

Cribl LogStream Documentation Manual

Version: v2.4.2

Generated: 2021-02-16 22:24:27

INTRODUCTION	12
Welcome	12
About Cribl LogStream	15
Basic Concepts	16
Getting Started Guide	19
DEPLOYMENT	38
Deployment Types	38
Single-Instance Deployment	39
Distributed Deployment	48
Bootstrap Workers from Master	63
Splunk App Deployment *	67
Kubernetes Master Deployment	71
Kubernetes Worker Deployment	80
Sizing and Scaling	89
Config Files	93
cribl.yml	95
inputs.yml	96
outputs.yml	97
licenses.yml	99
regexes.yml	100
breakers.yml	101
mappings.yml	102
instance.yml	103
Licensing	104
Access Management	110
Authentication	111
Local Users	117
Roles	119
Version Control	127
Persistent Queues	137
Securing	140
Monitoring	143
Internal Metrics	151
Upgrading	155
Uninstalling	160
WORKING WITH DATA	161

Event Model	161
Event Processing Order	163
Routes	165
Pipelines	170
Data Onboarding	175
Functions	179
Auto Timestamp	185
Aggregations	189
CEF Serializer	195
Clone	197
Comment	198
DNS Lookup	200
Drop	204
Dynamic Sampling	205
Eval	209
Flatten	212
GeoIP	214
Grok	217
JSON Unroll	219
Lookup	221
Mask	225
Numerify	231
Parser	233
Publish Metrics	241
Regex Extract	246
Redis	250
Regex Filter	253
Rename	254
Rollup Metrics	256
Sampling	258
Serialize	260
Suppress	263
Tee	266
Trim Timestamp	269
Unroll	270
XML Unroll	272
Prometheus Publisher (beta)	275
Reverse DNS (deprecated)	277

Sources		279
Splunk TCP		282
Splunk HEC		285
Syslog		288
Elasticsearch API		291
TCP JSON		295
TCP (RAW)		298
HTTP/S (Bulk API)		302
Raw HTTP/S		306
Kafka		309
Kinesis		313
Kinesis Firehose		316
Azure Event Hubs	;	319
Metrics		321
Prometheus		324
SQS		327
S3		331
Office 365 Service	es	338
Office 365 Activity	у	341
SNMP Trap		344
Datagens		346
Cribl Internal		348
AppScope		351
Collectors		354
Filesystem/NFS		358
S3		362
Script		366
REST / API Endpo	int	370
Scheduling and R	unning	378
Destinations		384
Output Router		388
Splunk Single Inst	tance	390
Splunk Load Bala	nced	394
Splunk HEC		402
S3 Compatible St	ores	406
Kinesis Streams		411
CloudWatch Logs		414
SQS		417

Filesystem/NFS	421
Elasticsearch	424
Honeycomb	427
TCP JSON	430
Syslog	433
Kafka	437
Azure Blob Storage	442
Azure Monitor Logs	445
Azure Event Hubs	448
Google Cloud Storage	451
StatsD	454
StatsD Extended	457
Graphite	460
SNMP Trap	463
InfluxDB	465
MinIO	468
New Relic	472
Wavefront	475
SignalFx	478
Sumo Logic	481
Datadog	484
DevNull	488
Default	489
Data Preview	490
Securing Data	495
Encryption	496
Decryption	501
Scripts	503
Using Datagens	505
CLI Reference	510
EXPRESSION REFERENCE	518
Introduction to Expression Syntax	518
Cribl Expressions	521
KNOWLEDGE	532
Regex Library	532
Grok Patterns Library	535
Event Breakers	537
Lookups Library	547

Parsers Library	551
Schema Library	553
Global Variables Library	555
USE CASES	557
Ingest-time Fields	557
Ingest-time Lookups	560
Sampling	564
Access Logs: Apache, ELB, CDN, S3, etc.	566
Firewall Logs: VPC Flow Logs, Cisco ASA, Etc.	569
Masking and Obfuscation	572
Lookups as Filters for Masks	576
Lookups and Regex Magic	579
Regex Filtering	588
Encrypting Sensitive Data	590
Syslog Data Reduction	595
Splunk to Elasticsearch	602
Reducing Windows XML Events	610
BEST PRACTICES	621
Managing Large Lookups	621
TROUBLESHOOTING	625
Diagnosing Issues	625
Known Issues	627
THIRD-PARTY SOFTWARE	633
Credits	633
Introduction	635
Welcome	635
About Cribl LogStream	638
Basic Concepts	639
Getting Started Guide	642
Deployment	661
Deployment Types	661
Single-Instance Deployment	662
Distributed Deployment	671
Bootstrap Workers from Master	686
Splunk App Deployment *	690
Kubernetes Master Deployment	694
Kubernetes Worker Deployment	703

Sizing and Scaling	/12
Config Files	716
cribl.yml	718
inputs.yml	719
outputs.yml	720
licenses.yml	722
regexes.yml	723
breakers.yml	724
mappings.yml	725
instance.yml	726
Licensing	727
Access Management	733
Authentication	734
Local Users	740
Roles	742
Version Control	750
Persistent Queues	760
Securing	763
Monitoring	766
Internal Metrics	774
Upgrading	778
Uninstalling	783
Working With Data	784
Event Model	784
Event Processing Order	786
Routes	788
Pipelines	793
Data Onboarding	798
Functions	802
Auto Timestamp	808
Aggregations	812
CEF Serializer	818
Clone	820
Comment	821
DNS Lookup	823
Drop	827
Dynamic Sampling	828
Eval	832

Flatten	835
GeoIP	837
Grok	840
JSON Unroll	842
Lookup	844
Mask	848
Numerify	854
Parser	856
Publish Metrics	864
Regex Extract	869
Redis	873
Regex Filter	876
Rename	877
Rollup Metrics	879
Sampling	881
Serialize	883
Suppress	886
Tee	889
Trim Timestamp	892
Unroll	893
XML Unroll	895
Prometheus Publisher (beta)	898
Reverse DNS (deprecated)	900
Sources	902
Splunk TCP	905
Splunk HEC	908
Syslog	911
Elasticsearch API	914
TCP JSON	918
TCP (RAW)	921
HTTP/S (Bulk API)	925
Raw HTTP/S	929
Kafka	932
Kinesis	936
Kinesis Firehose	939
Azure Event Hubs	942
Metrics	944
Prometheus	947

SQS	950
S3	954
Office 365 Services	961
Office 365 Activity	964
SNMP Trap	967
Datagens	969
Cribl Internal	971
AppScope	974
Collectors	977
Filesystem/NFS	981
S3	985
Script	989
REST / API Endpoint	993
Scheduling and Running	1001
Destinations	1007
Output Router	1011
Splunk Single Instance	1013
Splunk Load Balanced	1017
Splunk HEC	1025
S3 Compatible Stores	1029
Kinesis Streams	1034
CloudWatch Logs	1037
SQS	1040
Filesystem/NFS	1044
Elasticsearch	1047
Honeycomb	1050
TCP JSON	1053
Syslog	1056
Kafka	1060
Azure Blob Storage	1065
Azure Monitor Logs	1068
Azure Event Hubs	1071
Google Cloud Storage	1074
StatsD	1077
StatsD Extended	1080
Graphite	1083
SNMP Trap	1086
InfluxDB	1088

MinIO	1091
New Relic	1095
Wavefront	1098
SignalFx	1101
Sumo Logic	1104
Datadog	1107
DevNull	1111
Default	1112
Data Preview	1113
Securing Data	1118
Encryption	1119
Decryption	1124
Scripts	1126
Using Datagens	1128
CLI Reference	1133
Expression Reference	1141
Introduction to Expression Syntax	1141
Cribl Expressions	1144
Knowledge	1155
Regex Library	1155
Grok Patterns Library	1158
Event Breakers	1160
Lookups Library	1170
Parsers Library	1174
Schema Library	1176
Global Variables Library	1178
Use Cases	1180
Ingest-time Fields	1180
Ingest-time Lookups	1183
Sampling	1187
Access Logs: Apache, ELB, CDN, S3, etc.	1189
Firewall Logs: VPC Flow Logs, Cisco ASA, Etc.	1192
Masking and Obfuscation	1195
Lookups as Filters for Masks	1199
Lookups and Regex Magic	1202
Regex Filtering	1211
Encrypting Sensitive Data	1213
Syslog Data Reduction	1218

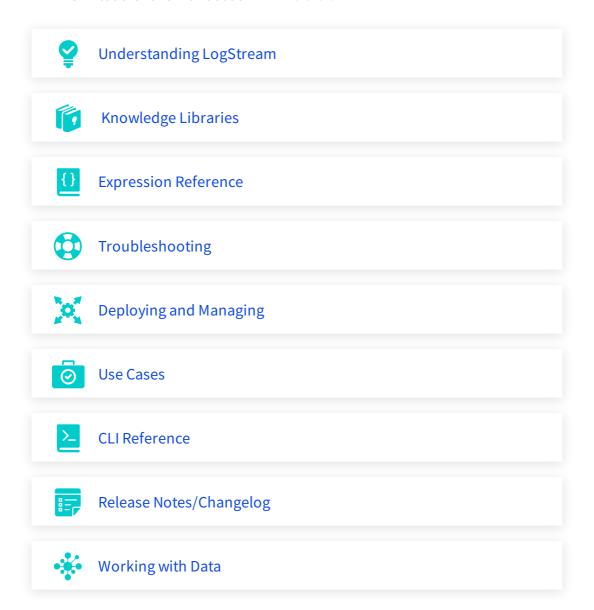
Splunk to Elasticsearch	1225
Reducing Windows XML Events	1233
Best Practices	1244
Managing Large Lookups	1244
Troubleshooting	1248
Diagnosing Issues	1248
Known Issues	1250
Third-Party Software	1256
Credits	1256

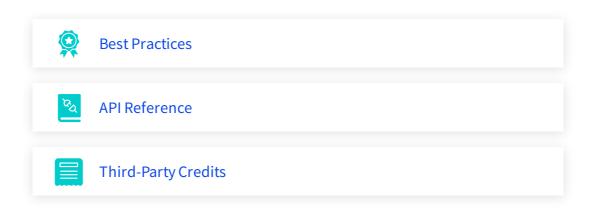
INTRODUCTION

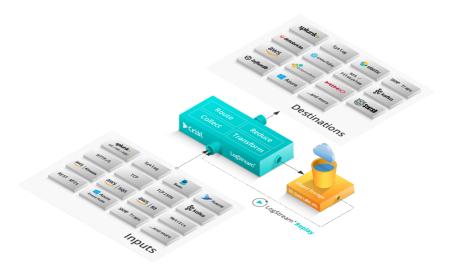
Welcome

Looking for Cribl LogStream documentation? We've goat you covered!

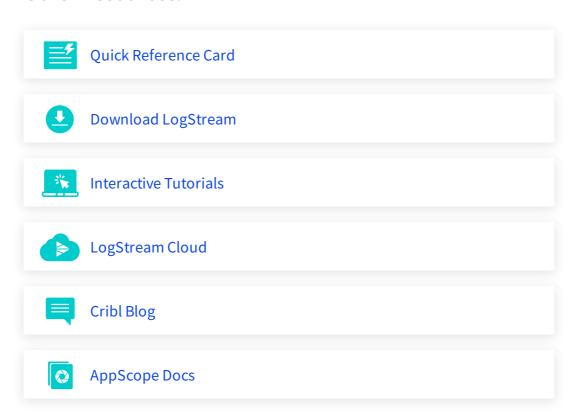
- Questions not answered here? We'd love to help you. Meet us in #Cribl Community Slack – sign up here.
- Download entire manual as PDF v.2.4.1.







Other Resources:



Page	14 0	f 1257
------	------	--------

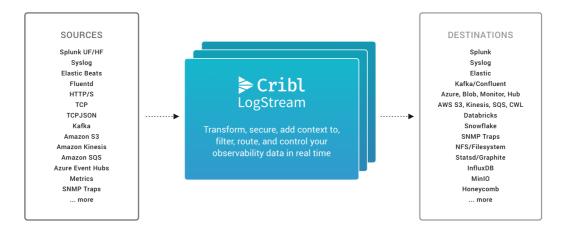
About Cribl LogStream

Getting started with Cribl LogStream

What Is Cribl LogStream?

Cribl LogStream helps you process machine data – logs, instrumentation data, application data, metrics, etc. – in real time, and deliver them to your analysis platform of choice. It allows you to:

- Add context to your data, by enriching it with information from external data sources.
- Help secure your data, by redacting, obfuscating, or encrypting sensitive fields
- Optimize your data, per your performance and cost requirements.



Cribl LogStream ships in a single, no-dependencies package. It provides a refreshing and modern interface for working with and transforming your data. It scales with – and works inline with – your existing infrastructure, and is transparent to your applications.

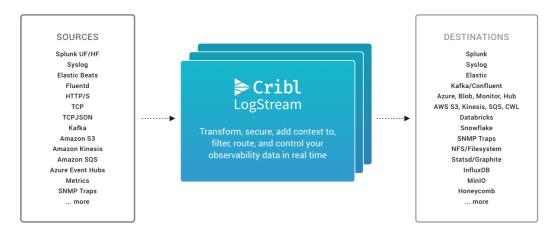
Who Is Cribl LogStream For?

Cribl LogStream is built for administrators, managers, and users of operational and security intelligence products and services.

Basic Concepts

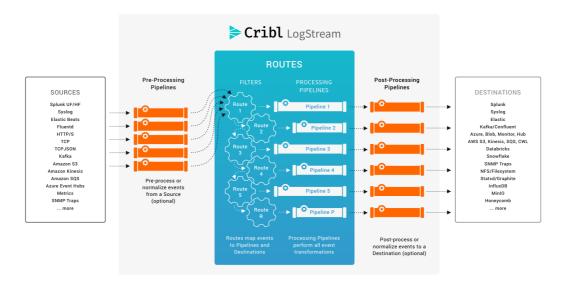
Notable features and concepts to get a fundamental understanding of Cribl LogStream

As we describe features and concepts, it helps to have a mental model of Cribl LogStream as a system that receives events from various sources, processes them, and then sends them to one or more destinations.



Sources, LogStream, destinations

Let's zoom in on the center of the above diagram, to take a closer look at the processing and transformation options that LogStream provides internally. The basic interface concepts to work with are Routes, which manage data flowing through Pipelines, which consist of Functions.



Routes, Pipelines, Functions

Routes

Routes evaluate incoming events against filter expressions to find the appropriate Pipeline to send them to. Routes are **evaluated in order**. A Route can be associated **with only one** Pipeline and one output. By default, a Route-Pipeline-Output tuple will consume matching events.

If the Route's Final flag is disabled, one or more event **clones** are sent down the associated Pipeline, while the original event continues down the rest of the Routes. This is very useful in cases where the same set of events needs to be processed in multiple ways and delivered to different destinations. For more details, see Routes.

Pipelines

A series of Functions is called a **Pipeline**, and the order in which the Functions are executed matters. Events are delivered to the beginning of a pipeline by a Route, and as they're processed by a Function, the events are passed to the next Function down the line.



Events only move forward – toward the end of the Pipeline, and eventually out of the system. For more details, see Pipelines.

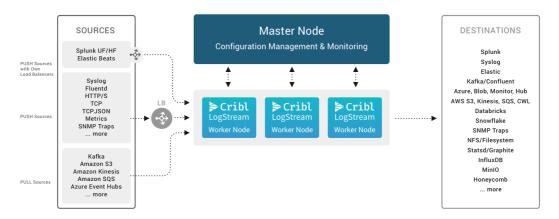
Functions

At its core, a **Function** is a piece of code that executes on an event, and that encapsulates the smallest amount of processing that can happen to that event. For instance, a very simple Function can be one that replaces the term foo with bar on each event. Another one can hash or encrypt bar. Yet another function can add a field – say, dc=jfk-42 – to any event with source=*us-nyc-application.log.

Functions process each event that passes through them. To help improve performance, functions can optionally be configured with filters, to limit their processing scope to matching events only. For more details, see Functions.

A Scalable Model

You can scale LogStream up to meet enterprise needs in a distributed deployment. Here, multiple LogStream Workers (instances) share the processing load. But as you can see in the preview schematic below, even complex deployments follow the same basic model outlined above.



Distributed deployment architecture

Getting Started Guide

This guide walks you through planning, installing, and configuring a single-instance deployment of Cribl LogStream. You'll capture some realistic sample log data, and then use LogStream's built-in Functions to redact, parse, refine, and shrink the data.

By the end of this guide, you'll have assembled all of LogStream's basic building blocks: a Source, Route, and Pipeline, several Functions, and a Destination. You can complete this tutorial using LogStream's included sample data, without connections to – or licenses on – any inbound or outbound services.

Assuming a cold start (from initial LogStream download and installation), this guide might take about an hour. But you can work through it in chunks, and LogStream will persist your work between sessions.

☐ If you've already downloaded, installed, and launched LogStream, skip ahead to Add a Source.

Requirements for this Tutorial

The minimum requirements for running this tutorial are the same as for a LogStream production single-instance deployment.

OS

• Linux: RedHat, CentOS, Ubuntu, Amazon Linux (64bit)

System

- +4 physical cores = +8 vCPUs; +8GB RAM all beyond your basic OS/VM requirements
- 5GB free disk space (more if persistent queuing is enabled)

i We assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). For details, see Recommended AWS, Azure, and GCP Instance Types.

Browser Support

• Firefox 65+, Chrome 70+, Safari 12+, Microsoft Edge

Network Ports

By default, LogStream listens on the following ports:

Component	Default Port
UI default	9000
HTTP Inbound, default	10080
User options	+ Other data ports as required.

You can override these defaults as needed.

Plan for Production

For higher processing volumes, users typically enable LogStream's Distributed Deployment option. While beyond the scope of this tutorial, that option has a few additional requirements, which we list here for planning purposes:

- Port 9000 or 4200 must be available on the Master Node for Workers' communications.
- Git (1.8.3.1 or higher) must be installed on the Master Node, to manage configuration changes.

See Sizing and Scaling for further details about configuring LogStream to handle large data streams.

Download and Install LogStream

Download the latest version of LogStream at https://cribl.io/download/.

Un-tar the resulting .tgz file in a directory of your choice (e.g., /opt/). Here's general syntax, and a specific example:

```
tar xvzf cribl~version><build><arch>.tgz
tar xvzf cribl-2.3.1-1d4e05c5-linux-x64.tgz
```

You'll now have LogStream installed in a cribl subdirectory, by default: /opt/cribl/. We'll refer to this cribl subdirectory throughout this documentation as \$CRIBL_HOME.

Run LogStream

In your terminal, switch to the \$CRIBL_HOME/bin directory (e.g,: /opt/cribl/bin). Here, you can start, top, and verify the LogStream server using these basic ./cribl CLI commands:

Start: ./cribl startStop: ./cribl stop

• Get status: ./cribl status

☐ For other available commands, see CLI Reference.

Next, in your browser, open http://<hostname>:9000 (e.g.,
http://localhost:9000) and log in with default credentials (admin ,
admin).

Register your copy of LogStream if desired.

After registering, you'll be prompted to change the default password.

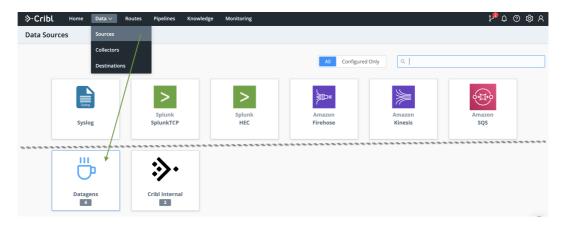
You are now ready to configure a working LogStream installation – with a Source, Destination, Pipeline, and Route – and to assemble several built-in Functions to refine sample log data.

Get Data Flowing

Add a Source

Each LogStream Source represents a data input. Options include Splunk, Elastic Beats, Kinesis, Kafka, syslog, HTTP, TCP JSON, and others.

For this tutorial, we'll enable a LogStream built-in datagen (i.e., data generator) that generates a stream of realistic sample log data.



Addiing a datagen Source

- 1. From LogStream's top menu, select **Data > Sources**.
- 2. From the **Data Sources** page's tiles or left menu, select **Datagens**.

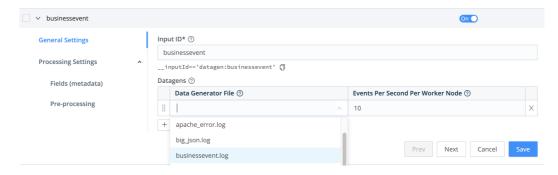
(You can use the search box to jump to the **Datagens** tile.)

- 3. Click Add New to open the New Datagen source pane.
- 4. In the Input ID field, name this Source businessevent.
- 5. In the Data Generator File drop-down, select businessevent.log.

This generates...log events for a business scenario. We'll look at their structure shortly, in Capture and Filter Sample Data.

6. Click Save.

The **On** slider in the **Enabled** column indicates that your datagen Source has started generating sample data.



Configuring a datagen Source

Add a Destination

Each LogStream Destination represents a data output. Options include Splunk, Kafka, Kinesis, InfluxDB, Snowflake, Databricks, TCP JSON, and others.

For this tutorial, we'll use LogStream's built-in **DevNull** Destination. This simply discards events – not very exciting! But it simulates a real output, so it provides a configuration-free quick start for testing LogStream setups. It's ideal for our purposes.

To verify that **DevNull** is enabled, let's walk through setting up a Destination, then setting it up as LogStream's default output:

- 1. From LogStream's top menu, select **Data > Destinations**.
- 2. Select **DevNull** from the **Data Destinations** page's tiles or left menu.

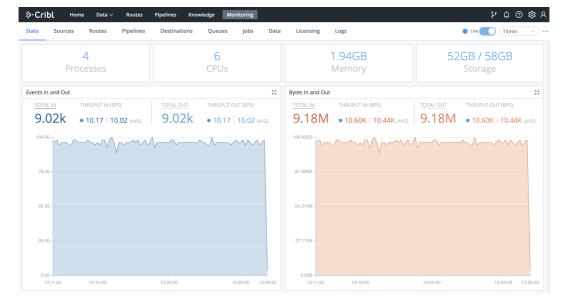
(You can use the search box to jump to the **DevNull** tile.)

- 3. On the resulting **devnull** row, look for the **Live** indicator under **Enabled**. This confirms that the **DevNull** Destination is ready to accept events.
- 4. From the **Data Destinations** page's left nav, select the **Default** Destination at the top.
- On the resulting Manage Default Destination page, verify that the Default Output ID drop-down points to the devnull Destination we just examined.

We've now set up data flow on both sides. Is data flowing? Let's check.

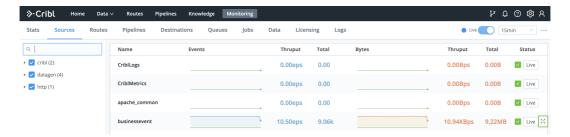
Monitor Data Throughput

From the top menu, select **Monitoring**. This opens a summary dashboard, where you should see a steady flow of data in and out of LogStream. The left graph shows events in/out. The right graph shows bytes in/out.



Monitoring dashboard

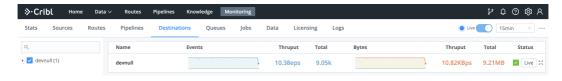
Monitoring displays data from the preceding 24 hours. You can use the Monitoring submenu to open detailed displays of LogStream components, collection jobs and tasks, and LogStream's own internallogs. Click Sources on the lower (white) submenu to switch to this view:



Monitoring Sources

This is a compact display of each Source's inbound events and bytes as a sparkline. You can click each Source's Expand button (highlighted at right) to zoom up detailed graphs.

Click **Destinations** on the lower submenu. This displays a similar sparklines view, where you can confirm data flow out to the devnull Destination:



Monitoring Destinations

With confidence that we've got data flowing, let's send it through a LogStream Pipeline, where we can add Functions to refine the raw data.

Create a Pipeline

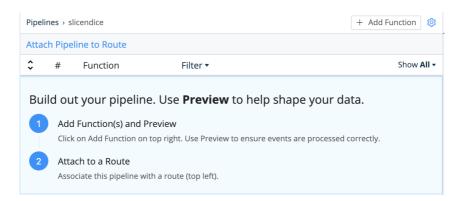
A Pipeline is a stack of LogStream Functions that process data. Pipelines are central to refining your data, and also provide a central LogStream workspace – so let's get one going.

1. From the top menu, select Pipelines.

You now have a two-pane view, with business on the left and party on the right a Pipelines list on the left and Sample Data controls on the right. (We'll capture some sample data momentarily.)

- 2. At the **Pipelines** pane's upper right, click **+ Add Pipeline**, then select **Create Pipeline**.
- 3. In the new Pipeline's **ID** field, enter a unique identifier. (For this tutorial, you might use slicendice.)
- 4. Optionally, enter a **Description** of this Pipeline's purpose.
- 5. Click Save.

Your empty Pipeline now prompts you to preview data, add Functions, and attach a Route. So let's capture some data to preview.



Pipeline prompt to add Functions

Capture and Filter Sample Data

The right **Sample Data** pane provides multiple tools for grabbing data from multiple places (inbound streams, copy/paste, and uploaded files); for

previewing and testing data transformations as you build them; and for saving and reloading sample files.

Since we've already got live (simulated) data flowing in from the datagen Source we built, let's grab some of that data.

Capture New Data

- 1. In the right pane, click Capture New.
- 2. In the **Capture Sample Data** modal, immediately change the generated **File Name** to a name you'll recognize, like be_raw.log.
- 3. Click Capture, then accept the drop-down's defaults click Start.
- 4. Click **Save as Sample File**. This saves to the **File Name** you entered above. You're now previewing the events in the right pane. (Note that this pane's **Preview Simple** tab now has focus.)
- 5. Click **Show more** to expand one or more events.

By skimming the key-value pairs within the data's _raw fields, you'll notice the scenario underlying this preview data (provided by the businessevents.log datagen): these are business logs from a mobile-phone provider.

To set up our next step, find at least one marketState K=V pair. Having captured and examined this raw data, let's use this K=V pair to crack open LogStream's most basic data-transformation tool, Filtering.

Filter Data and Manage Sample Files

- 1. Click the right pane's **Sample Data** tab.
- 2. Again click Capture New.
- 3. In the **Capture Sample Data** modal, replace the **Filter Expression** field's default true value with this simple regex:

```
_raw.match(/marketState=TX/)
```

We're going to Texas! If you type this in, rather than pasting it, notice how LogStream provides typeahead assist to complete a well-formed JavaScript expression.

You can also click the Expand button at the Filter Expression field's right

edge to open a modal to validate your expression. The adjacent dropdown enables you to restore previously used expressions



4. Click Capture, then Start.

Using the **Capture** drop-down's default limits of 10 seconds and 10 events, you'll notice that with this filter applied, it takes much longer for LogStream to capture 10 matching events.

- 5. Click **Cancel** to discard this filtered data and close the modal.
- 6. On the right pane's **Sample Data** tab, click **Simple** beside be_raw.log.

This restores our preview of our original, unfiltered capture. We're ready to transform this sample data in more interesting ways, by building out our Pipeline's Functions.

Refine Data with Functions

Functions are pieces of JavaScript code that LogStream invokes on each event that passes through them. By default, this means all events – each Function has a **Filter** field whose value defaults to true. As we just saw with data capture, you can replace this value with an expression that scopes the Function down to particular matching events.

In this Pipeline, we'll use some of LogStream's core Functions to:

- Redact (mask) sensitive data
- Extract (parse) the _raw field's key-value pairs as separate fields.
- Add a new field.
- Delete the original _raw field, now that we've extracted its contents.
- Rename a field for better legibility..

Mask: Redact Sensitive Data

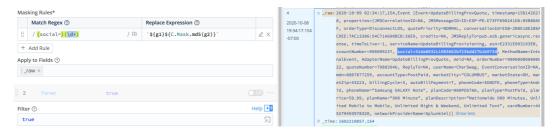
In the right **Preview** pane, notice each that event includes a **social** key, whose value is a (fictitious) raw Social Security number. Before this data goes any further through our Pipeline, let's use LogStream's Mask Function to swap in an md5 hash of each SSN.

- 1. In the left Pipelines pane, click + Add Function.
- 2. Search for Mask, then click it.
- 3. In the new Function's Masking Rules, click the into Match Regex field.
- 4. Enter or paste this regex, which simply looks for digits following social = : (social =)(\d+)
- 5. In **Replace Expression**, paste the following hash function. The backticks are literal: `\${g1}\${C.Mask.md5(g2)}`
- 6. Note that **Apply to Fields** defaults to _raw . This is what we want to target, so we'll accept this default.
- 7. Click Save.

You'll immediately notice some obvious changes:

- The **Preview** pane has switched from its **IN** to its **OUT** tab, to show you the outbound effect of the Pipeline you just saved.
- Each event's _raw field has changed color, to indicate that it's undergone some redactions.

Now locate at least one event's **Show more** link, and click to expand it. You can verify that the social values have now been hashed.



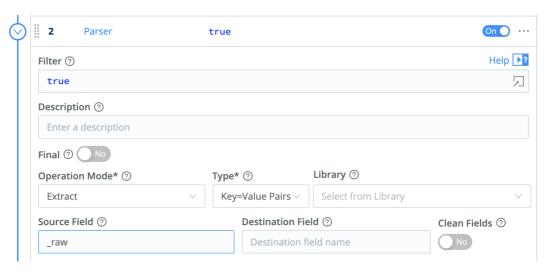
Mask Function and hashed result

Parser: Extract Events

Having redacted sensitive data, we'll next use a Parser function to lift up all the _raw field's key-value pairs as fields:

- 1. In the left Pipelines pane, click + Add Function.
- 2. Search for Parser, then click it.

- 3. Leave the Operation Mode set to its Extract default.
- 4. Set the Type to Key=Value Pairs.
- 5. Leave the **Source Field** set to its _raw default.
- 6. Click Save.



Parser configured to extract K=V pairs from _raw

You should see the **Preview** pane instantly light up with a lot more fields, parsed from _raw . You now have rich structured data, but not all of this data is particularly interesting: Note how many fields have NA ("Not Applicable") values. We can enhance the **Parser** Function to ignore fields with NA values.

 In the Function's Fields Filter Expression field (near the bottom), enter this negation expression: value≠'NA'

Note the single-quoted value. If you type (rather than paste) this expression, watch how typeahead matches the first quote you type.

2. Click **Save**, and watch the **Preview** pane.



Filtering the Parser Function to ignore fields with 'NA' values

Several fields should disappear – such as credits, EventConversationID, and ReplyTo. The remaining fields should display meaningful values.

Congratulations! Your log data is already starting to look better-organized and less bloated.

□ Missed It?

If you didn't see the fields change, slide the Parser Function **Off**, click **Save** below, and watch the **Preview** pane change. Using these toggles, you can preserve structure as you test and troubleshoot each Function's effect.

Note that each Function also has a **Final** toggle, defaulting to **Off**. Enabling **Final** anywhere in the Functions stack will prevent data from flowing to any Functions lower in the UI.

Be sure to toggle the Function back **On**, and click **Save** again, before you proceed!



Toggling a Function off and on

Next, let's add an extra field, and conditionally infer its value from existing values. We'll also remove the _raw field, now that it's redundant. To add and remove fields, the **Eval** Function is our pal.

Eval: Add and Remove Fields

Let's assume we want to enrich our data by identifying the manufacturer of a certain popular phone handset. We can infer this from the existing phoneType field that we've lifted up for each event.

Add Field (Enrich)

- 1. In the left **Pipelines** pane, click + Add Function .
- 2. Search for Eval, then click it.

3. Click into the new Function's Evaluate Fields table.

Here you add new fields to events, defining each field as a key-value pair. If we needed more key-value pairs, we could click + Add Field for more rows.

- 4. In Name, enter: phoneCompany.
- 5. In Value Expression, enter this JS ternary expression that tests
 phoneType 's value:
 phoneType.startsWith('iPhone') ? 'Apple' : 'Other' (Note the ?
 and : operators, and the single-quoted values.)
- 6. Click **Save**. Examine some events in the **Preview** pane, and each should now contain a phoneCompany field that matches its phoneType.



Adding a field to enrich data

Remove Field (Shrink Data)

Now that we've parsed out all of the _raw field's data – it can go. Deleting a (large) redundant field will give us cleaner events, and reduced load on downstream resources.

- 1. Still in the **Eval** Function, click into **Remove Fields**
- 2. Type: _raw and press **Tab** or **Enter**.
- 3. Click Save.

The **Preview** pane's diff view should now show each event's _raw field stripped out.



Removing a field to streamline data

Our log data has now been cleansed, structured, enriched, and slimmed-down. Let's next look at how to make it more legible, by giving fields simpler names.

Rename: Refine Field Names

1. In the left **Pipelines** pane, click + Add Function.

This rhythm should now be familiar to you.

- 2. Search for Rename, then click it.
- 3. Click into the new Function's **Rename Fields** table.

This has the same structure you saw above in Eval: Each row defines a key-value pair.

- 4. In **Current Name**, enter the longhaired existing field name: conversationId.
- 5. In New Name, enter the simplified field name: ID.
- 6. Watch any event's conversationId field in the **Preview** pane as you click **Save** at left. This field should change to ID in all events.

Drop: Remove Unneeded Events

We've already refined our data substantially. To further slim it down, a Pipeline can entirely remove events that aren't of interest for a particular downstream service.

As the "Pipeline" name implies, your LogStream installation can have multiple Pipelines, each configured to send out a data stream tailored to a particular Destination. This helps you get the right data in the right places most efficiently.

Here, let's drop all events for customers who use prepaid monthly phone service (i.e., **not** postpaid):

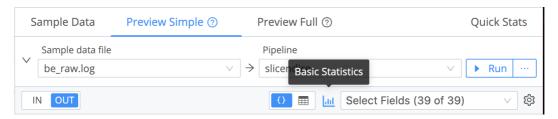
- 1. In the left **Pipelines** pane, click + Add Function.
- 2. Search for Drop, then click it.
- 3. Click into the new Function's Filter field.
- 4. Replace the default true value with this JS negation expression: accountType≠'PostPaid'
- 5. Click Save.

Now scroll through the right **Preview** pane. Depending on your data sample, you should now see multiple events struck out and faded – indicating that LogStream will drop them before forwarding the data.

A Second Look at Our Data

Torture the data enough, and it will confess. By what factor have our transformations refined our data's volume? Let's check.

In the right **Preview** pane, click the **Basic Statistics** button:



Displaying Basic Statistics

Even without the removal of the _raw field (back in Eval) and the dropped events, you should see a substantial % reduction in the Full Event Length.



Data reduction quantified

Woo hoo! Before we wrap up our configuration: If you're curious about individual Functions' independent contribution to the data reduction shown

here, you can test it now. Use the toggle **Off > Save > Basic Statistics** sequence to check various changes.

Add and Attach a Route

We've now built a complete, functional Pipeline. But so far, we've tested its effects only on the static data sample we captured earlier. To get dynamic data flowing through a Pipeline, we need to filter that data in, by defining a LogStream Route.

1. At the **Pipelines** page's top left, click **Attach Pipeline to Route**.

This displays the **Routes** page. It's structured very similarly to the **Pipelines** page, so the rhythm here should feel familiar.

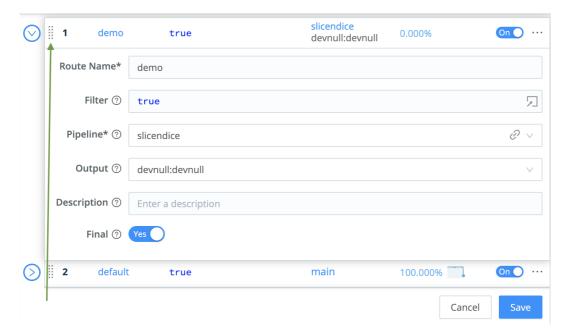
- 2. Click + Add Route.
- 3. Enter a unique, meaningful Route Name, like demo.
- 4. Leave the **Filter** field set to its true default, allowing it to deliver all events.

Because a Route delivers events to a Pipeline, it offers a first stage of filtering. In production, you'd typically configure each Route to filter events by appropriate source, sourcetype, index, host, _time, or other characteristics. The **Filter** field accepts JavaScript expressions, including AND (&&) and OR (||) operators.

- 5. Set the **Pipeline** drop-down to our configured slicendice Pipeline.
- 6. Set the Output drop-down to either devnull or default.

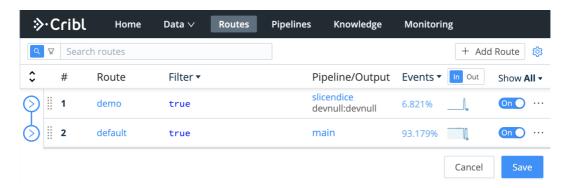
This doesn't matter, because we've set default as a pointer to devnull. In production, you'd set this carefully.

- 7. You can leave the **Description** empty, and leave **Final** set to **Yes**.
- 8. Grab the new Route by its left handle, and drag it above the default Route, so that our new Route will process events first. You should see something like the screenshot below.
- 9. Click **Save** to save the new Route to the routing table.



Configuring and adding a Route

The sparklines should immediately confirm that data is flowing through your new Route:



Live Routes

To confirm data flow through the whole system we've built, select **Monitoring > Routes** from LogStream's top menu and examine demo.



Monitoring data flow through Routes

Also select Monitoring > Pipelines and examine slicendice.



Monitoring data flow through Pipelines

What Have We Done?

Look at you! Give yourself a pat on the back! In this short, scenic tour – with no hit to your cloud-services charges – you've build a simple but complete LogStream system, exercising all of its basic components:

- Downloaded, installed, and run LogStream.
- Configured a Source to hook up an input.
- Configured a Destination to feed an output.
- Monitored data throughput, and checked it twice.
- Built a Pipeline.
- Configured LogStream Functions to redact, parse, enrich, trim, rename, and drop event data.
- Added and attached a Route to get data flowing through our Pipeline.

Next Steps

Interested in guided walk-throughs of more-advanced LogStream features? We suggest that next, you check out:

- LogStream Sandboxes: Work through general and specific scenarios in containers. with terminal access and free, hosted data inputs and outputs.
- Use Cases documentation: Bring your own services to build solutions to specific challenges.
- Cribl Concept: Pipelines Video showing how to build and use Pipelines at multiple LogStream stages.
- Cribl Concept: Routing Video about using Routes to send different data through different paths.

Cleaning Up

Oh yeah, you've still got the LogStream server running, with its businessevent.log datagen wtill firing events. If you'd like to shut these down for now, in reverse order:

- 1. Go to Data > Sources > Datagens.
- 2. Slide businessevent to **Off**, and click **Save**. (Refer back to the screenshot above.)
- 3. In your terminal's \$CRIBL_HOME/bin directory, shut down the server with: ./cribl stop

That's it! Enjoy using LogStream.

DEPLOYMENT

Deployment Types

Deployment guide to get you started with Cribl

There are at least **two** key factors that will determine the type of Cribl LogStream deployment in your environment:

- Amount of Incoming Data: This is defined as the amount of data planned to be ingested per unit of time. E.g. How many MB/s or GB/day?
- Amount of Data Processing: This is defined as the amount of processing that will happen on incoming data. E.g., is most data passing through and just being routed? Or are there a lot of transformations, regex extractions, field encryptions? Is there a need for heavy re-serialization?

Single Instance Deployment

When volume is low and/or amount of processing is light, you can get started with a single instance deployment.

Distributed Deployment

To accommodate increased load, we recommend scaling up and perhaps out with multiple instances.

Splunk App Deployment

If you have an existing Splunk Heavy Forwarder infrastructure that you want to use, you can deploy Cribl App for Splunk. See the note below before you plan.

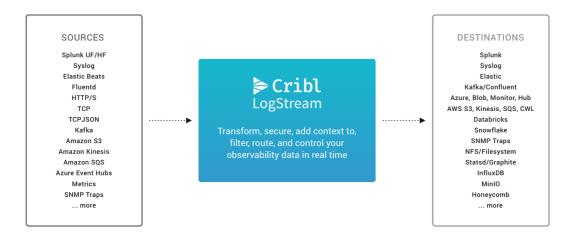
△ Cribl App for Splunk Deprecation Notice Click here.

Single-Instance Deployment

Getting started with Cribl LogStream on a single instance

For small-volume or light processing environments – or for test or evaluation use cases – a single instance of Cribl LogStream might be sufficient to serve all inputs, event processing, and outputs. This page outlines how to implement a single-instance deployment.

Architecture



Requirements

- OS:
 - Linux: Red Hat, CentOS, Ubuntu, Amazon Linux (64bit)
 - Mac OS is no longer supported as of v. 2.3, due to LogStream's incorporation of Linux-native components.
- System:
 - +4 physical cores, +8GB RAM
 - 5GB free disk space (more if persistent queuing is enabled)
 - i We assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). All quantities listed above are minimum requirements. To fulfill the above

requirements using cloud-based virtual machines, see Recommended AWS, Azure, and GCP Instance Types.

• Browser Support: Firefox 65+, Chrome 70+, Safari 12+, Microsoft Edge

Network Ports

By default, LogStream listens on the following ports:

Component	Default Port
UI	9000
HTTP In	10080
Splunk to Cribl LogStream data port	localhost:10000 (Cribl App for Splunk)
criblstream Splunk search command to Cribl LogStream	localhost:10420 (Cribl App for Splunk)
User options	+ Other data ports as required.

Overriding Default Ports

The above ports can be overridden in the following configuration files:

- Cribl UI port (9000): Default definitions for host, port, and other settings are set in \$CRIBL_HOME/default/cribl/cribl.yml, and can be overridden by defining alternatives in \$CRIBL_HOME/local/cribl/cribl.yml.
- Data Ports: HTTP In (10080), TCPJSON in (10420) Splunk to Cribl (10000): Default definitions for host, port and other settings are set in \$CRIBL_HOME/default/cribl/inputs.yml, and can be overridden by defining alternatives in \$CRIBL_HOME/local/cribl/inputs.yml.

Installing on Linux

- Install the package on your instance of choice. Download it here.
- Ensure that required ports are available (see Network Ports).
- Un-tar in a directory of choice, e.g., /opt/:
 - tar xvzf cribl ← version ← < build ← < arch > .tgz

Running

Go to the \$CRIBL_HOME/bin directory, where the package was extracted (e.g.: /opt/cribl/bin). Here, you can use ./cribl to:

- Start: ./cribl start
- Stop: ./cribl stop
- Reload: ./cribl reload
- Restart: ./cribl restart
- Get status: ./cribl status
- Switch a distributed deployment to single-instance mode:

```
./cribl mode-single (uses the default address:port 0.0.0.0:9000)
```

i Executing the restart or stop command cancels any currently running collection jobs. For other available commands, see CLI Reference.

Next, go to http://<hostname>:9000 and log in with default credentials (admin:admin). You can now start configuring Cribl LogStream with Sources and Destinations, or start creating Routes and Pipelines.

i In the case of an API port conflict, the process will retry binding for 10 minutes before exiting.

Enabling Start on Boot

Cribl LogStream ships with a CLI utility that can update your system's configuration to start LogStream at system boot time. The basic format to invoke this utility is:

[sudo] \$CRIBL_HOME/bin/cribl boot-start [enable|disable] [options] [args]

i Boot-start is supported only on Linux. For options and arguments, see the CLI Reference.

Newer systems use systemd to start processes at boot, while older ones use initd.

Using systemd

To **enable** Cribl LogStream to start at boot time with **systemd**, you need to run the boot-start command. Make sure you first create any user you want to specify to run LogStream. E.g., to run LogStream on boot as existing user cribl, you'd use:

```
sudo $CRIBL_HOME/bin/cribl boot-start enable -m systemd -u cribl
```

This will install a unit file (as below) and start Cribl LogStream at boot time as user cribl. A -configDir option can be used to specify where to install the unit file. If not specified, this location defaults to /etc/systemd/system.

If necessary, change ownership for the Cribl LogStream installation:

```
[sudo] chown -R cribl $CRIBL_HOME
```

Next, use the enable command to ensure that the service starts on system boot:

```
[sudo] systemctl enable cribl
```

To **disable** starting at boot time, run the following command:

```
sudo $CRIBL_HOME/bin/cribl boot-start disable
```

Note the file's default 65536 limit on maximum open files. Cribl recommends going no higher than this value.

```
Installed systemd File
```

```
[Unit]
Description=Systemd service file for Cribl LogStream.
After=network.target
```

[Service]
Type=forking
User=cribl
Restart=on-failure
RestartSec=5
LimitNOFILE=65536
PIDFile=/install/path/to/cribl/pid/cribl.pid
ExecStart=/install/path/to/cribl/bin/cribl start
ExecStop=/install/path/to/cribl/bin/cribl stop
ExecStopPost='/bin/rm -f /install/path/to/cribl/pid/cribl.pid'
ExecReload=/install/path/to/cribl/bin/cribl reload
TimeoutSec=60

```
[Install]
WantedBy=multi-user.target
```

Using initd

To **enable** Cribl LogStream to start at boot time with **initd**, you need to run the boot-start command. If the user that you want to run LogStreams does not exist, create it prior to executing. E.g., running LogStream as user cribl on boot:

```
sudo $CRIBL_HOME/bin/cribl boot-start enable -m initd -u cribl
```

This will install an init.d script in /etc/init.d/cribl.init.d, and start Cribl LogStream at boot time as user cribl.A -configDir option can be used to specify where to install the script. If not specified, this location defaults to /etc/init.d.

If necessary, change ownership for the Cribl LogStream installation:

```
[sudo] chown -R cribl $CRIBL_HOME
```

To **disable** starting at boot time, run the following command:

```
sudo $CRIBL_HOME/bin/cribl boot-start disable
```

△ Do NOT Run LogStream as Root!

If LogStream is required to listen on ports 1–1024, it will need privileged access. On a Linux system with POSIX capabilities, you can achieve this by adding the CAP_NET_BIND_SERVICE capability. For example: # setcap cap_net_bind_service=+ep \$CRIBL_HOME/bin/cribl

On some OS versions (such as CentOS), you must add an -i switch to the setcap command. For example: # setcap -i cap_net_bind_service=+ep \$CRIBL_HOME/bin/cribl

Upon starting the LogStream server, a Port xxx is already in use error might indicate that setcap did not successfully execute.

System Proxy Configuration

You can direct all outbound HTTP/S requests to go through proxy servers. Initial configuration and changing these variables requires restarting LogStream on the affected nodes if the application is already running when the changes are

applied. You do so by setting a few environment variables before starting LogStream, as follows:

Configure the HTTP_PROXY and HTTPS_PROXY environment variables either with your proxy's IP address, or with a DNS name that resolves to that IP address. Optionally, follow either convention with a colon and the port number to which you want to send queries.

```
HTTP_PROXY examples:

$ export HTTP_PROXY=http://10.15.20.25:1234
$ export HTTP_PROXY=http://proxy.example.com:1234

HTTPS_PROXY examples:
```

```
$ export HTTPS_PROXY=http://10.15.20.25:5678
$ export HTTPS_PROXY=http://proxy.example.com:5678
```

i Case Conflicts

The environment variables' names can be either uppercase or lowercase. However, if you set duplicate versions of the same name, the lowercase version takes precedence. E.g., if you've set both HTTPS_PROXY and https_proxy, the IP address specified in https_proxy will take effect.

Proxy Confguration with systemd

If you are proxying outbound traffic with systemd, list your proxy environment variables in the systemd unit file's [Service] section by adding statements of this form:

```
Installed systemd File

[Service]
...
Environment=https_proxy=<yourproxy>
Environment=https_proxy=http://proxy.example.com:1234
Environment=https_proxy=http://10.10.1.1:8080
```

This will prevent LogStream from throwing "failed to send anonymized telemetry metadata" errors.

Authenticating on Proxies

You can use HTTP Basic authentication on HTTP or HTTPS proxies. Specify the user name and password in the proxy URL. For example:

```
$ export HTTP_PROXY=http://username:password@proxy.example.com:1234
$ export HTTPS_PROXY=http://username:password@proxy.example.com:5678
```

Bypassing Proxies with NO_PROXY

If you've set the above environment variables, you can negate them for specified (or all) hosts. Set the NO_PROXY environment variable to identify URLs that should bypass the proxy server, and instead be sent as direct requests. Use the following format:

```
$ export NO_PROXY="<list of hosts/domains>"
```

Usage notes:

- Within the list, separate the host/domain names with commas or spaces.
- Optionally, each host/domain entry can be followed by a port. If specified, the port must match. If not specified, the protocol's default port is assumed.
- If specified, subdomain names must match. E.g.,
 NO_PROXY=foo.example.com will send requests directly to
 https://foo.example.com, but https://bar.example.com requests will go
 through the proxy.
- You can use leading wildcards like NO_PROXY="*.us, .org".
- NO_PROXY="*" disables all proxies.
- NO_PROXY with an empty list disables no proxies.

Where Proxies Apply

Proxy configuration is relevant to the following LogStream components that make outbound HTTP/S requests:

Destinations

S3 Compatible Stores

- AWS Kinesis Streams
- AWS CloudWatch Logs
- AWS SQS
- Azure Blob Storage
- Azure Event Hubs
- Azure Monitor Logs
- Elasticsearch
- Honeycomb
- Splunk HEC

Sources

- AWS Kinesis Streams
- AWS SQS
- AWS S3
- Azure Event Hubs

Collectors

S3 Collector

Proxying Multiple LogStream Instances in One Browser

LogStream stores authentication tokens based on each http header's URI scheme, host, and port information. Within a given browser, LogStream enforces a same-origin policy to isolate instances.

This means that if you want to run multiple proxied LogStream instances in one browser session, you must assign them different URI schemes, hosts, and/or ports. Otherwise, logging into an extra LogStream instance will expire the prior instance's session and log it out.

For example, assume that you've set up this pair of Apache proxy forward rules:

- https://web/cribla forwards to cribl_hosta:8001/cribla.
- https://web/criblb forwards to cribl_hostb:8001/criblb.

These two proxied addresses cannot be run simultaneously in the same browser session. However, this pair – which lead with separate URI schemes – could:

- https://web/cribla forwards to cribl_hosta:8001/cribla.
- https://web2/criblb forwards to cribl_hostb:8001/criblb.

Where separate instances **must** share URI formats, a workaround is to open the second instance in an incognito/private browsing window, or in a completely different browser.

Scaling Up

A single-instance installation can be configured to scale up and utilize as many resources on the host as required. See Sizing and Scaling for details.

Distributed Deployment

Getting started with Cribl LogStream on a distributed deployment

Distributed Deployment

To sustain higher incoming data volumes, and/or increased processing, you can scale from a single instance up to a multi-instance, distributed deployment. Instances in the deployment independently serve all inputs, process events, and send to outputs. The instances are managed centrally by a single Master Node, which is responsible for keeping configurations in sync, and for tracking and monitoring the instances' activity metrics.

Concepts

Single Instance – a single Cribl LogStream instance, running as a standalone (not distributed) installation on one server.

Master Node – a LogStream instance running in **Master** mode, used to centrally author configurations and monitor Worker Nodes in a distributed deployment.

Worker Node – a LogStream instance running as a **managed Worker**, whose configuration is fully managed by a Master Node. (By default, will poll the master for configuration changes every 10 seconds.)

Worker Group – a collection of Worker Nodes that share the same configuration. You map Nodes to a Worker Group using a Mapping Ruleset.

Worker Process – a Linux process within a Single Instance, or within Worker Nodes, that handles data inputs, processing, and output. The process count is constrained by the number of physical or virtual CPUs available; for details, see Sizing and Scaling.

Mapping Ruleset – an ordered list of filters, used to map Workers Nodes into Worker Groups.

△ A Worker Node's local running config can be manually overridden/changed, but changes won't persist on the filesystem. To

permanently modify a Worker Node's config, save, commit, and deploy it from the Master. See Deploying Configurations below.

LogStream 2.4 introduces role-based access control, at the Worker Group level. Users will be able to access Workers only within those Worker Groups on which they've been granted access.

Aggregating Workers

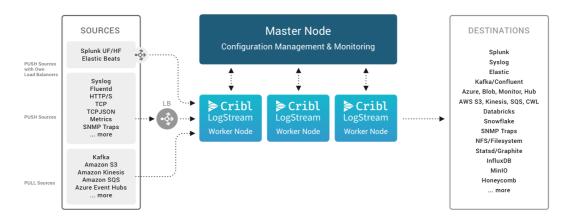
To clarify how the above concepts add up hierarchically, let's use a military metaphor involving toy soldiers:

- Worker Process = soldier.
- Worker Node = multiple Worker Processes = squad.
- Worker Group = multiple Worker Nodes = platoon.

Multiple Worker Groups are very useful in meeting organizational or geographic constraints reflected in configuration. E.g., you might have a U.S. Worker Group with certain TLS certificates and output settings, versus and APAC Worker Group and an EMEA Worker Group that each have distinct certs and settings.

Architecture

This is an overview of a distributed LogStream deployment's components.



Distributed deployment architecture

Master Node Requirements

OS:

- Linux: RedHat, CentOS, Ubuntu, AWS Linux (64bit)
- System:
 - +4 physical cores, +8GB RAM
 - 5GB free disk space
- Git: git must be available on the Master Node. See details below.
- Browser Support: Firefox 65+, Chrome 70+, Safari 12+, Microsoft Edge
 - i We assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). All quantities listed above are minimum requirements.
 - Mac OS is no longer supported as of v. 2.3, due to LogStream's incorporation of Linux-native features.

Worker Node Requirements

See Single-Instance Deployment for requirements and Sizing and Scaling for capacity planning details.

Network Ports - Master Node

In a distributed deployment, Workers communicate with the Master Node on these ports. Ensure that the Master is reachable on those ports from **all** Workers.

Component	Default Port
Heartbeat	4200

Network Ports - Worker Nodes

By default, all LogStream Worker instances listen on the following ports:

Component	Default Port
UI	9000

Installing on Linux

See Single-Instance Deployment, as the installation procedures are identical.

Version Control with git

LogStream requires git (version 1.8.3.1 or higher) to be available locally on the host where the Master Node will run. **Configuration changes must be committed to git before they're deployed.**

If you don't have git installed, check here for details on how to get started.

The Master node uses git to:

- Manage configuration versions across Worker Groups.
- Provide users with an audit trail of all configuration changes.
- Allow users to display diffs between current and previous config versions.

Setting up Master and Worker Nodes

1. Configuring a Master Node

You can configure a Master Node through the UI, through the instance.yml config file, or through the command line.

Using the UI

In Settings > Distributed Settings > Distributed Management > General Settings, select Mode Master. Supply the required Master settings (Address and Port). Customize the optional settings if desired. Then click Save to restart.

☐ Worker UI Access

If you enable the nearby **Distributed Settings > Master Settings > Worker UI access** option (enabledWorkerRemoteAccess key), you will be able to click through from the Master's **Manage Worker Nodes** screen to an authenticated view of each Worker's UI. An orange

header labeled **Viewing Worker:** <host/GUID> will appear to confirm that you are remotely viewing a Worker's UI.



Worker UI access

Using YAML Config File

In \$CRIBL_HOME/local/_system/instance.yml ,under the distributed
section, set mode to master:

```
$CRIBL_HOME/local/_system/instance.yml

distributed:
   mode: master
master:
   host: <IP or 0.0.0.0>
   port: 4200
   tls:
       disabled: true
   ipWhitelistRegex: /.*/
   authToken: <auth token>
   enabledWorkerRemoteAccess: false
   compression: none
   connectionTimeout: 5000
   writeTimeout: 10000
```

Using the Command Line

You can configure a Master Node using a CLI command of this form:

```
./cribl mode-master [options] [args]
```

For all options, see the CLI Reference.

2. Configuring a Worker Node

On each LogStream instance you designate as a Worker Node, you can configure the Worker through the UI, the instance.yml config file, environment variables, or the command line.

Using the UI

In Settings > Distributed Settings > Distributed Management > General Settings, select Mode Worker. Supply the required Master settings (Address and Port). Customize the optional settings if desired. Then click Save to restart.

Using YAML Config File

In \$CRIBL_HOME/local/_system/instance.yml ,under the distributed
section, set mode to worker:

```
$CRIBL_HOME/local/_system/instance.yml
distributed:
 mode: worker
  envRegex: /^CRIBL_/
 master:
   host: <master address>
    port: 4200
    authToken: <token here>
    compression: none
    tls:
      disabled: true
    connectionTimeout: 5000
    writeTimeout: 10000
  tags:
      - tag1
       - tag2
      - tag42
  group: teamsters
```

Using Environment Variables

You can configure Worker Nodes via environment variables, as in this example:

```
CRIBL_DIST_MASTER_URL=tcp://criblmaster@masterHostname:4203
./cribl start
```

See the Environment Variables section for more details.

Using the Command Line

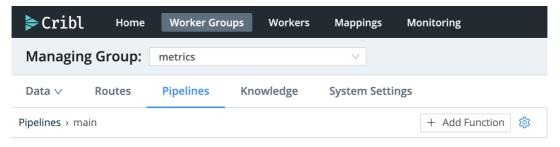
You can configure a Worker Node using CLI commands of this form:

```
./cribl mode-worker -H <master-hostname-or-IP> -p <port> [options] [args]
```

The -H and -p parameters are required. For other options, see the CLI Reference. Here is an example command:

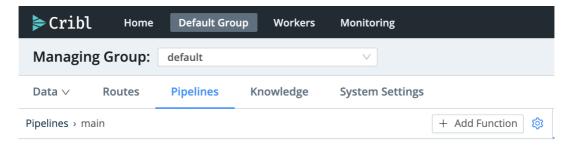
Menu Changes in Distributed Mode

Compared to a single-instance deployment, deploying in distributed mode changes LogStream's menu structure in a few ways. The top menu adds **Worker Groups**, **Workers**, and **Mappings** tabs – all to manage Workers and their assignments.



Distributed deployment: menu structure

If you have a LogStream Free or LogStream One license, the Worker Groups tab instead reads **Default Group**, because these license types allow only this single group. Therefore, throughout this documentation, interpret any reference to the "Worker Groups tab" as "Default Group tab" in your installation.

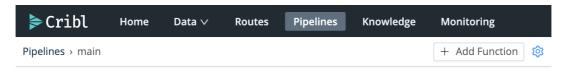


Distributed deployment with LogStream Free/One license

To access the **Data** (drop-down), **Routes**, **Pipelines**, and **Knowledge** items on the light-colored submenu shown above, click the **Worker Groups** tab, then click into your desired Worker Group to display its submenu. This submenu also adds a **System Settings** tab, through which you can manage configuration per Worker Group.

(With a LogStream Free or LogStream One license, you'd click the **Default Group** tab, whose **System Settings** submenu tab configures only that single group.)

For comparison, here is a single-instance deployment's single-level top menu:



Single-instance deployment: single-level menu

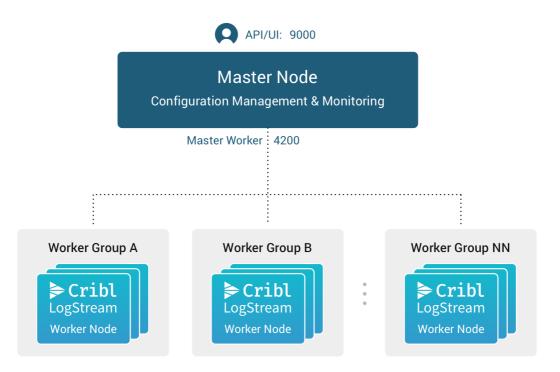
⚠ This repositioning of Data, Routes, Pipelines, and Knowledge tabs to the Worker Groups (or Default Group) submenu also applies to several instructions and screenshots that you'll see throughout this documentation.

Where procedures are written around a single-instance scenario, just click into your appropriate Worker Group to access the same tabs on its submenu.

How Do Workers and Master Work Together

The Master Node has two primary roles:

- 1. Serves as a central location for Workers' operational metrics. The Master ships with a monitoring console that has a number of dashboards, covering almost every operational aspect of the deployment.
- 2. Serves as a central location for authoring, validating, deploying, and synchronizing configurations across Worker Groups.



Master Node/Worker Nodes relationship

Network Port Requirements (Defaults)

- UI access to Master Node: TCP 9000.
- Worker Node to Master Node: TCP 4200 (Heartbeat/Metrics/other).

Master/Worker Node Communication

Workers will periodically (every 10 seconds) send a heartbeat to the Master. This heartbeat includes information about themselves, and a set of current system metrics. The heartbeat payload includes facts – such as hostname, IP address, GUID, tags, environment variables, current software/configuration version, etc. – that the Master tracks with the connection.

The failure of a Worker Node to successfully send two consecutive heartbeat messages to the Master will cause the respective Worker to be removed from the Workers page in the Master's UI until the Master receives a heartbeat message from the affected Worker.

When a Worker Node checks in with the Master:

- The Worker sends heartbeat to Master.
- The Master uses the Worker's facts and Mapping Rules to map it to a Worker Group.
- The Worker Node pulls its Group's updated configuration bundle, if necessary.

Config Bundle Management

Config bundles are compressed archives of all config files and associated data that a Worker needs to operate. The Master creates bundles upon Deploy, and manages them as follows:

- Bundles are wiped clean on startup.
- While running, at most 5 bundles per group are kept.
- Bundle cleanup is invoked when a new bundle is created.

The Worker pulls bundles from the Master and manages them as follows:

- Last 5 bundles and backup files are kept.
- At any point in time, all files created in the last 10 minutes are kept.
- Bundle cleanup is invoked after a reconfigure.

Worker Groups

Worker Groups facilitate authoring and management of configuration settings for a particular set of Workers. To create a new Worker Group, go to the **Worker Groups** top-level menu and click + **Add New**.

Configuring a Worker Group

Click on the newly created Group to display an interface for **authoring and validating** its configuration. You can configure everything for this Group as if it were a single Cribl LogStream instance – using exactly the same visual interface for Routes, Pipelines, Sources, Destinations and System Settings.

△ Can't Log into the Worker Node as Admin User?

To explicitly set passwords for Worker Groups, see User Authentication.

Mapping Workers to Worker Groups

Mapping Rulesets are used to map Workers to Worker Groups. Within a ruleset, a list of rules evaluate Filter expressions on the information that Workers send to the Master.

Only one Mapping Ruleset can be active at any one time, although a ruleset can contain multiple rules. At least one Worker Group should be defined and present in the system.

The ruleset behavior is similar to Routes, where the order matters, and the **Filter** section supports full JS expressions. The ruleset matching strategy is first-match, and one Worker can belong to only one Worker Group.

Creating a Mapping Ruleset

To create a Mapping Ruleset, start on the **Mappings** top-level menu, then click + **Add New**.

i The Mappings top-level menu appears only when you have started LogStream with Distributed Settings > Mode set to Master.

Click on the newly created item, and start adding rules by clicking on + Add Rule. While you build and refine rules, the Preview in the right pane will show which currently reporting and tracked workers map to which Worker Groups.

A ruleset must be activated before it can be used by the Master. To activate it, go to **Mappings** and click **Activate** on the required ruleset. The **Activate** button will then change to an **Active** toggle. Using the adjacent buttons, you can also **Configure** or **Delete** a ruleset, or **Clone** a ruleset if you'd like to work on it offline, test different filters, etc.

Although not required, Workers can be configured to send a Group with their payload. See below how this ranks in mapping priority.

Add a Mapping Rule – Example

Within a Mapping Ruleset, click + Add Rule to define a new rule. Assume that you want to define a rule for all hosts that satisfy this set of conditions:

- IP address starts with 10.10.42, AND:
- More than 6 CPUs OR CRIBL_HOME environment variable contains w0, AND:
- Belongs to Group 420.

Rule Configuration

• Rule Name: myFirstRule

- Filter: (conn_ip.startsWith('10.10.42.') &6 cpus > 6) || env.CRIBL_HOME.match('w0')
- Group: Group420

Default Worker Group and Mapping

When a LogStream instance runs as Master, the following are created automatically:

- A default Worker Group.
- A default Mapping Ruleset,
 - with a default Rule matching all (true).

Mapping Order of Priority

Priority for mapping to a group is as follows: Mapping Rules > Group sent by Worker > default Group.

- If a Filter matches, use that Group.
- Else, if a Worker has a Group defined, use that.
- Else, map to the default Group.

Deploying Configurations

Your typical workflow for deploying LogStream configurations is the following:

- 1. Work on configs.
- 2. Save your changes.
- 3. Commit (and optionally push).
- 4. Deploy.

Deployment is the last step after configuration changes have been saved and committed. **Deploying here means propagating updated configs to Workers.** You deploy new configurations at the Group level: Locate your desired Group and click on **Deploy**. Workers that belong to the group will start **pulling** updated configurations on their next check-in with the Master.

\triangle Can't Log into the the Worker Node as Admin User?

When a Worker Node pulls its first configs, the admin password will be randomized, unless specifically changed. This means that users won't be able to log in on the Worker Node with default credentials. For details, see User Authentication.

Configuration Files

On the Master, a Worker Group's configuration lives under: \$CRIBL_HOME/groups/<groupName>/local/cribl/.

On the managed Worker, after configs have been pulled, they're extracted under: \$CRIBL_HOME/local/cribl/.

Lookup Files

On the Master, a Group's lookup files live under: \$CRIBL_HOME/groups/<groupName>/data/lookups.

On the managed Worker, after configs have been pulled, lookups are extracted under: \$CRIBL_HOME/data/lookups. When deployed via the Master, lookup files are distributed to Workers as part of a configuration deployment.

If you want your lookup files to be part of the LogStream configuration's version control process, we recommended deploying using the Master Node. Otherwise, you can update your lookup file out-of-band on the individual Workers. The latter is especially useful for larger lookup files (> 10 MB, for example), or for lookup files maintained using some other mechanism, or for lookup files that are updated frequently.

For other options, see Managing Large Lookups.

i Some configuration changes will require restarts, while many others require only reloads. See here for details.

Restarts/reloads of each Worker Process are handled automatically by the Worker. Note that individual Worker Nodes might temporarily disappear from the Master's **Workers** tab while restarting.

Worker Process Rolling Restart

During a restart, to minimize ingestion disruption and increase availability of network ports, Worker Processes on a Worker Node are restarted in a rolling fashion. 20% of running processes – with a minimum of one process – are

restarted at a time. A Worker Process must come up and report as started before the next one is restarted. This rolling restart continues until all processes have restarted. If a Worker Process fails to restart, configurations will be rolled back.

Auto-Scaling Workers and Load-Balancing Incoming Data

If data flows in via Load Balancers, make sure to register all instances. Each Cribl LogStream node exposes a health endpoint that your Load Balancer can check to make a data/connection routing decision.

Health Check Endpoint	Healthy Response
curl http:// <host>:<port>/api/v1/health</port></host>	{"status":"healthy"}

Environment Variables

- CRIBL_DIST_MASTER_URL URL of the Master Node. Format:
 <tls|tcp>://<authToken>@host:port?
 group=defaultGroup&tag=tag1&tag=tag2&tls.<tls-settings below>.
 - tls.privKeyPath Private KeyPath.
 - tls.passphrase Key Passphrase.
 - tls.caPath CA Certificate Path.
 - tls.certPath Certificate Path.
 - tls.rejectUnauthorized Validate Client Certs. Boolean, defaults to false.
 - tls.requestCert Authenticate Client (mutual auth). Boolean, defaults to false.
 - tls.commonNameRegex Regex matching peer certificate > subject > common names allowed to connect. Used only if tls.requestCert is set to true.
- CRIBL_DIST_MODE worker | master.Defaults to worker iff
 CRIBL_DIST_MASTER_URL is present.
- CRIBL_HOME Auto setup on startup. Defaults to parent of bin directory.
- CRIBL_CONF_DIR Auto setup on startup. Defaults to parent of bin directory.
- CRIBL_NOAUTH Disables authentication. Careful here!!

- CRIBL_VOLUME_DIR Sets a directory that persists modified data between different containers or ephemeral instances.
- Deprecated variables: CRIBL_CONFIG_LOCATION, CRIBL_SCRIPTS_LOCATION

Workers GUID

When you install and first run the software, a GUID is generated and stored in a .dat file located in CRIBL_HOME/bin/, e.g.:

```
# cat CRIBL_HOME/bin/676f6174733432.dat
{"it":1570724418,"phf":0,"guid":"48f7b21a-0c03-45e0-a699-
01e0b7a1e061"}
```

When deploying Cribl LogStream as part of a host image or VM, be sure to remove this file, so that you don't end up with duplicate GUIDs. The file will be regenerated on next run.

Bootstrap Workers from Master

Boot fully provisioned workers

This feature of LogStream allows workers to completely provision themselves on initial boot, directly from the master. It allows a fleet of any number of nodes to launch. and be fully functional within the cluster, in seconds.

How Does It Work?

A LogStream Master Node (v2.2 or higher) provides a bootstrap API endpoint, at /init/install-worker.sh, which returns a shell script. You can run this shell script on any supported machine (see Restrictions below) without LogStream installed, fully provisioning the machine as a Worker Node.

Although you can specify the download URL when you execute the initial curl command, the LogStream package is not downloaded until the script is generated by the API, and then later executed.

△ Root Access or sudo

Note that the script will install LogStream into <code>/opt/cribl</code>, and will make system-level changes. For systems like Ubuntu, which don't allow direct root access, you'll need to use the <code>sudo</code> command when executing the script.

API Spec

Request Format

GET http://<master hostname or IP>:9000/init/install-worker.sh

Query Strings

String	Required?	Description
token	optional	Master Node's shared secret (authToken). By default, this is set to criblmaster. You can find this secret in the the Master Node's Distributed Settings section.
group	optional	Name of the cluster's work group. If not specified, falls back to default .
download_url	optional	Provide the complete URL to a Cribl LogStream installation binary. This is especially useful if the Worker Nodes don't have access to the Internet to download from cribl.io.

Example HTTP Request

HTTP

GET http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-554

Response

```
Shell
#!/bin/sh
### START CRIBL MASTER TEMPLATE SETTINGS ###
CRIBL MASTER HOST="<Master FQDN/IP>"
CRIBL_AUTH_TOKEN="<Auth token string>"
CRIBL_VERSION="<Version>"
CRIBL_GROUP="<Default group preference>"
CRIBL_MASTER_PORT="<Master heartbeat port>"
CRIBL_DOWNLOAD_URL="<download url>"
### END CRIBL MASTER TEMPLATE SETTINGS ###
# Set defaults
checkrun() { $1 --help >/dev/null 2>/dev/null; }
faildep() { [ $? -eq 127 ] & echo "$1 not found" & exit 1; }
[ -z "${CRIBL_MASTER_HOST}" ] && echo "CRIBL_MASTER_HOST not set" && exit 1
CRIBL_INSTALL_DIR="${CRIBL_INSTALL_DIR:-/opt/cribl}"
CRIBL MASTER PORT="${CRIBL MASTER PORT:-4200}"
CRIBL_AUTH_TOKEN="${CRIBL_AUTH_TOKEN:-criblmaster}"
CRIBL_GROUP="${CRIBL_GROUP:-default}"
if [ -z "${CRIBL_DOWNLOAD_URL}" ]; then
    FILE="cribl-${CRIBL_VERSION}-linux-x64.tgz"
    \label{lownload_url="https://cdn.cribl.io/dl/$(echo $\{CRIBL_VERSION\} \mid cut -d '-' -f 1)/$\{FIIA \mid CRIBL_VERSION\} \mid cut -d '-' -f 1/$\{FIIA \mid CRIBL_VERSION\} \mid cut -d '-' -f 1/$\}
fi
UBUNTU=0
CENTOS=0
AMAZON=0
echo "Checking dependencies"
checkrun curl & faildep curl
checkrun adduser & faildep adduser
checkrun usermod & faildep usermod
BOOTSTART=1
checkrun systemctl & [ $? -eq 127 ] & BOOTSTART=0
checkrun update-rc.d & [ \$? -eq 127 ] & BOOTSTART=0
echo "Checking OS version"
lsb_release -d 2>/dev/null | grep -i ubuntu & [ $? -eq 0 ] & UBUNTU=1
cat /etc/system-release 2>/dev/null | grep -i amazon & [ $? -eq 0 ] & AMAZON=1
echo "Creating cribl user"
if [ $UBUNTU -eq 1 ]; then
    adduser cribl --home /home/cribl --gecos "Cribl LogStream User" --disabled-password
fi
if [ CENTOS - eq 1 ] || [ AMAZON - eq 1 ]; then
    adduser cribl -d /home/cribl -c "Cribl LogStream User" -m
    usermod -aG wheel cribl
fi
```

curl Option

An easy way of wrapping HTTP methods is to use the curl command. Here is an example, which uses a GET operation by default, with the same URL used in the above HTTP example:

Shell

curl http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-5!

Chaining Script Execution

The GET and curl procedures above will only output the contents of the script that needs executing – the script will still need to be manually executed. However, you can automate that part, too, using the command below. This passes the script's contents to the sh shell to immediately execute. As noted above, on Ubuntu and similar systems, you might need to insert sudo before the sh.

Shell

curl http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-5!

Adding Download URL

We'll now graduate to the next level by adding more to the above commands. All the preceding commands excluded the download_url parameter so, by default, the script gets configured to download the LogStream package from the public Cribl repository.

To successfully execute the curl command while also specifying the download_url, you must enclose the URL in double quotes. The reason for this is that the & character that joins multiple HTTP parameters is interpreted by the shell as the operator to run commands in the background. Quoting the URL, as shown in this example, prevents this.

Shell

curl "http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-!

Status Codes

Status Code	Reason
200 – OK	All is well. You should have received the script as a response.
403 – Forbidden	Either the node is not configured as a Master, or the token provided is invalid.

Restrictions

Keep the following in mind when using this feature:

- Each Worker must normally have access to the internet in order to download the Cribl LogStream installation binary from cribl.io. Where this isn't feasible, you can use the download_url switch to point to a binary in a restricted location.
- By default, Worker Nodes communicate with the Master on port 4200. Ensure that access between all Workers and the Master is open on this port.
- TLS is not enabled by default. If enabled and configured, access to this feature will be over https instead of http.
- Red Hat, Ubuntu, CentOS, and Amazon Linux are the only supported Worker platforms.

User Data

For public-cloud customers, an easy way to use this feature is in an instance's user data. Simply use the following script (changing the command as needed. based on the information above). Upon launch, the Worker Node will reach out to the Master, download the script, download the LogStream package from the specified location, and then install and configure LogStream:

Shell

#!/bin/bash
curl http://<master-node-ip/host-address>:9000/init/install-worker.sh?token=<auth-token> | sh -

Splunk App Deployment *

Getting started with Cribl App for Splunk

△ Cribl App for Splunk for HFs Is Deprecated as of Cribl LogStream v.2.1

Cribl will continue to support this package, but **customers are advised to begin planning now for the eventual removal of support**.

See Single-Instance Deployment and Distributed Deployment for alternatives.

Deploying Cribl App for Splunk

In a Splunk environment, Cribl LogStream can be installed and configured as a Splunk app (Cribl App for Splunk). Depending on your requirements and architecture, it can run either on a Search Head or on a Heavy Forwarder. Cribl App for Splunk cannot be used in a Cribl LogStream Distributed Deployment as a Master or managed as Worker.

Running on a Search Head (SH)

When running on an SH, Cribl LogStream is set to **mode-searchhead**, the default mode for the app. It listens for **localhost traffic** generated by a custom command: | criblstream. The command is used to forward search results to the LogStream instance's TCP JSON input on port 10420, but it's also capable of sending to any other LogStream instance listening for TCP JSON.

Once received, data can be processed and forwarded to any of the supported Destinations. In addition, several out-of-the box saved searches are ready to run and send their results to Cribl with a single click.

Installing the Cribl App for Splunk on an SH

- Select an instance on which to install.
- Ensure that ports 10000, 10420, and 9000 are available. See the Requirements section for more info.
- Get the bits here, and install as a regular Splunk app.
- Restart the Splunk instance.
- Go to https://<instance>/en-US/app/cribl or https://<instance>:9000, and log in with Splunk admin role credentials.

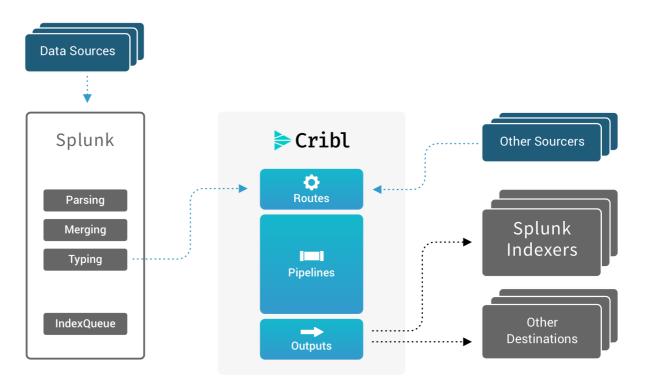
Typical Use Cases for Search Head Mode

• Working with search results in a Cribl LogStream pipeline.

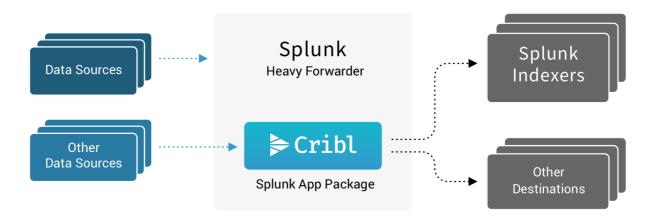
• Sending search results to any Destination supported by Cribl LogStream.

Running on a Heavy Forwarder (HF)

When running on an HF, Cribl LogStream is set to **mode-hwf**. It receives events from the local Splunk process per routing configurations in props.conf and transforms.conf. Data is parsed and processed first by Splunk pipelines, and then by LogStream. By default, all data except internal indexes is routed out right after the Typing pipeline.



Cribl LogStream is capable of accepting data **streams** (unbroken events) or **events** from other sources. In this case, the HF will deliver **events** locally to LogStream, which processes them and sends them to one or more destinations downstream. When receivers are Splunk indexers, LogStream can also load-balance across them.



Installing the Cribl App for Splunk on an HF

• Select an instance on which to install.

- Ensure that ports 10000, 10420, and 9000 are available. See here.
- Get the bits here, and install as a regular Splunk app.
- Set Cribl to mode-hwf: \$SPLUNK_HOME/etc/apps/cribl/bin/cribl mode-hwf.



- Restart the Splunk instance.
- Go to https://<instance>:9000 and log in with Splunk admin role credentials.

```
☐ Note About Splunk Warnings

If you come across messages similar to the following example, on startup or in logs, please ignore them. They are benign warnings.

Invalid value in stanza [route2criblQueue]/[hecCriblQueue] in /opt/splunk/etc/apps/cribl/default/transforms.conf, line 11: (key: DEST_KEY, value: criblQueue) / line 24: (key: DEST_KEY, value: $1)
```

Relevant configurations in Cribl App for Splunk on an HF

When Cribl App for Splunk is installed on an HF (in mode-hwf), below are the **relevant sections** in configuration files that enable Splunk to send data to Cribl LogStream:

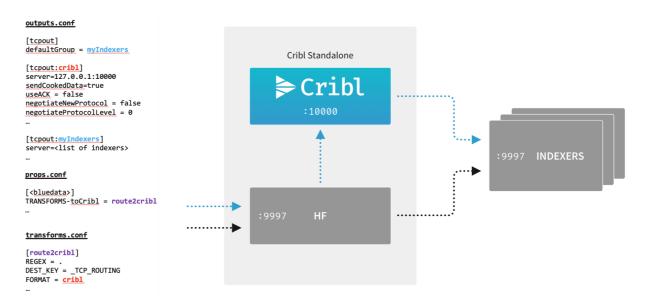
```
apps/cribl/default/outputs.conf
[tcpout]
disabled = false
defaultGroup = cribl
[tcpout:cribl]
server=127.0.0.1:10000
sendCookedData=true
useACK = false
negotiateNewProtocol = false
negotiateProtocolLevel = 0
apps/cribl/default/inputs.conf
[splunktcp]
route=has_key:_replicationBucketUUID:replicationQueue;has_key:_dstrx:typingQueue;has_key:__CRIBN
apps/cribl/default/transforms.conf
[route2cribl]
SOURCE_KEY = _MetaData:Index
REGEX = ^[^_]
DEST_KEY = _TCP_ROUTING
FORMAT = cribl
```

```
[route2criblQueue]
SOURCE_KEY = _MetaData:Index
REGEX = ^[^_]
DEST_KEY = queue
FORMAT = criblQueue

apps/cribl/default/props.conf
[default]
TRANSFORMS-cribl = route2criblQueue, route2cribl
```

Configuring Cribl LogStream with a Subset of Your Data

The props.conf stanza above will apply the above transforms to **everything**. Depending on your requirements, you might want to target only a subset of your sources, sourcetypes, or hosts. For example, the diagram below shows the **effective** configurations of outputs.conf, props.conf, and transforms.conf to send <bluedata> events through Cribl LogStream.



Configure Cribl LogStream to Send Data to Splunk Indexers

To send data from Cribl LogStream to a set of Splunk indexers, use the LogStream UI to go to **Destinations** > Splunk Load Balanced, then enter the required information.

Kubernetes Master Deployment

Boot a fully provisioned Master Node via Helm

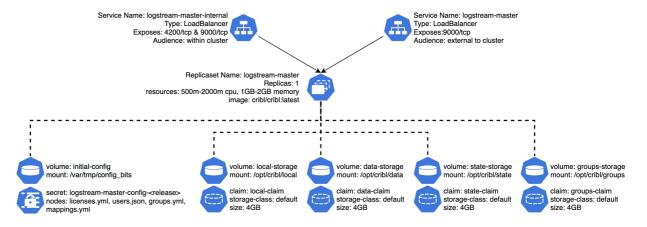
This page outlines how to deploy a Cribl LogStream Master Node (or single instance) to AWS via Kubernetes, using a Cribl-provided Helm chart.

<u>^</u>

This chart is a work in progress, provided as-is. Cribl expects to further develop and refine it.

Deployment

As built, Cribl's chart will deploy a Master Server for LogStream, consisting of a deployment, two services, and a number of persistent volumes.



Deployment schematic

Note that this chart creates two load-balanced services:

- The main one (named after the Helm release), which is intended as the primary service interface for users.
- The "internal" one (named <helm-release > internal), which is intended for the workergroup-to-master communication.
 - By default, this chart installs only a LogStream Master Node. To also deploy LogStream Worker Groups via Helm, you can use the Set Up Worker Groups/Mappings override described below

You can also use Cribl's separate logstream-workergroup chart. For details, see Kubernetes Deployment: Worker Group in this documentation.

AWS and Kubernetes Prerequisites

This section covers both general and specific prerequisites, with a bias toward the EKS-oriented approach that Cribl uses for its own deployments.

Set Up AWS CLI

Install the AWS CLI, version 2, according to AWS' instructions.

Next, create or modify your ~/.aws/config file to include (at least) a [profile] section with the following SSO (single-sign-on) details:

```
~/.aws/config

[profile <your-profile-name>]
sso_start_url = https://<your-domain>/start#/
sso_region = <your-AWS-SSO-region>
sso_account_id = <your-AWS-SSO-account-ID>
sso_role_name = <your-AWS-role-name>
region = <your-AWS-deployment-region>
```

Set Up kubectl

You will, of course, need kubectl set up on your local machine or VM. Follow Kubernetes' installation instructions.

Add a Cluster to Your kubeconfig File

You must modify your ~/.kube/config file to instruct kubectl what cluster (context) to work with.

1. Run a command of this form:

This should return a response like this:

```
Added new context arn:aws:eks:us-west-2:4242424242:cluster/<cluster-name> to /Users/<username>/.kube/config
```

2. In the resulting ~/.kube/config file's args section, as the new first child, insert the profile argument that you provided to the aws command. For example:

```
/.kube/config
args:
   --profile=<profile-name>
   --region
[ ... ]
```

3. Also change the command: aws pair to include the full path to the aws executable. This is usually in /usr/local/bin, in which case you'd insert: command: /usr/local/bin/aws.

This section of ~/.kube/config should now look something like this:

```
~/.kube/config

args:
    --profile=<profile-name>
    --region
    - us-west-2
    - eks
    - get-token
    --cluster-name
    - lab
    command: /usr/local/bin/aws
    env:
    - name: AWS_PROFILE
    value: <profile-name>
```

With these AWS and Kubernetes prerequisites completed, you're now set up to run kubectl commands against your cluster, as long as you have an active aws SSO login session.

Next, do the Helm setup.

Install Helm and Cribl Repo

- 1. You'll need Helm (preferably v.3.x) installed. Follow the instructions here.
- 2. Add Cribl's repo to Helm, using this command:
 helm repo add cribl https://criblio.github.io/helm-charts/

Persistent Storage

The chart requires persistent storage. It will use your default StorageClass, or (if you prefer) you can override config.scName with the name of a specific StorageClass to use.

Cribl has tested this chart primarily using AWS EBS storage, via the CSI EBS driver. The volumes are created as ReadWriteOnce claims. For details about storage classes, see Kubernetes' Storage Classes documentation.

AWS-Specific Notes

If you're running on EKS, Cribl highly recommends that you use Availability Zone–specific node groups. For details, see eksctl.io's Autoscaling documentation.

△ Do not allow a single node group to spans AZs. This can lead to trouble in mounting volumes, because EBS volumes are AZ-specific.

Configure the Chart's Values

You'll want to override some of the chart's default values. The easiest way is to copy this chart's default values. yaml file from our repo. save it locally, modify it, and install it in Helm:

- Copy the raw contents of: https://github.com/criblio/helm-charts/blob/master/helm-chart-sources/logstream-master/values.yaml
- 2. Save this as a local file, e.g.: /bar/values.yaml
- 3. Modify values as necessary (see Values to Override below).
- 4. Install your updated values to Helm, using this command: helm install -f /bar/values.yaml

Values to Override

This section covers the most likely values to override. To see the full scope of values available, run: helm show values cribl/logstream-master

Key	Туре	Default Value	Description
config.adminPassword	String	[No default]	The password you want to assign to the admin user.
config.token	String	[No default]	The auth key you want to set up for Worker access. If you set this value, the LogStream instance will be configured only as a Master server for a distributed deployment. (You can also configure this later via the LogStream UI, after launching the instance in single-instance mode.)
config.license	String	[No default]	The license for your LogStream instance. If you do not set this, it will default to the Free license. You can change this in the LogStream UI as well.
config.groups	List	[No default]	Array of Worker Group names to configure for the Master instance. This will create a mapping for each Group, which looks for the tag <groupname>, and will create the basic structure of each Group's configuration.</groupname>

config.scName	String		The StorageClass name for all of the persistent volumes.
config. rejectSelfSignedCerts	Number	0	Either 0 (allow self-signed certificates) or 1 (deny self-signed certs).
config.healthPort	number	9000	The port to use for health checks (readiness/live).
service.ports	Array of Maps	- name: api port: 9000 protocol: TCP external: true - name: mastercomm port: 4200 protocol: TCP external: false	The ports to make available, both in the deployment and in the service. Each "map" in this list needs the following values set: name A descriptive name, identifying what the port is being used for. port The container port to be made available. protocol The protocol in use for this port (UDP or TCP). external Set to true to expose the port on the external service, or false to not expose it.
service.annotations	String	[No default]	Annotations for the the service component – this is where you'll want to put load-balancer–specific configuration directives.
image.tag	String	latest	The container image tag to pull from. Cribl will increment this tag per LogStream version. By default, this will use a version equivalent to the chart's appVersion value. You can override this with latest to get the latest LogStream version, or with a specific LogStream version number (like "2.3.3").

Match Versions

Cribl recommends that you use the same LogStream version on Worker Nodes versus the Master Node. So if, for any reason, you're not yet upgrading your Workers to the version in the Master's default values.yaml > criblImage.tag, be sure to override that criblImage.tag value to match the version you're running on all Workers.

EKS-Specific Values

If you're deploying to EKS, many annotations are available for the load balancer. Set these as values for the service.annotations key. Internally, we typically use the annotations for logging to S3, like this:

```
values.yaml [excerpt]
```

```
service.beta.kubernetes.io/aws-load-balancer-access-log-enabled: "true"
   service.beta.kubernetes.io/aws-load-balancer-access-log-emit-interval: "5"
   service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-name: "<bucket name>"
   service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-prefix: "ELB"
```

For an exhaustive list of annotations you can use with AWS's Elastic Load Balancers, see the Kubernetes Service documentation.

Basic Chart Installation

With the above prerequisites and configuration completed, you're ready to install our chart to deploy a LogStream Master Node. Here are some example commands:

• To install the chart with the release name logstream-master:

```
helm install logstream-master cribl/logstream-master
```

• To install the chart using the storage class ebs-sc:

```
helm install logstream-master cribl/logstream-master --set config.scName='lebs-sc
```

Change the Configuration

If you don't override its default values, this Helm chart effectively creates a single-instance deployment of LogStream, using the standard container image. You can later configure distributed mode, licensing, user passwords, etc., all from the LogStream UI. However, you also have the option to change these configuration details upfront, by installing with value overrides. Here are some common examples.

Apply a License

If you have a standard or enterprise license, you can use the config.license parameter to add it as an override to your install:

```
helm install logstream-master cribl/logstream-master --set config.license="<long encoded license string redacted>"
```

Set the Admin Password

Normally, when you first install LogStream and log into the UI, it prompts you to change the default admin password. You can skip the password-change challenge by setting your admin password via the config.adminPassword parameter:

helm install logstream-master cribl/logstream-master --set config.adminPassword="
<new password>"

Set Up Worker Groups/Mappings

As mentioned above, the chart's default is to install a vanilla deployment of LogStream. If you are deploying as a Master, you can use the config.groups parameter to define the Worker Groups you want created and mapped. Each group in the list you provide will be created as a Worker Group, with a Mapping Rule to seek a tag with that Worker Group's name in it:

helm install logstream-master cribl/logstream-master --set config.groups= {group1,group2,group3}

The example above will create three Worker Groups – group1, group2, and group3 – and a Mapping Rule for each.

Upgrading from a Pre-2.4.0 Version

LogStream 2.4.0 introduces the \$CRIBL_VOLUME_DIR environment variable, which simplifies the persistent-storage requirement for logstream-master. Instead of maintaining four separate persistent volumes (one each for \$CRIBL_HOME/{data,state,local,groups}), you can consolidate all four into a single volume.

In the Helm chart, we handle this via the helm upgrade command. For the chart's 2.4.0 version, the upgrade option creates a new, larger volume, and consolidates the data from the original volumes to that volume. An initContainer handles the logistics. When the initContainer completes, the logstream-master pod comes back up with a single consolidated volume.

□ Back Up Your Data First

While we've tested this upgrade repeatedly, differences in environments can always cause problems. Therefore, we recommend that you back up your data before running the upgrade command. This is best done with a combination of kubectl and tar:

kubectl exec <pod name> -n <namespace> -- bash -c "cd /opt/cribl; tar cf - {state,data

This command executes the tar -based backup of all four volumes, and outputs it to a local .tar file (cribl_backup.tar).

Running the Upgrade

Helm makes upgrades easy. You simply need to run helm repo update to ensure you have the latest repo updates available, followed by helm upgrade to actually upgrade the containers.

For example, if you've installed the Helm charts in the logstream namespace, named your release ls-master, and set up your Helm repo according to the prerequisites section above (i.e., named it cribl), run the following:

```
helm repo update
helm upgrade ls-master -n logstream cribl/logstream-master
```

Upgrade Order of Operations

While there should be no major problems running a 2.4.0 master and 2.3.4 workers, Cribl does not recommend this. Instead, upgrade the master Helm chart to 2.4.0 first, and then upgrade the workers. (For details, see Kubernetes Worker Deployment.)

Idempotency of Upgrade

The upgrade operation performs a potentially destructive action in coalescing the 4 volumes to a single volume. But that operation proceeds only if the single volume has no data on it. Once the upgrade is performed the first time, any further upgrade operations will effectively skip that coalescence operation, without causing any additional issues.

Recovering from a Failed Upgrade

If the upgrade fails, the suggested recovery path is to remov the Helm chart, reinstall it, and then run this command to restore the data from the backup:

```
cat cribl_backup.tar| kubectl -n <namespace> exec --stdin <pod name> -- bash -c "cd /opt/cribl/c
```

This will restore the data into the "new" volume (which is mounted as /opt/cribl/config-volume). If you want to double-check that:

```
kubectl -n <namespace> exec <pod name> -- bash -c "ls -alR /opt/cribl/config-volume"
```

Uninstall the Infrastructure

To spin down deployed pods, use the helm uninstall command – where <release-name> is the namespace you assigned when you installed the chart:

```
helm uninstall <release-name>
```

You can append the --dry-run flag to verify which releases will be uninstalled before actually uninstalling them:

Known Issues

- Cribl's current architecture supports **only** TCP ports in Worker Groups' service > ports configuration. This restriction might be removed in future versions.
- The upgrade process creates an initContainer, which will run prior to any instance of the LogStream pod. Because the coalescence operation will not overwrite existing data, this is not a functional problem. But depending on your persistent-volume setup, the initContainer 's precedence might cause pod restarts to take additional time while waiting for the volume claims to release. The only upgrade path that will have this issue is 2.3.* -> 2.4.0. In the next iteration, we'll remove the initContainer from the upgrade path.
- The upgrade process leaves the old PersistentVolumes and PersistentVolumeClaims around. This is, unfortunately, necessary for this upgrade path. In follow-on versions, we will remove these volumes from the chart.

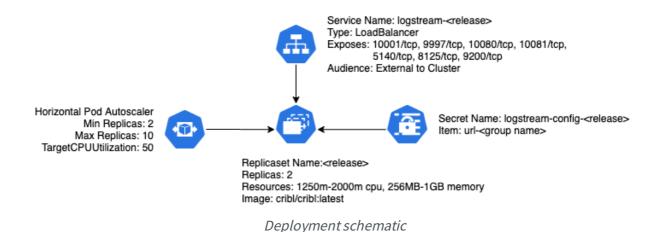
Kubernetes Worker Deployment

Boot a fully provisioned Worker Group via Helm

This page outlines how to deploy a Cribl LogStream Worker Group to AWS via Kubernetes, using a Cribl-provided Helm chart.

Deployment

As built, Cribl's chart will deploy a simple Worker Group for LogStream, consisting of a deployment, a service, a horizontal pod autoscaler configuration, and a secret used for configuration.



☐ This chart will deploy only a LogStream Worker Group. To deploy a LogStream Master Node, see Kubernetes Deployment: Master.

AWS and Kubernetes Prerequisites

This section covers both general and specific prerequisites, with a bias toward the EKS-oriented approach that Cribl uses for its own deployments.

Set Up AWS CLI

Install the AWS CLI, version 2, according to AWS' instructions.

Next, create or modify your ~/.aws/config file to include (at least) a [profile] section with the following SSO (single-sign-on) details:

~/.aws/config

```
[profile <your-profile-name>]
sso_start_url = https://<your-domain>/start#/
sso_region = <your-AWS-SSO-region>
sso_account_id = <your-AWS-SSO-account-ID>
sso_role_name = <your-AWS-role-name>
region = <your-AWS-deployment-region>
```

Set Up kubectl

You will, of course, need kubectl set up on your local machine or VM. Follow Kubernetes' installation instructions.

Add a Cluster to Your kubeconfig File

You must modify your ~/.kube/config file to instruct kubectl what cluster (context) to work with.

1. Run a command of this form:

```
aws --profile <profile-name> eks update-kubeconfig --name <cluster-name>
```

This should return a response like this:

```
Added new context arn:aws:eks:us-west-2:4242424242:cluster/<cluster-name> to /Users/<username>/.kube/config
```

2. In the resulting ~/.kube/config file's args section, as the new first child, insert the profile argument that you provided to the aws command. For example:

```
/.kube/config
args:
   ---profile=<profile-name>
   --region
[ ... ]
```

3. Also change the command: aws pair to include the full path to the aws executable. This is usually in /usr/local/bin, in which case you'd insert: command: /usr/local/bin/aws.

This section of ~/.kube/config should now look something like this:

```
-/.kube/config

args:
    - --profile=<profile-name>
    - -region
    - us-west-2
    - eks
    - get-token
    - --cluster-name
    - lab
    command: /usr/local/bin/aws
```

With these AWS and Kubernetes prerequisites completed, you're now set up to run kubectl commands against your cluster, as long as you have an active aws SSO login session.

Next, do the Helm setup.

Install Helm and Cribl Repo

- 1. You'll need Helm (preferably v.3.x) installed. Follow the instructions here.
- 2. Add Cribl's repo to Helm, using this command: helm repo add cribl https://criblio.github.io/helm-charts/
- 3. Display the default values available to configure Cribl's logstream-workergroup chart: helm show values cribl/logstream-workergroup

Configure the Chart's Values

You'll want to override some of the values you've just displayed. The easiest way is to copy this chart's default values.yaml file from our repo. save it locally, modify it, and install it in Helm:

- Copy the raw contents of: https://github.com/criblio/helm-charts/blob/master/helm-chart-sources/logstream-workergroup/values.yaml
- 2. Save this as a local file, e.g.: /foo/values.yaml
- 3. Modify values as necessary (see Values to Override below).
- 4. Install your updated values to Helm, using this command: helm install -f /foo/values.yaml

Values to Override

This section covers the most likely values to override.

Key	Туре	Default Value	Description
config.tag	String	kubernetes	The tag/group to include in the URL. (This is included as both a tag value

			and a group value.)
config.token	String	criblmaster	The authentication token for your LogStream Master.
config.host	String	logstream- master	The resolveable hostname of your LogStream Master.
<pre>config. rejectSelfSignedCerts</pre>	Number	0	One of: 0 – allow self- signed certs, or 1 – deny self- signed certs.
service.ports	Array of Maps	- name: tcpjson port: 10001 protocol: TCP - name: s2s port: 9997 protocol: TCP - name: http port: 10080 protocol: TCP - name: https port: 10081 protocol: TCP - name: syslog port: 5140 protocol: TCP - name: metrics port: 8125 protocol: TCP - name: elastic port: 9200 protocol: TCP	The ports to make available, both in the deployment and in the service. Each "map" in this list needs the following values set: name A descriptive name, identifying what the port is being used for. port The container port to be made available. protocol The protocol in use for this

			port (UDP or TCP).
service.annotations	String	[No default]	Annotations for the the service component – this is where you'll want to put load-balancer– specific configuration directives.
<pre>criblImage.tag</pre>	String	2.3.3	The container image tag to pull from. Cribl will increment this tag per LogStream version. By default, this will use a version equivalent to the chart's appVersion value. You can override this with latest to get the latest LogStream version, or with a specific LogStream version number.
autoscaling.minReplicas	Number	2	The minimum number of LogStream pods to run.
autoscaling.maxReplicas	Number	10	The maximum number of LogStream pods to scale up to.
autoscaling.targetCPUUtilizationPercentage	Number	50	The CPU utilization percentage that triggers scaling action.

rbac.create	Boolean	false	Enable Service Account, Cluster Role, and Role Binding creation.
rbac.resources	List	["pods"]	Set the resource boundary for the role being created (K8s resources).
rbac.verbs	List	["get", "list"]	Set the API verbs allowed the role (default: read ops).

Match Versions

Cribl recommends that you use the same LogStream version on Master Nodes versus Worker Group Nodes. So, if you're not yet upgrading your Master to the version in the current values.yaml > criblImage.tag, be sure to override that criblImage.tag value to match the version you're running on the Master.

EKS-Specific Values

If you're deploying to EKS, many annotations are available for the load balancer. Set these as values for the service.annotations key. Internally, we typically use the annotations for logging to S3, like this:

```
values.yaml[excerpt]

service:
    type: LoadBalancer
    annotations: {
        service.beta.kubernetes.io/aws-load-balancer-access-log-enabled: "true"
        service.beta.kubernetes.io/aws-load-balancer-access-log-emit-interval: "5"
        service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-name: "<bucket name>"
        service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-prefix: "ELB"
    }
```

For an exhaustive list of annotations you can use with AWS's Elastic Load Balancers, see the Kubernetes Service documentation.

Install the Chart

With the above prerequisites and configuration completed, you're ready to install our chart to deploy a LogStream Worker Group. Here are some example commands:

• To install the chart with the release name logstream-wg:

```
helm install logstream-wg cribl/logstream-workergroup
```

• To install the chart using the LogStream Master logstream.lab.cribl.io:

```
helm install logstream-wg cribl/logstream-workergroup --set config.host='logstream.lab.cribl.io
```

• To install the chart using the LogStream Master logstream.lab.cribl.io in the namespace cribl-helm:

```
helm install logstream-wg cribl/logstream-workergroup --set config.host='logstream.lab.cribl.io' -n cribl-helm
```

Upgrading

You upgrade using the helm upgrade command. But it's important to ensure that your Helm repository cache is up to date, so first issue this command:

```
helm repo update
```

After this step, invoke:

```
helm upgrade <release> -n <namespace> cribl/logstream-workergroup
```

For the example above, where the release is logstream—wg and is installed in the cribl—helm namespace, the command would be:

```
helm upgrade logstream-wg -n cribl-helm cribl/logstream-workergroup
```

This Helm chart's upgrade is idempotent, so you can use the upgrade mechanism to upgrade the chart, but you can also use it to change its configuration (as outlined in Change the Configuration).

Optional: Kubernetes API Access

Starting in version 2.4.0, the chart provides options for Worker Groups to access the Kubernetes API. The values.yaml file provides three relevant options:

- rbac.create enables the creation of a Service Account, Cluster Role, and Role Binding (which binds the first two together) for the release.
- rbac.resources specifies the Kubernetes API resources that will be available to the release.

• rbac.verbs - specifies the API verbs that will be available to the release.

For more information on the verbs and resources available, see Kubernetes' Using RBAC Authorization documentation.

Change the Configuration

Once you've installed a release, you can get its values.yaml file by using the helm get values command. For example, assuming a release name of logstream-wg, you could use this command:

```
helm get values logstream-wg -o yaml > values.yaml
```

This will retrieve a local values.yaml file containing the values in the running release, including any values that you overrode when you installed the release.

You can now make changes to this local values.yaml file, and then use the helm upgrade operation to "upgrade" the release with the new configuration.

For example, assume you wanted to add an additional TCP-based syslog port, listening on port 5141, to the existing logstream-wg release. In the values.yaml file's service > ports section, you'd add the three key-value pairs shown below:

```
values.yaml(excerpt)
service:
[ ... ]
ports:
[ ... ]
- name: syslog
port: 5141
protocol: TCP
```

Then you'd run:

helm upgrade logstream-wg cribl/logstream-workergroup -f values.yaml

Uninstall the Infrastructure

To spin down deployed pods, use the helm uninstall command – where <release-name> is the namespace you assigned when you installed the chart:

```
helm uninstall <release-name>
```

You can append the --dry-run flag to verify which releases will be uninstalled before actually uninstalling them:

Notes on This Example

- If you installed in a namespace, you'll need to include the -n <namespace> option in any helm command.
- In the above syslog example, you'd still need to configure a corresponding syslog Source in your LogStream Master, and then commit and deploy it to your Worker Group(s).

Known Issues

• The chart currently supports **only** TCP ports in service > ports for Worker Groups. This restriction might be removed in future versions.

Sizing and Scaling

A Cribl LogStream installation can be scaled **up** within a single instance and/or scaled **out** across multiple instances. Scaling allows for:

- Increased data volumes of any size.
- Increased processing complexity.
- Increased deployment availability.
- Increased number of destinations.

Scale Up

A LogStream installation can be configured to scale up and utilize as many resources on the host as required. In a single-instance deployment, you govern resource allocation through the **General Settings** > **Worker Processes Settings** section.

In a distributed deployment, you allocate resources per Worker Group.

Navigate to Worker Groups > Group Name > System Settings > Worker

Processes.

Either way, these controls are available:

Process count: Indicates the number of Worker Processes to spawn.
 Positive numbers specify an absolute number of Workers. Negative numbers specify the number of Workers relative to the number of CPUs in the system. like this:

{ number of CPUs available minus this setting }. A 0 setting is interpreted as 1 Worker Process. LogStream corrects for excessive negative offsets by guaranteeing at least 1 Process. Defaults to -2.

- Memory (MB): Amount of memory available to each Worker Process, in MB.
 Defaults to 2048.
 - i Throughout these guidelines, we assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). Each LogStream

instance requires the following resources to run, beyond those reserved for the vCPU's operating system:

- +4 physical cores, +8GB RAM
- 5GB free disk space (more if persistent queuing is enabled)

For example, assuming a Cribl LogStream system with 6 physical cores (12 vCPUs):

- If **Process count** is set to 4, then the system will spawn 4 processes, using up to 4 vCPUs, leaving 8 free.
- If **Process count** is set to -2, then the system will spawn 10 processes (12-2), using up to 10 vCPUs. This will leave 2 vCPUs free.
 - i LogStream incorporates guardrails that prevent spawning more processes than available vCPUs.

It's important to understand that worker processes operate in parallel, i.e., independently of each other. This means that:

- Data coming in on a single connection will be handled by a single worker process. To get the full benefits of multiple Worker Processes, data should come over multiple connections..
 - E.g., it's better to have 5 connections to TCP 514, each bringing in 200GB/day, than one at 1TB/day.
- 2. Each Worker Process will maintain and manage its own outputs. E.g., if an instance with 2 worker processes is configured with a Splunk output, then the Splunk destination will see 2 inbound connections.

Capacity and Performance Considerations

As with most data processing applications, Cribl LogStream's expected resource utilization will be commensurate with the type of processing that is occurring. For instance, a function that adds a static field on an event will likely perform faster than one that applies a regex to finding and replacing a string. At the time of this writing:

• A Worker Process will use up to 1 physical core, or 2 vCPUs.

- Processing performance is proportional to CPU clock speed.
- All processing happens in-memory.
- Processing does not require significant disk allocation.

Estimating Requirements

Current guidance for capacity planning is: Allocate 1 physical core for each 400GB/day of IN+OUT throughput. So, to estimate the number of cores needed: Sum your expected input and output volume, then divide by 400GB.

- Example 1: 100GB IN -> 100GB out to each of 3 destinations = 400GB total = 1 physical core.
- Example 2: 3TB IN -> 1TB out = 4TB total = 10 physical cores.
- Example 3: 4 TB IN -> full 4TB to Destination A, plus 2 TB to Destination B = 10TB total = 25 physical cores.

Recommended AWS, Azure, and GCP Instance Types

You could meet the requirement above with multiples of the following instances:

AWS – Compute Optimized Instances. For other options, see here.

Minimum	Recommended
c5d.2xlarge (4 physical cores, 8vCPUs)	c5d.4xlarge or higher (8 physical cores, 16vCPUs)
c5.2xlarge (4 physical cores, 8vCPUs)	c5.4xlarge or higher (8 physical cores, 16vCPUs)

Azure – Compute Optimized Instances

Minimum	Recommended
Standard_F8s_v2 (4 physical cores, 8vCPUs)	Standard_F16s_v2 or higher (8 physical cores, 16vCPUs)

GCP – Compute Optimized Instances

Minimum	Recommended
c2-standard-8 (4 physical cores, 8vCPUs)	c2-standard-16 or higher (8 physical cores, 16vCPUs)
n2-standard-8 (4 physical cores, 8vCPUs)	n2-standard-16 or higher (8 physical cores, 16vCPUs)

Scale Out

When data volume, processing needs, or other requirements exceed what a single instance can sustain, a Cribl LogStream deployment can span multiple nodes. This is known as a Distributed Deployment, and it can be configured and managed centrally by a single master instance. See Distributed Deployment for more details.

Config Files

Understanding Configuration Paths and Files

Even though all the Routes, Pipelines, and Functions can be managed from the UI, it's important to understand how the configuration works under the hood. At the time of this writing this is how configuration paths and files are laid on the filesystem.

\$CRIBL_HOME	Standalone Install: /path/to/install/cribl/
, <u>-</u>	Splunk App Install: \$SPLUNK_HOME/etc/apps/cribl/

All paths below are relative to \$CRIBL_HOME.

Default Configurations	default/cribl
Local Configurations	local/cribl
System Configuration	<pre>(default local)/cribl/cribl.yml See cribl.yml</pre>
API Configuration	(default local)/cribl/api.yml
Source Configuration	<pre>(default local)/cribl/inputs.yml See inputs.yml</pre>
Destination Configuration	<pre>(default local)/cribl/outputs.yml See outputs.yml</pre>
License Configuration	(default local)/cribl/licenses.yml
Regexes Configuration	(default local)/cribl/regexes.yml
Breakers Configuration	(default local)/cribl/breakers.yml
Limits Configuration	(default local)/cribl/limits.yml

Pipelines Configuration	<pre>(default local)/cribl/pipelines/<pname> Each pipeline's conf is contained therein</pname></pre>
Routes Configuration	(default local)/cribl/pipelines/routes.yml
Functions	<pre>(default local)/cribl/functions/<function_name> Each function's code, conf is contained therein</function_name></pre>
Functions Conf	<pre>(default local)/cribl/functions/<function_name>/ Each function's conf contained therein.</function_name></pre>

Configurations and Restart

- Any configuration changes resulting from UI interactions, for instance, changing the order of functions in a pipeline, or changing the order of routes, do not require restarts.
- All Cribl LogStream configuration file changes resulting from direct file manipulations in (bin|local|default)/cribl/... will require restarts.
- Worker Nodes might temporarily disappear from the Master's **Workers** tab while restarting.
- In the case of a Cribl App for Splunk, Splunk configurations file changes might or might not require restarts. Please check with recent Splunk docs.

Configuration Layering and Precedence

Similar to most *nix systems, Cribl configurations in local take precedence over those in default . There is no layering of configuration files.

△ Editing Configuration Files Manually

When config files ${\it must}$ be edited manually, save all changes in local .

cribl.yml

cribl.yml contains settings for configuring API and other system properties.

```
$CRIBL_HOME/default/cribl/cribl.yml
api:
  # Address to bind to. Default: 0.0.0.0
  host: 0.0.0.0
  # Port to listen to. Default: 9000
  port: 9000
  # Flag to enable/disable UI. Default: false
  disabled : false
  # SSL Settings
  ssl:
    # SSL is enabled by default
    disabled: false
    # Path to private key
    privKeyPath: /path/to/privkey.pem
    # Path to certificate
    certPath: /path/to/cert.pem
auth:
  # Type of authentication.
  type: splunk
  host: localhost
  port: 8089
  ssl: true
workers: # worker processes, memory in MB
  memory: 2048
kms.local:
  # Encryption key management system settings. Default type: local.
  type: local
crypto:
  # Crypto settings.
  keyPath: $CRIBL_HOME/local/cribl/auth/keys.json
system:
  # Upgradability options: api, auto, false
  upgrade: api
  # Restart options: api, false
  restart: api
  # installType options: standalone, splunk-app
  installType: standalone
  # Flag to enable/disable intercom. Default: true
  intercom: true
license:
  accepted: true
# distributed mode: master | worker | single
distributed:
  mode: master
```

inputs.yml

\$CRIBL_HOME/default/cribl/inputs.yml

inputs.yml contains settings for configuring inputs into Cribl.

inputs:
 # Input name
 local-splunk:
 # Input type
 type: splunk

Address to listen to for incoming events
host: localhost
Port to listen to for incoming events
port: 10000

• •

secureTCPJSON:
 type: tcpjson
 disabled: false
 host: 0.0.0.0
 port: 10002
 tls:
 disabled: false
 privKeyPath: /opt/privkey.pem
 certPath: /opt/cert.pem
 requestCert: false
 rejectUnauthorized: false
 ipWhitelistRegex: /.*/

authToken: ""

outputs.yml

outputs.yml contains settings for configuring outputs from Cribl. Also see Destinations for more info.

```
$CRIBL_HOME/default/cribl/outputs.yml
outputs:
  # Default output setting
  default:
   type: default
    defaultId: local-splunk
  # Output Name
  local-splunk:
    # Output type
    type: splunk
    # Output host address to send data from
    host: localhost
    # Output port to send data from
    port: 9999
  # Output name
  myFilesystemDestination:
    # Output type
    type: filesystem
    # Final destination path. Writable by Cribl.
    destPath: /path/to/destiation
    # Staging destination path. Writable by Cribl.
    stagePath: /tmp/foo
    # Partition schema for outputted files
    partitionExpr: >-
      `${host}/${sourcetype}`
    # Format of the output data
    format: json
    # The output filename prefix
    baseFileName: CriblOut
    # Compression options. None | Gzip
    compress: none
    # Maximum uncompressed output file size
    maxFileSizeMB: 32
    # Maximum amount of time to keep inactive files open.
    maxFileOpenTimeSec: 300
    # Maximum amount of time to keep inactive files open.
    maxFileIdleTimeSec: 30
    # Maximum number of files to keep open concurrently.
    maxOpenFiles: 100
  myS3Destination:
    # Output type
    type: s3
    # S3 bucket address
    bucket: s2.bucket.address.here
    # Prefix to append to files before uploading
    destPath: keyprefix
    # AWS API key, if not present will fallback on env.AWS_ACCESS_KEY_ID, or the meta-data endpo
    awsApiKey: key
    # AWS Secret Key. If left blank, Cribl will fallback on env.AWS_SECRET_ACCESS_KEY, or the r
    awsSecretKey: secretkey
    # Staging destination path. Writable by Cribl.
```

stagePath: /tmp/foo

licenses.yml

licenses.yml maintains a list of licenses for Cribl.

\$CRIBL_HOME/default/cribl/licenses.yml

licenses:

- # List of license keys
- eyJ0eXAiOiJKV1QiLCJhasdfasfasdfdasfasdfa-Abo2_ogVbR_5VKeAelZlTc5b-TKQax9R1ywnoOG8guis2RC0sSl

regexes.yml

regexes.yml maintains a list of regexes. Cribl's Regex Library ships under default.

```
$CRIBL_HOME/default/cribl/regexes.yml
```

```
"uuid":
 lib: cribl
 description: UUID/GUID
 regex: /[0-9a-f]\{8\}-[0-9a-f]\{4\}-[1-5][0-9a-f]\{3\}-[89ab][0-9a-f]\{3\}-[0-9a-f]\{12\}/gm
 sampleData: 9a50fa34-58b1-4a67-8b8d-ea9c0ae48c8f
   eb671525-2b9e-4140-ae21-a0a8a81b506e
 tags: uuid,guid
"aws_secret_key":
 description: AWS Secret Access Key
 regex: /(?<![A-Za-z0-9]/+=])[A-Za-z0-9]/+=]{40}(?![A-Za-z0-9]/+=])/gm
 lib: cribl
 sampleData: wJalrXUtnFEMI/K7MDENG/bPxRfiCYEXAMPLEKEY
 tags: aws,access,key,secret
"aws_access_key":
 lib: cribl
 description: AWS Access Key ID
 regex: /(A3T[A-Z0-9]|AKIA|AGPA|AIDA|AROA|AIPA|ANPA|ANVA|ASIA)[A-Z0-9]\{16\}(?![A-Za-z0-9\/+=])/(2)
 sampleData: >-2
    AKIAIOSFODNN7EXAMPLE
 tags: aws,access,key
"private_key":
 description: Private key block
 regex: /----BEGIN (DSA|RSA|EC|PGP|OPENSSH) PRIVATE KEY(\sBLOCK)?----[\s\S]*/gm
 lib: cribl
 tags: ssh,openssh,dsa,ec,rsa,private key
"slack_token":
 lib: cribl
 description: Slack Token
 regex: /xox[p|b|o|a][\s\S]*/g
 sampleData: xoxp-23984754863-2348975623103
   xoxa-23984754863-2348975623103
   xoxb-23984754863-2348975623103
   xoxo-23984754863-2348975623103
 tags: slack, token
```

breakers.yml

Cribl's default Event Breaker Library is located under \$CRIBL_HOME/default/cribl/breakers.yml.

```
$CRIBL_HOME/default/cribl/breakers.yml
AWS Ruleset:
     lib: cribl
      description: Event breaking rules for common AWS data sources
      tags: flowlogs,elb,alb,loadbalancer,cdn
      rules:
            - name: AWS VPC Flow
                  condition: /^\d+\s+\d+\s+eni-\w+.*(OK|NODATA|SKIPDATA)?$/.test(_raw) || sourcetype='aws:
                   eventBreakerRegex: /[\n\r]+/
                  timestampAnchorRegex: /(?=\d{10}\s\d{10})/
                  timestamp:
                        type: format
                       length: 150
                        format: "%s"
                  timestampTimezone: utc
                   maxEventBytes: 1024
             - name: AWS ALB
                   condition: /^{(:https:|h2|wss:)}\d+-\d+-\d+.*?arn:aws:elasticloadbalancing/.test(_raw) ||
                   eventBreakerRegex: /[\n\r]+/
                   timestampAnchorRegex: /\w+\s/
                   timestamp:
                        type: format
                        length: 150
                        format: "%Y-%m-%dT%H:%M:%S.%f%Z"
                  timestampTimezone: local
                  maxEventBytes: 4096
             - name: AWS ELB
                   condition: /^d+-d+-d+.*?(?:d+)_3/.test(_raw) || sourcetype='aws:elb:accesslossibset || sourcet
                   eventBreakerRegex: /[\n\r]+/
                   timestampAnchorRegex: /^/
                   timestamp:
                        type: format
                        length: 150
```

format: "%Y-%m-%dT%H:%M:%S.%f%Z"

timestampTimezone: local
maxEventBytes: 4096

mappings.yml

Mapping ruleset configurations are located under \$CRIBL_HOME/local/cribl/mappings.yml . \$CRIBL_HOME/default/cribl/mappings.yml

```
...
rulesets:
 default: # ruleset name
   conf:
     functions:
        - filter: env.CRIBL_HOME.match('w0') # filter to match
          description: w0 # rule name/id
          final: true
          conf:
           add:
              - name: groupId
               value: "'myGroup42'" # group to map to
        - filter: env.CRIBL_HOME.match('w1')
          id: eval
          description: w1
          final: true
          conf:
           add:
              - name: groupId
               value: "'NewGroup22'"
  newruleset: # another ruleset
   conf:
      functions:
        - filter: (cpus>12 & env.CRIBL_HOME.match('w0')) || release.startsWith('18')
          description: catch all
          final: true
          conf:
            add:
              - name: groupId
               value: "'NewGroup2'"
```

instance.yml

```
Instance configuration is located under
$CRIBL_HOME/local/_system/instance.yml.
 $CRIBL_HOME/local/_system/instance.yml
 distributed:
     # mode master | worker | single
   mode: master
   master:
     host: 0.0.0.0
     port: 4203
     tls:
       disabled: true
     ipWhitelistRegex: /.*/
     authToken: criblmaster
     compression: none
     connectionTimeout: 5000
     writeTimeout: 10000
   group: default
   envRegex: /^CRIBL_/
   tags:
        - tag1
        - tag2
        - tag42
```

Licensing

Every Cribl LogStream download package ships with a Free license that allows for processing of up to 1 TB/day. LogStream Free and LogStream One licenses require sending anonymized telemetry metadata to Cribl. (For details, see Telemetry Data below). Cribl does not require a separate license for sending data from LogStream to LogStream, such as sending from one worker group managed by master node A to a different worker group managed by master node B. The same license used on master node A can be used on master node B in that situation.

Enterprise, Standard, and Sales Trial licenses do **not** require sending telemetry metadata, and are entitled to a defined, per-license daily ingestion volume.

This page summarizes all these license types.

Managing Licenses

You can add and manage licenses in **Settings** > **Licensing**. Click + **Add License** to paste in a license key provided to you by Cribl.

☐ License Expiration and Renewal

For LogStream v.2.2 and earlier, the latest Free license expires on: 2020-12-15T00:00:00+00:00

For LogStream v.2.3 and later, Free licenses do not expire.

LogStream One and LogStream Standard licenses must be renewed annually.

License Types

Cribl offers five LogStream license types, summarized below.

i For a detailed comparison of what's included in each license type, please see Cribl Pricing.

Enterprise License

This is a license available for purchase.

- Up to unlimited data ingestion.
- Role-based access control.
- External authentication (via LDAP, Splunk, and OpenID Connect identity providers).
- Git remote backup.
- All other LogStream features included.

Contact Cribl Sales at sales@cribl.io for more information.

Standard License

This is a license available for purchase. Compared to an Enterprise license, it offers a cost discount, in exchange for some limitations (all data volumes below based on uncompressed data size):

- Daily ingestion up to 5 TB/day.
- Maximum 1 Worker Group.

Contact Cribl Sales at sales@cribl.io for more information.

Free License

Free licenses ship in the download package, and are permanent. They impose some limitations:

- Daily ingestion up to 1 TB/day.
- Maximum 10 Worker Processes.
- Maximum 1 Worker Group.

"One" License

LogStream One is a type of free license that allows for higher processing volume, but only to **one** Splunk (Single-Instance or Load-Balanced) or Elasticsearch Destination. This combination is designed to help users explore LogStream's value in routing large data volumes to these common services. Contact Cribl Sales at sales@cribl.io to convert a Free license to a LogStream One license, which must be renewed annually.

- Daily ingestion up to 5 TB/day, only to one of either Splunk or Elasticsearch outputs.
- Maximum 50 Worker Processes
- Maximum 1 Worker Group

Sales Trial License

A license type used when preparing a POC (proof of concept), or a pilot, with requirements that go beyond those afforded by the Free or One license. Contact Cribl Sales at sales@cribl.io for more information.

i LogStream Free and LogStream One licenses require sending of anonymized telemetry metadata to Cribl. These licenses will block inputs if sending fails after a grace period of 24 hours.

Combining License Types

Multiple license types can coexist on an instance. However, only a **single type** of license can be effective at any one time. When multiple types coexist, the following method of resolution is used:

- If there are any unexpired Enterprise or Standard licenses use only these licenses to compute the effective license.
- Else, if there are any Sales Trial licenses use only Sales Trial licenses to compute the effective license.
- Else, if there exists a Free or One license use only the Free or One license to compute the effective license.

When an Enterprise or Standard license expires, Cribl LogStream will fall back to the Sales Trial or Free/One types. However, an expired Sales Trial license cannot fall back to a Free/One license.

Upon expiration of a paid license, if there is no fallback license, LogStream will backpressure and block all incoming data.

Licensing in Distributed Deployments

LogStream 2.2.x or Earlier

In distributed deployments of LogStream versions through 2.2.x, licenses should be configured both on the Master Node and on each of the Worker Groups. This allows for different Worker Groups to have different licensing capacities.

- To configure the Master: **Settings** > **Licensing**.
- To configure Worker Groups: Worker Groups > [Select a Group] > System
 Settings > Licensing.

LogStream 2.3.x or Later

☐ As of LogStream 2.3, you no longer need to add licenses directly to Worker Groups. The Master will push license information down to Worker Groups as part of the heartbeat.

LogStream will attempt to balance (or rebalance) Worker Processes/threads as evenly as possible across all licensed Worker Nodes.

△ LogStream 2.3 changes licensing in other ways that might require you to update an existing LogStream configuration. Please see Upgrading to LogStream 2.3.

Telemetry Data

A Free or One license requires sharing of telemetry metadata with Cribl. Cribl uses this metadata to help us understand how to improve the product and prioritize new features. Telemetry payloads are sent to an endpoint located on https://cdn.cribl.io/telemetry/. (For versions prior to 2.2, this endpoint is 34.220.85.61:8000.)

i Testing telemetry endpoint connectivity

To manually test connectivity to the telemetry endpoint, especially if you are needing to configure a proxy, you can use the following

command:

```
$ curl https://cdn.cribl.io/telemetry/
Expected response:
cribl /// living the stream!
```

If you get a 302 response code then you may be missing the / at the end of the URL.

If you would like this feature disabled in order to deploy in your environment, please contact Cribl Sales at sales@cribl.io, and we will work with you to issue licenses on a case-by-case basis.

i Once you have received a license that removes the telemetry requirement, disable telemetry in LogStream's UI at Settings > System > General Settings > Upgrade & Share Settings > Sharing and Live Help. Toggle the slider to No.

Data Shared Per Interval (roughly, every minute):

- Version
- Instance's GUID
- License ID
- Earliest, Latest Time
- Number of Events In, Out
- Number of Bytes In, Out
- Number of Open, Closed, Active Connections
- Number of Routes
- Number of Pipelines

Licensing FAQ

How do I check my license type, restrictions, and/or expiration date?

Open LogStream's **Settings > Licensing** page to see these details.

How can I track my actual data ingestion volume over the last 30 days?

Forward Cribl Internal metrics to your Metrics Destination of choice, and run a report on cribl.total.in_bytes.

How does LogStream enforce license limits?

If your data throughput exceeds your license quota, Chuck Norris will track you down and make your life a living hell.

However, that will happen only in your nightmares. In the product itself:

- Free, One, and Standard licenses enforce data ingestion quotas through limits on the number of Worker Groups and Worker Processes.
- Enterprise license keys turn off all enforcement, between annual true-ups.
- When an Enterprise or Standard license expires, LogStream will attempt to fall back to a trial or free license, or only if that fails will block incoming data. For details, see Combining License Types.

I'm using LogStream 2.3.0 or higher, with its "permanent, Free" license. Why is LogStream claiming an expired license, and blocking inputs?

This can happen if you've upgraded from a LogStream version below 2.3.0, in which you previously entered this earlier version's Free (time-limited) license key. To remedy this, go to **Settings > Licensing**, click to select and expand your expired Free license, and then click **Delete license**. LogStream will fall back to the new, permanent Free license behavior, and will restore throughput.

If I pull data from compressed S3 buckets, is my license quota applied to the compressed or the uncompressed size of the file objects?

To measure license consumption, LogStream uses the uncompressed size.

Access Management

Cribl LogStream provides a range of access-management features for users with different security requirements. For details, see the following topics:

- Authentication: Authenticating users in LogStream.
- Local Users: Creating and managing users and their permissions.
- Roles: Managing roles and policies to assign to users.
 - i Role-based access control can be enabled only on distributed deployments with an Enterprise license.

Authentication

User authentication in LogStream

Cribl LogStream supports **local**, **Splunk**, **LDAP**, and **SSO/OpenID Connect** authentication methods.

Local Authentication

To set up local authentication, navigate to **Settings > General Settings > Authentication Settings** and select **Local**.

You can then manage users through the **Settings > Local Users** UI. All changes made to users are persisted in a file located at \$CRIBL HOME/local/cribl/auth/users.json.

Line format:

```
{"username":"user","first":"Elvis","last":"Bath","disabled":"false",
"passwd":"Yrt0MOD1w80zyMYB8WMcEleOtYESMwZw2qIZyTvueOE"}
```

The file is monitored for modifications every 60s, and will be reloaded if changes are detected.

Adding users through direct modification of the file is also supported, but not recommended.

Manual Password Replacement

To manually add, change, or restore a password, replace the affected user's passwd key-value pair with a password key, in this format: "password":" <newPlaintext>" . LogStream will hash all plaintext password(s), identified by the password key, during the next file reload, and will rename the plaintext password key.

Starting with the same users.json line above:

```
{"username":"user","first":"Elvis","last":"Bath","disabled":"false",
"passwd":"Yrt0MOD1w80zyMYB8WMcEleOtYESMwZw2qIZyTvueOE"}
```

...you'd modify the final key-value pair to something like:

```
{"username":"user","first":"Elvis","last":"Bath","disabled":"false",
"password":"V3ry53CuR&pW9"}
```

Within at most one minute after you save the file, LogStream will rename the password key back to passwd, and will hash its value, re-creating something resembling the original example.

Set Worker Passwords

In a distributed deployment, once a worker has been set to point to the Master Node, LogStream will set each Worker node's admin password with a randomized password which is different from the admin user's password on the Master Node. This is by design, as a security precaution, but may lead to situations where administrators cannot log into a Worker Node directly and must rely on accessing them via the Master

To explicitly push a known/new password to your Worker Node, set and push a new password to the Worker Group.

In the Master Node's UI:

- 1. From the top menu, select Worker Groups.
- 2. Select the desired Worker Group.
- 3. From the Worker Groups submenu, select System Settings.
- 4. Select **Local Users**, then expand the desired user.
- 5. Update the **Password** field and select **Save**.

Every 10 seconds, the Worker Nodes will request an update of configuration from the Master and new password settings will be reflected.

Authentication Controls

You can customize authentication behavior at **General Settings > API Server Settings > Advanced.** The options here include:

- Logout on Roles change: If role-based access control is enabled, determines whether users are automatically logged out of LogStream when their assigned Roles change. Defaults to Yes.
- Auth-token TTL: Sets authentication tokens' valid lifetime, in seconds. Defaults to 3600 (60 minutes).

- Login rate limit: Sets the number of login attempts allowed over a (selectable) unit of time. Defaults to 2/second.
- **HTTP header**: Enables you to specify one or more custom HTTP headers to be sent with every response.

The cribl.secret File

When Cribl LogStream first starts, it creates a \$CRIBL_HOME/local/cribl/auth/cribl.secret file. This file contains a key that is used to generate auth tokens for users, encrypt their passwords, and encrypt encryption keys.

Default local credentials are: admin/admin

Back up and secure access to this file by applying strict permissions – e.g., 600.

Splunk Authentication

Splunk authentication is very helpful when deploying in the same environment as Splunk, and requires the user to have Splunk admin role permissions. To set up Splunk authentication:

Navigate to **Settings > General Settings > Authentication Settings** and select **Splunk**.

- Host: Splunk hostname (typically a search head).
- Port: Splunk management port (defaults to 8089).
- SSL: Set to Yes if enabled.
- Fallback to local: Attempt local authentication if Splunk authentication is unsuccessful. Defaults to false.

Notes: The Splunk searchhead does not need to be locally installed on the LogStream instance. See also Role Mapping below.

LDAP Authentication

LDAP authentication is supported, and can be set up as follows:

Navigate to Settings > General Settings > Authentication Settings, and select LDAP.

- Secure: Enable to use a secure LDAP connections (ldaps://). Disable for an insecure (ldap://)connection.
- LDAP servers: List of LDAP servers. Each entry should contain host:port (e.g., localhost:389).
- Bind DN: Distinguished name of entity to authenticate with LDAP server. E.g., 'cn=admin,dc=example,dc=org'.
- Password: Distinguished Name password used to authenticate with LDAP server.
- User search base: Starting point to search LDAP for users, e.g., 'dc=example,dc=org'.
- Username field: LDAP user search field, e.g., cn or (cn (or uid).
- User search filter: LDAP search filter to apply when finding user, e.g., (& (group=admin)(!(department=123*))).Optional.
- **Group search base**: Starting point to search LDAP for groups, e.g., dc=example,dc=org.Optional.
- **Group member field**: LDAP group search field, e.g., member . Optional.
- **Group search filter**: LDAP search filter to apply when finding group, e.g., (& (cn=cribl*)(objectclass=group)).Optional.
- Fallback to local: Attempt local authentication if LDAP authentication is down or is mis-configured. Defaults to No.
- Connection timeout (ms): Defaults to 5000.
- Reject unauthorized: Valid for secure LDAP connections. Set to Yes to reject unauthorized server certificates.
- **Group name field**: LDAP group field, e.g., cn.

Note: See also Role Mapping below.

SSO/OpenID Connect Authentication

LogStream supports SSO/OpenID user authentication (login/password) and authorization (user's group membership, which you can map to Cribl Roles). Set this up as follows:

Navigate to **Settings > General Settings > Authentication Settings** and select **OpenID Connect**.

- **Provider name**: The name of the identity provider service. You can select **Google** or **Okta**, both supported natively. Manual entries are also allowed.
- Audience: The Audience from provider configuration. This will be the base URL, e.g.: https://yourDomain.com:9000.
- **Client ID**: The client_id from provider configuration.
- Client secret: The client_secret from provider configuration.
- **Scope**: Space-separated list of authentication scopes. The default list is: openid profile email.
- Authentication URL: The full path to the provider's authentication endpoint.

 Be sure to configure the callback URL at the provider as

 <yourDomainUrl>/api/v1/auth/authorization-code/callback, e.g.:

 https://yourDomain.com:9000/api/v1/auth/authorization-code/callback.
- Token URL: The full path to the provider's access token URL.
- Logout URL: The full path to the provider's logout URL. Leave blank if the provider does not support logout or token revocation.
- Validate certs: Whether to validate certificates. Defaults to Yes . Toggle to No to allow insecure self-signed certificates.
- **Filter type**: Select either **Email whitelist** or **User info filter**. This selection displays one of the following fields:
 - Email whitelist: Wildcard list of emails that are allowed access.
 - User info filter: JavaScript expression to filter against user profile attributes. E.g.: name.startsWith("someUser") &6 email.endsWith("domain.com")
- **Group name field**: Field on the id_token that contains the user groups. Defaults to groups.

Note the following details when filling in the form – for example, when using Okta:

- <Issuer URI> is the account at the identity provider.
- Audience is the URL of the host that will be connecting to the Issuer (e.g., https://localhost:9000). The issuer (Okta, in this example) will redirect back to this site upon authentication success or failure.

See also Role Mapping below.

Role Mapping

This section is displayed only on distributed deployments with an Enterprise license. For details on mapping your external identity provider's configured groups to corresponding LogStream user access Roles, see External Groups and LogStream Roles.

- **Default role**: Default LogStream Role to assign to all groups not explicitly mapped to a Role.
- Mapping: On each mapping row, enter an external group name on the left, and select the corresponding LogStream Role on the right drop-down list. Click + Add Mapping to add more rows.

Local Users

This page covers how to create and manage LogStream users, including their credentials and (where enabled) their access roles. These options apply if you're using the **Local** Authentication type, which is detailed here.

Creating and Managing Local Users

On the Master Node – or in a single-instance deployment – you manage users by selecting Settings > Access Management > Local Users.

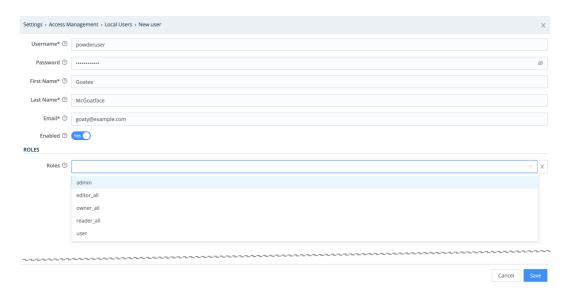
The resulting **Manage Local Users** page will initially show only the default admin user. You are operating as this user.



Managing users

To create a new LogStream user, click + Add New. To edit an existing user, click anywhere on its row. With either selection, you will see the modal shown below.

The first few fields are self-explanatory: they establish the user's credentials. If you want to establish or maintain a user's credentials on LogStream, but prevent them from currently logging in, you can toggle the **Enabled** slider to No.



Adding Roles

If you've enabled role-based access control you can use the modal's bottom **Roles** section to assign access Roles to this new or existing user.

i For details, see Roles. Role-based access control can be enabled only on distributed deployments with an Enterprise license.

Click + Add Role to assign each desired role to this user. The options on the Roles drop-down reflect the Roles you've configured in Settings > Access Management > Roles.

Note that when you assign multiple Roles to a user, the Roles' permissions are additive: This user is granted a superset of the highest permissions contained in all the assigned Roles.

When you've configured (or reconfigured) this user as desired, click **Save**.

By default, LogStream will log out a user upon a change in their assigned Roles. You can defeat this behavior at **General Settings > API Server Settings > Advanced > Logout on roles change.**

Roles

Define and manage access-control roles and policies

Cribl LogStream offers role-based access control (RBAC) to serve these common enterprise goals:

- **Security**: Limit the blast radius of inadvertent or intentional errors, by restricting each user's actions to their needed scope within the application.
- Accountability: Ensure compliance, by restricting read and write access to sensitive data.
- Operational efficiency: Match enterprise workflows, by delegating access over subsets of objects/resources to appropriate users and teams.
 - i Role-based access control is enabled only on distributed deployments with an Enterprise license.

RBAC Concepts

LogStream's RBAC mechanism is designed around the following concepts, which you manage in the UI:

- Roles: Logical entities that are associated with one or multiple Policies (groups of permissions). You use each Role to consistently apply these permissions to multiple LogStream users.
- **Policies**: A set of **permissions**. A Role that is granted a given Policy can access, or perform an action on, a specified LogStream object or objects.
- Permissions: Access rights to navigate to, view, change, or delete specified objects in LogStream.
- Users: You map Roles to LogStream users in the same way that you map user groups to users in LDAP and other common access-control frameworks.

Users are independent LogStream objects that you can configure even without RBAC enabled. For details, see Local Users.

How LogStream RBAC Works

LogStream RBAC is designed to grant arbitrary permissions over objects, attributes, and actions at arbitrary levels.

i As of v. 2.4, Roles are customizable only down to the Worker Group level. E.g., you can grant Edit permission on Worker Group WG1 to User A and User B, but cannot grant them finer-grained permissions on child objects such as Pipelines, Routes, etc.

LogStream's UI will be presented differently to users assigned Roles with access restrictions. Controls will be visible but disabled, and search and log results will be limited, depending on each user's permissions.

LogStream Roles can be integrated with external authorization/IAM mechanisms, such as LDAP and OIDC and mapped to their respective groups, tags, etc.

Using Roles

LogStream ships with a set of default Roles, which you can supplement.

Default Roles

These Roles ship with LogStream by default:

Name	Description
admin	Superusers – authorized to do anything and everything in the system.
owner_all	Read/write access to (and Deploy permission on) all Worker Groups.
editor_all	Read/write access to all Worker Groups.
reader_all	Read-only access to all Worker Groups.
user	Default role that gets only a home/landing page to authenticate. This

is a fallback for users who have not yet been assigned a higher role by an admin.

Cribl **strongly recommends** that you do not edit or delete these default roles. However, you can readily clone them (see **Clone Role** below), and modify the duplicates to meet your needs.

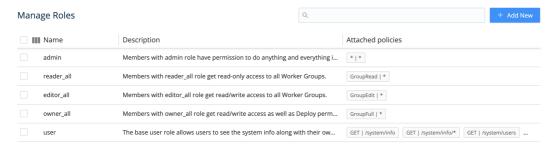
☐ Initial Installation or Upgrade

When you first install LogStream with the prerequisites to enable RBAC (Enterprise license and distributed deployment), you will be granted the **admin** role. Using this role, you can then define and apply additional roles for other users.

You will similarly be granted the **admin** role upon upgrading an existing LogStream installation from pre-2.4 versions to v. 2.4 or higher. This maintains backwards-compatible access to everything your organization has configured under the previous LogStream version's single role.

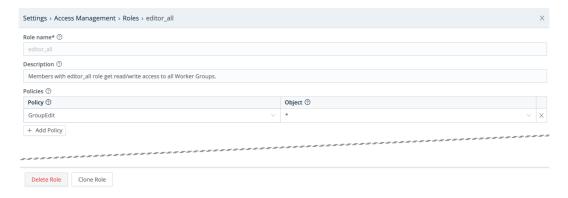
Adding and Modifying Roles

In a distributed environment, you manage Roles at the Master level, for the entire deployment. On the Master Node, select **Settings** > **Access Management** > **Roles**.



Manage Roles page

To add a new Role, click + Add New at the upper right. To edit an existing Role, click anywhere on its row. Here again, either way, the resulting modal offers basically the same options.



Add/edit Role modal

The options at the modal's top and bottom are nearly self-explanatory:

Role name: Unique name for this Role.

Description: Optional free-text description.

Delete Role: And...it's gone. (But first, there's a confirmation prompt.)

Clone Role: Opens a **New role** version of the modal, duplicating the **Description** and **Policies** of the Role you started with.

The modal's central **Policies** section (described below) is its real working area.

Adding and Modifying Policies

The **Policies** section is an expandable table. In each row, you select a Policy using the left drop-down, and apply that Policy to objects (i.e., assign permissions on those objects) using the right drop-down.

Let's highlight an example from the above screen capture of LogStream's built-in Roles: The editor_all Role has the GroupEdit Policy, with permission to exercise it on any and all Worker Groups (as indicated by the * wildcard).



Policies on the left, objects on the right

To add a new Policy to a Role:

- 1. Click + Add Policy to add a new row to the Policies table.
- 2. Select a Policy from the left column drop-down.

3. Accept the default object on the right; or select one from the drop-down.

To modify an already-assigned Policy, just edit its row's drop-downs in the **Policies** table.

To remove a Policy from the Role, click its close box at right.

In all cases, click **Save** to confirm your changes and close the modal.

Default Policies

In the **Policies** table's left column, the drop-down offers the following default Policies:

Name	Description
GroupRead	The most basic Worker Group-level permission. Enables users to view a Worker Group and/or its configuration.
GroupEdit	Building on GroupRead, grants the ability to also change and commit a Worker Group's configuration.
GroupFull	Building on GroupEdit, grants the ability to also deploy a Worker Group.
* (wildcard)	Grants all permissions on associated objects.

Objects and Permissions

In the **Policies** table's right column, use the drop-down to select the LogStream objects on which the left column's Policy will apply. (Remember that in v. 2.4, the objects available for selection are specific Worker Groups, or a wildcard representing all Worker Groups.) For example:

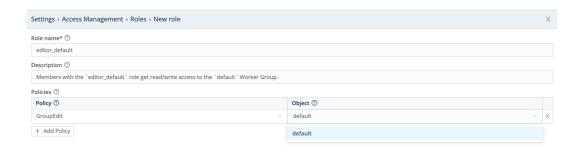
- Worker Group <id>
- NewGroup2
- default (Worker Group)
- * (all Worker Groups)

Extending Default Roles

Here's a basic example that ties together the above concepts and facilities. It demonstrates how to add a Role whose permissions are restricted to a particular Worker Group.

Here, we've cloned the editor_all Role that we unpacked above. We've named the clone editor_default .

We've kept the GroupEdit Policy from editor_all. But in the right column, we're restricting its object permissions to the default Worker Group that ships with LogStream.



Cloning a default Role

You can readily adapt this example to create a Role that has permissions on an arbitrarily named Worker Group of your own.

Roles and Users

Once you've defined a Role, you can associate it with LogStream users. On the Master Node, select **Settings** > **Access Management** > **Local Users**. For details, see Local Users.

Note that when you assign multiple Roles to a given user, the Roles' permissions are additive: This user is granted a superset of all the permissions contained in all the assigned Roles.

By default, LogStream will log out a user upon a change in their assigned Roles. You can defeat this behavior at **General Settings > API Server Settings > Advanced > Logout on roles change.**

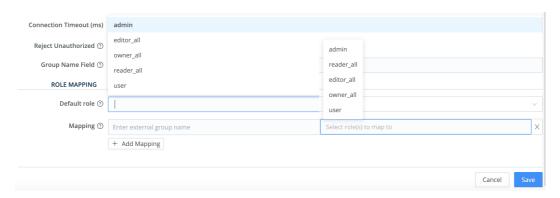
External Groups and LogStream Roles

You can map user groups from external identity providers (LDAP, Splunk, or OIDC) to LogStream Roles, as follows:

 On the Master Node, select Settings > Access Management > Authentication.

- 2. From the **Type** drop-down, select **LDAP**, **Splunk**, or **OpenID Connect**, according to your needs.
- 3. On the resulting **Authentication Settings** page, configure your identity provider's connection and other basics. (For configuration details, see the appropriate **Authentication** section.)
- 4. Under **Role Mapping**, first select a LogStream **Default role** to apply to external user groups that have no explicit LogStream mapping defined below.
- 5. Next, map external groups as you've configured them in your external identity provider (left field below) to LogStream Roles (right drop-down list below).
- 6. To map more user groups, click + Add Mapping.
- 7. When your configuration is complete, click **Save**.

Here's a composite showing the built-in Roles available on both the **Default Role** and the **Mapping** drop-downs:



Mapping external user groups to LogStream Roles

And here, we've set a conservative **Default Role** and one explicit **Mapping**:



External user groups mapped to LogStream Roles

Page 126 of 1257

Version Control

Tracking, backing up, and restoring configuration changes for single-instance and distributed deployments

Cribl LogStream integrates with Git clients and remote repositories to provide version control of LogStream's configuration. This integration offers backup and rollback for single-instance and distributed deployments.

These options are separate from the Git repo responsible for version control of Worker configurations, located on the Master Node in distributed deployments. We cover all these options and requirements below.

Git Installation (Local or Standalone/Single-Instance)

To verify that git is available, run:

git --version

The minimum version that LogStream requires is: **1.8.3.1.** If you don't have git installed, see the installation links here.

Git Required for Distributed Deployments

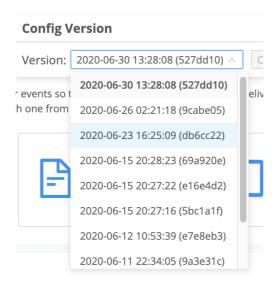
For distributed deployments, git must be installed and available locally on the host running the Master Node.

All configuration changes must be committed before they are deployed. The Master notifies Workers that a new configuration is available, and Workers pull the new configuration from the Master Node.

Reverting Commits

Once Git is installed, you can revert to a previous commit using the git CLI. You can also restore a Worker Group's previous commit using LogStream's UI:

Select the commit from the **Config Version** drop-down, as shown below. Then, in the resulting modal, verify the diff'ed configuration change and click **OK**.



Support For Remote Repositories

Git **remote** repositories are supported – but not required – for version control of all configuration changes. You can configure a Standalone Master Node with Git remote push capabilities through the LogStream CLI, or through the LogStream UI (via **Settings > Distributed Settings > Git Settings**).

To create a repo, see these tutorials:

- Setting Up a Repository (CLI instructions, host-agnostic, from Atlassian).
- Creating a New Repository (specific to GitHub's Web UI).
- Create a Repo (longer GitHub-specific tutorial, also covers committing changes).
 - i Currently, LogStream supports push and pull only against the master branch on each remote repo.

Several tutorial links and examples on this page point to GitHub, based on its wide adoption. The basic principles are the same for other Git repo providers, including private Git servers. GitHub's own UI and documentation periodically change, and linked tutorials' screenshots might differ from GitHub's current UI.

Remote Formats Supported

Remote URI schema patterns should match this regex:

 $(?:git|ssh|ftps?|file|https?|git@[-\w.]+):(\/\)?(.*?)(\.git\/?)?$

You can find a list of supported formats here.

For example:

- Local Gitservers: git://<host.xyz>:<port>/<user>/path/to/repo.git

Connecting to a Remote with a Personal Access Token over HTTPS (Recommended)

Cribl recommends connecting to a remote repo over HTTPS. The example below shows a token-based HTTPS connection to GitHub.

Example: Connecting to GitHub over HTTPS

- 1. Create a new GitHub repository.
 - i For best results, create a new empty repo, with no readme file and no commit history. This will prevent git push errors.

Note the user name and email with which you log into the repo provider.

- 2. Create a personal access token with repo scope.
- 3. Copy the token to your clipboard.
- 4. In Cribl LogStream, go to Settings > Distributed Settings > Git Settings.
- 5. Fill in the **Remote URL** field with your repo name, user name, and token (in place of a password). Use the format below, replacing both <username> placeholders with your user name on the repo provider:

https://<username>:<token>@github.com/<username>/<reponame>.git

For additional details, see GitHub's Creating a Personal Access Token tutorial.

▲ For GitHub repos specifically, use only personal access tokens in the **Remote URL** field. GitHub has announced that it will end support for plaintext passwords as of August 13, 2021.

Connecting to a Remote with SSH

You can set up SSH keys from the CLI, or upload keys via the UI. If you have a passphrase set, this functionality is available only through the CLI – see Encryption: Configuring Keys with the CLI. The example below outlines the UI steps.

Example: Connecting to GitHub with SSH

1. Create a new GitHub repository.

i For best results, create a new empty repo, with no readme file and no commit history. This will prevent git push errors.

Note the user name and email with which you log into the repo provider.

- 2. Add an SSH public key to your GitHub account.
- 3. In Cribl LogStream, go to **Settings > Distributed Settings > Git Settings**.
- 4. Fill in the remote repo URL and the SSH private key. In the example format below, replace <username> with your user name on the repo provider:

Remote URL: <protocol>://git@github.com:<username>/<reponame>.git

SSH private key: <ssh-private-key>

For GitHub specifically, the URL/protocol format must be:

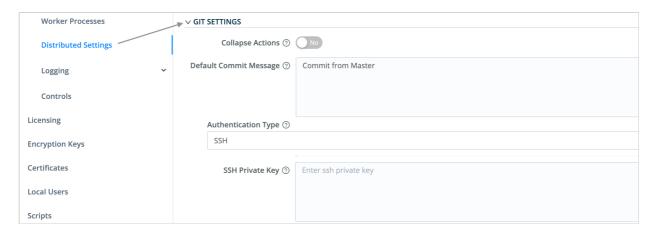
Remote URL: git@github.com:<user>/<reponame>.git

For example:

Remote URL: git@github.com:taylorswift/leadsheets.git

5. As the user running LogStream, run this command to add the GitHub keys to known_hosts: ssh-keyscan -H github.com >> ~/.ssh/known_hosts

For additional details, see GitHub's Connecting to GitHub with SSH tutorial.



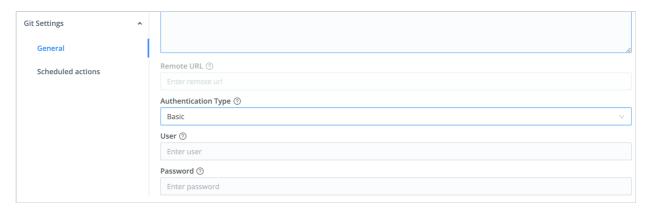
LogStream's Git settings

Additional Git Settings

On the **Git Settings** > **General** tab, you can change the **Authentication Type** from its **SSH** default to **Basic** authentication. This displays two additional fields:

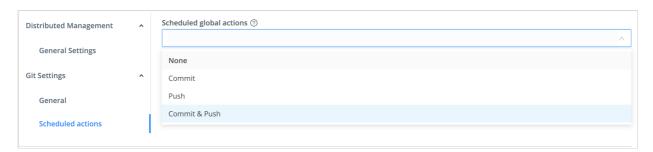
• User: Username on the repo.

• Password: Authentication password (e.g., a GitHub personal access token).



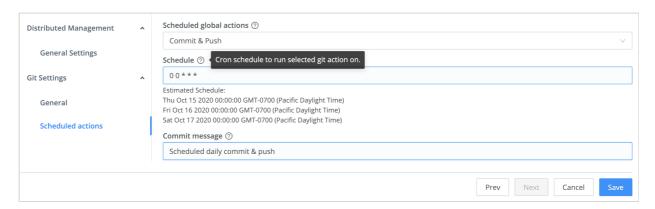
Git Authentication Type settings

On the **Git Settings** > **Scheduled Actions** tab, you can schedule a **Commit, Push**, or **Commit & Push** action to occur on a predefined interval.



Git Scheduled Actions selection

For the selected action type, you can define a [cron schedule] (cron schedule), and a commit message distinct from the **General** tab's **Default Commit Message**. Then click **Save**.

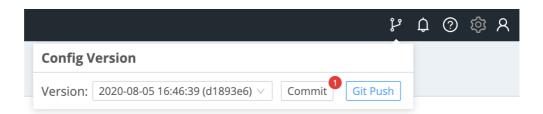


Saving a Git Scheduled Action

You can schedule only one type of action. To swap to a different type, select it from the **Scheduled global actions** drop-down, and resave. To turn off scheduled Git commands, select **None** from the drop-down, and resave.

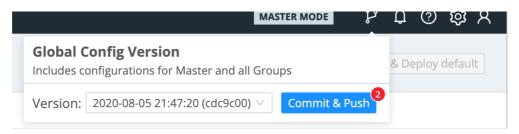
Pushing to a Remote Repo

Once you've configured a remote, a **Git Push** button appears in the Version Control overlay.



Git Push button

If you enabled the **Git Settings** > **Collapse Actions** option, you will instead see a combined **Commit & Push** button in the overlay.



Git combined actions button

Troubleshooting Push Errors

This section anticipates common errors you might see in LogStream's UI, or in the git CLI, when pushing a commit.

Failed to Push Some Refs

Your first push to a remote repo might fail with one of several failed to push some refs errors.

As a first step in debugging these errors, edit the \$CRIBL_HOME/.git/config file to make sure that its name and email key values match the credentials you've set on your repo provider or git server.

Also make sure that the remote "origin" key value matches the remote you set when you connected to the remote repo. This example shows all three keys, with placeholder values:

```
[user]
   name = <your-login-name>
   email = <email@example.com>
[remote "origin"]
   url = https://<user-name>:<token>@github.com/<username>/<repo-name>
```

Next, verify the remote repo from the command line, as follows:

```
cd $CRIBL_HOME/.git
git remote -v
```

In response, git should echo your configured remote twice – once for fetch and once for push operations.

If all of the above settings are correct, the push is very likely blocking because the remote repo has some commit history, or was simply created with a readme.md file. For command-line instructions to remedy this – by syncing your local repo to its remote – see GitHub's Dealing with Non-Fast-Forward Errors topic.

Large Files Detected

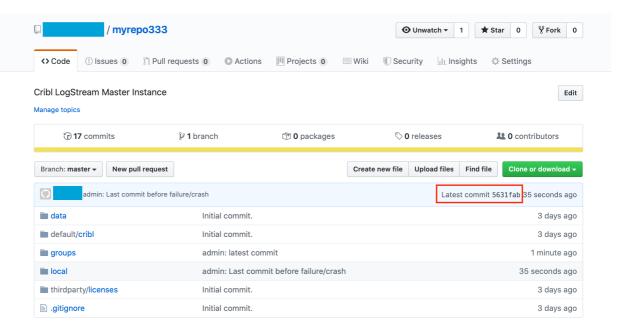
A push command might also trigger "large file" warnings or, more seriously, errors of this form (CLI/GitHub example):

Cribl recommends adding such large files to .gitignore , to exclude them from subsequent push commands. As the above examples show, typical culprits are large .csv or .mmdb lookup files. A simple option is to place these files in a \$CRIBL_HOME subdirectory that's already listed in .gitignore – for details, see Managing Large Lookups.

Other available workarounds include staging such files **outside** \$CRIBL_HOME, or using plugins to accommodate the large files. For GitHub-specific options, see Working with Large Files.

Restoring Master from a Remote Repo

If a remote repo is configured and has the latest known good Master configuration, this section outlines the general steps to restore the config from that repo.



Restoring from remote repo

Let's assume that the entire \$CRIBL_HOME directory of the Master is corrupted, or you're starting from scratch. Let's also assume that the remote repo has the form:

git@github.com:<username>/<reponame>.git.

- 1. **Important**: In a directory of choice, untar the **same Cribl LogStream version** that you're trying to restore, but do not start it.
- 2. If you are using SSH key authentication, specify the key using the following command:

```
GIT_SSH_COMMAND='ssh -i .key -o IdentitiesOnly=yes' git fetch origin
```

3. Ensure that you have proper access to the remote repo:

```
# git ls-remote git@github.com:/.git
56331fabb4822eaec4ca0ffd008d6e9974c1e419f HEAD
5631fabb4822eaec4ca0ffd008d6e9974c1e419f refs/heads/master
```

4. Change directory into \$CRIBL_HOME and initialize git:

```
# git init
```

5. Next, add/configure the remote:

```
# git remote add origin git@github.com:<username>/<reponame>.git
```

6. Now set up your local branch to exactly match the remote branch:

```
# git fetch origin
# git reset --hard origin/master
```

7. Finally, to confirm that the commits match, run this command while in \$CRIBL_HOME . Note the commit hash:

```
# git show --abbrev-commit commit 5631fab (HEAD \rightarrow master, origin/master)
```

```
Author: First Last
Date: Fri Jan 31 10:16:07 2020 -0500

admin: Last commit before failure/crash
```

That last step above pulls in all the latest configs from the remote repo, and you should be able to start the Master as normal. Once up and running, Workers should start checking in after about 60 seconds.

∆ Verify cribl.secret

The cribl.secret file-located at \$CRIBL_HOME/local/cribl/auth/cribl.secret - contains the secret key that is used to encrypt sensitive settings on configuration files (e.g., AWS Secret Access Key, etc.). Make sure this file is properly restored on the new Master, because it is required to make encrypted conf file settings usable again.

.gitignore File

A .gitignore file specifies files that git should ignore when tracking changes. Each line specifies a pattern, which should match a file path to be ignored. Cribl LogStream ships with a .gitgnore file containing a number of patterns/rules, under a section of the file labeled CRIBL SECTION.

```
.gitignore
# Do NOT REMOVE CRIBL and CUSTOM header lines!
# DO NOT REMOVE rules under the CRIBL section as they may be reintroduced on update.
# You can ONLY comment out rules in the CRIBL section.
# You can add new rules in the CUSTOM section.
### CRIBL SECTION -- DO NOT REMOVE ###
default/ui/**
default/data/ui/**
bin/**
log/**
pid/**
data/uploads/**
diag/**
**/state/**
### CUSTOM SECTION -- DO NOT REMOVE ###
<User defined patterns/rules go here>
```

CRIBL Section

☐ Do Not Remove CRIBL SECTION or CUSTOM SECTION Headers

The CRIBL SECTION is used by Cribl LogStream to define default patterns/rules that ship with every version. Do **not** add or remove any of the lines here, because Chuck Norris will easily find you!

Maslow's theory of higher needs does not apply to Chuck Norris. He has only two needs: killing people and finding people to kill. Seriously, do not remove them, as they will be overwritten on the next update. The only modifications that will survive updates are commented lines.

CUSTOM Section

User-defined, custom patterns/rules can be **safely defined** under the CUSTOM SECTION . Cribl LogStream will **not** modify the contents of CUSTOM SECTION .

Good candidates to add here include large lookup files – especially large binary database files. See Troubleshooting: Large Files Detected, above.

Files skipped with .gitignore

If you have files that are skipped with .gitgnore, you will need to back them up and restore them via means other than Git. E.g., you can periodically copy/rsync them to a backup destination, and then restore them to their original locations after you complete the steps above.

Persistent Queues

Persistent queuing (PQ) is a feature that helps minimize data loss if a downstream receiver is unreachable. Durability is provided by writing data to disk for the duration of the outage, and forwarding it upon recovery.

PQs are implemented on the outbound side, meaning that each Source can take advantage of a Destination's queue.

How Does Persistent Queueing Work

Each LogStream output has an in-memory queue that helps it absorb temporary imbalances between inbound and outbound data rates. E.g., if there is an inbound burst of data, the output will store events in the queue, and output them at the rate that the receiver can sync (as opposed to blocking or dropping them). Only when this queue is full will the output impose backpressure upstream.

Backpressure behavior can be configured to either **block** or **drop**. In block mode, the output will refuse to accept new data until the receiver is ready. The system will back propagate block "signals" all the way back to the sender (assuming it supports backpressure, too). In drop behavior, the output will discard new events until the receiver is ready.

In some environments, the in-memory queues and their block/drop behavior are acceptable. Persistent queues serve environments where more durability is required (e.g., outages last longer than memory queues can sustain), or where upstream senders do not support backpressure (e.g., ephemeral/network senders),

Engaging persistent queues in these scenarios can help minimize data loss. Once the in-memory queue is full, the LogStream output will write its data to disk. Then, when the receiver is ready, the output will start draining the queues in FIFO (first in, first out) fashion.

Persistent Queue Details and Constraints

Persistent queues are:

- Available at the output side (i.e., after processing).
- Engaged only when all of the receivers of that output exert blocking.
- Drained when at least one receiver can accept data.
- Not infinite in size. I.e., if data cannot be delivered out, you might run out of disk space.
- Not able to fully protect in cases of application failure. E.g., in-memory data might get lost if a crash occurs.
- Not able to protect in cases of hardware failure. E.g., disk failure, corruption, or machine/host loss.

Persistent Queue Support

The following LogStream Destinations support Persistent Queuing:

- Splunk Single Instance
- Splunk Load Balanced
- Splunk HEC
- Kinesis
- Cloudwatch Logs
- SQS
- Azure Monitor Logs
- Azure Event Hubs
- StatsD
- StatsD Extended
- Graphite
- TCP JSON
- Syslog
- Elasticsearch
- Honeycomb
- InfluxDB
- Wavefront
- SignalFx

Configuring Persistent Queueing

Persistent Queueing is configured individually for each output that supports it. To enable persistent queueing, go to the output's (Destination's) configuration

page and set the **Backpressure Behavior** control to **Persistent Queueing**. This exposes the following additional controls:

- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Sufficient disk space is required for queuing to operate properly. You configure the minimum disk space in **Settings > General Settings > Limits > Min Free Disk Space**. If available disk space falls below this threshold, LogStream will stop maintaining persistent queues, and data loss will begin. The default is 5GB. Be sure to set this on your worker nodes rather than on the master node when in distributed mode.

Securing

You can secure Cribl LogStream's API and UI access by configuring SSL. To do so, you can use your own certs and private keys, or you can generate a pair with OpenSSL, as shown here:

```
openssl req -nodes -new -x509 -newkey rsa:2048 -keyout myKey.pem -out myCert.pem -days 420
```

This command will generate both a self-signed cert (certified for 420 days), and an unencrypted, 2048-bit RSA private key.

In the LogStream UI, you can configure the cert via **Settings** > **Certificates**. You can configure the key via:

- Settings > Encryption Keys (single-instance deployments), or
- Worker Groups > <group-name> > System Settings > Encryption Keys (distributed deployments).

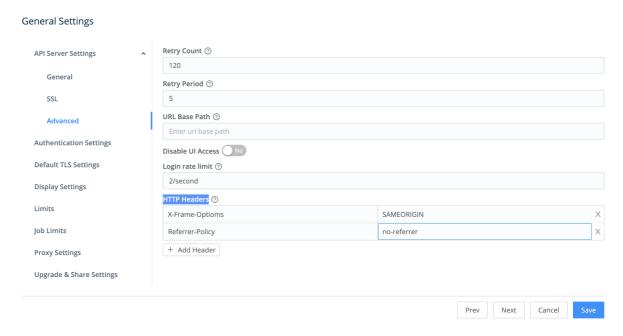
Alternatively, you can edit the local/cribl.yml file's api section to directly set the privKeyPath and certPath attributes. For example:

```
cribl.yml

api:
  host: 0.0.0.0
  port: 9000
  disabled: false
  ssl:
    disabled: false
    privKeyPath: /path/to/myKey.pem
    certPath: /path/to/myCert.pem
...
```

Custom HTTP Headers

You can encode custom, security-related HTTP headers, as needed. As shown in the examples below, you specify these at **Settings > General > API Server Settings > Advanced > HTTP Headers**. Click **+ Add Header** to display extra rows for new key-value pairs.



Custom HTTP headers

TLS Settings and Traffic Types

This table shows TLS client/server pairs, and encryption defaults, per traffic type.

Traffic Type	TLS Client	TLS Server	Encryption	Cert Auth	CN* Check
UI	Browser	Cribl LogStream	Default disabled	Default disabled	Default disabled
API	Worker	Master	Default disabled	Default disabled	Default disabled
Worker-to- Master	Worker	Master	Default disabled	Default disabled	Default disabled
Data	Any data sender	Cribl LogStream (Source)	Default disabled	Default disabled	Default disabled
Data	Cribl LogStream (Destination)	Any data receiver	Default disabled	Default disabled	Default disabled
Authentication					
• Local	Browser	Cribl LogStream	Default Disabled	N/A	N/A
• LDAP	Cribl LogStream	LDAP Provider	Custom	N/A	Default Disabled

• Splunk	Cribl LogStream	Splunk Search Head	Default Enabled	N/A	Default Disabled
• OIDC†/Okta	Browser and Cribl LogStream	Okta	Default Enabled	N/A	Enabled (Browser)
• OIDC/Google	Browser and Cribl LogStream	Google	Default Enabled	N/A	Enabled (Browser)

^{*} Common name

† OpenID Connect

You can configure advanced, system-wide TLS settings for versions, cipher lists, and ECDH Curve names via Settings > System > General Settings > Default TLS Settings.

CA Certificates and Environment Variables

Where LogStream Sources and Destinations support TLS, each Source's or Destination's configuration provides a **CA Certificate Path** field where you can point to corresponding Certificate Authority (CA) .pem file(s). However, you can also use environment variables to manage CAs globally. Here are some common scenarios:

1. How do I add a set of trusted root CAs to the list of trusted CAs that LogStream trusts?

Set this environment variable in each Worker's environment (e.g., in its systemd unit file): NODE_EXTRA_CA_CERTS=/path/to/file_with_certs.pem . For details, see nodejs docs.

2. How do I make LogStream trust all TLS certificates presented by any server it connects to?

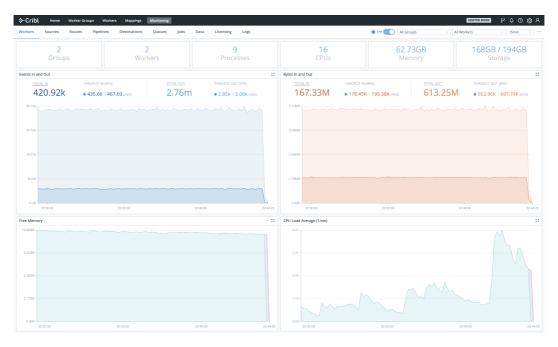
Set this environment variable: NODE_TLS_REJECT_UNAUTHORIZED=0 - for details, see nodejs docs.

Monitoring

To get an operational view of a Cribl LogStream deployment, you can consult the following resources.

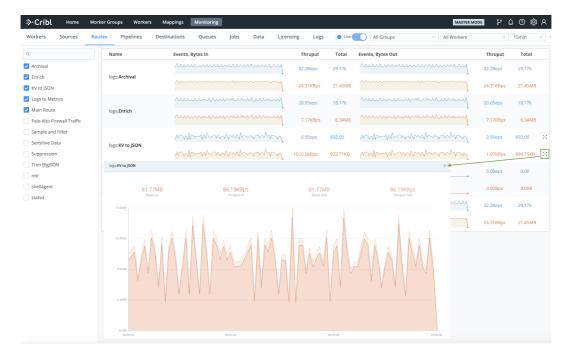
Monitoring Page

Select **Monitoring** from the top menu. This exposes information about traffic in and out of the system, as well as collection jobs and tasks. It tracks events, bytes, splits by data fields over time, and broader system metrics. Coverage is limited to the previous 24 hours. (Byte-related charts show the uncompressed size of processed data.)



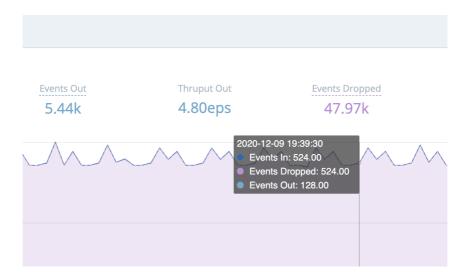
Monitoring page

Dense displays are condensed to sparklines for legibility. Hover over the right edge to display Maximize buttons that you can click to zoom these up to detailed graphs.



Sparklines and fly-out

You can hover over an expanded graph fly-out to display further details.



Throughput details

Internal Logs and Metrics

Select **Logs** from the **Monitoring** submenu. LogStream's **internal logs** and **internal metrics** provide comprehensive information about an instance's status/health, inputs, outputs, Pipelines, Routes, Functions, and traffic.

Health Endpoint

Query this endpoint on any instance to check the instance's health. (Details below.)

Types of Logs

LogStream provides the following log types, by originating process:

- API Server Logs These logs are emitted primarily by the API/main process. They correspond to the top-level cribl.log that shows up on the Diag page. Filesystem location: \$CRIBL_HOME/log/cribl.log
- Worker Process(es) Logs These logs are emitted by all the worker processes, and are very common in standalone instances or Worker Nodes. Filesystem location: \$CRIBL_HOME/log/worker/N/cribl.log
- Worker Group Logs These logs are emitted by all processes that help a
 Master Node configure Worker Groups. Filesystem location:
 \$CRIBL_HOME/log/group/GROUPNAME/cribl.log

LogStream rotates logs every 5 MB, keeping the most recent 5 logs. In a distributed deployment, all Workers forward their metrics to the Master Node, which then consolidates them to provide a deployment-wide view.

Forward Logs and Metrics Externally

LogStream supports forwarding internal logs and metrics to your preferred external monitoring solution. To send out internal data, go to **Data > Sources** and enable the **Cribl Internal** Source.

This will send all cribllog logs and internal metrics down through Routes and Pipelines, just like another data source. Both logs and metrics will have a field called source, set to the value cribl, which you can use in Route filters.

For recommendations about useful Cribl metrics to monitor, see Internal Metrics.

i CriblMetrics Override

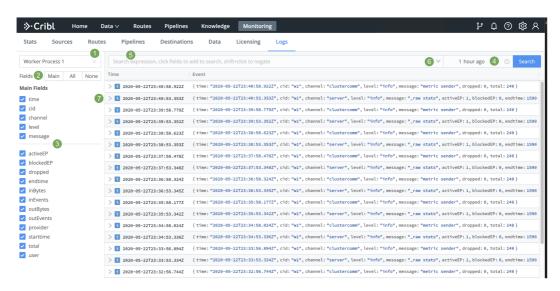
The **Disable field metrics** setting (in **Settings > System > General Settings > Limits**) applies only to metrics sent to the Master Node.
When the **Cribl Internal** Source is enabled, LogStream ignores this

Disable field metrics setting, and full-fidelity data will flow down the Routes.

Search Internal Logs

LogStream exists because logs are great and wonderful things! Using its **Monitoring > Logs** page, you can search all LogStream's internal logs at once – from a single location, for both Master and Worker Nodes. This enables you to query across all internal logs for strings of interest.

The labels on this screenshot highlight the key controls you can use (see the descriptions below):



Logs page (controls highlighted)

- 1. Log file selector: Choose the Node to view. In a Distributed Deployment, this list will be hierarchical, with Workers displayed inside their Master.
- 2. **Fields selector**: Click the **Main | All | None** toggles to quickly select or deselect multiple check boxes below.
- 3. **Fields**: Select or deselect these check boxes to determine which columns are displayed in the Results pane at right. (The upper **Main Fields** group will contain data for *every* event; other fields might not display data for all events.)
- 4. **Time range selector**: Select a standard or custom range of log data to display.

5. **Search box**: To limit the displayed results, enter a JavaScript expression here. An expression must evaluate to truthy to return results. You can press **Shift+Enter** to insert a newline.

Typeahead assist is available for expression completion:



Click a field in any event to add it to a query:



Click other fields to append them to a query:



Shift+click to negate a field:



- To modify the depth of information that is originally input to the Logs page, see Logging Settings.
- 6. Click the Search box's history arrow (right side) to retrieve recent queries:



7. The Results pane displays most-recent events first. Each event's icon is color-coded to match the event's severity level.

Click individual log events to unwrap an expanded view of their fields:



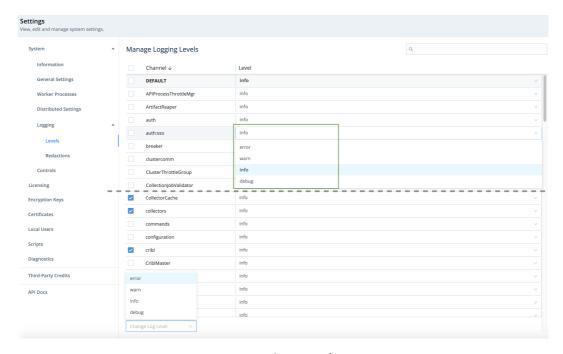
Logging Settings

Through LogStream's System Settings, you can adjust the level (verbosity) of internal logging data processed, per logging channel. You can also redact fields in customized ways.

Change Logging Levels

Select Settings > System > Logging > Levels to open the Manage Logging Levels page. Here, you can:

- Modify one channel by clicking its Level column. In the resulting dropdown, you can set a verbosity level ranging from error up to debug. (Top of composite screenshot below.)
- Modify multiple channels by selecting their check boxes, then clicking the
 Change log level drop-down at the bottom of the page. (Bottom of
 composite screenshot below.) You can select all channels at once by
 clicking the top check box. You can search for channels at top right.



Manage Logging Levels screen

Change Logging Redactions

Select **Settings > System > Logging > Redactions**: to open the **Redact Internal Log Fields** page. Here, you can customize the redaction of sensitive, verbose, or just ugly data within LogStream's internal logs.



Redact Internal Log Fields page

It's easiest to understand this page's fields from bottom to top:

- Default fields: LogStream always redacts these fields. You can't modify this list.
- Additional fields: Type or paste in the names of other fields you want to redact. Use a tab or hard return to confirm each entry.
- **Custom redact string:** Unless this field is empty, it defines a literal string that will override LogStream's default redaction pattern, explained below.

Default Redact String

By default, LogStream transforms this page's selected fields by applying the following redaction pattern:

- Echo the field value's first two characters.
- Replace all intermediate characters with a literal ... ellipsis.
- Echo the value's last two characters.

Anything you enter in the **Custom redact string** field will override this default ?? ... ?? pattern.

Health Endpoint

Each LogStream instance exposes a health endpoint – typically used in conjunction with a Load Balancer – that you can use to make operational decisions.

Health Check Endpoint	Healthy Response
<pre>curl http(s)://<host>:<port>/api/v1/health</port></host></pre>	{"status":"healthy"}

Internal Metrics

When sending LogStream metrics to a metric system of analysis, such as InfluxDB, Splunk or Elasticsearch, some metrics are particularly valuable. You can use these metrics to set up alerts when a Worker Node is having a problem, a Node is down, a Destination is down, a Source stops providing incoming data, etc.

LogStream reports its internal metrics within the LogStream UI (in the same way that it reports internal logs at Monitoring > Logs). To expose metrics for capture or routing, enable the Cribl Internal Source > CriblMetrics section.

By default, LogStream generates internal metrics every 2 seconds. To consume metrics at longer intervals, you can use or adapt the <code>cribl-metrics_rollup</code> Pipeline that ships with LogStream. Attach it to your <code>Cribl Internal</code> Source as a <code>pre-processing Pipeline</code>. The Pipeline's <code>Rollup Metrics</code> Function has a default <code>Time Window</code> of 30 seconds, which you can adjust to a different granularity as needed.

You can also use our public endpoints to automate monitoring using your own external tools.

Total Throughput Metrics

Five important metrics below are prefixed with total. These power the top of LogStream's **Monitoring** dashboard. The first two report on Sources, the remainder on Destinations.

- total.in_bytes
- total.in_events
- total.out_events
- total.out_bytes
- total.dropped_events (new in LogStream 2.4) helpful for discovering situations such as: you've disabled a Destination without noticing.

Interpreting Total Metrics

These total. metrics' values could reflect LogStream's health, but could also report low activity simply due to the Source system. E.g., logs from a store site will be low at low buying periods.

Also, despite the total. prefix, these metrics are each specific to the Worker Process that's generating them.

You can distinguish unique metrics by their #input=<id> dimension. For example, total.in_events|#input=foo would be one unique metric; total.in_events|#input=bar would be another.

System Health Metrics

Five specific metrics are most valuable for monitoring system health. The first two are LogStream composite metrics; the remaining three report on your hardware or VM infrastructure.

- health.inputs
- health.outputs see the JSON Examples below for both health. metrics.
- system.load_avg
- system.free_mem
- system.disk_used valuable if you know your disk size, especially for monitoring Persistent Queues. Here, a 0 value typically indicates that the disk-usage data provider has not yet provided the metric with data. (Getting the first value should take about one minute.)

All of the above metrics take these three values:

- 0 = green = healthy.
- 1 = yellow = warning.
- 2 = red = trouble.

Health Inputs/Outputs JSON Examples

The health.inputs metrics are reported per Source, and the health.outputs metrics per Destination. The health.inputs example below has two configured Sources, and two LogStream-internal inputs. The health.outputs example includes the built-in devnull Destination, and six user-configured Destinations.

Given all the 0 values here, everything is in good shape!

```
"health.inputs": [
    { "model": { "ci": "http:http", "input": "http:http" }}, "val": 0},
    { "model": { "ci": "cribl:CriblLogs", "input": "cribl:CriblLogs" }}, "val": 0},
    { "model": { "ci": "cribl:CriblMetrics", "input": "cribl:CriblMetrics" }}, "val": 0},
    { "model": { "ci": "datagen:DatagenWeblog", "input": "datagen:DatagenWeblog" }}, "val": 0
    ],
    "health.outputs": [
    { "model": { "output": "devnull:devnull" }}, "val": 0},
    { "model": { "output": "router:MyOut1" }}, "val": 0},
    { "model": { "output": "router:MyOut2" }}, "val": 0},
    { "model": { "output": "tcpjson:MyTcpOut2" }}, "val": 0},
    { "model": { "output": "router:MyOut3" }}, "val": 0},
    { "model": { "output": "router:MyOut3" }}, "val": 0},
    { "model": { "output": "router:MyOut4" }}, "val": 0},
    { "model": { "output": "router:MyOut4" }}, "val": 0},
}
```

Persistent Queue Metrics

Five metrics below are valuable for monitoring Persistent Queues' behavior:

• pq.queue_size

- pq.in_bytes
- pq.in_events
- pq.out_events
- pq.out_bytes

These are aggregate metrics. But you can distinguish unique metrics per queue Destination, using the #output=<id> dimension. For example, pq.out_events|#output=kafka would be one unique metric; pq.out_events|#output=camus would be another.

Other Metrics Endpoints and Dimensions

Below are basics on using the /system/metrics endpoint, the /system/info endpoint, and the cribl_wp dimension.

/system/metrics Endpoint

/system/metrics is LogStream's primary public metrics endpoint, which returns most internal metrics. Note that many of these retrieved metrics report configuration only, not runtime behavior. For details, see our API Docs.

/system/info Endpoint

/system/info generates the JSON displayed in the LogStream UI at Settings > Diagnostics > System Info. Its two most useful properties are loadavg and memory.

loadavg Example

```
"loadavg": [1.39599609375, 1.22265625, 1.31494140625],
```

This property is an array containing the 1-, 5-, and 15-minute load averages at the UNIX OS level. (On Windows, the return value is always [0, 0, 0].) For details, see the Node.js os.loadavg() documentation.

memory Example

```
"memory": { "free": 936206336, "total": 16672968704 },
```

Divide total / free to monitor memory pressure. If the result exceeds 90%, this indicates a risky situation: you're running out of memory.

cpus Alternative

The cpus metric returns an array of CPU/memory key-value pairs. This provides an alternative way of determining loadavg, but it requires you to query all your CPUs individually, and then average. In the example below, Idle = user + nice + sys.

```
"cpus": [{ "times": {
    "user": 19881260,
    "nice": 39370,
    "sys": 5250130,
    "idle": 76088790,
```

cribl_wp Metric Dimension

cribl_wp is a useful dimension that identifies the Worker Process that processed each event.

Upgrading

This page outlines how to upgrade Cribl LogStream's Single-Instance or Distributed Deployment packages along one of the following supported upgrade paths:

- v2.x ==> v2.x
- v1.7.x/v2.0.x ==> v2.x.x
- v1.6.x or below ==> v1.7.x ==> v2.x.x
 - △ See notes on Upgrading to LogStream 2.3 below.

LogStream does **not** support direct upgrades from a Beta to a GA version. To get the GA version running, you must perform a new install.

Standalone/Single-Instance

This path requires upgrading only the single/standalone node:

- 1. Stop Cribl LogStream.
- 2. Uncompress the new version on top of the old one.

On some Linux systems, tar might complain with: cribl/bin/cribl: Cannot open: File exists. In this case, please remove the cribl/bin/cribl directory if it's empty, and untar again. If you have custom functions in cribl/bin/cribl, please move them under \$CRIBL_HOME/local/cribl/functions/ before untarring again.

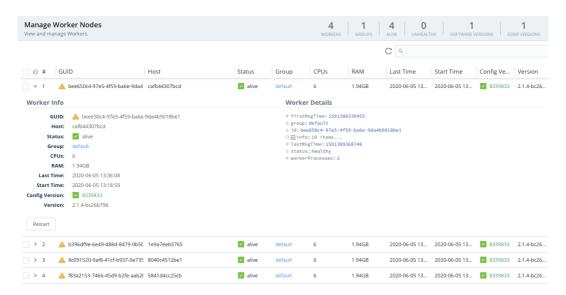
3. Restart LogStream.

Distributed Deployment

For a distributed deployment, the order of upgrade is: Upgrade first the Master Node, then upgrade the Worker Nodes, then commit and deploy the changes on the Master.

Upgrade the Master Node

- 1. Commit and deploy your desired last version. (This will be your most recent checkpoint.)
 - Optionally, git push to your configured remote repo.
- 2. Stop Cribl LogStream.
 - Optional but recommended: Back up the entire \$CRIBL_HOME directory.
 - Optional: Check that the Worker Nodes are still functioning as expected. In absence of the Master Node, they should continue to work with their last deployed configurations.
- 3. Uncompress the new LogStream version on top of the old one.
- 4. Restart LogStream and log back in.
- 5. Wait for all the Worker Nodes to report to the Master, and ensure that they are correctly reporting the last committed configuration version.
- i Workers' UI will not be available until the Worker version has been upgraded to match the version on the Master. Errors like those below will appear until the Worker nodes are upgraded.



Worker Node version mismatch

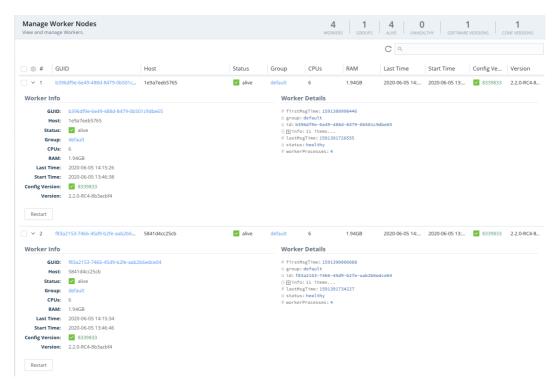
Upgrade the Worker Nodes

These are the same basic steps as when upgrading a Standalone Instance, above:

- 1. Stop Cribl LogStream on each Worker Node.
- 2. Uncompress the new version on top of the old one.
- 3. Restart LogStream.

Commit and Deploy Changes on the Master Node

- 1. Ensure that newly upgraded Worker Nodes report to the Master with their new software version.
- 2. Commit and deploy the newly updated configuration **only after all** Workers have upgraded.



Post-2.1.4 upgrade to 2.2

Upgrading to LogStream 2.3

As of version 2.3, LogStream Free and One licenses are permanent, but they enforce certain restrictions that especially affect distributed deployments:

- Even if you have more than one Worker Group defined, only one Worker Group will be visible and usable.
 - This will be the first Group listed in \$CRIBL_HOME/local/cribl/groups.yml - typically, the default Group. You can edit groups.yml to move the desired Group to the top.
- Your cluster will be limited to 10 Worker Processes across all Worker Nodes.
 - LogStream will balance (or rebalance) these Processes as evenly as possible across the Worker Nodes.
- Authentication will fall back to local authorization. You will not be able to authenticate via Splunk, LDAP, or SSO/OpenID.
- **Git Push** to remote repos will not be supported through the product.
 - ⚠ If you are upgrading LogStream Free or LogStream One from version 2.2.x or lower, these changes might require you to adjust your existing configuration and/or workflows.

See Licensing for details on all current license options.

As of LogStream 2.3, licenses no longer need to be deployed directly to Worker Groups. The Master will push license information down to Worker Groups as part of the heartbeat.

Splunk App Package Upgrade Steps

△ See Deprecation note for v.2.1.

Follow these steps to upgrade from v.1.7, or higher, of the Cribl App for Splunk:

- 1. Stop Splunk.
- 2. Untar/unzip the new app version on top of the old one.

On some Linux systems, tar might complain with: cribl/bin/cribl: Cannot open: File exists. In this case, please remove the cribl/bin/cribl directory if it's empty, and untar again. If you have
custom functions in cribl/bin/cribl , please move them under
\$CRIBL_HOME/local/cribl/functions/ before untarring again.

3. Restart Splunk.

Upgrading from Splunk App v.1.6 (or Lower)

As of v.1.7, contrary to prior versions, Cribl's Splunk App package defaults to Search Head Mode. If you have v.1.6 or earlier deployed as a Heavy Forwarder app, upgrading requires an extra step to restore this setting:

- 1. Stop Splunk.
- 2. Untar/unzip the new app version on top of the old one.
- 3. Convert to HF mode by running:
 \$\$PLUNK_HOME/etc/apps/cribl/bin/cribld mode-hwf
- 4. Restart Splunk.

Uninstalling

Uninstalling the Standalone Version

- Stop Cribl LogStream (stopping the main process).
- Back up necessary configurations/data.
- Remove the directory where Cribl LogStream is installed.

In a distributed deployment, repeat the above steps for the Master instance and all Worker instances.

Uninstalling the Splunk App Version

- Stop Splunk.
- Back up necessary configurations/data.
- Remove the Cribl App in \$SPLUNK_HOME/etc/apps.
- Remove the Cribl module in \$SPLUNK_HOME/etc/modules/cribl (some versions).

WORKING WITH DATA

Event Model

All data processing in Cribl LogStream is based on discrete data entities commonly known as **events**. An event is defined as a collection of key-value pairs (fields). Some Sources deliver events directly, while others might deliver bytestreams that need to be broken up by Event Breakers. Events travel from a Source through Pipelines' Functions, and on to Destinations.

The internal representation of a Cribl LogStream event is as follows:

```
Cribl LogStream Event Model
```

```
{
   "_raw": "<body of non-JSON parse-able event>",
   "_time": "<timestamp in UNIX epoch format>",
   "_inputId": "<Id/Name of Source that delivered the event>",
   "_other1": "<Internal field1>",
   "_other2": "<Internal field2>",
   "_otherN": "<Internal fieldN>",
   "key1": "<value1>",
   "key2": "<value2>",
   "keyN": "<valueN>",
   "...": "..."
}
```

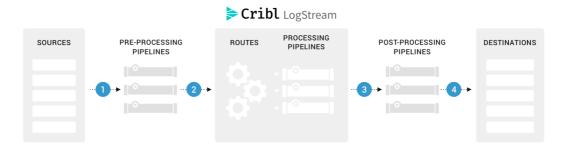
Some notes about these representative fields:

- Fields that start with a double-underscore are known as internal fields, and each Source can add one or many to each event. For example, Syslog adds both a __inputId and a __srcIpPort field. Internal fields are used only within Cribl LogStream, and are not passed down to Destinations.
- Upon arrival from a Source, if an event cannot be JSON-parsed, all of its content will be assigned to _raw .
- If a timestamp is not configured to be extracted, the current time (in UNIX epoch format) will be assigned to _time .

Using Capture

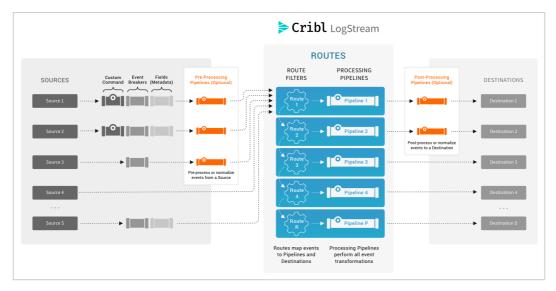
One way to see what an event looks like as it travels through the system is to use the **Capture** feature. While in **Preview** (right pane):

- 1. Click **Start a Capture**.
- 2. In the resulting modal, enter a **Filter expression** to narrow down the events of interest.
- 3. Click **Capture...** and (optionally) change the default Time and/or Event limits.
- 4. Select the desired Where to capture option. There are four options:
- 1. Before the pre-processing Pipeline Capture events right after they're delivered by the respective Input.
- 2. Before the Routes Capture events right after the pre-processing Pipeline, before they go down the Routes.
- 3. Before the post-processing Pipeline Capture events right after the Processing Pipeline that actually handled them, before any post-processing Pipeline.
- **4. Before the Destination** Capture events right after the post-processing Pipeline, before they go out to the configured Destination.



Event Processing Order

The expanded schematic below shows how all events in the Cribl LogStream ecosystem are processed linearly, from left to right.



LogStream in great detail

Here are the stages of event processing:

- Sources: Data arrives from your choice of external providers. (LogStream supports Splunk, HTTP/S, Elastic Beats, Amazon Kinesis/S3/SQS, Kafka, TCP raw or JSON, and many others.)
- 2. Custom command: Optionally, you can pass this input's data to an external command before the data continues downstream. This external command will consume the data via stdin, will process it and send its output via stdout.
- 3. Event Breakers can, optionally, break up incoming bytestreams into discrete events.
- 4. Fields/Metadata: Optionally, you can add these enrichments to each incoming event. You add fields by specifying key/value pairs, per Source, in a format similar to LogStream's Eval function. Each key defines a field name, and each value is a JavaScript expression (or constant) used to compute the field's value.

- 5. Pre-processing Pipeline: Optionally, you can use a single Pipeline to condition (normalize) data from this input before the data reaches the Routes.
- 6. Routes map incoming events to Processing Pipelines and Destinations. A Route can accept data from multiple Sources, but each Route can be associated with only one Pipeline and one Destination.
- 7. Processing Pipelines perform all event transformations. Within a Pipeline, you define these transformations as a linear series of Functions.

 A Function is an atomic piece of JavaScript code invoked on each event.
- 8. Post-processing Pipeline: Optionally, you can append a Pipeline a to condition (normalize) data from each Processing Pipeline before the data reaches its Destination.
- 9. Destinations: Each Route/Pipeline combination forwards processed data to your choice of streaming or storage Destination. (LogStream supports Splunk, Syslog, Elastic, Kafka/Confluent, Amazon S3, Filesystem/NFS, and many other options.)

i Pipelines Everywhere

All Pipelines have the same basic internal structure – they're a series of Functions. The three Pipeline types identified above differ only in their position in the system.

Routes

What Are Routes

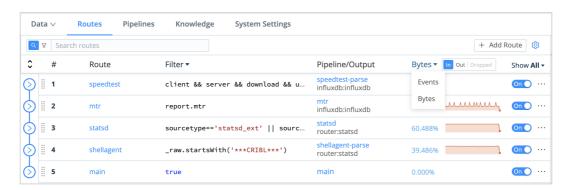
Before incoming events are transformed by a processing Pipeline, Cribl LogStream uses a set of filters to first select a **subset** of events to deliver to the correct Pipeline. This selection is made via Routes.

How Do Routes Work

Routes apply filter expressions on incoming events to send matching results to the appropriate Pipeline. Filters are JavaScript-syntax-compatible expressions that are configured with each Route. Examples are: true,

i There can be multiple Routes in the system, but each Route can be associated with only one Pipeline.

Routes are evaluated in their display order, top->down. The stats shown in the **Bytes/Events** (toggle) column are for the most-recent 15 minutes.

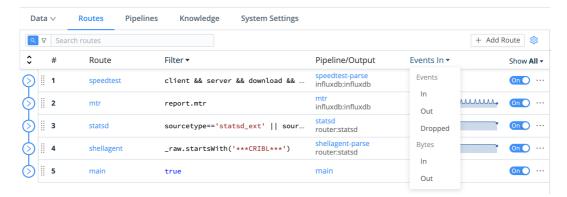


Routes and bytes

In the example above, incoming events will be evaluated first against the Route named **speedtest**, then against **mtr**, then against **statsd**, and so on. At the end, the **main** Route serves as a catch-all for any event that does not match any of the other Routes.

Above, note the selectors to toggle between displaying Events versus Bytes, and to display In versus Out.

When you condense the Routes page to a narrower viewport, LogStream consolidates the In/Out/Dropped selectors onto an expanded Bytes/Events drop-down menu, as shown below.



Routes and events (combined menu)

Managing the Routes Page

To apply a Route before another, simply drag it vertically. Use the sliders to turn Routes **On/Off** inline, as necessary, to facilitate development and debugging.

You can press the] (right-bracket) shortcut key to toggle between the Preview pane and the expanded Routes display shown above. (This works when no field has focus.)

Output Destination

You can configure each Route with an output Destination that denotes where to send events after they're processed by the Pipeline.

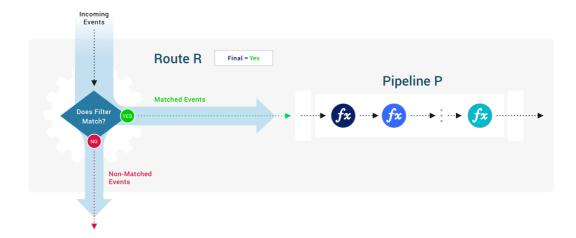
The Final Toggle

When an event that enters the system and matches a Route-Pipeline pair, it will usually be either:

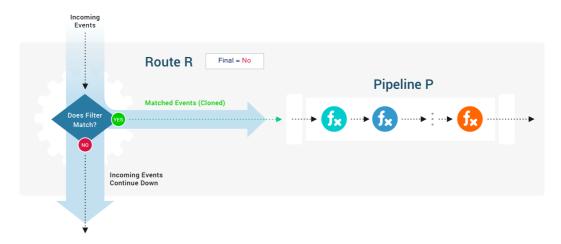
- Dropped by a function, or
- Transformed (optionally) and exit the system.

This behavior is ensured by the Final toggle in Route settings. It defaults to Yes, meaning that matched events will be **consumed** by that Route, and will

not be evaluated against any other Routes that sit below it.



If the Final toggle is set to No, clone(s) of the matching events will be processed by the configured Pipeline, and the original events will be allowed to continue their trip to be evaluated and/or processed by other Route-Pipeline pairs.



This is very useful in cases where the same set of events needs to be processed differently and delivered to different Destinations. Each clone can be decorated with key-value pairs as necessary.

Final Flag and Cloning Considerations

Depending on your cloning needs, you might want to follow a **most-specific first** or a **most-general first** processing strategy. The general goal is to minimize the number of filters/Routes an event gets evaluated against. For example:

If cloning is not needed at all (i.e., all Final toggles stay at default), then it
makes sense to start with the broadest expression at the top, so as to
consume as many events as early as possible.

• If cloning is needed on a narrow set of events, then it might make sense to do that upfront, and follow it with a Route that consumes those clones immediately after.

Route Groups

A Route group is a collection of consecutive Routes that can be moved up and down the Route stack together. Groups help with managing long lists of Routes. They are a UI visualization only: While Routes are in a group, those Routes maintain their global position order.

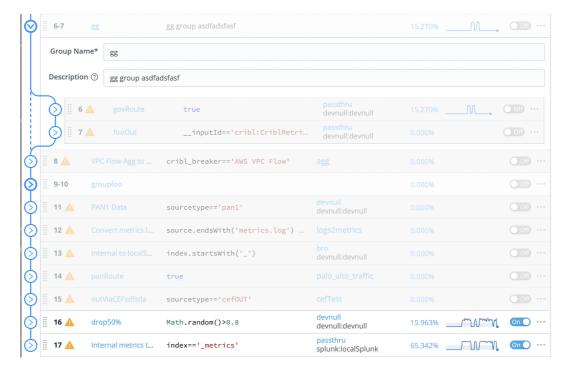
i Route groups work much like Function groups, offering similar UI controls and drag-and-drop options.

Unreachable Routes

Routes display an "unreachable" warning indicator (orange triangle) when data can't reach them. This will be true when, with your current configuration, any Route higher in the stack matches **all** three of these conditions:

- Previous Route is enabled (slider is set to On).
- Previous Route is final (Final slider is set to Yes).
- Previous Route's Filter expression evaluates to true, (e.g., true, 1 == 1, etc.).

Note that the third condition above can be triggered intermittently by a randomizing method like Math.random(). This might be included in a previous Route's own Filter expression, or in a Pipeline Function (such as one configured for random data sampling).



Unreachable Route warnings

Routing with Output Router

Output Router Destinations offer another way to route data. These function as meta-Destinations, in that they allow selection of actual Destinations based on rules. Rules are evaluated in order, top->down, with the first match being the winner.

Pipelines

What Are Pipelines

After your data has been matched by a Route, it gets delivered to a Pipeline. A Pipeline is a list of Functions that work on the data.

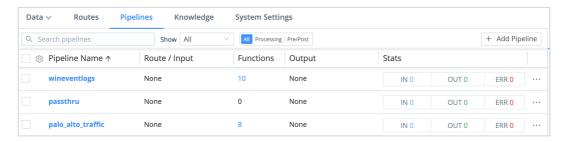
i As with Routes, the order in which the Functions are listed matters. A Pipeline's Functions are evaluated in order, top->down.

Accessing Pipelines

Select **Pipelines** from LogStream's (or a Worker Group's) top menu. To configure a new Pipeline, click **+ Add Pipeline**.

How Do Pipelines Work

Events are always delivered to the beginning of a Pipeline via a Route. The data in the **Stats** column shown below are for the last 15 minutes.



Pipelines and Route inputs

i You can press the] (right-bracket) shortcut key to toggle between the Preview pane and an expanded Pipelines display. This works when no field has focus.

Within the Pipeline, events are processed by each Function, in order. A Pipeline will always move events in the direction that points outside of the system. This

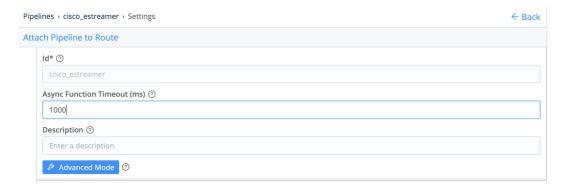
is on purpose, to keep the design simple and avoid potential loops.



Pipeline Functions

Pipeline Settings

Click the gear icon at top right to open the Pipeline's Settings. Here, you can attach the Pipeline to a Route. In the Settings' **Async function timeout (ms)** field, you can enter a buffer to adjust for Functions that might take much longer to execute than normal. (An example would be a Lookup Function processing a large lookup file.)



Pipeline Settings

Advanced Mode (JSON Editor)

Click **Advanced Mode** to edit the Pipeline's definition as JSON text. In this mode's editor, you can directly edit multiple values. You can also use the **Import** and **Export** buttons here to copy and modify existing Pipeline configurations.

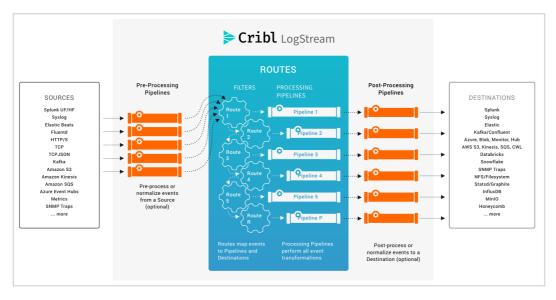
```
Pipelines > elastic > Advanced Settings Mode
                                                                                                      ← Back
Attached to Route: elastic
     Simple Mode
                                                                                       ⊥ Import
                                                                                                  ± Export
             "id": "elastic",
             "conf": {
       3
               "output": "default",
       4
       5
               "groups": {},
               "asyncFuncTimeout": 1000,
       6
               "functions": [
                   "id": "drop",
                   "filter": "host!='192.168.1.241'",
      10
                   "disabled": true,
      11
                   "conf": {}
      12
      13
      14
             }
      15
      16
```

Advanced Pipeline Editing

i You can streamline the above display by organizing related Functions into Function groups.

Types of Pipelines

You can apply various Pipeline types at different stages of data flow. All Pipelines have the same basic internal structure (a series of Functions) – the types below differ only in their position in the system.



Pre-processing, processing, and post-processing Pipelines

Pre-Processing Pipelines

These are Pipelines that are attached to a Source to condition (normalize) the events **before** they're delivered to a processing Pipeline. They're optional.

Typical use cases are event formatting, or applying Functions to **all** events of an input. (E.g., to extract a message field before pushing events to various processing Pipelines.)

You configure these Pipelines just like any other Pipeline, by selecting **Pipelines** from the top menu. You then attach your configured Pipeline to individual **Sources**, using the Source's **Pre-Processing > Pipeline** drop-down.

Fields extracted using pre-processing Pipelines are made available to Routes.

Processing Pipelines

These are "normal" event processing Pipelines, attached directly to Routes.

Post-Processing Pipelines

These Pipelines are attached to a Destination to normalize the events before they're sent out. A post-processing Pipeline's Functions apply to **all** events exiting to the attached Destination.

Typical use cases are applying Functions that transform or shape events per receiver requirements. (E.g., to ensure that a _time field exists for all events bound to a Splunk receiver.)

You configure these Pipelines as normal, by selecting **Pipelines** from the top menu. You then attach your configured Pipeline to individual **Destinations**, using the Destination's **Post-Processing > Pipeline** drop-down.

You can also use a Destination's **Post-Processing** options to add **System Fields** like <code>cribl_input</code>, identifying the LogStream Source that processed the events.

Best Practices for Pipelines

Functions in a Pipeline are equipped with their own filters. Even though filters are not required, we recommend using them as often as possible.

As with Routes, the general goal is to minimize extra work that a Function will do. The fewer events a Function has to operate on, the better the overall performance.

For example, if a Pipeline has two Functions, **f1** and **f2**, and if **f1** operates on source 'foo' and **f2** operates on source 'bar', it might make sense to apply source='foo' versus source='bar' filters on these two Functions, respectively.

Data Onboarding

Onboarding data into Cribl LogStream can vary in complexity, depending on your organization's needs, requirements, and constraints. Proper onboarding from all Sources is key to system performance, troubleshooting, and ultimately the quality of data and decisions both in LogStream and in downstream Destinations.

General Onboarding Steps

Typically, a data onboarding process revolves around these steps, both before and after turning on the Source:

- Create configuration settings.
- Verify that settings do the right thing.
- Iterate.

Below, we break down individual steps.

Before Turning On the Source

Cribl recommends that you take the following steps to verify and tune incoming data, before it starts flowing.

Preview Sample Data

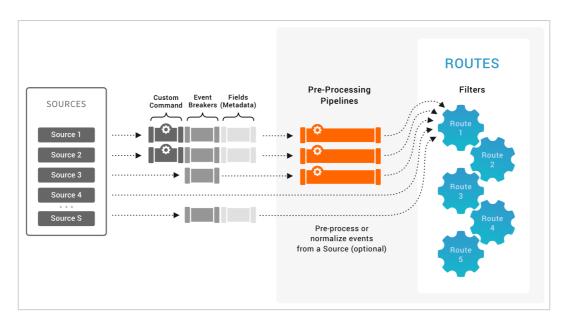
Use a sample of your real data in Data Preview. Sample data can come from a sample Source file that you upload or paste into LogStream.

You can also obtain sample data in a live data capture from a Source. One way to do this **before** going to production is to configure your Source with a **devnull** Pipeline (which just drops all events) as a pre-processing Pipeline. Then, let data flow in for just long enough to capture a sufficient sample.

Check the Processing Order

While events can be processed almost arbitrarily by functions in LogStream Pipelines, make sure you understand the event processing order. This is very

important, as it tells you exactly where certain processing steps occur. For instance, as we'll see just below, quite a few steps can be accomplished at the Source level, before data even hits LogStream Routes.



Source-level processing options

Custom Command

Where supported, data streams will be handled by **custom commands**. These are external system commands that can (optionally) be used to pre-process the data. You can specify any command, script, etc., that consumes via stdin and outputs via stdout.

Verify that such commands are doing what's expected, as they are the very **first** in a series of processing steps.

Event Breakers

Next, data streams are handled by Event Breakers, which:

- Convert data streams into discrete events.
- Extract and assign timestamps to each event.

If the resulting events do not look correct, feel free to use **non-default** breaking rules and timestamp recognition patterns. Downstream, you can use the Auto Timestamp function to modify _time as needed, if timestamps were not recognized properly. Examples of such errors are:

Timestamps too far out in the future or past

- Wrong timezone.
- Incorrect timestamp is selected from multiple timestamps present in the event.

Fields (Metadata)

Next, events can be enriched with Fields (Metadata). This is where you'd add static or dynamic fields to all events delivered by a particular Source.

Pre-Processing Pipeline

Next, you can optionally configure a pre-processing Pipeline on a particular Source. This is extremely useful in these cases:

- Drop non-useful events as early as possible (so as to save on CPU processing).
- Normalize events from this Source to conform a certain shape or structure.
- Fix/touch up events accordingly. E.g., if event breakers assigned the wrong timestamp, this is the best place to use the Auto Timestamp function to adjust _time .

We Can't Say This Enough

Verify, verify, verify your data's integrity before turning on the Source.

After Turning On the Source

Use data Destinations to verify that certain metrics of interest are accurate. This will depend significantly on the capabilities of each Destination, but here's a basic checklist list of things to ensure:

- Timestamps are correct.
- All necessary fields are assigned to events.
- All expected events show up correctly. (E.g., if a Drop or Suppress Function was configured, ensure that it's not dropping unintended events.)
- Throughput both in bytes and in events per second (EPS) is what's expected, or is within a certain tolerance.

Iterate

Iterate on the steps above as necessary. E.g., adjust fields values and timestamps as needed.			
	Remember that there is almost always a workaround. Any arbitrary event transformation that you need is likely just a Function or two away.		

Functions

What Are Functions

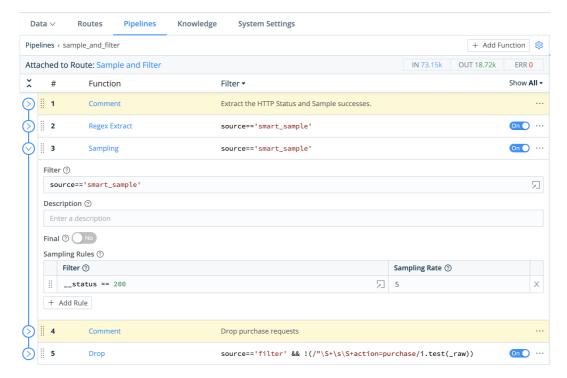
When events enter a Pipeline, they're processed by a series of Functions. At its core, a Function is code that executes on an event, and it encapsulates the smallest amount of processing that can happen to that event.

The term "processing" means a variety of possible options: string replacement, obfuscation, encryption, event-to-metrics conversions, etc. For example, a Pipeline can be composed of several Functions – one that replaces the term foo with bar, another one that hashes bar, and a final one that adds a field (say, dc=jfk-42) to any event that matches source='us-nyc-application.log'.

How Do They Work

Functions are atomic pieces of JavaScript code that are invoked on each event that passes through them. To help improve performance, Functions can be configured with filters to further scope their invocation to matching events only.

You can add as many Functions in a Pipeline as necessary, though the more you have, the longer it will take each event to pass through. Also, you can turn Functions **On/Off** within a Pipeline as necessary. This enables you to preserve structure as you optimize or debug.



Functions stack in a Pipeline

You can reposition Functions up or down the Pipeline stack to adjust their execution order. Use a Function's left grab handle to drag and drop it into place.

The Final Toggle

Similar to the Final toggle in Routes, the Final toggle here controls the flow of events at the Function level. Its states are:

- No (default): means that matching events processed by this Function will be passed down to the next Function.
- Yes: means that this Function is the last one that will be applied to
 matching events. All Functions further down the Pipeline will be skipped. A
 Function with Final set to Yes will display an Findicator in the Pipeline
 stack.

Out-of-the-Box Functions

Cribl LogStream ships with several Functions out-of-the-box, and you can chain them together to meet your requirements. For more details, see individual **Functions**, and the **Use Cases** section, within this documentation.

Custom Functions

For an overview of adding custom Functions to Cribl LogStream, see our blog post, Extending Cribl: Building Custom Functions.

What Functions to Use When

- Add, remove, update fields:
 Eval, Lookup, Regex Extract
- Find & Replace, including basic sed -like, obfuscate, redact, hash, etc.:
 Mask, Eval
- Add GeoIP information to events:
 GeoIP
- Extract fields:Regex Extract, Parser
- Extract timestamps: Auto Timestamp
- Drop events:
 Drop, Regex Filter, Sampling, Suppress, Dynamic Sampling
- Sample events (e.g, high-volume, low-value data):
 Sampling, Dynamic Sampling
- Suppress events (e.g, duplicates, etc.):
 Suppress
- Serialize events to CEF format (send to various SIEMs):
 CEF Serializer
- Serialize / change format (e.g., convert JSON to CSV):
 Serialize
- Convert JSON arrays into their own events: JSON Unroll, XML Unroll
- Flatten nested structures (e.g., nested JSON):
 Flatten
- Aggregate events in real-time (i.e. statistical aggregations):
 Aggregations

- Convert events to metrics format:
 Publish Metrics, Prometheus Publisher (beta)
- Resolve hostname from IP address: Reverse DNS (beta)
- Extract numeric values from event fields, converting them to type number:
 Numerify
- Send events out to a command or a local file, via stdin, from any point in a Pipeline:

Tee

- Convert an XML event's elements into individual events:
 XML Unroll
- Duplicate events in the same Pipeline, with optional added fields:
 Clone
- Add a text comment within a Pipeline's UI, to label steps without changing event data:

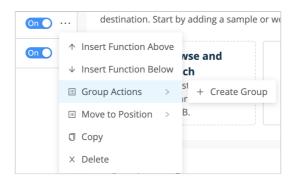
Comment

Function Groups

A Function group is a collection of consecutive Functions that can be moved up and down a Pipeline's Functions stack together. Groups help you manage long stacks of Functions by streamlining their display. They are a UI visualization only: While Functions are in a group, those Functions maintain their global position order in the Pipeline.

i Function groups work much like Route groups.

To build a group from any Function, click the Function's ••• (Options) menu, then select **Group Actions** > **Create Group**.



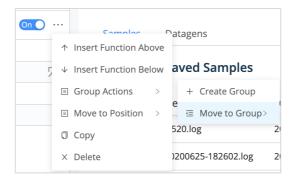
Creating a group

You'll need to enter a **Group Name** before you can save or resave the Pipeline. Optionally, enter a **Description**.



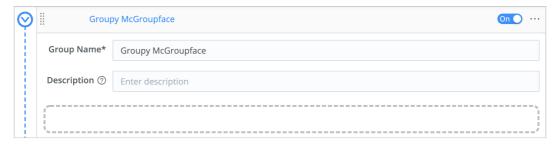
Naming a group

Once you've saved at least one group to a Pipeline, other Functions'
••• (Options) > Group Actions submenus will add options to Move to Group or Ungroup/Ungroup All.



Expanded Group Actions submenu

You can also use a Function's left grab handle to drag and drop it into, or out of, a group. A saved group that's empty displays a dashed target into which you can drag and drop Functions.



Drag-and-drop target

Auto Timestamp

Description

The Auto Timestamp Function extracts time to a destination field, given a source field in the event. By default, Auto Timestamp makes a first best effort and populates _time . When you add a sample (via paste or a local file), you should accomplish time and event breaking at the same time you add the data.

This Function allows fine-grained and powerful transformations to populate new time fields, or to edit existing time fields. You can use the Function's Additional timestamps section to create custom time fields using regex and custom JavaScript strptime functions.

The Auto Timestamp Function uses the same basic algorithm as the Event Breaker Function and the C.Time.timestampFinder() native method.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. The default true setting passes all events through the Function.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Source field: Field to search for a timestamp. Defaults to _raw .

Destination field: Field to place extracted timestamp in. Defaults to _time . Supports nested addressing.

Default timezone: Select a timezone to assign to timestamps that lack timezone info. Defaults to Local . (This drop-down includes support for legacy names: EST5EDT, CST6CDT, MST7MDT, and PST8PDT.)

Additional timestamps: Add Regex/Strptime pairs to extract additional timestamp formats.

- Regex: Regex, with first capturing group matching the timestamp.
- Strptime format: Select or enter the strptime format for the captured timestamp.

Click **Add timestamp** to add more rows.

Advanced Settings

Time expression: Expression with which to format extracted time. Current time, as a JavaScript Date object, is in global time. Defaults to time.getTime() / 1000.

For details about Cribl LogStream's Library (native) time methods, see: C.Time – Time Functions.

Start scan offset: How far into the string to look for a time string.

Max timestamp scan depth: Maximum string length at which to look for a timestamp.

Default time: How to set the time field if no timestamp is found. Defaults to **Current time.**

Two fields enable you to constrain (clamp) the parsed timestamp, to prevent the Function from mistakenly extracting non-time values as unrealistic timestamps:

- Earliest timestamp allowed: Enter a string that specifies the latest allowable timestamp, relative to now. (Sample value: -42years . Default value: -420weeks .) Parsed values earlier than this date will be set to the **Default time**.
- Future timestamp allowed: Enter a string that specifies the latest allowable timestamp, relative to now. (Sample value: +42days. Default value: +1week.) Parsed values after this date will be set to the Default time.

Format Reference

Referencing https://github.com/d3/d3-time-format#locale_format:

```
%a - abbreviated weekday name.*
%A - full weekday name.*
%b - abbreviated month name.*
%B - full month name.*
%c - the locale's date and time, such as %x, %X.*
%d - zero-padded day of the month as a decimal number [01,31].
%e - space-padded day of the month as a decimal number [ 1,31]; equivalent to %_d.
%f - microseconds as a decimal number [000000, 9999999].
%H - hour (24-hour clock) as a decimal number [00,23].
%I - hour (12-hour clock) as a decimal number [01,12].
%j - day of the year as a decimal number [001,366].
%m - month as a decimal number [01,12].
%M - minute as a decimal number [00,59].
%L - milliseconds as a decimal number [000, 999].
%p - either AM or PM.*
%Q - milliseconds since UNIX epoch.
%s - seconds since UNIX epoch.
%S - second as a decimal number [00,61].
%u - Monday-based (ISO 8601) weekday as a decimal number [1,7].
%U - Sunday-based week of the year as a decimal number [00,53].
%V - ISO 8601 week of the year as a decimal number [01, 53].
%w - Sunday-based weekday as a decimal number [0,6].
%W - Monday-based week of the year as a decimal number [00,53].
%x - the locale's date, such as %-m/%-d/%Y.*
%X - the locale's time, such as %-I:%M:%S %p.*
%y - year without century as a decimal number [00,99].
```

```
%Y - year with century as a decimal number.
%Z - time zone offset, such as -0700, -07:00, -07, or Z.
%% - a literal percent sign (%).
```

Directives marked with an asterisk (*) might be affected by the locale definition.

Complying with the Format

In order to use auto timestamping upon ingestion, the formatting used must match the %Z parameters above. E.g., this Function will automatically parse all of these formats:

```
• 2020/06/10T17:17:35.004-0700
```

- 2020/06/10T17:17:35.004-07:00
- 2020/06/10T17:17:35.004-07
- 2020/06/10T10:17:35.004Z
- 2020/06/10T11:17:35.004 EST

To parse other formats, you can use the Additional Timestamps section's internal Regex or Strptime Format operators.

Basic Example

Filter: name.startsWith('kumquats') & value='specific string here'

This will allow the Auto Timestamp Function to act only on events matching the specified parameters.

Sample:

```
Sep 20 12:03:55 PA-VM 1,2019/09/20 13:03:58,CRIBL,TRAFFIC,end,2049,2019/09/20 14:03:58,314.817.1
```

To add this sample (after creating an Auto Timestamp Function with the above **Filter** expression): Go to **Preview** > **Add a Sample** > **Paste a Sample**, and add the data snippet above. Do not make any changes to timestamping or line breaking, and select **Save as Sample File**.

By default, LogSteram will inspect the first 150 characters, and extract the first valid timestamp it sees. You can modify this character limit under **Advanced Settings** > **Max Timestamp Scan Depth**.

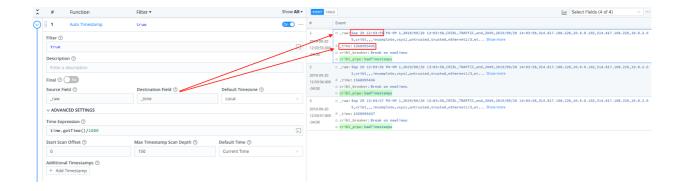
LogStream grabs the first part of the event, and settles on the first matching value to display for time:

```
_time 1569006235
```

GMT: Friday, 20 September 2019, 7:03:55 PM GMT

Your Local Time: Friday, 20 September 2019 PDT, 12:03:55 AM GMT -07:00

Because no explicit timezone has been set (under **Default Timezone**), _time inherits the **Local** timezone, which in this example is GMT -07:00.



i Timezone Dependencies and Details

LogStream uses ICU for timezone information. It does not query external files or the operating system. The bundled ICU is updated periodically.

For additional timezone details, see: https://www.iana.org/time-zones.

Advanced Settings Example

The datetime.strptime() method creates a datetime object from the string passed in by the Regex field.

Here, we'll use datetime.strptime() to match a timestamp in AM/PM format at the end of a line.

Sample:

This is a sample event that will push the datetime values further on inside the event. This is still a sample event and finally here is the datetime information!: Server_UTC_Timestamp="04/27/2020 2:30:15 PM"

Max timestamp scan depth: 200

Click to add Additional timestamps:

Regex: $(d{1,2})/(d{2})/(d{4})$ s $(d{1,2}):(d{2})$ s $(w{2})$

Strptime format: '%m/%d/%Y %H:%M:%S %p'

i Gnarly Details

This Function supports the %f (microseconds) directive, but LogStream will truncate it to millisecond resolution.

For further examples, see Extracting Timestamps from Messy Logs.

Aggregations

Description

The Aggregations Function performs aggregate statistics on event data.

Safeguarding Data

Upon shutdown, LogStream will attempt to flush the buffers that hold aggregated data, to avoid data loss. If you set a **Time window** greater than an hour, Cribl recommends adjusting the **Aggregation memory limit** and/or **Aggregation event limit** to prevent the system from running out of memory. This is especially necessary for high-cardinality data. (Both settings default to unlimited, but we recommend setting defined limits, based on testing.)

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Time window: The time span of the tumbling window for aggregating events. Must be a valid time string (e.g., 10s). Must match pattern \d+[sm]\$.

Aggregates: Aggregate function(s) to perform on events. E.g., sum(bytes).where(action='REJECT').as(TotalBytes). Expression format: aggFunction(<FieldExpression>).where(<FilterExpression>).as(<outputField>). See more examples below.

Note: When used without as(), the aggregate's output will be placed in a field labeled <aggFunction>_<fieldName>. If there are conflicts, the last aggregate wins.
 For example, given two aggregates - sum(bytes).where(action='REJECT') and sum(bytes) - the latter one (sum_bytes) is the winner.

Group by Fields: Fields to group aggregates by.

Evaluate fields: Set of key/value pairs to evaluate and add/set. Fields are added in the context of an aggregated event, before they're sent out. Does not apply to passthrough events.

Time Window Settings

Cumulative aggregations: Determines if the aggregations should be retained for cumulative aggregations, or reset to 0, when flushing out an aggregation table event. Defaults to No.

Lag tolerance: The lag tolerance represents the tumbling window tolerance to late events. Must be a valid time string (e.g., 10s). Must match pattern \d+[sm]\$.

Idle bucket time limit: The amount of time to wait before flushing a bucket that has not received events. Must be a valid time string (e.g., 10s). Must match pattern \d+[sm]\$.

Output Settings

Passthrough mode: Determines whether to pass through the original events along with the aggregation events. Defaults to No.

Metrics mode: Determines whether to output aggregates as metrics. Defaults to No, causing aggregates to be output as events.

Sufficient stats mode: Determines whether to output *only* statistics sufficient for the supplied aggregations. Defaults to No, meaning output richer statistics.

Output prefix: A prefix that is prepended to all of the fields output by this Aggregations Function.

Advanced Settings

Aggregation event limit: The maximum number of events to include in any given aggregation event. Defaults to unlimited. Must be at least 1.

Aggregation memory limit: The memory usage limit to impose upon aggregations. Defaults to unlimited (i.e., the amount of memory available in the system).

List of Aggregate Functions

avg(expr:FieldExpression): Returns the average of the values of the parameter.
count(expr:FieldExpression): Returns the number of occurrences of the values of the parameter.

dc(expr: FieldExpression, errorRate: number = 0.01): Returns the estimated
number of distinct values of the <expr> parameter, within a relative error rate.
distinct_count(expr: FieldExpression, errorRate: number = 0.01): Returns the
estimated number of distinct values of the <expr> parameter, within a relative error rate.
earliest(expr:FieldExpression): Returns the earliest (based on _time) observed
value of the parameter.

first(expr:FieldExpression): Returns the first observed value of the parameter.
last(expr:FieldExpression): Returns the last observed value of the parameter.
latest(expr:FieldExpression): Returns the latest (based on _time) observed value of the parameter.

max(expr:FieldExpression): Returns the maximum value of the parameter.
min(expr:FieldExpression): Returns the minimum value of the parameter.
per_second(expr:FieldExpression): Returns the per second rate (based on _time)
observed value of the parameter.

perc(level: number, expr: FieldExpression):Returns <level> percentile value of the numeric values of the <expr> parameter.

rate(expr:FieldExpression, timeString: string = '1s'):Returns the rate (based
on _time) observed value of the parameter.

stddev(expr:FieldExpression): Returns the sample standard deviation of the values of the parameter.

stddevp(expr:FieldExpression): Returns the population standard deviation of the values of the parameter.

sum(expr:FieldExpression) : Returns the sum of the values of the parameter.
sumsq(expr:FieldExpression) : Returns the sum of squares of the values of the
parameter.

variance(expr:FieldExpression): Returns the sample variance of the values of the parameter.

variancep(expr:FieldExpression): Returns the population variance of the values of the parameter.

How Do Time Window Settings Work?

Lag Tolerance

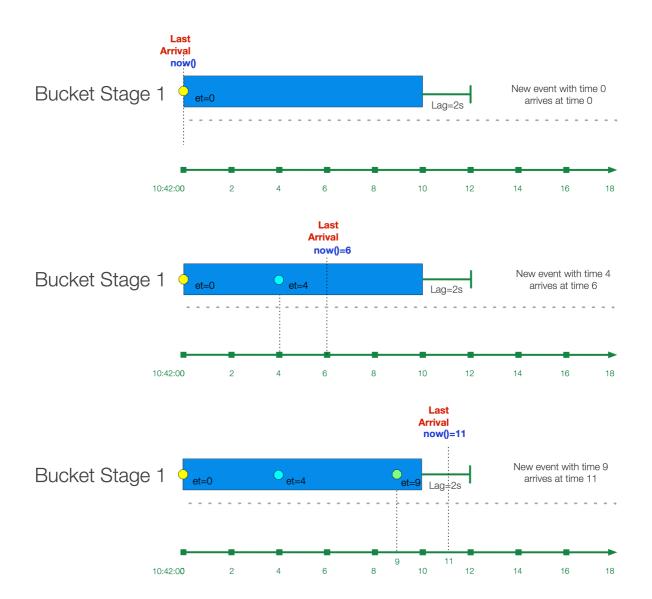
As events are aggregated into windows, there is a good chance that most will arrive later than their event time. For instance, given a 10s window (10:42:00 - 10:42:10), an event with timestamp 10:42:03 might come in 2 seconds later at 10:42:05.

In several cases, there will also be late, or lagging, events that will arrive **after** the latest time window boundary. For example, an event with timestamp 10:42:04 might arrive at

10:42:12 . Lag Tolerance is the setting that governs how long to wait – after the latest window boundary – and still accept late events.

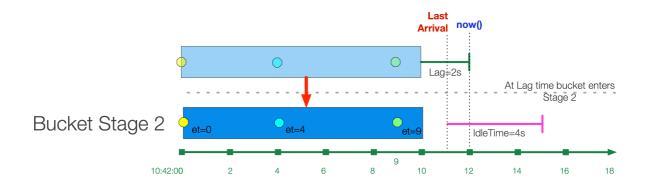
10s Window Aggregation

Settings: Lag=2s, IdleTime=4s



The "bucket" of events is said to be in Stage 1, where it's still accepting new events, but it's not yet finalized. Notice how in the third case, an event with event time 10:42:09 arrives 1 second past the window boundary at 10:42:11, but it's still accepted because it happens before the lag time expires.

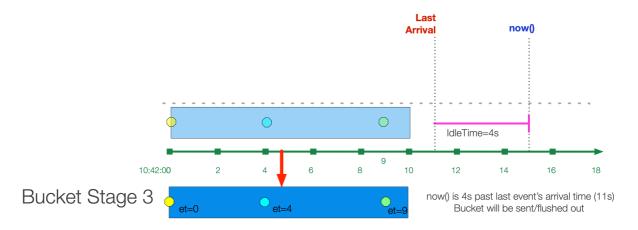
After the lag time expires, the bucket moves to Stage 2.



If the bucket is created from a historic stream, then the bucket is initiated in Stage 2. Lag time is not considered. A "historic" stream is one where the latest time of a bucket is before now(). E.g., if the window size is 10s, and now()=10:42:42, an event with event_time=10 will be placed in a Stage 2 bucket with range 10:42:10 - 10:42:20.

Idle Bucket Time Limit

While Lag Tolerance works with event time, Idle Bucket Time Limit works on arrival time (i.e., real time). It is defined as the amount of time to wait before flushing a bucket that has not received events.



After the Idle Time limit is reached, the bucket is "flushed" and sent out of the system.

Examples

Assume we're working with VPC Flowlog events that have the following structure:

version account_id interface_id srcaddr dstaddr srcport dstport protocol
packets bytes start end action log_status

For example:

2 99999XXXXX eni-02f03c2880e4aaa3 10.0.1.70 10.0.1.11 9999 63030 6 6556 262256 1554562460 1554562475 ACCEPT OK

2 496698360409 eni-08e66c4525538d10b 37.23.15.38 10.0.2.232 4373 8108 6 1 52 1554562456 1554562466 REJECT OK

Scenario A:

Every 10s, compute sum of bytes and output it in a field called Total Bytes.

Time Window: 10s

Aggregations: sum(bytes).as(TotalBytes)

Scenario B:

Every 10s, compute sum of bytes, outputitin a field called Total Bytes, group by srcaddr.

Time Window: 10s

Aggregations: sum(bytes).as(TotalBytes)

Group by Fields: srcaddr

Scenario C:

Every 10s, compute sum of $\,$ bytes $\,$ but only where action is $\,$ REJECT , output it in a field $\,$ called $\,$ TotalBytes , group by $\,$ srcaddr .

Time Window: 10s

Aggregations: sum(bytes).where(action='REJECT').as(TotalBytes)

Group by Fields: srcaddr

Scenario D:

Every 10s, compute sum of bytes but only where action is REJECT, output it in a field called TotalBytes. Also, compute distinct count of srcaddr.

Time Window: 10s

Aggregations:

sum(bytes).where(action='REJECT').as(TotalBytes)
distinct_count(srcaddr).where(action='REJECT')

i For further examples, see Engineering Deep Dive: Streaming Aggregations Part 2

– Memory Optimization

CEF Serializer

Description

The CEF Serializer takes a list of fields and/or values, and formats them in the Common Event Format (CEF) standard. CEF defines a syntax for log records. It is composed of a standard prefix, and a variable extension formatted as a series of key-value pairs.

Format

CEF:Version|Device Vendor|Device Product|Device Version|Device
Event Class ID|Name|Severity|[Extension]

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Output field: The field to which the CEF formatted event will be output. Nested addressing supported. Defaults to _raw .

Header Fields

CEF Header field definitions. The field values below will be written pipe (|) – delimited in the Output Field. Names cannot be changed. Values can be computed with JS expression, or can be constants.

- **cef_version**: Defaults to CEF:0.
- device_vendor: Defaults to Cribl.
- device_product: Defaults to Cribl.
- device_version: Defaults to C.version.
- device_event_class_id: Defaults to 420.

- name: Defaults to Cribl Event.
- severity: Defaults to 6.

Extension Fields

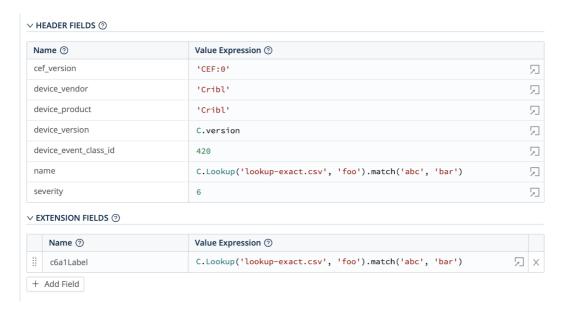
CEF Extension field definitions. Field names and values will be written in key=value format. Select each field's Name from the drop-down list. Values can be computed with JS expressions, or can be constants.

Example

For each CEF field, allowed values include strings, plus any custom Cribl function. For example, if using a lookup:

Name: Name
Value expression: C.Lookup('lookup-exact.csv', 'foo').match('abc', 'bar')

This can be used for any of the CEF Header Fields.



The resulting event has the following structure for an **Output Field** set to _CEF_out:

_CEF_out:CEF:0|Cribl|Cribl|42.0-61c12259|420|Business Group 6|6|c6a1Label=Colorado_Ext_Bldg7

Clone

Description

The Clone Function clones events, with optional added fields. Cloned events will be sent to the same Destination as the original event, because they are in the same Pipeline.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Clones: Create clones with the specified fields added and set.

Fields: Set of key-value pairs to add. Nested addressing is supported.

Examples

In this example, the Destination will receive a clone with an env field set to staging .

Field: env

Value: staging

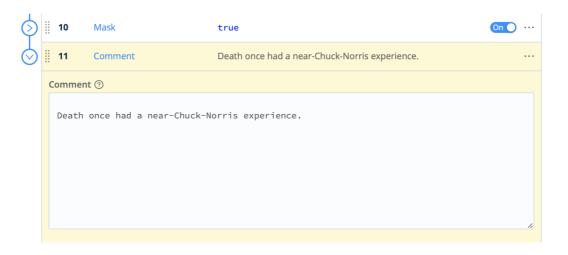
Comment

Description

The Comment Function adds a text comment in a Pipeline. It makes no changes to event data. The added comment is visible only within the Pipeline UI, where it is useful for labeling Pipeline steps.

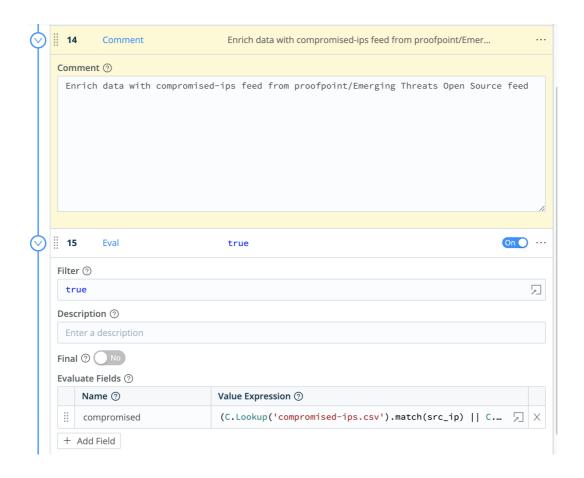
Usage

Comment: Add your comment as plain text in this field.



Examples

This comment labels the Pipeline's next function:



DNS Lookup

Description

The DNS Lookup Function offers two operations useful in enriching security and other data:

- DNS lookups based on host name as text, resolving to A record (IP address) or to other record types.
- Reverse DNS Lookup. (This duplicates LogStream's existing Reverse DNS Function, which is now deprecated.)

To reduce DNS lookups and minimize latency, the DNS Lookup Function incorporates a configurable DNS cache. If you need additional caching, consider enabling OS-level DNS caching on each LogStream Worker that will execute this Function. (OS-level caching options include DNSMasq, nscd, systemd-resolved, etc.)

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No .

DNS Lookup Fields Section

Lookup field name: Name of the field containing the domain to look up.

Resource record type: DNS record type (RR) to return. Defaults to A 'record.

Output field name: Lookup result(s) will be added to this field. Leave blank to overwrite the original field specified in **Lookup field name**.

Reverse DNS Lookup Field(s) Section

Lookup field name: Name of the field containing the IP address to look up.

⚠ If the field value is not in IPv4 or IPv6 format, the lookup is skipped.

Output field name: Name of the field in which to add the resolved hostname. Leave blank to overwrite the original field specified in **Lookup field name**.

Advanced Settings

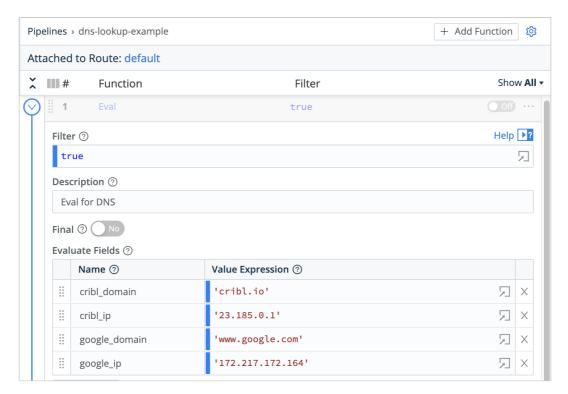
DNS server(s) overrides: IP address(es), in RFC 5952 format, of the DNS server(s) to use for resolution. IPv4 examples: 1.1.1.1, 4.2.2.2:53 . IPv6 examples: [2001:4860:4860::8888], [2001:4860:4860::8888]:1053. If this field is not specified, LogStream will use the system's DNS server.

Reload period (minutes): How often to refresh the DNS cache. Use 0 to disable refreshes. Defaults to 30 minutes.

Maximum cache size: Maximum number of DNS resolutions to cache locally. Before changing the default 5000, contact Cribl Support to understand the implications. Highest allowed value is 10000.

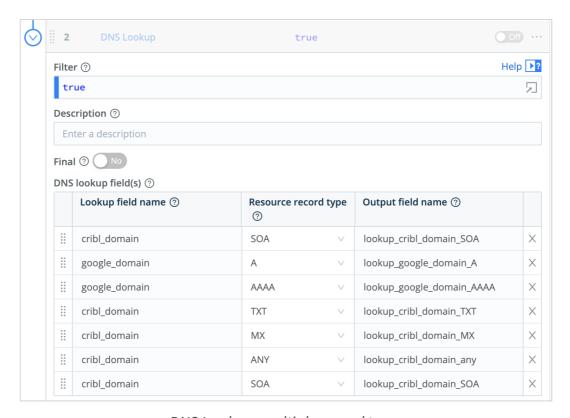
Example

This example Pipeline chains two Functions. First, we have an **Eval** Function that defines key-value pairs for two alphabetical domain names and two numeric IP addresses.



DNS Lookup: Eval Function

Next, the DNS Lookup Function looks up several record types for the two domain names, placing each retrieved record type in its own output field.



DNS Lookup: multiple record types

Finally, the same Function's **Reverse DNS lookup** section retrieves domain names for the two IP addresses.



DNS Lookup: reverse lookups

Drop

Description

The Drop Function will drop/delete any events that meet the Filter expression.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Example

Assume that we care only about errors, so we want to filter out any events that contain the word "success," regardless of case: "success," "SUCCESS," etc.

In our Drop Function, we'll use the JavaScript search() method to search the _raw field's contents for our target pattern. We know that search() returns a non-negative integer to indicate the starting position of the first match in the string, or -1 if no match. So we can evaluate the Function as true when the return value is >= 0.

Filter: _raw.search(/success/i) ≥ 0

Dynamic Sampling

Description

The Dynamic Sampling Function filters out events based on an expression, a sample mode, and events' volume. Your sample mode's configuration determines what percentage of incoming events will be passed along to the next step.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events passed into the Function will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Sample mode: Defines how sample rate will be derived. For formulas and usage details, see **Sample Modes** below. Supported methods:

- Logarithmic (the default): log(previousPeriodCount).
- Square root: sqrt(previousPeriodCount).

Sample group key: Expression used to derive sample group key. For example: \${domain}:\${httpCode} . Each sample group will have its own derived sampling rate, based on volume. Defaults to `\${host}`.

All events without a host field passing through the Function will be associated with the same group and sampled the same.

Advanced Settings

- Sample period Sec: How often (in seconds) sample rates will be adjusted.

 Defaults to 30.
- **Minimum events**: Minimum number of events that must be received, in previous sample period, for sampling mode to be applied to current

period. If the number of events received for a sample group is less than this minimum, a sample rate of 1:1 is used. Defaults to 30.

• Max sampling rate. Maximum sampling rate. If the computed sampling rate is above this value, the rate will be limited to this value.

How Does Dynamic Sampling Work

Compared to static sampling, where users must select a sample rate a priori, Dynamic Sampling allows for **automatically adjusting** sampling rates, based on incoming data volume per sample group. This Function allows users to set only the aggressiveness/coarseness of this adjustment. Square Root is more aggressive than Logarithmic mode.

As an event passes through the Function, it's evaluated against the Sample Group Key expression to determine the sample group it will be associated with. For example, given an event with these fields: ... ip=1.2.3.42, port=1234 ..., and a Sample Group Key of `\${ip}:\${port}`, the event will be associated with the 1.2.3.42:1234 sample group.

⚠ If the Sample Group Key is left at its `\${host}` default, all events without a host will be associated with the same group and sampled the same.

When a sample group is new, it will initially have a sample rate of 1:1 for Sample Period seconds (this value defaults to 30 seconds). Once Sample Period seconds have elapsed, a sample rate will be derived based on the configured Sample Mode, using the sample group's event volume during the **previous** sample period.

For example, assuming a Logarithmic Sample Mode:

Period 0 (first 30s): Number of events in sample group: 1000 , Sample Rate:
1:1 , Events allowed: ALL
Sample Rate calculation for next period: Math.ceil(Math.log(1000)) = 7

Period 1 (next 30s) -- Number of events in sample group: 4000 , Sample Rate: 7:1 : Events allowed: 572

Sample Rate calculation for next period: Math.ceil(Math.log(4000)) = 9

Period 2 (next 30s) -- Number of events in sample group: 12000, Sample Rate: 9:1: Events allowed: 1334

Sample Rate calculation for next period: Math.ceil(Math.log(12000)) =
10

Period 3 (next 30s) -- Number of events in sample group: 2000, Sample Rate:

10:1: Events allowed: 200

Sample Rate calculation for **next** period: Math.ceil(Math.log(2000)) = 8

•••

Sample Modes

- Logarithmic The sample rate is derived, for each sample group, using a natural log: Math.ceil(Math.log(lastPeriodVolume)). This mode is less aggressive, and drops fewer events.
- 2. Square Root The sample rate is derived, for each sample group, using:
 Math.ceil(Math.sqrt(lastPeriodVolume)) . This mode is more
 aggressive, and drops more events.

Example

Here's an example that illustrates the effectiveness of using the Square Root sample mode.

Settings:

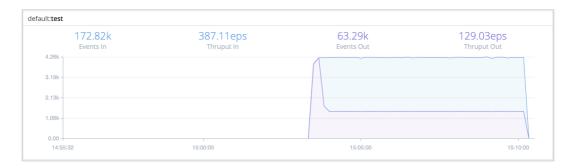
Sample Mode: Square Root

Sample Period (sec): 20

Minimum Events: 3
Max. Sampling Rate: 3

Results:

Events In: 4.23K Events Out: 1.41K



In this generic example, we reduced the incoming event volume from 4.23K to 1.41K. Your own results will vary depending on multiple parameters – the Sample Group Key, Sample Period, Minimum Events, Max Sampling Rate, and rate of incoming events.

i For further examples, see Getting Smart and Practical With Dynamic Sampling.

Eval

Description

The Eval Function adds or removes fields from events. (In Splunk, these are index-time fields.)

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Evaluate fields: Set of key/value pairs to add. The left-hand side input (**Name**) is the key name. The right-hand side input (**Value Expression**) is a JS expression to compute the value – this can be a constant. Nested addressing is supported. Strings intended to be used as values must be single- or double-quoted. (For details, see Introduction to Expression Syntax.)

Keep fields: List of fields to keep. Wildcards (*) and nested addressing are supported. Takes precedence over **Remove fields** (below). To reference a parent object and all children requires using the (*) wildcard. For example, if _raw is converted to an object then use _raw* to refer to itself and all children.

Remove fields: List of fields to remove. Wildcards (*) and nested addressing are supported. Cannot remove fields matching Keep fields. Cribl LogStream internal fields that start with __ (double underscore) cannot be removed via wildcard. Instead, they need to be specified individually. For example, __myField cannot be removed by specifying __myF*.

Using Keep and Remove

A field matching an entry in *both* **Keep** (wildcard or not) and **Remove** will *not* be removed. This is useful for implementing "remove all but" functionality. For

example, to keep only _time, _raw, source, sourcetype, host, we can specify them all in **Keep**, while specifying * in **Remove**.

Negated terms are supported in both **Keep fields** and **Remove fields**. The list is order-sensitive when negated terms are used. Examples:

- !foobar, foo* means "All fields that start with 'foo' except foobar."
- !foo*, * means "All fields except for those that start with 'foo!"

Examples

Scenario A: Create field myField with static value of value1:

- Name: myField
- Value Expression: 'value1'

Scenario B: Set field action to blocked if login=error:

- Name: action
- Value Expression: login='fail' ? 'blocked' : action

Scenario C: Create a multivalued field called myTags . (i.e., array):

- Name: myTags
- Value Expression: ['failed', 'blocked']

Scenario D: Add value error to the multivalued field myTags:

- Name: myTags
- Value Expression: login='error' ? [... myTags, 'error'] : myTags

(The above expression is literal, and uses JavaScript spread syntax.)

Scenario E: Rename an identification field to the shorter ID – copying over the original field's value, and removing the old field:

- Name: ID
- Value Expression: identification
- Remove Field: identification
 - i See Ingest-time Fields for more examples.

Advanced Usage Notes

Note 1

The Eval Function has the ability to execute expressions without assigning their value to the field of an event. You can do this by simply leaving the left-hand side input empty, and having the right-hand side do the assignment.

- Simple Example: Object.assign(foo, JSON.parse(bar), JSON.parse(baz)) on the right-hand side (and left-hand side empty) will JSON-parse the strings in bar and baz, merge them, and assign their value to foo, an already existing field.
- Another Example: To parse JSON, enter Object.assign(__e, JSON.parse(_raw)) on the right-hand side (and left-hand side empty).
 __e is a special variable that refers to the (context) event within a JS expression. In this case, content parsed from _raw is added at the top level of the event.

Note 2

You can also use the Eval Function to set and unset control fields (e.g., _TCP_ROUTING in Splunk), via this syntax: _ctrl.<name> . Control fields can be referenced only on the left-hand side of Add. (I.e., they cannot be read or used on the right-hand side, and cannot be referenced in Remove.)

To unset/delete a control field, set its value to undefined. These fields are normally not needed for event computations, and modifying them is suggested to be done only by experts. Please reach out to Cribl if you need help with this topic.

Flatten

Description

The Flatten Function is used to flatten fields out of a nested structure.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Fields: List of top-level fields to include for flattening. Supports * wildcards. Defaults to empty array, which means all fields.

Prefix: Prefix string for flattened field names. Defaults to empty.

Depth: Number representing the nested levels to consider for flattening. Minimum 1. Defaults to 5.

Delimiter: Delimiter to be used for flattening. Defaults to _ (underscore).

Example

Add the following test sample in **Preview** > **Paste a Sample**:

```
input
{ "accounting" : [ { "firstName" : "John", "lastName" : "Doe", "age" : 23 }, { "firstName" : "Mage" : "Mage" : 23 }, }
```

Under **Select Event Breaker**, choose **ndjson** (newline-delimited JSON), and click **Save as a Sample File**.

Here's sample output with all settings at default:

```
output
{
    "accounting_0_firstName": "John",
    "accounting_0_lastName": "Doe",
    "accounting_0_age": 23,
    "accounting_1_firstName": "Mary",
    "accounting_1_lastName": "Smith",
    "accounting_1_age": 32,
    "sales_0_firstName": "Sally",
    "sales_0_lastName": "Green",
```

```
"sales_0_age": 27,
"sales_1_firstName": "Jim",
"sales_1_lastName": "Galley",
"sales_1_age": 41,
}
```

Using the Flatten Function's default settings, we successfully create top-level fields from the nested JSON structure, as expected.

GeolP

Description

The GeoIP Function enriches events with geographic fields, given an IP address. It is optimized for binary databases such as MaxMind's GeoIP.

☐ For details on setting up MaxMind (and similar) databases, see Managing Large Lookups.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

GeoIP file (.mmdb): Path to a Maxmind database, in binary format, with .mmdb extension.

i If the database file is located within the lookup directory (\$CRIBL_HOME/data/lookups/), the GeoIP file does not need to be an absolute path.

In distributed deployments, ensure that the Maxmind database file is in the same location on both the Master and Worker Nodes.

IP field: Field name in which to find an IP to look up. Can be nested. Defaults to ip.

Result field: Field name in which to store the GeoIP lookup results. Defaults to geoip.

Examples

Assume that you are receiving SMTP logs, and need to see geolocation information associated with IPs using the SMTP service.

Here's a sample of our data, from IPSwitch IMail Server logs:

```
03:19 03:22 SMTPD(00180250) [192.168.1.131] connect 74.136.132.88 port 2539 03:19 03:22 SMTPD(00180250) [74.136.132.88] EHLO msnbc.com 03:19 03:22 SMTPD(00180250) [74.136.132.88] MAIL FROM: <info-jjgcdshx@test.us> 03:19 03:22 SMTPD(00180250) [74.136.132.88] RCPT To:<user@domain.com>
```

In this example, we'll chain together three Functions. First, we'll use a Regex Extract Function to isolate the host's IP. Next, we'll use the GeoIP Function to look up the extracted IP against our geoIP database, placing the returned info into a new __geoip field. Finally we'll use an Eval Function to parse that field's city, state, country, ZIP, latitude, and longitude.

Function 1 – Regex Extract

Regex: $[(?<ip>\S+)]$

Source field: _raw

Result: 74.136.132.88

Function 2 - GeoIP

Event's IP field: ip
Result field: __geoip

Function 3 - Eval

Name	Value Expression
City	geoip.city.names.en
Country	geoip.country.names.en
Zip	geoip.postal.code
Lat	geoip.location.latitude
Long	geoip.location.longitude

In the Eval Function's Remove fields setting, you could specify thegeoip field for removal, if desired. However, its prefix makes it an internal field	
· · · · · · · · · · · · · · · · · · ·	
anyway.	
☐ For a hosted tutorial on applying the GeoIP Function, see Cribl's GeoIP and Threat Feed Enrichment Sandbox.	

Grok

Description

The Grok Function extracts structured fields from unstructured log data, using modular regex patterns.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Optional description of this Function's purpose in this Pipeline. Defaults to empty.

Final: If toggled to Yes, stops data from being fed to downstream Functions. Defaults to No.

```
Pattern: Grok pattern to extract fields. Syntax supported: % {PATTERN_NAME: FIELD_NAME} .
```

Click + Add pattern to chain more patterns.

Source field: Field on which to perform Grok extractions. Defaults to _raw .

Management

You can add and edit Grok patterns via LogStream's UI by selecting **Knowledge > Grok Patterns.** Pattern files are located at: \$CRIBL_HOME/(default|local)/cribl/grokpatterns/

Example

Example event:

```
{"_raw": "2020-09-16T04:20:42.45+01:00 DEBUG This is a sample debug log message"}`

Pattern: %{TIMESTAMP_IS08601:event_time} %{LOGLEVEL:log_level} %

{GREEDYDATA:log_message}
```

Source Field: _raw

Event after extraction:

```
{"_raw": "2020-09-16T04:20:42.45+01:00 DEBUG This is a sample debug log message",
    "_time": 1600226442.045,
    "event_time": "2020-09-16T04:20:42.45+01:00",
    "log_level": "DEBUG",
    "log_message": "This is a sample debug log message",
}
```

Note the new fields added to the event: event_time, log_level, and log_message.

References

- Syntax for a Grok pattern is %{PATTERN_NAME:FIELD_NAME} . E.g.: %{IP:client} % {WORD:method} .
- Useful links for creating and testing Grok patterns: http://grokdebug.herokuapp.com and http://grokconstructor.appspot.com/.
- Additional patterns are available here: https://github.com/logstash-plugins/logstash-patterns-core/tree/master/patterns.

JSON Unroll

Description

The JSON Unroll Function accepts a JSON object string _raw field, unrolls/explodes an array of objects therein into individual events, while also inheriting top level fields. See example(s). Cribl highly recommends not using this JSON Unroll function for certain types of data. Instead, perform the unrolling using an event breaker for those inputs which support configuring an event breaker. Specifying the event breaker type JSON Array and toggling the JSON Extract Fields option to Yes will accomplish the same unrolling but much more efficiently. This is recommended, for example, for CloudTrail and Office635 events, which are collected as JSON arrays.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Path: Path to array to unroll, e.g., foo.0.bar.

New name: The name that the exploded array element will receive in each new event. Leave empty to expand the array element with its original name.

Example(s)

Assume you have an incoming event that has a _raw field as a JSON object string like this:

```
Sample_raw field
```

Settings:

Path: allCars
New Name: cars

Output Events:

```
Resulting Events

Event 1:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"Ford","models'

Event 2:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"GM","models":|

Event 3:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"Fiat","models'

Event 4:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"Blackberry","r
```

Each element under the original **allCars** array is now placed in a **cars** field in its own event, inheriting original top level fields; **date**, **name** and **age**

Lookup

Description

The Lookup Function enriches events with external fields, using lookup table files in CSV, compressed .csv.gz, or binary .mmdb format.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Lookup file path (.csv, .csv.gz): Path to the lookup file. Select an existing file that you've uploaded via LogStream's UI at Knowledge > Lookups Libary, or specify the path. You can reference environment variables via \$, e.g.: \$CRIBL_HOME/file.csv.

i When you configure this field via a distributed deployment's Master Node, LogStream will swap \$CRIBL_HOME/groups/<groupname>/ for \$CRIBL_HOME when validating whether the file exists. In this case, the default upload path changes from \$CRIBL_HOME/data/lookups (single-instance deployments) to \$CRIBL_HOME/groups/<groupname>/data/lookups/ (distributed deployments).

Match mode: Defines the format of the lookup file, and indicates the matching logic that will be performed. Defaults to Exact .

Match type: For CIDR and Regex Match modes, this attribute refines how to resolve multiple matches. First match will return the first matching entry. Most specific will scan all entries, finding the most specific match. All will return all matches in the output, as arrays. (Defaults to First match. Not displayed for Exact Match mode.)

Lookup fields (.csv): Field(s) that should be used to key into the lookup table.

- Lookup field name in event: Exact field name as it appears in events. Nested addressing supported.
- Corresponding field name in lookup: The field name as it appears in the lookup file. Defaults to the Lookup field name in event value. This input is optional.

△ Case-Sensitive / Multiple Matches

Lookups are case-sensitive by default. (See the **Ignore case** option below.)

If the lookup file contains duplicate key names with different values, all **Match mode**s of this Function will use **only** the value in the key's **final** instance, ignoring all preceding instances.

Output field(s): Field(s) to add to events after matching the lookup table. Defaults to **all** if not specified.

- Output field name from lookup: Field name, as it appears in the lookup file.
- Lookup field name in event: Field name to add to event. Defaults to the lookup field name. This input is optional. Nested addressing is supported.

Advanced Settings

Reload period (sec): Periodically check the underlying file for modtime changes, and reload if necessary. Use -1 to disable. Defaults to 60.

Ignore case: Ignore case when performing Match mode: Exact lookups. Defaults to No.

Add to raw event: Whether to append the looked-up values to the _raw field, as key=value pairs. Defaults to No.

Examples

Example 1: Regex Lookups

Assign a sourcetype field to events if their _raw field matches a particular regex.

```
regex, sourcetype
"^[^,]+,[^,]+,[^,]+,THREAT",pan:threat
"^[^,]+,[^,]+,[^,]+,TRAFFIC",pan:traffic
"^[^,]+,[^,]+,[^,]+,SYSTEM",pan:system
```

Match mode: Regex

Match type: First match

Lookup field name in event: _raw

Corresponding field name in lookup: regex

```
Events before and after

### BEFORE:

{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:58
{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,FOOBAR,THREAT,end,2049,2018/09/20 13:03:58
```

```
### AFTER:
```

```
{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:5
    "sourcetype": "pan:traffic"
    }
{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,FOOBAR,THREAT,end,2049,2018/09/20 13:03:58
    "sourcetype": "pan:threat"
}
```

Example 2: CIDR Lookups

Assign a location field to events if their destination_ip field matches a particular CIDR range.

```
paloaltoips.csv

range,location
10.0.0.0/24,San Francisco
10.0.0.0/16,California
10.0.0.0/8,US
```

Match mode: CIDR

Match type: See options below

Lookup field name in event: destination_ip

Corresponding field name in lookup: range

i In Match mode: CIDR with Match type: Most specific, the lookup will implicitly search for matches from most specific to least specific. There is no need to pre-sort data.

Note that Match mode: CIDR with Match type: First Match is likely the most performant with large lookups. This can be used as an alternative to Most specific, if the file is sorted with the most specific/relevant entries first. This mode still performs a table scan, top to bottom.

Events before and after

```
### BEFORE:
{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:
    "destination_ip": "10.0.0.102"
}
### AFTER with Match Type: First Match

{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:
    "destination_ip": "10.0.0.102",
    "location": "San Francisco"
}
### AFTER with Match Type: Most Specific
```

```
{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:
   "destination_ip": "10.0.0.102",
   "location": "San Francisco"
}

### AFTER with Match Type: All

{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:
   "destination_ip": "10.0.0.102",
   "location": [
        "San Francisco",
        "California",
        "US",
   ]}
```

More Examples and Scenarios

More examples:

- Ingest-time Lookups.
- Lookups and Regex Magic.
- Lookups as Filters for Masks.

See also:

- Managing Large Lookups to optimize file locations for large lookup files.
- Redis Function for faster lookups using a Redis integration.

Mask

Description

The Mask Function masks, or replaces, patterns in events. This is especially useful for redacting PII (personally identifiable information) and other sensitive data.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Masking rules: Match Regex and Replace Expression pairs. Defaults to empty. Each row has the following fields:

- Match regex: Pattern to replace. Supports capture groups. Use /g to replace all matches, e.g.: /foo(bar)/g
- **Replace expression**: A JavaScript expression or literal to replace all matching content.

To add more rows, click + Add Rule.

Apply to fields: Fields on which to apply the masking rules. Defaults to $_{\rm raw}$. Add more fields by typing in their names, separated by hard returns. Supports wildcards (\star) and nested addressing.

i Negated terms are also supported. When you negate field names, the fields list is order-sensitive. E.g., !foobar before foo* means "Apply to all fields that start with foo , except foobar." However, !foo* before * means "Apply to all fields, except for those that start with foo."

Advanced Settings

Evaluate fields: Optionally, specify fields to add to events in which one or more of the **Masking Rules** were matched. These fields can be useful in downstream processing and reporting. You specify the fields as key–value expression pairs, like those in the Eval Function.

- Name: Field name.
- Value Expression: JavaScript expression to compute the value (can be a constant).

Evaluating the Replace Expression

The **Replace expression** field accepts a full JS expression that evaluates to a value, so you're not necessarily limited to what's under C.Mask . For example, you can do conditional replacement: g1%2=1 ? `fieldA="odd"`: `fieldA="even"`

The Replace expression can reference other event fields as event. <fieldName> . For example, `\${g1}\${event.source}` . Note that this is slightly different from other expression inputs, where event fields are referenced without event. Here, we require the event. prefix for the following reasons:

- We don't expect this to be a common case.
- Expanding the event in the replace context would have a high performance hit on the common path.
- There is a slight chance that there might be a gN field in the event.

Examples

Example 1: Transform a String

Here, we'll simply search for the string dfhgdfgj, and replace that value (if found) with Trans AM. This will help close America's muscle-car gap:

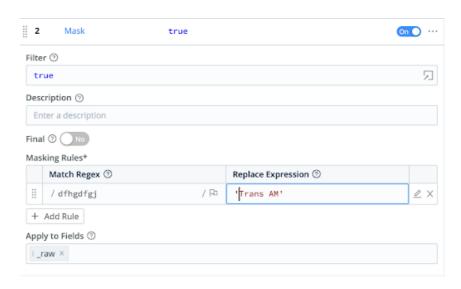


Event before masking

Configure the Mask Function > Masking Rules as follows:

Match Regex: dfhgdfgj

Replace Expression: Trans AM



Mask Function configuration

Result: Vroom vroom!

Event after masking

Example 2: Mask Sensitive Data

Assume that you're ingesting data whose <code>_raw</code> fields contain unredacted Social Security numbers in the Key=Value pattern <code>social=########</code>.

```
a_raw: 2020-07-22 05:22:43,330,Event [Event=UpdateBillingProvQuote, timestamp=1577371]
       0, properties={JMSCorrelationID=NA, JMSMessageID=ID:ESP-PD.C7A19FC656293:AB21BCF
       E, orderType=NewActivation, quotePriority=NORMAL, conversationId=ESB~BEBFAB927C87
      5E35:81E10EA8:47283ADA8A10:5568, credits=NA, JMSReplyTo=pub.esb.genericasync.resp
      onse, timeToLive=-1, serviceName=UpdateBillingProvisioning, esn=10D9C064A00987, a
      ccountNumber=900001336, social=518057110, MethodName=InternalEvent, AdapterName=U
      pdateBillingProvQuote, meid=NA, orderNumber=900000000002363, quoteNumber=4258319
      8, ReplyTo=NA, userName=yosem7, EventConversationID=NA, mdn=6248526355, accountTy
      pe=PostPaid, marketCity="JOLIET", marketState=IL, marketZip=60432, billingCycle=2
       4, autoBillPayment=T, phoneCode=SGS5, phoneType=Android, phoneName="Samsung GALAX
       Y S5", planCode=1400POST5L90, planType=PostPaid, planPrice=89.99, planName="1400"
       Minute Family", planDescription="Nationwide 1400 Minutes, Unlimited Mobile to Mob
       ile, Unlimited Night & Weekend, Unlimited Data", cardNumber=3569948084568945, net
       workProviderName=Splunktel}] Showless
# time: 1595395363.33
α host: 127.0.0.1
a index: cribl
a source: /opt/tibco/tra/apps/ESB/logs/business_event.log
α sourcetype: business_event
```

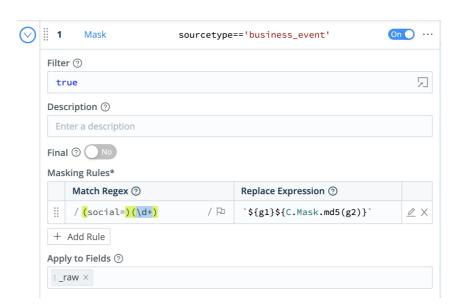
Event with unredacted SSNs

You can use a Mask Function to run an md5 hash of the social keys' numeric values, replacing the original values with the hashed values. Configure the Masking Rules as follows:

Match Regex: (social=)(\d+)
Replace Expression: `\${g1}\${C.Mask.md5(g2)}`

In the first example everything in the Match regex field was replaced by the Replace Expression. However if that isn't desired then you can use capture groups in the Match Regex to define individual string components for manipulation or, alternatively, use string literals in the Replace expression for retaining any static text. Any content matching the Match Regex that is not inserted into the Replace expression will not be retained.

In this example, social= is assigned to capture group g1 for later reference. The value of social= will be hashed by referencing it as g2 in the md5 function. If we didn't make social= its own capture group (or specified social= as a literal in the Replace Expression) then we cannot reference it using g1 in the Replace expression, the value of social= would instead be assigned to g1, and the entire social=######## string would be replaced with a hash of the social security number, which probably isn't desired because no one would know the value being hashed without a field name preceding it.



Mask Function configuration

Result: The sensitive values are replaced by their md5 hashes.

```
a _raw: 2020-07-22 05:22:43,330,Event [Event=UpdateBillingProvQuote, timestamp=1577371270, properties={JMSCorrel
      ationID=NA, JMSMessageID=ID:ESP-PD.C7A19FC656293:AB21BCFE, orderType=NewActivation, quotePriority=NORMA
       L, conversationId=ESB~BEBFAB927C875E35:81E10EA8:47283ADA8A10:5568, credits=NA, JMSReplyTo=pub.esb.generi
       casync.response, timeToLive=-1, serviceName=UpdateBillingProvisioning, esn=10D9C064A00987, accountNumber
       =900001336, social=dlce1763a8e5213781a30f8e7ba9172f, MethodName=InternalEvent, AdapterName=UpdateBilling
       ProvQuote, meid=NA, orderNumber=9000000000002363, quoteNumber=42583198, ReplyTo=NA, userName=yosem7, Eve
       ntConversationID=NA, mdn=6248526355, accountType=PostPaid, marketCity="JOLIET", marketState=IL, marketZi
       p=60432, billingCycle=24, autoBillPayment=T, phoneCode=SGS5, phoneType=Android, phoneName="Samsung GALAX
       Y S5", planCode=1400POST5L90, planType=PostPaid, planPrice=89.99, planName="1400 Minute Family", planDes
       cription="Nationwide 1400 Minutes, Unlimited Mobile to Mobile, Unlimited Night & Weekend, Unlimited Dat
       a", cardNumber=3569948084568945, networkProviderName=Splunktel}] Showless
# time: 1595395363.33
α cribl_pipe: business_event
α host: 127.0.0.1
a index: cribl
a source: /opt/tibco/tra/apps/ESB/logs/business_event.log
\alpha sourcetype: business_event
```

Event with hashed SSNs

- i In scenarios where you need to send unmodified values to certain Destinations (such as archival stores), you can narrow the Mask Function's scope by setting the associated Route's Output field.
 - For further masking examples, see Masking and Obfuscation.

Numerify

Description

The Numerify Function converts event fields that are numbers to type number.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Ignore Fields: Specify fields to **not** numerify. Type in field names, separated by hard returns. By default, this list is empty, and Numerify applies to **all** fields. Supports wildcards (*) and nested addressing.

i Double Negatives

Negated terms are also supported. When you negate field names, the fields list is order-sensitive. E.g., !foobar before foo* means "Ignore all fields that start with foo, except foobar." However, !foo* before * means "Ignore all fields, except for those that start with foo."

Format: Optionally, reformat or truncate the extracted numeric value. Select one of:

- None: Applies no reformatting (the default).
- **Floor**: Rounds the number down to the lower adjacent integer (truncates it).
- **Ceil**: Rounds the number up to the higher adjacent integer, removing decimal digits.

- Round: Rounds (truncates) the number to a specified number of digits. This option exposes an extra field:
 - **Digits**: Number of digits after the decimal point. Enter a value between 0 20; defaults to 2.

Examples

Scenario A:

Assume an event whose text contains a numeric value that must be extracted to perform some numeric analysis. The text looks like this:

```
version=11.5.0.0.1.1588476445
```

We can extract the numeric value by chaining together two Functions:

- 1. A Regex Extract Function. Set its **Regex** field to /version=(?<ver>\d+)/, to capture the first set of digits found in the event string.
- 2. Then use Numerify.

This captures the substring 11 and converts it to a numeric 11 value.

Scenario B:

Assume email transaction log events like the sample below. The final field is the message's size, in bytes. We want to extract this as a numeric value, for analysis in LogStream or downstream services:

```
03:19 03:22 SMTPD (00180250) [209.221.59.70] C:\IMail\spool\D28de0018025017cd.SMD 3827
```

Again, we can accomplish this with two Functions:

- 2. Then use Numerify.

Parser

Description

The Parser Function can be used to extract fields out of events, or to reserialize (rewrite) events with a subset of fields. Reserialization will maintain the format of the events.

For example: If an event contains comma-delimited fields, and fieldA and fieldB are filtered out, those fields' positions will be set to null, but not deleted completely.

Parser cannot remove fields that it did not create. A subsequent Eval Function can do so.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Operation mode: Extract will create new fields. **Reserialize** will extract, filter fields, and then reserialize.

Type: Parser/Formatter type to use. Options:

- CSV
- Extended Log File Format (ELFF)
- Common Log Format (CLF)
- K=V Pairs
- JSON
- Delimited Values

Setting **Type** to **Delimited Values** displays the following extra options:

- **Delimiter**: Delimiter character to split value. Defaults to comma (,). You can also specify pipe (|) or tab characters.
- Quote char: Character used to quote literal values. Defaults to ".
- * Escape char: Character used to escape delimiter or quote characters.
 Defaults to \ .
- **Null value**: Field value representing the null value. These fields will be omitted. Defaults to .

Library: Select an option from the Parsers Library.

Source field: Field that contains text to be parsed. Not usually needed in Serialize mode.

Destination field: Name of field in which to add extracted and serialized fields. If multiple new fields are created and this setting is configured then all new fields are created as elements of an array with the array name set to the name specified for this setting. If you want all new fields to be independent, rather than in an array, then specify them using **List of fields** below. (Extract and Serialize modes only.)

Clean fields: This option appears for **Type: K=V Pairs**. Toggle to Yes to clean field names by replacing non-alphanumeric characters with _ . This will also strip leading and trailing " symbols.

List of fields: Fields expected to be extracted, in order. If not specified, Parser will auto-generate fields.

Fields to keep: List of fields to keep. Supports wildcards (*). Takes precedence over **Fields to remove**. Nested addressing supported.

Fields to remove: List of fields to remove. Supports wildcards (*). Cannot remove fields matching **Fields to keep**. Nested addressing supported.

i Negated terms are supported in both Fields to remove and Fields to keep. When you use negated terms, the list is order-sensitive.
E.g., !foobar, foo* means "All fields that start with foo, except foobar." However, !foo*, * means "All fields, except for those that start with foo."

Fields filter expression: Expression to evaluate against {index, name, value} context of each field. Return truthy to keep, falsy to remove field. Index is zero-based.

How Fields Settings Interact

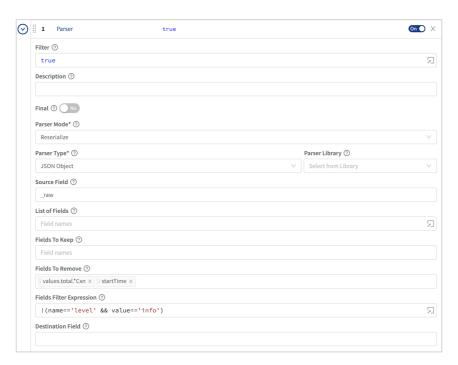
The **Fields to keep**, **Fields to remove**, and **Fields filter expression** settings interact as follows:

- Order of evaluation: Fields to keep > Fields to remove > Fields filter expression.
- If a field is in both Fields to keep and Fields to remove, Fields to keep takes precedence.
- If a field is in both Fields to remove and Fields filter expression, Fields to remove takes precedence.

Example 1

Insert the following sample, using Preview > Add a Sample > Paste a Sample: 2019/06/24 05:10:55 PM Z a=000,b=001,c=002,d=003,e=004,f=005,g1=006,g2=007,g3=008

Create the following test Parser Function (or import this Pipeline: https://raw.githubusercontent.com/weeb-cribl/cribl-samples/master/parser/functions/parser/parser_1.json).



Parser Function initial configuration

First, set the Parser type to Key=Value Pairs.

Scenario A:

Keep fields a, b, c. Drop the rest.

Expected result: a, b, c

• Fields to Keep: a, b, c

• Fields to Remove: *

• Fields Filter Expression:

Result: The event will gain four new fields and values, as follows.

• a: 000

• b: 001

• c: 002

• cribl_pipe: parser2



Scenario A result

You can check your stats by clicking the **Preview** pane's **Basic Statistics** (chart) button. In the resulting pop-up, the **Number of Fields** should have incremented ty four.

Now that you have the hang of it, try out the other simple scenarios below.

Scenario B:

Keep fields a, b, those that start with g. Drop the rest.

Expected result: a, b, g1, g2, g3

• Fields to keep: a, b

• Fields to remove: [empty]

• Fields filter expression: name.startsWith('g')

Scenario C:

Keep fields a, b, those that start with g but only if value is 007. Drop the rest.

Expected result: a, b, g2

- Fields to keep: a, b
- Fields to remove: [empty]
- Fields filter expression: name.startsWith('g') & value='007'

Scenario D:

Keep fields a, b, c, those that start with g, unless it's g1. Drop the rest.

Expected result: a, b, c, g2, g3

- Fields to keep: a, b, c
- Fields to remove: g1
- Fields filter expression: name.startsWith('g')

Scenario E:

Keep fields a, b, c, those that start with g but only if index is greater than 6. Drop the rest.

Expected result: a, b, c, g2, g3

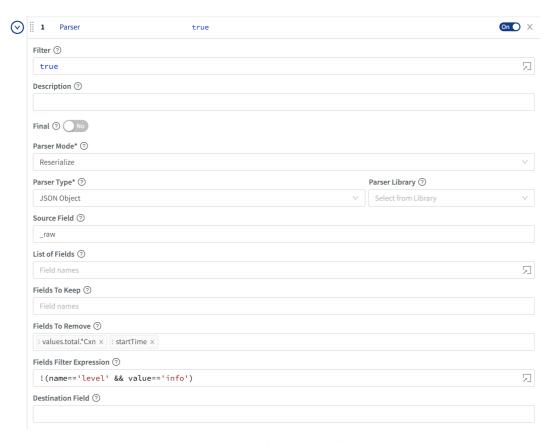
- Fields to keep: a, b, c
- Fields to remove: [empty]
- Fields filter expression: name.startsWith('g') & index>6
 - i The index refers to the location of a field in the array of all fields extracted by this Parser. It is zero-based. In the case above, g2 and g3 have index values of 7 and 8, respectively.

Example 2

Assume we have a JSON event that needs to be **reserialized**, given these requirements:

- 1. Remove the level field only if it's set to info.
- 2. Remove the startTime field, and all fields in the values.total. path that end in Cxn .

Parser Function configuration:



Parser Function configuration for Example 2

JSON event after being processed by the Function:

```
□ "_raw":{
□ "_raw":{
                                   "channel": "server"
   "channel": "server"
                                   "endTime":1549503300000
   "endTime":1549503300000
                                   "keyCount":0
   "keyCount":0
                                   "level": "info"
   "level":"info"
                                   "message":"_raw stats"
   "message":"_raw stats"
                                   "startTime": 1549503240000
   "startTime":1549503240000
                                   "time":1549503300401
   "time":1549503300401
                                   □ "values":{
  □ "values":{
                                      □ "total":{
     □ "total":{
                                         "activeCxn":2
        "activeCxn":2
                                         <del>"closeCxn"</del>:4
        "closeCxn":4
                                         "inBytes":61724
        "inBytes":61724
                                         "inEvents":210
        "inEvents":210
                                         "openCxn":4
        "openCxn":4
                                         "outBytes":61724
        "outBytes":61724
                                         "outEvents":210
        "outEvents":210
                                      }
                                   }
   }
                                }
}
```

Example 2 event transformation

Example 3

Insert the following sample, using **Preview > Add a Sample > Paste a Sample**:

```
2019/06/24 15:25:36 PM Z
a=000,b=001,c=002,d=003,e=004,f=005,g1=006,g2=007,g3=008,
```

For all scenarios below, first create a Parser Function to extract all fields, by setting the **Parser type** to Key=Value Pairs . Then add a second Parser Function with the configuration shown under **Parser 2**.

Scenario A:

```
Serialize fields a, b, c, d in CSV format.
```

Expected result: _raw field will have this value 000,001,002,003

Parser 2:

- Operation mode: Serialize
- Source field: [empty]
- Destination field: [empty]

- Type: CSV
- List of fields: a, b, c, d (needed for positional formats)

Scenario B:

Serialize fields a, b, c in JSON format, under a field called bar.

```
Expected result: bar field will be set to: {"a":"000","b":"001","c":"002","d":"003"}
```

Parser 2:

- Operation mode: Serialize
- Source field: [empty]
- Destination field: bar
- Type: JSON
- List of fields: [empty]
- Fields to keep: a, b, c, d

Publish Metrics

Description

The Publish Metrics Function extracts, formats, and outputs metrics from events.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Overwrite: If set to Yes, overwrite previous metric specs. Otherwise, append. Defaults to No.

Metrics

Add Metrics: List of metrics to extract from the event and format. Destinations can pass the formatted metrics to a metrics aggregation platform. Click **Add Metrics** to add new rows containing the following options:

- Event field name: The name of the field (in the event) that contains the metric value.
- Metric name expression: JavaScript expression to evaluate the metric field name. Defaults to the Event field name value.
 - The JavaScript expression will evaluate the metric field name only after the metrics are processed for transport to the Destination. While in the processing Pipeline, the metric name expression appears as a literal.
- Metric type: Select Counter, Timer, or Gauge (the default).

Remove Metrics: Optionally, enter a List of field names to look for when removing metrics. Where a metric's field name matches an element in this list, LogStream will remove that metric from the event.

Dimensions

Add Dimensions: Optional list of dimensions to associate with every extracted metric value. If this Function is used to process output from the Aggregations Function, leave this field blank, because dimensions will be automatically discovered. Defaults to !_* *.

Remove Dimensions: Optional list of dimensions to associate with every extracted metric value. Leave blank if this function is used to process output from the Aggregation function as dimensions will be automatically discovered. If this Function is used to process output from the Aggregations Function, leave this field blank, because dimensions will be automatically discovered.

i The Add Dimensions and Remove Dimensions fields support wildcards and negated terms. When you use negated terms, the list is order-sensitive. E.g., !foobar before foo* means "All fields that start with foo, except foobar." However, !foo* before * means "All fields, except for those that start with foo."

Overwrite: If set to Yes, overwrite previous metric specs. Otherwise, append. Defaults to No.

Fields Color Coding

On the right Preview pane's **OUT** tab, the Publish Metrics Function adds the following color codes to field labels:

```
1  m  a _metric: metrics_pool.num_metrics

2020-03-05  a _metric_type: gauge

18:42:52.707  # _value: 229

a _cribl_pipe: publish_some_metrics

# _earliest: 1583451770704

a _host: zebrastripes

# _latest: 1583451772704

a _source: _cribl
```

Dimension: purple | Value: cyan (light blue) | Info: dark blue

These are in addition to the color codes applied to field values, which are listed here.

Examples

Scenario A:

Assume we're working with AWS VPC Flowlog events that have the following structure:

version account_id interface_id srcaddr dstaddr srcport dstport protocol packets
bytes start end action log_status

For example:

```
2 99999XXXXX eni-02f03c2880e4aaa3 10.0.1.70 10.0.1.11 9999 63030 6 6556 262256 1554562460 1554562475 ACCEPT OK
```

... and we want to use values of packets and bytes as metrics across these dimensions: action, interface_id, and dstaddr.

To reference the packets and bytes fields by name, as 'packets' and 'bytes', our Pipeline will need a Parser Function before the Publish Metrics Function.

Parser Function

Filter: Set as needed
Operation mode: Extract

Type: Extended Log File Format (automatically set when specifying a library)

Library: AWS VPC Flow Logs

Source: _raw

(No need to specify any other fields.)

Publish Metrics Function

Below, the metric_name prefix was arbitrarily chosen. Because there is no JavaScript expression to evaluate – i.e. this is literal text – the strings specified for the Metric name expression will be identical to those in the final metrics data sent to the Destination. See Raw Output below.

Metrics

Event Field NaLme	Metric Name Expression	Metric Type
bytes	`metric_name.bytes`	Gauge
packets	`metric_name.packets`	Gauge

Dimensions

Dimensions	
action interface_id dstaddr	

All specified dimension names must align with those from the original event. When you preview the Function's output, the metrics and dimensions will all have special highlighting to separate them from other fields. Additional highlighting is used to differentiate the metrics from the dimensions. (If one or more metrics/dimensions are not highlighted as expected, check the Function's configuration.)

Raw Output

```
metric_name.bytes:262256|g#action:REJECT,interface_id:eni-
02f03c2880e4aaa3,dstaddr:10.0.1.11

metric_name.packets:6556|g#action:REJECT,interface_id:eni-
02f03c2880e4aaa3,dstaddr:10.0.1.11
```

i Compatible Destinations

All text after the # symbol represents the dimensions as key-value pairs. In order for dimension data to be included in metrics, the Destination type cannot be standard **StatsD**.

Formatted Output

```
{
  "action": "REJECT",
  "interface_id": "eni-02f03c2880e4aaa3",
  "dstaddr": "10.0.1.11",
  "metric_name.bytes": 262256,
  "metric_name.packets": 6556,
}
```

Scenario B:

Assume that we want to extract some metrics from specific fields in PANOS logs, whose events have the following structure:

future_use_0,receive_time, serial_number, type, threat_content_type, future_use_1, generated_time, source_ip, destination_ip, nat_source_ip, nat_destination_ip, rule_name, source_user, destination_user, application, virtual_system, source_zone, destination_zone, inbound_interface, outbound_interface, log_action, future_use_2, session_id, repeat_count, source_port, destination_port, nat_source_port, nat_destination_port, flags, protocol, action, bytes, bytes_sent, bytes_received, packets, start_time, elapsed_time, category, future_use_3, sequence_number, action_flags, source_location, destination_location, future_use_4, packets_sent, packets_received, session_end_reason, device_group_hierarchy_level_1, device_group_hierarchy_level_2, device_group_hierarchy_level_3, device_group_hierarchy_level_4, virtual_system_name, device_name, action_source, source_vm_uuid, destination_vm_uuid, tunnel_id_imsi, monitor_tag_imei, parent_session_id, parent_start_time, tunnel_type, sctp_association_id, sctp_chunks, sctp_chunks_sent, sctp_chunks_received

For example:

```
Jan 10 10:19:15 DMZ-internal.nsa.gov 1,2019/01/10
10:19:15,001234567890002,TRAFFIC,drop,2304,2019/01/10
10:19:15,209.118.103.150,160.177.222.249,0.0.0.0,0.0.0.0,InternalServer,,,not-applicable,vsys1,inside,z1-FW-Transit,ethernet1/2,,All traffic,2019/01/10
10:19:15,0,1,63712,443,0,0,0×0,udp,deny,60,60,0,1,2019/01/10
10:19:15,0,any,0,0123456789,0×0,Netherlands,10.0.0.0-10.255.255.255,0,1,0,policy-deny,0,0,0,0,TDMZ-internal,from-policy,,,0,,0,,N/A,0,0,0,0,1202585d-b4d5-5b4c-aaa2-d80d77ba456e,0
```

Our goal is to use the four values of bytes_sent, bytes_received, packets_sent, and packets_received as metrics across these dimensions: destination_ip, inbound_interface, outbound_interface, and destination_port.

Here again, our Pipeline will need a Parser Function before the Publish Metrics Function.

Parser Function

Filter: Set as needed Operation mode: Extract

Type: Extended Log File Format (automatically set when specifying a Library)

Library: Palo Alto Traffic

Source: _raw

(No need to specify any other fields.)

Publish Metrics Function

Set up the Publish Metrics Function as follows.

Metrics

Event Field Name	Metric Name Expression	Metric Type
bytes_sent	metric.\${host}.bytes_sent	Counter
bytes_received	metric.\${host}.bytes_rcvd	Counter
packets_sent	metric.\${host}.pkts_sent	Counter
packets_received	metric.\${host}.pkts_rcvd	Counter

Added Dimensions

destination_ip , inbound_interface , outbound_interface , destination_port

Raw Output

 $\label{lem:metric.10.10.12.192.bytes_sent:60|c|\#destination_ip:160.177.222.249, inbound_interface:e:ethernet1/2, destination_port:443$

 $\label{lem:metric.10.10.12.192.bytes_rcvd:0|c|\#destination_ip:160.177.222.249, inbound_interface: ethernet1/2, destination_port:443$

metric.10.10.12.192.pkts_sent:1|c|#destination_ip:160.177.222.249,inbound_interface:e
thernet1/2,destination_port:443

metric.10.10.12.192.pkts_rcvd:0|c|#destination_ip:160.177.222.249,inbound_interface:e
thernet1/2,destination_port:443

Here again, all text after the # symbol represents the dimensions as key-value pairs. (See the Compatible Destinations note above.) Unlike the first example, this example uses JavaScript expressions, which you can see evaluated in the raw output where the \${host} has been converted to 10.10.12.192.

Regex Extract

Description

The Regex Extract Function extracts fields using regex named groups. (In Splunk, these will be index-time fields). Fields that start with ___ (double underscore) are special fields in Cribl LogStream. They are ephemeral: they can be used by any Function downstream, but will not be added to events, and will not exit the Pipeline.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Regex: Regex literal. Must contain named capturing groups, e.g.: (?<foo>bar). Can contain special _NAME_N and _VALUE_N capturing groups, which extract both the name and value of a field, e.g.: (?<_NAME_0>[^\s=]+)=(?<_VALUE_0>[^\s]+). Defaults to empty. See Examples below.

Additional regex: Click + Add Regex to chain extra regex conditions.

Source field: Field on which to perform regex field extraction. Nested addressing is supported. Defaults to _raw .

Advanced Settings

Max exec: The maximum number of times to apply the Regex to the source field when the global flag is set, or when using <code>_NAME_N</code> and <code>_VALUE_N</code> capturing groups. Named capturing groups will always use a value of <code>1.Defaults</code> to <code>100</code>.

Field name format expression: JavaScript expression to format field names when <code>_NAME_n</code> and <code>_VALUE_n</code> capturing groups are used. E.g., to append XX to all field names, use: `\${name}_XX` (backticks are literal). If not specified, names will be sanitized using regex: /^[_0-9]+|[^a-zA-Z0-9_]+/g . The original field name is in the global <code>name</code> .

Overwrite existing fields: Whether to overwrite existing event fields with extracted values. If set to No (the default), existing fields will be converted to an array. If toggled to Yes, Regex Extract will create array fields if applied multiple times, or if fields exist. (E.g., if src_ip is extracted in an input Pipeline where it is assigned a value of 10.1.2.2, and is also in a processing Pipeline with a value of 10.2.3.3, then the resulting field is ["10.1.2.2", "10.2.3.3"].)

Examples

Example 1

Assume a simple event that looks like this: metric1=23 metric2=42 dc=23 abc=xyz

Extract only the metric1 field:

Regex: metric1=(?<metric1>\d+)

Result: metric1:"23"

Example 2

Use the first line of the sample here:

https://github.com/weeb-cribl/cribl-

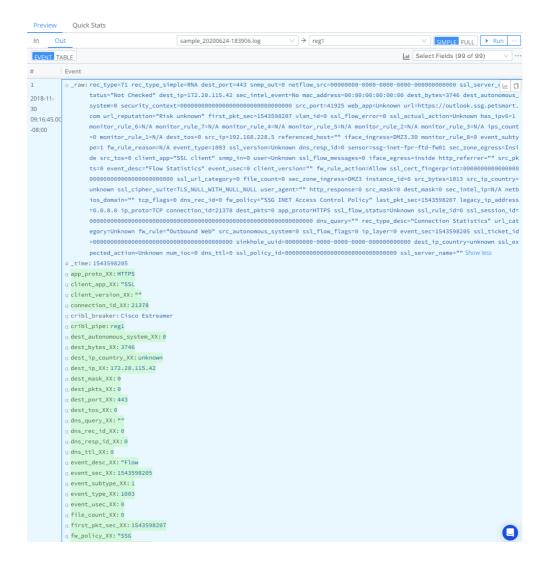
samples/blob/master/parser/functions/parser/cisco_estreamer.log

Extract all k=v pairs, and add an _XX suffix to the end of the extracted fields:

Regex: (?<_NAME_0>[^\s]+)=(?<_VALUE_0>[^\s]+)

Field Name Format Expression: `\${name}_XX`

Result:



Example 3

This example uses similar syntax as Example 2, but with a more complex event structure.

In the right **Sample Data** pane, click **Paste** and insert the following sample:

Sample Data

```
<134>1 2020-12-22T17:06:08Z CORP_INT_NLB CheckPoint 18160 - [action:"Accept"; conn_direction:"Ir
```

This event is from a CheckPoint Firewall CMA system. With this type of event structure, properly extracting each event field into a separate metadata field requires two-stage processing. So we'll use two Regex Extract Functions.

The first Regex Function splits the event to separate the actual data from the header information. We'll split after the CheckPoint 18160 string, by capturing everything between the [and]:

Regex; \[(?<__fields>.*)\]

Source: _raw

Next, add this second Regex Extract Function to extract all k=v pairs:

```
Regex: (?<_NAME_0>[^ :]+):(?<_VALUE_0>[^;]+);
```

Source: __fields

Results:

```
on:"Internal"; flags:"4606212"; ifdir:"inbound"; ifname:"bond2.1025"; logid:"0"; logui
                 d:"{0x5fe25889,0x0,0x80ad57cd,0xeb91c0c3}"; origin:"192.168.20.54"; originsicname:"CN=TST3
                 2-VSX0-FW-DC-01_tst302-shd,0=CORP-SEC-SHRD-CMA..t7xpcz"; sequencenum:"3"; time:"160865676
                 8"; version:"5"; __policy_id_tag:"product=VPN-1 & FireWall-1[db_tag={15E4B45A-663B-5B49-BD
                 59-CD9B9F21AA16};mgmt=SHRDFW01CON;date=1608236862;policy_name=TEST-SHRD-POL\]"; dst:"192.1
                 68.79.20"; log_delay:"1608656768"; layer_name:"TEST-SHRD-POL Security"; layer_uuid:"e914c2
                  f3-d7bd-4a77-8e7a-7a5e403447aa";\ match\_id:"1";\ parent\_rule:"0";\ rule\_action:"Accept";\ rule\_action:"Accept ";\ rule\_action:"Accept "
                  _uid:"001ab86d-d201-4b61-9b64-0fede1a9f059"; product:"VPN-1 & FireWall-1"; proto:"17"; s_p
                  ort:"45519"; service:"123"; service_id:"ntp-udp"; src:"192.168.79.22"; ] Showless
  # _time: 1608656768
 α action: "Accept"
 a conn_direction: "Internal"
 α cribl_breaker: Break on newlines
 α cribl_pipe: asfasfdasfd
a dst: "192.168.79.20"
α flags: "4606212"
a ifdir: "inbound"
α ifname: "bond2.1025"
α layer_name: "TEST-SHRD-POL Security"
α layer_uuid: "e914c2f3-d7bd-4a77-8e7a-7a5e403447aa"
α log_delay: "1608656768"
a logid: "0"
a loguid: "{0x5fe25889,0x0,0x80ad57cd,0xeb91c0c3}"
α match_id: "1"
α origin: "192.168.20.54"
α originsicname: "CN=TST32-VSX0-FW-DC-01_tst302-shd,0=CORP-SEC-SHRD-CMA..t7xpcz"
a policy_id_tag: "product=VPN-1 & FireWall-1[db_tag={15E4B45A-663B-5B49-BD59-CD9B9F21AA16}
α product: "VPN-1 & FireWall-1"
a proto: "17"
α rule_action: "Accept"
g rule uid: "001ab86d-d201-4b61-9b64-0fede1a9f059"
```

i For further examples, see Using Cribl to Analyze DNS Logs in Real Time – Part 2.

Redis

Description

The Redis Function interacts with Redis stores, setting and getting key-hash and key-value combinations. Redis' in-memory caching of these key pairs enables large lookup tables that would be cumbersome with a .CSV or binary lookup file.

You can use LogStream Collectors (e.g., a REST Collector) to retrieve reference data from desired endpoints, and then use this Function to store the data on Redis and retrieve it to enrich your production data. Note that LogStream does not cache the data returned from this Redis Function.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to

Result field: Name of the field in which to store the returned value. (Leave empty to discard the returned value.)

Command: Redis command to perform. Required. (A complete list of Redis commands is at: https://redis.io/commands.)

Key: A JavaScript expression to compute the value of the key to operate on. Can also be a constant, e.g.: username. This is a required field. Click the icon at right to open a validation modal.

Args: A JavaScript expression to compute arguments to the operation. Can return an array. Click the icon at right to open a validation modal.

Redis URL: Redis URL to connect to. The formatis: [redis[s]:]//[[user] [:password@]][host][:port][/db-number][?db=db-number[&password=bar[&option=value]]]

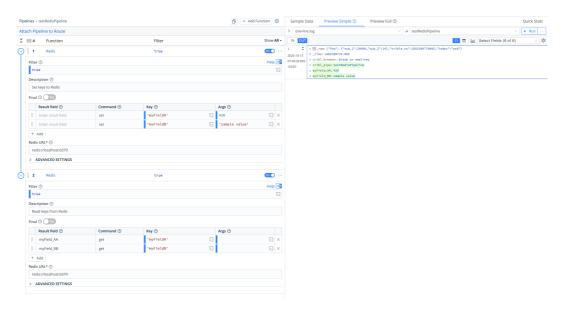
For example: redis://user:secret@localhost:6379/0?foo=bar&qux=baz

Advanced Settings

Max blocking time: Maximum amount of time (in seconds) before assuming that Redis is down and passing events through. Defaults to 60 seconds. Use 0 to disable timeouts.

Example

This Pipeline demonstrates the use of a pair of Redis Functions. The first Function sets two key-value pairs in Redis. The second Function their values, by key, into two corresponding new **Result fields**.



Redis set and get Functions

Redis Function #1

Description: Set keys to Redis

Command: set
Key: 'myFieldA'

Args: 420

Command: set
Key: 'myFieldB'

Args: 'sample value'

Redis Function #2

Description: Read keys from Redis

Result field: myField_AA

Command: get
Key: 'myFieldA'

Result field: myField_BB

Command: get
Key: 'myFieldB'

Regex Filter

Description

The Regex Filter Function filters out events based on regex matches.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Regex: Regex to test against. Defaults to empty.

Additional regex: Click + Add Regex to chain extra regex conditions.

Field: Name of the field to test against the regex. Defaults to <code>_raw</code> . Supports nested addressing.

Examples

See Regex Filtering for examples.

Rename

Description

The Rename Function is designed to change fields' names or reformat their names (e.g., by normalizing names to camelcase). You can use Rename to change specified fields (much like the Eval Function), or for bulk renaming based on a JavaScript expression (much like the Parser Function).

Compared to these alternatives, Rename offers a streamlined way to alter only field names, without other effects.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Optionally, enter a simple description of this step in the Pipeline. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Base fields: Enter one or more source field names to rename. If empty, rename will be performed on top-level fields.

Rename fields: Each row here is a key-value pair that defines how to rename fields. The current name is the key, and the new name is the value. Click + Add Field to add more rows.

- **Current name**: Original name of the field to rename. You must quote literal identifiers (non-alphanumeric characters such as spaces or hyphens).
- New name: New or reformatted name for the field. Here again, you must quote literals.

Renaming expression: An optional JavaScript expression (or literal) used to compute multiple fields' new names. This expression is evaluated against a {name, value} context, and the expression returns a value with which to rename fields.

You can use both Rename fields (to rename specified field names) and Renaming expression* (to globally rename fields) in the same Function. The Rename fields** strategy will execute first.

Example

Change the level field, and all fields that start with out, to all-uppercase.

Example event:

```
{"inEvents": 622,
    "level": "info",
    "outEvents": 311,
    "outBytes": 144030,
    "activeCxn": 0,
    "openCxn": 0,
    "closeCxn": 0,
    "activeEP": 105,
    "blockedEP": 0
}
```

Rename Fields:

```
Current Name: level
New Name: LEVEL
Renaming Expression: name.startsWith('out') ? name.toUpperCase():
name
```

Event after Rename:

```
{"inEvents": 622,
   "LEVEL": "info",
   "OUTEVENTS": 311,
   "OUTBYTES": 144030,
   "activeCxn": 0,
   "openCxn": 0,
   "closeCxn": 0,
   "activeEP": 105,
   "blockedEP": 0
}
```

Rollup Metrics

Description

The Rollup Metrics Function merges/rolls up frequently generated incoming metrics into more manageable time windows.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Optional description of this Function's purpose in this Pipeline. Defaults to empty.

Final: If toggled to Yes, stops data from being fed to downstream Functions. Defaults to No.

Dimensions: List of data dimensions across which to perform rollups. Supports wildcards. Defaults to * wildcard, meaning all original dimensions.

Time window: The time span over which to roll up (aggregate) metrics. Must be a valid time string (e.g., 10s). Must match pattern: \d+[sm]\$.

⚠ With high-cardinality data, beware of setting long time windows. Doing can cause high memory consumption and/or lost data, because memory is flushed upon restarts and redeployments.

Gauge update: The operation to use when rolling up gauge metrics. Defaults to Last; other options are Maximum, Minimum, or Average.

Examples

Scenario A:

Assume that you have metrics coming in at a rate that is too high. For example, LogStream's internal metrics come in at a 2s interval.

To roll up these metrics to 1-minute granularity, you would set up the Rollup Metrics Function with a **Time Window** value of 60s.

Scenario B:

Assume that you have metrics coming up with multiple dimensions – e.g. host, source, data_center, and application. You want to aggregate these metrics to eliminate some dimensions.

Here, you would configure Rollup Metrics Function with a **Time Window** value that matches the metrics' generation – e.g., 10s . In the **Dimensions** field, you would remove the default \star wildcard, and would specify only the dimensions you want to keep – e.g.: host , data_center .

Sampling

Description

The Sampling Function filters out events, based on an expression and a sampling rate.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

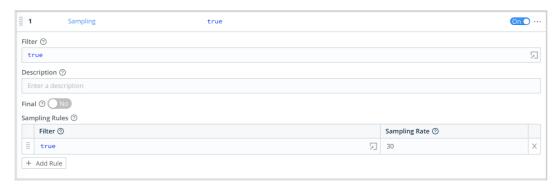
Final: If true, stops data from being fed to downstream Functions. Defaults to No .

Sampling rules: Events matching these rules will be sampled at the rates you specify:

- Filter: Filter expression matching events to be sampled. Use true to match all.
- Sampling rate: Enter an integer N. (Defaults to 1.) Sampling will pick 1/N events matching this rule.

How It Works

Setting this Function's **Sampling rate** to 30 would mean that only 1 of every 30 events would be kept.



Let's assume that we save this setting, and then capture data from a datagen Source by selecting Preview > Start a Capture > Capture. In the Capture Sample Data modal, select: 100 seconds, 100 events, and As they come in. Then start the capture, and Save as Sample File.

Next, in the **Preview** pane, click **Simple** beside the new file's name. If you then click the **Basic Statistics** (chart) button, you should see that we've kept about 4 of the original 100 events, or close to 1 in 30.

	Full Event Length ⑦	Number of Fields ⑦	Number of Events ⑦
IN	28.82KB	41	100
OUT	1.42KB	38	4
DIFF	↓ -95.08%	↓ -7.32%	↓ -96.00%

Examples

See Sampling for examples.

Serialize

Description

Use the Serialize Function to serialize an event's content into a predefined format.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Type: Data output format. Defaults to CSV.

Library: Browse Parser/Formatter library.

Fields to serialize: Required for CSV, ELFF, and CLF Types. (All other formats support wildcard field lists.)

Source field: Field containing the object to serialize. Leave blank to serialize top-level event fields.

Destination field: Field to serialize the data into. Defaults to _raw .

Examples

Scenario A: JSON to CSV

```
Assume a simple event that looks like this: {"time":"2019-08-
25T14:19:10.240Z", "channel": "input", "level": "info", "message": "initializing input", "type": "kafka"}
```

We want to serialize these fields: _time , channel , level , and type into a single string, in CSV format, stored in a new destination field called test .

To properly extract the key-value pairs from this event structure, we'll use a built-in Event Breaker:

- 1. Copy the above sample event to your clipboard.
- 2. In the **Preview** pane, select **Paste a Sample**, and paste in the sample event.
- 3. Under **Select Event Breaker**, choose **ndjson** (newline-delimited JSON), and click **Save as a Sample File**.

Now you're ready to configure the Serialize Function, using the settings below:

Type: CSV

Fields to Serialize: _time channel level type

Destination Field: test Source Field: [leave empty]

Result: test: 1566742750.24, input, info, kafka

In the new test field, you now see the time, channel, level, and type keys extracted as top-level fields.

Scenario B: CSV to JSON

Let's assume that a merchant wants to extract a subset of each customer order, to aggregate anonymized order statistics across their customer base. The transaction data is originally in CSV format, but the statistical data must be in JSON.

Here's a CSV header (which we don't want to process), followed by a row that represents one order:

```
orderID, custName, street, city, state, zip
20200622102822, john smith, 100 Main St., Anytown, AK, 99911
```

To convert to JSON, we'll need to first parse each field from the CSV to a manipulable field in the Pipeline, which the Serialize Function will be able to reference. In this example, the new manipulable field is message.

Use the Parser Function:

Filter: true

Operation mode: Extract

Type: CSV

Source field: _raw

Destination field: message

List of fields: orderID custName street city state zip

Now use the Serialize Function:

Filter: true Type: JSON

Fields to serialize: city state

Source field: message

Destination field: orderStats

Suppress

Description

The Suppress Function suppresses events over a time period, based on evaluating a key expression.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No .

Key expression: Suppression key expression used to uniquely identify events to suppress. For example, `\${ip}:\${port}` will use the fields ip and port from each event to generate the key.

Number to allow: The number of events to allow per time period. Defaults to 1.

Suppression period (sec): The number of seconds to suppress events after 'Number to allow' events are received. Defaults to 300.

Drop suppressed events: Specifies if suppressed events should be dropped, or just tagged with suppress=1. Defaults to Yes, meaning drop.

Advanced Settings

Maximum cache size: The maximum number of keys that can be cached before idle entries are removed. Before changing the default 50000, contact Cribl Support to understand the implications.

Suppression period timeout: The number of suppression periods of inactivity before a cache entry is considered idle. This defines a multiple of the **Suppression period (sec)** value. Before changing the default 2, contact Cribl Support to understand the implications.

Num events to trigger cache clean-up: Check cache for idle sessions every N events when cache size exceeds the **Maximum cache size**. Before changing the default 10000, contact Cribl Support to understand the implications.

Examples

In the examples below, **Filter** is the Function-level Filter expression:

1. Suppress by the value of the host field:

Filter: true

Key expression: host Number to allow: 1

Suppression period (sec): 30

Using a datagen sample as a source, generate at least 100 events over 2 minutes.

Result: One event per unique host value will be allowed in every 30s. Events without a host field will **not** be suppressed.

2. Suppress by the value of the host and port tuple:

Filter: true

Key expression: `\${host}:\${port}`

Number to allow: 1

Suppression period (sec): 300

Result: One event per unique host: port tuple value will be allowed in every 300s.

- ⚠ Suppression will **also** apply to events without a host or a port field. The reason is that if field is not present, `\${field}` results in the literal undefined.
- 3. To **guarantee** that suppression applies **only** to events with host and port, check for their presence using a Filter:

Filter: host≠undefined & port≠undefined

Key expression: `\${host}:\${port}`

Number to allow: 1

Suppression period (sec): 300

4. Decorate events that qualify for suppression:

Filter: true

Key expression: `\${host}:\${port}`

Number to allow: 1

Suppression period (sec): 300 Drop suppressed events: No

Result: No events will be suppressed. But all qualifying events will gain an added field suppress=1, which can be used downstream to further transform these events.

Tee

Description

The Tee Function tees events out to a command of choice, via stdin. The output is one JSON-formatted event per line. You can send the events to (for example) a local file on the LogStream worker. This can be useful in verifying the data being processed in a Pipeline.

The Filesystem/NFS Destination offers similar capability, but only after the data leaves the Pipeline. Tee, by comparison, can be inserted at any point in the Pipeline.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Command: Command to execute and receive events (via stdin) – one JSON-formatted event per line.

Args: Click + Add Arg to supply arguments to the command.

Restart on exit: Restart the process if it exits and/or we fail to write to it. Defaults to Yes.

Environment variables: Environment variables to set or overwrite. Click + Add Variable to add key/value pairs.

Communication Protocol

Data is passed to the command through its stdin, using the following protocol:

- First line: Metadata serialized in JSON, containing the following fields:
 - format: Serialization format for event. Defaults to JSON.
 - conf: Full Function configuration.
- Remaining: Payload.

Examples

Assume that we are parsing PANOS Traffic logs, and want to see how they look at a particular step in the processing Pipeline We'll assume that the Parser Function is already in place, so we'll insert the

Tee Function at any (arbitrary) later point in the Pipeline.

Scenario A:

The Tee Function itself requires only that we define the **Command** field. In this particular example, that **Command** will be tee itself.

We've also clicked + Add Arg, to specify a local output file in the resulting Args field. (A file path would normally be the first argument to a tee command executed from the command line. The LogStream user must have write permission on the specified file path.)

Command: tee

Args: /opt/cribl/foo.log

In this first scenario, assume that we have the Parser configured to parse, but not keep any fields. After changes are deployed and PANOS logs are received, if we tail foo.log, we'd see the following:

```
Line 1: {"format":"json","conf":{"restartOnExit":true,"env":
{},"command":"tee","args":["/opt/cribl/foo.log"]}

Line 2: {"_raw":"Oct 09 10:19:15 DMZ-internal.nsa.gov 1,2019/10/09

10:19:15,001234567890002,TRAFFIC,drop,2304,2019/10/09

10:19:15,209.118.103.150,160.177.222.249,0.0.0.0,0.0.0,InternalServer,,,not-applicable,vsys1,inside,z1-FW-Transit,ethernet1/2,,All traffic,2019/10/09

10:19:15,0,1,63712,443,0,0,0×0,udp,deny,60,60,0,1,2019/10/09

10:19:15,0,any,0,0123456789,0×0,Netherlands,10.0.0.0-10.255.255.255,0,1,0,policy-deny,0,0,0,0,,DMZ-internal,from-policy,,,0,,0,,N/A,0,0,0,0,1202585d-b4d5-5b4c-aaa2-d80d77ba456e,0","_time":1593185574.663,"host":"127.0.0.1"}
```

In Line 2 above, note that the _raw field makes up most of the contents, with only the _time and host fields added.

Scenario B:

Assume that we use the Tee Function, using the same **Command** and arguments, but we've modified the the Parser Function to retain five fields: receive_time, source_port, destination_port bytes_received, and packets_received.

This time, if we tail foolog, we'll see something like the following. If you compare this output to the previous output example, you'll notice the five fields appended to this event:

```
Line 3: {"_raw":"Oct 09 10:19:15 DMZ-internal.nsa.gov 1,2019/10/09 10:19:15,001234567890002,TRAFFIC,drop,2304,2019/10/09 10:19:15,209.118.103.150,160.177.222.249,0.0.0.0,0.0.0.0,InternalServer,,,not-applicable,vsys1,inside,z1-FW-Transit,ethernet1/2,,All traffic,2019/10/09 10:19:15,0,1,63712,443,0,0,0×0,udp,deny,60,60,0,1,2019/10/09 10:19:15,0,any,0,0123456789,0×0,Netherlands,10.0.0.0-10.255.255.255,0,1,0,policy-deny,0,0,0,0,TDMZ-internal,from-policy,,,0,,0,,N/A,0,0,0,0,1202585d-b4d5-5b4c-aaa2-
```

```
d80d77ba456e,0","_time":1593185606.965,"host":"127.0.0.1","receive_time":"2019/10/09
10:19:15","source_port":"63712","destination_port":"443","bytes_received":"0","packet
s_received":"0"}
```

i In this Function's **Command** field, you can specify commands other than tee itself. For example: By using nc as the command, and specifying localhost and a port number (as two separate arguments), you'll see event data being received via nc on the specified port.

Trim Timestamp

Description

The Trim Timestamp Function removes timestamp patterns from events, and (optionally) stores them in a specified field.

This Function looks for a timestamp pattern that exists between the characters indicated by numeric timestartpos and timeendpos fields. It removes timestartpos and timeendpos along with the timestamp pattern.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this step in the Pipeline. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Field name: Name of field in which to save the timestamp. (If empty, timestamp will not be saved to a field.)

Example

Remove the timestamp pattern (indicated by timestartpos and timeendpos) from _raw , and stash it in a field called time_field .

Example event before:

```
{"_raw": "Event [Event=UpdateBillingProvQuote, timestamp=1581426279, properties={JMSCorrelation:
"timestartpos":0,
"timestartpos":23
}```

**Field Name**: `time_field`

**Example Event after:**
```

```
{"_raw": "2020-05-22 16:32:11,359 Event [Event=UpdateBillingProvQuote, timestamp=1581426279, properties={JMSCorrelationID=NA, JMSMessageID=ID:ESP-PD.D2BB2D95F857B:FA323D61, orderType=RatePlanFeatureChange, quotePriority=NORMAL}", "time_field":"2020-05-22 16:32:11,359"
```

Unroll

Description

The Unroll Function accepts an array field – or an expression to evaluate an array field – and breaks/unrolls the array into individual events.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Source field expression: Field in which to find/calculate the array to unroll. E.g.: $_{raw}$, $_{raw.split(/\n/)}$. Defaults to $_{raw}$.

Destination field: Field (within the destination event) in which to place the unrolled value. Defaults to raw.

Example

Assume we want to break/unroll each line of this event:

Sample Event							
USER	PID %CI	PU %MEM	VSZ	RSS TTY	STAT	START	TIME COMMAND
root	1 0	.0 0.5	38000	5356 ?	Ss	2018	2:02 /lib/systemd/systemdsystem ·
root	2 0	0.0	0	0 ?	S	2018	0:00 [kthreadd]
root	3 0	0.0	0	0 ?	S	2018	1:51 [ksoftirqd/0]
root	5 0	0.0	0	0 ?	S<	2018	0:00 [kworker/0:0H]
root	7 0	0.0	0	0 ?	S	2018	3:55 [rcu_sched]
root	8 0	.0 0.0	0	0 ?	S	2018	0:00 [rcu_bh]

Settings

Source field expression: _raw.split(/\n/)

i The split() JavaScript method breaks _raw into an ordered set of substrings/values, puts these values into an array, and returns the array.

Destination field: _raw

Resulting Events

Event 1: USER	PID %CPU	%MEM VS	Z RSS	TTY	STAT	START	TIME	COMMAND
Event 2: root	1 0.0	0.5 3800	0 5356	?	Ss	2018	2:02	/lib/systemd/systemdsystem ·
Event 3: root	2 0.0	0.0	0 0	?	S	2018	0:00	[kthreadd]
Event 4: root	3 0.0	0.0	0 0	?	S	2018	1:51	[ksoftirqd/0]
Event 5: root	5 0.0	0.0	0 0	?	S<	2018	0:00	[kworker/0:0H]
Event 6: root	7 0.0	0.0	0 0	?	S	2018	3:55	[rcu_sched]
Event 7: root	8 0.0	0.0	0 0	?	S	2018	0:00	[rcu_bh]

XML Unroll

Description

The XML Unroll Function accepts a proper XML event with a set of elements, and converts the elements into individual events.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Unroll elements regex: Path to the array to unroll. E.g.:
 ^root\.child\.ElementToUnroll\$

Copy elements regex: Regex matching elements to copy into each unrolled event.

E.g.: ^root\.(childA|childB|childC)\$

Unroll index field: LogStream will add a field with this name, containing the 0-based index at which the element was located within the event. In Splunk, this will be an index-time field. Supports nested addressing. Name defaults to unroll_idx.

Pretty print: Whether to pretty print the output XML.

Examples

Assume that the following sample is ingested as a single event:

```
<state>NY</state>
        <city>New York</city>
    </Child>
    <Child>
        <state>NJ</state>
        <city>Edgewater</city>
    </Child>
    <Child>
       <state>CA</state>
        <city>Oakland</city>
    </Child>
    <Child>
        <state>CA</state>
        <city>San Francisco</city>
    </Child>
</Parent>
```

i If you insert this sample using Preview > Add a Sample > Paste a Sample, adjust Event Breaker settings to add the sample as a single event. One way to do this is to add a regex Event Breaker that (by design) will not match anything present in the sample. For example: /[\n\r]+donotbreak(?!\s)/. As of LogStream 2.3, you can also use the built-in Do Not Break Ruleset.

Set up the XML Unroll Function using these settings:

Unroll elements regex: ^Parent\.Child\$

Copy elements regex: ^Parent\.(myID|branchLocation)\$

Output 4 Events:

```
Resulting Events
# Event 1
<?xml version="1.0"?>
<Child>
  <myID>123456</myID>
  <branchLocation>US/branchLocation>
  <state>NY</state>
  <city>New York</city>
</Child>
# Event 2
<?xml version="1.0"?>
<Child>
  <myID>123456</myID>
  <branchLocation>US/branchLocation>
  <state>NJ</state>
  <city>Edgewater
</Child>
```

```
# Event 3
<?xml version="1.0"?>
<Child>
 <myID>123456</myID>
 <branchLocation>US/branchLocation>
 <state>CA</state>
 <city>Oakland</city>
</Child>
# Event 4
<?xml version="1.0"?>
<Child>
 <myID>123456</myID>
 <branchLocation>US</branchLocation>
 <state>CA</state>
 <city>San Francisco</city>
</Child>
```

Prometheus Publisher (beta)

Description

The Prometheus Publisher Function allows for metrics to be published to a Prometheus-compatible metrics endpoint. These can be upstream metrics received by LogStream, or metrics derived from the output of LogStream's Publish Metrics or Aggregation Functions. A Prometheus instance is responsible for collecting the metrics at that endpoint, and for performing its own processing of the metric data.

In the current LogStream version, the endpoint is: http://<worker_node_IP>:<apiport>/metrics. Within LogStream, that endpoint redirects from http://<worker_node_IP>:9000/metrics to http://<worker_node_IP>:9000/api/v1/metrics.

⚠ If used, this Function must follow any Publish Metrics or Aggregations Functions within the same Pipeline. This is to ensure that any data not originating from a metrics input is transformed into metrics format.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Fields to publish: Wildcard list of fields to publish to the Prometheus endpoint.

Advanced Settings

Batch write interval: How often, in milliseconds, the contents should be published. Defaults to 5000.

Passthrough mode: If set to No (the default), overrides the Final setting, and suppresses output to downstream Functions' Destinations. Toggle to Yes to allow events to flow to consumers beyond the Prometheus endpoint. In effect, when previewing the pipeline output what you'll see is your event fields will have strikethrough font applied to them. This does not mean the Prometheus function is not matching your events but rather indicative of the Passthrough being disabled.

Update mode: On the default No setting, suppresses output to downstream Functions' Destinations. (This overrides the Final setting.) Toggle to Yes to allow events to flow to consumers beyond the

Example

This example uses the same PANOS sample data as the Publish Metrics Function, and is similarly preceded in a Pipeline by a Parser Function that extracts fields from the PANOS log.

Filter: Set as appropriate.

Fields to publish: Set as appropriate. We'll use the default of * for this example. **Advanced settings**: Accept defaults.

After committing and deploying changes, you should be able to use a <code>curl</code> command (-L needed to follow the redirect mentioned above) to verify that metrics are being published, just a few seconds after data is ingested on an idle system.

```
curl output

$ curl -L http://<worker_node_IP>:9000/metrics

# TYPE perf_192_168_1_248_bytes_sent counter
metric_192_168_1_248_bytes_sent {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"

# TYPE perf_192_168_1_248_bytes_rcvd counter
metric_192_168_1_248_bytes_rcvd {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"

# TYPE perf_192_168_1_248_pkts_sent counter
metric_192_168_1_248_pkts_sent {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"

# TYPE perf_192_168_1_248_pkts_rcvd counter
metric_192_168_1_248_pkts_rcvd counter
metric_192_168_1_248_pkts_rcvd {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"
```

Now, we need to have Prometheus scrape the metrics. In this very basic example, you can add the target endpoint to the prometheus.yml file, under the scrape_configs -> static_configs section. Specify the endpoint in IP:port syntax, because Prometheus assumes (and requires) /metrics for all endpoints.

Restart Prometheus. Within just a few seconds, you should be able to use its query interface to retrieve metrics published by LogStream.

Reverse DNS (deprecated)

Description

The Reverse DNS Function resolves hostnames from a numeric IP address, using a reverse DNS lookup.

⚠ This Function is deprecated. Use the DNS Lookup Function's reverse lookup feature instead.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Lookup Fields

Lookup field name: Name of the field containing the IP address to look up.

⚠ If the field value is not in IPv4 or IPv6 format, the lookup is skipped.

Output field name: Name of the field in which to add the resolved hostname. Leave blank to overwrite the lookup field.

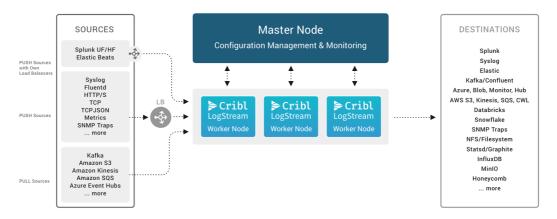
Reload period (minutes): How often to refresh the DNS cache. Use 0 to disable refreshes. Defaults to 60 minutes.

Example

Lookup field name: dest_ip Output field name: dest_host **Result**: See the dest_ip field, and the newly created dest_host field, in the events.

Sources

Cribl LogStream can receive data from various Sources, including Splunk, HTTP, Elastic Beats, Kinesis, Kafka, TCP JSON, and many others.



Push and Pull Sources

PUSH Sources

Supported data Sources that **send** to Cribl LogStream:

- Splunk TCP
- Splunk HEC
- Syslog
- Elasticsearch API
- TCP JSON
- TCP Raw
- HTTP/S
- Raw HTTP/S
- Kinesis Firehose
- SNMP Trap
- Metrics

Data from these Sources is normally sent to a set of LogStream Workers through a loadbalancer. Some Sources, such as Splunk forwarders, have native loadbalancing capabilities, so you should point these directly at LogStream.

PULL Sources

Supported Sources that Cribl LogStream fetches data from:

- Kafka
- Kinesis Streams
- Azure Event Hubs
- SQS
- S3
- Office 365 Services
- Office 365 Activity
- Prometheus

Internal Sources

Sources that are internal to Cribl LogStream:

- Datagens
- Cribl Internal

Configuring and Managing Sources

For each Source *type*, you can create multiple definitions, depending on your requirements.

To configure Sources, select **Data > Sources**, select the desired type from the tiles or the left menu, and then click **+ Add New**.

Preconfigured Sources

To accelerate your setup, LogStream ships with several common Sources configured for typical listening ports, but not switched on. Open, clone (if desired), modify, and enable any of these preconfigured Sources to get started quickly:

- Syslog TCP Port 9514, UDP Port 9514
- Splunk TCP Port 9997
- Splunk HEC Port 8088
- TCP JSON Port 10070
- TCP Port 10060

- HTTP Port 10080
- Elasticsearch API Port 9200
- SNMP Trap Port 9162
- Cribl Internal > CriblLogs Internal
- Cribl Internal > CriblMetrics Internal

Backpressure Behavior

On the Destination side, you can configure how each LogStream output will respond to a **backpressure** situation – a situation where its in-memory queue is overwhelmed with data.

All Destinations default to **Block** mode, in which they will refuse to accept new data until the downstream receiver is ready. Here, LogStream will backpropagate block signals through the Source, all the way back to the sender (if it supports backpressure, too).

All Destinations also support **Drop** mode, which will simply discard new events until the receiver is ready.

Several Destinations also support a **Persistent Queue** option to minimize data loss. Here, the Destination will write data to disk until the receiver is ready. Then it will drain the disk-buffered data in FIFO (first in, first out) order. See Persistent Queues for details about all three modes, and about **Persistent Queue** support.

Other BackPressure Options

The S3 Source provides a configurable **Advanced Settings > Socket timeout** option, to prevent data loss (partial downloading of logs) during backpressure delays.

Diagnosing Backpressure Errors

When backpressure affects HTTP Sources (Splunk HEC, HTTP/S, Raw HTTP/S, and Kinesis Firehose), LogStream internal logs will show a 503 error code.

Splunk TCP

Cribl LogStream supports receiving Splunk data from Universal or Heavy Forwarders.

Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive Splunk TCP Data

Select **Data > Sources**, then select **Splunk > Splunk TCP** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Splunk TCP > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a Splunk TCP Source preconfigured to listen on Port 9997. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Splunk Source definition.

Address: Enter hostname/IP to listen for Splunk data. E.g., localhost or 0.0.0.0.

Port: Enter port number.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection. Defaults to .* (i.e., all IPs).

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

* CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect.

Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol V1 or V2.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

• __inputId

Configuring a Splunk Forwarder

To configure a Splunk forwarder (UF, HF) use the following outputs.conf stanzas:

```
.../outputs.conf

[tcpout]
disabled = false
defaultGroup = cribl, <optional_clone_target_group>,

[tcpout:cribl]
server = [<cribl_ip>|<cribl_host>]:<port>, [<cribl_ip>|<cribl_host>]:<port>, ...
sendCookedData=true
# As of Splunk 6.5, using forceTimebasedAutoLB is no longer recommended. Ensure this is left at
# forceTimebasedAutoLB = false
```

Splunk HEC

Cribl LogStream supports receiving data over HTTP/S using the Splunk HEC (HTTP Event Collector).

i Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive Data over Splunk HEC

Select Data > Sources, then select Splunk > HEC from the Data Sources page's tiles or left menu. Click Add New to open the HEC > New Source modal, which provides the fields outlined below.

LogStream ships with a Splunk HEC Source preconfigured to listen on Port 8088. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Splunk HEC Source definition.

Address: Enter the hostname/IP on which to listen for HTTP(S) data. (E.g., localhost or 0.0.0.0.)

Port: Enter the port number.

Splunk HEC endpoint: Absolute path on which to listen for the Splunk HTTP Event Collector API requests. This input supports the /event and /raw endpoints. Defaults to /services/collector.

Allowed Indexes: List the values allowed in the HEC event index field. Allows wildcards. Leave blank to skip validation.

Splunk HEC acks: Whether to enable Splunk HEC acknowledgments. Defaults to No . Some sources may require HEC acks to be enabled and, as a result, may keep TCP connections open while waiting for an ack. This behavior can exhaust available file descriptors. Cribl does not maintain a comprehensive list of such sources. Refer to your source's documentation for more information.

Auth Tokens

Token: Shared secret to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access **will be permitted**.

Description: Optional description for this token.

Fields: Fields (metadata) to add to events referencing this token. Each field is a Name/Value pair.

These fields may be overridden by fields added at the request level.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

Common name: Regex matching peer certificate subject common names allowed to connect.
 Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol V1 or V2.

Processing Settings

Event Breakers

This section defines event breaking rulesets that will be applied, in order, on the /raw endpoint.

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

These fields may be overridden by fields added at the token or request level.

Pre-Processing Pipeline

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __hecToken

Format and Endpoint Example

- Configure Cribl LogStream to listen on port 10080 with an authToken of myToken42.
- Send a payload to your Cribl LogStream receiver.

Note: Token specification can be either Splunk <token> or <token> .

Splunk HEC Event Endpoint

```
curl -k http://<myCriblHost>:10080/services/collector/event -H 'Authorization: myToken42' -d '{'
curl -k http://<myCriblHost>:10080/services/collector -H 'Authorization: myToken42' -d '{"event'

# Multiple Events
curl -k http://<myCriblHost>:10080/services/collector -H 'Authorization: myToken42' -d '{"event'
```

Syslog

Cribl LogStream supports receiving of data over syslog.

Type: Push | TLS Support: YES | Event Breaker Support: No
This Syslog Source supports RFC 3164 and RFC 5424.

Configuring Cribl LogStream to Receive Data over Syslog

Select **Data > Sources**, then select **Syslog** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Syslog > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a Syslog Source preconfigured to listen for both UDP and TCP traffic on Port 9514. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Syslog Source definition.

Address: Enter the hostname/IP on which to listen for data., E.g. localhost or 0.0.0.0.

UDP port: Enter the UDP port number to listen on. Not required if listening on TCP.

TCP port: Enter the TCP port number to listen on. Not required if listening on UDP.

Fields to keep: List of fields from source data to retain and pass through. Supports wildcards. Defaults to * wildcard, meaning keep all fields. Fields **not** specified here (by wildcard or specific name) will be removed from the event.

TLS Settings (TCP Only)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

IP whitelist regex: Regex matching IP addresses that are allowed to send data. Defaults to .* (i.e., all IPs).

Max buffer size (events): Maximum number of events to buffer when downstream is blocking. The buffer is only in memory. (This setting is applicable only to UDP syslog.)

Default timezone: Timezone to assign to timestamps that omit timezone info. Accept the default Local value, or use the drop-down list to select a specific timezone by city name or GMT/UTC offset.

Single msg per UDP: Whether to treat UDP packet data received as a full syslog message. Defaults to No. (I.e., newlines in the packet will be treated as event delimiters.)

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but are accessible and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __srcIpPort
- __syslogFail: true for data that fails RFC 3164/5424 validation as syslog format.

Elasticsearch API

Cribl LogStream supports receiving data over HTTP/S using the Elasticsearch Bulk API. (See the Configuring Filebeat example below.)

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring LogStream to Receive Data over HTTP(S), Using the Elasticsearch Bulk API Protocol

Select Data > Sources, then select Elasticsearch API from the Data Sources page's tiles or left menu. Click Add New to open the Elasticsearch API > New Source modal, which provides the fields outlined below.

☐ LogStream ships with an Elasticsearch API Source preconfigured to listen on Port 9200. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Elasticsearch Source definition.

Address: Enter the hostname/IP on which to listen for Elasticsearch data. (E.g., localhost or 0.0.0.0.)

Port: Enter the port number.

Auth tokens: Shared secrets to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

Elasticsearch API endpoint (for Bulk API): Absolute path on which to listen for Elasticsearch API requests. Currently, the only supported option is the default /elastic, which LogStream expands as /elastic/_bulk. Other entries are faked as success. Use an empty string to disable.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Field Normalization

The Elasticsearch API input normalizes the following fields:

- atimestamp becomes _time atmillisecond resolution.
- host is set to host.name.
- Original object host is stored in __host .

The Elasticsearch Destination does the reverse, and it also recognizes the presence of __host .

Internal Settings

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __id
- _type
- __index
- __host

Configuring Filebeat

To set up Filebeat to send data to LogStream, use its Elasticsearch output. If an Auth Token is configured here, add it in Filebeat configuration under output.elasticsearch.headers, as in this example:

```
...filebeat.yml

output.elasticsearch:
    # Array of hosts to connect to.
    hosts: ["http://<LOGSTREAM_HOST>:9200/elastic"]
```

output.elasticsearch.headers:
 Authorization: "myToken42"

TCP JSON

Cribl LogStream supports receiving of data over TCP in JSON format (see protocol below).

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring Cribl LogStream to Receive TCP JSON Data

Select **Data > Sources**, then select **TCP JSON** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **TCP JSON > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a TCP JSON Source preconfigured to listen on Port 10070. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this TCP JSON Source definition.

Address: Enter hostname/IP to listen for TCP JSON data. E.g., localhost or 0.0.0.0.

Port: Enter the port number to listen on.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection. Defaults to .* (i.e., all IPs).

Shared secret (authToken): Shared secret to be provided by any client (in authToken header field). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

Common name: Regex matching peer certificate subject common names allowed to connect.
 Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

• Enable proxy protocol: Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

inputId

Format

LogStream expects TCP JSON events in newline-delimited JSON format:

1. A header line. Can be empty – e.g., {} . If authToken is enabled (see above) it should be included here as a field called authToken . When authToken is not set, the header line is optional. In this case, the first line will be treated as an event if does not look like a header record.

In addition, if events need to contain common fields, they can be included here under fields . In the example below, region and AZ will be automatically added to all events.

2. A JSON event/record per line.

Sample TCP JSON Events

```
{"authToken":"myToken42", "fields": {"region": "us-east-1", "AZ":"az1"}}

{"_raw":"this is a sample event ", "host":"myHost", "source":"mySource", "fieldA":"valueA", "fieldst":"myOtherHost", "source":"myOtherSource", "_raw": "{\"message\":\"Something informative | myOtherSource", "_raw": "fields": "source": "myOtherSource", "_raw": "fields": "source": "sour
```

TCP JSON Field Mapping to Splunk

If a TCP JSON Source is routed to a Splunk destination, fields within the JSON payload are mapped to Splunk fields. Fields that do not have corresponding (native) Splunk fields become index-time fields. For example, let's assume we have a TCP JSON event as below:

```
{"_time":1541280341, "host":"myHost", "source":"mySource", "_raw":"this is a sample
event ", "fieldA":"valueA"}
```

Here, _time , host , and source become their corresponding fields in Splunk. The value of _raw becomes the actual body of the event, and fieldA becomes an index-time field (fieldA::`valueA``).

Example

- 1. Configure Cribl LogStream to listen on port 10001 for TCP JSON. Set authToken to myToken42.
- 2. Create a file called test.json with the payload above.
- 3. Send it over to your Cribl LogStream host: cat test.json | nc <myCriblHost> 10001

TCP (RAW)

Cribl LogStream supports receiving of data over TCP. (See examples and header options below.)

i Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive TCP Data

Select **Data > Sources**, then select **TCP** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **TCP > New Source** modal, which provides the fields outlined below.

□ LogStream ships with a TCP Source preconfigured to listen on Port 10060. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this TCP Source definition.

Address: Enter hostname/IP to listen for raw TCP data. E.g., localhost or 0.0.0.0.

Port: Enter port number.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection. Defaults to .* (i.e., all IPs).

Enable Header: Toggle to Yes to indicate that client will pass a header record with every new connection. The header can contain an authToken, and an object with a list of fields and values to add to every event. These fields can be used to simplify Event Breaker selection, routing, etc. Header format:

\[
\[
\begin{align*}
\]
\[
\begin{align*}
\beg

```
{ "authToken" : "myToken", "fields": { "field1": "value1", "field2": "value2" }}.
```

• **Shared secret (authToken)**: Shared secret to be provided by any client (in authToken header field). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Custom Command

In this section, you can pass the data from this input to an external command for processing before the data continues downstream.

Enabled: Defaults to No. When toggled to Yes:

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + Add Argument to add each argument for the command. You can drag arguments vertically to resequence them.

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

• Name: Field name.

• Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol V1 or V2.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions can use them to make processing decisions.

Fields accessible for this Source:

- __inputId
- __srcIpPort
- __channel

TCP Source Example

Every new TCP connection may contain an **optional** header line, with an authToken and a list of fields and values to add to every event.

```
SamplerawTCPtest
{"authToken":"myToken42", "fields": {"region": "us-east-1", "AZ":"az1"}}
this is event number 1
this is event number 2
```

Enabling the Example

- Configure LogStream to listen on port 7777 for raw TCP. Set authToken to myToken42.
- 2. Create a file called test.raw, with the payload above.
- 3. Send it over to your Cribl LogStream host, using this command: cat
 test.raw | nc <myCriblHost> 7777

HTTP/S (Bulk API)

Cribl LogStream supports receiving data over HTTP/S using the Cribl Bulk API, Splunk HEC, or Elastic Bulk API.

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data over HTTP(S)

Select Data > Sources, then select HTTP from the Data Sources page's tiles or left menu. Click Add New to open the HTTP > New Source modal, which provides the fields outlined below.

☐ LogStream ships with an HTTP Source preconfigured to listen on Port 10080, and on several default endpoints. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this HTTP(S) Source definition.

Address: Enter the hostname/IP on which to listen for HTTP(S) data. (E.g., localhost or 0.0.0.0.)

Port: Enter the port number.

Auth tokens: Shared secrets to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access **will be permitted**.

Cribl HTTP event API: Absolute path on which to listen for Cribl HTTP API requests. Currently, the only supported option is the default <code>/cribl</code>, which LogStream expands as <code>/cribl/_bulk</code>. Use an empty string to disable. Maximum payload size is 2MB.

Elastic API endpoint (for Bulk API): Absolute path on which to listen for Elasticsearch API requests. Currently, the only supported option is the default /elastic , which LogStream expands as /elastic/_bulk . Other entries are faked as success. Use an empty string to disable.

i Cribl generally recommends that you use the dedicated Elasticsearch API Source instead of this endpoint. The Elastic API implementation here is provided for backward compatibility, and for users who want to ingest multiple inputs on one HTTP/S port.

Splunk HEC endpoint: Absolute path on which to listen for Splunk HTTP Event Collector (HEC) API requests. Use an empty string to disable. Default entry is /services/collector.

i This Splunk HEC implementation is an event (i.e., not raw) endpoint. For details, see Splunk's documentation. To send data to it from a HEC client, use either /services/collector or /services/collector/event . (See the examples below.)

Cribl generally recommends that you use the dedicated Splunk HEC Source instead of this endpoint. The Splunk HEC implementation here is provided for backward compatibility, and for users who want to ingest multiple inputs on one HTTP/S port.

Splunk HEC acks: Whether to enable Splunk HEC acknowledgements. Defaults to No.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

Common name: Regex matching peer certificate subject common names allowed to connect.
 Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

```
• __inputId
```

- __id (Elastic In)
- __type (Elastic In)
- __index (Elastic In)
- host (Elastic In)

Format and Endpoint

LogStream expects HTTP(S) events to be formatted as one JSON record per event. Here are two event records:

```
Sample Event Format
```

Note 1: Events can be sent as separate POSTs, but Cribl **highly** recommends combining multiple events in newline-delimited groups, and POSTing them together.

Note 2: If an HTTP(S) source is routed to a Splunk destination, fields within the JSON payload are mapped to Splunk fields. Fields that do not have corresponding (native) Splunk fields become indextime fields. For example, let's assume we have a HTTP(S) event like this:

```
{"_time":1541280341, "host":"myHost", "source":"mySource", "_raw":"this is a sample
event ", "fieldA":"valueA"}
```

Here, _time, host and source become their corresponding fields in Splunk. The value of _raw becomes the actual body of the event, and fieldA becomes an index-time field. (fieldA:: valueA).

Examples

For the following examples:

- 1. Configure Cribl to listen on port 10080 for HTTP (default). Set authToken to myToken42.
- 2. Send a payload to your Cribl LogStream receiver.

Cribl Endpoint - Single Event

Cribl Single Event Example:

```
curl -k http://<myCriblHost>:10080/cribl/_bulk -H 'Authorization: myToken42' -d '{"_raw":"this :
```

Cribl Endpoint - Multiple Events

```
Cribl Endpoint - Multiple Events
```

```
curl -k http://<myCriblHost>:10080/cribl/_bulk -H 'Authorization: myToken42' -d $'{"_raw":"this
```

Splunk HEC Event Endpoint

```
Splunk HEC Event Endpoint
```

```
curl -k http://<myCriblHost>:10080/services/collector/event -H 'Authorization: myToken42' -d '{'
curl -k http://<myCriblHost>:10080/services/collector -H 'Authorization: myToken42' -d '{"event"
```

 ${f i}$ For Splunk HEC, the token specification can be either Splunk <token> or <token> .

Raw HTTP/S

Cribl LogStream supports receiving raw HTTP data. The Raw HTTP Source listens on a specific port, captures every HTTP request to that port, and creates a corresponding event that it pushes to its configured Event Breakers.

i Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive Raw HTTP Data

Select **Data > Sources**, then select **Raw HTTP** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Raw HTTP > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Raw HTTP Source definition.

Address: Enter the address to bind on. Defaults to 0.0.0.0 (all addresses).

Port: Enter the port number to listen on.

Auth tokens: Shared secrets to be provided by any client. Click **Generate** to create a new secret. If empty, permits open access.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

- Name: Field name.
- Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Allowed URI paths: List of URI paths accepted by this input. Supports wildcards, e.g., /api/v*/hook . Defaults to *, which allows all paths.

Allowed HTTP methods: List of HTTP methods accepted by this input. Supports wildcards, e.g., P*, GET. Defaults to *, which allows all methods.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions can use them to make processing decisions.

Fields accessible for this Source:

- __inputId
- __srcIpPort
- __channel

Kafka

Cribl LogStream supports receiving data records from a Kafka cluster.

i Type: Pull | TLS Support: YES | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Kafka Topics

Select **Data > Sources**, then select **Kafka** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Kafka > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Brokers: List of Kafka brokers to use, e.g., localhost:9092.

Topics: List of topics to subscribe to.

Group ID: The name of the consumer group to which this Cribl LogStream instance belongs.

From beginning: Whether to start reading from the earliest available data. Relevant only during initial subscription. Defaults to Yes.

TLS Settings (Client Side)

Enabled: defaults to No . When toggled to Yes :

Autofill?: This setting is experimental.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Authentication

This section governs SASL (Simple Authentication and Security Layer) authentication.

Enabled: Defaults to No. When toggled to Yes:

SASL mechanism: Use this drop-down to select the SASL authentication mechanism to use.

Username: Enter the username for your account.

Password: Enter the account's password.

Schema Registry

This section governs Kafka Schema Registry Authentication for AVRO-encoded data with a schema stored in the Confluent Schema Registry.

Enabled: defaults to No. When toggled to Yes:

Schema registry URL: URL for access to the Confluent Schema Registry. (E.g., http://<hostname>:8081.)

TLS enabled: defaults to No. When toggled to Yes, displays the following TLS settings for the Schema Registry:

i These have the same format as the TLS Settings (Client Side) above.

TLS Settings (Schema Registry)

Validate server certs: Reject certificates that are not authorized by a CA specified in the **CA Certificate Path** field. Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __topicIn (indicates the Kafka topic that the event came from; see
 - __topicOut in our Kafka Destination documentation)
- __schemaId (when using Schema Registry)

Kinesis

Cribl LogStream supports receiving data records from Amazon Kinesis Streams.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Kinesis Streams

Select **Data > Sources**, then select **Kinesis** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Kinesis > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Kinesis Stream Source definition.

Stream name: Kinesis stream name (not ARN) to read data from.

Shard iterator start: Location at which to start reading a shard for the first time. Defaults to Earliest Record.

Record data format: Format of data inside the Kinesis Stream records. Gzip compression is automatically detected. Options include:

- Cribl (the default): Use this option if LogStream wrote data to Kinesis in this format. This is a type of NDJSON.
- Newline JSON: Use if the records contain newline-delimited JSON
 (NDJSON) events e.g., Kubernetes logs ingested through Kinesis. This is a
 good choice if you don't know the records' format.
- CloudWatch Logs: Use if you've configured CloudWatch to send logs to Kinesis.
- Event per line: NDJSON can use this format when it fails to parse lines as valid JSON.

Region: Region where the Kinesis stream is located. Required.

Authentication

Authentication Method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for Kinesis Streams: Whether to use Assume Role credentials to access Kinesis Streams. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

• Name: Field name.

• Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Shard selection expression: A JavaScript expression to be called with each shardId for the stream. The shard will be processed if the expression evaluates to a truthy value. Defaults to true.

Service Period: Time interval (in minutes) between consecutive service calls. Defaults to 1 minute.

Endpoint: Kinesis stream service endpoint. If empty, the endpoint will be automatically constructed from the region.

Signature version: Signature version to use for signing Kinesis Stream requests. Defaults to v4.

Verify KPL checksums: Enable this setting to verify Kinesis Producer Library (KPL) event checksums.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

• __inputId

Kinesis Firehose

Cribl LogStream supports receiving data from Amazon Kinesis Data Firehose delivery streams via Kinesis' **HTTP endpoint** destination option.

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring LogStream to Receive Data over HTTP(S) from Amazon Kinesis Firehose

Select **Data > Sources**, then select **Amazon > Firehose** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Firehose > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Address: Address to bind on. Defaults to 0.0.0.0 (all addresses).

Port: Enter the port number to listen on.

Auth tokens: Shared secrets to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

• Name: Field name.

• Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions

can use them to make processing decisions.

Fields accessible for this Source:

- __inputId
- __firehoseArn
- __firehoseReqId
- __firehoseEndpoint

Limitations/Troubleshooting

If you set the optional IntervalInSeconds and/or SizeInMBs parameters in the Kinesis Firehose BufferingHints API, beware of selecting extreme values (toward the ends of the API's supported ranges). These can send more bytes than LogStream can buffer, causing LogStream to send HTTP 500 error responses to Kinesis Firehose.

Azure Event Hubs

Cribl LogStream supports receiving data records from Azure Event Hubs.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Azure Event Hubs

Select **Data > Sources**, then select **Azure Event Hubs** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Azure Event Hubs > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this source definition.

Brokers: List of Event Hubs Kafka brokers to connect to, e.g., yourdomain.servicebus.windows.net:9093. Get the hostname from the host portion of the primary or secondary connection string in Shared Access Policies.

Event Hub name: The name of the Event Hub (a.k.a. Kafka Topic) to subscribe to.

Group ID: Specifies the name of the consumer group to which this Cribl LogStream instance belongs. Should always be \$Default for Event Hubs.

From beginning: Whether to start reading from the earliest available data. Relevant only during initial subscription. Defaults to Yes.

TLS Settings (Client Side)

Enabled: Defaults to Yes.

Validate server certs: Whether to reject connections to servers without signed certificates. Defaults to No – and for Event Hubs, must always be disabled.

Authentication

Enabled: Defaults to No. When toggled to Yes:

- **SASL** mechanism: SASL (Simple Authentication and Security Layer) authentication mechanism to use. Currently, PLAIN is the only mechanism supported for Event Hubs Kafka brokers.
- **Username**: The username for authentication. For Event Hubs, this should always be \$ConnectionString.
- **Password**: Connection-string primary key or connection-string secondary key from the Event Hub workspace.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Metrics

Cribl LogStream supports receiving metrics in these wire formats/protocols: StatsD, StatsD Extended, and Graphite. Automatic protocol detection will happen on the first line received over a TCP connection or a UDP packet. Lines not matching the detected protocol will be dropped.

i Type: Push | TLS Support: No | Event Breaker Support: No

Configuring Cribl LogStream to Receive Metrics

Select **Data > Sources**, then select **Metrics** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Metrics > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Address: Enter the hostname/IP to listen to. Defaults to 0.0.0.0.

UDP port: Enter the UDP port number to listen on. Not required if listening on TCP.

TCP port: Enter the TCP port number to listen on. Not required if listening on UDP.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

IP whitelist regex: Regex matching IP addresses that are allowed to send data. Defaults to .* (i.e., all IPs.)

Max buffer size (events): Maximum number of events to buffer when downstream is blocking. Defaults to 1000.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __srcIpPort
- __metricsInType

Metric Event Schema and Destination Support

Metric data is read into the following event schema:

```
Text

_metric - the metric name
_metric_type - the type of the metric (gauge, counter, timer)
_value - the value of the metric
_time - metric_time or Date.now()/1000
dim1 - value of dimension1
dim3 - value of dimension2
....
```

LogStream places sufficient information into a field called __criblMetric to enable these events to be properly serialized out to any metric outputs

(independent of the input type).

The following Destinations natively support the __criblMetric field:

- Splunk
- Splunk HEC
- InfluxDB
- Statsd
- Statsd Extended
- Graphite

Prometheus

Cribl LogStream supports receiving data from Prometheus.

i Type: Pull | TLS Support: No | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Prometheus

Select **Data > Sources**, then select **Prometheus** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Prometheus > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Extra dimensions: Dimensions to include in events. By default, host and source are included.

Discovery type: Target discovery mechanism. Use Static (the default) to manually enter a list of targets. Select DNS or AWS EC2 options enable dynamic discovery of endpoints to scrape. Your selection determines which fields are displayed lower in this section:

- Targets: Displayed for Discovery type: Static . List of Prometheus targets
 to pull metrics from, values can be in URL or host[:port] format. For
 example: http://localhost:9090/metrics, localhost:9090, or localhost. In the
 cases where just host[:port] are specified, the endpoint will resolve to
 'http://host[:port]/metrics'.
- **DNS names**: Displayed for **Discovery type**: DNS . Enter a list of DNS names to resolve.
- Record type: Displayed for Discovery type: DNS. Select the DNS record type to resolve. Defaults to SRV (Service record). Other options are A or AAA record.

• **Region**: Displayed for **Discovery type**: AWS EC2 . Select the AWS region in which to discover the EC2 instances with metrics endpoints to scrape.

Poll interval: How often (in minutes) to scrape targets for metrics. Defaults to 15. This value must be an integer that divides evenly into 60 minutes.

Log level: Set the verbosity level to one of debug, info (the default), warn, or error.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Keep alive time (seconds): How often workers should check in with the scheduler to keep job subscription alive. Defaults to 60 seconds.

Worker timeout (periods): How many **Keep alive time** periods before an inactive worker's job subscription will be revoked. Defaults to 3 periods.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

__source

- __isBroken
- __inputId
- __final
- __criblMetrics
- __channel
- __cloneCount

SQS

Cribl LogStream supports receiving events from Amazon Simple Queuing Service.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support:

Configuring Cribl LogStream to Receive Data from Amazon SQS

Select **Data > Sources**, then select **SQS** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **SQS > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this SQS Source definition.

Queue: The name, URL, or ARN of the SQS queue to read events from.

This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks. To specify a non-AWS URL, use the format: '{url}/<queueName>' . (E.g., ':port/<myQueueName>' .)

Create queue: Create queue if it does not exist.

Region: AWS Region where the SQS queue is located. Required, unless the **Queue** entry is a URL or ARN that includes a Region.

Authentication

Authentication Method: Select an AWS authentication method.

Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.

• Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

 API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.

 Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for SQS: Whether to use Assume Role credentials to access SQS. Defaults to No.

AWS account ID: SQS queue owner's AWS account ID. Leave empty if SQS queue is in same AWS account.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the external ID to use when assuming role.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Endpoint: SQS service endpoint. If empty, the endpoint will be automatically constructed from the AWS Region.

Signature version: Signature version to use for signing SQS requests. Defaults to v4.

Max messages: The maximum number of messages that SQS should return in a poll request. Amazon SQS never returns more messages than this value. (However, fewer messages might be returned.) Acceptable values: 1 to 10. Defaults to 10.

Visibility timeout seconds: The duration (in seconds) that the received messages are hidden from subsequent retrieve requests, after they're retrieved by a ReceiveMessage request. Defaults to 600.

Num receivers: The number of receiver processes to run. The higher the number, the better the throughput, at the expense of CPU overhead. Defaults to 3.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __sqsSysAttrs

The _sqsSysAttrs field can take on the following properties, which are reported to LogStream from SQS:

- __sqsSysAttrs.ApproximateFirstReceiveTimestamp: Returns the time (epoch time in milliseconds) the message was first received from the queue.
- __sqsSysAttrs.ApproximateReceiveCount: Returns the number of times a message has been received from the queue without being deleted.
- __sqsSysAttrs.SenderId: For an IAM user, returns the IAM user ID (e.g.: ABCDEFGHI1JKLMNOPQ23R). For an IAM role, returns the IAM role ID (e.g.: ABCDE1F2GH3I4JK5LMNOP:i-a123b456).
- __sqsSysAttrs.SentTimestamp: Returns the time (epoch time in milliseconds) the message was sent to the queue.

- __sqsSysAttrs.MessageDeduplicationId: Returns the value provided by the producer that calls the SendMessage action.
- __sqsSysAttrs.MessageGroupId: Returns the value provided by the producer that calls the SendMessage action messages with the same MessageGroupId are returned in sequence.
- __sqsSysAttrs.SequenceNumber: Returns the sequence-number value provided by Amazon SQS.
- __sqsSysAttrs.AWSTraceHeader: Returns the AWS X-Ray trace header string.

For background on these message properties, see AWS' ReceiveMessage > Request Parameters documentation.

SQS Permissions

The following permissions are needed on the SQS queue:

- sqs:ReceiveMessage
- sqs:DeleteMessage
- sqs:GetQueueAttributes
- sqs:GetQueueUrl
- sqs:CreateQueue (optional, if and only if you want LogStream to create the queue)

Troubleshooting Notes

▲ VPC endpoints for SQS might need to be set up in your account. Check with your administrator for details. Cribl LogStream supports receiving data from Amazon S3 buckets, using event notifications through SQS.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support: YES

S3 Setup Strategy

i The source S3 bucket must be configured to send s3:ObjectCreated:* events to an SQS queue, either directly (easiest) or via SNS (Amazon Simple Notification Service). See the event notification configuration guidelines below.

SQS messages will be deleted after they're read, unless an error occurs, in which case LogStream will retry. This means that although LogStream will ignore files not matching the **Filename Filter**, their SQS events/notifications will still be read, and then deleted from the queue (along with those from files that match).

These ignored files will no longer be available to other S3 Sources targeting the same SQS queue. If you still need to process these files, we suggest one of these alternatives:

- Using a different, dedicated SQS queue. (Preferred and recommended.)
- Applying a broad filter on a single Source, and then using preprocessing Pipelines an/or Route filters for further processing.

Configuring Cribl LogStream to Receive Data from Amazon S3

Select **Data > Sources**, then select **S3** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **S3 > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this S3 Source definition.

Queue: The name, URL, or ARN of the SQS queue to read events from. When specifying a non-AWS URL, you must use the format: {url}/<queueName>.(E.g., https://host:port/<queueName>.) This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks.

Filename filter: Regex matching file names to download and process. Defaults to .*, to match all characters.

Region: AWS Region where the S3 bucket and SQS queue are located. Required, unless the **Queue** entry is a URL or ARN that includes a Region.

Authentication

Authentication method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for S3: Whether to use Assume Role credentials to access S3. Defaults to Yes.

Enable for SQS: Whether to use Assume Role credentials when accessing SQS (Amazon Simple Queue Service). Defaults to No.

AWS account ID: SQS queue owner's AWS account ID. Leave empty if the SQS queue is in the same AWS account.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + **Add Argument** to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

This section defines event breaking rulesets that will be applied, in order.

Event Breaker Rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker Buffer Timeout: The amount of time (in milliseconds) that the Event Breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Endpoint: S3 service endpoint. If empty, defaults to AWS's region-specific endpoint. Otherwise, used to point to an S3-compatible endpoint.

Signature version: Signature version to use for signing SQS requests. Defaults to v4.

Num receivers: The number of receiver processes to run,. The higher the number, the better the throughput, at the expense of CPU overhead. Defaults to 1.

Max messages: The maximum number of messages that SQS should return in a poll request. Amazon SQS never returns more messages than this value. (However, fewer messages might be returned.) Acceptable values: 1 to 10. Defaults to 1.

Visibility timeout seconds: The duration (in seconds) that the received messages are hidden from subsequent retrieve requests, after being retrieved by a ReceiveMessage request. Defaults to 600.

i LogStream will automatically extend this timeout until the initial request's files have been processed – notably, in the case of large files that require additional processing time.

Socket timeout: Socket inactivity timeout (in seconds). Increase this value if retrievals time out during backpressure. Defaults to 300 seconds.

Skip file on error: Toggle to **Yes** to skip files that trigger a processing error. (E.g., corrupted files.) Defaults to **No**, which enables retries after a processing error.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __source

How to Configure S3 to Send Event Notifications to SQS

- i For step-by-step instructions, see AWS' Walkthrough: Configure a Bucket for Notifications (SNS Topic and SQS Queue).
- 1. Create a Standard SQS Queue. Note its ARN.
- 2. Replace its access policy with one similar to the examples below. To do so, select the queue; and then, in the **Permissions** tab, click: **Edit Policy Document (Advanced)**. (These examples differ only at line 9, showing public access to the SQS queue versus S3-only access to the queue.)
- 3. In the Amazon S3 console, add a notification configuration to publish events of the s3:0bjectCreated:* type to the SQS queue.

Permissive SQS access policy Restrictive SQS access policy "Version": "example-2020-04-20", "Id": "example-ID", "Statement": ["Sid": "<SID name>", "Effect": "Allow", "Principal": { "AWS":"*" }, "Action": ["SQS:SendMessage" "Resource": "example-SQS-queue-ARN", "Condition": { "ArnLike": { "aws:SourceArn": "arn:aws:s3:*:*:example-bucket-name" } }]

S3 and SQS Permissions

The following permissions are required on the S3 bucket:

- s3:GetObject
- s3:ListBucket

The following permissions are required on the SQS queue:

- sqs:ReceiveMessage
- sqs:DeleteMessage
- sqs:ChangeMessageVisibility
- sqs:GetQueueAttributes
- sqs:GetQueueUrl

Best Practices

- When LogStream instances are deployed on AWS, use IAM Roles whenever possible.
 - Not only is this safer, but it also makes the configuration simpler to maintain.
- Although optional, we highly recommend that you use a Filename Filter.
 - This will ensure that LogStream ingests only files of interest.
 - Ingesting only what's strictly needed improves latency, processing power, and data quality.
- If higher throughput is needed, increase Advanced Settings > Number of Receivers and/or Max messages. However, do note:
 - These are set at 1 by default. Which means, each Worker Process, in each LogStream Worker Node, will run 1 receiver consuming 1 message (i.e. S3 file) at a time.
 - Total S3 objects processed at a time per Worker Node = Worker Processes x Number of Receivers x Max Messages
 - Increased throughput implies additional CPU utilization.
- When ingesting large files, tune up the **Visibility Timeout**, or consider using smaller objects.

• The default value of 600s works well in most cases, and while you certainly can increase it, we suggest that you also consider using smaller S3 objects.

Troubleshooting Notes

- VPC endpoints for SQS and for S3 might need to be set up in your account. Check with your administrator for details.
- If you're having connectivity issues, but no problems with the CLI, see if the AWS CLI proxy is in use. Check with your administrator for details.

Office 365 Services

Cribl LogStream supports receiving data from the Office 365 Service Communications API. This facilitates analyzing the status and history of service incidents on multiple Microsoft cloud services, along with associated incident and Message Center communications.

Type: Pull | TLS Support: YES | Event Breaker Support: YES

TLS is enabled via the HTTPS protocol on this Source's underlying REST API.

Configuring Cribl LogStream to Receive Data from Office 365 Services

Select **Data > Sources**, then select **Office 365 > Services** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Services > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Office 365 Services definition.

Tenant ID: Enter the Office 365 Azure tenant ID.

App ID: Enter the Office 365 Azure application ID.

Client secret: Enter the Office 365 Azure client secret.

Content Types

Here, you can configure polling separately for the following types of data from the Office 365 Service Communications API:

- **Current Status**: Get a real-time view of current and ongoing service incidents.
- Messages: Find incident and Message Center communications.
- Historical Status: Get a historical view of service incidents.

As of this revision, this Microsoft API provides data for Office 365, Yammer, Dynamics CRM, and Microsoft Intune cloud services. For each of these content types, this section provides the following controls:

Enabled: Toggle this to Yes for each service that you want to poll.

Interval: Optionally, override the default polling interval. See About Polling Intervals below.

Log level: Set the verbosity level to one of debug, info (the default), warn, or error.

About Polling Intervals

To poll the Office 365 Service Communications API, LogStream uses the **Interval** field's value to establish the search date range and the cron schedule (e.g.: */\${interval} * * * *).

Therefore, intervals set in minutes – those for **Current Status** and **Historical Status** – must divide evenly into 60 minutes to create a predictable schedule. Dividing 60 by intervals like 1, 2, 3, 4, 5, 6, 10, 12, 15, 20, or 60 itself yields an integer, so you can enter any of these values.

LogStream will reject intervals like 23, 42, or 45, or 75 – which would yield non-integer results, meaning unpredictable schedules.

The **Historical Status** service polls only once per day. So here, the **Interval** field's value simply establishes the hour of the day at which to poll. (In distributed deployments, this time is set based on the Master Node's system time. In single-instance deployments, it is set based on the API server's time zone.)

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Keep Alive Time (seconds): How often Workers should check in with the scheduler to keep their job subscription alive. Defaults to 60.

Worker timeout (periods): The number of Keep Alive Time periods before an inactive Worker will have its job subscription revoked. Defaults to 3.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __final
- __inputId
- __isBroken
- __source

Office 365 Activity

Cribl LogStream supports receiving data from the Office 365 Management Activity API. This facilitates analyzing actions and events on Azure Active Directory, Exchange, and SharePoint, along with global auditing and Data Loss Prevention data.

Type: Pull | TLS Support: YES | Event Breaker Support: YES

TLS is enabled via the HTTPS protocol on this Source's underlying REST API.

Configuring Cribl LogStream to Receive Data from Office 365 Activity

Select **Data > Sources**, then select **Office 365 > Activity** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Activity > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Office 365 Activity definition.

Tenant ID: Enter the Office 365 Azure tenant ID.

App ID: Enter the Office 365 Azure application ID.

Client secret: Enter the Office 365 Azure client secret.

Subscription Plan: Select the Office 365 subscription plan for your organization. This is typically Enterprise and GCC Government Plan.

Content Types

Here, you can configure polling independently for the following types of audit data from the Office 365 Management Activity API:

- Active Directory
- Exchange

- SharePoint
- General: All workloads not included in the above content types
- DLP.All: Data Loss Prevention events only, for all workloads

For each of these content types, this section provides the following controls:

Enabled: Toggle this to Yes for each service that you want to poll.

Interval: Optionally, override the default polling interval. See About Polling Intervals below.

Log level: Set the verbosity level to one of debug, info (the default), warn, or error.

About Polling Intervals

To poll the Office 365 Management Activity API, LogStream uses the Interval field's value to establish the search date range and the cron schedule (e.g.: */\${interval} * * * *).

Therefore, intervals set in minutes must divide evenly into 60 minutes to create a predictable schedule. Dividing 60 by intervals like 1, 2, 3, 4, 5, 6, 10, 12, 15, 20, or 60 itself yields an integer, so you can enter any of these values.

LogStream will reject intervals like 23, 42, or 45, or 75 – which would yield non-integer results, meaning unpredictable schedules.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Keep Alive Time (seconds): How often Workers should check in with the scheduler to keep their job subscription alive. Defaults to 60.

Worker timeout (periods): The number of Keep Alive Time periods before an inactive Worker will have its job subscription revoked. Defaults to 3.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __final
- __inputId
- __isBroken
- __source

SNMP Trap

Cribl LogStream supports receiving data from SNMP Traps.

i Type: Push | TLS Support: NO | Event Breaker Support: No

Configuring Cribl LogStream to Receive SNMP Traps

Select **Data > Sources**, then select **SNMP Trap** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **SNMP Trap > New Source** pane, which provides the fields outlined below.

☐ LogStream ships with an SNMP Trap Source preconfigured to listen on Port 9162. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Address: Address to bind on. Defaults to 0.0.0.0 (all addresses).

UDP Port: Port on which to receive SNMP traps. Defaults to 162.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

IP whitelist regex: Regex matching IP addresses that are allowed to send data. Defaults to .* i.e. all IPs.

Max buffer size (events): Maximum number of events to buffer when downstream is blocking. Defaults to 1000.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __snmpVersion: Acceptable values are 0, 2, 3. Versions: 0 =v1, 2 =v2c, 3 =v3.
- __srcIpPort : <hostname>|port
- __snmpRaw : Buffer containing Raw SNMP packet

Considerations for Working with SNMP Trap Data

- It's possible to work with SNMP metadata (i.e., we'll decode the packet). Options include dropping, routing, etc.
- SNMP packets can be forwarded to other SNMP destinations. However, the
 contents of the incoming packet cannot be modified i.e., we'll forward
 the packets verbatim as they came in.
- SNMP packets can be forwarded to non-SNMP destinations (e.g., Splunk, Syslog, S3, etc.).
- Non-SNMP input data cannot be sent to SNMP destinations.

Datagens

Cribl LogStream supports generating of data from datagen files. See Using Datagens for more details.

i Type: Internal | TLS Support: N/A | Event Breaker Support: No

Configuring Cribl LogStream to Generate Sample Data

Select **Data > Sources**, then select **Datagens** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Datagens > New Source** pane, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Datagens: List of datagens.

- Data generator file: Name of the datagen file.
- Events per second per Worker Node: Maximum number of events to generate per second, per worker node. Defaults to 10.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

• __inputId

Cribl Internal

The Cribl Internal Source enables you to capture and send LogStream's own internal **logs** and **metrics** through Routes and Pipelines. In distributed mode, only Worker Node internal logs can be processed through this Source. (Logs on the Master remain on the Master, since the Master Node is not part of any processing path.)

i Type: Internal | TLS Support: N/A | Event Breaker Support: No

Configuring Cribl Internal Logs/Metrics to Behave as a Data Source

Select **Data > Sources**, then select **Cribl Internal** from the **Data Sources** page's tiles or left menu.

Next, on the **CriblLogs** and/or the **CriblMetrics** row, slide the **Enabled** slider to Yes. Confirm your choice in the resulting message box.

To proceed to the configuration options listed below, click anywhere on the **CriblLogs** or the **CriblMetrics** row.



Cribl Internal Sources - click to configure

CriblLogs Settings

General Settings

Enabled: This duplicates the parent page's **Enabled** slider. Keep it at Yes to enable Cribl logs as a Source.

Input ID: Enter a unique name to identify this CriblLogs Source definition.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

CriblMetrics Settings

General Settings

Enabled: This duplicates the parent page's **Enabled** slider. Keep it at Yes to enable Cribl metrics as a Source.

Input ID: Enter a unique name to identify this CriblMetrics Source definition.

Metric name prefix: Enter an optional prefix that will be applied to metrics provided by LogStream. The prefix defaults to cribl.logstream. .

i If LogStream detects source, sourcetype, host, or index fields in metrics from external sources, it copies their values into new dimensions with added event_ prefixes (e.g., event_sourcetype). This leaves the original dimensions (and their values) intact.

Note that you can disable metric collection for any or all of these four fields at System Settings > General Settings > Limits > Disable field metrics.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Reporting Metrics Less Frequently

By default, LogStream generates internal metrics every 2 seconds. To consume metrics at longer intervals, you can use or adapt the <code>cribl-metrics_rollup</code> Pipeline that ships with LogStream. Attach it to your <code>Cribl Internal</code> Source as a <code>pre-processing Pipeline</code>. The Pipeline's <code>Rollup Metrics</code> Function has a default <code>Time Window</code> of 30 seconds, which you can adjust to a different granularity as needed.

Internal Fields

The following fields will be added to all events/metrics:

source:setto cribl.

• host: set to the hostname of the Cribl instance.

Use these fields to guide these events/metrics through Cribl Routes.

△ All Cribl internal fields are subject to change and modification. Cribl provides them to assist with analytics and diagnostics, but does not guarantee that they will remain available.

AppScope

AppScope is an open-source instrumentation utility from Cribl. It offers visibility into any Linux command or application, regardless of runtime, with no code modification. For details about configuring the AppScope CLI, loader, and library, see: https://appscope.dev/docs.

Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring LogStream to Receive AppScope Data

Select **Data > Sources**, then select **AppScope** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **AppScope > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this AppScope Source definition.

Address: Enter the hostname/IP on which to listen for AppScope data. (E.g., localhost .) Defaults to 0.0.0.0, meaning all addresses.

Port: Enter the port number to listen on.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection.

Auth tokens: Shared secrets to be provided by any client (in authToken header field). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Enable if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

Internal Settings

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

• __inputId

^ Reviewers: This is placeholder copy/pasta from the Elasticsearch Source. Which internal fields are actually supported here?

Collectors

Collectors enable you to dispatch on-demand collection tasks that fetch data from local or remote locations. As of v.2.3, LogStream supports scheduled collection jobs. These recurring jobs can make batch collection of stored data more like continual processing of streaming data.

How Do Collectors Work

You can configure a LogStream Node to retrieve data from a remote system via **Data** > **Collectors**. Data collection is a multi-step process:

First, define a Collector instance. In this step, you configure **collector-specific settings** by selecting a Collector type and pointing it at a specific target. (E.g., the target will be a directory if the type is Filesystem, or an S3 bucket/path if the type is Amazon S3.)

Next, schedule or manually run the Collector. In this step, you configure scheduled-job-specific or run-specific settings – such as the run Mode (such as Discovery or Full Run), the Filter expression to match the data against, the time range, etc.

When a Node receives this configuration, it prepares the infrastructure to execute a collection job. A collection job is typically made up of one or more tasks that: discover the data to be fetched; fetch data that match the run filter; and finally, pass the results either through the Routes or (optionally) into a specific Pipeline and Destination.

i On the Manage Collectors page, click Job Inspector to see the results of recent collection runs. Select the Show system jobs check box to also display discovery jobs and collection jobs for the Office 365 System/Activity Sources.

Scheduled Collection Jobs

You might process data from inherently non-streaming sources, such as REST endpoints, blob stores, etc. Scheduled jobs enable you to emulate a data stream by scraping data from these sources in batches, on a set interval.

You can schedule a specific job to pick up new data from the source – data that hadn't been picked up in previous invocations of this scheduled job. This essentially transforms a non-streaming data source into a streaming data source.

Collectors in Distributed Deployments

In a distributed deployment, Collectors are set up at the Worker Group level, and the tasks are executed by Worker Nodes. The Master Node oversees the task distribution, and tries to maintain a fair balance across jobs.

When Workers ask for tasks, the Master will normally try to assign the next task from a job with the least tasks in progress. This is known as "Least-In-Flight Scheduling," and provides the fairest task distribution for most cases. Default behavior can be changed via **Settings > General Settings > Job Limits > Job Dispatching**.

Collector Types

Cribl LogStream currently provides the following Collector options:

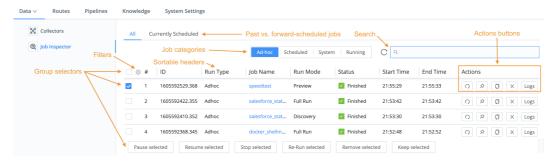
- Filesystem/NFS enables data collection from local or remote filesystem locations.
- S3 enables data collection from Amazon S3 buckets or S3-compatible stores.
- Script enables data collection via custom scripts.
- REST enables data collection via REST API calls. Provides four Discover options, to support progressively more complex (and dynamic) item enumerations.

Monitoring Collection Jobs

Select **Monitoring > Jobs** from the top menu to see a graphical display of inflight collection jobs and their tasks.

Inspecting Jobs

Select **Data > Collectors > Job Inspector** to view and manage pending, in-flight, and completed collection jobs.



Job Inspector: all the things

Here are the options available on the Job Inspector page:

- All vs. Currently Scheduled tabs: Click Currently Scheduled to see jobs foward-scheduled for future execution – including their cron schedule details, last execution, and next scheduled execution. Click All to see all jobs initiated in the past, regardless of completion status.
- Job categories (buttons): Select among Ad hoc, Scheduled, System, and Running. (At this level, Scheduled means scheduled jobs already running or finished.)
- Filters: Click the gear icon to open a drop-down with multiple options to filter the jobs shown within your selected category.
- Group selectors: Select one or more check boxes to display the Pause,
 Resume, etc., buttons shown along the bottom.
- Sortable headers: Click any column to reverse its sort direction.
- Search bar: Click to filter displayed jobs by arbitrary strings.
- Action buttons: For finished jobs, the icons (from left to right) indicate: Rerun; Keep job artifacts; Copy job artifacts; Delete job artifacts; and Display job logs in a modal. For running jobs, the options (again from left to right) are: Pause; Stop; Copy job artifacts; Delete job artifacts; and Live (show collection status in a modal).

What's Next

See the configuration instructions for the collector type you want to configure, Then proceed to instructions for scheduling and running collection jobs.

Filesystem/NFS
S3
Script
REST
Scheduling and Running

Filesystem/NFS

Cribl LogStream supports collecting data from a local or a remote filesystem location.

Configuring a Filesystem Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID: Unique ID for this Collector. E.g., DysonV11Roomba960.

Collector type: Defines the type of Collector to configure.

i Set this to Filesystem to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Auto-populate from: Select a Destination with which to auto-populate Collector settings. Useful when replaying data.

Directory: The directory from which to collect data. Templating is supported (e.g., /myDir/\${host}/\${year}/\${month}/). You can also use templating to specify (e.g.) a Splunk bucket from which to collect. Symlinks will not be followed. More on templates and Filters.

Recursive: If set to Yes (the default), data collection will recurse through subdirectories.

Destructive: If set to Yes, the Collector will delete files after collection. Defaults to No.

Max batch size (files): Maximum number of lines written to the discovery results files each time. To override this limit in the Collector's Schedule/Run modal, use Advanced Settings > Upper task bundle size.

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + Add Argument to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

- Pipeline: Select a Pipeline to process results.
- **Destination**: Select a Destination to receive results.
 - i You might disable **Send to Routes** when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. This keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

Cribl LogStream supports collecting data from Amazon S3 stores. This page covers how to configure the Collector.

For a step-by-step tutorial on using LogStream to replay data from an S3-compatible store, see our Data Collection & Replay sandbox. The sandbox takes about 30 minutes. It provides a hosted environment, with all inputs and outputs preconfigured for you.

Configuring an S3 Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID: Unique ID for this Collector. E.g., Attic42TreasureChest.

Collector type: Defines the type of Collector to configure.

i Set this to S3 to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Auto-populate from: Select a Destination with which to auto-populate Collector settings. Useful when replaying data.

S3 bucket: Simple Storage Service bucket from which to collect data.

Region: S3 Region from which to retrieve data.

Path: Path, within the bucket, from which to collect data. Templating is supported (e.g., /myDir/\${host}/\${year}/\${month}/). More on templates and Filters.

Recursive: If set to Yes (the default), data collection will recurse through subdirectories.

Max batch size (files): Maximum number of lines written to the discovery results files each time. To override this limit in the Collector's Schedule/Run modal, use Advanced Settings > Upper task bundle size.

Authentication

API key: Enter API key. If empty, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials. Optional when running on AWS.

Secret key: Enter secret key. if empty, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials. Optional when running on AWS.

Assume Role

Enable Assume Role: Slide to Yes to enable Assume Role behavior.

AssumeRole ARN: Amazon Resource Name (ARN) of the role to assume.

External ID: External ID to use when assuming role.

Additional Collector Settings

Endpoint: S3 service endpoint. If empty, LogStream will automatically construct the endpoint from the region.

Signature version: Signature version to use for signing S3 requests. Defaults to $\vee 4$.

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + **Add Argument** to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

• Pipeline: Select a Pipeline to process results.

• **Destination**: Select a Destination to receive results.

i You might disable **Send to Routes** when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. This keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data – potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

Scheduling and Running

Script

Cribl LogStream supports flexible data collection configured by your custom scripts.

Configuring a Script Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID: Unique ID for this Collector. E.g., sh2GetStuff.

Collector type: Defines the type of Collector to configure.

i Set this to Script to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Discover script: Script to discover which objects/files to collect. This script should output one task per line in stdout.

Collect script: Script to perform data collections. Pass in tasks from the Discover script as \$CRIBL_COLLECT_ARG . Should output results to stdout.

Shell: Shell in which to execute scripts. Defaults to /bin/bash.

Scripts will allow you to execute almost anything on the system where Cribl LogStream is running. Make sure you understand the impact of what you're executing before you do so! These scripts run

as the user running LogStream, so if you are running it as root, these commands will run with root user permissions. 🧟 🧟

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + Add Argument to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

- Pipeline: Select a Pipeline to process results.
- **Destination**: Select a Destination to receive results.
 - i You might disable **Send to Routes** when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. This keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

REST / API Endpoint

Cribl LogStream supports collecting data from REST endpoints.

Configuring a REST Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID

Unique ID for this Collector. E.g., rest42json.

Collector Type

Defines the type of Collector to configure.

i Set this to REST to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Discover Type

Once you've selected the REST Collector type above, this exposes a Discover type drop-down. Here you have four options, corresponding to different use cases. Each Discover type selection will expose a different set of Collector Settings fields. Below, we cover the Discover types from simplest to most-complex.

- **Discover type: None** matches cases where one simple API call will retrieve all the data you need. This suppresses the Discover stage. (Example: Collect a list of configured LogStream Pipelines.)
- **Discover type: Item List** matches cases where you want to enumerate a known list of items to retrieve. (Example: Collect network traffic data that's tagged with specific subnets.)
- Discover type: JSON Response provides a Discover result field where you can (optionally) define Discover tasks as a JSON array of objects. Each entry returned by Discover will generate a Collect task. (Example: Collect data for specific geo locations the National Weather Service API's stream of worldwide weather data. This API requires multiple parameters in the request URL latitude, longitude, etc. so an Item List would not work.)
- Discover type: HTTP Request matches cases where you need to dynamically discover what you can collect from a REST endpoint. This Discover type most fully exploits LogStream's Discover-Before-Collect architecture. (Example: Make a REST call to get a list of available log files, then run Collect against each of those files.)

Common Collector Settings / Discover Type: None

These remaining **Collector Settings** options appear for **Discover type**: None, as well as for all other **Discover type** selections:

i Time Range Variables

The following fields fields accept \${earliest} and \${latest} variables, which reference any **Time Range** values that have been set in manual or scheduled collection jobs:

- Collect URL, Collect parameters, Collect headers
- Discover URL, Discover parameters, Discover headers.

As an example, here is a **Collect URL** entry using these variables: http://localhost/path?from=\${earliest}&to=\${latest}

Both variables are formatted as UNIX epoch time, in seconds units. When using them in contexts that require milliseconds resolution, multiply them by 1,000 to convert to ms.

Collect URL: URL (constant or JavaScript expression) to use for the Collect operation.

i Any variables used in a URL (path or parameters) must be encoded using: C.Encode.uri(paramName).

As of v.2.3.2, URLs/expressions specified in this field follow redirects.

Collect method: Select the HTTP verb to use for the Collect operation – GET, POST, or POST with body.

Collect POST body: Template for POST body to send with the Collect request. (This field is displayed only when you set the Collect method to POST with body .) You can reference parameters from the Discover response using template params of the form: \${variable}.

Collect parameters: Optional HTTP request parameters to append to the request URL. These refine or narrow the request. Click **+ Add Parameter** to add parameters as key-value pairs:

- Name: Field name.
- Value: JavaScript expression to compute the field's value (can be a constant).

Collect headers:: Click + Add Header to (optionally) add collection request haaders as key-value pairs:

- Name: Header name.
- Value: JavaScript expression to compute the header's value (can be a constant).
 - i By adding the appropriate **Collect headers**, you can specify API Keybased authentication as an alternative to the Authentication: Basic or Login options below.

Authentication

In the **Authentication** drop-down, select an authentication method to use for discover and collect REST requests:

- None: Compatible with REST servers like AWS, where you embed a secret directly in the request URL.
- Basic: Compatible with Basic Authentication servers. Selecting Basic exposes additional fields in which you specify a Basic Auth zxUsername and Password.
- Login: Enables you to specify several credentials, then perform a POST to an endpoint during the Discover operation. The POST response returns a token, which LogStream uses for later Collect operations.

Selecting Login exposes the following additional fields:

- Login URL: URL for the login API call, which is expected to be a POST call.
- Username: Login username.
- Password: Login password.
- **POST Body**: Template for POST body to send with the login request. The \${username} and \${password} variables specify the corresponding credentials' locations in the message.
- **Token Attribute**: Path to the token attribute in the login response body. Supports nested attributes.
- Authorize Expression: JavaScript expression used to compute the Authorization header to pass in Discover and Collect calls. Uses \${token} to reference the token obtained from the login POST request.

Discover Type: Item List

Setting the **Discover type** to Item List exposes this additional field above the Common Collector Settings:

Discover Items: List of items to return from the Discover task. Each returned item will generate a Collect task, and can be referenced using \${id} in the **Collect URL**, the **Collect parameters**, or the **Collect headers**.

Discover Type: JSON Response

Setting the **Discover type** to JSON Response exposes these additional fields above the Common Collector Settings:

Discover result: Allows hard-coding the Discover result. Must be a JSON object. Works with the Discover data field.

Discover data field: Within the response JSON, name of the field or array element to pull results from. Leave blank if the result is an array of values. Sample entry: items, json: { items: [{id: 'first'},{id: 'second'}] }

Discover Type: HTTP Request

Setting the **Discover type** to HTTP Request exposes these additional fields above the Common Collector Settings:

Discover URL: Enter the URL to use for the Discover operation. This can be a constant URL, or a JavaScript expression to derive the URL.

i Any variables used in a URL (path or parameters) must be encoded using: C.Encode.uri(paramName).

As of v.2.3.2, URLs/expressions specified in this field follow redirects.

Discover method: Select the HTTP verb to use for the Discover operation – GET, POST, or POST with body.

Discover POST body: Template for POST body to send with the Discover request. (This field is displayed only when you set the **Discover method** to POST with body .)

Discover parameters: Optional HTTP request parameters to append to the Discover request URL. These refine or narrow the request. Click

- + Add Parameter to add parameters as key-value pairs:
 - Name: Parameter name.
 - Value: JavaScript expression to compute the parameter's value (can also be a constant).

Discover headers: Optional Discover request headers.: Click **+ Add Header** to add headers as key-value pairs:

- Name: Header name.
- Value: JavaScript expression to compute the header's value (can also be a constant).

Discover data field: Within the response JSON, name of the field that contains Discover results. Leave blank if the result is an array.

i The following sections describe the Collector Settings' remaining tabs, whose settings and content apply equally to all **Discover type** selections.

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + **Add Argument** to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

• Pipeline: Select a Pipeline to process results.

• **Destination**: Select a Destination to receive results.

i You might disable Send to Routes when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. One use case might be a REST Collector that gathers a known, simple type of data from a single endpoint. This approach keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data – potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

> Scheduling and Running

Scheduling and Running

Once you've configured a Collector, you can either run it immediately ("ad hoc") to collect data, or schedule it to run on a recurring interval. Scheduling requires some extra configuration upfront, so we cover this option first.

i For ad hoc collection, you can configure whether a job interrupted by a LogStream shutdown will automatically resume upon LogStream restart.

But regardless of this configuration, explicitly restarting or stopping LogStream (via ./cribl restart, ./cribl stop, or Settings > Controls > Restart) will cancel any currently running jobs.

A **scheduled** job will **not** resume upon restart.

Schedule Configuration

Click **Schedule** beside a configured Collector to display the **Schedule configuration** modal. This provides the following controls.

Enabled: Slide to Yes to enable this collection schedule.

⚠ The scheduled job will keep running on this schedule forever, unless you toggle **Enabled** back to Off. The Off setting preserves the schedule's configuration, but prevents its execution.

Cron schedule: A cron schedule on which to run this job.

• The **Estimated schedule** below this field shows the next few collection runs, as examples of the cron interval you've scheduled.

Skippable: Skippable jobs can be delayed up to their next run time if the system is hitting concurrency limits. Defaults to Yes.

Skippable Jobs and Concurrency Limits

If toggled to Yes, the **Skippable** option obliges these concurrency limits in **Settings** > **General Settings** > **Job Limits**:

- Concurrent Job Limit
- Concurrent Scheduled Job Limit

When the above limits delay a Skippable job:

- The Skippable job will be granted slightly higher priority than non-Skippable jobs.
- If the job receives resources to run before its next scheduled run, LogStream will run the delayed job, then snap back to the original cron schedule.
- If resources do not free up before the next scheduled run: LogStream will skip the delayed run, and snap back to the original cron schedule.

Set **Skippable** to No if you absolutely must have all your data, for compliance or other reasons. In this case, LogStream will build up a backlog of jobs to run.

You can think of **Skippable**: No as behaving more like the TCP protocol, with **Skippable**: Yes behaving more like UDP.

Max Concurrent Runs: Sets the maximum number of instances of this scheduled job that may simultaneously run.

⚠ All collection jobs are constrained by the following Settings > General Settings > Job Limits:

- Concurrent Task Limit
- Max Task Usage Percentage

Run Configuration and Shared Settings

Most of the remaining fields and options below are shared with the Run configuration modal, which you can open by clicking Run beside a configured Collector.

Mode

Depending on your requirements, you can schedule or run a collector in these modes:

- Preview default for Run, but not offered for Scheduled Jobs
- Discovery default for Scheduled Jobs
- Full Run

Preview

In the Preview mode, a collection job will return only a sample subset of matching results (e.g., 100 events). This is very useful in cases when users need a data sample to:

- Ensure that the correct data comes in.
- Iterate on filter expressions.
- Capture a sample to iterate on pipelines.

i Schedule configuration omits the Preview option, because Preview is designed for immediate analysis and decision making. To configure a Scheduled Job with high confidence, you might want to first manually run Preview jobs with the same Collector, to verify that you're collecting the data you expect.

Preview Settings

In Preview mode, you can optionally set these limits:

- Capture time (sec): Maximum time interval (in seconds) to collect data.
- Capture up to N events: Maximum number of events to capture.

Discovery

In Discovery mode, a collection job will return only **the list of objects/files** to be collected, but none of the data. This mode is typically used to ensure that the Filter expression and time range are correct before a Full Run job collects unintended data.

Send to Routes

In Discovery mode, this slider enables you to send discovery results to LogStream Routes. Defaults to $\,\mathrm{No}\,$.

i This setting overrides the Collector configuration's Result Routing > Send to Routes setting.

Full Run

In Full Run mode, the collection job is fully executed by Worker Nodes, and will return all data matching the Run configuration.

Time Range

Set an **Absolute** or **Relative** time range for data collection.

 $\begin{tabular}{ll} \hline \end{tabular} The \textbf{\it Relative} option is the default, and is particularly useful for configuring scheduled jobs. \\ \hline \end{tabular}$

Absolute

Select the **Absolute** button to set fixed collection boundaries in your local time. Next, use the **Earliest** and **Latest** controls to set the start date/time and end date/time.

Relative

Select the **Relative** button to set collection boundaries relative to the current time. Next, use the **Earliest** and **Latest** to set start and end times like these:

- Earliest example values: -1hr, -42m, -42m@h
- Latest example values: now, -20m, +42m@h

Relative Time Syntax

For Relative times, the **Earliest** and **Latest** controls accept the following syntax:

```
[+ |- ] < time_integer > < time_unit > @ < snap - to_time_unit >
```

To break down this syntax:

Syntax Element	Values Supported
Offset	Specify: - for times in the past, + for times in the future, or omit with now .
<time_integer></time_integer>	Specify any integer, or omit with now .
<time_unit></time_unit>	Specify the now constant, or one of the following abbreviations: s[econds], m[inutes], h[ours], d[ays], w[eeks], mon[ths], q[uarters], y[ears].
@ <snap- to_time_unit></snap- 	Optionally, you can append the @ modifier, followed by any of the above <time_unit> s, to round down to the nearest instance of that unit. (See the next section for details.)</time_unit>

LogStream validates relative time values using these rules:

- Earliest must not be later than Latest.
- Values without units get interpreted as seconds. (E.g., -1 = -1s.)

Snap-to-Time Syntax

The a snap modifier always rounds **down** (backwards) from any specified time. This is true even in relative time expressions with + (future) offsets. For example:

- ad snaps back to the beginning of today, 12:00 AM (midnight).
- +128mah looks forward 128 minutes, then snaps back to the nearest round hour. (If you specified this in the Latest field, and ran the Collector at 4:20 PM, collection would end at 6:00 PM. The expression would look forward to 6:28 PM, but snap back to 6:00 PM.)

Other options:

- aw or aw7 to snap back to the beginning of the week defined here as the preceding Sunday.
- To snap back to other days of a week, use w1 (Monday) through w6 (Saturday).
- am to snap back to the 1st of a month.

- aq to snap back to the beginning of the most recent quarter Jan. 1, Apr. 1, Jul. 1, or Oct. 1.
- ay to snap back to Jan. 1.

Filter

This is a JavaScript filter expression that is evaluated against token values in the provided collector path (see below), and against the events being collected. The **Filter** value defaults to true, which matches all data, but this value can be customized almost arbitrarily.

For example, if a Filesystem or S3 collector is run with this Filter:

```
host='myHost' & source.endsWith('.log') || source.endsWith('.txt')
```

...then only files/objects with .log or .txt extensions will be fetched. And, from those, only those events with host field myHost will be collected.

For more extensive options, see Tokens for Filtering below.

Advanced Settings

Log Level: Level at which to set task logging. More-verbose levels are useful for troubleshooting jobs and tasks, but use them sparingly.

Lower task bundle size: Limits the bundle size for small tasks. E.g., bundle five 200KB files into one 1MB task bundle. Defaults to 1MB.

Upper task bundle size: Limits the bundle size for files above the **Lower task bundle size**. E.g., bundle five 2MB files into one 10MB task bundle. Files greater than this size will be assigned to individual tasks. Defaults to 10MB.

Reschedule tasks: Whether to automatically reschedule tasks that failed with non-fatal errors. Defaults to Yes; does not apply to fatal errors.

Max task reschedule: Maximum number of times a task can be rescheduled. Defaults to 1.

Job timeout: Maximum time this job will be allowed to run. Units are seconds, if not specified. Sample values: 30, 45s, or 15m. Minimum granularity is 10 seconds, so a 45s value would round up to a 50-second timeout. Defaults to 0, meaning unlimited time (no timeout).

Tokens for Filtering

Let's look at the options for path-based (basic) and time-based token filtering.

Basic Tokens

In collectors with paths, such as Filesystem or S3, LogStream supports path filtering via token notation. Basic tokens' syntax follows that of JS template literals: \${<token_name>} - where token name is the field (name) of interest.

For example, if the path was set to \\var\log\\$\{\hostname\}\\$\{\sourcetype\}\, you could use a Filter such as \hostname='myHost' & sourcetype='mySourcetype' to collect data only from the \\var\log\myHost\mySourcetype\ subdirectory.

Time-based Tokens

In paths with time partitions, LogStream supports further filtering via time-based tokens. This has a direct effect with earliest and latest boundaries. When a job runs against a path with time partitions, the job traverses a minimal superset of the required directories to satisfy the time range, before subsequent event _time_filtering.

About Partitions and Tokens

LogStream processes time-based tokens as follows:

- For each path, time partitions must be notated in descending order. So Year/Month/Day order is supported, but Day/Month/Year is not.
- Paths may contain more than one partition. E.g., /my/path/2020-04/20/.
- In a given path, each time component can be used only once.
 So /my/path/\${_time:%Y}/\${_time:%m}/\${_time:%d}/... is a valid expression format, but /my/path/\${_time:%Y}/\${_time:%m}/\${host}/\${_time:%Y}/... (with a repeated Y) is not supported.
- For each path, all extracted dates/times are considered in UTC.

The following strptime format components are allowed:

- 'Yy', for years
- 'mBbj', for months
- 'dj', for days
- 'HI', for hours
- 'M', for minutes
- 'S', for seconds

Token Syntax

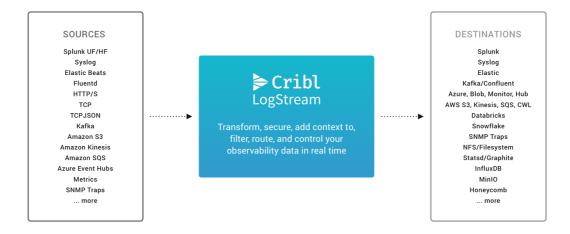
Time-based token syntax follows that of a slightly modified JS template literal:

\${_time: <some_strptime_format_component>}.Examples:

Filter	Matches
/my/path/\${_time:%Y}/\${_time:%m}/\${_time:%d}/	/my/path/2020/04/20/
/my/path/\${_time:year=%Y}/\${_time:month=%m}/\${_time:date=%d}/	/my/path/year=2020/month=
/my/path/\${_time:%Y-%m-%d}/	/my/path/2020-05-20/

Destinations

Cribl LogStream can send data to various Destinations, including Splunk, Kafka, Kinesis, InfluxDB, Snowflake, Databricks, TCP JSON, and many others.



Streaming Destinations

Destinations that accept events in real time are referred to as streaming Destinations:

- Splunk Single Instance
- Splunk Load Balanced
- Splunk HEC
- AWS Kinesis Streams
- AWS CloudWatch Logs
- AWS SQS
- Elasticsearch
- Honeycomb
- TCP JSON
- Syslog
- Kafka
- Azure Event Hubs
- Azure Monitor Logs
- StatsD
- StatsD Extended
- Graphite

- SNMP Trap
- InfluxDB
- New Relic
- Wavefront
- Sumo Logic
- Datadog
- SignalFx

Non-Streaming Destinations

Destinations that accept events in groups or batches are referred to as nonstreaming Destinations:

- S3 Compatible Stores
- Filesystem/NFS
- MinIO
- Azure Blob Storage
- Google Cloud Storage
 - i The S3 Compatible Stores Destination can be adapted to send data to downstream services like Databricks and Snowflake, for which LogStream currently has no preconfigured Destination. For details, please contact Cribl Support.

Other Destinations

LogStream also provides these special-purpose Destinations:

- Output Router: Flexible "meta-destination." Here, you can configure rules that route data to multiple outputs.
- DevNull: An output that simply drops events. Preconfigured and active when you install LogStream, so it requires no configuration. Useful for testing.
- Default: Here, you can specify a default output from among your configured Destinations.

How Does Non-Streaming Delivery Work

Cribl LogStream uses a staging directory in the local filesystem to format and write outputted events before sending them to configured Destinations. After a set of conditions is met – typically file size and number of files, further details below – data is compressed and then moved to the final Destination.

An inventory of open, or in-progress, files is kept in the staging directory's root, to avoid having to walk that directory at startup. This can get expensive if staging is also the final directory. At startup, Cribl LogStream will check for any leftover files in progress from prior sessions, and will ensure that they're moved to their final Destination. The process of moving to the final Destination is delayed after startup (default delay: 30 seconds). Processing of these files is paced at one file per service period (which defaults to 1 second).

Batching Conditions

Several conditions govern when files are closed and rolled out:

- 1. File reaches its configured maximum size.
- 2. File reaches its configured maximum open time.
- 3. File reaches its configured maximum idle time.

If a new file needs to be open, Cribl LogStream will enforce the maximum number of open files, by closing files in the order in which they were opened.

Data Delivery

Data is delivered to all Destinations on an at-least-once basis. When a Destination is unreachable, there are three possible behaviors:

- Block Cribl LogStream will block incoming events.
- **Drop** Cribl LogStream will drop events addressed to that Destination.
- Queue Cribl LogStream will Persistent-Queue events to that Destination.

You can configure the desired behavior through a Destination's **Backpressure Behavior** option. If this option is not present, Cribl LogStream's default behavior is to **Block**.

Configuring Destinations

For each Destination **type**, you can create multiple definitions, depending on your requirements.

To configure Destinations, select **Data > Destinations**, select the desired type from the tiles or the left menu, then click **+ Add New**.

Output Router

Output Routers are meta-destinations that allow for output selection based on rules. Rules are evaluated in order, top->down, with the first match being the winner.

Configuring Cribl LogStream to Send to an Output Router

Select **Data > Destinations**, then select **Output Router** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Output Router > New Destination** modal, which provides the following fields.

Router name: Enter a unique name to identify this Router definition.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Rules: A list of event routing rules. Each provides the following settings:

- Filter expression: JavaScript expression to select events to send to output.
- Output: Output to send matching events to.
- **Description**: Optionally, enter a description of this rule's purpose.
- **Final**: Flag that controls whether to stop the event from being checked against other rules lower in the stack. Defaults to Yes.

Notes

- An Output Router cannot reference another. This is by design, so as to avoid cycles.
- Events that do not match any of the rules are dropped. Use a catchall rule to change this behavior.

- No post-processing (conditioning) can be done here. Use Pre-Processing Pipelines on the Source tier.
- Data can be cloned by toggling the Final flag to No . (The default is Yes , i.e., no cloning.)

Example

Scenario:

- Send all events where host starts with 66 to Destination San Francisco.
- From the rest of the events:
 - Send all events with method field POST or GET to both Seattle and Los Angeles (i.e., clone).
- Send the remaining events to New York City.

Router Name: router66

Filter Expression	Output	Final
host.startsWith('66')	San Francisco	Yes
method='POST' method='GET	Seattle	No
method='POST' method='GET'	Los Angeles	Yes
true	New York	Yes

Splunk Single Instance

Splunk Enterprise is a streaming Destination type.

Configuring Cribl LogStream to Output to Splunk Destinations

Select **Data > Destinations**, then select **Splunk > Single Instance** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Single Instance > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Splunk Single Instance definition.

Address: Hostname of the Splunk receiver.

Port: The port number on the host.

Nested field serialization: Specifies how to serialize nested fields into indextime fields. Defaults to None.

Throttling: Throttle rate in bytes per second. Multiple byte units such as KB, MB, GB etc. are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttle engaged, excesses data will be dropped only if Backpressure Behavior is set to drop, and blocked for all other settings.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Output multi metrics: Toggle to Yes to output multiple-measurement metric data points. (Supported in Splunk 8.0 and above, this format enables sending multiple metrics in a single event, improving the efficiency of your Splunk capacity.)

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

i Single .pem File

If you have a **single** .pem file containing cacert, key, and cert sections, enter it in all of these fields above: **CA certificate path**, **Private key path** (**mutual auth**), and **Certificate path** (**mutual auth**).

Timeout Settings

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (in milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Notes about Forwarding to Splunk

Data sent to Splunk is not compressed.

- If events have a Cribl LogStream internal field called __criblMetrics , they'll be forwarded to Splunk as metric events.
- If events do **not** have a _raw field, they'll be serialized to JSON prior to sending to Splunk.

Splunk Load Balanced

Splunk is a streaming Destination type, and with the **Splunk Load Balanced** output, you can load-balance data out to multiple Splunk receivers.

How Does Load Balancing Work

Cribl LogStream will attempt to load-balance outbound data as fairly as possibly across all receivers. Data is sent by all worker processes to all receivers simultaneously, and the amount sent to each receiver depends on these parameters:

- 1. Respective destination weight.
- 2. Respective destination historical data.

By default, historical data is tracked for 300s. LogStream uses this data to influence the traffic sent to each destination, to ensure that differences decay over time, and that total ratios converge towards configured weights.

Example

Suppose we have two receivers, A and B, each with weight of 1 (i.e., they are configured to receive equal amounts of data). Suppose further that the load-balance stats period is set at the default 300s and – to make things easy – for each period, there are 200 events of equal size (Bytes) that need to be balanced.

Interval	Time Range	Events to be dispensed
1	time=0s> time=300s	200

Both A and B start this interval with 0 historical stats each.

Let's assume that, due to various circumstances, 200 events are "balanced" as follows:

A = 120 events and B = 80 events – a difference of 40 events and a ratio of 1.5:1.

Interval	Time Range	Events to be dispensed

At the beginning of interval 2, the load-balancing algorithm will look back to the previous interval stats and carry **half** of the receiving stats forward. I.e., receiver A will start the interval with **60** and receiver B with **40**. To determine how many events A and B will receive during this next interval, LogStream will use their weights and their stats as follows:

Total number of events: events to be dispensed + stats carried forward = 200 + 60 + 40 = 300.

Number of events per each destination (weighed): 300/2 = 150 (they're equal, due to equal weight).

Number of events to send to each destination A: 150 - 60 = 90 and B: 150 - 40 = 110.

Totals at end of interval 2: A=120+90=210, B=80+110=190, a difference of 20 events and a ratio of 1.1:1.

Over the subsequent intervals, the difference becomes exponentially less pronounced, and eventually insignificant. Thus, the load gets balanced fairly.

Configuring Cribl LogStream to Load-Balance to Multiple Splunk Destinations

To configure load balancing, first select **Data > Destinations**, then select **Splunk > Load Balanced** from the **Data Destinations** page's tiles or left menu. Then click **Add New** to open the **Load Balanced > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Splunk LB Destination definition.

Indexer Discovery: When toggled to Yes, enables automatic discovery of indexers in an indexer clustering environment. See Indexer Discovery for the resulting UI options displayed below. When set to No (the default), displays the Destinations section below.

Exclude current host IPs: Exclude all IPs of the current host from the list of any resolved hostnames. Defaults to Yes.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block. When toggled to Persistent Queue, adds the Persistent Queue Settings section (left tab) to the modal.

Destinations

The **Destinations** section appears only when **Indexer discovery** is set to No. Here, you specify a known set of Splunk receivers on which to load-balance data.

Click + Add Destination to specify more receivers on new rows. Each row provides the following fields:

- Address: Hostname of the Splunk receiver. Optionally, you can paste in a comma-separated list, in <host>:<port> format.
- Port: Port number to send data to.
- TLS: Whether to inherit TLS configs from group setting, or disable TLS.
 Defaults to inherit.
- TLS servername: Servername to use if establishing a TLS connection. If not specified, defaults to connection host (if not an IP). Otherwise, uses the global TLS settings.
- Load weight: The weight to apply to this Destination for load-balancing purposes.

Indexer Discovery

Toggling the **Indexer Discovery** toggle to Yes displays the following fields instead of the **Destinations** section:

Site: Clustering site from which indexers need to be discovered. In the case of a single site cluster, default is the default entry.

Cluster Master URI: Full URI of Splunk Cluster Master, in the format: scheme://host:port.(Worker Nodes normally access the Cluster Master on port 8089 to get the list of currently online indexers.)

Auth token: Authentication token required to authenticate to Cluster Master for indexer discovery.

Refresh period: Time interval (in seconds) between two consecutive fetches of indexer list from Cluster Master. Defaults to 60.

Enabling Cluster Master Authentication

To enable token authentication on the Splunk Cluster Master, you can find complete instructions in Splunk's Enable or Disable Token Authentication documentation. The following capabilites are required:

list_indexer_cluster and list_indexerdiscovery.

For details on creating the token, see Splunk's Create Authentication Tokens topic – especially its section on how to Configure Token Expiry and "Not Before" Settings.

⚠ Be sure to give the token an **Expiration** setting well in the future, whether you use **Relative Time** or **Absolute Time**. Otherwise, the token will inherit Splunk's default expiration time of +30d (30 days in the future), which will cause indexer discovery to fail.

If you have a failover site configured on Splunk's Cluster Master, Cribl respects this configuration, and forwards the data to the failover site in case of site failure.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

i Single PEM File

If you have a **single** .pem file containing cacert, key, and cert sections, enter this file's path in all of these fields above:

CA certificate path, Private key path (mutual auth), and Certificate path (mutual auth).

Timeout Settings

- **Connection timeout**: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000 ms.
- Write timeout: Amount of time (in milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000 ms.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Output Multi Metrics: Toggle this slider to Yes to output multiple-measurement metric data points. (Supported in Splunk 8.0 and above, this format enables sending multiple metrics in a single event, improving the efficiency of your Splunk capacity.)

Minimize in-flight data loss: If set to Yes (the default), LogStream will check whether the indexer is shutting down and, if so, stop sending data. This helps minimize data loss during shutdown.

DNS resolution period (seconds): Re-resolve any hostnames after each interval of this many seconds, and pick up destinations from A records. Defaults to 600 seconds.

Load balance stats period (seconds): Lookback traffic history period. Defaults to 300 seconds. (Note that If multiple receivers are behind a hostname – i.e., multiple A records – all resolved IPs will inherit the weight of the host, unless each IP is specified separately. In Cribl LogStream load balancing, IP settings take priority over those from hostnames.)

Max connections: Constrains the number of concurrent indexer connections, per Worker Process, to limit memory utilization. If set to a number > 0, then on every DNS resolution period (or indexer discovery), LogStream will randomly select this subset of discovered IPs to connect to. LogStream will rotate IPs in future resolution periods – monitoring weight and historical data, to ensure fair load balancing of events among IPs.

Nested field serialization: Specifies whether and how to serialize nested fields into index-time fields. Select None (the default) or JSON.

Throttling: Throttle rate, in bytes per second. Multiple byte units such as KB, MB, GB, etc., are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttling is engaged, excess data will be dropped only if **Backpressure behavior** is set to **Drop events**. (Data will be blocked for all other **Backpressure behavior** settings.)

SSL Configuration for Splunk Cloud – Special Note

To connect to Splunk Cloud, you **might** need to extract the private and public key from the Splunk-provided Splunk Cloud Certificate, which is typically bundled in an app. Use the following steps:

Step 1. Test connectivity to Splunk Cloud, using the Root CA certificate:

```
openssl s_client -CApath path_to_ca.pem -connect
hostnameToSplunkCloud:9997
```

Step 2. Extract the Private key from the Splunk Cloud Certificate. At the prompt, you will need the sslPassword value from the outputs.conf file bundled with the Splunk Cloud app. Using Elliptic Curve keys:

```
openssl ec -in path_to_server_cert.pem -out private.pem
```

If you are using RSA keys, instead use:

```
openssl rsa -in path_to_server_cert.pem -out private.pem
```

Step 3. Extract the Public Key for the Server Certificate:

```
openssl x509 -in path_to_server_cert.pem -out server.pem
```

Step 4. In the LogStream Destination's **TLS Settings (Client Side)** section, enter the following:

• CA Certificate Path: Path to CA Certificate.

- Private Key Path (mutual auth): Path to private.pem (Step 2 above).
- Certificate Path (mutual auth): Path to server.pem (Step 3 above).

Notes About Forwarding to Splunk

- Data sent to Splunk is not compressed.
- If events have a LogStream internal field called __criblMetrics , they'll be forwarded to Splunk as metric events.
- If events do **not** have a _raw field, they'll be serialized to JSON prior to sending to Splunk.

Splunk HEC

Splunk HEC is a streaming Destination type. In a typical deployment, Cribl LogStream will be installed/co-located in a Splunk heavy forwarder. If this output is enabled, it can send data out to a Splunk HEC (HTTP Event Collector) destination through the event endpoint.

Configuring Cribl LogStream to Output to Splunk HEC Destinations

Select **Data > Destinations**, then select **Splunk > HEC** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **HEC > New Destination** modeal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Splunk HEC definition.

Splunk HEC endpoint: URL of a Splunk HEC endpoint to send events to (e.g., http://myhost.example.com:8088/services/collector/event).

HEC auth token: Splunk HEC authentication token.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the Backpressure behavior is set to Persistent Queue.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5. Each request can potentially hit a different HEC receiver.

Max body size (KB): Maximum size, in KB, of the request body. Defaults to 4096. Lowering the size can potentially result in more parallel requests and also cause outbound requests to be made sooner.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured Max body size. Defaults to 1.

- Retries happen on this flush interval.
 - Any HTTP response code in the 2xx range is considered success.
 - Any response code in the 5xx range is considered a retryable error, which will not trigger Persistent Queue (PQ) usage.
 - Any other response code will trigger PQ (if PQ is configured as the Backpressure behavior).

Extra HTTP headers: Click + **Add Header** to add **Name/Value** pairs to pass as additional HTTP headers.

Next processing queue: Specify the next Splunk processing queue to send the events to, after HEC processing. Defaults to indexQueue.

Default_TCP_ROUTING: Specify the value of the _TCP_ROUTING field for events that do not have _ctrl._TCP_ROUTING set. Defaults to nowhere .

i This is useful only when you expect the HEC receiver to route this data on to another destination.

Output multi metrics: Toggle to Yes to output multiple-measurement metric data points. (Supported in Splunk 8.0 and above, this format enables sending multiple metrics in a single event, improving the efficiency of your Splunk capacity.)

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new connection will be reattempted. This is to prevent sticking to a particular Destination when there is a constant flow of events.
- If the server does not support keepalives or if the server closes a pooled connection while idle – a new connection will be established for next request.
- When resolving the Destination's hostname, LogStream will pick the first IP
 in the list for use in the next connection. Round-robin DNS would help with
 event balancing.

Page 405 of 1257

S3 Compatible Stores

S3 is a non-streaming Destination type. Cribl LogStream does **not** have to run on AWS in order to deliver data to S3.

Stores that are S3-compatible will also work with this Destination type. For example, the S3 Destination can be adapted to send data to services like Databricks and Snowflake, for which LogStream currently has no preconfigured Destination. For these integrations, please contact Cribl Support.

Configuring Cribl LogStream to Output to S3 Destinations

Select Data > Destinations, then select Amazon > S3 from the Data Destinations page's tiles or left menu. Click Add New to open the S3 > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this S3 definition.

S3 bucket name: Name of the destination S3 Bucket. This value can be a constant, or a JavaScript expression that can be evaluated only at init time. E.g., referencing a Global Variable: myBucket-\${C.vars.myVar}.

i Event-level variables are not available for JavaScript expressions. This is because the bucket name is evaluated only at Destination initialization. If you want to use event-level variables in file paths, Cribl recommends specifying them in the Partitioning Expression field (described below), because this is evaluated for each file.

Region: Region where the S3 bucket is located.

Staging location: Filesystem location in which to locally buffer files before compressing and moving to final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: When set to Yes (the default), adds the **Output ID** field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:

- Clone the Destination.
- Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with **Add output ID** enabled.

Key prefix: Prefix to add to files before uploading. This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. If left blank, Cribl LogStream will fall back to event.__partition . Defaults to `\${host}/\${sourcetype}`. Partitioning by time is also possible, e.g., `\${host}/\${C.Time.strftime(_time, '%Y-%m-%d')}/\${sourcetype}`

Data format: Format of the output data. Defaults to JSON.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Select the data compression format to use before moving data to final destination. Defaults to none. Cribl recommends setting this to gzip.

Backpressure behavior: Select whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

Authentication

Authentication method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for S3: Whether to use Assume Role credentials to access S3. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports <code>c*</code> wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Endpoint: S3 service endpoint. If empty, the endpoint will be automatically constructed from the region.

Object ACL: Object ACL (Access Control List) to assign to uploaded objects.

Storage class: Select a storage class for uploaded objects. Defaults to Standard.

Server side encryption: Server side encryption type for uploaded objects. Defaults to none. This is used for enabling encryption for data at rest.

Signature version: Signature version to use for signing S3 requests. Defaults to v4.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location.

Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Amazon S3 Permissions

The following permissions are needed to write to an Amazon S3 bucket:

s3:GetObject

s3:ListBucket

s3:GetBucketLocation

s3:PutObject

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

• __partition

Kinesis Streams

Cribl LogStream can output events to **Amazon Kinesis Data Streams** records of up to 1MB uncompressed. Cribl LogStream does **not** have to run on AWS in order to deliver data to a Kinesis Data Stream.

Configuring Cribl LogStream to Output to Amazon Kinesis Data Streams

Select **Data > Destinations**, then select **Amazon > Kinesis** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Kinesis > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Kinesis definition.

Stream name: Enter the name of the Kinesis Data Stream to which to send events.

Region: Select the AWS Region where the Kinesis Data Stream is located.

Endpoint: Kinesis Stream service endpoint. If empty, the endpoint will be automatically constructed from the region.

Signature version: Signature version to use for signing Kinesis stream requests. Defaults to v4.

Put request concurrency: Maximum number of ongoing put requests before blocking. Defaults to 5.

Max record size (KB, uncompressed): Maximum size of each individual record before compression. For non-compressible data, 1MB (the default) is the maximum recommended size.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the Backpressure behavior is set to Persistent Queue.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Authentication

Authentication method: Select an AWS authentication method.

- Auto: This default option uses the environment variables AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env. AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for Kinesis Streams: Whether to use Assume Role credentials to access Kinesis Streams. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes cribl_pipe (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Format

Currently, outputted events use the following record format:

- Header line containing information about the payload (currently supports one type, as shown below).
- Newline-Delimited JSON (that is, each Kinesis record will contain multiple events, in **ndjson** format).

Record payloads (including header and body) will be gzip-compressed, and then Kinesis will base64-encode them.

Sample Kinesis Record

```
{"format":"ndjson","count":8,"size":3960}
{"_raw":"07-03-2018 18:33:51.136 -0700 ERROR TcpOutputFd - Read error. Connection reset by peer'
{"_raw":"07-03-2018 18:33:51.136 -0700 INFO TcpOutputProc - Connection to 127.0.0.1:10000 close
```

CloudWatch Logs

Cribl LogStream supports sending data to Amazon CloudWatch Logs. This is a streaming Destination type. Cribl LogStream does **not** have to run on AWS in order to deliver data to CloudWatch Logs.

Configuring Cribl LogStream to Output to Amazon CloudWatch Logs

Select **Data > Destinations**, then select **Amazon > CloudWatch Logs** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **CloudWatch Logs > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this CloudWatch definition.

Log group name: CloudWatch log group to associate events with.

Log stream prefix: Prefix for CloudWatch log stream name. This prefix will be used to generate a unique log stream name per Cribl LogStream instance. (E.g., myStream_myHost_myOutputId .)

Region: AWS region where the CloudWatch Logs group is located.

Endpoint: CloudWatch Logs service endpoint. If empty, defaults to AWS' Region-specific endpoint. Otherwise, use this field to point to a CloudWatchLogs-compatible endpoint.

Signature version: Signature version to use for signing CloudWatch Logs requests. Defaults to $\ v4$.

Max queue size: Maximum number of queued batches before blocking. Defaults to 5.

Max record size (KB, uncompressed): Maximum size of each individual record before compression. For non-compressible data, 1MB (the default) is the maximum recommended size.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Authentication

Authentication Method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

 API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials. Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for CloudWatch Logs: Whether to use Assume Role credentials to access CloudWatch Logs. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

SQS

Cribl LogStream supports sending events to Amazon Simple Queuing Service.

Configuring Cribl LogStream to Send Data to Amazon SQS

Select **Data > Destinations**, then select **Amazon > SQS** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **SQS > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this SQS Destination.

Queue name: The name, URL, or ARN of the SQS queue to send events to.
This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks. To specify a non-AWS URL, use the format: '{url}/<queueName>' . (E.g., ':port/<myQueueName>' .)

Message group ID: This parameter applies only to queues of type FIFO. Enter the tag that specifies that a message belongs to a specific message group. (Messages belonging to the same message group are processed in FIFO order.) Defaults to cribl. Use event field __messageGroupId to override this value.

Create queue: Specifies whether to create the queue if it does not exist. Defaults to Yes.

Region: Region where SQS queue is located.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the Backpressure behavior is set to Persistent Queue. Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Authentication

Authentication Method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for SQS: Whether to use Assume Role credentials to access SQS. Defaults to No.

AWS account ID: Enter the SQS queue owner's AWS account ID. Leave empty if the SQS queue is in the same AWS account where this LogStream instance is located.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Endpoint: SQS service endpoint. If empty, the endpoint will be automatically constructed from the region.

Signature version: Signature version to use for signing SQS requests. Defaults to v4.

Max queue size: Maximum number of queued batches before blocking. Defaults to 100.

Max record size (KB): Maximum size of each individual record. Per the SQS spec, the maximum allowed value is 256 KB. (the default).

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Max concurrent requests: The maximum number of in-progress API requests before backpressure is applied. Defaults to 10.

SQS Permissions

The following permissions are needed to write to an SQS queue:

- sqs:ListQueues
- sqs:SendMessage
- sqs:SendMessageBatch
- sqs:CreateQueue
- sqs:GetQueueAttributes
- sqs:SetQueueAttributes
- sqs:GetQueueUrl

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions can use them to make processing decisions.

Fields for this Destination:

- __messageGroupId
- __sqsMsgAttrs
- _sqsSysAttrs

Filesystem/NFS

Filesystem is a non-streaming Destination type that Cribl LogStream can use to output files to a local or a network-attached filesystem (NFS).

Configuring Cribl LogStream to Output to Filesystem Destinations

Select **Data > Destinations**, then select **Filesystem** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Filesystem > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Filesystem definition.

Output location: Final destination for the output files.

Staging location: Local filesystem location in which to buffer files before compressing and moving them to the final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: When set to Yes (the default), adds the **Output ID** field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

- For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:
 - Clone the Destination.
 - Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with **Add output ID** enabled.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. Defaults to `\${host}/\${sourcetype}`. If left blank, Cribl LogStream will fall back to event.__partition . Partitioning by time is also possible, e.g.: `\${host}/\${C.Time.strftime(_time, '%Y-%m-%d')}/\${sourcetype}`

Data format: Format of the output data. Defaults to json.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Data compression format used before moving to final destination. Default none . It is recommended that gzip is used.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this will be closed and moved to final output location. Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Backpressure Behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

- __partition
 - i To export events from an intermediate stage within a Pipeline to a file, see the Tee Function.

Elasticsearch

Cribl LogStream can send events to an Elasticsearch cluster using the Bulk API.

Configuring Cribl LogStream to Output to Elasticsearch

Select Data > Destinations, then select Elasticsearch from the

Data Destinations page's tiles or left menu. Click Add New to open the

Elasticsearch > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Elasticsearch Destination definition.

Bulk API URL: URL of an Elasticsearch cluster to send events to.
(E.g., http://<myElasticCluster>:9200/_bulk .)

Index: Elasticsearch Index where to send events to. Note that this value can be overwritten by an event's __index field.

Type: Specify document type to use for events. Note that this value can be overwritten by an event's __type field.

Authentication enabled: Set to No by default. Toggle to Yes to enter a **Username** and **Password**.

Backpressure behavior: Specify whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (s): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Field Normalization

This Destination normalizes the following fields:

- _time becomes atimestamp atmillisecond resolution.
- host.name is set to host.

See also our Elasticsearch Source documentation's Field Normalization section.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this Destination:

- __id
- __type
- __index

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular destination when there is a constant flow of events.
- If the server does not support keepalives (or if the server closes a pooled connection while idle), a new connection will be established for the next request.
- When resolving the Destination's hostname, LogStream will pick the first IP
 in the list for use in the next connection. Round-robin DNS would help with
 event balancing.

Honeycomb

Cribl LogStream supports sending events to a Honeycomb dataset.

Configuring Cribl LogStream to Output to Honeycomb

Select **Data > Destinations**, then select **Honeycomb** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Honeycomb > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Honeycomb definition.

Dataset name: Name of the dataset to send events to. (E.g., iLoveObservabilityDataset .)

API Key: Team API Key to which the dataset belongs. (E.g., teamWilde .)

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to $\, 1 \, \, \text{MB} \,$.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id> . Defaults to

\$CRIBL HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular Destination when there is a constant flow of events.
- If the server does not support keepalives (or if the server closes a pooled connection while idle), a new connection will be established for the next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

TCP JSON

Cribl LogStream supports sending data over TCP in JSON format. **TCP JSON** is a streaming Destination type.

Configuring Cribl LogStream to Output in TCP JSON Format

Select **Data > Destinations**, then select **TCP JSON** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **TCP JSON > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Destination definition.

Address: Hostname of the receiver.

Port: Port number to connect to on the host.

Auth token: Optional authentication token to include as part of the connection header. Defaults to empty.

Compression: Codec to use to compress the data before sending. Defaults to None.

Throttling: Throttle rate in bytes per second. Multiple byte units such as KB, MB, GB etc. are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttle engaged, excesses data will be dropped only if Backpressure Behavior is set to drop, and blocked for all other settings.

Backpressure behavior: Specifies whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Autofill?: This setting is experimental.

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Timeout Settings

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish before retrying. Defaults to 10000.

Write timeout: Amount of time (in milliseconds) to wait for a write to complete before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Format

TCP JSON events are sent in newline-delimited JSON format, consisting of:

- 1. A header line. Can be empty, e.g.: {} . If **Auth Token** is enabled, the token will be included here as a field called authToken . In addition, if events contain common fields, they will be included here under fields .
- 2. A JSON event/record per line.

Example

See an example in our TCP JSON Source topic.

Syslog

Cribl LogStream supports sending of data over syslog via TCP. Syslog is a streaming Destination type.

i This Syslog Destination supports RFC 3164 and RFC 5424.

Configuring Cribl LogStream to output in **Syslog** format

Select **Data > Destinations**, then select **Syslog** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Syslog > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Syslog definition.

Protocol: The network protocol to use for sending out syslog messages. Defaults to TCP; UDP is also available.

Address: Address/hostname of the receiver.

Port: Port number to connect to on the host.

Facility: Default value for message facility. If set, will be overwritten by the value of __facility . Defaults to user .

Severity: Default value for message severity. If set, will be overwritten by the value of __severity . Defaults to notice .

App name: Default value for application name. If set, will be overwritten by the value of __appname . Defaults to Cribl .

Message format: The syslog message format supported by the receiver. Defaults to RFC3164.

Timestamp format: The timestamp format to use when serializing an event's time field. Defaults to Syslog.

Throttling: Throttle rate in bytes per second. Multiple byte units such as KB, MB, GB etc. are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttle engaged, excesses data will be dropped only if Backpressure Behavior is set to drop, and blocked for all other settings.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No . When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Timeout Settings

i These timeout settings apply only to the TCP protocol.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.

- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this destination:

- __priority
- __facility
- __severity
- __procid
- __appname
- __msgid
- __syslogout

Kafka

Cribl LogStream supports sending data to a Kafka topic. Kafka is a streaming Destination type.

Configuring Cribl LogStream to Output to Kafka

Select **Data > Destinations**, then select **Kafka** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Kafka > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Kafka definition.

Brokers: List of Kafka brokers to connect to. (E.g., localhost:9092.)

Topic: The topic on which to publish events. Can be overwritten using event's __topic field.

Acknowledgments: Select the number of required acknowledgments. Defaults to Leader.

Record data format: Format to use to serialize events before writing to Kafka. Defaults to JSON.

Compression: Codec to compress the data before sending to Kafka. Select None, Gzip, or Snappy.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Authentication

Authentication parameters to use when connecting to brokers. Using TLS is highly recommended.

Enabled: Defaults to No. When toggled to Yes:

- **SASL** mechanism: Select the SASL (Simple Authentication and Security Layer) authentication mechanism to use,
- Username: The username for authentication.
- Password: The password for authentication.

Schema Registry

This section governs Kafka Schema Registry Authentication for AVRO-encoded data with a schema stored in the Confluent Schema Registry.

Enabled: defaults to No. When toggled to Yes:

- Schema registry URL: URL for access to the Confluent Schema Registry. (E.g., http://<hostname>:8081.)
- Default key schema ID: Used when __keySchemaIdOut is not present to transform key values. Leave blank if key transformation is not required by default.
- Default value schema ID: Used when __valueSchemaIdOut not present to transform _raw . Leave blank if value transformation is not required by default.
- TLS enabled: defaults to No. When toggled to Yes, displays the following TLS settings for the Schema Registry:

TLS Settings (Schema Registry)

- i These have the same format as the TLS Settings (Client Side) above.
- Validate server certs: Require client to reject any connection that is not authorized by a CA in the CA certificate path, or by another trusted CA (e.g., the system's CA). Defaults to No.
- **Server name (SNI)**: Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.
- Certificate name: The name of the predefined certificate.
- CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.
- Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS. Use only if mutual auth is required.
- Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS. Use only if mutual auth is required.
- Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (KB, uncompressed): Maximum size (KB) of each record batch before compression. Setting should be < message.max.bytes settings in Kafka brokers. Defaults to 768.

Max events per batch: Maximum number of events in a batch before forcing a flush. Defaults to 1000.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this Destination:

- _topicOut
- __key
- __headers

- <u>__keySchemaIdOut</u>
- __valueSchemaIdOut

Azure Blob Storage

Azure Blob Storage is a non-streaming Destination type. Cribl LogStream does not have to run on Azure in order to deliver data to it. Azure Data Lake Storage Gen2 (hierarchical namespace) is also supported.

Configuring Cribl LogStream to Output to Azure Blob Storage

Select **Data > Destinations**, then select **Azure > Azure Blob Storage** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Blob Storage > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Destination definition.

Connection string: Enter your Azure Connection String. The format is:

DefaultEndpointsProtocol=[http|https];AccountName=<your-account-name>;AccountKey=
<your-account-key>

A fictitious example, using Microsoft's recommended HTTPS option:

DefaultEndpointsProtocol=https;AccountName=storagesample;AccountKey=12345678 ... 32

If left blank, LogStream will fall back to: env.AZURE_STORAGE_CONNECTION_STRING.

☐ The Connection string field replaces the Account name and Account key fields provided in this Destination's UI prior to LogStream 2.4. If you configured an Azure Blob Storage Destination before upgrading to LogStream 2.4, those fields' values are now concatenated into the Connection string value.

In case values were concatenated incorrectly, the original field keys and values are retained in LogStream's configuration files (e.g.,

\$CRIBL_HOME/groups/<group-name>/cribl/outputs.yml). However, Cribl recommends simply re-entering the correct Azure **Connection string** here.

For further details, see Microsoft's Configure Azure Storage Connection Strings documentation.

Container name: Enter the container name. (A container organizes a set of blobs, similar to a directory in a file system.)

Create container: Defaults to No . Toggle to Yes to create the configured container in Azure Blob Storage if it does not already exist.

Blob prefix: Root directory to prepend to path before uploading.

Staging location: Local filesystem location in which to buffer files before compressing and moving them to the final destination. Cribl recommends that this location be stable and highperformance.

Add Output ID: When set to Yes (the default), adds the Output ID field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

- ▲ For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:
 - Clone the Destination.
 - Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with Add output ID enabled.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. Defaults to `\${host}/\${sourcetype}` If left blank, Cribl LogStream will fall back to event.__partition.

Data format: Format of the output data. Defaults to json.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Data compression format used before moving to final destination. Defaults to none. Cribl recommends setting to gzip.

Backpressure behavior: Whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes cribl_pipe (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max file size (MB): Maximum uncompressed output file size. Files reaching this size will be closed and moved to the final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location. Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Default: 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Default: 100.

i LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

• __partition

Azure Monitor Logs

Cribl LogStream supports sending of data over to Azure Monitor Logs. This is a streaming Destination type.

Configuring Cribl LogStream to Output to Azure Monitor Logs

Select **Data > Destinations**, then select **Azure > Monitor Logs** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Monitor Logs > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Azure Monitor Logs definition.

Workspace ID: Enter the Azure Log Analytics Workspace ID. (See **Workspace->Advanced settings** in the Azure Dashboard.)

Workspace key: Enter the Azure Log Analytics Workspace Primary or Secondary Shared Key. (In the Azure Dashboard, see **Workspace->Advanced settings**.)

Log type: The Record Type of events sent to this LogAnalytics workspace. Defaults to Cribl.

Resource ID: Resource ID of the Azure resource to associate the data with. This populates the _ResourceId property, and allows the data to be included in resource-centric queries. (Optional, but if this field is not specified, the data will not be included in resource-centric queries.)

Backpressure behavior: Whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular Destination when there is a constant flow of events.
- If keepalives are not supported by the server (or if the server closes a
 pooled connection while idle), a new connection will be established for the
 next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

Azure Event Hubs

Cribl LogStream supports sending data to Azure Event Hubs. This is a streaming Destination type.

Configuring Cribl LogStream to Output to Azure Event Hubs

Select Data > Destinations, then select Azure > Event Hubs from the Data Destinations page's tiles or left menu. Click Add New to open the Event Hubs > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Azure Event Hubs definition.

Brokers: List of Event Hub Kafka brokers to connect to. (E.g., yourdomain.servicebus.windows.net:9093.) Find the hostname in Shared Access Policies, in the host portion of the primary or secondary connection string.

Event Hub name: The name of the Event Hub (a.k.a., Kafka Topic) on which to publish events. Can be overwritten using the __topicOut field.

Acknowledgments: Control the number of required acknowledgments. Defaults to Leader.

Record data format: Format to use to serialize events before writing to the Event Hub Kafka brokers. Defaults to JSON.

Compression: Codec to use to compress the data before sending it to Event Hub Kafka brokers. Defaults to Gzip.

Backpressure behavior: Whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled Defaults to Yes.

Validate server certs: Defaults to No – and for Event Hubs, this must always be disabled.

Authentication

Authentication parameters to use when connecting to brokers. Using TLS is highly recommended.

Enabled: Defaults to Yes. (Toggling to No hides the remaining settings in this group.)

SASL mechanism: SASL (Simple Authentication and Security Layer) authentication mechanism to use, PLAIN is the only mechanism currently supported for Event Hub Kafka brokers.

Username: The username for authentication. For Event Hub, this should always be \$ConnectionString.

Password: Connection-string primary key or connection-string secondary key from the Event Hub workspace.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (KB, uncompressed): Maximum size (KB) of each record batch before compression. Setting should be < message.max.bytes settings in Kafka brokers. Defaults to 768.

Max events per batch: Maximum number of events in a batch before forcing a flush. Defaults to 1000.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this Destination:

•	topicOut
•	key
•	headers
•	keySchemaIdOut
•	valueSchemaIdOut

Google Cloud Storage

Google Cloud Storage is a non-streaming Destination type.

Configuring Cribl LogStream to Output to Google Cloud Storage Destinations

Select **Data > Destinations**, then select **Google Cloud > Cloud Storage** from the **Data Destinations** page's tiles or left menu.

Next, click **Add New** to open the **Cloud Storage > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Cloud Storage definition.

Bucket name: Name of the destination bucket. This value can be a constant. or a JavaScript expression that can be evaluated only at init time. E.g., referencing a Global Variable: myBucket-\${C.vars.myVar}.

Region: Region where the bucket is located.

Staging location: Filesystem location in which to locally buffer files before compressing and moving to final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: Whether to append output's ID to staging location. Defaults to Yes .

Key prefix: Prefix to add to files before uploading. This value can be a constant or a JavaScript expression.

Partitioning expression: JavaScript expression to define how files are
partitioned and organized. If left blank, Cribl LogStream will fall back to
 event.__partition . Defaults to `\${host}/\${sourcetype}` . Partitioning
by time is also possible, e.g., `\${host}/\${C.Time.strftime(_time, '%Y-%m-%d')}/\${sourcetype}`

Data format: Format of the output data. Defaults to JSON.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Select the data compression format to use before moving data to final destination. Defaults to none. Cribl recommends setting this to gzip.

Backpressure behavior: Select whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

Authentication

Authentication is via HMAC (Hash-based Message Authentication Code). To create a key and secret, see Google Cloud's Managing HMAC Keys for Service Accounts documentation

Access key: Enter the HMAC access key.

Secret key: Enter the HMAC secret.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports <code>c*</code> wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Endpoint: The Google Cloud Storage service endpoint. Typically, there is no reason to change the default https://storage.googleapis.com endpoint.

Object ACL: Select an Access Control List to assign to uploaded objects. Defaults to private.

Storage class: Select a storage class for uploaded objects.

Signature version: Signature version to use for signing requests. Defaults to $\,^{\vee4}$.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location.

Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

__partition

StatsD

Cribl LogStream supports sending data to a StatsD Destination. This is a streaming Destination type.

Configuring Cribl LogStream to Output via StatsD

While on the **Data Destinations** page, select **Metrics > StatsD** from the tiles or the left menu, then click **Add New**. The resulting **New StatsD destination** pane contains the following fields.

Select Data > Destinations, then select Metrics > StatsD from the

Data Destinations page's tiles or left menu. Click Add New to open the StatsD >

New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this StatsD definition.

Destination protocol: Protocol to use when communicating with the Destination. Defaults to UDP.

Host: The hostname of the Destination.

Port: Destination port. Defaults to 8125.

i The next two settings apply only to the TCP protocol, and are not displayed for UDP.

Throttling: Rate (in bytes per second) at which at which to throttle while writing to an output. Also takes numerical values in multiples of bytes (KB, MB, GB, etc.). Default value of 0 indicates no throttling.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

- i This section is displayed only for TCP, and only when the Backpressure behavior is set to Persistent Queue.
- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Timeout Settings

i These timeout settings apply only to the TCP protocol, and are not displayed for UDP.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

• cribl_host - LogStream Node that processed the event.

- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (bytes): Used when Protocol is UDP. Specifies the maximum size of packets sent to the Destination. (Also known as the MTU – maximum transmission unit – for the network path to the Destination system.) Defaults to 512.

Flush period (sec): Used when Protocol is TCP. Specifies how often buffers should be flushed, sending records to the Destination. Defaults to 1.

StatsD Extended

Cribl LogStream supports sending data to a StatsD Destination. This is a streaming Destination type.

Configuring Cribl LogStream to Output via StatsD Extended

Select Data > Destinations, then select Metrics > StatsD Extended from the Data Destinations page's tiles or left menu. Click Add New to open the StatsD Extended > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this StatsD Extended definition.

Destination protocol: Protocol to use when communicating with the Destination. Defaults to UDP.

Host: The hostname of the Destination.

Port: Destination port. Defaults to 8125.

i The next two settings apply only to the TCP protocol, and are not displayed for UDP.

Throttling: Rate (in bytes per second) at which at which to throttle while writing to an output. Also takes numerical values in multiples of bytes (KB, MB, GB, etc.). Default value of 0 indicates no throttling.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

- i This section is displayed only for TCP, and only when the Backpressure behavior is set to Persistent Queue.
- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Timeout Settings

i These timeout settings apply only to the TCP protocol, and are not displayed for UDP.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

cribl_host - LogStream Node that processed the event.

- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (bytes): Used when Protocol is UDP. Specifies the maximum size of packets sent to the Destination. (Also known as the MTU – maximum transmission unit – for the network path to the Destination system.) Defaults to 512.

Flush period (sec): Used when Protocol is TCP. Specifies how often buffers should be flushed, sending records to the Destination. Defaults to 1.

Graphite

Cribl LogStream supports sending data to a Graphite backend Destination. This is a streaming Destination type.

Configuring Cribl LogStream to Output to a Graphite Backend

Select **Data > Destinations**, then select **Metrics > Graphite** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Graphite > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Graphite definition.

Destination protocol: Protocol to use when communicating with the Destination. Defaults to UDP.

Host: The hostname of the Destination.

Port: Destination port. Defaults to 8125.

i The next two settings apply only to the TCP protocol, and are not displayed for UDP.

Throttling: Rate (in bytes per second) at which at which to throttle while writing to an output. Also takes numerical values in multiples of bytes (KB, MB, GB, etc.). Default value of 0 indicates no throttling.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

- i This section is displayed only for TCP, and only when the Backpressure behavior is set to Persistent Queue.
- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Timeout Settings

i These timeout settings apply only to the TCP protocol, and are not displayed for UDP.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

cribl_host - LogStream Node that processed the event.

- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (bytes): Used when Protocol is UDP. Specifies the maximum size of packets sent to the Destination. (Also known as the MTU – maximum transmission unit – for the network path to the destination system.) Defaults to 512.

Flush period (sec): Used when Protocol is TCP. Specifies how often buffers should be flushed, sending records to the Destination. Defaults to 1.

SNMP Trap

Cribl LogStream supports forwarding of SNMP Traps out.

Configuring Cribl LogStream to Forward to SNMP Traps

While on the **Data Destinations** page, select **SNMP Trap** from the tiles or the left menu, then click **Add New**. The resulting **SNMP Trap > New Destination** modal, which provides the following fields.

Select **Data > Destinations**, then select **SNMP Trap** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **New SNMP destination** pane, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this SNMP Trap definition.

SNMP Trap destinations: One or more SNMP destinations to forward traps to.

- Address: Destination host.
- Port: Destination port. Defaults to 162.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.

• cribl_output - LogStream Destination that processed the event.

Considerations for Working with SNMP Traps Data

- It's possible to work with SNMP metadata (i.e., we'll decode the packet).
 Options include dropping, routing, etc. However, packets cannot be modified and sent to another SNMP Destination.
- SNMP packets can be forwarded to non-SNMP Destinations (e.g., Splunk, Syslog, S3, etc.).
- SNMP packets can be forwarded to other SNMP Destinations. However, the contents of the incoming packet cannot be modified i.e., we'll forward the packets verbatim as they came in.
- Non-SNMP input data cannot be sent to SNMP Destinations.

InfluxDB

Cribl LogStream supports sending data to InfluxDB.

Configuring Cribl LogStream to Output to InfluxDB

Select **Data > Destinations**, then select **InfluxDB** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **InfluxDB > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this InfluxDB definition.

Write API URL: URL of an InfluxDB cluster to send events to. (E.g., http://localhost:8086/write.)

Database name: The database on which to write data points.

Timestamp precision: Sets the precision for the supplied UNIX time values. Defaults to Milliseconds .

Dynamic value fields: When enabled, LogStream will pull the value field from the metric name. (E.g., db.query.user will use db.query as the measurement and user as the value field). Defaults to Yes.

Value field name: Name of the field in which to store the metric when sending to InfluxDB. This will be used as a fallback if dynamic name generation is enabled but fails. Defaults to value.

Authentication enabled: Set to No by default. Toggle to Yes to enter a **Username** and **Password**.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

MinIO

MinIO is a non-streaming Destination type, to which Cribl LogStream can output objects.

Configuring Cribl LogStream to Output to MinIO Destinations.

Select **Data > Destinations**, then select **MinIO** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **MinIO > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this MinIO definition.

MinIO endpoint: MinIO service URL (e.g., http://minioHost:9000).

MinIO bucket name: Name of the destination MinIO bucket. This value can be a constant, or a JavaScript expression that can be evaluated only at init time. E.g., referencing a Global Variable: myBucket-\${C.vars.myVar} . Ensure that the bucket already exists, otherwise MinIO will generate "bucket does not exist" errors.

API key: If left blank, LogStream will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.

Secret key: If left blank, Cribl LogStream will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Staging location: Filesystem location in which to locally buffer files before compressing and moving to final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: When set to Yes (the default), adds the **Output ID** field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

- For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:
 - Clone the Destination.
 - Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with **Add output ID** enabled.

Key prefix: Prefix to apply to files/objects before uploading to the specified bucket. MinIO will display key prefixes as folders.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. If left blank, Cribl LogStream will fall back to event.__partition. Defaults to `\${host}/\${sourcetype}`.

i LogStream's internal __partition field can be populated in multiple ways. The precedence order is: explicit Partitioning expression value -> \${host}/\${sourcetype} (default) Partitioning expression value -> user-defined event.__partition, set with an Eval Function (takes effect only where this Partitioning expression field is blank).

Data format: Format of the output data. Defaults to json.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Select the data compression format to use before moving data to final destination. Defaults to none. Cribl recommends setting this to gzip.

Backpressure behavior: Select whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

How MinIO Composes File Names

The full path to a file consists of:

```
<bucket_name>/<keyprefix><partition_expression | __partition>
<file_name_prefix><filename>.<extension>
```

As an example, assume that the MinIO bucket name is bucket1, the Key prefix is aws, the Partitioning expression is `\${host}/\${sourcetype}`, the source is undefined, the File name prefix is the default CriblOut, and the Data format is json. Here, the full path as displayed in MinIO would have this form: /bucket1/aws/192.168.1.241/undefined/CriblOut-<randomstring>0.json

i Although MinIO will display the Key prefix and Partitioning expression values as folders, both are actually just part of the overall key name, along with the file name.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Region: Region where the MinIO service/cluster is located. Leave blank when using a containerized MinIO.

Object ACL: ACL (Access Control List) to assign to uploaded objects. Defaults to Private .

Storage class: Select a storage class for uploaded objects. Defaults to Standard.

Server side encryption: Server side encryption type for uploaded objects. Defaults to none.

Signature version: Signature version to use for signing MinIO requests. Defaults to v4.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location.

Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the

Max file open time (sec) conditions is met.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

• __partition

New Relic

Cribl LogStream supports sending events to the New Relic Log API and the New Relic Metric API.

Configuring Cribl LogStream to Output to New Relic

Select Data > Destinations, then select New Relic from the Data Destinations page's tiles or left menu. Click Add New to open the New Relic > New destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this New Relic definition.

API key: Enter your New Relic Insert API key, as created on New Relic's **Insights API keys** page.

Region: Select which New Relic region endpoint to use.

Log type: Name of the logType to send with events. E.g., observability or access_log.

i This sets a default. Where a sourcetype is specified in an event, it will override this value.

Log message field: Name of the field to send as the log message value. If not specified, the event will be serialized and sent as JSON.

Fields: Additional metadata fields to (optionally) add, as **Name-Value** pairs.

- Name: Enter the metadata field name.
- Value: JavaScript expression to compute field's value, enclosed in single quotes, double quotes, or backticks. (Can evaluate to a constant.)
- Add Field: Click to add more metadata Name-Value pairs.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

 $\begin{tabular}{ll} \textbf{Compress}: Whether to compress the payload body before sending. Defaults to No . \end{tabular}$

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 1000 KB.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured **Max body size**. Defaults to 1 second.

Extra HTTP headers: Click + Add Header to insert extra headers as Name/Value pairs.

Wavefront

Cribl LogStream supports sending events to Wavefront analytics.

Configuring Cribl LogStream to Output to Wavefront

Select **Data > Destinations**, then select **Wavefront** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Wavefront > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Wavefront definition.

Auth token: Wavefront API authentication token. For details, see Wavefront's Generating an API Token topic. Required.

Domain name: WaveFront domain name, e.g., longboard. Required.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to $\, 1 \, \, \text{MB} \,$.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id> . Defaults to

\$CRIBL HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

 $\begin{tabular}{ll} \textbf{Compress}: Whether to compress the payload body before sending. Defaults to No . \end{tabular}$

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured **Max body size**. Defaults to 1 second.

Extra HTTP headers: Click + **Add Header** to insert extra headers as **Name/Value** pairs.

Notes About Wavefront

For details on integrating with Wavefront, see these Wavefront resources:

- Direct Data Ingestion, and adjacent topics on Wavefront Proxies.
- Wavefront Data Format.

SignalFx

Cribl LogStream supports sending events to SignalFx.

Configuring Cribl LogStream to Output to SignalFx

Select **Data > Destinations**, then select **SignalFx** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **SignalFx > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this SignalFx definition.

Auth token: SignalFx API access token. For details, see SignalFx's Manage Tokens topic. Required.

Realm: SignalFx realm name (e.g., us0). Required.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Whether to compress the payload body before sending. Defaults to No .

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured **Max body size**. Defaults to 1 second.

Extra HTTP headers: Click + **Add Header** to insert extra headers as **Name/Value** pairs.

Notes About SignalFx

For details on integrating with SignalFx, see the SignalFx Developers Guide, with particular reference to the SignalFx HTTP Send Metrics Reference.

Sumo Logic

Cribl LogStream can send log and metric events to Sumo Logic over HTTP.

Configuring Cribl LogStream to Output to Sumo Logic

Select **Data > Destinations**, then select **Sumo Logic** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Sumo Logic > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Sumo Logic Destination definition.

API URL: Enter the URL of the Sumo Logic HTTP collector to which events should be sent. (E.g.,

https://endpoint6.collection.us2.sumologic.com/receiver/v1/http/<longhash>.)

Custom source name: Optionally, override the source name configured on the Sumo Logic HTTP collector. This value will be sent with events via the X-Sumo-Name HTTP header.

Custom source category: Optionally, override the source category configured on the Sumo Logic HTTP collector. This value will be sent with events via the X-Sumo-Category HTTP header.

Backpressure behavior: Optionally, override the source category configured on the Sumo Logic HTTP collector.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

If an event contains the internal field __criblMetrics , LogStream will send it to Sumo Logic as a metric event. Otherwise, LogStream will send it as a log event.

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular destination when there is a constant flow of events.
- If the server does not support keepalives (or if the server closes a pooled connection while idle), a new connection will be established for the next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

Datadog

Cribl LogStream can send log and metric events to Datadog. (Datadog supports metrics only of type gauge, counter, and rate via its REST API.)

LogStream sends events to the following Datadog endpoints in the US region. Use a DNS lookup to discover and include the corresponding IP addresses in your firewall rules' allowlist.

- Logs: https://http-intake.logs.datadoghq.com/v1/input
- Metrics: https://api.datadoghq.com/api/v1/series

Configuring Cribl LogStream to Output to Datadog

Select **Data > Destinations**, then select **Datadog** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Datadog > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Destination definition.

API Key: Enter an API key available in your Datadog profile.

Send logs as: Specify the content type to use when sending logs. Defaults to application/json, where each log message is represented by a JSON object. The alternative text/plain option sends one message per line, with newline \n delimiters.

Message field: Name of the event field that contains the message to send. If not specified, LogStream sends a JSON representation of the whole event (regardless of whether **Send logs as** is set to JSON or plain text).

Source: Name of the source to send with logs. If you're sending logs as JSON objects (i.e., you've selected **Send logs as**: application/json), the event's source field (if set) will override this value.

Host: Name of the host to send with logs. If you're sending logs as JSON objects, the event's host field (if set) will override this value.

Service: Name of the service to send with logs. If you're sending logs as JSON objects, the event's __service field (if set) will override this value.

Tags: List of tags to send with logs (e.g., env:prod, env_staging:east).

Severity: Default value for message severity. If you're sending logs as JSON objects, the event's __severity field (if set) will override this value. Defaults to info; the drop-down offers many other severity options.

i Datadog uses the above five fields (source , host , __service , __severity , and tags) to enhance searches and UX.

Backpressure behavior: Specify whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None . Select Gzip to enable compression.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress log events' payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (s): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

If an event contains the internal field __criblMetrics , LogStream will send it to Datadog as a metric event. Otherwise, LogStream will send it as a log event.

You can use these fields to override outbound event values for log events:

- __service
- __severity

No internal fields are supported for metric events.

For More Information

You might find these Datadog references helpful:

- Submit Metrics
- Send Logs
- Metrics Types

DevNull

The DevNull Destination simply drops events. Cribl provides this as a basic output to test Pipelines and Routes.

Configuring Cribl LogStream to Forward to DevNull

DevNull requires no configuration: A DevNull Destination is preconfigured and active as soon as you install Cribl LogStream.

To verify this, select **Data > Destinations** from the top menu. On the resulting **Data Destinations** page, select **DevNull** from the tiles or the left menu. Look for the **Live** indicator at the top right.

Default

The **Default** Destination simply enables you to specify a default output from among your configured Destinations.

Configuring Cribl LogStream's Default Destination

Select **Data > Destinations**, then select **Default** from the **Data Destinations** page's tiles or left menu. From the resulting **Manage Default Destination** page, click anywhere on the default row to proceed.



Default Destination - click to configure

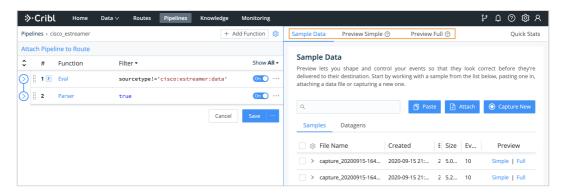
In the resulting **Destinations > Default** modal, use the **Default Output ID** dropdown to select one of your configured Destinations as LogStream's default.

The only other field here is the Output ID, whose value is locked to default.

Data Preview

Sample Data Preview is a LogStream feature that allows for visual inspection of events as they make their trip into a Pipeline. It helps you shape and control events before they're delivered to a Destination, as well as assisting with troubleshooting LogStream Functions.

Preview works by taking a set of Sample events, passing them through the Pipeline, and displaying the result in a separate pane. Any time a Function is modified, added, or removed, the Pipeline changes, and so does its displayed output.

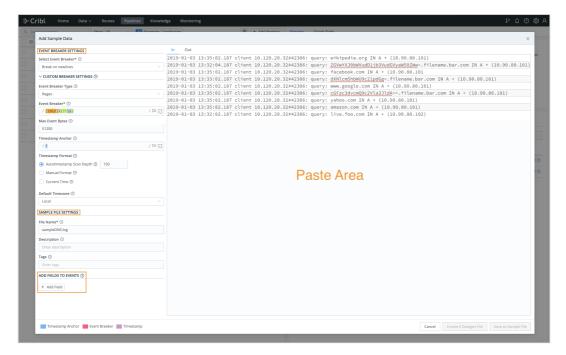


Preview options

While you're in a Pipeline, you can add samples through one of the supported options: Paste, Attach, or Capture New. The Paste and Attach options work with content that needs to be broken into events, while the Capture New option works with events only.

Adding Sample Data (Using Paste as an Example)

When you click on the corresponding option, you'll be presented with a modal like the one shown below.



Add Sample Data modal

i The Capture New modal is slightly different – it does not require event breaking.

Paste Area

This is where the content of the paste (or uploaded file) is displayed.

Event Breaker Settings

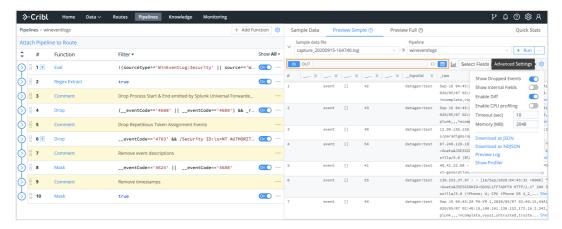
An Event Breaker is a regular expression that tells Cribl LogStream how to break the file or pasted content into events. Breaking will occur at the **start** of the match. Cribl LogStream ships with several common breaker patterns out of the box, but you can also configure custom breakers. The UI here is interactive, and you can iterate until you find the exact pattern.

Fields

The Fields section enables users to add, or overwrite. key/value pairs on the sample.

IN Tab: Displaying Samples on the Way IN to the Pipeline

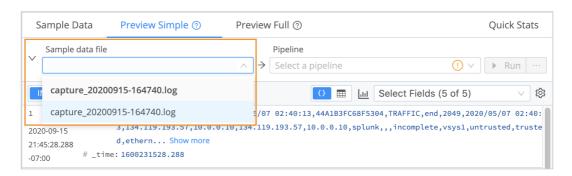
The Preview pane offers two display options for the event: Event and Table. (You can also download data as JSON or NDJSON, using the **Advanced Settings** menu at the top right.) Each format can be useful, depending on the type of data you are previewing.



Event, Table, and Advanced options

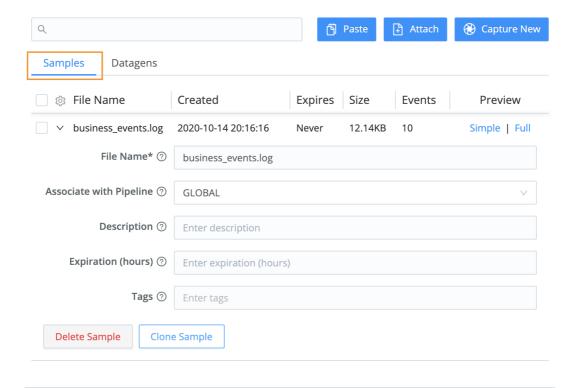
In the Advanced Settings menu's Timeout (sec) and Memory (MB) fields, you can increase the defaults to adjust for cases where very large data samples fail to load. For example, you might increase the Timeout (sec) to 30 and the Memory (MB) to 3048.

As you add more samples to your system, you can easily access them via the Samples drop-down near the top right, under Quick Stats.



Selecting an existing sample

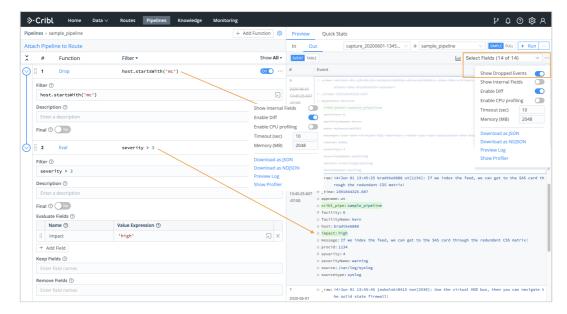
You can also manage, clone/modify, and delete samples via the **Samples** tab below.



OUT Tab: Displaying Samples on the Way OUT of the Pipeline

As data traverses Functions in a Pipeline, events can be modified, and some might be dropped altogether. The **OUT** tab indicates changes using this color coding:

- **Dropped events**: When events are dropped, the **OUT** tab displays them as grayed-out text, with strikethrough. You can control their display using the **Advanced Settings** menu's **Show Dropped Events** slider.
- Added fields: When LogStream's processing adds new fields, these fields are highlighted green. You can control these fields' display using the Select Fields drop-down.
- Redacted fields: These fields are highlighted amber.
- **Deleted fields**: These fields are highlighted red.



Dropped and added fields in a Pipeline's output

Securing Data

Cribl LogStream can be used to encrypt sensitive data in real time and route it to an end system. Decrypted retrieval can be implemented on a per-system basis. Currently, decryption is supported only when Splunk is the end system.

- Data Encryption
- Data Decryption

Encryption

Encryption of Data in Motion

With Cribl LogStream, you can encrypt fields or patterns within events in real time, by using C.Crypto.encrypt() in a Mask function. The Mask function accepts multiple replacement rules and multiple fields to apply them to.

A Match Regex defines the pattern of content to be replaced. The Replace Expression is a JS expression or literal to replace matched content. The C.Crypto.encrypt() method can be used here to generate an encrypted string from a value passed to it.

i C.Crypto.encrypt() Syntax

(method) Crypto.encrypt(value: any, keyclass: number,
keyId?: string, defaultVal?: string): string
Encrypt the given value with the keyId or a keyId picked up
automatically based on keyclass

- @param {string | Buffer} value what to encrypt
- @param keyclass if keyld isn't specified, pick one at the given key class
- @param keyld encryption keyld, takes precedence over keyclass
- @param defaultVal what to return if encryptions fails for any reason, if unspecified the original value is returned
- @returns - if encryption succeeds the encrypted value, otherwise defaultVal if specified, otherwise value.

Encryption Keys

Symmetric keys can be configured through the CLI or UI. Users are free to define as many keys as required. Each key is characterized by the following:

- keyId: ID of the key.
- algorithm: Algorithm used with the key
- keyclass: Cribl Key Class (below) that the key belongs to.
- kms: Key management system for the key. Defaults to local.

- created: Time (epoch) when key was generated.
- expires: Time (epoch) after which the key is invalid. Useful for key rotation.
- useIV: Flag that indicates whether or not an initialization vector was used.

Key Classes

Key Classes in Cribl LogStream are collections of keys that can be used to implement multiple levels of access control. Users (or groups of users) with access to data with encrypted patterns can be associated with key classes, for even more granular, pattern-level compartmentalized access.

Example

Users U0, U1 have been given access to keyclass 0 which contains key IDs 0 and 1. These keys are used to encrypt certain patterns in datasetA. Even though users U0, U1, U2 have access to read this dataset, only U0 and U1 can decrypt its encrypted patterns.

Key Class	Dataset
keyclass: 0 Keys: keyId: 0, keyId: 1 Users: U0, U1	datasetA Users: U0, U1, U2

User U1 has been given access to an **additional** keyclass, 1, which contains key IDs 11 and 22. These keys are used to encrypt certain **other** patterns in dataset A. Even though users U0, U1, U2 have access to read this dataset – same as above – only U1 can decrypt the additional encrypted patterns.

Key Class	Dataset
keyclass: 1 Keys: keyId: 11, keyId: 22 Users: U1	datasetA Users: U0, U1, U2

Configuring Keys with CLI

When using the local key management system, encryption keys in Cribl LogStream are encrypted with

\$CRIBL_HOME/local/cribl/auth/cribl.secret and stored in

\$CRIBL_HOME/local/cribl/auth/keys.json.Cribl monitors the keys.json file for changes every 60 seconds.

i When installed as a Splunk app, \$CRIBL_HOME is \$SPLUNK_HOME/etc/apps/cribl.

Listing Keys

Keys are added and listed using the keys command:

\$CRIBL_HOME/bin/cribl keys list -g <workerGroupID>

Sample Command Output

keyId	algorithm	keyclass	kms	created	expires	useIV
1	aes-256-cbc	0	local	1544906269.316	0	false
2	aes-256-cbc	_		1544906272.452	_	false
3	aes-256-cbc	2	local	1544906275.948	1545906275	true
4	aes-256-cbc	3	local	1544906278.026	0	false

Adding Keys

```
Displaying --help:

$CRIBL_HOME/bin/cribl keys add --help

SampleCommandOutput

Add encryption keys
Usage: [options] [args]

Options:

[-c <keyclass>] - key class to set for the key

[-k <kms>] - KMS to use, must be configured, see cribl.yml

[-e <expires>] - expiration time, epoch time

[-i] - use an initialization vector

-g <group> - Group ID
```

Adding a key to keyclass 1, with no expiration date, on the default Worker Group:

```
$CRIBL_HOME/bin/cribl keys add -c 1 -i -g default
```

Sample Command Output

Adding key: success. Key count=1

Listing keys to verify key generation:

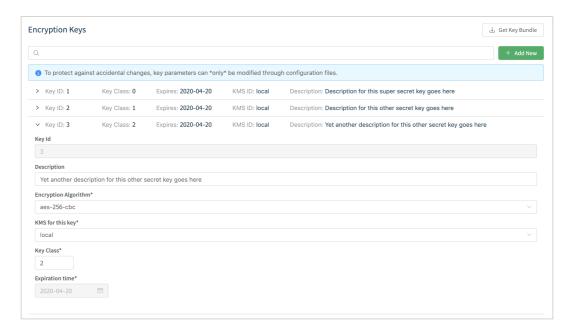
\$CRIBL_HOME/bin/cribl keys list -g default

Sample Command Output

keyId	algorithm	keyclass	kms	created	expires	useIV
1	aes-256-cbc	1	local	1545243364.342	0	true

Configuring Keys with the UI

The key management interface can be accessed through **Settings > Encryption Keys**. Here, you can list and add new keys. To protect against accidental changes, a key's parameters, once saved, can be edited only through configuration files.



List or create keys through LogStream's UI

Sync auth/(cribl.secret|keys.json)

To successfully decrypt data, the decrypt command will need access to the same keys that were used to encrypt. The cribl.secret and keys.json files in \$CRIBL_HOME/local/cribl/auth (in the Cribl instance where encryption happened) should be synced/copied over to the ones on the

Search Head/decrypting side. When using the UI, these files can be downloaded through the **Get Key Bundle** button.

Decryption

Decryption of Data

Currently, Cribl LogStream supports decryption only when Splunk is the end system. In Splunk, decryption is available to users of any role with permissions to run the decrypt command that ships with Cribl App for Splunk. Further restrictions can be applied with Splunk capabilities. This page provides details.

Decrypting in Splunk

Decryption in Splunk is implemented via a custom command called decrypt. To use the command, users must belong to a Splunk role that has permissions to execute it. Capabilities, which are aligned to Cribl Key Classes, can be associated with a particular role to further control the scope of decrypt.

i Decrypt Command Is Search Head ONLY

To ensure that keys don't get distributed to all search peers – including peers that your search head can search, but you don't have full control over – decrypt is scoped to run locally on the installed search head.

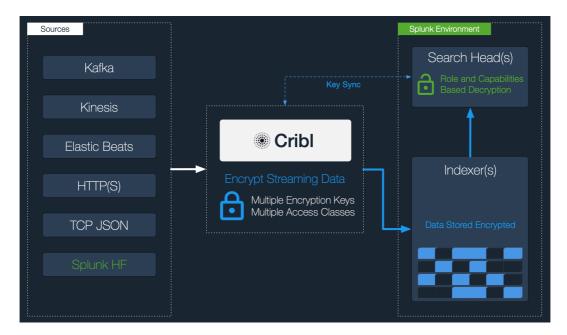
Restricting Access with Splunk Capabilities

In Splunk, capability names should follow the format cribl_keyclass_N, where N is the Cribl Key Class. For example, a role with capability cribl_keyclass_1 has access to all key IDs associated with key class_1.

Capability Name	Corresponding Cribl Key Class
cribl_keyclass_1	1
cribl_keyclass_2	2
•••	
cribl_keyclass_N	N

Configuring Splunk Search Head to Decrypt Data

You set up decryption in Splunk according to this schematic:



1. Install the Cribl App for Splunk on your Splunk search head.

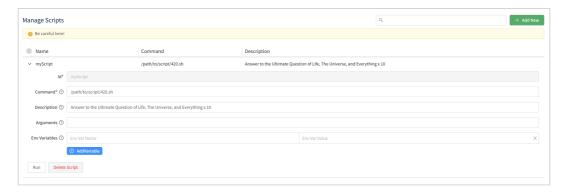
As of LogStream v1.7, the app will run in search head mode by default. If the app has previously been installed and later modified, you can convert it to search head mode with the command: \$CRIBL_HOME/bin/cribld mode-searchhead. (When installed as a Splunk app, \$CRIBL_HOME is \$SPLUNK_HOME/etc/apps/cribl.)

- 2. Assign permissions to the decrypt command, per your requirements.
- 3. Assign capabilities to your roles, per your requirements. If you'd like to create more capabilities, ensure that they follow the naming convention defined above.
- 4. Sync auth/(cribl.secret|keys.json). To successfully decrypt data, the decrypt command will need access to the same keys that were used to encrypt. The cribl.secret and keys.json files in \$CRIBL_HOME/local/cribl/auth which must be in the same Cribl instance where encryption happened should be synced/copied over to the files on the Search Head/decrypting side. When using the UI, these files can be downloaded through the Get Key Bundle button.

Scripts

Admins can run scripts (e.g., shell scripts) from within Cribl LogStream by configuring and executing them thru **Settings** > **Scripts**. Scripts are typically used to call custom automation jobs or, more generally, to trigger tasks on demand. For example, you can use Scripts to run an Ansible job, or to place a call to another automation system, when Cribl LogStream configs are updated.

Scripts will allow you to execute almost anything on the system where Cribl LogStream is running. Make sure you understand the impact of what you're executing before you do so!



Settings > Manage Scripts page

The Manage Scripts page provides the following tields:

- ID: Unique ID for this script.
- **Command**: Command to execute for this script.
- **Description**: Brief description about this script. Optional.
- Arguments: Arguments to pass when executing this script.
- Env variables: Extra environment variables to set when executing script.

i Scripts in Distributed Deployments

• Scripts can be deployed from Master Node, but can be **run** only locally from each Worker Node.

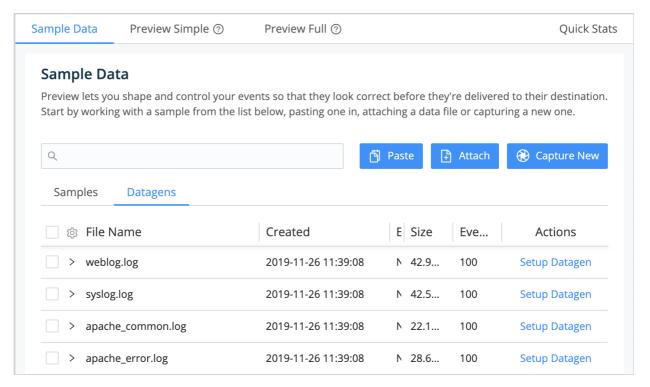
If the Script command is referencing a file (e.g., 420.sh), that file
must exist on the Cribl LogStream instance. In other words, the
Script management interface cannot be used to upload or
manage script files.

Using Datagens

Data generators for testing and troubleshooting

Cribl LogStream's Datagens feature enables you to generate sample data for the purposes of troubleshooting Routes, Pipelines, Functions, and general connectivity.

Several Datagen template files ship with the product, out of the box. You can create others from sample files or live captures.



Preview pane – add samples via paste, attach/upload file, or live capture

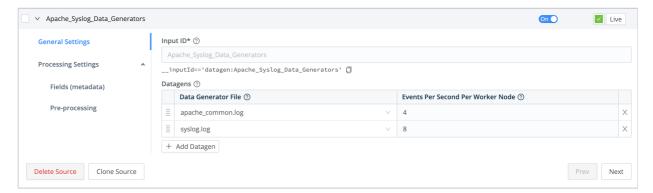
As outlined in the following tutorial: Once you've created a template, you can configure a Datagen Source to use the template to generate real-time data at a given EPS (events per second) rate.

Enabling a Datagen

To see how Datagens work, start by enabling a pair of LogStream's out-of-the-box generators:

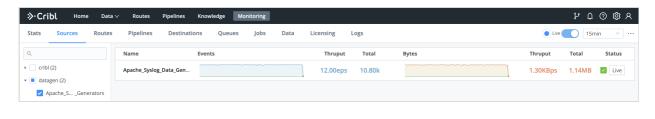
Navigate to **Sources** > **Datagens** and click **Add New**.

Select a Data Generator File (e.g., apache_common.log) and set it at 4 EPS/worker process. Select another Data Generator File (e.g., syslog.log) and set it at 8 EPS/worker process. Hit Save.



Selecting Datagens files and event rates

On the **Monitoring** page, under **Sources**, search for datagen and confirm that the Source is generating data.



Creating a Datagen Template from a Sample File

To convert a sample into a template:

Go to Preview > Paste a Sample, and add a sample like the AWS VPC Flow logs below:

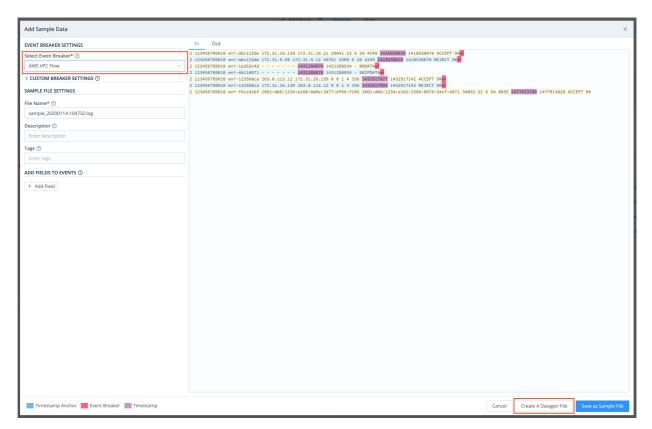
Sample VPC Flow Logs

```
2 123456789010 eni-abc123de 172.31.16.139 172.31.16.21 20641 22 6 20 4249 1418530010 1418530070 2 123456789010 eni-abc123de 172.31.9.69 172.31.9.12 49761 3389 6 20 4249 1418530010 1418530070 F 2 123456789010 eni-1a2b3c4d - - - - - - 1431280876 1431280934 - NODATA 2 123456789010 eni-4b118871 - - - - - - 1431280876 1431280934 - SKIPDATA 2 123456789010 eni-1235b8ca 203.0.113.12 172.31.16.139 0 0 1 4 336 1432917027 1432917142 ACCEPT 2 123456789010 eni-1235b8ca 172.31.16.139 203.0.113.12 0 0 1 4 336 1432917094 1432917142 REJECT 2 123456789010 eni-f41c42bf 2001:db8:1234:a100:8d6e:3477:df66:f105 2001:db8:1234:a102:3304:8879
```

From the **Event Breaker** drop-down, select **AWS VPC Flow** to ensure that:

- The pasted text gets broken properly into individual events (notice the Event Breaker on newlines).
- Timestamps are extracted correctly (text highlighted purple below).

Once you've verified these results, click Create a Datagen File.



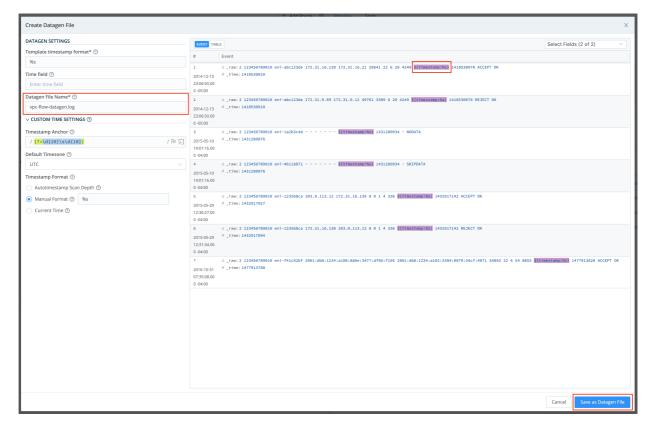
Creating a Datagen template

On the resulting Create Datagen File screen:

- Enter a file name, e.g.: vpc-flow-datagen.log
- Ensure that the timestamp template format is correct: \${timestamp: %s}

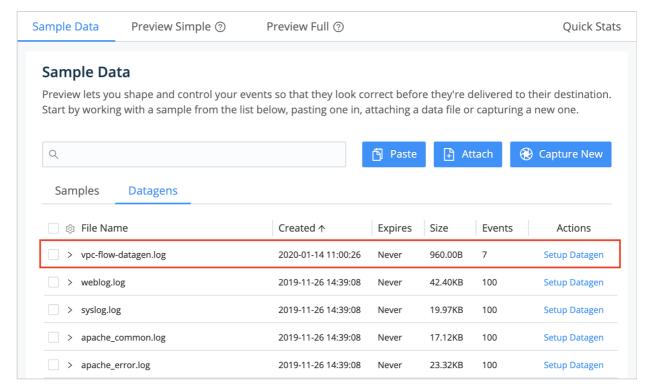
\${timestamp: <format>} is a template that the datagen engine uses to insert the current time – in each newly generated event – using the given format. In this case, %s is the desired strftime format for the timestamp (i.e., the epoch).

Once you've verified these results, click Save as Datagen File.



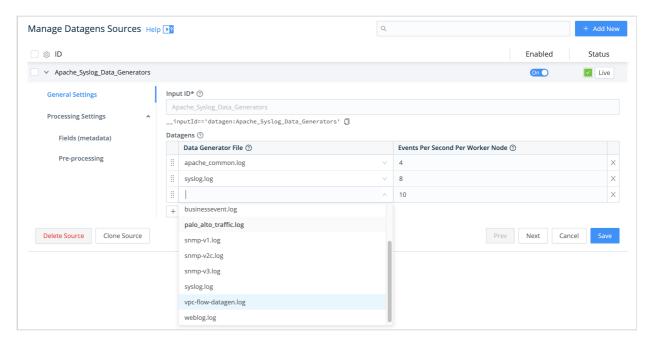
Saving a named Datagen template

To confirm that the Datagen file has been created, check **Preview > Datagens**.



Verifying Datagen file creation

Now, to start using your newly created Datagen file, go back to **Sources > Datagens**. Add it using the drop-down shown below.



Adding new template file to Datagens Source

CLI Reference

Command line interface basics

In addition to starting and stopping the Cribl LogStream server, LogStream's command line interface enables you to initiate many configuration and administrative tasks directly from your terminal.

Command Syntax

To execute CLI commands, the basic syntax is:

```
cd $CRIBL_HOME/bin
./cribl <command> <sub-command> <options> <arguments>
```

Avoiding Surprises

Immediate Execution

As indicated in the sample output below, some commands take effect immediately.

Commands that require further input will echo the sub-commands, options, and arguments they expect.

Persistent Volumes

If you start LogStream with the CRIBL_VOLUME_DIR variable, all subsequent CLI commands should have this variable defined. Otherwise, those commands will apply LogStream's default directories, yielding misleading results.

You can set CRIBL_VOLUME_DIR as an environment variable, or you can explicitly include it in each command, as in this example:

```
CRIBL_VOLUME_DIR /opt/cribl/bin/cribl status.
```

Commands Available

To see a list of available commands, enter ./cribl alone (or the equivalent ./cribl help). To execute a command, or to see its required parameters, enter ./cribl <command> .

help

Displays help (commands list).

```
Cribl LogStream - N.n.n ≺build no.>
Usage: [sub-command] [options] [args]
Commands:
```

- Display help help mode-master - Configure to a master instance

mode-single - Configure to a single instance mode-worker - Configure to a worker instance

reload - Reload Cribl LogStream
restart - Restart Cribl LogStream
start - Start Cribl LogStream
status - Status of Cribl LogStream
stop - Stop Cribl LogStream
version - Print Cribl LogStream version and installation type

auth - Cribl LogStream Auth
boot-start - Enable/Disable Cribl LogStream boot-start
diag - Manage diagnostics bundles
groups - Manage Worker Groups
keys - Manage encryption keys

mode-searchhead - Configure Cribl LogStream to run on a Splunk Search Head - Listen on a port for traffic and output stats and data

node — Execute a JavaScript file pipe — Feed stdin to a pipeline splunk-decrypt - Splunk decrypt search command task - Run Cribl LogStream task vars - Manage global variables

mode-master

Configures Cribl LogStream as a Master instance.

Options

```
[-H <host>] - Host (defaults to 0.0.0.0).
[-p <port>] - Port (defaults to 4200).
[-n <certName>] - Name of saved certificate.
[-k <privKeyPath>] - Server path containing the private key (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use.
   \hbox{[-i < ipWhitelistRegex>]} - \hbox{Regex matching IP addresses that are allowed to establish a connection}
```

Sample Response

```
Settings updated.
You will need to restart LogStream before your changes take full effect.
```

mode-single

Configures Cribl LogStream as a single-instance deployment.

Sample Response

```
Settings updated.
You will need to restart LogStream before your changes take full effect.
```

mode-worker

Configures Cribl LogStream as a Worker instance.

Usage

```
./cribl mode-worker -H <host> -p <port>
```

The -H <host> -p <port> parameters are required.

Options

```
-H <host> - Master Node's Hostname or IP address.
-p <port> - Master Node's cluster communications port (defaults to 4200).

[-n <certName>] - Name of saved certificate.

[-k <privKeyPath>] - Server path containing the private key (in PEM format) to use. Can reference certPath>] - Server path containing certificates (in PEM format) to use. Can reference certPath>] - Authentication token to include as part of the connection header. By default certificates (in PEM format) to use. Can reference certificates (in PE
```

Sample Response

```
Settings updated.
You will need to restart LogStream before your changes take full effect.
```

reload

Reloads Cribl LogStream. Executes immediately.

```
Reload request submitted to Cribl LogStream
```

restart

Restarts Cribl LogStream. Executes immediately.

 \triangle Executing this command cancels any running collection jobs.

```
Stopping Cribl LogStream, process 56572 ......

Cribl LogStream is not running
Starting Cribl LogStream...
```

```
..
Cribl LogStream started with pid 57233
API Server is available at http://192.168.0.100:9000
```

start

Starts Cribl LogStream. Executes immediately. Upon first run, echoes LogStream's default login credentials.

```
Starting Cribl LogStream...
...
Cribl LogStream started with pid 57279
API Server is available at http://192.168.0.100:9000
```

status

Displays status of Cribl LogStream, including the API Server address, instance's mode (Master or Worker), process ID, and GUID (fictitious example below). Executes immediately.

Cribl LogStream Status

Address: http://192.168.0.100:9000

Mode: worker Status: Up

Software Version: 42.0-7f4c190a

Master: localhost:4200

PID: 3859

GUID: 76-ea411263a64b9-e419daee4-ef-dd2e2f

stop

Stops Cribl LogStream. Executes immediately.

♠ Executing this command cancels any running collection jobs.

```
Stopping Cribl LogStream, process 57233 ......
Cribl LogStream is not running
```

version

Displays Cribl LogStream version and installation type. Executes immediately.

Version: 2.2-0####x##

Installation type: standalone

☐ The version command echoes standalone for both single-instance and distributed deployments. This simply confirms that you're running a freestanding Cribl LogStream server, not the Cribl App for Splunk.

auth

Log into or out of Cribl LogStream.

```
Commands:
login - Log in to Cribl LogStream, args:
  [-h <host>] - Host URL (e.g. http://localhost:9000)
  [-u <username>] - Username
  [-p <password>] - Password
  [-f <file>] - File with credentials
logout - Log out from Cribl LogStream
```

Login Examples

Launch interactive login:

```
$CRIBL_HOME/bin/cribl auth login
```

Append credentials as command arguments:

```
$CRIBL_HOME/bin/cribl auth login -h <url> -u <username> -p <password>
```

i All -h and host arguments are optional, provided that the API host and port are listed in the cribl.yml file's api: section

Provide credentials in environment variables:

```
CRIBL_HOST=<url> CRIBL_USERNAME=<username> CRIBL_PASSWORD=<password>
$CRIBL_HOME/bin/cribl auth login
```

Provide credentials in a file:

```
$CRIBL_HOME/bin/cribl auth login -f <path/to/file>
```

--

Corresponding file contents:

```
host=<url>
username=<username>
password=<password>
```

boot-start

Enables or disables Cribl LogStream boot-start.

```
Usage: [sub-command] [options] [args]

Commands:
disable - Disable Cribl LogStream boot-start, args:
    [-m <manager>] - Init manager (systemd|initd)
    [-c <configDir>] - Config directory for the init manager
enable - Enable Cribl LogStream boot-start, args:
    [-m <manager>] - Init manager (systemd|initd)
    [-u <user>] - User to run Cribl LogStream as
    [-c <configDir>] - Config directory for the init manager
```

diag

Manages diagnostic bundles.

```
create - Creates diagnostic bundle for Cribl LogStream

list - List existing Cribl LogStream diagnostic bundles

send - Send LogStream diagnostics bundle to Cribl Support, args:
   -c <caseNumber> - Cribl Support Case Number
[-p <path>] - Diagnostic bundle path (if empty then new bundle will be created)
```

groups

Manages Worker Groups.

```
Usage: [sub-command] [options] [args]

Commands:
commit - Commit, args:
    [-g <group>] - Group ID
    [-m <message>] - Commit message
commit-deploy - Commit & Deploy, args:
    -g <group> - Group ID
    [-m <message>] - Commit message
deploy - Deploy, args:
    -g <group> - Group ID
    [-v <version>] - Deploy version
list - List Worker Groups
```

keys

Manages encryption keys. You must append the -g <group> argument to specify a Worker Group. As a fallback, append the argument -g default, e.g.:

```
./cribl keys list -g default
```

```
Usage: [sub-command] [options] [args]

Commands:
add - Add encryption keys, args:
    -g <group> - Group ID
    [-c <keyclass>] - key class to set for the key
    [-k <kms>] - KMS to use, must be configured, see cribl.yml
    [-e <expires>] - expiration time, epoch time
    [-i ] - use an initialization vector

list - List encryption keys
    -g <group> - Group ID
```

mode-searchhead

Configures Cribl LogStream to run on a Splunk Search Head.

nc

Listens on a port for traffic, and outputs stats and data. (Netcat-like utility.)

node

Executes a JavaScript file. Displays a command prompt for path/filename input, as shown here:

>

pipe

Feeds stdin to a pipeline. Examples:

```
cat sample.log | ./cribl pipe -p <pipelineName>
cat sample.log | ./cribl pipe -p <pipelineName> 2>/dev/null
```

scope

Greps your apps by the syscalls. Executes immediately.

```
splunk-decrypt
```

task

Runs a Cribl LogStream task. Requires definitions for the dir, executor, and path properties.

vars

Manages LogStream Global Variables.

```
Usage: [sub-command] [options] [args]
Commands:
add - Add global variable, args:
  -i <id> - Global variable ID
-t <type> - Type
-v <value> - Value
[-a <args>] - Arguments
  [-d <description>] - Description
 [-c <tags>] - Custom Tags (comma separated list)
[-g <group>] - Group ID
get - List encryption keys, args:
  [-i <id>] - Global variable ID
  [-g <group>] - Group ID
remove - Remove global variable, args:
   -i <id> - Global variable ID
  [-g <group>] - Group ID
update - Update global variable, args:
  -i <id> - Global variable ID

-t <type> - Type

-v <value> - Value

[-a <args>] - Arguments
  [-d <description>] - Description
  [-c <tags>] - Custom Tags (comma separated list)
[-g <group>] - Group ID
```

EXPRESSION REFERENCE

Introduction to Expression Syntax

As data travels through a Cribl LogStream Pipeline, it is operated on by a series of Functions. Functions are fundamentally JavaScript code.

Functions that ship with Cribl LogStream are configurable via a set of inputs. Some of these configuration options are literals, such as field names, and others can be JavaScript expressions.

Expressions are **valid units** of code that resolve to a value. Every syntactically valid expression resolves to some value, but conceptually, there are two types of expressions: those that **assign** value to a variable (a.k.a., with side effects), and those that **evaluate** to a value.

Assigning a value	Evaluating to a value
42	(Math.random() * 42)
x = 42	3 + 4
<pre>newFoo = foo.slice(30)</pre>	'foobar'
	'42'

Filters and Value Expressions

Filters

Filters are used in Routes to select a stream of the data flow, and in Functions to scope or narrow down the applicability of a Function. Filters are expressions that must evaluate to either true (or truthy) or false (or falsy). Keep this in mind when creating Routes or Functions. For example:

```
    sourcetype='access_combined' & host.startsWith('web')
```

```
source.endsWith('.log') ||
sourcetype='aws:cloudwatchlogs:vpcflow'
```

Truthy	Falsy
true	false
42	null
-42	undefined
3.14	0
"foo"	NaN
Infinity	11
-Infinity	11 11

Value Expressions

Value expressions are typically used in Functions to assign a value – for example, to a new field. For example:

```
Math.floor(_time/3600)
```

source.replace(/.{3}/, 'XXX')

Considerations and Best Practices for Creating Predictable Expressions

- In a value expression, ensure that the source variable is not null, undefined, or empty. For example, assume you want to have a field called len, to be assigned the length of a second field called employeeID. But you're not sure if employeeID exists. Instead of employeeID.length, you can use a safer shorthand, such as: (employeeID || '').length.
- If a field does not exist (undefined), and you're doing a comparison with its properties, then the boolean expression will always evaluate to false. For example, if employeeID is undefined, then both of these expressions will evaluate to false: employeeID.length > 10 , and employeeID.length < 10 .
- = means "equal to," while == means "equal value and equal type." For example, 5 = 5 evaluates to true, while 5 == "5" evaluates to false.
- A ternary operator is a very powerful way to create conditional values. For example, if you wanted to assign either minor or adult to a field groupAge, based on the value of age, you could do: (age ≥ 18)?
 'adult': 'minor'.

Expressions Using Fields with Non-Alphanumeric Characters

If there are fields with non-alphanumeric characters – e.g., @timestamp or user-agent or kubernetes.namespace_name – you can access them using __e['<field-name-here>']. (Note the single quotes.) More details here. In any other place where the field is referenced – e.g., in the Eval function's field names – you should use a single-quoted literal, of the form: '<field-name-here>'.

Wildcard Lists

Wildcard Lists are used throughout the product, especially in various Functions, such as Eval, Mask, Publish Metrics, Parser, etc.

Wildcard Lists, as their name implies, accept strings with asterisks (\star) to represent one or more terms. They also accept strings that start with an exclamation mark (!) to **negate** one or more terms.

Wildcard Lists are order-sensitive **only** when negated terms are used. This allows for implementing any combination of whitelists and blacklists.

For example:

Wildcard List	Value	Meaning
List 1	!foobar, foo*	All terms that start with foo , except foobar .
List 2	!foo*, *	All terms, except for those that start with foo .

Cribl Expressions

Native Cribl LogStream function methods can be found under C.*, and can be invoked from any Function that allows for expression evaluations. For example, to create a field that is the SHA1 of a another field's value, you can use the Eval Function with this **Evaluate Fields** pair:

Name	Value Expression
myNewField	C.Mask.sha1(myOtherField)

C.Crypto – Data Encryption and Decryption Functions

```
C.Crypto.decrypt
```

(method) Crypto.decrypt(value: string): string

Decrypt all occurrences of ciphers in the given value. Instances that cannot be decrypted (for any reason) are left intact.

@param - value - string in which to look for ciphers

@returns - value with ciphers decrypted

C.Crypto.encrypt

(method) Crypto.encrypt(value: any, keyclass: number, keyId?:
string, defaultVal?: string): string
Encrypt the given value with the keyId, or with a keyId picked up

Encrypt the given value with the keyId, or with a keyId picked up automatically based on keyclass.

- @param {string | Buffer} value what to encrypt.
- @param keyclass if keyId isn't specified, pick one at the given keyclass.
- @param keyId encryption keyld, takes precedence over keyclass.
- @param defaultVal what to return if encryption fails for any reason; if unspecified, the original value is returned.
- @returns if encryption succeeds, the encrypted value; otherwise, defaultVal if specified; otherwise, value.

C.Decode - Data Decoding Functions

C.Decode.base64

(method) Decode.base64(val: string, resultEnc?: string): any

```
Performs base64 decoding of the given string. Returns a string or Buffer,
depending on the resultEnc value, which defaults to 'utf8'.
@param - val - value to base64-decode
@param - resultEnc - encoding to use to convert the binary data to a string.
Defaults to 'utf8'. Use 'utf8-valid' to validate that result is valid UTF8;
use 'buffer' if you need the binary data in a Buffer.
C.Decode.gzip
(method) Decode.gzip(value: any, encoding?: string): string
Gunzip the supplied value.
@param - value - the value to gunzip.
@param - encoding - encoding of value, for example: 'base64', 'hex',
'utf-8', 'binary'. Default is 'base64'. If data is received as Buffer (from
gzip with encoding: 'none' ), decoding is skipped.
C.Decode.hex
(method) Decode.hex(val: string): number
Performs hex to number conversion. (Returns NaN if value cannot be
converted to a number.)
@param - val - hex string to parse to a number (e.g., "0xcafe").
C.Decode.uri
(method) Decode.uri(val: string): string
Performs URI-decoding of the given string.
@param - val - value to URI-decode.
C.Encode – Data Encoding Functions
C.Encode.base64
(method) Encode.base64(val: any, trimTrailEq?: boolean): string
Returns a base64 representation of the given string or Buffer.
@param - val - value to base64-encode.
@param - trimTrailEq - whether to trim any trailing = .
C.Encode.gzip
(method) Encode.gzip(value: string, encoding?: string): any
```

@param - encoding - encoding of value, for example: 'base64', 'hex',
'utf-8', 'binary', 'none'. Defaultis 'base64'. If 'none' is specified,

Gzip, and optionally base64-encode, the supplied value.

@param - value - the value to gzip.

data will be returned as a Buffer.

```
C.Encode.hex
(method) Encode.hex(val: string | number): string
Rounds the number to an integer and returns its hex representation
(lowercase). If a string is provided, it will be parsed into a number or NaN.
@param - val - value to convert to hex.

C.Encode.uri
(method) Encode.uri(val: string): string
Returns the URI-encoded representation of the given string.
@param - val - value to uri encode.
```

C.env – Environment

```
C.env
(property) env: {[key: string]: string;}
Returns an object containing Cribl LogStream's environment variables.
```

C.Lookup – Inline Lookup Functions

```
C.Lookup - Exact Lookup
(property) Lookup: (file: string, primaryKey?: string,
otherFields?: string[], ignoreCase?: boolean) ⇒ InlineLookup
Returns an instance of a lookup to use inline.
Example invocation:
C.Lookup('lookup_name.csv',
'IP_field_name_in_lookup_file').match(host)
C.LookupCIDR - CIDR Lookup
(property) LookupCIDR: (file: string, primaryKey?: string,
otherFields?: string[]) ⇒ InlineLookup
Returns an instance of a CIDR lookup to use inline.
C.LookupIgnoreCase - Case-insensitive Lookup
(property) LookupIgnoreCase: (file: string, primaryKey?: string,
otherFields?: string[]) ⇒ InlineLookup
Returns an instance of a lookup (ignoring case) to use inline. Works identically
to C.Lookup, except ignores the case of lookup values. (Equivalent to calling
C.Lookup with its fourth ignoreCase? parameter set to true ).
C.[LookupRegex](http://google.com) - Regex Lookup
(property) LookupRegex: (file: string, primaryKey?: string,
```

```
otherFields?: string[]) ⇒ InlineLookup
Returns an instance of a Regex lookup to use inline.
(method) InlineLookup.match(value: string, fieldToReturn?:
string): any
@param - value - the value to look up.
@param - fieldToReturn - name of the lookup file > field to return.
E.g., C.Lookup('lookup-exact.csv', 'foo').match('abc', 'bar')
Return the value of field bar in the lookup table if field foo matches abc.
Example 1: C.LookupCIDR('lookup-cidr.csv',
'foo').match('192.168.1.1', 'bar')
Return the value of field bar in the lookup table if the CIDR range in foo
includes 192.168.1.1.
Example 2: C.LookupCIDR('lookup-cidr.csv', 'cidr').match(hostIP,
'location')
Example 3: C.LookupRegex('lookup-regex.csv',
'foo').match('manchester', 'bar')
Return the value of field bar in the lookup table if the regex in foo matches the
string manchester.
```

C.Mask - Data Masking Functions

string): string

```
C.Mask.CC
(method) Mask.CC(value: string, unmasked?: number, maskChar?:
string): string
Check whether a value could be a valid credit card number, and mask a subset
of the value. By default, all digits except the last 4 will be replaced with X.
@param - value - a string whose digits to mask IFF it could be a valid credit
card number.
@param - unmasked - number of digits to leave unmasked: positive for left,
negative for right, 0 for none.
@param - maskChar - a string/char to replace a digit with.

C.Mask.IMEI
```

value. By default, all digits except the last 4 will be replaced with X.

@param - value - a string whose digits to mask IFF it could be a valid IMEI

Check whether a value could be a vlaid IMEI number, and mask a subset of the

(method) Mask.IMEI(value: string, unmasked?: number, maskChar?:

number.

@param – unmasked – number of digits to leave unmasked: positive for left, negative for right, 0 for none.

@param - maskChar - a string/char to replace a digit with.

C.Mask.isCC

(method) Mask.isCC(value: string): boolean

Checks whether the given value could be a valid credit card number, by computing the string's Lunh's checksum modulo 10 == 0.

@param - value - a string to check for being a valid credit card number.

C.Mask.isIMEI

(method) Mask.isIMEI(value: string): boolean

Checks whether the given value could be a valid IMEI number, by computing the string's Lunh's checksum modulo 10 == 0.

@param - value - a string to check for being a valid IMEI number

C.Mask.luhn

(method) Mask.luhn(value: string, unmasked?: number, maskChar?:
string): string

Check that value Lunh's checksum mod 10 is $\,0\,$, and mask a subset of the value. By default, all digits except the last 4 will be replaced with $\,X\,$. If the value's Lunh's checksum mod 10 is not $\,0\,$, then the value is returned unmodified.

@param - value - a string whose digits to mask IFF the value's Lunh's checksum mod 10 is 0.

@param – unmasked – number of digits to leave unmasked: positive for left, negative for right, 0 for none.

@param - maskChar - a string/char to replace a digit with.

C.Mask.LUHN_SUB

(property) Mask.LUHN_SUB: any

C.Mask.luhnChecksum

(method) Mask.luhnChecksum(value: string, mod?: number): number Generates the Luhn checksum (used to validate certain credit card numbers, IMEIs, etc.). By default, the mod 10 of the checksum is returned. Pass mod = 0 to get the actual checksum.

@param – value – a string whose digits you want to perform the Lunh checksum on.

@param - mod - return checksum modulo this number. If 0, skip modulo. Default is 10.

C.Mask.md5

(method) Mask.md5(value: string, len?: string | number): string
Generate MD5 hash of a given value.

@param - value - compute the hash of this.

@param – len – length of hash to return: 0 for full hash, a +number for left or a -number for right substring. If a string is passed it's length will be used.

C.Mask.random

(method) Mask.random(len?: string | number): string
Generates a random alphanumeric string.

@param – len – a number indicating the length of the result; or, if a string, use its length.

C.Mask.REDACTED

(property) Mask.REDACTED: string

The literal 'REDACTED'.

C.Mask.repeat

(method) Mask.repeat(len?: string | number, char?: string):
string

Generates a repeating char/string pattern, e.g., XXXX.

@param – len – a number indicating the length of the result; or, if a string, use its length.

@param - char - pattern to repeat len times.

C.Mask.sha1

(method) Mask.sha1(value: string, len?: string | number): string
Generate SHA1 hash of given value.

@param - value - compute the hash of this.

@param – len - length of hash to return: 0 for full hash, a +number for left, or a -number for right.

substring. If a string is passed, its length will be used

C.Misc – Miscellaneous Utility Functions

C.Misc.zip()

(method) Misc.zip(keys: string[], values: any[], dest?: any): any Set the given keys to the corresponding values on the given dest object. If dest is not provided, a new object will be constructed.

@param - keys - field names corresponding to keys.

@param - values - values corresponding to values.

```
@param - dest - object on which to set field values.
@returns - object on which the fields were set.

E.g., people = C.Misc.zip(titles, names)

Sample data: titles=['ceo', 'svp', 'vp'], names=['foo', 'bar', 'baz']

Create an object called people, with key names from elements in titles, and with corresponding values from elements in names.

Result: "people": {"ceo": "foo", "svp": "bar", "vp": "baz"}
```

C.Net – Network Functions

```
C.Net.cidrMatch()
(method) Net.cidrMatch(cidrIpRange: string, ipAddress: string):
boolean
Determines whether the supplied IPv4 ipAddress is inside the range of
addresses identified by cidrIpRange . For example: C.Net.cidrMatch
('10.0.0.0/24', '10.0.0.100') returns true.
@param - cidrIpRange - IPv4 address range in CIDR format. E.g.,
10.0.0.0/24.
@param - ipAddress - The IPv4 IP address to test for inclusion in
cidrIpRange.
C.Net.ipv6Normalize()
(method) Net.ipv6Normalize(address: string): string
Normalize an IPV6 address based on RFC draft-ietf-6man-text-addr-
representation-04.
@param - address - the IPV6 address to normalize.
C.Net.isPrivate()
(method) Net.isPrivate(address: string): string
Determine whether the supplied IPv4 address is in the range of private
addresses per RFC1819.
@param - address - address to test.
```

C.os – System Functions

C.confVersion

Returns Cribl LogStream config version.

C.os.hostname()

Returns hostname of the system running this Cribl LogStream instance.

C.Schema - Schema Functions

```
C.Schema()
(property) Schema: (id: string) ⇒ SchemaValidator
(method) SchemaValidator.validate(data: any): boolean
Validates the given object against the schema.
@param - data - object to be validated.
@returns - true when schema is valid; otherwise, false.
Example: C.Schema('schema1').validate(myField) will validate if
myField object conforms to schema1.
See Schema Library for more details.
C.Text - Text Functions
C.Text.entropy()
(method) Text.entropy(bytes: any): number
Computes the Shannon entropy of the given buffer or string.
@param - bytes - value to undergo Shannon entropy computation.
@returns – the entropy value; or –1 in case of an error.
C.Text.hashCode()
(method) Text.hashCode(val: string | Buffer | number): number
Computes hashcode (djb2) of the given value.
@param - val - value to be hashed.
@returns - hashcode value.
C.Text.isASCII()
(method) Text.isASCII(bytes: any): boolean
Checks whether all bytes or chars are in the ASCII printable range.
@param - bytes - value to check for character range.
@returns - true if all chars/bytes are within ASCII printable range; otherwise,
false.
C.Text.isUTF8()
(method) Text.isUTF8(bytes: any): boolean
Checks whether the given Buffer contains valid UTF8.
@param - bytes - bytes to check.
@returns - true if bytes are UTF8; otherwise, false.
```

(method) C.Text.parseWinEvent(xml: string, nonValues?: string[]):

C.Text.parseWinEvent

any

Parses an XML string representing a Windows event into a compact, prettified JSON object. Works like C.Text.parseXml, but with Windows events, produces more-compact output. For a usage example, see Reducing Windows XML Events.

@param - xml - an XML string; or an event field containing the XML.

@param - nonValues - array of string values. Elements whose value equals any of the values in this array will be omitted from the returned object.

Defaults to ['-'], meaning that elements whose value equals - will be discarded.

@returns – an object representing the parsed Windows Event; or undefined if the input could not be parsed.

C.Text.parseXml

(method) C.Text.parseXml(xml:string, keepAttr?:boolean,
keepMetadata?:boolean, nonValues?:string[]): any
Parses an XML string and returns a JSON object. Can be used with Eval
Function to parse XML fields contained in an event, or with ad hoc XML.
@param - xml - XML string, or an event field containing the XML.
@param - keepAttr - whether or not to include attributes in the returned
object. Defaults to true.

@param - keepMetadata - whether or not to include metadata found in the XML. The keepAttr parameter must be set to true for this to work. Defaults to false. (Eligible metadata includes namespace definitions and prefixes, and XML declaration attributes such as encoding, version, etc.)

@param - nonValues - array of string values. Elements whose value equals any of the values in this array will be omitted from the returned object.

Defaults to [] (empty array), meaning discard no elements.

@returns - an object representing the parsed XML; or undefined if the input could not be parsed. An input collection of elements will be parsed into an

C.Text.relativeEntropy()

array of objects.

(method) Text.relativeEntropy(bytes: any, modelName?: string):
number

Computes the relative entropy of the given buffer or string.

@param - bytes - value whose relative entropy to compute.

@param - modelName - Name of the model to test the string with.

@returns – the relative entropy value, or –1 in case of an error.

C.Time - Time Functions

```
(method) Time.adjustTZ(epochTime: number, tzTo: string, tzFrom?:
string): number
Adjust a timestamp from one timezone to another.
@param - epochTime - UNIX epoch time.
@param - tzTo - timezone to adjust to.
@param - tzFrom - optional timezone of the timestamp.
C.Time.clamp()
(method) Time.clamp(date, earliest, latest, defaultDate?): number
Constrains an event's parsed timestamp to realistic earliest and latest
boundaries.
@param - date - Timestamp originally parsed from event, in UNIX epoch
time (ms) or JavaScript Date format.
@param - earliest - earliest allowable timestamp, in UNIX epoch time (ms)
or JS Date format.
@param - latest - latest allowable timestamp, in UNIX epoch time (ms) or JS
Date format.
@param - defaultDate - optional default date, in UNIX epoch time (ms) or
JS Date format, to substitute for values outside the earliest or latest
boundaries.
C.Time.strftime()
(method) Time.strftime(date: number | Date, format: string, utc?:
boolean): string
Format a Date object or number as a time string, using strftime specifier.
@param - date - Date object or number (seconds since epoch) to format.
@param - format - specifier to use to format the date.
@param - utc - whether to output the time in UTC, rather than in local
timezone.
@returns – representation of the given date.
C.Time.strptime()
(method) Time.strptime(str: string, format: string, utc?: boolean,
strict?: boolean): Date
Extract time from a string using strptime specifier.
@param - str - string to parse to a timestamp (see strict flag).
@param - format - strptime specifier.
@param - utc - whether to interpret times as UTC, rather than as local time.
@param - strict - whether to return null if there are any extra characters
after timestamp.
```

C.Time.adjustTZ()

@returns – a parsed Date object, if successful; otherwise, null if the specifier did not match.

C.Time.timestampFinder()

(method) Time.timestampFinder(utc?: boolean).find(<source-field>):
AutoTimeParser

Extract time from the specified <source-field>, using the same algorithm as the Auto Timestamp Function and the Event Breaker Function.

- @param utc whether to output the time in UTC, rather than in local timezone.
- @param <source-field> the field in which to search for the time.
- @returns representation of the extracted time.

C.vars - Global Variables

See Global Variables Library for more details.

C.version – Cribl LogStream Version

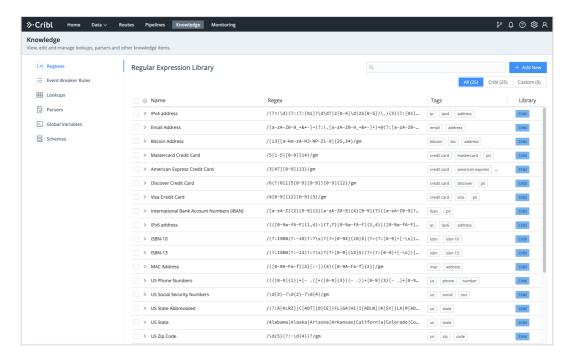
(property) version: string Cribl LogStream Version.

KNOWLEDGE

Regex Library

What Is the Regex Library

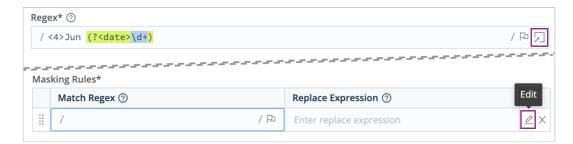
Cribl LogStream ships with a Regex Library that contains a set of pre-built common regex patterns. This library serves as an easily accessible repository of regular expressions. The Library is searchable, and you can assign tags to each pattern for further organization or categorization. The Library is located under **Knowledge** > **Regex Library**.



Regular Expression Library

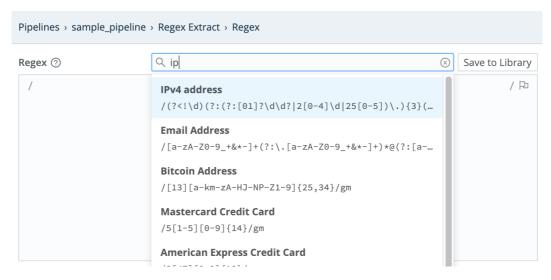
Using Library Patterns

As of this version, the Library contains 25 patterns shipped by Cribl LogStream. To insert a pattern into a Function's regex field, first click the pop-out or Edit icon beside that field.



Opening a Regex modal

In the resulting Regex or Rules modal, Regex Library patterns will appear as typeahead options. Click a pattern to paste it in. You can then use the pattern as-is, or modify it as necessary.



Inserting a pattern from the Regex Library

Adding Patterns to the Library

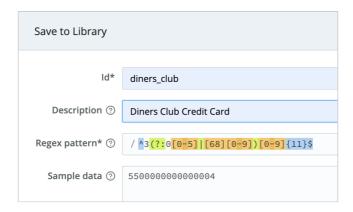
You can also add new, custom patterns to the Library. In the same modal, once you've built your pattern, click the **Save to Library** button.



Adding a custom pattern to the Regex Library from a Function's Regex modal

In the resulting modal, give your custom pattern a unique ID. Optionally, you can also provide a **Description** (name) and groom the **Sample data**. Then click

Save.



Identifying the custom pattern

Your custom pattern will now reside in the Regex Library. It will be available to Functions using the same typeahead assist as Cribl's pre-built patterns.

Cribl vs. Custom and Priority

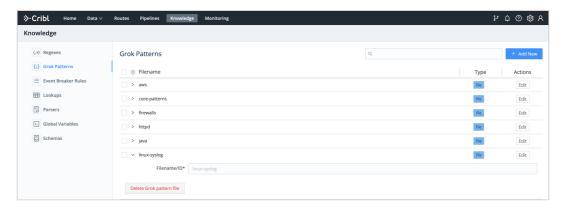
Within the Library, patterns shipped by Cribl will be listed under the **Cribl** tab, while those built by users will be found under **Custom**. Over time, Cribl LogStream will ship more patterns, and this distinction allows for both sets to grow independently.

In the case of an ID/Name conflict, the Custom pattern takes priority in listings and search. For example, if a Cribl-provided pattern and a Custom one are both named <code>ipv4</code>, the one from Cribl will not be displayed or delivered as a search result.

Grok Patterns Library

What Is the Grok Patterns Library

Cribl LogStream ships with a Grok Patterns Library that contains a set of prebuilt common patterns, organized as files.



Grok Patterns Library

Managing Library Patterns

You can access the Grok Patterns Library in the UI by selecting **Knowledge** > **Grok Patterns**. The library contains several pattern files that Cribl provides for basic Grok scenarios, and is searchable.

To edit a pattern file, click **Edit** in its **Actions** column.

To create a new pattern file, click + Add New. In the resulting Create Grok Patterns modal, assign a unique Filename/ID, populate the file with patterns, then click Save.



Adding Grok patterns

i Pattern files reside in: \$CRIBL_HOME/(default|local)/cribl/grok-patterns/

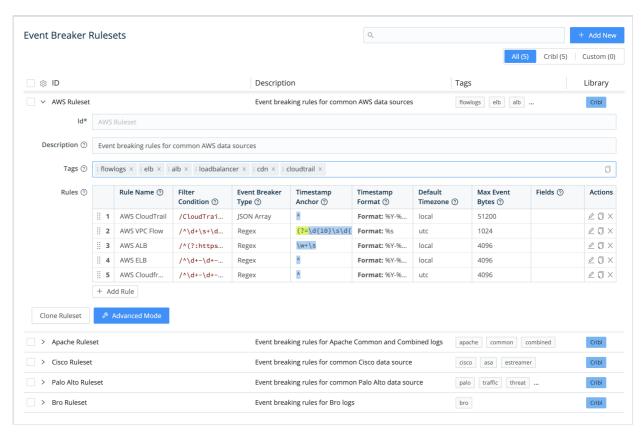
Using Grok Patterns

In the current LogStream version, you apply Grok patterns by inserting a Grok Function into a Pipeline, then manually typing or pasting patterns into the Pattern field(s).

Event Breakers

What Are Event Breakers

Event Breakers help break incoming streams of data into discrete events. You access the Event Breakers management interface under **Knowledge** > **Event Breakers**. On the resulting **Event Breaker Rulesets** page, you can edit, add, delete, search, and tag Event Breaker rules and rulesets, as necessary.



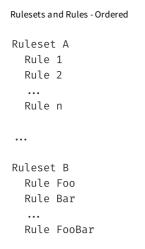
Event Breaker Rulesets page

Event Breaker Rulesets

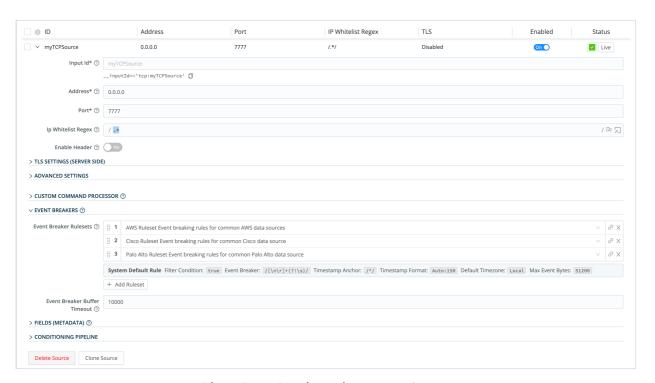
Rulesets are **collections of Event Breaker rules** that are associated with **Sources**. Rules define configurations needed to break down a stream of data into events.

i Event Breakers are accessible only on Sources that require incoming events to be broken into a better-defined format. Check individual LogStream Sources' documentation for Event Breaker support.

Rules within a ruleset are ordered and evaluated top->down. One or more rulesets can be associated with a Source, and these rulesets are also evaluated top->down. For a stream from a given Source, the first matching rule goes into effect.



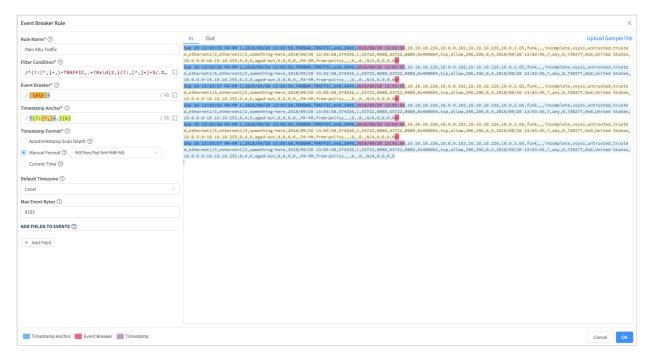
An example of multiple rulesets associated with a Source:



Three Event Breaker rulesets on a Source

Rule Example

This rule breaks on newlines and uses Manual timestamping **after** the sixth comma, as indicated by this pattern: $^(?:[^,]*,){6}$.



An Event Breaker rule

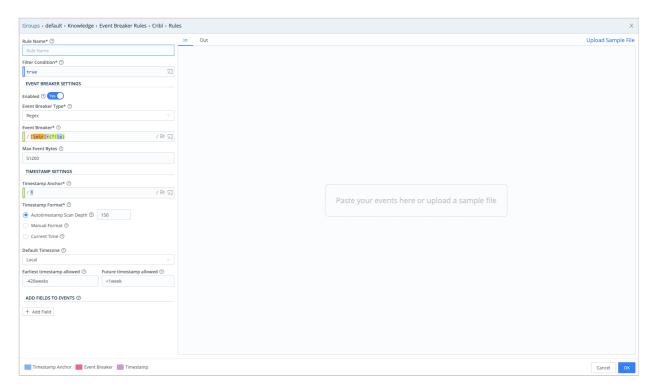
System Default Rule

The system default rule functionally sits at the bottom of the ruleset/rule hierarchy (but is built-in and not displayed on the Event Breakers page), and goes into effect if there are no matching rules:

- Filter Condition defaults to true
- Event Breaker to [\n\r]+(?!\s)
- Timestamp anchor to ^
- Timestamp format to Auto and a scan depth of 150 bytes
- Max Event Bytes to 51200
- Default Timezone to Local

How Do Event Breakers Work

On the **Event Breaker Rulesets** page (see screenshot above), click + **Add New** to create a new Event Breaker ruleset. Click + **Add Rule** within a ruleset to add a new Event Breaker.



Adding a new Event Breaker rule

Each Event Breaker includes the following components, which you configure from top to bottom in the above **Event Breaker Rule** modal:

Filter Condition

As a stream of data moves into the engine, a rule's filter expression is applied. If the expression evaluates to true, the rule configurations are engaged for the entire duration of that stream. Else, the next rule down the line is evaluated.

Event Breaker Type

After a breaker pattern has been selected, it will apply on the stream **continuously**. See below for specific information on different Event Breaker Types.

Timestamp Settings

After events are synthesized out of streams, LogStream will attempt timestamping. First, a timestamp anchor will be located inside the event. Next, starting there, the engine will try to do one of the following:

- Scan up to a configurable depth into the event and autotimestamp, or
- Timestamp using a manually supplied strptime format, or
- Timestamp the event with the current time.

The closer an anchor is to the timestamp pattern, the better the performance and accuracy – especially if multiple timestamps exist within an event. For the manually supplied option, the anchor

must lead the engine right before the timestamp pattern begins.



Anchors preceding timestamps

☐ This timestamping executes the same basic algorithm as the Auto Timestamp Function and the C.Time.timestampFinder() native method.

Add Fields to Events

After events have been timestamped, one or more fields can be added here as key-value pairs. In each field's **Value Expression**, you can fully evaluate the field value using JavaScript expressions.

Event Breaker Types

Several types of Event Breaker can applied to incoming data streams:

- Regex
- File Header
- JSON Array
- JSON New Line Delimited
- Timestamp
- CSV

Regex

The Regex breaker uses regular expressions to find breaking points in data streams.

After a breaker regex pattern has been selected, it will apply on the stream **continuously**. Breaking will occur at the beginning of the match, and the matched content will be consumed/thrown away.

If necessary, you can use a positive lookahead regex to keep the content – e.g.: (?=pattern)

Capturing groups are **not allowed** to be used anywhere in the Event Breaker pattern, as they will further break the stream – which is often undesirable. Breaking will also occur if **Max Event Bytes** has been reached.

Example

Break after a newline or carriage return, but only if followed by a timestamp pattern:

Event Breaker: $[\n\r]+(?=\d+-\d+\s\d+:\d+:\d+)$

```
Sample Event - Multiline
--- input ---
2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can program the aux:
   Try to connect the FTP sensor, maybe it will connect the digital bus!
   Try to navigate the AGP panel, maybe it will quantify the mobile alarm!
2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can program the aux:
   Try to connect the FTP sensor, maybe it will connect the digital bus!
   Try to navigate the AGP panel, maybe it will quantify the mobile alarm!
--- output event 1 ---
  "_raw": "2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can prog
  "_time": 1589920332
}
--- output event 2 ---
  "_raw": "2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can prog
   _time": 1589920332
}
```

File Header

You can use the File Header breaker to break files with headers, such as IIS or Bro logs. This type of breaker relies on a header section that lists field names. The header section is typically present at the top of the file, and can be single-line or greater.

After the file has been broken into events, fields will also be extracted, as follows:

- Header Line: Regex matching a file header line. For example, ^# .
- Field Delimiter: Field delimiter regex. For example, \s+.
- **Field Regex**: Regex with one capturing group, capturing all the fields to be broken by field delimiter. For example, ^#[Ff]ields[:]?\s+(.*)
- Null Values: Representation of a null value. Null fields are not added to events.
- Clean Fields: Whether to clean up field names by replacing non <code>[a-zA-Z0-9]</code> characters with

Example

Using the values above, let's see how this sample file breaks up:

```
--- output event 1 ---
{
 " raw": "1331904608.080000
                                 - 192.168.204.59 137 192.168.204.255 137
                                                                                      udp",
 "ts": "1331904608.080000",
 "id_orig_h": "192.168.204.59",
 "id_orig_p": "137",
 "id_resp_h": "192.168.204.255",
 "id_resp_p": "137",
  "proto": "udp",
  time": 1331904608.08
--- output event 2 ---
 "_raw": "1331904609.190000
                                - 192.168.202.83 48516 192.168.207.4 53
                                                                                       udp",
 "ts": "1331904609.190000",
 "id_orig_h": "192.168.202.83",
 "id_orig_p": "48516",
 "id_resp_h": "192.168.207.4",
 "id_resp_p": "53",
 "proto": "udp",
 "_time": 1331904609.19
```

JSON Array

You can use the JSON Array to extract events from an array in a JSON document (e.g., an Amazon CloudTrail file).

- Array Field: Optional path to array in a JSON event with records to extract. For example,
- Timestamp Field: Optional path to timestamp field in extracted events. For example, eventTime or level1.level2.eventTime.
- **JSON Extract Fields**: Enable this slider to auto-extract fields from JSON events. If disabled, only _raw and time will be defined on extracted events.
- Timestamp Format: If JSON Extract Fields is set to No, you must set this to Autotimestamp or Current Time. If JSON Extract Fields is set to Yes, you can select any option here.

Example

Using the values above, let's see how this sample file breaks up:

```
--- output event 2 --- {
    "_raw": "{\"eventVersion\":\"1.05\",\"eventTime\":\"2020-04-08T01:35:56Z\",\"eventSource\":\"6
    "_time": 1586309756,
    "cribl_breaker": "j-array"
}
```

JSON New Line Delimited

You can use the JSON New Line Delimited breaker to break and extract fields in newline-delimited JSON streams.

Example

Using default values, let's see how this sample stream breaks up:

```
{\sf Sample \ Event \ - Newline \ Delimited \ JSON \ Breaker}
```

```
--- input ---
{"time":"2020-05-25T18:00:54.201Z","cid":"w1","channel":"clustercomm","level":"info","message":'
 \{ "time": "2020-05-25T18:00:54.246Z", "cid": "w0", "channel": "clustercomm", "level": "info", "message": "compared to the compared to the c
 --- output event 1 ---
                "\_raw": "{\ ''time\ '':\ ''2020-05-25T18:00:54.201Z\ '',\ ''cid\ '':\ ''w1\ '',\ ''channel\ '':\ ''clustercomm\ '',\ ''lustercomm'',\ ''lust
              "time": "2020-05-25T18:00:54.201Z",
                "cid": "w1",
                "channel": "clustercomm",
                 "level": "info",
                 "message": "metric sender",
                "total": 720,
                "dropped": 0,
                    "_time": 1590429654.201,
}
 --- output event 2 ---
              "_raw": "{\"time\":\"2020-05-25T18:00:54.246Z\",\"cid\":\"w0\",\"channel\":\"clustercomm\",\"left and the control of the co
                "time": "2020-05-25T18:00:54.246Z",
                "cid": "w0",
                "channel": "clustercomm",
                "level": "info",
                "message": "metric sender",
                "total": 720,
                "dropped": 0,
                  "_time": 1590429654.246,
```

Timestamp

You can use the Timestamp breaker to break events at the beginning of any line in which LogStream finds a timestamp. This type enables breaking on lines whose timestamp pattern is not known ahead of time.

Example

Using default values, let's see how this sample stream breaks up:

Sample Event - Timestamp Based Breaker --- input --{"level": "debug", "ts": "2021-02-02T10:38:46.365Z", "caller": "sdk/sync.go:42", "msg": "Handle ENIConfevel": "debug", "ts": "2021-02-02T10:38:56.365Z", "caller": "sdk/sync.go:42", "msg": "Handle ENIConfevel": "sdk/sync.go:42", "msg": "Handle ENIConfevel": "sdk/sync.go:42", "msg": "Handle ENIConfevel": "sdk/sync.go:42", "caller": \"sdk/sync.go:42" --- output event 1 --{ "_raw": "{\"level\":\"debug\",\"ts\":\"2021-02-02T10:38:46.365Z\",\"caller\":\"sdk/sync.go:42\" "_raw": "{\"level\":\"debug\",\"ts\":\"2021-02-02T10:38:56.365Z\",\"caller\":\"sdk/sync.go:42\" "_time": 1612262336.365 }

CSV

The CSV breaker extracts fields in CSV streams that include a header line. Selecting this type exposes these extra fields:

- Delimiter: Delimiter character to use to split values. Defaults to: ,
- Quote Char: Character used to quote literal values. Defaults to: "
- Escape Char: Character used to escape the quote character in field values. Defaults to: "

Example: Using default values, let's see how this sample stream breaks up:

Example

Using default values, let's see how this sample stream breaks up:

```
Sample Event - CSV Breaker
--- input ---
time,host,source,model,serial,bytes_in,bytes_out,cpu
1611768713, "myHost1", "anet", "cisco", "ASN4204269", 11430, 43322, 0.78
1611768714, "myHost2", "anet", "cisco", "ASN420423", 345062, 143433, 0.28
--- output event 1 ---
 "time": "1611768713",
  "host": "myHost1",
  "source": "anet",
  "model": "cisco",
  "serial": "ASN4204269",
 "bytes_in": "11430",
 "bytes_out": "43322",
  "cpu": "0.78",
  "_time": 1611768713
```

```
}
--- output event 2 ---
{
    "_raw": "\"1611768714\",\"myHost2\",\"anet\",\"cisco\",\"ASN420423\",\"345062\",\"143433\",\"(
    "time": "1611768714",
    "host": "myHost2",
    "source": "anet",
    "model": "cisco",
    "serial": "ASN420423",
    "bytes_in": "345062",
    "bytes_out": "143433",
    "cpu": "0.28",
    "_time": 1611768714
}
```

With **Type: CSV** selected, an Event Breaker will properly add quotes around all values, regardless of their initial state.

Cribl versus Custom Rulesets

Event Breaker rulesets shipped by Cribl will be listed under the **Cribl** tag, while user-built rulesets will be listed under **Custom**. Over time, Cribl will ship more patterns, so this distinction allows for both sets to grow independently. In the case of an ID/Name conflict, the Custom pattern takes priority in listings and search.

Lookups Library

What Are Lookups

Lookups are data tables that can be used in Cribl LogStream to enrich events as they're processed by the Lookup Function. You can access the Lookups library under **Knowledge** > **Lookups**, which provides a management interface for all lookups.

This library is searchable, and each lookup can be tagged as necessary. There's full support for .csv files. Compressed files are supported, but must be in gzip format (.gz extension). You can add files in multimedia database (.mmdb) binary format, but you cannot edit these binary files through LogStream's UI.



Lookups Library

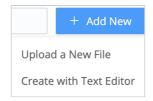
How Does the Library Work

In single-instance deployments, all files handled by the interface are stored in \$CRIBL_HOME/data/lookups . In distributed deployments, the storage path on the Master Node is \$CRIBL_HOME/groups/<groupname>/data/lookups/for each Worker Group.

i For large and/or frequently replicated lookup files, you might want to bypass the Lookups Library UI and instead manually place the files in a different location. This can both reduce deploy traffic and prevent errors with LogStream's default Git integration. For details, see Managing Large Lookups.

Adding Lookup Files

To upload or create a new lookup file, click + Add New, then click the appropriate option from the drop-down.

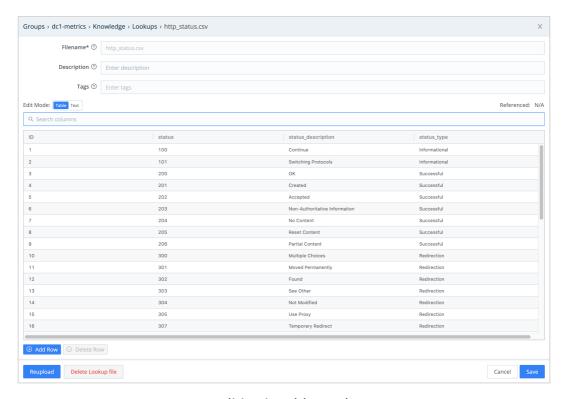


Adding a lookup file

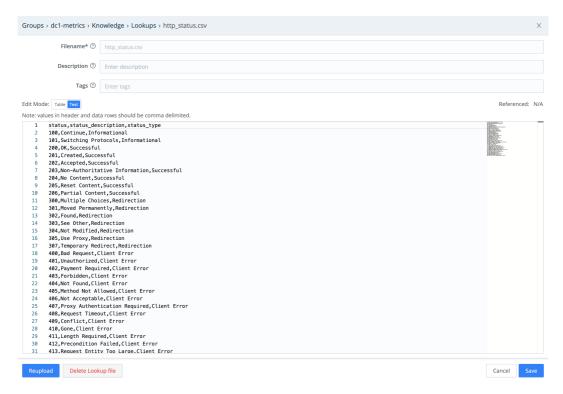
Modifying Lookup Files

To re-upload, expand, edit, or delete an existing .csv or .gz lookup file, click its row on the Lookups page. (No editing option is available for .mmdb files; you can only re-upload or delete these.)

In the resulting modal, you can edit files in **Table** or **Text** mode. However, **Text** mode is disabled for files larger than 1 MB.



Editing in table mode



Editing in text mode

Memory Sizing for Large Lookups

For large lookup files, you'll need to provide extra memory beyond basic requirements for LogStream and the OS. To determine how much extra memory to add to a Worker Node for a lookup, use this formula:

```
Lookup file's uncompressed size (MB) * 2.25 (to 2.75) * numWorkerProcesses = Extra RAM required for lookup
```

For example, if you have a lookup file that is 1 GB (1,000 MB) on disk, and three Worker Processes, you could use an average 2.50 as the multiplier:

```
1,000 * 2.50 * 3 = 7,500
```

In this case, the Node's server or VM would need an extra 7,500 MB (7.5 GB) to accommodate the lookup file across all three worker processes.

What's with That Multiplier?

We've cited a squishy range of 2.25–2.75 for the multiplier, because we've found that it varies inversely with the number of columns in the lookup file:

• The fewer columns, the higher the extra memory overhead (2.75 multiplier).

• The more columns, the lower the overhead (2.25 multiplier).

In Cribl's testing:

- 5 columns required a multiplier of 2.75
- 10 columns required a multiplier of only 2.25.

These are general (not exact) guidelines, and this multiplier depends only on the lookup table's number of columns. The memory overhead imposed by each additional row appears to decline when more columns are present in the data.

Other Scenarios

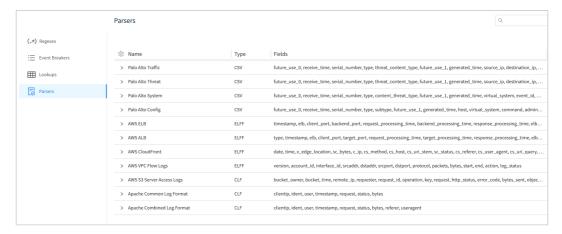
See also:

- Lookup Function.
- Ingest-time Lookups use case.
- Managing Large Lookups use case.
- Redis Function for faster lookups using a Redis integration (bypasses the Lookups Library).

Parsers Library

What Are Parsers

Parsers are definitions and configurations for the Parser Function. You can find the library under **Knowledge** > **Parsers**, and its purpose is to provide an interface for creating and editing Parsers. The library is searchable, and each parser can be tagged as necessary.



Parsers Library

Parsers can be used to **extract** or **reserialize** events. See **Parser Function** page for examples.

Supported Parser Types:

- CSV Parse and reserialize comma-separated values.
- ELFF Parse and reserialize events in Extended Log File Format.
- CLF Parse and reserialize events in Common Log Format.

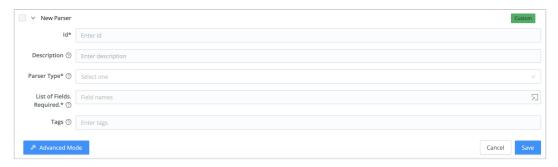
Creating a Parser

To create a parser, follow these steps:

- 1. Go to Knowledge > Parsers and click Add New.
- 2. Enter a unique ID.
- 3. Optionally, enter a **Description**.

- 4. Select a **Parser type** (see the supported types above).
- 5. Enter the **List of fields** expected to be extracted, in order.

 Click this field's Maximize icon (far right) if you'd like to open a modal where you can work with sample data and iterate on results.
- 6. Optionally, enter any desired Tags.



Adding a new parser

Schema Library

What Are Schemas

Schemas are JSON definitions that are used to validate of JSON events. They're based on the popular JSON Schema standard, and all schemas matching draft version 2019-09 are supported. You can find the library under **Knowledge** > **Schemas**. Its purpose is to provide an interface for creating, editing, and maintaining Schemas.

You validate a schema using the C.Schema('<schema name>').validate(<object field>) built-in method. This function can be called anywhere in Cribl LogStream that JavaScript expressions are supported.

Typical use cases for Schema validation:

- Making a decision before sending an event down to a destination.
- Making a decision before accepting an event. (E.g., drop an event if invalid.)
- Making a decision to route an event based on the result of validation.

Example

To add this example JSON Schema, go to **Knowledge** > **Schemas** and click **Add New**. Enter the following:

- ID: schema1.
- Description: (Enter your own description here.)
- Schema: Paste the following schema.

JSON Schema-Sample

{
 "\$id": "https://example.com/person.schema.json",
 "\$schema": "http://json-schema.org/draft-07/schema#",
 "title": "Person",
 "type": "object",
 "required": ["firstName", "lastName", "age"],
 "properties": {
 "firstName": {
 "type": "string",
 "description": "The person's first name."
 },
 "lastName": {
 "type": "string",
 "description": "The person's first name."
 },
 "lastName": {
 "type": "string",
 }
}

```
"description": "The person's last name."
},
"age": {
   "description": "Age in years which must be equal to or greater than zero.",
   "type": "integer",
   "minimum": 0,
   "maximum": 42
}
}
}
```

Assume that events look like this:

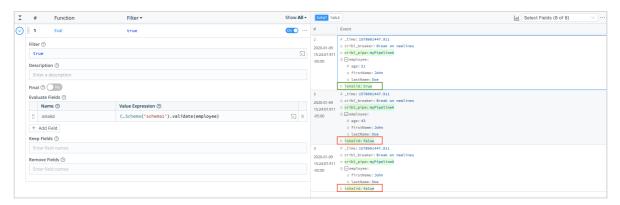
```
Events
{"employee":{"firstName": "John", "lastName": "Doe", "age": 21}}
{"employee":{"firstName": "John", "lastName": "Doe", "age": 43}}
{"employee":{"firstName": "John", "lastName": "Doe"}}
```

To validate whether the employee field is valid per schema1, we can use the following:

```
C.Schema('schema1').validate(employee)
```

Results:

- First event is valid.
- Second event is **not valid** because age is greater than the maximum of 42.
- Third event is **not valid** because age is missing.



Schema validation results for the above events

Global Variables Library

What Are Global Variables

Global Variables are reusable JavaScript expressions that can be accessed in Functions in any Pipeline. You can access the library under Knowledge > Global Variables.

Typical use cases for Global Variables include:

- Storing a constant that you can reference from any Function in any Pipeline.
- Storing a relatively long value expression, or one that uses one or more arguments.

Global Variables can be of the following types:

- Number
- String
- Boolean
- Object
- Array
- Expression

Global Variables can be accessed via C.vars. – which can be called anywhere in Cribl LogStream that JS expressions are supported. Typeahead is provided. More on Cribl Expressions here.

Examples

Scenario 1:

Assign field foo the value in the Answer Global Variable.

- Global Variable Name: theAnswer <-- ships with Cribl LogStream by default.
- Global Variable Value: 42
- Sample Eval Function: foo = C.vars.theAnswer

Scenario 2:

Assign field nowEpoch the current time, in epoch format.

- Global Variable Name: epoch <-- ships with Cribl LogStream by default.
- Global Variable Value: Date.now()/1000
- **Sample Eval Function:** nowEpoch = C.vars.epoch()

Scenario 3:

Create a new field called storage, by converting the value of event field size to human-readable format.

- Global Variable Name: convertBytes <-- ships with Cribl LogStream by default
- Global Variable Value: `\${Math.round(bytes / Math.pow(1024, (Math.floor(Math.log(bytes) / Math.log(1024)))),
 2)}\${['Bytes', 'KiB', 'MiB', 'GiB', 'TiB', 'PiB', 'EiB', 'ZiB', 'YiB'][(Math.floor(Math.log(bytes) / Math.log(1024)))]}`
- Global Variable Argument: bytes
- **Sample Eval Function:** storage = C.vars.convertBytes(size)

Note the use of bytes here as an argument.

USE CASES

Ingest-time Fields

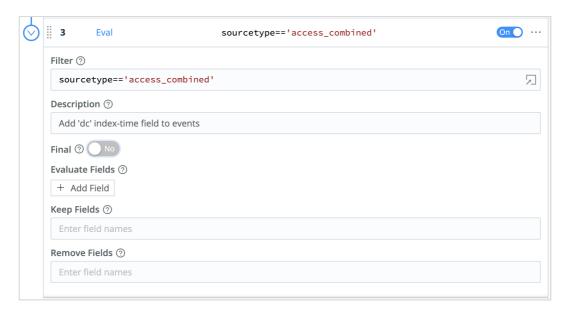
Adding Fields to Data in Motion

To add new fields to any event, we use the out-of-the-box **Eval** Function. We can either apply a Filter to select the events, or we can use the default true Filter expression to apply the Function to all incoming events.

Adding Fields Example

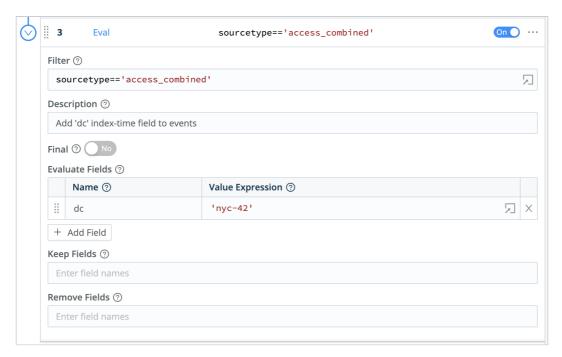
Let's see how we add dc::nyc-42 to all events with sourcetype='access_combined':

- First make sure you have a Route and Pipeline configured to match desired events.
- Next, let's add a **Eval** function to it:



Defining the Eval Function's filter expression

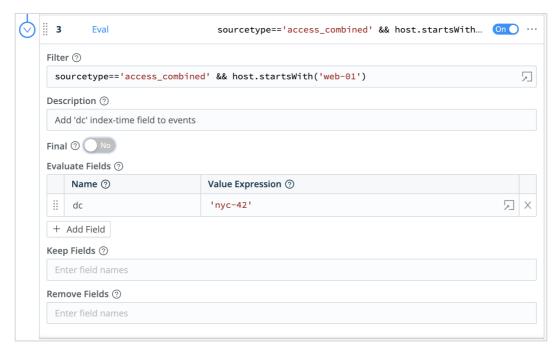
• Next, let's click on + Add Field, add our dc field, and click Save.



Adding the dc field

To confirm, verify that this search returns results: sourcetype="access_combined" dc::nyc-42

• You can add more conditions to the filter, if you'd like. For example, to limit the field to only events from hosts that start with web-01, we can change the filter input as below:



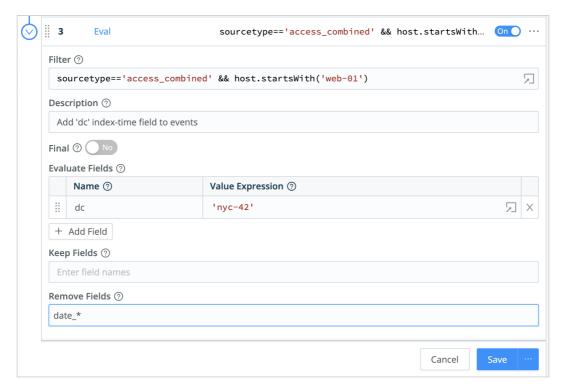
Refining the filter

This is a **very** powerful method to change incoming events in real time. In addition to providing the right context at the right time, users can further benefit substantially by using tstats for **faster** analytics.

Removing Fields

You can remove fields by listing and/or wildcarding field names. Let's see how we can remove all fields that start with date_ ::

- First, make sure you have a Route and Pipeline configured to match desired events.
- Next, let's add a **Eval** function to it (as above).
- Next, in **Remove Fields**, add date_* and hit Save.



Goodbye date_ field

To confirm, verify that this search: sourcetype="access_combined" date_minute=* will soon stop returning results. Enjoy a more efficient Splunk!

Ingest-time Lookups

Enriching Data in Motion

To enrich events with new fields from external sources (such as .csv files), we use LogStream's out-of-the-box Lookup Function. Ingestion-time lookups are not only great for normalizing field names and values, but also ideal for use cases where:

- Fast access via the looked-up value is required. For example, when you don't have a datacenter field in your events, but you do have a host-to-datacenter map, and you need to search by datacenter.
- Looked-up information must be temporally correct. For example, assume that you have a highly dynamic infrastructure, and you need to resolve a resource name (e.g., a container name) to its address. You can't afford to defer this to search time/runtime, as the resource and its records might no longer exist.
 - i To use large binary databases (like GeoIP .mmdb files) for LogStream lookups, see
 Managing Large Lookups. To achieve faster lookups, use LogStream's Redis Function.

Working with Lookups - Example 1

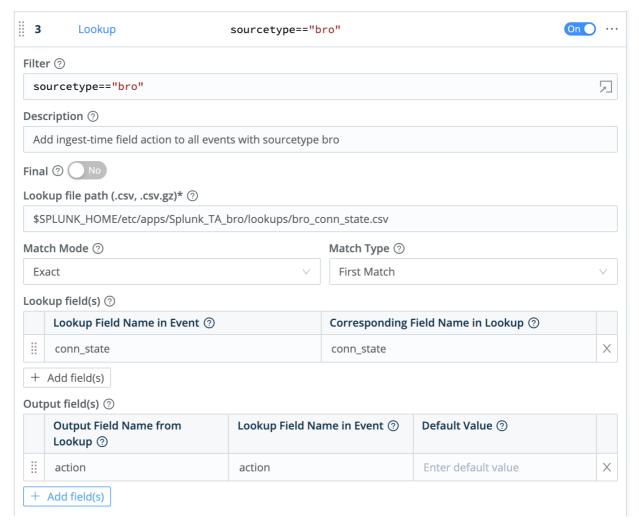
Let's assume we have the following lookup file. Given the field <code>conn_state</code> in an event, we would like to add a corresponding ingestion-time field called <code>action</code>.

```
action, "conn_state", "conn_state_meaning"
dropped,S0, "Connection attempt seen, no reply."
allowed,S1, "Connection established, not terminated."
allowed,SF, "Normal establishment and termination."
blocked,REJ, "Connection attempt rejected."
allowed,S2, "Connection established and close attempt by originator seen (but no reply from respondence), "Connection established and close attempt by responder seen (but no reply from originallowed,RSTO, "Connection established, originator aborted (sent a RST)."
allowed,RSTR, "Established, responder aborted."
dropped,RSTOSO, "Originator sent a SYN followed by a RST, we never saw a SYN-ACK from the respondence, "Conped, SHR, "Responder sent a SYN ACK followed by a FIN, we never saw a SYN ACK from the respondence, SHR, "Responder sent a SYN followed by a FIN, we never saw a SYN from the originator allowed,OTH, "No SYN seen, just midstream traffic (a 'partial connection' that was not later close
```

First, make sure you have a Route and Pipeline configured to match desired events.

Next, let's add a **Lookup** function to the Pipeline, with these settings:

- Lookup file path: \$SPLUNK_HOME/etc/apps/Splunk_TA_bro/lookups/bro_conn_state.csv (note that Environment variables are allowed in the path).
- Lookup Field Name in Event set to conn_state .
- Corresponding Field Name in Lookup set to conn_state.
- Output Field Name from Lookup set to action .
- Lookup Field Name in Event set to action.



Lookup Function to add action field

To confirm success, verify that this search returns expected results: sourcetype="bro" action::allowed.Change the action value as necessary.

Working with Lookups – Example 2

Let's assume we have the following lookup file, and given **both** the fields impact and priority in an event, we would like to add a corresponding ingestion-time field called severity.

```
cisco_sourcefire_severity.csv
impact,priority,severity
1,high,critical
2,high,critical
```

```
3, high, high
4, high, high
0,high,high
"*",high,high
"*", medium, medium
1,low,medium
2,low,medium
3,low,low
4,low,low
0,low,low
"*",low,low
1, none, low
2, none, low
3, none, informational
4, none, informational
0, none, informational
"*", none, informational
```

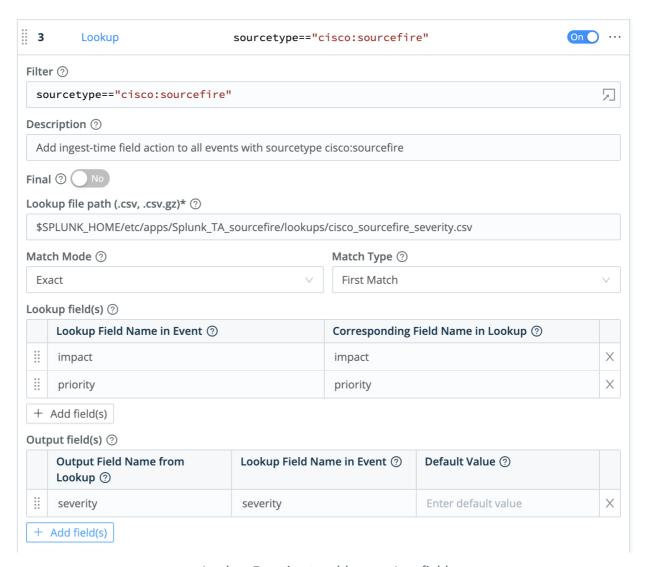
First, make sure you have a Route and Pipeline configured to match desired events.

Next, let's add a **Lookup** function to the Pipeline, with these settings:

• Lookup file path:

\$SPLUNK_HOME/etc/apps/Splunk_TA_sourcefire/lookups/cisco_sourcefire_severity.csv (note that Environment variables are allowed in the path).

- Lookup Field Name(s) in Event set to impact and priority.
- Corresponding Field Name(s) in Lookup setto impact and priority.
- Output Field Name from Lookup set to severity.
- Lookup Field Name in Event set to severity.



Lookup Function to add severity field

To confirm success, verify that this search returns expected results:

sourcetype="cisco:sourcefire" severity::medium.Changethe severity value as necessary.

Sampling

Sampling at Ingest-Time

Let's say that you wanted to analyze and troubleshoot with **highly verbose/voluminous** data – for example, CDN logs, ELB Access Logs, or VPC Flows – but you were concerned about storage requirements and search performance. With Sampling, you can bring in enough samples that your analysis remains statistically significant, and also do all the necessary troubleshooting.

See the example below, or see more details in Access Logs and Firewall Logs.

Sampling Example

Let's use the out-of-the-box **Sampling** function to sample all events from sourcetype='access_combined' where status is '200'. We'll sample these at 5:1 (and all other events at 1:1). This should lower the volume of all verbose successes (200 s), but still bring in **all** potentially erroneous events (400 s, 500`s, etc.) that can be used for troubleshooting.

- First, make sure you have a Route and Pipeline configured to match desired events.
- Next, let's add a Regex Extract Function to extract the status field from
 _raw , and let's call the resulting field __status . Remember, fields that
 start with __ are special fields in Cribl LogStream, and can be used
 anywhere in a Pipeline.



Extracting the __status field

Next, let's add a **Sampling** function, and scope it to all events where sourcetype='access_combined'. Let's apply a filter condition of __status = 200, and a Sample Rate of 5.



Sampling success responses

To confirm that sampling works, compare the event count of all 200 s before and after.

i Each time an event goes through the **Sampling** function, an indextime sampled::<rate> field is added to it. You can use this field in your statistical functions, as necessary.

Access Logs: Apache, ELB, CDN, S3, etc.

Recipe for Sampling Access Logs

Access logs are extremely common. They're often emitted by web servers or similar/related technologies (proxies, loadbalancers, etc.), and tend to be highly voluminous. Typical examples include Apache access logs, and CDN logs such as those from Amazon Cloudfront, Amazon S3 Server Access Logs, AWS ELB Access Logs, etc.

For large installations, bringing everything into an analytics tool is often so cost-prohibitive (storage, resources, license, etc.) that most users don't even bother. However, some of the logs contain relevant information when looked at individually (e.g., errors). The much larger majority contains relevant information when looked at in the aggregate (e.g., successes to determine traffic patterns, etc.).

It would be great if we could find a middle ground. With the Sampling Function, you can! Specifically, you can:

- Ingest enough sample events from the majority category that your aggregate analysis remains statistically significant.
- Ingest *all* events from the minority categories, and perform troubleshooting and introspection with full-fidelity data.

Using status as the Sampling Condition

Most of the access logs (including the ones mentioned above) have very similar formats. One quick way to sample is to look at the value of the status field.

2XX s indicate success and tend to be, by far, the most common ones – with

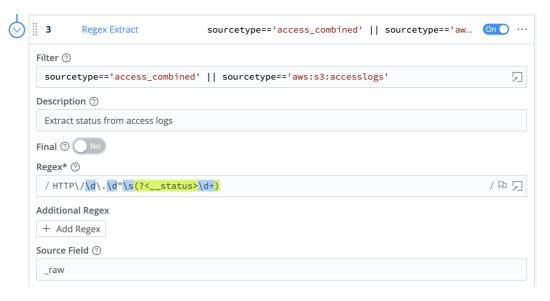
200 being the top. **Therefore**, 200 is the perfect candidate for sampling. All other statuses occur much less frequently, indicate conditions that often need to be looked at, and can be brought in with full fidelity.

Sample Status 200 at 5:1

1. Add a Regex Extract Function that looks at these sourcetypes:

```
sourcetype='access_combined' ||
sourcetype='aws:s3:accesslogs'
```

2. Configure that Function to extract a field called __status with this
 regex: /HTTP\/\d\.\d"\s(?<__status>\d+)/



Defining the Regex Extract Function

- 3. Add a Sampling Function to sample 5:1 when __status=200.
- 4. Save.



Sampling success reponses

Note About Sampling

Each time an event goes through the **Sampling** Function, an index-time sampled::<rate> field is added to it. Use this field in your statistical Functions, as necessary.

Other Sourcetypes

Examples of other sourcetypes that will benefit from sampling, but might need a different __status extraction regex:

Sourcetype	Filter Expression
Amazon Cloudfront Access Logs	<pre>sourcetype='aws:cloudfront:accesslogs'</pre>
Amazon ELB Access Logs	sourcetype='aws:elb:accesslogs'

Firewall Logs: VPC Flow Logs, Cisco ASA, Etc.

Recipe for Sampling Firewall Logs

Firewall logs are another source of important operational (and security) data. Typical examples include Amazon VPC Flow Logs, Cisco ASA Logs, and other technologies such as Juniper, Checkpoint, pfSense, etc.

As with Access Logs, bringing in everything for operational analysis might be cost-prohibitive. But sampling with Cribl LogStream can help you:

- Ingest enough sample events from the majority category that your aggregate analysis remains statistically significant. E.g., sample all ACCEPT s at 5:1.
- Ingest all events from the minority categories, and perform troubleshooting and introspection with full-fidelity data. E.g., bring in all REJECT s.

Sampling VPC Flow Logs

AWS' VPC Flow Logs feature enables you to capture information about the IP traffic going to and from network interfaces in your VPC. Flow Log data can be published to Amazon CloudWatch Logs and Amazon S3.

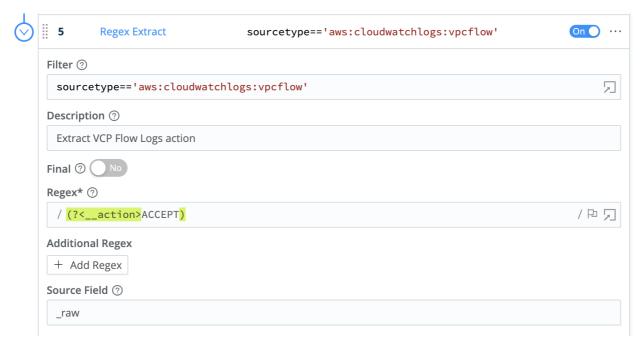
Typical VPC Flow Logs look like this:

Flow Log Records for Accepted and Rejected Traffic

```
2 123456789010 eni-abc123de 172.31.16.139 172.31.16.21 20641 22 6 20 4249 1418530010 1418530070 2 123456789010 eni-abc123de 172.31.9.69 172.31.9.12 49761 3389 6 20 4249 1418530010 1418530070 F
```

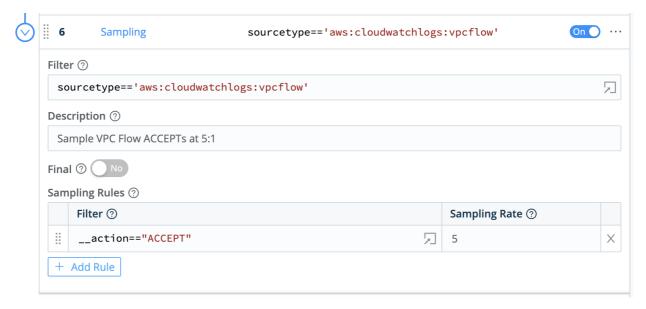
Let's use a very simple Filter condition and only look for ACCEPT events:

- 1. Add a Regex Extract Function that looks at: sourcetype='aws:cloudwatchlogs:vpcflow'
- 2. Configure that Function to extract a field called __action with this regex: /(?
 <__action>ACCEPT)/



Extracting the __action field

- 3. Add a **Sampling** Function to sample 5:1 when __action="ACCEPT".
- 4. Save.



Sampling ACCEPT events

Note About Sampling

Each time an event goes through the Sampling Function, an index-time field is added to it: sampled: <rate> . It's advisable that you use that in your statistical functions, as necessary.

Other Sourcetypes

Other sourcetypes that will benefit from sampling, but might need a different __action extraction regex:

Sourcetype	Filter Expression
Cisco ASA Logs	sourcetype='cisco:asa'
Related sourcetypes to consider sampling:	<pre>sourcetype='cisco:fwsm' sourcetype='cisco:pix'</pre>

Masking and Obfuscation

Masking and Anonymization of Data in Motion

To mask patterns in real time, we use the out-of-the-box Mask Function. This is similar to sed, but with much more powerful functionality.

Masking Capabilities

The Mask Function accepts multiple replacement rules, and accepts multiple fields to apply them to.

Match Regex is a JS regex pattern that describes the content to be replaced. It can optionally contain matching groups. By default, it will stop after the first match, but using /g will make the Function replace all matches.

Replace Expression is a JS expression or literal to replace matched content.

Matching groups can be referenced in the Replace Expression as g1, g2 ... gN, and the entire match as g0.

There are several masking methods that are available under C.Mask.:

- C.Mask.random: Generates a random alphanumeric string
- C.Mask.repeat: Generates a repeating char/string pattern, e.g., XXXX
- C.Mask.REDACTED: The literal 'REDACTED'
- C.Mask.md5: Generates a MD5 hash of given value
- C.Mask.sha1: Generates a SHA1 hash of given value
- C.Mask.sha256: Generates a SHA256 hash of given value

Almost all methods have an optional len parameter which can be used to control the length of the replacement. len can be either a number or string. If it's a string, its length will be used. For example:



Defining the replacement length

Masking Examples

Let's look at the various ways that we can mask a string like this one: cardNumber=214992458870391. The **Regex Match** we'll use is: $/(cardNumber=)(\d+)/g$. In this example:

- g0 = cardNumber = 214992458870391
- g1 = cardNumber=
- g2 = 214992458870391

Random Masking with default character length (4):

- Replace Expression: `\${g1}\${C.Mask.random()}`
- Result: cardNumber=HRhc

Random Masking with defined character length:

- Replace Expression: `\${g1}\${C.Mask.random(7)}`
- Result: cardNumber=neNSm8r

Random Masking with length preserving replacement:

- Replace Expression: `\${g1}\${C.Mask.random(g2)}`
- **Result**: cardNumber=DroJ73qmyaro51u3

Repeat Masking with default character length (4):

- Replace Expression: `\${g1}\${C.Mask.repeat()}`
- Result: Result: cardNumber=XXXX

Repeat Masking with defined character choice and length:

- Replace Expression: `\${g1}\${C.Mask.repeat(6, 'Y')}`
- **Result**: cardNumber=YYYYYY

Repeat Masking with length preserving replacement:

- Replace Expression: `\${g1}\${C.Mask.repeat(g2)}`
- **Result**: cardNumber=XXXXXXXXXXXXXXXX

Literal **REDACTED** masking:

- **Replace Expression**: `\${g1}\${C.Mask.REDACTED}`
- **Result**: cardNumber=REDACTED

Hash Masking (applies to: md5, sha1 and sha256):

- Replace Expression: `\${g1}\${C.Mask.md5(g2)}`
- **Result**: cardNumber=f5952ec7e6da54579e6d76feb7b0d01f

Hash Masking with left N-length* substring (applies to: md5, sha1 and sha256):

- Replace Expression: `\${g1}\${C.Mask.md5(g2, 12)}`
- Result: cardNumber=d65a3ddb2749
 - *Replacement length will **not** exceed that of the hash algorithm output; MD5: 32 chars, SHA1: 40 chars, SHA256: 64 chars.

Hash Masking with right N-length* substring (applies to: **md5**, **sha1** and **sha256**):

- Replace Expression: `\${g1}\${C.Mask.md5(g2, -12)}`
- Result: cardNumber= 933bfcebf992
 - *Replacement length will **not** exceed that of the hash algorithm output; MD5: 32 chars, SHA1: 40 chars, SHA256: 64 chars.

Hash Masking with length* preserving replacement (applies to: md5, sha1 and sha256):

- Replace Expression: `\${g1}\${C.Mask.md5(g2, g2)}`
- Result: cardNumber= d65a3ddb27493f5

 *Replacement length will not exceed that of the hash algorithm output;

 MD5: 32 chars, SHA1: 40 chars, SHA256: 64 chars.

Lookups as Filters for Masks

Overview

You can make your data architecture more maintainable by using Lookups to route and transform events within Cribl LogStream. This use case demonstrates an unusual solution, but one that served one Cribl customer's particular goals (which might overlap with yours):

- Ingest many hundreds of different sourcetype / index field combinations.
- Send all this data through a common Pipeline.
- Stack four Mask Functions in the Pipeline.
- Evaluate and process each sourcetype / index field combination only within its applicable Mask Functions – either two or three Masks per combination.
 - ⚠ This last restriction reduces latency, by preventing Mask Functions from evaluating non-applicable events, simply to ignore them.

Just to reiterate, this use case outlined here responded to this customer's requirements – one Pipeline combining multiple Mask Functions, for many sourcetype / index combinations.

More typically, you'd use multiple Pipelines to process different sourcetype / index combinations.

To enable this approach, the example below centralizes masking logic for multiple conditions in a Lookup table and corresponding Lookup Functions. The Lookup's output filters events to the applicable Mask Functions. Specifically, we'll show how to instruct LogStream to:

- Check for a particular index / sourcetype combination in each event,
 and
- Based on that combination, determine which Masks to apply to that event.

Design the Lookup

To use a lookup as a filter, you'd start by creating a comma-separated lookup table in this format, and adding it to LogStream:

```
index_tracker.csv

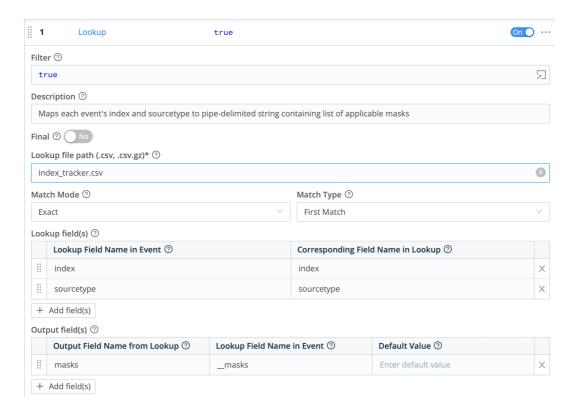
index, sourcetype, masks
apache_common, sourcetypec, ssn|credit_card|auth_token
syslog, sourcetypeb, ssn|auth_token
weblog, sourcetypea, auth_token|bearer_token
```

Below the header, each row specifies an index, a sourcetype, and (in the third column) a pipe-delimited list of applicable masks.

To make this example work, the table must have only **one** row for each index/sourcetype combination. (This unusual restriction is particular to this scenario.) So, as you build out the lookup table, you cannot add new masks for **existing** index/sourcetype combinations by appending new rows. Instead, you must modify the third column of the existing rows.

Configure the Pipeline

Create a LogStream Pipeline with a Lookup Function configured like this, pointing to your lookup table:



Lookup Function's configuration

This Function keys against both the index and sourcetype fields. When it finds a matching combination, it adds a new key-value pair to your event for future filtering.

The key of that key-value pair (namely, __masks) starts with a double underscore, to make it a LogStream internal field. This convention ensures that the key-value pair will **not** get passed along to the Destination.

However, you might prefer to export the key-value pair. For example, you might want a Splunk Destination to index the list of masks applied to a given event, alongside that event. (This approach applies to many forensic use cases.) If so, remove the double underscore from the above Function's **Lookup Field Name** in Event value, and from the subsequent Filter expressions for each Mask Function.

Each Mask Function has a JavasScript Filter that breaks the pipe-delimited string into an array, and determines whether the tag for that type of mask (e.g., bearer_auth) is in the __masks key-value pair. If so, it applies the mask processing. If not, the event moves on to the Pipeline's next Mask Function.

Here are the four Mask Functions below the Lookup Function:



Mask Functions

In this particular example, the pipe-delimited mask tags in the lookup table's third column match the Mask Functions' names, as well as matching their **Filter** conditions. This is just for simplicity – the Functions could have any names, as long as the **Filter** expressions match the tags.

Lookups and Regex Magic

Regular expressions are not just for field extractions – they can also be used inside lookup tables, and in Functions, to replace and manipulate values within fields. Let's walk through a Pipeline that demonstrates four different ways to leverage regular expressions in LogStream.

Why Lookup Tables Matter

When organizations use host naming standards, it is easy to understand things like regions, availability zones (AZs), IP addresses, and more. For example, consider an Amazon host called:

```
ec2-35-162-133-145.us-west1-a.compute.amazonaws.com
```

This is an EC2 host with a (dashed) IP address 35-162-133-145, in the us-west1 region, in Availability Zone a . You can also see the domain: compute.amazonaws.com.

While we can understand the enriched host names, we don't know which indexes to route the data to, nor which sourcetypes to assign to the events, without looking up this information from another source. Doing so is often a huge challenge for organizations. To solve this challenge, let's combine Regex Extract, Lookup, and Eval Functions with some sample events to demonstrate the power of LogStream.

Sample Events

The events below have timestamps broken out, but no indexes, sourcetypes, or other details have been assigned yet:

```
a_raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                 [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α cribl_breaker: Break on newlines
         a_raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Loo
                  king for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α cribl_breaker: Break on newlines
-06:00
          α_raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.py
                 [DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           α cribl_breaker: Break on newlines
         α_raw: Feb 06 2021 12:59:44.232 GMT: ec2-76-187-246-132.europe-west3-b.compute.amazonaws.com: cloud-init[2929]: stages.py[DEBUG]: R
                 unning module power-state-change (<module 'cloudinit.config.cc_power_stat... Show more
2021-02-06
06:59:44.232 # _time: 1612616384.232
           α cribl breaker: Break on newlines
         a _raw: Feb 06 2021 17:04:16.921 GMT: ec2-67-205-202-104.northamerica-northeast1-c.compute.amazonaws.com: cloud-init[2929]: util.py
                 [DEBUG]: Running command ['lxc-is-container'] with allowed return codes [0... Show more
2021-02-06
11:04:16.921 # _time: 1612631056.921
           a cribl breaker: Break on newlines
           \alpha_raw: Feb 06 2021 19:45:47.687 GMT: ec2-87-209-176-201.southamerica-east1-a.compute.amazonaws.com: DataSourceEc2.py[DEBUG]: Remove
                 d the following from metadata urls: ['http://instance-data.:8773']
2021-02-06
13:45:47.687 # _time: 1612640747.687
           α cribl_breaker: Break on newlines
-06:00
```

The Regex Extract Function

Before we can assign an index or sourcetype, we need to extract the host, region, az, and domain fields from the events. We can use a Regex Extract Function with this regular expression to extract all four fields:

```
GMT:\s+(?<host>[^.]+)\.(?<region>\w+-\w+\d+)-(?<az>[^.]+)\.(?<domain>[^:]+):
```

Here's that Regex Extract in a LogStream Pipeline:



Results of the Regex Extract Function

On the **OUT** tab of LogStream's Preview pane, the extracted fields of az, domain, host, and region now appear below the _raw event. You can use these extracted fields for searching in your preferred search solution.

```
IN OUT
                                                                                     Select Fields (8 of 8)
            a_raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                 [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
           a host: ec2-35-162-133-145
           a region: us-west1
          a_raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                  oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α az: b
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           α host: ec2-48-169-111-182
         α region: us-east2
3
            a _raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.p
                  y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
00:29:11.841 # _time: 1612592951.841
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           α host: ec2-21-187-232-201
           α region: asia-northeast3
```

Lookups

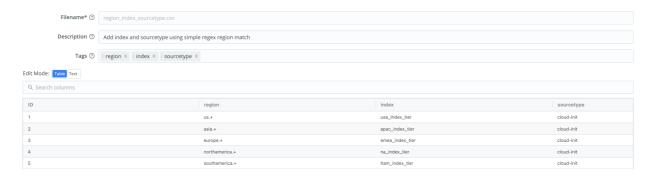
We still need to determine the index and sourcetype. LogStream's Lookup Function enriches events with external fields. We'll use it with the newly extracted region field to assign an index and sourcetype to these events.

Lookup File

First, we need to create a lookup table for the Function to reference. For this, we'll use regex again.

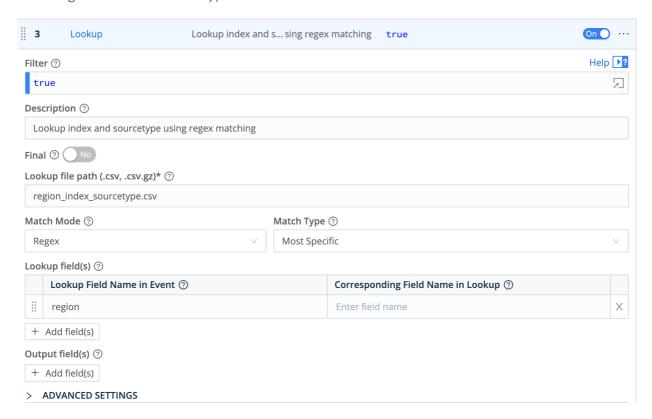
In the table below, five simple regular expressions map the extracted region field to the appropriate index and sourcetype. For example, the region us-west1-a starts with us, so it matches the first regular expression: us.+

We use this lookup table's first row to assign an index of usa_index_tier, and a sourcetype of cloud-init, to each matching event. The region patterns in the table's four remaining rows work the same way.



Lookup Function

With the lookup table saved as region_index_sourcetype.csv , the Lookup Function below matches the events' extracted region field against the regular expressions, and returns the matching index and sourcetype .



There's actually more here than meets the eye. Note that we've specified no **Output Fields**. From the Lookup Function's documentation, we know this means that the Function will default to outputting all fields in the lookup table. So we get the contents of both remaining columns in the table we saw above: index and sourcetype.

Results of the Lookup Function

With the Lookup Function added to our Pipeline, the Preview pane's **OUT** tab shows that the index and sourcetype are now added to each event.

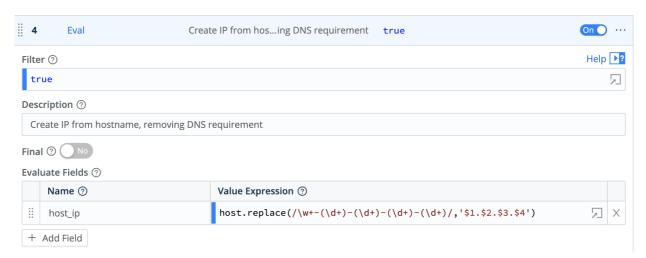
```
a _raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                   [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://145.133.162.36/latest/api/token', 'headers': {'X-a
2021-02-05
                   ws-ec2-metadata-token-ttl-seconds': '21600', 'User-Agent': 'Cloud-Init/19.3-5.amzn2'}, 'allow_redirects': True, 'method':
20:18:31.286
                   'PUT', 'timeout': 1.0} configuration Showless
-06:00
            # _time: 1612577911.286
            α az: a
            α cribl_breaker: Break on newlines
            α cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
            α host: ec2-35-162-133-145
            α index: usa_index_tier
            a region: us-west1
            α sourcetype: cloud-init
           a _raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                   oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α cribl_breaker: Break on newlines
            a cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
           α host: ec2-48-169-111-182
            α index: usa_index_tier
            α region: us-east2
           a sourcetype: cloud-init
            a _raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.p
                   y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
            α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
           a host: ec2-21-187-232-201
            α index: apac_index_tier
            α region: asia-northeast3
            a sourcetype: cloud-init
```

Getting Host IP Address from Host Name

Since the IP address is present in the host field, we can create the host_ip field using an Eval Function with this replace method:

```
host.replace(/\w+-(\d+)-(\d+)-(\d+)/,'$1.$2.$3.$4')
```

This regular expression uses capture groups and pulls the four IP octets present in the hostname to build the <code>host_ip</code>. These four capture groups are noted as \$1.\$2.\$3.\$4, respectively. This method is very fast, and removes the need to perform a DNS lookup from the <code>host</code> field to get the host's IP address. Magic!



Results of the Eval Function and Replace Method

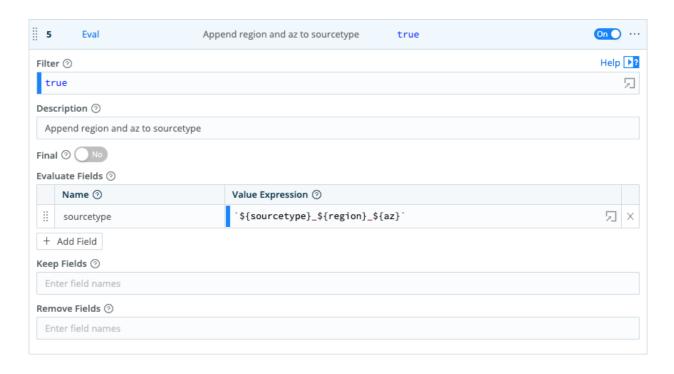
The host_ip field is now added to the events, displayed below host:

```
IN OUT
                                                                                      Select Fields (11 of 11)
           a _raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                  [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           \alpha domain: compute.amazonaws.com
           a host: ec2-35-162-133-145
           a host_ip: 35.162.133.145
           α index: usa_index_tier
           a region: us-west1
           α sourcetype: cloud-init
           a _raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                  oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α az: b
           α cribl_breaker: Break on newlines
           \alpha cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           a host: ec2-48-169-111-182
           α host_ip: 48.169.111.182
           α index: usa index tier
           α region: us-east2
         α sourcetype: cloud-init
           a_raw: Feb 06 2021 06;29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]; atomic_helper.p
                  y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           a host: ec2-21-187-232-201
           a host_ip: 21.187.232.201
           α index: apac_index_tier
           α region: asia-northeast3
           a sourcetype: cloud-init
```

Customizing the Sourcetype

Finally, let's put some sense into the <code>sourcetype</code> field, using another Eval Function. By combining the values of the $\{sourcetype\}_{\{region\}_{\{az\}}, the sourcetype becomes cloud-init_us-west1_a - so now you can understand much more about the sourcetype at a glance.$

Examine this Eval Function's value expression, taking careful note of the backticks (` `) and braces ({ }) that surround the field names, and the underscore (_) that separates them.



Results of the Eval Function to Combine Values

Take a look at the updated sourcetypes, and enjoy exploring LogStream with your new knowledge!

```
Select Fields (11 of 11)
IN OUT
1
            α_raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                   [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α cribl breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           a host: ec2-35-162-133-145
           a host_ip: 35.162.133.145
           \alpha index: usa_index_tier
           a region: us-west1
           a sourcetype: cloud-init_us-west1_a
         a _raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                  oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α az: b
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           α host: ec2-48-169-111-182
           α host_ip: 48.169.111.182
           α index: usa_index_tier
           α region: us-east2
           α sourcetype: cloud-init_us-east2_b
           a _raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.p
                  y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           a az: a
           α cribl_breaker: Break on newlines
           a cribl_pipe: setting_index_by_region_availability_zone
           a host: ec2-21-187-232-201
           a host_ip: 21.187.232.201
           α index: apac_index_tier
            α region: asia-northeast3
            α sourcetype: cloud-init_asia-northeast3_a
```

Try This at Home

Below you'll find the the lookup table, Pipeline, and sample events demonstrated in this use case. Create the lookup file first, and then import the Pipeline. (The order matters, because the Pipeline import depends on the lookup table's presence.)

Lookup Table

To create the lookup table in LogStream's UI, select **Knowledge > Lookups**, then click + **Add New** and select **Create with Text Editor**.

Copy and paste in the header and rows listed below, then save the result as:

```
region_index_sourcetype
region,index,sourcetype
us.+,usa_index_tier,cloud-init
asia.+,apac_index_tier,cloud-init
europe.+,emea_index_tier,cloud-init
northamerica.+,na_index_tier,cloud-init
```

southamerica.+,ltam index tier,cloud-init

region_index_sourcetype.csv.

Pipeline

Below is an export of the whole LogStream Pipeline presented here. Import this JSON to get a Pipeline named: setting_index_by_region_availability_zone.json.

```
Magic Pipeline
  "id": "setting_index_by_region_availability_zone",
  "conf": {
    "output": "default",
    "groups": {},
    "asyncFuncTimeout": 1000,
    "functions": [
        "filter": "true",
        "conf": {
          "comment": "This pipeline demonstrates four different ways to leverage regular express
        "id": "comment"
      },
        "filter": "true",
        "conf": {
          "source": "_raw",
          "iterations": 100,
          "overwrite": false,
          "regex": "/GMT:\\s+(?<host>[^.]+)\\.(?<region>\\w+\\d+)-(?<az>[^.]+)\\.(?<domain)</pre>
        },
        "id": "regex_extract",
        "disabled": false,
        "description": "Extract host, region, availability_zone, and domain"
      },
      {
```

```
"filter": "true",
  "conf": {
   "matchMode": "regex",
    "matchType": "specific",
    "reloadPeriodSec": 60,
    "addToEvent": false,
    "inFields": [
        "eventField": "region"
   ],
    "ignoreCase": false,
    "file": "region_index_sourcetype.csv"
  },
  "id": "lookup",
  "disabled": false,
  "description": "Lookup index and sourcetype using regex matching"
},
{
  "filter": "true",
  "conf": {
    "add": [
      {
        "name": "host_ip",
```

Sample Events

And here's a sample of raw events that you can upload or copy/paste into LogStream's Preview pane to test the Pipeline's Functions:

Sample events

```
Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[29] Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[29] Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-Feb 06 2021 12:59:44.232 GMT: ec2-76-187-246-132.europe-west3-b.compute.amazonaws.com: cloud-init Feb 06 2021 17:04:16.921 GMT: ec2-67-205-202-104.northamerica-northeast1-c.compute.amazonaws.com Feb 06 2021 19:45:47.687 GMT: ec2-87-209-176-201.southamerica-east1-a.compute.amazonaws.com: Data
```

From here, modify the sample data, lookup table, and Functions to adapt this approach to your own needs!

Regex Filtering

Regex Filtering of Data in Motion

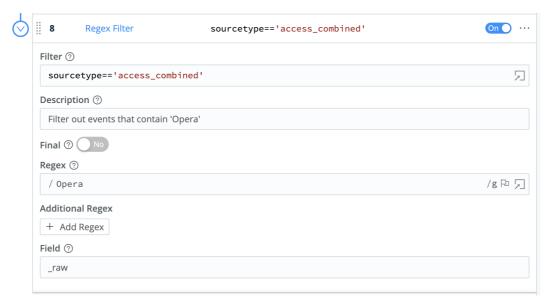
To filter events in real time, we use the out-of-the-box **Regex Filter** Function. This is similar to nullqueueing with TRANSFORMS in Splunk, but the matching condition is way more flexible.

Regex Filtering Example

Let's see how we can filter out any sourcetype—'access_combined' events whose _raw field contains the pattern Opera:

First, make sure you have a Route and Pipeline configured to match desired events.

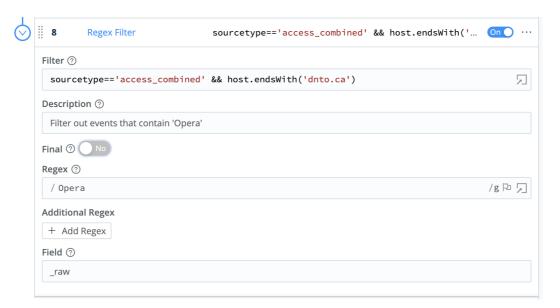
Next, let's add a Regex Filter Function to it:



Defining the Regex Filter Function

Next, verify that this search does **not** return any results: sourcetype="access_combined" Opera

You can add more conditions to the Filter input field. For example, to further limit the filtering to only events from hosts with domain <code>dnto.ca</code>, change the filter input as shown below:



Filtering by host

This is a very flexible method for filtering incoming events in real time, on virtually any arbitrary conditions.

Encrypting Sensitive Data

Encryption at Ingest-Time and Decryption in Splunk

With Cribl LogStream, you can encrypt your sensitive data in real time before it's forwarded to and stored at a destination. Using the out-of-the-box Mask function, you can define patterns to encrypt with specific key IDs or key classes. To decrypt in Splunk, you will need to install Cribl App for Splunk on your search head. (The app will default to mode-searchhead.)

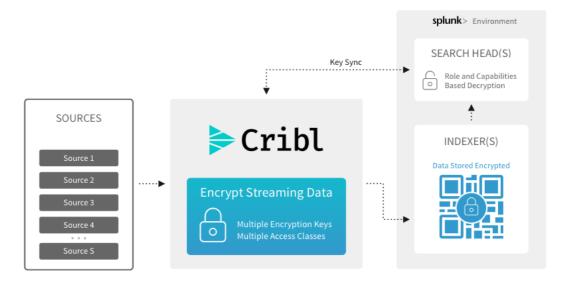
Keys and Key Classes

Symmetric encryption keys can be configured through the CLI or the UI. They're used to encrypt the patterns, and users are free to define as many keys as required.

Key classes are collections of keys that can be used to implement multiple levels of access control. Users (or groups of users) that have access to data with encrypted patterns can be associated with key classes. You can use these classes to provide more-granular access rights, such as read versus decryption permissions on a dataset.

Encrypting in Cribl LogStream and Decrypting in Splunk

- 1. Define one or more keys and key classes on Cribl LogStream.
- 2. Sync auth with the decryption side (Splunk Search Head)
- 3. Apply the Mask function to patterns of interest, using C.Crypto.encrypt().
- 4. Decrypt on the Splunk search head, using Role-Based Access Control on the decrypt command.

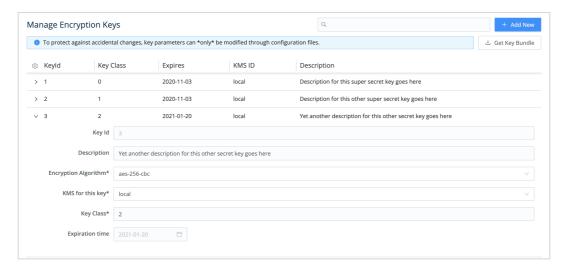


Encrypting in LogStream, decrypting in Splunk

Example

Encryption Side

• Generate keys via the UI, in **Settings** > **Encryption Keys**:



Adding a new encryption key

• Or generate one or more keys via the CLI, as follows:

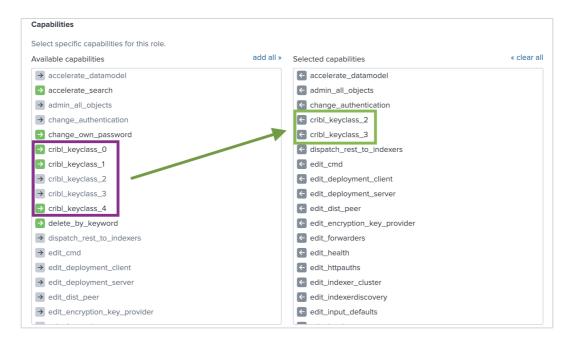
```
$CRIBL_HOME/bin/cribl keys add -c 1 -i
...
$CRIBL_HOME/bin/cribl keys add -c <N> -i
```

Add -e <epoch> if you'd like to set expiration for your keys.

• Sync auth/(cribl.secret|keys.json). To decrypt data, the decrypt command will need access to these keys. The cribl.secret and keys.json files in \$CRIBL_HOME/local/cribl/auth should be synced/copied over to the search head (decryption side).

Decryption Side

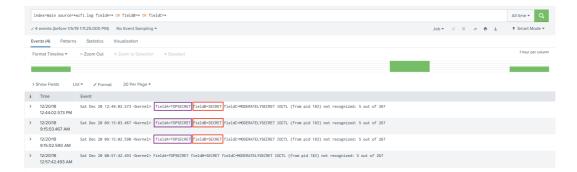
- Install Cribl App for Splunk on your search head. It will default to modesearchhead.
- Assign permissions to the decrypt command, per your requirements.
- Assign capabilities to your Roles, per your requirements. Capability names should follow the format cribl_keyclass_N, where N is the Cribl Key Class. For example, a role with capability cribl_keyclass_1 has access to all key IDs associated with key class 1. You can use more capabilities, as long as they follow this naming convention.



Selecting capabilities

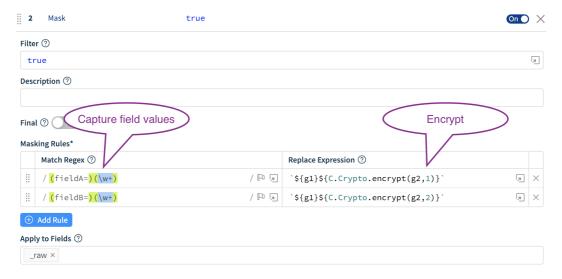
Usage

Before Encryption: Sample un-encrypted events. Notice the values of fieldA and fieldB.



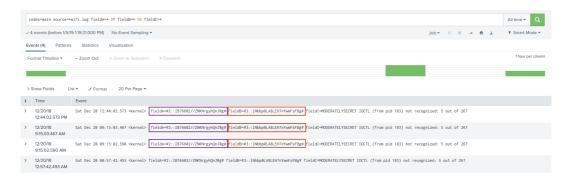
Events before encryption

Next, encrypt fieldA values with key class 1, and fieldB with key class 2.



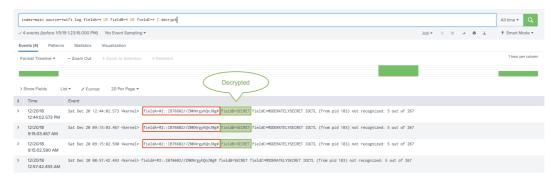
Encrypting two fields with separate key classes

After Encryption: again, notice the values of fieldA and fieldB.



Both fields encrypted

Here, we've decrypted fieldB but not fieldA. This is because the logged-in user has been assigned the capability cribl_keyclass_2, but not cribl_keyclass_1.



One field decrypted

Syslog Data Reduction

When ingesting data from syslog senders, Cribl LogStream can readily trim data volume by 30% or more, optimizing infrastructure for downstream services like Splunk or Elasticsearch. Here, we outline some best practices for replacing your syslog server with LogStream.

Syslog Event Parsing

By default, a LogStream Syslog Source will produce the following fields:
_time, appname, facility and facilityName, host, message, and
severity and severityName.

Parsed syslog event

This output is much more readable, but we haven't saved any volume – we now have redundant pairs of fields (numeric versus text) representing the facility and severity.

Below, we'll outline how to streamline syslog events to something more like this:

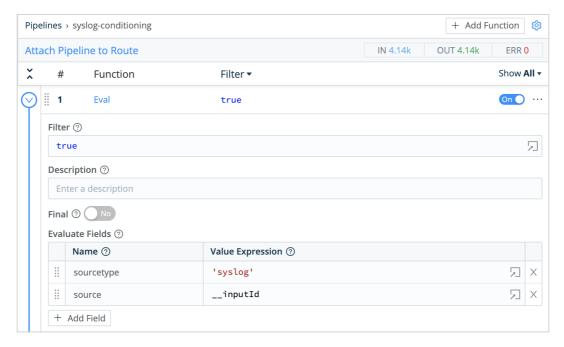
Parsed and redacted syslog event

This extracts the essentials, removes the redundancies, adds one new field that identifies the LogStream Pipeline we're about to build, and shrinks the outbound _raw payload to just its message component. For still further efficiencies, we'll look at how to drop or downsample frequent events, and how to balance high-volume syslog inputs across LogStream worker processes.

Create Pre-Processing Pipeline

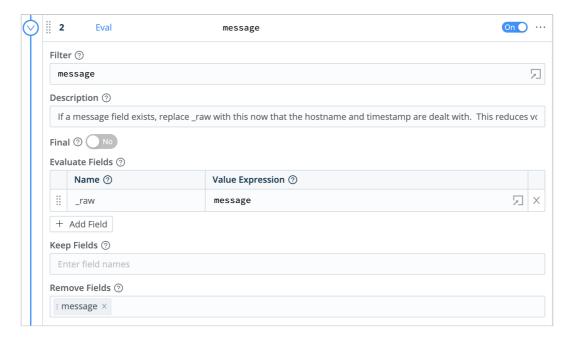
Even before setting up a syslog Source, our first step is to create a preprocessing Pipeline that will be available to normalize incoming events on all syslog Sources, reducing data volume as shown above.

The Pipeline starts with an **Eval** Function, whose **Evaluate Fields** section tags syslog events with two new fields indicating their origin: sourcetype: 'syslog' and source: __inputId . Because this Pipeline is designed only to condition all incoming syslog data, we leave **Filter** set to its default true value, to process all events.



Eval Function to tag syslog events' origin

A second **Eval** Function filters for the presence of a message field. If found, it overwrites the _raw field with message , and then deletes message as redundant. This function alone typically reduces syslog data volume by 15–20%.



Eval Function to rewrite message as _raw

⚠ Before using this Pipeline in production, preview sample data to verify that you're not dropping any essential information.

This third **Eval** Function deletes two redundant fields. Its **Filter** condition makes sure both of these string fields exist and contain values: severity \neq null & facility \neq null . If so, it removes their corresponding numeric fields, severity and facility.



Eval Function to remove redundant numeric fields

To further shrink the output, this fourth **Eval** Function removes procid fields that contain only a dash – meaning "no value extracted." We set **Filter**: procid='-' and **Remove Fields**: procid.



Eval Function to remove empty procid fields

With these four Functions enabled, the **Preview** pane's **Basic Statistics** confirm that we've reduced the _raw field's length by more than 30%.



Data reduction example

Dropping Noisy Data

With some syslog senders, like VMware ESX/ESXi, 80–90% of incoming events can be of debug severity. To further reduce volume, you could use this final **Drop** Function to drop all these events. Its **Filter** is simply severityName='debug'.



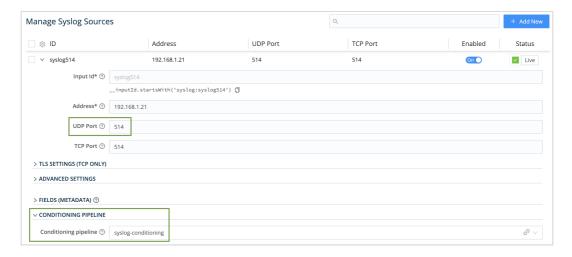
Drop Function to remove debug events

Enabling this Function could up our volume savings to about 40%. Depending on your needs, you might prefer to:

- Use a Function like this in your Route's processing Pipeline, rather than in this upfront Pipeline.
- Also drop info events.
- Instead use a Sampling Function to sample debug events (at a ratio like 1:10), or a Dynamic Sampling function.
- Instead use a Suppress Function to clamp down the frequency of debug events.

Create Syslog Source

Once we've built and saved the pre-processing Pipeline, our next step is to add a Syslog Source.



Syslog Source configured for UDP and pre-processing Pipeline

Specify the UDP Port where you want this Source to listen for syslog data.

Then attach the pre-processing Pipeline that you created above, and save the Source.

i Cribl generally recommends selecting UDP, rather than TCP, for high-volume syslog senders. This facilitates efficient load balancing by not continuously tying such senders to any one LogStream Worker Process. See Sizing and Scaling for more details.

Fields/Metadata

In the pre-processing Pipeline, we tagged *all* incoming syslog events with new sourcetype and source fields to indicate their origin. Alternatively, you could use the Source's **Fields/Metadata** section to define similar key-value pairs, specific to each of your Syslog Sources.

Create Route(s)

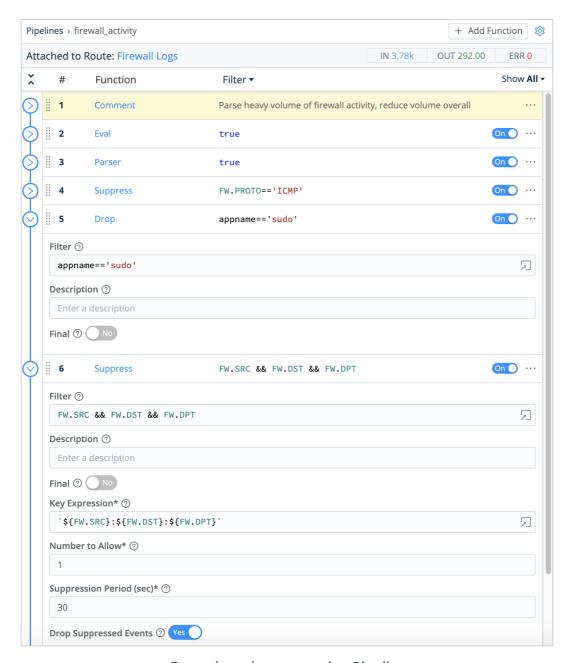
Create Routes, as needed, for each of your Syslog Sources. Give each Route a corresponding **Filter** expression, and set its **Output** to the Destination where you want to send its processed syslog data.



Example syslog Route

Processing Pipelines, and Next Steps

For any or all syslog Routes, you can define and attach a processing Pipeline. These can apply more-granular Filters and Functions to further reduce volume, using techniques like Sampling, Dynamic Sampling, or (as shown below) Drop and Suppression. Your most-verbose Syslog Sources are ideal targets for such further processing.



Example syslog processing Pipeline

Splunk to Elasticsearch

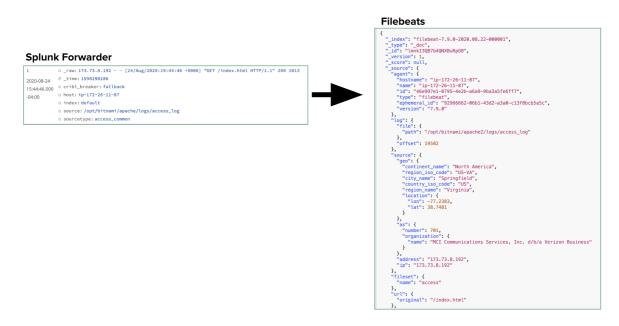
To route data from existing Splunk infrastructure to Elasticsearch services, you might face a daunting task: re-architecting your entire forwarding tier. This could require retooling lots of servers – up to hundreds, or thousands – to uninstall their Splunk forwarders, and swap in Elastic-compatible agents.

Cribl LogStream can reduce this effort to just a few hours: Configure one Splunk outputs.conf stanza to output to LogStream, and propagate that across all your Splunk servers. Done!

Next, you can easily configure LogStream to listen for Splunk data on one port, and to route that data to all the Elasticsearch destinations you want to feed.

Transforming Data from Splunk to Elastic Format

Also, in LogStream's core, you can easily design a Pipeline that modifies the original Splunk event into Elastic's Common Schema – making it look exactly like an event generated by an Elastic agent. These transformations help you make the most of Elastic's offerings, like Filebeats, etc.

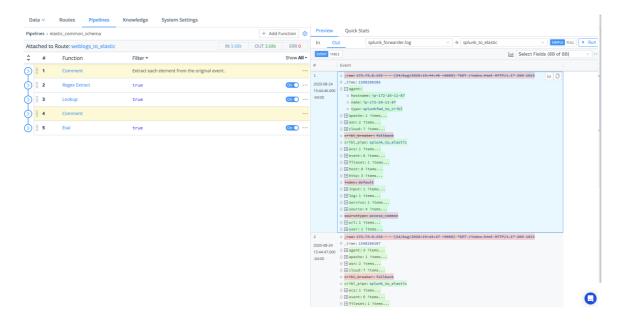


Transforming to Elastic Common Schema

Some of the LogStream Functions useful in transforming Splunk-generated events into Elastic's format are:

- Regex Extract: Extract a portion of the raw event, and place it into a specified field.
- Lookup: key off the host IP to add fields like hostname, name, id, and type.
- Eval: Turn key-value pairs into nested JSON objects.
- GeoIP: Correlate source IP to a geographic database.

We'll show all four in our example Pipeline below, although you might need only a subset.



LogStream Pipeline and Data Preview

Goat Rid of Some Guesswork

LogStream will offer you further time savings as you configure the internal data transformation. LogStream's Data Preview features enable you to test transformations' results as you build your Pipeline, before you commit or run it.

This eliminates blind guesswork in Splunk configuration files to specify source -> index transformations, check the results, and then start all over again. In particular, LogStream's Regex Extract Function provides a regex101-like UI, to facilitate precisely designing and debugging your regex expressions.

Let's goat started on the example.

Configure Splunk Forwarder

First, in a Splunk App, configure a Splunk forwarder (UF or HF) to specify your Cribl Workers as destinations. Use outputs.conf stanzas of this form:

```
.../outputs.conf

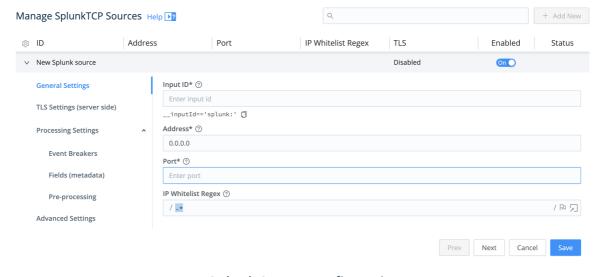
[tcpout]
disabled = false
defaultGroup = cribl, <optional_clone_target_group>,

[tcpout:cribl]
server = [<cribl_ip>|<cribl_host>]:<port>, [<cribl_ip>|<cribl_host>]:<port>, ...
sendCookedData=true
```

Push the app using the deployment server.

Configure Splunk Source in LogStream

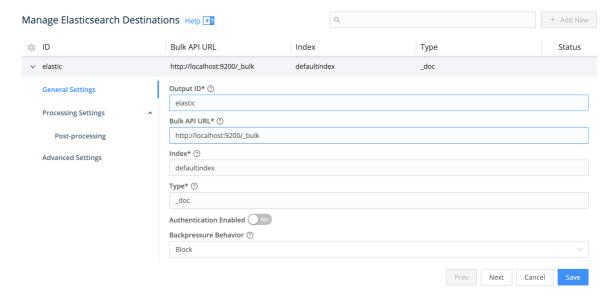
Next, in LogStream, configure a Splunk Source. The key requirement here is to set the **Port** to listen on. (Optionally, you can also configure TLS, Event Breakers, metadata fields, and/or a pre-processing Pipeline.)



Splunk Source configuration

Configure Elasticsearch Destination

To configure LogStream's output, set up an Elasticsearch Destination by specifying the **Bulk API URL** and **Index**.



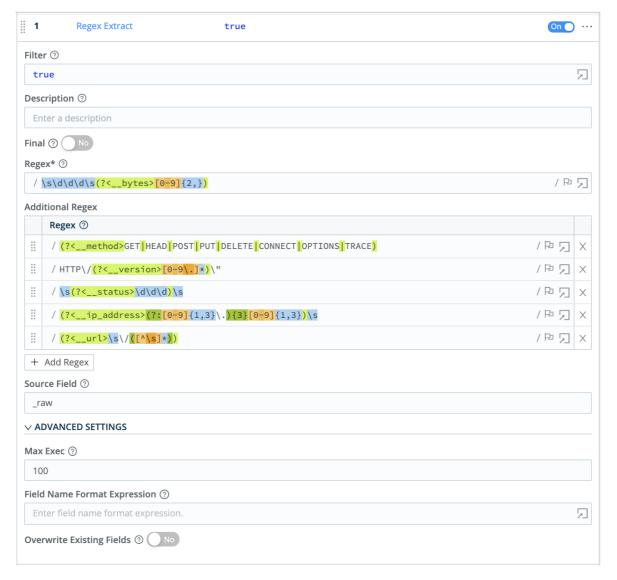
Elasticsearch Destination configuration

Configure Pipeline

Next, this section shows several Functions that you can assemble into a Pipeline to transform incoming Splunk events to match the Elastic Common Schema.

Regex Extract Function

First, use a Regex Extract Function to break the Splunk events into fields. Try the sample configuration shown below:



Regex Extract Function

Here are the six rows of regex in this example:

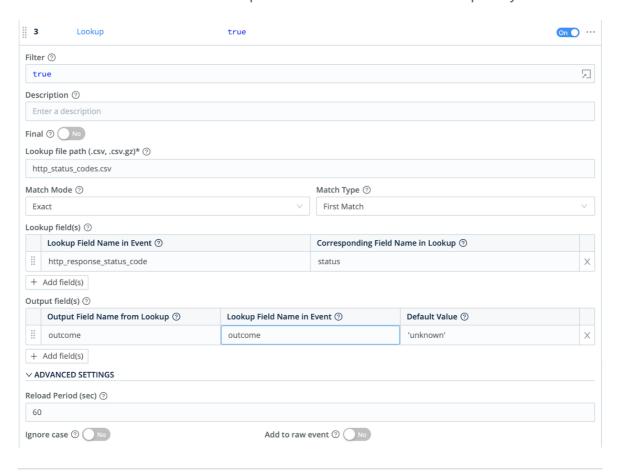
```
Regex; Additional Regex

/\s\d\d\s(?<_bytes>[0-9]{2,})/
/(?<_method>GET|HEAD|POST|PUT|DELETE|CONNECT|OPTIONS|TRACE)/
/HTTP\/(?<_version>[0-9\.]*)\"/
/\s(?<_status>\d\d\d)\s/
/(?<_ip_address>(?:[0-9]{1,3}\.){3}[0-9]{1,3})\s/
/(?<_url>\s\/([^\s]*))/
```

As you refine your expression, capture a sample of incoming Splunk data to test your regex's results in LogStream's right Preview pane.

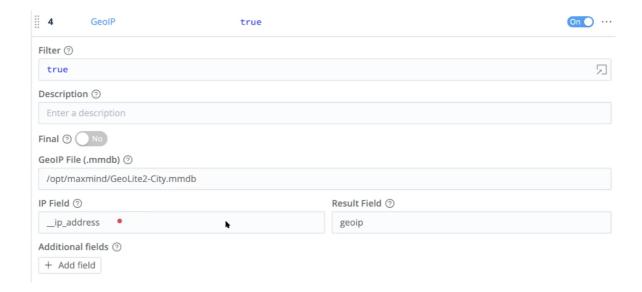
Lookup Function

In this example, we next add a Lookup Function, to translate HTTP error codes to readable text. Note this Function's optional **Reload Period** field, in which you can define a reload interval for a lookup file whose contents refresh frequently.



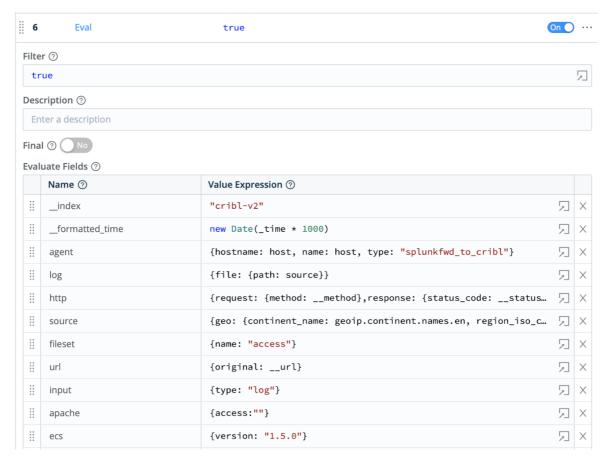
GeoIP Function

To enrich the Splunk data, we next use a GeoIP Function. This a specialized lookup against a database of IP addresses by geographic location. This Function's output can provide Elasticsearch with location fields like lat and long.



Eval Function

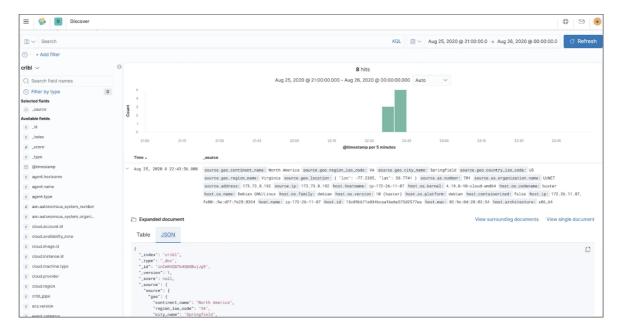
Finally, to further enrich the outbound events, the Pipeline uses an Eval Function. This adds multiple key-value pairs that define and populate fields conforming to the Elastic Common Schema.



Eval Function

Results

After attaching your Pipeline to a Route, here's an an exported event, all happy in Elasticsearch with nested JSON.



Event as exported to Elasticsearch

For More Info

For additional details on configuring Splunk forwarders for LogStream, see this related documentation:

- Configuring a Splunk (TCP) Forwarder
- Configuring Cribl App for Splunk on an HF

Reducing Windows XML Events

Here, we demonstrate how to use just a few LogStream Functions to parse WindowsXML events and reduce their volume by 34–70%, dramatically reducing your downstream infrastructure requirements.

Using Eval and C.Text.parseWinEvent

LogStream's internal C.Text.parseWinEvent method parses Windows XML strings and returns a prettified JSON object. You can use this function within an Eval Function to parse an event or an ad hoc XML string. It works like C.Text.parseXml, but with Windows events, it produces more-compact output.

As you can see from its signature, C.Text.parseWinEvent accepts an optional nonValues parameter that can further reduce an event's size by discarding characters that you specify as redundant.

i (method) C.Text.parseWinEvent(xml: string, nonValues?: string[]): any

@param - xml - an XML string; or an event field containing the XML.

@param - nonValues - array of string values. Elements whose value equals any of the values in this array will be omitted from the returned object. Defaults to ['-'], meaning that elements whose value equals - will be discarded.

@returns – an object representing the parsed Windows Event; or undefined if the input could not be parsed.

XML: Threat or Menace?

When working with XML, an anonymous Reddit user's quote sums up the challenge: "Some languages can be read by human, but not by machines, while others can be read by machines but not by humans. XML solves this problem by being readable to neither." An example of a Windows XML event only reinforces this quote:



This Windows XML raw event is 1.36KB in size:

	_raw Length ②	Full Event Length ③	Number of Fields ③	Number of Events ②
IN	1.36KB	1.43KB	3	1
OUT	1.36KB	1.47KB	4	1
DIFF	0.00%	个 2.73%	↑ 33.33%	0.00%

Eval to the Rescue

In our initial Eval Function below, the **Value Expression** uses C.Text.parseWinEvent to simply parse the _raw Windows XML event and turn it into a prettified JSON object:



The resulting JSON event is now down to 921.00B in size, a 34.07% reduction of the event:



Removing Unnecessary Fields with a Better Eval

But we can do better. The fields containing essentially null values ('0', $'0\times0'$, or '-') bloat events, demanding extra infrastructure and storage:

```
(} = _raw:
  {} = Event:
     ⊕ EventData:
       {} □ Data:
          α AuthenticationPackageName: NTLM
           α FailureReason: 9%2313
           a IpAddress: 64.203.230.138
           a IpPort: 0 ←
           a KeyLength: 0 ◀
           α LmPackageName: - ◀
           a LogonProcessName: NtLmSsp
           a LogonType: 3
           a ProcessId: 0x0 ←
           α ProcessName: - ◀
           a Status: 0xc000006d
           α SubjectDomainName: - ←
           a SubjectLogonId: 0x0 ←
           a SubjectUserSid: S-1-0-0
           a SubStatus: 0xc000006a
           a TargetUserName: ADMINISTRATOR
                                                                    Field Values: ['0','0x0','-']
           a TargetUserSid: S-1-0-0

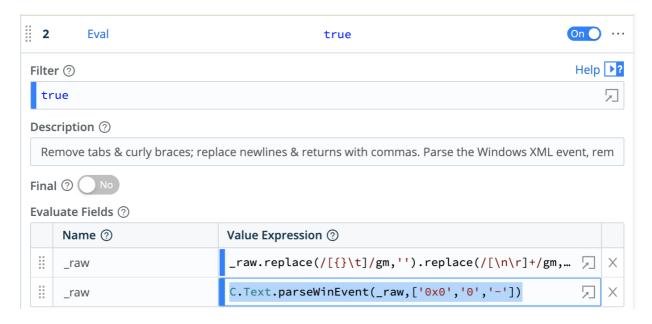
    □ TransmittedServices: - 
    □

           α WorkstationName: - ◀
     {} □ System:
        a Channel: Security
        a Computer: EC2AMAZ-CPMK6J5.cribl.poc
        a EventID: 4625
        a EventRecordID: 1598720
        {} = Execution:
          a ProcessID: 740
           α ThreadID: 948
        a Keywords: 0x8010000000000000
        α Level: Θ ◀
        α Opcode: 0 ◀
        {} ☐ Provider:
          a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}

    Name: Microsoft-Windows-Security-Auditing

        α Task: 12544
        {} ☐ TimeCreated:
           a SystemTime: 2020-11-12T17:27:33.608596000Z
```

Let's amplify the reduction by removing all of the fields whose values are in the set: ['0','0×0','-']. This improved version of our Eval Function parses the Windows XML event, and discards ['0','0×0','-'] values. (Its preceding row also tidies up events by removing tabs and curly braces, and replacing newlines and returns with commas.) The result is an even **smaller** prettified JSON object:



If you compare this Preview-pane screenshot to the Preview screenshot above, you can confirm that the fields with values matching $['0', '0\times0', '-']$ are removed:

```
(} □ _raw:
   {} 	☐ EventData:
         {} Data:

    AuthenticationPackageName: NTLM

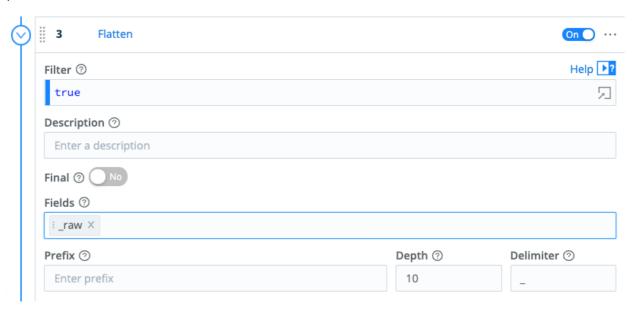
            α FailureReason: %%2313
            a IpAddress: 64.203.230.138
            a LogonProcessName: NtLmSsp
            α LogonType: 3
            a Status: 0xc000006d
            a SubjectUserSid: S-1-0-0
            α SubStatus: 0xc000006a
            a TargetUserName: ADMINISTRATOR
            a TargetUserSid: S-1-0-0
      {} - System:
         α Channel: Security
         a Computer: EC2AMAZ-CPMK6J5.cribl.poc
         a EventID: 4625
         a EventRecordID: 1598720
         {} = Execution:
            a ProcessID: 740
            α ThreadID: 948
         a Keywords: 0x8010000000000000
         ⊕ Provider:
            a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
            α Name: Microsoft-Windows-Security-Auditing
         α Task: 12544
         {} 	☐ TimeCreated:
            a SystemTime: 2020-11-12T17:27:33.608596000Z
```

The event is now down to 678.00B in size, translating to a 51.47% reduction from the original event:

	_raw Length ⑦	Full Event Length ⑦	Number of Fields ⑦	Number of Events ⑦
IN	1.36KB	1.43KB	3	1
OUT	678.00B	785.00B	4	1
DIFF	↓ -51.47%	↓ -46.42%	↑ 33.33%	0.00%

Flatten Function

LogStream's Flatten Function is designed to flatten fields out of a nested structure. Let's flatten the JSON object within <code>_raw</code> , to see if we can further reduce the event's size before we send it to our preferred destinations:



Using Flatten, we've successfully created top-level fields from the nested JSON structure:

```
a _raw: <Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider
       y-Auditing' Guid='{54849625-5478-4994-A5BA-3E3B0328C30D}'/><EventID>4625</Event... Show more
α _raw_Event_EventData_Data_AuthenticationPackageName: NTLM
a _raw_Event_EventData_Data_FailureReason: %%2313
a _raw_Event_EventData_Data_IpAddress: 64.203.230.138
a _raw_Event_EventData_Data_LogonProcessName: NtLmSsp
a _raw_Event_EventData_Data_LogonType: 3
a _raw_Event_EventData_Data_Status: 0xc000006d
a _raw_Event_EventData_Data_SubjectUserSid: S-1-0-0
a _raw_Event_EventData_Data_SubStatus: 0xc0000006a
a _raw_Event_EventData_Data_TargetUserName: ADMINISTRATOR
a _raw_Event_EventData_Data_TargetUserSid: S-1-0-0
a _raw_Event_System_Channel: Security
a _raw_Event_System_Computer: EC2AMAZ-CPMK6J5.cribl.poc
a _raw_Event_System_EventID: 4625
α _raw_Event_System_EventRecordID: 1598720
a _raw_Event_System_Execution_ProcessID: 740
a _raw_Event_System_Execution_ThreadID: 948
α _raw_Event_System_Keywords: 0x8010000000000000
a _raw_Event_System_Provider_Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
{\tt \alpha\_raw\_Event\_System\_Provider\_Name: Microsoft-Windows-Security-Auditing}
α _raw_Event_System_Task: 12544
a _raw_Event_System_TimeCreated_SystemTime: 2020-11-12T17:27:33.608596000Z
# time: 1609852209.078
α cribl breaker: Break on newlines
α cribl_pipe: Parse_Windows_XML_Events
```

Don't worry about the _raw field's deletion (red strikeout). This is the Flatten Function's default behavior. We'll restore _raw after we clean and reduce the event even more.

The flattened field names are extracted from _raw and delimited with _ . These field names are quite long. We can optimize them using the Rename Function.

Rename Function

Rename is designed to change fields' names, or to reformat their names (e.g., to normalize names to camelcase). You can use Rename to change specified fields (much like Eval), or to accomplish bulk renaming based on a JavaScript expression (much like the Parser Function). But Rename offers a streamlined way to alter only field names, without other effects.

Let's use Rename to remove any unnecessary prefixes from the field names, to further shrink our events. In the **Renaming Expression**, we build a a JavaScript expression to match the field names' prefixes (up to the underscore):



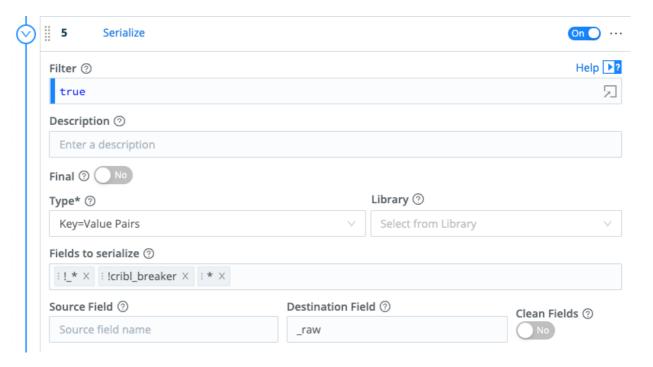
The resulting field names are now much more compact, and easier to work with and manage:

```
a _raw: <Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider
       y-Auditing' Guid='{54849625-5478-4994-A5BA-3E3B0328C30D}'/><EventID>4625</Event... Show more
# _time: 1609852209.078
α AuthenticationPackageName: NTLM
α Channel: Security
α Computer: EC2AMAZ-CPMK6J5.cribl.poc
a cribl_breaker: Break on newlines
α cribl_pipe: Parse_Windows_XML_Events
a EventID: 4625
a EventRecordID: 1598720
α FailureReason: %%2313
a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
a IpAddress: 64.203.230.138
a Keywords: 0x8010000000000000
α LogonProcessName: NtLmSsp
a LogonType: 3
α Name: Microsoft-Windows-Security-Auditing
a ProcessID: 740
a Status: 0xc000006d
a SubjectUserSid: S-1-0-0
a SubStatus: 0xc000006a
a SystemTime: 2020-11-12T17:27:33.608596000Z
a TargetUserName: ADMINISTRATOR
α TargetUserSid: S-1-0-0
α Task: 12544
α ThreadID: 948
```

Serialize Function

We started with bloated Windows XML data, and we've parsed and prettified it into JSON. Next, we'll extract key-value pairs. We'll use the Serialize Function, which serializes an event's content into a predefined format.

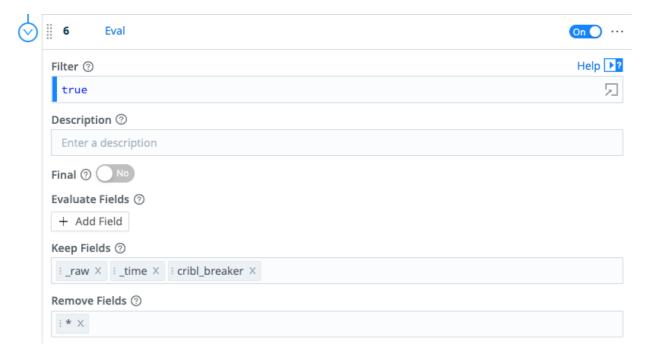
We set Serialize to change the **Type** to key-value pairs. (The Function's other supported target Types include JSON Object and CSV.) Here, Serialize takes the extracted fields and puts them back into _raw:



In the Preview pane, the _raw field is now back, serialized into compact, tidy key-value pairs:

```
α _raw: Name=Microsoft-Windows-Security-Auditing Guid={54849625-5478-4994-A58A-3E3B0328C30D} EventID=4625 Task=12544 Keywords=
        0x801000000000000 SystemTime=2020-11-12T17:27:33.608596000Z EventRecordID=1598720 ProcessID=740 ThreadID=948 Channel=
        Security Computer=EC2AMAZ-CPMK6J5.cribl.poc SubjectUserSid=S-1-0-0 TargetUserSid=S-1-0-0 TargetUserName=ADMINISTRATOR
        Status=0xc000006d FailureReason=%%2313 SubStatus=0xc000006a LogonType=3 LogonProcessName="NtLmSsp " AuthenticationPack
        ageName=NTLM IpAddress=64.203.230.138 Show less
# _time: 1609852209.078
\alpha AuthenticationPackageName: NTLM
α Channel: Security
α Computer: EC2AMAZ-CPMK6J5.cribl.poc
a cribl breaker: Break on newlines
α cribl_pipe: Parse_Windows_XML_Events
a EventID: 4625
a EventRecordID: 1598720
α FailureReason: %%2313
a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
a IpAddress: 64.203.230.138
α Keywords: 0x8010000000000000
α LogonProcessName: NtLmSsp
a LogonType: 3
\alpha Name: Microsoft-Windows-Security-Auditing
α ProcessID: 740
α Status: 0xc000006d
a SubjectUserSid: S-1-0-0
a SubStatus: 0xc000006a
a SystemTime: 2020-11-12T17:27:33.608596000Z
α TargetUserName: ADMINISTRATOR
α TargetUserSid: S-1-0-0
α Task: 12544
α ThreadID: 948
```

The last step is to remove any of the extracted fields you don't need before sending events to your destinations. We'll again call on the Eval Function, which adds or removes fields in events. (For a Splunk destination, these are index-time fields.) This final Eval Function looks like this:



To sum up, we've successfully transformed the original Windows XML event into key-value pairs:

```
a _raw: Name=Microsoft-Windows-Security-Auditing Guid={54849625-5478-4994-A5BA-3E3B0328C30D} EventID=4625 Task=12544 Keywords=
0x8010000000000000000 SystemTime=2020-11-12T17:27:33.608596000Z EventRecordID=1598720 ProcessID=740 ThreadID=948 Channel=
Security Computer=EC2AMAZ-CPMK6J5.cribl.poc SubjectUserSid=S-1-0-0 TargetUserSid=S-1-0-0 TargetUserName=ADMINISTRATOR
Status=0xc000006d FailureReason=%%2313 SubStatus=0xc000006a LogonType=3 LogonProcessName="NtLmSsp" AuthenticationPack
ageName=NTLM IpAddress=64.203.230.138 Show less

# _time: 1609852209.078
cribl_breaker: Break on newlines
cribl_pipe: Parse_Windows_XML_Events
```

And we've dramatically reduced the event's size, while retaining all of the necessary fields. The event is now down to 513.00B in size, translating to a 63.28% reduction from the original Windows XML:

	_raw Length ③	Full Event Length ③	Number of Fields ⑦	Number of Events ①
IN	1.36KB	1.43KB	3	1
OUT	513.00B	621.00B	4	1
DIFF	↓ -63.28%	↓ -57.61%	↑ 33.33%	0.00%

Try This at Home

Below is an export of the whole LogStream Pipeline presented here. Import this JSON to experiment with it and modify it to match your own needs:

```
Win XML Pipeline
{
    "id": "Windows_Security_Events",
    "conf": {
        "output": "default",
        "groups": {},
        "asyncFuncTimeout": 1000,
```

```
"functions": [
 {
    "filter": "true",
    "conf": {
     "comment": "This LogStream Pipeline reduces Microsoft Windows XML events, retains full
    "id": "comment"
 },
    "filter": "true",
    "conf": {
      "add": [
       {
          "name": "_raw",
          "value": "_raw.replace(/[{}\\t]/gm,'').replace(/[\\n\\r]+/gm,',')"
       },
        {
          "name": "_raw",
          "value": "C.Text.parseWinEvent(_raw,['0×0','0','-'])"
       }
     ]
    },
    "id": "eval",
    "disabled": false,
    "description": "Remove tabs & curly braces; replace newlines & returns with commas. Pars
 },
 {
   "filter": "true",
    "conf": {
     "fields": [
       "_raw"
     ],
     "prefix": "",
     "depth": 5,
     "delimiter": " "
   },
    "id": "flatten",
    "disabled": false,
    "description": "Flatten the object into key value fields"
 },
 {
   "filter": "true",
    "conf": {
      "baseFields": [],
      "renameExpr": "name.replace(/_raw_Event_\\w+_/,'')",
```

Finally, here's a sample of Windows XML events that you can upload to LogStream's Preview pane to try this out:

```
Sample data
```

```
<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic
<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic
SeBackupPrivilege
SeRestorePrivilege
SeTakeOwnershipPrivilege
SeDebugPrivilege
SeDebugPrivilege
SeSystemEnvironmentPrivilege
SeLoadDriverPrivilege
SeImpersonatePrivilege
SeImpersonatePrivilege
SeDelegateSessionUserImpersonatePrivilege
SeEnableDelegationPrivilege</pre>
SeEnableDelegationPrivilege
SeVentData

VeventData

Event>
```

<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic
<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic
<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic</pre>

BEST PRACTICES

Managing Large Lookups

This page offers a general approach to managing large lookup files. While LogStream's Git integration normally helps manage configuration changes, large lookups are exceptions. In many cases, you might want to exclude these files from Git, to reduce excessive deploy traffic. This approach can also prevent Git Push commands from encountering large file errors.

Good scenarios for this approach are:

- Large binary files like databases which don't benefit from Git's typical efficient storage of only the deltas between versions. (With binary files, Git must replace the whole file for each new version.)
- Files updated frequently and/or files updated independent of LogStream.
- Files replicated on many Worker Nodes.

About the MaxMind GeoLite Example

We'll illustrate this with an example that often combines all three conditions: setting up the free, popular MaxMind GeoLite2 City database to support LogStream's GeoIP lookup Function. This example anticipates a LogStream production distributed deployment, where the GeoLite database is updated nightly across multiple Workers.

This example includes complete instructions for this particular setup. However, you can generalize the example to other MaxMind databases, and to other large lookup files – including large .csv 's that similarly receive frequent updates.

Reducing Deploy Traffic

The general approach for handling large lookups is:

- Do not place these files in the standard \$CRIBL_HOME/data/lookups.
- Instead, place them in a \$CRIBL_HOME subdirectory that's excluded from Git version control, through inclusion in the \$CRIBL_HOME/.gitignore file. Deploying the files to the Master Node and all desired Workers will require a manual procedure and will be required for the initial deployment as well as subsequent updates.

The example below uses \$CRIBL_HOME/state subdirectory, which is already listed in the default .gitignore file that ships with LogStream.

i If you prefer, you can use a different path, including a path outside \$CRIBL_HOME . If you choose this alternative, be sure to add that path to .gitignore .

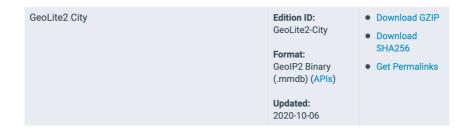
However, Cribl recommends using a \$CRIBL_HOME subdirectory like \$CRIBL_HOME/state, because this inherits appropriate permissions and simplifies backup/restore operations.

Let's move on to the MaxMind GeoLite specifics.

Download and Extract the Database

To enable the GeoIP Function using the MaxMind GeoLite 2 City database, your first steps are:

- 1. Create a free MaxMind account, at the page linked above.
- 2. Log in to your MaxMind account portal and select **Download Databases**.
- 3. On the Download page, look for the database you want. (In this example, you'd locate the **GeoLite2 City** section.) Note the **Format: GeoIP2 Binary**, and select **Download GZIP**.



- 4. Extract the archive to your local system.
- 5. Change to the directory created when you extracted the archive. This directory's name will correspond to the date you downloaded the file, so in the above 2020-10-06 example, you would use: \$ cd GeoLite2-City_20201006

Copy the Database File to the Master and Worker Nodes (Recommended)

In distributed deployments, Cribl recommends copying the MaxMind database separately to the Master and all Worker Nodes, e.g.. placing it in the \$CRIBL_HOME/state path. This will minimize the Git commit/deploy overhead around nightly updates to the binary database file.

Once in the database's directory, execute commands of this form:

```
$ scp *.mmdb <user>@<master-node>:
$ scp *.mmdb <user>@<worker-node>:
```

△ Copy the file to each Worker in the Worker Group where you intend to use LogStream's GeoIP Function.

The above commands copy the .mmdb database file into your user's home directory on each Node. Next, we'll move it to \$CRIBL_HOME/state on each Node. Execute these commands on both the Master and Worker Nodes:

```
$ sudo mv ~/*.mmdb <$CRIBL_HOME>/state/
$ sudo chown -R cribl:cribl <$CRIBL_HOME>/state/
```

Now that the database is in place, your Pipelines can use the GeoIP Function to enrich data. In the Function's **GeoIP file (.mmdb)** field, insert the complete \$CRIBL_HOME/state/<filename>.mmdb file path.

Copy the Database File Only to the Master (Alternative)

In smaller deployments, you might choose to copy this MaxMind database only to Master Node, and to let Workers receive updates via Git commit/deploy. In this case, the final commands above might look like this:

Shell

```
$ sudo cp ~/*.mmdb /opt/cribl/groups/<group-name>/data/lookups/
$ cd /opt/cribl/groups/<group-name>/data/lookups/
$ sudo chown cribl:cribl *.mmdb
```

Automatic Updates to the MaxMind Database

To set up automatic updates, see MaxMind's Automatic Updates for GeoIP2 and GeoIP Legacy Databases documentation. You'll need two modifications specific to LogStream:

- This must be set up on the Master, and on each Worker in any Group using GeoIP lookups.

Memory Considerations

Storage aside, large lookup files can also require additional RAM on each Worker Node that processes the lookups. For details, see Memory Sizing for Large Lookups.

TROUBLESHOOTING

Diagnosing Issues

To help diagnose LogStream problems, you can share a diagnostic bundle with Cribl Support. The bundle contains a snapshot of configuration files and logs at the time the bundle was created, and gives troubleshooters insights into how LogStream was configured and operating at that time.

What's in the Diagnostic Bundle

The following directories (and their contents) off of \$CRIBL_HOME are included:

- /default/*
- /local/*
- /log/*
- /groups/*
- /state/jobs/* includes only the latest 10 task from the latest 10 jobs.

Creating and Exporting a Diagnostic Bundle

Users can create and **securely** share bundles with Cribl Support either from the UI or from the CLI. In either case, you'll need outbound internet access to https://diag-upload.cribl.io and a valid Support Case number.

Using the UI

To create a bundle, go to **Settings > Diagnostics > Diagnostic Bundle** and click **Create diagnostic bundle**.

- To download the bundle locally to your machine, click **Export**.
- To share the bundle with Cribl Support, toggle **Send to Cribl Support** to **Yes**, enter your case number, and then click **Export**.

You can create a bundle from individual workers if you have the Worker UI access setting enabled. Go to Workers > <worker-name> > System Settings > Diagnostics > Diagnostic Bundle, and click Create Diagnostic Bundle.

Previously created bundles are stored in $\c SCRIBL_HOME/diag$. They're also listed in the UI, where you can re-download them or share them with Cribl Support.

Using the CLI

To create a bundle using the CLI, use the diag command.

```
diag command CLI
# $CRIBL_HOME/bin/cribl diag
Usage: [sub-command] [options] [args]
Commands:
get - List existing Cribl LogStream diagnostic bundles
create - Creates diagnostic bundle for Cribl LogStream
send - Send LogStream diagnostic bundle to Cribl Support, args:
  -c <caseNumber> - Cribl Case Number
  [-p <path>]
              - Diagnostic bundle path (if empty, then new bundle will be created)
## Creating a diagnostic bundle
# $CRIBL_HOME/bin/cribl diag create
Created Cribl LogStream diagnostic bundle at /opt/cribl/diag/cribl-logstream ← hostname ← <datetir
## Creating and sending a diagnostic bundle
# $CRIBL_HOME/bin/cribl diag send -c 420420
Sent LogStream diagnostic bundle to Cribl Support
## Sending a previously created diagnostic bundle
# $CRIBL_HOME/bin/cribl diag send -p /opt/cribl/diag/cribl-logstream—hostname>-<datetime>.tar.§
Sent LogStream diagnostic bundle to Cribl Support
```

Known Issues

2021-02-10 – v.2.4.0, 2.4.1 – With Splunk HEC Source, JSON payloads containing embedded objects trigger high CPU usage

Problem: Splunk HEC JSON payloads containing nested objects trigger high CPU usage, due to a flaw in JSON parsing.

Workaround: If you encounter this problem, rolling back to LogStream 2.3.4 resolves it. **Fix**: In LogStream 2.4.2.

2021-01-30 - v.2.4.0 - Worker Nodes cannot connect to Master

Problem: Worker Nodes cannot connect to the Master after the Master is upgraded to v.2.4.0.

Workaround: Disable compression on the Workers. You can do so through the Workers' UI at System Settings > Distributed Settings > Master Settings > Compression, or by commenting out this line in each Worker's cribl.yml config file:

compression: gzip

Fix: In LogStream 2.4.1.

2021-01-25 – v.2.4.0 – S3 collection stops working due to auth secret key issues.

Problem: S3 collection stops after upgrade to 2.4.0 due to secret key re-encryption. **Workaround:** Re-configure S3, save and re-deploy.

Fix: In LogStream 2.4.1.

2021-01-14 – v.2.4.0 – Google Cloud Storage Destination Needs Extra Endpoint to Initialize

Problem: The Google Cloud Storage Destination fails to initialize, displaying an error of the form: Bucket does not exist!

Workaround: In the outputs.yml file, under your cribl-gcp-bucket key endpoint, add: https://storage.googleapis.com.(in a single-instance deployment, locate this file at \$CRIBL_HOME/local/cribl/outputs.yml. In a distributed deployment, locate it

at \$CRIBL_HOME/groups/<group name>/local/cribl/outputs.yml.)
Fix: In LogStream 2.4.1.

2021-01-14 – v.2.4.0 – Worker Groups' Settings > Access Management Is Absent from UI

Problem: In this release, the Worker Groups > <group-name> > System Settings UI did not display the expected Access Management, Authentication, and Local Users sections. Workaround: Manually edit the users.json file.

Fix: In LogStream 2.4.1.

2021-01-13 – v.2.4.0 – Route Filters Aren't Copied to Capture Modal

Problem: On the **Routes** page, selecting **Capture New** in the right pane does not copy custom **Filter** expressions to the resulting **Capture Sample Data** modal. That modal's **Filter Expression** field always defaults to true.

Workarounds: 1. Bypass the Capture New button. Instead, from the Route's own
••• (Options) menu, select Capture. This initiates a capture with the Filter Expression
correctly populated. 2. Copy/paste the expression into the Capture Sample Data modal's
Filter Expression field. Or, if the expression is displayed in that field's history drop-down,
retrieve it.

Fix: In LogStream 2.4.1.

2021-01-13 – v.2.4.0 – Destinations' Documentation Doesn't Render from UI

Problem: Clicking the **Help Help D?** link in a Destination's configuration modal displays the error message: "Unable to load docs. Please check LogStream's online documentation instead."

Workarounds: 1. Go directly to the online Destinations docs, starting here. 2. Follow the UI link to the docs landing page, click through to open or download the current PDF, and scroll to its Destinations section.

Fix: In LogStream 2.4.1.

2021-01-13 – v.2.4.0 – Esc Key Doesn't Consistently Close Modals

Problem: Pressing Esc with focus on a modal's drop-down or slider doesn't close the modal as expected. (Pressing Esc with focus on a free-text field, combo box, or nothing does close the modal – displaying a confirmation dialog first, if you have unsaved changes.)

Workarounds: Click the **X** close box at upper right, or click **Cancel** at lower right. **Fix**: Planned for LogStream 2.4.2.

2020-12-17 – v.2.3.0+ – Free-License Expiration Notice, Blocked Inputs

Problem: LogStream reports an expired Free license, and blocks inputs, even though Free licenses in v.2.3.0 do not expire.

Workaround: This is caused by time-limited Free license key originally entered in a LogStream version prior to 2.3.0. Go to **Settings > Licensing**, click to select and expand your expired Free license, and click **Delete license**. LogStream will recognize the new, permanent Free license, and will restore throughput.

Fix: In LogStream 2.4.1.

2020-11-14 – v.2.3.3 – Null Fields Redacted in Preview, but Still Forwarded

Problem: Where event fields have null values, LogStream (by default) displays them as struck-out in the right Preview pane. The preview is misleading, because the events are still sent to the output.

Workaround: If you do want to prevent fields with null values from reaching the output, use an Eval Function, with an appropriate Filter expression, to remove them.

Fix: Preview corrected in LogStream 2.3.4.

2020-10-27 – v.2.3.2 – Cannot Name or Save New Event Breaker Rule

Problem: After clicking **Add Rule** in a new or existing Event Breaker Ruleset, the **Event Breaker Rule** modal's **Rule Name** field is disabled. Because **Rule Name** is mandatory field, this also disables saving the Rule via the **OK** button. **Fix**: In LogStream 2.3.3.

2020-10-12 – v.2.3.1 – Deleting One Function Deletes Others in Same Group

Problem: After inserting a new Function into a group and saving the Pipeline, deleting the Function also deletes other Functions lower down in the same group.

Fix: In LogStream 2.3.2.

Workaround: Move the target Function out of the group, resave the Pipeline, and only then delete the Function.

2020-09-27 – v.2.3.1 – Enabling Boot Start as Different User Fails

Problem: When a root user tries to enable boot-start as a different user (e.g., using /opt/cribl/bin/cribl boot-start enable -u <some-username>), they receive an error of this form:

error: found user=0 as owner for path=/opt/cribl, expected uid=NaN.

Please make sure CRIBL_HOME and its contents are owned by the uid=NaN by running:

[sudo] chown -R NaN:[\$group] /opt/cribl

Fix: In LogStream 2.3.2.

Workaround: Install LogStream 2.2.3 (which you can download here), then upgrade to 2.3.1.

2020-09-17 – v.2.3.0 – Worker Groups menu tab hidden after upgrade to LogStream 2.3.0

Problem: Upon upgrading an earlier, licensed LogStream installation to v. 2.3.0, the **Worker Groups** tab might be absent from the Master Node's top menu.

Fix: In LogStream 2.3.1.

Workaround: Click the **Home > Worker Groups** tile to access Worker Groups.

2020-09-17 – v.2.3.0 – Cannot Start LogStream 2.3.0 on RHEL 6, RHEL 7

Problem: Upon upgrading to v. 2.3.0, LogStream might fail to start on RHEL 6 or 7, with an error message of the following form. This occurs when the user running LogStream doesn't match the LogStream binary's owner. LogStream 2.3.0 applies a restrictive permissions check using id -un <uid>, which does not work with the version of id that ships with these RHEL releases.

id: 0: No such user

ERROR: Cannot run command because user=root with uid=0 does not own executable

Fix: In LogStream 2.3.1.

Workaround: Update your RHEL environment's id version, if possible.

2020-09-17 – v.2.3.0 – Cannot Start LogStream 2.3.0 with OpenId Connect

Problem: Upon upgrading an earlier LogStream installation to v. 2.3.0, OIDC users might be unable to restart the LogStream server.

Fix: In LogStream 2.3.1.

Workaround: Edit \$CRIBL_HOME/default/cribl/cribl.yml to add the following lines to its the auth section:

filter_type: email_whitelist
scope: openid profile email

2020-06-11 – v.2.1.x – Can't switch from Worker to Master Mode

Problem: In a Distributed deployment, attempting to switch Distributed Settings from Worker to Master Mode blocks with a spurious "Git not available...Please install and try again" error message.

Fix: In LogStream 2.3.0.

Workaround: To initialize git, switch first from Worker to Single mode, and then from Single to Master mode.

2020-05-19 - v.2.1.x - Login page blocks

Problem: Entering valid credentials on the login page (e.g., http://localhost:9000/login) yields only a spinner.

Fix: In LogStream 2.3.0.

Workaround: Trim /login from the URL.

2020-02-22 - v.2.1.x - Deleting resources in default/

Problem: In a Distributed deployment, deleting resources in default/ causes them to reappear on restart.

Workaround/Fix: In progress.

2019-10-22 – v. 2.0 – In-product upgrade issue on v2.0

Problem: Using in-product upgrade feature in v.1.7 (or earlier) fails to upgrade to v2.0, due to package-name convention change.

Workaround/Fix: Download the new version and upgrade per steps laid out here.

2019-08-27 – v.1.7 – In-product upgrade issue on v1.7

Problem: Using in-product upgrade feature in v1.6 (or earlier) fails to upgrade to v1.7 due to package name convention change.

Workaround/Fix: Download the new package and upgrade per steps laid out here.

2019-03-21 - v.1.4 - S3 stagePath issue on upgrade to v.1.4+

Problem: When upgrading from v1.2 with a S3 output configured, stagePath was allowed to be undefined. In v.1.4+, stagePath is a required field. This might causing

schema violations when upgrading older configs.

Workaround/Fix: Reconfigure the output with a valid stagePath filesystem path.

THIRD-PARTY SOFTWARE

Credits

Various components in Cribl LogStream are built and enhanced with software under free or open source licenses. We thank those projects' contributors!

@azure-storage-blob-10.3.0 ag-grid-community-19.1.2

ag-grid-react-19.1.2

ajv-6.9.2

ajv-errors-1.0.1

antd-3.13.0

as-table-1.0.36

avsc-5.4.9

aws-sdk-2.323.0

cidr-matcher-1.0.5

classnames-2.2.6

color-hash-1.0.3

d3-time-format-2.1.3

date-fns-1.29.0

diff-3.5.0

diff2html-2.11.3

echarts-4.3.0

echarts-4.6.0

escodegen-1.11.1

esprima-4.0.1

express-4.16.3

fast-bitset-1.3.2

file-saver-1.3.8

http-proxy-agent-3.0.0

https-proxy-agent-4.0.0

jwt-simple-0.5.6

kafkajs-1.11.0

kafkajs-1.4.5

kafkajs-snappy-1.1.0

ldapts-1.10.0

limiter-1.1.4

lodash-4.17.15

lz4js-0.2.0

maxmind-3.1.2

node-cache-4.2.0

node-uuid-1.4.8

numeral-2.0.6

pako-1.0.10

papaparse-5.0.0-beta.0

pbf-3.2.1

proxy-from-env-1.0.0

query-string-6.1.0

react-16.7.0

react-dom-16.7.0

react-grid-layout-0.16.6

react-router-dom-4.3.1

react-sortable-hoc-0.8.3

react-split-pane-0.1.82

regexpp-2.0.0

requirejs-2.3.6

resize-observer-polyfill-1.5.0

rxjs-6.2.2

saxen-8.1.0

simple-git-1.126.0

snappyjs-0.6.0

snmp-native-1.2.0

streamcount-1.0.1

tar-stream-1.6.1

url-0.11.0

winston-3.0.0

xmlbuilder-10.0.0

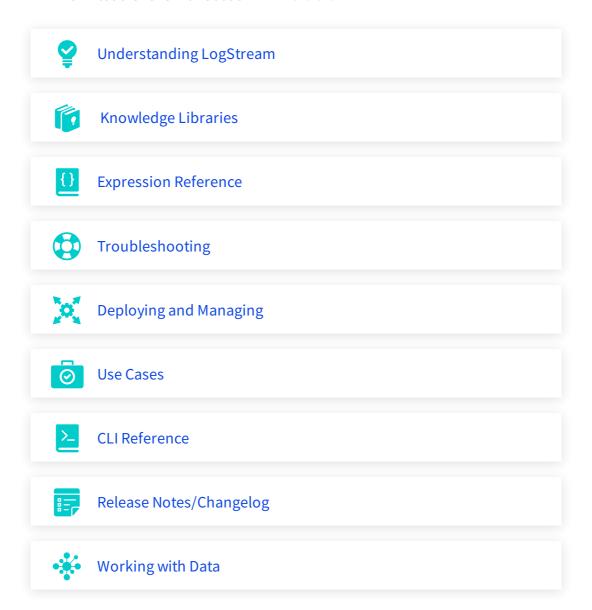
yaml-1.3.2

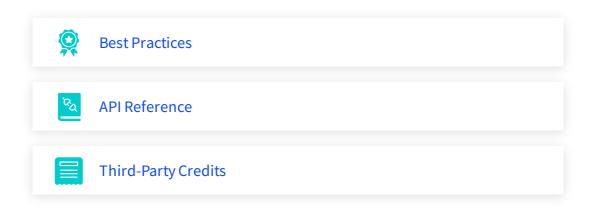
Introduction

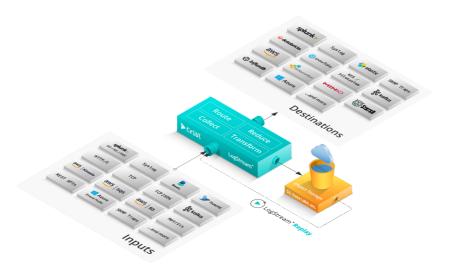
Welcome

Looking for Cribl LogStream documentation? We've goat you covered!

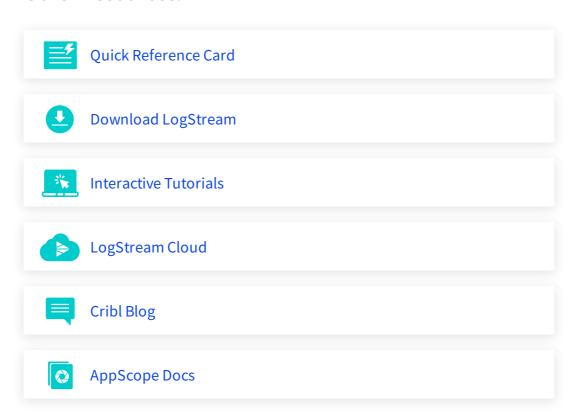
- Questions not answered here? We'd love to help you. Meet us in #Cribl Community Slack – sign up here.
- Download entire manual as PDF v.2.4.1.







Other Resources:



Page 637 of 1257

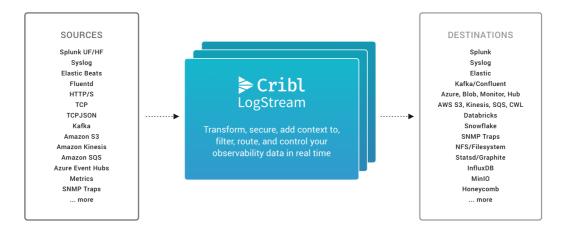
About Cribl LogStream

Getting started with Cribl LogStream

What Is Cribl LogStream?

Cribl LogStream helps you process machine data – logs, instrumentation data, application data, metrics, etc. – in real time, and deliver them to your analysis platform of choice. It allows you to:

- Add context to your data, by enriching it with information from external data sources.
- Help secure your data, by redacting, obfuscating, or encrypting sensitive fields
- Optimize your data, per your performance and cost requirements.



Cribl LogStream ships in a single, no-dependencies package. It provides a refreshing and modern interface for working with and transforming your data. It scales with – and works inline with – your existing infrastructure, and is transparent to your applications.

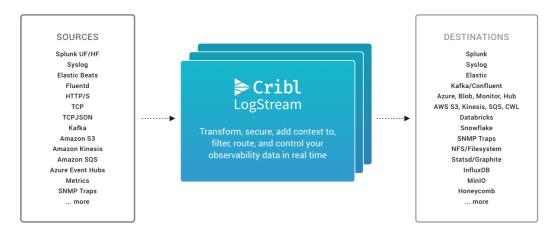
Who Is Cribl LogStream For?

Cribl LogStream is built for administrators, managers, and users of operational and security intelligence products and services.

Basic Concepts

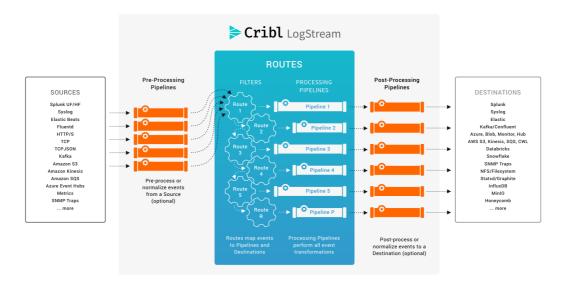
Notable features and concepts to get a fundamental understanding of Cribl LogStream

As we describe features and concepts, it helps to have a mental model of Cribl LogStream as a system that receives events from various sources, processes them, and then sends them to one or more destinations.



Sources, LogStream, destinations

Let's zoom in on the center of the above diagram, to take a closer look at the processing and transformation options that LogStream provides internally. The basic interface concepts to work with are Routes, which manage data flowing through Pipelines, which consist of Functions.



Routes, Pipelines, Functions

Routes

Routes evaluate incoming events against filter expressions to find the appropriate Pipeline to send them to. Routes are **evaluated in order**. A Route can be associated **with only one** Pipeline and one output. By default, a Route-Pipeline-Output tuple will consume matching events.

If the Route's Final flag is disabled, one or more event **clones** are sent down the associated Pipeline, while the original event continues down the rest of the Routes. This is very useful in cases where the same set of events needs to be processed in multiple ways and delivered to different destinations. For more details, see Routes.

Pipelines

A series of Functions is called a **Pipeline**, and the order in which the Functions are executed matters. Events are delivered to the beginning of a pipeline by a Route, and as they're processed by a Function, the events are passed to the next Function down the line.



Events only move forward – toward the end of the Pipeline, and eventually out of the system. For more details, see Pipelines.

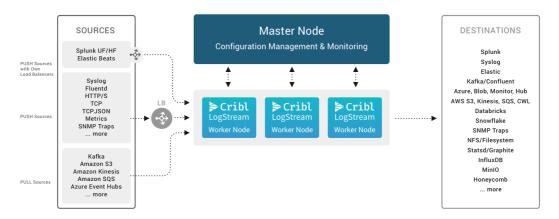
Functions

At its core, a **Function** is a piece of code that executes on an event, and that encapsulates the smallest amount of processing that can happen to that event. For instance, a very simple Function can be one that replaces the term foo with bar on each event. Another one can hash or encrypt bar. Yet another function can add a field – say, dc=jfk-42 – to any event with source=*us-nyc-application.log.

Functions process each event that passes through them. To help improve performance, functions can optionally be configured with filters, to limit their processing scope to matching events only. For more details, see Functions.

A Scalable Model

You can scale LogStream up to meet enterprise needs in a distributed deployment. Here, multiple LogStream Workers (instances) share the processing load. But as you can see in the preview schematic below, even complex deployments follow the same basic model outlined above.



Distributed deployment architecture

Getting Started Guide

This guide walks you through planning, installing, and configuring a single-instance deployment of Cribl LogStream. You'll capture some realistic sample log data, and then use LogStream's built-in Functions to redact, parse, refine, and shrink the data.

By the end of this guide, you'll have assembled all of LogStream's basic building blocks: a Source, Route, and Pipeline, several Functions, and a Destination. You can complete this tutorial using LogStream's included sample data, without connections to – or licenses on – any inbound or outbound services.

Assuming a cold start (from initial LogStream download and installation), this guide might take about an hour. But you can work through it in chunks, and LogStream will persist your work between sessions.

☐ If you've already downloaded, installed, and launched LogStream, skip ahead to Add a Source.

Requirements for this Tutorial

The minimum requirements for running this tutorial are the same as for a LogStream production single-instance deployment.

OS

• Linux: RedHat, CentOS, Ubuntu, Amazon Linux (64bit)

System

- +4 physical cores = +8 vCPUs; +8GB RAM all beyond your basic OS/VM requirements
- 5GB free disk space (more if persistent queuing is enabled)

i We assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). For details, see Recommended AWS, Azure, and GCP Instance Types.

Browser Support

• Firefox 65+, Chrome 70+, Safari 12+, Microsoft Edge

Network Ports

By default, LogStream listens on the following ports:

Component	Default Port
UI default	9000
HTTP Inbound, default	10080
User options	+ Other data ports as required.

You can override these defaults as needed.

Plan for Production

For higher processing volumes, users typically enable LogStream's Distributed Deployment option. While beyond the scope of this tutorial, that option has a few additional requirements, which we list here for planning purposes:

- Port 9000 or 4200 must be available on the Master Node for Workers' communications.
- Git (1.8.3.1 or higher) must be installed on the Master Node, to manage configuration changes.

See Sizing and Scaling for further details about configuring LogStream to handle large data streams.

Download and Install LogStream

Download the latest version of LogStream at https://cribl.io/download/.

Un-tar the resulting .tgz file in a directory of your choice (e.g., /opt/). Here's general syntax, and a specific example:

```
tar xvzf cribl → version → <build → <arch > .tgz
tar xvzf cribl -2.3.1-1d4e05c5-linux-x64.tgz
```

You'll now have LogStream installed in a cribl subdirectory, by default: /opt/cribl/. We'll refer to this cribl subdirectory throughout this documentation as \$CRIBL_HOME.

Run LogStream

In your terminal, switch to the \$CRIBL_HOME/bin directory (e.g,: /opt/cribl/bin). Here, you can start, top, and verify the LogStream server using these basic ./cribl CLI commands:

Start: ./cribl startStop: ./cribl stop

• Get status: ./cribl status

☐ For other available commands, see CLI Reference.

Next, in your browser, open http://chostname:9000 (e.g., http://localhost:9000) and log in with default credentials (admin, admin).

Register your copy of LogStream if desired.

After registering, you'll be prompted to change the default password.

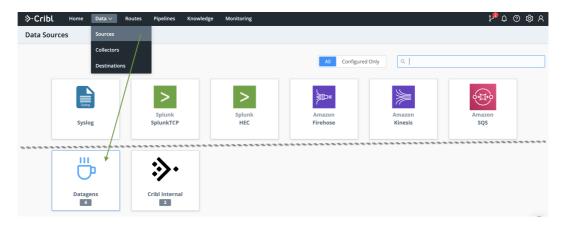
You are now ready to configure a working LogStream installation – with a Source, Destination, Pipeline, and Route – and to assemble several built-in Functions to refine sample log data.

Get Data Flowing

Add a Source

Each LogStream Source represents a data input. Options include Splunk, Elastic Beats, Kinesis, Kafka, syslog, HTTP, TCP JSON, and others.

For this tutorial, we'll enable a LogStream built-in datagen (i.e., data generator) that generates a stream of realistic sample log data.



Addiing a datagen Source

- 1. From LogStream's top menu, select **Data > Sources**.
- 2. From the **Data Sources** page's tiles or left menu, select **Datagens**.

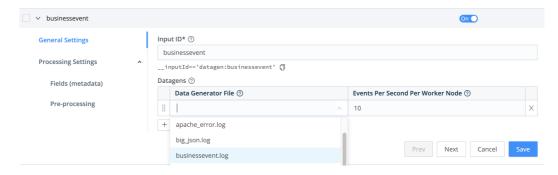
(You can use the search box to jump to the **Datagens** tile.)

- 3. Click Add New to open the New Datagen source pane.
- 4. In the Input ID field, name this Source businessevent.
- 5. In the Data Generator File drop-down, select businessevent.log.

This generates...log events for a business scenario. We'll look at their structure shortly, in Capture and Filter Sample Data.

6. Click Save.

The **On** slider in the **Enabled** column indicates that your datagen Source has started generating sample data.



Configuring a datagen Source

Add a Destination

Each LogStream Destination represents a data output. Options include Splunk, Kafka, Kinesis, InfluxDB, Snowflake, Databricks, TCP JSON, and others.

For this tutorial, we'll use LogStream's built-in **DevNull** Destination. This simply discards events – not very exciting! But it simulates a real output, so it provides a configuration-free quick start for testing LogStream setups. It's ideal for our purposes.

To verify that **DevNull** is enabled, let's walk through setting up a Destination, then setting it up as LogStream's default output:

- 1. From LogStream's top menu, select **Data > Destinations**.
- 2. Select **DevNull** from the **Data Destinations** page's tiles or left menu.

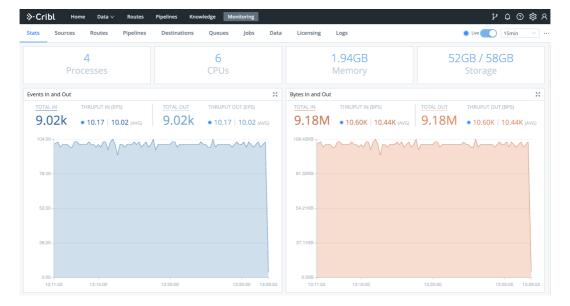
(You can use the search box to jump to the **DevNull** tile.)

- 3. On the resulting **devnull** row, look for the **Live** indicator under **Enabled**. This confirms that the **DevNull** Destination is ready to accept events.
- 4. From the **Data Destinations** page's left nav, select the **Default** Destination at the top.
- 5. On the resulting Manage Default Destination page, verify that the Default Output ID drop-down points to the devnull Destination we just examined.

We've now set up data flow on both sides. Is data flowing? Let's check.

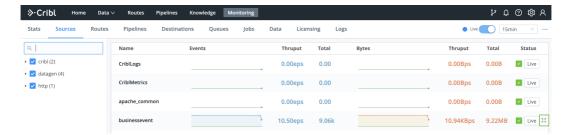
Monitor Data Throughput

From the top menu, select **Monitoring**. This opens a summary dashboard, where you should see a steady flow of data in and out of LogStream. The left graph shows events in/out. The right graph shows bytes in/out.



Monitoring dashboard

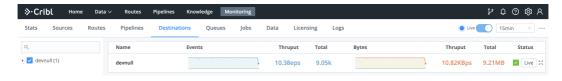
Monitoring displays data from the preceding 24 hours. You can use the Monitoring submenu to open detailed displays of LogStream components, collection jobs and tasks, and LogStream's own internallogs. Click Sources on the lower (white) submenu to switch to this view:



Monitoring Sources

This is a compact display of each Source's inbound events and bytes as a sparkline. You can click each Source's Expand button (highlighted at right) to zoom up detailed graphs.

Click **Destinations** on the lower submenu. This displays a similar sparklines view, where you can confirm data flow out to the devnull Destination:



Monitoring Destinations

With confidence that we've got data flowing, let's send it through a LogStream Pipeline, where we can add Functions to refine the raw data.

Create a Pipeline

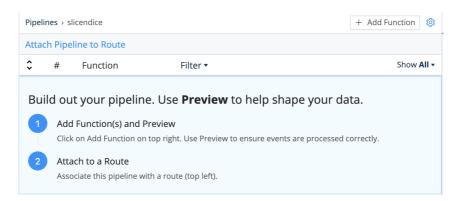
A Pipeline is a stack of LogStream Functions that process data. Pipelines are central to refining your data, and also provide a central LogStream workspace – so let's get one going.

1. From the top menu, select **Pipelines**.

You now have a two-pane view, with business on the left and party on the right a Pipelines list on the left and Sample Data controls on the right. (We'll capture some sample data momentarily.)

- 2. At the **Pipelines** pane's upper right, click **+ Add Pipeline**, then select **Create Pipeline**.
- 3. In the new Pipeline's **ID** field, enter a unique identifier. (For this tutorial, you might use slicendice.)
- 4. Optionally, enter a **Description** of this Pipeline's purpose.
- 5. Click Save.

Your empty Pipeline now prompts you to preview data, add Functions, and attach a Route. So let's capture some data to preview.



Pipeline prompt to add Functions

Capture and Filter Sample Data

The right **Sample Data** pane provides multiple tools for grabbing data from multiple places (inbound streams, copy/paste, and uploaded files); for

previewing and testing data transformations as you build them; and for saving and reloading sample files.

Since we've already got live (simulated) data flowing in from the datagen Source we built, let's grab some of that data.

Capture New Data

- 1. In the right pane, click Capture New.
- 2. In the **Capture Sample Data** modal, immediately change the generated **File Name** to a name you'll recognize, like be_raw.log.
- 3. Click Capture, then accept the drop-down's defaults click Start.
- 4. Click **Save as Sample File**. This saves to the **File Name** you entered above. You're now previewing the events in the right pane. (Note that this pane's **Preview Simple** tab now has focus.)
- 5. Click **Show more** to expand one or more events.

By skimming the key-value pairs within the data's _raw fields, you'll notice the scenario underlying this preview data (provided by the businessevents.log datagen): these are business logs from a mobile-phone provider.

To set up our next step, find at least one marketState K=V pair. Having captured and examined this raw data, let's use this K=V pair to crack open LogStream's most basic data-transformation tool, Filtering.

Filter Data and Manage Sample Files

- 1. Click the right pane's **Sample Data** tab.
- 2. Again click Capture New.
- 3. In the **Capture Sample Data** modal, replace the **Filter Expression** field's default true value with this simple regex:

```
_raw.match(/marketState=TX/)
```

We're going to Texas! If you type this in, rather than pasting it, notice how LogStream provides typeahead assist to complete a well-formed JavaScript expression.

You can also click the Expand button at the Filter Expression field's right

edge to open a modal to validate your expression. The adjacent dropdown enables you to restore previously used expressions



4. Click Capture, then Start.

Using the **Capture** drop-down's default limits of 10 seconds and 10 events, you'll notice that with this filter applied, it takes much longer for LogStream to capture 10 matching events.

- 5. Click Cancel to discard this filtered data and close the modal.
- 6. On the right pane's **Sample Data** tab, click **Simple** beside be_raw.log.

This restores our preview of our original, unfiltered capture. We're ready to transform this sample data in more interesting ways, by building out our Pipeline's Functions.

Refine Data with Functions

Functions are pieces of JavaScript code that LogStream invokes on each event that passes through them. By default, this means all events – each Function has a **Filter** field whose value defaults to true. As we just saw with data capture, you can replace this value with an expression that scopes the Function down to particular matching events.

In this Pipeline, we'll use some of LogStream's core Functions to:

- Redact (mask) sensitive data
- Extract (parse) the _raw field's key-value pairs as separate fields.
- Add a new field.
- Delete the original _raw field, now that we've extracted its contents.
- Rename a field for better legibility..

Mask: Redact Sensitive Data

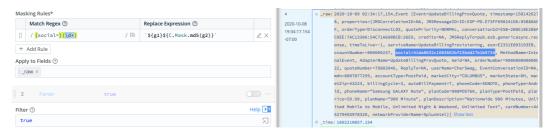
In the right **Preview** pane, notice each that event includes a **social** key, whose value is a (fictitious) raw Social Security number. Before this data goes any further through our Pipeline, let's use LogStream's Mask Function to swap in an md5 hash of each SSN.

- 1. In the left Pipelines pane, click + Add Function.
- 2. Search for Mask, then click it.
- 3. In the new Function's **Masking Rules**, click the into **Match Regex** field.
- 4. Enter or paste this regex, which simply looks for digits following social = : (social =)(\d+)
- 5. In **Replace Expression**, paste the following hash function. The backticks are literal: `\${g1}\${C.Mask.md5(g2)}`
- 6. Note that **Apply to Fields** defaults to _raw . This is what we want to target, so we'll accept this default.
- 7. Click Save.

You'll immediately notice some obvious changes:

- The **Preview** pane has switched from its **IN** to its **OUT** tab, to show you the outbound effect of the Pipeline you just saved.
- Each event's _raw field has changed color, to indicate that it's undergone some redactions.

Now locate at least one event's **Show more** link, and click to expand it. You can verify that the social values have now been hashed.



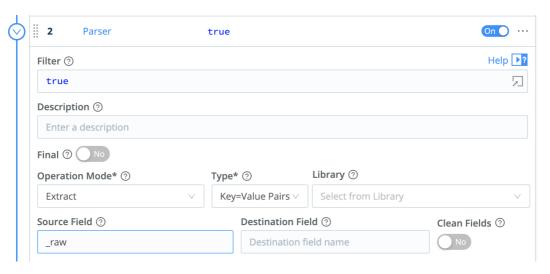
Mask Function and hashed result

Parser: Extract Events

Having redacted sensitive data, we'll next use a Parser function to lift up all the _raw field's key-value pairs as fields:

- 1. In the left Pipelines pane, click + Add Function.
- 2. Search for Parser, then click it.

- 3. Leave the Operation Mode set to its Extract default.
- 4. Set the Type to Key=Value Pairs.
- 5. Leave the **Source Field** set to its _raw default.
- 6. Click Save.



Parser configured to extract K=V pairs from _raw

You should see the **Preview** pane instantly light up with a lot more fields, parsed from _raw . You now have rich structured data, but not all of this data is particularly interesting: Note how many fields have NA ("Not Applicable") values. We can enhance the **Parser** Function to ignore fields with NA values.

 In the Function's Fields Filter Expression field (near the bottom), enter this negation expression: value≠'NA'

Note the single-quoted value. If you type (rather than paste) this expression, watch how typeahead matches the first quote you type.

2. Click **Save**, and watch the **Preview** pane.



Filtering the Parser Function to ignore fields with 'NA' values

Several fields should disappear – such as credits, EventConversationID, and ReplyTo. The remaining fields should display meaningful values.

Congratulations! Your log data is already starting to look better-organized and less bloated.

□ Missed It?

If you didn't see the fields change, slide the Parser Function **Off**, click **Save** below, and watch the **Preview** pane change. Using these toggles, you can preserve structure as you test and troubleshoot each Function's effect.

Note that each Function also has a **Final** toggle, defaulting to **Off**. Enabling **Final** anywhere in the Functions stack will prevent data from flowing to any Functions lower in the UI.

Be sure to toggle the Function back **On**, and click **Save** again, before you proceed!



Toggling a Function off and on

Next, let's add an extra field, and conditionally infer its value from existing values. We'll also remove the _raw field, now that it's redundant. To add and remove fields, the **Eval** Function is our pal.

Eval: Add and Remove Fields

Let's assume we want to enrich our data by identifying the manufacturer of a certain popular phone handset. We can infer this from the existing phoneType field that we've lifted up for each event.

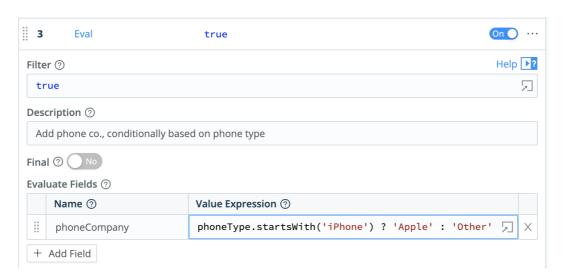
Add Field (Enrich)

- 1. In the left **Pipelines** pane, click + Add Function.
- 2. Search for Eval, then click it.

3. Click into the new Function's Evaluate Fields table.

Here you add new fields to events, defining each field as a key-value pair. If we needed more key-value pairs, we could click + Add Field for more rows.

- 4. In Name, enter: phoneCompany.
- 5. In Value Expression, enter this JS ternary expression that tests
 phoneType 's value:
 phoneType.startsWith('iPhone') ? 'Apple' : 'Other' (Note the ?
 and : operators, and the single-quoted values.)
- 6. Click **Save**. Examine some events in the **Preview** pane, and each should now contain a phoneCompany field that matches its phoneType.



Adding a field to enrich data

Remove Field (Shrink Data)

Now that we've parsed out all of the _raw field's data – it can go. Deleting a (large) redundant field will give us cleaner events, and reduced load on downstream resources.

- 1. Still in the **Eval** Function, click into **Remove Fields**
- 2. Type: _raw and press **Tab** or **Enter**.
- 3. Click Save.

The **Preview** pane's diff view should now show each event's _raw field stripped out.



Removing a field to streamline data

Our log data has now been cleansed, structured, enriched, and slimmed-down. Let's next look at how to make it more legible, by giving fields simpler names.

Rename: Refine Field Names

1. In the left **Pipelines** pane, click + Add Function.

This rhythm should now be familiar to you.

- 2. Search for Rename, then click it.
- 3. Click into the new Function's **Rename Fields** table.

This has the same structure you saw above in Eval: Each row defines a key-value pair.

- 4. In **Current Name**, enter the longhaired existing field name: conversationId.
- 5. In New Name, enter the simplified field name: ID.
- 6. Watch any event's conversationId field in the **Preview** pane as you click **Save** at left. This field should change to ID in all events.

Drop: Remove Unneeded Events

We've already refined our data substantially. To further slim it down, a Pipeline can entirely remove events that aren't of interest for a particular downstream service.

As the "Pipeline" name implies, your LogStream installation can have multiple Pipelines, each configured to send out a data stream tailored to a particular Destination. This helps you get the right data in the right places most efficiently.

Here, let's drop all events for customers who use prepaid monthly phone service (i.e., **not** postpaid):

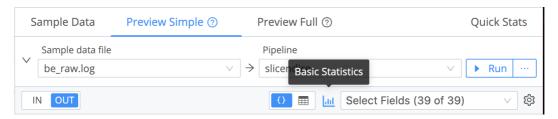
- 1. In the left **Pipelines** pane, click + Add Function.
- 2. Search for Drop, then click it.
- 3. Click into the new Function's Filter field.
- 4. Replace the default true value with this JS negation expression: accountType≠'PostPaid'
- 5. Click Save.

Now scroll through the right **Preview** pane. Depending on your data sample, you should now see multiple events struck out and faded – indicating that LogStream will drop them before forwarding the data.

A Second Look at Our Data

Torture the data enough, and it will confess. By what factor have our transformations refined our data's volume? Let's check.

In the right **Preview** pane, click the **Basic Statistics** button:



Displaying Basic Statistics

Even without the removal of the _raw field (back in Eval) and the dropped events, you should see a substantial % reduction in the Full Event Length.



Data reduction quantified

Woo hoo! Before we wrap up our configuration: If you're curious about individual Functions' independent contribution to the data reduction shown

here, you can test it now. Use the toggle **Off > Save > Basic Statistics** sequence to check various changes.

Add and Attach a Route

We've now built a complete, functional Pipeline. But so far, we've tested its effects only on the static data sample we captured earlier. To get dynamic data flowing through a Pipeline, we need to filter that data in, by defining a LogStream Route.

1. At the **Pipelines** page's top left, click **Attach Pipeline to Route**.

This displays the **Routes** page. It's structured very similarly to the **Pipelines** page, so the rhythm here should feel familiar.

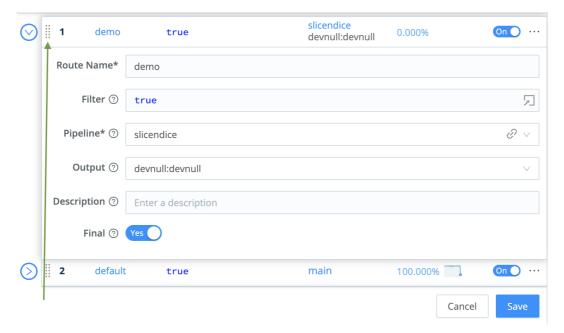
- 2. Click + Add Route.
- 3. Enter a unique, meaningful Route Name, like demo.
- 4. Leave the **Filter** field set to its true default, allowing it to deliver all events.

Because a Route delivers events to a Pipeline, it offers a first stage of filtering. In production, you'd typically configure each Route to filter events by appropriate source, sourcetype, index, host, _time, or other characteristics. The **Filter** field accepts JavaScript expressions, including AND (&&) and OR (||) operators.

- 5. Set the **Pipeline** drop-down to our configured slicendice Pipeline.
- 6. Set the Output drop-down to either devnull or default.

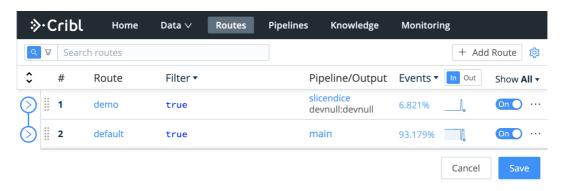
This doesn't matter, because we've set default as a pointer to devnull. In production, you'd set this carefully.

- 7. You can leave the **Description** empty, and leave **Final** set to **Yes**.
- 8. Grab the new Route by its left handle, and drag it above the default Route, so that our new Route will process events first. You should see something like the screenshot below.
- 9. Click **Save** to save the new Route to the routing table.



Configuring and adding a Route

The sparklines should immediately confirm that data is flowing through your new Route:



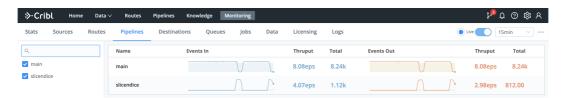
Live Routes

To confirm data flow through the whole system we've built, select **Monitoring > Routes** from LogStream's top menu and examine demo.



Monitoring data flow through Routes

Also select Monitoring > Pipelines and examine slicendice.



Monitoring data flow through Pipelines

What Have We Done?

Look at you! Give yourself a pat on the back! In this short, scenic tour – with no hit to your cloud-services charges – you've build a simple but complete LogStream system, exercising all of its basic components:

- Downloaded, installed, and run LogStream.
- Configured a Source to hook up an input.
- Configured a Destination to feed an output.
- Monitored data throughput, and checked it twice.
- Built a Pipeline.
- Configured LogStream Functions to redact, parse, enrich, trim, rename, and drop event data.
- Added and attached a Route to get data flowing through our Pipeline.

Next Steps

Interested in guided walk-throughs of more-advanced LogStream features? We suggest that next, you check out:

- LogStream Sandboxes: Work through general and specific scenarios in containers. with terminal access and free, hosted data inputs and outputs.
- Use Cases documentation: Bring your own services to build solutions to specific challenges.
- Cribl Concept: Pipelines Video showing how to build and use Pipelines at multiple LogStream stages.
- Cribl Concept: Routing Video about using Routes to send different data through different paths.

Cleaning Up

Oh yeah, you've still got the LogStream server running, with its businessevent.log datagen wtill firing events. If you'd like to shut these down for now, in reverse order:

- 1. Go to Data > Sources > Datagens.
- 2. Slide businessevent to **Off**, and click **Save**. (Refer back to the screenshot above.)
- 3. In your terminal's \$CRIBL_HOME/bin directory, shut down the server with: ./cribl stop

That's it! Enjoy using LogStream.

Deployment

Deployment Types

Deployment guide to get you started with Cribl

There are at least **two** key factors that will determine the type of Cribl LogStream deployment in your environment:

- Amount of Incoming Data: This is defined as the amount of data planned to be ingested per unit of time. E.g. How many MB/s or GB/day?
- Amount of Data Processing: This is defined as the amount of processing that will happen on incoming data. E.g., is most data passing through and just being routed? Or are there a lot of transformations, regex extractions, field encryptions? Is there a need for heavy re-serialization?

Single Instance Deployment

When volume is low and/or amount of processing is light, you can get started with a single instance deployment.

Distributed Deployment

To accommodate increased load, we recommend scaling up and perhaps out with multiple instances.

Splunk App Deployment

If you have an existing Splunk Heavy Forwarder infrastructure that you want to use, you can deploy Cribl App for Splunk. See the note below before you plan.

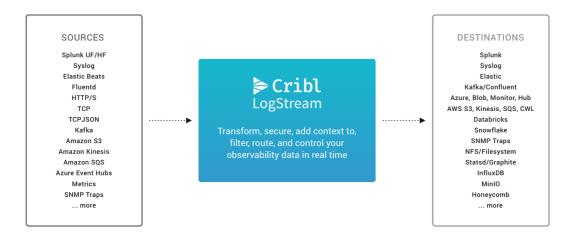
△ Cribl App for Splunk Deprecation Notice
Click here.

Single-Instance Deployment

Getting started with Cribl LogStream on a single instance

For small-volume or light processing environments – or for test or evaluation use cases – a single instance of Cribl LogStream might be sufficient to serve all inputs, event processing, and outputs. This page outlines how to implement a single-instance deployment.

Architecture



Requirements

- OS:
 - Linux: Red Hat, CentOS, Ubuntu, Amazon Linux (64bit)
 - Mac OS is no longer supported as of v. 2.3, due to LogStream's incorporation of Linux-native components.
- System:
 - +4 physical cores, +8GB RAM
 - 5GB free disk space (more if persistent queuing is enabled)
 - i We assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). All quantities listed above are minimum requirements. To fulfill the above

requirements using cloud-based virtual machines, see Recommended AWS, Azure, and GCP Instance Types.

• Browser Support: Firefox 65+, Chrome 70+, Safari 12+, Microsoft Edge

Network Ports

By default, LogStream listens on the following ports:

Component	Default Port
UI	9000
HTTP In	10080
Splunk to Cribl LogStream data port	localhost:10000 (Cribl App for Splunk)
criblstream Splunk search command to Cribl LogStream	localhost:10420 (Cribl App for Splunk)
User options	+ Other data ports as required.

Overriding Default Ports

The above ports can be overridden in the following configuration files:

- Cribl UI port (9000): Default definitions for host, port, and other settings are set in \$CRIBL_HOME/default/cribl/cribl.yml, and can be overridden by defining alternatives in \$CRIBL_HOME/local/cribl/cribl.yml.
- Data Ports: HTTP In (10080), TCPJSON in (10420) Splunk to Cribl (10000): Default definitions for host, port and other settings are set in \$CRIBL_HOME/default/cribl/inputs.yml, and can be overridden by defining alternatives in \$CRIBL_HOME/local/cribl/inputs.yml.

Installing on Linux

- Install the package on your instance of choice. Download it here.
- Ensure that required ports are available (see Network Ports).
- Un-tar in a directory of choice, e.g., /opt/:

Running

Go to the \$CRIBL_HOME/bin directory, where the package was extracted (e.g.: /opt/cribl/bin). Here, you can use ./cribl to:

- Start: ./cribl start
- Stop: ./cribl stop
- Reload: ./cribl reload
- Restart: ./cribl restart
- Get status: ./cribl status
- Switch a distributed deployment to single-instance mode:
 - ./cribl mode-single (uses the default address:port 0.0.0.0:9000)
 - i Executing the restart or stop command cancels any currently running collection jobs. For other available commands, see CLI Reference.

Next, go to http://<hostname>:9000 and log in with default credentials (admin:admin). You can now start configuring Cribl LogStream with Sources and Destinations, or start creating Routes and Pipelines.

i In the case of an API port conflict, the process will retry binding for 10 minutes before exiting.

Enabling Start on Boot

Cribl LogStream ships with a CLI utility that can update your system's configuration to start LogStream at system boot time. The basic format to invoke this utility is:

[sudo] \$CRIBL_HOME/bin/cribl boot-start [enable|disable] [options] [args]

i Boot-start is supported only on Linux. For options and arguments, see the CLI Reference.

Newer systems use systemd to start processes at boot, while older ones use initd.

Using systemd

To **enable** Cribl LogStream to start at boot time with **systemd**, you need to run the boot-start command. Make sure you first create any user you want to specify to run LogStream. E.g., to run LogStream on boot as existing user cribl, you'd use:

```
sudo $CRIBL_HOME/bin/cribl boot-start enable -m systemd -u cribl
```

This will install a unit file (as below) and start Cribl LogStream at boot time as user cribl. A -configDir option can be used to specify where to install the unit file. If not specified, this location defaults to /etc/systemd/system.

If necessary, change ownership for the Cribl LogStream installation:

```
[sudo] chown -R cribl $CRIBL_HOME
```

Next, use the enable command to ensure that the service starts on system boot:

```
[sudo] systemctl enable cribl
```

To **disable** starting at boot time, run the following command:

```
sudo $CRIBL_HOME/bin/cribl boot-start disable
```

Note the file's default 65536 limit on maximum open files. Cribl recommends going no higher than this value.

```
Installed systemd File
```

```
[Unit]
Description=Systemd service file for Cribl LogStream.
After=network.target
```

[Service]
Type=forking
User=cribl
Restart=on-failure
RestartSec=5
LimitNOFILE=65536
PIDFile=/install/path/to/cribl/pid/cribl.pid
ExecStart=/install/path/to/cribl/bin/cribl start
ExecStop=/install/path/to/cribl/bin/cribl stop
ExecStopPost='/bin/rm -f /install/path/to/cribl/pid/cribl.pid'
ExecReload=/install/path/to/cribl/bin/cribl reload
TimeoutSec=60

```
[Install]
WantedBy=multi-user.target
```

Using initd

To **enable** Cribl LogStream to start at boot time with **initd**, you need to run the boot-start command. If the user that you want to run LogStreams does not exist, create it prior to executing. E.g., running LogStream as user cribl on boot:

```
sudo $CRIBL_HOME/bin/cribl boot-start enable -m initd -u cribl
```

This will install an init.d script in /etc/init.d/cribl.init.d, and start Cribl LogStream at boot time as user cribl.A -configDir option can be used to specify where to install the script. If not specified, this location defaults to /etc/init.d.

If necessary, change ownership for the Cribl LogStream installation:

```
[sudo] chown -R cribl $CRIBL_HOME
```

To **disable** starting at boot time, run the following command:

```
sudo $CRIBL_HOME/bin/cribl boot-start disable
```

△ Do NOT Run LogStream as Root!

If LogStream is required to listen on ports 1–1024, it will need privileged access. On a Linux system with POSIX capabilities, you can achieve this by adding the CAP_NET_BIND_SERVICE capability. For example: # setcap cap_net_bind_service=+ep \$CRIBL_HOME/bin/cribl

On some OS versions (such as CentOS), you must add an -i switch to the setcap command. For example: # setcap -i cap_net_bind_service=+ep \$CRIBL_HOME/bin/cribl

Upon starting the LogStream server, a Port xxx is already in use error might indicate that setcap did not successfully execute.

System Proxy Configuration

You can direct all outbound HTTP/S requests to go through proxy servers. Initial configuration and changing these variables requires restarting LogStream on the affected nodes if the application is already running when the changes are

applied. You do so by setting a few environment variables before starting LogStream, as follows:

Configure the HTTP_PROXY and HTTPS_PROXY environment variables either with your proxy's IP address, or with a DNS name that resolves to that IP address. Optionally, follow either convention with a colon and the port number to which you want to send queries.

```
HTTP_PROXY examples:

$ export HTTP_PROXY=http://10.15.20.25:1234
$ export HTTP_PROXY=http://proxy.example.com:1234

HTTPS_PROXY examples:
```

```
$ export HTTPS_PROXY=http://10.15.20.25:5678
$ export HTTPS_PROXY=http://proxy.example.com:5678
```

i Case Conflicts

The environment variables' names can be either uppercase or lowercase. However, if you set duplicate versions of the same name, the lowercase version takes precedence. E.g., if you've set both HTTPS_PROXY and https_proxy, the IP address specified in https_proxy will take effect.

Proxy Confguration with systemd

If you are proxying outbound traffic with systemd, list your proxy environment variables in the systemd unit file's [Service] section by adding statements of this form:

```
Installed systemd File

[Service]
...
Environment=https_proxy=<yourproxy>
Environment=https_proxy=http://proxy.example.com:1234
Environment=https_proxy=http://10.10.1.1:8080
```

This will prevent LogStream from throwing "failed to send anonymized telemetry metadata" errors.

Authenticating on Proxies

You can use HTTP Basic authentication on HTTP or HTTPS proxies. Specify the user name and password in the proxy URL. For example:

```
$ export HTTP_PROXY=http://username:password@proxy.example.com:1234
$ export HTTPS_PROXY=http://username:password@proxy.example.com:5678
```

Bypassing Proxies with NO_PROXY

If you've set the above environment variables, you can negate them for specified (or all) hosts. Set the NO_PROXY environment variable to identify URLs that should bypass the proxy server, and instead be sent as direct requests. Use the following format:

```
$ export NO_PROXY="<list of hosts/domains>"
```

Usage notes:

- Within the list, separate the host/domain names with commas or spaces.
- Optionally, each host/domain entry can be followed by a port. If specified, the port must match. If not specified, the protocol's default port is assumed.
- If specified, subdomain names must match. E.g.,
 NO_PROXY=foo.example.com will send requests directly to
 https://foo.example.com, but https://bar.example.com requests will go
 through the proxy.
- You can use leading wildcards like NO_PROXY="*.us, .org".
- NO_PROXY="*" disables all proxies.
- NO_PROXY with an empty list disables no proxies.

Where Proxies Apply

Proxy configuration is relevant to the following LogStream components that make outbound HTTP/S requests:

Destinations

S3 Compatible Stores

- AWS Kinesis Streams
- AWS CloudWatch Logs
- AWS SQS
- Azure Blob Storage
- Azure Event Hubs
- Azure Monitor Logs
- Elasticsearch
- Honeycomb
- Splunk HEC

Sources

- AWS Kinesis Streams
- AWS SQS
- AWS S3
- Azure Event Hubs

Collectors

S3 Collector

Proxying Multiple LogStream Instances in One Browser

LogStream stores authentication tokens based on each http header's URI scheme, host, and port information. Within a given browser, LogStream enforces a same-origin policy to isolate instances.

This means that if you want to run multiple proxied LogStream instances in one browser session, you must assign them different URI schemes, hosts, and/or ports. Otherwise, logging into an extra LogStream instance will expire the prior instance's session and log it out.

For example, assume that you've set up this pair of Apache proxy forward rules:

- https://web/cribla forwards to cribl_hosta:8001/cribla.
- https://web/criblb forwards to cribl_hostb:8001/criblb.

These two proxied addresses cannot be run simultaneously in the same browser session. However, this pair – which lead with separate URI schemes – could:

- https://web/cribla forwards to cribl_hosta:8001/cribla.
- https://web2/criblb forwards to cribl_hostb:8001/criblb.

Where separate instances **must** share URI formats, a workaround is to open the second instance in an incognito/private browsing window, or in a completely different browser.

Scaling Up

A single-instance installation can be configured to scale up and utilize as many resources on the host as required. See Sizing and Scaling for details.

Distributed Deployment

Getting started with Cribl LogStream on a distributed deployment

Distributed Deployment

To sustain higher incoming data volumes, and/or increased processing, you can scale from a single instance up to a multi-instance, distributed deployment. Instances in the deployment independently serve all inputs, process events, and send to outputs. The instances are managed centrally by a single Master Node, which is responsible for keeping configurations in sync, and for tracking and monitoring the instances' activity metrics.

Concepts

Single Instance – a single Cribl LogStream instance, running as a standalone (not distributed) installation on one server.

Master Node – a LogStream instance running in **Master** mode, used to centrally author configurations and monitor Worker Nodes in a distributed deployment.

Worker Node – a LogStream instance running as a **managed Worker**, whose configuration is fully managed by a Master Node. (By default, will poll the master for configuration changes every 10 seconds.)

Worker Group – a collection of Worker Nodes that share the same configuration. You map Nodes to a Worker Group using a Mapping Ruleset.

Worker Process – a Linux process within a Single Instance, or within Worker Nodes, that handles data inputs, processing, and output. The process count is constrained by the number of physical or virtual CPUs available; for details, see Sizing and Scaling.

Mapping Ruleset – an ordered list of filters, used to map Workers Nodes into Worker Groups.

⚠ A Worker Node's local running config can be manually overridden/changed, but changes won't persist on the filesystem. To

permanently modify a Worker Node's config, save, commit, and deploy it from the Master. See Deploying Configurations below.

LogStream 2.4 introduces role-based access control, at the Worker Group level. Users will be able to access Workers only within those Worker Groups on which they've been granted access.

Aggregating Workers

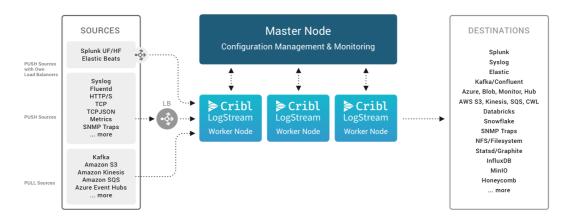
To clarify how the above concepts add up hierarchically, let's use a military metaphor involving toy soldiers:

- Worker Process = soldier.
- Worker Node = multiple Worker Processes = squad.
- Worker Group = multiple Worker Nodes = platoon.

Multiple Worker Groups are very useful in meeting organizational or geographic constraints reflected in configuration. E.g., you might have a U.S. Worker Group with certain TLS certificates and output settings, versus and APAC Worker Group and an EMEA Worker Group that each have distinct certs and settings.

Architecture

This is an overview of a distributed LogStream deployment's components.



Distributed deployment architecture

Master Node Requirements

OS:

- Linux: RedHat, CentOS, Ubuntu, AWS Linux (64bit)
- System:
 - +4 physical cores, +8GB RAM
 - 5GB free disk space
- Git: git must be available on the Master Node. See details below.
- Browser Support: Firefox 65+, Chrome 70+, Safari 12+, Microsoft Edge
 - i We assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). All quantities listed above are minimum requirements.
 - Mac OS is no longer supported as of v. 2.3, due to LogStream's incorporation of Linux-native features.

Worker Node Requirements

See Single-Instance Deployment for requirements and Sizing and Scaling for capacity planning details.

Network Ports - Master Node

In a distributed deployment, Workers communicate with the Master Node on these ports. Ensure that the Master is reachable on those ports from **all** Workers.

Component	Default Port
Heartbeat	4200

Network Ports - Worker Nodes

By default, all LogStream Worker instances listen on the following ports:

Component	Default Port
UI	9000

Installing on Linux

See Single-Instance Deployment, as the installation procedures are identical.

Version Control with git

LogStream requires git (version 1.8.3.1 or higher) to be available locally on the host where the Master Node will run. **Configuration changes must be committed to git before they're deployed.**

If you don't have git installed, check here for details on how to get started.

The Master node uses git to:

- Manage configuration versions across Worker Groups.
- Provide users with an audit trail of all configuration changes.
- Allow users to display diffs between current and previous config versions.

Setting up Master and Worker Nodes

1. Configuring a Master Node

You can configure a Master Node through the UI, through the instance.yml config file, or through the command line.

Using the UI

In Settings > Distributed Settings > Distributed Management > General Settings, select Mode Master. Supply the required Master settings (Address and Port). Customize the optional settings if desired. Then click Save to restart.

☐ Worker UI Access

If you enable the nearby **Distributed Settings > Master Settings > Worker UI access** option (enabledWorkerRemoteAccess key), you will be able to click through from the Master's **Manage Worker Nodes** screen to an authenticated view of each Worker's UI. An orange

header labeled **Viewing Worker:** <host/GUID> will appear to confirm that you are remotely viewing a Worker's UI.



Worker UI access

Using YAML Config File

In \$CRIBL_HOME/local/_system/instance.yml ,under the distributed
section, set mode to master:

```
$CRIBL_HOME/local/_system/instance.yml

distributed:
   mode: master

master:
   host: <IP or 0.0.0.0>
   port: 4200
   tls:
       disabled: true
   ipWhitelistRegex: /.*/
   authToken: <auth token>
   enabledWorkerRemoteAccess: false
   compression: none
   connectionTimeout: 5000
   writeTimeout: 10000
```

Using the Command Line

You can configure a Master Node using a CLI command of this form:

```
./cribl mode-master [options] [args]
```

For all options, see the CLI Reference.

2. Configuring a Worker Node

On each LogStream instance you designate as a Worker Node, you can configure the Worker through the UI, the instance.yml config file, environment variables, or the command line.

Using the UI

In Settings > Distributed Settings > Distributed Management > General Settings, select Mode Worker. Supply the required Master settings (Address and Port). Customize the optional settings if desired. Then click Save to restart.

Using YAML Config File

In \$CRIBL_HOME/local/_system/instance.yml ,under the distributed
section, set mode to worker:

```
$CRIBL_HOME/local/_system/instance.yml
distributed:
 mode: worker
  envRegex: /^CRIBL_/
 master:
   host: <master address>
    port: 4200
    authToken: <token here>
    compression: none
    tls:
      disabled: true
    connectionTimeout: 5000
    writeTimeout: 10000
  tags:
      - tag1
       - tag2
      - tag42
  group: teamsters
```

Using Environment Variables

You can configure Worker Nodes via environment variables, as in this example:

```
CRIBL_DIST_MASTER_URL=tcp://criblmaster@masterHostname:4203
./cribl start
```

See the Environment Variables section for more details.

Using the Command Line

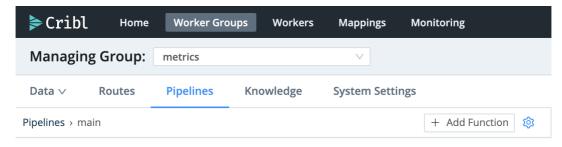
You can configure a Worker Node using CLI commands of this form:

```
./cribl mode-worker -H <master-hostname-or-IP> -p <port> [options] [args]
```

The -H and -p parameters are required. For other options, see the CLI Reference. Here is an example command:

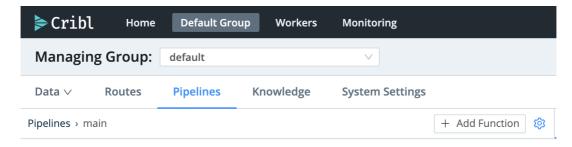
Menu Changes in Distributed Mode

Compared to a single-instance deployment, deploying in distributed mode changes LogStream's menu structure in a few ways. The top menu adds **Worker Groups**, **Workers**, and **Mappings** tabs – all to manage Workers and their assignments.



Distributed deployment: menu structure

If you have a LogStream Free or LogStream One license, the Worker Groups tab instead reads **Default Group**, because these license types allow only this single group. Therefore, throughout this documentation, interpret any reference to the "Worker Groups tab" as "Default Group tab" in your installation.

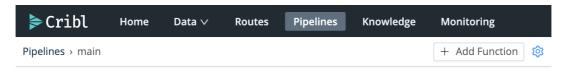


Distributed deployment with LogStream Free/One license

To access the **Data** (drop-down), **Routes**, **Pipelines**, and **Knowledge** items on the light-colored submenu shown above, click the **Worker Groups** tab, then click into your desired Worker Group to display its submenu. This submenu also adds a **System Settings** tab, through which you can manage configuration per Worker Group.

(With a LogStream Free or LogStream One license, you'd click the **Default Group** tab, whose **System Settings** submenu tab configures only that single group.)

For comparison, here is a single-instance deployment's single-level top menu:



Single-instance deployment: single-level menu

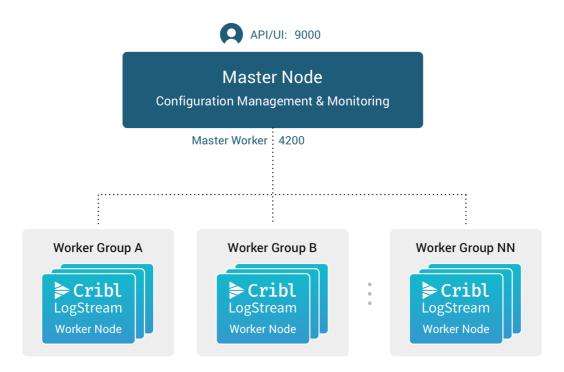
⚠ This repositioning of Data, Routes, Pipelines, and Knowledge tabs to the Worker Groups (or Default Group) submenu also applies to several instructions and screenshots that you'll see throughout this documentation.

Where procedures are written around a single-instance scenario, just click into your appropriate Worker Group to access the same tabs on its submenu.

How Do Workers and Master Work Together

The Master Node has two primary roles:

- 1. Serves as a central location for Workers' operational metrics. The Master ships with a monitoring console that has a number of dashboards, covering almost every operational aspect of the deployment.
- 2. Serves as a central location for authoring, validating, deploying, and synchronizing configurations across Worker Groups.



Master Node/Worker Nodes relationship

Network Port Requirements (Defaults)

- UI access to Master Node: TCP 9000.
- Worker Node to Master Node: TCP 4200 (Heartbeat/Metrics/other).

Master/Worker Node Communication

Workers will periodically (every 10 seconds) send a heartbeat to the Master. This heartbeat includes information about themselves, and a set of current system metrics. The heartbeat payload includes facts – such as hostname, IP address, GUID, tags, environment variables, current software/configuration version, etc. – that the Master tracks with the connection.

The failure of a Worker Node to successfully send two consecutive heartbeat messages to the Master will cause the respective Worker to be removed from the Workers page in the Master's UI until the Master receives a heartbeat message from the affected Worker.

When a Worker Node checks in with the Master:

- The Worker sends heartbeat to Master.
- The Master uses the Worker's facts and Mapping Rules to map it to a Worker Group.
- The Worker Node pulls its Group's updated configuration bundle, if necessary.

Config Bundle Management

Config bundles are compressed archives of all config files and associated data that a Worker needs to operate. The Master creates bundles upon Deploy, and manages them as follows:

- Bundles are wiped clean on startup.
- While running, at most 5 bundles per group are kept.
- Bundle cleanup is invoked when a new bundle is created.

The Worker pulls bundles from the Master and manages them as follows:

- Last 5 bundles and backup files are kept.
- At any point in time, all files created in the last 10 minutes are kept.
- Bundle cleanup is invoked after a reconfigure.

Worker Groups

Worker Groups facilitate authoring and management of configuration settings for a particular set of Workers. To create a new Worker Group, go to the **Worker Groups** top-level menu and click + **Add New**.

Configuring a Worker Group

Click on the newly created Group to display an interface for **authoring and validating** its configuration. You can configure everything for this Group as if it were a single Cribl LogStream instance – using exactly the same visual interface for Routes, Pipelines, Sources, Destinations and System Settings.

△ Can't Log into the Worker Node as Admin User?

To explicitly set passwords for Worker Groups, see User Authentication.

Mapping Workers to Worker Groups

Mapping Rulesets are used to map Workers to Worker Groups. Within a ruleset, a list of rules evaluate Filter expressions on the information that Workers send to the Master.

Only one Mapping Ruleset can be active at any one time, although a ruleset can contain multiple rules. At least one Worker Group should be defined and present in the system.

The ruleset behavior is similar to Routes, where the order matters, and the **Filter** section supports full JS expressions. The ruleset matching strategy is first-match, and one Worker can belong to only one Worker Group.

Creating a Mapping Ruleset

To create a Mapping Ruleset, start on the **Mappings** top-level menu, then click + **Add New**.

i The Mappings top-level menu appears only when you have started LogStream with Distributed Settings > Mode set to Master.

Click on the newly created item, and start adding rules by clicking on + Add Rule. While you build and refine rules, the Preview in the right pane will show which currently reporting and tracked workers map to which Worker Groups.

A ruleset must be activated before it can be used by the Master. To activate it, go to **Mappings** and click **Activate** on the required ruleset. The **Activate** button will then change to an **Active** toggle. Using the adjacent buttons, you can also **Configure** or **Delete** a ruleset, or **Clone** a ruleset if you'd like to work on it offline, test different filters, etc.

Although not required, Workers can be configured to send a Group with their payload. See below how this ranks in mapping priority.

Add a Mapping Rule – Example

Within a Mapping Ruleset, click + Add Rule to define a new rule. Assume that you want to define a rule for all hosts that satisfy this set of conditions:

- IP address starts with 10.10.42, AND:
- More than 6 CPUs OR CRIBL_HOME environment variable contains w0, AND:
- Belongs to Group 420.

Rule Configuration

• Rule Name: myFirstRule

- Filter: (conn_ip.startsWith('10.10.42.') &6 cpus > 6) || env.CRIBL_HOME.match('w0')
- Group: Group420

Default Worker Group and Mapping

When a LogStream instance runs as Master, the following are created automatically:

- A default Worker Group.
- A default Mapping Ruleset,
 - with a default Rule matching all (true).

Mapping Order of Priority

Priority for mapping to a group is as follows: Mapping Rules > Group sent by Worker > default Group.

- If a Filter matches, use that Group.
- Else, if a Worker has a Group defined, use that.
- Else, map to the default Group.

Deploying Configurations

Your typical workflow for deploying LogStream configurations is the following:

- 1. Work on configs.
- 2. Save your changes.
- 3. Commit (and optionally push).
- 4. Deploy.

Deployment is the last step after configuration changes have been saved and committed. **Deploying here means propagating updated configs to Workers.** You deploy new configurations at the Group level: Locate your desired Group and click on **Deploy**. Workers that belong to the group will start **pulling** updated configurations on their next check-in with the Master.

△ Can't Log into the the Worker Node as Admin User?

When a Worker Node pulls its first configs, the admin password will be randomized, unless specifically changed. This means that users won't be able to log in on the Worker Node with default credentials. For details, see User Authentication.

Configuration Files

On the Master, a Worker Group's configuration lives under: \$CRIBL_HOME/groups/<groupName>/local/cribl/.

On the managed Worker, after configs have been pulled, they're extracted under: \$CRIBL_HOME/local/cribl/.

Lookup Files

On the Master, a Group's lookup files live under: \$CRIBL_HOME/groups/<groupName>/data/lookups.

On the managed Worker, after configs have been pulled, lookups are extracted under: \$CRIBL_HOME/data/lookups. When deployed via the Master, lookup files are distributed to Workers as part of a configuration deployment.

If you want your lookup files to be part of the LogStream configuration's version control process, we recommended deploying using the Master Node. Otherwise, you can update your lookup file out-of-band on the individual Workers. The latter is especially useful for larger lookup files (> 10 MB, for example), or for lookup files maintained using some other mechanism, or for lookup files that are updated frequently.

For other options, see Managing Large Lookups.

i Some configuration changes will require restarts, while many others require only reloads. See here for details.

Restarts/reloads of each Worker Process are handled automatically by the Worker. Note that individual Worker Nodes might temporarily disappear from the Master's **Workers** tab while restarting.

Worker Process Rolling Restart

During a restart, to minimize ingestion disruption and increase availability of network ports, Worker Processes on a Worker Node are restarted in a rolling fashion. 20% of running processes – with a minimum of one process – are

restarted at a time. A Worker Process must come up and report as started before the next one is restarted. This rolling restart continues until all processes have restarted. If a Worker Process fails to restart, configurations will be rolled back.

Auto-Scaling Workers and Load-Balancing Incoming Data

If data flows in via Load Balancers, make sure to register all instances. Each Cribl LogStream node exposes a health endpoint that your Load Balancer can check to make a data/connection routing decision.

Health Check Endpoint	Healthy Response
curl http:// <host>:<port>/api/v1/health</port></host>	{"status":"healthy"}

Environment Variables

- CRIBL_DIST_MASTER_URL URL of the Master Node. Format:
 <tls|tcp>://<authToken>@host:port?
 group=defaultGroup&tag=tag1&tag=tag2&tls.<tls-settings below>.
 - tls.privKeyPath Private KeyPath.
 - tls.passphrase Key Passphrase.
 - tls.caPath CA Certificate Path.
 - tls.certPath Certificate Path.
 - tls.rejectUnauthorized Validate Client Certs. Boolean, defaults to false.
 - tls.requestCert Authenticate Client (mutual auth). Boolean, defaults to false.
 - tls.commonNameRegex Regex matching peer certificate > subject > common names allowed to connect. Used only if tls.requestCert is set to true.
- CRIBL_DIST_MODE worker | master.Defaults to worker iff
 CRIBL_DIST_MASTER_URL is present.
- CRIBL_HOME Auto setup on startup. Defaults to parent of bin directory.
- CRIBL_CONF_DIR Auto setup on startup. Defaults to parent of bin directory.
- CRIBL_NOAUTH Disables authentication. Careful here!!

- CRIBL_VOLUME_DIR Sets a directory that persists modified data between different containers or ephemeral instances.
- Deprecated variables: CRIBL_CONFIG_LOCATION, CRIBL_SCRIPTS_LOCATION

Workers GUID

When you install and first run the software, a GUID is generated and stored in a .dat file located in CRIBL_HOME/bin/, e.g.:

```
# cat CRIBL_HOME/bin/676f6174733432.dat
{"it":1570724418,"phf":0,"guid":"48f7b21a-0c03-45e0-a699-
01e0b7a1e061"}
```

When deploying Cribl LogStream as part of a host image or VM, be sure to remove this file, so that you don't end up with duplicate GUIDs. The file will be regenerated on next run.

Bootstrap Workers from Master

Boot fully provisioned workers

This feature of LogStream allows workers to completely provision themselves on initial boot, directly from the master. It allows a fleet of any number of nodes to launch. and be fully functional within the cluster, in seconds.

How Does It Work?

A LogStream Master Node (v2.2 or higher) provides a bootstrap API endpoint, at /init/install-worker.sh, which returns a shell script. You can run this shell script on any supported machine (see Restrictions below) without LogStream installed, fully provisioning the machine as a Worker Node.

Although you can specify the download URL when you execute the initial curl command, the LogStream package is not downloaded until the script is generated by the API, and then later executed.

△ Root Access or sudo

Note that the script will install LogStream into <code>/opt/cribl</code>, and will make system-level changes. For systems like Ubuntu, which don't allow direct root access, you'll need to use the <code>sudo</code> command when executing the script.

API Spec

Request Format

GET http://<master hostname or IP>:9000/init/install-worker.sh

Query Strings

String	Required?	Description
token	optional	Master Node's shared secret (authToken). By default, this is set to criblmaster. You can find this secret in the Master Node's Distributed Settings section.
group	optional	Name of the cluster's work group. If not specified, falls back to default .
download_url	optional	Provide the complete URL to a Cribl LogStream installation binary. This is especially useful if the Worker Nodes don't have access to the Internet to download from cribl.io.

Example HTTP Request

HTTP

GET http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-554

Response

```
Shell
#!/bin/sh
### START CRIBL MASTER TEMPLATE SETTINGS ###
CRIBL MASTER HOST="<Master FQDN/IP>"
CRIBL_AUTH_TOKEN="<Auth token string>"
CRIBL_VERSION="<Version>"
CRIBL_GROUP="<Default group preference>"
CRIBL_MASTER_PORT="<Master heartbeat port>"
CRIBL_DOWNLOAD_URL="<download url>"
### END CRIBL MASTER TEMPLATE SETTINGS ###
# Set defaults
checkrun() { $1 --help >/dev/null 2>/dev/null; }
faildep() { [ $? -eq 127 ] & echo "$1 not found" & exit 1; }
[ -z "${CRIBL_MASTER_HOST}" ] && echo "CRIBL_MASTER_HOST not set" && exit 1
CRIBL_INSTALL_DIR="${CRIBL_INSTALL_DIR:-/opt/cribl}"
CRIBL MASTER PORT="${CRIBL MASTER PORT:-4200}"
CRIBL_AUTH_TOKEN="${CRIBL_AUTH_TOKEN:-criblmaster}"
CRIBL_GROUP="${CRIBL_GROUP:-default}"
if [ -z "${CRIBL_DOWNLOAD_URL}" ]; then
    FILE="cribl-${CRIBL_VERSION}-linux-x64.tgz"
    \label{lownload_url="https://cdn.cribl.io/dl/$(echo $\{CRIBL_VERSION\} \mid cut -d '-' -f 1)/$\{FIIA \mid CRIBL_VERSION\} \mid cut -d '-' -f 1/$\{FIIA \mid CRIBL_VERSION\} \mid cut -d '-' -f 1/$\}
fi
UBUNTU=0
CENTOS=0
AMAZON=0
echo "Checking dependencies"
checkrun curl & faildep curl
checkrun adduser & faildep adduser
checkrun usermod & faildep usermod
BOOTSTART=1
checkrun systemctl & [ $? -eq 127 ] & BOOTSTART=0
checkrun update-rc.d & [ \$? -eq 127 ] & BOOTSTART=0
echo "Checking OS version"
lsb_release -d 2>/dev/null | grep -i ubuntu & [ $? -eq 0 ] & UBUNTU=1
cat /etc/system-release 2>/dev/null | grep -i amazon & [ $? -eq 0 ] & AMAZON=1
echo "Creating cribl user"
if [ $UBUNTU -eq 1 ]; then
    adduser cribl --home /home/cribl --gecos "Cribl LogStream User" --disabled-password
fi
if [ CENTOS - eq 1 ] || [ AMAZON - eq 1 ]; then
    adduser cribl -d /home/cribl -c "Cribl LogStream User" -m
    usermod -aG wheel cribl
fi
```

curl Option

An easy way of wrapping HTTP methods is to use the curl command. Here is an example, which uses a GET operation by default, with the same URL used in the above HTTP example:

Shell

curl http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-5!

Chaining Script Execution

The GET and curl procedures above will only output the contents of the script that needs executing – the script will still need to be manually executed. However, you can automate that part, too, using the command below. This passes the script's contents to the sh shell to immediately execute. As noted above, on Ubuntu and similar systems, you might need to insert sudo before the sh.

Shell

curl http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-5!

Adding Download URL

We'll now graduate to the next level by adding more to the above commands. All the preceding commands excluded the download_url parameter so, by default, the script gets configured to download the LogStream package from the public Cribl repository.

To successfully execute the curl command while also specifying the download_url, you must enclose the URL in double quotes. The reason for this is that the & character that joins multiple HTTP parameters is interpreted by the shell as the operator to run commands in the background. Quoting the URL, as shown in this example, prevents this.

Shell

curl "http://<master hostname or IP>:9000/init/install-worker.sh?token=79364d6e-dead-beef-4c6e-5

Status Codes

Status Code	Reason
200 – OK	All is well. You should have received the script as a response.
403 – Forbidden	Either the node is not configured as a Master, or the token provided is invalid.

Restrictions

Keep the following in mind when using this feature:

- Each Worker must normally have access to the internet in order to download the Cribl LogStream installation binary from cribl.io. Where this isn't feasible, you can use the download_url switch to point to a binary in a restricted location.
- By default, Worker Nodes communicate with the Master on port 4200. Ensure that access between all Workers and the Master is open on this port.
- TLS is not enabled by default. If enabled and configured, access to this feature will be over https: instead of http.
- Red Hat, Ubuntu, CentOS, and Amazon Linux are the only supported Worker platforms.

User Data

For public-cloud customers, an easy way to use this feature is in an instance's user data. Simply use the following script (changing the command as needed. based on the information above). Upon launch, the Worker Node will reach out to the Master, download the script, download the LogStream package from the specified location, and then install and configure LogStream:

Shell

#!/bin/bash
curl http://<master-node-ip/host-address>:9000/init/install-worker.sh?token=<auth-token> | sh -

Splunk App Deployment *

Getting started with Cribl App for Splunk

△ Cribl App for Splunk for HFs Is Deprecated as of Cribl LogStream v.2.1

Cribl will continue to support this package, but **customers are advised to begin planning now for the eventual removal of support**.

See Single-Instance Deployment and Distributed Deployment for alternatives.

Deploying Cribl App for Splunk

In a Splunk environment, Cribl LogStream can be installed and configured as a Splunk app (Cribl App for Splunk). Depending on your requirements and architecture, it can run either on a Search Head or on a Heavy Forwarder. Cribl App for Splunk cannot be used in a Cribl LogStream Distributed Deployment as a Master or managed as Worker.

Running on a Search Head (SH)

When running on an SH, Cribl LogStream is set to **mode-searchhead**, the default mode for the app. It listens for **localhost traffic** generated by a custom command: | criblstream. The command is used to forward search results to the LogStream instance's TCP JSON input on port 10420, but it's also capable of sending to any other LogStream instance listening for TCP JSON.

Once received, data can be processed and forwarded to any of the supported Destinations. In addition, several out-of-the box saved searches are ready to run and send their results to Cribl with a single click.

Installing the Cribl App for Splunk on an SH

- Select an instance on which to install.
- Ensure that ports 10000, 10420, and 9000 are available. See the Requirements section for more info.
- Get the bits here, and install as a regular Splunk app.
- Restart the Splunk instance.
- Go to https://<instance>/en-US/app/cribl or https://<instance>:9000, and log in with Splunk admin role credentials.

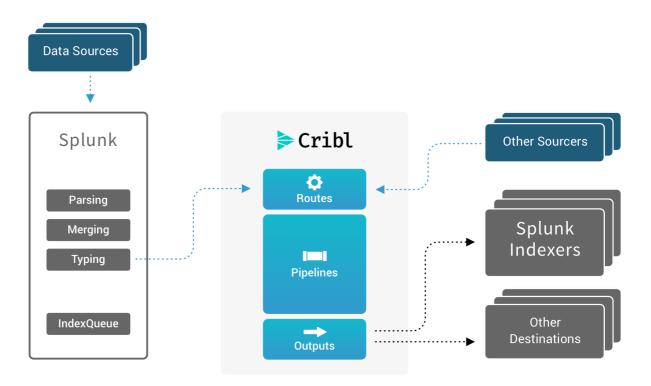
Typical Use Cases for Search Head Mode

• Working with search results in a Cribl LogStream pipeline.

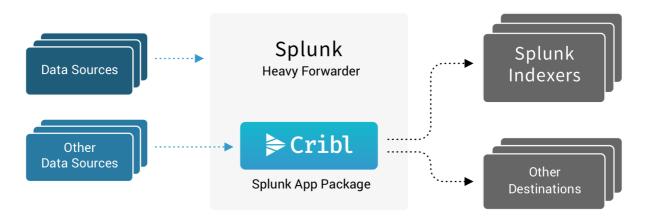
• Sending search results to any Destination supported by Cribl LogStream.

Running on a Heavy Forwarder (HF)

When running on an HF, Cribl LogStream is set to **mode-hwf**. It receives events from the local Splunk process per routing configurations in props.conf and transforms.conf. Data is parsed and processed first by Splunk pipelines, and then by LogStream. By default, all data except internal indexes is routed out right after the Typing pipeline.



Cribl LogStream is capable of accepting data **streams** (unbroken events) or **events** from other sources. In this case, the HF will deliver **events** locally to LogStream, which processes them and sends them to one or more destinations downstream. When receivers are Splunk indexers, LogStream can also load-balance across them.



Installing the Cribl App for Splunk on an HF

• Select an instance on which to install.

- Ensure that ports 10000, 10420, and 9000 are available. See here.
- Get the bits here, and install as a regular Splunk app.
- Set Cribl to mode-hwf: \$SPLUNK_HOME/etc/apps/cribl/bin/cribl mode-hwf.



- Restart the Splunk instance.
- Go to https://<instance>:9000 and log in with Splunk admin role credentials.

```
Note About Splunk Warnings
If you come across messages similar to the following example, on startup or in logs, please ignore them. They are benign warnings.
Invalid value in stanza [route2criblQueue]/[hecCriblQueue] in /opt/splunk/etc/apps/cribl/default/transforms.conf, line 11: (key: DEST_KEY, value: criblQueue) / line 24: (key: DEST_KEY, value: $1)
```

Relevant configurations in Cribl App for Splunk on an HF

When Cribl App for Splunk is installed on an HF (in mode-hwf), below are the **relevant sections** in configuration files that enable Splunk to send data to Cribl LogStream:

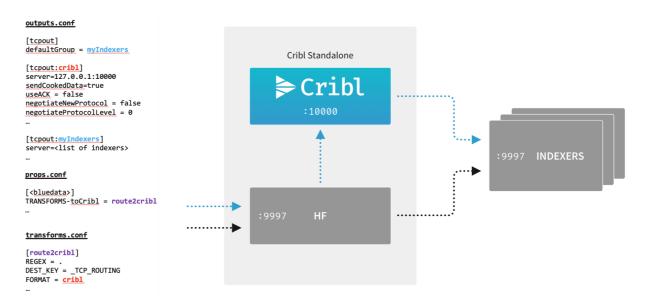
```
apps/cribl/default/outputs.conf
[tcpout]
disabled = false
defaultGroup = cribl
[tcpout:cribl]
server=127.0.0.1:10000
sendCookedData=true
useACK = false
negotiateNewProtocol = false
negotiateProtocolLevel = 0
apps/cribl/default/inputs.conf
[splunktcp]
route=has_key:_replicationBucketUUID:replicationQueue;has_key:_dstrx:typingQueue;has_key:__CRIBN
apps/cribl/default/transforms.conf
[route2cribl]
SOURCE_KEY = _MetaData:Index
REGEX = ^[^_]
DEST_KEY = _TCP_ROUTING
FORMAT = cribl
```

```
[route2criblQueue]
SOURCE_KEY = _MetaData:Index
REGEX = ^[^_]
DEST_KEY = queue
FORMAT = criblQueue

apps/cribl/default/props.conf
[default]
TRANSFORMS-cribl = route2criblQueue, route2cribl
```

Configuring Cribl LogStream with a Subset of Your Data

The props.conf stanza above will apply the above transforms to **everything**. Depending on your requirements, you might want to target only a subset of your sources, sourcetypes, or hosts. For example, the diagram below shows the **effective** configurations of outputs.conf, props.conf, and transforms.conf to send <bluedata> events through Cribl LogStream.



Configure Cribl LogStream to Send Data to Splunk Indexers

To send data from Cribl LogStream to a set of Splunk indexers, use the LogStream UI to go to **Destinations** > Splunk Load Balanced, then enter the required information.

Kubernetes Master Deployment

Boot a fully provisioned Master Node via Helm

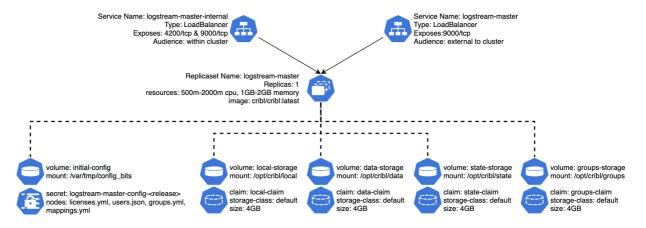
This page outlines how to deploy a Cribl LogStream Master Node (or single instance) to AWS via Kubernetes, using a Cribl-provided Helm chart.

<u>^</u>

This chart is a work in progress, provided as-is. Cribl expects to further develop and refine it.

Deployment

As built, Cribl's chart will deploy a Master Server for LogStream, consisting of a deployment, two services, and a number of persistent volumes.



Deployment schematic

Note that this chart creates two load-balanced services:

- The main one (named after the Helm release), which is intended as the primary service interface for users.
- The "internal" one (named <helm-release > internal), which is intended for the workergroup-to-master communication.
 - By default, this chart installs only a LogStream Master Node. To also deploy LogStream Worker Groups via Helm, you can use the Set Up Worker Groups/Mappings override described below

You can also use Cribl's separate logstream-workergroup chart. For details, see Kubernetes Deployment: Worker Group in this documentation.

AWS and Kubernetes Prerequisites

This section covers both general and specific prerequisites, with a bias toward the EKS-oriented approach that Cribl uses for its own deployments.

Set Up AWS CLI

Install the AWS CLI, version 2, according to AWS' instructions.

Next, create or modify your ~/.aws/config file to include (at least) a [profile] section with the following SSO (single-sign-on) details:

```
~/.aws/config

[profile <your-profile-name>]

sso_start_url = https://<your-domain>/start#/
sso_region = <your-AWS-SSO-region>
sso_account_id = <your-AWS-SSO-account-ID>
sso_role_name = <your-AWS-role-name>
region = <your-AWS-deployment-region>
```

Set Up kubectl

You will, of course, need kubectl set up on your local machine or VM. Follow Kubernetes' installation instructions.

Add a Cluster to Your kubeconfig File

You must modify your ~/.kube/config file to instruct kubectl what cluster (context) to work with.

1. Run a command of this form:

This should return a response like this:

```
Added new context arn:aws:eks:us-west-2:4242424242:cluster/<cluster-name> to /Users/<username>/.kube/config
```

2. In the resulting ~/.kube/config file's args section, as the new first child, insert the profile argument that you provided to the aws command. For example:

```
/.kube/config
args:
   --profile=<profile-name>
   --region
[ ... ]
```

3. Also change the command: aws pair to include the full path to the aws executable. This is usually in /usr/local/bin, in which case you'd insert: command: /usr/local/bin/aws.

This section of ~/.kube/config should now look something like this:

```
~/.kube/config

args:
    --profile=<profile-name>
    --region
    - us-west-2
    - eks
    - get-token
    --cluster-name
    - lab
    command: /usr/local/bin/aws
    env:
    - name: AWS_PROFILE
    value: <profile-name>
```

With these AWS and Kubernetes prerequisites completed, you're now set up to run kubectl commands against your cluster, as long as you have an active aws SSO login session.

Next, do the Helm setup.

Install Helm and Cribl Repo

- 1. You'll need Helm (preferably v.3.x) installed. Follow the instructions here.
- 2. Add Cribl's repo to Helm, using this command:
 helm repo add cribl https://criblio.github.io/helm-charts/

Persistent Storage

The chart requires persistent storage. It will use your default StorageClass, or (if you prefer) you can override config.scName with the name of a specific StorageClass to use.

Cribl has tested this chart primarily using AWS EBS storage, via the CSI EBS driver. The volumes are created as ReadWriteOnce claims. For details about storage classes, see Kubernetes' Storage Classes documentation.

AWS-Specific Notes

If you're running on EKS, Cribl highly recommends that you use Availability Zone–specific node groups. For details, see eksctl.io's Autoscaling documentation.

△ Do not allow a single node group to spans AZs. This can lead to trouble in mounting volumes, because EBS volumes are AZ-specific.

Configure the Chart's Values

You'll want to override some of the chart's default values. The easiest way is to copy this chart's default values. yaml file from our repo. save it locally, modify it, and install it in Helm:

- Copy the raw contents of: https://github.com/criblio/helm-charts/blob/master/helm-chart-sources/logstream-master/values.yaml
- 2. Save this as a local file, e.g.: /bar/values.yaml
- 3. Modify values as necessary (see Values to Override below).
- 4. Install your updated values to Helm, using this command: helm install -f /bar/values.yaml

Values to Override

This section covers the most likely values to override. To see the full scope of values available, run: helm show values cribl/logstream-master

Key	Туре	Default Value	Description
config.adminPassword	String	[No default]	The password you want to assign to the admin user.
config.token	String	[No default]	The auth key you want to set up for Worker access. If you set this value, the LogStream instance will be configured only as a Master server for a distributed deployment. (You can also configure this later via the LogStream UI, after launching the instance in single-instance mode.)
config.license	String	[No default]	The license for your LogStream instance. If you do not set this, it will default to the Free license. You can change this in the LogStream UI as well.
config.groups	List	[No default]	Array of Worker Group names to configure for the Master instance. This will create a mapping for each Group, which looks for the tag <groupname>, and will create the basic structure of each Group's configuration.</groupname>

config.scName	String		The StorageClass name for all of the persistent volumes.
<pre>config. rejectSelfSignedCerts</pre>	Number	0	Either 0 (allow self-signed certificates) or 1 (deny self-signed certs).
config.healthPort	number	9000	The port to use for health checks (readiness/live).
service.ports	Array of Maps	- name: api port: 9000 protocol: TCP external: true - name: mastercomm port: 4200 protocol: TCP external: false	The ports to make available, both in the deployment and in the service. Each "map" in this list needs the following values set: name A descriptive name, identifying what the port is being used for. port The container port to be made available. protocol The protocol in use for this port (UDP or TCP). external Set to true to expose the port on the external service, or false to not expose it.
service.annotations	String	[No default]	Annotations for the the service component – this is where you'll want to put load-balancer–specific configuration directives.
image.tag	String	latest	The container image tag to pull from. Cribl will increment this tag per LogStream version. By default, this will use a version equivalent to the chart's appVersion value. You can override this with latest to get the latest LogStream version, or with a specific LogStream version number (like "2.3.3").

Match Versions

Cribl recommends that you use the same LogStream version on Worker Nodes versus the Master Node. So if, for any reason, you're not yet upgrading your Workers to the version in the Master's default values.yaml > criblImage.tag, be sure to override that criblImage.tag value to match the version you're running on all Workers.

EKS-Specific Values

If you're deploying to EKS, many annotations are available for the load balancer. Set these as values for the service.annotations key. Internally, we typically use the annotations for logging to S3, like this:

```
values.yaml [excerpt]
```

```
service.beta.kubernetes.io/aws-load-balancer-access-log-enabled: "true"
   service.beta.kubernetes.io/aws-load-balancer-access-log-emit-interval: "5"
   service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-name: "<bucket name>"
   service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-prefix: "ELB"
```

For an exhaustive list of annotations you can use with AWS's Elastic Load Balancers, see the Kubernetes Service documentation.

Basic Chart Installation

With the above prerequisites and configuration completed, you're ready to install our chart to deploy a LogStream Master Node. Here are some example commands:

• To install the chart with the release name logstream-master:

```
helm install logstream-master cribl/logstream-master
```

• To install the chart using the storage class ebs-sc:

```
helm install logstream-master cribl/logstream-master --set config.scName='lebs-sc
```

Change the Configuration

If you don't override its default values, this Helm chart effectively creates a single-instance deployment of LogStream, using the standard container image. You can later configure distributed mode, licensing, user passwords, etc., all from the LogStream UI. However, you also have the option to change these configuration details upfront, by installing with value overrides. Here are some common examples.

Apply a License

If you have a standard or enterprise license, you can use the config.license parameter to add it as an override to your install:

```
helm install logstream-master cribl/logstream-master --set config.license="<long encoded license string redacted>"
```

Set the Admin Password

Normally, when you first install LogStream and log into the UI, it prompts you to change the default admin password. You can skip the password-change challenge by setting your admin password via the config.adminPassword parameter:

helm install logstream-master cribl/logstream-master --set config.adminPassword="
<new password>"

Set Up Worker Groups/Mappings

As mentioned above, the chart's default is to install a vanilla deployment of LogStream. If you are deploying as a Master, you can use the config.groups parameter to define the Worker Groups you want created and mapped. Each group in the list you provide will be created as a Worker Group, with a Mapping Rule to seek a tag with that Worker Group's name in it:

helm install logstream-master cribl/logstream-master --set config.groups= {group1,group2,group3}

The example above will create three Worker Groups – group1, group2, and group3 – and a Mapping Rule for each.

Upgrading from a Pre-2.4.0 Version

LogStream 2.4.0 introduces the \$CRIBL_VOLUME_DIR environment variable, which simplifies the persistent-storage requirement for logstream-master. Instead of maintaining four separate persistent volumes (one each for \$CRIBL_HOME/{data,state,local,groups}), you can consolidate all four into a single volume.

In the Helm chart, we handle this via the helm upgrade command. For the chart's 2.4.0 version, the upgrade option creates a new, larger volume, and consolidates the data from the original volumes to that volume. An initContainer handles the logistics. When the initContainer completes, the logstream-master pod comes back up with a single consolidated volume.

☐ Back Up Your Data First

While we've tested this upgrade repeatedly, differences in environments can always cause problems. Therefore, we recommend that you back up your data before running the upgrade command. This is best done with a combination of kubectl and tar:

kubectl exec <pod name> -n <namespace> -- bash -c "cd /opt/cribl; tar cf - {state,data

This command executes the tar -based backup of all four volumes, and outputs it to a local .tar file (cribl_backup.tar).

Running the Upgrade

Helm makes upgrades easy. You simply need to run helm repo update to ensure you have the latest repo updates available, followed by helm upgrade to actually upgrade the containers.

For example, if you've installed the Helm charts in the logstream namespace, named your release ls-master, and set up your Helm repo according to the prerequisites section above (i.e., named it cribl), run the following:

```
helm repo update
helm upgrade ls-master -n logstream cribl/logstream-master
```

Upgrade Order of Operations

While there should be no major problems running a 2.4.0 master and 2.3.4 workers, Cribl does not recommend this. Instead, upgrade the master Helm chart to 2.4.0 first, and then upgrade the workers. (For details, see Kubernetes Worker Deployment.)

Idempotency of Upgrade

The upgrade operation performs a potentially destructive action in coalescing the 4 volumes to a single volume. But that operation proceeds only if the single volume has no data on it. Once the upgrade is performed the first time, any further upgrade operations will effectively skip that coalescence operation, without causing any additional issues.

Recovering from a Failed Upgrade

If the upgrade fails, the suggested recovery path is to remov the Helm chart, reinstall it, and then run this command to restore the data from the backup:

```
cat cribl_backup.tar| kubectl -n <namespace> exec --stdin <pod name> -- bash -c "cd /opt/cribl/c
```

This will restore the data into the "new" volume (which is mounted as <code>/opt/cribl/config-volume</code>). If you want to double-check that:

```
kubectl -n <namespace> exec <pod name> -- bash -c "ls -alR /opt/cribl/config-volume"
```

Uninstall the Infrastructure

To spin down deployed pods, use the helm uninstall command – where <release-name> is the namespace you assigned when you installed the chart:

```
helm uninstall <release-name>
```

You can append the --dry-run flag to verify which releases will be uninstalled before actually uninstalling them:

Known Issues

- Cribl's current architecture supports **only** TCP ports in Worker Groups' service > ports configuration. This restriction might be removed in future versions.
- The upgrade process creates an initContainer, which will run prior to any instance of the LogStream pod. Because the coalescence operation will not overwrite existing data, this is not a functional problem. But depending on your persistent-volume setup, the initContainer 's precedence might cause pod restarts to take additional time while waiting for the volume claims to release. The only upgrade path that will have this issue is 2.3.* -> 2.4.0. In the next iteration, we'll remove the initContainer from the upgrade path.
- The upgrade process leaves the old PersistentVolumes and PersistentVolumeClaims around. This is, unfortunately, necessary for this upgrade path. In follow-on versions, we will remove these volumes from the chart.

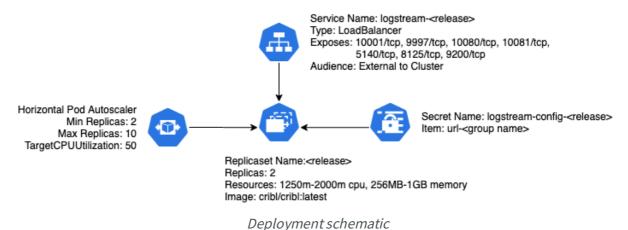
Kubernetes Worker Deployment

Boot a fully provisioned Worker Group via Helm

This page outlines how to deploy a Cribl LogStream Worker Group to AWS via Kubernetes, using a Cribl-provided Helm chart.

Deployment

As built, Cribl's chart will deploy a simple Worker Group for LogStream, consisting of a deployment, a service, a horizontal pod autoscaler configuration, and a secret used for configuration.



Deployment schematic

☐ This chart will deploy only a LogStream Worker Group. To deploy a LogStream Master Node, see Kubernetes Deployment: Master.

AWS and Kubernetes Prerequisites

This section covers both general and specific prerequisites, with a bias toward the EKS-oriented approach that Cribl uses for its own deployments.

Set Up AWS CLI

Install the AWS CLI, version 2, according to AWS' instructions.

Next, create or modify your ~/.aws/config file to include (at least) a [profile] section with the following SSO (single-sign-on) details:

~/.aws/config

```
[profile <your-profile-name>]
sso_start_url = https://<your-domain>/start#/
sso_region = <your-AWS-SSO-region>
sso_account_id = <your-AWS-SSO-account-ID>
sso_role_name = <your-AWS-role-name>
region = <your-AWS-deployment-region>
```

Set Up kubectl

You will, of course, need kubectl set up on your local machine or VM. Follow Kubernetes' installation instructions.

Add a Cluster to Your kubeconfig File

You must modify your ~/.kube/config file to instruct kubectl what cluster (context) to work with.

1. Run a command of this form:

```
aws --profile <profile-name> eks update-kubeconfig --name <cluster-name>
```

This should return a response like this:

```
Added new context arn:aws:eks:us-west-2:4242424242:cluster/<cluster-name> to /Users/<username>/.kube/config
```

2. In the resulting ~/.kube/config file's args section, as the new first child, insert the profile argument that you provided to the aws command. For example:

```
/.kube/config

args:
    ---profile=<profile-name>
    --region
[...]
```

3. Also change the command: aws pair to include the full path to the aws executable. This is usually in /usr/local/bin, in which case you'd insert: command: /usr/local/bin/aws.

This section of ~/.kube/config should now look something like this:

```
~/.kube/config

args:
    - --profile=<profile-name>
    - --region
    - us-west-2
    - eks
    - get-token
    - --cluster-name
    - lab
    command: /usr/local/bin/aws
```

With these AWS and Kubernetes prerequisites completed, you're now set up to run kubectl commands against your cluster, as long as you have an active aws SSO login session.

Next, do the Helm setup.

Install Helm and Cribl Repo

- 1. You'll need Helm (preferably v.3.x) installed. Follow the instructions here.
- 2. Add Cribl's repo to Helm, using this command: helm repo add cribl https://criblio.github.io/helm-charts/
- 3. Display the default values available to configure Cribl's logstream-workergroup chart: helm show values cribl/logstream-workergroup

Configure the Chart's Values

You'll want to override some of the values you've just displayed. The easiest way is to copy this chart's default values.yaml file from our repo. save it locally, modify it, and install it in Helm:

- Copy the raw contents of: https://github.com/criblio/helm-charts/blob/master/helm-chart-sources/logstream-workergroup/values.yaml
- 2. Save this as a local file, e.g.: /foo/values.yaml
- 3. Modify values as necessary (see Values to Override below).
- 4. Install your updated values to Helm, using this command: helm install -f /foo/values.yaml

Values to Override

This section covers the most likely values to override.

Key	Туре	Default Value	Description
config.tag	String	kubernetes	The tag/group to include in the URL. (This is included as both a tag value

			and a group value.)
config.token	String	criblmaster	The authentication token for your LogStream Master.
config.host	String	logstream- master	The resolveable hostname of your LogStream Master.
config. rejectSelfSignedCerts	Number	0	One of: 0 – allow self- signed certs, or 1 – deny self- signed certs.
service.ports	Array of Maps	- name: tcpjson port: 10001 protocol: TCP - name: s2s port: 9997 protocol: TCP - name: http port: 10080 protocol: TCP - name: https port: 10081 protocol: TCP - name: syslog port: 5140 protocol: TCP - name: metrics port: 8125 protocol: TCP - name: elastic port: 9200 protocol: TCP	The ports to make available, both in the deployment and in the service. Each "map" in this list needs the following values set: name A descriptive name, identifying what the port is being used for. port The container port to be made available. protocol The protocol in use for this

			port (UDP or TCP).
service.annotations	String	[No default]	Annotations for the the service component – this is where you'll want to put load-balancer– specific configuration directives.
criblImage.tag	String	2.3.3	The container image tag to pull from. Cribl will increment this tag per LogStream version. By default, this will use a version equivalent to the chart's appVersion value. You can override this with latest to get the latest LogStream version, or with a specific LogStream version number.
autoscaling.minReplicas	Number	2	The minimum number of LogStream pods to run.
autoscaling.maxReplicas	Number	10	The maximum number of LogStream pods to scale up to.
autoscaling.targetCPUUtilizationPercentage	Number	50	The CPU utilization percentage that triggers scaling action.

rbac.create	Boolean	false	Enable Service Account, Cluster Role, and Role Binding creation.
rbac.resources	List	["pods"]	Set the resource boundary for the role being created (K8s resources).
rbac.verbs	List	["get", "list"]	Set the API verbs allowed the role (default: read ops).

Match Versions

Cribl recommends that you use the same LogStream version on Master Nodes versus Worker Group Nodes. So, if you're not yet upgrading your Master to the version in the current values.yaml > criblImage.tag, be sure to override that criblImage.tag value to match the version you're running on the Master.

EKS-Specific Values

If you're deploying to EKS, many annotations are available for the load balancer. Set these as values for the service.annotations key. Internally, we typically use the annotations for logging to S3, like this:

```
values.yaml[excerpt]

service:
    type: LoadBalancer
    annotations: {
        service.beta.kubernetes.io/aws-load-balancer-access-log-enabled: "true"
        service.beta.kubernetes.io/aws-load-balancer-access-log-emit-interval: "5"
        service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-name: "<bucket name>"
        service.beta.kubernetes.io/aws-load-balancer-access-log-s3-bucket-prefix: "ELB"
    }
```

For an exhaustive list of annotations you can use with AWS's Elastic Load Balancers, see the Kubernetes Service documentation.

Install the Chart

With the above prerequisites and configuration completed, you're ready to install our chart to deploy a LogStream Worker Group. Here are some example commands:

• To install the chart with the release name logstream-wg:

```
helm install logstream-wg cribl/logstream-workergroup
```

• To install the chart using the LogStream Master logstream.lab.cribl.io:

```
helm install logstream-wg cribl/logstream-workergroup --set config.host='logstream.lab.cribl.io
```

• To install the chart using the LogStream Master logstream.lab.cribl.io in the namespace cribl-helm:

```
helm install logstream-wg cribl/logstream-workergroup --set config.host='logstream.lab.cribl.io' -n cribl-helm
```

Upgrading

You upgrade using the helm upgrade command. But it's important to ensure that your Helm repository cache is up to date, so first issue this command:

```
helm repo update
```

After this step, invoke:

```
helm upgrade <release> -n <namespace> cribl/logstream-workergroup
```

For the example above, where the release is logstream—wg and is installed in the cribl—helm namespace, the command would be:

```
helm upgrade logstream-wg -n cribl-helm cribl/logstream-workergroup
```

This Helm chart's upgrade is idempotent, so you can use the upgrade mechanism to upgrade the chart, but you can also use it to change its configuration (as outlined in Change the Configuration).

Optional: Kubernetes API Access

Starting in version 2.4.0, the chart provides options for Worker Groups to access the Kubernetes API. The values.yaml file provides three relevant options:

- rbac.create enables the creation of a Service Account, Cluster Role, and Role Binding (which binds the first two together) for the release.
- rbac.resources specifies the Kubernetes API resources that will be available to the release.

• rbac.verbs - specifies the API verbs that will be available to the release.

For more information on the verbs and resources available, see Kubernetes' Using RBAC Authorization documentation.

Change the Configuration

Once you've installed a release, you can get its values.yaml file by using the helm get values command. For example, assuming a release name of logstream-wg, you could use this command:

```
helm get values logstream-wg -o yaml > values.yaml
```

This will retrieve a local values.yaml file containing the values in the running release, including any values that you overrode when you installed the release.

You can now make changes to this local values.yaml file, and then use the helm upgrade operation to "upgrade" the release with the new configuration.

For example, assume you wanted to add an additional TCP-based syslog port, listening on port 5141, to the existing logstream-wg release. In the values.yaml file's service > ports section, you'd add the three key-value pairs shown below:

```
values.yaml(excerpt)
service:
[ ... ]
ports:
[ ... ]
- name: syslog
port: 5141
protocol: TCP
```

Then you'd run:

helm upgrade logstream-wg cribl/logstream-workergroup -f values.yaml

Uninstall the Infrastructure

To spin down deployed pods, use the helm uninstall command – where <release-name> is the namespace you assigned when you installed the chart:

```
helm uninstall <release-name>
```

You can append the --dry-run flag to verify which releases will be uninstalled before actually uninstalling them:

helm uninstall <release-name> --dry-run

Notes on This Example

- If you installed in a namespace, you'll need to include the -n <namespace> option in any helm command.
- In the above syslog example, you'd still need to configure a corresponding syslog Source in your LogStream Master, and then commit and deploy it to your Worker Group(s).

Known Issues

• The chart currently supports **only** TCP ports in service > ports for Worker Groups. This restriction might be removed in future versions.

Sizing and Scaling

A Cribl LogStream installation can be scaled **up** within a single instance and/or scaled **out** across multiple instances. Scaling allows for:

- Increased data volumes of any size.
- Increased processing complexity.
- Increased deployment availability.
- Increased number of destinations.

Scale Up

A LogStream installation can be configured to scale up and utilize as many resources on the host as required. In a single-instance deployment, you govern resource allocation through the **General Settings** > **Worker Processes Settings** section.

In a distributed deployment, you allocate resources per Worker Group.

Navigate to Worker Groups > Group Name > System Settings > Worker

Processes.

Either way, these controls are available:

Process count: Indicates the number of Worker Processes to spawn.
 Positive numbers specify an absolute number of Workers. Negative numbers specify the number of Workers relative to the number of CPUs in the system. like this:

{ number of CPUs available minus this setting }. A 0 setting is interpreted as 1 Worker Process. LogStream corrects for excessive negative offsets by guaranteeing at least 1 Process. Defaults to -2.

- Memory (MB): Amount of memory available to each Worker Process, in MB.
 Defaults to 2048.
 - i Throughout these guidelines, we assume that 1 physical core is equivalent to 2 virtual/hyperthreaded CPUs (vCPUs). Each LogStream

instance requires the following resources to run, beyond those reserved for the vCPU's operating system:

- +4 physical cores, +8GB RAM
- 5GB free disk space (more if persistent queuing is enabled)

For example, assuming a Cribl LogStream system with 6 physical cores (12 vCPUs):

- If **Process count** is set to 4, then the system will spawn 4 processes, using up to 4 vCPUs, leaving 8 free.
- If **Process count** is set to -2, then the system will spawn 10 processes (12-2), using up to 10 vCPUs. This will leave 2 vCPUs free.
 - i LogStream incorporates guardrails that prevent spawning more processes than available vCPUs.

It's important to understand that worker processes operate in parallel, i.e., independently of each other. This means that:

- Data coming in on a single connection will be handled by a single worker process. To get the full benefits of multiple Worker Processes, data should come over multiple connections..
 - E.g., it's better to have 5 connections to TCP 514, each bringing in 200GB/day, than one at 1TB/day.
- 2. Each Worker Process will maintain and manage its own outputs. E.g., if an instance with 2 worker processes is configured with a Splunk output, then the Splunk destination will see 2 inbound connections.

Capacity and Performance Considerations

As with most data processing applications, Cribl LogStream's expected resource utilization will be commensurate with the type of processing that is occurring. For instance, a function that adds a static field on an event will likely perform faster than one that applies a regex to finding and replacing a string. At the time of this writing:

A Worker Process will use up to 1 physical core, or 2 vCPUs.

- Processing performance is proportional to CPU clock speed.
- All processing happens in-memory.
- Processing does not require significant disk allocation.

Estimating Requirements

Current guidance for capacity planning is: Allocate 1 physical core for each 400GB/day of IN+OUT throughput. So, to estimate the number of cores needed: Sum your expected input and output volume, then divide by 400GB.

- Example 1: 100GB IN -> 100GB out to each of 3 destinations = 400GB total = 1 physical core.
- Example 2: 3TB IN -> 1TB out = 4TB total = 10 physical cores.
- Example 3: 4 TB IN -> full 4TB to Destination A, plus 2 TB to Destination B = 10TB total = 25 physical cores.

Recommended AWS, Azure, and GCP Instance Types

You could meet the requirement above with multiples of the following instances:

AWS – Compute Optimized Instances. For other options, see here.

Minimum	Recommended
c5d.2xlarge (4 physical cores, 8vCPUs)	c5d.4xlarge or higher (8 physical cores, 16vCPUs)
c5.2xlarge (4 physical cores, 8vCPUs)	c5.4xlarge or higher (8 physical cores, 16vCPUs)

Azure – Compute Optimized Instances

Minimum	Recommended
Standard_F8s_v2 (4 physical cores, 8vCPUs)	Standard_F16s_v2 or higher (8 physical cores, 16vCPUs)

GCP – Compute Optimized Instances

Minimum	Recommended
c2-standard-8 (4 physical cores, 8vCPUs)	c2-standard-16 or higher (8 physical cores, 16vCPUs)
n2-standard-8 (4 physical cores, 8vCPUs)	n2-standard-16 or higher (8 physical cores, 16vCPUs)

Scale Out

When data volume, processing needs, or other requirements exceed what a single instance can sustain, a Cribl LogStream deployment can span multiple nodes. This is known as a Distributed Deployment, and it can be configured and managed centrally by a single master instance. See Distributed Deployment for more details.

Config Files

Understanding Configuration Paths and Files

Even though all the Routes, Pipelines, and Functions can be managed from the UI, it's important to understand how the configuration works under the hood. At the time of this writing this is how configuration paths and files are laid on the filesystem.

\$CRIBL_HOME	Standalone Install: /path/to/install/cribl/
, <u>-</u>	Splunk App Install: \$SPLUNK_HOME/etc/apps/cribl/

All paths below are relative to \$CRIBL_HOME.

Default Configurations	default/cribl
Local Configurations	local/cribl
System Configuration	<pre>(default local)/cribl/cribl.yml See cribl.yml</pre>
API Configuration	(default local)/cribl/api.yml
Source Configuration	<pre>(default local)/cribl/inputs.yml See inputs.yml</pre>
Destination Configuration	<pre>(default local)/cribl/outputs.yml See outputs.yml</pre>
License Configuration	(default local)/cribl/licenses.yml
Regexes Configuration	(default local)/cribl/regexes.yml
Breakers Configuration	(default local)/cribl/breakers.yml
Limits Configuration	(default local)/cribl/limits.yml

Pipelines Configuration	<pre>(default local)/cribl/pipelines/<pname> Each pipeline's conf is contained therein</pname></pre>
Routes Configuration	(default local)/cribl/pipelines/routes.yml
Functions	<pre>(default local)/cribl/functions/<function_name> Each function's code, conf is contained therein</function_name></pre>
Functions Conf	<pre>(default local)/cribl/functions/<function_name>/ Each function's conf contained therein.</function_name></pre>

Configurations and Restart

- Any configuration changes resulting from UI interactions, for instance, changing the order of functions in a pipeline, or changing the order of routes, do not require restarts.
- All Cribl LogStream configuration file changes resulting from direct file manipulations in (bin|local|default)/cribl/... will require restarts.
- Worker Nodes might temporarily disappear from the Master's **Workers** tab while restarting.
- In the case of a Cribl App for Splunk, Splunk configurations file changes might or might not require restarts. Please check with recent Splunk docs.

Configuration Layering and Precedence

Similar to most *nix systems, Cribl configurations in local take precedence over those in default . There is no layering of configuration files.

△ Editing Configuration Files Manually

When config files ${\it must}$ be edited manually, save all changes in local .

cribl.yml

cribl.yml contains settings for configuring API and other system properties.

```
$CRIBL_HOME/default/cribl/cribl.yml
  # Address to bind to. Default: 0.0.0.0
  host: 0.0.0.0
  # Port to listen to. Default: 9000
  port: 9000
  # Flag to enable/disable UI. Default: false
  disabled : false
  # SSL Settings
  ssl:
    # SSL is enabled by default
    disabled: false
    # Path to private key
    privKeyPath: /path/to/privkey.pem
    # Path to certificate
    certPath: /path/to/cert.pem
auth:
  # Type of authentication.
  type: splunk
  host: localhost
  port: 8089
  ssl: true
workers: # worker processes, memory in MB
  memory: 2048
kms.local:
  # Encryption key management system settings. Default type: local.
  type: local
crypto:
  # Crypto settings.
  keyPath: $CRIBL_HOME/local/cribl/auth/keys.json
system:
  # Upgradability options: api, auto, false
  upgrade: api
  # Restart options: api, false
  restart: api
  # installType options: standalone, splunk-app
  installType: standalone
  # Flag to enable/disable intercom. Default: true
  intercom: true
license:
  accepted: true
# distributed mode: master | worker | single
distributed:
  mode: master
```

inputs.yml

inputs.yml contains settings for configuring inputs into Cribl.

\$CRIBL_HOME/default/cribl/inputs.yml inputs: # Input name local-splunk: # Input type type: splunk # Address to listen to for incoming events host: localhost # Port to listen to for incoming events port: 10000 secureTCPJSON: type: tcpjson disabled: false host: 0.0.0.0 port: 10002 tls: disabled: false privKeyPath: /opt/privkey.pem certPath: /opt/cert.pem requestCert: false rejectUnauthorized: false

ipWhitelistRegex: /.*/

authToken: ""

outputs.yml

outputs.yml contains settings for configuring outputs from Cribl. Also see Destinations for more info.

```
$CRIBL_HOME/default/cribl/outputs.yml
outputs:
  # Default output setting
  default:
   type: default
    defaultId: local-splunk
  # Output Name
  local-splunk:
    # Output type
    type: splunk
    # Output host address to send data from
    host: localhost
    # Output port to send data from
    port: 9999
  # Output name
  myFilesystemDestination:
    # Output type
    type: filesystem
    # Final destination path. Writable by Cribl.
    destPath: /path/to/destiation
    # Staging destination path. Writable by Cribl.
    stagePath: /tmp/foo
    # Partition schema for outputted files
    partitionExpr: >-
      `${host}/${sourcetype}`
    # Format of the output data
    format: json
    # The output filename prefix
    baseFileName: CriblOut
    # Compression options. None | Gzip
    compress: none
    # Maximum uncompressed output file size
    maxFileSizeMB: 32
    # Maximum amount of time to keep inactive files open.
    maxFileOpenTimeSec: 300
    # Maximum amount of time to keep inactive files open.
    maxFileIdleTimeSec: 30
    # Maximum number of files to keep open concurrently.
    maxOpenFiles: 100
  myS3Destination:
    # Output type
    type: s3
    # S3 bucket address
    bucket: s2.bucket.address.here
    # Prefix to append to files before uploading
    destPath: keyprefix
    # AWS API key, if not present will fallback on env.AWS_ACCESS_KEY_ID, or the meta-data endpo
    awsApiKey: key
    # AWS Secret Key. If left blank, Cribl will fallback on env.AWS_SECRET_ACCESS_KEY, or the r
    awsSecretKey: secretkey
```

Staging destination path. Writable by Cribl.

stagePath: /tmp/foo

licenses.yml

licenses.yml maintains a list of licenses for Cribl.

\$CRIBL_HOME/default/cribl/licenses.yml

licenses:

- # List of license keys
- eyJ0eXAiOiJKV1QiLCJhasdfasfasdfdasfasdfa-Abo2_ogVbR_5VKeAelZlTc5b-TKQax9R1ywnoOG8guis2RC0sSl

regexes.yml

\$CRIBL_HOME/default/cribl/regexes.yml

regexes.yml maintains a list of regexes. Cribl's Regex Library ships under default.

```
...
"uuid":
lib: cribl
```

```
description: UUID/GUID
 regex: /[0-9a-f]\{8\}-[0-9a-f]\{4\}-[1-5][0-9a-f]\{3\}-[89ab][0-9a-f]\{3\}-[0-9a-f]\{12\}/gm
 sampleData: 9a50fa34-58b1-4a67-8b8d-ea9c0ae48c8f
   eb671525-2b9e-4140-ae21-a0a8a81b506e
 tags: uuid,guid
"aws_secret_key":
 description: AWS Secret Access Key
 regex: /(?<![A-Za-z0-9]/+=])[A-Za-z0-9]/+=]{40}(?![A-Za-z0-9]/+=])/gm
 lib: cribl
 sampleData: wJalrXUtnFEMI/K7MDENG/bPxRfiCYEXAMPLEKEY
 tags: aws,access,key,secret
"aws_access_key":
 lib: cribl
 description: AWS Access Key ID
 regex: /(A3T[A-Z0-9]|AKIA|AGPA|AIDA|AROA|AIPA|ANPA|ANVA|ASIA)[A-Z0-9]\{16\}(?![A-Za-z0-9\/+=])/(2)
 sampleData: >-2
    AKIAIOSFODNN7EXAMPLE
 tags: aws,access,key
"private_key":
 description: Private key block
 regex: /----BEGIN (DSA|RSA|EC|PGP|OPENSSH) PRIVATE KEY(\sBLOCK)?----[\s\S]*/gm
 lib: cribl
 tags: ssh,openssh,dsa,ec,rsa,private key
"slack_token":
 lib: cribl
 description: Slack Token
 regex: /xox[p|b|o|a][\s\S]*/g
 sampleData: xoxp-23984754863-2348975623103
   xoxa-23984754863-2348975623103
   xoxb-23984754863-2348975623103
   xoxo-23984754863-2348975623103
 tags: slack, token
```

breakers.yml

timestampTimezone: local
maxEventBytes: 4096

Cribl's default Event Breaker Library is located under \$CRIBL_HOME/default/cribl/breakers.yml.

```
$CRIBL_HOME/default/cribl/breakers.yml
AWS Ruleset:
     lib: cribl
      description: Event breaking rules for common AWS data sources
      tags: flowlogs,elb,alb,loadbalancer,cdn
      rules:
            - name: AWS VPC Flow
                  condition: /^\d+\s+\d+\s+eni-\w+.*(OK|NODATA|SKIPDATA)?$/.test(_raw) || sourcetype='aws:
                   eventBreakerRegex: /[\n\r]+/
                  timestampAnchorRegex: /(?=\d{10}\s\d{10})/
                  timestamp:
                       type: format
                      length: 150
                       format: "%s"
                  timestampTimezone: utc
                   maxEventBytes: 1024
             - name: AWS ALB
                   condition: /^{(:https:|h2|wss:)}\d+-\d+-\d+.*?arn:aws:elasticloadbalancing/.test(_raw) ||
                   eventBreakerRegex: /[\n\r]+/
                   timestampAnchorRegex: /\w+\s/
                   timestamp:
                       type: format
                        length: 150
                        format: "%Y-%m-%dT%H:%M:%S.%f%Z"
                   timestampTimezone: local
                  maxEventBytes: 4096
             - name: AWS ELB
                   condition: /^d+-d+-d+.*?(?:d+)_3/.test(_raw) || sourcetype='aws:elb:accesslossibset || sourcet
                   eventBreakerRegex: /[\n\r]+/
                   timestampAnchorRegex: /^/
                   timestamp:
                       type: format
                        length: 150
                        format: "%Y-%m-%dT%H:%M:%S.%f%Z"
```

mappings.yml

Mapping ruleset configurations are located under \$CRIBL_HOME/local/cribl/mappings.yml . \$CRIBL_HOME/default/cribl/mappings.yml

```
...
rulesets:
 default: # ruleset name
   conf:
     functions:
        - filter: env.CRIBL_HOME.match('w0') # filter to match
          description: w0 # rule name/id
          final: true
          conf:
           add:
              - name: groupId
               value: "'myGroup42'" # group to map to
        - filter: env.CRIBL_HOME.match('w1')
          id: eval
          description: w1
          final: true
          conf:
           add:
              - name: groupId
               value: "'NewGroup22'"
  newruleset: # another ruleset
   conf:
      functions:
        - filter: (cpus>12 & env.CRIBL_HOME.match('w0')) || release.startsWith('18')
          description: catch all
          final: true
          conf:
            add:
              - name: groupId
               value: "'NewGroup2'"
```

instance.yml

```
Instance configuration is located under
$CRIBL_HOME/local/_system/instance.yml.
 $CRIBL_HOME/local/_system/instance.yml
 distributed:
     # mode master | worker | single
   mode: master
   master:
     host: 0.0.0.0
     port: 4203
     tls:
       disabled: true
     ipWhitelistRegex: /.*/
     authToken: criblmaster
     compression: none
     connectionTimeout: 5000
     writeTimeout: 10000
   group: default
   envRegex: /^CRIBL_/
   tags:
        - tag1
        - tag2
        - tag42
```

Licensing

Every Cribl LogStream download package ships with a Free license that allows for processing of up to 1 TB/day. LogStream Free and LogStream One licenses require sending anonymized telemetry metadata to Cribl. (For details, see Telemetry Data below). Cribl does not require a separate license for sending data from LogStream to LogStream, such as sending from one worker group managed by master node A to a different worker group managed by master node B. The same license used on master node A can be used on master node B in that situation.

Enterprise, Standard, and Sales Trial licenses do **not** require sending telemetry metadata, and are entitled to a defined, per-license daily ingestion volume.

This page summarizes all these license types.

Managing Licenses

You can add and manage licenses in **Settings** > **Licensing**. Click + **Add License** to paste in a license key provided to you by Cribl.

☐ License Expiration and Renewal

For LogStream v.2.2 and earlier, the latest Free license expires on: 2020-12-15T00:00:00+00:00

For LogStream v.2.3 and later, Free licenses do not expire.

LogStream One and LogStream Standard licenses must be renewed annually.

License Types

Cribl offers five LogStream license types, summarized below.

i For a detailed comparison of what's included in each license type, please see Cribl Pricing.

Enterprise License

This is a license available for purchase.

- Up to unlimited data ingestion.
- Role-based access control.
- External authentication (via LDAP, Splunk, and OpenID Connect identity providers).
- Git remote backup.
- All other LogStream features included.

Contact Cribl Sales at sales@cribl.io for more information.

Standard License

This is a license available for purchase. Compared to an Enterprise license, it offers a cost discount, in exchange for some limitations (all data volumes below based on uncompressed data size):

- Daily ingestion up to 5 TB/day.
- Maximum 1 Worker Group.

Contact Cribl Sales at sales@cribl.io for more information.

Free License

Free licenses ship in the download package, and are permanent. They impose some limitations:

- Daily ingestion up to 1 TB/day.
- Maximum 10 Worker Processes.
- Maximum 1 Worker Group.

"One" License

LogStream One is a type of free license that allows for higher processing volume, but only to **one** Splunk (Single-Instance or Load-Balanced) or Elasticsearch Destination. This combination is designed to help users explore LogStream's value in routing large data volumes to these common services. Contact Cribl Sales at sales@cribl.io to convert a Free license to a LogStream One license, which must be renewed annually.

- Daily ingestion up to 5 TB/day, only to one of either Splunk or Elasticsearch outputs.
- Maximum 50 Worker Processes
- Maximum 1 Worker Group

Sales Trial License

A license type used when preparing a POC (proof of concept), or a pilot, with requirements that go beyond those afforded by the Free or One license. Contact Cribl Sales at sales@cribl.io for more information.

i LogStream Free and LogStream One licenses require sending of anonymized telemetry metadata to Cribl. These licenses will block inputs if sending fails after a grace period of 24 hours.

Combining License Types

Multiple license types can coexist on an instance. However, only a **single type** of license can be effective at any one time. When multiple types coexist, the following method of resolution is used:

- If there are any unexpired Enterprise or Standard licenses use only these licenses to compute the effective license.
- Else, if there are any Sales Trial licenses use only Sales Trial licenses to compute the effective license.
- Else, if there exists a Free or One license use only the Free or One license to compute the effective license.

When an Enterprise or Standard license expires, Cribl LogStream will fall back to the Sales Trial or Free/One types. However, an expired Sales Trial license cannot fall back to a Free/One license.

Upon expiration of a paid license, if there is no fallback license, LogStream will backpressure and block all incoming data.

Licensing in Distributed Deployments

LogStream 2.2.x or Earlier

In distributed deployments of LogStream versions through 2.2.x, licenses should be configured both on the Master Node and on each of the Worker Groups. This allows for different Worker Groups to have different licensing capacities.

- To configure the Master: **Settings** > **Licensing**.
- To configure Worker Groups: Worker Groups > [Select a Group] > System
 Settings > Licensing.

LogStream 2.3.x or Later

☐ As of LogStream 2.3, you no longer need to add licenses directly to Worker Groups. The Master will push license information down to Worker Groups as part of the heartbeat.

LogStream will attempt to balance (or rebalance) Worker Processes/threads as evenly as possible across all licensed Worker Nodes.

△ LogStream 2.3 changes licensing in other ways that might require you to update an existing LogStream configuration. Please see Upgrading to LogStream 2.3.

Telemetry Data

A Free or One license requires sharing of telemetry metadata with Cribl. Cribl uses this metadata to help us understand how to improve the product and prioritize new features. Telemetry payloads are sent to an endpoint located on https://cdn.cribl.io/telemetry/. (For versions prior to 2.2, this endpoint is 34.220.85.61:8000.)

i Testing telemetry endpoint connectivity

To manually test connectivity to the telemetry endpoint, especially if you are needing to configure a proxy, you can use the following

command:

```
$ curl https://cdn.cribl.io/telemetry/
Expected response:
```

cribl /// living the stream!

If you get a 302 response code then you may be missing the / at the end of the URL.

If you would like this feature disabled in order to deploy in your environment, please contact Cribl Sales at sales@cribl.io, and we will work with you to issue licenses on a case-by-case basis.

i Once you have received a license that removes the telemetry requirement, disable telemetry in LogStream's UI at Settings > System > General Settings > Upgrade & Share Settings > Sharing and Live Help. Toggle the slider to No.

Data Shared Per Interval (roughly, every minute):

- Version
- Instance's GUID
- License ID
- Earliest, Latest Time
- Number of Events In, Out
- Number of Bytes In, Out
- Number of Open, Closed, Active Connections
- Number of Routes
- Number of Pipelines

Licensing FAQ

How do I check my license type, restrictions, and/or expiration date?

Open LogStream's **Settings > Licensing** page to see these details.

How can I track my actual data ingestion volume over the last 30 days?

Forward Cribl Internal metrics to your Metrics Destination of choice, and run a report on cribl.total.in_bytes.

How does LogStream enforce license limits?

If your data throughput exceeds your license quota, Chuck Norris will track you down and make your life a living hell.

However, that will happen only in your nightmares. In the product itself:

- Free, One, and Standard licenses enforce data ingestion quotas through limits on the number of Worker Groups and Worker Processes.
- Enterprise license keys turn off all enforcement, between annual true-ups.
- When an Enterprise or Standard license expires, LogStream will attempt to fall back to a trial or free license, or only if that fails will block incoming data. For details, see Combining License Types.

I'm using LogStream 2.3.0 or higher, with its "permanent, Free" license. Why is LogStream claiming an expired license, and blocking inputs?

This can happen if you've upgraded from a LogStream version below 2.3.0, in which you previously entered this earlier version's Free (time-limited) license key. To remedy this, go to **Settings > Licensing**, click to select and expand your expired Free license, and then click **Delete license**. LogStream will fall back to the new, permanent Free license behavior, and will restore throughput.

If I pull data from compressed S3 buckets, is my license quota applied to the compressed or the uncompressed size of the file objects?

To measure license consumption, LogStream uses the uncompressed size.

Access Management

Cribl LogStream provides a range of access-management features for users with different security requirements. For details, see the following topics:

- Authentication: Authenticating users in LogStream.
- Local Users: Creating and managing users and their permissions.
- Roles: Managing roles and policies to assign to users.
 - i Role-based access control can be enabled only on distributed deployments with an Enterprise license.

Authentication

User authentication in LogStream

Cribl LogStream supports **local**, **Splunk**, **LDAP**, and **SSO/OpenID Connect** authentication methods.

Local Authentication

To set up local authentication, navigate to **Settings > General Settings > Authentication Settings** and select **Local**.

You can then manage users through the **Settings > Local Users** UI. All changes made to users are persisted in a file located at \$CRIBL HOME/local/cribl/auth/users.json.

Line format:

```
{"username":"user","first":"Elvis","last":"Bath","disabled":"false",
"passwd":"Yrt0MOD1w80zyMYB8WMcEleOtYESMwZw2qIZyTvueOE"}
```

The file is monitored for modifications every 60s, and will be reloaded if changes are detected.

Adding users through direct modification of the file is also supported, but not recommended.

Manual Password Replacement

To manually add, change, or restore a password, replace the affected user's passwd key-value pair with a password key, in this format: "password":" <newPlaintext>" . LogStream will hash all plaintext password(s), identified by the password key, during the next file reload, and will rename the plaintext password key.

Starting with the same users.json line above:

```
{"username":"user","first":"Elvis","last":"Bath","disabled":"false",
"passwd":"Yrt0MOD1w80zyMYB8WMcEleOtYESMwZw2qIZyTvueOE"}
```

...you'd modify the final key-value pair to something like:

```
{"username":"user","first":"Elvis","last":"Bath","disabled":"false",
"password":"V3ry53CuR&pW9"}
```

Within at most one minute after you save the file, LogStream will rename the password key back to passwd, and will hash its value, re-creating something resembling the original example.

Set Worker Passwords

In a distributed deployment, once a worker has been set to point to the Master Node, LogStream will set each Worker node's admin password with a randomized password which is different from the admin user's password on the Master Node. This is by design, as a security precaution, but may lead to situations where administrators cannot log into a Worker Node directly and must rely on accessing them via the Master

To explicitly push a known/new password to your Worker Node, set and push a new password to the Worker Group.

In the Master Node's UI:

- 1. From the top menu, select Worker Groups.
- 2. Select the desired Worker Group.
- 3. From the Worker Groups submenu, select System Settings.
- 4. Select **Local Users**, then expand the desired user.
- 5. Update the **Password** field and select **Save**.

Every 10 seconds, the Worker Nodes will request an update of configuration from the Master and new password settings will be reflected.

Authentication Controls

You can customize authentication behavior at **General Settings > API Server Settings > Advanced.** The options here include:

- Logout on Roles change: If role-based access control is enabled, determines whether users are automatically logged out of LogStream when their assigned Roles change. Defaults to Yes.
- Auth-token TTL: Sets authentication tokens' valid lifetime, in seconds. Defaults to 3600 (60 minutes).

- Login rate limit: Sets the number of login attempts allowed over a (selectable) unit of time. Defaults to 2/second.
- **HTTP header**: Enables you to specify one or more custom HTTP headers to be sent with every response.

The cribl.secret File

When Cribl LogStream first starts, it creates a \$CRIBL_HOME/local/cribl/auth/cribl.secret file. This file contains a key that is used to generate auth tokens for users, encrypt their passwords, and encrypt encryption keys.

Default local credentials are: admin/admin

Back up and secure access to this file by applying strict permissions – e.g., 600.

Splunk Authentication

Splunk authentication is very helpful when deploying in the same environment as Splunk, and requires the user to have Splunk admin role permissions. To set up Splunk authentication:

Navigate to **Settings > General Settings > Authentication Settings** and select **Splunk**.

- Host: Splunk hostname (typically a search head).
- Port: Splunk management port (defaults to 8089).
- SSL: Set to Yes if enabled.
- Fallback to local: Attempt local authentication if Splunk authentication is unsuccessful. Defaults to false.

Notes: The Splunk searchhead does not need to be locally installed on the LogStream instance. See also Role Mapping below.

LDAP Authentication

LDAP authentication is supported, and can be set up as follows:

Navigate to Settings > General Settings > Authentication Settings, and select LDAP.

- Secure: Enable to use a secure LDAP connections (ldaps://). Disable for an insecure (ldap://)connection.
- LDAP servers: List of LDAP servers. Each entry should contain host:port (e.g., localhost:389).
- Bind DN: Distinguished name of entity to authenticate with LDAP server. E.g., 'cn=admin,dc=example,dc=org'.
- Password: Distinguished Name password used to authenticate with LDAP server.
- User search base: Starting point to search LDAP for users, e.g., 'dc=example,dc=org'.
- Username field: LDAP user search field, e.g., cn or (cn (or uid).
- User search filter: LDAP search filter to apply when finding user, e.g., (& (group=admin)(!(department=123*))).Optional.
- **Group search base**: Starting point to search LDAP for groups, e.g., dc=example,dc=org.Optional.
- **Group member field**: LDAP group search field, e.g., member . Optional.
- **Group search filter**: LDAP search filter to apply when finding group, e.g., (& (cn=cribl*)(objectclass=group)).Optional.
- Fallback to local: Attempt local authentication if LDAP authentication is down or is mis-configured. Defaults to No.
- Connection timeout (ms): Defaults to 5000.
- Reject unauthorized: Valid for secure LDAP connections. Set to Yes to reject unauthorized server certificates.
- **Group name field**: LDAP group field, e.g., cn.

Note: See also Role Mapping below.

SSO/OpenID Connect Authentication

LogStream supports SSO/OpenID user authentication (login/password) and authorization (user's group membership, which you can map to Cribl Roles). Set this up as follows:

Navigate to **Settings > General Settings > Authentication Settings** and select **OpenID Connect**.

- **Provider name**: The name of the identity provider service. You can select **Google** or **Okta**, both supported natively. Manual entries are also allowed.
- Audience: The Audience from provider configuration. This will be the base URL, e.g.: https://yourDomain.com:9000.
- **Client ID**: The client_id from provider configuration.
- Client secret: The client_secret from provider configuration.
- **Scope**: Space-separated list of authentication scopes. The default list is: openid profile email.
- Authentication URL: The full path to the provider's authentication endpoint.

 Be sure to configure the callback URL at the provider as

 <yourDomainUrl>/api/v1/auth/authorization-code/callback, e.g.:

 https://yourDomain.com:9000/api/v1/auth/authorization-code/callback.
- Token URL: The full path to the provider's access token URL.
- Logout URL: The full path to the provider's logout URL. Leave blank if the provider does not support logout or token revocation.
- Validate certs: Whether to validate certificates. Defaults to Yes . Toggle to No to allow insecure self-signed certificates.
- **Filter type**: Select either **Email whitelist** or **User info filter**. This selection displays one of the following fields:
 - Email whitelist: Wildcard list of emails that are allowed access.
 - User info filter: JavaScript expression to filter against user profile attributes. E.g.: name.startsWith("someUser") &6 email.endsWith("domain.com")
- **Group name field**: Field on the id_token that contains the user groups. Defaults to groups.

Note the following details when filling in the form – for example, when using Okta:

- <Issuer URI> is the account at the identity provider.
- Audience is the URL of the host that will be connecting to the Issuer (e.g., https://localhost:9000). The issuer (Okta, in this example) will redirect back to this site upon authentication success or failure.

See also Role Mapping below.

Role Mapping

This section is displayed only on distributed deployments with an Enterprise license. For details on mapping your external identity provider's configured groups to corresponding LogStream user access Roles, see External Groups and LogStream Roles.

- **Default role**: Default LogStream Role to assign to all groups not explicitly mapped to a Role.
- Mapping: On each mapping row, enter an external group name on the left, and select the corresponding LogStream Role on the right drop-down list. Click + Add Mapping to add more rows.

Local Users

This page covers how to create and manage LogStream users, including their credentials and (where enabled) their access roles. These options apply if you're using the **Local** Authentication type, which is detailed here.

Creating and Managing Local Users

On the Master Node – or in a single-instance deployment – you manage users by selecting Settings > Access Management > Local Users.

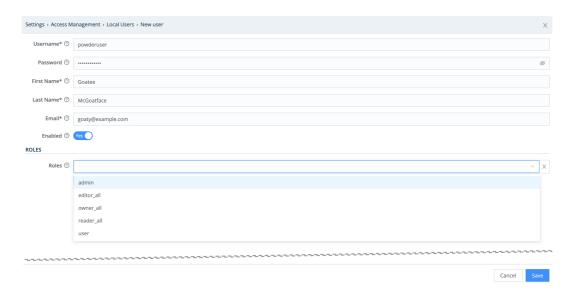
The resulting **Manage Local Users** page will initially show only the default admin user. You are operating as this user.



Managing users

To create a new LogStream user, click + Add New. To edit an existing user, click anywhere on its row. With either selection, you will see the modal shown below.

The first few fields are self-explanatory: they establish the user's credentials. If you want to establish or maintain a user's credentials on LogStream, but prevent them from currently logging in, you can toggle the **Enabled** slider to No.



Adding Roles

If you've enabled role-based access control you can use the modal's bottom **Roles** section to assign access Roles to this new or existing user.

i For details, see Roles. Role-based access control can be enabled only on distributed deployments with an Enterprise license.

Click + Add Role to assign each desired role to this user. The options on the Roles drop-down reflect the Roles you've configured in Settings > Access Management > Roles.

Note that when you assign multiple Roles to a user, the Roles' permissions are additive: This user is granted a superset of the highest permissions contained in all the assigned Roles.

When you've configured (or reconfigured) this user as desired, click **Save**.

By default, LogStream will log out a user upon a change in their assigned Roles. You can defeat this behavior at **General Settings > API Server Settings > Advanced > Logout on roles change.**

Roles

Define and manage access-control roles and policies

Cribl LogStream offers role-based access control (RBAC) to serve these common enterprise goals:

- **Security**: Limit the blast radius of inadvertent or intentional errors, by restricting each user's actions to their needed scope within the application.
- Accountability: Ensure compliance, by restricting read and write access to sensitive data.
- Operational efficiency: Match enterprise workflows, by delegating access over subsets of objects/resources to appropriate users and teams.
 - i Role-based access control is enabled only on distributed deployments with an Enterprise license.

RBAC Concepts

LogStream's RBAC mechanism is designed around the following concepts, which you manage in the UI:

- Roles: Logical entities that are associated with one or multiple Policies (groups of permissions). You use each Role to consistently apply these permissions to multiple LogStream users.
- **Policies**: A set of **permissions**. A Role that is granted a given Policy can access, or perform an action on, a specified LogStream object or objects.
- Permissions: Access rights to navigate to, view, change, or delete specified objects in LogStream.
- Users: You map Roles to LogStream users in the same way that you map user groups to users in LDAP and other common access-control frameworks.

Users are independent LogStream objects that you can configure even without RBAC enabled. For details, see Local Users.

How LogStream RBAC Works

LogStream RBAC is designed to grant arbitrary permissions over objects, attributes, and actions at arbitrary levels.

i As of v. 2.4, Roles are customizable only down to the Worker Group level. E.g., you can grant Edit permission on Worker Group WG1 to User A and User B, but cannot grant them finer-grained permissions on child objects such as Pipelines, Routes, etc.

LogStream's UI will be presented differently to users assigned Roles with access restrictions. Controls will be visible but disabled, and search and log results will be limited, depending on each user's permissions.

LogStream Roles can be integrated with external authorization/IAM mechanisms, such as LDAP and OIDC and mapped to their respective groups, tags, etc.

Using Roles

LogStream ships with a set of default Roles, which you can supplement.

Default Roles

These Roles ship with LogStream by default:

Name	Description
admin	Superusers – authorized to do anything and everything in the system.
owner_all	Read/write access to (and Deploy permission on) all Worker Groups.
editor_all	Read/write access to all Worker Groups.
reader_all	Read-only access to all Worker Groups.
user	Default role that gets only a home/landing page to authenticate. This

is a fallback for users who have not yet been assigned a higher role by an admin.

Cribl **strongly recommends** that you do not edit or delete these default roles. However, you can readily clone them (see **Clone Role** below), and modify the duplicates to meet your needs.

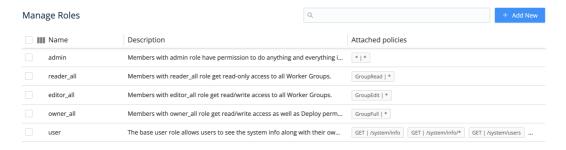
☐ Initial Installation or Upgrade

When you first install LogStream with the prerequisites to enable RBAC (Enterprise license and distributed deployment), you will be granted the **admin** role. Using this role, you can then define and apply additional roles for other users.

You will similarly be granted the **admin** role upon upgrading an existing LogStream installation from pre-2.4 versions to v. 2.4 or higher. This maintains backwards-compatible access to everything your organization has configured under the previous LogStream version's single role.

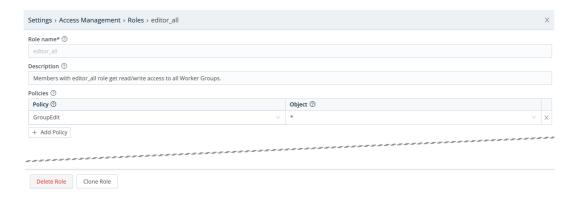
Adding and Modifying Roles

In a distributed environment, you manage Roles at the Master level, for the entire deployment. On the Master Node, select **Settings** > **Access Management** > **Roles**.



Manage Roles page

To add a new Role, click + Add New at the upper right. To edit an existing Role, click anywhere on its row. Here again, either way, the resulting modal offers basically the same options.



Add/edit Role modal

The options at the modal's top and bottom are nearly self-explanatory:

Role name: Unique name for this Role.

Description: Optional free-text description.

Delete Role: And...it's gone. (But first, there's a confirmation prompt.)

Clone Role: Opens a **New role** version of the modal, duplicating the **Description** and **Policies** of the Role you started with.

The modal's central **Policies** section (described below) is its real working area.

Adding and Modifying Policies

The **Policies** section is an expandable table. In each row, you select a Policy using the left drop-down, and apply that Policy to objects (i.e., assign permissions on those objects) using the right drop-down.

Let's highlight an example from the above screen capture of LogStream's built-in Roles: The editor_all Role has the GroupEdit Policy, with permission to exercise it on any and all Worker Groups (as indicated by the * wildcard).



Policies on the left, objects on the right

To add a new Policy to a Role:

- 1. Click + Add Policy to add a new row to the Policies table.
- 2. Select a Policy from the left column drop-down.

3. Accept the default object on the right; or select one from the drop-down.

To modify an already-assigned Policy, just edit its row's drop-downs in the **Policies** table.

To remove a Policy from the Role, click its close box at right.

In all cases, click **Save** to confirm your changes and close the modal.

Default Policies

In the **Policies** table's left column, the drop-down offers the following default Policies:

Name	Description
GroupRead	The most basic Worker Group-level permission. Enables users to view a Worker Group and/or its configuration.
GroupEdit	Building on GroupRead, grants the ability to also change and commit a Worker Group's configuration.
GroupFull	Building on GroupEdit, grants the ability to also deploy a Worker Group.
* (wildcard)	Grants all permissions on associated objects.

Objects and Permissions

In the **Policies** table's right column, use the drop-down to select the LogStream objects on which the left column's Policy will apply. (Remember that in v. 2.4, the objects available for selection are specific Worker Groups, or a wildcard representing all Worker Groups.) For example:

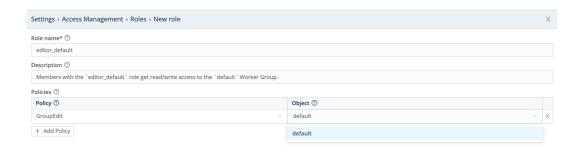
- Worker Group <id>
- NewGroup2
- default (Worker Group)
- * (all Worker Groups)

Extending Default Roles

Here's a basic example that ties together the above concepts and facilities. It demonstrates how to add a Role whose permissions are restricted to a particular Worker Group.

Here, we've cloned the editor_all Role that we unpacked above. We've named the clone editor_default .

We've kept the GroupEdit Policy from editor_all. But in the right column, we're restricting its object permissions to the default Worker Group that ships with LogStream.



Cloning a default Role

You can readily adapt this example to create a Role that has permissions on an arbitrarily named Worker Group of your own.

Roles and Users

Once you've defined a Role, you can associate it with LogStream users. On the Master Node, select **Settings** > **Access Management** > **Local Users**. For details, see Local Users.

Note that when you assign multiple Roles to a given user, the Roles' permissions are additive: This user is granted a superset of all the permissions contained in all the assigned Roles.

By default, LogStream will log out a user upon a change in their assigned Roles. You can defeat this behavior at **General Settings > API Server Settings > Advanced > Logout on roles change.**

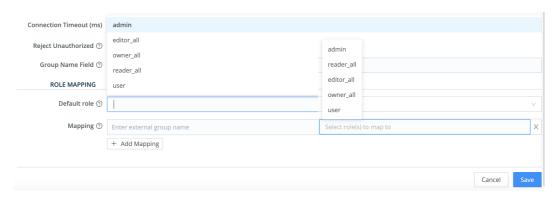
External Groups and LogStream Roles

You can map user groups from external identity providers (LDAP, Splunk, or OIDC) to LogStream Roles, as follows:

 On the Master Node, select Settings > Access Management > Authentication.

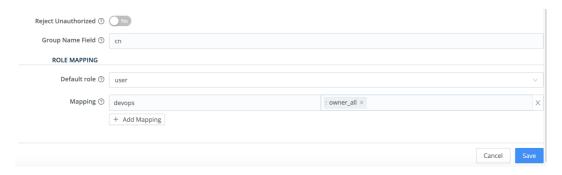
- 2. From the **Type** drop-down, select **LDAP**, **Splunk**, or **OpenID Connect**, according to your needs.
- 3. On the resulting **Authentication Settings** page, configure your identity provider's connection and other basics. (For configuration details, see the appropriate **Authentication** section.)
- 4. Under **Role Mapping**, first select a LogStream **Default role** to apply to external user groups that have no explicit LogStream mapping defined below.
- 5. Next, map external groups as you've configured them in your external identity provider (left field below) to LogStream Roles (right drop-down list below).
- 6. To map more user groups, click + Add Mapping.
- 7. When your configuration is complete, click **Save**.

Here's a composite showing the built-in Roles available on both the **Default Role** and the **Mapping** drop-downs:



Mapping external user groups to LogStream Roles

And here, we've set a conservative **Default Role** and one explicit **Mapping**:



External user groups mapped to LogStream Roles

Page 749 of 1257		

Version Control

Tracking, backing up, and restoring configuration changes for single-instance and distributed deployments

Cribl LogStream integrates with Git clients and remote repositories to provide version control of LogStream's configuration. This integration offers backup and rollback for single-instance and distributed deployments.

These options are separate from the Git repo responsible for version control of Worker configurations, located on the Master Node in distributed deployments. We cover all these options and requirements below.

Git Installation (Local or Standalone/Single-Instance)

To verify that git is available, run:

git --version

The minimum version that LogStream requires is: **1.8.3.1.** If you don't have git installed, see the installation links here.

Git Required for Distributed Deployments

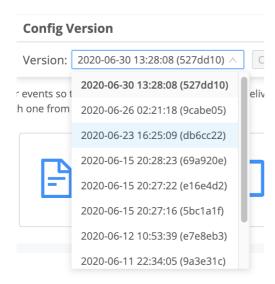
For distributed deployments, git must be installed and available locally on the host running the Master Node.

All configuration changes must be committed before they are deployed. The Master notifies Workers that a new configuration is available, and Workers pull the new configuration from the Master Node.

Reverting Commits

Once Git is installed, you can revert to a previous commit using the git CLI. You can also restore a Worker Group's previous commit using LogStream's UI:

Select the commit from the **Config Version** drop-down, as shown below. Then, in the resulting modal, verify the diff'ed configuration change and click **OK**.



Support For Remote Repositories

Git **remote** repositories are supported – but not required – for version control of all configuration changes. You can configure a Standalone Master Node with Git remote push capabilities through the LogStream CLI, or through the LogStream UI (via **Settings > Distributed Settings > Git Settings**).

To create a repo, see these tutorials:

- Setting Up a Repository (CLI instructions, host-agnostic, from Atlassian).
- Creating a New Repository (specific to GitHub's Web UI).
- Create a Repo (longer GitHub-specific tutorial, also covers committing changes).
 - i Currently, LogStream supports push and pull only against the master branch on each remote repo.

Several tutorial links and examples on this page point to GitHub, based on its wide adoption. The basic principles are the same for other Git repo providers, including private Git servers. GitHub's own UI and documentation periodically change, and linked tutorials' screenshots might differ from GitHub's current UI.

Remote Formats Supported

Remote URI schema patterns should match this regex:

```
(?:git|ssh|ftps?|file|https?|git@[-\w.]+):(\/\)?(.*?)(\.git\/?)?
```

You can find a list of supported formats here.

For example:

- Local Gitservers: git://<host.xyz>:<port>/<user>/path/to/repo.git

Connecting to a Remote with a Personal Access Token over HTTPS (Recommended)

Cribl recommends connecting to a remote repo over HTTPS. The example below shows a token-based HTTPS connection to GitHub.

Example: Connecting to GitHub over HTTPS

- 1. Create a new GitHub repository.
 - i For best results, create a new empty repo, with no readme file and no commit history. This will prevent git push errors.

Note the user name and email with which you log into the repo provider.

- 2. Create a personal access token with repo scope.
- 3. Copy the token to your clipboard.
- 4. In Cribl LogStream, go to Settings > Distributed Settings > Git Settings.
- 5. Fill in the **Remote URL** field with your repo name, user name, and token (in place of a password). Use the format below, replacing both <username> placeholders with your user name on the repo provider:

https://<username>:<token>@github.com/<username>/<reponame>.git

For additional details, see GitHub's Creating a Personal Access Token tutorial.

▲ For GitHub repos specifically, use only personal access tokens in the **Remote URL** field. GitHub has announced that it will end support for plaintext passwords as of August 13, 2021.

Connecting to a Remote with SSH

You can set up SSH keys from the CLI, or upload keys via the UI. If you have a passphrase set, this functionality is available only through the CLI – see Encryption: Configuring Keys with the CLI. The example below outlines the UI steps.

Example: Connecting to GitHub with SSH

1. Create a new GitHub repository.

i For best results, create a new empty repo, with no readme file and no commit history. This will prevent git push errors.

Note the user name and email with which you log into the repo provider.

- 2. Add an SSH public key to your GitHub account.
- 3. In Cribl LogStream, go to **Settings > Distributed Settings > Git Settings**.
- 4. Fill in the remote repo URL and the SSH private key. In the example format below, replace <username> with your user name on the repo provider:

Remote URL: <protocol>://git@github.com:<username>/<reponame>.git

SSH private key: <ssh-private-key>

For GitHub specifically, the URL/protocol format must be:

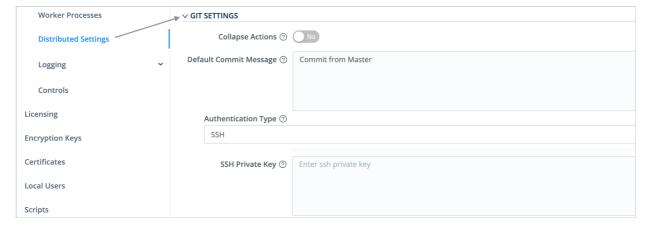
Remote URL: git@github.com:<user>/<reponame>.git

For example:

Remote URL: git@github.com:taylorswift/leadsheets.git

5. As the user running LogStream, run this command to add the GitHub keys to known_hosts: ssh-keyscan -H github.com >> ~/.ssh/known_hosts

For additional details, see GitHub's Connecting to GitHub with SSH tutorial.



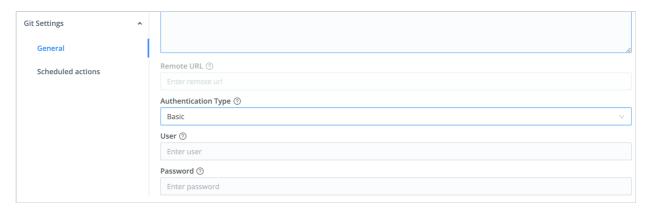
LogStream's Git settings

Additional Git Settings

On the **Git Settings** > **General** tab, you can change the **Authentication Type** from its **SSH** default to **Basic** authentication. This displays two additional fields:

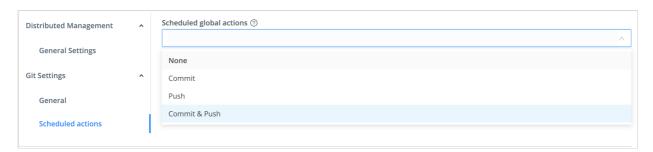
• User: Username on the repo.

• Password: Authentication password (e.g., a GitHub personal access token).



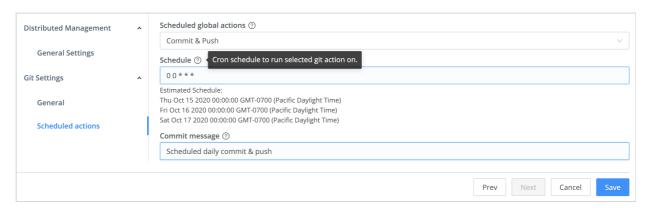
Git Authentication Type settings

On the **Git Settings** > **Scheduled Actions** tab, you can schedule a **Commit, Push**, or **Commit & Push** action to occur on a predefined interval.



Git Scheduled Actions selection

For the selected action type, you can define a [cron schedule] (cron schedule), and a commit message distinct from the **General** tab's **Default Commit Message**. Then click **Save**.

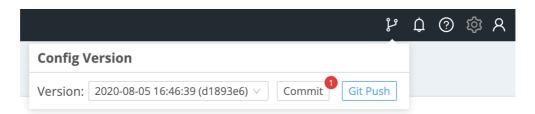


Saving a Git Scheduled Action

You can schedule only one type of action. To swap to a different type, select it from the **Scheduled global actions** drop-down, and resave. To turn off scheduled Git commands, select **None** from the drop-down, and resave.

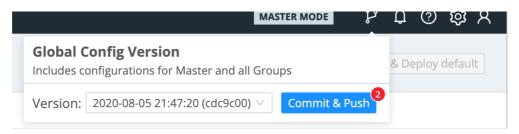
Pushing to a Remote Repo

Once you've configured a remote, a **Git Push** button appears in the Version Control overlay.



Git Push button

If you enabled the **Git Settings** > **Collapse Actions** option, you will instead see a combined **Commit & Push** button in the overlay.



Git combined actions button

Troubleshooting Push Errors

This section anticipates common errors you might see in LogStream's UI, or in the git CLI, when pushing a commit.

Failed to Push Some Refs

Your first push to a remote repo might fail with one of several failed to push some refs errors.

As a first step in debugging these errors, edit the \$CRIBL_HOME/.git/config file to make sure that its name and email key values match the credentials you've set on your repo provider or git server.

Also make sure that the remote "origin" key value matches the remote you set when you connected to the remote repo. This example shows all three keys, with placeholder values:

```
[user]
   name = <your-login-name>
   email = <email@example.com>
[remote "origin"]
   url = https://<user-name>:<token>@github.com/<username>/<repo-name>
```

Next, verify the remote repo from the command line, as follows:

```
cd $CRIBL_HOME/.git
git remote -v
```

In response, git should echo your configured remote twice – once for fetch and once for push operations.

If all of the above settings are correct, the push is very likely blocking because the remote repo has some commit history, or was simply created with a readme.md file. For command-line instructions to remedy this – by syncing your local repo to its remote – see GitHub's Dealing with Non-Fast-Forward Errors topic.

Large Files Detected

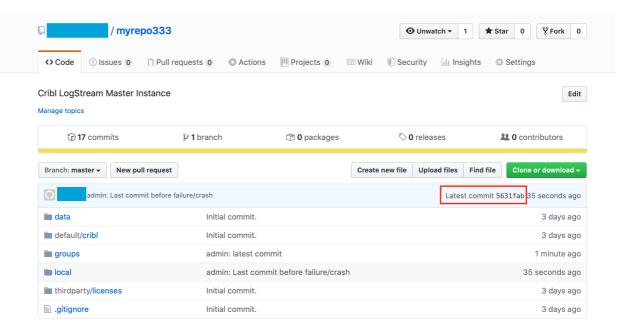
A push command might also trigger "large file" warnings or, more seriously, errors of this form (CLI/GitHub example):

Cribl recommends adding such large files to .gitignore , to exclude them from subsequent push commands. As the above examples show, typical culprits are large .csv or .mmdb lookup files. A simple option is to place these files in a \$CRIBL_HOME subdirectory that's already listed in .gitignore – for details, see Managing Large Lookups.

Other available workarounds include staging such files **outside** \$CRIBL_HOME, or using plugins to accommodate the large files. For GitHub-specific options, see Working with Large Files.

Restoring Master from a Remote Repo

If a remote repo is configured and has the latest known good Master configuration, this section outlines the general steps to restore the config from that repo.



Restoring from remote repo

Let's assume that the entire \$CRIBL_HOME directory of the Master is corrupted, or you're starting from scratch. Let's also assume that the remote repo has the form:

git@github.com:<username>/<reponame>.git.

- 1. **Important**: In a directory of choice, untar the **same Cribl LogStream version** that you're trying to restore, but do not start it.
- 2. If you are using SSH key authentication, specify the key using the following command:

```
GIT_SSH_COMMAND='ssh -i .key -o IdentitiesOnly=yes' git fetch origin
```

3. Ensure that you have proper access to the remote repo:

```
# git ls-remote git@github.com:/.git
56331fabb4822eaec4ca0ffd008d6e9974c1e419f HEAD
5631fabb4822eaec4ca0ffd008d6e9974c1e419f refs/heads/master
```

4. Change directory into \$CRIBL_HOME and initialize git:

```
# git init
```

5. Next, add/configure the remote:

```
# git remote add origin git@github.com:<username>/<reponame>.git
```

6. Now set up your local branch to exactly match the remote branch:

```
# git fetch origin
# git reset --hard origin/master
```

7. Finally, to confirm that the commits match, run this command while in \$CRIBL_HOME . Note the commit hash:

```
# git show --abbrev-commit commit 5631fab (HEAD \rightarrow master, origin/master)
```

```
Author: First Last
Date: Fri Jan 31 10:16:07 2020 -0500

admin: Last commit before failure/crash
.....
```

That last step above pulls in all the latest configs from the remote repo, and you should be able to start the Master as normal. Once up and running, Workers should start checking in after about 60 seconds.

∆ Verify cribl.secret

The cribl.secret file-located at \$CRIBL_HOME/local/cribl/auth/cribl.secret - contains the secret key that is used to encrypt sensitive settings on configuration files (e.g., AWS Secret Access Key, etc.). Make sure this file is properly restored on the new Master, because it is required to make encrypted conf file settings usable again.

.gitignore File

A .gitignore file specifies files that git should ignore when tracking changes. Each line specifies a pattern, which should match a file path to be ignored. Cribl LogStream ships with a .gitgnore file containing a number of patterns/rules, under a section of the file labeled CRIBL SECTION.

```
.gitignore
# Do NOT REMOVE CRIBL and CUSTOM header lines!
# DO NOT REMOVE rules under the CRIBL section as they may be reintroduced on update.
# You can ONLY comment out rules in the CRIBL section.
# You can add new rules in the CUSTOM section.
### CRIBL SECTION -- DO NOT REMOVE ###
default/ui/**
default/data/ui/**
bin/**
log/**
pid/**
data/uploads/**
diag/**
**/state/**
### CUSTOM SECTION -- DO NOT REMOVE ###
<User defined patterns/rules go here>
```

CRIBL Section

☐ Do Not Remove CRIBL SECTION or CUSTOM SECTION Headers

The CRIBL SECTION is used by Cribl LogStream to define default patterns/rules that ship with every version. Do **not** add or remove any of the lines here, because Chuck Norris will easily find you!

Maslow's theory of higher needs does not apply to Chuck Norris. He has only two needs: killing people and finding people to kill. Seriously, do not remove them, as they will be overwritten on the next update. The only modifications that will survive updates are commented lines.

CUSTOM Section

User-defined, custom patterns/rules can be **safely defined** under the CUSTOM SECTION . Cribl LogStream will **not** modify the contents of CUSTOM SECTION .

Good candidates to add here include large lookup files – especially large binary database files. See Troubleshooting: Large Files Detected, above.

Files skipped with .gitignore

If you have files that are skipped with .gitgnore, you will need to back them up and restore them via means other than Git. E.g., you can periodically copy/rsync them to a backup destination, and then restore them to their original locations after you complete the steps above.

Persistent Queues

Persistent queuing (PQ) is a feature that helps minimize data loss if a downstream receiver is unreachable. Durability is provided by writing data to disk for the duration of the outage, and forwarding it upon recovery.

PQs are implemented on the outbound side, meaning that each Source can take advantage of a Destination's queue.

How Does Persistent Queueing Work

Each LogStream output has an in-memory queue that helps it absorb temporary imbalances between inbound and outbound data rates. E.g., if there is an inbound burst of data, the output will store events in the queue, and output them at the rate that the receiver can sync (as opposed to blocking or dropping them). Only when this queue is full will the output impose backpressure upstream.

Backpressure behavior can be configured to either **block** or **drop**. In block mode, the output will refuse to accept new data until the receiver is ready. The system will back propagate block "signals" all the way back to the sender (assuming it supports backpressure, too). In drop behavior, the output will discard new events until the receiver is ready.

In some environments, the in-memory queues and their block/drop behavior are acceptable. Persistent queues serve environments where more durability is required (e.g., outages last longer than memory queues can sustain), or where upstream senders do not support backpressure (e.g., ephemeral/network senders).

Engaging persistent queues in these scenarios can help minimize data loss. Once the in-memory queue is full, the LogStream output will write its data to disk. Then, when the receiver is ready, the output will start draining the queues in FIFO (first in, first out) fashion.

Persistent Queue Details and Constraints

Persistent queues are:

- Available at the output side (i.e., after processing).
- Engaged only when all of the receivers of that output exert blocking.
- Drained when at least one receiver can accept data.
- Not infinite in size. I.e., if data cannot be delivered out, you might run out of disk space.
- Not able to fully protect in cases of application failure. E.g., in-memory data might get lost if a crash occurs.
- Not able to protect in cases of hardware failure. E.g., disk failure, corruption, or machine/host loss.

Persistent Queue Support

The following LogStream Destinations support Persistent Queuing:

- Splunk Single Instance
- Splunk Load Balanced
- Splunk HEC
- Kinesis
- Cloudwatch Logs
- SQS
- Azure Monitor Logs
- Azure Event Hubs
- StatsD
- StatsD Extended
- Graphite
- TCP JSON
- Syslog
- Elasticsearch
- Honeycomb
- InfluxDB
- Wavefront
- SignalFx

Configuring Persistent Queueing

Persistent Queueing is configured individually for each output that supports it. To enable persistent queueing, go to the output's (Destination's) configuration

page and set the **Backpressure Behavior** control to **Persistent Queueing**. This exposes the following additional controls:

- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Sufficient disk space is required for queuing to operate properly. You configure the minimum disk space in **Settings > General Settings > Limits > Min Free Disk Space**. If available disk space falls below this threshold, LogStream will stop maintaining persistent queues, and data loss will begin. The default is 5GB. Be sure to set this on your worker nodes rather than on the master node when in distributed mode.

Securing

You can secure Cribl LogStream's API and UI access by configuring SSL. To do so, you can use your own certs and private keys, or you can generate a pair with OpenSSL, as shown here:

```
openssl req -nodes -new -x509 -newkey rsa:2048 -keyout myKey.pem -out myCert.pem -days 420
```

This command will generate both a self-signed cert (certified for 420 days), and an unencrypted, 2048-bit RSA private key.

In the LogStream UI, you can configure the cert via **Settings** > **Certificates**. You can configure the key via:

- Settings > Encryption Keys (single-instance deployments), or
- Worker Groups > <group-name> > System Settings > Encryption Keys (distributed deployments).

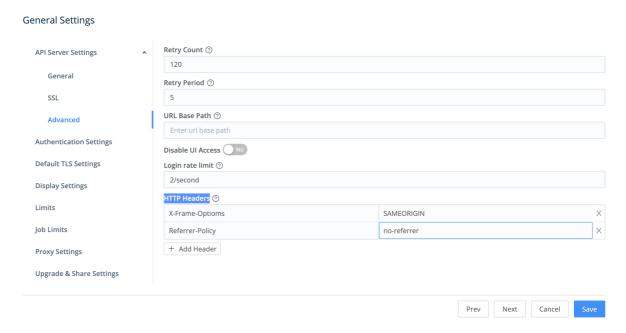
Alternatively, you can edit the local/cribl.yml file's api section to directly set the privKeyPath and certPath attributes. For example:

```
cribl.yml

api:
  host: 0.0.0.0
  port: 9000
  disabled: false
  ssl:
    disabled: false
    privKeyPath: /path/to/myKey.pem
    certPath: /path/to/myCert.pem
...
```

Custom HTTP Headers

You can encode custom, security-related HTTP headers, as needed. As shown in the examples below, you specify these at **Settings > General > API Server Settings > Advanced > HTTP Headers**. Click **+ Add Header** to display extra rows for new key-value pairs.



Custom HTTP headers

TLS Settings and Traffic Types

This table shows TLS client/server pairs, and encryption defaults, per traffic type.

Traffic Type	TLS Client	TLS Server	Encryption	Cert Auth	CN* Check
UI	Browser	Cribl LogStream	Default disabled	Default disabled	Default disabled
API	Worker	Master	Default disabled	Default disabled	Default disabled
Worker-to- Master	Worker	Master	Default disabled	Default disabled	Default disabled
Data	Any data sender	Cribl LogStream (Source)	Default disabled	Default disabled	Default disabled
Data	Cribl LogStream (Destination)	Any data receiver	Default disabled	Default disabled	Default disabled
Authentication					
• Local	Browser	Cribl LogStream	Default Disabled	N/A	N/A
• LDAP	Cribl LogStream	LDAP Provider	Custom	N/A	Default Disabled

• Splunk	Cribl LogStream	Splunk Search Head	Default Enabled	N/A	Default Disabled
• OIDC†/Okta	Browser and Cribl LogStream	Okta	Default Enabled	N/A	Enabled (Browser)
• OIDC/Google	Browser and Cribl LogStream	Google	Default Enabled	N/A	Enabled (Browser)

^{*} Common name

You can configure advanced, system-wide TLS settings for versions, cipher lists, and ECDH Curve names via Settings > System > General Settings > Default TLS Settings.

CA Certificates and Environment Variables

Where LogStream Sources and Destinations support TLS, each Source's or Destination's configuration provides a **CA Certificate Path** field where you can point to corresponding Certificate Authority (CA) .pem file(s). However, you can also use environment variables to manage CAs globally. Here are some common scenarios:

1. How do I add a set of trusted root CAs to the list of trusted CAs that LogStream trusts?

Set this environment variable in each Worker's environment (e.g., in its systemd unit file): NODE_EXTRA_CA_CERTS=/path/to/file_with_certs.pem . For details, see nodejs docs.

2. How do I make LogStream trust all TLS certificates presented by any server it connects to?

Set this environment variable: NODE_TLS_REJECT_UNAUTHORIZED=0 - for details, see nodejs docs.

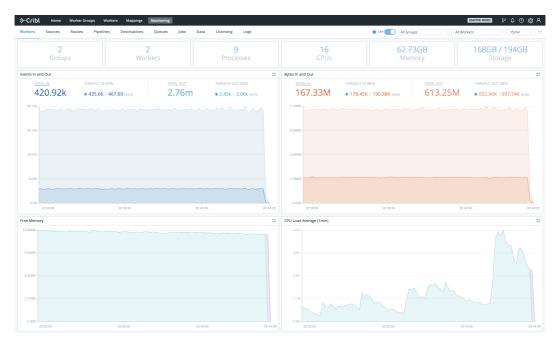
[†] OpenID Connect

Monitoring

To get an operational view of a Cribl LogStream deployment, you can consult the following resources.

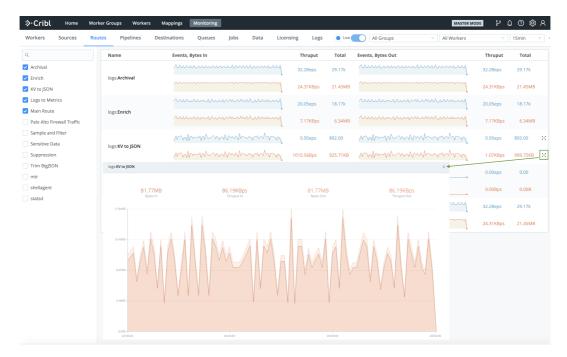
Monitoring Page

Select **Monitoring** from the top menu. This exposes information about traffic in and out of the system, as well as collection jobs and tasks. It tracks events, bytes, splits by data fields over time, and broader system metrics. Coverage is limited to the previous 24 hours. (Byte-related charts show the uncompressed size of processed data.)



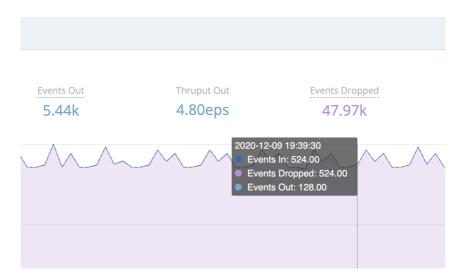
Monitoring page

Dense displays are condensed to sparklines for legibility. Hover over the right edge to display Maximize buttons that you can click to zoom these up to detailed graphs.



Sparklines and fly-out

You can hover over an expanded graph fly-out to display further details.



Throughput details

Internal Logs and Metrics

Select **Logs** from the **Monitoring** submenu. LogStream's **internal logs** and **internal metrics** provide comprehensive information about an instance's status/health, inputs, outputs, Pipelines, Routes, Functions, and traffic.

Health Endpoint

Query this endpoint on any instance to check the instance's health. (Details below.)

Types of Logs

LogStream provides the following log types, by originating process:

- API Server Logs These logs are emitted primarily by the API/main process. They correspond to the top-level cribl.log that shows up on the Diag page. Filesystem location: \$CRIBL_HOME/log/cribl.log
- Worker Process(es) Logs These logs are emitted by all the worker processes, and are very common in standalone instances or Worker Nodes. Filesystem location: \$CRIBL_HOME/log/worker/N/cribl.log
- Worker Group Logs These logs are emitted by all processes that help a
 Master Node configure Worker Groups. Filesystem location:
 \$CRIBL_HOME/log/group/GROUPNAME/cribl.log

LogStream rotates logs every 5 MB, keeping the most recent 5 logs. In a distributed deployment, all Workers forward their metrics to the Master Node, which then consolidates them to provide a deployment-wide view.

Forward Logs and Metrics Externally

LogStream supports forwarding internal logs and metrics to your preferred external monitoring solution. To send out internal data, go to **Data > Sources** and enable the **Cribl Internal** Source.

This will send all cribl.log logs and internal metrics down through Routes and Pipelines, just like another data source. Both logs and metrics will have a field called source, set to the value cribl, which you can use in Route filters.

For recommendations about useful Cribl metrics to monitor, see Internal Metrics.

i CriblMetrics Override

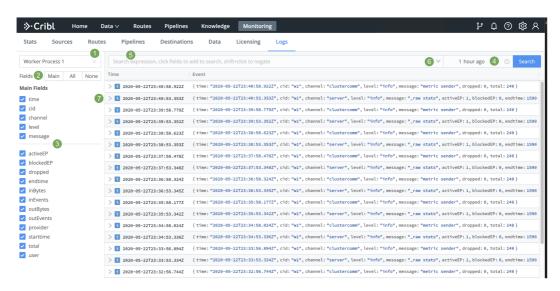
The **Disable field metrics** setting (in **Settings > System > General Settings > Limits**) applies only to metrics sent to the Master Node.
When the **Cribl Internal** Source is enabled, LogStream ignores this

Disable field metrics setting, and full-fidelity data will flow down the Routes.

Search Internal Logs

LogStream exists because logs are great and wonderful things! Using its **Monitoring > Logs** page, you can search all LogStream's internal logs at once – from a single location, for both Master and Worker Nodes. This enables you to query across all internal logs for strings of interest.

The labels on this screenshot highlight the key controls you can use (see the descriptions below):



Logs page (controls highlighted)

- 1. **Log file selector**: Choose the Node to view. In a **Distributed Deployment**, this list will be hierarchical, with Workers displayed inside their Master.
- 2. **Fields selector**: Click the **Main | All | None** toggles to quickly select or deselect multiple check boxes below.
- 3. **Fields**: Select or deselect these check boxes to determine which columns are displayed in the Results pane at right. (The upper **Main Fields** group will contain data for *every* event; other fields might not display data for all events.)
- 4. **Time range selector**: Select a standard or custom range of log data to display.

5. **Search box**: To limit the displayed results, enter a JavaScript expression here. An expression must evaluate to truthy to return results. You can press **Shift+Enter** to insert a newline.

Typeahead assist is available for expression completion:



Click a field in any event to add it to a query:



Click other fields to append them to a query:



Shift+click to negate a field:



- i To modify the depth of information that is originally input to the Logs page, see Logging Settings.
- 6. Click the Search box's history arrow (right side) to retrieve recent queries:



7. The Results pane displays most-recent events first. Each event's icon is color-coded to match the event's severity level.

Click individual log events to unwrap an expanded view of their fields:



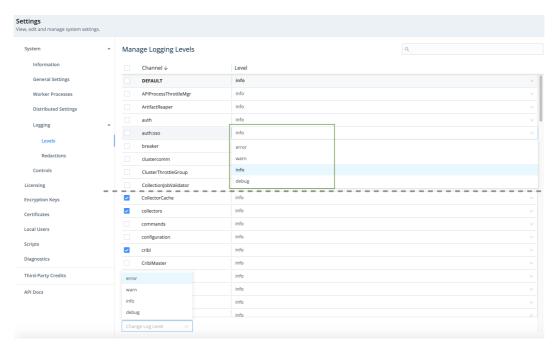
Logging Settings

Through LogStream's System Settings, you can adjust the level (verbosity) of internal logging data processed, per logging channel. You can also redact fields in customized ways.

Change Logging Levels

Select Settings > System > Logging > Levels to open the Manage Logging Levels page. Here, you can:

- Modify one channel by clicking its Level column. In the resulting dropdown, you can set a verbosity level ranging from error up to debug. (Top of composite screenshot below.)
- Modify multiple channels by selecting their check boxes, then clicking the
 Change log level drop-down at the bottom of the page. (Bottom of
 composite screenshot below.) You can select all channels at once by
 clicking the top check box. You can search for channels at top right.



Manage Logging Levels screen

Change Logging Redactions

Select **Settings > System > Logging > Redactions**: to open the **Redact Internal Log Fields** page. Here, you can customize the redaction of sensitive, verbose, or just ugly data within LogStream's internal logs.



Redact Internal Log Fields page

It's easiest to understand this page's fields from bottom to top:

- **Default fields**: LogStream always redacts these fields. You can't modify this list.
- Additional fields: Type or paste in the names of other fields you want to redact. Use a tab or hard return to confirm each entry.
- **Custom redact string:** Unless this field is empty, it defines a literal string that will override LogStream's default redaction pattern, explained below.

Default Redact String

By default, LogStream transforms this page's selected fields by applying the following redaction pattern:

- Echo the field value's first two characters.
- Replace all intermediate characters with a literal ... ellipsis.
- Echo the value's last two characters.

Anything you enter in the **Custom redact string** field will override this default ?? ... ?? pattern.

Health Endpoint

Each LogStream instance exposes a health endpoint – typically used in conjunction with a Load Balancer – that you can use to make operational decisions.

Health Check Endpoint		Healthy Response		
	<pre>curl http(s)://<host>:<port>/api/v1/health</port></host></pre>	{"status":"healthy"}		

Internal Metrics

When sending LogStream metrics to a metric system of analysis, such as InfluxDB, Splunk or Elasticsearch, some metrics are particularly valuable. You can use these metrics to set up alerts when a Worker Node is having a problem, a Node is down, a Destination is down, a Source stops providing incoming data, etc.

LogStream reports its internal metrics within the LogStream UI (in the same way that it reports internal logs at Monitoring > Logs). To expose metrics for capture or routing, enable the Cribl Internal Source > CriblMetrics section.

By default, LogStream generates internal metrics every 2 seconds. To consume metrics at longer intervals, you can use or adapt the <code>cribl-metrics_rollup</code> Pipeline that ships with LogStream. Attach it to your <code>Cribl Internal</code> Source as a <code>pre-processing Pipeline</code>. The Pipeline's <code>Rollup Metrics</code> Function has a default <code>Time Window</code> of 30 seconds, which you can adjust to a different granularity as needed.

You can also use our public endpoints to automate monitoring using your own external tools.

Total Throughput Metrics

Five important metrics below are prefixed with total. These power the top of LogStream's **Monitoring** dashboard. The first two report on Sources, the remainder on Destinations.

- total.in_bytes
- total.in_events
- total.out_events
- total.out_bytes
- total.dropped_events (new in LogStream 2.4) helpful for discovering situations such as: you've disabled a Destination without noticing.

Interpreting Total Metrics

These total. metrics' values could reflect LogStream's health, but could also report low activity simply due to the Source system. E.g., logs from a store site will be low at low buying periods.

Also, despite the total. prefix, these metrics are each specific to the Worker Process that's generating them.

You can distinguish unique metrics by their #input=<id> dimension. For example, total.in_events|#input=foo would be one unique metric; total.in_events|#input=bar would be another.

System Health Metrics

Five specific metrics are most valuable for monitoring system health. The first two are LogStream composite metrics; the remaining three report on your hardware or VM infrastructure.

- health.inputs
- health.outputs see the JSON Examples below for both health. metrics.
- system.load_avg
- system.free_mem
- system.disk_used valuable if you know your disk size, especially for monitoring Persistent Queues. Here, a 0 value typically indicates that the disk-usage data provider has not yet provided the metric with data. (Getting the first value should take about one minute.)

All of the above metrics take these three values:

- 0 = green = healthy.
- 1 = yellow = warning.
- 2 = red = trouble.

Health Inputs/Outputs JSON Examples

The health.inputs metrics are reported per Source, and the health.outputs metrics per Destination. The health.inputs example below has two configured Sources, and two LogStream-internal inputs. The health.outputs example includes the built-in devnull Destination, and six user-configured Destinations.

Given all the 0 values here, everything is in good shape!

```
"health.inputs": [
    { "model": { "ci": "http:http", "input": "http:http" }}, "val": 0},
    { "model": { "ci": "cribl:CriblLogs", "input": "cribl:CriblLogs" }}, "val": 0},
    { "model": { "ci": "cribl:CriblMetrics", "input": "cribl:CriblMetrics" }}, "val": 0},
    { "model": { "ci": "datagen:DatagenWeblog", "input": "datagen:DatagenWeblog" }}, "val": 0
    ],
    "health.outputs": [
    { "model": { "output": "devnull:devnull" }}, "val": 0},
    { "model": { "output": "router:MyOut1" }}, "val": 0},
    { "model": { "output": "router:MyOut2" }}, "val": 0},
    { "model": { "output": "tcpjson:MyTcpOut2" }}, "val": 0},
    { "model": { "output": "router:MyOut3" }}, "val": 0},
    { "model": { "output": "router:MyOut3" }}, "val": 0},
    { "model": { "output": "router:MyOut4" }}, "val": 0},
    { "model": { "output": "router:MyOut4" }}, "val": 0},
}
```

Persistent Queue Metrics

Five metrics below are valuable for monitoring Persistent Queues' behavior:

• pq.queue_size

- pq.in_bytes
- pq.in_events
- pq.out_events
- pq.out_bytes

These are aggregate metrics. But you can distinguish unique metrics per queue Destination, using the #output=<id> dimension. For example, pq.out_events|#output=kafka would be one unique metric; pq.out_events|#output=camus would be another.

Other Metrics Endpoints and Dimensions

Below are basics on using the /system/metrics endpoint, the /system/info endpoint, and the cribl_wp dimension.

/system/metrics Endpoint

/system/metrics is LogStream's primary public metrics endpoint, which returns most internal metrics. Note that many of these retrieved metrics report configuration only, not runtime behavior. For details, see our API Docs.

/system/info Endpoint

/system/info generates the JSON displayed in the LogStream UI at **Settings > Diagnostics > System Info.** Its two most useful properties are loadayg and memory.

loadavg Example

```
"loadavg": [1.39599609375, 1.22265625, 1.31494140625],
```

This property is an array containing the 1-, 5-, and 15-minute load averages at the UNIX OS level. (On Windows, the return value is always [0, 0, 0].) For details, see the Node.js os.loadavg() documentation.

memory Example

```
"memory": { "free": 936206336, "total": 16672968704 },
```

Divide total / free to monitor memory pressure. If the result exceeds 90%, this indicates a risky situation: you're running out of memory.

cpus Alternative

The cpus metric returns an array of CPU/memory key-value pairs. This provides an alternative way of determining loadavg, but it requires you to query all your CPUs individually, and then average. In the example below, Idle = user + nice + sys.

```
"cpus": [{ "times": {
    "user": 19881260,
    "nice": 39370,
    "sys": 5250130,
    "idle": 76088790,
```

cribl_wp Metric Dimension

cribl_wp is a useful dimension that identifies the Worker Process that processed each event.

Upgrading

This page outlines how to upgrade Cribl LogStream's Single-Instance or Distributed Deployment packages along one of the following supported upgrade paths:

- v2.x ==> v2.x
- v1.7.x/v2.0.x ==> v2.x.x
- v1.6.x or below ==> v1.7.x ==> v2.x.x
 - △ See notes on Upgrading to LogStream 2.3 below.

LogStream does **not** support direct upgrades from a Beta to a GA version. To get the GA version running, you must perform a new install.

Standalone/Single-Instance

This path requires upgrading only the single/standalone node:

- 1. Stop Cribl LogStream.
- 2. Uncompress the new version on top of the old one.

On some Linux systems, tar might complain with: cribl/bin/cribl: Cannot open: File exists. In this case, please remove the cribl/bin/cribl directory if it's empty, and untar again. If you have custom functions in cribl/bin/cribl, please move them under \$CRIBL_HOME/local/cribl/functions/ before untarring again.

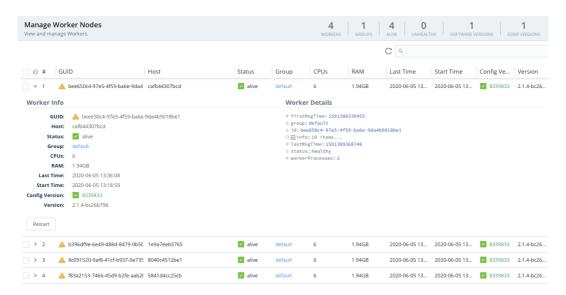
3. Restart LogStream.

Distributed Deployment

For a distributed deployment, the order of upgrade is: Upgrade first the Master Node, then upgrade the Worker Nodes, then commit and deploy the changes on the Master.

Upgrade the Master Node

- 1. Commit and deploy your desired last version. (This will be your most recent checkpoint.)
 - Optionally, git push to your configured remote repo.
- 2. Stop Cribl LogStream.
 - Optional but recommended: Back up the entire \$CRIBL_HOME directory.
 - Optional: Check that the Worker Nodes are still functioning as expected. In absence of the Master Node, they should continue to work with their last deployed configurations.
- 3. Uncompress the new LogStream version on top of the old one.
- 4. Restart LogStream and log back in.
- 5. Wait for all the Worker Nodes to report to the Master, and ensure that they are correctly reporting the last committed configuration version.
 - i Workers' UI will not be available until the Worker version has been upgraded to match the version on the Master. Errors like those below will appear until the Worker nodes are upgraded.



Worker Node version mismatch

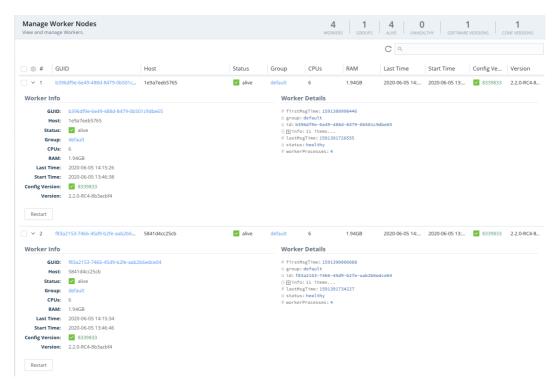
Upgrade the Worker Nodes

These are the same basic steps as when upgrading a Standalone Instance, above:

- 1. Stop Cribl LogStream on each Worker Node.
- 2. Uncompress the new version on top of the old one.
- 3. Restart LogStream.

Commit and Deploy Changes on the Master Node

- 1. Ensure that newly upgraded Worker Nodes report to the Master with their new software version.
- 2. Commit and deploy the newly updated configuration **only after all** Workers have upgraded.



Post-2.1.4 upgrade to 2.2

Upgrading to LogStream 2.3

As of version 2.3, LogStream Free and One licenses are permanent, but they enforce certain restrictions that especially affect distributed deployments:

- Even if you have more than one Worker Group defined, only one Worker Group will be visible and usable.
 - This will be the first Group listed in \$CRIBL_HOME/local/cribl/groups.yml - typically, the default Group. You can edit groups.yml to move the desired Group to the top.
- Your cluster will be limited to 10 Worker Processes across all Worker Nodes.
 - LogStream will balance (or rebalance) these Processes as evenly as possible across the Worker Nodes.
- Authentication will fall back to local authorization. You will not be able to authenticate via Splunk, LDAP, or SSO/OpenID.
- **Git Push** to remote repos will not be supported through the product.
 - ⚠ If you are upgrading LogStream Free or LogStream One from version 2.2.x or lower, these changes might require you to adjust your existing configuration and/or workflows.

See Licensing for details on all current license options.

As of LogStream 2.3, licenses no longer need to be deployed directly to Worker Groups. The Master will push license information down to Worker Groups as part of the heartbeat.

Splunk App Package Upgrade Steps

△ See Deprecation note for v.2.1.

Follow these steps to upgrade from v.1.7, or higher, of the Cribl App for Splunk:

- 1. Stop Splunk.
- 2. Untar/unzip the new app version on top of the old one.

On some Linux systems, tar might complain with: cribl/bin/cribl: Cannot open: File exists. In this case, please remove the cribl/bin/cribl directory if it's empty, and untar again. If you have
custom functions in cribl/bin/cribl , please move them under
\$CRIBL_HOME/local/cribl/functions/ before untarring again.

3. Restart Splunk.

Upgrading from Splunk App v.1.6 (or Lower)

As of v.1.7, contrary to prior versions, Cribl's Splunk App package defaults to Search Head Mode. If you have v.1.6 or earlier deployed as a Heavy Forwarder app, upgrading requires an extra step to restore this setting:

- 1. Stop Splunk.
- 2. Untar/unzip the new app version on top of the old one.
- 3. Convert to HF mode by running:
 \$\$PLUNK_HOME/etc/apps/cribl/bin/cribld mode-hwf
- 4. Restart Splunk.

Uninstalling

Uninstalling the Standalone Version

- Stop Cribl LogStream (stopping the main process).
- Back up necessary configurations/data.
- Remove the directory where Cribl LogStream is installed.

In a distributed deployment, repeat the above steps for the Master instance and all Worker instances.

Uninstalling the Splunk App Version

- Stop Splunk.
- Back up necessary configurations/data.
- Remove the Cribl App in \$SPLUNK_HOME/etc/apps.
- Remove the Cribl module in \$SPLUNK_HOME/etc/modules/cribl (some versions).

Working With Data

Event Model

All data processing in Cribl LogStream is based on discrete data entities commonly known as **events**. An event is defined as a collection of key-value pairs (fields). Some Sources deliver events directly, while others might deliver bytestreams that need to be broken up by Event Breakers. Events travel from a Source through Pipelines' Functions, and on to Destinations.

The internal representation of a Cribl LogStream event is as follows:

```
Cribl LogStream Event Model
```

```
{
   "_raw": "<body of non-JSON parse-able event>",
   "_time": "<timestamp in UNIX epoch format>",
   "_inputId": "<Id/Name of Source that delivered the event>",
   "_other1": "<Internal field1>",
   "_other2": "<Internal field2>",
   "_otherN": "<Internal fieldN>",
   "key1": "<value1>",
   "key2": "<value2>",
   "keyN": "<valueN>",
   "...": "..."
}
```

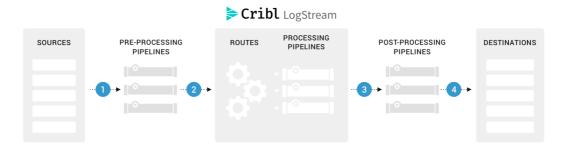
Some notes about these representative fields:

- Fields that start with a double-underscore are known as internal fields, and each Source can add one or many to each event. For example, Syslog adds both a __inputId and a __srcIpPort field. Internal fields are used only within Cribl LogStream, and are not passed down to Destinations.
- Upon arrival from a Source, if an event cannot be JSON-parsed, all of its content will be assigned to _raw .
- If a timestamp is not configured to be extracted, the current time (in UNIX epoch format) will be assigned to _time .

Using Capture

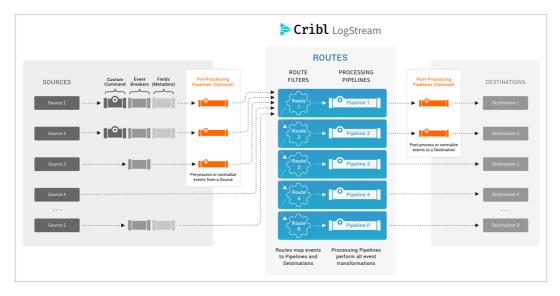
One way to see what an event looks like as it travels through the system is to use the **Capture** feature. While in **Preview** (right pane):

- 1. Click **Start a Capture**.
- 2. In the resulting modal, enter a **Filter expression** to narrow down the events of interest.
- 3. Click **Capture...** and (optionally) change the default Time and/or Event limits.
- 4. Select the desired Where to capture option. There are four options:
- 1. Before the pre-processing Pipeline Capture events right after they're delivered by the respective Input.
- 2. Before the Routes Capture events right after the pre-processing Pipeline, before they go down the Routes.
- 3. Before the post-processing Pipeline Capture events right after the Processing Pipeline that actually handled them, before any post-processing Pipeline.
- **4. Before the Destination** Capture events right after the post-processing Pipeline, before they go out to the configured Destination.



Event Processing Order

The expanded schematic below shows how all events in the Cribl LogStream ecosystem are processed linearly, from left to right.



LogStream in great detail

Here are the stages of event processing:

- Sources: Data arrives from your choice of external providers. (LogStream supports Splunk, HTTP/S, Elastic Beats, Amazon Kinesis/S3/SQS, Kafka, TCP raw or JSON, and many others.)
- 2. Custom command: Optionally, you can pass this input's data to an external command before the data continues downstream. This external command will consume the data via stdin, will process it and send its output via stdout.
- 3. Event Breakers can, optionally, break up incoming bytestreams into discrete events.
- 4. Fields/Metadata: Optionally, you can add these enrichments to each incoming event. You add fields by specifying key/value pairs, per Source, in a format similar to LogStream's Eval function. Each key defines a field name, and each value is a JavaScript expression (or constant) used to compute the field's value.

- 5. Pre-processing Pipeline: Optionally, you can use a single Pipeline to condition (normalize) data from this input before the data reaches the Routes.
- 6. Routes map incoming events to Processing Pipelines and Destinations. A Route can accept data from multiple Sources, but each Route can be associated with only one Pipeline and one Destination.
- 7. Processing Pipelines perform all event transformations. Within a Pipeline, you define these transformations as a linear series of Functions.

 A Function is an atomic piece of JavaScript code invoked on each event.
- 8. Post-processing Pipeline: Optionally, you can append a Pipeline a to condition (normalize) data from each Processing Pipeline before the data reaches its Destination.
- 9. Destinations: Each Route/Pipeline combination forwards processed data to your choice of streaming or storage Destination. (LogStream supports Splunk, Syslog, Elastic, Kafka/Confluent, Amazon S3, Filesystem/NFS, and many other options.)

i Pipelines Everywhere

All Pipelines have the same basic internal structure – they're a series of Functions. The three Pipeline types identified above differ only in their position in the system.

Routes

What Are Routes

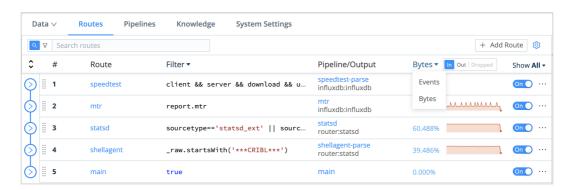
Before incoming events are transformed by a processing Pipeline, Cribl LogStream uses a set of filters to first select a **subset** of events to deliver to the correct Pipeline. This selection is made via Routes.

How Do Routes Work

Routes apply filter expressions on incoming events to send matching results to the appropriate Pipeline. Filters are JavaScript-syntax-compatible expressions that are configured with each Route. Examples are: true,

i There can be multiple Routes in the system, but each Route can be associated with only **one** Pipeline.

Routes are evaluated in their display order, top->down. The stats shown in the **Bytes/Events** (toggle) column are for the most-recent 15 minutes.

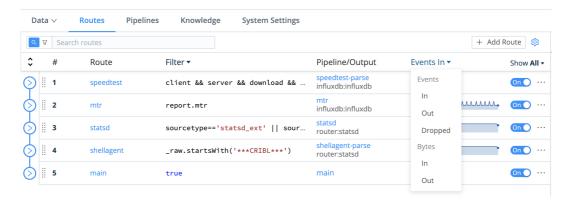


Routes and bytes

In the example above, incoming events will be evaluated first against the Route named **speedtest**, then against **mtr**, then against **statsd**, and so on. At the end, the **main** Route serves as a catch-all for any event that does not match any of the other Routes.

Above, note the selectors to toggle between displaying Events versus Bytes, and to display In versus Out.

When you condense the Routes page to a narrower viewport, LogStream consolidates the In/Out/Dropped selectors onto an expanded Bytes/Events drop-down menu, as shown below.



Routes and events (combined menu)

Managing the Routes Page

To apply a Route before another, simply drag it vertically. Use the sliders to turn Routes **On/Off** inline, as necessary, to facilitate development and debugging.

You can press the] (right-bracket) shortcut key to toggle between the Preview pane and the expanded Routes display shown above. (This works when no field has focus.)

Output Destination

You can configure each Route with an output Destination that denotes where to send events after they're processed by the Pipeline.

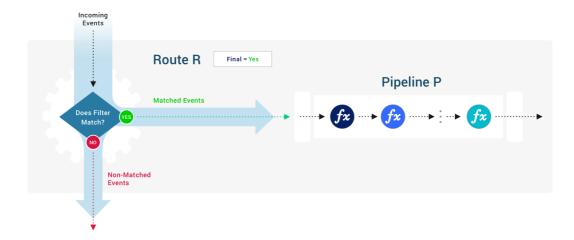
The Final Toggle

When an event that enters the system and matches a Route-Pipeline pair, it will usually be either:

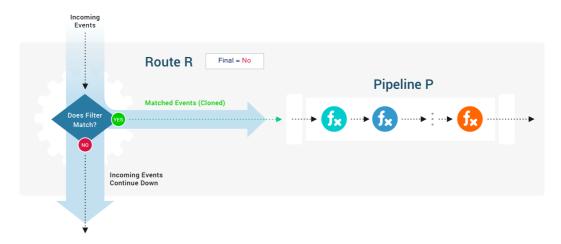
- Dropped by a function, or
- Transformed (optionally) and exit the system.

This behavior is ensured by the Final toggle in Route settings. It defaults to Yes, meaning that matched events will be **consumed** by that Route, and will

not be evaluated against any other Routes that sit below it.



If the Final toggle is set to No, clone(s) of the matching events will be processed by the configured Pipeline, and the original events will be allowed to continue their trip to be evaluated and/or processed by other Route-Pipeline pairs.



This is very useful in cases where the same set of events needs to be processed differently and delivered to different Destinations. Each clone can be decorated with key-value pairs as necessary.

Final Flag and Cloning Considerations

Depending on your cloning needs, you might want to follow a **most-specific first** or a **most-general first** processing strategy. The general goal is to minimize the number of filters/Routes an event gets evaluated against. For example:

If cloning is not needed at all (i.e., all Final toggles stay at default), then it
makes sense to start with the broadest expression at the top, so as to
consume as many events as early as possible.

• If cloning is needed on a narrow set of events, then it might make sense to do that upfront, and follow it with a Route that consumes those clones immediately after.

Route Groups

A Route group is a collection of consecutive Routes that can be moved up and down the Route stack together. Groups help with managing long lists of Routes. They are a UI visualization only: While Routes are in a group, those Routes maintain their global position order.

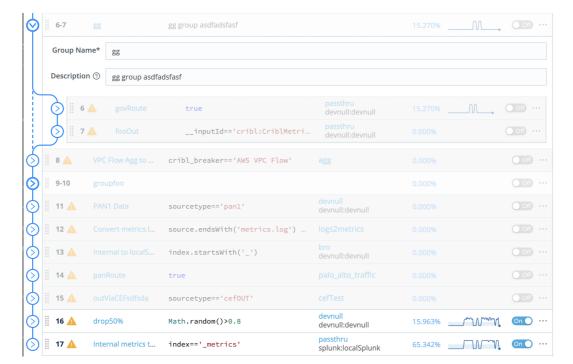
i Route groups work much like Function groups, offering similar UI controls and drag-and-drop options.

Unreachable Routes

Routes display an "unreachable" warning indicator (orange triangle) when data can't reach them. This will be true when, with your current configuration, any Route higher in the stack matches **all** three of these conditions:

- Previous Route is enabled (slider is set to On).
- Previous Route is final (Final slider is set to Yes).
- Previous Route's Filter expression evaluates to true, (e.g., true, 1 == 1, etc.).

Note that the third condition above can be triggered intermittently by a randomizing method like Math.random(). This might be included in a previous Route's own Filter expression, or in a Pipeline Function (such as one configured for random data sampling).



Unreachable Route warnings

Routing with Output Router

Output Router Destinations offer another way to route data. These function as meta-Destinations, in that they allow selection of actual Destinations based on rules. Rules are evaluated in order, top->down, with the first match being the winner.

Pipelines

What Are Pipelines

After your data has been matched by a Route, it gets delivered to a Pipeline. A Pipeline is a list of Functions that work on the data.

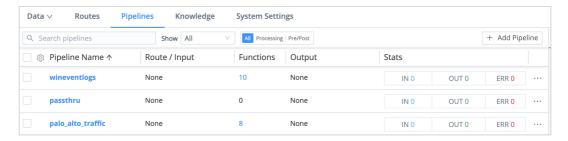
i As with Routes, the order in which the Functions are listed matters. A Pipeline's Functions are evaluated in order, top->down.

Accessing Pipelines

Select **Pipelines** from LogStream's (or a Worker Group's) top menu. To configure a new Pipeline, click **+ Add Pipeline**.

How Do Pipelines Work

Events are always delivered to the beginning of a Pipeline via a Route. The data in the **Stats** column shown below are for the last 15 minutes.



Pipelines and Route inputs

i You can press the] (right-bracket) shortcut key to toggle between the Preview pane and an expanded Pipelines display. This works when no field has focus.

Within the Pipeline, events are processed by each Function, in order. A Pipeline will always move events in the direction that points outside of the system. This

is on purpose, to keep the design simple and avoid potential loops.



Pipeline Functions

Pipeline Settings

Click the gear icon at top right to open the Pipeline's Settings. Here, you can attach the Pipeline to a Route. In the Settings' **Async function timeout (ms)** field, you can enter a buffer to adjust for Functions that might take much longer to execute than normal. (An example would be a Lookup Function processing a large lookup file.)



Pipeline Settings

Advanced Mode (JSON Editor)

Click **Advanced Mode** to edit the Pipeline's definition as JSON text. In this mode's editor, you can directly edit multiple values. You can also use the **Import** and **Export** buttons here to copy and modify existing Pipeline configurations.

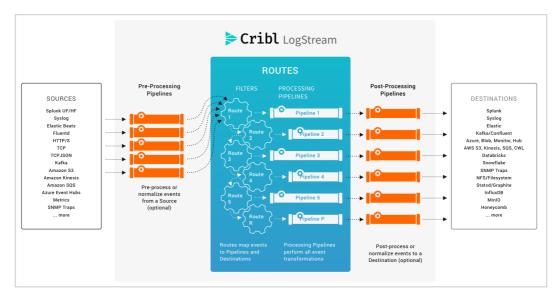
```
Pipelines > elastic > Advanced Settings Mode
                                                                                                      ← Back
Attached to Route: elastic
     Simple Mode
                                                                                       ⊥ Import
                                                                                                  ± Export
             "id": "elastic",
             "conf": {
       3
               "output": "default",
       4
       5
               "groups": {},
               "asyncFuncTimeout": 1000,
       6
               "functions": [
                   "id": "drop",
                   "filter": "host!='192.168.1.241'",
      10
                   "disabled": true,
      11
                   "conf": {}
      12
      13
      14
             }
      15
      16
```

Advanced Pipeline Editing

i You can streamline the above display by organizing related Functions into Function groups.

Types of Pipelines

You can apply various Pipeline types at different stages of data flow. All Pipelines have the same basic internal structure (a series of Functions) – the types below differ only in their position in the system.



Pre-processing, processing, and post-processing Pipelines

Pre-Processing Pipelines

These are Pipelines that are attached to a Source to condition (normalize) the events **before** they're delivered to a processing Pipeline. They're optional.

Typical use cases are event formatting, or applying Functions to **all** events of an input. (E.g., to extract a message field before pushing events to various processing Pipelines.)

You configure these Pipelines just like any other Pipeline, by selecting **Pipelines** from the top menu. You then attach your configured Pipeline to individual **Sources**, using the Source's **Pre-Processing > Pipeline** drop-down.

Fields extracted using pre-processing Pipelines are made available to Routes.

Processing Pipelines

These are "normal" event processing Pipelines, attached directly to Routes.

Post-Processing Pipelines

These Pipelines are attached to a Destination to normalize the events before they're sent out. A post-processing Pipeline's Functions apply to **all** events exiting to the attached Destination.

Typical use cases are applying Functions that transform or shape events per receiver requirements. (E.g., to ensure that a _time field exists for all events bound to a Splunk receiver.)

You configure these Pipelines as normal, by selecting **Pipelines** from the top menu. You then attach your configured Pipeline to individual **Destinations**, using the Destination's **Post-Processing > Pipeline** drop-down.

You can also use a Destination's **Post-Processing** options to add **System Fields** like <code>cribl_input</code>, identifying the LogStream Source that processed the events.

Best Practices for Pipelines

Functions in a Pipeline are equipped with their own filters. Even though filters are not required, we recommend using them as often as possible.

As with Routes, the general goal is to minimize extra work that a Function will do. The fewer events a Function has to operate on, the better the overall performance.

For example, if a Pipeline has two Functions, **f1** and **f2**, and if **f1** operates on source 'foo' and **f2** operates on source 'bar', it might make sense to apply source='foo' versus source='bar' filters on these two Functions, respectively.

Data Onboarding

Onboarding data into Cribl LogStream can vary in complexity, depending on your organization's needs, requirements, and constraints. Proper onboarding from all Sources is key to system performance, troubleshooting, and ultimately the quality of data and decisions both in LogStream and in downstream Destinations.

General Onboarding Steps

Typically, a data onboarding process revolves around these steps, both before and after turning on the Source:

- Create configuration settings.
- Verify that settings do the right thing.
- Iterate.

Below, we break down individual steps.

Before Turning On the Source

Cribl recommends that you take the following steps to verify and tune incoming data, before it starts flowing.

Preview Sample Data

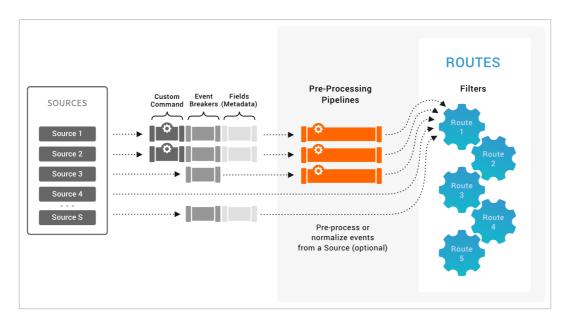
Use a sample of your real data in Data Preview. Sample data can come from a sample Source file that you upload or paste into LogStream.

You can also obtain sample data in a live data capture from a Source. One way to do this **before** going to production is to configure your Source with a **devnull** Pipeline (which just drops all events) as a pre-processing Pipeline. Then, let data flow in for just long enough to capture a sufficient sample.

Check the Processing Order

While events can be processed almost arbitrarily by functions in LogStream Pipelines, make sure you understand the event processing order. This is very

important, as it tells you exactly where certain processing steps occur. For instance, as we'll see just below, quite a few steps can be accomplished at the Source level, before data even hits LogStream Routes.



Source-level processing options

Custom Command

Where supported, data streams will be handled by **custom commands**. These are external system commands that can (optionally) be used to pre-process the data. You can specify any command, script, etc., that consumes via stdin and outputs via stdout.

Verify that such commands are doing what's expected, as they are the very **first** in a series of processing steps.

Event Breakers

Next, data streams are handled by Event Breakers, which:

- Convert data streams into discrete events.
- Extract and assign timestamps to each event.

If the resulting events do not look correct, feel free to use **non-default** breaking rules and timestamp recognition patterns. Downstream, you can use the Auto Timestamp function to modify _time as needed, if timestamps were not recognized properly. Examples of such errors are:

Timestamps too far out in the future or past

- Wrong timezone.
- Incorrect timestamp is selected from multiple timestamps present in the event.

Fields (Metadata)

Next, events can be enriched with Fields (Metadata). This is where you'd add static or dynamic fields to all events delivered by a particular Source.

Pre-Processing Pipeline

Next, you can optionally configure a pre-processing Pipeline on a particular Source. This is extremely useful in these cases:

- Drop non-useful events as early as possible (so as to save on CPU processing).
- Normalize events from this Source to conform a certain shape or structure.
- Fix/touch up events accordingly. E.g., if event breakers assigned the wrong timestamp, this is the best place to use the Auto Timestamp function to adjust _time.

We Can't Say This Enough

Verify, verify, verify your data's integrity before turning on the Source.

After Turning On the Source

Use data Destinations to verify that certain metrics of interest are accurate. This will depend significantly on the capabilities of each Destination, but here's a basic checklist list of things to ensure:

- Timestamps are correct.
- All necessary fields are assigned to events.
- All expected events show up correctly. (E.g., if a Drop or Suppress Function was configured, ensure that it's not dropping unintended events.)
- Throughput both in bytes and in events per second (EPS) is what's expected, or is within a certain tolerance.

Iterate

Iterate on the steps above as necessary. E.g., adjust fields values and timestamps as needed.	
	member that there is almost always a workaround. Any arbitrary ent transformation that you need is likely just a Function or two ay.

Functions

What Are Functions

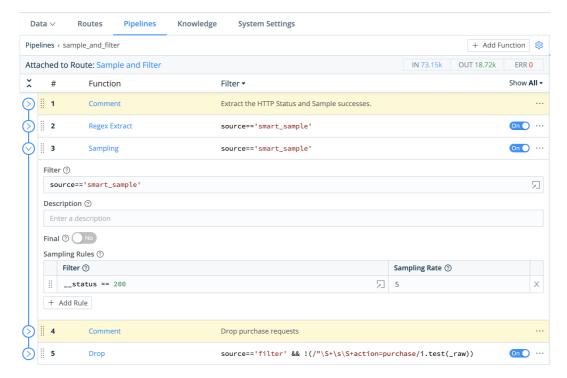
When events enter a Pipeline, they're processed by a series of Functions. At its core, a Function is code that executes on an event, and it encapsulates the smallest amount of processing that can happen to that event.

The term "processing" means a variety of possible options: string replacement, obfuscation, encryption, event-to-metrics conversions, etc. For example, a Pipeline can be composed of several Functions – one that replaces the term foo with bar, another one that hashes bar, and a final one that adds a field (say, dc=jfk-42) to any event that matches source='us-nyc-application.log'.

How Do They Work

Functions are atomic pieces of JavaScript code that are invoked on each event that passes through them. To help improve performance, Functions can be configured with filters to further scope their invocation to matching events only.

You can add as many Functions in a Pipeline as necessary, though the more you have, the longer it will take each event to pass through. Also, you can turn Functions **On/Off** within a Pipeline as necessary. This enables you to preserve structure as you optimize or debug.



Functions stack in a Pipeline

You can reposition Functions up or down the Pipeline stack to adjust their execution order. Use a Function's left grab handle to drag and drop it into place.

The Final Toggle

Similar to the Final toggle in Routes, the Final toggle here controls the flow of events at the Function level. Its states are:

- No (default): means that matching events processed by this Function will be passed down to the next Function.
- Yes: means that this Function is the last one that will be applied to
 matching events. All Functions further down the Pipeline will be skipped. A
 Function with Final set to Yes will display an Findicator in the Pipeline
 stack.

Out-of-the-Box Functions

Cribl LogStream ships with several Functions out-of-the-box, and you can chain them together to meet your requirements. For more details, see individual **Functions**, and the **Use Cases** section, within this documentation.

Custom Functions

For an overview of adding custom Functions to Cribl LogStream, see our blog post, Extending Cribl: Building Custom Functions.

What Functions to Use When

- Add, remove, update fields:
 Eval, Lookup, Regex Extract
- Find & Replace, including basic sed -like, obfuscate, redact, hash, etc.:
 Mask, Eval
- Add GeoIP information to events:
 GeoIP
- Extract fields:Regex Extract, Parser
- Extract timestamps: Auto Timestamp
- Drop events:
 Drop, Regex Filter, Sampling, Suppress, Dynamic Sampling
- Sample events (e.g, high-volume, low-value data):
 Sampling, Dynamic Sampling
- Suppress events (e.g, duplicates, etc.):
 Suppress
- Serialize events to CEF format (send to various SIEMs):
 CEF Serializer
- Serialize / change format (e.g., convert JSON to CSV):
 Serialize
- Convert JSON arrays into their own events: JSON Unroll, XML Unroll
- Flatten nested structures (e.g., nested JSON):
 Flatten
- Aggregate events in real-time (i.e. statistical aggregations):
 Aggregations

- Convert events to metrics format:
 Publish Metrics, Prometheus Publisher (beta)
- Resolve hostname from IP address: Reverse DNS (beta)
- Extract numeric values from event fields, converting them to type number:
 Numerify
- Send events out to a command or a local file, via stdin, from any point in a Pipeline:

Tee

- Convert an XML event's elements into individual events:
 XML Unroll
- Duplicate events in the same Pipeline, with optional added fields:
 Clone
- Add a text comment within a Pipeline's UI, to label steps without changing event data:

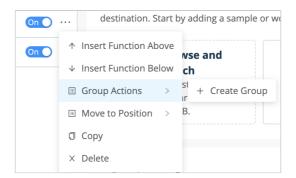
Comment

Function Groups

A Function group is a collection of consecutive Functions that can be moved up and down a Pipeline's Functions stack together. Groups help you manage long stacks of Functions by streamlining their display. They are a UI visualization only: While Functions are in a group, those Functions maintain their global position order in the Pipeline.

i Function groups work much like Route groups.

To build a group from any Function, click the Function's ••• (Options) menu, then select **Group Actions** > **Create Group**.



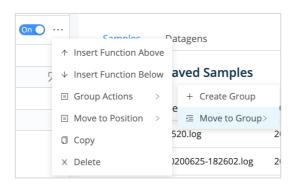
Creating a group

You'll need to enter a **Group Name** before you can save or resave the Pipeline. Optionally, enter a **Description**.



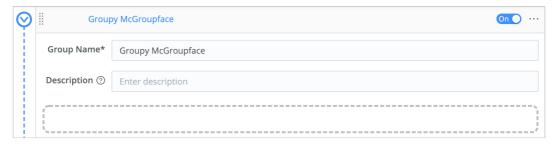
Naming a group

Once you've saved at least one group to a Pipeline, other Functions'
••• (Options) > Group Actions submenus will add options to Move to Group or Ungroup/Ungroup All.



Expanded Group Actions submenu

You can also use a Function's left grab handle to drag and drop it into, or out of, a group. A saved group that's empty displays a dashed target into which you can drag and drop Functions.



Drag-and-drop target

Auto Timestamp

Description

The Auto Timestamp Function extracts time to a destination field, given a source field in the event. By default, Auto Timestamp makes a first best effort and populates _time . When you add a sample (via paste or a local file), you should accomplish time and event breaking at the same time you add the data.

This Function allows fine-grained and powerful transformations to populate new time fields, or to edit existing time fields. You can use the Function's Additional timestamps section to create custom time fields using regex and custom JavaScript strptime functions.

The Auto Timestamp Function uses the same basic algorithm as the Event Breaker Function and the C.Time.timestampFinder() native method.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. The default true setting passes all events through the Function.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Source field: Field to search for a timestamp. Defaults to _raw .

Destination field: Field to place extracted timestamp in. Defaults to _time . Supports nested addressing.

Default timezone: Select a timezone to assign to timestamps that lack timezone info. Defaults to Local . (This drop-down includes support for legacy names: EST5EDT, CST6CDT, MST7MDT, and PST8PDT.)

Additional timestamps: Add Regex/Strptime pairs to extract additional timestamp formats.

- Regex: Regex, with first capturing group matching the timestamp.
- Strptime format: Select or enter the strptime format for the captured timestamp.

Click **Add timestamp** to add more rows.

Advanced Settings

Time expression: Expression with which to format extracted time. Current time, as a JavaScript Date object, is in global time. Defaults to time.getTime() / 1000.

For details about Cribl LogStream's Library (native) time methods, see: C.Time – Time Functions.

Start scan offset: How far into the string to look for a time string.

Max timestamp scan depth: Maximum string length at which to look for a timestamp.

Default time: How to set the time field if no timestamp is found. Defaults to **Current time.**

Two fields enable you to constrain (clamp) the parsed timestamp, to prevent the Function from mistakenly extracting non-time values as unrealistic timestamps:

- Earliest timestamp allowed: Enter a string that specifies the latest allowable timestamp, relative to now. (Sample value: -42years . Default value: -420weeks .) Parsed values earlier than this date will be set to the **Default time**.
- Future timestamp allowed: Enter a string that specifies the latest allowable timestamp, relative to now. (Sample value: +42days. Default value: +1week.) Parsed values after this date will be set to the Default time.

Format Reference

Referencing https://github.com/d3/d3-time-format#locale_format:

```
%a - abbreviated weekday name.*
%A - full weekday name.*
%b - abbreviated month name.*
%B - full month name.*
%c - the locale's date and time, such as %x, %X.*
%d - zero-padded day of the month as a decimal number [01,31].
%e - space-padded day of the month as a decimal number [ 1,31]; equivalent to %_d.
%f - microseconds as a decimal number [000000, 9999999].
%H - hour (24-hour clock) as a decimal number [00,23].
%I - hour (12-hour clock) as a decimal number [01,12].
%j - day of the year as a decimal number [001,366].
%m - month as a decimal number [01,12].
%M - minute as a decimal number [00,59].
%L - milliseconds as a decimal number [000, 999].
%p - either AM or PM.*
%Q - milliseconds since UNIX epoch.
%s - seconds since UNIX epoch.
%S - second as a decimal number [00,61].
%u - Monday-based (ISO 8601) weekday as a decimal number [1,7].
%U - Sunday-based week of the year as a decimal number [00,53].
%V - ISO 8601 week of the year as a decimal number [01, 53].
%w - Sunday-based weekday as a decimal number [0,6].
%W - Monday-based week of the year as a decimal number [00,53].
%x - the locale's date, such as %-m/%-d/%Y.*
%X - the locale's time, such as %-I:%M:%S %p.*
%y - year without century as a decimal number [00,99].
```

```
%Y - year with century as a decimal number.
%Z - time zone offset, such as -0700, -07:00, -07, or Z.
%% - a literal percent sign (%).
```

Directives marked with an asterisk (*) might be affected by the locale definition.

Complying with the Format

In order to use auto timestamping upon ingestion, the formatting used must match the %Z parameters above. E.g., this Function will automatically parse all of these formats:

```
• 2020/06/10T17:17:35.004-0700
```

- 2020/06/10T17:17:35.004-07:00
- 2020/06/10T17:17:35.004-07
- 2020/06/10T10:17:35.004Z
- 2020/06/10T11:17:35.004 EST

To parse other formats, you can use the Additional Timestamps section's internal Regex or Strptime Format operators.

Basic Example

Filter: name.startsWith('kumquats') & value='specific string here'

This will allow the Auto Timestamp Function to act only on events matching the specified parameters.

Sample:

```
Sep 20 12:03:55 PA-VM 1,2019/09/20 13:03:58,CRIBL,TRAFFIC,end,2049,2019/09/20 14:03:58,314.817.1
```

To add this sample (after creating an Auto Timestamp Function with the above **Filter** expression): Go to **Preview** > **Add a Sample** > **Paste a Sample**, and add the data snippet above. Do not make any changes to timestamping or line breaking, and select **Save as Sample File**.

By default, LogSteram will inspect the first 150 characters, and extract the first valid timestamp it sees. You can modify this character limit under **Advanced Settings** > **Max Timestamp Scan Depth**.

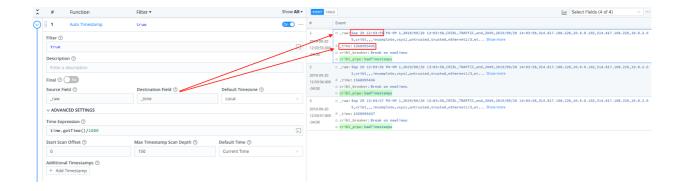
LogStream grabs the first part of the event, and settles on the first matching value to display for time:

```
_time 1569006235
```

GMT: Friday, 20 September 2019, 7:03:55 PM GMT

Your Local Time: Friday, 20 September 2019 PDT, 12:03:55 AM GMT -07:00

Because no explicit timezone has been set (under **Default Timezone**), _time inherits the **Local** timezone, which in this example is GMT -07:00.



i Timezone Dependencies and Details

LogStream uses ICU for timezone information. It does not query external files or the operating system. The bundled ICU is updated periodically.

For additional timezone details, see: https://www.iana.org/time-zones.

Advanced Settings Example

The datetime.strptime() method creates a datetime object from the string passed in by the Regex field.

Here, we'll use datetime.strptime() to match a timestamp in AM/PM format at the end of a line.

Sample:

This is a sample event that will push the datetime values further on inside the event. This is still a sample event and finally here is the datetime information!: Server_UTC_Timestamp="04/27/2020 2:30:15 PM"

Max timestamp scan depth: 200

Click to add Additional timestamps:

Regex: $(d{1,2})/(d{2})/(d{4})$ s $(d{1,2}):(d{2})$ s $(w{2})$

Strptime format: '%m/%d/%Y %H:%M:%S %p'

i Gnarly Details

This Function supports the %f (microseconds) directive, but LogStream will truncate it to millisecond resolution.

For further examples, see Extracting Timestamps from Messy Logs.

Aggregations

Description

The Aggregations Function performs aggregate statistics on event data.

Safeguarding Data

Upon shutdown, LogStream will attempt to flush the buffers that hold aggregated data, to avoid data loss. If you set a **Time window** greater than an hour, Cribl recommends adjusting the **Aggregation memory limit** and/or **Aggregation event limit** to prevent the system from running out of memory. This is especially necessary for high-cardinality data. (Both settings default to unlimited, but we recommend setting defined limits, based on testing.)

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Time window: The time span of the tumbling window for aggregating events. Must be a valid time string (e.g., 10s). Must match pattern \d+[sm]\$.

Aggregates: Aggregate function(s) to perform on events. E.g., sum(bytes).where(action='REJECT').as(TotalBytes). Expression format: aggFunction(<FieldExpression>).where(<FilterExpression>).as(<outputField>). See more examples below.

Note: When used without as(), the aggregate's output will be placed in a field labeled <aggFunction>_<fieldName>. If there are conflicts, the last aggregate wins.
 For example, given two aggregates - sum(bytes).where(action='REJECT') and sum(bytes) - the latter one (sum_bytes) is the winner.

Group by Fields: Fields to group aggregates by.

Evaluate fields: Set of key/value pairs to evaluate and add/set. Fields are added in the context of an aggregated event, before they're sent out. Does not apply to passthrough events.

Time Window Settings

Cumulative aggregations: Determines if the aggregations should be retained for cumulative aggregations, or reset to 0, when flushing out an aggregation table event. Defaults to No.

Lag tolerance: The lag tolerance represents the tumbling window tolerance to late events. Must be a valid time string (e.g., 10s). Must match pattern \d+[sm]\$.

Idle bucket time limit: The amount of time to wait before flushing a bucket that has not received events. Must be a valid time string (e.g., 10s). Must match pattern \d+[sm]\$.

Output Settings

Passthrough mode: Determines whether to pass through the original events along with the aggregation events. Defaults to No.

Metrics mode: Determines whether to output aggregates as metrics. Defaults to No, causing aggregates to be output as events.

Sufficient stats mode: Determines whether to output *only* statistics sufficient for the supplied aggregations. Defaults to No, meaning output richer statistics.

Output prefix: A prefix that is prepended to all of the fields output by this Aggregations Function.

Advanced Settings

Aggregation event limit: The maximum number of events to include in any given aggregation event. Defaults to unlimited. Must be at least 1.

Aggregation memory limit: The memory usage limit to impose upon aggregations. Defaults to unlimited (i.e., the amount of memory available in the system).

List of Aggregate Functions

avg(expr:FieldExpression): Returns the average of the values of the parameter.
count(expr:FieldExpression): Returns the number of occurrences of the values of the parameter.

dc(expr: FieldExpression, errorRate: number = 0.01): Returns the estimated
number of distinct values of the <expr> parameter, within a relative error rate.
distinct_count(expr: FieldExpression, errorRate: number = 0.01): Returns the
estimated number of distinct values of the <expr> parameter, within a relative error rate.
earliest(expr:FieldExpression): Returns the earliest(based on _time) observed
value of the parameter.

first(expr:FieldExpression): Returns the first observed value of the parameter.
last(expr:FieldExpression): Returns the last observed value of the parameter.
latest(expr:FieldExpression): Returns the latest (based on _time) observed value of the parameter.

max(expr:FieldExpression): Returns the maximum value of the parameter.
min(expr:FieldExpression): Returns the minimum value of the parameter.
per_second(expr:FieldExpression): Returns the per second rate (based on _time)
observed value of the parameter.

perc(level: number, expr: FieldExpression):Returns <level> percentile value of the numeric values of the <expr> parameter.

rate(expr:FieldExpression, timeString: string = '1s'):Returns the rate (based
on _time) observed value of the parameter.

stddev(expr:FieldExpression): Returns the sample standard deviation of the values of the parameter.

stddevp(expr:FieldExpression): Returns the population standard deviation of the values of the parameter.

sum(expr:FieldExpression) : Returns the sum of the values of the parameter.
sumsq(expr:FieldExpression) : Returns the sum of squares of the values of the
parameter.

variance(expr:FieldExpression): Returns the sample variance of the values of the parameter.

variancep(expr:FieldExpression): Returns the population variance of the values of the parameter.

How Do Time Window Settings Work?

Lag Tolerance

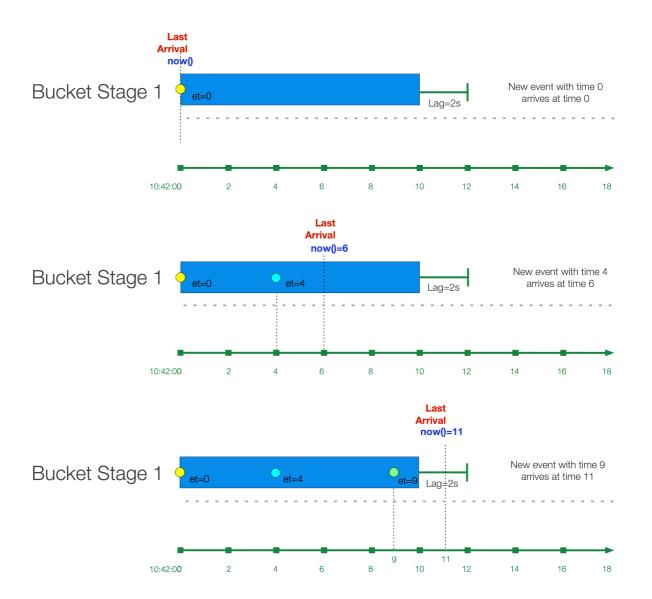
As events are aggregated into windows, there is a good chance that most will arrive later than their event time. For instance, given a 10s window (10:42:00 - 10:42:10), an event with timestamp 10:42:03 might come in 2 seconds later at 10:42:05.

In several cases, there will also be late, or lagging, events that will arrive **after** the latest time window boundary. For example, an event with timestamp 10:42:04 might arrive at

10:42:12 . Lag Tolerance is the setting that governs how long to wait – after the latest window boundary – and still accept late events.

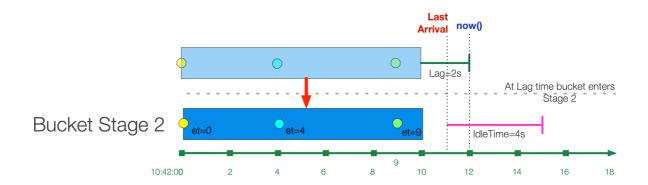
10s Window Aggregation

Settings: Lag=2s, IdleTime=4s



The "bucket" of events is said to be in Stage 1, where it's still accepting new events, but it's not yet finalized. Notice how in the third case, an event with event time 10:42:09 arrives 1 second past the window boundary at 10:42:11, but it's still accepted because it happens before the lag time expires.

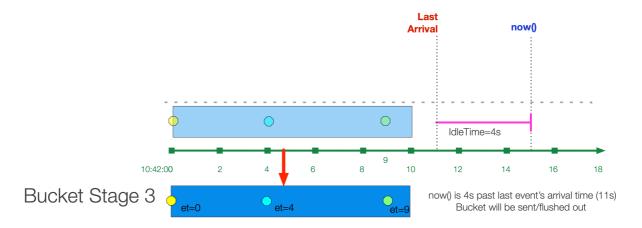
After the lag time expires, the bucket moves to Stage 2.



If the bucket is created from a historic stream, then the bucket is initiated in Stage 2. Lag time is not considered. A "historic" stream is one where the latest time of a bucket is before now(). E.g., if the window size is 10s, and now()=10:42:42, an event with event_time=10 will be placed in a Stage 2 bucket with range 10:42:10 - 10:42:20.

Idle Bucket Time Limit

While Lag Tolerance works with event time, Idle Bucket Time Limit works on arrival time (i.e., real time). It is defined as the amount of time to wait before flushing a bucket that has not received events.



After the Idle Time limit is reached, the bucket is "flushed" and sent out of the system.

Examples

Assume we're working with VPC Flowlog events that have the following structure:

version account_id interface_id srcaddr dstaddr srcport dstport protocol
packets bytes start end action log_status

For example:

2 99999XXXXX eni-02f03c2880e4aaa3 10.0.1.70 10.0.1.11 9999 63030 6 6556 262256 1554562460 1554562475 ACCEPT OK

2 496698360409 eni-08e66c4525538d10b 37.23.15.38 10.0.2.232 4373 8108 6 1 52 1554562456 1554562466 REJECT OK

Scenario A:

Every 10s, compute sum of bytes and output it in a field called Total Bytes.

Time Window: 10s

Aggregations: sum(bytes).as(TotalBytes)

Scenario B:

Every 10s, compute sum of bytes, outputitin a field called Total Bytes, group by srcaddr.

Time Window: 10s

Aggregations: sum(bytes).as(TotalBytes)

Group by Fields: srcaddr

Scenario C:

Every 10s, compute sum of bytes but only where action is REJECT, output it in a field called TotalBytes, group by srcaddr.

Time Window: 10s

Aggregations: sum(bytes).where(action='REJECT').as(TotalBytes)

Group by Fields: srcaddr

Scenario D:

Every 10s, compute sum of bytes but only where action is REJECT, output it in a field called TotalBytes. Also, compute distinct count of srcaddr.

Time Window: 10s

Aggregations:

sum(bytes).where(action='REJECT').as(TotalBytes)
distinct_count(srcaddr).where(action='REJECT')

i For further examples, see Engineering Deep Dive: Streaming Aggregations Part 2Memory Optimization

CEF Serializer

Description

The CEF Serializer takes a list of fields and/or values, and formats them in the Common Event Format (CEF) standard. CEF defines a syntax for log records. It is composed of a standard prefix, and a variable extension formatted as a series of key-value pairs.

Format

CEF:Version|Device Vendor|Device Product|Device Version|Device
Event Class ID|Name|Severity|[Extension]

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Output field: The field to which the CEF formatted event will be output. Nested addressing supported. Defaults to _raw .

Header Fields

CEF Header field definitions. The field values below will be written pipe (|) – delimited in the Output Field. Names cannot be changed. Values can be computed with JS expression, or can be constants.

- **cef_version**: Defaults to CEF:0.
- device_vendor: Defaults to Cribl.
- device_product: Defaults to Cribl.
- device_version: Defaults to C.version.
- device_event_class_id: Defaults to 420.

- name: Defaults to Cribl Event.
- severity: Defaults to 6.

Extension Fields

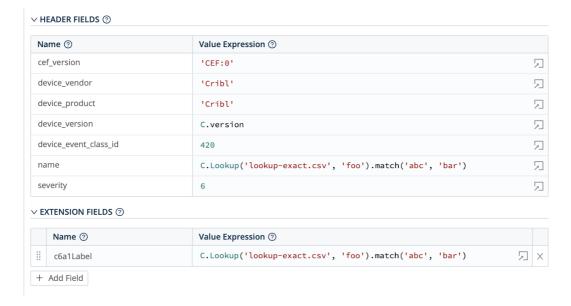
CEF Extension field definitions. Field names and values will be written in key=value format. Select each field's Name from the drop-down list. Values can be computed with JS expressions, or can be constants.

Example

For each CEF field, allowed values include strings, plus any custom Cribl function. For example, if using a lookup:

Name: Name
Value expression: C.Lookup('lookup-exact.csv', 'foo').match('abc', 'bar')

This can be used for any of the CEF Header Fields.



The resulting event has the following structure for an **Output Field** set to _CEF_out:

_CEF_out:CEF:0|Cribl|Cribl|42.0-61c12259|420|Business Group 6|6|c6a1Label=Colorado_Ext_Bldg7

Clone

Description

The Clone Function clones events, with optional added fields. Cloned events will be sent to the same Destination as the original event, because they are in the same Pipeline.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Clones: Create clones with the specified fields added and set.

Fields: Set of key-value pairs to add. Nested addressing is supported.

Examples

In this example, the Destination will receive a clone with an env field set to staging .

Field: env

Value: staging

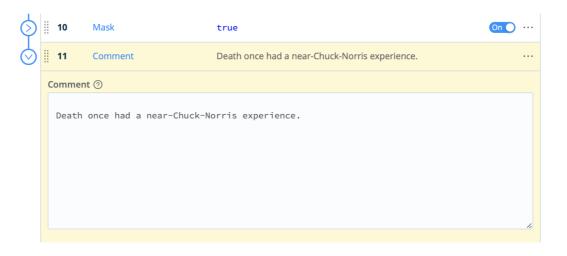
Comment

Description

The Comment Function adds a text comment in a Pipeline. It makes no changes to event data. The added comment is visible only within the Pipeline UI, where it is useful for labeling Pipeline steps.

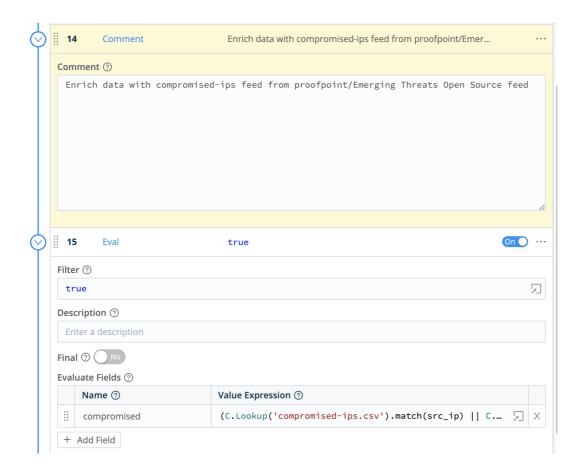
Usage

Comment: Add your comment as plain text in this field.



Examples

This comment labels the Pipeline's next function:



DNS Lookup

Description

The DNS Lookup Function offers two operations useful in enriching security and other data:

- DNS lookups based on host name as text, resolving to A record (IP address) or to other record types.
- Reverse DNS Lookup. (This duplicates LogStream's existing Reverse DNS Function, which is now deprecated.)

To reduce DNS lookups and minimize latency, the DNS Lookup Function incorporates a configurable DNS cache. If you need additional caching, consider enabling OS-level DNS caching on each LogStream Worker that will execute this Function. (OS-level caching options include DNSMasq, nscd, systemd-resolved, etc.)

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No .

DNS Lookup Fields Section

Lookup field name: Name of the field containing the domain to look up.

Resource record type: DNS record type (RR) to return. Defaults to A 'record.

Output field name: Lookup result(s) will be added to this field. Leave blank to overwrite the original field specified in **Lookup field name**.

Reverse DNS Lookup Field(s) Section

Lookup field name: Name of the field containing the IP address to look up.

⚠ If the field value is not in IPv4 or IPv6 format, the lookup is skipped.

Output field name: Name of the field in which to add the resolved hostname. Leave blank to overwrite the original field specified in **Lookup field name**.

Advanced Settings

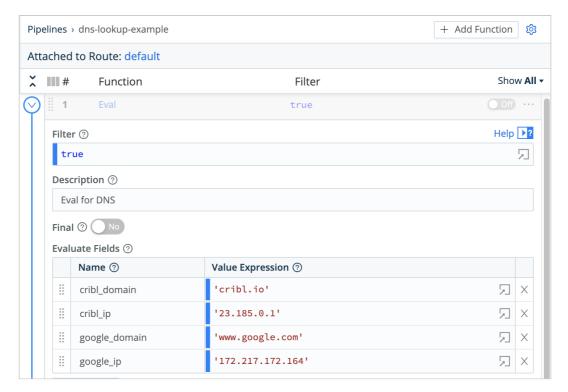
DNS server(s) overrides: IP address(es), in RFC 5952 format, of the DNS server(s) to use for resolution. IPv4 examples: 1.1.1.1, 4.2.2.2:53 . IPv6 examples: [2001:4860:4860::8888], [2001:4860:4860::8888]:1053. If this field is not specified, LogStream will use the system's DNS server.

Reload period (minutes): How often to refresh the DNS cache. Use 0 to disable refreshes. Defaults to 30 minutes.

Maximum cache size: Maximum number of DNS resolutions to cache locally. Before changing the default 5000, contact Cribl Support to understand the implications. Highest allowed value is 10000.

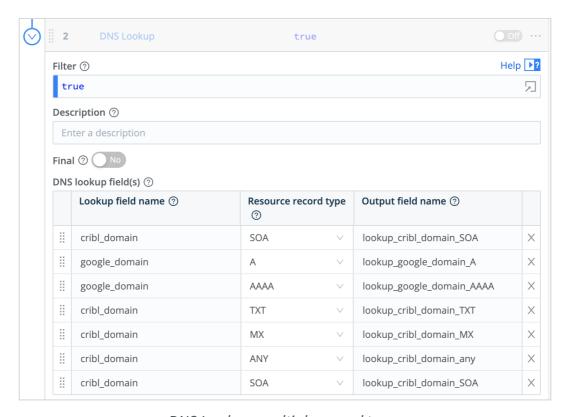
Example

This example Pipeline chains two Functions. First, we have an **Eval** Function that defines key-value pairs for two alphabetical domain names and two numeric IP addresses.



DNS Lookup: Eval Function

Next, the DNS Lookup Function looks up several record types for the two domain names, placing each retrieved record type in its own output field.



DNS Lookup: multiple record types

Finally, the same Function's **Reverse DNS lookup** section retrieves domain names for the two IP addresses.



DNS Lookup: reverse lookups

Drop

Description

The Drop Function will drop/delete any events that meet the Filter expression.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Example

Assume that we care only about errors, so we want to filter out any events that contain the word "success," regardless of case: "success," "SUCCESS," etc.

In our Drop Function, we'll use the JavaScript search() method to search the _raw field's contents for our target pattern. We know that search() returns a non-negative integer to indicate the starting position of the first match in the string, or -1 if no match. So we can evaluate the Function as true when the return value is >= 0.

Filter: _raw.search(/success/i) ≥ 0

Dynamic Sampling

Description

The Dynamic Sampling Function filters out events based on an expression, a sample mode, and events' volume. Your sample mode's configuration determines what percentage of incoming events will be passed along to the next step.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events passed into the Function will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Sample mode: Defines how sample rate will be derived. For formulas and usage details, see **Sample Modes** below. Supported methods:

- Logarithmic (the default): log(previousPeriodCount).
- Square root: sqrt(previousPeriodCount).

Sample group key: Expression used to derive sample group key. For example: \${domain}:\${httpCode} . Each sample group will have its own derived sampling rate, based on volume. Defaults to `\${host}`.

All events without a host field passing through the Function will be associated with the same group and sampled the same.

Advanced Settings

- Sample period Sec: How often (in seconds) sample rates will be adjusted.

 Defaults to 30.
- **Minimum events**: Minimum number of events that must be received, in previous sample period, for sampling mode to be applied to current

period. If the number of events received for a sample group is less than this minimum, a sample rate of 1:1 is used. Defaults to 30.

• Max sampling rate. Maximum sampling rate. If the computed sampling rate is above this value, the rate will be limited to this value.

How Does Dynamic Sampling Work

Compared to static sampling, where users must select a sample rate a priori, Dynamic Sampling allows for **automatically adjusting** sampling rates, based on incoming data volume per sample group. This Function allows users to set only the aggressiveness/coarseness of this adjustment. Square Root is more aggressive than Logarithmic mode.

As an event passes through the Function, it's evaluated against the Sample Group Key expression to determine the sample group it will be associated with. For example, given an event with these fields: ... ip=1.2.3.42, port=1234 ..., and a Sample Group Key of `\${ip}:\${port}`, the event will be associated with the 1.2.3.42:1234 sample group.

⚠ If the Sample Group Key is left at its `\${host}` default, all events without a host will be associated with the same group and sampled the same.

When a sample group is new, it will initially have a sample rate of 1:1 for Sample Period seconds (this value defaults to 30 seconds). Once Sample Period seconds have elapsed, a sample rate will be derived based on the configured Sample Mode, using the sample group's event volume during the **previous** sample period.

For example, assuming a Logarithmic Sample Mode:

```
Period 0 (first 30s): Number of events in sample group: 1000 , Sample Rate:
1:1 , Events allowed: ALL
Sample Rate calculation for next period: Math.ceil(Math.log(1000)) = 7
```

Period 1 (next 30s) -- Number of events in sample group: 4000 , Sample Rate:
7:1: Events allowed: 572
Sample Rate calculation for next period: Math.ceil(Math.log(4000)) = 9

Period 2 (next 30s) -- Number of events in sample group: 12000, Sample Rate: 9:1: Events allowed: 1334

Sample Rate calculation for next period: Math.ceil(Math.log(12000)) =
10

Period 3 (next 30s) -- Number of events in sample group: 2000, Sample Rate:

10:1: Events allowed: 200

Sample Rate calculation for **next** period: Math.ceil(Math.log(2000)) = 8

• • •

Sample Modes

- Logarithmic The sample rate is derived, for each sample group, using a natural log: Math.ceil(Math.log(lastPeriodVolume)). This mode is less aggressive, and drops fewer events.
- 2. Square Root The sample rate is derived, for each sample group, using:
 Math.ceil(Math.sqrt(lastPeriodVolume)) . This mode is more
 aggressive, and drops more events.

Example

Here's an example that illustrates the effectiveness of using the Square Root sample mode.

Settings:

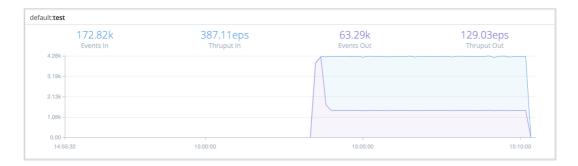
Sample Mode: Square Root

Sample Period (sec): 20

Minimum Events: 3
Max. Sampling Rate: 3

Results:

Events In: 4.23K Events Out: 1.41K



In this generic example, we reduced the incoming event volume from 4.23K to 1.41K. Your own results will vary depending on multiple parameters – the Sample Group Key, Sample Period, Minimum Events, Max Sampling Rate, and rate of incoming events.

i For further examples, see Getting Smart and Practical With Dynamic Sampling.

Eval

Description

The Eval Function adds or removes fields from events. (In Splunk, these are index-time fields.)

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Evaluate fields: Set of key/value pairs to add. The left-hand side input (**Name**) is the key name. The right-hand side input (**Value Expression**) is a JS expression to compute the value – this can be a constant. Nested addressing is supported. Strings intended to be used as values must be single- or double-quoted. (For details, see Introduction to Expression Syntax.)

Keep fields: List of fields to keep. Wildcards (*) and nested addressing are supported. Takes precedence over **Remove fields** (below). To reference a parent object and all children requires using the (*) wildcard. For example, if _raw is converted to an object then use _raw* to refer to itself and all children.

Remove fields: List of fields to remove. Wildcards (*) and nested addressing are supported. Cannot remove fields matching Keep fields. Cribl LogStream internal fields that start with __ (double underscore) cannot be removed via wildcard. Instead, they need to be specified individually. For example, __myField cannot be removed by specifying __myF*.

Using Keep and Remove

A field matching an entry in *both* **Keep** (wildcard or not) and **Remove** will *not* be removed. This is useful for implementing "remove all but" functionality. For

example, to keep only _time, _raw, source, sourcetype, host, we can specify them all in **Keep**, while specifying * in **Remove**.

Negated terms are supported in both **Keep fields** and **Remove fields**. The list is order-sensitive when negated terms are used. Examples:

- !foobar, foo* means "All fields that start with 'foo' except foobar."
- !foo*, * means "All fields except for those that start with 'foo!"

Examples

Scenario A: Create field myField with static value of value1:

- Name: myField
- Value Expression: 'value1'

Scenario B: Set field action to blocked if login=error:

- Name: action
- Value Expression: login='fail' ? 'blocked' : action

Scenario C: Create a multivalued field called myTags . (i.e., array):

- Name: myTags
- Value Expression: ['failed', 'blocked']

Scenario D: Add value error to the multivalued field myTags:

- Name: myTags
- Value Expression: login='error' ? [... myTags, 'error'] : myTags

(The above expression is literal, and uses JavaScript spread syntax.)

Scenario E: Rename an identification field to the shorter ID – copying over the original field's value, and removing the old field:

- Name: ID
- Value Expression: identification
- Remove Field: identification
 - i See Ingest-time Fields for more examples.

Advanced Usage Notes

Note 1

The Eval Function has the ability to execute expressions without assigning their value to the field of an event. You can do this by simply leaving the left-hand side input empty, and having the right-hand side do the assignment.

- Simple Example: Object.assign(foo, JSON.parse(bar), JSON.parse(baz)) on the right-hand side (and left-hand side empty) will JSON-parse the strings in bar and baz, merge them, and assign their value to foo, an already existing field.
- Another Example: To parse JSON, enter Object.assign(__e, JSON.parse(_raw)) on the right-hand side (and left-hand side empty).
 __e is a special variable that refers to the (context) event within a JS expression. In this case, content parsed from _raw is added at the top level of the event.

Note 2

You can also use the Eval Function to set and unset control fields (e.g., _TCP_ROUTING in Splunk), via this syntax: _ctrl.<name> . Control fields can be referenced only on the left-hand side of Add. (I.e., they cannot be read or used on the right-hand side, and cannot be referenced in Remove.)

To unset/delete a control field, set its value to undefined. These fields are normally not needed for event computations, and modifying them is suggested to be done only by experts. Please reach out to Cribl if you need help with this topic.

Flatten

Description

The Flatten Function is used to flatten fields out of a nested structure.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Fields: List of top-level fields to include for flattening. Supports * wildcards. Defaults to empty array, which means all fields.

Prefix: Prefix string for flattened field names. Defaults to empty.

Depth: Number representing the nested levels to consider for flattening. Minimum 1. Defaults to 5.

Delimiter: Delimiter to be used for flattening. Defaults to _ (underscore).

Example

Add the following test sample in **Preview** > **Paste a Sample**:

```
input
{ "accounting" : [ { "firstName" : "John", "lastName" : "Doe", "age" : 23 }, { "firstName" : "Mage" : "Mage" : 23 }, }
```

Under **Select Event Breaker**, choose **ndjson** (newline-delimited JSON), and click **Save as a Sample File**.

Here's sample output with all settings at default:

```
output
{
    "accounting_0_firstName": "John",
    "accounting_0_lastName": "Doe",
    "accounting_0_age": 23,
    "accounting_1_firstName": "Mary",
    "accounting_1_lastName": "Smith",
    "accounting_1_age": 32,
    "sales_0_firstName": "Sally",
    "sales_0_lastName": "Green",
```

```
"sales_0_age": 27,
"sales_1_firstName": "Jim",
"sales_1_lastName": "Galley",
"sales_1_age": 41,
}
```

Using the Flatten Function's default settings, we successfully create top-level fields from the nested JSON structure, as expected.

GeoIP

Description

The GeoIP Function enriches events with geographic fields, given an IP address. It is optimized for binary databases such as MaxMind's GeoIP.

☐ For details on setting up MaxMind (and similar) databases, see Managing Large Lookups.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

GeoIP file (.mmdb): Path to a Maxmind database, in binary format, with .mmdb extension.

i If the database file is located within the lookup directory (\$CRIBL_HOME/data/lookups/), the GeoIP file does not need to be an absolute path.

In distributed deployments, ensure that the Maxmind database file is in the same location on both the Master and Worker Nodes.

IP field: Field name in which to find an IP to look up. Can be nested. Defaults to ip.

Result field: Field name in which to store the GeoIP lookup results. Defaults to geoip.

Examples

Assume that you are receiving SMTP logs, and need to see geolocation information associated with IPs using the SMTP service.

Here's a sample of our data, from IPSwitch IMail Server logs:

```
03:19 03:22 SMTPD(00180250) [192.168.1.131] connect 74.136.132.88 port 2539 03:19 03:22 SMTPD(00180250) [74.136.132.88] EHLO msnbc.com 03:19 03:22 SMTPD(00180250) [74.136.132.88] MAIL FROM: <info-jjgcdshx@test.us> 03:19 03:22 SMTPD(00180250) [74.136.132.88] RCPT To:<user@domain.com>
```

In this example, we'll chain together three Functions. First, we'll use a Regex Extract Function to isolate the host's IP. Next, we'll use the GeoIP Function to look up the extracted IP against our geoIP database, placing the returned info into a new __geoip field. Finally we'll use an Eval Function to parse that field's city, state, country, ZIP, latitude, and longitude.

Function 1 – Regex Extract

Regex: $[(?<ip>\S+)]$

Source field: _raw

Result: 74.136.132.88

Function 2 - GeoIP

Event's IP field: ip
Result field: __geoip

Function 3 - Fval

Name	Value Expression
City	geoip.city.names.en
Country	geoip.country.names.en
Zip	geoip.postal.code
Lat	geoip.location.latitude
Long	geoip.location.longitude

In the Eval Function's Remove fields setting, you could specify thegeoip		
field for removal, if desired. However, its prefix makes it an internal field		
anyway.		
☐ For a hosted tutorial on applying the GeoIP Function, see Cribl's GeoIP and Threat Feed Enrichment Sandbox.		

Grok

Description

The Grok Function extracts structured fields from unstructured log data, using modular regex patterns.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Optional description of this Function's purpose in this Pipeline. Defaults to empty.

Final: If toggled to Yes, stops data from being fed to downstream Functions. Defaults to No.

```
Pattern: Grok pattern to extract fields. Syntax supported: % {PATTERN_NAME: FIELD_NAME}.
```

Click + Add pattern to chain more patterns.

Source field: Field on which to perform Grok extractions. Defaults to _raw .

Management

You can add and edit Grok patterns via LogStream's UI by selecting **Knowledge > Grok**Patterns. Pattern files are located at: \$CRIBL_HOME/(default|local)/cribl/grokpatterns/

Example

Example event:

```
{"_raw": "2020-09-16T04:20:42.45+01:00 DEBUG This is a sample debug log message"}`

Pattern: %{TIMESTAMP_IS08601:event_time} %{LOGLEVEL:log_level} %

{GREEDYDATA:log_message}
```

Source Field: _raw

Event after extraction:

```
{"_raw": "2020-09-16T04:20:42.45+01:00 DEBUG This is a sample debug log message",
    "_time": 1600226442.045,
    "event_time": "2020-09-16T04:20:42.45+01:00",
    "log_level": "DEBUG",
    "log_message": "This is a sample debug log message",
}
```

Note the new fields added to the event: event_time, log_level, and log_message.

References

- Syntax for a Grok pattern is %{PATTERN_NAME:FIELD_NAME} . E.g.: %{IP:client} % {WORD:method} .
- Useful links for creating and testing Grok patterns: http://grokdebug.herokuapp.com and http://grokconstructor.appspot.com/.
- Additional patterns are available here: https://github.com/logstash-plugins/logstash-patterns-core/tree/master/patterns.

JSON Unroll

Description

The JSON Unroll Function accepts a JSON object string _raw field, unrolls/explodes an array of objects therein into individual events, while also inheriting top level fields. See example(s). Cribl highly recommends not using this JSON Unroll function for certain types of data. Instead, perform the unrolling using an event breaker for those inputs which support configuring an event breaker. Specifying the event breaker type JSON Array and toggling the JSON Extract Fields option to Yes will accomplish the same unrolling but much more efficiently. This is recommended, for example, for CloudTrail and Office635 events, which are collected as JSON arrays.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Path: Path to array to unroll, e.g., foo.0.bar.

New name: The name that the exploded array element will receive in each new event. Leave empty to expand the array element with its original name.

Example(s)

Assume you have an incoming event that has a _raw field as a JSON object string like this:

```
Sample_raw field
```

Settings:

Path: allCars
New Name: cars

Output Events:

```
Resulting Event 1:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"Ford","models'

Event 2:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"GM","models":|

Event 3:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"Fiat","models'

Event 4:
{"_raw":"{"date":"9/25/18 9:10:13.000 PM","name":"Amrit","age":42,"cars":{"name":"Blackberry","r
```

Each element under the original **allCars** array is now placed in a **cars** field in its own event, inheriting original top level fields; **date**, **name** and **age**

Lookup

Description

The Lookup Function enriches events with external fields, using lookup table files in CSV, compressed .csv.gz, or binary .mmdb format.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Lookup file path (.csv, .csv.gz): Path to the lookup file. Select an existing file that you've uploaded via LogStream's UI at Knowledge > Lookups Libary, or specify the path. You can reference environment variables via \$, e.g.: \$CRIBL_HOME/file.csv.

i When you configure this field via a distributed deployment's Master Node, LogStream will swap \$CRIBL_HOME/groups/<groupname>/ for \$CRIBL_HOME when validating whether the file exists. In this case, the default upload path changes from \$CRIBL_HOME/data/lookups (single-instance deployments) to \$CRIBL_HOME/groups/<groupname>/data/lookups/ (distributed deployments).

Match mode: Defines the format of the lookup file, and indicates the matching logic that will be performed. Defaults to Exact .

Match type: For CIDR and Regex Match modes, this attribute refines how to resolve multiple matches. First match will return the first matching entry. Most specific will scan all entries, finding the most specific match. All will return all matches in the output, as arrays. (Defaults to First match. Not displayed for Exact Match mode.)

Lookup fields (.csv): Field(s) that should be used to key into the lookup table.

- Lookup field name in event: Exact field name as it appears in events. Nested addressing supported.
- Corresponding field name in lookup: The field name as it appears in the lookup file. Defaults to the Lookup field name in event value. This input is optional.

△ Case-Sensitive / Multiple Matches

Lookups are case-sensitive by default. (See the **Ignore case** option below.)

If the lookup file contains duplicate key names with different values, all **Match mode**s of this Function will use **only** the value in the key's **final** instance, ignoring all preceding instances.

Output field(s): Field(s) to add to events after matching the lookup table. Defaults to **all** if not specified.

- Output field name from lookup: Field name, as it appears in the lookup file.
- Lookup field name in event: Field name to add to event. Defaults to the lookup field name. This input is optional. Nested addressing is supported.

Advanced Settings

Reload period (sec): Periodically check the underlying file for modtime changes, and reload if necessary. Use -1 to disable. Defaults to 60.

Ignore case: Ignore case when performing Match mode: Exact lookups. Defaults to No.

Add to raw event: Whether to append the looked-up values to the _raw field, as key=value pairs. Defaults to No.

Examples

Example 1: Regex Lookups

Assign a sourcetype field to events if their _raw field matches a particular regex.

```
regex, sourcetype
"^[^,]+,[^,]+,[^,]+,THREAT",pan:threat
"^[^,]+,[^,]+,[^,]+,TRAFFIC",pan:traffic
"^[^,]+,[^,]+,[^,]+,SYSTEM",pan:system
```

Match mode: Regex

Match type: First match

Lookup field name in event: _raw

Corresponding field name in lookup: regex

```
Events before and after

### BEFORE:

{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:58
{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,FOOBAR,THREAT,end,2049,2018/09/20 13:03:58
```

```
### AFTER:
```

```
{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,F00BAR,TRAFFIC,end,2049,2018/09/20 13:03:5
    "sourcetype": "pan:traffic"
    }
{"_raw": "Sep 20 13:03:55 PA-VM 1,2018/09/20 13:03:58,F00BAR,THREAT,end,2049,2018/09/20 13:03:55
    "sourcetype": "pan:threat"
}
```

Example 2: CIDR Lookups

Assign a location field to events if their destination_ip field matches a particular CIDR range.

```
paloaltoips.csv

range,location
10.0.0.0/24,San Francisco
10.0.0.0/16,California
10.0.0.0/8,US
```

Match mode: CIDR

Match type: See options below

Lookup field name in event: destination ip

Corresponding field name in lookup: range

i In Match mode: CIDR with Match type: Most specific, the lookup will implicitly search for matches from most specific to least specific. There is no need to pre-sort data.

Note that Match mode: CIDR with Match type: First Match is likely the most performant with large lookups. This can be used as an alternative to Most specific, if the file is sorted with the most specific/relevant entries first. This mode still performs a table scan, top to bottom.

Events before and after

```
### BEFORE:
{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,F00BAR,TRAFFIC,end,2049,2018/09/20 13:03
   "destination_ip": "10.0.0.102"
}
### AFTER with Match Type: First Match
{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,F00BAR,TRAFFIC,end,2049,2018/09/20 13:03
   "destination_ip": "10.0.0.102",
   "location": "San Francisco"
}
### AFTER with Match Type: Most Specific
```

```
{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:
   "destination_ip": "10.0.0.102",
   "location": "San Francisco"
}

### AFTER with Match Type: All

{"_raw": "Sep 20 13:03:55 PA-VM 1, 2018/09/20 13:03:58,FOOBAR,TRAFFIC,end,2049,2018/09/20 13:03:
   "destination_ip": "10.0.0.102",
   "location": [
        "San Francisco",
        "California",
        "US",
   ]}
```

More Examples and Scenarios

More examples:

- Ingest-time Lookups.
- Lookups and Regex Magic.
- Lookups as Filters for Masks.

See also:

- Managing Large Lookups to optimize file locations for large lookup files.
- Redis Function for faster lookups using a Redis integration.

Mask

Description

The Mask Function masks, or replaces, patterns in events. This is especially useful for redacting PII (personally identifiable information) and other sensitive data.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Masking rules: Match Regex and Replace Expression pairs. Defaults to empty. Each row has the following fields:

- Match regex: Pattern to replace. Supports capture groups. Use /g to replace all matches, e.g.: /foo(bar)/g
- **Replace expression**: A JavaScript expression or literal to replace all matching content.

To add more rows, click + Add Rule.

Apply to fields: Fields on which to apply the masking rules. Defaults to $_{\rm raw}$. Add more fields by typing in their names, separated by hard returns. Supports wildcards (\star) and nested addressing.

i Negated terms are also supported. When you negate field names, the fields list is order-sensitive. E.g., !foobar before foo* means "Apply to all fields that start with foo , except foobar." However, !foo* before * means "Apply to all fields, except for those that start with foo."

Advanced Settings

Evaluate fields: Optionally, specify fields to add to events in which one or more of the **Masking Rules** were matched. These fields can be useful in downstream processing and reporting. You specify the fields as key–value expression pairs, like those in the Eval Function.

- Name: Field name.
- Value Expression: JavaScript expression to compute the value (can be a constant).

Evaluating the Replace Expression

The **Replace expression** field accepts a full JS expression that evaluates to a value, so you're not necessarily limited to what's under C.Mask . For example, you can do conditional replacement: g1%2=1 ? `fieldA="odd"`: `fieldA="even"`

The Replace expression can reference other event fields as event. <fieldName> . For example, `\${g1}\${event.source}` . Note that this is slightly different from other expression inputs, where event fields are referenced without event. Here, we require the event. prefix for the following reasons:

- We don't expect this to be a common case.
- Expanding the event in the replace context would have a high performance hit on the common path.
- There is a slight chance that there might be a gN field in the event.

Examples

Example 1: Transform a String

Here, we'll simply search for the string dfhgdfgj, and replace that value (if found) with Trans AM. This will help close America's muscle-car gap:

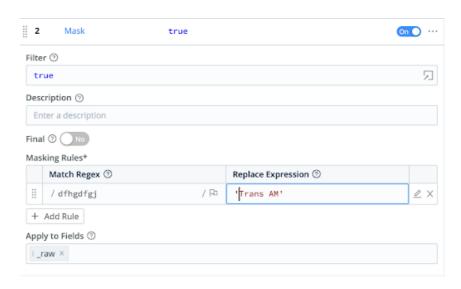


Event before masking

Configure the Mask Function > Masking Rules as follows:

Match Regex: dfhgdfgj

Replace Expression: Trans AM



Mask Function configuration

Result: Vroom vroom!

Event after masking

Example 2: Mask Sensitive Data

Assume that you're ingesting data whose <code>_raw</code> fields contain unredacted Social Security numbers in the Key=Value pattern <code>social=########</code>.

```
a_raw: 2020-07-22 05:22:43,330,Event [Event=UpdateBillingProvQuote, timestamp=1577371]
       0, properties={JMSCorrelationID=NA, JMSMessageID=ID:ESP-PD.C7A19FC656293:AB21BCF
       E, orderType=NewActivation, quotePriority=NORMAL, conversationId=ESB~BEBFAB927C87
      5E35:81E10EA8:47283ADA8A10:5568, credits=NA, JMSReplyTo=pub.esb.genericasync.resp
      onse, timeToLive=-1, serviceName=UpdateBillingProvisioning, esn=10D9C064A00987, a
      ccountNumber=900001336, social=518057110, MethodName=InternalEvent, AdapterName=U
      pdateBillingProvQuote, meid=NA, orderNumber=900000000002363, quoteNumber=4258319
      8, ReplyTo=NA, userName=yosem7, EventConversationID=NA, mdn=6248526355, accountTy
      pe=PostPaid, marketCity="JOLIET", marketState=IL, marketZip=60432, billingCycle=2
       4, autoBillPayment=T, phoneCode=SGS5, phoneType=Android, phoneName="Samsung GALAX
       Y S5", planCode=1400POST5L90, planType=PostPaid, planPrice=89.99, planName="1400"
       Minute Family", planDescription="Nationwide 1400 Minutes, Unlimited Mobile to Mob
       ile, Unlimited Night & Weekend, Unlimited Data", cardNumber=3569948084568945, net
       workProviderName=Splunktel}] Show less
# time: 1595395363.33
α host: 127.0.0.1
a index: cribl
a source: /opt/tibco/tra/apps/ESB/logs/business_event.log
α sourcetype: business_event
```

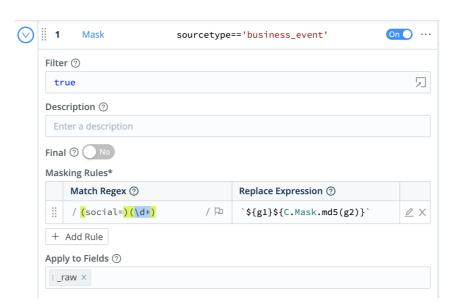
Event with unredacted SSNs

You can use a Mask Function to run an md5 hash of the social keys' numeric values, replacing the original values with the hashed values. Configure the Masking Rules as follows:

Match Regex: (social=)(\d+)
Replace Expression: `\${g1}\${C.Mask.md5(g2)}`

In the first example everything in the Match regex field was replaced by the Replace Expression. However if that isn't desired then you can use capture groups in the Match Regex to define individual string components for manipulation or, alternatively, use string literals in the Replace expression for retaining any static text. Any content matching the Match Regex that is not inserted into the Replace expression will not be retained.

In this example, social= is assigned to capture group g1 for later reference. The value of social= will be hashed by referencing it as g2 in the md5 function. If we didn't make social= its own capture group (or specified social= as a literal in the Replace Expression) then we cannot reference it using g1 in the Replace expression, the value of social= would instead be assigned to g1, and the entire social=####### string would be replaced with a hash of the social security number, which probably isn't desired because no one would know the value being hashed without a field name preceding it.



Mask Function configuration

Result: The sensitive values are replaced by their md5 hashes.

```
a _raw: 2020-07-22 05:22:43,330,Event [Event=UpdateBillingProvQuote, timestamp=1577371270, properties={JMSCorrel
      ationID=NA, JMSMessageID=ID:ESP-PD.C7A19FC656293:AB21BCFE, orderType=NewActivation, quotePriority=NORMA
       L, conversationId=ESB~BEBFAB927C875E35:81E10EA8:47283ADA8A10:5568, credits=NA, JMSReplyTo=pub.esb.generi
       casync.response, timeToLive=-1, serviceName=UpdateBillingProvisioning, esn=10D9C064A00987, accountNumber
       =900001336, social=dlce1763a8e5213781a30f8e7ba9172f, MethodName=InternalEvent, AdapterName=UpdateBilling
       ProvQuote, meid=NA, orderNumber=9000000000002363, quoteNumber=42583198, ReplyTo=NA, userName=yosem7, Eve
       ntConversationID=NA, mdn=6248526355, accountType=PostPaid, marketCity="JOLIET", marketState=IL, marketZi
       p=60432, billingCycle=24, autoBillPayment=T, phoneCode=SGS5, phoneType=Android, phoneName="Samsung GALAX
       Y S5", planCode=1400POST5L90, planType=PostPaid, planPrice=89.99, planName="1400 Minute Family", planDes
       cription="Nationwide 1400 Minutes, Unlimited Mobile to Mobile, Unlimited Night & Weekend, Unlimited Dat
       a", cardNumber=3569948084568945, networkProviderName=Splunktel}] Showless
# time: 1595395363.33
α cribl_pipe: business_event
α host: 127.0.0.1
a index: cribl
a source: /opt/tibco/tra/apps/ESB/logs/business_event.log
\alpha sourcetype: business_event
```

Event with hashed SSNs

- i In scenarios where you need to send unmodified values to certain Destinations (such as archival stores), you can narrow the Mask Function's scope by setting the associated Route's Output field.
 - For further masking examples, see Masking and Obfuscation.

Numerify

Description

The Numerify Function converts event fields that are numbers to type number.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Ignore Fields: Specify fields to **not** numerify. Type in field names, separated by hard returns. By default, this list is empty, and Numerify applies to **all** fields. Supports wildcards (*) and nested addressing.

i Double Negatives

Negated terms are also supported. When you negate field names, the fields list is order-sensitive. E.g., !foobar before foo* means "Ignore all fields that start with foo, except foobar." However, !foo* before * means "Ignore all fields, except for those that start with foo."

Format: Optionally, reformat or truncate the extracted numeric value. Select one of:

- None: Applies no reformatting (the default).
- **Floor**: Rounds the number down to the lower adjacent integer (truncates it).
- **Ceil**: Rounds the number up to the higher adjacent integer, removing decimal digits.

- Round: Rounds (truncates) the number to a specified number of digits. This option exposes an extra field:
 - **Digits**: Number of digits after the decimal point. Enter a value between 0 20; defaults to 2.

Examples

Scenario A:

Assume an event whose text contains a numeric value that must be extracted to perform some numeric analysis. The text looks like this:

```
version=11.5.0.0.1.1588476445
```

We can extract the numeric value by chaining together two Functions:

- 1. A Regex Extract Function. Set its **Regex** field to /version=(?<ver>\d+)/, to capture the first set of digits found in the event string.
- 2. Then use Numerify.

This captures the substring 11 and converts it to a numeric 11 value.

Scenario B:

Assume email transaction log events like the sample below. The final field is the message's size, in bytes. We want to extract this as a numeric value, for analysis in LogStream or downstream services:

```
03:19 03:22 SMTPD (00180250) [209.221.59.70] C:\IMail\spool\D28de0018025017cd.SMD 3827
```

Again, we can accomplish this with two Functions:

- 2. Then use Numerify.

Parser

Description

The Parser Function can be used to extract fields out of events, or to reserialize (rewrite) events with a subset of fields. Reserialization will maintain the format of the events.

For example: If an event contains comma-delimited fields, and fieldA and fieldB are filtered out, those fields' positions will be set to null, but not deleted completely.

Parser cannot remove fields that it did not create. A subsequent Eval Function can do so.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Operation mode: Extract will create new fields. **Reserialize** will extract, filter fields, and then reserialize.

Type: Parser/Formatter type to use. Options:

- CSV
- Extended Log File Format (ELFF)
- Common Log Format (CLF)
- K=V Pairs
- JSON
- Delimited Values

Setting **Type** to **Delimited Values** displays the following extra options:

- **Delimiter**: Delimiter character to split value. Defaults to comma (,). You can also specify pipe (|) or tab characters.
- Quote char: Character used to quote literal values. Defaults to ".
- * Escape char: Character used to escape delimiter or quote characters.
 Defaults to \ .
- **Null value**: Field value representing the null value. These fields will be omitted. Defaults to .

Library: Select an option from the Parsers Library.

Source field: Field that contains text to be parsed. Not usually needed in Serialize mode.

Destination field: Name of field in which to add extracted and serialized fields. If multiple new fields are created and this setting is configured then all new fields are created as elements of an array with the array name set to the name specified for this setting. If you want all new fields to be independent, rather than in an array, then specify them using **List of fields** below. (Extract and Serialize modes only.)

Clean fields: This option appears for **Type: K=V Pairs**. Toggle to Yes to clean field names by replacing non-alphanumeric characters with _ . This will also strip leading and trailing " symbols.

List of fields: Fields expected to be extracted, in order. If not specified, Parser will auto-generate fields.

Fields to keep: List of fields to keep. Supports wildcards (*). Takes precedence over **Fields to remove**. Nested addressing supported.

Fields to remove: List of fields to remove. Supports wildcards (*). Cannot remove fields matching **Fields to keep**. Nested addressing supported.

Negated terms are supported in both Fields to remove and Fields to keep. When you use negated terms, the list is order-sensitive.
E.g., !foobar, foo* means "All fields that start with foo, except foobar." However, !foo*, * means "All fields, except for those that start with foo."

Fields filter expression: Expression to evaluate against {index, name, value} context of each field. Return truthy to keep, falsy to remove field. Index is zero-based.

How Fields Settings Interact

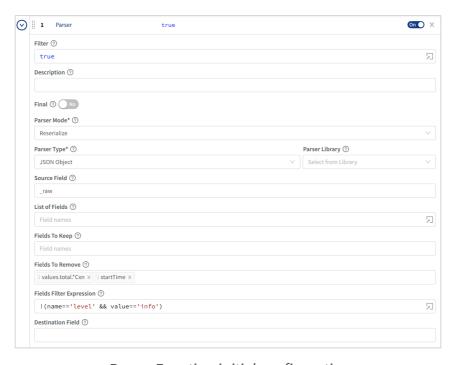
The **Fields to keep**, **Fields to remove**, and **Fields filter expression** settings interact as follows:

- Order of evaluation: Fields to keep > Fields to remove > Fields filter expression.
- If a field is in both Fields to keep and Fields to remove, Fields to keep takes precedence.
- If a field is in both Fields to remove and Fields filter expression, Fields to remove takes precedence.

Example 1

Insert the following sample, using Preview > Add a Sample > Paste a Sample: 2019/06/24 05:10:55 PM Z a=000,b=001,c=002,d=003,e=004,f=005,g1=006,g2=007,g3=008

Create the following test Parser Function (or import this Pipeline: https://raw.githubusercontent.com/weeb-cribl/cribl-samples/master/parser/functions/parser/parser_1.json).



Parser Function initial configuration

First, set the Parser type to Key=Value Pairs.

Scenario A:

Keep fields a, b, c. Drop the rest.

Expected result: a, b, c

• Fields to Keep: a, b, c

• Fields to Remove: *

• Fields Filter Expression:

Result: The event will gain four new fields and values, as follows.

• a: 000

• b: 001

• c: 002

• cribl_pipe: parser2



Scenario A result

You can check your stats by clicking the **Preview** pane's **Basic Statistics** (chart) button. In the resulting pop-up, the **Number of Fields** should have incremented ty four.

Now that you have the hang of it, try out the other simple scenarios below.

Scenario B:

Keep fields a, b, those that start with g. Drop the rest.

Expected result: a, b, g1, g2, g3

• Fields to keep: a, b

• Fields to remove: [empty]

• Fields filter expression: name.startsWith('g')

Scenario C:

Keep fields a, b, those that start with g but only if value is 007. Drop the rest.

Expected result: a, b, g2

- Fields to keep: a, b
- Fields to remove: [empty]
- Fields filter expression: name.startsWith('g') & value='007'

Scenario D:

Keep fields a, b, c, those that start with g, unless it's g1. Drop the rest.

Expected result: a, b, c, g2, g3

- Fields to keep: a, b, c
- Fields to remove: g1
- Fields filter expression: name.startsWith('g')

Scenario E:

Keep fields a, b, c, those that start with g but only if index is greater than 6. Drop the rest.

Expected result: a, b, c, g2, g3

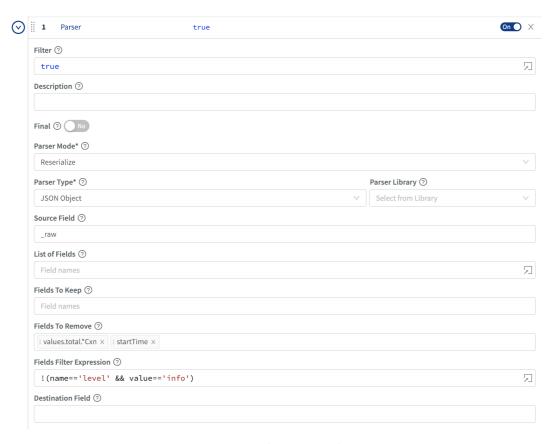
- Fields to keep: a, b, c
- Fields to remove: [empty]
- Fields filter expression: name.startsWith('g') & index>6
 - i The index refers to the location of a field in the array of all fields extracted by this Parser. It is zero-based. In the case above, g2 and g3 have index values of 7 and 8, respectively.

Example 2

Assume we have a JSON event that needs to be **reserialized**, given these requirements:

- 1. Remove the level field only if it's set to info.
- 2. Remove the startTime field, and all fields in the values.total. path that end in Cxn .

Parser Function configuration:



Parser Function configuration for Example 2

JSON event after being processed by the Function:

```
□ "_raw":{
□ "_raw":{
                                   "channel": "server"
   "channel": "server"
                                   "endTime":1549503300000
   "endTime":1549503300000
                                   "keyCount":0
   "keyCount":0
                                   "level": "info"
   "level":"info"
                                   "message":"_raw stats"
   "message":"_raw stats"
                                   "startTime": 1549503240000
   "startTime":1549503240000
                                   "time":1549503300401
   "time":1549503300401
                                   □ "values":{
  □ "values":{
                                      □ "total":{
     □ "total":{
                                         "activeCxn":2
        "activeCxn":2
                                         <del>"closeCxn"</del>:4
        "closeCxn":4
                                         "inBytes":61724
        "inBytes":61724
                                         "inEvents":210
        "inEvents":210
                                         "openCxn":4
        "openCxn":4
                                         "outBytes":61724
        "outBytes":61724
                                         "outEvents":210
        "outEvents":210
                                      }
                                   }
   }
                                }
}
```

Example 2 event transformation

Example 3

Insert the following sample, using **Preview > Add a Sample > Paste a Sample**:

```
2019/06/24 15:25:36 PM Z
a=000,b=001,c=002,d=003,e=004,f=005,g1=006,g2=007,g3=008,
```

For all scenarios below, first create a Parser Function to extract all fields, by setting the **Parser type** to Key=Value Pairs . Then add a second Parser Function with the configuration shown under **Parser 2**.

Scenario A:

Serialize fields a, b, c, d in CSV format.

Expected result: _raw field will have this value 000,001,002,003

Parser 2:

- Operation mode: Serialize
- Source field: [empty]
- Destination field: [empty]

- Type: CSV
- List of fields: a, b, c, d (needed for positional formats)

Scenario B:

Serialize fields a, b, c in JSON format, under a field called bar.

```
Expected result: bar field will be set to:
```

```
{"a":"000","b":"001","c":"002","d":"003"}
```

Parser 2:

- Operation mode: Serialize
- Source field: [empty]
- Destination field: bar
- Type: JSON
- List of fields: [empty]
- Fields to keep: a, b, c, d

Publish Metrics

Description

The Publish Metrics Function extracts, formats, and outputs metrics from events.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Overwrite: If set to Yes, overwrite previous metric specs. Otherwise, append. Defaults to No.

Metrics

Add Metrics: List of metrics to extract from the event and format. Destinations can pass the formatted metrics to a metrics aggregation platform. Click **Add Metrics** to add new rows containing the following options:

- Event field name: The name of the field (in the event) that contains the metric value.
- Metric name expression: JavaScript expression to evaluate the metric field name. Defaults to the Event field name value.
 - The JavaScript expression will evaluate the metric field name only after the metrics are processed for transport to the Destination. While in the processing Pipeline, the metric name expression appears as a literal.
- Metric type: Select Counter, Timer, or Gauge (the default).

Remove Metrics: Optionally, enter a List of field names to look for when removing metrics. Where a metric's field name matches an element in this list, LogStream will remove that metric from the event.

Dimensions

Add Dimensions: Optional list of dimensions to associate with every extracted metric value. If this Function is used to process output from the Aggregations Function, leave this field blank, because dimensions will be automatically discovered. Defaults to !_* *.

Remove Dimensions: Optional list of dimensions to associate with every extracted metric value. Leave blank if this function is used to process output from the Aggregation function as dimensions will be automatically discovered. If this Function is used to process output from the Aggregations Function, leave this field blank, because dimensions will be automatically discovered.

i The Add Dimensions and Remove Dimensions fields support wildcards and negated terms. When you use negated terms, the list is order-sensitive. E.g., !foobar before foo* means "All fields that start with foo, except foobar." However, !foo* before * means "All fields, except for those that start with foo."

Overwrite: If set to Yes, overwrite previous metric specs. Otherwise, append. Defaults to No.

Fields Color Coding

On the right Preview pane's **OUT** tab, the Publish Metrics Function adds the following color codes to field labels:

```
1  m  a _metric: metrics_pool.num_metrics

2020-03-05  a _metric_type: gauge

18:42:52.707  # _value: 229

a _cribl_pipe: publish_some_metrics

# _earliest: 1583451770704

a _host: zebrastripes

# _latest: 1583451772704

a _source: _cribl
```

Dimension: purple | Value: cyan (light blue) | Info: dark blue

These are in addition to the color codes applied to field values, which are listed here.

Examples

Scenario A:

Assume we're working with AWS VPC Flowlog events that have the following structure:

version account_id interface_id srcaddr dstaddr srcport dstport protocol packets
bytes start end action log_status

For example:

```
2 99999XXXXX eni-02f03c2880e4aaa3 10.0.1.70 10.0.1.11 9999 63030 6 6556 262256 1554562460 1554562475 ACCEPT OK
```

... and we want to use values of packets and bytes as metrics across these dimensions: action, interface_id, and dstaddr.

To reference the packets and bytes fields by name, as 'packets' and 'bytes', our Pipeline will need a Parser Function before the Publish Metrics Function.

Parser Function

Filter: Set as needed Operation mode: Extract

Type: Extended Log File Format (automatically set when specifying a library)

Library: AWS VPC Flow Logs

Source: _raw

(No need to specify any other fields.)

Publish Metrics Function

Below, the metric_name prefix was arbitrarily chosen. Because there is no JavaScript expression to evaluate – i.e. this is literal text – the strings specified for the Metric name expression will be identical to those in the final metrics data sent to the Destination. See Raw Output below.

Metrics

Event Field NaLme	Metric Name Expression	Metric Type
bytes	`metric_name.bytes`	Gauge
packets	`metric_name.packets`	Gauge

Dimensions

Dimensions	
action interface_id dstaddr	

All specified dimension names must align with those from the original event. When you preview the Function's output, the metrics and dimensions will all have special highlighting to separate them from other fields. Additional highlighting is used to differentiate the metrics from the dimensions. (If one or more metrics/dimensions are not highlighted as expected, check the Function's configuration.)

Raw Output

```
metric_name.bytes:262256|g#action:REJECT,interface_id:eni-
02f03c2880e4aaa3,dstaddr:10.0.1.11

metric_name.packets:6556|g#action:REJECT,interface_id:eni-
02f03c2880e4aaa3,dstaddr:10.0.1.11
```

i Compatible Destinations

All text after the # symbol represents the dimensions as key-value pairs. In order for dimension data to be included in metrics, the Destination type cannot be standard **StatsD**.

Formatted Output

```
{
  "action": "REJECT",
  "interface_id": "eni-02f03c2880e4aaa3",
  "dstaddr": "10.0.1.11",
  "metric_name.bytes": 262256,
  "metric_name.packets": 6556,
}
```

Scenario B:

Assume that we want to extract some metrics from specific fields in PANOS logs, whose events have the following structure:

future_use_0,receive_time, serial_number, type, threat_content_type, future_use_1, generated_time, source_ip, destination_ip, nat_source_ip, nat_destination_ip, rule_name, source_user, destination_user, application, virtual_system, source_zone, destination_zone, inbound_interface, outbound_interface, log_action, future_use_2, session_id, repeat_count, source_port, destination_port, nat_source_port, nat_destination_port, flags, protocol, action, bytes, bytes_sent, bytes_received, packets, start_time, elapsed_time, category, future_use_3, sequence_number, action_flags, source_location, destination_location, future_use_4, packets_sent, packets_received, session_end_reason, device_group_hierarchy_level_1, device_group_hierarchy_level_2, device_group_hierarchy_level_3, device_group_hierarchy_level_4, virtual_system_name, device_name, action_source, source_vm_uuid, destination_vm_uuid, tunnel_id_imsi, monitor_tag_imei, parent_session_id, parent_start_time, tunnel_type, sctp_association_id, sctp_chunks, sctp_chunks_sent, sctp_chunks_received

For example:

```
Jan 10 10:19:15 DMZ-internal.nsa.gov 1,2019/01/10
10:19:15,001234567890002,TRAFFIC,drop,2304,2019/01/10
10:19:15,209.118.103.150,160.177.222.249,0.0.0.0,0.0.0.0,InternalServer,,not-applicable,vsys1,inside,z1-FW-Transit,ethernet1/2,,All traffic,2019/01/10
10:19:15,0,1,63712,443,0,0,0×0,udp,deny,60,60,0,1,2019/01/10
10:19:15,0,any,0,0123456789,0×0,Netherlands,10.0.0.0-10.255.255.255,0,1,0,policy-deny,0,0,0,0,0,TMZ-internal,from-policy,,,0,,0,,N/A,0,0,0,0,1202585d-b4d5-5b4c-aaa2-d80d77ba456e,0
```

Our goal is to use the four values of bytes_sent, bytes_received, packets_sent, and packets_received as metrics across these dimensions: destination_ip, inbound_interface, outbound_interface, and destination_port.

Here again, our Pipeline will need a Parser Function before the Publish Metrics Function.

Parser Function

Filter: Set as needed Operation mode: Extract

Type: Extended Log File Format (automatically set when specifying a Library)

Library: Palo Alto Traffic

Source: _raw

(No need to specify any other fields.)

Publish Metrics Function

Set up the Publish Metrics Function as follows.

Metrics

Event Field Name	Metric Name Expression	Metric Type
bytes_sent	metric.\${host}.bytes_sent	Counter
bytes_received	metric.\${host}.bytes_rcvd	Counter
packets_sent	metric.\${host}.pkts_sent	Counter
packets_received	metric.\${host}.pkts_rcvd	Counter

Added Dimensions

destination_ip , inbound_interface , outbound_interface , destination_port

Raw Output

 $\label{lem:metric.10.10.12.192.bytes_sent:60|c|\#destination_ip:160.177.222.249, inbound_interface:e:ethernet1/2, destination_port:443$

metric.10.10.12.192.bytes_rcvd:0|c|#destination_ip:160.177.222.249,inbound_interface:
ethernet1/2,destination_port:443

metric.10.10.12.192.pkts_sent:1|c|#destination_ip:160.177.222.249,inbound_interface:e
thernet1/2,destination_port:443

metric.10.10.12.192.pkts_rcvd:0|c|#destination_ip:160.177.222.249,inbound_interface:e
thernet1/2,destination_port:443

Here again, all text after the # symbol represents the dimensions as key-value pairs. (See the Compatible Destinations note above.) Unlike the first example, this example uses JavaScript expressions, which you can see evaluated in the raw output where the \${host} has been converted to 10.10.12.192.

Regex Extract

Description

The Regex Extract Function extracts fields using regex named groups. (In Splunk, these will be index-time fields). Fields that start with ___ (double underscore) are special fields in Cribl LogStream. They are ephemeral: they can be used by any Function downstream, but will not be added to events, and will not exit the Pipeline.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Regex: Regex literal. Must contain named capturing groups, e.g.: (?<foo>bar). Can contain special _NAME_N and _VALUE_N capturing groups, which extract both the name and value of a field, e.g.: (?<_NAME_0>[^\s=]+)=(?<_VALUE_0>[^\s]+). Defaults to empty. See Examples below.

Additional regex: Click + Add Regex to chain extra regex conditions.

Source field: Field on which to perform regex field extraction. Nested addressing is supported. Defaults to $_raw$.

Advanced Settings

Max exec: The maximum number of times to apply the Regex to the source field when the global flag is set, or when using <code>_NAME_N</code> and <code>_VALUE_N</code> capturing groups. Named capturing groups will always use a value of <code>1.Defaults</code> to <code>100</code>.

Field name format expression: JavaScript expression to format field names when _NAME_n and _VALUE_n capturing groups are used. E.g., to append XX to all field names, use: `\${name}_XX` (backticks are literal). If not specified, names will be sanitized using regex: /^[_0-9]+|[^a-zA-Z0-9_]+/g . The original field name is in the global _name .

Overwrite existing fields: Whether to overwrite existing event fields with extracted values. If set to No (the default), existing fields will be converted to an array. If toggled to Yes, Regex Extract will create array fields if applied multiple times, or if fields exist. (E.g., if src_ip is extracted in an input Pipeline where it is assigned a value of 10.1.2.2, and is also in a processing Pipeline with a value of 10.2.3.3, then the resulting field is ["10.1.2.2", "10.2.3.3"].)

Examples

Example 1

Assume a simple event that looks like this: metric1=23 metric2=42 dc=23 abc=xyz

Extract only the metric1 field:

Regex: metric1=(?<metric1>\d+)

Result: metric1:"23"

Example 2

Use the first line of the sample here:

https://github.com/weeb-cribl/cribl-

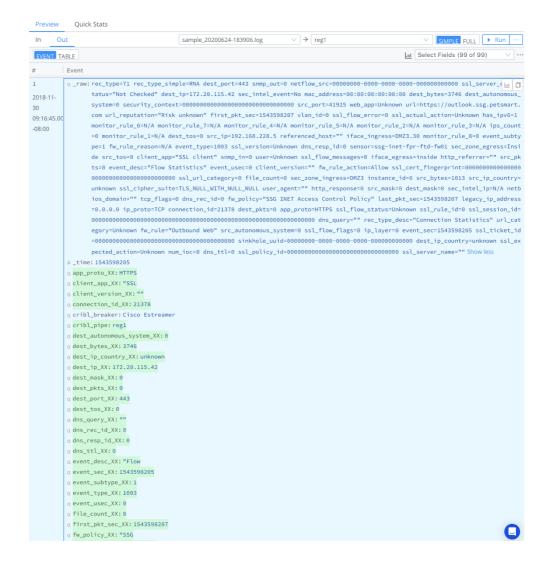
samples/blob/master/parser/functions/parser/cisco_estreamer.log

Extract all k=v pairs, and add an _XX suffix to the end of the extracted fields:

Regex: (?<_NAME_0>[^\s]+)=(?<_VALUE_0>[^\s]+)

Field Name Format Expression: `\${name}_XX`

Result:



Example 3

This example uses similar syntax as Example 2, but with a more complex event structure.

In the right **Sample Data** pane, click **Paste** and insert the following sample:

Sample Data

```
<134>1 2020-12-22T17:06:08Z CORP_INT_NLB CheckPoint 18160 - [action:"Accept"; conn_direction:"Ir
```

This event is from a CheckPoint Firewall CMA system. With this type of event structure, properly extracting each event field into a separate metadata field requires two-stage processing. So we'll use two Regex Extract Functions.

The first Regex Function splits the event to separate the actual data from the header information. We'll split after the CheckPoint 18160 string, by capturing everything between the [and]:

Regex; \[(?<__fields>.*)\]

Source: _raw

Next, add this second Regex Extract Function to extract all k=v pairs:

```
Regex: (?<_NAME_0>[^ :]+):(?<_VALUE_0>[^;]+);
```

Source: __fields

Results:

```
on:"Internal"; flags:"4606212"; ifdir:"inbound"; ifname:"bond2.1025"; logid:"0"; logui
                 d:"{0x5fe25889,0x0,0x80ad57cd,0xeb91c0c3}"; origin:"192.168.20.54"; originsicname:"CN=TST3
                 2-VSX0-FW-DC-01_tst302-shd,0=CORP-SEC-SHRD-CMA..t7xpcz"; sequencenum:"3"; time:"160865676
                 8"; version:"5"; __policy_id_tag:"product=VPN-1 & FireWall-1[db_tag={15E4B45A-663B-5B49-BD
                 59-CD9B9F21AA16};mgmt=SHRDFW01CON;date=1608236862;policy_name=TEST-SHRD-POL\]"; dst:"192.1
                 68.79.20"; log_delay:"1608656768"; layer_name:"TEST-SHRD-POL Security"; layer_uuid:"e914c2
                  f3-d7bd-4a77-8e7a-7a5e403447aa";\ match\_id:"1";\ parent\_rule:"0";\ rule\_action:"Accept";\ rule\_action:"Accept ";\ rule\_action:"Accept "
                  _uid:"001ab86d-d201-4b61-9b64-0fede1a9f059"; product:"VPN-1 & FireWall-1"; proto:"17"; s_p
                  ort:"45519"; service:"123"; service_id:"ntp-udp"; src:"192.168.79.22"; ] Showless
  # _time: 1608656768
 α action: "Accept"
 a conn_direction: "Internal"
 α cribl_breaker: Break on newlines
 α cribl_pipe: asfasfdasfd
a dst: "192.168.79.20"
α flags: "4606212"
a ifdir: "inbound"
α ifname: "bond2.1025"
α layer_name: "TEST-SHRD-POL Security"
α layer_uuid: "e914c2f3-d7bd-4a77-8e7a-7a5e403447aa"
α log_delay: "1608656768"
a logid: "0"
a loguid: "{0x5fe25889,0x0,0x80ad57cd,0xeb91c0c3}"
α match_id: "1"
α origin: "192.168.20.54"
α originsicname: "CN=TST32-VSX0-FW-DC-01_tst302-shd,0=CORP-SEC-SHRD-CMA..t7xpcz"
a policy_id_tag: "product=VPN-1 & FireWall-1[db_tag={15E4B45A-663B-5B49-BD59-CD9B9F21AA16}
α product: "VPN-1 & FireWall-1"
a proto: "17"
α rule_action: "Accept"
g rule uid: "001ab86d-d201-4b61-9b64-0fede1a9f059"
```

i For further examples, see Using Cribl to Analyze DNS Logs in Real Time – Part 2.

Redis

Description

The Redis Function interacts with Redis stores, setting and getting key-hash and key-value combinations. Redis' in-memory caching of these key pairs enables large lookup tables that would be cumbersome with a .CSV or binary lookup file.

You can use LogStream Collectors (e.g., a REST Collector) to retrieve reference data from desired endpoints, and then use this Function to store the data on Redis and retrieve it to enrich your production data. Note that LogStream does not cache the data returned from this Redis Function.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to

Result field: Name of the field in which to store the returned value. (Leave empty to discard the returned value.)

Command: Redis command to perform. Required. (A complete list of Redis commands is at: https://redis.io/commands.)

Key: A JavaScript expression to compute the value of the key to operate on. Can also be a constant, e.g.: username. This is a required field. Click the icon at right to open a validation modal.

Args: A JavaScript expression to compute arguments to the operation. Can return an array. Click the icon at right to open a validation modal.

Redis URL: Redis URL to connect to. The formatis: [redis[s]:]//[[user] [:password@]][host][:port][/db-number][?db=db-number[&password=bar[&option=value]]]

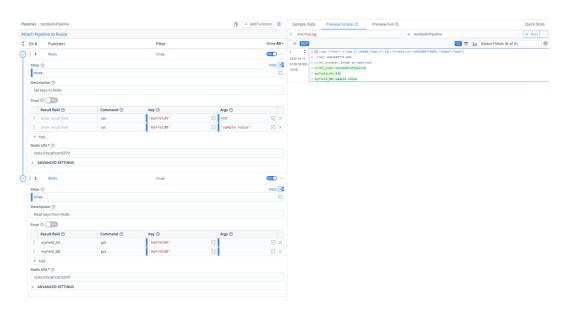
For example: redis://user:secret@localhost:6379/0?foo=bar&qux=baz

Advanced Settings

Max blocking time: Maximum amount of time (in seconds) before assuming that Redis is down and passing events through. Defaults to 60 seconds. Use 0 to disable timeouts.

Example

This Pipeline demonstrates the use of a pair of Redis Functions. The first Function sets two key-value pairs in Redis. The second Function their values, by key, into two corresponding new **Result fields**.



Redis set and get Functions

Redis Function #1

Description: Set keys to Redis

Command: set
Key: 'myFieldA'

Args: 420

Command: set
Key: 'myFieldB'
Args: 'sample value'

Redis Function #2

Description: Read keys from Redis

Result field: myField_AA

Command: get
Key: 'myFieldA'

Result field: myField_BB

Command: get
Key: 'myFieldB'

Regex Filter

Description

The Regex Filter Function filters out events based on regex matches.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Regex: Regex to test against. Defaults to empty.

Additional regex: Click + Add Regex to chain extra regex conditions.

Field: Name of the field to test against the regex. Defaults to <code>_raw</code> . Supports nested addressing.

Examples

See Regex Filtering for examples.

Rename

Description

The Rename Function is designed to change fields' names or reformat their names (e.g., by normalizing names to camelcase). You can use Rename to change specified fields (much like the Eval Function), or for bulk renaming based on a JavaScript expression (much like the Parser Function).

Compared to these alternatives, Rename offers a streamlined way to alter only field names, without other effects.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Optionally, enter a simple description of this step in the Pipeline. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Base fields: Enter one or more source field names to rename. If empty, rename will be performed on top-level fields.

Rename fields: Each row here is a key-value pair that defines how to rename fields. The current name is the key, and the new name is the value. Click + Add Field to add more rows.

- **Current name**: Original name of the field to rename. You must quote literal identifiers (non-alphanumeric characters such as spaces or hyphens).
- New name: New or reformatted name for the field. Here again, you must quote literals.

Renaming expression: An optional JavaScript expression (or literal) used to compute multiple fields' new names. This expression is evaluated against a {name, value} context, and the expression returns a value with which to rename fields.

You can use both Rename fields (to rename specified field names) and Renaming expression* (to globally rename fields) in the same Function. The Rename fields** strategy will execute first.

Example

Change the level field, and all fields that start with out, to all-uppercase.

Example event:

```
{"inEvents": 622,
    "level": "info",
    "outEvents": 311,
    "outBytes": 144030,
    "activeCxn": 0,
    "openCxn": 0,
    "closeCxn": 0,
    "activeEP": 105,
    "blockedEP": 0
}
```

Rename Fields:

```
Current Name: level
New Name: LEVEL
Renaming Expression: name.startsWith('out') ? name.toUpperCase():
name
```

Event after Rename:

```
{"inEvents": 622,
  "LEVEL": "info",
  "OUTEVENTS": 311,
  "OUTBYTES": 144030,
  "activeCxn": 0,
  "openCxn": 0,
  "closeCxn": 0,
  "activeEP": 105,
  "blockedEP": 0
}
```

Rollup Metrics

Description

The Rollup Metrics Function merges/rolls up frequently generated incoming metrics into more manageable time windows.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Optional description of this Function's purpose in this Pipeline. Defaults to empty.

Final: If toggled to Yes, stops data from being fed to downstream Functions. Defaults to No.

Dimensions: List of data dimensions across which to perform rollups. Supports wildcards. Defaults to * wildcard, meaning all original dimensions.

Time window: The time span over which to roll up (aggregate) metrics. Must be a valid time string (e.g., 10s). Must match pattern: \d+[sm]\$.

⚠ With high-cardinality data, beware of setting long time windows. Doing can cause high memory consumption and/or lost data, because memory is flushed upon restarts and redeployments.

Gauge update: The operation to use when rolling up gauge metrics. Defaults to Last; other options are Maximum, Minimum, or Average.

Examples

Scenario A:

Assume that you have metrics coming in at a rate that is too high. For example, LogStream's internal metrics come in at a 2s interval.

To roll up these metrics to 1-minute granularity, you would set up the Rollup Metrics Function with a **Time Window** value of 60s.

Scenario B:

Assume that you have metrics coming up with multiple dimensions – e.g. host, source, data_center, and application. You want to aggregate these metrics to eliminate some dimensions.

Here, you would configure Rollup Metrics Function with a **Time Window** value that matches the metrics' generation – e.g., 10s . In the **Dimensions** field, you would remove the default \star wildcard, and would specify only the dimensions you want to keep – e.g.: host , data_center .

Sampling

Description

The Sampling Function filters out events, based on an expression and a sampling rate.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

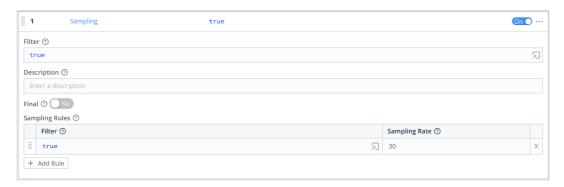
Final: If true, stops data from being fed to downstream Functions. Defaults to No .

Sampling rules: Events matching these rules will be sampled at the rates you specify:

- **Filter**: Filter expression matching events to be sampled. Use true to match all.
- Sampling rate: Enter an integer N. (Defaults to 1.) Sampling will pick 1/N events matching this rule.

How It Works

Setting this Function's **Sampling rate** to 30 would mean that only 1 of every 30 events would be kept.



Let's assume that we save this setting, and then capture data from a datagen Source by selecting Preview > Start a Capture > Capture. In the Capture Sample Data modal, select: 100 seconds, 100 events, and As they come in. Then start the capture, and Save as Sample File.

Next, in the **Preview** pane, click **Simple** beside the new file's name. If you then click the **Basic Statistics** (chart) button, you should see that we've kept about 4 of the original 100 events, or close to 1 in 30.

	Full Event Length ⑦	Number of Fields ⑦	Number of Events ⑦
IN	28.82KB	41	100
OUT	1.42KB	38	4
DIFF	↓ -95.08%	↓ -7.32%	↓ -96.00%

Examples

See Sampling for examples.

Serialize

Description

Use the Serialize Function to serialize an event's content into a predefined format.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Type: Data output format. Defaults to CSV.

Library: Browse Parser/Formatter library.

Fields to serialize: Required for CSV, ELFF, and CLF Types. (All other formats support wildcard field lists.)

Source field: Field containing the object to serialize. Leave blank to serialize top-level event fields.

Destination field: Field to serialize the data into. Defaults to _raw .

Examples

Scenario A: JSON to CSV

```
Assume a simple event that looks like this: {"time":"2019-08-
25T14:19:10.240Z", "channel": "input", "level": "info", "message": "initializing input", "type": "kafka"}
```

We want to serialize these fields: _time , channel , level , and type into a single string, in CSV format, stored in a new destination field called test .

To properly extract the key-value pairs from this event structure, we'll use a built-in Event Breaker:

- 1. Copy the above sample event to your clipboard.
- 2. In the **Preview** pane, select **Paste a Sample**, and paste in the sample event.
- 3. Under **Select Event Breaker**, choose **ndjson** (newline-delimited JSON), and click **Save as a Sample File**.

Now you're ready to configure the Serialize Function, using the settings below:

Type: CSV

Fields to Serialize: _time channel level type

Destination Field: test Source Field: [leave empty]

Result: test: 1566742750.24, input, info, kafka

In the new test field, you now see the time, channel, level, and type keys extracted as top-level fields.

Scenario B: CSV to JSON

Let's assume that a merchant wants to extract a subset of each customer order, to aggregate anonymized order statistics across their customer base. The transaction data is originally in CSV format, but the statistical data must be in JSON.

Here's a CSV header (which we don't want to process), followed by a row that represents one order:

```
orderID, custName, street, city, state, zip
20200622102822, john smith, 100 Main St., Anytown, AK, 99911
```

To convert to JSON, we'll need to first parse each field from the CSV to a manipulable field in the Pipeline, which the Serialize Function will be able to reference. In this example, the new manipulable field is message.

Use the Parser Function:

Filter: true

Operation mode: Extract

Type: CSV

Source field: _raw

Destination field: message

List of fields: orderID custName street city state zip

Now use the Serialize Function:

Filter: true Type: JSON

Fields to serialize: city state

Source field: message

Destination field: orderStats

Suppress

Description

The Suppress Function suppresses events over a time period, based on evaluating a key expression.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No .

Key expression: Suppression key expression used to uniquely identify events to suppress. For example, `\${ip}:\${port}` will use the fields ip and port from each event to generate the key.

Number to allow: The number of events to allow per time period. Defaults to 1.

Suppression period (sec): The number of seconds to suppress events after 'Number to allow' events are received. Defaults to 300.

Drop suppressed events: Specifies if suppressed events should be dropped, or just tagged with suppress=1. Defaults to Yes, meaning drop.

Advanced Settings

Maximum cache size: The maximum number of keys that can be cached before idle entries are removed. Before changing the default 50000, contact Cribl Support to understand the implications.

Suppression period timeout: The number of suppression periods of inactivity before a cache entry is considered idle. This defines a multiple of the **Suppression period (sec)** value. Before changing the default 2, contact Cribl Support to understand the implications.

Num events to trigger cache clean-up: Check cache for idle sessions every N events when cache size exceeds the **Maximum cache size**. Before changing the default 10000, contact Cribl Support to understand the implications.

Examples

In the examples below, **Filter** is the Function-level Filter expression:

1. Suppress by the value of the host field:

Filter: true

Key expression: host Number to allow: 1

Suppression period (sec): 30

Using a datagen sample as a source, generate at least 100 events over 2 minutes.

Result: One event per unique host value will be allowed in every 30s. Events without a host field will **not** be suppressed.

2. Suppress by the value of the host and port tuple:

Filter: true

Key expression: `\${host}:\${port}`

Number to allow: 1

Suppression period (sec): 300

Result: One event per unique host: port tuple value will be allowed in every 300s.

- ⚠ Suppression will **also** apply to events without a host or a port field. The reason is that if field is not present, `\${field}` results in the literal undefined.
- 3. To **guarantee** that suppression applies **only** to events with host and port, check for their presence using a Filter:

Filter: host≠undefined & port≠undefined

Key expression: `\${host}:\${port}`

Number to allow: 1

Suppression period (sec): 300

4. Decorate events that qualify for suppression:

Filter: true

Key expression: `\${host}:\${port}`

Number to allow: 1

Suppression period (sec): 300 Drop suppressed events: No

Result: No events will be suppressed. But all qualifying events will gain an added field suppress=1, which can be used downstream to further transform these events.

Tee

Description

The Tee Function tees events out to a command of choice, via stdin. The output is one JSON-formatted event per line. You can send the events to (for example) a local file on the LogStream worker. This can be useful in verifying the data being processed in a Pipeline.

The Filesystem/NFS Destination offers similar capability, but only after the data leaves the Pipeline. Tee, by comparison, can be inserted at any point in the Pipeline.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Command: Command to execute and receive events (via stdin) – one JSON-formatted event per line.

Args: Click + Add Arg to supply arguments to the command.

Restart on exit: Restart the process if it exits and/or we fail to write to it. Defaults to Yes.

Environment variables: Environment variables to set or overwrite. Click + Add Variable to add key/value pairs.

Communication Protocol

Data is passed to the command through its stdin, using the following protocol:

- First line: Metadata serialized in JSON, containing the following fields:
 - format: Serialization format for event. Defaults to JSON.
 - conf: Full Function configuration.
- Remaining: Payload.

Examples

Assume that we are parsing PANOS Traffic logs, and want to see how they look at a particular step in the processing Pipeline We'll assume that the Parser Function is already in place, so we'll insert the

Tee Function at any (arbitrary) later point in the Pipeline.

Scenario A:

The Tee Function itself requires only that we define the **Command** field. In this particular example, that **Command** will be tee itself.

We've also clicked + Add Arg, to specify a local output file in the resulting Args field. (A file path would normally be the first argument to a tee command executed from the command line. The LogStream user must have write permission on the specified file path.)

Command: tee

Args: /opt/cribl/foo.log

In this first scenario, assume that we have the Parser configured to parse, but not keep any fields. After changes are deployed and PANOS logs are received, if we tail foo.log, we'd see the following:

```
Line 1: {"format":"json","conf":{"restartOnExit":true,"env":
{},"command":"tee","args":["/opt/cribl/foo.log"]}

Line 2: {"_raw":"Oct 09 10:19:15 DMZ-internal.nsa.gov 1,2019/10/09

10:19:15,001234567890002,TRAFFIC,drop,2304,2019/10/09

10:19:15,209.118.103.150,160.177.222.249,0.0.0.0,0.0.0,InternalServer,,,not-applicable,vsys1,inside,z1-FW-Transit,ethernet1/2,,All traffic,2019/10/09

10:19:15,0,1,63712,443,0,0,0×0,udp,deny,60,60,0,1,2019/10/09

10:19:15,0,any,0,0123456789,0×0,Netherlands,10.0.0.0-10.255.255.255,0,1,0,policydeny,0,0,0,0,0,0,Titlernal,from-policy,,,0,,0,,N/A,0,0,0,0,1202585d-b4d5-5b4c-aaa2-d80d77ba456e,0","_time":1593185574.663,"host":"127.0.0.1"}
```

In Line 2 above, note that the _raw field makes up most of the contents, with only the _time and host fields added.

Scenario B:

Assume that we use the Tee Function, using the same **Command** and arguments, but we've modified the the Parser Function to retain five fields: receive_time, source_port, destination_port bytes_received, and packets_received.

This time, if we tail foolog, we'll see something like the following. If you compare this output to the previous output example, you'll notice the five fields appended to this event:

```
Line 3: {"_raw":"Oct 09 10:19:15 DMZ-internal.nsa.gov 1,2019/10/09 10:19:15,001234567890002,TRAFFIC,drop,2304,2019/10/09 10:19:15,209.118.103.150,160.177.222.249,0.0.0.0,0.0.0.0,InternalServer,,,not-applicable,vsys1,inside,z1-FW-Transit,ethernet1/2,,All traffic,2019/10/09 10:19:15,0,1,63712,443,0,0,0×0,udp,deny,60,60,0,1,2019/10/09 10:19:15,0,any,0,0123456789,0×0,Netherlands,10.0.0.0-10.255.255.255,0,1,0,policy-deny,0,0,0,0,TDMZ-internal,from-policy,,,0,,0,,N/A,0,0,0,0,1202585d-b4d5-5b4c-aaa2-
```

```
d80d77ba456e,0","_time":1593185606.965,"host":"127.0.0.1","receive_time":"2019/10/09
10:19:15","source_port":"63712","destination_port":"443","bytes_received":"0","packet
s_received":"0"}
```

i In this Function's **Command** field, you can specify commands other than tee itself. For example: By using nc as the command, and specifying localhost and a port number (as two separate arguments), you'll see event data being received via nc on the specified port.

Trim Timestamp

Description

The Trim Timestamp Function removes timestamp patterns from events, and (optionally) stores them in a specified field.

This Function looks for a timestamp pattern that exists between the characters indicated by numeric timestartpos and timeendpos fields. It removes timestartpos and timeendpos along with the timestamp pattern.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description about this step in the Pipeline. Defaults to empty.

Final: If true, stops data from being fed to the downstream Functions. Defaults to No.

Field name: Name of field in which to save the timestamp. (If empty, timestamp will not be saved to a field.)

Example

Remove the timestamp pattern (indicated by timestartpos and timeendpos) from _raw , and stash it in a field called time_field .

Example event before:

```
{"_raw": "Event [Event=UpdateBillingProvQuote, timestamp=1581426279, properties={JMSCorrelation:
"timestartpos":0,
"timestartpos":23
}```
**Field Name**: `time_field`
**Example Event after:**
```

```
{"_raw": "2020-05-22 16:32:11,359 Event [Event=UpdateBillingProvQuote, timestamp=1581426279, properties={JMSCorrelationID=NA, JMSMessageID=ID:ESP-PD.D2BB2D95F857B:FA323D61, orderType=RatePlanFeatureChange, quotePriority=NORMAL}", "time_field":"2020-05-22 16:32:11,359" }```
```

Unroll

Description

The Unroll Function accepts an array field – or an expression to evaluate an array field – and breaks/unrolls the array into individual events.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Source field expression: Field in which to find/calculate the array to unroll. E.g.: $_{raw}$, $_{raw.split(/\n/)}$. Defaults to $_{raw}$.

Destination field: Field (within the destination event) in which to place the unrolled value. Defaults to raw .

Example

Assume we want to break/unroll each line of this event:

em ·
. •

Settings

Source field expression: _raw.split(/\n/)

i The split() JavaScript method breaks _raw into an ordered set of substrings/values, puts these values into an array, and returns the array.

Destination field: _raw

Resulting Events

Event 1: USER	PID %CPU %MEM	VSZ RSS TTY	STAT START	TIME COMMAND
Event 2: root	1 0.0 0.5	38000 5356 ?	Ss 2018	2:02 /lib/systemd/systemdsystem ·
Event 3: root	2 0.0 0.0	0 0?	S 2018	0:00 [kthreadd]
Event 4: root	3 0.0 0.0	0 0?	S 2018	1:51 [ksoftirqd/0]
Event 5: root	5 0.0 0.0	0 0?	S< 2018	0:00 [kworker/0:0H]
Event 6: root	7 0.0 0.0	0 0?	S 2018	3:55 [rcu_sched]
Event 7:	8 0.0 0.0	0 0 ?	S 2018	0:00 [rcu bh]

XML Unroll

Description

The XML Unroll Function accepts a proper XML event with a set of elements, and converts the elements into individual events.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No .

Unroll elements regex: Path to the array to unroll. E.g.:
 ^root\.child\.ElementToUnroll\$

Copy elements regex: Regex matching elements to copy into each unrolled event.

E.g.: ^root\.(childA|childB|childC)\$

Unroll index field: LogStream will add a field with this name, containing the 0-based index at which the element was located within the event. In Splunk, this will be an index-time field. Supports nested addressing. Name defaults to unroll_idx.

Pretty print: Whether to pretty print the output XML.

Examples

Assume that the following sample is ingested as a single event:

```
<state>NY</state>
        <city>New York</city>
    </Child>
    <Child>
        <state>NJ</state>
        <city>Edgewater</city>
    </Child>
    <Child>
       <state>CA</state>
        <city>Oakland</city>
    </Child>
    <Child>
        <state>CA</state>
        <city>San Francisco</city>
    </Child>
</Parent>
```

i If you insert this sample using Preview > Add a Sample > Paste a Sample, adjust Event Breaker settings to add the sample as a single event. One way to do this is to add a regex Event Breaker that (by design) will not match anything present in the sample. For example: /[\n\r]+donotbreak(?!\s)/. As of LogStream 2.3, you can also use the built-in Do Not Break Ruleset.

Set up the XML Unroll Function using these settings:

Unroll elements regex: ^Parent\.Child\$

Copy elements regex: ^Parent\.(myID|branchLocation)\$

Output 4 Events:

```
Resulting Events
# Event 1
<?xml version="1.0"?>
<Child>
  <myID>123456</myID>
  <branchLocation>US/branchLocation>
  <state>NY</state>
  <city>New York</city>
</Child>
# Event 2
<?xml version="1.0"?>
<Child>
  <myID>123456</myID>
  <branchLocation>US/branchLocation>
  <state>NJ</state>
  <city>Edgewater
</Child>
```

```
# Event 3
<?xml version="1.0"?>
<Child>
 <myID>123456</myID>
 <branchLocation>US/branchLocation>
 <state>CA</state>
 <city>Oakland</city>
</Child>
# Event 4
<?xml version="1.0"?>
<Child>
 <myID>123456</myID>
 <branchLocation>US</branchLocation>
 <state>CA</state>
 <city>San Francisco</city>
</Child>
```

Prometheus Publisher (beta)

Description

The Prometheus Publisher Function allows for metrics to be published to a Prometheus-compatible metrics endpoint. These can be upstream metrics received by LogStream, or metrics derived from the output of LogStream's Publish Metrics or Aggregation Functions. A Prometheus instance is responsible for collecting the metrics at that endpoint, and for performing its own processing of the metric data.

In the current LogStream version, the endpoint is: http://<worker_node_IP>:<apiport>/metrics. Within LogStream, that endpoint redirects from http://<worker_node_IP>:9000/metrics to http://<worker_node_IP>:9000/api/v1/metrics.

⚠ If used, this Function must follow any Publish Metrics or Aggregations Functions within the same Pipeline. This is to ensure that any data not originating from a metrics input is transformed into metrics format.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Fields to publish: Wildcard list of fields to publish to the Prometheus endpoint.

Advanced Settings

Batch write interval: How often, in milliseconds, the contents should be published. Defaults to 5000.

Passthrough mode: If set to No (the default), overrides the Final setting, and suppresses output to downstream Functions' Destinations. Toggle to Yes to allow events to flow to consumers beyond the Prometheus endpoint. In effect, when previewing the pipeline output what you'll see is your event fields will have strikethrough font applied to them. This does not mean the Prometheus function is not matching your events but rather indicative of the Passthrough being disabled.

Update mode: On the default No setting, suppresses output to downstream Functions' Destinations. (This overrides the Final setting.) Toggle to Yes to allow events to flow to consumers beyond the

Example

This example uses the same PANOS sample data as the Publish Metrics Function, and is similarly preceded in a Pipeline by a Parser Function that extracts fields from the PANOS log.

Filter: Set as appropriate.

Fields to publish: Set as appropriate. We'll use the default of * for this example. **Advanced settings**: Accept defaults.

After committing and deploying changes, you should be able to use a <code>curl</code> command (-L needed to follow the redirect mentioned above) to verify that metrics are being published, just a few seconds after data is ingested on an idle system.

```
curl output

$ curl -L http://<worker_node_IP>:9000/metrics

# TYPE perf_192_168_1_248_bytes_sent counter
metric_192_168_1_248_bytes_sent {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"

# TYPE perf_192_168_1_248_bytes_rcvd counter
metric_192_168_1_248_bytes_rcvd {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"

# TYPE perf_192_168_1_248_pkts_sent counter
metric_192_168_1_248_pkts_sent {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"

# TYPE perf_192_168_1_248_pkts_rcvd counter
metric_192_168_1_248_pkts_rcvd counter
metric_192_168_1_248_pkts_rcvd {destination_ip="160.177.222.249",inbound_interface="ethernet1/2"
```

Now, we need to have Prometheus scrape the metrics. In this very basic example, you can add the target endpoint to the prometheus.yml file, under the scrape_configs -> static_configs section. Specify the endpoint in IP:port syntax, because Prometheus assumes (and requires) /metrics for all endpoints.

Restart Prometheus. Within just a few seconds, you should be able to use its query interface to retrieve metrics published by LogStream.

Reverse DNS (deprecated)

Description

The Reverse DNS Function resolves hostnames from a numeric IP address, using a reverse DNS lookup.

⚠ This Function is deprecated. Use the DNS Lookup Function's reverse lookup feature instead.

Usage

Filter: Filter expression (JS) that selects data to be fed through the Function. Defaults to true, meaning that all events will be evaluated.

Description: Simple description of this Function. Defaults to empty.

Final: If true, stops data from being fed to downstream Functions. Defaults to No.

Lookup Fields

Lookup field name: Name of the field containing the IP address to look up.

⚠ If the field value is not in IPv4 or IPv6 format, the lookup is skipped.

Output field name: Name of the field in which to add the resolved hostname. Leave blank to overwrite the lookup field.

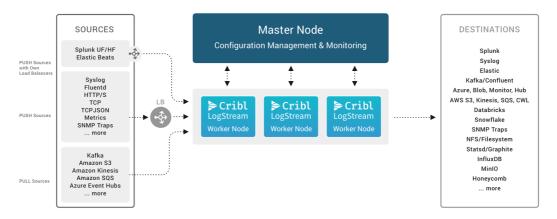
Reload period (minutes): How often to refresh the DNS cache. Use 0 to disable refreshes. Defaults to 60 minutes.

Example

Lookup field name: dest_ip Output field name: dest_host **Result**: See the dest_ip field, and the newly created dest_host field, in the events.

Sources

Cribl LogStream can receive data from various Sources, including Splunk, HTTP, Elastic Beats, Kinesis, Kafka, TCP JSON, and many others.



Push and Pull Sources

PUSH Sources

Supported data Sources that send to Cribl LogStream:

- Splunk TCP
- Splunk HEC
- Syslog
- Elasticsearch API
- TCP JSON
- TCP Raw
- HTTP/S
- Raw HTTP/S
- Kinesis Firehose
- SNMP Trap
- Metrics

Data from these Sources is normally sent to a set of LogStream Workers through a loadbalancer. Some Sources, such as Splunk forwarders, have native loadbalancing capabilities, so you should point these directly at LogStream.

PULL Sources

Supported Sources that Cribl LogStream fetches data from:

- Kafka
- Kinesis Streams
- Azure Event Hubs
- SQS
- S3
- Office 365 Services
- Office 365 Activity
- Prometheus

Internal Sources

Sources that are internal to Cribl LogStream:

- Datagens
- Cribl Internal

Configuring and Managing Sources

For each Source *type*, you can create multiple definitions, depending on your requirements.

To configure Sources, select **Data > Sources**, select the desired type from the tiles or the left menu, and then click **+ Add New**.

Preconfigured Sources

To accelerate your setup, LogStream ships with several common Sources configured for typical listening ports, but not switched on. Open, clone (if desired), modify, and enable any of these preconfigured Sources to get started quickly:

- Syslog TCP Port 9514, UDP Port 9514
- Splunk TCP Port 9997
- Splunk HEC Port 8088
- TCP JSON Port 10070
- TCP Port 10060

- HTTP Port 10080
- Elasticsearch API Port 9200
- SNMP Trap Port 9162
- Cribl Internal > CriblLogs Internal
- Cribl Internal > CriblMetrics Internal

Backpressure Behavior

On the Destination side, you can configure how each LogStream output will respond to a **backpressure** situation – a situation where its in-memory queue is overwhelmed with data.

All Destinations default to **Block** mode, in which they will refuse to accept new data until the downstream receiver is ready. Here, LogStream will backpropagate block signals through the Source, all the way back to the sender (if it supports backpressure, too).

All Destinations also support **Drop** mode, which will simply discard new events until the receiver is ready.

Several Destinations also support a **Persistent Queue** option to minimize data loss. Here, the Destination will write data to disk until the receiver is ready. Then it will drain the disk-buffered data in FIFO (first in, first out) order. See Persistent Queues for details about all three modes, and about **Persistent Queue** support.

Other BackPressure Options

The S3 Source provides a configurable **Advanced Settings > Socket timeout** option, to prevent data loss (partial downloading of logs) during backpressure delays.

Diagnosing Backpressure Errors

When backpressure affects HTTP Sources (Splunk HEC, HTTP/S, Raw HTTP/S, and Kinesis Firehose), LogStream internal logs will show a 503 error code.

Splunk TCP

Cribl LogStream supports receiving Splunk data from Universal or Heavy Forwarders.

Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive Splunk TCP Data

Select **Data > Sources**, then select **Splunk > Splunk TCP** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Splunk TCP > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a Splunk TCP Source preconfigured to listen on Port 9997. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Splunk Source definition.

Address: Enter hostname/IP to listen for Splunk data. E.g., localhost or 0.0.0.0.

Port: Enter port number.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection. Defaults to .* (i.e., all IPs).

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

* CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

Common name: Regex matching peer certificate subject common names allowed to connect.
 Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol V1 or V2.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

• __inputId

Configuring a Splunk Forwarder

To configure a Splunk forwarder (UF, HF) use the following outputs.conf stanzas:

```
.../outputs.conf

[tcpout]
disabled = false
defaultGroup = cribl, <optional_clone_target_group>,

[tcpout:cribl]
server = [<cribl_ip>|<cribl_host>]:<port>, [<cribl_ip>|<cribl_host>]:<port>, ...
sendCookedData=true
# As of Splunk 6.5, using forceTimebasedAutoLB is no longer recommended. Ensure this is left at
# forceTimebasedAutoLB = false
```

Splunk HEC

Cribl LogStream supports receiving data over HTTP/S using the Splunk HEC (HTTP Event Collector).

i Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive Data over Splunk HEC

Select Data > Sources, then select Splunk > HEC from the Data Sources page's tiles or left menu. Click Add New to open the HEC > New Source modal, which provides the fields outlined below.

LogStream ships with a Splunk HEC Source preconfigured to listen on Port 8088. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Splunk HEC Source definition.

Address: Enter the hostname/IP on which to listen for HTTP(S) data. (E.g., localhost or 0.0.0.0.)

Port: Enter the port number.

Splunk HEC endpoint: Absolute path on which to listen for the Splunk HTTP Event Collector API requests. This input supports the /event and /raw endpoints. Defaults to /services/collector.

Allowed Indexes: List the values allowed in the HEC event index field. Allows wildcards. Leave blank to skip validation.

Splunk HEC acks: Whether to enable Splunk HEC acknowledgments. Defaults to No . Some sources may require HEC acks to be enabled and, as a result, may keep TCP connections open while waiting for an ack. This behavior can exhaust available file descriptors. Cribl does not maintain a comprehensive list of such sources. Refer to your source's documentation for more information.

Auth Tokens

Token: Shared secret to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access **will be permitted**.

Description: Optional description for this token.

Fields: Fields (metadata) to add to events referencing this token. Each field is a Name/Value pair.

These fields may be overridden by fields added at the request level.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

Common name: Regex matching peer certificate subject common names allowed to connect.
 Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol V1 or V2.

Processing Settings

Event Breakers

This section defines event breaking rulesets that will be applied, in order, on the /raw endpoint.

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

These fields may be overridden by fields added at the token or request level.

Pre-Processing Pipeline

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __hecToken

Format and Endpoint Example

- Configure Cribl LogStream to listen on port 10080 with an authToken of myToken42.
- Send a payload to your Cribl LogStream receiver.

Note: Token specification can be either Splunk <token> or <token> .

Splunk HEC Event Endpoint

```
curl -k http://<myCriblHost>:10080/services/collector/event -H 'Authorization: myToken42' -d '{'
curl -k http://<myCriblHost>:10080/services/collector -H 'Authorization: myToken42' -d '{"event'

# Multiple Events
curl -k http://<myCriblHost>:10080/services/collector -H 'Authorization: myToken42' -d '{"event'
```

Syslog

Cribl LogStream supports receiving of data over syslog.

Type: Push | TLS Support: YES | Event Breaker Support: No
This Syslog Source supports RFC 3164 and RFC 5424.

Configuring Cribl LogStream to Receive Data over Syslog

Select **Data > Sources**, then select **Syslog** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Syslog > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a Syslog Source preconfigured to listen for both UDP and TCP traffic on Port 9514. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Syslog Source definition.

Address: Enter the hostname/IP on which to listen for data., E.g. localhost or 0.0.0.0.

UDP port: Enter the UDP port number to listen on. Not required if listening on TCP.

TCP port: Enter the TCP port number to listen on. Not required if listening on UDP.

Fields to keep: List of fields from source data to retain and pass through.

Supports wildcards. Defaults to * wildcard, meaning keep all fields. Fields not specified here (by wildcard or specific name) will be removed from the event.

TLS Settings (TCP Only)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

IP whitelist regex: Regex matching IP addresses that are allowed to send data. Defaults to .* (i.e., all IPs).

Max buffer size (events): Maximum number of events to buffer when downstream is blocking. The buffer is only in memory. (This setting is applicable only to UDP syslog.)

Default timezone: Timezone to assign to timestamps that omit timezone info. Accept the default Local value, or use the drop-down list to select a specific timezone by city name or GMT/UTC offset.

Single msg per UDP: Whether to treat UDP packet data received as a full syslog message. Defaults to No. (I.e., newlines in the packet will be treated as event delimiters.)

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but are accessible and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __srcIpPort
- _syslogFail: true for data that fails RFC 3164/5424 validation as syslog format.

Elasticsearch API

Cribl LogStream supports receiving data over HTTP/S using the Elasticsearch Bulk API. (See the Configuring Filebeat example below.)

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring LogStream to Receive Data over HTTP(S), Using the Elasticsearch Bulk API Protocol

Select Data > Sources, then select Elasticsearch API from the Data Sources page's tiles or left menu. Click Add New to open the Elasticsearch API > New Source modal, which provides the fields outlined below.

☐ LogStream ships with an Elasticsearch API Source preconfigured to listen on Port 9200. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Elasticsearch Source definition.

Address: Enter the hostname/IP on which to listen for Elasticsearch data. (E.g., localhost or 0.0.0.0.)

Port: Enter the port number.

Auth tokens: Shared secrets to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access **will be permitted**.

Elasticsearch API endpoint (for Bulk API): Absolute path on which to listen for Elasticsearch API requests. Currently, the only supported option is the default /elastic, which LogStream expands as /elastic/_bulk. Other entries are faked as success. Use an empty string to disable.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Field Normalization

The Elasticsearch API input normalizes the following fields:

- atimestamp becomes _time atmillisecond resolution.
- host is set to host.name.
- Original object host is stored in __host .

The Elasticsearch Destination does the reverse, and it also recognizes the presence of __host .

Internal Settings

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __id
- __type
- __index
- __host

Configuring Filebeat

To set up Filebeat to send data to LogStream, use its Elasticsearch output. If an Auth Token is configured here, add it in Filebeat configuration under output.elasticsearch.headers, as in this example:

```
...filebeat.yml

output.elasticsearch:
    # Array of hosts to connect to.
    hosts: ["http://<LOGSTREAM_HOST>:9200/elastic"]
```

output.elasticsearch.headers:
 Authorization: "myToken42"

TCP JSON

Cribl LogStream supports receiving of data over TCP in JSON format (see protocol below).

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring Cribl LogStream to Receive TCP JSON Data

Select **Data > Sources**, then select **TCP JSON** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **TCP JSON > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a TCP JSON Source preconfigured to listen on Port 10070. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this TCP JSON Source definition.

Address: Enter hostname/IP to listen for TCP JSON data. E.g., localhost or 0.0.0.0.

Port: Enter the port number to listen on.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection. Defaults to .* (i.e., all IPs).

Shared secret (authToken): Shared secret to be provided by any client (in authToken header field). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

Common name: Regex matching peer certificate subject common names allowed to connect.
 Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

• Enable proxy protocol: Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

inputId

Format

LogStream expects TCP JSON events in newline-delimited JSON format:

1. A header line. Can be empty – e.g., {} . If **authToken** is enabled (see above) it should be included here as a field called authToken . When authToken is **not** set, the header line is **optional**. In this case, the first line will be treated as an event if does not look like a header record.

In addition, if events need to contain common fields, they can be included here under fields . In the example below, region and AZ will be automatically added to all events.

2. A JSON event/record per line.

Sample TCP JSON Events

```
{"authToken":"myToken42", "fields": {"region": "us-east-1", "AZ":"az1"}}

{"_raw":"this is a sample event ", "host":"myHost", "source":"mySource", "fieldA":"valueA", "fieldst":"myOtherHost", "source":"myOtherSource", "_raw": "{\"message\":\"Something informative | myOtherSource", "_raw": "fields": "source": "myOtherSource", "_raw": "fields": "source": "sour
```

TCP JSON Field Mapping to Splunk

If a TCP JSON Source is routed to a Splunk destination, fields within the JSON payload are mapped to Splunk fields. Fields that do not have corresponding (native) Splunk fields become index-time fields. For example, let's assume we have a TCP JSON event as below:

```
{"_time":1541280341, "host":"myHost", "source":"mySource", "_raw":"this is a sample
event ", "fieldA":"valueA"}
```

Here, _time , host , and source become their corresponding fields in Splunk. The value of _raw becomes the actual body of the event, and fieldA becomes an index-time field (fieldA::`valueA``).

Example

- 1. Configure Cribl LogStream to listen on port 10001 for TCP JSON. Set authToken to myToken42.
- 2. Create a file called test.json with the payload above.
- 3. Send it over to your Cribl LogStream host: cat test.json | nc <myCriblHost> 10001

TCP (RAW)

Cribl LogStream supports receiving of data over TCP. (See examples and header options below.)

i Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive TCP Data

Select **Data > Sources**, then select **TCP** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **TCP > New Source** modal, which provides the fields outlined below.

☐ LogStream ships with a TCP Source preconfigured to listen on Port 10060. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this TCP Source definition.

Address: Enter hostname/IP to listen for raw TCP data. E.g., localhost or 0.0.0.0.

Port: Enter port number.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection. Defaults to .* (i.e., all IPs).

Enable Header: Toggle to Yes to indicate that client will pass a header record with every new connection. The header can contain an authToken, and an object with a list of fields and values to add to every event. These fields can be used to simplify Event Breaker selection, routing, etc. Header format:

```
{ "authToken" : "myToken", "fields": { "field1": "value1", "field2": "value2" }}.
```

• **Shared secret (authToken)**: Shared secret to be provided by any client (in authToken header field). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Custom Command

In this section, you can pass the data from this input to an external command for processing before the data continues downstream.

Enabled: Defaults to No. When toggled to Yes:

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + Add Argument to add each argument for the command. You can drag arguments vertically to resequence them.

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

• Name: Field name.

• Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol V1 or V2.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions can use them to make processing decisions.

Fields accessible for this Source:

- __inputId
- __srcIpPort
- __channel

TCP Source Example

Every new TCP connection may contain an **optional** header line, with an authToken and a list of fields and values to add to every event.

```
SamplerawTCPtest
{"authToken":"myToken42", "fields": {"region": "us-east-1", "AZ":"az1"}}
this is event number 1
this is event number 2
```

Enabling the Example

- Configure LogStream to listen on port 7777 for raw TCP. Set authToken to myToken42.
- 2. Create a file called test.raw, with the payload above.
- 3. Send it over to your Cribl LogStream host, using this command: cat
 test.raw | nc <myCriblHost> 7777

HTTP/S (Bulk API)

Cribl LogStream supports receiving data over HTTP/S using the Cribl Bulk API, Splunk HEC, or Elastic Bulk API.

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data over HTTP(S)

Select Data > Sources, then select HTTP from the Data Sources page's tiles or left menu. Click Add New to open the HTTP > New Source modal, which provides the fields outlined below.

☐ LogStream ships with an HTTP Source preconfigured to listen on Port 10080, and on several default endpoints. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this HTTP(S) Source definition.

Address: Enter the hostname/IP on which to listen for HTTP(S) data. (E.g., localhost or 0.0.0.0.)

Port: Enter the port number.

Auth tokens: Shared secrets to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access **will be permitted**.

Cribl HTTP event API: Absolute path on which to listen for Cribl HTTP API requests. Currently, the only supported option is the default <code>/cribl</code>, which LogStream expands as <code>/cribl/_bulk</code>. Use an empty string to disable. Maximum payload size is 2MB.

Elastic API endpoint (for Bulk API): Absolute path on which to listen for Elasticsearch API requests. Currently, the only supported option is the default /elastic , which LogStream expands as /elastic/_bulk . Other entries are faked as success. Use an empty string to disable.

i Cribl generally recommends that you use the dedicated Elasticsearch API Source instead of this endpoint. The Elastic API implementation here is provided for backward compatibility, and for users who want to ingest multiple inputs on one HTTP/S port.

Splunk HEC endpoint: Absolute path on which to listen for Splunk HTTP Event Collector (HEC) API requests. Use an empty string to disable. Default entry is /services/collector.

i This Splunk HEC implementation is an event (i.e., not raw) endpoint. For details, see Splunk's documentation. To send data to it from a HEC client, use either /services/collector or /services/collector/event . (See the examples below.)

Cribl generally recommends that you use the dedicated Splunk HEC Source instead of this endpoint. The Splunk HEC implementation here is provided for backward compatibility, and for users who want to ingest multiple inputs on one HTTP/S port.

Splunk HEC acks: Whether to enable Splunk HEC acknowledgements. Defaults to No.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No. When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect.

Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __id (Elastic In)
- __type (Elastic In)
- __index (Elastic In)
- host (Elastic In)

Format and Endpoint

LogStream expects HTTP(S) events to be formatted as one JSON record per event. Here are two event records:

Sample Event Format

Note 1: Events can be sent as separate POSTs, but Cribl **highly** recommends combining multiple events in newline-delimited groups, and POSTing them together.

Note 2: If an HTTP(S) source is routed to a Splunk destination, fields within the JSON payload are mapped to Splunk fields. Fields that do not have corresponding (native) Splunk fields become indextime fields. For example, let's assume we have a HTTP(S) event like this:

```
{"_time":1541280341, "host":"myHost", "source":"mySource", "_raw":"this is a sample
event ", "fieldA":"valueA"}
```

Here, _time, host and source become their corresponding fields in Splunk. The value of _raw becomes the actual body of the event, and fieldA becomes an index-time field. (fieldA:: valueA).

Examples

For the following examples:

- 1. Configure Cribl to listen on port 10080 for HTTP (default). Set authToken to myToken42.
- 2. Send a payload to your Cribl LogStream receiver.

Cribl Endpoint - Single Event

Cribl Single Event Example:

```
curl -k http://<myCriblHost>:10080/cribl/_bulk -H 'Authorization: myToken42' -d '{"_raw":"this :
```

Cribl Endpoint - Multiple Events

```
Cribl Endpoint - Multiple Events
```

```
curl -k http://<myCriblHost>:10080/cribl/_bulk -H 'Authorization: myToken42' -d $'{"_raw":"this
```

Splunk HEC Event Endpoint

```
Splunk HEC Event Endpoint
```

```
curl -k http://<myCriblHost>:10080/services/collector/event -H 'Authorization: myToken42' -d '{'
curl -k http://<myCriblHost>:10080/services/collector -H 'Authorization: myToken42' -d '{"event"
```

Raw HTTP/S

Cribl LogStream supports receiving raw HTTP data. The Raw HTTP Source listens on a specific port, captures every HTTP request to that port, and creates a corresponding event that it pushes to its configured Event Breakers.

i Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring Cribl LogStream to Receive Raw HTTP Data

Select **Data > Sources**, then select **Raw HTTP** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Raw HTTP > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Raw HTTP Source definition.

Address: Enter the address to bind on. Defaults to 0.0.0.0 (all addresses).

Port: Enter the port number to listen on.

Auth tokens: Shared secrets to be provided by any client. Click **Generate** to create a new secret. If empty, permits open access.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

- Name: Field name.
- Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Allowed URI paths: List of URI paths accepted by this input. Supports wildcards, e.g., /api/v*/hook . Defaults to *, which allows all paths.

Allowed HTTP methods: List of HTTP methods accepted by this input. Supports wildcards, e.g., P*, GET. Defaults to *, which allows all methods.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions can use them to make processing decisions.

Fields accessible for this Source:

- __inputId
- __srcIpPort
- __channel

Kafka

Cribl LogStream supports receiving data records from a Kafka cluster.

i Type: Pull | TLS Support: YES | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Kafka Topics

Select **Data > Sources**, then select **Kafka** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Kafka > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Brokers: List of Kafka brokers to use, e.g., localhost:9092.

Topics: List of topics to subscribe to.

Group ID: The name of the consumer group to which this Cribl LogStream instance belongs.

From beginning: Whether to start reading from the earliest available data. Relevant only during initial subscription. Defaults to Yes.

TLS Settings (Client Side)

Enabled: defaults to No . When toggled to Yes :

Autofill?: This setting is experimental.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Authentication

This section governs SASL (Simple Authentication and Security Layer) authentication.

Enabled: Defaults to No. When toggled to Yes:

SASL mechanism: Use this drop-down to select the SASL authentication mechanism to use.

Username: Enter the username for your account.

Password: Enter the account's password.

Schema Registry

This section governs Kafka Schema Registry Authentication for AVRO-encoded data with a schema stored in the Confluent Schema Registry.

Enabled: defaults to No. When toggled to Yes:

Schema registry URL: URL for access to the Confluent Schema Registry. (E.g., http://<hostname>:8081.)

TLS enabled: defaults to No . When toggled to Yes, displays the following TLS settings for the Schema Registry:

i These have the same format as the TLS Settings (Client Side) above.

TLS Settings (Schema Registry)

Validate server certs: Reject certificates that are not authorized by a CA specified in the **CA Certificate Path** field. Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __topicIn (indicates the Kafka topic that the event came from; see
 - __topicOut in our Kafka Destination documentation)
- __schemaId (when using Schema Registry)

Kinesis

Cribl LogStream supports receiving data records from Amazon Kinesis Streams.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Kinesis Streams

Select **Data > Sources**, then select **Kinesis** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Kinesis > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Kinesis Stream Source definition.

Stream name: Kinesis stream name (not ARN) to read data from.

Shard iterator start: Location at which to start reading a shard for the first time. Defaults to Earliest Record.

Record data format: Format of data inside the Kinesis Stream records. Gzip compression is automatically detected. Options include:

- Cribl (the default): Use this option if LogStream wrote data to Kinesis in this format. This is a type of NDJSON.
- Newline JSON: Use if the records contain newline-delimited JSON
 (NDJSON) events e.g., Kubernetes logs ingested through Kinesis. This is a
 good choice if you don't know the records' format.
- CloudWatch Logs: Use if you've configured CloudWatch to send logs to Kinesis.
- Event per line: NDJSON can use this format when it fails to parse lines as valid JSON.

Region: Region where the Kinesis stream is located. Required.

Authentication

Authentication Method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for Kinesis Streams: Whether to use Assume Role credentials to access Kinesis Streams. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

• Name: Field name.

• Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Shard selection expression: A JavaScript expression to be called with each shardId for the stream. The shard will be processed if the expression evaluates to a truthy value. Defaults to true.

Service Period: Time interval (in minutes) between consecutive service calls. Defaults to 1 minute.

Endpoint: Kinesis stream service endpoint. If empty, the endpoint will be automatically constructed from the region.

Signature version: Signature version to use for signing Kinesis Stream requests. Defaults to v4.

Verify KPL checksums: Enable this setting to verify Kinesis Producer Library (KPL) event checksums.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

• __inputId

Kinesis Firehose

Cribl LogStream supports receiving data from Amazon Kinesis Data Firehose delivery streams via Kinesis' **HTTP endpoint** destination option.

i Type: Push | TLS Support: YES | Event Breaker Support: No

Configuring LogStream to Receive Data over HTTP(S) from Amazon Kinesis Firehose

Select **Data > Sources**, then select **Amazon > Firehose** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Firehose > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Address: Address to bind on. Defaults to 0.0.0.0 (all addresses).

Port: Enter the port number to listen on.

Auth tokens: Shared secrets to be provided by any client (Authorization: <token>). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No. When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

• Name: Field name.

• Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions

can use them to make processing decisions.

Fields accessible for this Source:

- __inputId
- __firehoseArn
- __firehoseReqId
- __firehoseEndpoint

Limitations/Troubleshooting

If you set the optional IntervalInSeconds and/or SizeInMBs parameters in the Kinesis Firehose BufferingHints API, beware of selecting extreme values (toward the ends of the API's supported ranges). These can send more bytes than LogStream can buffer, causing LogStream to send HTTP 500 error responses to Kinesis Firehose.

Azure Event Hubs

Cribl LogStream supports receiving data records from Azure Event Hubs.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Azure Event Hubs

Select **Data > Sources**, then select **Azure Event Hubs** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Azure Event Hubs > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this source definition.

Brokers: List of Event Hubs Kafka brokers to connect to, e.g., yourdomain.servicebus.windows.net:9093. Get the hostname from the host portion of the primary or secondary connection string in Shared Access Policies.

Event Hub name: The name of the Event Hub (a.k.a. Kafka Topic) to subscribe to.

Group ID: Specifies the name of the consumer group to which this Cribl LogStream instance belongs. Should always be \$Default for Event Hubs.

From beginning: Whether to start reading from the earliest available data. Relevant only during initial subscription. Defaults to Yes.

TLS Settings (Client Side)

Enabled: Defaults to Yes.

Validate server certs: Whether to reject connections to servers without signed certificates. Defaults to No – and for Event Hubs, must always be disabled.

Authentication

Enabled: Defaults to No. When toggled to Yes:

- **SASL** mechanism: SASL (Simple Authentication and Security Layer) authentication mechanism to use. Currently, PLAIN is the only mechanism supported for Event Hubs Kafka brokers.
- **Username**: The username for authentication. For Event Hubs, this should always be \$ConnectionString.
- **Password**: Connection-string primary key or connection-string secondary key from the Event Hub workspace.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Metrics

Cribl LogStream supports receiving metrics in these wire formats/protocols: StatsD, StatsD Extended, and Graphite. Automatic protocol detection will happen on the first line received over a TCP connection or a UDP packet. Lines not matching the detected protocol will be dropped.

Type: Push | TLS Support: No | Event Breaker Support: No

Configuring Cribl LogStream to Receive Metrics

Select **Data > Sources**, then select **Metrics** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Metrics > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Address: Enter the hostname/IP to listen to. Defaults to 0.0.0.0.

UDP port: Enter the UDP port number to listen on. Not required if listening on TCP.

TCP port: Enter the TCP port number to listen on. Not required if listening on UDP.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Defaults to No . Toggle to Yes if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

IP whitelist regex: Regex matching IP addresses that are allowed to send data. Defaults to .* (i.e., all IPs.)

Max buffer size (events): Maximum number of events to buffer when downstream is blocking. Defaults to 1000.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __srcIpPort
- __metricsInType

Metric Event Schema and Destination Support

Metric data is read into the following event schema:

```
Text

_metric - the metric name
_metric_type - the type of the metric (gauge, counter, timer)
_value - the value of the metric
_time - metric_time or Date.now()/1000
dim1 - value of dimension1
dim3 - value of dimension2
....
```

LogStream places sufficient information into a field called __criblMetric to enable these events to be properly serialized out to any metric outputs

(independent of the input type).

The following Destinations natively support the __criblMetric field:

- Splunk
- Splunk HEC
- InfluxDB
- Statsd
- Statsd Extended
- Graphite

Prometheus

Cribl LogStream supports receiving data from Prometheus.

i Type: Pull | TLS Support: No | Event Breaker Support: No

Configuring Cribl LogStream to Receive Data from Prometheus

Select **Data > Sources**, then select **Prometheus** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Prometheus > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Extra dimensions: Dimensions to include in events. By default, host and source are included.

Discovery type: Target discovery mechanism. Use Static (the default) to manually enter a list of targets. Select DNS or AWS EC2 options enable dynamic discovery of endpoints to scrape. Your selection determines which fields are displayed lower in this section:

- Targets: Displayed for Discovery type: Static . List of Prometheus targets to pull metrics from, values can be in URL or host[:port] format. For example: http://localhost:9090/metrics, localhost:9090, or localhost. In the cases where just host[:port] are specified, the endpoint will resolve to 'http://host[:port]/metrics'.
- **DNS names**: Displayed for **Discovery type**: DNS . Enter a list of DNS names to resolve.
- Record type: Displayed for Discovery type: DNS. Select the DNS record type to resolve. Defaults to SRV (Service record). Other options are A or AAA record.

• **Region**: Displayed for **Discovery type**: AWS EC2 . Select the AWS region in which to discover the EC2 instances with metrics endpoints to scrape.

Poll interval: How often (in minutes) to scrape targets for metrics. Defaults to 15. This value must be an integer that divides evenly into 60 minutes.

Log level: Set the verbosity level to one of debug, info (the default), warn, or error.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Keep alive time (seconds): How often workers should check in with the scheduler to keep job subscription alive. Defaults to 60 seconds.

Worker timeout (periods): How many **Keep alive time** periods before an inactive worker's job subscription will be revoked. Defaults to 3 periods.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

__source

- __isBroken
- __inputId
- __final
- __criblMetrics
- __channel
- __cloneCount

SQS

Cribl LogStream supports receiving events from Amazon Simple Queuing Service.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support:

Configuring Cribl LogStream to Receive Data from Amazon SQS

Select **Data > Sources**, then select **SQS** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **SQS > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this SQS Source definition.

Queue: The name, URL, or ARN of the SQS queue to read events from.

This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks. To specify a non-AWS URL, use the format: '{url}/<queueName>' . (E.g., ':port/<myQueueName>' .)

Create queue: Create queue if it does not exist.

Region: AWS Region where the SQS queue is located. Required, unless the **Queue** entry is a URL or ARN that includes a Region.

Authentication

Authentication Method: Select an AWS authentication method.

Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.

• Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

 API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.

 Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for SQS: Whether to use Assume Role credentials to access SQS. Defaults to No.

AWS account ID: SQS queue owner's AWS account ID. Leave empty if SQS queue is in same AWS account.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the external ID to use when assuming role.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Endpoint: SQS service endpoint. If empty, the endpoint will be automatically constructed from the AWS Region.

Signature version: Signature version to use for signing SQS requests. Defaults to v4.

Max messages: The maximum number of messages that SQS should return in a poll request. Amazon SQS never returns more messages than this value. (However, fewer messages might be returned.) Acceptable values: 1 to 10. Defaults to 10.

Visibility timeout seconds: The duration (in seconds) that the received messages are hidden from subsequent retrieve requests, after they're retrieved by a ReceiveMessage request. Defaults to 600.

Num receivers: The number of receiver processes to run. The higher the number, the better the throughput, at the expense of CPU overhead. Defaults to 3.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __sqsSysAttrs

The _sqsSysAttrs field can take on the following properties, which are reported to LogStream from SQS:

- __sqsSysAttrs.ApproximateFirstReceiveTimestamp: Returns the time (epoch time in milliseconds) the message was first received from the queue.
- __sqsSysAttrs.ApproximateReceiveCount: Returns the number of times a message has been received from the queue without being deleted.
- __sqsSysAttrs.SenderId: For an IAM user, returns the IAM user ID (e.g.: ABCDEFGHI1JKLMNOPQ23R). For an IAM role, returns the IAM role ID (e.g.: ABCDE1F2GH3I4JK5LMNOP:i-a123b456).
- __sqsSysAttrs.SentTimestamp: Returns the time (epoch time in milliseconds) the message was sent to the queue.

- __sqsSysAttrs.MessageDeduplicationId: Returns the value provided by the producer that calls the SendMessage action.
- __sqsSysAttrs.MessageGroupId: Returns the value provided by the producer that calls the SendMessage action messages with the same MessageGroupId are returned in sequence.
- __sqsSysAttrs.SequenceNumber: Returns the sequence-number value provided by Amazon SQS.
- __sqsSysAttrs.AWSTraceHeader: Returns the AWS X-Ray trace header string.

For background on these message properties, see AWS' ReceiveMessage > Request Parameters documentation.

SQS Permissions

The following permissions are needed on the SQS queue:

- sqs:ReceiveMessage
- sqs:DeleteMessage
- sqs:GetQueueAttributes
- sqs:GetQueueUrl
- sqs:CreateQueue (optional, if and only if you want LogStream to create the queue)

Troubleshooting Notes

▲ VPC endpoints for SQS might need to be set up in your account. Check with your administrator for details. Cribl LogStream supports receiving data from Amazon S3 buckets, using event notifications through SQS.

i Type: Pull | TLS Support: YES (secure API) | Event Breaker Support: YES

S3 Setup Strategy

i The source S3 bucket must be configured to send s3:ObjectCreated:* events to an SQS queue, either directly (easiest) or via SNS (Amazon Simple Notification Service). See the event notification configuration guidelines below.

SQS messages will be deleted after they're read, unless an error occurs, in which case LogStream will retry. This means that although LogStream will ignore files not matching the **Filename Filter**, their SQS events/notifications will still be read, and then deleted from the queue (along with those from files that match).

These ignored files will no longer be available to other S3 Sources targeting the same SQS queue. If you still need to process these files, we suggest one of these alternatives:

- Using a different, dedicated SQS queue. (Preferred and recommended.)
- Applying a broad filter on a single Source, and then using preprocessing Pipelines an/or Route filters for further processing.

Configuring Cribl LogStream to Receive Data from Amazon S3

Select **Data > Sources**, then select **S3** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **S3 > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this S3 Source definition.

Queue: The name, URL, or ARN of the SQS queue to read events from. When specifying a non-AWS URL, you must use the format: {url}/<queueName>.(E.g., https://host:port/<queueName>.) This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks.

Filename filter: Regex matching file names to download and process. Defaults to .*, to match all characters.

Region: AWS Region where the S3 bucket and SQS queue are located. Required, unless the **Queue** entry is a URL or ARN that includes a Region.

Authentication

Authentication method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for S3: Whether to use Assume Role credentials to access S3. Defaults to Yes.

Enable for SQS: Whether to use Assume Role credentials when accessing SQS (Amazon Simple Queue Service). Defaults to No.

AWS account ID: SQS queue owner's AWS account ID. Leave empty if the SQS queue is in the same AWS account.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + **Add Argument** to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

This section defines event breaking rulesets that will be applied, in order.

Event Breaker Rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker Buffer Timeout: The amount of time (in milliseconds) that the Event Breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Endpoint: S3 service endpoint. If empty, defaults to AWS's region-specific endpoint. Otherwise, used to point to an S3-compatible endpoint.

Signature version: Signature version to use for signing SQS requests. Defaults to v4.

Num receivers: The number of receiver processes to run,. The higher the number, the better the throughput, at the expense of CPU overhead. Defaults to 1.

Max messages: The maximum number of messages that SQS should return in a poll request. Amazon SQS never returns more messages than this value. (However, fewer messages might be returned.) Acceptable values: 1 to 10. Defaults to 1.

Visibility timeout seconds: The duration (in seconds) that the received messages are hidden from subsequent retrieve requests, after being retrieved by a ReceiveMessage request. Defaults to 600.

i LogStream will automatically extend this timeout until the initial request's files have been processed – notably, in the case of large files that require additional processing time.

Socket timeout: Socket inactivity timeout (in seconds). Increase this value if retrievals time out during backpressure. Defaults to 300 seconds.

Skip file on error: Toggle to **Yes** to skip files that trigger a processing error. (E.g., corrupted files.) Defaults to **No**, which enables retries after a processing error.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __source

How to Configure S3 to Send Event Notifications to SQS

- i For step-by-step instructions, see AWS' Walkthrough: Configure a Bucket for Notifications (SNS Topic and SQS Queue).
- 1. Create a Standard SQS Queue. Note its ARN.
- 2. Replace its access policy with one similar to the examples below. To do so, select the queue; and then, in the **Permissions** tab, click: **Edit Policy Document (Advanced)**. (These examples differ only at line 9, showing public access to the SQS queue versus S3-only access to the queue.)
- 3. In the Amazon S3 console, add a notification configuration to publish events of the s3:0bjectCreated:* type to the SQS queue.

Permissive SQS access policy Restrictive SQS access policy "Version": "example-2020-04-20", "Id": "example-ID", "Statement": ["Sid": "<SID name>", "Effect": "Allow", "Principal": { "AWS":"*" }, "Action": ["SQS:SendMessage" "Resource": "example-SQS-queue-ARN", "Condition": { "ArnLike": { "aws:SourceArn": "arn:aws:s3:*:*:example-bucket-name" } }]

S3 and SQS Permissions

The following permissions are required on the S3 bucket:

- s3:GetObject
- s3:ListBucket

The following permissions are required on the SQS queue:

- sqs:ReceiveMessage
- sqs:DeleteMessage
- sqs:ChangeMessageVisibility
- sqs:GetQueueAttributes
- sqs:GetQueueUrl

Best Practices

- When LogStream instances are deployed on AWS, use IAM Roles whenever possible.
 - Not only is this safer, but it also makes the configuration simpler to maintain.
- Although optional, we highly recommend that you use a Filename Filter.
 - This will ensure that LogStream ingests only files of interest.
 - Ingesting only what's strictly needed improves latency, processing power, and data quality.
- If higher throughput is needed, increase Advanced Settings > Number of Receivers and/or Max messages. However, do note:
 - These are set at 1 by default. Which means, each Worker Process, in each LogStream Worker Node, will run 1 receiver consuming 1 message (i.e. S3 file) at a time.
 - Total S3 objects processed at a time per Worker Node = Worker Processes x Number of Receivers x Max Messages
 - Increased throughput implies additional CPU utilization.
- When ingesting large files, tune up the **Visibility Timeout**, or consider using smaller objects.

• The default value of 600s works well in most cases, and while you certainly can increase it, we suggest that you also consider using smaller S3 objects.

Troubleshooting Notes

- VPC endpoints for SQS and for S3 might need to be set up in your account. Check with your administrator for details.
- If you're having connectivity issues, but no problems with the CLI, see if the AWS CLI proxy is in use. Check with your administrator for details.

Office 365 Services

Cribl LogStream supports receiving data from the Office 365 Service Communications API. This facilitates analyzing the status and history of service incidents on multiple Microsoft cloud services, along with associated incident and Message Center communications.

Type: Pull | TLS Support: YES | Event Breaker Support: YES

TLS is enabled via the HTTPS protocol on this Source's underlying REST API.

Configuring Cribl LogStream to Receive Data from Office 365 Services

Select **Data > Sources**, then select **Office 365 > Services** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Services > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Office 365 Services definition.

Tenant ID: Enter the Office 365 Azure tenant ID.

App ID: Enter the Office 365 Azure application ID.

Client secret: Enter the Office 365 Azure client secret.

Content Types

Here, you can configure polling separately for the following types of data from the Office 365 Service Communications API:

- **Current Status**: Get a real-time view of current and ongoing service incidents.
- Messages: Find incident and Message Center communications.
- Historical Status: Get a historical view of service incidents.

As of this revision, this Microsoft API provides data for Office 365, Yammer, Dynamics CRM, and Microsoft Intune cloud services. For each of these content types, this section provides the following controls:

Enabled: Toggle this to Yes for each service that you want to poll.

Interval: Optionally, override the default polling interval. See About Polling Intervals below.

Log level: Set the verbosity level to one of debug, info (the default), warn, or error.

About Polling Intervals

To poll the Office 365 Service Communications API, LogStream uses the **Interval** field's value to establish the search date range and the cron schedule (e.g.: */\${interval} * * * *).

Therefore, intervals set in minutes – those for **Current Status** and **Historical Status** – must divide evenly into 60 minutes to create a predictable schedule. Dividing 60 by intervals like 1, 2, 3, 4, 5, 6, 10, 12, 15, 20, or 60 itself yields an integer, so you can enter any of these values.

LogStream will reject intervals like 23, 42, or 45, or 75 – which would yield non-integer results, meaning unpredictable schedules.

The **Historical Status** service polls only once per day. So here, the **Interval** field's value simply establishes the hour of the day at which to poll. (In distributed deployments, this time is set based on the Master Node's system time. In single-instance deployments, it is set based on the API server's time zone.)

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Keep Alive Time (seconds): How often Workers should check in with the scheduler to keep their job subscription alive. Defaults to 60.

Worker timeout (periods): The number of Keep Alive Time periods before an inactive Worker will have its job subscription revoked. Defaults to 3.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __final
- __inputId
- __isBroken
- __source

Office 365 Activity

Cribl LogStream supports receiving data from the Office 365 Management Activity API. This facilitates analyzing actions and events on Azure Active Directory, Exchange, and SharePoint, along with global auditing and Data Loss Prevention data.

Type: Pull | TLS Support: YES | Event Breaker Support: YES

TLS is enabled via the HTTPS protocol on this Source's underlying REST API.

Configuring Cribl LogStream to Receive Data from Office 365 Activity

Select **Data > Sources**, then select **Office 365 > Activity** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Activity > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Office 365 Activity definition.

Tenant ID: Enter the Office 365 Azure tenant ID.

App ID: Enter the Office 365 Azure application ID.

Client secret: Enter the Office 365 Azure client secret.

Subscription Plan: Select the Office 365 subscription plan for your organization. This is typically Enterprise and GCC Government Plan.

Content Types

Here, you can configure polling independently for the following types of audit data from the Office 365 Management Activity API:

- Active Directory
- Exchange

- SharePoint
- General: All workloads not included in the above content types
- DLP.All: Data Loss Prevention events only, for all workloads

For each of these content types, this section provides the following controls:

Enabled: Toggle this to Yes for each service that you want to poll.

Interval: Optionally, override the default polling interval. See About Polling Intervals below.

Log level: Set the verbosity level to one of debug, info (the default), warn, or error.

About Polling Intervals

To poll the Office 365 Management Activity API, LogStream uses the Interval field's value to establish the search date range and the cron schedule (e.g.: */\${interval} * * * *).

Therefore, intervals set in minutes must divide evenly into 60 minutes to create a predictable schedule. Dividing 60 by intervals like 1, 2, 3, 4, 5, 6, 10, 12, 15, 20, or 60 itself yields an integer, so you can enter any of these values.

LogStream will reject intervals like 23, 42, or 45, or 75 – which would yield non-integer results, meaning unpredictable schedules.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Keep Alive Time (seconds): How often Workers should check in with the scheduler to keep their job subscription alive. Defaults to 60.

Worker timeout (periods): The number of Keep Alive Time periods before an inactive Worker will have its job subscription revoked. Defaults to 3.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __final
- __inputId
- __isBroken
- __source

SNMP Trap

Cribl LogStream supports receiving data from SNMP Traps.

i Type: Push | TLS Support: NO | Event Breaker Support: No

Configuring Cribl LogStream to Receive SNMP Traps

Select **Data > Sources**, then select **SNMP Trap** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **SNMP Trap > New Source** pane, which provides the fields outlined below.

☐ LogStream ships with an SNMP Trap Source preconfigured to listen on Port 9162. You can clone or directly modify this Source to further configure it, and then enable it.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Address: Address to bind on. Defaults to 0.0.0.0 (all addresses).

UDP Port: Port on which to receive SNMP traps. Defaults to 162.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

IP whitelist regex: Regex matching IP addresses that are allowed to send data. Defaults to .* i.e. all IPs.

Max buffer size (events): Maximum number of events to buffer when downstream is blocking. Defaults to 1000.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

- __inputId
- __snmpVersion: Acceptable values are 0, 2, 3. Versions: 0 =v1, 2 =v2c, 3 =v3.
- __srcIpPort : <hostname>|port
- __snmpRaw : Buffer containing Raw SNMP packet

Considerations for Working with SNMP Trap Data

- It's possible to work with SNMP metadata (i.e., we'll decode the packet). Options include dropping, routing, etc.
- SNMP packets can be forwarded to other SNMP destinations. However, the contents of the incoming packet **cannot** be modified i.e., we'll forward the packets verbatim as they came in.
- SNMP packets can be forwarded to non-SNMP destinations (e.g., Splunk, Syslog, S3, etc.).
- Non-SNMP input data **cannot** be sent to SNMP destinations.

Datagens

Cribl LogStream supports generating of data from datagen files. See Using Datagens for more details.

i Type: Internal | TLS Support: N/A | Event Breaker Support: No

Configuring Cribl LogStream to Generate Sample Data

Select **Data > Sources**, then select **Datagens** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **Datagens > New Source** pane, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this Source definition.

Datagens: List of datagens.

- Data generator file: Name of the datagen file.
- Events per second per Worker Node: Maximum number of events to generate per second, per worker node. Defaults to 10.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Fields for this Source:

• __inputId

Cribl Internal

The Cribl Internal Source enables you to capture and send LogStream's own internal **logs** and **metrics** through Routes and Pipelines. In distributed mode, only Worker Node internal logs can be processed through this Source. (Logs on the Master remain on the Master, since the Master Node is not part of any processing path.)

i Type: Internal | TLS Support: N/A | Event Breaker Support: No

Configuring Cribl Internal Logs/Metrics to Behave as a Data Source

Select **Data > Sources**, then select **Cribl Internal** from the **Data Sources** page's tiles or left menu.

Next, on the **CriblLogs** and/or the **CriblMetrics** row, slide the **Enabled** slider to Yes. Confirm your choice in the resulting message box.

To proceed to the configuration options listed below, click anywhere on the **CriblLogs** or the **CriblMetrics** row.



Cribl Internal Sources - click to configure

CriblLogs Settings

General Settings

Enabled: This duplicates the parent page's **Enabled** slider. Keep it at Yes to enable Cribl logs as a Source.

Input ID: Enter a unique name to identify this CriblLogs Source definition.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

CriblMetrics Settings

General Settings

Enabled: This duplicates the parent page's **Enabled** slider. Keep it at Yes to enable Cribl metrics as a Source.

Input ID: Enter a unique name to identify this CriblMetrics Source definition.

Metric name prefix: Enter an optional prefix that will be applied to metrics provided by LogStream. The prefix defaults to cribl.logstream. .

i If LogStream detects source, sourcetype, host, or index fields in metrics from external sources, it copies their values into new dimensions with added event_ prefixes (e.g., event_sourcetype). This leaves the original dimensions (and their values) intact.

Note that you can disable metric collection for any or all of these four fields at System Settings > General Settings > Limits > Disable field metrics.

Processing Settings

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Reporting Metrics Less Frequently

By default, LogStream generates internal metrics every 2 seconds. To consume metrics at longer intervals, you can use or adapt the <code>cribl-metrics_rollup</code> Pipeline that ships with LogStream. Attach it to your <code>Cribl Internal</code> Source as a <code>pre-processing Pipeline</code>. The Pipeline's <code>Rollup Metrics</code> Function has a default <code>Time Window</code> of 30 seconds, which you can adjust to a different granularity as needed.

Internal Fields

The following fields will be added to all events/metrics:

source:setto cribl.

• host: set to the hostname of the Cribl instance.

Use these fields to guide these events/metrics through Cribl Routes.

All Cribl internal fields are subject to change and modification. Cribl provides them to assist with analytics and diagnostics, but does not guarantee that they will remain available.

AppScope

AppScope is an open-source instrumentation utility from Cribl. It offers visibility into any Linux command or application, regardless of runtime, with no code modification. For details about configuring the AppScope CLI, loader, and library, see: https://appscope.dev/docs.

Type: Push | TLS Support: YES | Event Breaker Support: YES

Configuring LogStream to Receive AppScope Data

Select **Data > Sources**, then select **AppScope** from the **Data Sources** page's tiles or left menu. Click **Add New** to open the **AppScope > New Source** modal, which provides the following fields.

General Settings

Input ID: Enter a unique name to identify this AppScope Source definition.

Address: Enter the hostname/IP on which to listen for AppScope data. (E.g., localhost .) Defaults to 0.0.0.0, meaning all addresses.

Port: Enter the port number to listen on.

IP whitelist regex: Regex matching IP addresses that are allowed to establish a connection.

Auth tokens: Shared secrets to be provided by any client (in authToken header field). Click **Generate** to create a new secret. If empty, unauthenticated access will be permitted.

TLS Settings (Server Side)

Enabled defaults to No . When toggled to Yes:

Certificate name: Name of the predefined certificate.

Private key path: Path on server where to find the private key to use in PEM format. Path can reference \$ENV_VARS.

Passphrase: Passphrase to use to decrypt private key.

Certificate path: Server path at which to find certificates (in PEM format) to use. Path can reference \$ENV_VARS.

CA certificate path: Server path at which to find CA certificates (in PEM format) to use. Path can reference \$ENV_VARS.

Validate client certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Authenticate client (mutual auth): Require clients to present their certificates. Used to perform mutual authentication using SSL certs. Defaults to No . When toggled to Yes:

• Common name: Regex matching peer certificate subject common names allowed to connect. Defaults to .*.

Minimum TLS version: Optionally, select the minimum TLS version to accept from connections.

Maximum TLS version: Optionally, select the maximum TLS version to accept from connections.

Processing Settings

Event Breakers

Event Breaker rulesets: A list of event breaking rulesets that will be applied to the input data stream before the data is sent through the Routes. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute field's value (can be a constant).

Pre-Processing

In this section's **Pipeline** drop-down list, you can select a single existing Pipeline to process data from this input before the data is sent through the Routes.

Advanced Settings

Enable proxy protocol: Enable if the connection is proxied by a device that supports Proxy Protocol v1 or v2.

Internal Settings

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and Functions can use them to make processing decisions.

Field for this Source:

• __inputId

^ Reviewers: This is placeholder copy/pasta from the Elasticsearch Source. Which internal fields are actually supported here?

Collectors

Collectors enable you to dispatch on-demand collection tasks that fetch data from local or remote locations. As of v.2.3, LogStream supports scheduled collection jobs. These recurring jobs can make batch collection of stored data more like continual processing of streaming data.

How Do Collectors Work

You can configure a LogStream Node to retrieve data from a remote system via **Data** > **Collectors**. Data collection is a multi-step process:

First, define a Collector instance. In this step, you configure **collector-specific settings** by selecting a Collector type and pointing it at a specific target. (E.g., the target will be a directory if the type is Filesystem, or an S3 bucket/path if the type is Amazon S3.)

Next, schedule or manually run the Collector. In this step, you configure scheduled-job-specific or run-specific settings – such as the run Mode (such as Discovery or Full Run), the Filter expression to match the data against, the time range, etc.

When a Node receives this configuration, it prepares the infrastructure to execute a collection job. A collection job is typically made up of one or more tasks that: discover the data to be fetched; fetch data that match the run filter; and finally, pass the results either through the Routes or (optionally) into a specific Pipeline and Destination.

i On the Manage Collectors page, click Job Inspector to see the results of recent collection runs. Select the Show system jobs check box to also display discovery jobs and collection jobs for the Office 365 System/Activity Sources.

Scheduled Collection Jobs

You might process data from inherently non-streaming sources, such as REST endpoints, blob stores, etc. Scheduled jobs enable you to emulate a data stream by scraping data from these sources in batches, on a set interval.

You can schedule a specific job to pick up new data from the source – data that hadn't been picked up in previous invocations of this scheduled job. This essentially transforms a non-streaming data source into a streaming data source.

Collectors in Distributed Deployments

In a distributed deployment, Collectors are set up at the Worker Group level, and the tasks are executed by Worker Nodes. The Master Node oversees the task distribution, and tries to maintain a fair balance across jobs.

When Workers ask for tasks, the Master will normally try to assign the next task from a job with the least tasks in progress. This is known as "Least-In-Flight Scheduling," and provides the fairest task distribution for most cases. Default behavior can be changed via **Settings > General Settings > Job Limits > Job Dispatching**.

Collector Types

Cribl LogStream currently provides the following Collector options:

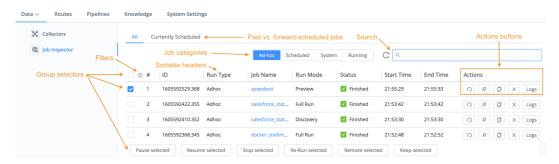
- Filesystem/NFS enables data collection from local or remote filesystem locations.
- S3 enables data collection from Amazon S3 buckets or S3-compatible stores.
- Script enables data collection via custom scripts.
- REST enables data collection via REST API calls. Provides four Discover options, to support progressively more complex (and dynamic) item enumerations.

Monitoring Collection Jobs

Select **Monitoring > Jobs** from the top menu to see a graphical display of inflight collection jobs and their tasks.

Inspecting Jobs

Select **Data > Collectors > Job Inspector** to view and manage pending, in-flight, and completed collection jobs.



Job Inspector: all the things

Here are the options available on the Job Inspector page:

- All vs. Currently Scheduled tabs: Click Currently Scheduled to see jobs foward-scheduled for future execution – including their cron schedule details, last execution, and next scheduled execution. Click All to see all jobs initiated in the past, regardless of completion status.
- Job categories (buttons): Select among Ad hoc, Scheduled, System, and Running. (At this level, Scheduled means scheduled jobs already running or finished.)
- Filters: Click the gear icon to open a drop-down with multiple options to filter the jobs shown within your selected category.
- Group selectors: Select one or more check boxes to display the Pause,
 Resume, etc., buttons shown along the bottom.
- Sortable headers: Click any column to reverse its sort direction.
- Search bar: Click to filter displayed jobs by arbitrary strings.
- Action buttons: For finished jobs, the icons (from left to right) indicate: Rerun; Keep job artifacts; Copy job artifacts; Delete job artifacts; and Display job logs in a modal. For running jobs, the options (again from left to right) are: Pause; Stop; Copy job artifacts; Delete job artifacts; and Live (show collection status in a modal).

What's Next

See the configuration instructions for the collector type you want to configure, Then proceed to instructions for scheduling and running collection jobs.

Filesystem/NFS
S3
Script
REST
Scheduling and Running

Filesystem/NFS

Cribl LogStream supports collecting data from a local or a remote filesystem location.

Configuring a Filesystem Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID: Unique ID for this Collector. E.g., DysonV11Roomba960.

Collector type: Defines the type of Collector to configure.

i Set this to Filesystem to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Auto-populate from: Select a Destination with which to auto-populate Collector settings. Useful when replaying data.

Directory: The directory from which to collect data. Templating is supported (e.g., /myDir/\${host}/\${year}/\${month}/). You can also use templating to specify (e.g.) a Splunk bucket from which to collect. Symlinks will not be followed. More on templates and Filters.

Recursive: If set to Yes (the default), data collection will recurse through subdirectories.

Destructive: If set to Yes, the Collector will delete files after collection. Defaults to No.

Max batch size (files): Maximum number of lines written to the discovery results files each time. To override this limit in the Collector's Schedule/Run modal, use Advanced Settings > Upper task bundle size.

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + Add Argument to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

- Pipeline: Select a Pipeline to process results.
- **Destination**: Select a Destination to receive results.
 - i You might disable **Send to Routes** when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. This keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

Cribl LogStream supports collecting data from Amazon S3 stores. This page covers how to configure the Collector.

☐ For a step-by-step tutorial on using LogStream to replay data from an S3-compatible store, see our Data Collection & Replay sandbox. The sandbox takes about 30 minutes. It provides a hosted environment, with all inputs and outputs preconfigured for you.

Configuring an S3 Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID: Unique ID for this Collector. E.g., Attic42TreasureChest.

Collector type: Defines the type of Collector to configure.

i Set this to S3 to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Auto-populate from: Select a Destination with which to auto-populate Collector settings. Useful when replaying data.

S3 bucket: Simple Storage Service bucket from which to collect data.

Region: S3 Region from which to retrieve data.

Path: Path, within the bucket, from which to collect data. Templating is supported (e.g., /myDir/\${host}/\${year}/\${month}/). More on templates and Filters.

Recursive: If set to Yes (the default), data collection will recurse through subdirectories.

Max batch size (files): Maximum number of lines written to the discovery results files each time. To override this limit in the Collector's Schedule/Run modal, use Advanced Settings > Upper task bundle size.

Authentication

API key: Enter API key. If empty, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials. Optional when running on AWS.

Secret key: Enter secret key. if empty, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials. Optional when running on AWS.

Assume Role

Enable Assume Role: Slide to Yes to enable Assume Role behavior.

AssumeRole ARN: Amazon Resource Name (ARN) of the role to assume.

External ID: External ID to use when assuming role.

Additional Collector Settings

Endpoint: S3 service endpoint. If empty, LogStream will automatically construct the endpoint from the region.

Signature version: Signature version to use for signing S3 requests. Defaults to $\vee 4$.

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + **Add Argument** to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

• Pipeline: Select a Pipeline to process results.

• **Destination**: Select a Destination to receive results.

i You might disable **Send to Routes** when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. This keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data – potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

Scheduling and Running

Script

Cribl LogStream supports flexible data collection configured by your custom scripts.

Configuring a Script Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID: Unique ID for this Collector. E.g., sh2GetStuff.

Collector type: Defines the type of Collector to configure.

i Set this to Script to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Discover script: Script to discover which objects/files to collect. This script should output one task per line in stdout.

Collect script: Script to perform data collections. Pass in tasks from the Discover script as \$CRIBL_COLLECT_ARG . Should output results to stdout.

Shell: Shell in which to execute scripts. Defaults to /bin/bash.

Scripts will allow you to execute almost anything on the system where Cribl LogStream is running. Make sure you understand the impact of what you're executing before you do so! These scripts run

as the user running LogStream, so if you are running it as root, these commands will run with root user permissions. 🧟 🧟

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + Add Argument to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the Routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

- Pipeline: Select a Pipeline to process results.
- **Destination**: Select a Destination to receive results.
 - i You might disable **Send to Routes** when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. This keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

REST / API Endpoint

Cribl LogStream supports collecting data from REST endpoints.

Configuring a REST Collector

From the top menu, select **Data > Collectors**. On the resulting **Manage Collectors page**, click **Add New**. The resulting **New Collector** modal displays the following options and fields.

Collector Settings

The Collector Settings determine how data is collected before processing.

Collector ID

Unique ID for this Collector. E.g., rest42json.

Collector Type

Defines the type of Collector to configure.

i Set this to REST to configure the Collector as shown below.

The sections described below are spread across several tabs. Click the tab links at left, or the **Next** and **Prev** buttons, to navigate among tabs. Click **Save** when you've configured your Collector.

Discover Type

Once you've selected the REST Collector type above, this exposes a Discover type drop-down. Here you have four options, corresponding to different use cases. Each Discover type selection will expose a different set of Collector Settings fields. Below, we cover the Discover types from simplest to most-complex.

- **Discover type: None** matches cases where one simple API call will retrieve all the data you need. This suppresses the Discover stage. (Example: Collect a list of configured LogStream Pipelines.)
- **Discover type: Item List** matches cases where you want to enumerate a known list of items to retrieve. (Example: Collect network traffic data that's tagged with specific subnets.)
- Discover type: JSON Response provides a Discover result field where you can (optionally) define Discover tasks as a JSON array of objects. Each entry returned by Discover will generate a Collect task. (Example: Collect data for specific geo locations the National Weather Service API's stream of worldwide weather data. This API requires multiple parameters in the request URL latitude, longitude, etc. so an Item List would not work.)
- Discover type: HTTP Request matches cases where you need to
 dynamically discover what you can collect from a REST endpoint. This
 Discover type most fully exploits LogStream's Discover-Before-Collect
 architecture. (Example: Make a REST call to get a list of available log files,
 then run Collect against each of those files.)

Common Collector Settings / Discover Type: None

These remaining **Collector Settings** options appear for **Discover type**: None, as well as for all other **Discover type** selections:

i Time Range Variables

The following fields fields accept \${earliest} and \${latest} variables, which reference any **Time Range** values that have been set in manual or scheduled collection jobs:

- Collect URL, Collect parameters, Collect headers
- Discover URL, Discover parameters, Discover headers.

As an example, here is a **Collect URL** entry using these variables: http://localhost/path?from=\${earliest}&to=\${latest}

Both variables are formatted as UNIX epoch time, in seconds units. When using them in contexts that require milliseconds resolution, multiply them by 1,000 to convert to ms.

Collect URL: URL (constant or JavaScript expression) to use for the Collect operation.

i Any variables used in a URL (path or parameters) must be encoded using: C.Encode.uri(paramName).

As of v.2.3.2, URLs/expressions specified in this field follow redirects.

Collect method: Select the HTTP verb to use for the Collect operation – GET, POST, or POST with body.

Collect POST body: Template for POST body to send with the Collect request. (This field is displayed only when you set the **Collect method** to POST with body .) You can reference parameters from the Discover response using template params of the form: \${variable}.

Collect parameters: Optional HTTP request parameters to append to the request URL. These refine or narrow the request. Click + **Add Parameter** to add parameters as key-value pairs:

- Name: Field name.
- Value: JavaScript expression to compute the field's value (can be a constant).

Collect headers:: Click + Add Header to (optionally) add collection request haaders as key-value pairs:

- Name: Header name.
- Value: JavaScript expression to compute the header's value (can be a constant).
 - i By adding the appropriate **Collect headers**, you can specify API Keybased authentication as an alternative to the Authentication: Basic or Login options below.

Authentication

In the **Authentication** drop-down, select an authentication method to use for discover and collect REST requests:

- None: Compatible with REST servers like AWS, where you embed a secret directly in the request URL.
- Basic: Compatible with Basic Authentication servers. Selecting Basic exposes additional fields in which you specify a Basic Auth zxUsername and Password.
- Login: Enables you to specify several credentials, then perform a POST to an endpoint during the Discover operation. The POST response returns a token, which LogStream uses for later Collect operations.

Selecting Login exposes the following additional fields:

- Login URL: URL for the login API call, which is expected to be a POST call.
- Username: Login username.
- Password: Login password.
- **POST Body**: Template for POST body to send with the login request. The \${username} and \${password} variables specify the corresponding credentials' locations in the message.
- **Token Attribute**: Path to the token attribute in the login response body. Supports nested attributes.
- Authorize Expression: JavaScript expression used to compute the Authorization header to pass in Discover and Collect calls. Uses \${token} to reference the token obtained from the login POST request.

Discover Type: Item List

Setting the **Discover type** to Item List exposes this additional field above the Common Collector Settings:

Discover Items: List of items to return from the Discover task. Each returned item will generate a Collect task, and can be referenced using \${id} in the **Collect URL**, the **Collect parameters**, or the **Collect headers**.

Discover Type: JSON Response

Setting the **Discover type** to JSON Response exposes these additional fields above the Common Collector Settings:

Discover result: Allows hard-coding the Discover result. Must be a JSON object. Works with the Discover data field.

Discover data field: Within the response JSON, name of the field or array element to pull results from. Leave blank if the result is an array of values. Sample entry: items, json: { items: [{id: 'first'},{id: 'second'}] }

Discover Type: HTTP Request

Setting the **Discover type** to HTTP Request exposes these additional fields above the Common Collector Settings:

Discover URL: Enter the URL to use for the Discover operation. This can be a constant URL, or a JavaScript expression to derive the URL.

i Any variables used in a URL (path or parameters) must be encoded using: C.Encode.uri(paramName).

As of v.2.3.2, URLs/expressions specified in this field follow redirects.

Discover method: Select the HTTP verb to use for the Discover operation – GET, POST, or POST with body.

Discover POST body: Template for POST body to send with the Discover request. (This field is displayed only when you set the **Discover method** to POST with body .)

Discover parameters: Optional HTTP request parameters to append to the Discover request URL. These refine or narrow the request. Click

- + Add Parameter to add parameters as key-value pairs:
 - Name: Parameter name.
 - Value: JavaScript expression to compute the parameter's value (can also be a constant).

Discover headers: Optional Discover request headers.: Click **+ Add Header** to add headers as key-value pairs:

- Name: Header name.
- Value: JavaScript expression to compute the header's value (can also be a constant).

Discover data field: Within the response JSON, name of the field that contains Discover results. Leave blank if the result is an array.

i The following sections describe the Collector Settings' remaining tabs, whose settings and content apply equally to all **Discover type** selections.

Result Settings

The Result Settings determine how LogStream transforms and routes the collected data.

Custom Command

In this section, you can pass the data from this input to an external command for processing, before the data continues downstream.

Enabled: Defaults to No . Toggle to Yes to enable the custom command.

Command: Enter the command that will consume the data (via stdin) and will process its output (via stdout).

Arguments: Click + **Add Argument** to add each argument to the command. You can drag arguments vertically to resequence them.

Event Breakers

In this section, you can apply event breaking rules to convert data streams to discrete events.

Event Breaker rulesets: A list of event breaking rulesets that will be applied, in order, to the input data stream. Defaults to System Default Rule.

Event Breaker buffer timeout: The amount of time (in milliseconds) that the event breaker will wait for new data to be sent to a specific channel, before flushing out the data stream, as-is, to the routes. Defaults to 10000.

Fields (Metadata)

In this section, you can add fields/metadata to each event, using Eval-like functionality.

Name: Field name.

Value: JavaScript expression to compute the field's value (can be a constant).

Result Routing

Send to Routes: If set to Yes (the default), events will be sent to normal routing and event processing. Toggle to No to select a specific Pipeline/Destination combination. The No setting exposes these two additional fields:

- Pipeline: Select a Pipeline to process results.
- **Destination**: Select a Destination to receive results.
 - i You might disable Send to Routes when configuring a Collector that will connect data from a specific Source to a specific Pipeline and Destination. One use case might be a REST Collector that gathers a known, simple type of data from a single endpoint. This approach keeps the Collector's configuration self-contained and separate from LogStream's routing table for live data potentially simplifying the Routes structure.

Preprocess Pipeline: Pipeline to process results before sending to Routes. Optional.

Throttling: Rate (in bytes per second) to throttle while writing to an output. Also takes values with multiple-byte units, such as KB, MB, GB, etc. (Example: 42 MB.) Default value of 0 indicates no throttling.

Advanced Settings

Advanced Settings enable you to customize post-processing and administrative options.

Time to live: How long to keep the job's artifacts on disk after job completion. This also affects how long a job is listed in **Job Inspector**. Defaults to 4h.

Remove Discover fields: List of fields to remove from the Discover results. This is useful when discovery returns sensitive fields that should not be exposed in the Jobs user interface. You can specify wildcards (such as aws*).

Resume job on boot: Toggle to Yes to resume ad hoc collection jobs if LogStream restarts during the jobs' execution.

What's Next

> Scheduling and Running

Scheduling and Running

Once you've configured a Collector, you can either run it immediately ("ad hoc") to collect data, or schedule it to run on a recurring interval. Scheduling requires some extra configuration upfront, so we cover this option first.

i For ad hoc collection, you can configure whether a job interrupted by a LogStream shutdown will automatically resume upon LogStream restart.

But regardless of this configuration, explicitly restarting or stopping LogStream (via ./cribl restart, ./cribl stop, or Settings > Controls > Restart) will cancel any currently running jobs.

A scheduled job will not resume upon restart.

Schedule Configuration

Click **Schedule** beside a configured Collector to display the **Schedule configuration** modal. This provides the following controls.

Enabled: Slide to Yes to enable this collection schedule.

⚠ The scheduled job will keep running on this schedule forever, unless you toggle **Enabled** back to Off. The Off setting preserves the schedule's configuration, but prevents its execution.

Cron schedule: A cron schedule on which to run this job.

• The **Estimated schedule** below this field shows the next few collection runs, as examples of the cron interval you've scheduled.

Skippable: Skippable jobs can be delayed up to their next run time if the system is hitting concurrency limits. Defaults to Yes.

Skippable Jobs and Concurrency Limits

If toggled to Yes, the **Skippable** option obliges these concurrency limits in **Settings** > **General Settings** > **Job Limits**:

- Concurrent Job Limit
- Concurrent Scheduled Job Limit

When the above limits delay a Skippable job:

- The Skippable job will be granted slightly higher priority than non-Skippable jobs.
- If the job receives resources to run before its next scheduled run, LogStream will run the delayed job, then snap back to the original cron schedule.
- If resources do not free up before the next scheduled run: LogStream will skip the delayed run, and snap back to the original cron schedule.

Set **Skippable** to No if you absolutely must have all your data, for compliance or other reasons. In this case, LogStream will build up a backlog of jobs to run.

You can think of Skippable: No as behaving more like the TCP protocol, with Skippable: Yes behaving more like UDP.

Max Concurrent Runs: Sets the maximum number of instances of this scheduled job that may simultaneously run.

⚠ All collection jobs are constrained by the following Settings > General Settings > Job Limits:

- Concurrent Task Limit
- Max Task Usage Percentage

Run Configuration and Shared Settings

Most of the remaining fields and options below are shared with the Run configuration modal, which you can open by clicking Run beside a configured Collector.

Mode

Depending on your requirements, you can schedule or run a collector in these modes:

- Preview default for Run, but not offered for Scheduled Jobs
- Discovery default for Scheduled Jobs
- Full Run

Preview

In the Preview mode, a collection job will return only a sample subset of matching results (e.g., 100 events). This is very useful in cases when users need a data sample to:

- Ensure that the correct data comes in.
- Iterate on filter expressions.
- Capture a sample to iterate on pipelines.

i Schedule configuration omits the Preview option, because Preview is designed for immediate analysis and decision making. To configure a Scheduled Job with high confidence, you might want to first manually run Preview jobs with the same Collector, to verify that you're collecting the data you expect.

Preview Settings

In Preview mode, you can optionally set these limits:

- Capture time (sec): Maximum time interval (in seconds) to collect data.
- Capture up to N events: Maximum number of events to capture.

Discovery

In Discovery mode, a collection job will return only **the list of objects/files** to be collected, but none of the data. This mode is typically used to ensure that the Filter expression and time range are correct before a Full Run job collects unintended data.

Send to Routes

In Discovery mode, this slider enables you to send discovery results to LogStream Routes. Defaults to $\,\mathrm{No}\,$.

i This setting overrides the Collector configuration's Result Routing > Send to Routes setting.

Full Run

In Full Run mode, the collection job is fully executed by Worker Nodes, and will return all data matching the Run configuration.

Time Range

Set an **Absolute** or **Relative** time range for data collection.

 $\begin{tabular}{ll} \hline \end{tabular} The \textbf{\it Relative} option is the default, and is particularly useful for configuring scheduled jobs. \\ \hline \end{tabular}$

Absolute

Select the **Absolute** button to set fixed collection boundaries in your local time. Next, use the **Earliest** and **Latest** controls to set the start date/time and end date/time.

Relative

Select the **Relative** button to set collection boundaries relative to the current time. Next, use the **Earliest** and **Latest** to set start and end times like these:

- Earliest example values: -1hr, -42m, -42m@h
- Latest example values: now, -20m, +42m@h

Relative Time Syntax

For Relative times, the **Earliest** and **Latest** controls accept the following syntax:

```
[+ \vdash] < time_integer > < time_unit > 0 < snap-to_time_unit > 0 < snap-to_tim
```

To break down this syntax:

Syntax Element	Values Supported
Offset	Specify: - for times in the past, + for times in the future, or omit with now .
<time_integer></time_integer>	Specify any integer, or omit with now .
<time_unit></time_unit>	Specify the now constant, or one of the following abbreviations: s[econds], m[inutes], h[ours], d[ays], w[eeks], mon[ths], q[uarters], y[ears].
@ <snap- to_time_unit></snap- 	Optionally, you can append the @ modifier, followed by any of the above <time_unit> s, to round down to the nearest instance of that unit. (See the next section for details.)</time_unit>

LogStream validates relative time values using these rules:

- Earliest must not be later than Latest.
- Values without units get interpreted as seconds. (E.g., -1 = -1s.)

Snap-to-Time Syntax

The a snap modifier always rounds **down** (backwards) from any specified time. This is true even in relative time expressions with + (future) offsets. For example:

- ad snaps back to the beginning of today, 12:00 AM (midnight).
- +128mah looks forward 128 minutes, then snaps back to the nearest round hour. (If you specified this in the Latest field, and ran the Collector at 4:20 PM, collection would end at 6:00 PM. The expression would look forward to 6:28 PM, but snap back to 6:00 PM.)

Other options:

- aw or aw7 to snap back to the beginning of the week defined here as the preceding Sunday.
- To snap back to other days of a week, use w1 (Monday) through w6 (Saturday).
- am to snap back to the 1st of a month.

- aq to snap back to the beginning of the most recent quarter Jan. 1, Apr. 1, Jul. 1, or Oct. 1.
- ay to snap back to Jan. 1.

Filter

This is a JavaScript filter expression that is evaluated against token values in the provided collector path (see below), and against the events being collected. The **Filter** value defaults to true, which matches all data, but this value can be customized almost arbitrarily.

For example, if a Filesystem or S3 collector is run with this Filter:

```
host='myHost' & source.endsWith('.log') || source.endsWith('.txt')
```

...then only files/objects with .log or .txt extensions will be fetched. And, from those, only those events with host field myHost will be collected.

For more extensive options, see Tokens for Filtering below.

Advanced Settings

Log Level: Level at which to set task logging. More-verbose levels are useful for troubleshooting jobs and tasks, but use them sparingly.

Lower task bundle size: Limits the bundle size for small tasks. E.g., bundle five 200KB files into one 1MB task bundle. Defaults to 1MB.

Upper task bundle size: Limits the bundle size for files above the **Lower task bundle size**. E.g., bundle five 2MB files into one 10MB task bundle. Files greater than this size will be assigned to individual tasks. Defaults to 10MB.

Reschedule tasks: Whether to automatically reschedule tasks that failed with non-fatal errors. Defaults to Yes; does not apply to fatal errors.

Max task reschedule: Maximum number of times a task can be rescheduled. Defaults to 1.

Job timeout: Maximum time this job will be allowed to run. Units are seconds, if not specified. Sample values: 30, 45s, or 15m. Minimum granularity is 10 seconds, so a 45s value would round up to a 50-second timeout. Defaults to 0, meaning unlimited time (no timeout).

Tokens for Filtering

Let's look at the options for path-based (basic) and time-based token filtering.

Basic Tokens

In collectors with paths, such as Filesystem or S3, LogStream supports path filtering via token notation. Basic tokens' syntax follows that of JS template literals: \${<token_name>} - where token_name is the field (name) of interest.

For example, if the path was set to \\var\log\\$\{\hostname\}\\$\{\sourcetype\}\, you could use a Filter such as \hostname='myHost' & sourcetype='mySourcetype' to collect data only from the \\var\log\myHost\mySourcetype\ subdirectory.

Time-based Tokens

In paths with time partitions, LogStream supports further filtering via time-based tokens. This has a direct effect with earliest and latest boundaries. When a job runs against a path with time partitions, the job traverses a minimal superset of the required directories to satisfy the time range, before subsequent event _time_filtering.

About Partitions and Tokens

LogStream processes time-based tokens as follows:

- For each path, time partitions must be notated in descending order. So Year/Month/Day order is supported, but Day/Month/Year is not.
- Paths may contain more than one partition. E.g., /my/path/2020-04/20/.
- In a given path, each time component can be used only once.
 So /my/path/\${_time:%Y}/\${_time:%m}/\${_time:%d}/... is a valid expression format, but /my/path/\${_time:%Y}/\${_time:%m}/\${host}/\${_time:%Y}/... (with a repeated Y) is not supported.
- For each path, all extracted dates/times are considered in UTC.

The following strptime format components are allowed:

- 'Yy', for years
- 'mBbj', for months
- 'dj', for days
- 'HI', for hours
- 'M', for minutes
- 'S', for seconds

Token Syntax

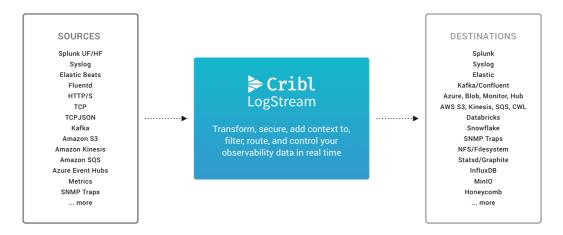
Time-based token syntax follows that of a slightly modified JS template literal:

\${_time: <some_strptime_format_component>} . Examples:

Filter	Matches
/my/path/\${_time:%Y}/\${_time:%m}/\${_time:%d}/	/my/path/2020/04/20/
/my/path/\${_time:year=%Y}/\${_time:month=%m}/\${_time:date=%d}/	/my/path/year=2020/month=
/my/path/\${_time:%Y-%m-%d}/	/my/path/2020-05-20/

Destinations

Cribl LogStream can send data to various Destinations, including Splunk, Kafka, Kinesis, InfluxDB, Snowflake, Databricks, TCP JSON, and many others.



Streaming Destinations

Destinations that accept events in real time are referred to as streaming Destinations:

- Splunk Single Instance
- Splunk Load Balanced
- Splunk HEC
- AWS Kinesis Streams
- AWS CloudWatch Logs
- AWS SQS
- Elasticsearch
- Honeycomb
- TCP JSON
- Syslog
- Kafka
- Azure Event Hubs
- Azure Monitor Logs
- StatsD
- StatsD Extended
- Graphite

- SNMP Trap
- InfluxDB
- New Relic
- Wavefront
- Sumo Logic
- Datadog
- SignalFx

Non-Streaming Destinations

Destinations that accept events in groups or batches are referred to as nonstreaming Destinations:

- S3 Compatible Stores
- Filesystem/NFS
- MinIO
- Azure Blob Storage
- Google Cloud Storage
 - i The S3 Compatible Stores Destination can be adapted to send data to downstream services like Databricks and Snowflake, for which LogStream currently has no preconfigured Destination. For details, please contact Cribl Support.

Other Destinations

LogStream also provides these special-purpose Destinations:

- Output Router: Flexible "meta-destination." Here, you can configure rules that route data to multiple outputs.
- DevNull: An output that simply drops events. Preconfigured and active when you install LogStream, so it requires no configuration. Useful for testing.
- Default: Here, you can specify a default output from among your configured Destinations.

How Does Non-Streaming Delivery Work

Cribl LogStream uses a staging directory in the local filesystem to format and write outputted events before sending them to configured Destinations. After a set of conditions is met – typically file size and number of files, further details below – data is compressed and then moved to the final Destination.

An inventory of open, or in-progress, files is kept in the staging directory's root, to avoid having to walk that directory at startup. This can get expensive if staging is also the final directory. At startup, Cribl LogStream will check for any leftover files in progress from prior sessions, and will ensure that they're moved to their final Destination. The process of moving to the final Destination is delayed after startup (default delay: 30 seconds). Processing of these files is paced at one file per service period (which defaults to 1 second).

Batching Conditions

Several conditions govern when files are closed and rolled out:

- 1. File reaches its configured maximum size.
- 2. File reaches its configured maximum open time.
- 3. File reaches its configured maximum idle time.

If a new file needs to be open, Cribl LogStream will enforce the maximum number of open files, by closing files in the order in which they were opened.

Data Delivery

Data is delivered to all Destinations on an at-least-once basis. When a Destination is unreachable, there are three possible behaviors:

- **Block** Cribl LogStream will block incoming events.
- **Drop** Cribl LogStream will drop events addressed to that Destination.
- Queue Cribl LogStream will Persistent-Queue events to that Destination.

You can configure the desired behavior through a Destination's **Backpressure Behavior** option. If this option is not present, Cribl LogStream's default behavior is to **Block**.

Configuring Destinations

For each Destination **type**, you can create multiple definitions, depending on your requirements.

To configure Destinations, select **Data > Destinations**, select the desired type from the tiles or the left menu, then click **+ Add New**.

Output Router

Output Routers are meta-destinations that allow for output selection based on rules. Rules are evaluated in order, top->down, with the first match being the winner.

Configuring Cribl LogStream to Send to an Output Router

Select **Data > Destinations**, then select **Output Router** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Output Router > New Destination** modal, which provides the following fields.

Router name: Enter a unique name to identify this Router definition.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Rules: A list of event routing rules. Each provides the following settings:

- Filter expression: JavaScript expression to select events to send to output.
- Output: Output to send matching events to.
- **Description**: Optionally, enter a description of this rule's purpose.
- **Final**: Flag that controls whether to stop the event from being checked against other rules lower in the stack. Defaults to Yes.

Notes

- An Output Router cannot reference another. This is by design, so as to avoid cycles.
- Events that do not match any of the rules are dropped. Use a catchall rule to change this behavior.

- No post-processing (conditioning) can be done here. Use Pre-Processing Pipelines on the Source tier.
- Data can be cloned by toggling the Final flag to No . (The default is Yes , i.e., no cloning.)

Example

Scenario:

- Send all events where host starts with 66 to Destination San Francisco.
- From the rest of the events:
 - Send all events with method field POST or GET to both Seattle and Los Angeles (i.e., clone).
- Send the remaining events to New York City.

Router Name: router66

Filter Expression	Output	Final
host.startsWith('66')	San Francisco	Yes
method='POST' method='GET	Seattle	No
method='POST' method='GET'	Los Angeles	Yes
true	New York	Yes

Splunk Single Instance

Splunk Enterprise is a streaming Destination type.

Configuring Cribl LogStream to Output to Splunk Destinations

Select **Data > Destinations**, then select **Splunk > Single Instance** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Single Instance > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Splunk Single Instance definition.

Address: Hostname of the Splunk receiver.

Port: The port number on the host.

Nested field serialization: Specifies how to serialize nested fields into indextime fields. Defaults to None.

Throttling: Throttle rate in bytes per second. Multiple byte units such as KB, MB, GB etc. are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttle engaged, excesses data will be dropped only if Backpressure Behavior is set to drop, and blocked for all other settings.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Output multi metrics: Toggle to Yes to output multiple-measurement metric data points. (Supported in Splunk 8.0 and above, this format enables sending multiple metrics in a single event, improving the efficiency of your Splunk capacity.)

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

i Single .pem File

If you have a **single** .pem file containing cacert, key, and cert sections, enter it in all of these fields above: **CA certificate path**, **Private key path** (**mutual auth**), and **Certificate path** (**mutual auth**).

Timeout Settings

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (in milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Notes about Forwarding to Splunk

Data sent to Splunk is not compressed.

- If events have a Cribl LogStream internal field called __criblMetrics , they'll be forwarded to Splunk as metric events.
- If events do **not** have a _raw field, they'll be serialized to JSON prior to sending to Splunk.

Splunk Load Balanced

Splunk is a streaming Destination type, and with the **Splunk Load Balanced** output, you can load-balance data out to multiple Splunk receivers.

How Does Load Balancing Work

Cribl LogStream will attempt to load-balance outbound data as fairly as possibly across all receivers. Data is sent by all worker processes to all receivers simultaneously, and the amount sent to each receiver depends on these parameters:

- 1. Respective destination weight.
- 2. Respective destination historical data.

By default, historical data is tracked for 300s. LogStream uses this data to influence the traffic sent to each destination, to ensure that differences decay over time, and that total ratios converge towards configured weights.

Example

Suppose we have two receivers, A and B, each with weight of 1 (i.e., they are configured to receive equal amounts of data). Suppose further that the load-balance stats period is set at the default 300s and – to make things easy – for each period, there are 200 events of equal size (Bytes) that need to be balanced.

Interval	Time Range	Events to be dispensed
1	time=0s> time=300s	200

Both A and B start this interval with 0 historical stats each.

Let's assume that, due to various circumstances, 200 events are "balanced" as follows:

A = 120 events and B = 80 events – a difference of 40 events and a ratio of 1.5:1.

Interval	Time Range	Events to be dispensed

At the beginning of interval 2, the load-balancing algorithm will look back to the previous interval stats and carry **half** of the receiving stats forward. I.e., receiver A will start the interval with **60** and receiver B with **40**. To determine how many events A and B will receive during this next interval, LogStream will use their weights and their stats as follows:

Total number of events: events to be dispensed + stats carried forward = 200 + 60 + 40 = 300.

Number of events per each destination (weighed): 300/2 = 150 (they're equal, due to equal weight).

Number of events to send to each destination A: 150 - 60 = 90 and B: 150 - 40 = 110.

Totals at end of interval 2: A=120+90=210, B=80+110=190, a difference of 20 events and a ratio of 1.1:1.

Over the subsequent intervals, the difference becomes exponentially less pronounced, and eventually insignificant. Thus, the load gets balanced fairly.

Configuring Cribl LogStream to Load-Balance to Multiple Splunk Destinations

To configure load balancing, first select **Data > Destinations**, then select **Splunk > Load Balanced** from the **Data Destinations** page's tiles or left menu. Then click **Add New** to open the **Load Balanced > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Splunk LB Destination definition.

Indexer Discovery: When toggled to Yes, enables automatic discovery of indexers in an indexer clustering environment. See Indexer Discovery for the resulting UI options displayed below. When set to No (the default), displays the Destinations section below.

Exclude current host IPs: Exclude all IPs of the current host from the list of any resolved hostnames. Defaults to Yes.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block . When toggled to Persistent Queue , adds the Persistent Queue Settings section (left tab) to the modal.

Destinations

The **Destinations** section appears only when **Indexer discovery** is set to No. Here, you specify a known set of Splunk receivers on which to load-balance data.

Click + Add Destination to specify more receivers on new rows. Each row provides the following fields:

- Address: Hostname of the Splunk receiver. Optionally, you can paste in a comma-separated list, in <host>:<port> format.
- Port: Port number to send data to.
- TLS: Whether to inherit TLS configs from group setting, or disable TLS.
 Defaults to inherit .
- TLS servername: Servername to use if establishing a TLS connection. If not specified, defaults to connection host (if not an IP). Otherwise, uses the global TLS settings.
- Load weight: The weight to apply to this Destination for load-balancing purposes.

Indexer Discovery

Toggling the **Indexer Discovery** toggle to Yes displays the following fields instead of the **Destinations** section:

Site: Clustering site from which indexers need to be discovered. In the case of a single site cluster, default is the default entry.

Cluster Master URI: Full URI of Splunk Cluster Master, in the format: scheme://host:port.(Worker Nodes normally access the Cluster Master on port 8089 to get the list of currently online indexers.)

Auth token: Authentication token required to authenticate to Cluster Master for indexer discovery.

Refresh period: Time interval (in seconds) between two consecutive fetches of indexer list from Cluster Master. Defaults to 60.

Enabling Cluster Master Authentication

To enable token authentication on the Splunk Cluster Master, you can find complete instructions in Splunk's Enable or Disable Token Authentication documentation. The following capabilites are required:

list_indexer_cluster and list_indexerdiscovery.

For details on creating the token, see Splunk's Create Authentication Tokens topic – especially its section on how to Configure Token Expiry and "Not Before" Settings.

⚠ Be sure to give the token an **Expiration** setting well in the future, whether you use **Relative Time** or **Absolute Time**. Otherwise, the token will inherit Splunk's default expiration time of +30d (30 days in the future), which will cause indexer discovery to fail.

If you have a failover site configured on Splunk's Cluster Master, Cribl respects this configuration, and forwards the data to the failover site in case of site failure.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

i Single PEM File

If you have a **single** .pem file containing cacert, key, and cert sections, enter this file's path in all of these fields above:

CA certificate path, Private key path (mutual auth), and Certificate path (mutual auth).

Timeout Settings

- **Connection timeout**: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000 ms.
- Write timeout: Amount of time (in milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000 ms.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Output Multi Metrics: Toggle this slider to Yes to output multiple-measurement metric data points. (Supported in Splunk 8.0 and above, this format enables sending multiple metrics in a single event, improving the efficiency of your Splunk capacity.)

Minimize in-flight data loss: If set to Yes (the default), LogStream will check whether the indexer is shutting down and, if so, stop sending data. This helps minimize data loss during shutdown.

DNS resolution period (seconds): Re-resolve any hostnames after each interval of this many seconds, and pick up destinations from A records. Defaults to 600 seconds.

Load balance stats period (seconds): Lookback traffic history period. Defaults to 300 seconds. (Note that If multiple receivers are behind a hostname – i.e., multiple A records – all resolved IPs will inherit the weight of the host, unless each IP is specified separately. In Cribl LogStream load balancing, IP settings take priority over those from hostnames.)

Max connections: Constrains the number of concurrent indexer connections, per Worker Process, to limit memory utilization. If set to a number > 0, then on every DNS resolution period (or indexer discovery), LogStream will randomly select this subset of discovered IPs to connect to. LogStream will rotate IPs in future resolution periods – monitoring weight and historical data, to ensure fair load balancing of events among IPs.

Nested field serialization: Specifies whether and how to serialize nested fields into index-time fields. Select None (the default) or JSON.

Throttling: Throttle rate, in bytes per second. Multiple byte units such as KB, MB, GB, etc., are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttling is engaged, excess data will be dropped only if **Backpressure behavior** is set to **Drop events**. (Data will be blocked for all other **Backpressure behavior** settings.)

SSL Configuration for Splunk Cloud – Special Note

To connect to Splunk Cloud, you **might** need to extract the private and public key from the Splunk-provided Splunk Cloud Certificate, which is typically bundled in an app. Use the following steps:

Step 1. Test connectivity to Splunk Cloud, using the Root CA certificate:

```
openssl s_client -CApath path_to_ca.pem -connect
hostnameToSplunkCloud:9997
```

Step 2. Extract the Private key from the Splunk Cloud Certificate. At the prompt, you will need the sslPassword value from the outputs.conf file bundled with the Splunk Cloud app. Using Elliptic Curve keys:

```
openssl ec -in path_to_server_cert.pem -out private.pem
```

If you are using RSA keys, instead use:

```
openssl rsa -in path_to_server_cert.pem -out private.pem
```

Step 3. Extract the Public Key for the Server Certificate:

```
openssl x509 -in path_to_server_cert.pem -out server.pem
```

Step 4. In the LogStream Destination's **TLS Settings (Client Side)** section, enter the following:

• CA Certificate Path: Path to CA Certificate.

- Private Key Path (mutual auth): Path to private.pem (Step 2 above).
- Certificate Path (mutual auth): Path to server.pem (Step 3 above).

Notes About Forwarding to Splunk

- Data sent to Splunk is not compressed.
- If events have a LogStream internal field called __criblMetrics , they'll be forwarded to Splunk as metric events.
- If events do **not** have a _raw field, they'll be serialized to JSON prior to sending to Splunk.

Splunk HEC

Splunk HEC is a streaming Destination type. In a typical deployment, Cribl LogStream will be installed/co-located in a Splunk heavy forwarder. If this output is enabled, it can send data out to a Splunk HEC (HTTP Event Collector) destination through the event endpoint.

Configuring Cribl LogStream to Output to Splunk HEC Destinations

Select **Data > Destinations**, then select **Splunk > HEC** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **HEC > New Destination** modeal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Splunk HEC definition.

Splunk HEC endpoint: URL of a Splunk HEC endpoint to send events to (e.g., http://myhost.example.com:8088/services/collector/event).

HEC auth token: Splunk HEC authentication token.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the Backpressure behavior is set to Persistent Queue.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5. Each request can potentially hit a different HEC receiver.

Max body size (KB): Maximum size, in KB, of the request body. Defaults to 4096. Lowering the size can potentially result in more parallel requests and also cause outbound requests to be made sooner.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured Max body size. Defaults to 1.

- Retries happen on this flush interval.
 - Any HTTP response code in the 2xx range is considered success.
 - Any response code in the 5xx range is considered a retryable error, which will not trigger Persistent Queue (PQ) usage.
 - Any other response code will trigger PQ (if PQ is configured as the Backpressure behavior).

Extra HTTP headers: Click + **Add Header** to add **Name/Value** pairs to pass as additional HTTP headers.

Next processing queue: Specify the next Splunk processing queue to send the events to, after HEC processing. Defaults to indexQueue.

Default_TCP_ROUTING: Specify the value of the _TCP_ROUTING field for events that do not have _ctrl._TCP_ROUTING set. Defaults to nowhere .

i This is useful only when you expect the HEC receiver to route this data on to another destination.

Output multi metrics: Toggle to Yes to output multiple-measurement metric data points. (Supported in Splunk 8.0 and above, this format enables sending multiple metrics in a single event, improving the efficiency of your Splunk capacity.)

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new connection will be reattempted. This is to prevent sticking to a particular Destination when there is a constant flow of events.
- If the server does not support keepalives or if the server closes a pooled connection while idle – a new connection will be established for next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

Page 1028 of 1257

S3 Compatible Stores

S3 is a non-streaming Destination type. Cribl LogStream does **not** have to run on AWS in order to deliver data to S3.

Stores that are S3-compatible will also work with this Destination type. For example, the S3 Destination can be adapted to send data to services like Databricks and Snowflake, for which LogStream currently has no preconfigured Destination. For these integrations, please contact Cribl Support.

Configuring Cribl LogStream to Output to S3 Destinations

Select Data > Destinations, then select Amazon > S3 from the Data Destinations page's tiles or left menu. Click Add New to open the S3 > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this S3 definition.

S3 bucket name: Name of the destination S3 Bucket. This value can be a constant, or a JavaScript expression that can be evaluated only at init time. E.g., referencing a Global Variable: myBucket-\${C.vars.myVar}.

i Event-level variables are not available for JavaScript expressions. This is because the bucket name is evaluated only at Destination initialization. If you want to use event-level variables in file paths, Cribl recommends specifying them in the Partitioning Expression field (described below), because this is evaluated for each file.

Region: Region where the S3 bucket is located.

Staging location: Filesystem location in which to locally buffer files before compressing and moving to final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: When set to Yes (the default), adds the Output ID field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:

- Clone the Destination.
- Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with **Add output ID** enabled.

Key prefix: Prefix to add to files before uploading. This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. If left blank, Cribl LogStream will fall back to event.__partition . Defaults to `\${host}/\${sourcetype}`. Partitioning by time is also possible, e.g., `\${host}/\${C.Time.strftime(_time, '%Y-%m-%d')}/\${sourcetype}`

Data format: Format of the output data. Defaults to JSON.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Select the data compression format to use before moving data to final destination. Defaults to none. Cribl recommends setting this to gzip.

Backpressure behavior: Select whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

Authentication

Authentication method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for S3: Whether to use Assume Role credentials to access S3. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports <code>c*</code> wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Endpoint: S3 service endpoint. If empty, the endpoint will be automatically constructed from the region.

Object ACL: Object ACL (Access Control List) to assign to uploaded objects.

Storage class: Select a storage class for uploaded objects. Defaults to Standard.

Server side encryption: Server side encryption type for uploaded objects. Defaults to none. This is used for enabling encryption for data at rest.

Signature version: Signature version to use for signing S3 requests. Defaults to v4.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location.

Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Amazon S3 Permissions

The following permissions are needed to write to an Amazon S3 bucket:

s3:GetObject

s3:ListBucket

s3:GetBucketLocation

s3:PutObject

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

• __partition

Kinesis Streams

Cribl LogStream can output events to **Amazon Kinesis Data Streams** records of up to 1MB uncompressed. Cribl LogStream does **not** have to run on AWS in order to deliver data to a Kinesis Data Stream.

Configuring Cribl LogStream to Output to Amazon Kinesis Data Streams

Select **Data > Destinations**, then select **Amazon > Kinesis** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Kinesis > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Kinesis definition.

Stream name: Enter the name of the Kinesis Data Stream to which to send events.

Region: Select the AWS Region where the Kinesis Data Stream is located.

Endpoint: Kinesis Stream service endpoint. If empty, the endpoint will be automatically constructed from the region.

Signature version: Signature version to use for signing Kinesis stream requests. Defaults to v4.

Put request concurrency: Maximum number of ongoing put requests before blocking. Defaults to 5.

Max record size (KB, uncompressed): Maximum size of each individual record before compression. For non-compressible data, 1MB (the default) is the maximum recommended size.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the Backpressure behavior is set to Persistent Queue.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Authentication

Authentication method: Select an AWS authentication method.

- Auto: This default option uses the environment variables AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env. AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for Kinesis Streams: Whether to use Assume Role credentials to access Kinesis Streams. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes cribl_pipe (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Format

Currently, outputted events use the following record format:

- Header line containing information about the payload (currently supports one type, as shown below).
- Newline-Delimited JSON (that is, each Kinesis record will contain multiple events, in **ndjson** format).

Record payloads (including header and body) will be gzip-compressed, and then Kinesis will base64-encode them.

Sample Kinesis Record

```
{"format":"ndjson","count":8,"size":3960}
{"_raw":"07-03-2018 18:33:51.136 -0700 ERROR TcpOutputFd - Read error. Connection reset by peer'
{"_raw":"07-03-2018 18:33:51.136 -0700 INFO TcpOutputProc - Connection to 127.0.0.1:10000 close
```

CloudWatch Logs

Cribl LogStream supports sending data to Amazon CloudWatch Logs. This is a streaming Destination type. Cribl LogStream does **not** have to run on AWS in order to deliver data to CloudWatch Logs.

Configuring Cribl LogStream to Output to Amazon CloudWatch Logs

Select **Data > Destinations**, then select **Amazon > CloudWatch Logs** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **CloudWatch Logs > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this CloudWatch definition.

Log group name: CloudWatch log group to associate events with.

Log stream prefix: Prefix for CloudWatch log stream name. This prefix will be used to generate a unique log stream name per Cribl LogStream instance. (E.g., myStream_myHost_myOutputId .)

Region: AWS region where the CloudWatch Logs group is located.

Endpoint: CloudWatch Logs service endpoint. If empty, defaults to AWS' Region-specific endpoint. Otherwise, use this field to point to a CloudWatchLogs-compatible endpoint.

Signature version: Signature version to use for signing CloudWatch Logs requests. Defaults to $\ v4$.

Max queue size: Maximum number of queued batches before blocking. Defaults to 5.

Max record size (KB, uncompressed): Maximum size of each individual record before compression. For non-compressible data, 1MB (the default) is the maximum recommended size.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Authentication

Authentication Method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

 API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials. Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for CloudWatch Logs: Whether to use Assume Role credentials to access CloudWatch Logs. Defaults to No.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

SQS

Cribl LogStream supports sending events to Amazon Simple Queuing Service.

Configuring Cribl LogStream to Send Data to Amazon SQS

Select **Data > Destinations**, then select **Amazon > SQS** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **SQS > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this SQS Destination.

Queue name: The name, URL, or ARN of the SQS queue to send events to.
This value must be a JavaScript expression (which can evaluate to a constant), enclosed in single quotes, double quotes, or backticks. To specify a non-AWS URL, use the format: '{url}/<queueName>' . (E.g., ':port/<myQueueName>' .)

Message group ID: This parameter applies only to queues of type FIFO. Enter the tag that specifies that a message belongs to a specific message group. (Messages belonging to the same message group are processed in FIFO order.) Defaults to cribl. Use event field __messageGroupId to override this value.

Create queue: Specifies whether to create the queue if it does not exist. Defaults to Yes.

Region: Region where SQS queue is located.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the Backpressure behavior is set to Persistent Queue. Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Authentication

Authentication Method: Select an AWS authentication method.

- Auto: This default option uses the environment variables
 AWS_ACCESS_KEY_ID and AWS_SECRET_ACCESS_KEY, or the attached IAM
 role. Works only when running on AWS.
- Manual: You must select this option when not running on AWS.

The Manual option exposes these corresponding additional fields:

- API key: Enter your AWS API key. If not present, will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.
- Secret key: Enter your AWS secret key. If not present, will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Assume Role

Enable for SQS: Whether to use Assume Role credentials to access SQS. Defaults to No.

AWS account ID: Enter the SQS queue owner's AWS account ID. Leave empty if the SQS queue is in the same AWS account where this LogStream instance is located.

AssumeRole ARN: Enter the Amazon Resource Name (ARN) of the role to assume.

External ID: Enter the External ID to use when assuming role.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Endpoint: SQS service endpoint. If empty, the endpoint will be automatically constructed from the region.

Signature version: Signature version to use for signing SQS requests. Defaults to v4.

Max queue size: Maximum number of queued batches before blocking. Defaults to 100.

Max record size (KB): Maximum size of each individual record. Per the SQS spec, the maximum allowed value is 256 KB. (the default).

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Max concurrent requests: The maximum number of in-progress API requests before backpressure is applied. Defaults to 10.

SQS Permissions

The following permissions are needed to write to an SQS queue:

- sqs:ListQueues
- sqs:SendMessage
- sqs:SendMessageBatch
- sqs:CreateQueue
- sqs:GetQueueAttributes
- sqs:SetQueueAttributes
- sqs:GetQueueUrl

Internal Fields

Cribl LogStream uses a set of internal fields to assist in handling of data. These "meta" fields are **not** part of an event, but they are accessible, and functions can use them to make processing decisions.

Fields for this Destination:

- __messageGroupId
- __sqsMsgAttrs
- _sqsSysAttrs

Filesystem/NFS

Filesystem is a non-streaming Destination type that Cribl LogStream can use to output files to a local or a network-attached filesystem (NFS).

Configuring Cribl LogStream to Output to Filesystem Destinations

Select **Data > Destinations**, then select **Filesystem** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Filesystem > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Filesystem definition.

Output location: Final destination for the output files.

Staging location: Local filesystem location in which to buffer files before compressing and moving them to the final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: When set to Yes (the default), adds the **Output ID** field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

- ▲ For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:
 - Clone the Destination.
 - Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with **Add output ID** enabled.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. Defaults to `\${host}/\${sourcetype}`. If left blank, Cribl LogStream will fall back to event.__partition . Partitioning by time is also possible, e.g.: `\${host}/\${C.Time.strftime(_time, '%Y-%m-%d')}/\${sourcetype}`

Data format: Format of the output data. Defaults to json.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Data compression format used before moving to final destination. Default none . It is recommended that gzip is used.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this will be closed and moved to final output location. Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Backpressure Behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

- __partition
 - i To export events from an intermediate stage within a Pipeline to a file, see the Tee Function.

Elasticsearch

Cribl LogStream can send events to an Elasticsearch cluster using the Bulk API.

Configuring Cribl LogStream to Output to Elasticsearch

Select **Data > Destinations**, then select **Elasticsearch** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Elasticsearch > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Elasticsearch Destination definition.

Bulk API URL: URL of an Elasticsearch cluster to send events to.
(E.g., http://<myElasticCluster>:9200/_bulk .)

Index: Elasticsearch Index where to send events to. Note that this value can be overwritten by an event's __index field.

Type: Specify document type to use for events. Note that this value can be overwritten by an event's __type field.

Authentication enabled: Set to No by default. Toggle to Yes to enter a **Username** and **Password**.

Backpressure behavior: Specify whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (s): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Field Normalization

This Destination normalizes the following fields:

- _time becomes atimestamp atmillisecond resolution.
- host.name is set to host.

See also our Elasticsearch Source documentation's Field Normalization section.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this Destination:

- __id
- __type
- __index

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular destination when there is a constant flow of events.
- If the server does not support keepalives (or if the server closes a pooled connection while idle), a new connection will be established for the next request.
- When resolving the Destination's hostname, LogStream will pick the first IP
 in the list for use in the next connection. Round-robin DNS would help with
 event balancing.

Honeycomb

Cribl LogStream supports sending events to a Honeycomb dataset.

Configuring Cribl LogStream to Output to Honeycomb

Select **Data > Destinations**, then select **Honeycomb** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Honeycomb > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Honeycomb definition.

Dataset name: Name of the dataset to send events to. (E.g., iLoveObservabilityDataset .)

API Key: Team API Key to which the dataset belongs. (E.g., teamWilde .)

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id> . Defaults to

\$CRIBL HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular Destination when there is a constant flow of events.
- If the server does not support keepalives (or if the server closes a pooled connection while idle), a new connection will be established for the next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

TCP JSON

Cribl LogStream supports sending data over TCP in JSON format. **TCP JSON** is a streaming Destination type.

Configuring Cribl LogStream to Output in TCP JSON Format

Select **Data > Destinations**, then select **TCP JSON** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **TCP JSON > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Destination definition.

Address: Hostname of the receiver.

Port: Port number to connect to on the host.

Auth token: Optional authentication token to include as part of the connection header. Defaults to empty.

Compression: Codec to use to compress the data before sending. Defaults to None.

Throttling: Throttle rate in bytes per second. Multiple byte units such as KB, MB, GB etc. are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttle engaged, excesses data will be dropped only if Backpressure Behavior is set to drop, and blocked for all other settings.

Backpressure behavior: Specifies whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Autofill?: This setting is experimental.

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Timeout Settings

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish before retrying. Defaults to 10000.

Write timeout: Amount of time (in milliseconds) to wait for a write to complete before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Format

TCP JSON events are sent in newline-delimited JSON format, consisting of:

- 1. A header line. Can be empty, e.g.: {} . If **Auth Token** is enabled, the token will be included here as a field called authToken . In addition, if events contain common fields, they will be included here under fields .
- 2. A JSON event/record per line.

Example

See an example in our TCP JSON Source topic.

Syslog

Cribl LogStream supports sending of data over syslog via TCP. Syslog is a streaming Destination type.

i This Syslog Destination supports RFC 3164 and RFC 5424.

Configuring Cribl LogStream to output in **Syslog** format

Select **Data > Destinations**, then select **Syslog** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Syslog > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Syslog definition.

Protocol: The network protocol to use for sending out syslog messages. Defaults to TCP; UDP is also available.

Address: Address/hostname of the receiver.

Port: Port number to connect to on the host.

Facility: Default value for message facility. If set, will be overwritten by the value of __facility . Defaults to user .

Severity: Default value for message severity. If set, will be overwritten by the value of __severity . Defaults to notice .

App name: Default value for application name. If set, will be overwritten by the value of __appname . Defaults to Cribl .

Message format: The syslog message format supported by the receiver. Defaults to RFC3164.

Timestamp format: The timestamp format to use when serializing an event's time field. Defaults to Syslog.

Throttling: Throttle rate in bytes per second. Multiple byte units such as KB, MB, GB etc. are also allowed. E.g., 42 MB. Default value of 0 indicates no throttling. When throttle engaged, excesses data will be dropped only if Backpressure Behavior is set to drop, and blocked for all other settings.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled defaults to No . When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Timeout Settings

i These timeout settings apply only to the TCP protocol.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.

- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this destination:

- __priority
- __facility
- __severity
- __procid
- __appname
- __msgid
- __syslogout

Kafka

Cribl LogStream supports sending data to a Kafka topic. Kafka is a streaming Destination type.

Configuring Cribl LogStream to Output to Kafka

Select **Data > Destinations**, then select **Kafka** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Kafka > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Kafka definition.

Brokers: List of Kafka brokers to connect to. (E.g., localhost:9092.)

Topic: The topic on which to publish events. Can be overwritten using event's __topic field.

Acknowledgments: Select the number of required acknowledgments. Defaults to Leader.

Record data format: Format to use to serialize events before writing to Kafka. Defaults to JSON.

Compression: Codec to compress the data before sending to Kafka. Select None, Gzip, or Snappy.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

TLS Settings (Client Side)

Enabled defaults to No. When toggled to Yes:

Validate server certs: Reject certificates that are not authorized by a CA in the **CA certificate path**, or by another trusted CA (e.g., the system's CA). Defaults to No.

Server name (SNI): Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.

Certificate name: The name of the predefined certificate.

CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.

Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS . Use only if mutual auth is required.

Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS . **Use only if mutual auth is required**.

Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Authentication

Authentication parameters to use when connecting to brokers. Using TLS is highly recommended.

Enabled: Defaults to No. When toggled to Yes:

- **SASL** mechanism: Select the SASL (Simple Authentication and Security Layer) authentication mechanism to use,
- Username: The username for authentication.
- **Password**: The password for authentication.

Schema Registry

This section governs Kafka Schema Registry Authentication for AVRO-encoded data with a schema stored in the Confluent Schema Registry.

Enabled: defaults to No. When toggled to Yes:

- Schema registry URL: URL for access to the Confluent Schema Registry. (E.g., http://<hostname>:8081.)
- Default key schema ID: Used when __keySchemaIdOut is not present to transform key values. Leave blank if key transformation is not required by default.
- **Default value schema ID**: Used when __valueSchemaIdOut not present to transform _raw . Leave blank if value transformation is not required by default.
- TLS enabled: defaults to No . When toggled to Yes, displays the following TLS settings for the Schema Registry:

TLS Settings (Schema Registry)

- i These have the same format as the TLS Settings (Client Side) above.
- Validate server certs: Require client to reject any connection that is not authorized by a CA in the CA certificate path, or by another trusted CA (e.g., the system's CA). Defaults to No.
- **Server name (SNI)**: Server name for the SNI (Server Name Indication) TLS extension. This must be a host name, not an IP address.
- Certificate name: The name of the predefined certificate.
- CA certificate path: Path on client containing CA certificates (in PEM format) to use to verify the server's cert. Path can reference \$ENV_VARS.
- Private key path (mutual auth): Path on client containing the private key (in PEM format) to use. Path can reference \$ENV_VARS. Use only if mutual auth is required.
- Certificate path (mutual auth): Path on client containing certificates in (PEM format) to use. Path can reference \$ENV_VARS. Use only if mutual auth is required.
- Passphrase: Passphrase to use to decrypt private key.

Minimum TLS version: Optionally, select the minimum TLS version to use when connecting.

Maximum TLS version: Optionally, select the maximum TLS version to use when connecting.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (KB, uncompressed): Maximum size (KB) of each record batch before compression. Setting should be < message.max.bytes settings in Kafka brokers. Defaults to 768.

Max events per batch: Maximum number of events in a batch before forcing a flush. Defaults to 1000.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this Destination:

- __topicOut
- __key
- __headers

- __keySchemaIdOut
- __valueSchemaIdOut

Azure Blob Storage

Azure Blob Storage is a non-streaming Destination type. Cribl LogStream does not have to run on Azure in order to deliver data to it. Azure Data Lake Storage Gen2 (hierarchical namespace) is also supported.

Configuring Cribl LogStream to Output to Azure Blob Storage

Select **Data > Destinations**, then select **Azure > Azure Blob Storage** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Blob Storage > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Destination definition.

Connection string: Enter your Azure Connection String. The format is:

DefaultEndpointsProtocol=[http|https];AccountName=<your-account-name>;AccountKey=
<your-account-key>

A fictitious example, using Microsoft's recommended HTTPS option:

DefaultEndpointsProtocol=https;AccountName=storagesample;AccountKey=12345678 ... 32

If left blank, LogStream will fall back to: env.AZURE_STORAGE_CONNECTION_STRING.

☐ The Connection string field replaces the Account name and Account key fields provided in this Destination's UI prior to LogStream 2.4. If you configured an Azure Blob Storage Destination before upgrading to LogStream 2.4, those fields' values are now concatenated into the Connection string value.

In case values were concatenated incorrectly, the original field keys and values are retained in LogStream's configuration files (e.g.,

\$CRIBL_HOME/groups/<group-name>/cribl/outputs.yml). However, Cribl recommends simply re-entering the correct Azure **Connection string** here.

For further details, see Microsoft's Configure Azure Storage Connection Strings documentation.

Container name: Enter the container name. (A container organizes a set of blobs, similar to a directory in a file system.)

Create container: Defaults to No . Toggle to Yes to create the configured container in Azure Blob Storage if it does not already exist.

Blob prefix: Root directory to prepend to path before uploading.

Staging location: Local filesystem location in which to buffer files before compressing and moving them to the final destination. Cribl recommends that this location be stable and highperformance.

Add Output ID: When set to Yes (the default), adds the Output ID field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

- ▲ For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:
 - Clone the Destination.
 - Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with Add output ID enabled.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. Defaults to `\${host}/\${sourcetype}` If left blank, Cribl LogStream will fall back to event.__partition.

Data format: Format of the output data. Defaults to json.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Data compression format used before moving to final destination. Defaults to none. Cribl recommends setting to gzip.

Backpressure behavior: Whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes cribl_pipe (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max file size (MB): Maximum uncompressed output file size. Files reaching this size will be closed and moved to the final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location. Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Default: 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Default: 100.

i LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

• __partition

Azure Monitor Logs

Cribl LogStream supports sending of data over to Azure Monitor Logs. This is a streaming Destination type.

Configuring Cribl LogStream to Output to Azure Monitor Logs

Select **Data > Destinations**, then select **Azure > Monitor Logs** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Monitor Logs > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Azure Monitor Logs definition.

Workspace ID: Enter the Azure Log Analytics Workspace ID. (See **Workspace->Advanced settings** in the Azure Dashboard.)

Workspace key: Enter the Azure Log Analytics Workspace Primary or Secondary Shared Key. (In the Azure Dashboard, see **Workspace->Advanced settings**.)

Log type: The Record Type of events sent to this LogAnalytics workspace. Defaults to Cribl.

Resource ID: Resource ID of the Azure resource to associate the data with. This populates the _ResourceId property, and allows the data to be included in resource-centric queries. (Optional, but if this field is not specified, the data will not be included in resource-centric queries.)

Backpressure behavior: Whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Oueue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues .

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular Destination when there is a constant flow of events.
- If keepalives are not supported by the server (or if the server closes a
 pooled connection while idle), a new connection will be established for the
 next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

Azure Event Hubs

Cribl LogStream supports sending data to Azure Event Hubs. This is a streaming Destination type.

Configuring Cribl LogStream to Output to Azure Event Hubs

Select Data > Destinations, then select Azure > Event Hubs from the Data Destinations page's tiles or left menu. Click Add New to open the Event Hubs > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Azure Event Hubs definition.

Brokers: List of Event Hub Kafka brokers to connect to. (E.g., yourdomain.servicebus.windows.net:9093.) Find the hostname in Shared Access Policies, in the host portion of the primary or secondary connection string.

Event Hub name: The name of the Event Hub (a.k.a., Kafka Topic) on which to publish events. Can be overwritten using the __topicOut field.

Acknowledgments: Control the number of required acknowledgments. Defaults to Leader.

Record data format: Format to use to serialize events before writing to the Event Hub Kafka brokers. Defaults to JSON.

Compression: Codec to use to compress the data before sending it to Event Hub Kafka brokers. Defaults to Gzip.

Backpressure behavior: Whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

TLS Settings (Client Side)

Enabled Defaults to Yes.

Validate server certs: Defaults to No – and for Event Hubs, this must always be disabled.

Authentication

Authentication parameters to use when connecting to brokers. Using TLS is highly recommended.

Enabled: Defaults to Yes. (Toggling to No hides the remaining settings in this group.)

SASL mechanism: SASL (Simple Authentication and Security Layer) authentication mechanism to use, PLAIN is the only mechanism currently supported for Event Hub Kafka brokers.

Username: The username for authentication. For Event Hub, this should always be \$ConnectionString.

Password: Connection-string primary key or connection-string secondary key from the Event Hub workspace.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (KB, uncompressed): Maximum size (KB) of each record batch before compression. Setting should be < message.max.bytes settings in Kafka brokers. Defaults to 768.

Max events per batch: Maximum number of events in a batch before forcing a flush. Defaults to 1000.

Flush period (sec): Maximum time between requests. Low settings could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Fields for this Destination:

•	topicOut
•	key
•	headers
•	keySchemaIdOut
•	valueSchemaIdOut

Google Cloud Storage

Google Cloud Storage is a non-streaming Destination type.

Configuring Cribl LogStream to Output to Google Cloud Storage Destinations

Select **Data > Destinations**, then select **Google Cloud > Cloud Storage** from the **Data Destinations** page's tiles or left menu.

Next, click **Add New** to open the **Cloud Storage > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Cloud Storage definition.

Bucket name: Name of the destination bucket. This value can be a constant. or a JavaScript expression that can be evaluated only at init time. E.g., referencing a Global Variable: myBucket-\${C.vars.myVar}.

Region: Region where the bucket is located.

Staging location: Filesystem location in which to locally buffer files before compressing and moving to final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: Whether to append output's ID to staging location. Defaults to Yes .

Key prefix: Prefix to add to files before uploading. This value can be a constant or a JavaScript expression.

Partitioning expression: JavaScript expression to define how files are
partitioned and organized. If left blank, Cribl LogStream will fall back to
 event.__partition . Defaults to `\${host}/\${sourcetype}` . Partitioning
by time is also possible, e.g., `\${host}/\${C.Time.strftime(_time, '%Y-%m-%d')}/\${sourcetype}`

Data format: Format of the output data. Defaults to JSON.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Select the data compression format to use before moving data to final destination. Defaults to none. Cribl recommends setting this to gzip.

Backpressure behavior: Select whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

Authentication

Authentication is via HMAC (Hash-based Message Authentication Code). To create a key and secret, see Google Cloud's Managing HMAC Keys for Service Accounts documentation

Access key: Enter the HMAC access key.

Secret key: Enter the HMAC secret.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports <code>c*</code> wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Endpoint: The Google Cloud Storage service endpoint. Typically, there is no reason to change the default https://storage.googleapis.com endpoint.

Object ACL: Select an Access Control List to assign to uploaded objects. Defaults to private.

Storage class: Select a storage class for uploaded objects.

Signature version: Signature version to use for signing requests. Defaults to $\,^{\vee4}$.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location.

Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the Max file open time (sec) conditions are met.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

__partition

StatsD

Cribl LogStream supports sending data to a StatsD Destination. This is a streaming Destination type.

Configuring Cribl LogStream to Output via StatsD

While on the **Data Destinations** page, select **Metrics > StatsD** from the tiles or the left menu, then click **Add New**. The resulting **New StatsD destination** pane contains the following fields.

Select Data > Destinations, then select Metrics > StatsD from the

Data Destinations page's tiles or left menu. Click Add New to open the StatsD >

New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this StatsD definition.

Destination protocol: Protocol to use when communicating with the Destination. Defaults to UDP.

Host: The hostname of the Destination.

Port: Destination port. Defaults to 8125.

i The next two settings apply only to the TCP protocol, and are not displayed for UDP.

Throttling: Rate (in bytes per second) at which at which to throttle while writing to an output. Also takes numerical values in multiples of bytes (KB, MB, GB, etc.). Default value of 0 indicates no throttling.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

- i This section is displayed only for TCP, and only when the Backpressure behavior is set to Persistent Queue.
- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Timeout Settings

i These timeout settings apply only to the TCP protocol, and are not displayed for UDP.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

• cribl_host - LogStream Node that processed the event.

- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (bytes): Used when Protocol is UDP. Specifies the maximum size of packets sent to the Destination. (Also known as the MTU – maximum transmission unit – for the network path to the Destination system.) Defaults to 512.

Flush period (sec): Used when Protocol is TCP. Specifies how often buffers should be flushed, sending records to the Destination. Defaults to 1.

StatsD Extended

Cribl LogStream supports sending data to a StatsD Destination. This is a streaming Destination type.

Configuring Cribl LogStream to Output via StatsD Extended

Select Data > Destinations, then select Metrics > StatsD Extended from the Data Destinations page's tiles or left menu. Click Add New to open the StatsD Extended > New Destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this StatsD Extended definition.

Destination protocol: Protocol to use when communicating with the Destination. Defaults to UDP.

Host: The hostname of the Destination.

Port: Destination port. Defaults to 8125.

i The next two settings apply only to the TCP protocol, and are not displayed for UDP.

Throttling: Rate (in bytes per second) at which at which to throttle while writing to an output. Also takes numerical values in multiples of bytes (KB, MB, GB, etc.). Default value of 0 indicates no throttling.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

- i This section is displayed only for TCP, and only when the Backpressure behavior is set to Persistent Queue.
- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Timeout Settings

i These timeout settings apply only to the TCP protocol, and are not displayed for UDP.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

cribl_host - LogStream Node that processed the event.

- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (bytes): Used when Protocol is UDP. Specifies the maximum size of packets sent to the Destination. (Also known as the MTU – maximum transmission unit – for the network path to the Destination system.) Defaults to 512.

Flush period (sec): Used when Protocol is TCP. Specifies how often buffers should be flushed, sending records to the Destination. Defaults to 1.

Graphite

Cribl LogStream supports sending data to a Graphite backend Destination. This is a streaming Destination type.

Configuring Cribl LogStream to Output to a Graphite Backend

Select **Data > Destinations**, then select **Metrics > Graphite** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Graphite > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Graphite definition.

Destination protocol: Protocol to use when communicating with the Destination. Defaults to UDP.

Host: The hostname of the Destination.

Port: Destination port. Defaults to 8125.

i The next two settings apply only to the TCP protocol, and are not displayed for UDP.

Throttling: Rate (in bytes per second) at which at which to throttle while writing to an output. Also takes numerical values in multiples of bytes (KB, MB, GB, etc.). Default value of 0 indicates no throttling.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

- i This section is displayed only for TCP, and only when the Backpressure behavior is set to Persistent Queue.
- Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.
- Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.
- Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.
- Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Timeout Settings

i These timeout settings apply only to the TCP protocol, and are not displayed for UDP.

Connection timeout: Amount of time (in milliseconds) to wait for the connection to establish, before retrying. Defaults to 10000.

Write timeout: Amount of time (milliseconds) to wait for a write to complete, before assuming connection is dead. Defaults to 60000.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

cribl_host - LogStream Node that processed the event.

- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Max record size (bytes): Used when Protocol is UDP. Specifies the maximum size of packets sent to the Destination. (Also known as the MTU – maximum transmission unit – for the network path to the destination system.) Defaults to 512.

Flush period (sec): Used when Protocol is TCP. Specifies how often buffers should be flushed, sending records to the Destination. Defaults to 1.

SNMP Trap

Cribl LogStream supports forwarding of SNMP Traps out.

Configuring Cribl LogStream to Forward to SNMP Traps

While on the **Data Destinations** page, select **SNMP Trap** from the tiles or the left menu, then click **Add New**. The resulting **SNMP Trap > New Destination** modal, which provides the following fields.

Select **Data > Destinations**, then select **SNMP Trap** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **New SNMP destination** pane, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this SNMP Trap definition.

SNMP Trap destinations: One or more SNMP destinations to forward traps to.

- Address: Destination host.
- Port: Destination port. Defaults to 162.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.

• cribl_output - LogStream Destination that processed the event.

Considerations for Working with SNMP Traps Data

- It's possible to work with SNMP metadata (i.e., we'll decode the packet).
 Options include dropping, routing, etc. However, packets cannot be modified and sent to another SNMP Destination.
- SNMP packets can be forwarded to non-SNMP Destinations (e.g., Splunk, Syslog, S3, etc.).
- SNMP packets can be forwarded to other SNMP Destinations. However, the contents of the incoming packet cannot be modified i.e., we'll forward the packets verbatim as they came in.
- Non-SNMP input data cannot be sent to SNMP Destinations.

InfluxDB

Cribl LogStream supports sending data to InfluxDB.

Configuring Cribl LogStream to Output to InfluxDB

Select **Data > Destinations**, then select **InfluxDB** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **InfluxDB > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this InfluxDB definition.

Write API URL: URL of an InfluxDB cluster to send events to. (E.g., http://localhost:8086/write.)

Database name: The database on which to write data points.

Timestamp precision: Sets the precision for the supplied UNIX time values. Defaults to Milliseconds .

Dynamic value fields: When enabled, LogStream will pull the value field from the metric name. (E.g., db.query.user will use db.query as the measurement and user as the value field). Defaults to Yes.

Value field name: Name of the field in which to store the metric when sending to InfluxDB. This will be used as a fallback if dynamic name generation is enabled but fails. Defaults to value.

Authentication enabled: Set to No by default. Toggle to Yes to enter a **Username** and **Password**.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to $1 \, \text{MB}$.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

MinIO

MinIO is a non-streaming Destination type, to which Cribl LogStream can output objects.

Configuring Cribl LogStream to Output to MinIO Destinations.

Select **Data > Destinations**, then select **MinIO** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **MinIO > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this MinIO definition.

MinIO endpoint: MinIO service URL (e.g., http://minioHost:9000).

MinIO bucket name: Name of the destination MinIO bucket. This value can be a constant, or a JavaScript expression that can be evaluated only at init time. E.g., referencing a Global Variable: myBucket-\${C.vars.myVar} . Ensure that the bucket already exists, otherwise MinIO will generate "bucket does not exist" errors.

API key: If left blank, LogStream will fall back to env.AWS_ACCESS_KEY_ID, or to the metadata endpoint for IAM credentials.

Secret key: If left blank, Cribl LogStream will fall back to env.AWS_SECRET_ACCESS_KEY, or to the metadata endpoint for IAM credentials.

Staging location: Filesystem location in which to locally buffer files before compressing and moving to final destination. Cribl recommends that this location be stable and high-performance.

Add Output ID: When set to Yes (the default), adds the **Output ID** field's value to the staging location's file path. This ensures that each Destination's logs will write to its own bucket.

- ⚠ For a Destination originally configured in a LogStream version below 2.4.0, the Add Output ID behavior will be switched off on the backend, regardless of this slider's state. This is to avoid losing any files pending in the original staging directory, upon LogStream upgrade and restart. To enable this option for such Destinations, Cribl's recommended migration path is:
 - Clone the Destination.
 - Redirect the Routes referencing the original Destination to instead reference the new, cloned Destination.

This way, the original Destination will process pending files (after an idle timeout), and the new, cloned Destination will process newly arriving events with **Add output ID** enabled.

Key prefix: Prefix to apply to files/objects before uploading to the specified bucket. MinIO will display key prefixes as folders.

Partitioning expression: JavaScript expression to define how files are partitioned and organized. If left blank, Cribl LogStream will fall back to event.__partition. Defaults to `\${host}/\${sourcetype}`.

i LogStream's internal __partition field can be populated in multiple ways. The precedence order is: explicit Partitioning expression value -> \${host}/\${sourcetype} (default) Partitioning expression value -> user-defined event.__partition, set with an Eval Function (takes effect only where this Partitioning expression field is blank).

Data format: Format of the output data. Defaults to json.

File name prefix expression: The output filename prefix. Must be a JavaScript expression (which can evaluate to a constant), enclosed in quotes or backticks. Defaults to CriblOut.

Compress: Select the data compression format to use before moving data to final destination. Defaults to none. Cribl recommends setting this to gzip.

Backpressure behavior: Select whether to block or drop events when all receivers in this group are exerting backpressure. Defaults to Block.

How MinIO Composes File Names

The full path to a file consists of:

```
<bucket_name>/<keyprefix><partition_expression | __partition>
<file_name_prefix><filename>.<extension>
```

As an example, assume that the MinIO bucket name is bucket1, the Key prefix is aws, the Partitioning expression is `\${host}/\${sourcetype}`, the source is undefined, the File name prefix is the default CriblOut, and the Data format is json. Here, the full path as displayed in MinIO would have this form: /bucket1/aws/192.168.1.241/undefined/CriblOut-<randomstring>0.json

i Although MinIO will display the Key prefix and Partitioning expression values as folders, both are actually just part of the overall key name, along with the file name.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Region: Region where the MinIO service/cluster is located. Leave blank when using a containerized MinIO.

Object ACL: ACL (Access Control List) to assign to uploaded objects. Defaults to Private .

Storage class: Select a storage class for uploaded objects. Defaults to Standard.

Server side encryption: Server side encryption type for uploaded objects. Defaults to none.

Signature version: Signature version to use for signing MinIO requests. Defaults to v4.

Max file size (MB): Maximum uncompressed output file size. Files of this size will be closed and moved to final output location. Defaults to 32.

Max file open time (sec): Maximum amount of time to write to a file. Files open for longer than this limit will be closed and moved to final output location.

Defaults to 300.

Max file idle time (sec): Maximum amount of time to keep inactive files open. Files open for longer than this limit will be closed and moved to final output location. Defaults to 30.

Max open files: Maximum number of files to keep open concurrently. When exceeded, the oldest open files will be closed and moved to final output location. Defaults to 100.

i Cribl LogStream will close files when either of the Max file size (MB) or the

Max file open time (sec) conditions is met.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

Field for this Destination:

• __partition

New Relic

Cribl LogStream supports sending events to the New Relic Log API and the New Relic Metric API.

Configuring Cribl LogStream to Output to New Relic

Select Data > Destinations, then select New Relic from the Data Destinations page's tiles or left menu. Click Add New to open the New Relic > New destination modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this New Relic definition.

API key: Enter your New Relic Insert API key, as created on New Relic's **Insights API keys** page.

Region: Select which New Relic region endpoint to use.

Log type: Name of the logType to send with events. E.g., observability or access_log.

i This sets a default. Where a sourcetype is specified in an event, it will override this value.

Log message field: Name of the field to send as the log message value. If not specified, the event will be serialized and sent as JSON.

Fields: Additional metadata fields to (optionally) add, as **Name-Value** pairs.

- Name: Enter the metadata field name.
- Value: JavaScript expression to compute field's value, enclosed in single quotes, double quotes, or backticks. (Can evaluate to a constant.)
- Add Field: Click to add more metadata Name-Value pairs.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

 $\begin{tabular}{ll} \textbf{Compress}: Whether to compress the payload body before sending. Defaults to No . \end{tabular}$

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 1000 KB.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured **Max body size**. Defaults to 1 second.

Extra HTTP headers: Click **+ Add Header** to insert extra headers as **Name/Value** pairs.

Wavefront

Cribl LogStream supports sending events to Wavefront analytics.

Configuring Cribl LogStream to Output to Wavefront

Select **Data > Destinations**, then select **Wavefront** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Wavefront > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Wavefront definition.

Auth token: Wavefront API authentication token. For details, see Wavefront's Generating an API Token topic. Required.

Domain name: WaveFront domain name, e.g., longboard. Required.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id> . Defaults to

\$CRIBL HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

 $\begin{tabular}{ll} \textbf{Compress}: Whether to compress the payload body before sending. Defaults to No . \end{tabular}$

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured **Max body size**. Defaults to 1 second.

Extra HTTP headers: Click + **Add Header** to insert extra headers as **Name/Value** pairs.

Notes About Wavefront

For details on integrating with Wavefront, see these Wavefront resources:

- Direct Data Ingestion, and adjacent topics on Wavefront Proxies.
- Wavefront Data Format.

SignalFx

Cribl LogStream supports sending events to SignalFx.

Configuring Cribl LogStream to Output to SignalFx

Select **Data > Destinations**, then select **SignalFx** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **SignalFx > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this SignalFx definition.

Auth token: SignalFx API access token. For details, see SignalFx's Manage Tokens topic. Required.

Realm: SignalFx realm name (e.g., us0). Required.

Backpressure behavior: Select whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Whether to compress the payload body before sending. Defaults to No.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values can cause the payload size to be smaller than the configured **Max body size**. Defaults to 1 second.

Extra HTTP headers: Click + **Add Header** to insert extra headers as **Name/Value** pairs.

Notes About SignalFx

For details on integrating with SignalFx, see the SignalFx Developers Guide, with particular reference to the SignalFx HTTP Send Metrics Reference.

Sumo Logic

Cribl LogStream can send log and metric events to Sumo Logic over HTTP.

Configuring Cribl LogStream to Output to Sumo Logic

Select **Data > Destinations**, then select **Sumo Logic** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Sumo Logic > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Sumo Logic Destination definition.

API URL: Enter the URL of the Sumo Logic HTTP collector to which events should be sent. (E.g.,

https://endpoint6.collection.us2.sumologic.com/receiver/v1/http/<longhash>.)

Custom source name: Optionally, override the source name configured on the Sumo Logic HTTP collector. This value will be sent with events via the X-Sumo-Name HTTP header.

Custom source category: Optionally, override the source category configured on the Sumo Logic HTTP collector. This value will be sent with events via the X-Sumo-Category HTTP header.

Backpressure behavior: Optionally, override the source category configured on the Sumo Logic HTTP collector.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None; Gzip is also available.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress the payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (sec): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

If an event contains the internal field __criblMetrics , LogStream will send it to Sumo Logic as a metric event. Otherwise, LogStream will send it as a log event.

Notes on HTTP-based Outputs

- Cribl LogStream will attempt to use keepalives to reuse a connection for multiple requests. After 2 minutes of the first use, the connection will be thrown away, and a new one will be reattempted. This is to prevent sticking to a particular destination when there is a constant flow of events.
- If the server does not support keepalives (or if the server closes a pooled connection while idle), a new connection will be established for the next request.
- When resolving the Destination's hostname, LogStream will pick the first IP in the list for use in the next connection. Round-robin DNS would help with event balancing.

Datadog

Cribl LogStream can send log and metric events to Datadog. (Datadog supports metrics only of type gauge, counter, and rate via its REST API.)

LogStream sends events to the following Datadog endpoints in the US region. Use a DNS lookup to discover and include the corresponding IP addresses in your firewall rules' allowlist.

- Logs: https://http-intake.logs.datadoghq.com/v1/input
- Metrics: https://api.datadoghq.com/api/v1/series

Configuring Cribl LogStream to Output to Datadog

Select **Data > Destinations**, then select **Datadog** from the **Data Destinations** page's tiles or left menu. Click **Add New** to open the **Datadog > New Destination** modal, which provides the following fields.

General Settings

Output ID: Enter a unique name to identify this Destination definition.

API Key: Enter an API key available in your Datadog profile.

Send logs as: Specify the content type to use when sending logs. Defaults to application/json, where each log message is represented by a JSON object. The alternative text/plain option sends one message per line, with newline \n delimiters.

Message field: Name of the event field that contains the message to send. If not specified, LogStream sends a JSON representation of the whole event (regardless of whether **Send logs as** is set to JSON or plain text).

Source: Name of the source to send with logs. If you're sending logs as JSON objects (i.e., you've selected **Send logs as**: application/json), the event's source field (if set) will override this value.

Host: Name of the host to send with logs. If you're sending logs as JSON objects, the event's host field (if set) will override this value.

Service: Name of the service to send with logs. If you're sending logs as JSON objects, the event's __service field (if set) will override this value.

Tags: List of tags to send with logs (e.g., env:prod, env_staging:east).

Severity: Default value for message severity. If you're sending logs as JSON objects, the event's __severity field (if set) will override this value. Defaults to info; the drop-down offers many other severity options.

i Datadog uses the above five fields (source , host , __service , __severity , and tags) to enhance searches and UX.

Backpressure behavior: Specify whether to block, drop, or queue events when all receivers in this group are exerting backpressure. Defaults to Block.

Persistent Queue Settings

i This section is displayed when the **Backpressure behavior** is set to **Persistent Queue**.

Max file size: The maximum size to store in each queue file before closing it. Enter a numeral with units of KB, MB, etc. Defaults to 1 MB.

Max queue size: The maximum amount of disk space the queue is allowed to consume. Once this limit is reached, queueing is stopped, and data blocking is applied. Enter a numeral with units of KB, MB, etc.

Queue file path: The location for the persistent queue files. This will be of the form: your/path/here/<worker-id>/<output-id>. Defaults to \$CRIBL_HOME/state/queues.

Compression: Codec to use to compress the persisted data, once a file is closed. Defaults to None . Select Gzip to enable compression.

Processing Settings

Post-Processing

Pipeline: Pipeline to process data before sending the data out using this output.

System fields: A list of fields to automatically add to events that use this output. By default, includes <code>cribl_pipe</code> (identifying the LogStream Pipeline that processed the event). Supports wildcards. Other options include:

- cribl_host LogStream Node that processed the event.
- cribl_wp LogStream Worker Process that processed the event.
- cribl_input LogStream Source that processed the event.
- cribl_output LogStream Destination that processed the event.

Advanced Settings

Compress: Toggle this slider to Yes to compress log events' payload body before sending.

Request timeout: Amount of time (in seconds) to wait for a request to complete before aborting it. Defaults to 30.

Request concurrency: Maximum number of concurrent requests before blocking. This is set per Worker Process. Defaults to 5.

Max body size (KB): Maximum size of the request body. Defaults to 4096 KB.

Flush period (s): Maximum time between requests. Low values could cause the payload size to be smaller than its configured maximum. Defaults to 1.

Extra HTTP headers: Name/Value pairs to pass as additional HTTP headers.

Internal Fields

Cribl LogStream uses a set of internal fields to assist in forwarding data to a Destination.

If an event contains the internal field __criblMetrics , LogStream will send it to Datadog as a metric event. Otherwise, LogStream will send it as a log event.

You can use these fields to override outbound event values for log events:

- __service
- __severity

No internal fields are supported for metric events.

For More Information

You might find these Datadog references helpful:

- Submit Metrics
- Send Logs
- Metrics Types

DevNull

The DevNull Destination simply drops events. Cribl provides this as a basic output to test Pipelines and Routes.

Configuring Cribl LogStream to Forward to DevNull

DevNull requires no configuration: A DevNull Destination is preconfigured and active as soon as you install Cribl LogStream.

To verify this, select **Data > Destinations** from the top menu. On the resulting **Data Destinations** page, select **DevNull** from the tiles or the left menu. Look for the **Live** indicator at the top right.

Default

The **Default** Destination simply enables you to specify a default output from among your configured Destinations.

Configuring Cribl LogStream's Default Destination

Select **Data > Destinations**, then select **Default** from the **Data Destinations** page's tiles or left menu. From the resulting **Manage Default Destination** page, click anywhere on the default row to proceed.



Default Destination - click to configure

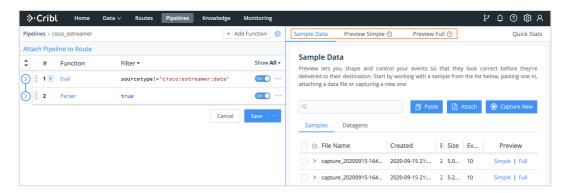
In the resulting **Destinations > Default** modal, use the **Default Output ID** dropdown to select one of your configured Destinations as LogStream's default.

The only other field here is the Output ID, whose value is locked to default.

Data Preview

Sample Data Preview is a LogStream feature that allows for visual inspection of events as they make their trip into a Pipeline. It helps you shape and control events before they're delivered to a Destination, as well as assisting with troubleshooting LogStream Functions.

Preview works by taking a set of Sample events, passing them through the Pipeline, and displaying the result in a separate pane. Any time a Function is modified, added, or removed, the Pipeline changes, and so does its displayed output.

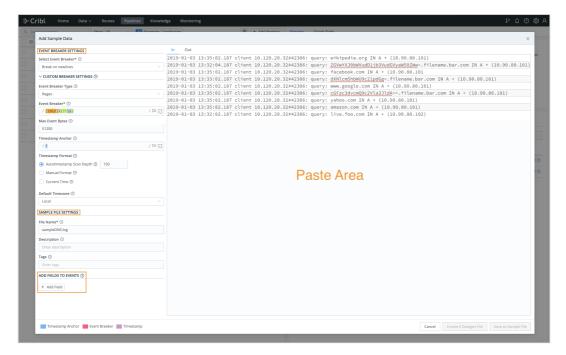


Preview options

While you're in a Pipeline, you can add samples through one of the supported options: Paste, Attach, or Capture New. The Paste and Attach options work with content that needs to be broken into events, while the Capture New option works with events only.

Adding Sample Data (Using Paste as an Example)

When you click on the corresponding option, you'll be presented with a modal like the one shown below.



Add Sample Data modal

i The Capture New modal is slightly different – it does not require event breaking.

Paste Area

This is where the content of the paste (or uploaded file) is displayed.

Event Breaker Settings

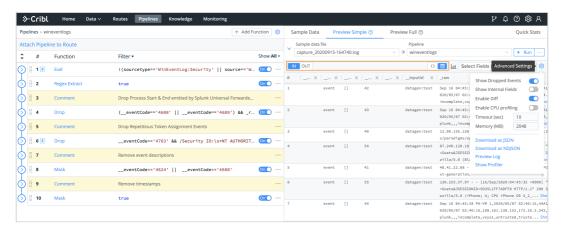
An Event Breaker is a regular expression that tells Cribl LogStream how to break the file or pasted content into events. Breaking will occur at the **start** of the match. Cribl LogStream ships with several common breaker patterns out of the box, but you can also configure custom breakers. The UI here is interactive, and you can iterate until you find the exact pattern.

Fields

The Fields section enables users to add, or overwrite. key/value pairs on the sample.

IN Tab: Displaying Samples on the Way IN to the Pipeline

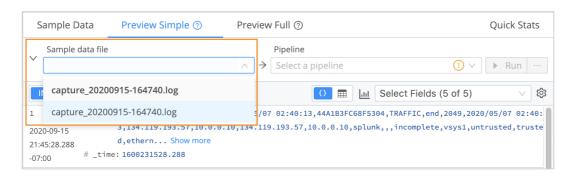
The Preview pane offers two display options for the event: Event and Table. (You can also download data as JSON or NDJSON, using the **Advanced Settings** menu at the top right.) Each format can be useful, depending on the type of data you are previewing.



Event, Table, and Advanced options

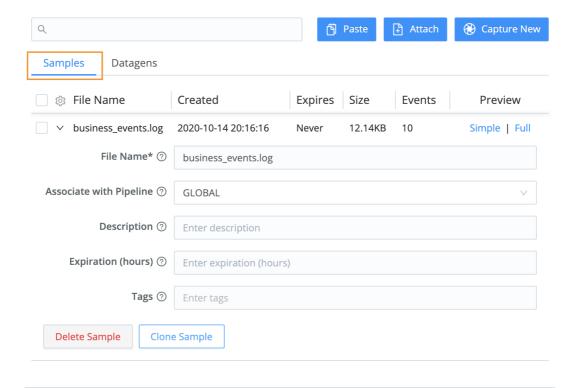
In the Advanced Settings menu's Timeout (sec) and Memory (MB) fields, you can increase the defaults to adjust for cases where very large data samples fail to load. For example, you might increase the Timeout (sec) to 30 and the Memory (MB) to 3048.

As you add more samples to your system, you can easily access them via the Samples drop-down near the top right, under Quick Stats.



Selecting an existing sample

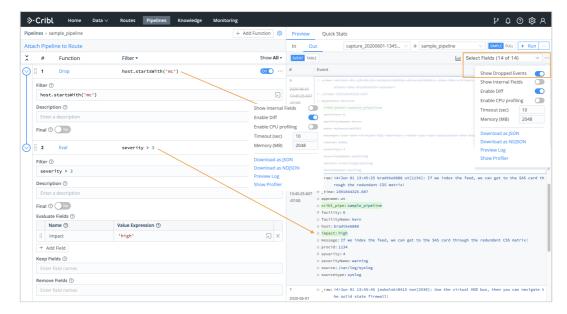
You can also manage, clone/modify, and delete samples via the **Samples** tab below.



OUT Tab: Displaying Samples on the Way OUT of the Pipeline

As data traverses Functions in a Pipeline, events can be modified, and some might be dropped altogether. The **OUT** tab indicates changes using this color coding:

- **Dropped events**: When events are dropped, the **OUT** tab displays them as grayed-out text, with strikethrough. You can control their display using the **Advanced Settings** menu's **Show Dropped Events** slider.
- Added fields: When LogStream's processing adds new fields, these fields are highlighted green. You can control these fields' display using the Select Fields drop-down.
- Redacted fields: These fields are highlighted amber.
- **Deleted fields**: These fields are highlighted red.



Dropped and added fields in a Pipeline's output

Securing Data

Cribl LogStream can be used to encrypt sensitive data in real time and route it to an end system. Decrypted retrieval can be implemented on a per-system basis. Currently, decryption is supported only when Splunk is the end system.

- Data Encryption
- Data Decryption

Encryption

Encryption of Data in Motion

With Cribl LogStream, you can encrypt fields or patterns within events in real time, by using C.Crypto.encrypt() in a Mask function. The Mask function accepts multiple replacement rules and multiple fields to apply them to.

A Match Regex defines the pattern of content to be replaced. The Replace Expression is a JS expression or literal to replace matched content. The C.Crypto.encrypt() method can be used here to generate an encrypted string from a value passed to it.

i C.Crypto.encrypt() Syntax

(method) Crypto.encrypt(value: any, keyclass: number,
keyId?: string, defaultVal?: string): string
Encrypt the given value with the keyId or a keyId picked up
automatically based on keyclass

- @param {string | Buffer} value what to encrypt
- @param keyclass if keyld isn't specified, pick one at the given key class
- @param keyld encryption keyld, takes precedence over keyclass
- @param defaultVal what to return if encryptions fails for any reason, if unspecified the original value is returned
- @returns - if encryption succeeds the encrypted value, otherwise defaultVal if specified, otherwise value.

Encryption Keys

Symmetric keys can be configured through the CLI or UI. Users are free to define as many keys as required. Each key is characterized by the following:

- keyId: ID of the key.
- algorithm: Algorithm used with the key
- keyclass: Cribl Key Class (below) that the key belongs to.
- kms: Key management system for the key. Defaults to local.

- created: Time (epoch) when key was generated.
- expires: Time (epoch) after which the key is invalid. Useful for key rotation.
- useIV: Flag that indicates whether or not an initialization vector was used.

Key Classes

Key Classes in Cribl LogStream are collections of keys that can be used to implement multiple levels of access control. Users (or groups of users) with access to data with encrypted patterns can be associated with key classes, for even more granular, pattern-level compartmentalized access.

Example

Users U0, U1 have been given access to keyclass 0 which contains key IDs 0 and 1. These keys are used to encrypt certain patterns in datasetA. Even though users U0, U1, U2 have access to read this dataset, only U0 and U1 can decrypt its encrypted patterns.

Key Class	Dataset
<pre>keyclass: 0 Keys: keyId: 0, keyId: 1 Users: U0, U1</pre>	datasetA Users: U0, U1, U2

User U1 has been given access to an **additional** keyclass, 1, which contains key IDs 11 and 22. These keys are used to encrypt certain **other** patterns in datasetA. Even though users U0, U1, U2 have access to read this dataset – same as above – only U1 can decrypt the additional encrypted patterns.

Key Class	Dataset
keyclass: 1 Keys: keyId: 11, keyId: 22 Users: U1	datasetA Users: U0, U1, U2

Configuring Keys with CLI

When using the local key management system, encryption keys in Cribl LogStream are encrypted with

\$CRIBL_HOME/local/cribl/auth/cribl.secret and stored in

\$CRIBL_HOME/local/cribl/auth/keys.json.Cribl monitors the keys.json file for changes every 60 seconds.

i When installed as a Splunk app, \$CRIBL_HOME is \$SPLUNK_HOME/etc/apps/cribl.

Listing Keys

Keys are added and listed using the keys command:

\$CRIBL_HOME/bin/cribl keys list -g <workerGroupID>

Sample Command Output

keyId	algorithm	keyclass	kms	created	expires	useIV
1	aes-256-cbc	0	local	1544906269.316	0	false
2	aes-256-cbc	1	local	1544906272.452	0	false
3	aes-256-cbc	2	local	1544906275.948	1545906275	true
4	aes-256-cbc	3	local	1544906278.026	0	false

Adding Keys

```
Displaying --help:

$CRIBL_HOME/bin/cribl keys add --help

SampleCommandOutput

Add encryption keys
Usage: [options] [args]

Options:

[-c <keyclass>] - key class to set for the key

[-k <kms>] - KMS to use, must be configured, see cribl.yml

[-e <expires>] - expiration time, epoch time
```

Adding a key to keyclass 1, with no expiration date, on the default Worker Group:

```
$CRIBL_HOME/bin/cribl keys add -c 1 -i -g default
```

[-i] - use an initialization vector

Sample Command Output

-g <group> - Group ID

Adding key: success. Key count=1

Listing keys to verify key generation:

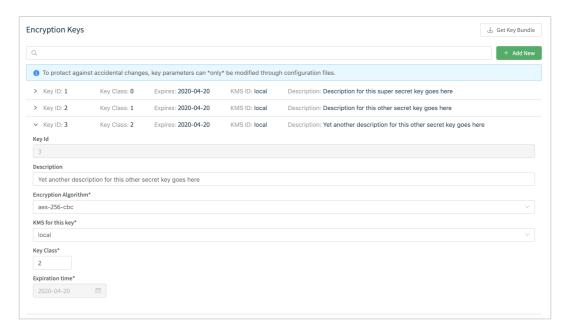
\$CRIBL_HOME/bin/cribl keys list -g default

Sample Command Output

keyId	algorithm	keyclass	kms	created	expires	useIV
1	aes-256-cbc	1	local	1545243364.342	0	true

Configuring Keys with the UI

The key management interface can be accessed through **Settings > Encryption Keys**. Here, you can list and add new keys. To protect against accidental changes, a key's parameters, once saved, can be edited only through configuration files.



List or create keys through LogStream's UI

Sync auth/(cribl.secret|keys.json)

To successfully decrypt data, the decrypt command will need access to the same keys that were used to encrypt. The cribl.secret and keys.json files in \$CRIBL_HOME/local/cribl/auth (in the Cribl instance where encryption happened) should be synced/copied over to the ones on the

Search Head/decrypting side. When using the UI, these files can be downloaded through the **Get Key Bundle** button.

Decryption

Decryption of Data

Currently, Cribl LogStream supports decryption only when Splunk is the end system. In Splunk, decryption is available to users of any role with permissions to run the decrypt command that ships with Cribl App for Splunk. Further restrictions can be applied with Splunk capabilities. This page provides details.

Decrypting in Splunk

Decryption in Splunk is implemented via a custom command called decrypt. To use the command, users must belong to a Splunk role that has permissions to execute it. Capabilities, which are aligned to Cribl Key Classes, can be associated with a particular role to further control the scope of decrypt.

i Decrypt Command Is Search Head ONLY

To ensure that keys don't get distributed to all search peers – including peers that your search head can search, but you don't have full control over – decrypt is scoped to run locally on the installed search head.

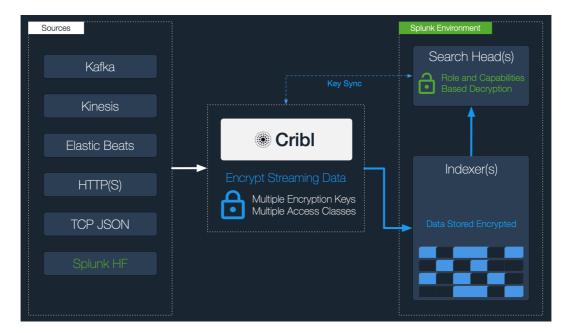
Restricting Access with Splunk Capabilities

In Splunk, capability names should follow the format cribl_keyclass_N, where N is the Cribl Key Class. For example, a role with capability cribl_keyclass_1 has access to all key IDs associated with key class_1.

Capability Name	Corresponding Cribl Key Class
cribl_keyclass_1	1
cribl_keyclass_2	2
cribl_keyclass_N	N

Configuring Splunk Search Head to Decrypt Data

You set up decryption in Splunk according to this schematic:



1. Install the Cribl App for Splunk on your Splunk search head.

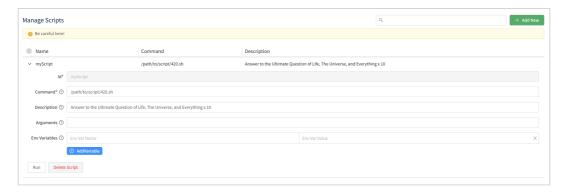
As of LogStream v1.7, the app will run in search head mode by default. If the app has previously been installed and later modified, you can convert it to search head mode with the command: \$CRIBL_HOME/bin/cribld mode-searchhead. (When installed as a Splunk app, \$CRIBL_HOME is \$SPLUNK_HOME/etc/apps/cribl.)

- 2. Assign permissions to the decrypt command, per your requirements.
- 3. Assign capabilities to your roles, per your requirements. If you'd like to create more capabilities, ensure that they follow the naming convention defined above.
- 4. Sync auth/(cribl.secret|keys.json). To successfully decrypt data, the decrypt command will need access to the same keys that were used to encrypt. The cribl.secret and keys.json files in \$CRIBL_HOME/local/cribl/auth which must be in the same Cribl instance where encryption happened should be synced/copied over to the files on the Search Head/decrypting side. When using the UI, these files can be downloaded through the Get Key Bundle button.

Scripts

Admins can run scripts (e.g., shell scripts) from within Cribl LogStream by configuring and executing them thru **Settings** > **Scripts**. Scripts are typically used to call custom automation jobs or, more generally, to trigger tasks on demand. For example, you can use Scripts to run an Ansible job, or to place a call to another automation system, when Cribl LogStream configs are updated.

Scripts will allow you to execute almost anything on the system where Cribl LogStream is running. Make sure you understand the impact of what you're executing before you do so!



Settings > Manage Scripts page

The Manage Scripts page provides the following tields:

- ID: Unique ID for this script.
- Command: Command to execute for this script.
- **Description**: Brief description about this script. Optional.
- Arguments: Arguments to pass when executing this script.
- Env variables: Extra environment variables to set when executing script.

i Scripts in Distributed Deployments

• Scripts can be deployed from Master Node, but can be **run** only locally from each Worker Node.

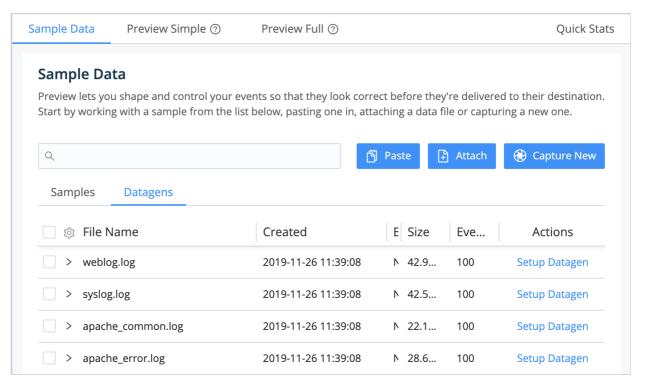
If the Script command is referencing a file (e.g., 420.sh), that file
must exist on the Cribl LogStream instance. In other words, the
Script management interface cannot be used to upload or
manage script files.

Using Datagens

Data generators for testing and troubleshooting

Cribl LogStream's Datagens feature enables you to generate sample data for the purposes of troubleshooting Routes, Pipelines, Functions, and general connectivity.

Several Datagen template files ship with the product, out of the box. You can create others from sample files or live captures.



Preview pane – add samples via paste, attach/upload file, or live capture

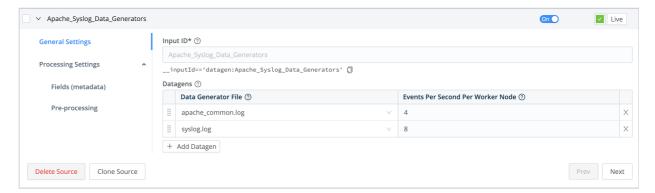
As outlined in the following tutorial: Once you've created a template, you can configure a Datagen Source to use the template to generate real-time data at a given EPS (events per second) rate.

Enabling a Datagen

To see how Datagens work, start by enabling a pair of LogStream's out-of-the-box generators:

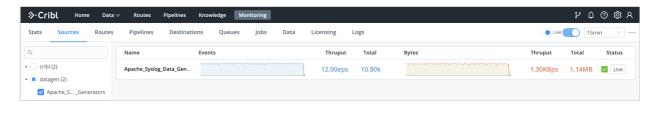
Navigate to **Sources** > **Datagens** and click **Add New**.

Select a Data Generator File (e.g., apache_common.log) and set it at 4 EPS/worker process. Select another Data Generator File (e.g., syslog.log) and set it at 8 EPS/worker process. Hit Save.



Selecting Datagens files and event rates

On the **Monitoring** page, under **Sources**, search for datagen and confirm that the Source is generating data.



Creating a Datagen Template from a Sample File

To convert a sample into a template:

Go to Preview > Paste a Sample, and add a sample like the AWS VPC Flow logs below:

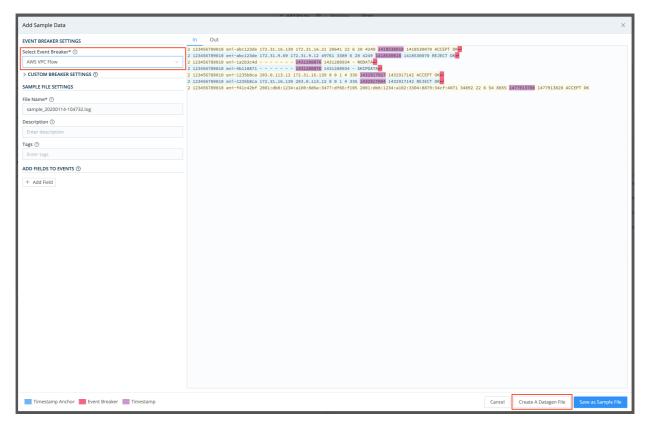
Sample VPC Flow Logs

```
2 123456789010 eni-abc123de 172.31.16.139 172.31.16.21 20641 22 6 20 4249 1418530010 1418530070 2 123456789010 eni-abc123de 172.31.9.69 172.31.9.12 49761 3389 6 20 4249 1418530010 1418530070 F 2 123456789010 eni-1a2b3c4d - - - - - - 1431280876 1431280934 - NODATA 2 123456789010 eni-4b118871 - - - - - - 1431280876 1431280934 - SKIPDATA 2 123456789010 eni-1235b8ca 203.0.113.12 172.31.16.139 0 0 1 4 336 1432917027 1432917142 ACCEPT 2 123456789010 eni-1235b8ca 172.31.16.139 203.0.113.12 0 0 1 4 336 1432917094 1432917142 REJECT 2 123456789010 eni-f41c42bf 2001:db8:1234:a100:8d6e:3477:df66:f105 2001:db8:1234:a102:3304:8879
```

From the **Event Breaker** drop-down, select **AWS VPC Flow** to ensure that:

- The pasted text gets broken properly into individual events (notice the Event Breaker on newlines).
- Timestamps are extracted correctly (text highlighted purple below).

Once you've verified these results, click Create a Datagen File.



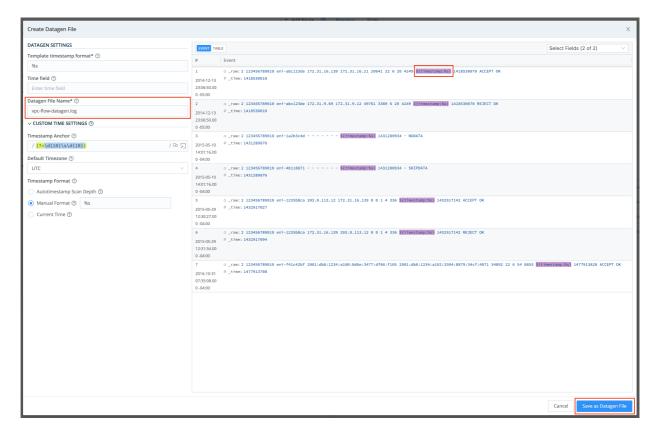
Creating a Datagen template

On the resulting Create Datagen File screen:

- Enter a file name, e.g.: vpc-flow-datagen.log
- Ensure that the timestamp template format is correct: \${timestamp: %s}

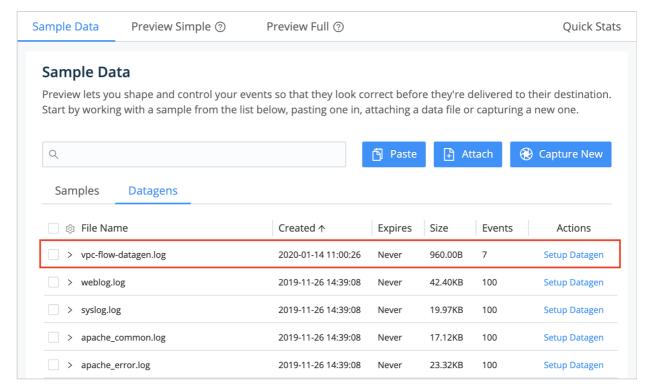
\${timestamp: <format>} is a template that the datagen engine uses to insert the current time – in each newly generated event – using the given format. In this case, %s is the desired strftime format for the timestamp (i.e., the epoch).

Once you've verified these results, click Save as Datagen File.



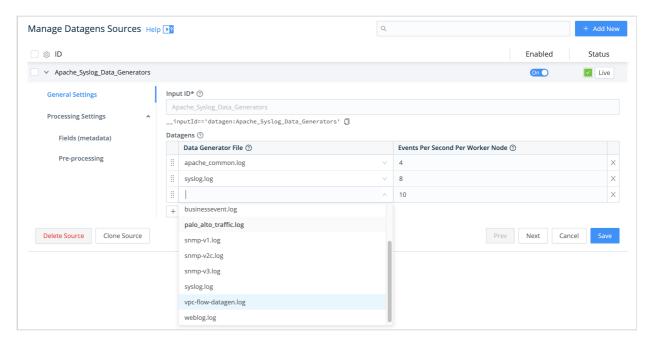
Saving a named Datagen template

To confirm that the Datagen file has been created, check **Preview > Datagens**.



Verifying Datagen file creation

Now, to start using your newly created Datagen file, go back to **Sources > Datagens**. Add it using the drop-down shown below.



Adding new template file to Datagens Source

CLI Reference

Command line interface basics

In addition to starting and stopping the Cribl LogStream server, LogStream's command line interface enables you to initiate many configuration and administrative tasks directly from your terminal.

Command Syntax

To execute CLI commands, the basic syntax is:

```
cd $CRIBL_HOME/bin
./cribl <command> <sub-command> <options> <arguments>
```

Avoiding Surprises

Immediate Execution

As indicated in the sample output below, some commands take effect immediately.

Commands that require further input will echo the sub-commands, options, and arguments they expect.

Persistent Volumes

If you start LogStream with the CRIBL_VOLUME_DIR variable, all subsequent CLI commands should have this variable defined. Otherwise, those commands will apply LogStream's default directories, yielding misleading results.

You can set CRIBL_VOLUME_DIR as an environment variable, or you can explicitly include it in each command, as in this example:

```
CRIBL_VOLUME_DIR /opt/cribl/bin/cribl status.
```

Commands Available

To see a list of available commands, enter ./cribl alone (or the equivalent ./cribl help). To execute a command, or to see its required parameters, enter ./cribl <command> .

help

Displays help (commands list).

```
Cribl LogStream - N.n.n ≺build no.>
Usage: [sub-command] [options] [args]
Commands:
```

- Display help help

mode-master - Configure to a master instance mode-single - Configure to a single instance mode-worker - Configure to a worker instance

reload - Reload Cribl LogStream
restart - Restart Cribl LogStream
start - Start Cribl LogStream
status - Status of Cribl LogStream
stop - Stop Cribl LogStream
version - Print Cribl LogStream version and installation type

auth - Cribl LogStream Auth
boot-start - Enable/Disable Cribl LogStream boot-start
diag - Manage diagnostics bundles
groups - Manage Worker Groups
keys - Manage encryption keys

mode-searchhead - Configure Cribl LogStream to run on a Splunk Search Head - Listen on a port for traffic and output stats and data

node — Execute a JavaScript file pipe — Feed stdin to a pipeline splunk-decrypt - Splunk decrypt search command task - Run Cribl LogStream task vars - Manage global variables

mode-master

Configures Cribl LogStream as a Master instance.

Options

```
[-H <host>] - Host (defaults to 0.0.0.0).
[-p <port>] - Port (defaults to 4200).
[-n <certName>] - Name of saved certificate.
[-k <privKeyPath>] - Server path containing the private key (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use. Can refered to the containing certificates (in PEM format) to use.
   \hbox{[-i < ipWhitelistRegex>]} - \hbox{Regex matching IP addresses that are allowed to establish a connection}
```

Sample Response

```
Settings updated.
You will need to restart LogStream before your changes take full effect.
```

mode-single

Configures Cribl LogStream as a single-instance deployment.

Sample Response

```
Settings updated.
You will need to restart LogStream before your changes take full effect.
```

mode-worker

Configures Cribl LogStream as a Worker instance.

Usage

```
./cribl mode-worker -H <host> -p <port>
```

The -H <host> -p <port> parameters are required.

Options

```
-H <host> - Master Node's Hostname or IP address.
-p <port> - Master Node's cluster communications port (defaults to 4200).

[-n <certName>] - Name of saved certificate.

[-k <privKeyPath>] - Server path containing the private key (in PEM format) to use. Can reference certPath>]

[-u <authToken>] - Authentication token to include as part of the connection header. By defaul certificates (in PEM format) to use. Can reference certPath>]

[-e <envRegex>] - Regex that selects environment variables to report to Master.

[-t <tags>] - Tag values to report to master.

[-g <group>] - Worker Group to report to master.
```

Sample Response

```
Settings updated.
You will need to restart LogStream before your changes take full effect.
```

reload

Reloads Cribl LogStream. Executes immediately.

```
Reload request submitted to Cribl LogStream
```

restart

Restarts Cribl LogStream. Executes immediately.

 \triangle Executing this command cancels any running collection jobs.

```
Stopping Cribl LogStream, process 56572 ......

Cribl LogStream is not running
Starting Cribl LogStream...
```

```
..
Cribl LogStream started with pid 57233
API Server is available at http://192.168.0.100:9000
```

start

Starts Cribl LogStream. Executes immediately. Upon first run, echoes LogStream's default login credentials.

```
Starting Cribl LogStream...
...
Cribl LogStream started with pid 57279
API Server is available at http://192.168.0.100:9000
```

status

Displays status of Cribl LogStream, including the API Server address, instance's mode (Master or Worker), process ID, and GUID (fictitious example below). Executes immediately.

```
Cribl LogStream Status
```

Address: http://192.168.0.100:9000

Mode: worker Status: Up

Software Version: 42.0-7f4c190a

Master: localhost:4200

PID: 3859

GUID: 76-ea411263a64b9-e419daee4-ef-dd2e2f

stop

Stops Cribl LogStream. Executes immediately.

Executing this command cancels any running collection jobs.

```
Stopping Cribl LogStream, process 57233 ......
Cribl LogStream is not running
```

version

Displays Cribl LogStream version and installation type. Executes immediately.

Version: 2.2-0####x##

Installation type: standalone

☐ The version command echoes standalone for both single-instance and distributed deployments. This simply confirms that you're running a freestanding Cribl LogStream server, not the Cribl App for Splunk.

auth

Log into or out of Cribl LogStream.

```
Commands:
login - Log in to Cribl LogStream, args:
  [-h <host>] - Host URL (e.g. http://localhost:9000)
  [-u <username>] - Username
  [-p <password>] - Password
  [-f <file>] - File with credentials
logout - Log out from Cribl LogStream
```

Login Examples

Launch interactive login:

```
$CRIBL_HOME/bin/cribl auth login
```

Append credentials as command arguments:

```
$CRIBL_HOME/bin/cribl auth login -h <url> -u <username> -p <password>
```

i All -h and host arguments are optional, provided that the API host and port are listed in the cribl.yml file's api: section

Provide credentials in environment variables:

```
CRIBL_HOST=<url> CRIBL_USERNAME=<username> CRIBL_PASSWORD=<password>
$CRIBL_HOME/bin/cribl auth login
```

Provide credentials in a file:

```
$CRIBL_HOME/bin/cribl auth login -f <path/to/file>
```

--

Corresponding file contents:

```
host=<url>
username=<username>
password=<password>
```

boot-start

Enables or disables Cribl LogStream boot-start.

```
Usage: [sub-command] [options] [args]

Commands:
disable - Disable Cribl LogStream boot-start, args:
    [-m <manager>] - Init manager (systemd|initd)
    [-c <configDir>] - Config directory for the init manager
enable - Enable Cribl LogStream boot-start, args:
    [-m <manager>] - Init manager (systemd|initd)
    [-u <user>] - User to run Cribl LogStream as
    [-c <configDir>] - Config directory for the init manager
```

diag

Manages diagnostic bundles.

```
create - Creates diagnostic bundle for Cribl LogStream

list - List existing Cribl LogStream diagnostic bundles

send - Send LogStream diagnostics bundle to Cribl Support, args:
   -c <caseNumber> - Cribl Support Case Number
[-p <path>] - Diagnostic bundle path (if empty then new bundle will be created)
```

groups

Manages Worker Groups.

```
Usage: [sub-command] [options] [args]

Commands:
commit - Commit, args:
    [-g <group>] - Group ID
    [-m <message>] - Commit message
commit-deploy - Commit & Deploy, args:
    -g <group> - Group ID
    [-m <message>] - Commit message
deploy - Deploy, args:
    -g <group> - Group ID
    [-v <version>] - Deploy version
list - List Worker Groups
```

keys

Manages encryption keys. You must append the -g <group> argument to specify a Worker Group. As a fallback, append the argument -g default, e.g.:

```
./cribl keys list -g default
```

```
Usage: [sub-command] [options] [args]

Commands:
add - Add encryption keys, args:
    -g <group> - Group ID
    [-c <keyclass>] - key class to set for the key
    [-k <kms>] - KMS to use, must be configured, see cribl.yml
    [-e <expires>] - expiration time, epoch time
    [-i ] - use an initialization vector

list - List encryption keys
    -g <group> - Group ID
```

mode-searchhead

Configures Cribl LogStream to run on a Splunk Search Head.

nc

Listens on a port for traffic, and outputs stats and data. (Netcat-like utility.)

node

Executes a JavaScript file. Displays a command prompt for path/filename input, as shown here:

>

pipe

Feeds stdin to a pipeline. Examples:

```
cat sample.log | ./cribl pipe -p <pipelineName>
cat sample.log | ./cribl pipe -p <pipelineName> 2>/dev/null
```

scope

Greps your apps by the syscalls. Executes immediately.

```
splunk-decrypt
```

task

Runs a Cribl LogStream task. Requires definitions for the dir, executor, and path properties.

vars

Manages LogStream Global Variables.

```
Usage: [sub-command] [options] [args]
Commands:
add - Add global variable, args:
  -i <id> - Global variable ID
-t <type> - Type
-v <value> - Value
[-a <args>] - Arguments
  [-d <description>] - Description
 [-c <tags>] - Custom Tags (comma separated list)
[-g <group>] - Group ID
get - List encryption keys, args:
  [-i <id>] - Global variable ID
  [-g <group>] - Group ID
remove - Remove global variable, args:
   -i <id> - Global variable ID
  [-g <group>] - Group ID
update - Update global variable, args:
  -i <id> - Global variable ID

-t <type> - Type

-v <value> - Value

[-a <args>] - Arguments
  [-d <description>] - Description
  [-c <tags>] - Custom Tags (comma separated list)
[-g <group>] - Group ID
```

Expression Reference

Introduction to Expression Syntax

As data travels through a Cribl LogStream Pipeline, it is operated on by a series of Functions. Functions are fundamentally JavaScript code.

Functions that ship with Cribl LogStream are configurable via a set of inputs. Some of these configuration options are literals, such as field names, and others can be JavaScript expressions.

Expressions are **valid units** of code that resolve to a value. Every syntactically valid expression resolves to some value, but conceptually, there are two types of expressions: those that **assign** value to a variable (a.k.a., with side effects), and those that **evaluate** to a value.

Assigning a value	Evaluating to a value
<pre>x = 42 newFoo = foo.slice(30)</pre>	(Math.random() * 42) 3 + 4 'foobar' '42'

Filters and Value Expressions

Filters

Filters are used in Routes to select a stream of the data flow, and in Functions to scope or narrow down the applicability of a Function. Filters are expressions that must evaluate to either true (or truthy) or false (or falsy). Keep this in mind when creating Routes or Functions. For example:

- sourcetype='access_combined' & host.startsWith('web')
- source.endsWith('.log') ||
 sourcetype='aws:cloudwatchlogs:vpcflow'

Truthy	Falsy
true	false
42	null
-42	undefined
3.14	0
"foo"	NaN
Infinity	
-Infinity	н н

Value Expressions

Value expressions are typically used in Functions to assign a value – for example, to a new field. For example:

```
• Math.floor(_time/3600)
```

source.replace(/.{3}/, 'XXX')

Considerations and Best Practices for Creating Predictable Expressions

- In a value expression, ensure that the source variable is not null, undefined, or empty. For example, assume you want to have a field called len, to be assigned the length of a second field called employeeID. But you're not sure if employeeID exists. Instead of employeeID.length, you can use a safer shorthand, such as: (employeeID || '').length.
- If a field does not exist (undefined), and you're doing a comparison with its properties, then the boolean expression will always evaluate to false. For example, if employeeID is undefined, then both of these expressions will evaluate to false: employeeID.length > 10 , and employeeID.length < 10 .
- = means "equal to," while == means "equal value and equal type." For example, 5 = 5 evaluates to true, while 5 == "5" evaluates to false.
- A ternary operator is a very powerful way to create conditional values. For example, if you wanted to assign either minor or adult to a field groupAge, based on the value of age, you could do: (age ≥ 18)?
 'adult': 'minor'.

Expressions Using Fields with Non-Alphanumeric Characters

If there are fields with non-alphanumeric characters – e.g., @timestamp or user-agent or kubernetes.namespace_name – you can access them using __e['<field-name-here>']. (Note the single quotes.) More details here. In any other place where the field is referenced – e.g., in the Eval function's field names – you should use a single-quoted literal, of the form: '<field-name-here>'.

Wildcard Lists

Wildcard Lists are used throughout the product, especially in various Functions, such as Eval, Mask, Publish Metrics, Parser, etc.

Wildcard Lists, as their name implies, accept strings with asterisks (\star) to represent one or more terms. They also accept strings that start with an exclamation mark (!) to **negate** one or more terms.

Wildcard Lists are order-sensitive **only** when negated terms are used. This allows for implementing any combination of whitelists and blacklists.

For example:

Wildcard List	Value	Meaning
List 1	!foobar, foo*	All terms that start with foo , except foobar .
List 2	!foo*, *	All terms, except for those that start with foo .

Cribl Expressions

Native Cribl LogStream function methods can be found under C.*, and can be invoked from any Function that allows for expression evaluations. For example, to create a field that is the SHA1 of a another field's value, you can use the Eval Function with this **Evaluate Fields** pair:

Name	Value Expression
myNewField	C.Mask.sha1(myOtherField)

C.Crypto – Data Encryption and Decryption Functions

```
C.Crypto.decrypt
```

(method) Crypto.decrypt(value: string): string

Decrypt all occurrences of ciphers in the given value. Instances that cannot be decrypted (for any reason) are left intact.

@param - value - string in which to look for ciphers

@returns - value with ciphers decrypted

C.Crypto.encrypt

(method) Crypto.encrypt(value: any, keyclass: number, keyId?:
string, defaultVal?: string): string
Encrypt the given value with the keyId, or with a keyId picked up.

Encrypt the given value with the keyId, or with a keyId picked up automatically based on keyclass.

- @param {string | Buffer} value what to encrypt.
- @param keyclass if keyId isn't specified, pick one at the given keyclass.
- @param keyId encryption keyId, takes precedence over keyclass.
- @param defaultVal what to return if encryption fails for any reason; if unspecified, the original value is returned.
- @returns if encryption succeeds, the encrypted value; otherwise, defaultVal if specified; otherwise, value.

C.Decode - Data Decoding Functions

C.Decode.base64

(method) Decode.base64(val: string, resultEnc?: string): any

```
Performs base64 decoding of the given string. Returns a string or Buffer,
depending on the resultEnc value, which defaults to 'utf8'.
@param - val - value to base64-decode
@param - resultEnc - encoding to use to convert the binary data to a string.
Defaults to 'utf8'. Use 'utf8-valid' to validate that result is valid UTF8;
use 'buffer' if you need the binary data in a Buffer.
C.Decode.gzip
(method) Decode.gzip(value: any, encoding?: string): string
Gunzip the supplied value.
@param - value - the value to gunzip.
@param - encoding - encoding of value, for example: 'base64', 'hex',
'utf-8', 'binary'. Default is 'base64'. If data is received as Buffer (from
gzip with encoding: 'none'), decoding is skipped.
C.Decode.hex
(method) Decode.hex(val: string): number
Performs hex to number conversion. (Returns NaN if value cannot be
converted to a number.)
@param - val - hex string to parse to a number (e.g., "0xcafe").
C.Decode.uri
(method) Decode.uri(val: string): string
Performs URI-decoding of the given string.
@param - val - value to URI-decode.
C.Encode – Data Encoding Functions
C.Encode.base64
(method) Encode.base64(val: any, trimTrailEq?: boolean): string
Returns a base64 representation of the given string or Buffer.
@param - val - value to base64-encode.
@param - trimTrailEq - whether to trim any trailing = .
C.Encode.gzip
(method) Encode.gzip(value: string, encoding?: string): any
```

@param - encoding - encoding of value, for example: 'base64', 'hex',
'utf-8', 'binary', 'none'. Defaultis 'base64'. If 'none' is specified,

Gzip, and optionally base64-encode, the supplied value.

@param - value - the value to gzip.

data will be returned as a Buffer.

```
C.Encode.hex
(method) Encode.hex(val: string | number): string
Rounds the number to an integer and returns its hex representation
(lowercase). If a string is provided, it will be parsed into a number or NaN.
@param - val - value to convert to hex.

C.Encode.uri
(method) Encode.uri(val: string): string
Returns the URI-encoded representation of the given string.
@param - val - value to uri encode.
```

C.env - Environment

```
C.env
(property) env: {[key: string]: string;}
Returns an object containing Cribl LogStream's environment variables.
```

C.Lookup – Inline Lookup Functions

```
C.Lookup - Exact Lookup
(property) Lookup: (file: string, primaryKey?: string,
otherFields?: string[], ignoreCase?: boolean) ⇒ InlineLookup
Returns an instance of a lookup to use inline.
Example invocation:
C.Lookup('lookup_name.csv',
'IP_field_name_in_lookup_file').match(host)
C.LookupCIDR - CIDR Lookup
(property) LookupCIDR: (file: string, primaryKey?: string,
otherFields?: string[]) ⇒ InlineLookup
Returns an instance of a CIDR lookup to use inline.
C.LookupIgnoreCase - Case-insensitive Lookup
(property) LookupIgnoreCase: (file: string, primaryKey?: string,
otherFields?: string[]) ⇒ InlineLookup
Returns an instance of a lookup (ignoring case) to use inline. Works identically
to C.Lookup, except ignores the case of lookup values. (Equivalent to calling
C.Lookup with its fourth ignoreCase? parameter set to true ).
C.[LookupRegex](http://google.com) - Regex Lookup
(property) LookupRegex: (file: string, primaryKey?: string,
```

```
otherFields?: string[]) ⇒ InlineLookup
Returns an instance of a Regex lookup to use inline.
(method) InlineLookup.match(value: string, fieldToReturn?:
string): any
@param - value - the value to look up.
@param - fieldToReturn - name of the lookup file > field to return.
E.g., C.Lookup('lookup-exact.csv', 'foo').match('abc', 'bar')
Return the value of field bar in the lookup table if field foo matches abc.
Example 1: C.LookupCIDR('lookup-cidr.csv',
'foo').match('192.168.1.1', 'bar')
Return the value of field bar in the lookup table if the CIDR range in foo
includes 192.168.1.1.
Example 2: C.LookupCIDR('lookup-cidr.csv', 'cidr').match(hostIP,
'location')
Example 3: C.LookupRegex('lookup-regex.csv',
'foo').match('manchester', 'bar')
Return the value of field bar in the lookup table if the regex in foo matches the
string manchester.
```

C.Mask - Data Masking Functions

```
C.Mask.CC
(method) Mask.CC(value: string, unmasked?: number, maskChar?:
string): string
Check whether a value could be a valid credit card number, and mask a subset
of the value. By default, all digits except the last 4 will be replaced with X.
@param - value - a string whose digits to mask IFF it could be a valid credit
card number.
@param - unmasked - number of digits to leave unmasked: positive for left,
negative for right, 0 for none.
@param - maskChar - a string/char to replace a digit with.

C.Mask.IMEI
```

string): string
Check whether a value could be a vlaid IMEI number, and mask a subset of the value. By default, all digits except the last 4 will be replaced with X.

@param - value - a string whose digits to mask IFF it could be a valid IMEI

(method) Mask.IMEI(value: string, unmasked?: number, maskChar?:

number.

@param – unmasked – number of digits to leave unmasked: positive for left, negative for right, 0 for none.

@param - maskChar - a string/char to replace a digit with.

C.Mask.isCC

(method) Mask.isCC(value: string): boolean

Checks whether the given value could be a valid credit card number, by computing the string's Lunh's checksum modulo 10 == 0.

@param - value - a string to check for being a valid credit card number.

C.Mask.isIMEI

(method) Mask.isIMEI(value: string): boolean

Checks whether the given value could be a valid IMEI number, by computing the string's Lunh's checksum modulo $10 == \emptyset$.

@param - value - a string to check for being a valid IMEI number

C.Mask.luhn

(method) Mask.luhn(value: string, unmasked?: number, maskChar?:
string): string

Check that value Lunh's checksum mod 10 is $\,0\,$, and mask a subset of the value. By default, all digits except the last 4 will be replaced with $\,X\,$. If the value's Lunh's checksum mod 10 is not $\,0\,$, then the value is returned unmodified.

@param - value - a string whose digits to mask IFF the value's Lunh's checksum mod 10 is 0.

@param – unmasked – number of digits to leave unmasked: positive for left, negative for right, 0 for none.

@param - maskChar - a string/char to replace a digit with.

C.Mask.LUHN_SUB

(property) Mask.LUHN_SUB: any

C.Mask.luhnChecksum

(method) Mask.luhnChecksum(value: string, mod?: number): number Generates the Luhn checksum (used to validate certain credit card numbers, IMEIs, etc.). By default, the mod 10 of the checksum is returned. Pass mod = 0 to get the actual checksum.

@param – value – a string whose digits you want to perform the Lunh checksum on.

@param - mod - return checksum modulo this number. If 0, skip modulo. Default is 10.

C.Mask.md5

(method) Mask.md5(value: string, len?: string | number): string
Generate MD5 hash of a given value.

@param - value - compute the hash of this.

@param – len – length of hash to return: 0 for full hash, a +number for left or a -number for right substring. If a string is passed it's length will be used.

C.Mask.random

(method) Mask.random(len?: string | number): string
Generates a random alphanumeric string.

@param – len – a number indicating the length of the result; or, if a string, use its length.

C.Mask.REDACTED

(property) Mask.REDACTED: string

The literal 'REDACTED'.

C.Mask.repeat

(method) Mask.repeat(len?: string | number, char?: string):
string

Generates a repeating char/string pattern, e.g., XXXX.

@param – len – a number indicating the length of the result; or, if a string, use its length.

@param - char - pattern to repeat len times.

C.Mask.sha1

(method) Mask.sha1(value: string, len?: string | number): string
Generate SHA1 hash of given value.

@param - value - compute the hash of this.

@param – len - length of hash to return: 0 for full hash, a +number for left, or a -number for right.

substring. If a string is passed, its length will be used

C.Misc – Miscellaneous Utility Functions

C.Misc.zip()

(method) Misc.zip(keys: string[], values: any[], dest?: any): any Set the given keys to the corresponding values on the given dest object. If dest is not provided, a new object will be constructed.

@param - keys - field names corresponding to keys.

@param - values - values corresponding to values.

```
@param - dest - object on which to set field values.
@returns - object on which the fields were set.

E.g., people = C.Misc.zip(titles, names)

Sample data: titles=['ceo', 'svp', 'vp'], names=['foo', 'bar', 'baz']

Create an object called people, with key names from elements in titles, and with corresponding values from elements in names.

Result: "people": {"ceo": "foo", "svp": "bar", "vp": "baz"}
```

C.Net – Network Functions

```
C.Net.cidrMatch()
(method) Net.cidrMatch(cidrIpRange: string, ipAddress: string):
boolean
Determines whether the supplied IPv4 ipAddress is inside the range of
addresses identified by cidrIpRange . For example: C.Net.cidrMatch
('10.0.0.0/24', '10.0.0.100') returns true.
@param - cidrIpRange - IPv4 address range in CIDR format. E.g.,
10.0.0.0/24.
@param - ipAddress - The IPv4 IP address to test for inclusion in
cidrIpRange.
C.Net.ipv6Normalize()
(method) Net.ipv6Normalize(address: string): string
Normalize an IPV6 address based on RFC draft-ietf-6man-text-addr-
representation-04.
@param - address - the IPV6 address to normalize.
C.Net.isPrivate()
(method) Net.isPrivate(address: string): string
Determine whether the supplied IPv4 address is in the range of private
addresses per RFC1819.
@param - address - address to test.
```

C.os – System Functions

C.confVersion

Returns Cribl LogStream config version.

C.os.hostname()

Returns hostname of the system running this Cribl LogStream instance.

C.Schema - Schema Functions

```
C.Schema()
(property) Schema: (id: string) ⇒ SchemaValidator
(method) SchemaValidator.validate(data: any): boolean
Validates the given object against the schema.
@param - data - object to be validated.
@returns - true when schema is valid; otherwise, false.
Example: C.Schema('schema1').validate(myField) will validate if
myField object conforms to schema1.
See Schema Library for more details.
C.Text - Text Functions
C.Text.entropy()
(method) Text.entropy(bytes: any): number
Computes the Shannon entropy of the given buffer or string.
@param - bytes - value to undergo Shannon entropy computation.
@returns – the entropy value; or –1 in case of an error.
C.Text.hashCode()
(method) Text.hashCode(val: string | Buffer | number): number
Computes hashcode (djb2) of the given value.
@param - val - value to be hashed.
@returns - hashcode value.
C.Text.isASCII()
(method) Text.isASCII(bytes: any): boolean
Checks whether all bytes or chars are in the ASCII printable range.
@param - bytes - value to check for character range.
@returns - true if all chars/bytes are within ASCII printable range; otherwise,
false.
C.Text.isUTF8()
(method) Text.isUTF8(bytes: any): boolean
Checks whether the given Buffer contains valid UTF8.
@param - bytes - bytes to check.
@returns - true if bytes are UTF8; otherwise, false.
```

(method) C.Text.parseWinEvent(xml: string, nonValues?: string[]):

C.Text.parseWinEvent

any

Parses an XML string representing a Windows event into a compact, prettified JSON object. Works like C.Text.parseXml, but with Windows events, produces more-compact output. For a usage example, see Reducing Windows XML Events.

@param - xml - an XML string; or an event field containing the XML.

@param - nonValues - array of string values. Elements whose value equals any of the values in this array will be omitted from the returned object.

Defaults to ['-'], meaning that elements whose value equals - will be discarded.

@returns – an object representing the parsed Windows Event; or undefined if the input could not be parsed.

C.Text.parseXml

(method) C.Text.parseXml(xml:string, keepAttr?:boolean,
keepMetadata?:boolean, nonValues?:string[]): any
Parses an XML string and returns a JSON object. Can be used with Eval
Function to parse XML fields contained in an event, or with ad hoc XML.
@param - xml - XML string, or an event field containing the XML.
@param - keepAttr - whether or not to include attributes in the returned
object. Defaults to true.

@param - keepMetadata - whether or not to include metadata found in the XML. The keepAttr parameter must be set to true for this to work. Defaults to false. (Eligible metadata includes namespace definitions and prefixes, and XML declaration attributes such as encoding, version, etc.)

@param - nonValues - array of string values. Elements whose value equals any of the values in this array will be omitted from the returned object.

Defaults to [] (empty array), meaning discard no elements.

@returns - an object representing the parsed XML; or undefined if the input

@returns – an object representing the parsed XML; or undefined if the input could not be parsed. An input collection of elements will be parsed into an array of objects.

C.Text.relativeEntropy()

(method) Text.relativeEntropy(bytes: any, modelName?: string):
number

Computes the relative entropy of the given buffer or string.

- @param bytes value whose relative entropy to compute.
- @param modelName Name of the model to test the string with.
- @returns the relative entropy value, or –1 in case of an error.

C.Time – Time Functions

```
C.Time.adjustTZ()
(method) Time.adjustTZ(epochTime: number, tzTo: string, tzFrom?:
string): number
Adjust a timestamp from one timezone to another.
@param - epochTime - UNIX epoch time.
@param - tzTo - timezone to adjust to.
@param - tzFrom - optional timezone of the timestamp.
C.Time.clamp()
(method) Time.clamp(date, earliest, latest, defaultDate?): number
Constrains an event's parsed timestamp to realistic earliest and latest
boundaries.
@param - date - Timestamp originally parsed from event, in UNIX epoch
time (ms) or JavaScript Date format.
@param - earliest - earliest allowable timestamp, in UNIX epoch time (ms)
or JS Date format.
@param - latest - latest allowable timestamp, in UNIX epoch time (ms) or JS
Date format.
@param - defaultDate - optional default date, in UNIX epoch time (ms) or
JS Date format, to substitute for values outside the earliest or latest
boundaries.
C.Time.strftime()
(method) Time.strftime(date: number | Date, format: string, utc?:
boolean): string
Format a Date object or number as a time string, using strftime specifier.
@param - date - Date object or number (seconds since epoch) to format.
@param - format - specifier to use to format the date.
@param - utc - whether to output the time in UTC, rather than in local
timezone.
@returns – representation of the given date.
C.Time.strptime()
(method) Time.strptime(str: string, format: string, utc?: boolean,
strict?: boolean): Date
Extract time from a string using strptime specifier.
@param - str - string to parse to a timestamp (see strict flag).
@param - format - strptime specifier.
@param - utc - whether to interpret times as UTC, rather than as local time.
@param - strict - whether to return null if there are any extra characters
after timestamp.
```

@returns – a parsed Date object, if successful; otherwise, null if the specifier did not match.

C.Time.timestampFinder()

(method) Time.timestampFinder(utc?: boolean).find(<source-field>):
AutoTimeParser

Extract time from the specified <source-field>, using the same algorithm as the Auto Timestamp Function and the Event Breaker Function.

- @param utc whether to output the time in UTC, rather than in local timezone.
- @param <source-field> the field in which to search for the time.
- @returns representation of the extracted time.

C.vars - Global Variables

See Global Variables Library for more details.

C.version – Cribl LogStream Version

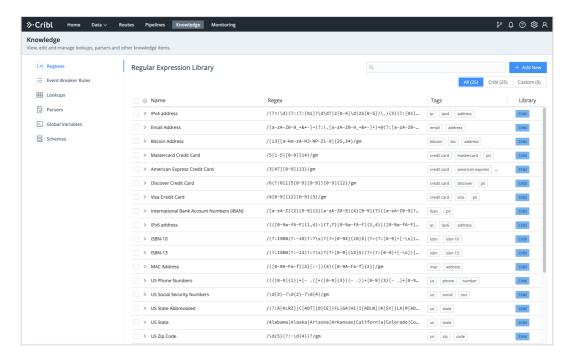
(property) version: string
Cribl LogStream Version.

Knowledge

Regex Library

What Is the Regex Library

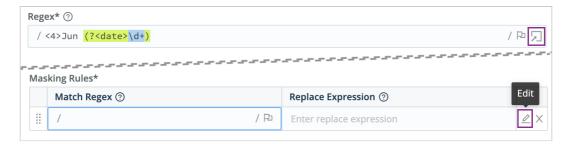
Cribl LogStream ships with a Regex Library that contains a set of pre-built common regex patterns. This library serves as an easily accessible repository of regular expressions. The Library is searchable, and you can assign tags to each pattern for further organization or categorization. The Library is located under **Knowledge** > **Regex Library**.



Regular Expression Library

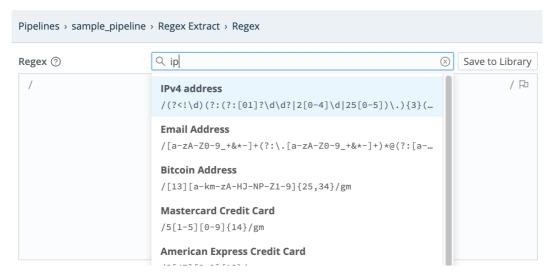
Using Library Patterns

As of this version, the Library contains 25 patterns shipped by Cribl LogStream. To insert a pattern into a Function's regex field, first click the pop-out or Edit icon beside that field.



Opening a Regex modal

In the resulting Regex or Rules modal, Regex Library patterns will appear as typeahead options. Click a pattern to paste it in. You can then use the pattern as-is, or modify it as necessary.



Inserting a pattern from the Regex Library

Adding Patterns to the Library

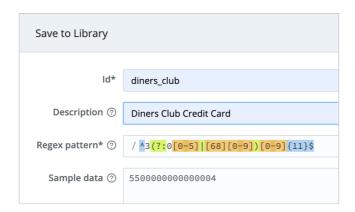
You can also add new, custom patterns to the Library. In the same modal, once you've built your pattern, click the **Save to Library** button.



Adding a custom pattern to the Regex Library from a Function's Regex modal

In the resulting modal, give your custom pattern a unique ID. Optionally, you can also provide a **Description** (name) and groom the **Sample data**. Then click

Save.



Identifying the custom pattern

Your custom pattern will now reside in the Regex Library. It will be available to Functions using the same typeahead assist as Cribl's pre-built patterns.

Cribl vs. Custom and Priority

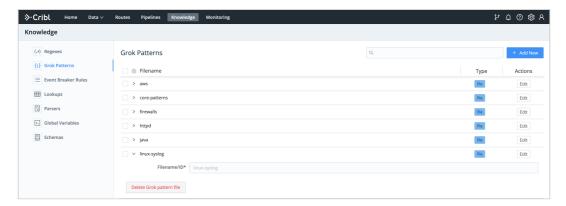
Within the Library, patterns shipped by Cribl will be listed under the **Cribl** tab, while those built by users will be found under **Custom**. Over time, Cribl LogStream will ship more patterns, and this distinction allows for both sets to grow independently.

In the case of an ID/Name conflict, the Custom pattern takes priority in listings and search. For example, if a Cribl-provided pattern and a Custom one are both named <code>ipv4</code>, the one from Cribl will not be displayed or delivered as a search result.

Grok Patterns Library

What Is the Grok Patterns Library

Cribl LogStream ships with a Grok Patterns Library that contains a set of prebuilt common patterns, organized as files.



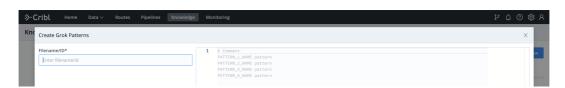
Grok Patterns Library

Managing Library Patterns

You can access the Grok Patterns Library in the UI by selecting **Knowledge** > **Grok Patterns**. The library contains several pattern files that Cribl provides for basic Grok scenarios, and is searchable.

To edit a pattern file, click **Edit** in its **Actions** column.

To create a new pattern file, click + Add New. In the resulting Create Grok Patterns modal, assign a unique Filename/ID, populate the file with patterns, then click Save.



Adding Grok patterns

i Pattern files reside in: \$CRIBL_HOME/(default|local)/cribl/grok-patterns/

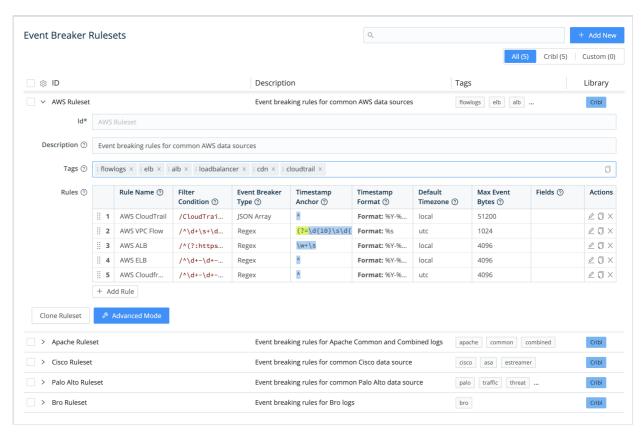
Using Grok Patterns

In the current LogStream version, you apply Grok patterns by inserting a Grok Function into a Pipeline, then manually typing or pasting patterns into the Pattern field(s).

Event Breakers

What Are Event Breakers

Event Breakers help break incoming streams of data into discrete events. You access the Event Breakers management interface under **Knowledge** > **Event Breakers**. On the resulting **Event Breaker Rulesets** page, you can edit, add, delete, search, and tag Event Breaker rules and rulesets, as necessary.



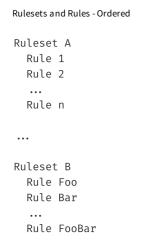
Event Breaker Rulesets page

Event Breaker Rulesets

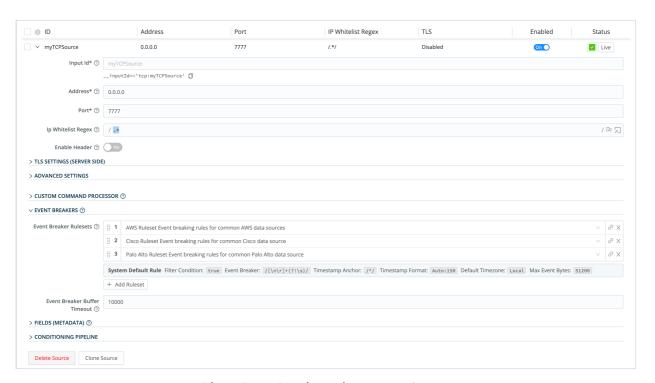
Rulesets are **collections of Event Breaker rules** that are associated with **Sources**. Rules define configurations needed to break down a stream of data into events.

i Event Breakers are accessible only on Sources that require incoming events to be broken into a better-defined format. Check individual LogStream Sources' documentation for Event Breaker support.

Rules within a ruleset are ordered and evaluated top->down. One or more rulesets can be associated with a Source, and these rulesets are also evaluated top->down. For a stream from a given Source, the first matching rule goes into effect.



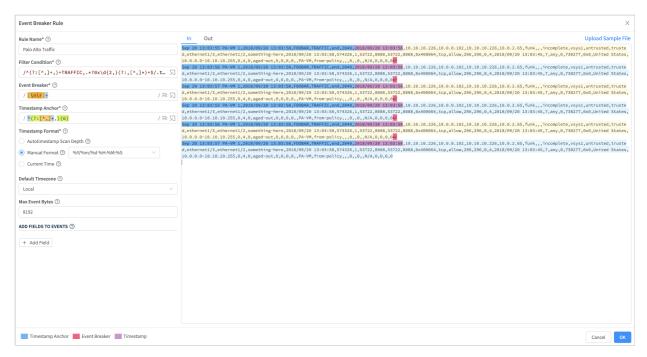
An example of multiple rulesets associated with a Source:



Three Event Breaker rulesets on a Source

Rule Example

This rule breaks on newlines and uses Manual timestamping **after** the sixth comma, as indicated by this pattern: $^(?:[^,]*,){6}$.



An Event Breaker rule

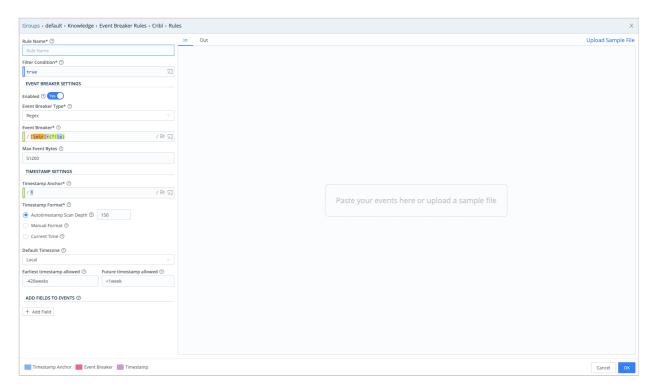
System Default Rule

The system default rule functionally sits at the bottom of the ruleset/rule hierarchy (but is built-in and not displayed on the Event Breakers page), and goes into effect if there are no matching rules:

- Filter Condition defaults to true
- Event Breaker to [\n\r]+(?!\s)
- Timestamp anchor to ^
- Timestamp format to Auto and a scan depth of 150 bytes
- Max Event Bytes to 51200
- Default Timezone to Local

How Do Event Breakers Work

On the **Event Breaker Rulesets** page (see screenshot above), click + **Add New** to create a new Event Breaker ruleset. Click + **Add Rule** within a ruleset to add a new Event Breaker.



Adding a new Event Breaker rule

Each Event Breaker includes the following components, which you configure from top to bottom in the above **Event Breaker Rule** modal:

Filter Condition

As a stream of data moves into the engine, a rule's filter expression is applied. If the expression evaluates to true, the rule configurations are engaged for the entire duration of that stream. Else, the next rule down the line is evaluated.

Event Breaker Type

After a breaker pattern has been selected, it will apply on the stream **continuously**. See below for specific information on different Event Breaker Types.

Timestamp Settings

After events are synthesized out of streams, LogStream will attempt timestamping. First, a timestamp anchor will be located inside the event. Next, starting there, the engine will try to do one of the following:

- Scan up to a configurable depth into the event and autotimestamp, or
- Timestamp using a manually supplied strptime format, or
- Timestamp the event with the current time.

The closer an anchor is to the timestamp pattern, the better the performance and accuracy – especially if multiple timestamps exist within an event. For the manually supplied option, the anchor

must lead the engine right before the timestamp pattern begins.



Anchors preceding timestamps

☐ This timestamping executes the same basic algorithm as the Auto Timestamp Function and the C.Time.timestampFinder() native method.

Add Fields to Events

After events have been timestamped, one or more fields can be added here as key-value pairs. In each field's **Value Expression**, you can fully evaluate the field value using JavaScript expressions.

Event Breaker Types

Several types of Event Breaker can applied to incoming data streams:

- Regex
- File Header
- JSON Array
- JSON New Line Delimited
- Timestamp
- CSV

Regex

The Regex breaker uses regular expressions to find breaking points in data streams.

After a breaker regex pattern has been selected, it will apply on the stream **continuously**. Breaking will occur at the beginning of the match, and the matched content will be consumed/thrown away.

If necessary, you can use a positive lookahead regex to keep the content – e.g.: (?=pattern)

Capturing groups are **not allowed** to be used anywhere in the Event Breaker pattern, as they will further break the stream – which is often undesirable. Breaking will also occur if **Max Event Bytes** has been reached.

Example

Break after a newline or carriage return, but only if followed by a timestamp pattern:

Event Breaker: $[\n\r]+(?=\d+-\d+\s\d+:\d+:\d+)$

```
Sample Event - Multiline
--- input ---
2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can program the aux:
   Try to connect the FTP sensor, maybe it will connect the digital bus!
   Try to navigate the AGP panel, maybe it will quantify the mobile alarm!
2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can program the aux:
   Try to connect the FTP sensor, maybe it will connect the digital bus!
   Try to navigate the AGP panel, maybe it will quantify the mobile alarm!
--- output event 1 ---
  "_raw": "2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can prog
  "_time": 1589920332
}
--- output event 2 ---
  "_raw": "2020-05-19 16:32:12 moen3628 ipsum[5213]: Use the mobile TCP feed, then you can prog
   _time": 1589920332
}
```

File Header

You can use the File Header breaker to break files with headers, such as IIS or Bro logs. This type of breaker relies on a header section that lists field names. The header section is typically present at the top of the file, and can be single-line or greater.

After the file has been broken into events, fields will also be extracted, as follows:

- Header Line: Regex matching a file header line. For example, ^# .
- Field Delimiter: Field delimiter regex. For example, \s+.
- **Field Regex**: Regex with one capturing group, capturing all the fields to be broken by field delimiter. For example, ^#[Ff]ields[:]?\s+(.*)
- Null Values: Representation of a null value. Null fields are not added to events.
- Clean Fields: Whether to clean up field names by replacing non <code>[a-zA-Z0-9]</code> characters with

Example

1331904609.190000

Using the values above, let's see how this sample file breaks up:

- 192.168.202.83 48516 192.168.207.4 53

udp

```
--- output event 1 ---
{
 " raw": "1331904608.080000
                                 - 192.168.204.59 137 192.168.204.255 137
                                                                                       udp",
 "ts": "1331904608.080000",
 "id_orig_h": "192.168.204.59",
 "id_orig_p": "137",
 "id_resp_h": "192.168.204.255",
 "id_resp_p": "137",
  "proto": "udp",
  _time": 1331904608.08
--- output event 2 ---
 "_raw": "1331904609.190000
                                - 192.168.202.83 48516 192.168.207.4 53
                                                                                       udp",
 "ts": "1331904609.190000",
 "id_orig_h": "192.168.202.83",
 "id_orig_p": "48516",
 "id_resp_h": "192.168.207.4",
 "id_resp_p": "53",
 "proto": "udp",
 "_time": 1331904609.19
```

JSON Array

You can use the JSON Array to extract events from an array in a JSON document (e.g., an Amazon CloudTrail file).

- Array Field: Optional path to array in a JSON event with records to extract. For example,
- Timestamp Field: Optional path to timestamp field in extracted events. For example, eventTime or level1.level2.eventTime.
- **JSON Extract Fields**: Enable this slider to auto-extract fields from JSON events. If disabled, only _raw and time will be defined on extracted events.
- Timestamp Format: If JSON Extract Fields is set to No, you must set this to Autotimestamp or Current Time. If JSON Extract Fields is set to Yes, you can select any option here.

Example

Using the values above, let's see how this sample file breaks up:

```
--- output event 2 --- {
    "_raw": "{\"eventVersion\":\"1.05\",\"eventTime\":\"2020-04-08T01:35:56Z\",\"eventSource\":\"6
    "_time": 1586309756,
    "cribl_breaker": "j-array"
}
```

JSON New Line Delimited

You can use the JSON New Line Delimited breaker to break and extract fields in newline-delimited JSON streams.

Example

Using default values, let's see how this sample stream breaks up:

```
Sample Event - Newline Delimited JSON Breaker
```

```
--- input ---
{"time":"2020-05-25T18:00:54.201Z","cid":"w1","channel":"clustercomm","level":"info","message":'
 \{ "time": "2020-05-25T18:00:54.246Z", "cid": "w0", "channel": "clustercomm", "level": "info", "message": "compared to the compared to the c
 --- output event 1 ---
                "\_raw": "{\ ''time\ '':\ ''2020-05-25T18:00:54.201Z\ '',\ ''cid\ '':\ ''w1\ '',\ ''channel\ '':\ ''clustercomm\ '',\ ''lustercomm'',\ ''lust
              "time": "2020-05-25T18:00:54.201Z",
                "cid": "w1",
                "channel": "clustercomm",
                 "level": "info",
                 "message": "metric sender",
                "total": 720,
                "dropped": 0,
                    "_time": 1590429654.201,
}
 --- output event 2 ---
              "_raw": "{\"time\":\"2020-05-25T18:00:54.246Z\",\"cid\":\"w0\",\"channel\":\"clustercomm\",\"left and the control of the co
                "time": "2020-05-25T18:00:54.246Z",
                "cid": "w0",
                "channel": "clustercomm",
                "level": "info",
                "message": "metric sender",
                "total": 720,
                "dropped": 0,
                  "_time": 1590429654.246,
```

Timestamp

You can use the Timestamp breaker to break events at the beginning of any line in which LogStream finds a timestamp. This type enables breaking on lines whose timestamp pattern is not known ahead of time.

Example

Using default values, let's see how this sample stream breaks up:

Sample Event - Timestamp Based Breaker --- input --{"level":"debug","ts":"2021-02-02T10:38:46.365Z","caller":"sdk/sync.go:42","msg":"Handle ENIConffglevel":"debug","ts":"2021-02-02T10:38:56.365Z","caller":"sdk/sync.go:42","msg":"Handle ENIConffglevel":"debug","ts":"2021-02-02T10:38:46.365Z\",\"caller\":\"sdk/sync.go:42\""_time": 1612262326.365 } --- output event 2 --{ "_raw": "{\"level\":\"debug\",\"ts\":\"2021-02-02T10:38:56.365Z\",\"caller\":\"sdk/sync.go:42\""_time": 1612262336.365 }

CSV

The CSV breaker extracts fields in CSV streams that include a header line. Selecting this type exposes these extra fields:

- Delimiter: Delimiter character to use to split values. Defaults to: ,
- Quote Char: Character used to quote literal values. Defaults to: "
- Escape Char: Character used to escape the quote character in field values. Defaults to: "

Example: Using default values, let's see how this sample stream breaks up:

Example

Using default values, let's see how this sample stream breaks up:

```
Sample Event - CSV Breaker
--- input ---
time,host,source,model,serial,bytes_in,bytes_out,cpu
1611768713, "myHost1", "anet", "cisco", "ASN4204269", 11430, 43322, 0.78
1611768714, "myHost2", "anet", "cisco", "ASN420423", 345062, 143433, 0.28
--- output event 1 ---
 "time": "1611768713",
  "host": "myHost1",
  "source": "anet",
  "model": "cisco",
  "serial": "ASN4204269",
 "bytes_in": "11430",
 "bytes_out": "43322",
  "cpu": "0.78",
  "_time": 1611768713
```

```
}
--- output event 2 ---
{
    "_raw": "\"1611768714\",\"myHost2\",\"anet\",\"cisco\",\"ASN420423\",\"345062\",\"143433\",\"(
    "time": "1611768714",
    "host": "myHost2",
    "source": "anet",
    "model": "cisco",
    "serial": "ASN420423",
    "bytes_in": "345062",
    "bytes_out": "143433",
    "cpu": "0.28",
    "_time": 1611768714
}
```

☐ With **Type: CSV** selected, an Event Breaker will properly add quotes around all values, regardless of their initial state.

Cribl versus Custom Rulesets

Event Breaker rulesets shipped by Cribl will be listed under the **Cribl** tag, while user-built rulesets will be listed under **Custom**. Over time, Cribl will ship more patterns, so this distinction allows for both sets to grow independently. In the case of an ID/Name conflict, the Custom pattern takes priority in listings and search.

Lookups Library

What Are Lookups

Lookups are data tables that can be used in Cribl LogStream to enrich events as they're processed by the Lookup Function. You can access the Lookups library under **Knowledge** > **Lookups**, which provides a management interface for all lookups.

This library is searchable, and each lookup can be tagged as necessary. There's full support for .csv files. Compressed files are supported, but must be in gzip format (.gz extension). You can add files in multimedia database (.mmdb) binary format, but you cannot edit these binary files through LogStream's UI.



Lookups Library

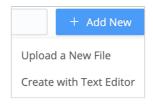
How Does the Library Work

In single-instance deployments, all files handled by the interface are stored in \$CRIBL_HOME/data/lookups . In distributed deployments, the storage path on the Master Node is \$CRIBL_HOME/groups/<groupname>/data/lookups/for each Worker Group.

i For large and/or frequently replicated lookup files, you might want to bypass the Lookups Library UI and instead manually place the files in a different location. This can both reduce deploy traffic and prevent errors with LogStream's default Git integration. For details, see Managing Large Lookups.

Adding Lookup Files

To upload or create a new lookup file, click + Add New, then click the appropriate option from the drop-down.

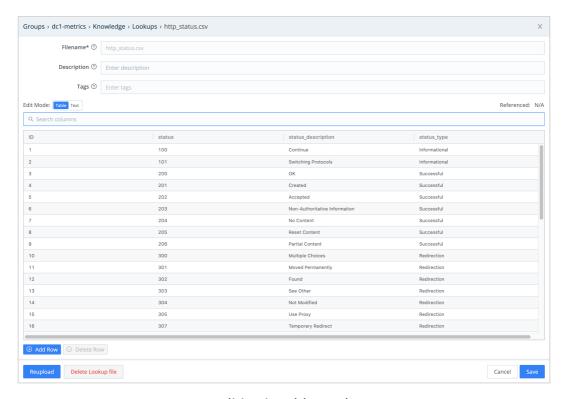


Adding a lookup file

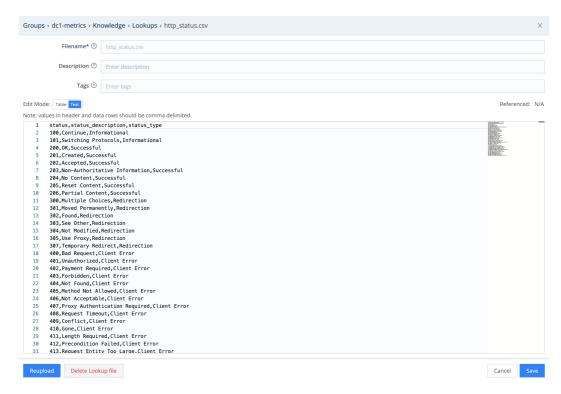
Modifying Lookup Files

To re-upload, expand, edit, or delete an existing .csv or .gz lookup file, click its row on the Lookups page. (No editing option is available for .mmdb files; you can only re-upload or delete these.)

In the resulting modal, you can edit files in **Table** or **Text** mode. However, **Text** mode is disabled for files larger than 1 MB.



Editing in table mode



Editing in text mode

Memory Sizing for Large Lookups

For large lookup files, you'll need to provide extra memory beyond basic requirements for LogStream and the OS. To determine how much extra memory to add to a Worker Node for a lookup, use this formula:

```
Lookup file's uncompressed size (MB) * 2.25 (to 2.75) * numWorkerProcesses = Extra RAM required for lookup
```

For example, if you have a lookup file that is 1 GB (1,000 MB) on disk, and three Worker Processes, you could use an average 2.50 as the multiplier:

```
1,000 * 2.50 * 3 = 7,500
```

In this case, the Node's server or VM would need an extra 7,500 MB (7.5 GB) to accommodate the lookup file across all three worker processes.

What's with That Multiplier?

We've cited a squishy range of 2.25–2.75 for the multiplier, because we've found that it varies inversely with the number of columns in the lookup file:

• The fewer columns, the higher the extra memory overhead (2.75 multiplier).

• The more columns, the lower the overhead (2.25 multiplier).

In Cribl's testing:

- 5 columns required a multiplier of 2.75
- 10 columns required a multiplier of only 2.25.

These are general (not exact) guidelines, and this multiplier depends only on the lookup table's number of columns. The memory overhead imposed by each additional row appears to decline when more columns are present in the data.

Other Scenarios

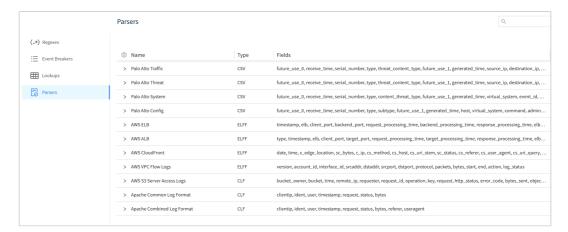
See also:

- Lookup Function.
- Ingest-time Lookups use case.
- Managing Large Lookups use case.
- Redis Function for faster lookups using a Redis integration (bypasses the Lookups Library).

Parsers Library

What Are Parsers

Parsers are definitions and configurations for the Parser Function. You can find the library under **Knowledge** > **Parsers**, and its purpose is to provide an interface for creating and editing Parsers. The library is searchable, and each parser can be tagged as necessary.



Parsers Library

Parsers can be used to **extract** or **reserialize** events. See **Parser Function** page for examples.

Supported Parser Types:

- CSV Parse and reserialize comma-separated values.
- ELFF Parse and reserialize events in Extended Log File Format.
- CLF Parse and reserialize events in Common Log Format.

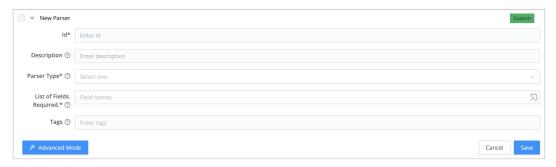
Creating a Parser

To create a parser, follow these steps:

- 1. Go to Knowledge > Parsers and click Add New.
- 2. Enter a unique ID.
- 3. Optionally, enter a **Description**.

- 4. Select a **Parser type** (see the supported types above).
- 5. Enter the **List of fields** expected to be extracted, in order.

 Click this field's Maximize icon (far right) if you'd like to open a modal where you can work with sample data and iterate on results.
- 6. Optionally, enter any desired Tags.



Adding a new parser

Schema Library

What Are Schemas

Schemas are JSON definitions that are used to validate of JSON events. They're based on the popular JSON Schema standard, and all schemas matching draft version 2019-09 are supported. You can find the library under **Knowledge** > **Schemas**. Its purpose is to provide an interface for creating, editing, and maintaining Schemas.

You validate a schema using the C.Schema('<schema name>').validate(<object field>) built-in method. This function can be called anywhere in Cribl LogStream that JavaScript expressions are supported.

Typical use cases for Schema validation:

- Making a decision before sending an event down to a destination.
- Making a decision before accepting an event. (E.g., drop an event if invalid.)
- Making a decision to route an event based on the result of validation.

Example

To add this example JSON Schema, go to **Knowledge** > **Schemas** and click **Add New**. Enter the following:

- ID: schema1.
- Description: (Enter your own description here.)
- Schema: Paste the following schema.

JSON Schema-Sample

{
 "\$id": "https://example.com/person.schema.json",
 "\$schema": "http://json-schema.org/draft-07/schema#",
 "title": "Person",
 "type": "object",
 "required": ["firstName", "lastName", "age"],
 "properties": {
 "firstName": {
 "type": "string",
 "description": "The person's first name."
 },
 "lastName": {
 "type": "string",
 "description": "The person's first name."
 },
 "lastName": {
 "type": "string",
 }
}

```
"description": "The person's last name."
},
"age": {
   "description": "Age in years which must be equal to or greater than zero.",
   "type": "integer",
   "minimum": 0,
   "maximum": 42
}
}
}
```

Assume that events look like this:

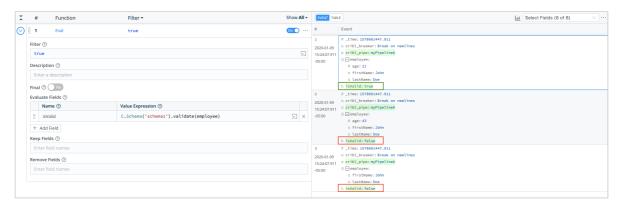
```
Events
{"employee":{"firstName": "John", "lastName": "Doe", "age": 21}}
{"employee":{"firstName": "John", "lastName": "Doe", "age": 43}}
{"employee":{"firstName": "John", "lastName": "Doe"}}
```

To validate whether the employee field is valid per schema1, we can use the following:

```
C.Schema('schema1').validate(employee)
```

Results:

- First event is valid.
- Second event is **not valid** because age is greater than the maximum of 42.
- Third event is **not valid** because age is missing.



Schema validation results for the above events

Global Variables Library

What Are Global Variables

Global Variables are reusable JavaScript expressions that can be accessed in Functions in any Pipeline. You can access the library under Knowledge > Global Variables.

Typical use cases for Global Variables include:

- Storing a constant that you can reference from any Function in any Pipeline.
- Storing a relatively long value expression, or one that uses one or more arguments.

Global Variables can be of the following types:

- Number
- String
- Boolean
- Object
- Array
- Expression

Global Variables can be accessed via C.vars. – which can be called anywhere in Cribl LogStream that JS expressions are supported. Typeahead is provided. More on Cribl Expressions here.

Examples

Scenario 1:

Assign field foo the value in the Answer Global Variable.

- Global Variable Name: theAnswer <-- ships with Cribl LogStream by default.
- Global Variable Value: 42
- Sample Eval Function: foo = C.vars.theAnswer

Scenario 2:

Assign field nowEpoch the current time, in epoch format.

- Global Variable Name: epoch <-- ships with Cribl LogStream by default.
- Global Variable Value: Date.now()/1000
- **Sample Eval Function:** nowEpoch = C.vars.epoch()

Scenario 3:

Create a new field called storage, by converting the value of event field size to human-readable format.

- Global Variable Name: convertBytes <-- ships with Cribl LogStream by default
- Global Variable Value: `\${Math.round(bytes / Math.pow(1024, (Math.floor(Math.log(bytes) / Math.log(1024)))),
 2)}\${['Bytes', 'KiB', 'MiB', 'GiB', 'TiB', 'PiB', 'EiB', 'ZiB', 'YiB'][(Math.floor(Math.log(bytes) / Math.log(1024)))]}`
- Global Variable Argument: bytes
- **Sample Eval Function:** storage = C.vars.convertBytes(size)

Note the use of bytes here as an argument.

Use Cases

Ingest-time Fields

Adding Fields to Data in Motion

To add new fields to any event, we use the out-of-the-box **Eval** Function. We can either apply a Filter to select the events, or we can use the default true Filter expression to apply the Function to all incoming events.

Adding Fields Example

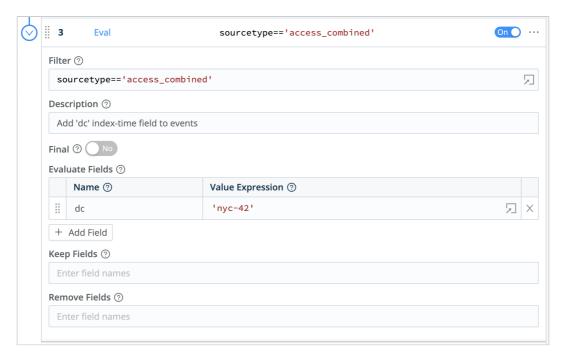
Let's see how we add dc::nyc-42 to all events with sourcetype='access_combined':

- First make sure you have a Route and Pipeline configured to match desired events.
- Next, let's add a **Eval** function to it:



Defining the Eval Function's filter expression

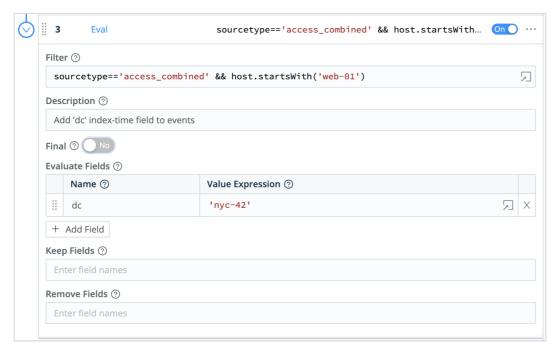
• Next, let's click on + Add Field, add our dc field, and click Save.



Adding the dc field

To confirm, verify that this search returns results: sourcetype="access_combined" dc::nyc-42

• You can add more conditions to the filter, if you'd like. For example, to limit the field to only events from hosts that start with web-01, we can change the filter input as below:



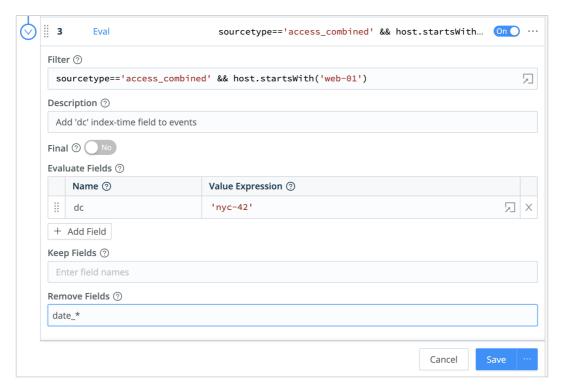
Refining the filter

This is a **very** powerful method to change incoming events in real time. In addition to providing the right context at the right time, users can further benefit substantially by using tstats for **faster** analytics.

Removing Fields

You can remove fields by listing and/or wildcarding field names. Let's see how we can remove all fields that start with date_ ::

- First, make sure you have a Route and Pipeline configured to match desired events.
- Next, let's add a **Eval** function to it (as above).
- Next, in **Remove Fields**, add date_* and hit Save.



Goodbye date_ field

To confirm, verify that this search: sourcetype="access_combined" date_minute=* will soon stop returning results. Enjoy a more efficient Splunk!

Ingest-time Lookups

Enriching Data in Motion

To enrich events with new fields from external sources (such as .csv files), we use LogStream's out-of-the-box Lookup Function. Ingestion-time lookups are not only great for normalizing field names and values, but also ideal for use cases where:

- Fast access via the looked-up value is required. For example, when you don't have a datacenter field in your events, but you do have a host-to-datacenter map, and you need to search by datacenter.
- Looked-up information must be temporally correct. For example, assume that you have a highly dynamic infrastructure, and you need to resolve a resource name (e.g., a container name) to its address. You can't afford to defer this to search time/runtime, as the resource and its records might no longer exist.
 - i To use large binary databases (like GeoIP .mmdb files) for LogStream lookups, see
 Managing Large Lookups. To achieve faster lookups, use LogStream's Redis Function.

Working with Lookups - Example 1

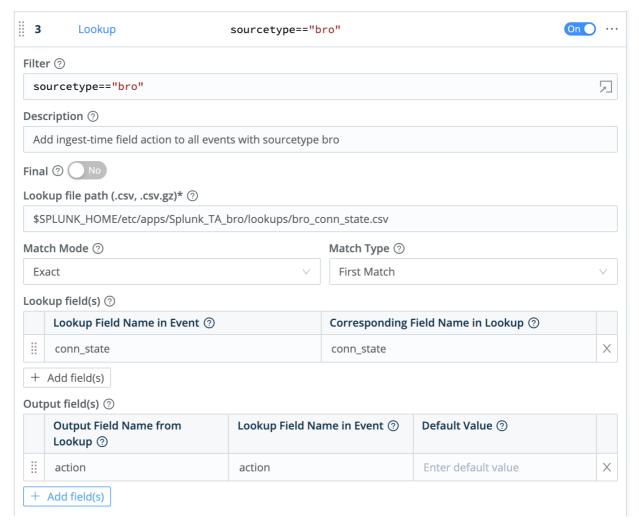
Let's assume we have the following lookup file. Given the field <code>conn_state</code> in an event, we would like to add a corresponding ingestion-time field called <code>action</code>.

```
action, "conn_state", "conn_state_meaning"
dropped,S0, "Connection attempt seen, no reply."
allowed,S1, "Connection established, not terminated."
allowed,SF, "Normal establishment and termination."
blocked,REJ, "Connection attempt rejected."
allowed,S2, "Connection established and close attempt by originator seen (but no reply from respondence), "Connection established and close attempt by responder seen (but no reply from originallowed,RSTO, "Connection established, originator aborted (sent a RST)."
allowed,RSTR, "Established, responder aborted."
dropped,RSTOSO, "Originator sent a SYN followed by a RST, we never saw a SYN-ACK from the respondence, "Conped, SHR, "Responder sent a SYN ACK followed by a FIN, we never saw a SYN ACK from the respondence, SHR, "Responder sent a SYN followed by a FIN, we never saw a SYN from the originator allowed,OTH, "No SYN seen, just midstream traffic (a 'partial connection' that was not later close
```

First, make sure you have a Route and Pipeline configured to match desired events.

Next, let's add a **Lookup** function to the Pipeline, with these settings:

- Lookup file path: \$SPLUNK_HOME/etc/apps/Splunk_TA_bro/lookups/bro_conn_state.csv (note that Environment variables are allowed in the path).
- Lookup Field Name in Event set to conn_state .
- Corresponding Field Name in Lookup set to conn_state.
- Output Field Name from Lookup set to action .
- Lookup Field Name in Event set to action.



Lookup Function to add action field

To confirm success, verify that this search returns expected results: sourcetype="bro" action::allowed.Change the action value as necessary.

Working with Lookups – Example 2

Let's assume we have the following lookup file, and given **both** the fields impact and priority in an event, we would like to add a corresponding ingestion-time field called severity.

```
cisco_sourcefire_severity.csv
impact,priority,severity
1,high,critical
2,high,critical
```

```
3, high, high
4, high, high
0,high,high
"*",high,high
"*", medium, medium
1,low,medium
2,low,medium
3,low,low
4,low,low
0,low,low
"*",low,low
1, none, low
2, none, low
3, none, informational
4, none, informational
0, none, informational
"*", none, informational
```

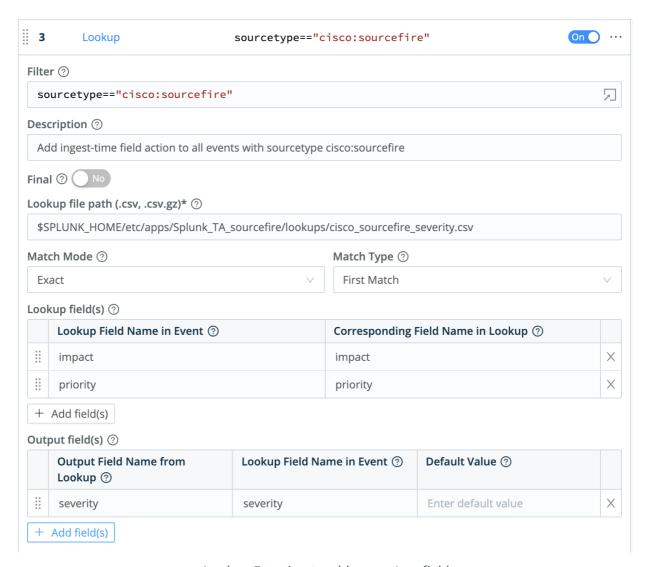
First, make sure you have a Route and Pipeline configured to match desired events.

Next, let's add a **Lookup** function to the Pipeline, with these settings:

• Lookup file path:

\$SPLUNK_HOME/etc/apps/Splunk_TA_sourcefire/lookups/cisco_sourcefire_severity.csv (note that Environment variables are allowed in the path).

- Lookup Field Name(s) in Event set to impact and priority.
- Corresponding Field Name(s) in Lookup setto impact and priority.
- Output Field Name from Lookup set to severity.
- Lookup Field Name in Event set to severity.



Lookup Function to add severity field

To confirm success, verify that this search returns expected results:

sourcetype="cisco:sourcefire" severity::medium.Changethe severity value as necessary.

Sampling

Sampling at Ingest-Time

Let's say that you wanted to analyze and troubleshoot with **highly verbose/voluminous** data – for example, CDN logs, ELB Access Logs, or VPC Flows – but you were concerned about storage requirements and search performance. With Sampling, you can bring in enough samples that your analysis remains statistically significant, and also do all the necessary troubleshooting.

See the example below, or see more details in Access Logs and Firewall Logs.

Sampling Example

Let's use the out-of-the-box **Sampling** function to sample all events from sourcetype='access_combined' where status is '200'. We'll sample these at 5:1 (and all other events at 1:1). This should lower the volume of all verbose successes (200 s), but still bring in **all** potentially erroneous events (400 s, 500`s, etc.) that can be used for troubleshooting.

- First, make sure you have a Route and Pipeline configured to match desired events.
- Next, let's add a Regex Extract Function to extract the status field from
 _raw , and let's call the resulting field __status . Remember, fields that
 start with __ are special fields in Cribl LogStream, and can be used
 anywhere in a Pipeline.



Extracting the __status field

Next, let's add a **Sampling** function, and scope it to all events where sourcetype='access_combined'. Let's apply a filter condition of __status = 200, and a Sample Rate of 5.



Sampling success responses

To confirm that sampling works, compare the event count of all 200 s before and after.

i Each time an event goes through the **Sampling** function, an indextime sampled::<rate> field is added to it. You can use this field in your statistical functions, as necessary.

Access Logs: Apache, ELB, CDN, S3, etc.

Recipe for Sampling Access Logs

Access logs are extremely common. They're often emitted by web servers or similar/related technologies (proxies, loadbalancers, etc.), and tend to be highly voluminous. Typical examples include Apache access logs, and CDN logs such as those from Amazon Cloudfront, Amazon S3 Server Access Logs, AWS ELB Access Logs, etc.

For large installations, bringing everything into an analytics tool is often so cost-prohibitive (storage, resources, license, etc.) that most users don't even bother. However, some of the logs contain relevant information when looked at individually (e.g., errors). The much larger majority contains relevant information when looked at in the aggregate (e.g., successes to determine traffic patterns, etc.).

It would be great if we could find a middle ground. With the Sampling Function, you can! Specifically, you can:

- Ingest enough sample events from the majority category that your aggregate analysis remains statistically significant.
- Ingest *all* events from the minority categories, and perform troubleshooting and introspection with full-fidelity data.

Using status as the Sampling Condition

Most of the access logs (including the ones mentioned above) have very similar formats. One quick way to sample is to look at the value of the status field.

2XX s indicate success and tend to be, by far, the most common ones – with

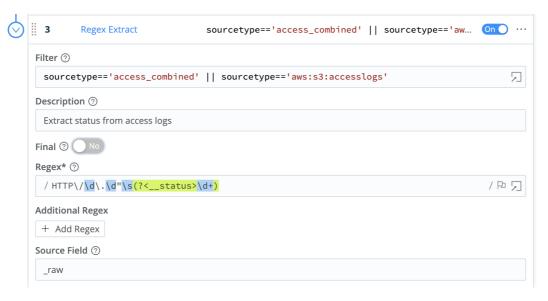
200 being the top. **Therefore, 200 is the perfect candidate for sampling.** All other statuses occur much less frequently, indicate conditions that often need to be looked at, and can be brought in with full fidelity.

Sample Status 200 at 5:1

1. Add a Regex Extract Function that looks at these sourcetypes:

```
sourcetype='access_combined' ||
sourcetype='aws:s3:accesslogs'
```

2. Configure that Function to extract a field called __status with this
 regex: /HTTP\/\d\.\d"\s(?<__status>\d+)/



Defining the Regex Extract Function

- 3. Add a Sampling Function to sample 5:1 when __status=200.
- 4. Save.



Sampling success reponses

Note About Sampling

Each time an event goes through the **Sampling** Function, an index-time sampled::<rate> field is added to it. Use this field in your statistical Functions, as necessary.

Other Sourcetypes

Examples of other sourcetypes that will benefit from sampling, but might need a different __status extraction regex:

Sourcetype	Filter Expression
Amazon Cloudfront Access Logs	<pre>sourcetype='aws:cloudfront:accesslogs'</pre>
Amazon ELB Access Logs	sourcetype='aws:elb:accesslogs'

Firewall Logs: VPC Flow Logs, Cisco ASA, Etc.

Recipe for Sampling Firewall Logs

Firewall logs are another source of important operational (and security) data. Typical examples include Amazon VPC Flow Logs, Cisco ASA Logs, and other technologies such as Juniper, Checkpoint, pfSense, etc.

As with Access Logs, bringing in everything for operational analysis might be cost-prohibitive. But sampling with Cribl LogStream can help you:

- Ingest enough sample events from the majority category that your aggregate analysis remains statistically significant. E.g., sample all ACCEPT s at 5:1.
- Ingest all events from the minority categories, and perform troubleshooting and introspection with full-fidelity data. E.g., bring in all REJECT s.

Sampling VPC Flow Logs

AWS' VPC Flow Logs feature enables you to capture information about the IP traffic going to and from network interfaces in your VPC. Flow Log data can be published to Amazon CloudWatch Logs and Amazon S3.

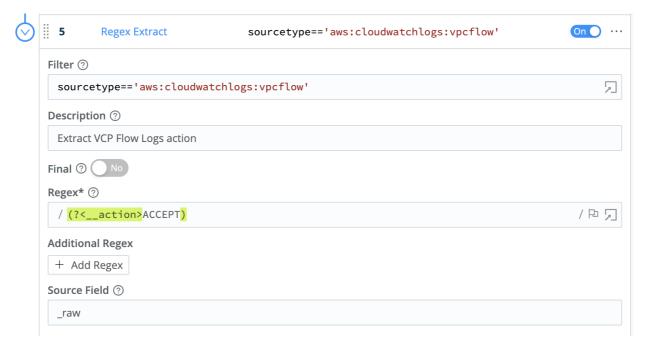
Typical VPC Flow Logs look like this:

Flow Log Records for Accepted and Rejected Traffic

```
2 123456789010 eni-abc123de 172.31.16.139 172.31.16.21 20641 22 6 20 4249 1418530010 1418530070 2 123456789010 eni-abc123de 172.31.9.69 172.31.9.12 49761 3389 6 20 4249 1418530010 1418530070 F
```

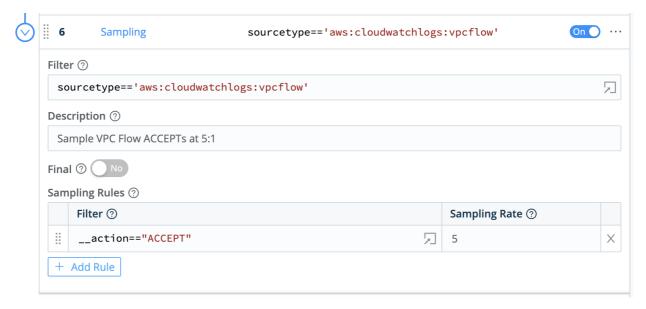
Let's use a very simple Filter condition and only look for ACCEPT events:

- 1. Add a Regex Extract Function that looks at: sourcetype='aws:cloudwatchlogs:vpcflow'
- 2. Configure that Function to extract a field called __action with this regex: /(?
 <__action>ACCEPT)/



Extracting the __action field

- 3. Add a **Sampling** Function to sample 5:1 when __action="ACCEPT".
- 4. Save.



Sampling ACCEPT events

Note About Sampling

Each time an event goes through the Sampling Function, an index-time field is added to it: sampled: <rate> . It's advisable that you use that in your statistical functions, as necessary.

Other Sourcetypes

Other sourcetypes that will benefit from sampling, but might need a different __action extraction regex:

Sourcetype	Filter Expression
Cisco ASA Logs	sourcetype='cisco:asa'
Related sourcetypes to consider sampling:	<pre>sourcetype='cisco:fwsm' sourcetype='cisco:pix'</pre>

Masking and Obfuscation

Masking and Anonymization of Data in Motion

To mask patterns in real time, we use the out-of-the-box Mask Function. This is similar to sed, but with much more powerful functionality.

Masking Capabilities

The Mask Function accepts multiple replacement rules, and accepts multiple fields to apply them to.

Match Regex is a JS regex pattern that describes the content to be replaced. It can optionally contain matching groups. By default, it will stop after the first match, but using /g will make the Function replace all matches.

Replace Expression is a JS expression or literal to replace matched content.

Matching groups can be referenced in the Replace Expression as g1, g2 ... gN, and the entire match as g0.

There are several masking methods that are available under C.Mask.:

- C.Mask.random: Generates a random alphanumeric string
- C.Mask.repeat: Generates a repeating char/string pattern, e.g., XXXX
- C.Mask.REDACTED: The literal 'REDACTED'
- C.Mask.md5: Generates a MD5 hash of given value
- C.Mask.sha1: Generates a SHA1 hash of given value
- C.Mask.sha256: Generates a SHA256 hash of given value

Almost all methods have an optional len parameter which can be used to control the length of the replacement. len can be either a number or string. If it's a string, its length will be used. For example:



Defining the replacement length

Masking Examples

Let's look at the various ways that we can mask a string like this one: cardNumber=214992458870391. The **Regex Match** we'll use is: $/(cardNumber=)(\d+)/g$. In this example:

- g0 = cardNumber = 214992458870391
- g1 = cardNumber=
- g2 = 214992458870391

Random Masking with default character length (4):

- Replace Expression: `\${g1}\${C.Mask.random()}`
- Result: cardNumber=HRhc

Random Masking with defined character length:

- Replace Expression: `\${g1}\${C.Mask.random(7)}`
- Result: cardNumber=neNSm8r

Random Masking with length preserving replacement:

- Replace Expression: `\${g1}\${C.Mask.random(g2)}`
- **Result**: cardNumber=DroJ73qmyaro51u3

Repeat Masking with default character length (4):

- Replace Expression: `\${g1}\${C.Mask.repeat()}`
- Result: Result: cardNumber=XXXX

Repeat Masking with defined character choice and length:

- Replace Expression: `\${g1}\${C.Mask.repeat(6, 'Y')}`
- **Result**: cardNumber=YYYYYY

Repeat Masking with length preserving replacement:

- Replace Expression: `\${g1}\${C.Mask.repeat(g2)}`
- **Result**: cardNumber=XXXXXXXXXXXXXXXX

Literal **REDACTED** masking:

- **Replace Expression**: `\${g1}\${C.Mask.REDACTED}`
- **Result**: cardNumber=REDACTED

Hash Masking (applies to: md5, sha1 and sha256):

- Replace Expression: `\${g1}\${C.Mask.md5(g2)}`
- **Result**: cardNumber=f5952ec7e6da54579e6d76feb7b0d01f

Hash Masking with left N-length* substring (applies to: md5, sha1 and sha256):

- Replace Expression: `\${g1}\${C.Mask.md5(g2, 12)}`
- Result: cardNumber=d65a3ddb2749
 - *Replacement length will **not** exceed that of the hash algorithm output; MD5: 32 chars, SHA1: 40 chars, SHA256: 64 chars.

Hash Masking with right N-length* substring (applies to: **md5**, **sha1** and **sha256**):

- Replace Expression: `\${g1}\${C.Mask.md5(g2, -12)}`
- Result: cardNumber= 933bfcebf992
 - *Replacement length will **not** exceed that of the hash algorithm output; MD5: 32 chars, SHA1: 40 chars, SHA256: 64 chars.

Hash Masking with length* preserving replacement (applies to: md5, sha1 and sha256):

- Replace Expression: `\${g1}\${C.Mask.md5(g2, g2)}`
- Result: cardNumber= d65a3ddb27493f5

 *Replacement length will not exceed that of the hash algorithm output;

 MD5: 32 chars, SHA1: 40 chars, SHA256: 64 chars.

Lookups as Filters for Masks

Overview

You can make your data architecture more maintainable by using Lookups to route and transform events within Cribl LogStream. This use case demonstrates an unusual solution, but one that served one Cribl customer's particular goals (which might overlap with yours):

- Ingest many hundreds of different sourcetype / index field combinations.
- Send all this data through a common Pipeline.
- Stack four Mask Functions in the Pipeline.
- Evaluate and process each sourcetype / index field combination only within its applicable Mask Functions – either two or three Masks per combination.
 - ⚠ This last restriction reduces latency, by preventing Mask Functions from evaluating non-applicable events, simply to ignore them.

Just to reiterate, this use case outlined here responded to this customer's requirements – one Pipeline combining multiple Mask Functions, for many sourcetype / index combinations.

More typically, you'd use multiple Pipelines to process different sourcetype / index combinations.

To enable this approach, the example below centralizes masking logic for multiple conditions in a Lookup table and corresponding Lookup Functions. The Lookup's output filters events to the applicable Mask Functions. Specifically, we'll show how to instruct LogStream to:

- Check for a particular index / sourcetype combination in each event,
 and
- Based on that combination, determine which Masks to apply to that event.

Design the Lookup

To use a lookup as a filter, you'd start by creating a comma-separated lookup table in this format, and adding it to LogStream:

```
index_tracker.csv

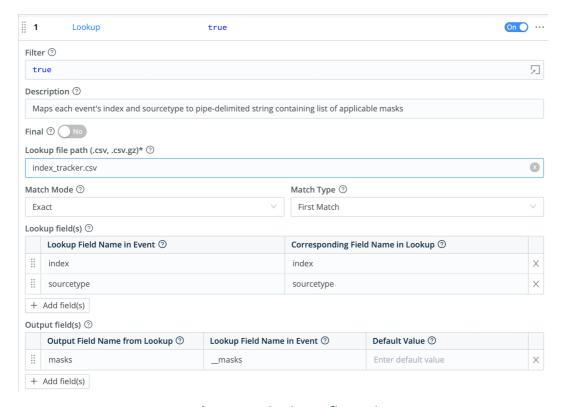
index,sourcetype,masks
apache_common, sourcetypec, ssn|credit_card|auth_token
syslog,sourcetypeb,ssn|auth_token
weblog,sourcetypea,auth_token|bearer_token
```

Below the header, each row specifies an index, a sourcetype, and (in the third column) a pipe-delimited list of applicable masks.

To make this example work, the table must have only **one** row for each index/sourcetype combination. (This unusual restriction is particular to this scenario.) So, as you build out the lookup table, you cannot add new masks for **existing** index/sourcetype combinations by appending new rows. Instead, you must modify the third column of the existing rows.

Configure the Pipeline

Create a LogStream Pipeline with a Lookup Function configured like this, pointing to your lookup table:



Lookup Function's configuration

This Function keys against both the index and sourcetype fields. When it finds a matching combination, it adds a new key-value pair to your event for future filtering.

The key of that key-value pair (namely, __masks) starts with a double underscore, to make it a LogStream internal field. This convention ensures that the key-value pair will **not** get passed along to the Destination.

However, you might prefer to export the key-value pair. For example, you might want a Splunk Destination to index the list of masks applied to a given event, alongside that event. (This approach applies to many forensic use cases.) If so, remove the double underscore from the above Function's **Lookup Field Name** in Event value, and from the subsequent Filter expressions for each Mask Function.

Each Mask Function has a JavasScript Filter that breaks the pipe-delimited string into an array, and determines whether the tag for that type of mask (e.g., bearer_auth) is in the __masks key-value pair. If so, it applies the mask processing. If not, the event moves on to the Pipeline's next Mask Function.

Here are the four Mask Functions below the Lookup Function:



Mask Functions

In this particular example, the pipe-delimited mask tags in the lookup table's third column match the Mask Functions' names, as well as matching their **Filter** conditions. This is just for simplicity – the Functions could have any names, as long as the **Filter** expressions match the tags.

Lookups and Regex Magic

Regular expressions are not just for field extractions – they can also be used inside lookup tables, and in Functions, to replace and manipulate values within fields. Let's walk through a Pipeline that demonstrates four different ways to leverage regular expressions in LogStream.

Why Lookup Tables Matter

When organizations use host naming standards, it is easy to understand things like regions, availability zones (AZs), IP addresses, and more. For example, consider an Amazon host called:

```
ec2-35-162-133-145.us-west1-a.compute.amazonaws.com
```

This is an EC2 host with a (dashed) IP address 35-162-133-145, in the us-west1 region, in Availability Zone a . You can also see the domain: compute.amazonaws.com.

While we can understand the enriched host names, we don't know which indexes to route the data to, nor which sourcetypes to assign to the events, without looking up this information from another source. Doing so is often a huge challenge for organizations. To solve this challenge, let's combine Regex Extract, Lookup, and Eval Functions with some sample events to demonstrate the power of LogStream.

Sample Events

The events below have timestamps broken out, but no indexes, sourcetypes, or other details have been assigned yet:

```
a_raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                 [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α cribl_breaker: Break on newlines
         a_raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Loo
                  king for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α cribl_breaker: Break on newlines
-06:00
          α_raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.py
                 [DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           α cribl_breaker: Break on newlines
         α_raw: Feb 06 2021 12:59:44.232 GMT: ec2-76-187-246-132.europe-west3-b.compute.amazonaws.com: cloud-init[2929]: stages.py[DEBUG]: R
                 unning module power-state-change (<module 'cloudinit.config.cc_power_stat... Show more
2021-02-06
06:59:44.232 # _time: 1612616384.232
           α cribl breaker: Break on newlines
         a _raw: Feb 06 2021 17:04:16.921 GMT: ec2-67-205-202-104.northamerica-northeast1-c.compute.amazonaws.com: cloud-init[2929]: util.py
                 [DEBUG]: Running command ['lxc-is-container'] with allowed return codes [0... Show more
2021-02-06
11:04:16.921 # _time: 1612631056.921
           a cribl breaker: Break on newlines
           \alpha_raw: Feb 06 2021 19:45:47.687 GMT: ec2-87-209-176-201.southamerica-east1-a.compute.amazonaws.com: DataSourceEc2.py[DEBUG]: Remove
                 d the following from metadata urls: ['http://instance-data.:8773']
2021-02-06
13:45:47.687 # _time: 1612640747.687
           α cribl_breaker: Break on newlines
-06:00
```

The Regex Extract Function

Before we can assign an index or sourcetype, we need to extract the host, region, az, and domain fields from the events. We can use a Regex Extract Function with this regular expression to extract all four fields:

```
GMT:\s+(?<host>[^.]+)\.(?<region>\w+-\w+\d+)-(?<az>[^.]+)\.(?<domain>[^:]+):
```

Here's that Regex Extract in a LogStream Pipeline:



Results of the Regex Extract Function

On the **OUT** tab of LogStream's Preview pane, the extracted fields of az, domain, host, and region now appear below the _raw event. You can use these extracted fields for searching in your preferred search solution.

```
IN OUT
                                                                                     Select Fields (8 of 8)
            a_raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                 [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
           a host: ec2-35-162-133-145
           a region: us-west1
          a_raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                  oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α az: b
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           α host: ec2-48-169-111-182
         α region: us-east2
3
            a _raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.p
                  y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
00:29:11.841 # _time: 1612592951.841
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           α host: ec2-21-187-232-201
           α region: asia-northeast3
```

Lookups

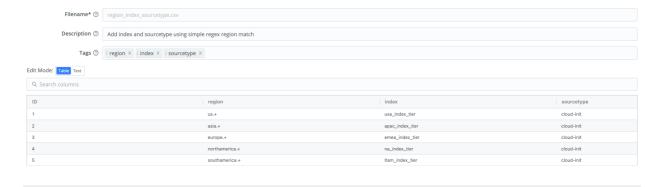
We still need to determine the index and sourcetype. LogStream's Lookup Function enriches events with external fields. We'll use it with the newly extracted region field to assign an index and sourcetype to these events.

Lookup File

First, we need to create a lookup table for the Function to reference. For this, we'll use regex again.

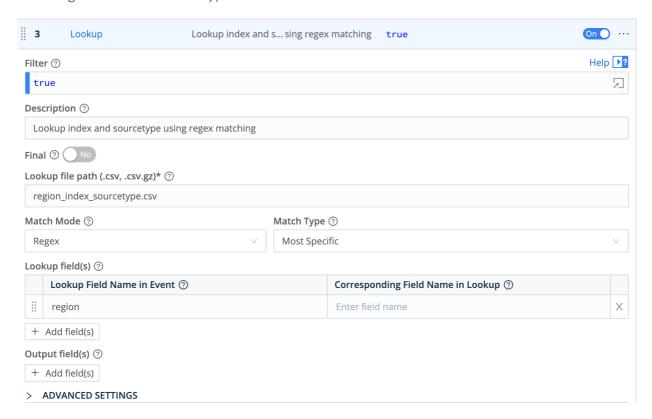
In the table below, five simple regular expressions map the extracted region field to the appropriate index and sourcetype. For example, the region us-west1-a starts with us, so it matches the first regular expression: us.+

We use this lookup table's first row to assign an index of usa_index_tier, and a sourcetype of cloud-init, to each matching event. The region patterns in the table's four remaining rows work the same way.



Lookup Function

With the lookup table saved as region_index_sourcetype.csv , the Lookup Function below matches the events' extracted region field against the regular expressions, and returns the matching index and sourcetype .



There's actually more here than meets the eye. Note that we've specified no **Output Fields**. From the Lookup Function's documentation, we know this means that the Function will default to outputting all fields in the lookup table. So we get the contents of both remaining columns in the table we saw above: index and sourcetype.

Results of the Lookup Function

With the Lookup Function added to our Pipeline, the Preview pane's **OUT** tab shows that the index and sourcetype are now added to each event.

```
a _raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                   [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://145.133.162.36/latest/api/token', 'headers': {'X-a
2021-02-05
                   ws-ec2-metadata-token-ttl-seconds': '21600', 'User-Agent': 'Cloud-Init/19.3-5.amzn2'}, 'allow_redirects': True, 'method':
20:18:31.286
                   'PUT', 'timeout': 1.0} configuration Showless
-06:00
            # _time: 1612577911.286
            α az: a
            α cribl_breaker: Break on newlines
            α cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
            α host: ec2-35-162-133-145
            α index: usa_index_tier
            a region: us-west1
            a sourcetype: cloud-init
           a _raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                   oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α cribl_breaker: Break on newlines
            a cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
           α host: ec2-48-169-111-182
           α index: usa_index_tier
            α region: us-east2
           a sourcetype: cloud-init
            a _raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.p
                   y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
            α domain: compute.amazonaws.com
           a host: ec2-21-187-232-201
            α index: apac_index_tier
            α region: asia-northeast3
            a sourcetype: cloud-init
```

Getting Host IP Address from Host Name

Since the IP address is present in the host field, we can create the host_ip field using an Eval Function with this replace method:

```
host.replace(/\w+-(\d+)-(\d+)-(\d+)/,'$1.$2.$3.$4')
```

This regular expression uses capture groups and pulls the four IP octets present in the hostname to build the <code>host_ip</code>. These four capture groups are noted as \$1.\$2.\$3.\$4, respectively. This method is very fast, and removes the need to perform a DNS lookup from the <code>host</code> field to get the host's IP address. Magic!



Results of the Eval Function and Replace Method

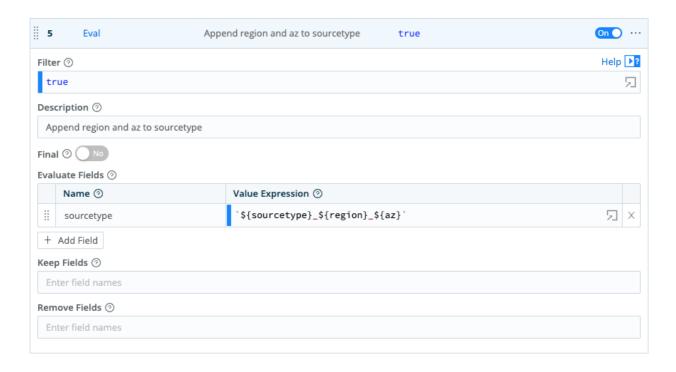
The host_ip field is now added to the events, displayed below host:

```
IN OUT
                                                                                      Select Fields (11 of 11)
           a _raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                  [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           \alpha domain: compute.amazonaws.com
           a host: ec2-35-162-133-145
           a host_ip: 35.162.133.145
           α index: usa_index_tier
           a region: us-west1
           α sourcetype: cloud-init
           a _raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                  oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α az: b
           α cribl_breaker: Break on newlines
           \alpha cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           a host: ec2-48-169-111-182
           α host_ip: 48.169.111.182
           α index: usa index tier
           α region: us-east2
         \alpha sourcetype: cloud-init
           a_raw: Feb 06 2021 06;29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]; atomic_helper.p
                  y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           α az: a
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           a host: ec2-21-187-232-201
           a host_ip: 21.187.232.201
           α index: apac_index_tier
           α region: asia-northeast3
           a sourcetype: cloud-init
```

Customizing the Sourcetype

Finally, let's put some sense into the sourcetype field, using another Eval Function. By combining the values of the $\{\text{sourcetype}\}_{\{\text{region}\}_{\{az\}}\}$, the sourcetype becomes cloud-init_us-west1_a - so now you can understand much more about the sourcetype at a glance.

Examine this Eval Function's value expression, taking careful note of the backticks (` `) and braces ({ }) that surround the field names, and the underscore (_) that separates them.



Results of the Eval Function to Combine Values

Take a look at the updated sourcetypes, and enjoy exploring LogStream with your new knowledge!

```
Select Fields (11 of 11)
IN OUT
1
            α_raw: Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[2929]: url_helper.py[DEBUG]:
                   [0/1] open 'http://145.133.162.42/latest/api/token' with {'url': 'http://1... Show more
2021-02-05
20:18:31.286 # _time: 1612577911.286
           α cribl breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           a host: ec2-35-162-133-145
           a host_ip: 35.162.133.145
           \alpha index: usa_index_tier
           a region: us-west1
           a sourcetype: cloud-init_us-west1_a
         a _raw: Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[2929]: __init__.py[DEBUG]: Lo
                  oking for data source in: ['Ec2', 'None'], via packages ['', u'cloudinit.s... Show more
2021-02-05
21:33:30.302 # _time: 1612582410.302
           α az: b
-06:00
           α cribl_breaker: Break on newlines
           α cribl_pipe: setting_index_by_region_availability_zone
           α domain: compute.amazonaws.com
           α host: ec2-48-169-111-182
           α host_ip: 48.169.111.182
           α index: usa_index_tier
           α region: us-east2
           α sourcetype: cloud-init_us-east2_b
           a _raw: Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-init[2929]: atomic_helper.p
                  y[DEBUG]: Atomically writing to file /var/lib/cloud/data/status.json (via ... Show more
2021-02-06
00:29:11.841 # _time: 1612592951.841
           a az: a
           α cribl_breaker: Break on newlines
           a cribl_pipe: setting_index_by_region_availability_zone
           a host: ec2-21-187-232-201
           a host_ip: 21.187.232.201
           α index: apac_index_tier
            α region: asia-northeast3
            α sourcetype: cloud-init_asia-northeast3_a
```

Try This at Home

Below you'll find the the lookup table, Pipeline, and sample events demonstrated in this use case. Create the lookup file first, and then import the Pipeline. (The order matters, because the Pipeline import depends on the lookup table's presence.)

Lookup Table

To create the lookup table in LogStream's UI, select **Knowledge > Lookups**, then click + **Add New** and select **Create with Text Editor**.

Copy and paste in the header and rows listed below, then save the result as:

```
region_index_sourcetype
us.+,usa_index_tier,cloud-init
asia.+,apac_index_tier,cloud-init
europe.+,emea_index_tier,cloud-init
northamerica.+,na_index_tier,cloud-init
```

southamerica.+,ltam index tier,cloud-init

region_index_sourcetype.csv.

Pipeline

Below is an export of the whole LogStream Pipeline presented here. Import this JSON to get a Pipeline named: setting_index_by_region_availability_zone.json.

```
Magic Pipeline
  "id": "setting_index_by_region_availability_zone",
  "conf": {
    "output": "default",
    "groups": {},
    "asyncFuncTimeout": 1000,
    "functions": [
        "filter": "true",
        "conf": {
          "comment": "This pipeline demonstrates four different ways to leverage regular express
        "id": "comment"
      },
        "filter": "true",
        "conf": {
          "source": "_raw",
          "iterations": 100,
          "overwrite": false,
          "regex": "/GMT:\\s+(?<host>[^.]+)\\.(?<region>\\w+\\d+)-(?<az>[^.]+)\\.(?<domain)</pre>
        },
        "id": "regex_extract",
        "disabled": false,
        "description": "Extract host, region, availability_zone, and domain"
      },
      {
```

```
"filter": "true",
  "conf": {
   "matchMode": "regex",
   "matchType": "specific",
    "reloadPeriodSec": 60,
    "addToEvent": false,
    "inFields": [
        "eventField": "region"
   ],
    "ignoreCase": false,
    "file": "region_index_sourcetype.csv"
  },
  "id": "lookup",
  "disabled": false,
  "description": "Lookup index and sourcetype using regex matching"
},
{
  "filter": "true",
  "conf": {
    "add": [
      {
        "name": "host_ip",
```

Sample Events

And here's a sample of raw events that you can upload or copy/paste into LogStream's Preview pane to test the Pipeline's Functions:

Sample events

```
Feb 06 2021 02:18:31.286 GMT: ec2-35-162-133-145.us-west1-a.compute.amazonaws.com: cloud-init[29] Feb 06 2021 03:33:30.302 GMT: ec2-48-169-111-182.us-east2-b.compute.amazonaws.com: cloud-init[29] Feb 06 2021 06:29:11.841 GMT: ec2-21-187-232-201.asia-northeast3-a.compute.amazonaws.com: cloud-Feb 06 2021 12:59:44.232 GMT: ec2-76-187-246-132.europe-west3-b.compute.amazonaws.com: cloud-init Feb 06 2021 17:04:16.921 GMT: ec2-67-205-202-104.northamerica-northeast1-c.compute.amazonaws.com Feb 06 2021 19:45:47.687 GMT: ec2-87-209-176-201.southamerica-east1-a.compute.amazonaws.com: Data
```

From here, modify the sample data, lookup table, and Functions to adapt this approach to your own needs!

Regex Filtering

Regex Filtering of Data in Motion

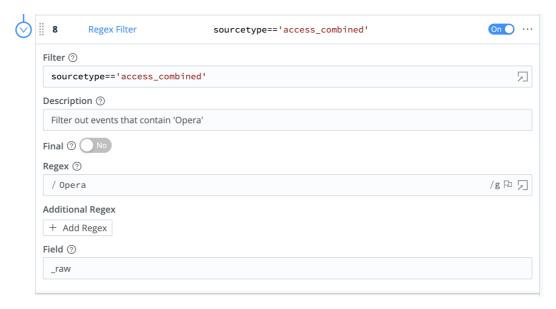
To filter events in real time, we use the out-of-the-box **Regex Filter** Function. This is similar to nullqueueing with TRANSFORMS in Splunk, but the matching condition is way more flexible.

Regex Filtering Example

Let's see how we can filter out any sourcetype—'access_combined' events whose _raw field contains the pattern Opera:

First, make sure you have a Route and Pipeline configured to match desired events.

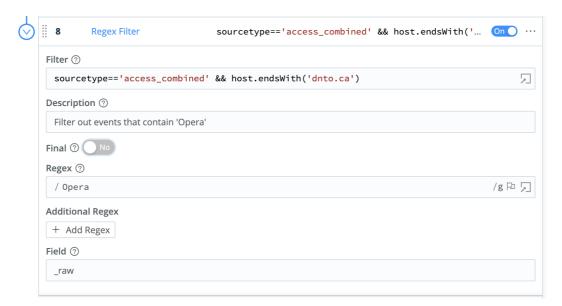
Next, let's add a Regex Filter Function to it:



Defining the Regex Filter Function

Next, verify that this search does **not** return any results: sourcetype="access_combined" Opera

You can add more conditions to the Filter input field. For example, to further limit the filtering to only events from hosts with domain <code>dnto.ca</code>, change the filter input as shown below:



Filtering by host

This is a very flexible method for filtering incoming events in real time, on virtually any arbitrary conditions.

Encrypting Sensitive Data

Encryption at Ingest-Time and Decryption in Splunk

With Cribl LogStream, you can encrypt your sensitive data in real time before it's forwarded to and stored at a destination. Using the out-of-the-box Mask function, you can define patterns to encrypt with specific key IDs or key classes. To decrypt in Splunk, you will need to install Cribl App for Splunk on your search head. (The app will default to mode-searchhead.)

Keys and Key Classes

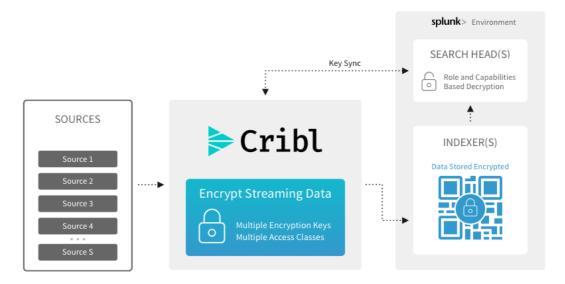
Symmetric encryption keys can be configured through the CLI or the UI.

They're used to encrypt the patterns, and users are free to define as many keys as required.

Key classes are collections of keys that can be used to implement multiple levels of access control. Users (or groups of users) that have access to data with encrypted patterns can be associated with key classes. You can use these classes to provide more-granular access rights, such as read versus decryption permissions on a dataset.

Encrypting in Cribl LogStream and Decrypting in Splunk

- 1. Define one or more keys and key classes on Cribl LogStream.
- 2. Sync auth with the decryption side (Splunk Search Head)
- 3. Apply the Mask function to patterns of interest, using C.Crypto.encrypt().
- 4. Decrypt on the Splunk search head, using Role-Based Access Control on the decrypt command.

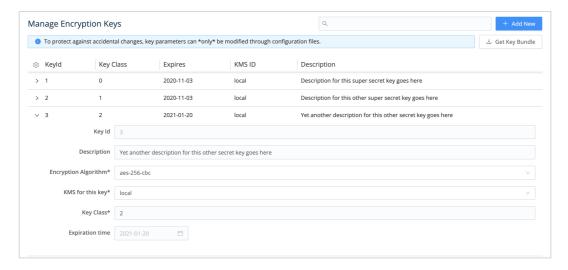


Encrypting in LogStream, decrypting in Splunk

Example

Encryption Side

• Generate keys via the UI, in **Settings** > **Encryption Keys**:



Adding a new encryption key

• Or generate one or more keys via the CLI, as follows:

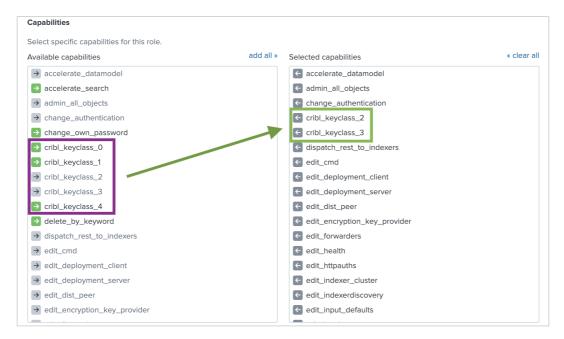
```
$CRIBL_HOME/bin/cribl keys add -c 1 -i
...
$CRIBL_HOME/bin/cribl keys add -c <N> -i
```

Add -e <epoch> if you'd like to set expiration for your keys.

• Sync auth/(cribl.secret|keys.json). To decrypt data, the decrypt command will need access to these keys. The cribl.secret and keys.json files in \$CRIBL_HOME/local/cribl/auth should be synced/copied over to the search head (decryption side).

Decryption Side

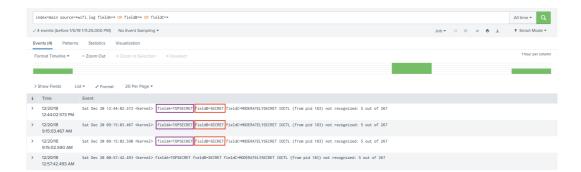
- Install Cribl App for Splunk on your search head. It will default to modesearchhead.
- Assign permissions to the decrypt command, per your requirements.
- Assign capabilities to your Roles, per your requirements. Capability names should follow the format cribl_keyclass_N, where N is the Cribl Key Class. For example, a role with capability cribl_keyclass_1 has access to all key IDs associated with key class 1. You can use more capabilities, as long as they follow this naming convention.



Selecting capabilities

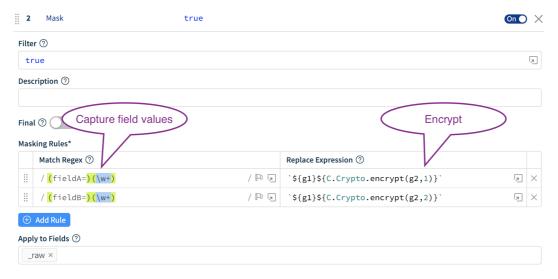
Usage

Before Encryption: Sample un-encrypted events. Notice the values of fieldA and fieldB.



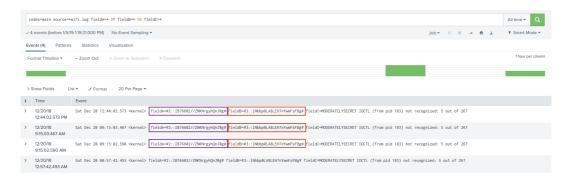
Events before encryption

Next, encrypt fieldA values with key class 1, and fieldB with key class 2.



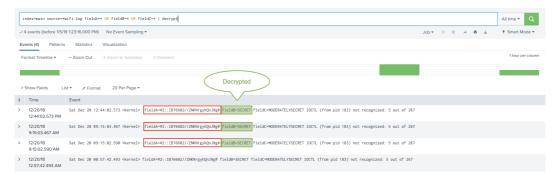
Encrypting two fields with separate key classes

After Encryption: again, notice the values of fieldA and fieldB.



Both fields encrypted

Here, we've decrypted fieldB but not fieldA. This is because the logged-in user has been assigned the capability cribl_keyclass_2, but not cribl_keyclass_1.



One field decrypted

Syslog Data Reduction

When ingesting data from syslog senders, Cribl LogStream can readily trim data volume by 30% or more, optimizing infrastructure for downstream services like Splunk or Elasticsearch. Here, we outline some best practices for replacing your syslog server with LogStream.

Syslog Event Parsing

By default, a LogStream Syslog Source will produce the following fields:
_time, appname, facility and facilityName, host, message, and
severity and severityName.

Parsed syslog event

This output is much more readable, but we haven't saved any volume – we now have redundant pairs of fields (numeric versus text) representing the facility and severity.

Below, we'll outline how to streamline syslog events to something more like this:

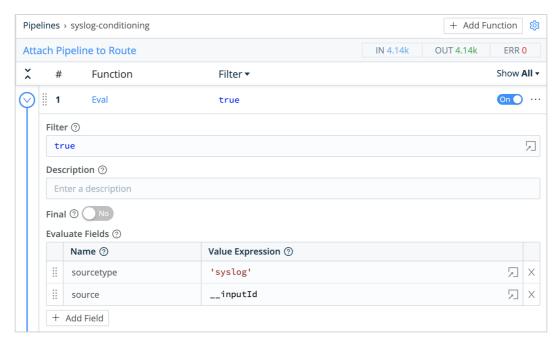
Parsed and redacted syslog event

This extracts the essentials, removes the redundancies, adds one new field that identifies the LogStream Pipeline we're about to build, and shrinks the outbound _raw payload to just its message component. For still further efficiencies, we'll look at how to drop or downsample frequent events, and how to balance high-volume syslog inputs across LogStream worker processes.

Create Pre-Processing Pipeline

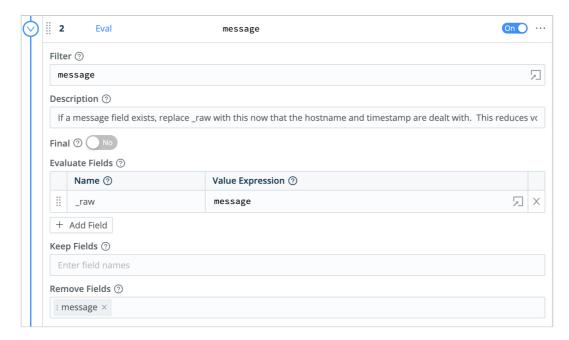
Even before setting up a syslog Source, our first step is to create a preprocessing Pipeline that will be available to normalize incoming events on all syslog Sources, reducing data volume as shown above.

The Pipeline starts with an **Eval** Function, whose **Evaluate Fields** section tags syslog events with two new fields indicating their origin: sourcetype: 'syslog' and source: __inputId . Because this Pipeline is designed only to condition all incoming syslog data, we leave **Filter** set to its default true value, to process all events.



Eval Function to tag syslog events' origin

A second **Eval** Function filters for the presence of a message field. If found, it overwrites the _raw field with message , and then deletes message as redundant. This function alone typically reduces syslog data volume by 15–20%.



Eval Function to rewrite message as _raw

⚠ Before using this Pipeline in production, preview sample data to verify that you're not dropping any essential information.

This third **Eval** Function deletes two redundant fields. Its **Filter** condition makes sure both of these string fields exist and contain values: severity \neq null & facility \neq null . If so, it removes their corresponding numeric fields, severity and facility.



Eval Function to remove redundant numeric fields

To further shrink the output, this fourth **Eval** Function removes procid fields that contain only a dash – meaning "no value extracted." We set **Filter**: procid='-' and **Remove Fields**: procid.



Eval Function to remove empty procid fields

With these four Functions enabled, the **Preview** pane's **Basic Statistics** confirm that we've reduced the _raw field's length by more than 30%.



Data reduction example

Dropping Noisy Data

With some syslog senders, like VMware ESX/ESXi, 80–90% of incoming events can be of debug severity. To further reduce volume, you could use this final **Drop** Function to drop all these events. Its **Filter** is simply severityName='debug'.



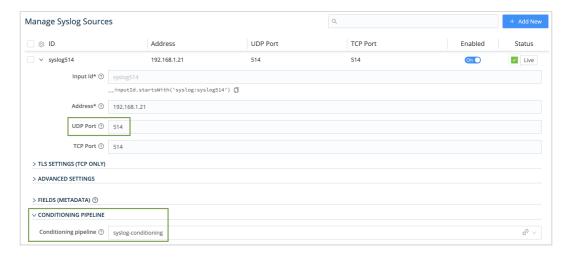
Drop Function to remove debug events

Enabling this Function could up our volume savings to about 40%. Depending on your needs, you might prefer to:

- Use a Function like this in your Route's processing Pipeline, rather than in this upfront Pipeline.
- Also drop info events.
- Instead use a Sampling Function to sample debug events (at a ratio like 1:10), or a Dynamic Sampling function.
- Instead use a Suppress Function to clamp down the frequency of debug events.

Create Syslog Source

Once we've built and saved the pre-processing Pipeline, our next step is to add a Syslog Source.



Syslog Source configured for UDP and pre-processing Pipeline

Specify the UDP Port where you want this Source to listen for syslog data.

Then attach the pre-processing Pipeline that you created above, and save the Source.

i Cribl generally recommends selecting UDP, rather than TCP, for high-volume syslog senders. This facilitates efficient load balancing by not continuously tying such senders to any one LogStream Worker Process. See Sizing and Scaling for more details.

Fields/Metadata

In the pre-processing Pipeline, we tagged *all* incoming syslog events with new sourcetype and source fields to indicate their origin. Alternatively, you could use the Source's **Fields/Metadata** section to define similar key-value pairs, specific to each of your Syslog Sources.

Create Route(s)

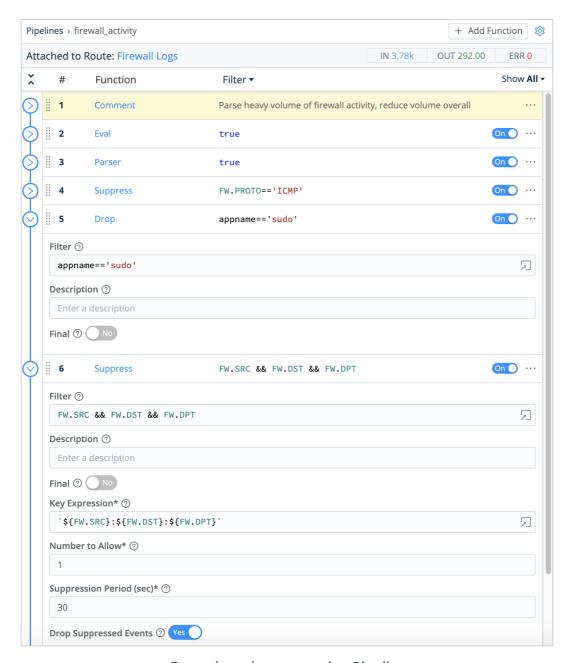
Create Routes, as needed, for each of your Syslog Sources. Give each Route a corresponding **Filter** expression, and set its **Output** to the Destination where you want to send its processed syslog data.



Example syslog Route

Processing Pipelines, and Next Steps

For any or all syslog Routes, you can define and attach a processing Pipeline. These can apply more-granular Filters and Functions to further reduce volume, using techniques like Sampling, Dynamic Sampling, or (as shown below) Drop and Suppression. Your most-verbose Syslog Sources are ideal targets for such further processing.



Example syslog processing Pipeline

Splunk to Elasticsearch

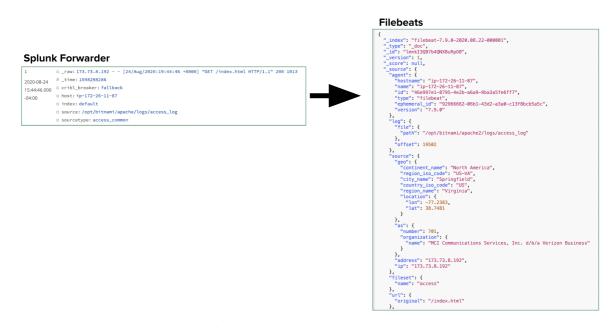
To route data from existing Splunk infrastructure to Elasticsearch services, you might face a daunting task: re-architecting your entire forwarding tier. This could require retooling lots of servers – up to hundreds, or thousands – to uninstall their Splunk forwarders, and swap in Elastic-compatible agents.

Cribl LogStream can reduce this effort to just a few hours: Configure one Splunk outputs.conf stanza to output to LogStream, and propagate that across all your Splunk servers. Done!

Next, you can easily configure LogStream to listen for Splunk data on one port, and to route that data to all the Elasticsearch destinations you want to feed.

Transforming Data from Splunk to Elastic Format

Also, in LogStream's core, you can easily design a Pipeline that modifies the original Splunk event into Elastic's Common Schema – making it look exactly like an event generated by an Elastic agent. These transformations help you make the most of Elastic's offerings, like Filebeats, etc.

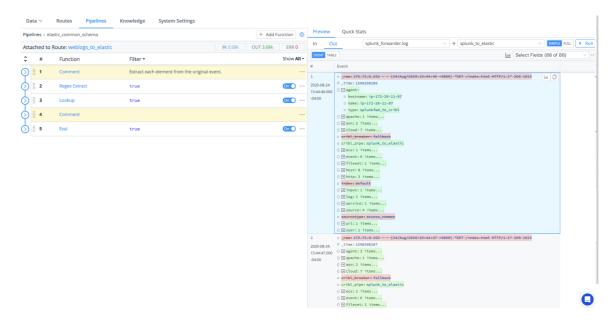


Transforming to Elastic Common Schema

Some of the LogStream Functions useful in transforming Splunk-generated events into Elastic's format are:

- Regex Extract: Extract a portion of the raw event, and place it into a specified field.
- Lookup: key off the host IP to add fields like hostname, name, id, and type.
- Eval: Turn key-value pairs into nested JSON objects.
- GeoIP: Correlate source IP to a geographic database.

We'll show all four in our example Pipeline below, although you might need only a subset.



LogStream Pipeline and Data Preview

Goat Rid of Some Guesswork

LogStream will offer you further time savings as you configure the internal data transformation. LogStream's Data Preview features enable you to test transformations' results as you build your Pipeline, before you commit or run it.

This eliminates blind guesswork in Splunk configuration files to specify source -> index transformations, check the results, and then start all over again. In particular, LogStream's Regex Extract Function provides a regex101-like UI, to facilitate precisely designing and debugging your regex expressions.

Let's goat started on the example.

Configure Splunk Forwarder

First, in a Splunk App, configure a Splunk forwarder (UF or HF) to specify your Cribl Workers as destinations. Use outputs.conf stanzas of this form:

```
.../outputs.conf

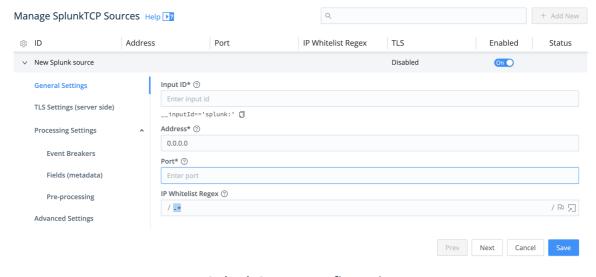
[tcpout]
disabled = false
defaultGroup = cribl, <optional_clone_target_group>,

[tcpout:cribl]
server = [<cribl_ip>|<cribl_host>]:<port>, [<cribl_ip>|<cribl_host>]:<port>, ...
sendCookedData=true
```

Push the app using the deployment server.

Configure Splunk Source in LogStream

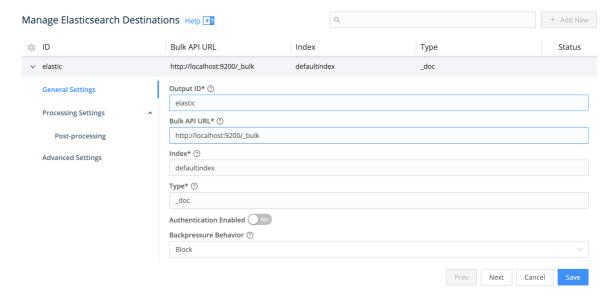
Next, in LogStream, configure a Splunk Source. The key requirement here is to set the **Port** to listen on. (Optionally, you can also configure TLS, Event Breakers, metadata fields, and/or a pre-processing Pipeline.)



Splunk Source configuration

Configure Elasticsearch Destination

To configure LogStream's output, set up an Elasticsearch Destination by specifying the **Bulk API URL** and **Index**.



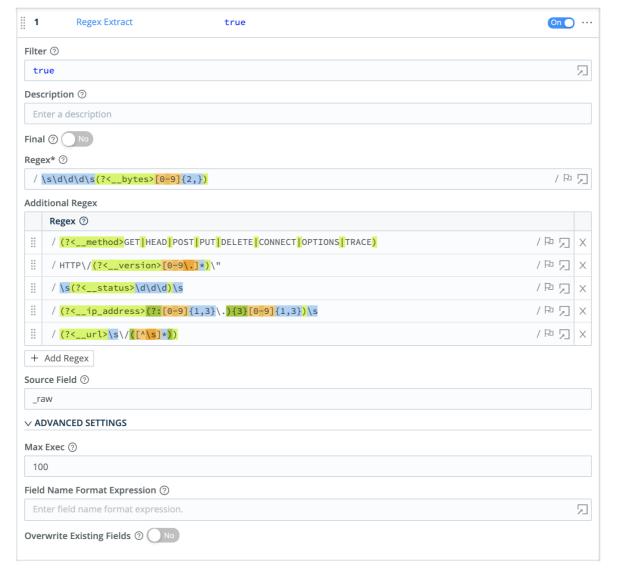
Elasticsearch Destination configuration

Configure Pipeline

Next, this section shows several Functions that you can assemble into a Pipeline to transform incoming Splunk events to match the Elastic Common Schema.

Regex Extract Function

First, use a Regex Extract Function to break the Splunk events into fields. Try the sample configuration shown below:



Regex Extract Function

Here are the six rows of regex in this example:

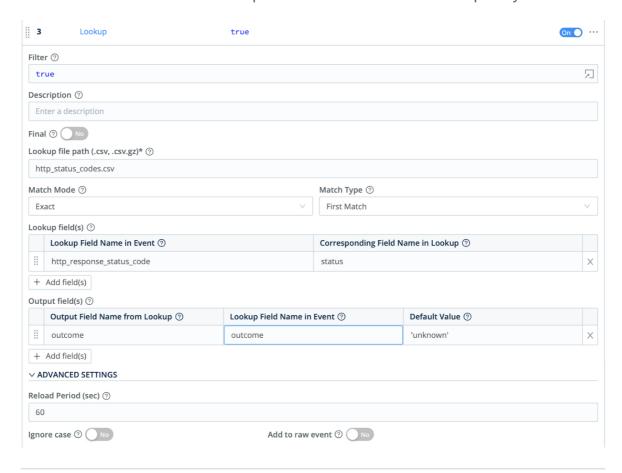
```
Regex; Additional Regex

/\s\d\d\s(?<_bytes>[0-9]{2,})/
/(?<_method>GET|HEAD|POST|PUT|DELETE|CONNECT|OPTIONS|TRACE)/
/HTTP\/(?<_version>[0-9\.]*)\"/
/\s(?<_status>\d\d\d)\s/
/(?<_ip_address>(?:[0-9]{1,3}\.){3}[0-9]{1,3})\s/
/(?<_url>\s\/([^\s]*))/
```

As you refine your expression, capture a sample of incoming Splunk data to test your regex's results in LogStream's right Preview pane.

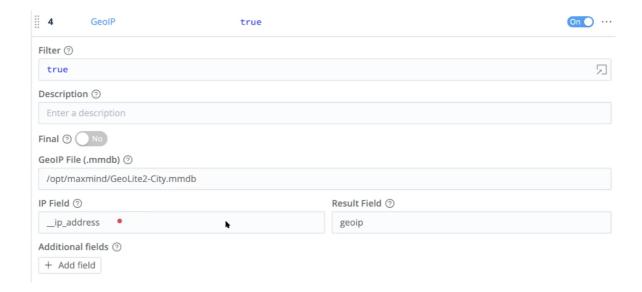
Lookup Function

In this example, we next add a Lookup Function, to translate HTTP error codes to readable text. Note this Function's optional **Reload Period** field, in which you can define a reload interval for a lookup file whose contents refresh frequently.



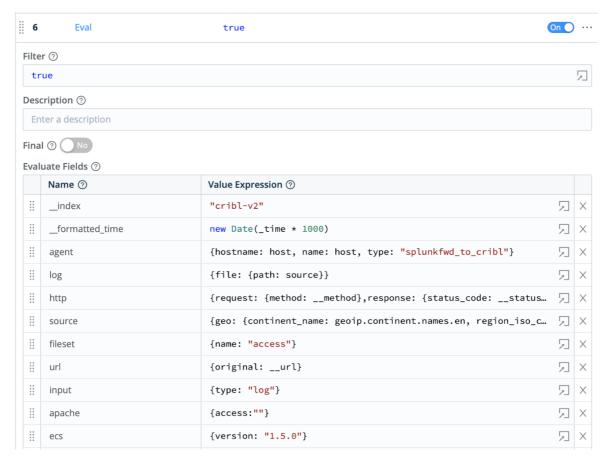
GeoIP Function

To enrich the Splunk data, we next use a GeoIP Function. This a specialized lookup against a database of IP addresses by geographic location. This Function's output can provide Elasticsearch with location fields like lat and long.



Eval Function

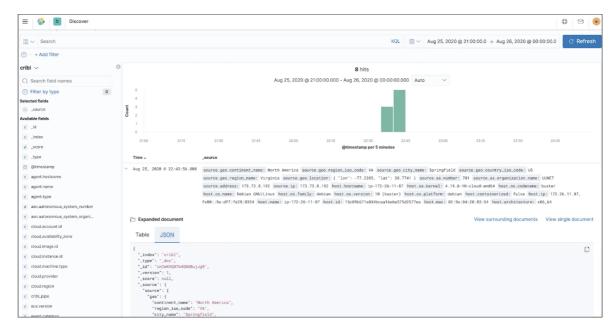
Finally, to further enrich the outbound events, the Pipeline uses an Eval Function. This adds multiple key-value pairs that define and populate fields conforming to the Elastic Common Schema.



Eval Function

Results

After attaching your Pipeline to a Route, here's an an exported event, all happy in Elasticsearch with nested JSON.



Event as exported to Elasticsearch

For More Info

For additional details on configuring Splunk forwarders for LogStream, see this related documentation:

- Configuring a Splunk (TCP) Forwarder
- Configuring Cribl App for Splunk on an HF

Reducing Windows XML Events

Here, we demonstrate how to use just a few LogStream Functions to parse WindowsXML events and reduce their volume by 34–70%, dramatically reducing your downstream infrastructure requirements.

Using Eval and C.Text.parseWinEvent

LogStream's internal C.Text.parseWinEvent method parses Windows XML strings and returns a prettified JSON object. You can use this function within an Eval Function to parse an event or an ad hoc XML string. It works like C.Text.parseXml, but with Windows events, it produces more-compact output.

As you can see from its signature, C.Text.parseWinEvent accepts an optional nonValues parameter that can further reduce an event's size by discarding characters that you specify as redundant.

- i (method) C.Text.parseWinEvent(xml: string, nonValues?: string[]): any
 - @param xml an XML string; or an event field containing the XML.
 - @param nonValues array of string values. Elements whose value equals any of the values in this array will be omitted from the returned object. Defaults to ['-'], meaning that elements whose value equals will be discarded.
 - @returns an object representing the parsed Windows Event; or undefined if the input could not be parsed.

XML: Threat or Menace?

When working with XML, an anonymous Reddit user's quote sums up the challenge: "Some languages can be read by human, but not by machines, while others can be read by machines but not by humans. XML solves this problem by being readable to neither." An example of a Windows XML event only reinforces this quote:



This Windows XML raw event is 1.36KB in size:

	_raw Length ②	Full Event Length ③	Number of Fields ③	Number of Events ②
IN	1.36KB	1.43KB	3	1
OUT	1.36KB	1.47KB	4	1
DIFF	0.00%	↑ 2.73%	↑ 33.33%	0.00%

Eval to the Rescue

In our initial Eval Function below, the **Value Expression** uses C.Text.parseWinEvent to simply parse the _raw Windows XML event and turn it into a prettified JSON object:



The resulting JSON event is now down to 921.00B in size, a 34.07% reduction of the event:



Removing Unnecessary Fields with a Better Eval

But we can do better. The fields containing essentially null values ('0', $'0\times0'$, or '-') bloat events, demanding extra infrastructure and storage:

```
(} = _raw:
  {} = Event:
     ⊕ EventData:
       {} □ Data:
          α AuthenticationPackageName: NTLM
           α FailureReason: 9%2313
           a IpAddress: 64.203.230.138
           a IpPort: 0 ←
           a KeyLength: 0 ◀
           α LmPackageName: - ◀
           a LogonProcessName: NtLmSsp
           a LogonType: 3
           a ProcessId: 0x0 ←
           α ProcessName: - ◀
           a Status: 0xc000006d
           α SubjectDomainName: - ←
           a SubjectLogonId: 0x0 ←
           a SubjectUserSid: S-1-0-0
           a SubStatus: 0xc000006a
           a TargetUserName: ADMINISTRATOR
                                                                    Field Values: ['0','0x0','-']
           a TargetUserSid: S-1-0-0

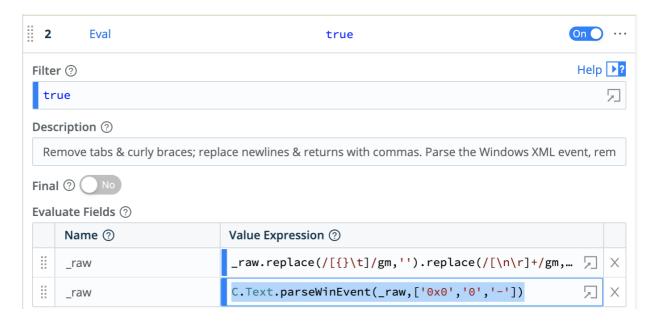
    □ TransmittedServices: - 
    □

           α WorkstationName: - ◀
     {} □ System:
        a Channel: Security
        a Computer: EC2AMAZ-CPMK6J5.cribl.poc
        a EventID: 4625
        a EventRecordID: 1598720
        {} = Execution:
          a ProcessID: 740
           α ThreadID: 948
        a Keywords: 0x8010000000000000
        α Level: Θ ◀
        α Opcode: 0 ◀
        {} ☐ Provider:
          a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}

    Name: Microsoft-Windows-Security-Auditing

        α Task: 12544
        {} ☐ TimeCreated:
           a SystemTime: 2020-11-12T17:27:33.608596000Z
```

Let's amplify the reduction by removing all of the fields whose values are in the set: ['0','0×0','-']. This improved version of our Eval Function parses the Windows XML event, and discards ['0','0×0','-'] values. (Its preceding row also tidies up events by removing tabs and curly braces, and replacing newlines and returns with commas.) The result is an even **smaller** prettified JSON object:



If you compare this Preview-pane screenshot to the Preview screenshot above, you can confirm that the fields with values matching $['0', '0\times0', '-']$ are removed:

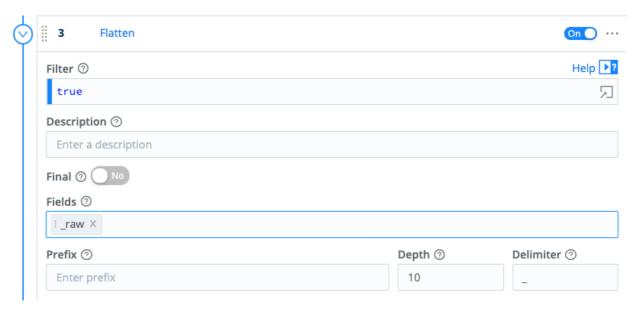


The event is now down to 678.00B in size, translating to a 51.47% reduction from the original event:

	_raw Length ⑦	Full Event Length ⑦	Number of Fields ⑦	Number of Events ⑦
IN	1.36KB	1.43KB	3	1
OUT	678.00B	785.00B	4	1
DIFF	↓ -51.47%	↓ -46.42%	↑ 33.33%	0.00%

Flatten Function

LogStream's Flatten Function is designed to flatten fields out of a nested structure. Let's flatten the JSON object within <code>_raw</code> , to see if we can further reduce the event's size before we send it to our preferred destinations:



Using Flatten, we've successfully created top-level fields from the nested JSON structure:

```
a _raw: <Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider
       y-Auditing' Guid='{54849625-5478-4994-A5BA-3E3B0328C30D}'/><EventID>4625</Event... Show more
α _raw_Event_EventData_Data_AuthenticationPackageName: NTLM
a _raw_Event_EventData_Data_FailureReason: %%2313
a _raw_Event_EventData_Data_IpAddress: 64.203.230.138
a _raw_Event_EventData_Data_LogonProcessName: NtLmSsp
a _raw_Event_EventData_Data_LogonType: 3
a _raw_Event_EventData_Data_Status: 0xc000006d
a _raw_Event_EventData_Data_SubjectUserSid: S-1-0-0
a _raw_Event_EventData_Data_SubStatus: 0xc0000006a
a _raw_Event_EventData_Data_TargetUserName: ADMINISTRATOR
a _raw_Event_EventData_Data_TargetUserSid: S-1-0-0
a _raw_Event_System_Channel: Security
a _raw_Event_System_Computer: EC2AMAZ-CPMK6J5.cribl.poc
a _raw_Event_System_EventID: 4625
α _raw_Event_System_EventRecordID: 1598720
a _raw_Event_System_Execution_ProcessID: 740
a _raw_Event_System_Execution_ThreadID: 948
α _raw_Event_System_Keywords: 0x8010000000000000
a _raw_Event_System_Provider_Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
{\tt \alpha\_raw\_Event\_System\_Provider\_Name: Microsoft-Windows-Security-Auditing}
α _raw_Event_System_Task: 12544
a _raw_Event_System_TimeCreated_SystemTime: 2020-11-12T17:27:33.608596000Z
# time: 1609852209.078
α cribl breaker: Break on newlines
α cribl_pipe: Parse_Windows_XML_Events
```

Don't worry about the _raw field's deletion (red strikeout). This is the Flatten Function's default behavior. We'll restore _raw after we clean and reduce the event even more.

The flattened field names are extracted from _raw and delimited with _ . These field names are quite long. We can optimize them using the Rename Function.

Rename Function

Rename is designed to change fields' names, or to reformat their names (e.g., to normalize names to camelcase). You can use Rename to change specified fields (much like Eval), or to accomplish bulk renaming based on a JavaScript expression (much like the Parser Function). But Rename offers a streamlined way to alter only field names, without other effects.

Let's use Rename to remove any unnecessary prefixes from the field names, to further shrink our events. In the **Renaming Expression**, we build a a JavaScript expression to match the field names' prefixes (up to the underscore):



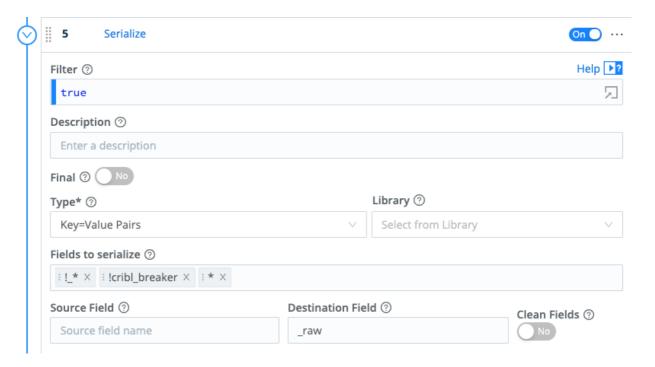
The resulting field names are now much more compact, and easier to work with and manage:

```
a _raw: <Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider
       y-Auditing' Guid='{54849625-5478-4994-A5BA-3E3B0328C30D}'/><EventID>4625</Event... Show more
# _time: 1609852209.078
α AuthenticationPackageName: NTLM
α Channel: Security
α Computer: EC2AMAZ-CPMK6J5.cribl.poc
a cribl_breaker: Break on newlines
α cribl_pipe: Parse_Windows_XML_Events
a EventID: 4625
a EventRecordID: 1598720
α FailureReason: %%2313
a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
a IpAddress: 64.203.230.138
a Keywords: 0x8010000000000000
α LogonProcessName: NtLmSsp
a LogonType: 3
α Name: Microsoft-Windows-Security-Auditing
a ProcessID: 740
a Status: 0xc000006d
a SubjectUserSid: S-1-0-0
a SubStatus: 0xc000006a
a SystemTime: 2020-11-12T17:27:33.608596000Z
α TargetUserName: ADMINISTRATOR
α TargetUserSid: S-1-0-0
α Task: 12544
α ThreadID: 948
```

Serialize Function

We started with bloated Windows XML data, and we've parsed and prettified it into JSON. Next, we'll extract key-value pairs. We'll use the Serialize Function, which serializes an event's content into a predefined format.

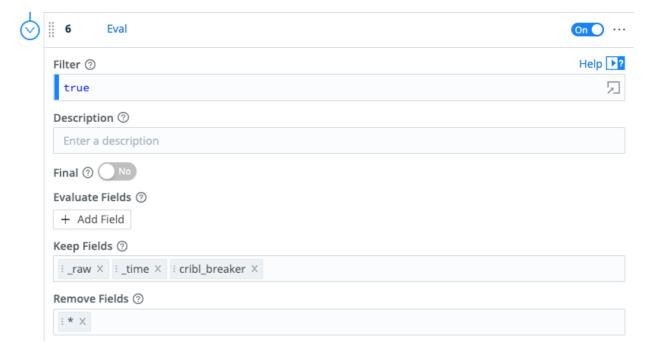
We set Serialize to change the **Type** to key-value pairs. (The Function's other supported target Types include JSON Object and CSV.) Here, Serialize takes the extracted fields and puts them back into _raw:



In the Preview pane, the _raw field is now back, serialized into compact, tidy key-value pairs:

```
α _raw: Name=Microsoft-Windows-Security-Auditing Guid={54849625-5478-4994-A58A-3E3B0328C30D} EventID=4625 Task=12544 Keywords=
        0x801000000000000 SystemTime=2020-11-12T17:27:33.608596000Z EventRecordID=1598720 ProcessID=740 ThreadID=948 Channel=
        Security Computer=EC2AMAZ-CPMK6J5.cribl.poc SubjectUserSid=S-1-0-0 TargetUserSid=S-1-0-0 TargetUserName=ADMINISTRATOR
        Status=0xc000006d FailureReason=%%2313 SubStatus=0xc000006a LogonType=3 LogonProcessName="NtLmSsp " AuthenticationPack
        ageName=NTLM IpAddress=64.203.230.138 Show less
# _time: 1609852209.078
\alpha AuthenticationPackageName: NTLM
α Channel: Security
α Computer: EC2AMAZ-CPMK6J5.cribl.poc
a cribl breaker: Break on newlines
α cribl_pipe: Parse_Windows_XML_Events
a EventID: 4625
a EventRecordID: 1598720
α FailureReason: %%2313
a Guid: {54849625-5478-4994-A5BA-3E3B0328C30D}
a IpAddress: 64.203.230.138
a Keywords: 0x8010000000000000
α LogonProcessName: NtLmSsp
a LogonType: 3
\alpha Name: Microsoft-Windows-Security-Auditing
α ProcessID: 740
α Status: 0xc000006d
a SubjectUserSid: S-1-0-0
a SubStatus: 0xc000006a
a SystemTime: 2020-11-12T17:27:33.608596000Z
α TargetUserName: ADMINISTRATOR
α TargetUserSid: S-1-0-0
α Task: 12544
α ThreadID: 948
```

The last step is to remove any of the extracted fields you don't need before sending events to your destinations. We'll again call on the Eval Function, which adds or removes fields in events. (For a Splunk destination, these are index-time fields.) This final Eval Function looks like this:



To sum up, we've successfully transformed the original Windows XML event into key-value pairs:

```
a_raw: Name=Microsoft-Windows-Security-Auditing Guid={54849625-5478-4994-A5BA-3E3B0328C30D} EventID=4625 Task=12544 Keywords=
0x80100000000000000 SystemTime=2020-11-12T17:27:33.608596000Z EventRecordID=1598720 ProcessID=740 ThreadID=948 Channel=
Security Computer=EC2AMAZ-CPMK6J5.cribl.poc SubjectUserSid=S-1-0-0 TargetUserSid=S-1-0-0 TargetUserName=ADMINISTRATOR
Status=0xc000006d FailureReason=%%2313 SubStatus=0xc000006a LogonType=3 LogonProcessName="NtLmSsp" AuthenticationPack
ageName=NTLM IpAddress=64.203.230.138 Show less

# _time: 1609852209.078
cribl_breaker: Break on newlines
cribl_pipe: Parse_Windows_XML_Events
```

And we've dramatically reduced the event's size, while retaining all of the necessary fields. The event is now down to 513.00B in size, translating to a 63.28% reduction from the original Windows XML:

	_raw Length ②	Full Event Length ③	Number of Fields ③	Number of Events ②
IN	1.36KB	1.43KB	3	1
OUT	513.00B	621.00B	4	1
DIFF	↓ -63.28%	↓ -57.61%	↑ 33.33%	0.00%

Try This at Home

Below is an export of the whole LogStream Pipeline presented here. Import this JSON to experiment with it and modify it to match your own needs:

```
Win XML Pipeline
{
    "id": "Windows_Security_Events",
    "conf": {
        "output": "default",
        "groups": {},
        "asyncFuncTimeout": 1000,
```

```
"functions": [
 {
    "filter": "true",
    "conf": {
     "comment": "This LogStream Pipeline reduces Microsoft Windows XML events, retains full
    "id": "comment"
 },
    "filter": "true",
    "conf": {
      "add": [
       {
          "name": "_raw",
          "value": "_raw.replace(/[{}\\t]/gm,'').replace(/[\\n\\r]+/gm,',')"
       },
        {
          "name": "_raw",
          "value": "C.Text.parseWinEvent(_raw,['0×0','0','-'])"
       }
     ]
    },
    "id": "eval",
    "disabled": false,
    "description": "Remove tabs & curly braces; replace newlines & returns with commas. Pars
 },
 {
   "filter": "true",
    "conf": {
     "fields": [
       "_raw"
     ],
     "prefix": "",
     "depth": 5,
     "delimiter": " "
   },
    "id": "flatten",
    "disabled": false,
    "description": "Flatten the object into key value fields"
 },
 {
   "filter": "true",
    "conf": {
      "baseFields": [],
      "renameExpr": "name.replace(/_raw_Event_\\w+_/,'')",
```

Finally, here's a sample of Windows XML events that you can upload to LogStream's Preview pane to try this out:

```
Sample data
```

<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic
<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic
<Event xmlns='http://schemas.microsoft.com/win/2004/08/events/event'><System><Provider Name='Mic</pre>

Best Practices

Managing Large Lookups

This page offers a general approach to managing large lookup files. While LogStream's Git integration normally helps manage configuration changes, large lookups are exceptions. In many cases, you might want to exclude these files from Git, to reduce excessive deploy traffic. This approach can also prevent Git Push commands from encountering large file errors.

Good scenarios for this approach are:

- Large binary files like databases which don't benefit from Git's typical efficient storage of only the deltas between versions. (With binary files, Git must replace the whole file for each new version.)
- Files updated frequently and/or files updated independent of LogStream.
- Files replicated on many Worker Nodes.

About the MaxMind GeoLite Example

We'll illustrate this with an example that often combines all three conditions: setting up the free, popular MaxMind GeoLite2 City database to support LogStream's GeoIP lookup Function. This example anticipates a LogStream production distributed deployment, where the GeoLite database is updated nightly across multiple Workers.

This example includes complete instructions for this particular setup. However, you can generalize the example to other MaxMind databases, and to other large lookup files – including large .csv 's that similarly receive frequent updates.

Reducing Deploy Traffic

The general approach for handling large lookups is:

- Do not place these files in the standard \$CRIBL_HOME/data/lookups.
- Instead, place them in a \$CRIBL_HOME subdirectory that's excluded from
 Git version control, through inclusion in the \$CRIBL_HOME/.gitignore
 file. Deploying the files to the Master Node and all desired Workers will
 require a manual procedure and will be required for the initial deployment
 as well as subsequent updates.

The example below uses \$CRIBL_HOME/state subdirectory, which is already listed in the default .gitignore file that ships with LogStream.

i If you prefer, you can use a different path, including a path outside \$CRIBL_HOME . If you choose this alternative, be sure to add that path to .gitignore .

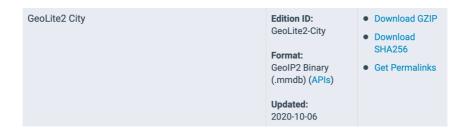
However, Cribl recommends using a \$CRIBL_HOME subdirectory like \$CRIBL_HOME/state, because this inherits appropriate permissions and simplifies backup/restore operations.

Let's move on to the MaxMind GeoLite specifics.

Download and Extract the Database

To enable the GeoIP Function using the MaxMind GeoLite 2 City database, your first steps are:

- 1. Create a free MaxMind account, at the page linked above.
- 2. Log in to your MaxMind account portal and select **Download Databases**.
- 3. On the Download page, look for the database you want. (In this example, you'd locate the **GeoLite2 City** section.) Note the **Format: GeoIP2 Binary**, and select **Download GZIP**.



- 4. Extract the archive to your local system.
- 5. Change to the directory created when you extracted the archive. This directory's name will correspond to the date you downloaded the file, so in the above 2020-10-06 example, you would use: \$ cd GeoLite2-City_20201006

Copy the Database File to the Master and Worker Nodes (Recommended)

In distributed deployments, Cribl recommends copying the MaxMind database separately to the Master and all Worker Nodes, e.g.. placing it in the \$CRIBL_HOME/state path. This will minimize the Git commit/deploy overhead around nightly updates to the binary database file.

Once in the database's directory, execute commands of this form:

```
$ scp *.mmdb <user>@<master-node>:
$ scp *.mmdb <user>@<worker-node>:
```

△ Copy the file to each Worker in the Worker Group where you intend to use LogStream's GeoIP Function.

The above commands copy the .mmdb database file into your user's home directory on each Node. Next, we'll move it to \$CRIBL_HOME/state on each Node. Execute these commands on both the Master and Worker Nodes:

```
$ sudo mv ~/*.mmdb <$CRIBL_HOME>/state/
$ sudo chown -R cribl:cribl <$CRIBL_HOME>/state/
```

Now that the database is in place, your Pipelines can use the GeoIP Function to enrich data. In the Function's **GeoIP file (.mmdb)** field, insert the complete \$CRIBL_HOME/state/<filename>.mmdb file path.

Copy the Database File Only to the Master (Alternative)

In smaller deployments, you might choose to copy this MaxMind database only to Master Node, and to let Workers receive updates via Git commit/deploy. In this case, the final commands above might look like this:

Shell

```
$ sudo cp ~/*.mmdb /opt/cribl/groups/<group-name>/data/lookups/
$ cd /opt/cribl/groups/<group-name>/data/lookups/
$ sudo chown cribl:cribl *.mmdb
```

Automatic Updates to the MaxMind Database

To set up automatic updates, see MaxMind's Automatic Updates for GeoIP2 and GeoIP Legacy Databases documentation. You'll need two modifications specific to LogStream:

- This must be set up on the Master, and on each Worker in any Group using GeoIP lookups.

Memory Considerations

Storage aside, large lookup files can also require additional RAM on each Worker Node that processes the lookups. For details, see Memory Sizing for Large Lookups.

Troubleshooting

Diagnosing Issues

To help diagnose LogStream problems, you can share a diagnostic bundle with Cribl Support. The bundle contains a snapshot of configuration files and logs at the time the bundle was created, and gives troubleshooters insights into how LogStream was configured and operating at that time.

What's in the Diagnostic Bundle

The following directories (and their contents) off of \$CRIBL_HOME are included:

- /default/*
- /local/*
- /log/*
- /groups/*
- /state/jobs/* includes only the latest 10 task from the latest 10 jobs.

Creating and Exporting a Diagnostic Bundle

Users can create and **securely** share bundles with Cribl Support either from the UI or from the CLI. In either case, you'll need outbound internet access to https://diag-upload.cribl.io and a valid Support Case number.

Using the UI

To create a bundle, go to **Settings > Diagnostics > Diagnostic Bundle** and click **Create diagnostic bundle**.

- To download the bundle locally to your machine, click Export.
- To share the bundle with Cribl Support, toggle **Send to Cribl Support** to **Yes**, enter your case number, and then click **Export**.

You can create a bundle from individual workers if you have the Worker UI access setting enabled. Go to Workers > <worker-name> > System Settings > Diagnostics > Diagnostic Bundle, and click Create Diagnostic Bundle.

Previously created bundles are stored in $\c SCRIBL_HOME/diag$. They're also listed in the UI, where you can re-download them or share them with Cribl Support.

Using the CLI

To create a bundle using the CLI, use the diag command.

```
diag command CLI
# $CRIBL_HOME/bin/cribl diag
Usage: [sub-command] [options] [args]
Commands:
get - List existing Cribl LogStream diagnostic bundles
create - Creates diagnostic bundle for Cribl LogStream
send - Send LogStream diagnostic bundle to Cribl Support, args:
  -c <caseNumber> - Cribl Case Number
  [-p <path>]
              - Diagnostic bundle path (if empty, then new bundle will be created)
## Creating a diagnostic bundle
# $CRIBL_HOME/bin/cribl diag create
Created Cribl LogStream diagnostic bundle at /opt/cribl/diag/cribl-logstream ← hostname ← <datetir
## Creating and sending a diagnostic bundle
# $CRIBL_HOME/bin/cribl diag send -c 420420
Sent LogStream diagnostic bundle to Cribl Support
## Sending a previously created diagnostic bundle
# $CRIBL_HOME/bin/cribl diag send -p /opt/cribl/diag/cribl-logstream—hostname>-<datetime>.tar.§
Sent LogStream diagnostic bundle to Cribl Support
```

Known Issues

2021-02-10 – v.2.4.0, 2.4.1 – With Splunk HEC Source, JSON payloads containing embedded objects trigger high CPU usage

Problem: Splunk HEC JSON payloads containing nested objects trigger high CPU usage, due to a flaw in JSON parsing.

Workaround: If you encounter this problem, rolling back to LogStream 2.3.4 resolves it. **Fix**: In LogStream 2.4.2.

2021-01-30 - v.2.4.0 - Worker Nodes cannot connect to Master

Problem: Worker Nodes cannot connect to the Master after the Master is upgraded to v.2.4.0.

Workaround: Disable compression on the Workers. You can do so through the Workers' UI at System Settings > Distributed Settings > Master Settings > Compression, or by commenting out this line in each Worker's cribl.yml config file:

compression: gzip

Fix: In LogStream 2.4.1.

2021-01-25 – v.2.4.0 – S3 collection stops working due to auth secret key issues.

Problem: S3 collection stops after upgrade to 2.4.0 due to secret key re-encryption. **Workaround:** Re-configure S3, save and re-deploy.

Fix: In LogStream 2.4.1.

2021-01-14 – v.2.4.0 – Google Cloud Storage Destination Needs Extra Endpoint to Initialize

Problem: The Google Cloud Storage Destination fails to initialize, displaying an error of the form: Bucket does not exist!

Workaround: In the outputs.yml file, under your cribl-gcp-bucket key endpoint, add: https://storage.googleapis.com.(in a single-instance deployment, locate this file at \$CRIBL_HOME/local/cribl/outputs.yml. In a distributed deployment, locate it

at \$CRIBL_HOME/groups/<group name>/local/cribl/outputs.yml.)
Fix: In LogStream 2.4.1.

2021-01-14 – v.2.4.0 – Worker Groups' Settings > Access Management Is Absent from UI

Problem: In this release, the **Worker Groups** > **system Settings** UI did not display the expected **Access Management**, **Authentication**, and **Local Users** sections. **Workaround**: **Manually edit** the users.json file.

Fix: In LogStream 2.4.1.

2021-01-13 – v.2.4.0 – Route Filters Aren't Copied to Capture Modal

Problem: On the **Routes** page, selecting **Capture New** in the right pane does not copy custom **Filter** expressions to the resulting **Capture Sample Data** modal. That modal's **Filter Expression** field always defaults to true.

Workarounds: 1. Bypass the Capture New button. Instead, from the Route's own
••• (Options) menu, select Capture. This initiates a capture with the Filter Expression
correctly populated. 2. Copy/paste the expression into the Capture Sample Data modal's
Filter Expression field. Or, if the expression is displayed in that field's history drop-down,
retrieve it.

Fix: In LogStream 2.4.1.

2021-01-13 – v.2.4.0 – Destinations' Documentation Doesn't Render from UI

Problem: Clicking the **Help Help D?** link in a Destination's configuration modal displays the error message: "Unable to load docs. Please check LogStream's online documentation instead."

Workarounds: 1. Go directly to the online Destinations docs, starting here. 2. Follow the UI link to the docs landing page, click through to open or download the current PDF, and scroll to its Destinations section.

Fix: In LogStream 2.4.1.

2021-01-13 – v.2.4.0 – Esc Key Doesn't Consistently Close Modals

Problem: Pressing Esc with focus on a modal's drop-down or slider doesn't close the modal as expected. (Pressing Esc with focus on a free-text field, combo box, or nothing does close the modal – displaying a confirmation dialog first, if you have unsaved changes.)

Workarounds: Click the **X** close box at upper right, or click **Cancel** at lower right. **Fix**: Planned for LogStream 2.4.2.

2020-12-17 – v.2.3.0+ – Free-License Expiration Notice, Blocked Inputs

Problem: LogStream reports an expired Free license, and blocks inputs, even though Free licenses in v.2.3.0 do not expire.

Workaround: This is caused by time-limited Free license key originally entered in a LogStream version prior to 2.3.0. Go to **Settings > Licensing**, click to select and expand your expired Free license, and click **Delete license**. LogStream will recognize the new, permanent Free license, and will restore throughput.

Fix: In LogStream 2.4.1.

2020-11-14 – v.2.3.3 – Null Fields Redacted in Preview, but Still Forwarded

Problem: Where event fields have null values, LogStream (by default) displays them as struck-out in the right Preview pane. The preview is misleading, because the events are still sent to the output.

Workaround: If you do want to prevent fields with null values from reaching the output, use an Eval Function, with an appropriate Filter expression, to remove them.

Fix: Preview corrected in LogStream 2.3.4.

2020-10-27 – v.2.3.2 – Cannot Name or Save New Event Breaker Rule

Problem: After clicking **Add Rule** in a new or existing Event Breaker Ruleset, the **Event Breaker Rule** modal's **Rule Name** field is disabled. Because **Rule Name** is mandatory field, this also disables saving the Rule via the **OK** button. **Fix**: In LogStream 2.3.3.

2020-10-12 – v.2.3.1 – Deleting One Function Deletes Others in Same Group

Problem: After inserting a new Function into a group and saving the Pipeline, deleting the Function also deletes other Functions lower down in the same group.

Fix: In LogStream 2.3.2.

Workaround: Move the target Function out of the group, resave the Pipeline, and only then delete the Function.

2020-09-27 – v.2.3.1 – Enabling Boot Start as Different User Fails

Problem: When a root user tries to enable boot-start as a different user (e.g., using /opt/cribl/bin/cribl boot-start enable -u <some-username>), they receive an error of this form:

error: found user=0 as owner for path=/opt/cribl, expected uid=NaN.

Please make sure CRIBL_HOME and its contents are owned by the uid=NaN by running:

[sudo] chown -R NaN:[\$group] /opt/cribl

Fix: In LogStream 2.3.2.

Workaround: Install LogStream 2.2.3 (which you can download here), then upgrade to 2.3.1.

2020-09-17 – v.2.3.0 – Worker Groups menu tab hidden after upgrade to LogStream 2.3.0

Problem: Upon upgrading an earlier, licensed LogStream installation to v. 2.3.0, the **Worker Groups** tab might be absent from the Master Node's top menu.

Fix: In LogStream 2.3.1.

Workaround: Click the **Home > Worker Groups** tile to access Worker Groups.

2020-09-17 – v.2.3.0 – Cannot Start LogStream 2.3.0 on RHEL 6, RHEL 7

Problem: Upon upgrading to v. 2.3.0, LogStream might fail to start on RHEL 6 or 7, with an error message of the following form. This occurs when the user running LogStream doesn't match the LogStream binary's owner. LogStream 2.3.0 applies a restrictive permissions check using id -un <uid>, which does not work with the version of id that ships with these RHEL releases.

id: 0: No such user

ERROR: Cannot run command because user=root with uid=0 does not own executable

Fix: In LogStream 2.3.1.

Workaround: Update your RHEL environment's id version, if possible.

2020-09-17 – v.2.3.0 – Cannot Start LogStream 2.3.0 with OpenId Connect

Problem: Upon upgrading an earlier LogStream installation to v. 2.3.0, OIDC users might be unable to restart the LogStream server.

Fix: In LogStream 2.3.1.

Workaround: Edit \$CRIBL_HOME/default/cribl/cribl.yml to add the following lines to its the auth section:

filter_type: email_whitelist
scope: openid profile email

2020-06-11 – v.2.1.x – Can't switch from Worker to Master Mode

Problem: In a Distributed deployment, attempting to switch Distributed Settings from Worker to Master Mode blocks with a spurious "Git not available...Please install and try again" error message.

Fix: In LogStream 2.3.0.

Workaround: To initialize git, switch first from Worker to Single mode, and then from Single to Master mode.

2020-05-19 - v.2.1.x - Login page blocks

Problem: Entering valid credentials on the login page (e.g., http://localhost:9000/login) yields only a spinner.

Fix: In LogStream 2.3.0.

Workaround: Trim /login from the URL.

2020-02-22 - v.2.1.x - Deleting resources in default/

Problem: In a Distributed deployment, deleting resources in default/ causes them to reappear on restart.

Workaround/Fix: In progress.

2019-10-22 – v. 2.0 – In-product upgrade issue on v2.0

Problem: Using in-product upgrade feature in v.1.7 (or earlier) fails to upgrade to v2.0, due to package-name convention change.

Workaround/Fix: Download the new version and upgrade per steps laid out here.

2019-08-27 – v.1.7 – In-product upgrade issue on v1.7

Problem: Using in-product upgrade feature in v1.6 (or earlier) fails to upgrade to v1.7 due to package name convention change.

Workaround/Fix: Download the new package and upgrade per steps laid out here.

2019-03-21 - v.1.4 - S3 stagePath issue on upgrade to v.1.4+

Problem: When upgrading from v1.2 with a S3 output configured, stagePath was allowed to be undefined. In v.1.4+, stagePath is a required field. This might causing

schema violations when upgrading older configs.

Workaround/Fix: Reconfigure the output with a valid stagePath filesystem path.

Third-Party Software

Credits

Various components in Cribl LogStream are built and enhanced with software under free or open source licenses. We thank those projects' contributors!

```
@azure-storage-blob-10.3.0
ag-grid-community-19.1.2
ag-grid-react-19.1.2
ajv-6.9.2
ajv-errors-1.0.1
antd-3.13.0
as-table-1.0.36
avsc-5.4.9
aws-sdk-2.323.0
cidr-matcher-1.0.5
classnames-2.2.6
color-hash-1.0.3
d3-time-format-2.1.3
date-fns-1.29.0
diff-3.5.0
diff2html-2.11.3
echarts-4.3.0
echarts-4.6.0
escodegen-1.11.1
esprima-4.0.1
express-4.16.3
fast-bitset-1.3.2
file-saver-1.3.8
http-proxy-agent-3.0.0
https-proxy-agent-4.0.0
jwt-simple-0.5.6
kafkajs-1.11.0
kafkajs-1.4.5
```

kafkajs-snappy-1.1.0

ldapts-1.10.0

limiter-1.1.4

lodash-4.17.15

lz4js-0.2.0

maxmind-3.1.2

node-cache-4.2.0

node-uuid-1.4.8

numeral-2.0.6

pako-1.0.10

papaparse-5.0.0-beta.0

pbf-3.2.1

proxy-from-env-1.0.0

query-string-6.1.0

react-16.7.0

react-dom-16.7.0

react-grid-layout-0.16.6

react-router-dom-4.3.1

react-sortable-hoc-0.8.3

react-split-pane-0.1.82

regexpp-2.0.0

requirejs-2.3.6

resize-observer-polyfill-1.5.0

rxjs-6.2.2

saxen-8.1.0

simple-git-1.126.0

snappyjs-0.6.0

snmp-native-1.2.0

streamcount-1.0.1

tar-stream-1.6.1

url-0.11.0

winston-3.0.0

xmlbuilder-10.0.0

yaml-1.3.2