

Cloudera Runtime 7.2.17

Kafka Connect

Date published: 2019-08-22

Date modified: 2023-06-27

CLOUDERA

<https://docs.cloudera.com/>

Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Kafka Connect Overview.....	5
Kafka Connect Setup.....	5
Using Kafka Connect.....	6
Configuring the Kafka Connect Role.....	7
Managing, Deploying and Monitoring Connectors.....	7
Writing Kafka data to Ozone with Kafka Connect.....	8
Writing data to Ozone in an unsecured cluster with Kafka Connect.....	8
Writing data to Ozone in a Kerberos and TLS/SSL enabled cluster with Kafka Connect.....	10
Using the AvroConverter.....	12
Enabling and configuring the AvroConverter.....	14
Configuring AvroConverter for Debezium connectors with Schema Registry integration.....	16
Reading Debezium Avro records from Kafka.....	17
Configuring EOS for source connectors.....	18
Securing Kafka Connect.....	19
Kafka Connect to Kafka broker security.....	20
Configuring TLS/SSL encryption for the Kafka Connect role.....	20
Configuring Kerberos authentication for the Kafka Connect role.....	20
Kafka Connect REST API security.....	21
Kafka Connect REST API authentication.....	21
Kafka Connect Rest API authorization.....	22
Kafka Connect connector configuration security.....	24
Kafka Connect Secrets Storage.....	24
Configuring connector JAAS configuration and Kerberos principal overrides.....	32
Configuring a Nexus repository allow list.....	33
Single Message Transforms.....	33
Configuring an SMT chain.....	34
ConvertFromBytes.....	36
ConvertToBytes.....	38
Connectors.....	40
Installing Kafka Connect connectors.....	43
Installing custom developed Kafka Connect connectors.....	44
Installing FileStream connectors.....	44
Setting up the Debezium Db2 Source connector [Technical preview].....	47
Setting up the Debezium MySQL Source connector.....	47
Setting up the Debezium Oracle Source connector.....	48
Setting up the Debezium PostgreSQL Source connector.....	49
Setting up the Debezium SQL Server Source connector.....	50
HTTP Source connector.....	51
JDBC Source connector.....	53

JMS Source connector.....	56
MQTT Source connector.....	58
SFTP Source connector.....	59
Stateless NiFi Source and Sink.....	62
Dataflow development best practices for Stateless NiFi.....	63
Kafka Connect worker assignment and Stateless NiFi.....	70
Kafka Connect log files and Stateless NiFi.....	73
Kafka Connect tasks and Stateless NiFi.....	73
Developing a dataflow for Stateless NiFi.....	76
Deploying a dataflow using Stateless NiFi.....	79
Downloading and viewing the predefined Stateless NiFi dataflows shipped in CDP.....	83
Configuring flow.snapshot for Stateless NiFi connectors.....	84
Tutorial: developing and deploying a JDBC Source dataflow in Kafka Connect using Stateless NiFi.....	86
Syslog TCP Source connector.....	96
Syslog UDP Source connector.....	99
ADLS Sink connector.....	101
Amazon S3 Sink.....	104
Configuration example.....	105
HDFS Sink.....	106
Configuration example for writing data to HDFS.....	106
Configuration example for writing data to Ozone FS.....	108
HDFS Stateless Sink connector.....	109
HTTP Sink connector.....	115
InfluxDB Sink connector.....	117
JDBC Sink connector.....	119
Kudu Sink connector.....	121
S3 Sink connector.....	123

Kafka Connect Overview

Get started with Kafka Connect.

Kafka Connect is a tool for streaming data between Apache Kafka and other systems in a reliable and scalable fashion. Kafka Connect makes it simple to quickly define connectors that move large collections of data into and out of Kafka. Source connectors can ingest entire databases or collect metrics from all your application servers into Kafka topics, making the data available for stream processing with low latency. Sink connectors can deliver data from Kafka topics into secondary storage and query systems or into batch systems for offline analysis.

Kafka Connect in CDP is shipped with many different Cloudera developed as well as publicly available sink and source connectors. Each of which cover a specific use case for streaming data. In addition to the connectors available by default, installing custom developed or third-party connectors is also possible. All connectors can be deployed, managed, and monitored using the Streams Messaging Manager UI (recommended), Streams Messaging Manager REST API, or Kafka Connect REST API.

See the *Related Information* section for more information regarding how to set up, configure, use, manage, and monitor Kafka Connect and its connectors.

Related Information

[Kafka Connect Setup](#)

[Using Kafka Connect](#)

[Securing Kafka Connect](#)

[Connectors](#)

[Streams Messaging Reference](#)

Kafka Connect Setup

Learn how you can set up Kafka Connect in CDP Public Cloud with the Streams Messaging cluster templates.

Kafka Connect in CDP is provided in the form of a Kafka service role. The role is called Kafka Connect. In CDP Public Cloud, Kafka Connect can be deployed with the Streams Messaging Light Duty, Heavy Duty, and High Availability templates available in the Data Hub service. However, there are some differences between how Kafka Connect is deployed with each of the templates.

Light Duty

In the Streams Messaging Light Duty template, Kafka Connect roles share a host group and nodes with other service roles. Specifically, Kafka Connect is deployed in the Broker and Core_broker host groups and shares nodes with Kafka broker, ZooKeeper, and SRM Driver roles. The Core_broker is a mandatory host group in Light Duty clusters and has an instance (node) count of three. This means that by default, all Light Duty cluster you deploy will include a minimum of three Kafka Connect roles. If required, Kafka Connect can also be scaled in Light Duty clusters by scaling the Broker host group.

Heavy Duty and High Availability

In the Heavy Duty and High Availability templates, Kafka Connect roles are deployed in a dedicated host group named Connect. The instance count of this host group is set to zero by default. This means that, unless specifically provisioned, Kafka Connect roles are **not** deployed in Heavy Duty or High Availability clusters by default. The instance (node) count of the Connect host group can be configured in Data Hub during cluster provisioning. Additionally, you can add Kafka Connect to already running Heavy Duty or High Availability cluster by scaling the Connect host group.

For more information regarding the Streams Messaging cluster templates, scaling, and cluster deployment with Data Hub, see the *Related Information* section below.

Streams Messaging Manager integration

Kafka Connect deployed in CDP Public Cloud is by default configured to integrate with the Streams Messaging Manager (SMM) instance running in the cluster. This means that, if you provisioned your cluster with Kafka Connect, you will be able to manage, monitor, and deploy Kafka Connect connectors with SMM as soon as the cluster is deployed with no additional configuration required. Kafka Connect can be managed in SMM on the Connect page. For more information, see *Monitoring Kafka Connect using Streams Messaging Manager*.

Related Information

[Setting up your Streams Messaging cluster](#)

[Streams Messaging cluster layout](#)

[Scaling Kafka Connect](#)

[Monitoring Kafka Connect using Streams Messaging Manager](#)

Using Kafka Connect

Learn how you can manage, monitor and configure Kafka Connect.

There are multiple ways that you can manage, monitor and configure Kafka connect in CDP. It is important to make a distinction between the Kafka Connect roles and Kafka Connect connectors as they are managed through different interfaces.

The Kafka Connect Role

Kafka Connect roles are roles that you can deploy within a Kafka service, they represent Kafka Connect workers. Because these roles are deployed with and managed by Cloudera Manager, their management and configuration is done through Cloudera Manager. In other words the connect workers you deploy will be managed and configured in Cloudera Manager. For more information see *Configuring Kafka Connect and the Cloudera Manager Documentation*.

Kafka Connect Connectors

Connectors are not managed by Cloudera Manager, instead multiple other interfaces can be used to interact with them. These are the following:

- Streams Messaging Manager UI

The Streams Messaging Manager UI is a recommended interface in CDP to manage connectors. Comprehensive documentation is provided on the usage of the UI in *Monitoring Kafka Connect*.

- The Streams Messaging Manager REST API

The Streams Messaging Manager REST API is a recommended interface in CDP to manage connectors. For more information, see the *SMM REST API Reference*.



Note: Both the Streams Messaging Manager UI and Streams Messaging Manager REST API support the same management actions.

- The Kafka Connect REST API

The Kafka connect REST API can be used to manage connectors. However, its usage is not recommended by Cloudera, nor is separate documentation provided. For more information, see the upstream Apache Kafka documentation.

Related Information

[Monitoring Kafka Connect](#)

[SMM REST API Reference](#)

[Apache Kafka Documentation](#)

Configuring the Kafka Connect Role

Learn more about how you can configure the Kafka Connect role as well as its notable properties.

The Kafka Connect role is deployed with and managed by Cloudera Manager. You can view a list of configuration properties as well as configure them by going to Kafka serviceConfiguration and selecting the Kafka Connect filter in the Filters pane. Most of the properties available are the standard worker properties defined in upstream Apache Kafka. However, there are a number of notable properties that you should be familiar with when using Kafka Connect on CDP. The notable properties are the following:

Broker List for Kafka Connect

These are the brokers Kafka Connect should connect to. You specify the brokers by adding IP:port or hostname:port pairs. The brokers you specify here should be the brokers that are running in the same cluster that Kafka Connect is deployed on. While you can specify a single broker, Cloudera recommends that you specify multiple for high availability.

group.id

The group ID is a unique string that Kafka Connect roles use to form a cluster of connect workers. This property should always be set to the same string for all Kafka Connect roles that are deployed in the same cluster. Kafka Connect will not function properly if there is a mismatch between the group IDs. Therefore, Cloudera recommends that you use the default values for all roles that you deploy.

plugin.path

The directory where the Kafka Connect connector plugins are stored. This is the directory where JAR files for custom connectors can be placed. Cloudera recommends that you use the default path which is /var/lib/kafka.

kafka.connect.rest.port and kafka.connect.secure.rest.port

These are the ports that Kafka Connect API will accept requests on. If the Kafka Connect role is TLS/SSL enabled, it will use the port specified in kafka.connect.secure.rest.port, if TLS/SSL is not enabled, it will use the port specified in kafka.connect.rest.port. Take note of the ports as they are needed when you set up SMM to manage and monitor Kafka Connect.

A comprehensive list of all properties available for the Kafka Connect role is available in Kafka Properties in Cloudera Runtime.

Related Information

[Kafka Properties in Cloudera Runtime](#)

Managing, Deploying and Monitoring Connectors

Learn more about managing connectors.

You can manage, monitor, deploy and interact with Kafka Connect and Kafka Connect connectors either through the SMM UI or SMM REST API. For more information see [Monitoring Kafka Connect](#), as well as the [SMM REST API Reference](#).

Related Information

[Monitoring Kafka Connect](#)

[SMM REST API Reference](#)

Writing Kafka data to Ozone with Kafka Connect

You can use the HDFS Sink Connector developed by Cloudera and shipped with Cloudera Runtime to write Kafka data to the Ozone filesystem. This can be done on both secure and unsecure clusters by deploying and configuring a new connector with the SMM UI.

Writing data to Ozone in an unsecured cluster with Kafka Connect

You can use the HDFS Sink Connector developed by Cloudera, i to write Kafka topic data to Ozone in an unsecure cluster. Connector deployment and configuration is done using the SMM UI.

About this task

The following steps walk you through how the Cloudera-developed HDFS Sink Connector can be set up and deployed using the SMM UI.

In addition to the connector setup, these steps also describe how you can create a test topic and populate it with data using Kafka command line tools. If you already have a topic that is ready for use and do not want to create a test topic, you can skip steps 1 through 3. These steps deal with topic creation, message consumption, and message production. These steps are optional.

Before you begin

An unsecure CDP PvC Base cluster with Kafka, SMM, and Ozone is set up and configured.

Procedure

1. Create a Kafka topic:

- a) SSH into one of the hosts in your cluster.

```
ssh [***USER***]@[***MY-CLUSTER-HOST.COM***]
```

- b) Create a topic with the kafka-topics tool.

```
kafka-topics --create --bootstrap-server [***MY-CLUSTER-HOST.COM:9092***] --replication-factor 1 --partitions 1 --topic [***TOPIC***]
```

If an out of memory exception is displayed while running this command, increase the JVM heap with the following command and try again:

```
export KAFKA_OPTS="-Xmx1g -Xms1g"
```

- c) Verify that the topic is created.

```
kafka-topics --list --bootstrap-server [***MY-CLUSTER-HOST.COM:9092***]
```

2. Produce messages to your topic with the kafka-console-producer.

```
kafka-console-producer --broker-list [***MY-CLUSTER-HOST.COM:9092***] --topic [***TOPIC***]
```

Enter test messages once the tool is running.

```
>my first message  
>my second message
```


3. Consume messages:

- a) Open a new terminal session and log in to one of the hosts in your cluster.
- b) Consume messages with the kafka-console-consumer.

```
kafka-console-consumer --from-beginning --bootstrap-server [***MY-CLUSTER-HOST.COM:9092***] --topic [***TOPIC***]
```

The messages you produced with kafka-console-producer appear. In addition, you can switch to the terminal session that is running kafka-console-producer and produce additional messages. These new messages appear in real time in the session running kafka-console-consumer.

4. Deploy and configure an HDFS Sink Connector:

- a) In Cloudera Manager, select the Streams Messaging Manager service.
- b) Click Streams Messaging Manager Web UI.
- c) Click the Connect option in the left-side menu.
- d) Click + New Connector to add a new connector.
- e) Go to the Sink Connectors tab and select the HDFS Sink Connector.

On the UI the HDFS Sink Connector is represented by its class name, which is `com.cloudera.dim.kafka.connector.hdfs.HdfsSinkConnector`.

- f) Enter a name for the connector.
- g) Configure the connector.

Use the following example as a template:

```
{
  "connector.class": "com.cloudera.dim.kafka.connect.hdfs.HdfsSinkConnector",
  "hdfs.uri": "ofs://ozonel/voll/bucket1/",
  "hdfs.output": "/topics_output/",
  "tasks.max": "1",
  "topics": "testTopic",
  "hadoop.conf.path": "file:///etc/hadoop/conf",
  "output.writer": "com.cloudera.dim.kafka.connect.partition.writers.txt.TxtPartitionWriter",
  "value.converter": "org.apache.kafka.connect.storage.StringConverter",
  "output.storage": "com.cloudera.dim.kafka.connect.hdfs.HdfsPartitionStorage",
  "hdfs.kerberos.authentication": "false"
}
```

Ensure that you replace the values of `hdfs.uri` and `hdfs.output` with valid Ozone paths. The template gives an example of how these paths should look like. Replace any other values depending on your cluster configuration and requirements.

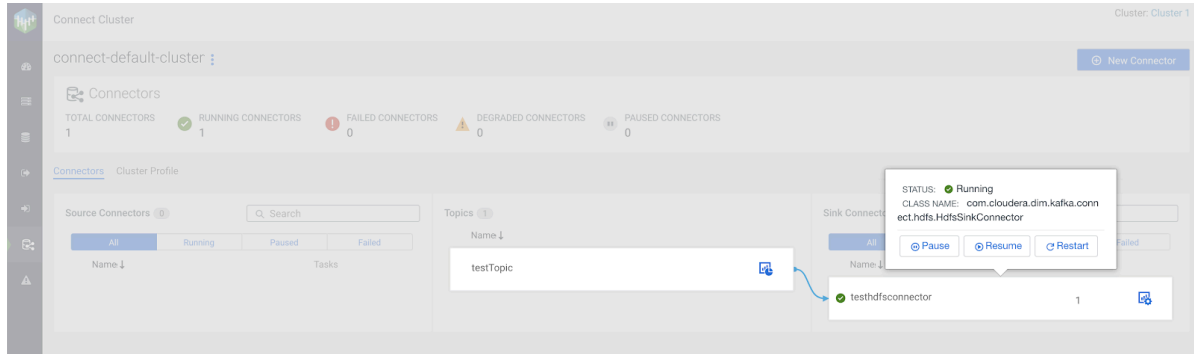
- h) Click Validate.

The validator displays any JSON errors in your configuration. Fix any errors that are displayed. If your JSON is valid, the JSON is valid message is displayed in the validator.

- i) Click Next.
- j) Review your connector configuration.
- k) Click Deploy .

5. Verify that connector deployment is successful:
 - a) In the SMM UI, click the Connect option in the left-side menu.
 - b) Click on either the topic or the connector you created.

If connector deployment is successful, a flow is displayed between the topic you specified and the connector you created.



6. Verify that topic data is written to Ozone.

You can do this by listing the files under the `ofs://` location you specified in the connector configuration.

Writing data to Ozone in a Kerberos and TLS/SSL enabled cluster with Kafka Connect

You can use the HDFS Sink Connector developed by Cloudera to write Kafka topic data to Ozone in a secure cluster. Connector deployment and configuration is done using the SMM UI.

About this task

The following steps walk you through how the Cloudera-developed HDFS Sink Connector can be set up and deployed using the SMM UI.

In addition to the connector setup, these steps also describe how you can create a test topic and populate it with data using Kafka command line tools. If you already have a topic that is ready for use and do not want to create a test topic, you can skip steps 1 through 3. These steps deal with topic creation, message consumption, and message production. They are not necessary to carry out.

Before you begin

- Ensure that a Kerberos and TLS/SSL enabled CDP PvC Base cluster with Kafka, SMM, and Ozone is set up and configured.
- To create a test topic with the Kafka console tools, you must ensure that a `.properties` client configuration file is available for use.

You can create one using the following example as a template:

```
sasl.jaas.config=com.sun.security.auth.module.Krb5LoginModule required u
seKeyTab=true keyTab="[*]PATH TO KEYTAB FILE[*]" principal="[*]KERBEROS
PRINCIPAL[*]";
security.protocol=SASL_SSL
sasl.mechanism=GSSAPI
sasl.kerberos.service.name=kafka
ssl.truststore.location=[*]TRUSTSTORE LOCATION[*]
```

Procedure

1. Create a Kafka topic:

- a) SSH into one of the hosts in your cluster.

```
ssh [***USER***]@[***MY-CLUSTER-HOST.COM***]
```

- b) Create a topic with the kafka-topics tool.

```
kafka-topics --create --bootstrap-server [***MY-CLUSTER-
HOST.COM:9093***] --replication-factor 1 --partitions 1 --topic [***T
OPIC***] --command-config [***CLIENT CONFIG FILE***]
```

If an out of memory exception is displayed, increase the JVM heap with the following command and try again:

```
export KAFKA_OPTS="-Xmx1g -Xms1g"
```

- c) Verify that the topic is created.

```
kafka-topics --list --bootstrap-server [***MY-CLUSTER-HOST.COM:9093***]
--command-config [***CLIENT CONFIG FILE***]
```

2. Produce messages to your topic with the kafka-console-producer.

```
kafka-console-producer --broker-list [***MY-CLUSTER-HOST.COM:9093***] --t
opic [***TOPIC***] --producer.config [***CLIENT CONFIG FILE***]
```

Enter test messages once the tool is running.

```
>my first message
>my second message
```

3. Consume messages:

- a) Open a new terminal session and log in to one of the hosts in your cluster.
- b) Consume messages with the kafka-console-consumer.

```
kafka-console-consumer --from-beginning --bootstrap-server [***MY-
CLUSTER-HOST.COM:9093***] --topic [***TOPIC***] --consumer.conf
ig [***CLIENT CONFIG FILE***]
```

The messages you produced with kafka-console-producer appear. You can switch to the terminal session that is running kafka-console-producer and produce additional messages. These new messages appear in real time in the session running kafka-console-consumer.

4. Deploy and configure an HDFS Sink Connector:

- a) In Cloudera Manager, select the Streams Messaging Manager service.
- b) Click Streams Messaging Manager Web UI to log in to the UI.
If prompted, enter a valid username as well as a password and click SIGN IN.
- c) Click the Connect option in the left-side menu.
- d) Click + New Connector to add a new connector.
- e) Go to the Sink Connectors tab and select the HDFS Sink Connector.

On the UI the HDFS Sink Connector is represented by its class name, which is `com.cloudera.dim.kafka.connect.hdfs.HdfsSinkConnector`.

- f) Enter a name for the connector.
- g) Configure the connector.

Use the following example as a template:

```
{
```

```

"connector.class": "com.cloudera.dim.kafka.connect.hdfs.HdfsSinkConnector",
"hdfs.uri": "ofs://ozonel/voll/bucket1/",
"tasks.max": "1",
"topics": "testTopic",
"hdfs.kerberos.authentication": "true",
"hdfs.kerberos.user.principal": "${cm-agent:ENV:kafka_connect_service_principal}",
"hdfs.kerberos.keytab.path": "${cm-agent:keytab}",
"hdfs.kerberos.namenode.principal": "hdfs/_HOST@REALM",
"hdfs.output": "/topics_output/",
"hadoop.conf.path": "file:///etc/hadoop/conf",
"output.writer": "com.cloudera.dim.kafka.connect.partition.writers.txt.TxtPartitionWriter",
"value.converter": "org.apache.kafka.connect.storage.StringConverter",
"output.storage": "com.cloudera.dim.kafka.connect.hdfs.HdfsPartitionStorage"
}

```

Ensure that you replace the values of `hdfs.uri` and `hdfs.output` with valid Ozone paths. The template gives an example of how these paths should look like. Replace other values depending on your cluster configuration and requirements.

h) Click Validate.

The validator displays any JSON errors in your configuration. Fix any errors that are displayed. If your JSON is valid, the JSON is valid message is displayed.

i) Click Next.

j) Review your connector configuration.

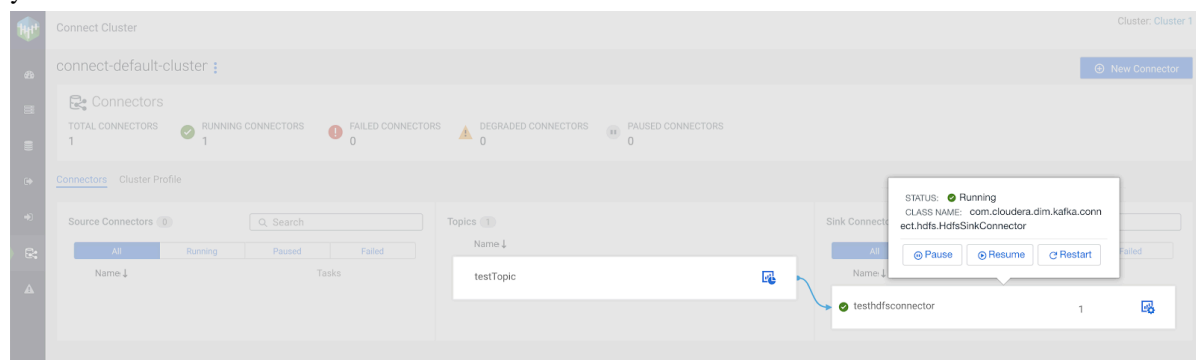
k) Click Deploy.

5. Verify that connector deployment is successful:

a) In the SMM UI, click the Connect option in the left-side menu.

b) Click on either the topic or the connector you created.

If connector deployment is successful, a flow is displayed between the topic you specified and the connector you created.



6. Verify that topic data is written to Ozone.

You can do this by listing the files under the `ofs://` location you specified in the connector configuration.

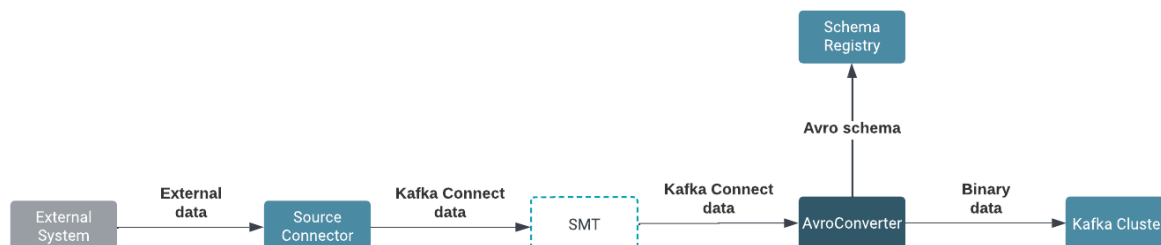
Using the AvroConverter

The `AvroConverter` is a Kafka Connect converter shipped with Cloudera Runtime that enables Kafka Connect connectors to serialize or deserialize Kafka messages, consisting of key and value pairs, from or to Avro. Using the Avro format can make your messages more compact compared to using for example JSON format. `AvroConverter` is able to convert from Avro format to Kafka Connect internal data format, and from Kafka Connect internal data to Avro format, and is able to handle Avro schemas integrated with Schema Registry automatically.

Using AvroConverter with a source connector

When using the AvroConverter with a Kafka Connect source connector, the source connector fetches the data from an external system. The fetched data is converted into Kafka Connect internal data format. If Single Message Transforms (SMT) are configured, the messages are converted based on the SMT conditions as well. If the AvroConverter is enabled, the messages are converted from the internal data format to Avro. At this step, the AvroConverter automatically also creates a new schema version in Schema Registry for the Avro data, if the schema has not been already created. In case the schema already exists in Schema Registry, schema creation is skipped and the existing ID will be assigned to the payload during the serialization. The converted data in Avro is written to the Kafka cluster with its corresponding schema ID from Schema Registry.

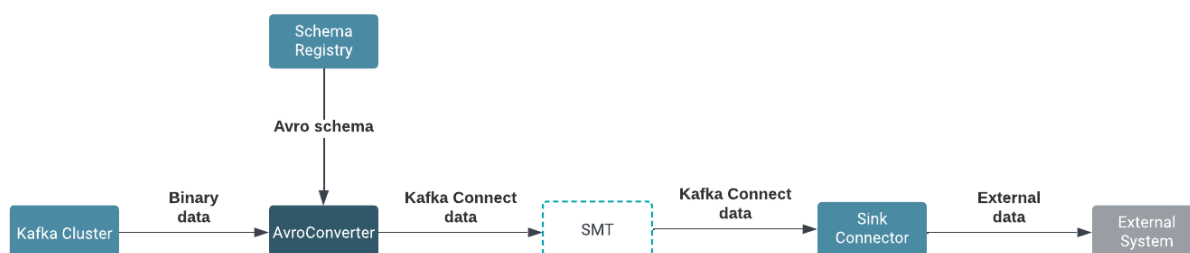
The following illustration shows the steps how AvroConverter works with a Kafka Connect source connector:



Using AvroConverter with a sink connector

When using the AvroConverter with a Kafka Connect sink connector, the data is fetched from the Kafka cluster, which includes its schema ID from Schema Registry. The AvroConverter extracts the schema ID and gets the corresponding schema from Schema Registry. The data is converted from Avro format to the Kafka Connect internal data format. If SMT are configured, the messages are converted based on the SMT conditions as well. The converted data is written by the Kafka Connect sink connector to the external system.

The following illustration shows the steps how AvroConverter works with a Kafka Connect sink connector:



Note: If the SMTs contain steps to change the schema, the schema will not be updated in Schema Registry. The schema is handled automatically by the AvroConverter before any SMTs are applied.

Mapping logical data types

AvroConverter is capable of mapping the data from the Kafka Connect internal data format to Avro. The mapping happens automatically during the conversion, which means that no property needs be configured. The AvroConverter also handles Kafka Connect specific logical types and maps them to Avro ones and vice versa.

The corresponding logical types for Kafka Connect and Avro with their direction is shown in the following table:

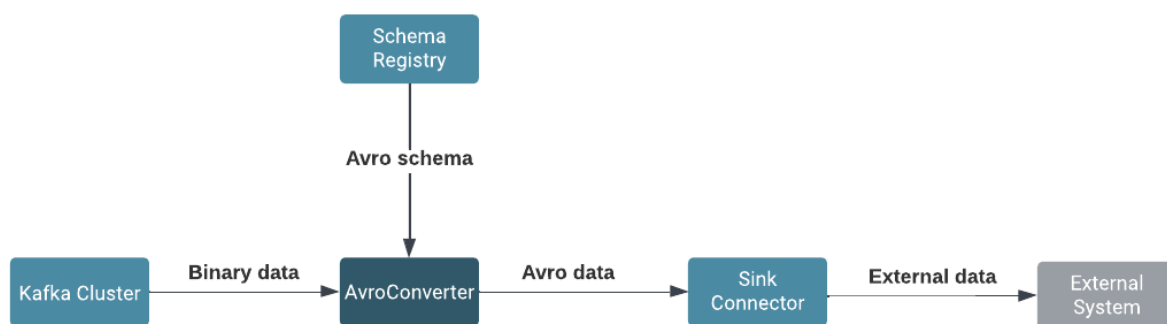
Table 1: Logical Type Mapping

Kafka Connect Logical Type	Direction	Avro Logical Type
Decimal	##	Decimal
Date	##	Date
Time	##	TimeMillis
Time	#	TimeMicros
Timestamp	##	TimestampMillis
Timestamp	#	TimestampMicros

Handling passthrough

When passthrough is enabled, messages are not converted between the Avro format and the Kafka Connect internal data format. In this case, the data is passed through unchanged to the Kafka Connect sink connector. If the passthrough is disabled, messages are converted from Avro to the Kafka Connect internal data format, which can make your input and output formats to be different.

The following illustration shows what happens when passthrough is enabled in the AvroConverter for a sink connector:



Enabling passthrough is recommended when both the input and output formats of your messages are Avro, and it can save any performance overhead generated by the conversion process. You need to ensure that the passthrough is only enabled if the connector can handle Avro objects. Cloudera provides the following connectors that are capable of handling passthrough:

- HDFS Sink connector
- S3 Sink connector

Any other third party connectors can make use of the passthrough if the connectors are implemented to handle the passthrough of the message.



Warning: Passthrough is enabled by default in AvroConverter. You need to make sure that the connector is able to handle passthrough before using AvroConverter or you must disable the passthrough. In case the connector cannot handle the passthrough, the connector will fail.

You can enable and disable passthrough by setting the `passthrough.enabled` property to true or false when configuring the connector. For more information, see the *KafkaAvroSerializer properties reference* documentation under the *Related information*

Enabling and configuring the AvroConverter

When the AvroConverter is enabled as a key or value converter, you also need to provide further configuration properties that affect the behavior of the AvroConverter. All configuration properties must be prefixed with either the

value.converter. or key.converter. prefixes. These prefixes control whether a specific property is applied to the key or value converter.

Enabling AvroConverter

You can enable AvroConverter for any connector using the key.converter and value.converter properties to the connector configuration. The value of these properties must be set to the class name of the AvroConverter.

For example, if you want to use the AvroConverter for both key and value conversion, you must add the following properties to the connector configuration:

```
"key.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
"value.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
```

Conversion is enabled separately for key or value conversion, which enables you to use different converters for the key and value. For example, if you want to use the AvroConverter to convert values, but use a different converter, such as the StringConverter, to convert keys, you must add the following properties to the connector configuration:

```
"key.converter": "org.apache.kafka.connect.storage.StringConverter",
"value.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
```

Configuring passthrough

The passthrough configuration in AvroConverter controls whether or not the data is converted into the Kafka Connect internal data format before it is written into an output file or if data is passed through the connector without conversion. Passthrough is enabled by the passthrough.enabled property, which is set to true by default.

The passthrough can be set with different prefixes for key and value converters as shown in the following example configuration:

```
"key.converter.passthrough.enabled": "false",
"value.converter.passthrough.enabled": "false",
```



Warning:

Passthrough is enabled by default in AvroConverter. You need to make sure that the connector is able to handle passthrough before using AvroConverter or you must disable the passthrough. In case the connector cannot handle the passthrough, the connector will fail. For more information, see the *Handling passthrough* section.

Enabling SerDes

To transform messages, the AvroConverter uses an embedded KafkaAvroSerializer and KafkaAvroDeserializer, which need to be configured with the AvroConverter to enable further control over the behavior.

You must provide the URL of the Schema Registry server to be able to use the KafkaAvroSerializer and KafkaAvroDeserializer, because both the serializer and the deserializer contain an embedded Schema Registry client. You can specify the Schema Registry server URL with the schema.registry.url property as shown in the following example:

```
"key.converter.schema.registry.url": "http://registryserver.example.com:9090/api/v1"
"value.converter.schema.registry.url": "http://registryserver.example.com:9090/api/v1"
```

You can configure additional properties related to serialization, deserialization and the Schema Registry client using the property reference documentations under the *Related information*.

Configuration example

The following example configuration shows a connector configuration with AvroConverter set as key and value converter with connection to a remote Schema Registry server.

```
{
  "value.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
  // use AvroConverter for the conversion of the Kafka message value
  "value.converter.passthrough.enabled": "false", // disable passthrough
  "value.converter.schema.registry.url": "http://registryserver.example.com:9090/api/v1", // URL of the schema registry server that the Schema Registry client inside the value converter will connect to
  "value.converter.schema.registry.auth.username": ["**USERNAM**"], // username to be used for authentication by the value converter's embedded Schema Registry client
  "value.converter.schema.registry.auth.password": ["**PASSWORD**"], // password to be used for authentication by the value converter's embedded Schema Registry client
  "key.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
  "key.converter.schema.registry.auth.username": ["**USERNAM**"],
  "key.converter.schema.registry.auth.password": ["**PASSWORD**"],
  "key.converter.passthrough.enabled": "false",
  "key.converter.schema.registry.url": "http://registryserver.example.com:9090/api/v1"
}
```

Configuring AvroConverter for Debezium connectors with Schema Registry integration

Learn about how to configure the AvroConverter for Debezium connector with Schema Registry integration.

When you want to use the AvroConverter for the Debezium connectors, you need to configure the passthrough, schema compatibility, schema name value and key prefix properties.

passthrough.enabled

You must disable the passthrough to have automatically created schemas for CDC records in Kafka as shown in the following example:

```
"value.converter.passthrough.enabled": "false"
"key.converter.passthrough.enabled": "false"
```

This will ensure that the CDC event fields are loaded into the internal schema of Kafka Connect, from which the field will be converted into Avro format by the AvroConverter before writing it to Kafka. When doing this conversion, the schemas are automatically created in Schema Registry as well.

schema.compatibility

Whenever a schema is created in Schema Registry, compatibility is also set to it by default. Based on the compatibility, if the schema evolves, the new version will be validated against old ones. For more information about the compatibility in Schema Registry, see the *Compatibility policies* documentation under the *Related information*.

The compatibility must be set before the first CDC event is processed by the connector, because the first schema will be created based on the compatibility policies. If backward incompatible changes could happen in the database, consider setting the parameters to NONE, because in this case the Avro schemas will also be incompatible.

The default compatibility can be set based on the following example:

```
"value.converter.schema.compatibility": "NONE"
"key.converter.schema.compatibility": "NONE"
```


**Note:**

The connector task will stop and fail if compatibility issues happen during schema evolution. In this case manual intervention is needed to continue the task.

schema.name.key.suffix and schema.name.value.suffix

When a new schema is created, the name of the schema consists of the name of the topic and an appended value, if set. However, it is possible to set additional suffixes for schema names. The default suffix is empty for the values, but :k is added as a prefix for keys. Key and value suffixes can be set on converters according the following example:

```
"value.converter.schema.name.value.suffix", "value_suffix"
"key.converter.schema.name.key.suffix", "key_suffix"
```

Additional configuration properties for Debezium connectors and Schema Registry integration with AvroConverter

The `bigint.unsigned.handling.mode` in a Debezium connector handles how bigint data types are mapped, which have the following values that can affect the behavior of the logical data type conversion:

precise

The bigint unsigned values will be stored as byte arrays in Avro with decimal logical type. You can get a byte array or a `BigDecimal` automatically when you would like to read this data. The object type you get depends on if logical type conversion is enabled or not.

long

Whatever is stored in the database, it will be stored in Avro as long, and you get a Java long object in Java Virtual Machine (JVM). In this case, precision of the mapping is smaller, but no need for logical type handling.

Check `schema.name.adjustment.mode` and `field.name.adjustment.mode` to ensure that there are no conflicts between your database related names and the Avro specific naming convention.

Check `decimal.handling.mode` to decide how precise numbers you need in different cases, and if you need to use logical type for higher precision

Check `time.precision.mode` to decide how precise dates and times you need in different cases, and if you need to use logical type for higher precision.

Reading Debezium Avro records from Kafka

Learn how to read Debezium CDC Avro records from Kafka using the `AvroConverter` functions.

Avro records can be consumed from Kafka with the Schema Registry integration as described in the *Integrating Kafka with Schema Registry* documentation.

If any of your connectors produce Avro records with logical types, you do not need to write the conversion from raw Avro data to proper Java objects as you can enable automatic conversion. To enable automatic logical type conversion, you need to set the `logical.type.conversion.enabled` property in the consumer configuration as shown with the following in the following example:

```
"logical.type.conversion.enabled": "true"
```

This conversion is handled by the internal `SerDes` classes in `KafkaAvroSerializer` and `KafkaAvroDeserializer`.

Logical type conversion automatically happens in Debezium connectors when they produce the Avro events to Kafka. Ensure that the `logical.type.conversion.enabled` property is not set as true for the key/value `SerDes` classes of the Avro Converter. The `logical.type.conversion.enabled` property only needs to be enabled on the consumer/producer side when the consumer/producer is not part of a Kafka Connect connector.

Related Information

[Single Message Transforms](#)

[KafkaAvroSerializer properties reference](#)
[KafkaAvroDeserializer properties reference](#)
[Schema Registry Client properties reference](#)
[Compatibility policies](#)
[Integrating Kafka with Schema Registry](#)

Configuring EOS for source connectors

Learn how to enable and configure exactly-once semantics (EOS) for Kafka Connect source connectors.

About this task

EOS is a framework that enables Kafka and Kafka applications to guarantee that each message is delivered precisely once without it being duplicated or lost. EOS is supported by some Kafka Connect source connectors. To enable EOS, you must enable EOS for the Kafka Connect service roles (Kafka Connect workers) and configure your connectors to use EOS.

EOS is enabled for the Kafka Connect service roles in Cloudera Manager with the Exactly-once support for source connectors property. Connectors can be configured using the SMM UI (recommended), SMM API, or Kafka Connect API.

Before you begin

- Enabling EOS for Kafka Connect service roles only enables the use of EOS on the level of the service. For individual connectors to make use of EOS, the connector must be specifically configured to use EOS.
- Not all connectors shipped in CDP support EOS. The following list collects the source connectors shipped in CDP that support EOS.
 - JDBC Source
 - SFTP Source
 - NiFi Stateless Source
 - FileStreamSource



Important: The NiFi Stateless Source connector is used to deploy your own, custom developed dataflows as Kafka Connect connectors. In order for your connector to support EOS, the dataflow that you deploy must include a source processor that supports state management on the scope of the cluster. For more information see the [official NiFi documentation](#).

- If you are using a third-party connector, review the documentation of the connector to learn whether EOS is supported.

Procedure

1. Configure Kafka Connect service roles:
 - a) In Cloudera Manager, select the Kafka service.
 - b) Go to Configuration.
 - c) Set Exactly-once support for source connectors in the cluster to preparing.
 - d) Restart all Kafka Connect service roles.
 - e) Set the Exactly-once support for source connectors in the cluster to enabled.
 - f) Restart all Kafka Connect service roles.
2. Configure source connectors.

The following table collects the common source connector properties that you can use to enable and configure EOS. Configure the following properties based on your requirements. Cloudera recommends that you use the

SMM UI to configure and deploy connectors. For more information on connector configuration and deployment see *Deploying and managing Kafka Connect connectors in SMM*

Table 2: EOS specific Kafka Connect source connector properties

Name	Default value	Description
exactly.once.support	requested	Permitted values are requested and required. If set to required, forces a preflight check for the connector to ensure that it can provide exactly-once delivery with the given configuration. Some connectors may be capable of providing exactly-once delivery but not signal to Kafka Connect that they support this. In this case, review the documentation for the connector before connector deployment and set this property to requested. Additionally, if the value is set to required but the worker that performs preflight validation does not have exactly-once support enabled for source connectors, requests to create or validate the connector will fail.
transaction.boundary	poll	Permitted values are poll, connector, and interval. If set to poll, a new producer transaction is started and committed for every batch of records that each task from this connector provides to Kafka Connect. If set to connector, relies on connector-defined transaction boundaries; note that not all connectors are capable of defining their own transaction boundaries, and in that case, attempts to create them with this property set to connector will fail. If set to interval, commits transactions only after a user-defined time interval has passed.
offsets.storage.topic	null	The name of a separate offsets topic to use for this connector. If left empty or not specified, the worker's global offsets topic name is used. If specified, the offsets topic is created if it does not already exist on the Kafka cluster targeted by this connector (which may be different from the one used for the worker's global offsets topic if the bootstrap.servers property of the connector's producer has been overridden from the worker's).
transaction.boundary.interval.ms	null	If transaction.boundary is set to interval, determines the interval for producer transaction commits by connector tasks. If unset, defaults to the value of the worker-level offset.flush.interval.ms property.

Related Information

[Deploying a Kafka Connect connector in SMM](#)

Securing Kafka Connect

Learn about the different options and features available for securing your Kafka Connect deployment.

Kafka Connect to Kafka broker security

Learn about the options you have for securing the connection between Kafka Connect and Kafka brokers.

Configuring TLS/SSL encryption for the Kafka Connect role

Kafka Connect roles inherit the TLS/SSL configuration of the parent Kafka service. If you are deploying Kafka Connect roles under a Kafka service that already has TLS/SSL enabled, Cloudera Manager will automatically enable TLS/SSL for Connect as well. If required however, you can manually enable or disable TLS/SSL.

Procedure

1. In Cloudera Manager, select the Kafka service.
2. Go to Configuration.
3. Find and configure the following properties based on your cluster and requirements.

Table 3:

Cloudera Manager Property	Description
Enable TLS/SSL for Kafka Connect ssl_enabled	Encrypt communication between clients and Kafka Connect using Transport Layer Security (TLS) (formerly known as Secure Socket Layer (SSL)).
Kafka Connect TLS/SSL Server JKS Keystore File Location ssl_server_keystore_location	The path to the TLS/SSL keystore file containing the server certificate and private key used for TLS/SSL. Used when Kafka Connect is acting as a TLS/SSL server.
Kafka Connect TLS/SSL Server JKS Keystore File Password ssl_server_keystore_password	The password for the Kafka Connect keystore file.
Kafka Connect TLS/SSL Server JKS Keystore Key Password ssl_server_keystore_keypassword	The password that protects the private key contained in the JKS keystore used when Kafka Connect is acting as a TLS/SSL server.
Kafka Connect TLS/SSL Trust Store File ssl_client_truststore_location	The location on disk of the trust store, used to confirm the authenticity of TLS/SSL servers that Kafka Connect might connect to. This trust store must contain the certificate(s) used to sign the service(s) connected to. If this parameter is not provided, the default list of well-known certificate authorities is used instead.
Kafka Connect TLS/SSL Trust Store Password ssl_client_truststore_password	The password for the Kafka Connect TLS/SSL Trust Store File. This password is not required to access the trust store; this field can be left blank. This password provides optional integrity checking of the file. The contents of trust stores are certificates, and certificates are public information.
ssl.client.auth	Client authentication mode for SSL connections. This configuration has three valid values, required, requested, and none. If set to required, client authentication is required. If set to requested, client authentication is requested and clients without certificates can still connect. If set to none, which is the default value, no client authentication is required.

4. Click Save Changes.
5. Restart the service.

Results

TLS/SSL encryption is configured for the Kafka Connect role.

Configuring Kerberos authentication for the Kafka Connect role

Learn how Kerberos authentication is configured for the Kafka Connect role.

Kafka Connect roles inherit the Kerberos configuration of the parent Kafka service. If you are deploying Kafka Connect roles under a Kafka service that already has Kerberos enabled, Cloudera Manager will automatically enable Kerberos for Kafka Connect as well. Other than making sure that the Kafka service's Kerberos configuration is correctly set, no additional user action is required.

Kafka Connect REST API security

Learn about the options you have for securing the Kafka Connect REST API.

Kafka Connect REST API authentication

Learn about the authentication options that are available for the Kafka Connect REST API.

Configuring TLS/SSL client authentication for the Kafka Connect REST API

Learn how to configure TLS/SSL client authentication for the Kafka Connect REST API.

About this task

You can secure the Kafka Connect API by configuring the Kafka Connect roles to require SSL Client authentication. This can be done by setting the SSL Client Authentication property to required. When set to required, only clients that pass SSL client authentication will be able to access the Kafka Connect API. As a result, any client that you would like to give access to should have its certificate added to the Kafka Connect truststore. This includes Streams Messaging Manager (SMM) as well. Cloudera recommends that in secure environments only SMM is given access to the Kafka Connect API.

In addition to setting client authentication to required, you may also want to consider setting up a firewall using third party tools to further secure access to the Kafka Connect API. Note however, that even with a firewall in place and SSL authentication set to required, if SMM is given access to the Kafka Connect API, then any user that has access to SMM will be able to interact with the Kafka Connect API. This is due to SMM not enforcing authorization checks when users are accessing Kafka Connect functionality within SMM. This is true for both the SMM UI and SMM REST API. As a result, caution is advised even if the Kafka Connect API itself is secured.

Complete the following steps to set SSL Client Authentication to required.

Procedure

1. Select the Kafka Service.
2. Go to Configuration.
3. Find the SSL Client Authentication property.
4. Set the property to required.



Important:

The SSL Client Authentication property is available for both Broker and Connect roles. Make sure that you are configuring the property for the Kafka Connect role or for the role group that includes your Kafka Connect roles.

5. Click Save Changes.
6. Restart the service.

Results

Only authenticated clients are allowed to connect to the Kafka Connect API.

What to do next

If you are using SMM to manage and monitor Kafka Connect, and you are not using auto TLS, add SMM's certificate to the Kafka Connect truststore.

Configuring SPNEGO Authentication and trusted proxies for the Kafka Connect REST API

Learn how you enable SPNEGO authentication and configure trusted proxies for the Kafka Connect REST API.

About this task



Note: SPNEGO authentication is automatically enabled if Kerberos is enabled for the Kafka service. As a result, SPNEGO authentication might already be enabled. However, if you upgraded your cluster from a version that did not have SPNEGO authentication available to a version that includes SPNEGO authentication, SPNEGO is not automatically enabled.

You can secure the Kafka Connect REST API by enabling SPNEGO authentication. This can be done with the Enable SPNEGO Authentication For Kafka Connect property in Cloudera Manager. If SPNEGO authentication is enabled, only users authenticated with Kerberos are able to access and use the REST API. Additionally, if Ranger authorization is enabled for the Kafka service, authenticated users are only able perform the operations that they are authorized for. If Ranger is not enabled, by default all authenticated users are able to perform all operations. Because users authenticate using Kerberos, securing the REST API using SPNEGO requires that Kerberos is enabled for the Kafka service.

In addition, the authentication framework supports trusted proxies. If the authenticated principal is recognized as a trusted proxy principal, then the authenticator accepts the doAs query parameter in the REST URL to specify the acting user. For example:

```
https://[***KAFKA_CONNECT_HOST***]:28085/connector-permissions?doAs=systest
```

In this example, systest is specified as the acting user. The doAs parameter is only accepted if the authenticated principal is recognized as a trusted proxy. By default, the Knox and SMM service principals (specifically, the short names of their Kerberos service principals) are recognized as trusted proxies. Trusted proxies can be configured with the List Of Trusted Proxy Services Cloudera Manager property.

Before you begin

Ensure that Kerberos is enabled for the Kafka service. For more information, see [Enable Kerberos authentication](#).

Procedure

1. In Cloudera Manager select the Kafka service.
2. Go to Configuration.
3. Find and enable the Enable SPNEGO Authentication For Kafka Connect property.
4. Find and configure the List Of Trusted Proxy Services property.
Add the Kerberos principal short names of the services that you want to act as trusted proxies.
5. Click Save Changes.
6. Restart the Kafka service.

Kafka Connect Rest API authorization

Learn about the authorization model shipped with Kafka Connect as well as how the model is integrated in Ranger.

Kafka Connect authorization model

Learn about the authorization model shipped with Kafka Connect.

Kafka Connect comes with an authorization model. Implementations are pluggable and it is up to the implementation how the capabilities of the model are utilized. The authorization model is implemented by default in Ranger. As a result, Cloudera recommends that you use Ranger.

Resource types

Cluster

Cluster is a singleton resource representing the cluster. This resource has the following operations:

- View: view connector plugins and the root path (basic cluster information).
- Validate: validate connector plugin configuration.
- Manage: manage cluster level resources such as loggers.

Connector

Connector is the resource representing a connector and its tasks. This resource has the following operations:

- View: view connector and task configuration as well as status.
- Manage: pause, resume, and restart connectors.
- Create: create connectors.
- Edit: edit connector configuration.
- Delete: delete connectors.

Error codes

When enabled, the Kafka Connect authentication and authorization framework can return the following HTTP errors in case of an authentication or authorization failure:

401 Unauthorized

The user attempts to connect the Kafka Connect REST server without authentication.

403 Forbidden

The user attempts an operation on a resource that is not permitted, but the user is entitled to view the resource. For example, the user tries to restart a connector that they only have view access to.

404 Not Found

The user attempts an operation on a resource that is not permitted and the resource is not visible by the user.

API endpoint to report user rights on a connector

The Kafka Connect version shipped with Cloudera Runtime includes the `/connector-permissions` endpoint. This endpoint is unique to Kafka Connect in Cloudera Runtime. It returns all connectors the user has permission to view along with access rights for each connector. The purpose of the endpoint is to allow the Streams Messaging Manager (SMM) UI to enable and disable UI controls based on what the user can do with each connector. This endpoint is not meant to check access rights in the client, these are always checked on the backend. The endpoint responds in the following format:

```
{ "createAllowed": true,
  "permissions": {
    "connectorName-1": { "editable": true, "deletable": true, "manageable": true },
    "connectorName-2": { "editable": true, "deletable": true, "manageable": true }
  }
}
```

Semantics for listing type endpoints

Listing endpoints (`/connectors` and `/connector-permissions`) only return the connectors that the user is permitted to view. In case the authenticated user has no rights to view any connectors it returns an empty list. Similarly to other endpoints, it responds with 401 Unauthorized if there is no authenticated user.

Kafka Connect Ranger integration

Learn about how the Kafka Connect authorization model is integrated with Ranger.

If Ranger authorization is enabled for the Kafka service and you have Kerberos and SPNEGO authentication enabled for Kafka Connect, you have the ability to set up fine grained access to the Kafka Connect REST API using the Ranger console.

This means that you can create policies and limit what operations an authenticated user is permitted to carry out. Policies that are set up are valid for any user that accesses the Kafka Connect REST API. This includes direct REST API access as well as access through Streams Messaging Manager (SMM).

If a user tries to carry out an operation on a resource that they do not have permission for, the REST API returns an error. For more information on the specific error messages, see *Kafka Connect authorization model*. In SMM, the policies control what UI operations are available for the user. For example, if a user does not have access rights to view select connectors, those connectors are not shown when logged in to the SMM UI. Similarly, if a user for example does not have the rights to create connectors, the New Connector option is not be available for the user.

Related Information

[Kafka Connect authorization model](#)

[Configuring SPNEGO Authentication and trusted proxies for the Kafka Connect REST API](#)

Kafka Connect connector configuration security

Learn about the options you have for hardening the security of connector configurations.

Kafka Connect Secrets Storage

Learn about the Kafka Connect Secrets Storage feature which can be used to hide and securely store sensitive data found in connector configurations.

In most production environments, the Kafka Connect connectors that you deploy connect to and move data either from or into secured services. Because of this, connectors that you deploy might require various credentials to access services. Credentials like passwords, access keys, or any other sensitive information must be provided to the connectors in their configuration. By default, the connector configurations can be easily read, their contents are unencrypted. Additionally, if you are using Streams Messaging Manager (SMM) to manage and deploy connectors, all sensitive information might be visible to anyone who has access to SMM with appropriate authorization rights.

Kafka Connect Secrets Storage is a feature within Kafka Connect that makes it possible for you to hide any sensitive values found in connector configurations and to store these values in a secure manner. For example, assume that you want to deploy a connector that connects to a secured service that requires a username and password for access. By default the password could be easily read by anyone with access to the connector configuration as the configuration is stored in plaintext. This poses a security risk. To circumvent this, you can mark the password as a secret. Once you mark the password as a secret, the Kafka Connect Secret Storage feature ensures that the password is stored and handled in a secure manner by doing the following:

- The property is replaced by a reference in the configuration. The reference resolves to the actual password at runtime. That is, the value is no longer readable if the configuration is accessed.
- The actual value of the password is encrypted and stored in an internal Kafka topic used by the feature.
- The value gets hidden (locked) on the SMM UI.

Kafka Connect Secrets Storage terms and concepts

Review the list of terms and concepts related to the Kafka Connect Secrets Storage feature.

Secure Connector

A secure connector is any connector intended to be used in CDP Kafka Connect whose configuration has a sensitive property.

Secret

Part of a secure connector configuration that must be protected. A secret can be any value used in the connector configuration that you mark as confidential data. For example, a secret can be a password, a JAAS credential, an SSL certificate, and so on. When you mark a property as a secret it becomes protected. Specifically, the value of the property is no longer stored in plaintext within the

configuration. Rather, it is replaced by a secret reference. The actual value is encrypted and stored in an internal Kafka topic (secrets topic).

Understand however, that this protection covers the configuration only. If you mark a property a secret that is otherwise public information, while the value of that property will be hidden in the configuration, the value could be looked up by other means. For example, sink connectors have a topics property. If auto topic creation is enabled, topics will be created by the connector automatically. This means that even if the topics property is marked as a secret, the value of the property could be retrieved by other means. For example, by listing the topics of the Kafka service.

Secrets topic

The secrets topic is an internal Kafka topic. The Kafka Connect Secrets Storage feature stores its internal state, configuration, secrets, and so on in this internal topic.

Secret reference

Part of a secure connector configuration that is already protected. Instead of the configuration's original value, a placeholder (reference) is used, which refers to the original value using a connector name and a bundle ID. The Kafka Connect Secrets Storage feature resolves these references and substitutes them with the original values stored in the secrets topic when connectors are used.

Encryption key

The encryption key is a randomly generated key. Its established value is unique at a certain point in time. However, it can be changed. It is used to encrypt and decrypt hashed secrets protecting them against tampering and retrieval.

Global key

Global key is a unique cryptographic key that is derived from a user supplied password. The global key protects the encryption key when the encryption key is persisted to the disk or exchanged on the network. The encryption key is hashed and encrypted with the global key. As a result, it is impossible to tamper with it, or retrieve the encryption key from wiretapping the network traffic in a man-in-the-middle style attack.

Bundle and Bundle ID

The set of secrets of a secure connector constitute a bundle. The bundle is identified by its bundle ID, which is maintained by the Kafka Connect Secrets Storage feature. If a secret is added, edited, or deleted, a new bundle is created, with a new bundle ID.

Managing secrets using the REST API

Learn how you can use the Kafka Connect Secrets Storage feature to manage secrets in connector configurations using the Kafka Connect REST API.

You can use the Kafka Connect Secrets Storage feature with the Kafka Connect REST API. The following sections walk you through how you can mark properties as secrets and how you can manage already existing secrets.



Important: The Kafka Connect Secrets Storage feature can also be used when managing connectors using the Streams Messaging Manager (SMM) UI. For more information, see [Deploying and managing Kafka Connect connectors in SMM](#).

Mark a configuration value as a secret in a connector configuration

To mark a property value in a connector configuration as a secret, you must:

- Add the property to the configuration.
- Add the secret.properties property. The value of this property is a comma separated list of the property keys that you want to mark as secrets.

For example, assume that you have two sensitive properties in your configuration that you want to mark as secrets, password.for.connector and some.confidential.config.

```
{  
  "name" = "Custom",
```

```

"connector.class" = "my.custom.secure.connector" ,
"secret.properties" = "password.for.connector,some.confidential.config" ,
"topics" = "...",
"password.for.connector" = "mypassword" ,
"some.confidential.config" = "sensitive information"
}

```

If the `secret.properties` is present in the configuration, the Kafka Connect Secrets Storage feature gets enabled. Once the connector configuration is submitted, the feature processes the configuration and makes a number of changes.

```

{
  "name" = "Custom" ,
  "connector.class" = "my.custom.secure.connector" ,
  "secret.properties" = "password.for.connector,some.confidential.config" ,
  "secret.bundle.id" = "18alb3a6-1e2f" ,
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/18alb3a6-1e2f:password.for.connector}" ,
  "some.confidential.config" = "${secret:Custom/18alb3a6-1e2f:some.confidential.config}"
}

```

Notice the following changes:

- A `secret.bundle.id` property is added. The bundle is the set of secrets of a secure connector. The bundle is identified by its bundle ID, which is maintained by the Kafka Connect Secrets Storage feature. If a secret is added, deleted, or edited, a new bundle is created with a new bundle ID.
- The `password.for.connector` and `some.confidential.config` properties had their values replaced with secret references. These references are placeholder values consisting of the connector name, bundle ID, and property key. If a reference like this is present in the configuration, it means that the value for that property is securely stored in an internal Kafka topic (secrets topic). The Kafka Connect Secrets Storage feature resolves these references and substitutes them with the original values stored in the secrets topic when the connectors are used.

Update a connector configuration with a new secret

If you want to update an already existing configuration with a new secret property, you must add the new property to the configuration as well as the `secret.properties` property. For example, assume that you have a configuration where `password.for.connector` and `some.confidential.config` are already marked as secrets. Now, you want to add the new `secret.config` property to the configuration and mark it as a secret as well.

```

{
  "name" = "Custom" ,
  "connector.class" = "my.custom.secure.connector" ,
  "secret.properties" = "password.for.connector,some.confidential.config,new.secret.config" ,
  "secret.bundle.id" = "18alb3a6-1e2f" ,
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/18alb3a6-1e2f:password.for.connector}" ,
  "some.confidential.config" = "${secret:Custom/18alb3a6-1e2f:some.confidential.config}" ,
  "new.secret.config" = "[***A NEW SECRET VALUE***]"
}

```

Once submitted and processed, the configuration will look like the following example:

```

{
  "name" = "Custom" ,
  "connector.class" = "my.custom.secure.connector" ,

```

```

"secret.properties" = "password.for.connector,some.confidential.config,new.secret.config",
"secret.bundle.id" = "c8a753d6-8242",
"topics" = "...",
"password.for.connector" = "${secret:Custom/c8a753d6-8242:password.for.connector}",
"some.confidential.config" = "${secret:Custom/c8a753d6-8242:some.confidential.config}",
"new.secret.config" = "${secret:Custom/c8a753d6-8242:new.secret.config}"
}

```

Notice the following changes:

- The `secret.bundle.id` value has changed.
This is done because adding a new secret creates a new set of secrets resulting in a new bundle being created.
- The values (secret references) of `password.for.connector` and `some.confidential.config` are updated with the new bundle ID.
- The value of `new.secret.config` is replaced by a secret reference.
This means that the value is now securely stored in the secrets topic.

Update a connector configuration by changing a secret

Updating already existing secret values can be done by updating the appropriate property values. For example, assume that you have three properties marked as secrets, `password.for.connector`, `some.confidential.config`, and `new.secret.config`. You want to update the value of `new.secret.config`.

```

{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.properties" = "password.for.connector,some.confidential.config,new.secret.config",
  "secret.bundle.id" = "c8a753d6-8242",
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/c8a753d6-8242:password.for.connector}",
  "some.confidential.config" = "${secret:Custom/c8a753d6-8242:some.confidential.config}",
  "new.secret.config" = "[***CHANGED VALUE***]"
}

```

Once submitted and processed, the configuration will look like the following example:

```

{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.properties" = "password.for.connector,some.confidential.config,new.secret.config",
  "secret.bundle.id" = "14b2f85f-1dcb",
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/14b2f85f-1dcb:password.for.connector}",
  "some.confidential.config" = "${secret:Custom/14b2f85f-1dcb:some.confidential.config}",
  "new.secret.config" = "${secret:Custom/14b2f85f-1dcb:new.secret.config}"
}

```

Notice the following changes:

- The `secret.bundle.id` value has changed.

This is done because deleting a secret creates a new set of secrets resulting in a new bundle being created. In this case the deleted secret is removed from the bundle.

- The values (secret references) of `password.for.connector` and `some.confidential.config` are updated with the new bundle ID.

This is done because updating a secret creates a new set of secrets resulting in a new bundle being created.

- The value of `new.secret.config` is replaced by a secret reference.

This means that the new value is now securely stored in the secrets topic.

Remove an existing secret from an existing configuration

Secrets can be removed from a configuration by updating the configuration and deleting the unnecessary properties from both the configuration and `secret.properties`. For example, assume you have three properties marked as secrets, `password.for.connector`, `some.confidential.config`, and `new.secret.config`.

```
{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.properties" = "password.for.connector,some.confidential.config,new.secret.config",
  "secret.bundle.id" = "14b2f85f-1dcb",
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/14b2f85f-1dcb:password.for.connector}",
  "some.confidential.config" = "${secret:Custom/14b2f85f-1dcb:some.confidential.config}",
  "new.secret.config" = "${secret:Custom/14b2f85f-1dcb:new.secret.config}"
}
```

Assume you want to delete `some.confidential.config`. To do this, you need to delete the property from the configuration and remove the property key from `secret.properties`.

```
{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.properties" = "password.for.connector,new.secret.config",
  "secret.bundle.id" = "14b2f85f-1dcb",
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/14b2f85f-1dcb:password.for.connector}",
  "new.secret.config" = "${secret:Custom/14b2f85f-1dcb:new.secret.config}"
}
```

Once submitted and processed, the configuration will look like the following example:

```
{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.properties" = "password.for.connector,new.secret.config",
  "secret.bundle.id" = "15641fca-483b",
  "topics" = "...",
  "password.for.connector" = "${secret:Custom/15641fca-483b:password.for.connector}",
  "new.secret.config" = "${secret:Custom/15641fca-483b:new.secret.config}"
}
```

Notice the following changes:

- The `secret.bundle.id` value has changed.

This is done because deleting a secret creates a new set of secrets resulting in a new bundle being created. In this case the deleted secret is removed from the bundle.

- The values (secret references) of `password.for.connector` and `some.confidential.config` are updated with the new bundle ID.

Remove all secrets

To remove all secrets from a configuration, you must remove all properties that were marked as secrets and also remove the `secret.properties` property. However, you must leave `secret.bundle.id` intact. For example, assume you had a single property marked as a secret, `password.for.connector`.

```
{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.properties" = "password.for.connector,new.secret.config",
  "secret.bundle.id" = "c8a753d6-8242",
  "topics" = "...",
  "password.for.connector"= "${secret:Custom/c8a753d6-8242:password.for.connector}
}
```

In a case like this, you need to remove `password.for.connector` and `secret.properties` from the configuration. However, you must preserve `secret.bundle.id`. The reason why you need to leave this property intact is because the encrypted values for the secret properties must be deleted from the internal secrets topic. The Kafka Connect Secrets Storage feature can only do this if it has a bundle ID for reference.

```
{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "secret.bundle.id" = "c8a753d6-8242",
  "topics" = "..."
}
```

Once the configuration is submitted and processed, `secret.bundle.id` is removed. The connector configuration at this point is considered unsecured.

```
{
  "name" = "Custom",
  "connector.class" = "my.custom.secure.connector",
  "topics" = "..."
}
```

Re-encrypting secrets

Using the Kafka Connect Secrets Storage feature places encrypted confidential information in an internal Kafka topic (secrets topic) used by the feature. If you want to change the encryption key used to encrypt the information placed in the topic, all secrets must be re-encrypted and must be migrated to a new topic. This is done using the `connect-secret-storage-migration` tool that is shipped with Cloudera Runtime.

About this task

The `connect-secret-storage-migration` tool migrates existing secrets from the currently used secrets topic (source) to a newly created one (target). It achieves this by reading through the source topic, decrypting the secrets stored within it using the existing global key, and re-encrypts them with a new, randomly generated encryption key and a new global key that you configure. After encryption is finished, newly encrypted secrets are produced into the target topic. The new encryption key will have a different security context (different global password and salt) than the currently used one.

The tool requires three configuration files to function. These configuration files provide the tool with information regarding the Kafka service as well as the source and target topics. The configuration files must be created manually.

Before you begin

In the configuration files you prepare, you must specify the current global password and PBE salt used by the feature. These values are specified in Cloudera Manager with the Kafka Connect Secrets Storage PBE Salt and Kafka Connect Secrets Storage Global Password properties. The default values for these properties are randomly generated and are hidden, but can be changed at any time. As a result, if you do not know the global password or PBE salt value, you must change both values, note down the values you configured, and restart all Kafka Connect roles before you complete the following steps.

Procedure

1. Stop all Kafka Connect roles.
2. Prepare the configuration files required for the tool.

You need to prepare three different configuration files. A Kafka client configuration file that the tool can use to access the Kafka service. Additionally, two configuration files that contain information regarding the source and target topics. All three files must be deployed on the cluster host where you run the tool.

- a) Create a Kafka client configuration file.

This configuration file must contain all client properties that the tool requires to access the Kafka service. The properties you need to specify in this file depend on the security configuration of the Kafka service. The following tabs contain examples for some of the most commonly used security configurations.

For SASL_SSL Kerberos

```
bootstrap.servers=[***HOST***]:[***PORT***]
security.protocol=SASL_SSL
sasl.mechanism=GSSAPI
sasl.kerberos.service.name=kafka
sasl.jaas.config=com.sun.security.auth.module.Krb5LoginModule required
useKeyTab=true storeKey=true keyTab="[***PATH TO KEYTAB***]" principal="[***KERBEROS PRINCIPAL***]";
ssl.truststore.location=[***PATH TO TRUSTSTORE.JKS***]
ssl.truststore.password=[***TRUSTSTORE PASSWORD***]
```

For SASL_SSL PLAIN

```
bootstrap.servers=[***HOST***]:[***PORT***]
security.protocol=SASL_SSL
sasl.mechanism=PLAIN
sasl.jaas.config=org.apache.kafka.common.security.plain.PlainLoginModule required username="[***USERNAME***]" password="[***PASSWORD***]";
ssl.truststore.location=[***PATH TO TRUSTSTORE.JKS***]
ssl.truststore.password=[***TRUSTSTORE PASSWORD***]
```

- b) Create the configuration file identifying the source topic and its security context. For example:

```
kafka.connect.secret.storage.topic : [***SOURCE TOPIC NAME***]
kafka.connect.secret.global.password : [***CURRENT GLOBAL PASSWORD***]
```

```
kafka.connect.secret.pbe.salt : [***CURRENT PBE SALT***]
```

Replace [***SOURCE TOPIC NAME***] with the name of the currently used secrets topic. The name of the current topic is specified in Cloudera Manager in the Kafka Connect Secrets Storage Topic Name property.

Replace [***GLOBAL PASSWORD***] with the currently used global password.

Replace [***PBE SALT***] with the currently used PBE salt value.

You can choose to not specify `kafka.connect.secret.global.password` in the configuration. If the property is not added to the configuration, you are prompted to enter the password when you run the tool.

- c) Create the configuration file identifying the target topic and its security context. For example:

```
kafka.connect.secret.storage.topic : [***TARGET TOPIC NAME***]
kafka.connect.secret.global.password : [***NEW GLOBAL PASSWORD***]
kafka.connect.secret.pbe.salt : [***NEW PBE SALT***]
kafka.connect.secret.storage.topic.replication.factor : [***REPLICATION
  FACTOR***]
kafka.connect.secret.storage.topic.configs.min.insync.replicas : [***IN-
  SYNC REPLICICA NUMBER***]
```

You can choose to not specify `kafka.connect.secret.global.password` in the configuration. If the property is not added to the configuration, you are prompted to enter the password when you run the tool.

3. Run the tool.

You must specify the three configuration files you created as parameters. The order in which you specify the files is fixed. The first file must be the Kafka client configuration, the second must be the file for the source topic, and the third file must be the file for the target topic. For example:

```
connect-secret-storage-migration [***KAFKA CLIENT
  CONFIGURATION***] [***SOURCE TOPIC CONFIGURATION***] [***TARGET TOPIC
  CONFIGURAITON***]
```

4. If prompted, enter the current and new global password.

The tool only prompts you to enter the passwords if they were not specified in the configuration files. The tool does not echo the passwords that you type in the console.

5. The tool prompts you to ensure that the Kafka Connect cluster (all Connect roles) are stopped. Press Enter to continue.

6. Wait until the tool is finished with re-encryption.

The tool reports its progress. Upon successful completion, the tool exits. Upon failure, an error message is printed. In such a case, fix any issues that the tool reports and run the tool again.

7. In Cloudera Manager, select the Kafka service.

8. Go to Configuration.

9. Find and configure the following properties:

- Kafka Connect Secrets Storage Topic Name
- Kafka Connect Secrets Storage Topic Replication Factor
- Kafka Connect Secrets Storage Topic Minimum In-Sync Replicas
- Kafka Connect Secrets Storage Global Password
- Kafka Connect Secrets Storage PBE Salt

Add the values you specified in the configuration file identifying the target topic.

10. Restart all Kafka Connect roles.

11. Verify that Kafka Connect is working as expected.

This can be done by verifying that all connectors that have properties marked as secrets are running and are in a healthy state. That is, the connectors are resolving all referenced secrets correctly. If you are experiencing issues, you can revert the configuration changes in Cloudera Manager and restart the re-encryption process.

12. Delete the configuration files you created.
13. Delete the source secrets topic.



Important: Ensure that you verified that Kafka Connect and your connectors are working as expected before deleting the topic. When you delete the topic, the secrets stored within it are lost. As a result, if there are issues with Kafka Connect, you might not be able to revert the cluster to a working state.

Configuring connector JAAS configuration and Kerberos principal overrides

Learn how to configure Kafka Connect connectors to override the default JAAS configuration and Kerberos principal that they use to connect to the Kafka service. Configuring these overrides enables you to set up fine grained access rules for the connectors that you deploy.

About this task

By default Kafka Connect connectors, specifically the clients managed by Connect, use the JAAS configuration and Kerberos principal of the Connect worker (Kafka Connect role) to establish a connection with the Kafka service. As a result, all connectors that you deploy authenticate with the same principal. Additionally, if authorization is enabled, authorization also happens using that same principal. Because the connectors use the same principal as the worker and because the worker has all permission on the Kafka service, all connectors have full access to Kafka by default.

This behavior might not be desirable and you might want to configure fine grained access rules for individual connectors that you deploy. For example, you might want to have your connectors to be able to manage specific Kafka resources only, but not have access to other resources. To achieve this, you must configure connectors to be required to specify unique connection credentials for the Kafka service connection. Configuration is done in Cloudera Manager using Connector Kafka Client Configuration Override Policy and Require Connectors To Override Kafka Client JAAS Configuration Kafka service properties.

Procedure

1. In Cloudera Manager select the Kafka service.
2. Go to Configuration.
3. Find and configure the following properties.

Table 4:

Property Name	Description
Connector Kafka Client Configuration Override Policy	<p>Specifies what client configuration can be overridden by the connector. Its value is the class name or alias of the ConnectorClientConfigOverridePolicy implementation. The framework includes three implementations (policies). These are None, All, and Principal.</p> <ul style="list-style-type: none"> • If set to None, configurations cannot be overridden. This is the default setting. • If set to Principal, SASL configurations can be overridden • If set to All, all configurations can be overridden
Require Connectors To Override Kafka Client JAAS Configuration	<p>Specifies whether Connector configurations should override the sasl .jaas.config property of the Kafka clients used by the connector. If set to true, connector configurations must override the client's sasl .jaas.config property. The override must also not match the Kerberos credentials of the Connect worker. Additionally, if set to true, a UnionCompositePolicy is installed under the connector.client.config.override.policy configuration. If Connector Kafka Client Configuration Override Policy is set to None, it is upgraded to Principal. Otherwise, the existing policy is also applied by the composite policy.</p>

4. Click Save Changes.
5. Restart the Kafka service.

What to do next

- Update the configuration of all existing connectors to include the `sasl.jaas.config` property with a valid JAAS configuration.
- Ensure that the principal specified for the connector has appropriate permissions assigned to it in your authorization service.

Configuring a Nexus repository allow list

To harden the security of your Kafka Connect deployment, you can configure a Nexus repository allow list. If an allow list is configured, Stateless NiFi Source and Sink connectors are only able to fetch NiFi extensions from the specified list of Nexus repositories.

About this task

NiFi dataflows within Kafka Connect are deployed using the Stateless NiFi Source and Sink connectors. When you deploy your custom developed dataflows with the Stateless NiFi connectors, you can configure the `nexus.url` property. This property specifies the URL of a Nexus repository that hosts the NiFi extensions that are not bundled with the connectors but are required by the dataflow that is run within the connector.

By default connectors you deploy are allowed to connect to any Nexus repository that you specify. However, this may pose a security risk. Because of this, you can configure Kafka Connect to only allow its connectors to connect to specific Nexus repositories. That is, you can configure a Nexus server allow list. If an allow list is configured, connectors are only allowed to connect and fetch NiFi extensions from the servers that are included in the list.

The allow list is configured in Cloudera Manager using the List Of Allowed Nexus Repository Urls property.

Procedure

1. In Cloudera Manager, select the Kafka service.
2. Go to Configuration.
3. Find and configure the List Of Allowed Nexus Repository Urls property.
Add the URLs of the Nexus repositories that you want to allow your connectors to fetch extensions from.
4. Click Save Changes.
5. Restart the Kafka service.

Results

The Nexus URLs specified in the `nexus.url` property of the Stateless NiFi Source and Sink connectors are validated against the list of configured URLs. Connections are only allowed to the repositories configured in the list.

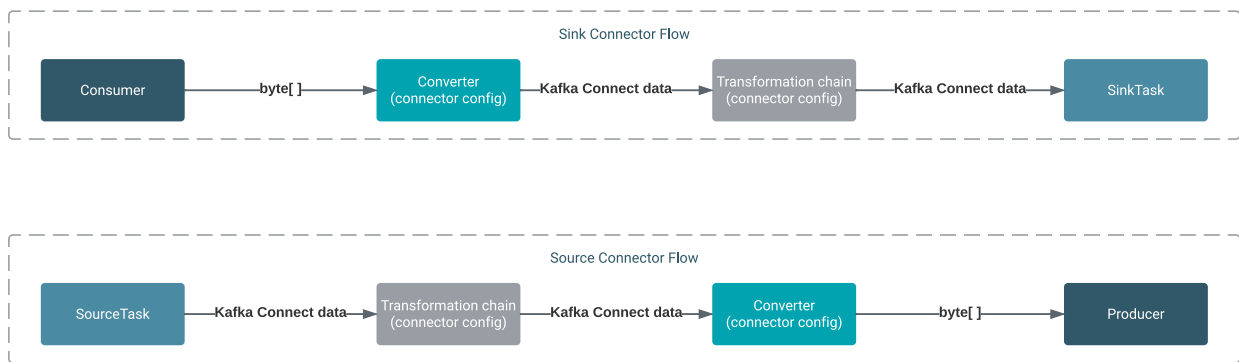
Single Message Transforms

Single Message Transforms (SMT) is a message transformation framework that you can deploy on top of Kafka Connect connectors to apply message transformations and filtering. Learn about the SMT framework as well as the transformation plugins available in CDP.

Kafka Connect connectors provide ready-to-use tools to integrate between Kafka and external data stores. Still, in many use cases, the data moved by the connectors require some sanitization and transformation. To provide extra flexibility built on top of connectors, Kafka Connect also supports an SMT framework.

The SMT framework installs a transformation chain on top of connectors that modifies and filters the data on a single message basis. An SMT chain consists of transform and predicate plugins. Transform plugins are used to modify the data. For example, you can insert, replace, mask as well as perform various other modifications on the messages moved by connectors. Predicate plugins are used to add additional logic to your chain so that the transformation chain is only applied to messages which satisfy specified conditions.

The SMT framework requires that data is converted to the Kafka Connect internal data format. This data format is specific to Kafka Connect and consists of a structure and schema descriptor (SchemaAndValue) specific to Connect.



Supported SMT plugins

Kafka in CDP ships with and supports the following SMT plugins:

- All Apache Kafka plugins. For more information, see [Transformations](#) in the Apache Kafka documentation.
- All Debezium plugins. For more information, see [Transformations](#) in the Debezium documentation.
- The following Cloudera specific plugins:
 - [ConvertFromBytes](#) on page 36
 - [ConvertToBytes](#) on page 38



Note: The ConvertFromBytes and ConvertToBytes transformation plugins transform binary data to and from the Kafka Connect internal data format. These plugins are specifically developed to enable the use of the SMT framework with connectors that only support binary data.

Configuring an SMT chain

Learn how to configure a Single Message Transformation (SMT) chain for Kafka Connect connectors.

SMT chains can be configured within the configuration of a Kafka Connect connector using SMT specific configuration properties. To set up a chain, you first define your transformation chain with the `transforms` property and optionally define your predicates using the `predicates` property. Afterward, you use `transforms.*` and `predicates.*` to configure the plugins in the chain. For example, the following configuration snippet sets up a transformation chain that filters messages based on their header and removes a specified field from messages.

```

transforms=FilterAudit,MaskField,

transforms.MaskField.type=org.apache.kafka.connect.transforms.MaskField$Value
transforms.MaskField.fields=CreditCardNumber

transforms.FilterAudit.type=org.apache.kafka.connect.transforms.Filter
transforms.FilterAudit.predicate=IsAudit
transforms.FilterAudit.negate=false

predicates=IsAudit
predicates.IsAudit.type=org.apache.kafka.connect.transforms.predicates.HasHeaderKey
predicates.IsAudit.name=Audit
  
```

The following sections go through the properties in this example and give an overview on how to set up a transformation chain.

Configuring transforms

The transforms property contains a comma-separated list of transformation aliases. Each alias represents one step in the transformation chain. The aliases you add to the property are arbitrary names, they are used in other properties to configure that particular transformation step. For example, the following defines a two step transformation chain.

```
transforms=FilterAudit,MaskField
```

The transforms.[***ALIAS***].type property specifies which transformation plugin should be used in a transformation step. [***ALIAS***] is one of the aliases that you specified in transforms. The value of the property is the fully qualified name of the transformation plugin that should be used in the step. For example, the following line specifies org.apache.kafka.connect.transforms.MaskField\$Value as the plugin for the MaskField step.

```
transforms.MaskField.type=org.apache.kafka.connect.transforms.MaskField$Value
```

Many transformation plugins support changing both the key and the value of a record. For these plugins, typically, a nested value or key class can be referenced as the type.

The transforms.[***ALIAS***].[***KEY***] property is used to configure the transformation plugins in your chain. This property is passed to the transformation plugin itself with transforms.[***ALIAS***] stripped from the property key. [***ALIAS***] is the alias of a plugin you specified in transforms. [***KEY***] is a property key that the plugin accepts. For example, the MaskField plugin has a fields property that specifies which fields should be removed from the structure.

```
transforms.MaskField.fields=CreditCardNumber
```

Configuring predicates

Predicates are a separate set of plugins. You use them to conditionally enable certain steps in the transformation chain. Predicates are configured in a similar way to transforms. You must specify the predicate aliases, associate the aliases with a plugin, and set plugin specific properties.

```
predicates=IsAudit
predicates.IsAudit.type=org.apache.kafka.connect.transforms.predicates.HasHeaderKey
predicates.IsAudit.name=Audit
```

In this example the IsAudit predicate is an instance of the HasHeaderKey predicate plugin. This predicate returns true for records where a specific header key is present. predicates.IsAudit.name=Audit configures the predicate to look for the Audit header in the records.

After a predicate is set up, you can associate the predicate with any number of transformation steps using the predicate property. If a predicate is associated with a transformation, that transformation step is only applied to the messages that match the condition specified in the predicate.

A good example for using a predicate is the Filter transformation plugin. This is because Filter filters (drops) all messages by default. Therefore, it must be used together with predicates to specify filtering logic. For example, the following configuration instructs the SMT framework that the FilterAudit step should only be invoked for messages where the IsAudit predicate returns true. That is, all messages with the Audit header will be filtered and will not be transformed by any subsequent steps in the transformation chain.

```
transforms.FilterAudit.predicate=IsAudit
transforms.FilterAudit.negate=false
```

The condition of a predicate can be inverted using negate. If negate is set to true, the SMT framework applies the transformation to any record that does **not** match the condition. For example, the following configuration instructs the

SMT framework that the FilterAudit step should only be invoked for messages where the IsAudit predicate returns false.

```
transforms.FilterAudit.predicate=IsAudit
transforms.FilterAudit.negate=true
```

ConvertFromBytes

ConvertFromBytes is a Cloudera specific transformation plugin that converts binary data to the Kafka Connect internal data format. You can use this plugin to make connectors that only support binary data compatible with the Single Message Transforms (SMT) framework.

Fully qualified names

- `com.cloudera.dim.kafka.connect.transformations.convert.ConvertFromBytes$Key`
- `com.cloudera.dim.kafka.connect.transformations.convert.ConvertFromBytes$Value`

Description



Important: Ensure that you have an in-depth understanding about the following aspects of the connector that you plan on using with ConvertFromBytes.

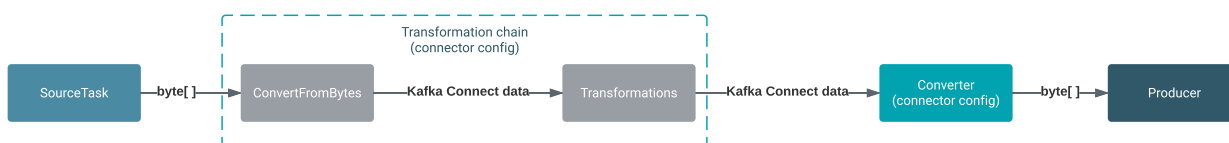
- The type of the connector
- The converter type used by the connector
- The header converter type used by the connector

These aspects of a connector heavily influence how the plugin and how the transformation chain that includes this plugin must be configured.

The ConvertFromBytes transformation plugin accepts binary data and converts it into the Kafka Connect internal data format with a nested converter that transforms binary data. To support header based converter logic, for example when using the AvroConverter with the schema ID encoded in the record header, the plugin requires a header converter to correctly transform record headers when interacting with the converter. This plugin supports both key and value conversion.

Using this plugin with connectors that only support binary data makes the connector fully compatible with the SMT framework. On their own, connectors that only support binary data have limited compatibility with transformations even if the binary data is structured. This is because transformations are only fully supported on data that is in the Kafka Connect internal data format. Binary only connectors like the Mirror Source (MirrorSourceConnector) or the Stateless NiFi Source connectors emit data that has the BYTES schema and do not provide conversion to the Kafka Connect internal data format by default. When you use a connector such as these with the ConvertFromBytes plugin, the binary data is parsed into a compatible structure, which can then be further processed with the transformation chain.

Figure 1: Source connector example flow with ConvertFromBytes



Example

The following configuration example adds a ConvertFromBytes transformation as a first step of the transformation chain. The conversion uses a schemaless JSON transformation to parse the binary data. The transformation steps, the

connector, or the converter, whichever comes directly after FromBytes, receives a properly structured record instead of binary data.

```
{
  "transforms" : "FromBytes,...",
  "transforms.FromBytes.type" : "com.cloudera.dim.kafka.connect.transformations.convert.ConvertFromBytes$Value",
  "transforms.FromBytes.converter" : "org.apache.kafka.connect.json.JsonConverter",
  "transforms.FromBytes.converter.schemas.enable" : "false"
}
```

Configuration properties

Table 5: ConvertFromBytes properties reference

Property	Default Value	Required	Description
converter		True	The fully qualified name of the converter implementation to use. For example: com.cloudera.dim.kafka.connect.converts.AvroConverter
header.converter	org.apache.kafka.connect.storage.SimpleHeaderConverter	True	The fully qualified name of the header converter implementation to use. This converter must match the header converter of the connector.
converter.		False	A configuration prefix. Use this prefix to configure the properties of the converter specified in converter. Property keys and values specified with the prefix are passed directly to the converter with the prefix stripped. For example: <pre>transformer.[***TRANSFORM ALIAS***].converter.[***CONVERTER PROPERTY KEY***]:[***CONVERTER PROPERTY VALUE***]</pre>

Property	Default Value	Required	Description
header.converter.		False	<p>A configuration prefix. Use this prefix to configure the properties of the header converter specified in header.converter. Property keys and values specified with the prefix are passed directly to the header converter with the prefix stripped. For example:</p> <pre> transform ms.[***TRANSFORM ALIAS***].con vert er.[***HEADER CONVERTER PROPERTY KEY***]:[***HEADER CONVERTER PROPERTY VALUE***] </pre>

ConvertToBytes

ConvertToBytes is a Cloudera specific transformation plugin that converts Kafka Connect internal data to binary data. You can use this plugin to make connectors that only support binary data compatible with the Single Message Transforms (SMT) framework.

Fully qualified names

- com.cloudera.dim.kafka.connect.transformations.convert.ConvertToBytes\$Key
- com.cloudera.dim.kafka.connect.transformations.convert.ConvertToBytes\$Value

Description



Important: Ensure that you have an in-depth understanding about the following aspects of the connector that you plan on using with ConvertToBytes.

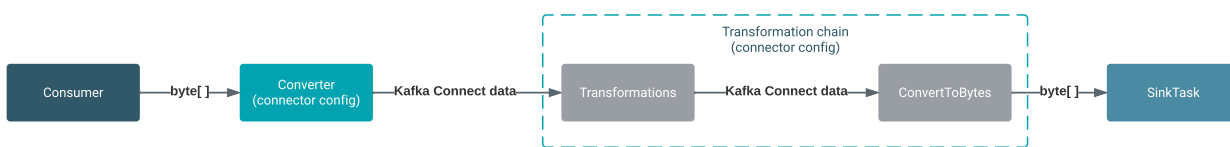
- The type of the connector
- The converter type used by the connector
- The header converter type used by the connector

These aspects of a connector heavily influence how the plugin and how the transformation chain that includes this plugin must be configured.

The ConvertToBytes transformation plugin accepts data in the Kafka Connect internal data format and converts it to binary data with a nested converter. To support header based converter logic, for example when using the AvroConverter with the schema ID encoded in the record header, the plugin requires a header converter to correctly transform record headers when interacting with the converter. This plugin supports both key and value conversion.

Using this plugin with connectors that only support binary data makes the connector fully compatible with the SMT framework. On their own, connectors that only support binary data, for example the Stateless NiFi Sink connector, have limited compatibility with transformations even if the binary data is structured. This is because the format of the data after transformations are carried out is normally the Kafka Connect internal data format. Binary only connectors, however, expect data that has the BYTES schema and do not provide conversion from the Kafka Connect internal data format by default. When you use the ConvertToBytes plugin with a binary only connector, the structured data is converted to binary format, which can then be picked up by the connector.

Figure 2: Source connector example flow with ConvertToBytes



Example

The following configuration example adds a `ConvertToBytes` transformation as the last step of the transformation chain. The conversion uses a schemaless JSON transformation to serialize the structured data. The transformation steps, the connector, or the converter, whichever comes directly after `ToBytes`, receives a properly structured record instead of binary data.

```
{
  "transforms" : "... ,ToBytes",
  "transforms.ToBytes.type" : "com.cloudera.dim.kafka.connect.transfor
mations.convert.ConvertToBytes$Value",
  "transforms.ToBytes.converter" : "org.apache.kafka.connect.json.Json
Converter",
  "transforms.ToBytes.converter.schemas.enable" : "false"
}
```

Configuration properties

Table 6: `ConvertToBytes` properties reference

Property	Default Value	Required	Description
<code>converter</code>		True	The fully qualified name of the converter implementation to use. For example: <code>com.cloudera.dim.kafka.connect.converts.AvroConverter</code>
<code>header.converter</code>	<code>org.apache.kafka.connect.storage.SimpleHeaderConverter</code>	True	The fully qualified name of the header converter implementation to use. This converter must match the header converter of the connector.
<code>converter.</code>		False	A configuration prefix. Use this prefix to configure the properties of the converter specified in <code>converter</code> . Property keys and values specified with the prefix are passed directly to the converter with the prefix stripped. For example: <div style="border: 1px solid #ccc; padding: 5px; margin-top: 10px;"> <pre>transfor ms.[***TRANSFORM ALIAS***].con vert er.[***CONVERTER PROPERTY KEY***]:[***CONVERTER PROPERTY VALUE***]</pre> </div>

Property	Default Value	Required	Description
header.converter.		False	<p>A configuration prefix. Use this prefix to configure the properties of the header converter specified in header.converter. Property keys and values specified with the prefix are passed directly to the header converter with the prefix stripped. For example:</p> <pre>transform.ms.[***TRANSFORM ALIAS***].converter.[***HEADER CONVERTER PROPERTY KEY***]:[***HEADER CONVERTER PROPERTY VALUE***]</pre>

Connectors

Learn what Kafka Connect connectors are shipped with Cloudera Runtime.

Cloudera Runtime comes prepackaged with a number of Cloudera developed Kafka Connect connectors. In addition, connectors that come packaged with the version of Apache Kafka that is included in Cloudera Runtime are also available for use. Manually installing and using your own custom connectors is also possible. The following collects the connectors shipped in Cloudera Runtime.

Debezium connectors

Debezium connectors capture changes from a wide variety of databases and produce the captured changes into Kafka. Using Debezium connectors makes it possible for your applications to consume and respond to change events regardless of where the changes originated from.

Table 7: Debezium connectors

Connector	Description
Debezium Db2 Source [Technical Preview]	The Debezium Db2 Source connector captures row-level changes in Db2 database tables and transfers the changes to Kafka topics.
Debezium MySQL Source	The Debezium MySQL Source connector reads the binary log (binlog) of a MySQL server, produces change events for row-level INSERT, UPDATE, and DELETE operations, and transfers the changes to Kafka topics.
Debezium Oracle Source	The Debezium Oracle Source connector captures and records row-level changes that occur in databases on an Oracle server, including tables that are added while the connector is running, and transfers the changes to Kafka topics.
Debezium PostgreSQL Source	The Debezium PostgreSQL Source connector captures row-level INSERT, UPDATE, and DELETE operations, produces change events for each change, and transfers the changes to Kafka topics.

Connector	Description
Debezium SQL Server Source	The Debezium SQL Server Source connector captures row-level INSERT, UPDATE, and DELETE operations that occur in the schemas of a SQL Server database, produces change events for each change, and transfer the changes to Kafka topics.



Note: Cloudera does not provide documentation on how to use or configure Debezium connectors. Documentation only covers the setup steps that are required before you can deploy the connectors in a CDP cluster. To learn more about these connectors, see the appropriate version of the [Debezium Documentation](#).

Stateless NiFi connectors

The Stateless NiFi Source and Sink connectors allow you to run NiFi dataflows within Kafka Connect. Using these connectors can grant you access to a number of NiFi features without having the need to deploy or maintain NiFi on your cluster.

Stateless NiFi connectors fall in to two categories. You have the base Stateless NiFi Source (`StatelessNiFiSource`) and Stateless NiFi Sink (`StatelessNiFiSink`) connectors. In addition, there are a number of ready-to-use connectors based on Stateless NiFi Source and Sink. These Stateless NiFi-based connectors run predefined dataflows developed by Cloudera and cover common data movement use cases.

Table 8: Base connectors

Connector	Description
Stateless NiFi Source	The Stateless NiFi Source and Sink connectors allow you to run NiFi dataflows within Kafka Connect. Using these connectors can grant you access to a number of NiFi features without having the need to deploy or maintain NiFi on your cluster.
Stateless NiFi Sink	

Table 9: Predefined flows/Stateless NiFi-based connectors

Connector	Description
HTTP Source	The HTTP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The HTTP Source connector listens on a port for HTTP POST requests and transfers the request body to a Kafka topic.
JDBC Source	The JDBC Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The JDBC Source connector reads records from a database table and transfers each record to Kafka in Avro or JSON format.
JMS Source	The JMS Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The JMS Source Connector consumes messages from a JMS broker and transfers the message body to Kafka.
MQTT Source	The MQTT Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The MQTT Source connector consumes messages from an MQTT broker and transfers them to Kafka.
SFTP Source	The SFTP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The SFTP Source connector obtains files from an SFTP server and transfers them to Kafka.
Syslog TCP Source	The Syslog TCP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The Syslog TCP Source connector listens on a port for syslog messages over TCP and transfers them to Kafka.

Connector	Description
Syslog UDP Source	The Syslog UDP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The Syslog UDP Source connector listens on a port for syslog messages over UDP and transfers them to Kafka.
ADLS Sink	The ADLS Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka connect framework. The ADLS Sink connector fetches messages from Kafka and uploads them to ADLS.
HDFS Stateless Sink	The HDFS Stateless Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect Framework. The HDFS Stateless Sink Connector fetches messages from Kafka and stores them in HDFS.
HTTP Sink	The HTTP Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The HTTP Sink connector obtains messages from a Kafka topic and transfers their content in a HTTP POST requests to a specified endpoint.
InfluxDB Sink	The InfluxDB Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect Framework. The InfluxDB Sink Connector fetches messages from Kafka and loads them to InfluxDB.
JDBC Sink	The JDBC Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The JDBC Sink connector fetches messages from Kafka and loads them into a database table.
Kudu Sink	The Kudu Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The Kudu Sink connector fetches messages from Kafka and loads them into a table in Kudu.
S3 Sink	The S3 Sink Connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. The S3 Sink connector fetches messages from Kafka and uploads them to AWS S3.

Standard connectors

The following are either Cloudera developed connectors or connectors that come packaged with Apache Kafka.

Table 10: Standard connectors

Connector	Description
Amazon S3 Sink	The Amazon S3 Sink connector is Cloudera developed connector that consumes data from Kafka topics and streams the data to an S3 bucket.
HDFS Sink	The HDFS Sink connector is a Cloudera developed connector that transfer data from Kafka topics to files on HDFS clusters.
MirrorSourceConnector	The <code>MirrorSourceConnector</code> is a connector used internally by Streams Replication Manager (SRM). Within SRM, this connector is responsible for replicating topics between the source and target cluster. Standalone use of this connector is not recommended by Cloudera.
MirrorHeartbeatConnector	The <code>MirrorHeartbeatConnector</code> is a connector used internally by Streams Replication Manager (SRM). Within SRM, this connector is responsible for creating the heartbeats topic in the target cluster. It also periodically produces heartbeats into the heartbeats topic. Standalone use of this connector is not recommended by Cloudera.

Connector	Description
MirrorCheckpointConnector	The <code>MirrorCheckpointConnector</code> is a connector used internally by Streams Replication Manager (SRM). Within SRM, this connector is responsible for replicating the committed group offsets between the source and target cluster. In addition, the connector is also capable of periodically applying the offsets to the consumer groups in the target cluster. Standalone use of this connector is not recommended by Cloudera.



Note: Standalone use of the `MirrorCheckpointConnector`, `MirrorHeartbeatConnector`, and `MirrorSourceConnector` is not recommended by Cloudera. Additionally Cloudera does not provide documentation on how to use or configure these connectors. However, you can review [Streams Replication Manager driver architecture](#) to learn more about how Streams Replication Manager (SRM) uses these connectors internally.

Example connectors (non-production)

The `FileStream Source` and `FileStream Sink` connectors are example connectors packaged with Apache Kafka. Both `FileStream` connectors are meant to be used to demonstrate the capabilities of Kafka Connect and are not production ready. As a result, even though these connectors are shipped with Runtime, they must be installed before they can be deployed. Prior to connector installation, the connectors are also not selectable for deployment on the SMM UI. See [Installing FileStream connectors](#) on page 44 for installation steps.

Table 11: Non-production example connectors

Connector	Description
<code>FileStream Sink</code>	The <code>FileStream Sink</code> connector reads data from Kafka and transfers that data to a local file.
<code>FileStream Source</code>	The <code>FileStream Source</code> connector reads data from a file and transfers the data to Kafka.



Note: Cloudera does not provide documentation regarding how to use or configure `FileStream` connectors. For more information on connector use, see the [Apache Kafka Documentation](#).



Note: Information presented on this page contains content that is modified and adapted from the [Debezium Documentation](#) by the Debezium Community, which is licensed under [CC BY 3.0](#).

Installing Kafka Connect connectors

Learn how to install custom developed (third party) connectors as well as the `FileStream` connectors in CDP.

Kafka Connect connectors are distributed as:

- A directory of JAR files:
The directory includes the JAR for the connector itself, as well as all its dependencies.
- An uber JAR/FAT JAR/JAR with dependencies file:

This is a single JAR file that contains the connector, as well as its dependencies.

In CDP, all connectors that do not come prepackaged with the Runtime distribution, as well as the `FileStream` connectors (`FileStreamSourceConnector` and `FileStreamsSinkConnector`), must be installed manually. This is done by making the connector JAR files available on all cluster hosts in a specific location. After installation is complete, you will be able to deploy the connector using the Streams Messaging Manager UI (recommended), Streams Messaging Manager API, or Kafka Connect API.

The location of the JAR files is determined by the Kafka Connect role's `plugin.path` property. Kafka Connect discovers connectors by looking at this directory path on the host machines. By default, the `plugin.path` property is set to `/var/lib/kafka`. This means that, by default, any connector placed in this directory will be discovered by Kafka Connect. Cloudera recommends that you use the default path.

The installation steps differ for custom developed connectors and the FileStream connectors. This is because the JAR file for the FileStream connectors is by default available on CDP cluster hosts. Additionally, FileStream connectors can be installed with an alternate installation method that involves the usage of an advanced configuration snippet.

Installing custom developed Kafka Connect connectors

Learn how to install custom developed (third party) connectors in CDP.

Before you begin

The following steps assume that `plugin.path` is set to `/var/lib/kafka`, which is the default path.



Important: Steps 1 and 2 must be carried out on all hosts in the cluster that have Kafka Connect roles deployed on them.

Procedure

1. Log in to a host that is running a Kafka Connect role.
2. Make the connector files available in or readable from `/var/lib/kafka`.

How you choose to complete this step will largely depend on your cluster environment. For example:

- You can download or copy the files directly to `/var/lib/kafka`.
- You can choose to place connector files in a location different from `/var/lib/kafka` and create symlinks that point to the location where the connector files are available.

Regardless of what method you choose, this step is considered complete once the JAR files are readable from `/var/lib/kafka`.

3. Restart all Kafka Connect roles:
 - a) In Cloudera Manager, select the Kafka service.
 - b) Go to Instances.
 - c) Select all Kafka Connect instances by checking the checkbox next to each instance.
 - d) Click Actions for selectedRestart.
 - e) Click Restart to confirm.

The roles are restarted once a **Finished** status is displayed.

Results

The connector is installed and is available for deployment. You are now able to deploy and manage the new connector from the Streams Messaging Manager (SMM) UI.

What to do next

Deploy an instance of the connector using SMM. For more information, see *Deploying and managing Kafka Connect connectors in SMM*.

Related Information

[Deploying and managing Kafka Connect connectors in SMM](#)

Installing FileStream connectors

Learn how to install the FileStream example connectors (FileStreamSourceConnector and FileStreamSinkConnector) that are shipped with Cloudera Runtime but are not installed by default. You can choose between two installation methods.

About this task

The JAR file for the FileStream connectors is shipped with Cloudera Runtime and is readily available on the cluster hosts. However, the file is not added to the Kafka Connect `plugin.path` directory by default. This is because the connectors are meant to be used for demonstrating the capabilities of Kafka Connect and are not production ready. As a result, even though these connectors do come packaged with Runtime, they must be installed before they can be deployed. The JAR file is located at `/opt/cloudera/parcels/CDH/jars/connect-file-[***KAFKA COMPONENT VERSION***].jar`.

You have two options when installing the FileStream connectors. You can install the connectors by copying or symlinking the JAR files to the `plugin.path` directory. Alternatively, you can add the location of the JAR file to the Kafka Connect role's `CLASSPATH` environment variable using an advanced configuration snippet.

The main difference between the two installation methods is that copying or symlinking the file requires that you log in to each Kafka Connect host in your cluster. Using an advanced configuration snippet, on the other hand, enables you to install the connector on all hosts by changing a single property in Cloudera Manager.

Although using an advanced configuration snippet is more convenient than copying or symlinking, be aware that setting an advanced configuration snippet is considered an advanced configuration practice. Therefore, Cloudera advises caution if you choose to install the FileStream connectors using an advanced configuration snippet.

Before you begin

- Go to [Cloudera Runtime component versions](#) and note down the component version of Apache Kafka. You need to specify the version during installation.

The version is made up of three parts. It contains the upstream Apache Kafka version (first three digits), the Runtime version (digits four to six), and the Runtime build number (last three digits denominated with a dash). For example: 3.1.1.7.1.8.0-801.

- The following steps assume that `plugin.path` is set to `/var/lib/kafka`, which is the default path.

Procedure

1. Install the FileStream connectors:

For Copy or symlink

Complete the following actions on all Kafka Connect hosts.

- a. Log in to a host that is running a Kafka Connect role.
- b. Copy or symlink the JAR file of the connectors to /var/lib/kafka.
 - To copy the file directly to /var/lib/kafka use the following command:

```
cp /opt/cloudera/parcels/CDH/jars/connect-file-[***KAFKA COMPONENT VERSION***].jar /var/lib/kafka/
```

- To create a symlink in /var/lib/kafka that points to /opt/cloudera/parcels/CDH/jars/connect-file-**[***KAFKA COMPONENT VERSION***]**.jar use the following command:

```
ln -sf /opt/cloudera/parcels/CDH/jars/connect-file-[***KAFKA COMPONENT VERSION***].jar /var/lib/kafka/connect-file-[***KAFKA COMPONENT VERSION***].jar
```

Replace **[***KAFKA COMPONENT VERSION***]** with the component version of Apache Kafka. Regardless of what method you choose, this step is considered complete once the JAR file is readable from /var/lib/kafka.

For Using an advanced configuration snippet

- a. In Cloudera Manager, select the Kafka service.
 - b. Go to Configuration.
 - c. Find the Kafka Connect Environment Advanced Configuration Snippet (Safety Valve) property.
 - d. Click **+** to add a property and enter the following key and value pair:
 - Key: CLASSPATH
 - Value: /opt/cloudera/parcels/CDH/jars/connect-file-**[***KAFKA COMPONENT VERSION***]**.jar
- Replace **[***KAFKA COMPONENT VERSION***]** with the component version of Apache Kafka.
- e. Click Save Changes.

2. Restart all Kafka Connect roles:

- a) In Cloudera Manager, select the Kafka service.
- b) Go to Instances.
- c) Select all Kafka Connect instances by checking the checkbox next to each instance.
- d) Click Actions for selected Restart.
- e) Click Restart to confirm.

The roles are restarted once a **Finished** status is displayed.

Results

The FileStream connectors are installed and available for deployment. You are now able to deploy and manage the FileStream connectors from the Streams Messaging Manager (SMM) UI.

What to do next

Deploy an instance of the FileStream Source or FileStream Sink connector using SMM. For more information, see *Deploying and managing Kafka Connect connectors in SMM*.

Related Information

[Setting an Advanced Configuration Snippet for a Cloudera Runtime Service](#)

[Deploying and managing Kafka Connect connectors in SMM](#)

Setting up the Debezium Db2 Source connector [Technical preview]

Learn about the CDP specific setup steps required before you can deploy the Debezium Db2 Source connector.

About this task



Important: This is a technical preview feature and considered under development. Do not use this in your production systems. If you have feedback, contact Support by logging a case on our Cloudera Support Portal at <https://my.cloudera.com/support.html>. Technical preview features are not guaranteed troubleshooting and fixes.

In CDP, before deploying an instance of the Debezium Db2 Source connector, you must download and deploy the Db2 JDBC driver on all Kafka Connect hosts. Otherwise, you will not be able to deploy the connector. The following list of steps walks you through this process.

For more information regarding how the Debezium DB2 connector works as well as its configuration and properties, see *Debezium connector for Db2* in the Debezium documentation.



Important: In addition to the following configuration steps, appropriately configuring your Db2 tables is also required. Ensure that you review *Setting up Db2* in the Debezium documentation and complete all setup tasks either before or after completing the following steps.

Procedure

1. Download the Db2 JDBC driver.

The JDBC driver can be downloaded from *Maven Central*. Cloudera recommends that you use version 11.5.0.0.

2. Copy and rename the driver to the location specified in the JDBC Driver Jars Path For Debezium Kafka property. The default location is `/var/lib/kafka_connect_jdbc`. The JDBC driver must be renamed to `db2-connector-java.jar`. For example:

```
cp ./jcc-[***VERSION***].jar var/lib/kafka_connect_jdbc/db2-connector-java.jar
```

3. Set correct file permissions. For example:

```
chmod 644 /var/lib/kafka_connect_jdbc/db2-connector-java.jar
```

4. In Cloudera Manager, restart the Kafka service.

Results

The Db2 JDBC driver is deployed on all Kafka Connect hosts. Deploying an instance of the Debezium Db2 Source connector in your cluster is now possible.

What to do next

Deploy an instance of the Debezium Db2 Source connector. Additionally, review *Debezium connector for Db2* in the Debezium documentation to learn more about the connector and its features.

Related Information

[Debezium connector for Db2](#)

[Setting up Db2](#)

[Maven central](#)

[Deploying and managing Kafka Connect connectors in SMM](#)

Setting up the Debezium MySQL Source connector

Learn about the CDP specific setup steps required before you can deploy the Debezium MySQL Source connector.

About this task

In CDP, before deploying an instance of the Debezium MySQL Source connector, you must download and deploy the MySQL JDBC driver on all Kafka Connect hosts. Otherwise, you will not be able to deploy the connector. The following list of steps walks you through this process.

For more information regarding how the Debezium MySQL connector works as well as its configuration and properties, see the *Debezium connector for MySQL* in the Debezium documentation.



Important: In addition to the following configuration steps, appropriately configuring your MySQL server is also required. Ensure that you review *Setting up MySQL* in the Debezium documentation and complete all setup tasks either before or after completing the following steps.

Procedure

1. Download the MySQL JDBC driver.

The JDBC driver can be downloaded from the *MySQL Community Downloads* page. Download the platform independent version. Cloudera recommends that you use version 8.0.27.

2. Extract the files from the downloaded archive.

You can use the tar command or any other archive manager to extract the archive. For example:

```
tar -xf mysql-connector-java-[*VERSION*].tar.gz
```

3. Navigate to the directory where the driver was extracted.

```
cd mysql-connector-java-[*VERSION*]
```

4. Copy and rename the driver to the location specified in the JDBC Driver Jars Path For Debezium Kafka property.

The default location is `/var/lib/kafka_connect_jdbc`. The JDBC driver must be renamed to `mysql-connector-java.jar`. For example:

```
cp ./mysql-connector-java-[*VERSION*].jar /var/lib/kafka_connect_jdbc/  
mysql-connector-java.jar
```

5. Set correct file permissions. For example:

```
chmod 644 /var/lib/kafka_connect_jdbc/mysql-connector-java.jar
```

Results

The MySQL JDBC driver is deployed on all Kafka Connect hosts. Deploying an instance of the Debezium MySQL Source connector in your cluster is now possible.

What to do next

Deploy an instance of the Debezium MySQL Source connector. Additionally, review *Debezium connector for MySQL* in the Debezium documentation to learn more about the connector and its features.

Related Information

[Debezium connector for MySQL](#)

[Setting up MySQL](#)

[MySQL Community Downloads](#)

Setting up the Debezium Oracle Source connector

Learn about the CDP specific setup steps required before you can deploy the Debezium Oracle Source connector.

About this task

In CDP, before deploying an instance of the Debezium Oracle Source connector, you must download and deploy the Oracle JDBC driver on all Kafka Connect hosts. Otherwise, you will not be able to deploy the connector. The following list of steps walks you through this process.

For more information regarding how the Debezium Oracle Source connector works as well as its configuration and properties, see the *Debezium Connector for Oracle* in the Debezium documentation.



Note: In addition to the following configuration steps, appropriately configuring your Oracle server is also required. Ensure that you review the *Setting up Oracle* in the Debezium documentation and complete all setup tasks either before or after completing the following steps.

Procedure

1. Download the Oracle JDBC driver.

The JDBC driver can be downloaded from *Maven Central*. Cloudera recommends that you use version 12.2.0.1, 19.8.0.0, or 21.1.0.0.

2. Copy and rename the driver to the location specified in the JDBC Driver Jars Path For Debezium Kafka property.
3. The default location is `/var/lib/kafka_connect_jdbc`. The JDBC driver must be renamed to `oracle-connector-java.jar`. For example:

```
cp ./ojdbc8-[*VERSION*].jar /var/lib/kafka_connect_jdbc/oracle-connector-java.jar
```

4. Set correct file permissions. For example:

```
chmod 644 /var/lib/kafka_connect_jdbc/oracle-connector-java.jar
```

Results

The Oracle JDBC driver is deployed on all Kafka Connect hosts. Deploying an instance of the Debezium Oracle Source connector in your cluster is now possible.

What to do next

Deploy an instance of the Debezium Oracle Source connector. Additionally, review *Debezium Connector for Oracle* in the Debezium documentation to learn more about the connector and its features.

Related Information

[Debezium Connector for Oracle](#)

[Setting up Oracle](#)

[Maven Central](#)

Setting up the Debezium PostgreSQL Source connector

Learn about the CDP specific setup steps required before you can deploy the Debezium PostgreSQL Source connector.

About this task

In CDP, before deploying an instance of the Debezium PostgreSQL Source connector, you must download and deploy the PostgreSQL JDBC driver on all Kafka Connect hosts. Otherwise, you will not be able to deploy the connector. The following list of steps walks you through this process.

For more information regarding how the Debezium PostgreSQL connector works as well as its configuration and properties, see the *Debezium connector for PostgreSQL* in the Debezium documentation.



Important: In addition to the following configuration steps, appropriately configuring your PostgreSQL server is also required. Ensure that you review the *Set up* section in the Debezium documentation for the PostgreSQL connector and complete all setup tasks either before or after completing the following steps.

Procedure

1. Download the PostgreSQL JDBC driver.

The JDBC driver can be downloaded from *Maven Central*. Cloudera recommends that you use version 42.2.24.

2. Copy and rename the driver to the location specified in the JDBC Driver Jars Path For Debezium Kafka property. The default location is `/var/lib/kafka_connect_jdbc`. The JDBC driver must be renamed to `postgresql-connector-java.jar`. For example:

```
cp ./postgresql-[***VERSION***].jar /var/lib/kafka_connect_jdbc/postgresql-connector-java.jar
```

3. Set correct file permissions. For example:

```
chmod 644 /var/lib/kafka_connect_jdbc/postgresql-connector-java.jar
```

Results

The PostgreSQL JDBC driver is deployed on all Kafka Connect hosts. Deploying instance of the Debezium PostgreSQL Source connector in your cluster is now possible

What to do next

Deploy an instance of the Debezium PostgreSQL Source connector. Additionally, review *Debezium connector for PostgreSQL* in the Debezium documentation to learn more about the connector and its features.

Related Information

[Debezium connector for PostgreSQL](#)

[Set up \(PostgreSQL connector\)](#)

[Maven Central](#)

Setting up the Debezium SQL Server Source connector

Learn about the CDP specific setup steps required before you can deploy the Debezium SQL Server Source connector.

About this task

In CDP, before deploying an instance of the Debezium SQL Server Source connector, you must download and deploy the SQL Server JDBC driver on all Kafka Connect hosts. Otherwise, you will not be able to deploy the connector. The following list of steps walks you through this process.

For more information regarding how the Debezium PostgreSQL connector works as well as its configuration and properties, see the *Debezium connector for SQL Server* in the Debezium documentation.



Important: In addition to the following configuration steps, appropriately configuring your SQL Server database is also required. Ensure that you review the *Setting up SQL Server* in the Debezium documentation and complete all setup tasks either before or after completing the following steps.

Procedure

1. Download the SQL Server JDBC driver.

The JDBC driver can be downloaded from the Microsoft JDBC Driver for SQL Server downloads page. Cloudera recommends that you use version 7.2.2.jre8. Download the zipped tar file, do not download the executable version.

2. Extract the files from the downloaded archive.

You can use the tar command or any other archive manager to extract the archive. For example:

```
tar -xf sqljdbc_***VERSION***_enu.tar.gz
```

3. Navigate to the directory where the driver was extracted.

```
cd sqljdbc_***VERSION***/enu
```

4. Copy and rename the driver to the location specified in the JDBC Driver Jars Path For Debezium Kafka property. The default location is `/var/lib/kafka_connect_jdbc`. The JDBC driver must be renamed to `mssql-connector-java.jar`. For example:

```
cp ./mssql-jdbc-***VERSION***.jre8.jar /var/lib/kafka_connect_jdbc/mssql-connector-java.jar
```

5. Set correct file permissions. For example:

```
chmod 644 /var/lib/kafka_connect_jdbc/mssql-connector-java.jar
```

Results

The SQL Server JDBC driver is deployed on all Kafka Connect hosts. Deploying instance of the Debezium SQL Server Source connector in your cluster is now possible.

What to do next

Deploy an instance of the Debezium SQL Server Source connector. Additionally, review *Debezium connector for SQL Server* in the Debezium documentation to learn more about the connector and its features.

Related Information

[Debezium connector for SQL Server](#)

[Setting up SQL Server](#)

[Microsoft JDBC Driver for SQL Server](#)

HTTP Source connector

The HTTP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The HTTP Source connector listens on a port for HTTP POST requests and transfers the request body to a Kafka topic. The Kafka topic this connector transfers messages to is determined by the value of the `topics` property. The connector does not perform record processing, the messages are transferred to Kafka as they are received in the HTTP request. The connector uses HTTPS for secure communication. Optionally, client certificate authentication (mutual TLS) can also be configured.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *HTTP Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- The HTTP Source connector must use at least one-way SSL, it cannot be used without SSL.
- The ListenHTTP processor in the NiFi flow keeps running when the connector is paused. Due to this, the connector accepts requests, but ultimately times out and sends an error response to the client.
- The tasks.max property must be 1 to avoid port conflict due to executing multiple tasks on the same Kafka Connect node. Cloudera recommends that you use a load balancer to keep track of the actual node which runs the connector.

Configuration example

In this example, the connector accepts client requests sent to the connector's URL. The connector's URL is based on the configuration of multiple properties. For example assume the following:

- The connector is running on kc-host-1.
- Listening Port is set to port 8080.
- Base Path is set to data-ingest.

In a case like this, the connector's URL will be the following:

```
https://kc-host-1:8080/data-ingest
```



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case and therefore can be removed from the configuration JSON.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
  "meta.smm.predefined.flow.name": "HTTP Source",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.HTTP Source Parameters:Listening Port": "[***PORT***]",
  "parameter.HTTP Source Parameters:Base Path": "[***BASE PATH***]",
}
```

```

"parameter.HTTP Source Parameters:Keystore Filename": "[**PATH TO
KEYSTORE**]",
"parameter.HTTP Source Parameters:Keystore Key Password": "[**KEYSTORE KEY
PASSWORD**]",
"parameter.HTTP Source Parameters:Keystore Password": "[**KEYSTORE
PASSWORD**]",
"parameter.HTTP Source Parameters:Keystore Type": "[**KEYSTORE TYPE**]"
}

```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic that the connector sends messages to.

Listening Port

The port that the client needs to use in the URL.

Base Path

The path that the client needs to use in the URL.

Keystore *

The properties for accessing the keystore used for HTTPS.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[HTTP Source properties reference](#)

JDBC Source connector

The JDBC Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The JDBC Source connector reads records from a database table and transfers each record to Kafka in Avro or JSON format. The Kafka topic this connector transfers messages to is determined by the value of the topics parameter in the configuration. The schema of the records can be inherited from the schema of the database table, or it can be a predefined schema retrieved from Schema Registry.

The connector supports incremental loading based on a table column containing increasing values (typically an ID sequence or timestamp). When the connector is started for the first time with incremental loading turned on, it can load all the existing data from the source table, or it can start from the current maximum value of the incrementing column.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These are the properties that are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *JDBC Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.
- The tasks.max property must be set to 1. Setting the property to any other value will not have an effect. This is because this connector performs incremental loading and must run as a single task.
- If Initial Load Strategy is set to Start at Beginning and “Maximum-value Columns is not set, the connector tries to load the entire content of the source database table to the target Kafka topic with every execution of the connector. Set Maximum-value Columns to avoid this.

Configuration example

In this example, the connector connects to a PostgreSQL database using username/password authentication, queries the records from a table, converts each record to JSON format, and sends them to a Kafka topic. The schema of the records is determined from the schema of the database table. The data loading is incremental, only new records are loaded based on a column containing incrementing values in the table.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
  "meta.smm.predefined.flow.name": "JDBC Source",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "key.attribute": "kafka.message.key",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.JDBC Source Parameters:Database Connection URL": "[***JDBC URL***]",
  "parameter.JDBC Source Parameters:Database Driver Location": "[***PATH TO JDBC DRIVER***]",
  "parameter.JDBC Source Parameters:Database Driver Class Name": "org.postgresql.Driver",
  "parameter.JDBC Source Parameters:Database Type": "PostgreSQL",
  "parameter.JDBC Source Parameters:Database User Name": "[***USERNAME***]",
```

```

"parameter.JDBC Source Parameters:Database User Password":
" [***PASSWORD***] ",
"parameter.JDBC Source Parameters:Database Table Name": " [***TABLE
NAME***] ",
"parameter.JDBC Source Parameters:Kafka Message Key Column": " [***COLUMN
NAME***] ",
"parameter.JDBC Source Parameters:Maximum-value Columns": " [***COLUMN
NAME***] ",
"parameter.JDBC Source Parameters:Kafka Message Data Format": "JSON",
"parameter.JDBC Source Parameters:Schema Access Strategy": "Inherit Schema"
}

```

The following list collects the properties from the configuration example that must be customized for this use case. **topics**

The name of the Kafka topic that the connector sends messages to.

Database Connection URL

The JDBC URL of the PostgreSQL database. For example, jdbc:postgresql://myhost:5432/mydb.

Database Driver Location

A comma-separated list of files or folders containing the JDBC client libraries.

Database Driver Class Name

The Java class name of the PostgreSQL Driver implementation.

Database Type

The type of the database. Because this is a PostgreSQL example, the property is set to PostgreSQL.

Database User Name

The username used for authenticating to the database.

Database User Password

The password of the user.

Database Table Name

The name of the database table to query data from.

Kafka Message Key Column

Specifies a database table column. The value of the column specified is used as the key of the Kafka message.

Maximum-value Columns

The name of the column used for incremental loading. The column value must be increasing, for example, coming from a sequence or system timestamp.

Kafka Message Data Format

The format of the messages the connector sends to Kafka.

Schema Access Strategy

Specifies the strategy used for determining the schema of the database record. In this example, this property is set to Inherit Schema, meaning that the schema is determined (inherited) from the database table.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[JDBC Source properties reference](#)

JMS Source connector

The JMS Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The JMS Source Connector consumes messages from a JMS broker and transfers the message body to Kafka. The Kafka topic this connector transfers messages to is determined by the value of the topics parameter in the configuration. The JMS Source connector does not perform record processing, the messages are transferred to Kafka as they are consumed from the JMS broker.

TLS can be used to establish a secure connection between the connector and the JMS broker. The keystore and truststore files necessary for securing the connection must be present on the cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *JMS Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.

Configuration example

In this example, the connector connects to an ActiveMQ Message Broker using a secured connection and username/password authentication.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
```



```

"meta.smm.predefined.flow.name": "JMS Source",
"meta.smm.predefined.flow.version": "1.0.0",
"key.converter": "org.apache.kafka.connect.storage.StringConverter",
"value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter"
,
,
"tasks.max": "1"
"nexus.url": "https://repository.cloudera.com/artifactory/repo",
"extensions.directory": "/tmp/nifi-stateless-extensions",
"working.directory": "/tmp/nifi-stateless-working",
"topics": "[***TOPIC NAME***]"
"parameter.JMS Source Parameters:JMS Broker URI": "[***BROKER URI***]",
"parameter.JMS Source Parameters:JMS Client Libraries": "[***PATH TO CLIENT LIBRARIES***]",
"parameter.JMS Source Parameters:JMS Connection Factory Class Name": "org.apache.activemq.ActiveMQSslConnectionFactory",
"parameter.JMS Source Parameters:JMS Destination Name": "[***QUEUE NAME***]",
"parameter.JMS Source Parameters:JMS Destination Type": "QUEUE",
"parameter.JMS Source Parameters:JMS User Name": "[***USERNAME***]",
"parameter.JMS Source Parameters:JMS User Password": "[***PASSWORD***]",
"parameter.JMS Source Parameters:Truststore Filename": "[***PATH TO TRUSTSTORE***]",
"parameter.JMS Source Parameters:Truststore Password": "[***TRUSTSTORE PASSWORD***]",
"parameter.JMS Source Parameters:Truststore Type": "[***TRUSTSTORE TYPE***]"
}

```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic that the connector sends messages to.

JMS Broker URI

The URI of the ActiveMQ Message Broker. For example, `ssl://activemq-hostname:61617`.

JMS Client Libraries

The path to the directory containing ActiveMQ client libraries.

JMS Connection Factory Class Name

The Java class name of the ActiveMQ ConnectionFactory implementation.

JMS Destination Name

The name of the destination (in this case a queue) on the ActiveMQ Message Broker that the messages are received from.

JMS Destination Type

The type of the destination. In this example, the type is `QUEUE`.

JMS User Name

The username used for authenticating to the ActiveMQ Message Broker.

JMS User Password

The password of the user.

Truststore *

These are the properties for accessing the truststore containing the ActiveMQ Message Broker's certificate used for secure communication.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[JMS Source properties reference](#)

MQTT Source connector

The MQTT Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The MQTT Source connector consumes messages from an MQTT broker and transfers them to Kafka. The Kafka topic this connector transfers messages to is determined by the value of the topics property in the configuration. This connector does not perform record processing. The messages are transferred to Kafka as they are consumed from the MQTT broker. TLS can be used to establish a secure connection between the connector and the MQTT broker. The keystore and truststore files necessary for securing the connection must be present on the cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *MQTT Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- If the keystore-related properties are removed and an empty truststore is provided, the connector does not use TLS for connecting to the MQTT broker. TLS is used if a truststore is provided that has a key in it.

Configuration example

In this example, the connector fetches messages from an MQTT broker and transfers them to a Kafka topic.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
```

```

"connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
"meta.smm.predefined.flow.name": "MQTT Source",
"meta.smm.predefined.flow.version": "1.0.0",
"key.converter": "org.apache.kafka.connect.storage.StringConverter",
"value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
"tasks.max": "1",
"nexus.url": "https://repository.cloudera.com/artifactory/repo",
"extensions.directory": "/tmp/nifi-stateless-extensions",
"working.directory": "/tmp/nifi-stateless-working",
"topics": "[***KAFKA TOPIC NAME***]",
"parameter.MQTT Source Parameters:MQTT Broker URI": "tcp://[***HOST***]:[***PORT***]",
"parameter.MQTT Source Parameters:MQTT Quality of Service": "0",
"parameter.MQTT Source Parameters:MQTT Topics": "[***MQTT TOPIC NAME***]"
}

```

The following list collects the properties from the configuration example that must be customized for this use case: **topics**

The name of the Kafka topic that the connector sends messages to.

MQTT Broker URI

The URI of the MQTT broker. In this example, the connection is not secure. As a result, the URI starts with `tcp://`. The port number you typically use in a case like this is 1883. If TLS is used, the URI should start with `ssl://`, and typically you use port number 8883.

MQTT Topics

The comma-separated list of MQTT topics that the connector fetches messages from.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[MQTT Source properties reference](#)

SFTP Source connector

The SFTP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The SFTP Source connector obtains files from an SFTP server and transfers them to Kafka. The SFTP Source connector supports record processing. In case record processing is enabled, input file is split into records and the records are transferred to Kafka. If record processing is not enabled, the connector forwards files to Kafka as they are fetched from the SFTP server.

If record processing is enabled, the input file can contain records in either JSON, CSV or Grok format. If the input records are in JSON or CSV format, the connector can either infer the schema of the records based on the data or read the schema from Schema Registry. If the input records are in Grok format, the connector can either derive the schema using the field names from the value of the Grok Expression property or read the schema from Schema Registry.

In addition, if record processing is enabled, the connector writes the records to Kafka in Avro format. The record schema gets embedded in these Avro messages if the schema is set to be either inferred or determined based on the Grok Expression property. If Schema Registry is used for getting the schema of the records, then the schema does not get embedded in the Avro messages.

If Schema Registry is used, and it is on a Kerberized cluster, the `krb5.file` property must point to the `krb5.conf` file that provides access to the cluster on which Schema Registry is present. This means that the `krb5.conf` file must be on the same cluster node that the connector runs on. The Kerberos keytab used to access Schema Registry must also be on the same cluster node that the connector runs on. The connection to Schema Registry can be secured by TLS. The truststore file necessary for securing the connection must also be on the same cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *SFTP Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- The Schema Registry URL property is mandatory even if Schema Registry is not used. If Schema Registry is not used, use the default value, or completely remove the property from the configuration JSON.
- Schema Branch and Schema Version can not be specified at the same time.
- The value of the Schema Access Strategy property is not independent of the value of the Input Data Format property. As a result, you must exercise caution when configuring Schema Access Strategy.
 - If the value of Input Data Format is JSON, the possible values for Schema Access Strategy are Schema Registry or Infer Schema.
 - If the value of Input Data Format is CSV, the possible values for Schema Access Strategy are Schema Registry or Infer Schema.
 - If the value of Input Data Format is GROK, the possible values for Schema Access Strategy are Schema Registry or Field Names From Grok Expression.
- If the Enable Record Processing property is set to false (record processing disabled), the entire input file is transferred to Kafka as one message.
- If the Enable Record Processing property is set to true (record processing enabled), the output is always in Avro format. However, the contents of the output message transferred to Kafka depends on the value of the Schema Access Strategy property.
 - If Schema Access Strategy is set to Schema Registry, the output Avro message does not contain the schema (the schema is not embedded in the output).
 - If Schema Access Strategy is set to either Infer Schema or Field Names From Grok Expression, the output Avro message contains the schema (the schema is embedded in the output).

Configuration example

In this example, the connector fetches files from an SFTP server that only requires Basic Authentication (username and password) and transfers these files to a Kafka topic.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
  "meta.smm.predefined.flow.name": "SFTP Source",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "topics": "[*** KAFKA TOPIC NAME***]",
  "parameter.SFTP Source Parameters:Hostname": "[***SFTP Server Hostname***]",
  "parameter.SFTP Source Parameters:Port": "22",
  "parameter.SFTP Source Parameters:Username": "[***SFTP USERNAME***]",
  "parameter.SFTP Source Parameters>Password": "[***SFTP PASSWORD***]",
  "parameter.SFTP Source Parameters:Remote Path": "[***FOLDER PATH***]",
  "parameter.SFTP Source Parameters:File Filter Regex": ".*",
  "parameter.SFTP Source Parameters:Enable Record Processing": "false"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

- topics**
 - The name of the Kafka topic that the connector sends messages to.
- Hostname**
 - The hostname of the remote system where the SFTP server runs. For example, my.sftp-server.com.
- Port**
 - The port that the remote system is listening on for file transfers.
- Username**
 - The username for connecting to the SFTP server.
- Password**
 - The password used to authenticate the user towards the SFTP server. In this example, the SFTP server requires Basic Authentication.
- Remote Path**
 - The path on the remote system that points to the directory from which to pull files. For example, / uploads.
- File Filter Regex**
 - The Java regular expression to use for filtering filenames. Only files whose names match the regular expression are fetched.
- Enable Record Processing**
 - Determines whether the contents of a file get parsed as records before sending them to Kafka. In this example, the property is set to false, meaning that the entire file gets forwarded to Kafka as one message.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[SFTP Source properties reference](#)

Stateless NiFi Source and Sink

The Stateless NiFi Source and Sink connectors allow you to run NiFi dataflows within Kafka Connect. Using these connectors can grant you access to a number of NiFi features without having the need to deploy or maintain NiFi on your cluster. Learn more about what Stateless NiFi is and how it shipped in CDP.

Apache NiFi is a powerful tool for authoring and running dataflows. It provides many capabilities that are necessary for large-scale enterprise deployments, such as data persistence and resilience, data lineage and traceability, as well as multi-tenancy. This, however, requires an administrator who ensures that this process is running and operational. Additionally, in most cases, adding more capabilities to your deployments results in more complexity.

There are times, however, when users do not need all the power of NiFi, and running it in a much simpler form factor is sufficient. A common use case is to use NiFi to pull data from many different sources, perform manipulations (for example, convert JSON to Avro), filter some records, and then publish the data to Apache Kafka. Another common use case is to pull data from Apache Kafka, perform manipulations and filtering, and then publish the data elsewhere.

For deployments where NiFi acts only as a bridge into and out of Kafka, it can be simpler to operationalize such a deployment by running the dataflow within Kafka Connect. The Stateless NiFi Source and Sink connectors allow users to do just that.

Stateless NiFi

Dataflows within Kafka Connect run using the Stateless NiFi dataflow engine. Stateless NiFi differs from the traditional NiFi engine in the following ways:

- Stateless NiFi is an engine that is designed to be embedded. This makes it convenient to run it within the Kafka Connect framework.
- Stateless NiFi does not provide a user interface (UI) or a REST API. Additionally, it does not support modifying the dataflow while it is running.
- Stateless NiFi does not persist flowfile content to disk. Instead, it holds the content in memory.
- Stateless NiFi does not use data prioritizers. Instead, it operates on data in a first-in-first-out order. Dataflows built for Stateless NiFi must have a single source and a single destination. The only exception to this is when the data is routed to exactly one of multiple destinations, such as a failure destination or a success destination.
- Stateless NiFi does not currently provide access to data lineage/provenance.
- Stateless NiFi does not support cyclic graphs. While it is common and desirable in traditional NiFi to have a failure relationship from a processor route back to the same processor, this can result in a `StackOverflowException` in Stateless NiFi. The preferred approach in Stateless NiFi is to create an output port for failures and route the data to that output port.

Stateless NiFi in CPD

Stateless NiFi in CDP comes in two flavors. You have the base Stateless NiFi Source (`StatelessNiFiSource`) and Stateless NiFi Sink (`StatelessNiFiSink`) connectors. In addition, there are a number of ready-to-use connectors based on Stateless NiFi Source and Sink.

Stateless NiFi Source and Sink connectors

The Stateless NiFi Source and Sink connectors are generic connectors that act as the base for Stateless NiFi in Kafka Connect. The idea behind these connectors is that you develop your own dataflow in NiFi, and deploy that dataflow as a Kafka Connect connector using either Stateless NiFi Source or Sink connectors. This way, you can create virtually any connector without having to write a single line of code. All of the heavy lifting is done on the NiFi UI.

Predefined dataflows and Stateless NiFi-based connectors

In addition to the base Stateless NiFi Source and Sink connectors, Cloudera also ships many connectors that are based on Stateless NiFi. Specifically, Cloudera ships a set of predefined dataflows with CDP. These dataflows are developed and maintained by Cloudera. Each dataflow covers a common data movement use case. These predefined flows are presented as standalone connectors on the Streams Messaging Manager (SMM) UI. Using these dataflows/connectors, you can immediately start using Stateless NiFi in CDP without having to develop your own dataflow.

You can find a full list of the Stateless NiFi-based connectors in [Connectors](#) on page 40. In addition, you can also use the SMM UI to check which of the connectors are Stateless NiFi-based. When adding a new connector using SMM, you are presented with a number of cards that represents the connectors available for deployment. Each card includes the following information about the connector:

- The connector's fully qualified class name.
- The connector's display name. If no display name is available, the card includes the unqualified classname.
- The version of the connector.

Notice how the fully qualified class name for many of the default available connectors is either `org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector` or `org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector`. Even though these are standalone connectors and have a unique display name, they are all Stateless NiFi connectors. Each of them runs a predefined flow designed for a specific use case. For example, JDBC Source runs a dataflow that reads records from a database table and transfers each record to Kafka.



Note: Unlike the base Stateless NiFi connectors, connectors running predefined flows **do not** accept a custom dataflow with `flow.snapshot`. The dataflow they use is specified by the `meta.smm.predefined.flow.name` and `meta.smm.predefined.flow.version` properties. While configuring these properties is possible, Cloudera does not recommend that you do so.

Dataflow development best practices for Stateless NiFi

By leveraging Stateless NiFi, you can build your own dataflows and use those as Kafka Connect connectors. However, there are some principles to keep in mind when building such a dataflow in NiFi.

General recommendations and criteria

Learn about the general recommendations and criteria that you need to follow when building a dataflow that you will run as a Kafka Connect connector.

- If your dataflow depends on external resources (like a JDBC driver), the resource must be made available on all Kafka Connect role instance nodes (Connect workers). Additionally, the resource must be in the same location on all nodes with proper file permissions. Cloudera also recommends parameterizing the file path so that you can specify it when deploying the connector.
- The dataflow should be designed in its dedicated process group and should have a parameter context assigned to its process group. This way, you have the ability to parameterize your connector, which enables you to provide custom property values when deploying the connector using Streams Messaging Manager (SMM).
- A source connector flow should always have a single source. This means that the flow should have only one processor with no incoming connections.
- A sink connector flow should always have a single destination. If you want to transfer the data consumed from the Kafka topic to two different destinations at the same time, it is preferred to have two distinct sink connectors with different consumer groups. If using a single sink connector, you must take into consideration the possibility of creating duplicates. If one destination is working as expected but not the other, the Kafka message is not acknowledged and consumed again at a later time while the data was still successfully sent to one destination.

Dataflow execution and scheduling

Learn how Stateless NiFi Source and Sink connectors handle dataflow execution and scheduling.

Stateless NiFi does not consider the scheduling settings of processors. When a Stateless NiFi connector is deployed, Kafka Connect starts an instance of the Stateless NiFi engine that executes the dataflow. The following sections discuss in detail how the dataflow is executed by both source and sink connectors

Source connectors

A Kafka Connect source connector obtains data from an external (source) data store and delivers the data to a Kafka service. The architecture of a source connector is as follows.

Figure 3: Source Connector Architecture

When a Stateless NiFi Source connector instance is deployed, a number of tasks are started for the connector instance. The number of tasks is determined by the `tasks.max` property, which can be set during connector deployment. Each task consists of the following:

- `SourceTask`

In Stateless NiFi connectors, the `SourceTask` is responsible for executing one instance of the connector's dataflow.

- `SMT Chain`

The Single Message Transforms (SMT) chain can modify and filter the data on a single message basis.

- `Converter`

The converter is responsible for converting Kafka message keys and values with a configured conversion mechanism. However, Stateless NiFi connectors currently only use keys of string type and values of byte array type. As a result, the converter does not perform any meaningful conversion.

- `Kafka Producer`

The Kafka producer is responsible for delivering messages to the configured Kafka topics.

With the Stateless NiFi Source connector, the `SourceTask` is the component executing an instance of the dataflow specified in the connector's `flow.snapshot` property. Any source connector flow you build should not attempt to deliver directly to Kafka itself using a processor such as `PublishKafka`. Instead, the data should be routed to a NiFi output port in the dataflow. Any flowfile delivered to that output port is obtained by the connector and delivered to Kafka.

The dataflow is triggered continuously in one thread within the `SourceTask`. This means that in one iteration, the entire flow is executed (one processor after another), and the data it fetches from the source data store in that execution ends up on the flow's output port in the form of flowfiles. These flowfiles are forwarded by the Kafka

Connect framework to the task's producer as messages (one flowfile equals one message). At this point the flow's execution is finished, and the flow is triggered again.

The producer is a regular Kafka producer that is responsible for delivering the messages to the Kafka topics. The producer is running on a separate thread. This means that while the dataflow is running, the producer can work on delivering the messages it received at the end of the flow's previous execution.

Each flowfile is delivered as a single Kafka message. Therefore, if a flowfile contains thousands of JSON records, totaling 5 MB for example, it is important to split those flowfiles into individual records using a `SplitRecord` processor before transferring the data to the output port. Otherwise, this results in a single Kafka message that is 5 MB while the default message size for Kafka is 1 MB.

Scheduling settings of the processors in the dataflow are ignored. The dataflow is executed in one thread from the first processor through the output port. Afterward, the flow is retriggered. One execution of the dataflow takes as long as necessary to fetch and process the data that accumulated in the source data store since the previous execution. However, the execution has a default timeout of one minute. This timeout can be configured on a connector level with the `dataflow.timeout` property.

Transactional sources

Unlike with traditional NiFi, Stateless NiFi keeps the contents of flowfiles in memory. As long as the source of the data is replayable and transactional, there is no concern over data loss. This is handled by treating the entire dataflow as a single transaction. Once data is obtained from a source component, the data is then transferred to the next processor in the flow. At this point, in traditional NiFi, the processor would acknowledge the data and NiFi would take ownership of that data, having persisted it to disk.

With Stateless NiFi, however, the data is not yet acknowledged. Instead, data is transferred from one processor to the next with each processor performing and completing its task. This is repeated until all data is queued up at the output port. At this point, the NiFi connector provides the data to Kafka. Only after Kafka has acknowledged the records does the connector acknowledge this to NiFi. Only at that point is the session committed, allowing the source processor to acknowledge receipt of the data.

As a result, if NiFi is restarted while processing some piece of data, the source processor will not have acknowledged the data and is able to replay the data, resulting in no data loss.

Initial data load

Stateless NiFi connectors are intended to run continuously. This means that the dataflow is retriggered immediately after a single execution finishes. This also means that the dataflow is only supposed to process as many messages at a time that accumulated in the source data store since the previous execution.

Depending on the connector, if there is already a significant amount of data in the source data store and the connector is not configured properly, you might run into problems with the dataflow's initial execution. In such a case, the dataflow tries to fetch all the data from the source data store during the initial execution. If there is too much data, one of the following happens:

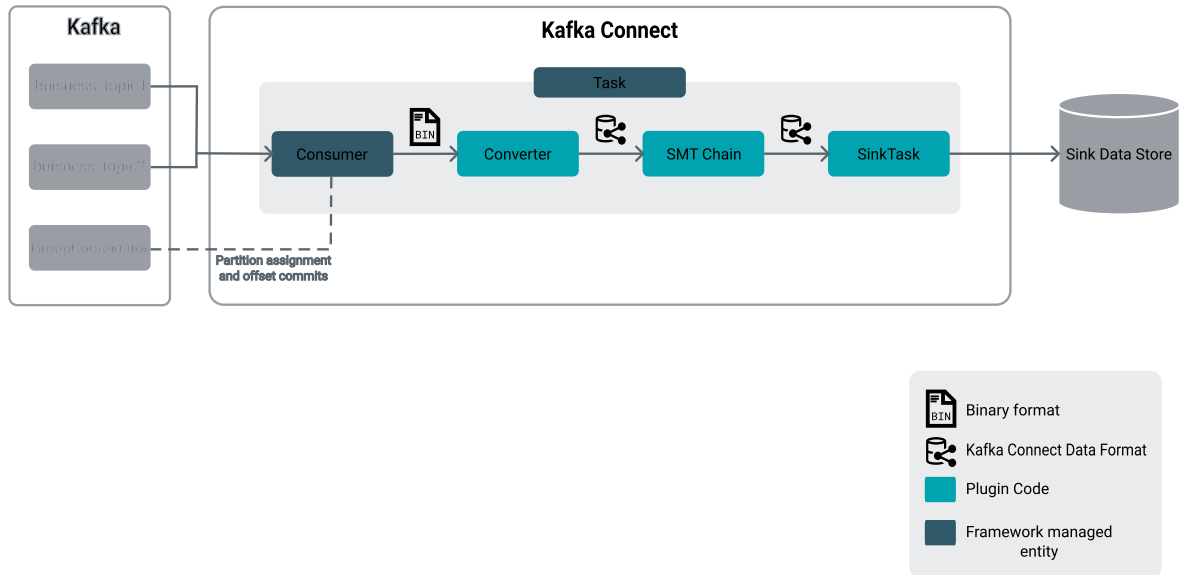
- No data appears in the target Kafka topics, and the log file includes an `InterruptedException` entry.
This means that loading and processing the data took more time than the timeout configured in `dataflow.timeout`.
- Not data appears in the target Kafka topics, and the connector throws an `OutOfMemory` error.
This can happen because the Stateless NiFi engine, which runs the dataflow, stores all data in memory. Trying to read too much data can exhaust the available memory capacity.

If you need to handle a large initial data load, Cloudera recommends that you process and transfer the initial load using another tool, for example regular NiFi. Afterward, you can configure and start the Kafka connector to handle subsequent changes in the source system.

Sink connectors

A Kafka Connect sink connector obtains data from a Kafka service and delivers that data to an external (sink) data store. The architecture of a sink connector is as follows.

Figure 4: Sink Connector Architecture



When a Stateless NiFi Sink connector instance is deployed, a number of tasks are started for the connector instance. The number of tasks is determined by the `tasks.max` property, which can be set during connector deployment. Each task consists of the following:

- **Kafka Consumer**

The consumer is responsible for fetching the messages from the Kafka topics and delivering them to the converter.

- **SMT Chain**

The Single Message Transforms (SMT) chain can modify and filter the data on a single message basis.

- **Converter**

The converter is responsible for converting Kafka message keys and values with a configured conversion mechanism. However, Stateless NiFi connectors currently only use keys of string type and values of byte array type. As a result, the converter does not perform any meaningful conversion.

- **SinkTask**

In Stateless NiFi connectors, the `SinkTask` is responsible for executing one instance of the connector's dataflow.

The consumer is a regular Kafka consumer reading messages from the configured topics and delivering them to the `SinkTask`. Messages are delivered through the converter and the SMT chain.

In a Stateless NiFi Sink connector, the `SinkTask` is the one executing an instance of the dataflow configured in the connector's `flow.snapshot` property. Any sink connector flow you build should not attempt to source data directly from Kafka using a processor such as `ConsumeKafka`. Instead, the data should be received from an input port within the dataflow. Each Kafka record is queued into the outbound connection of the input port so that the next processor in the flow is able to process it. The entire task (consumer + `SinkTask`) is executed in one thread. There are two important timeout parameters relevant for sink connectors.

- **Offset Flush Interval (`offset.flush.interval.ms`)**

This is the interval, in milliseconds, at which Kafka Connect attempts to commit task offsets.

- **`dataflow.timeout`**

Specifies the maximum amount of time to wait for the dataflow to complete.



Important: The configuration scope of the two properties is different. Offset Flush Interval (`offset.flush.interval.ms`) is a Kafka Connect role property that you set in Cloudera Manager. The value set for this property is **global** and affects **all** connectors in your deployment. On the other hand, `dataflow.timeout` is a connector level property. You set it during connector deployment and is configured on a per-connector basis.

The execution of a `SinkTask` is the following:

1. The consumer consumes messages from the Kafka topics for the period of time determined by Offset Flush Interval (default is one minute), and puts the messages on the input of the `SinkTask`, that is, the input port of the connector's dataflow.
2. After Offset Flush Interval is reached, the consumer stops consuming messages and the flow is triggered. Once triggered, the flow processes all messages that accumulated on its input port. Processing time is not infinite. The flow will time out once the value set in `dataflow.timeout` is reached.
3. After flow is finished running, the process restarts with step 1.

Notice that the flow is not triggered for each Kafka message individually. Instead, it is triggered for a batch of messages that the consumer put on its input port in the time specified by Offset Flush Interval. One execution (trigger) of the flow is one transaction. This means that the flow either manages to process all messages within the batch successfully (delivers them to the target system and commits the transaction), or rolls back the session. If the session is rolled back, the consumer is going to put the same batch of messages on the input of the flow in the next execution of step 1. Therefore, in the next execution of step 2, the flow is going to try and process the same batch of messages.

Each Kafka message is delivered to the target system as a single flowfile. Depending on the destination, Cloudera recommends that you consider merging your Kafka records together into bigger chunks before delivering them. For example, assume you are delivering data to HDFS. Because HDFS is not built for accessing and handling many small files, it is not recommended to send each individual Kafka message to HDFS as a separate file.

Merging considerations for sink connectors

Learn about the limitations and considerations when using `MergeContent` and `MergeRecord` processors with Stateless NiFi.



Tip: Before you continue, ensure that you read [Dataflow execution and scheduling](#) first. You must have an understanding of how dataflow execution works to understand the examples and settings discussed here.

NiFi supports many processors that can be used as sinks for Kafka data. If the destination service operates well with many small messages, then the data from Kafka can often be delivered directly to a sink processor. For example, a `PublishJMS` processor can easily handle many small messages. However, other services, like S3 or HDFS, perform much better if the data is first batched or merged together. For this reason, the `MergeContent` and `MergeRecord` processors are extremely popular in NiFi. They make it possible to merge many small flowfiles into one larger flowfile.

With `MergeContent` and `MergeRecord` in traditional NiFi, you can set a minimum and maximum size for the merged data along with a timeout. However, with Stateless NiFi and Kafka Connect, this might not work. This is because only a limited number of flowfiles are made available to the processor. You can still use these processors to merge the data together, but some limitations apply.

If `MergeContent` or `MergeRecord` is triggered but not enough flowfiles are present to create a batch, the processor does nothing. Therefore, Cloudera recommends that you set `Minimum Number of Entries` to 1 and `Minimum Group Size` to 0 B (these are default values). With these settings, the `MergeContent` and `MergeRecord` processors create a merged flowfile from the input flowfiles available in the given execution of the connector and do not wait for more data to arrive. Configuring `Max Bin Age` is not required in this case.

The `Maximum Number of Entries` and `Maximum Group Size` processor properties can also be used to limit the maximum size of the merged flowfile. For example, in a NiFi Stateless Sink connector the consumer running inside the connector puts a group of messages on the input flow. Afterward, the flow is executed. This execution of the flow can only work with the amount of messages that it has on its input when it gets triggered. If it had 115 messages on its input, and `MergeContent` processor's `Maximum Number of Entries` is 100, the flow is going to produce two outgoing messages, one with 100 entries and another one with 15.

To control how many messages are placed on the input of the flow for one execution, you can configure the Offset Flush Interval (`offset.flush.interval.ms`) Kafka Connect role property. The lower the interval, the less time the consumer has to put messages on the input of the flow.

NiFi Stateless Sink merging example

Let's look at a merging example with a NiFi Stateless Sink connector. Assume that the dataflow contains the following elements in the specified order:

```
input port -> MergeContent -> PutHDFS
```

Additionally, assume the following configurations:

- Offset Flush Interval is set to one minute (default). In this timeframe the internal consumer puts 1100 messages on the input port of the dataflow.
- `dataflow.timeout` is 90 seconds.
- The MergeContent processor has the following properties:
 - Maximum Number of Entries = 1
 - Minimum Group Size = 0B
 - Maximum Number of Entries = 100



Important: Offset Flush Interval is a Kafka Connect role property. It is configured in Cloudera Manager. `dataflow.timeout` is a connector property. It is configured during connector deployment, for example in Streams Messaging Manager (SMM). MergeContent's properties are processor properties configured in NiFi when you build the dataflow.

In this case, MergeContent creates 11 flowfiles, and PutHDFS has to upload these 11 files to HDFS within what is left of `dataflow.timeout` (90 seconds). Assume that PutHDFS manages to upload six files when `dataflow.timeout` is up. Since the dataflow did not finish execution within the dedicated timeframe, it gets interrupted, and the transaction is rolled back. The consumer starts putting messages on the input of the dataflow for a minute again.

However, since the previous transaction was rolled back, the messages that the dataflow tried to process were not committed. Therefore, the consumer puts the same 1100 messages on the input port of the dataflow. This behavior can result in an endless loop. Additionally, HDFS was not notified that the six files that were uploaded successfully should be removed. This results in these files being reuploaded in the next execution. You can resolve an issue like this by doing the following:

1. Decrease Offset Flush Interval.

For example, decreasing this interval to 30 seconds (from one minute) halves the number of messages that are put on the input port of the dataflow. Fewer messages on the input port means that the dataflow needs to process fewer messages in one execution. 550 instead of 1100.

2. Set Maximum Number of Entries to a value that is larger than the number of messages the you expect the dataflow will receive.

Increasing this value to, for example, 1,000,000,000 ensures that all messages are put in one HDFS file. This way all data is either successfully uploaded to HDFS, or the upload is interrupted and rolled back, but without leftover bits in the target system.

Considerations for listen-type source connectors

There are a number of limitations that apply to listen-type connectors that you should be aware of.

A listen-type source connector consists of a NiFi dataflow in which the first processor is a ListenX processor. These processors listen on specific interfaces and specific ports on the host where they are running. For example, a flow starting with a ListenHTTP processor listens on a specific port for HTTP requests made by external clients.

The Kafka Connect framework does not provide the concept of task affinity mapping, meaning that if a connector is deployed with multiple tasks, it is not possible to ensure that tasks are assigned to different Kafka Connect workers (Kafka Connect role instances). Two tasks started on the same worker in a listen-type connector results in a failure.

This happens because both connectors try to bind to the same port. For this reason, you must ensure that listen-type source connectors are deployed with a single task. That is, you must set `tasks.max=1` when configuring the connector.

A common deployment model is to have a load balancer in front of the Kafka Connect workers and a mapping rule linking to the port on each Kafka Connect worker. This way, the load balancer FQDN can be provided to the client without assuming which worker the task is going to be deployed on.

If multiple tasks are required for performance, you can choose to have multiple standalone Kafka Connect workers and deploy the connector independently on each worker. However, in such a case, Cloudera recommends using Cloudera Flow Management and NiFi as they provide a significantly more powerful environment for this type of use cases.

Mapping of NiFi features

There are some features that exist in NiFi that translate seamlessly in the context of Kafka Connect. The following discusses some of these NiFi features and how they can be utilized in Kafka Connect.

State management

In NiFi, a processor is capable of storing state about the work that it has accomplished. This is particularly important for source components such as the ListS3 processor. This processor keeps state about the data that it has already seen so that it does not list the same files repeatedly. State in NiFi can be stored local to the node (local state), or across all nodes in a cluster (cluster-wide state).

When using the Stateless NiFi Source connector, the state that is stored by these processors is provided to Kafka Connect and is stored within Kafka itself as source offsets and source partition information. This allows a tasks to restart and resume where they left off. If a processor stores local state in NiFi, it will be stored in Kafka using a source partition that corresponds directly to that task. As a result, each task is analogous to a node in a NiFi cluster. If the processor stores cluster-wide state in NiFi, the state will be stored in Kafka using a source partition that corresponds to a cluster-wide state.



Note: Kafka Connect does not allow for state to be stored for a sink task.

Primary node

NiFi provides several processors that are expected to run only on a single node in the cluster. This type of execution is configured by setting Execution to Primary Node Only in the SCHEDULING tab on the NiFi UI when configuring a processor.

When using the Stateless NiFi Source connector, if any source processor in the dataflow is set to run on a primary node only, only a single task will ever run, even if the `tasks.max` connector property is set to a value higher than one. If you attempt to use multiple tasks for a dataflow that has Primary Node Only set, a warning is logged.

Because processors should only be scheduled with Primary Node Only if they are sources of data, this setting is ignored for all sink tasks and for any processor in a source task that has incoming connections.

Processor yielding

When a processor determines that it is not capable of performing any work (for example, because the system that the processor is pulling from has no more data available), it might choose to yield. This means that the processor stops running for some amount of time. You can configure how long the processors stops running in a dataflow by configuring Processor Configure SETTINGS Yield Duration in NiFi. When using a source connector, if a processor chooses to yield, the source connector pauses for the configured amount of time before triggering the dataflow to run again.

Kafka Connect worker assignment and Stateless NiFi

Learn how NiFi Stateless connectors and tasks get assigned to Connect workers, in what cases you need to be aware of worker assignments, and how you can look up worker assignments.



When running Kafka Connect in a cluster environment, each cluster node in the Kafka Connect cluster has a Kafka Connect worker running on it. In CDP, when the cluster is managed by Cloudera Manager, Kafka Connect workers are represented by Kafka Connect role instances, which are deployed under your Kafka service. This means that the number of workers available in the cluster is identical to the number of Kafka Connect role instances.

When you deploy a connector, the connector is assigned to a worker. Additionally the task or tasks of the connector are also assigned to workers. The workers that the connector and tasks are assigned to might be the same, but might also be different. If the connector has multiple tasks (`tasks.max>1`), each task might be assigned to the same worker, or they might be assigned to different workers. In addition, if a Kafka Connect worker is stopped, the connectors and tasks that were assigned to it are moved to different workers in the Connect cluster. In general, you have no agency over worker assignment, it is decided by the Connect Framework, and it is not configurable. Additionally, worker assignments might change dynamically, however, this is a normal and expected behavior of Kafka Connect.



Note: In a test environment you can force deploy connectors and tasks to a specific worker by stopping all other workers in the Connect cluster.

With Stateless NiFi connectors, knowing exact worker assignments is required when you are running a listen-type connector (for example, HTTP Source, Syslog TCP Source, or Syslog UDP Source). With these types of connectors you must ensure that you send messages to the host of the task's worker and not the connector's worker.

You can look up worker assignments in Streams Messaging Manager (SMM) on the  Connector Profile page of the connector. The worker assigned to the connector is displayed in the **Connector Profile** section under **ASSIGNED WORKER** (1). The worker assigned to the task is displayed in the  **Tasks** section under **Worker ID** (2).



Connect Cluster / Connector Profile

sink-demo

[Connector Profile](#) Connector Settings

Connector Profile

CLASSNAME

org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector



STATUS

RUNNING

TOTAL TASKS

1

Tasks

Search by host

Rest



Status

2

Worker ID ↓



172.18.0.6

Running Ratio



RUNNING



Related Information

[Kafka Connect log files and Stateless NiFi](#)

Kafka Connect log files and Stateless NiFi

Learn how and where Stateless NiFi connectors produce their logs.

Each cluster host that runs a Kafka Connect worker (Kafka Connect role) has a log file specific to Kafka Connect. The logs are stored by default in `/var/log/kafka`. The Connect related logs are named `kafka-connect-[***HOSTNAME***].log`. The log directory path is specified by the Kafka Connect Log Directory Kafka Connect role property. This property is configured in Cloudera Manager.

Every log entry produced by a NiFi dataflow that is running in a Kafka Connect connector is directed to the Kafka Connect log file. Specifically, the dataflow log entries are written to the log file located on the cluster host that is running the Kafka Connect worker that is executing the connector's task.

For example, if your connector has a single task, which is executed by worker-1 running on host-1, then the log for that connector will be located on host-1. If a connector has multiple tasks, each running on a different worker, then each task's dataflow writes logs to the respective worker's log file.

If a Kafka Connect worker is stopped, its tasks are moved to other workers. As a result, it is possible that a specific connector's log entries first appear in the log file located on one host, but later on appear in another log file located on a different host.

Related Information

[Kafka connect worker assignment and Stateless NiFi](#)

Kafka Connect tasks and Stateless NiFi

Learn which Stateless NiFi connectors you can or cannot run with multiple Kafka Connect tasks (`task.max` \geq 2).

By default all connectors that you deploy start with a single Kafka Connect task. The number of tasks deployed for a connector can be configured using the `tasks.max` connector property, but not all types of connectors are meant to run using more than a single task.

Specifically with Stateless NiFi, there are a handful of source connectors that either cannot run with multiple tasks or function incorrectly when multiple tasks are configured. The following provides information on which connectors can or cannot be run with multiple tasks and what effects configuring multiple tasks has on both sink and source connectors.

Source connectors

Notice that the connector's task is made up of multiple components, these are as follows:

- `SourceTask`

In Stateless NiFi connectors, the `SourceTask` is responsible for executing one instance of the connector's dataflow.

- SMT Chain

The Single Message Transforms (SMT) chain can modify and filter the data on a single message basis.

- Converter

The converter is responsible for converting Kafka message keys and values with a configured conversion mechanism. However, Stateless NiFi connectors currently only use keys of string type and values of byte array type. As a result, the converter does not perform any meaningful conversion.

- Kafka Producer

The Kafka producer is responsible for delivering messages to the configured Kafka topics.

If you deploy a source connector with more than one task, then each task is going to have its own instance of the task components. Because the dataflow in the `SourceTask` is responsible for fetching the data, rather than the source data store providing the data to the `SourceTask`, it is not always possible to start a connector with more than one task.

For example, starting a JDBC Source connector, which is a Stateless NiFi-based connector, with two tasks would mean that two instances of the dataflow are running in parallel. Both are fetching data from the same table of the same database. These two tasks would be delivering the messages to the same Kafka topic. However, the two dataflow instances running in parallel would have no way of knowing whether a certain record was already been fetched by the other instance. This results in the same piece of data appearing in multiple Kafka messages.

Because of this, there are a number of limitations that you must keep in mind when using the Stateless NiFi Source connectors. These limitations apply when you are running a custom dataflow with a Stateless NiFi Source connector, as well when running any of the Stateless NiFi-based source connectors that Cloudera develops.

Listen-type connectors

A listen-type connector consists of a NiFi dataflow in which the first processor is a ListenX processor. Listen-type processors start a server with one or more background threads, listening on a defined network port. Stateless NiFi-based connectors of this category are:

- HTTP Source
- Syslog TCP Source
- Syslog UDP Source

For example, the HTTP Source connector's flow uses a ListenHTTP NiFi processor to start an HTTP server listening on a configured port. Deploying this connector with two tasks would mean that two ListenHTTP processors are running simultaneously, both listening on the same network port. Since there is no user-defined mapping between tasks and Kafka Connect workers, it is possible that these two instances get deployed to the same worker. This means that two HTTP servers would be listening on the same port on the same cluster node. This results in port collision. To avoid issues like this, you must ensure that listen-type connectors are always deployed with exactly one task.

Primary node only connectors

Connectors in this category are connectors that use a dataflow with a NiFi processor that is configured to run on the primary node only if the flow was running on a regular NiFi cluster. For example, the QueryDatabaseTableRecord processor is a processor that is intended to run on primary nodes only. This processor is used in the JDBC Source connector. Stateless NiFi-based connectors of this category are:

- JDBC Source
- SFTP Source

For these connectors, the configured value of the `tasks.max` property is ignored. The connector is always started with exactly one task.

All other connectors

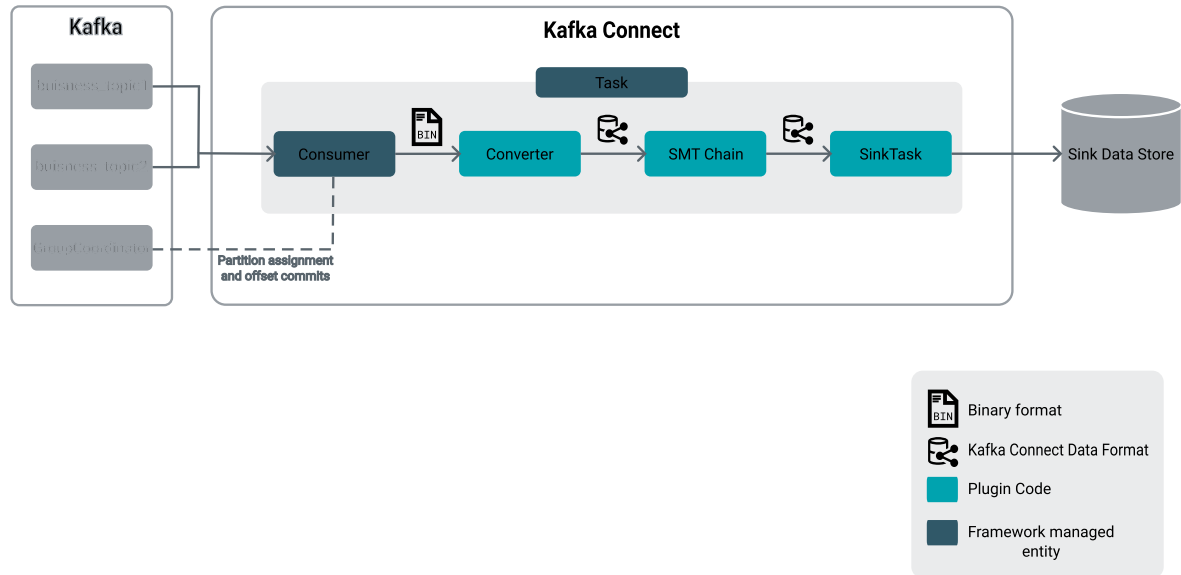
Connectors in this category are connectors that are neither listen-type or primary node only. Stateless NiFi-based connectors of this category are:

- JMS Source
- MQTT Source

There is no limitation to the number of tasks for these connectors. In the case of these connectors, it is the responsibility of the source data store system to make sure that one piece of data is only given to one of the connector's tasks

Sink Connectors

Figure 5: Sink Connector Architecture



Notice that the connector's task is made up of multiple components, these are as follows:

- Kafka Consumer

The consumer is responsible for fetching the messages from the Kafka topics and delivering them to the converter.

- SMT Chain

The Single Message Transforms (SMT) chain can modify and filter the data on a single message basis.

- Converter

The converter is responsible for converting Kafka message keys and values with a configured conversion mechanism. However, Stateless NiFi connectors currently only use keys of string type and values of byte array type. As a result, the converter does not perform any meaningful conversion.

- SinkTask

In Stateless NiFi connectors, the SinkTask is responsible for executing one instance of the connector's dataflow.

If you deploy a sink connector with more than one task, each task is going to have its own instance of the four task components.

With sink connectors, there is no limitation to using more than one task for a connector. When a connector is started with multiple tasks, each task has its own Kafka consumer, but these consumers belong to the same Kafka consumer group. Kafka makes sure that a message fetched by one consumer is not given to another consumer within the same consumer group.

There is no user-defined mapping between tasks and Kafka Connect workers, meaning that once deployed, the tasks of a connector may or may not be executed on the same worker.

Developing a dataflow for Stateless NiFi

Learn about the recommended process of building a dataflow that you can deploy with the Stateless NiFi Sink or Source connectors. This process involves building and designing a parameterized dataflow within a process group and then downloading the dataflow as a flow definition.

About this task



The general steps for building a dataflow are identical for both source and sink connector flows. For ease of understanding a dataflow example for a simple MQTT Source and MQTT Sink connector is provided without going into details (what processors to use, what parameters and properties to set, what relationships to define and so on).

Before you begin

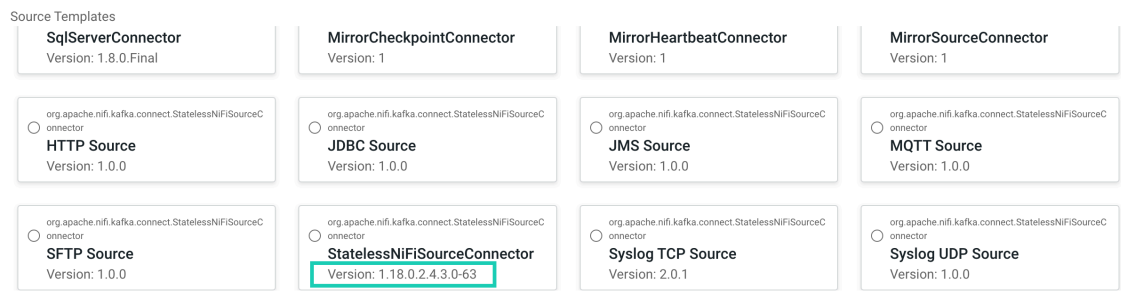
- Ensure that you reviewed [Dataflow development best practices for Stateless NiFi](#).
- You have access to a running instance of NiFi. The NiFi instance does not need to run on a cluster. A standalone instance running on **any machine**, like your own computer, is sufficient.
- Ensure that the version of your NiFi instance matches the versions of the Stateless NiFi plugin used by Kafka Connect on your cluster.

You can look up the Stateless NiFi plugin version either by using the Streams Messaging Manager (SMM) UI, or by logging into a Kafka Connect host and checking the Connect plugin directory. In addition, ensure that you note down the plugin version. You will need to manually edit the flow definition JSON and replace the NiFi version with the plugin version.

For Use SMM UI

1. Access the SMM UI, and click  **Connect** in the navigation sidebar.
2. Click the  New Connector option.
3. Locate the StatelessNiFiSourceConnector or StatelessNiFiSinkConnector cards. The version is located on the card.

The version is made up of multiple digits. The first three represent the NiFi version. For example, if the version on the card is 1.18.0.2.4.3.0-63, then you should use NiFi 1.18.0 to build your flow.

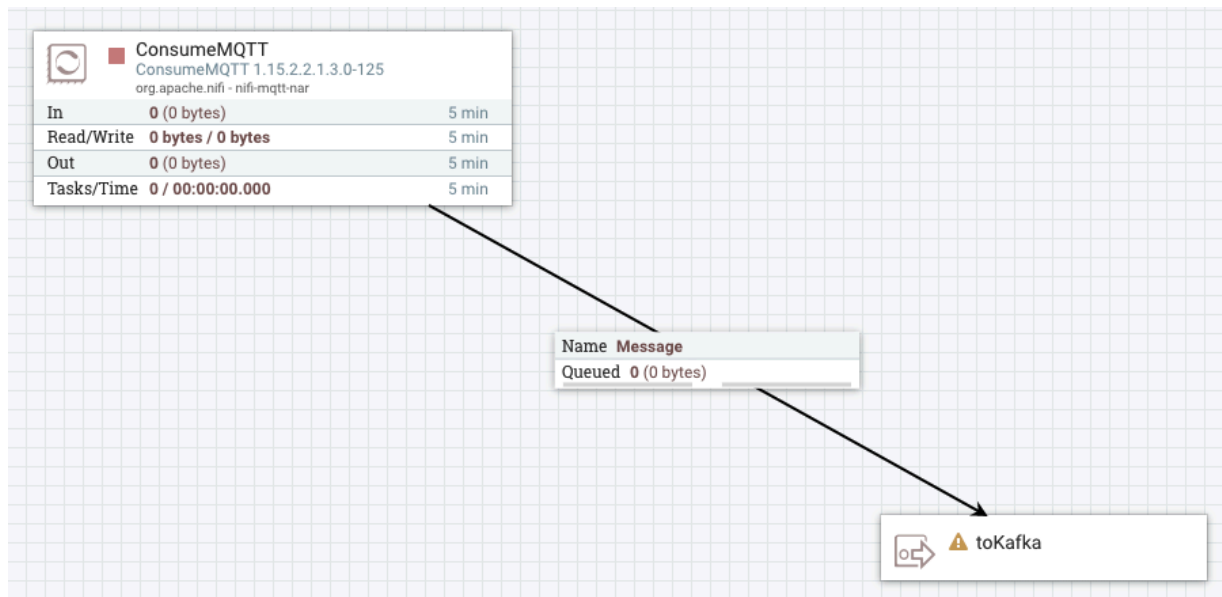


Procedure

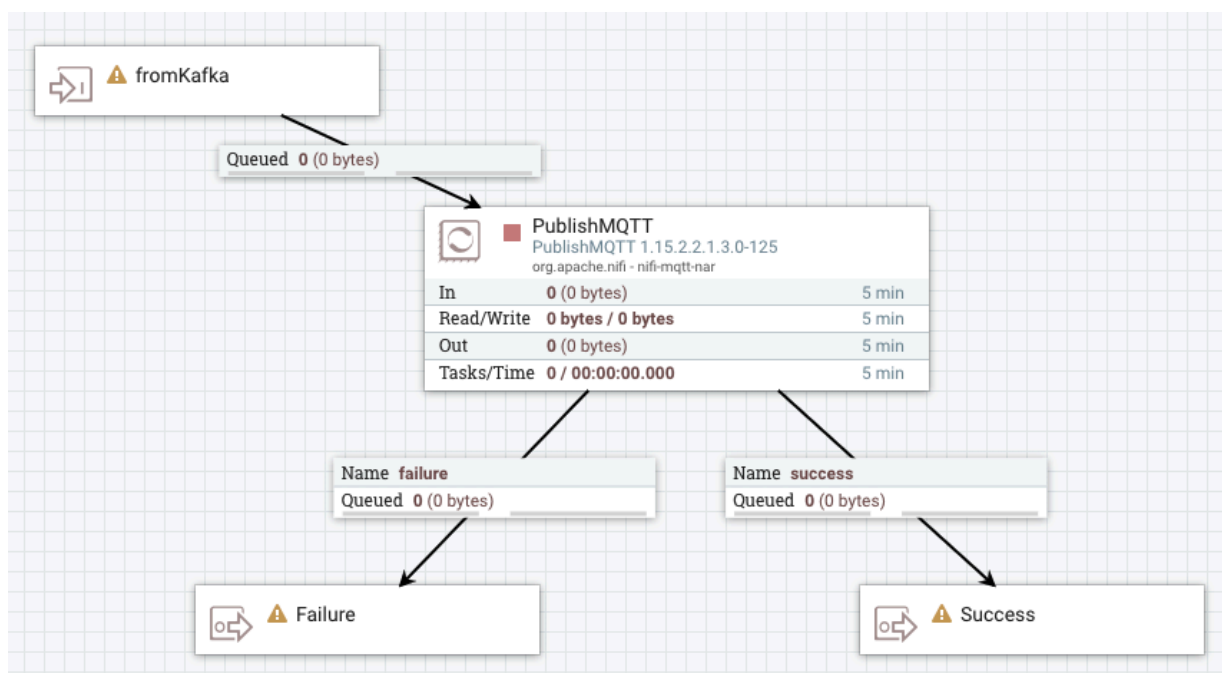
1. Access the NiFi UI.
2. Create a new process group.
3. Right-click the process group and select Configure.
4. Create a new parameter context and assign it to the process group:
 - a) Select Process Group Parameter Context Create new parameter context...
 - b) Enter a name for the parameter context.
 - c) Click Apply.
5. Click Apply and then click OK.
6. Close the configuration dialog.
7. Double-click the process group you created.

8. Design and parameterize your dataflow.

For example, you can build a dataflow that you deploy as a source connector. The following is an example dataflow for a simple MQTT source connector. It consists of a ConsumeMQTT processor and an output port (representing a Kafka topic).



Alternatively, you can also build a dataflow to deploy as a sink connector. The following is an example dataflow for a simple MQTT sink connector. It consists of an input port (representing a Kafka topic), a PublishMQTT processor, and two output ports.



In the case of this example, a failure output port is required. This is done so that in case the destination is not available, the session can be rolled back and the Kafka message can be declined. The failure port is specified in the configuration of the connector when deploying it through SMM.

When designing and building your dataflow, ensure that you update your parameter context and reference the parameters in your components. For example:

Keep in mind that sensitive properties in processors or controller services must always be parametrized. That is, they must get their values from a parameter in the process group's parameter context. This is because the values of sensitive properties are not exported with the process group, therefore, they must be configurable when deploying the connector.

Additionally, sensitive processor or controller service properties can only reference sensitive parameter context parameters. That is, the parameter must be marked as sensitive during creation.

For more information on building dataflows, parameters, and referencing parameters, see the *Flow Management* library.

9. Exit the process group.

10. Right-click the process group and select **Download flow definition Without external service**.

If you followed these steps, your process group is not referencing controller services that are outside of the process group. As a result, you can choose to download the flow without external services included.

11. Open the JSON file you downloaded in a text editor, and replace all occurrences of the NiFi version with the Stateless NiFi plugin version used in your cluster.

For example, if the NiFi version is 1.18.0 and your plugin version is 1.18.0.2.4.3.0-63, then every occurrence of 1.18.0 in the JSON must be replaced with 1.18.0.2.4.3.0-63.

Results

The dataflow is downloaded as a flow definition (JSON file).

What to do next

Deploy the dataflow as a Kafka Connect connector using the Stateless NiFi Source or Sink connectors. Continue with *Deploying a dataflow using stateless NiFi*.

Related Information

[Deploying a dataflow using Stateless NiFi](#)

[Flow Management](#)

Deploying a dataflow using Stateless NiFi

A custom developed NiFi dataflow can be deployed with the Stateless NiFi Source or Sink connectors using Streams Messaging Manager (SMM).

About this task


After building and downloading your dataflow, you can deploy it in Kafka Connect as a source or sink connector. This can be done using SMM. The following list of steps walk you through the process of deploying a Stateless NiFi Sink or Source connector and provide examples on how dataflows are configured using SMM.

These steps are for the SMM UI, however, the actions described here can also be completed using the SMM REST API. For more information, see *Cloudera Streams Messaging Manager REST API Reference*.

Before you begin

- A CDP cluster containing SMM and Kafka is available.
- You have access to the SMM UI.
- The flow definition of the dataflow that you want to deploy is available to you. For more information on how to design, build, and export a dataflow using the NiFi UI, see [Developing a dataflow for Stateless NiFi](#).
- The properties of Kafka consumers and producers used internally by the connector (Connect-managed clients) can be configured on a per-connector basis. The following steps give an example of how this can be done. However, in order for consumer and producer property overrides to take effect, you must ensure that the Connector Kafka Client Configuration Override Policy Kafka service property is set to All in Cloudera Manager. By default this property is set to None, which means that no overrides are allowed.

Procedure

1. Access the SMM UI.
2. Click  (Connect) on the navigation sidebar.
3. Click + New Connector to add a new connector.
4. Select a connector.

Which connector you select depends on the type of dataflow you want to run. Select `StatelessNiFiSourceConnector` found on the **Source Templates** tab if your dataflow collects (sources) data from systems and publishes it to Kafka. Select the `StatelessNiFiSinkConnector` from the **Sink Templates** tab if your dataflow delivers (sinks) Kafka data into other systems.

5. Enter a name for the connector in the Enter Name field.

Ensure that you add a unique and easily identifiable name. The name of the connector cannot be changed once the connector is deployed.

6. Configure the connector

Connector configuration can be broken into three steps. First, you configure the properties available in the default configuration template. Second, you configure other connector properties. Finally, you specify the flow parameters specific to your dataflow. These are the parameters that you set up in the parameter context when you built the flow in NiFi. Even though the following breaks the process up into multiple substeps, all properties are configured on the same page in the SMM UI.

- a) Configure the properties available in the default template.

The following properties from the default template are the ones that you should configure. Other properties in the template that are not highlighted here have working default values configured. Configuring the properties that are present in the template but not highlighted in the following list is not recommended by Cloudera.

flow.snapshot

Specifies the dataflow to run. You have multiple options on how to specify the value of this property. You can copy and paste the contents of the flow definition, upload the flow definition, or reference a file path. For more information on `flow.snapshot` configuration, see *Configuring flow.snapshot for Stateless NiFi connectors*.

input.port

The name of the input port in the NiFi dataflow that Kafka records are sent to by the connector's Kafka consumer. This property is specific to the Stateless NiFi Sink connector and is a port in NiFi terminology, not a network port. If the dataflow contains exactly one input port, this property is

optional and can be omitted. However, if the dataflow contains multiple Input Ports, this property must be specified.

output.port

The name of the output port in the NiFi dataflow that is used to deliver the messages from the dataflow to the connector's Kafka producer. This property is specific to the Stateless NiFi Source connector and is a port in NiFi terminology, and not a network port. If the dataflow contains exactly one port, this can be omitted. However, if the dataflow contains multiple ports (for example, a success and a failure port), this property must be specified. If any flowfile is sent to any port other than the specified port, it is considered as a failure. The session is rolled back and no data is collected.

topics

For source connectors, this is the name of the topic to deliver data to. For sink connectors, this can be a comma-separated list of topics to fetch data from.

krb5.file

Specifies the krb5.conf file to use if the dataflow interacts with any services that are secured using Kerberos. This property is optional and is only required if a connection is established with a Kerberized service. Use the default value if an appropriate krb5.conf is located in /etc.

b) Add other connector properties.

There are many other properties that you can set. These include the properties specific to the Stateless NiFi connectors as well as the Kafka Connect properties that come from the base framework. Which of these you require depends on your dataflow and use case.

In Stateless NiFi connectors, the Kafka message's key is always a string, its value is always a byte array. As a result, the key and value converters of the connector should not be modified. Use the default values, which are as follows:

```
"key.converter": "org.apache.kafka.connect.storage.StringConverter"  
"value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter"
```

A comprehensive list of all other properties that the connectors accept can be found in *Stateless NiFi Source connector reference*, *Stateless NiFi Sink Connector reference*, and the *Apache Kafka documentation*.

c) Configure dataflow-specific parameters.

If you followed the recommendations provided in *Dataflow development best practices for Stateless NiFi* and *Developing a dataflow for Stateless NiFi*, your dataflow is parameterized. The parameters that you defined within the dataflow can be assigned values using the configuration pane in SMM. This is done by adding a

configuration entry for each of the dataflow parameters that you want to configure. Dataflow parameters can be added to the configuration as follows:



Tip: Dataflow-specific parameters can also be automatically imported by clicking the Save and Enhance option when specifying the `flow.snapshot` property. For more information, see *Configuring flow.snapshot for Stateless NiFi connectors*.

```
"parameter.[***PARAMETER NAME***]": "[***VALUE***]"
```

For example, assume that you have a the following configuration entry:

```
"parameter.Directory": "/mydir"
```

In a case like this, any parameter context in the dataflow that has a parameter named `Directory` gets the specified value (`/mydir`). If the dataflow has child process groups, and those child process groups have their own parameter contexts, the value is used for all parameter contexts that contain a parameter named `Directory`.

You can also apply a parameter to a specific parameter context. This is done by prefixing the parameter name with the name of the parameter context followed by a colon.

```
"parameter.[***PARAMETER CONTEXT NAME***]:[***PARAMETER NAME***]": "[***VALUE***]"
```

For example, assume you have the following configuration entry:

```
"parameter.My Context:Directory": "/mydir"
```

In a case like this, only the parameter context called `My Context` gets the specified value for the `Directory` parameter.

Once all properties are configured, your configuration should look similar to the following example:

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
  "flow.snapshot": "[***FLOW DEFINITION JSON***]",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.MQTT Source Parameters:MQTT Broker URI": "tcp://[***HOST***]:[***PORT***]",
  "parameter.MQTT Source Parameters:MQTT Quality of Service": "0",
  "parameter.MQTT Source Parameters:MQTT Topics": "[***MQTT TOPIC NAME***]"
}
```

This example is for a custom built MQTT source connector. The example uses additional Connect properties such as `key.converter` and `value.converter`. The example also specifies a number of flow parameters such as `MQTT Broker URI` and `MQTT Topics`.

d) Configure Kafka producer and consumer properties.

Source connectors use internal Kafka producers to write message to topic, sink connectors use internal consumers to fetch messages from topics. The properties of these producers and consumers can be modified on

a per-connector basis. This is done by adding the required properties with the `producer.override` and `consumer.override` prefixes to the connector configuration.

For example, if you want to set the consumer of your sink connector to use `auto.offset.reset=earliest`, you need to add the following configuration entry to the configuration of the connector.

```
consumer.override.auto.offset.reset=earliest
```

7. Click Validate.

The validator displays any JSON errors in your configuration. Fix any errors that are displayed. If your JSON is valid, the **JSON is valid** message is displayed in the validator.

8. Click Next.

9. Review your connector configuration.

10. Click Deploy.

Related Information

[Cloudera Streams Messaging Manager REST API Reference](#)

[Stateless NiFi Source properties reference](#)

[Stateless NiFi Sink properties reference](#)

[Apache Kafka documentation](#)

[Dataflow development best practices for Stateless NiFi](#)

[Developing a dataflow for Stateless NiFi](#)

[Configuring flow.snapshot for Stateless NiFi connectors](#)

Downloading and viewing the predefined Stateless NiFi dataflows shipped in CDP

Learn how to download and view the predefined dataflows used by the Stateless NiFi-based connectors shipped in CDP. Predefined dataflows are downloaded from CDP cluster hosts. The downloaded files can be uploaded to NiFi to examine the dataflow's structure.

About this task

A number connectors shipped for use with CDP are based on Stateless NiFi and run predefined dataflows (dataflows developed by Cloudera). For example, the JDBC Source connector is a Stateless NiFi-based connector that runs a predefined dataflow.

To better understand how Stateless NiFi-based connectors like JDBC Source work, and how they move and manipulate Kafka data, it can be useful to view the dataflow that the connectors use.

These dataflows are installed on every cluster host that is running a Kafka Connect service role. You can fetch the dataflows from your cluster hosts, load them into NiFi, and look at their flow structure.

The flow definition JSON files of these dataflows are present on the cluster hosts at `/etc/kafka_connect_ext/flow-definitions`. However, you cannot simply load these flow definition files as is into NiFi. They must be edited first.

Before you begin

- You have access to a running CDP cluster that has a Kafka service and Kafka Connect service roles deployed.
- You have access to a running instance of NiFi. The NiFi instance does not need to run on a cluster. A standalone instance running on **any machine**, like your own computer, is sufficient.

Procedure

1. Fetch the flow definition of the chosen connector from the cluster.

The flow definition files are located on the cluster hosts at `/etc/kafka_connect_ext/flow-definitions`. How you complete this step depends on your cluster environment and what utilities are available to you. For example, if using `scp`, you can run the following command.

```
scp [***USER***]@[***CLUSTER_HOST***]:/etc/kafka_connect_ext/flow-definitions/[***CONNECTOR_NAME***]-[***VERSION***].json .
```

2. Open `.json` file in a text editor.
3. Copy contents of the "flow" JSON element to a separate `.json` file.
4. In the `.json` you created, replace the value of each "version" element with the version of NiFi you will use to view the dataflow.
5. In the NiFi UI, drag a new process group to the canvas.
- 6.



In the Add Process Group modal, click  and browse for your `.json` file.

7. Click Add.

Results

The dataflow is imported to the process group with all of its parameters.

What to do next

After the dataflow is imported, you can view the structure of the flow to gain a better understanding of how the dataflow operates. If you want to, you can make modifications and redeploy the modified dataflow as a custom connector.

Related Information

[Developing a dataflow for Stateless NiFi](#)

[Deploying a dataflow for Stateless NiFi](#)

Configuring `flow.snapshot` for Stateless NiFi connectors

Learn about the options you have in Streams Messaging Manager (SMM) when configuring the `flow.snapshot` Stateless NiFi connector property as well as the Save and Enhance feature that can be used to automatically import NiFi parameters into the connector configuration.



Important: Configuring the `flow.snapshot` property is only required when you are deploying your custom developed dataflows with the Stateless NiFi Source or Sink connectors. Cloudera developed connectors powered by Stateless NiFi do not require `flow.snapshot` to be set. If the connector you are deploying does not by default have a `flow.snapshot` entry in its configuration template in SMM, you do not need to configure it.

NiFi dataflows within Kafka Connect are deployed using the Stateless NiFi Source and Sink connectors. When you deploy your custom developed dataflows with the Stateless NiFi connectors, you must configure the `flow.snapshot` property. This property specifies the dataflow that is run within the connector.

Because the value of this property is the dataflow itself, the value you add to the property is not a simple string or integer. Instead, the value is the full contents of a flow definition JSON file, which depending on the dataflow, can be considerable in size.

As a result of this, the property is handled in a unique way by SMM. You can use multiple configuration methods to specify value. Additionally, you can choose to automatically enhance (amend) the configuration with the NiFi parameters that are available within the flow definition.

Configuration methods

You can use the following configuration methods when specifying a value for `flow.snapshot`. The following instructions assume that a flow definition JSON file is available to you and that you are logged in to the Connect tab of the SMM UI.



Note:

If the flow definition is larger than 64 KB, the following message appears:

```
Too large config, to view flow.snapshot click the edit button...
```

In a case like this, the dataflow is added, but you are not able to edit it directly in the **Properties** section.

Instead, you must click and edit or review the flow definition in the modal window that appears.

Copy and paste the contents of the flow definition

This method is similar to configuring any other connector property. The only difference is that the value you specify is considerably longer than other property value.

1. Open the flow definition in an editor.
2. Copy the full contents of the file.
3. Paste the contents of the file in the value field of the `flow.snapshot` property.

Upload the flow definition

This method involves uploading the flow definition from your own machine. This method is convenient if the flow definition JSON is available on the machine you are using to access the SMM UI.

1. Click next to `flow.snapshot`.
2. Click Browse... in the modal window that appears.
3. Select the flow definition file that you want to use.
4. Click Save or Save and Enhance.

Reference a file path

The `flow.snapshot` property accepts file paths. This means that instead of copy-pasting or uploading the full contents of the flow definition JSON, you can also specify a location on disk where the file is available. Cloudera recommends that you use this method if the flow definition is larger than 1 MB.

1. Deploy the flow definition file on all Kafka Connect hosts.



Important: Ensure that the file is deployed to all hosts at the same location and that the user that Kafka Connect is running as, which defaults to `kafka`, has read access to the file.

2. Set the file path as the value of the `flow.snapshot` property.

Using Save and Enhance

NiFi dataflows that you run within Kafka Connect are in the majority of cases parameterized. You can assign values to the parameters available within the dataflow on the SMM UI during connector configuration. This is done by adding each parameter, with specific prefixes, in the **Properties** section. However, this can take a considerable amount of time if done manually. This is especially true if the dataflow has many parameters and multiple parameter contexts.

In cases like this, you can choose to enhance the configuration by clicking the Save and Enhance option after a value for the `flow.snapshot` property is added. If this option is used, SMM parses the dataflow specified in `flow.snapshot`, extracts all parameters that are available in the dataflow, and adds them to properties of the connector. The Save and

Enhance option is available on the modal window that appears after you click next to the property.



Tutorial: developing and deploying a JDBC Source dataflow in Kafka Connect using Stateless NiFi

A step-by step tutorial that walks you through how you can create a JDBC Source dataflow and how to deploy the dataflow as a Kafka Connect connector using the Stateless NiFi Source connector. The connector/dataflow presented in this tutorial reads records from an Oracle database table and forwards them to Kafka in JSON format.

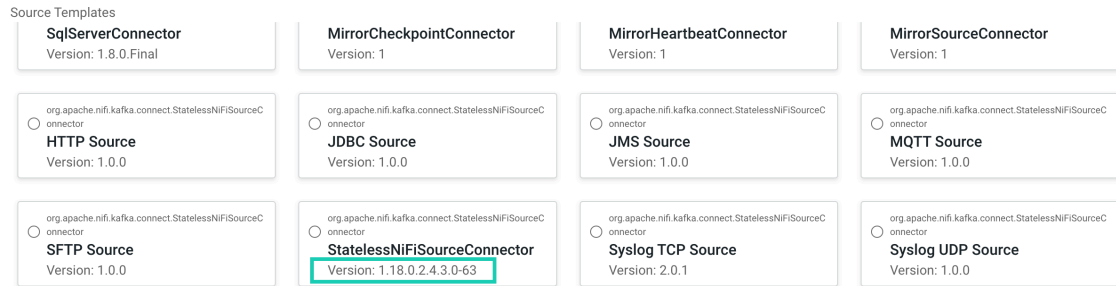
Before you begin

- Look up the Stateless NiFi plugin version either by using the Streams Messaging Manager (SMM) UI, or by logging into a Kafka Connect host and checking the Connect plugin directory.

For Use SMM UI

1. Access the SMM UI, and click  **Connect** in the navigation sidebar.
2. Click the  New Connector option.
3. Locate the StatelessNiFiSourceConnector or StatelessNiFiSinkConnector cards. The version is located on the card.

The version is made up of multiple digits. The first three represent the NiFi version. For example, if the version on the card is 1.18.0.2.4.3.0-63, then you should use NiFi 1.18.0 to build your flow.



- Download and start NiFi. You can download NiFi from <https://archive.apache.org/dist/nifi/>. This example uses NiFi 1.18.0 (nifi-1.18.0-bin.zip).
- The connector/dataflow developed in this tutorial requires the Oracle JDBC driver to function. Ensure that the driver JAR is deployed and available on every Kafka Connect host under the same path with correct file permissions. Note down the location where you deploy the driver, you will need to set the location as a property value during connector deployment. For example:

```
cp ./ojdbc8-[***VERSION***].jar /var/lib/kafka_connect_jdbc/oracle-connector-java.jar
```

```
chmod 644 /var/lib/kafka_connect_jdbc/oracle-connector-java.jar
```

Procedure

1. In NiFi, create a process group and give it a name.

The screenshot shows the NiFi web console interface. At the top, there is a navigation bar with the 'nifi' logo and various icons. Below the navigation bar, a status bar displays '0 / 0 bytes' and '14:50:06 CEST'. The main workspace is a grid where a process group box is visible. The process group is titled 'My JDBC Source' and contains a table with the following data:

My JDBC Source		
Queued	0 (0 bytes)	
In	0 (0 bytes) → 0	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 → 0 (0 bytes)	5 min

At the bottom of the interface, the text 'NiFi Flow' is visible.

2. Add a QueryDatabaseTableRecord processor and an output port to the process group.

Note down the name of the output port, you will need to set the name as the value of a connector property during connector deployment.

The screenshot shows the NiFi web console interface. At the top, there is a navigation bar with the 'nifi' logo and various icons. Below the navigation bar, a status bar displays '0 / 0 bytes' and '14:53:36 CEST'. The main workspace is a grid where a process group box is visible. The process group is titled 'My JDBC Source' and contains a table with the following data:

QueryDatabaseTableRecord		
QueryDatabaseTableRecord 1.18.0 org.apache.nifi - nifi-standard-nar		
In	0 (0 bytes)	5 min
Read/Write	0 bytes / 0 bytes	5 min
Out	0 (0 bytes)	5 min
Tasks/Time	0 / 00:00:00.000	5 min

Below the processor box, there is an output port box with the text 'Name success' and 'Queued 0 (0 bytes)'. An arrow points from the output port to a connector box titled 'Output to Kafka'.

At the bottom of the interface, the text 'NiFi Flow » My JDBC Source' is visible.

3. Add a DBCPConnectionPool and a JsonRecordSetWriter controller service to the process group.

The controller services are required for the QueryDatabaseTableRecord processor to access the database and to convert records.

My JDBC Source Configuration

GENERAL | CONTROLLER SERVICES

Name	Type	Bundle	State	Scope
DBCPConnectionP...	DBCPConnectionPool 1.18.0	org.apache.nifi - nif...	Invalid	My JDBC Source
JsonRecordSetWri...	JsonRecordSetWriter 1.18.0	org.apache.nifi - nif...	Disabled	My JDBC Source

Last updated: 14:57:04 CEST

Listed services are available to all descendant Processors and services of this Process Group.

NiFi Flow » My JDBC Source

4. Examine the properties of the QueryDatabaseTable processor and the two controller services. Collect the properties that you want to parameterize.

Parameterized properties can be set when deploying the connector. In this example, the following properties are parameterized:

- Database Type
- Table Name
- Maximum-value Columns
- Initial Load Strategy
- Database Connection URL
- Database Driver Class Name
- Database Driver Location
- Database User
- Database Password

5. Create a parameter context, give it a name, and add the above mentioned parameters to the parameter context.

The screenshot shows a dialog box titled "Add Parameter Context". It has three tabs: "SETTINGS", "PARAMETERS", and "INHERITANCE". The "PARAMETERS" tab is selected. The "Name" field contains the text "My JDBC Source Parameters" and has a "Referencing Components" icon to its right. The "Description" field is empty. At the bottom right, there are "CANCEL" and "APPLY" buttons.

You can assign values to the parameters in this step. The values you assign here serve as default values for the parameters when the flow is deployed as a connector. Assigning values is not mandatory as all parameters will

be configurable during connector deployment. However, assigning sensible default values is a good practice and makes the dataflow easier to deploy in Kafka Connect.

In this example, the connector will connect to an Oracle database. Therefore, values are specified for Database Driver Class Name and Database Type parameters. Additionally, Initial Load Strategy is set to Start at Beginning.

Ensure that you mark the Database Password parameter as a Sensitive Value during creation. This is because only sensitive parameters can be referenced from sensitive processor or controller service properties.

Add Parameter Context

SETTINGS PARAMETERS INHERITANCE

Name ^	Value	
Database Connection URL	No value set	
Database Driver Class Name	oracle.jdbc.driver.OracleDriver	
Database Driver Location	No value set	
Database Password	No value set	
Database Type	Oracle	
Database User	No value set	
Initial Load Strategy	Start at Beginning	
Maximum-value Columns	No value set	
Table Name	No value set	

Parameter Database Password

Referencing Components Pending Apply

CANCEL APPLY

6. Assign the parameter context to the process group.

My JDBC Source Configuration

GENERAL CONTROLLER SERVICES

Process Group Name
My JDBC Source

Process Group Parameter Context
My JDBC Source Parameters

Process Group Comments

Last updated: 15:21:35 CEST

NiFi Flow » My JDBC Source

7. Reference the parameters from the processor and the controller services.

The image shows two screenshots from the NiFi user interface. The top screenshot is the 'Configure Processor' dialog for the 'QueryDatabaseTableRecord' processor (version 1.18.0). The 'PROPERTIES' tab is active, showing a table of properties. The bottom screenshot is the 'Configure Controller Service' dialog for the 'DBCPConnectionPool' controller service (version 1.18.0). The 'PROPERTIES' tab is active, showing a table of properties.

Configure Processor: QueryDatabaseTableRecord 1.18.0

Property	Value
Database Connection Pooling Service	DBCPConnectionPool
Database Type	#{Database Type}
Table Name	#{Table Name}
Columns to Return	No value set
Additional WHERE clause	No value set
Custom Query	No value set
Record Writer	JsonRecordSetWriter
Maximum-value Columns	#{Maximum-value Columns}
Initial Load Strategy	#{Initial Load Strategy}
Max Wait Time	0 seconds
Fetch Size	0

Configure Controller Service: DBCPConnectionPool 1.18.0

Property	Value
Database Connection URL	#{Database Connection URL}
Database Driver Class Name	#{Database Driver Class Name}
Database Driver Location(s)	#{Database Driver Location}
Kerberos User Service	No value set
Kerberos Credentials Service	No value set
Kerberos Principal	No value set
Kerberos Password	No value set
Database User	#{Database User}
Password	Sensitive value set
Max Wait Time	500 millis
Max Total Connections	8
Validation query	No value set

At this point, one controller service is invalid, and the other is disabled, but that is expected. Neither processors, nor controller services need to be valid, started, or enabled at this point.

- Open the configuration of the QueryDatabaseTableRecord processor, and go to the SCHEDULING tab.

The screenshot shows the 'Configure Processor' dialog for the 'QueryDatabaseTableRecord 1.18.0' processor. The 'SCHEDULING' tab is selected. The 'Execution' dropdown menu is highlighted with a red box and is set to 'Primary node'. Other settings visible include 'Scheduling Strategy' set to 'Timer driven', 'Concurrent Tasks' set to '1', and 'Run Schedule' set to '0 sec'. The 'INVALID' status is shown at the top left. At the bottom right, there are 'CANCEL' and 'APPLY' buttons.

Notice that Execution is set to Primary node (which is the only option for this particular processor). For Stateless NiFi Kafka connectors, this setting means that the connector is going to be deployed with one task, regardless of the value of the tasks.max connector property.

- At this point, your dataflow is in the following state.

The screenshot shows the NiFi dataflow canvas. The 'QueryDatabaseTableRecord' processor is highlighted, showing its status and metrics. Below it, the 'Name success' processor is connected to the 'Output to Kafka' processor. The dataflow is not running, and there are warnings present. The status bar at the top shows '0 / 0 bytes' and '2' warnings. The bottom status bar shows 'NiFi Flow » My JDBC Source'.

Processor	In	Read/Write	Out	Tasks/Time
QueryDatabaseTableRecord	0 (0 bytes)	0 bytes / 0 bytes	0 (0 bytes)	0 / 00:00:00.000

Notice how the dataflow is not running, and there are warnings present. This is expected. The dataflow does not need to run in order for you to be able to download the flow definition. The properties you must configure to make the dataflow function will be set later on during connector deployment.

10. Right-click the process group and select **Download flow definition Without external service**.


If you followed these steps, your process group is not referencing controller services that are outside of the process group. As a result, you can choose to download the flow without external services included.

11. Open the downloaded JSON file, and replace **every occurrence** of the version number (1.18.0) with the version number of your Stateless NiFi plugin.

```
{
  "flowContents": {
    "identifier": "7f4e524d-e155-3d50-af2a-d31a49e87a50",
    "instanceIdentifier": "a3db7014-0187-1000-0e07-ae819736fda3",
    "name": "My JDBC Source",
    "comments": "",
    "position": {
      "x": 160.0,
      "y": 48.0
    },
    "processGroups": [],
    "remoteProcessGroups": [],
    "processors": [
      {
        "identifier": "d19b2cf9-f9ef-39c4-8591-51147e432fcd",
        "instanceIdentifier": "a3df4caa-0187-1000-9066-5fa475f8206d",
        "name": "QueryDatabaseTableRecord",
        "comments": "",
        "position": {
          "x": 64.0,
          "y": -56.0
        },
        "type": "org.apache.nifi.processors.standard.QueryDatabaseTableRecord",
        "bundle": {
          "group": "org.apache.nifi",
          "artifact": "nifi-standard-nar",
          "version": "1.18.0"
        },
        "properties": {
          "Table Name": "#{Table Name}"
        }
      }
    ]
  }
}
```

Replace with Stateless NiFi plugin version

12. Access the SMM UI.

13. Click  **Connect** in the navigation sidebar.

14. Click the  New Connector option.

15. Select StatelessNiFiSourceConnector from the available connectors.

The Connector Configuration page appears.

16. Configure the following connector properties.

- name=My Custom JDBC Connector
- output.port=Output to Kafka



This property must be set to the name of the output port in the dataflow. The output port name was specified in Step 2 on page 87.

- topics=jdbc_example

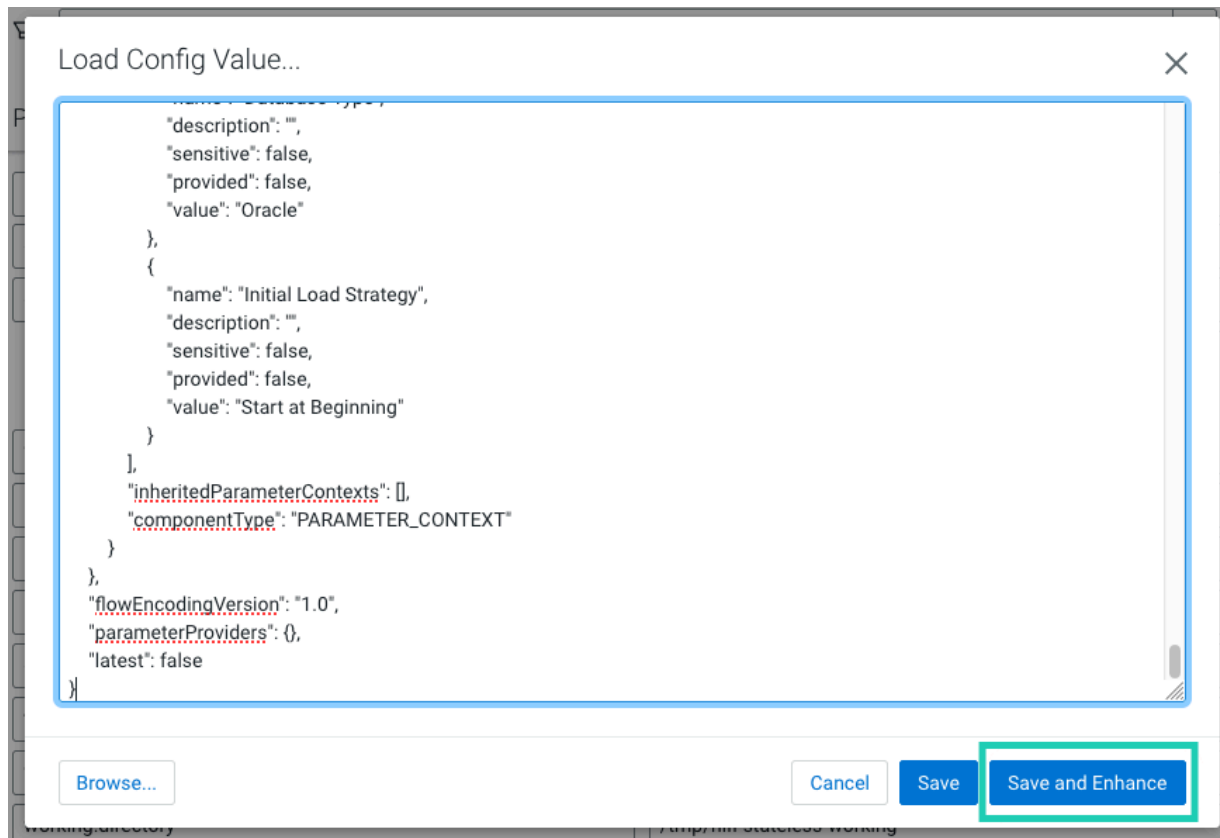
This is the name of the target topic that the connector writes data to.

17. Configure the flow.snapshot property.

The flow.snapshot property specifies the dataflow that the connector runs. This is the dataflow that you developed and exported as flow definition JSON in the previous steps. While you have multiple options when configuring this property, this tutorial uses the Import and Enhance option in SMM.

- a) Click  found next to flow.snapshot, and then click  Edit.
- b) Paste the contents of the flow definition JSON, or click Browse... to upload the file from your machine.
- c) Click Save and Enhance.

Save and Enhance is a unique SMM feature only available for this property. If this option is used, SMM parses the dataflow specified in flow.snapshot, extracts all parameters that are available in the dataflow, and adds them to connector configuration.

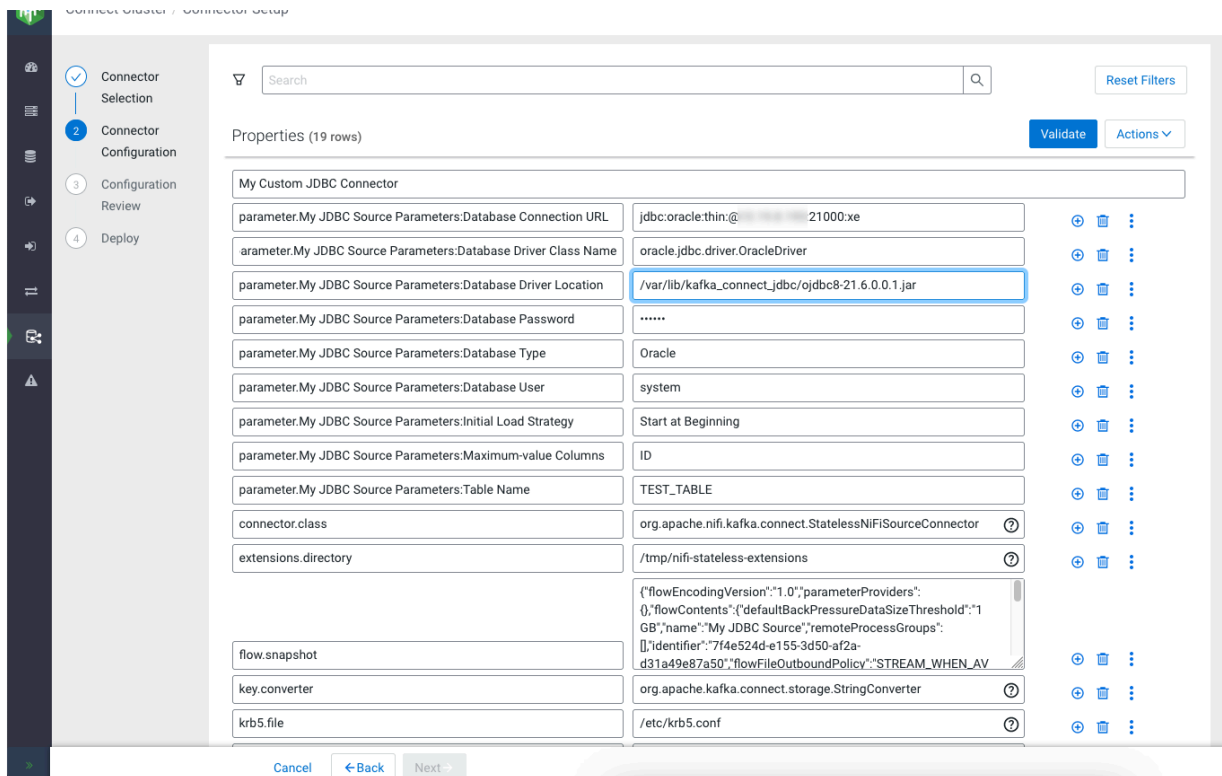


18. Configure the dataflow parameters.

These are the properties that have the parameter. prefix. In this specific example, values for Database Class Name, Database Type, and Initial Load Strategy are loaded automatically because values were specified within the dataflow in Step 5 on page 89. If needed, you can change the values.

Set Database Driver Location to the path where you deployed the JDBC driver.

19. Review your configuration.



Connector Selection

Connector Configuration

Configuration Review

Deploy

Search

Reset Filters

Validate

Actions

Properties (19 rows)

My Custom JDBC Connector

parameter.My JDBC Source Parameters:Database Connection URL	jdbc:oracle:thin:@...:21000:xe	
parameter.My JDBC Source Parameters:Database Driver Class Name	oracle.jdbc.driver.OracleDriver	
parameter.My JDBC Source Parameters:Database Driver Location	/var/lib/kafka_connect_jdbc/ojdbc8-21.6.0.0.1.jar	
parameter.My JDBC Source Parameters:Database Password	
parameter.My JDBC Source Parameters:Database Type	Oracle	
parameter.My JDBC Source Parameters:Database User	system	
parameter.My JDBC Source Parameters:Initial Load Strategy	Start at Beginning	
parameter.My JDBC Source Parameters:Maximum-value Columns	ID	
parameter.My JDBC Source Parameters:Table Name	TEST_TABLE	
connector.class	org.apache.nifi.kafka.connect.StatelessNIFISourceConnector	?
extensions.directory	/tmp/nifi-stateless-extensions	?
flow.snapshot	{ "flowEncodingVersion": "1.0", "parameterProviders": {}, "flowContents": { "defaultBackPressureDataSizeThreshold": "1 GB", "name": "My JDBC Source", "remoteProcessGroups": [], "identifier": "7f4e524d-e155-3d50-af2a-d31a49e87a50", "flowFileOutboundPolicy": "STREAM_WHEN_AV" } }	
key.converter	org.apache.kafka.connect.storage.StringConverter	?
krb5.file	/etc/krb5.conf	?

Cancel

← Back

Next →

20. Click Validate.

21. Click Next.


22. Click Deploy to deploy the connector.

Connector deployment is not instantaneous. You might need to wait a few seconds for deployment to finish.

Results

Your custom developed flow is deployed in Kafka Connect and is running as a Kafka Connect connector.

What to do next

Once deployment finishes, navigate to the  Connector Profile page of the connector. On this page you can view various details regarding the connector. Most importantly, you can view what the connector's status is, and which workers the connector and its task were assigned to.



Note: It might take some time for tasks to appear on this page.

Connect Cluster / Connector Profile Cluster: KAFKA-1

My Custom JDBC Connector Pause Resume Restart Delete New Connector

Connector Profile Connector Settings

Connector Profile

CLASSNAME: org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector
 ASSIGNED WORKER: -1. :28083

STATUS: RUNNING
 TOTAL TASKS: 1
 RUNNING TASKS: 1
 FAILED TASKS: 0
 PAUSED TASKS: 0

Tasks

Search by host Restart

Status	Worker ID	Task ID	Poll Batch Avg Time	Source Record Write Rate	Source Record Write Total
+	✓	-3.	0	516	0

< 1 >

In this particular case, the connector was assigned to the worker running on cluster host 1, whereas the task was assigned to worker running on cluster host 3. This means that dataflow log entries will be present in the Kafka Connect log file on host 3.

In addition to monitoring the connector, you can also check the contents of the topic that the connector is writing to.

This can be done in SMM by navigating to the Topics page, searching for the topic, and clicking . This opens the SMM Data Explorer in a modal window. You can use the Data Explorer to sample the data that is flowing through the topic. The database table's records appear in the topic in JSON format.

Related Information

[Kafka Connect tasks and Stateless NiFi](#)

Syslog TCP Source connector

The Syslog TCP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The Syslog TCP Source connector listens on a port for syslog messages over TCP and transfers them to Kafka. The connector accepts messages in one of the following formats: Syslog 3164, Syslog 5424, or Grok. If the input messages are in Grok format, the connector can either derive the schema using the field names from the value of the Grok Expression property or read the schema from Schema Registry.

The connector can write messages into Kafka in one of the following formats: Avro, JSON, or text. If the output format is text, the raw message is transferred to Kafka. If the output format is JSON, the input message is processed and converted into JSON format before it is transferred to Kafka. If the output format is Avro, the input message is processed and converted into Avro format before it is transferred to Kafka. The record schema is embedded into Avro messages if the schema is not fetched from Schema Registry. If Schema Registry is used for getting the schema of the records, then the schema does not get embedded in the Avro messages.

If Schema Registry is used, and it is on a Kerberized cluster, the `krb5.file` property must point to the `krb5.conf` file that provides access to the cluster on which Schema Registry is present. This means that the `krb5.conf` file must be on the same cluster node that the connector runs on. The Kerberos keytab that is used to access Schema Registry must also be on the same cluster node that the connector runs on.

The connection to Schema Registry can be secured by TLS. The truststore file necessary for securing the connection must be on the same cluster node that the connector runs on. Mutual TLS for securing the communication between the message sender and the connector itself is also supported. The keystore and truststore files necessary for securing the connection must be on the same cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *Syslog TCP Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- The Syslog TCP Source connector must use at least one-way SSL. It cannot be used without SSL.
- The Schema Registry URL property is mandatory even if Schema Registry is not used. If Schema Registry is not used, use the default value, or completely remove the property from the configuration JSON.
- Schemas are only read from Schema Registry if Input Data Format is set to GROK and Schema Access Strategy is set to Schema Registry.
- If Output Format is AVRO, the schema of the records can be embedded in the output data. Whether the schema is embedded is determined by the Schema Access Strategy property. If this property is set to Schema Registry, the schema is not embedded in the output messages. The schema is embedded in all other cases.
- Schema Branch and Schema Version can not be specified at the same time.

Configuration example

In this example, the connector uses mutual TLS to receive data in Syslog 3164 format which is then transferred to Kafka in JSON format.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
  "meta.smm.predefined.flow.name": "Syslog TCP Source",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.Syslog TCP Source Parameters:Port": "[***PORT***]",
  "parameter.Syslog TCP Source Parameters:Input Data Format": "Syslog 3164",
  "parameter.Syslog TCP Source Parameters:Output Format": "JSON",
  "parameter.Syslog TCP Source Parameters:Client Authentication": "REQUIRED",
  "parameter.Syslog TCP Source Parameters:SSL Keystore Filename": "[***THE FULLY-QUALIFIED FILENAME OF THE KEYSTORE***]",
  "parameter.Syslog TCP Source Parameters:SSL Keystore Key Password": "[***KEYSTORE KEY PASSWORD***]",
  "parameter.Syslog TCP Source Parameters:SSL Keystore Password": "[***KEYSTORE PASSWORD***]",
  "parameter.Syslog TCP Source Parameters:SSL Keystore Type": "[***KEYSTORE TYPE***]",
  "parameter.Syslog TCP Source Parameters:SSL Truststore Filename": "[***THE FULLY-QUALIFIED FILENAME OF THE TRUSTSTORE***]",
  "parameter.Syslog TCP Source Parameters:SSL Truststore Password": "[***TRUSTSTORE PASSWORD***]",
  "parameter.Syslog TCP Source Parameters:SSL Truststore Type": "[***TRUSTSTORE TYPE***]"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic that the connector sends messages to.

Port

The port that the connector listens on for incoming messages.

Input Data Format

Determines what format incoming messages are expected in.

Output Format

Determines the format in which messages are transferred to Kafka.

Client Authentication

Determines if one-way or two-way SSL is used. In this example, this property is set to **REQUIRED**, meaning that two-way SSL is used.

Keystore *

These are the properties for accessing the keystore containing the keypair used for secure communication.

Truststore *

These are the parameters for accessing the truststore containing the message sender's certificate used for secure communication.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[Syslog TCP Source properties reference](#)

Syslog UDP Source connector

The Syslog UDP Source connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The Syslog UDP Source connector listens on a port for syslog messages over UDP and transfers them to Kafka. The connector accepts messages in one of the following formats: Syslog 3164, Syslog 5424, or Grok. If the input messages are in Grok format, the connector can either derive the schema using the field names from the value of the Grok Expression property or read the schema from Schema Registry.

The connector can write messages into Kafka in one of the following formats: Avro, JSON, or text. If the output format is text, the raw message is transferred to Kafka. If the output format is JSON, the input message is processed and converted into JSON format before it is transferred to Kafka. If the output format is Avro, the input message is processed and converted into Avro format before it is transferred to Kafka. The record schema is embedded into Avro messages if the schema is not fetched from Schema Registry. If Schema Registry is used for getting the schema of the records, then the schema does not get embedded in the Avro messages.

If Schema Registry is used, and it is on a Kerberized cluster, the `krb5.file` property must point to the `krb5.conf` file that provides access to the cluster on which Schema Registry is present. This means that the `krb5.conf` file must be on the same cluster node that the connector runs on.

The Kerberos keytab that is used to access Schema Registry must also be on the same cluster node that the connector runs on. The connection to Schema Registry can be secured by TLS. The truststore file necessary for securing the connection must also be on the same cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Source properties

These are the properties that are specific to the Stateless NiFi Source connector. All Stateless NiFi Source connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Source property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *Syslog UDP Source properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- The Schema Registry URL property is mandatory even if Schema Registry is not used. If Schema Registry is not used, use the default value, or completely remove the property from the configuration JSON.
- Schemas are only read from Schema Registry if Input Data Format is set to GROK and Schema Access Strategy is set to Schema Registry.
- If Output Format is AVRO, the schema of the records can be embedded in the output data. Whether the schema is embedded is determined by the Schema Access Strategy property. If this property is set to Schema Registry, the schema is not embedded in the output messages. The schema is embedded in all other cases.
- Schema Branch and Schema Version can not be specified at the same time.

Configuration example

In this example, the connector receives data in Syslog 3164 format and transfers it to Kafka in JSON format.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSourceConnector",
  "meta.smm.predefined.flow.name": "Syslog UDP Source",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.Syslog UDP Source Parameters:Port": "[***PORT***]",
  "parameter.Syslog UDP Source Parameters:Input Data Format": "Syslog 3164",
  "parameter.Syslog UDP Source Parameters:Output Format": "JSON",
  "parameter.Syslog UDP Source Parameters:Output Grouping for JSON": "output-online"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic that the connector sends messages to.

Port

The port that the connector listens on for incoming messages.

Input Data Format

Determines what format incoming messages are expected in.

Output Format

Determines the format in which messages are transferred to Kafka.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Source properties reference](#)

[Syslog UDP Source properties reference](#)

ADLS Sink connector

The ADLS Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka connect framework. Learn about the connector, its properties, and configuration.

The ADLS Sink connector fetches messages from Kafka and uploads them to ADLS. The topic this connector receives messages from is determined by the value of the topics property in the configuration. The messages can contain unstructured (character or binary) data or they can be in Avro or JSON format.

If the input is unstructured data, record processing is disabled. In a case like this, multiple messages can be concatenated into a single output file on ADLS using a demarcator (for example, newline for text messages). Merging is optional, large binary data can be forwarded to Azure at a 1:1 ratio, one Kafka message equals a single ADLS file.

If the input is either Avro or JSON, record processing is enabled. In a case like this, the schema of the records can be a predefined schema retrieved from Schema Registry (for both Avro and JSON data), it can be embedded in the Avro data, or inferred from the JSON data. The output can be Avro, JSON or CSV. Multiple records are typically merged into a single output file before uploading to ADLS.

The connector can authenticate to Azure using an Account Key, SAS Token, Service Principal, or a Managed Identity.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *ADLS Sink properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.

Configuration example for fetching unstructured data

In this example, the connector fetches unstructured messages containing single line character data from Kafka. The connector concatenates the messages using newline characters as the demarcator. The files that the connector creates and uploads to ADLS are maximum 10 MB in size. The files are named according to the following pattern: `messages_***UUID***`. The connector authenticates to Azure using an Account Key.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "ADLS Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "PutAzureDataLakeStorage Failure",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.ADLS Sink Parameters:ADLS Account Name": "[***ACCOUNT NAME***]",
  "parameter.ADLS Sink Parameters:ADLS Account Key": "[***ACCOUNT KEY***]",
  "parameter.ADLS Sink Parameters:ADLS Filesystem Name": "[***FILESYSTEM NAME***]",
  "parameter.ADLS Sink Parameters:ADLS Directory Name": "[***DIRECTORY NAME***]",
  "parameter.ADLS Sink Parameters:Maximum File Size": "10 MB",
  "parameter.ADLS Sink Parameters:Kafka Message Data Format": "Raw",
  "parameter.ADLS Sink Parameters:Output File Demarcator": "\n",
  "parameter.ADLS Sink Parameters:Output Filename Pattern": "messages_${filename.uuid}"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics The name of the Kafka topic the connector fetches messages from.

ADLS Account Name

The Storage Account Name to use for authentication to Azure. This is also the target account where the output file is uploaded.

ADLS Account Key

The Storage Account Key to use for authentication to Azure.

ADLS Filesystem Name

The filesystem (or container) in the Storage Account to upload data to.

ADLS Directory Name

The target directory in the filesystem.

Maximum File Size

The maximum size of the output data file. In this example, the maximum size is 10 MB.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to Raw. This means that the connector expects unstructured text or binary data.

Output File Demarcator

Specifies the character sequence for demarcating (delimiting) message boundaries when multiple Kafka messages are ingested into an output file as raw messages. In this example, the property is set to `\n` (newline). This means that the newline character is used to separate the single line character data of the Kafka messages.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to `messages_{filename.uuid}`.

`messages_` is a string literal. `{filename.uuid}` is an expression that inserts a generated UUID in the filename. As a result, the files are named according to the following pattern: `messages_{***UUID***}`.

Configuration example for fetching JSON data

In this example, the connector fetches messages in JSON format from Kafka. The connector parses the JSON records and converts them to Avro format using the schema inferred from the JSON data. The schema is also embedded in the output file. The connector merges multiple records (all messages that are available for a single execution of the connector) into an Avro file, generates a name for the file using the following pattern `data_{***TIMESTAMP***}_{***SEQUENCE***}.avro`, and uploads the file to ADLS. The connector authenticates to Azure using a Service Principal.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "ADLS Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "PutAzureDataLakeStorage Failure",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.ADLS Sink Parameters:ADLS Account Name": "[***ACCOUNT NAME***]",
  "parameter.ADLS Sink Parameters:Azure Service Principal Tenant ID": "[***SERVICE PRINCIPAL TENANT ID***]",
  "parameter.ADLS Sink Parameters:Azure Service Principal Client ID": "[***SERVICE PRINCIPAL CLIENT ID***]",
  "parameter.ADLS Sink Parameters:Azure Service Principal Client Secret": "[***SERVICE PRINCIPAL CLIENT SECRET***]",
  "parameter.ADLS Sink Parameters:ADLS Filesystem Name": "[***FILESYSTEM NAME***]",
  "parameter.ADLS Sink Parameters:ADLS Directory Name": "[***DIRECTORY NAME***]",
  "parameter.ADLS Sink Parameters:Kafka Message Data Format": "JSON",
  "parameter.ADLS Sink Parameters:Schema Access Strategy": "Infer Schema",
  "parameter.ADLS Sink Parameters:Schema Write Strategy": "Embed Avro Schema",
  "parameter.ADLS Sink Parameters:Output File Data Format": "Avro",
  "parameter.ADLS Sink Parameters:Output Filename Pattern": "data_{filename.timestamp}_{filename.sequence}.avro",
  "parameter.ADLS Sink Parameters:Output Filename Timestamp Format": "yyyyMMdd_HHmss_SSS"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

ADLS Account Name

The Storage Account Name to use for authentication to Azure. This is also the target account where the output file is uploaded.

Service Principal *

The properties of the Service Principal to use for authentication to Azure.

ADLS Filesystem Name

The filesystem (or container) in the Storage Account to upload data to.

ADLS Directory Name

The target directory in the filesystem.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to JSON. This means that the connector expects JSON data.

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, this property is set to Infer Schema, meaning that the schema is determined (inferred) from the JSON data.

Schema Write Strategy

Specifies whether the record schema is written to the output data file. In this example, this property is set to Embed Avro Schema, meaning that the schema is embedded in the output Avro file.

Output File Data Format

Specifies the format of the records written to the output file. In this example, this property is set to Avro, meaning that the output file format is Avro.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to `data_${filename.timestamp}_${filename.sequence}.avro`.

`data_`, the underscore (`_`), and `.avro` are string literals. `${filename.timestamp}` and `${filename.sequence}` are expressions. `${filename.timestamp}` inserts the current timestamp in the filename. `${filename.sequence}` inserts an incrementing sequence value in the filename. As a result, the files are named according to the following pattern: `data_[***TIMESTAMP***]_[***SEQUENCE***].avro`.

Output Filename Timestamp Format

Specifies the timestamp format used in the `${filename.timestamp}` expression. The expression is used when generating the output filename.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[ADLS Sink properties reference](#)

Amazon S3 Sink

Learn more about the Amazon S3 Sink connector

The Amazon S3 Sink connector allows users to stream Kafka data into S3 buckets.

Configuration example

A simple configuration example for the Amazon S3 Sink connector.

The following is a simple configuration example for the Amazon S3 Sink connector. Short descriptions of the properties set in this example are also provided. For a full properties reference, see the *Amazon S3 Sink properties reference*.

```
{
  "aws.s3.bucket": "bring-me-the-bucket",
  "aws.s3.service_endpoint": "http://myendpoint:9090/",
  "aws.access_key_id": "EXAMPLEID",
  "aws.secret_access_key": "EXAMPLEKEY",
  "connector.class": "com.cloudera.dim.kafka.connect.s3.S3SinkConnector",
  "tasks.max": 1,
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
  "value.converter.passthrough.enabled": true,
  "value.converter.schema.registry.url": "http://schema-registry:9090/api/v1",
  "topics": "avro_topic",
  "output.storage": "com.cloudera.dim.kafka.connect.s3.S3PartitionStorage",
  "output.writer": "com.cloudera.dim.kafka.connect.partition.writers.avro.AvroPartitionWriter",
  "output.avro.passthrough.enabled": true
}
```

aws.s3.bucket

Target S3 bucket name.

aws.s3.service_endpoint

Target S3 host and port.

aws.access_key_id

The AWS secret key ID used for authentication.

aws.secret_access_key

The AWS secret access key used for authentication.

connector.class

Class name of the Amazon S3 Sink connector.

tasks.max

Maximum number of tasks.

key.converter

The converter capable of understanding the data format of the key of each record on this topic.

value.converter

The converter capable of understanding the data format of the value of each record on this topic.



Note: When the AvroConverter is used, you can specify Schema Registry properties to be used by the AvroConverter's Schema Registry client. This is done by adding the required Schema Registry property as a suffix to the value.converter property. For example, value.converter.schema.registry.url. Properties defined this way are passed on to the Schema Registry client used by the AvroConverter.

value.converter.passthrough.enabled

This property controls whether or not data is converted into the Kafka Connect intermediate data format before writing into an output file. Because in this example the input and output format is the same, the property is set to true, that is, data is not converted.

value.converter.schema.registry.url

The URL to Schema Registry. This is a mandatory property if the topic has records encoded in Avro format.

topics

List of topics to consume data from.

output.storage

The S3 storage implementation class.

output.writer

Determines the output file format. Because in this example the output format is Avro, AvroPartitionWriter is used.

output.avro.passthrough.enabled

This property has to match the configuration of the value.converter.passthrough.enabled property because both the input and output formats are Avro.

Related Information

[Amazon S3 Sink Connector Properties Reference](#)

HDFS Sink

Learn more about the HDFS Sink connector.

The HDFS Sink connector can be used to transfer data from Kafka topics to files on HDFS clusters. Each partition of every topic results in a collection of files named in the following pattern:

```
{topic name}_{partition number}_{end_offset}.{file extension}
```

For example, running the HDFS Sink connector on partition 0 of a topic named sourceTopic can yield the following series of files:

```
sourceTopic_0_50.avro - for record 0 ~ 50
sourceTopic_0_79.avro - holding record 51 ~ 79
...
```

The HDFS Sink connector periodically commits records to final result files. Each commit results in a separate "chunk" file.

Configuration example for writing data to HDFS

A simple configuration example for the HDFS Sink connector.

The following is a simple configuration example for the HDFS Sink connector. Short descriptions of the properties set in this example are also provided. For a full properties reference, see the *HDFS Sink properties reference*.

```
{
  "connector.class": "com.cloudera.dim.kafka.connect.hdfs.HdfsSinkConnector",
  "tasks.max": 1,
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
  "value.converter.passthrough.enabled": true,
  "value.converter.schema.registry.url": "http://localhost:9090/api/v1",
  "topics": "avro_topic",
```

```

"hdfs.uri": "hdfs://my-host.my-realm.com:8020",
"hdfs.output": "/topics_output/",
"output.writer": "com.cloudera.dim.kafka.connect.partition.writers.avro
.AvroPartitionWriter",
"output.avro.passthrough.enabled": true,
"hdfs.kerberos.authentication": true,
"hdfs.kerberos.user.principal": "user_account@MY-REALM.COM",
"hdfs.kerberos.keytab.path": "/path/to/user_account.keytab",
"hdfs.kerberos.namenode.principal": "hdfs/_HOST@MY-REALM.COM",
"hadoop.conf.path": "/etc/hadoop/"
}

```

connector.class

Class name of the HDFS Sink connector.

key.converter

The converter capable of understanding the data format of the key of each record on this topic.

value.converter

The converter capable of understanding the data format of the value of each record on this topic.



Note: When the AvroConverter is used, you can specify Schema Registry properties to be used by the AvroConverter's Schema Registry client. This is done by adding the required Schema Registry property as a suffix to the value.converter property. For example, value.converter.schema.registry.url. Properties defined this way are passed on to the Schema Registry client used by the AvroConverter.

value.converter.passthrough.enabled

This property controls whether or not data is converted into the Kafka Connect intermediate data format before writing into an output file. Because in this example the input and output format is the same, the property is set to true, that is, data is not converted.

value.converter.schema.registry.url

The URL to Schema Registry. This is a mandatory property if the topic has records encoded in Avro format.

topics

List of topics to consume data from.

hdfs.uri

The URI to the namenode of the HDFS cluster.

hdfs.output

The destination folder on the HDFS cluster where output files will reside.

output.writer

Determines the output file format. Because in this example the output format is Avro, AvroPartitionWriter is used.

output.avro.passthrough.enabled

This property has to match the configuration of the value.converter.passthrough.enabled property because both the input and output formats are Avro.

hdfs.kerberos.authentication

Enables or disables kerberos authentication.

hdfs.kerberos.user.principal

The user principal that the Kafka Connect role will use.

hdfs.kerberos.keytab.path

The path to the kerberos keytab file.

hdfs.kerberos.namenode.principal

The Kerberos principal used by the namenode. This is necessary when the HDFS cluster has data encryption turned on.

hadoop.conf.path

The path to the hadoop configuration files. This is necessary when the HDFS cluster has data encryption turned on.

Related Information

[HDFS Sink Connector Properties Reference](#)

Configuration example for writing data to Ozone FS

A simple configuration example for the HDFS Sink connector that writes data to the Ozone FS.

The following is a simple configuration example for the HDFS Sink connector. In this example data is written to the Ozone FS. Short descriptions of the properties set in this example are also provided. For a full properties reference, see the *HDFS Sink properties reference*.

```
{
  "connector.class": "com.cloudera.dim.kafka.connect.hdfs.HdfsSinkConnector",
  "tasks.max": 1,
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "com.cloudera.dim.kafka.connect.converts.AvroConverter",
  "value.converter.passthrough.enabled": true,
  "value.converter.schema.registry.url": "http://localhost:9090/api/v1",
  "topics": "avro_topic",
  "hdfs.uri": "ofs://ozone1/volume1/bucket1/",
  "hdfs.output": "/topics_output/",
  "output.writer": "com.cloudera.dim.kafka.connect.hdfs.avro.AvroPartitionWriter",
  "output.avro.passthrough.enabled": true,
  "hdfs.kerberos.authentication": true,
  "hdfs.kerberos.user.principal": "user_account@MY-REALM.COM",
  "hdfs.kerberos.keytab.path": "/path/to/user_account.keytab",
  "hadoop.conf.path": "/etc/hadoop/"
}
```

connector.class

Class name of the HDFS Sink connector.

key.converter

The converter capable of understanding the data format of the key of each record on this topic.

value.converter

The converter capable of understanding the data format of the value of each record on this topic.



Note: When the AvroConverter is used, you can specify Schema Registry properties to be used by the AvroConverter's Schema Registry client. This is done by adding the required Schema Registry property as a suffix to the value.converter property. For example, value.converter.schema.registry.url. Properties defined this way are passed on to the Schema Registry client used by the AvroConverter.

value.converter.passthrough.enabled

This property controls whether or not data is converted into the Kafka Connect intermediate data format before writing into an output file. Because in this example the input and output format is the same, the property is set to true, that is, data is not converted.

value.converter.schema.registry.url

The URL to Schema Registry. This is a mandatory property if the topic has records encoded in Avro format.

topics

List of topics to consume data from.

hdfs.uri

The Ozone FS (ofs) URI.

hdfs.output

The destination folder on the HDFS cluster where output files will reside.

output.writer

Determines the output file format. Because in this example the output format is Avro, AvroPartitionWriter is used.

output.avro.passthrough.enabled

This property has to match the configuration of the `value.converter.passthrough.enabled` property because both the input and output formats are Avro.

hdfs.kerberos.authentication

Enables or disables kerberos authentication.

hdfs.kerberos.user.principal

The user principal that the Kafka Connect role will use.

hdfs.kerberos.keytab.path

The path to the kerberos keytab file.

hadoop.conf.path

The path to the hadoop configuration files. This is necessary when the HDFS cluster has data encryption turned on.

HDFS Stateless Sink connector

The HDFS Stateless Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect Framework. Learn about the connector, its properties, and configuration.

The HDFS Stateless Sink Connector fetches messages from Kafka and stores them in HDFS. The topic this connector receives messages from is determined by the value of the `topics` parameter in the configuration. The messages can contain unstructured data (character or binary) or they can be in Avro or JSON format.

If the input is unstructured data, record processing is disabled. Multiple messages can be concatenated into a single output file on HDFS using a demarcator (for example newline for text messages). Merging is optional, large binary data can be forwarded to HDFS at a 1:1 ratio, one Kafka message equals a single HDFS file. However, in case of large files, this usage is encouraged only to avoid small file problems in HDFS.

If the input is either Avro or JSON, record processing is enabled. In a case like this, the schema of the records can be a predefined schema retrieved from Schema Registry (for both Avro and JSON data), it can be embedded in the Avro data, or inferred from the JSON data. The output can be Avro, JSON, CSV, or Parquet. Multiple records are typically merged into a single output file to store in HDFS.

The connector supports connecting to Kerberos enabled HDFS and Schema Registry services.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *HDFS Stateless Sink properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.
- If HDFS is running on a different cluster than the connector, the HDFS configuration resource files must be copied to the Kafka Connect node that the connector runs on.

Configuration example for unstructured data

In this example, the connector fetches unstructured messages from Kafka. The connector concatenates the messages using newline characters as the demarcator. The files that the connector creates and stores in HDFS are maximum 10 MB in size. The files are named according to the following pattern: `messages_[***UUID***]`.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "HDFS Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo/",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "",
  "topics": "[***TOPIC NAME***]",
  "parameter.HDFS Sink Parameters:Hadoop Configuration Resources": "/etc/hadoop/conf/core-site.xml,/etc/hadoop/conf/hdfs-site.xml",
  "parameter.HDFS Sink Parameters:Kafka Message Data Format": "Raw",
  "parameter.HDFS Sink Parameters:Output Directory Pattern": "[***OUTPUT DIRECTORY PATH***]",
  "parameter.HDFS Sink Parameters:Output Filename Pattern": "messages_{filename.uuid}"
}
```

```
"parameter.HDFS Sink Parameters:Maximum File Size": "10 MB",
"parameter.HDFS Sink Parameters:Output File Demarcator": " \n"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

Hadoop Configuration Resources

A comma delimited list of Hadoop configuration files (core-site.xml, hdfs-site.xml, and so on). If HDFS is running on a different cluster than the connector, ensure that you copy the configuration files to the Kafka Connect node that the connector runs on.



Note: In the majority of cases, these files can be found in `in/etc/hadoop/conf`.

Output Directory Pattern

Specifies the path to the output directory in HDFS.

Maximum File Size

The maximum size of the output data file. In this example, the maximum size is 10 MB.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to Raw. This means that the connector expects unstructured text or binary data and that record processing is disabled.

Output File Demarcator

Specifies the character sequence for demarcating (delimiting) message boundaries when multiple Kafka messages are ingested into an output file as raw messages. In this example, the property is set to `\n` (newline). This means that the newline character is used to separate the single line character data of the Kafka messages.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to `messages_${filename.uuid}`. The word `messages` and the underscore (`_`) are string literals. `${filename.uuid}` is an expression, which inserts a Universal Unique Identifier (UUID). As a result, the files are named according to the following pattern: `messages_***UUID***`.

Configuration example for Avro input and JSON output

In this example, the connector fetches messages in Avro format from Kafka. The connector parses the Avro records and converts them to JSON format using the schema definition stored in Schema Registry. The Avro file does not contain any schema reference. This is the reason why the schema must be specified in the connector parameters. The connector merges multiple records (all messages that are available for a single execution of the connector) into a JSON file, generates a name for the file using the following pattern: `messages_***UUID***.json`, and stores the file in HDFS.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "HDFS Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
}
```

```

"tasks.max": "1",
"nexus.url": "https://repository.cloudera.com/artifactory/repo/",
"working.directory": "/tmp/nifi-stateless-working",
"failure.ports": "",
"topics": "[***TOPIC NAME***]",
"parameter.HDFS Sink Parameters:Hadoop Configuration Resources": "/etc/hadoop/conf/core-site.xml,/etc/hadoop/conf/hdfs-site.xml",
"parameter.HDFS Sink Parameters:Kafka Message Data Format": "Avro",
"parameter.HDFS Sink Parameters:Output Directory Pattern": "[***OUTPUT DIRECTORY NAME***]",
"parameter.HDFS Sink Parameters:Output File Data Format": "JSON",
"parameter.HDFS Sink Parameters:Output Filename Pattern": "messages_{filename.uid}.json",
"parameter.HDFS Sink Parameters:Schema Access Strategy": "Schema Registry",
"parameter.HDFS Sink Parameters:Schema Name": "[***SCHEMA NAME***]",
"parameter.HDFS Sink Parameters:Schema Registry URL": "[***SCHEMA REGISTRY SERVER URL***]"
}

```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

Hadoop Configuration Resources

A comma delimited list of Hadoop configuration files (core-site.xml, hdfs-site.xml, and so on). If HDFS is running on a different cluster than the connector, ensure that you copy the configuration files to the Kafka Connect node that the connector runs on.



Note: In the majority of cases, these files can be found in/etc/hadoop/conf.

Output Directory Pattern

Specifies the path to the output directory in HDFS.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to Avro. This means that the connector expects Avro data and that record processing is enabled.

Output File Data Format

Specifies the format of the records written to the output file. In this example, this property is set to JSON, meaning that the output file format is JSON.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to messages_{filename.uid}.json. The word messages, the underscore (_), and .json are string literals. \${filename.uid} is an expression, which inserts a Universal Unique Identifier (UUID). As a result, the files are named according to the following pattern: messages_{***UUID***}.json.

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, this property is set to Schema Registry, meaning that the schema definition used to parse input data is taken from the Schema Registry.

Schema Name

Specifies the schema name to look up in Schema Registry.

Schema Registry URL

The URL of the Schema Registry server.



Note: Although not shown in this example, if Schema Registry is deployed on a secure cluster, the following Kerberos and/or TLS/SSL properties must also be set:

- Truststore Filename for Schema Registry
- Truststore Password for Schema Registry
- Truststore Type for Schema Registry
- Kerberos Keytab for Schema Registry
- Kerberos Principal for Schema Registry

Configuration example for JSON input and Parquet output

In this example, the connector fetches messages in JSON format from Kafka. The connector parses the JSON records and converts them to Parquet format using the schema inferred from the JSON data. The connector merges multiple records (all messages that are available for a single execution of the connector) into a Parquet file, generates a name for the file using the following pattern: `messages_***TIMESTAMP***_***SEQUENCE***.parquet.gz`, and stores the files in HDFS. The output directory path contains a fixed name and a timestamp appended to it: `/outputDir/***TIMESTAMP***`. The connector authenticates to HDFS using Kerberos.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "HDFS Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo/",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "",
  "topics": "[***TOPIC NAME***]",
  "parameter.HDFS Sink Parameters:Hadoop Configuration Resources": "/etc/hadoop/conf/core-site.xml,/etc/hadoop/conf/hdfs-site.xml",
  "parameter.HDFS Sink Parameters:Kafka Message Data Format": "JSON",
  "parameter.HDFS Sink Parameters:Output Directory Pattern": "/outputDir/${directory.timestamp}",
  "parameter.HDFS Sink Parameters:Output Directory Timestamp Format": "yyMMdd",
  "parameter.HDFS Sink Parameters:Output File Data Format": "Parquet",
  "parameter.HDFS Sink Parameters:Output Filename Pattern": "messages_${filename.timestamp}_${filename.sequence}.parquet.gz",
  "parameter.HDFS Sink Parameters:Output Filename Timestamp Format": "HHmmss",
  "parameter.HDFS Sink Parameters:Schema Access Strategy": "Infer Schema",
  "parameter.HDFS Sink Parameters:Kerberos Keytab for HDFS": "[***KEYTAB LOCATION***]",
  "parameter.HDFS Sink Parameters:Kerberos Principal for HDFS": "[***PRINCIPAL NAME***]",
  "parameter.HDFS Sink Parameters:Compression Codec for Parquet": "GZIP"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

Hadoop Configuration Resources

A comma delimited list of Hadoop configuration files (`core-site.xml`, `hdfs-site.xml`, and so on). If HDFS is running on a different cluster than the connector, ensure that you copy the configuration files to the Kafka Connect node that the connector runs on.



Note: In the majority of cases, these files can be found in `in/etc/hadoop/conf`.

Output Directory Pattern

Specifies the path to the output directory in HDFS. In this example, the path is set to `/outputDir/${directory.timestamp}`. As a result, the directory path in this example resolves to `/outputDir/[***TIMESTAMP***]`. The timestamp format is configured in the Output Directory Timestamp Format property.

Output Directory Timestamp Format

Specifies the timestamp format used in the `${directory.timestamp}` expression, which can be used in the Output Directory Pattern property's value. Directory Time Based Partitioning can be implemented using this expression. For example, if Output Directory Pattern is set to `/outputDir/${directory.timestamp}` and the timestamp format specified is `yyMMdd`, the output files will be generated in the `/outputDir/220101/` directory if the date is January 1, 2022.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to JSON. This means that the connector expects JSON data and that record processing is enabled.

Output File Data Format

Specifies the format of the records written to the output file. In this example, this property is set to Parquet, meaning that the output file format is Parquet.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to `messages_${filename.timestamp}_${filename.sequence}.parquet.gz`. The word `messages`, the underscores (`_`), and `parquet.gz` are string literals. `${filename.timestamp}` and `${filename.sequence}` are expressions. `${filename.timestamp}` inserts the current timestamp in the filename. `${filename.sequence}` inserts an incrementing sequence value in the filename. As a result, the files are named according to the following pattern: `messages_[***TIMESTAMP***]_[***SEQUENCE***].parquet.gz`.

Output Filename Timestamp Format

Specifies the timestamp format used in the `${filename.timestamp}` expression. The expression is used when generating the output filename. In this example, the format is `HHmmss`.

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, this property is set to Infer Schema, meaning that the schema is determined (inferred) from the JSON data.

Kerberos Keytab for HDFS

Specifies the keytab used to authenticate to HDFS with Kerberos.

Kerberos Principal for HDFS

Specifies the principal used to authenticate to HDFS with Kerberos.

Compression Codec

Specifies the codec used to compress files and store them in compressed format in HDFS, in the example the GZIP codec is used.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[HDFS Stateless Sink properties reference](#)

HTTP Sink connector

The HTTP Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The HTTP Sink connector obtains messages from a Kafka topic and transfers their content in a HTTP POST requests to a specified endpoint. The topic the connector receives messages from is determined by the value of the `topics` parameter in the configuration. The connector can forward the data it reads from Kafka as is (raw data) or can be configured to execute record processing. When record processing is enabled, the connector expects the incoming data in Avro or JSON format. In case of Avro, the connector can either read the record schema from the Avro file it receives (provided that the schema is embedded) or it can fetch the schema from Schema Registry. In case of JSON, it can either infer the schema or fetch it from Schema Registry. What strategy is used to retrieve the schema is determined by the `Schema Access Strategy` property.

If Schema Registry is used, and it is on a Kerberized cluster, the `krb5.file` property must point to the `krb5.conf` file that provides access to the cluster on which Schema Registry is present. This means that the `krb5.conf` file must be on the same cluster node that the connector runs on. The connections to Schema Registry and the HTTP server can be secured by TLS. The keystore and truststore files necessary for securing these connections must also be on the same cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *HTTP Sink properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.

- The value of the Schema Access Strategy property is not independent of the value of the Kafka Message Data Format property. As a result, you must exercise caution when configuring Schema Access Strategy.
 - If the value of Kafka Message Data Format is AVRO, the possible values for Schema Access Strategy are Schema Registry or Embedded Schema.
 - If the value of Kafka Message Data Format is JSON, the possible values for Schema Access Strategy are Schema Registry or Infer Schema.

Configuration example

In this example, the connector receives data in any format and forwards the raw data as the content of an HTTP POST request.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "HTTP Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "input.port": "Input from Kafka",
  "failure.ports": "Failure",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.HTTP Sink Parameters:Forward Raw Data": "true",
  "parameter.HTTP Sink Parameters:Remote URL": "http://[***SERVER HOSTNAME***]:[***PORT***]/[***PATH***]"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics The name of the Kafka topic the connector fetches messages from.

Forward Raw Data

Specifies whether messages from Kafka should be forwarded as is or converted to JSON. In this example, the property is set to true, meaning that the connector does not process any records. It forwards incoming data as is.

Remote URL

Identifies the HTTP endpoint that receives the messages sent by this connector. In this example, SSL is not used. As a result, the URL starts with http not https. For example, http://my.http-server.com:22000/contentListener.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[HTTP Sink properties reference](#)

InfluxDB Sink connector

The InfluxDB Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect Framework. Learn about the connector, its properties, and configuration.



Important: The InfluxDB Sink connector only supports InfluxDB v1.

The InfluxDB Sink Connector fetches messages from Kafka and loads them to InfluxDB. The topic this connector receives messages from is determined by the value of the topics parameter in the configuration. The messages the connector receives from Kafka can be in either Avro or JSON format.

If the connector input is in Avro format, then it can either read the schema from the Avro file it receives (provided that the schema is embedded) or it can fetch the schema from Schema Registry. If the connector's input is in JSON format, then it can either infer the schema or fetch it from Schema Registry. The strategy that is used to retrieve the schema is determined by the Schema Access Strategy property.

The connector can authenticate to InfluxDB using username and password.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *InfluxDB Sink properties reference*.

Notes and limitations

- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.
- The value of the Schema Access Strategy property is not independent of the value of the Kafka Message Data Format property. As a result, you must exercise caution when configuring Schema Access Strategy.
 - If the value of Kafka Message Data Format is AVRO, the possible values for Schema Access Strategy are Schema Registry or Embedded Schema.
 - If the value of Kafka Message Data Format is JSON, the possible values for Schema Access Strategy are Schema Registry or Infer Schema.
- The Line Protocol Query parameter contains a query expression which creates a string from input record fields. The query depends on the record schema. The result of the query must comply with the [line protocol format](#).

Configuration example

In this example, the connector receives data in JSON format (which is the default setting), transforms the data to line protocol format, and inserts it to InfluxDB. Schema Registry is not used in the example. As a result, this example does not include Schema Registry key and truststore configurations.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "InfluxDB Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "extensions.directory": "/tmp/stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo/",
  "failure.ports": "Retry from PutInfluxDB",
  "topics": "[***TOPIC NAME***]",
  "parameter.InfluxDB Connection URL": "http://[**SERVER HOSTNAME**]:[**PORT**]",
  "parameter.InfluxDB Database Name": "[***DATABASE NAME***]",
  "parameter.InfluxDB User Name": "[***USER NAME***]",
  "parameter.InfluxDB User Password": "[***USER PASSWORD***]",
  "parameter.Line Protocol Query": "select 'temp,device_id=' || device_id ||
  ',device_state=' || device_state || ',part_id=' || part_id || ',code=' ||
  code || ',result=' || result_code || ' temperature=' || temperature || ' '
  || ts || '000000' as payload from FLOWFILE",
  "parameter.Schema Access Strategy": "Infer Schema"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

InfluxDB Connection URL

The URL used to connect to InfluxDB.

InfluxDB Database Name

The name of the InfluxDB database.

InfluxDB User Name

The username used to connect to InfluxDB.

InfluxDB User Password

The password used to connect to InfluxDB.

Line Protocol Query

Specifies a query based on the Avro schema of the records and the measurement used.

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, the property is set to Infer Schema. This means that the schema is determined (inferred) from the JSON data.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[InfluxDB Sink properties reference](#)

JDBC Sink connector

The JDBC Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The JDBC Sink connector fetches messages from Kafka and loads them into a database table. The topic this connector receives messages from is determined by the value of the topics property in the configuration. The messages the connector receives from Kafka can be in either Avro or JSON format and must contain records that can be inserted into the database table. The schema of the records can be a predefined schema retrieved from Schema Registry (for both Avro and JSON data), it can be embedded in the Avro data, or inferred from the JSON data. The strategy that is used to retrieve the schema is determined by the Schema Access Strategy property.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR_NAME***] Parameters:
```

For a comprehensive list of these properties, see the *JDBC Sink properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.

Configuration example

In this example, the connector fetches messages in JSON format from Kafka, translates it to a SQL INSERT statement in order to load the data into a PostgreSQL database. The schema of the records must match the schema of the database table. The record schema in this example is inferred from the input JSON.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "JDBC Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "Retry from PutDatabaseRecord",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.JDBC Sink Parameters:Database Connection URL": "[***JDBC URL***]",
  "parameter.JDBC Sink Parameters:Database Driver Location": "[***PATH TO JDBC DRIVER***]",
  "parameter.JDBC Sink Parameters:Database Driver Class Name": "org.postgresql.Driver",
  "parameter.JDBC Sink Parameters:Database Type": "PostgreSQL",
  "parameter.JDBC Sink Parameters:Database User Name": "[***USERNAME***]",
  "parameter.JDBC Sink Parameters:Database User Password": "[***PASSWORD***]",
  "parameter.JDBC Sink Parameters:Database Table Name": "[***TABLE NAME***]",
  "parameter.JDBC Sink Parameters:Kafka Message Data Format": "JSON",
  "parameter.JDBC Sink Parameters:Schema Access Strategy": "Infer Schema"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

Database Connection URL

The JDBC URL of the PostgreSQL database. For example, jdbc:postgresql://myhost:5432/mydb.

Database Driver Location

The path to the PostgreSQL JDBC driver JAR file (or the directory containing the JAR).

Database Driver Class Name

The Java class name of the PostgreSQL Driver implementation.

Database Type

The type of the database. Because this is a PostgreSQL example, this property is set to PostgreSQL.

Database User Name

The username used for authenticating to the database.

Database User Password

The name of the database table to load data to.

Kafka Message Data Format

The format of the messages the connector receives from Kafka.

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, the property is set to Infer Schema. This means that the schema is determined (inferred) from the JSON data.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[JDBC Sink properties reference](#)

Kudu Sink connector

The Kudu Sink connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The Kudu Sink connector fetches messages from Kafka and loads them into a table in Kudu. The topic this connector receives messages from is determined by the value of the `topics` property in the configuration. The messages the connector receives from Kafka can be in either Avro or JSON format and must contain records that can be inserted into a Kudu table. If the connector input is in Avro format, then it can either read the schema from the Avro file it receives (provided that the schema is embedded) or it can fetch the schema from Schema Registry. If the connector's input is in JSON format, then it can either infer the schema, or fetch it from Schema Registry. The strategy that is used to retrieve the schema is determined by the `Schema Access Strategy` property.

Kudu is expected to be on a Kerberized cluster. Additionally, Schema Registry should be on the same cluster as Kudu. In the connector's configuration the `krb5.file` property must point to the `krb5.conf` file that provides access to the cluster that Kudu and Schema Registry are on. The keytab and truststore files that are necessary to access Kudu and Schema Registry must be present on the cluster node that the connector runs on.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *Kudu Sink properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.

- Schema Branch and Schema Version can not be specified at the same time.
- Schema Registry URL is mandatory even if no registry is used. In a case like this, use the default value, or remove the property completely from the configuration JSON. Alternatively, if your input data is JSON, you can use the Infer Schema value to infer the schema from the input, or if your input data is Avro, you can use the Embedded Schema value to read the schema embedded in the input.
- Truststore parameters are mandatory even if connecting to a non-secure Schema Registry or no registry is used. In that case the default (empty) truststore can be used.
- Kerberos parameters for Schema Registry are mandatory even if connecting to a non-secure Schema Registry or no registry is used. In that case the default (empty) keytab can be used.
- Kudu is expected to be on a Kerberized cluster. As a result, Kerberos parameters for Kudu are mandatory.

Configuration example

In this example, the connector receives data in JSON format from a Kafka topic and forwards it to a Kudu table. The record schema is inferred based on the data.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "Kudu Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "tasks.max": "1",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "input.port": "Input from Kafka",
  "failure.ports": "Failure",
  "krb5.file": "[***PATH TO KRB5 FILE***]",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.Kudu Sink Parameters:Kafka Message Data Format": "JSON",
  "parameter.Kudu Sink Parameters:Kerberos Keytab for Kudu": "[***PATH TO KUDU KEYTAB***]",
  "parameter.Kudu Sink Parameters:Kerberos Principal for Kudu": "[***KUDU PRINCIPAL***]",
  "parameter.Kudu Sink Parameters:Kudu Masters": "[***KUDU MASTER 1 WITH PORT***],[***KUDU MASTER 2 WITH PORT***],[***KUDU MASTER 3 WITH PORT***]",
  "parameter.Kudu Sink Parameters:Schema Access Strategy": "Infer Schema",
  "parameter.Kudu Sink Parameters:Table Name With Schema": "[***KUDU SCHEMA***].[***KUDU TABLE***]"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

krb5.file

The path to a krb5.conf file that contains the information needed to access the Kudu service on a Kerberized cluster

topics

The name of the Kafka topic that the connector fetches messages from.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to JSON. This means that the connector expects data in JSON format.

Kerberos Keytab for Kudu

The path to the keytab file that enables access to the Kudu service on a Kerberized cluster.

Kerberos Principal for Kudu

The principal used to access Kudu

Kudu Masters

A comma-separated list of server URLs that identify the Kudu service's master nodes. For example: localhost:7051,localhost:7151,localhost:7251

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, the property is set to Infer Schema. This means that the record schema is inferred based on the data the connector receives.

Table Name With Schema

The schema and name of the Kudu table in which the records get inserted.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[Kudu Sink properties reference](#)

S3 Sink connector

The S3 Sink Connector is a Stateless NiFi dataflow developed by Cloudera that is running in the Kafka Connect framework. Learn about the connector, its properties, and configuration.

The S3 Sink connector fetches messages from Kafka and uploads them to AWS S3. The topic this connector receives messages from is determined by the value of the topics property in the configuration. The messages can contain unstructured (character or binary) data or they can be in Avro or JSON format.

If the input is unstructured data, record processing is disabled. In a case like this, multiple messages can be concatenated into a single output file on S3 using a demarcator (for example, newline for text messages). Merging is optional, large binary data can be forwarded to S3 at a 1:1 ratio, one Kafka message equals a single S3 object.

If the input is either Avro or JSON, record processing is enabled. In a case like this, the schema of the records can be a predefined schema retrieved from Schema Registry (for both Avro and JSON data), it can be embedded in the Avro data, or inferred from the JSON data. The output can be Avro, JSON or CSV. Multiple records are typically merged into a single output file before uploading to S3.

The connector can authenticate to AWS using an Access Key.

Properties and configuration

Configuration is passed to the connector in a JSON file during creation. The properties of the connector can be categorized into three groups. These are as follows:

Common connector properties

These are the properties of the Kafka Connect framework that are accepted by all connectors. For a comprehensive list of these properties, see the *Apache Kafka documentation*.

Stateless NiFi Sink properties

These are the properties that are specific to the Stateless NiFi Sink connector. All Stateless NiFi Sink connectors share and accept these properties. For a comprehensive list of these properties, see the *Stateless NiFi Sink property reference*.

Connector/dataflow-specific properties

These properties are unique to this specific connector. Or to be more precise, unique to the dataflow running within the connector. These properties use the following prefix:

```
parameter.[***CONNECTOR NAME***] Parameters:
```

For a comprehensive list of these properties, see the *ADLS Sink properties reference*.



Tip: When creating a new connector using the SMM UI, all valid properties for this connector are presented in the default configuration template. You can use the template in the UI as a quick reference to view all accepted properties.

Notes and limitations

- Required properties must be assigned a valid value even if they are not used in the particular configuration. If a required property is not used, either leave its default value, or completely remove the property from the configuration JSON.
- If a property that has a default value is completely removed from the configuration JSON, the system uses the default value.
- Properties not marked as required must be completely removed from the configuration JSON if not set.
- Schema Branch and Schema Version can not be specified at the same time.

Configuration example for fetching unstructured data

In this example, the connector fetches unstructured messages containing single line character data from Kafka. The connector concatenates the messages using newline characters as the demarcator. The files that the connector creates and uploads to ADLS are maximum 10 MB in size. The files are named according to the following pattern: `messages_{lenid}.uuid`. The connector authenticates to AWS using an Access Key.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "S3 Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "PutS3Object Failure",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.S3 Sink Parameters:AWS Access Key ID": "[***ACCESS KEY ID***]",
  "parameter.S3 Sink Parameters:AWS Secret Access Key": "[***SECRET ACCESS KEY***]",
  "parameter.S3 Sink Parameters:S3 Region": "[***S3 REGION***]",
  "parameter.S3 Sink Parameters:S3 Bucket": "[***S3 Bucket***]",
  "parameter.S3 Sink Parameters:Maximum File Size": "10 MB",
  "parameter.S3 Sink Parameters:Kafka Message Data Format": "Raw",
  "parameter.S3 Sink Parameters:Output File Demarcator": "\n",
  "parameter.ADLS Sink Parameters:Output Filename Pattern": "messages_{filename.uuid}"
}
```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

AWS Access Key ID

The Access Key ID to use for authentication to AWS.

AWS Secret Access Key

The Secret Access Key to use for authentication to AWS.

S3 Region

The AWS Region of the S3 Bucket to upload data to.

S3 Bucket

The S3 Bucket to upload data to.

Maximum File Size

The maximum size of the output data file. In this example, the maximum size is 10 MB.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to Raw. This means that the connector expects unstructured text or binary data.

Output File Demarcator

Specifies the character sequence for demarcating (delimiting) message boundaries when multiple Kafka messages are ingested into an output file as raw messages. In this example, the property is set to `\n` (newline). This means that the newline character is used to separate the single line character data of the Kafka messages.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to `messages_${filename.uuid}`.

`messages_` is a string literal, `${filename.uuid}` is an expression that inserts a generated UUID in the filename. As a result, the files are named according to the following pattern: `messages_***UUID***`.

Configuration example for fetching JSON data

In this example, the connector fetches messages in JSON format from Kafka. The connector parses the JSON records and converts them to Avro format using the schema inferred from the JSON data. The schema is also embedded in the output file. The connector merges multiple records (all messages that are available for a single execution of the connector) into an Avro file, generates a name for the file using the following pattern `data_***TIMESTAMP***_***SEQUENCE***.avro`, and uploads the file to S3. The connector authenticates to AWS using an Access Key.



Note: All properties that are not present in this example have suitable default values or are unnecessary for this specific use case.

```
{
  "connector.class": "org.apache.nifi.kafka.connect.StatelessNiFiSinkConnector",
  "meta.smm.predefined.flow.name": "S3 Sink",
  "meta.smm.predefined.flow.version": "1.0.0",
  "key.converter": "org.apache.kafka.connect.storage.StringConverter",
  "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
  "nexus.url": "https://repository.cloudera.com/artifactory/repo",
  "extensions.directory": "/tmp/nifi-stateless-extensions",
  "working.directory": "/tmp/nifi-stateless-working",
  "failure.ports": "PutS3Object Failure",
  "topics": "[***KAFKA TOPIC NAME***]",
  "parameter.S3 Sink Parameters:AWS Access Key ID": "[***ACCESS KEY ID***]",
  "parameter.S3 Sink Parameters:AWS Secret Access Key": "[***SECRET ACCESS KEY***]"
}
```

```

"parameter.S3 Sink Parameters:S3 Region": "[***S3 REGION***]",
"parameter.S3 Sink Parameters:S3 Bucket": "[***S3 BUCKET***]"
"parameter.ADL Sink Parameters:Kafka Message Data Format": "JSON",
"parameter.ADL Sink Parameters:Schema Access Strategy": "Infer Schema",
"parameter.ADL Sink Parameters:Schema Write Strategy": "Embed Avro Schema",
,
"parameter.ADL Sink Parameters:Output File Data Format": "Avro",
"parameter.ADL Sink Parameters:Output Filename Pattern": "data_${filename.timestamp}_${filename.sequence}.avro",
"parameter.ADL Sink Parameters:Output Filename Timestamp Format": "yyyyMMdd_HH:mm:ss_SSS"
}

```

The following list collects the properties from the configuration example that must be customized for this use case.

topics

The name of the Kafka topic the connector fetches messages from.

AWS Access Key ID

The Access Key ID to use for authentication to AWS.

AWS Secret Access Key

The Secret Access Key to use for authentication to AWS.

S3 Region

The AWS Region of the S3 Bucket to upload data to.

S3 Bucket

The S3 Bucket to upload data to.

Kafka Message Data Format

The format of the messages the connector receives from Kafka. In this example, this property is set to JSON. This means that the connector expects JSON data.

Schema Access Strategy

Specifies the strategy used for determining the schema of the Kafka record. In this example, this property is set to Infer Schema, meaning that the schema is determined (inferred) from the JSON data.

Schema Write Strategy

Specifies whether the record schema is written to the output data file. In this example, this property is set to Embed Avro Schema, meaning that the schema is embedded in the output Avro file.

Output File Data Format

Specifies the format of the records written to the output file. In this example, this property is set to Avro, meaning that the output file format is Avro.

Output Filename Pattern

Specifies the structure of the name of the output file. This property accepts string literals (fixed text) as well as various expressions. In this example, the property is set to `data_${filename.timestamp}_${filename.sequence}.avro`.

`data_`, the underscore (`_`), and `.avro` are string literals. `${filename.timestamp}` and `${filename.sequence}` are expressions. `${filename.timestamp}` inserts the current timestamp in the filename. `${filename.sequence}` inserts an incrementing sequence value in the filename. As a result, the files are named according to the following pattern: `data_[***TIMESTAMP***]_[***SEQUENCE***].avro`.

Output Filename Timestamp Format

Specifies the timestamp format used in the `${filename.timestamp}` expression. The expression is used when generating the output filename.

Related Information

[Apache Kafka documentation](#)

[Stateless NiFi Sink properties reference](#)

[S3 Sink properties reference](#)