# Kafka Connect API
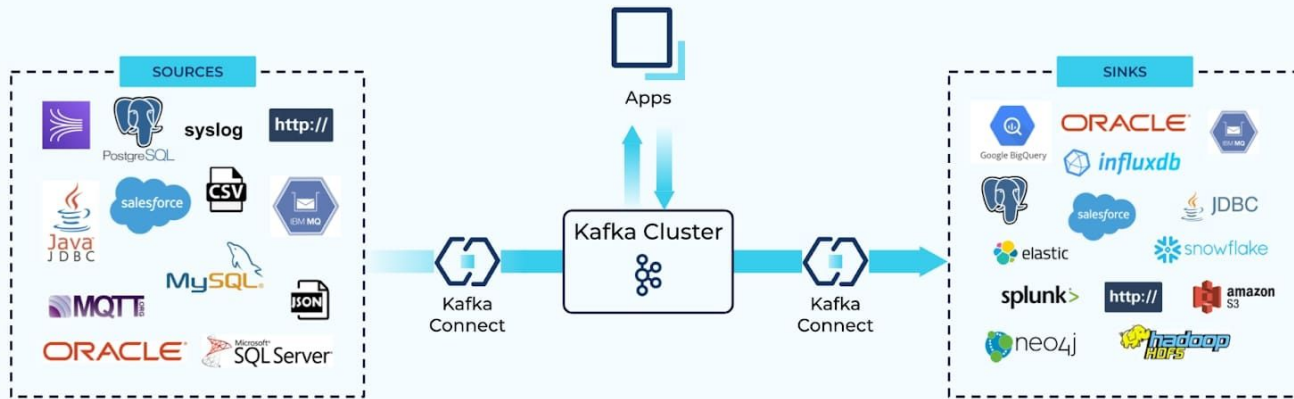
# Ingest Data from Upstream Systems

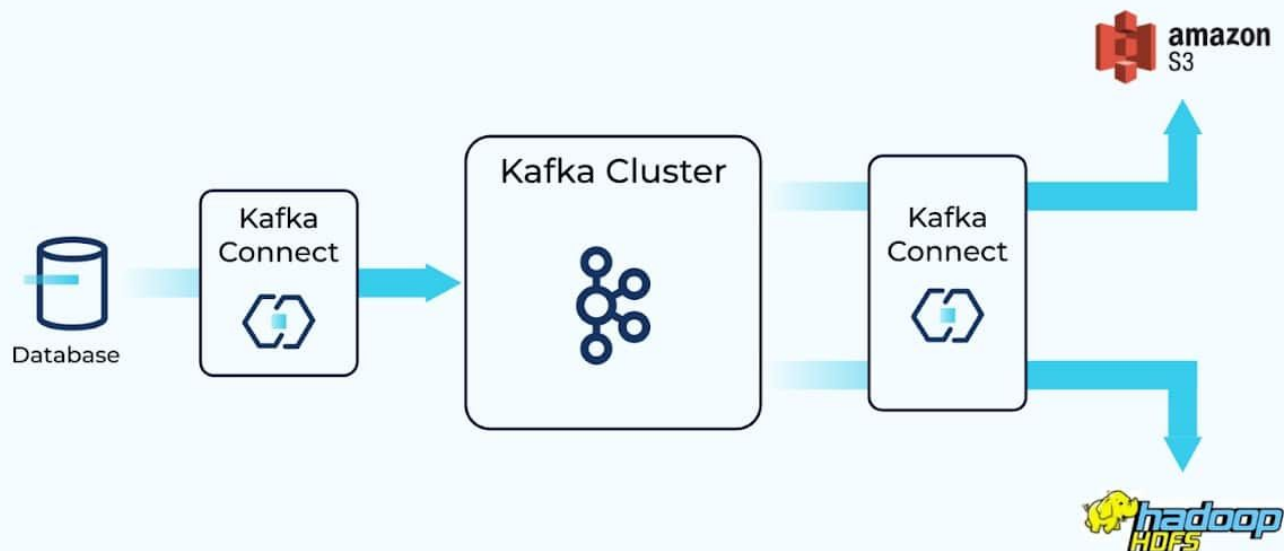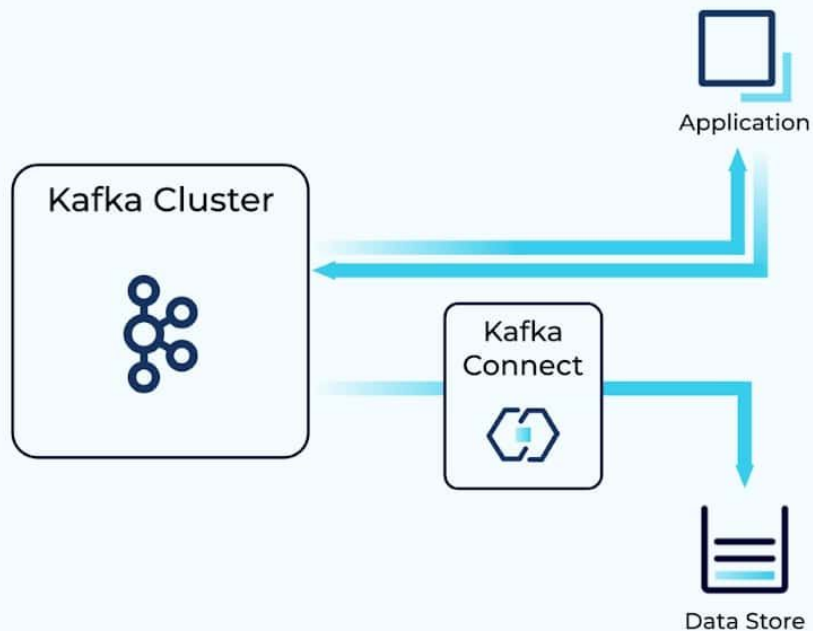# How Kafka Connect Works



```
{
    "connector.class":
        "io.confluent.connect.jdbc.JdbcSourceConnector",
    "connection.url":
        "jdbc:mysql://asgard:3306/demo",
    "table.whitelist":
        "sales,orders,customers"
}
```

@TheDanicaFine | developer.confluent.io

# Streaming Pipelines

# Writing to Datastores from Kafka



@TheDanicaFine | developer.confluent.io

# Evolve Processing from Old Systems to New

CONFLUENT
Developer

Existing
Application

Database

Kafka
Connect

Kafka Cluster

New
Application

@TheDanicaFine | developer.confluent.io

# Make Systems Real Time

# Why Not Write Your Own Integrations?

# Connectors

# Kafka Connect Workers

Ultimately, Kafka Connect workers are just JVM processes that you can deploy on bare metal or containers.

A few options present themselves:

- You're free to run a bare-metal, on-premises install of Confluent Platform
- For those leveraging infrastructure as a service, you may install Confluent Platform on those resources
- Terraform is an option on a couple cloud providers
- And of course, there's Docker which you can use for both on-prem and cloud-based installations

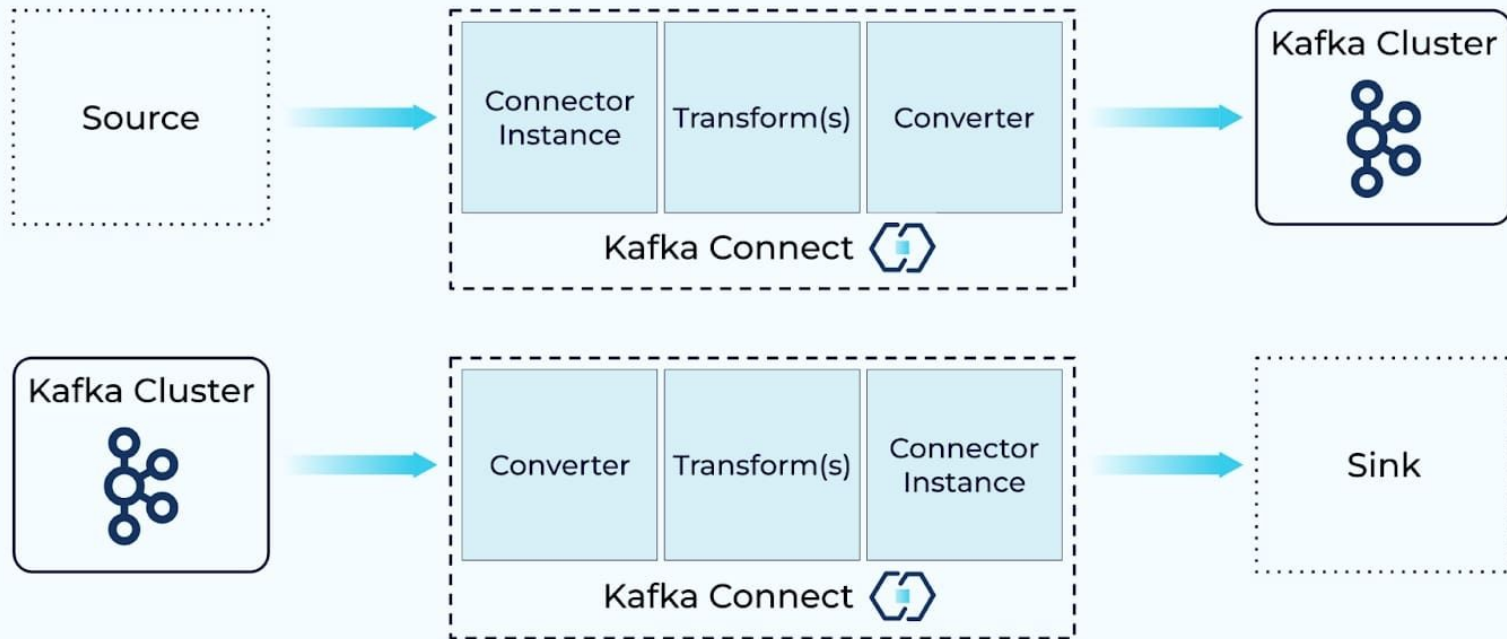# Managing a Kafka Connect Cluster

Once your Kafka Connect cluster is up and running, there's a bit of management that needs to be done:

- Connect workers have a number of default configuration settings that you may need to alter
- Depending on the needs of your systems, you might need to scale the Connect cluster up or down to suit demand changes
- And of course, you'll be monitoring for problems and fixing those that occur

# Inside Kafka Connect

# *Connectors*

CONFLUENT
**Developer**

| Source | → | Kafka Connect | → | Kafka Cluster |

**Kafka Connect** (top flow):
Source → [ **Connector Instance** | Transform(s) | Converter ] → Kafka Cluster

**Kafka Connect** (bottom flow):
Kafka Cluster → [ Converter | Transform(s) | **Connector Instance** ] → Sink

# For Example

- The **Debezium MySQL source connector** uses the MySQL bin log to read events from the database and stream these to Kafka Connect
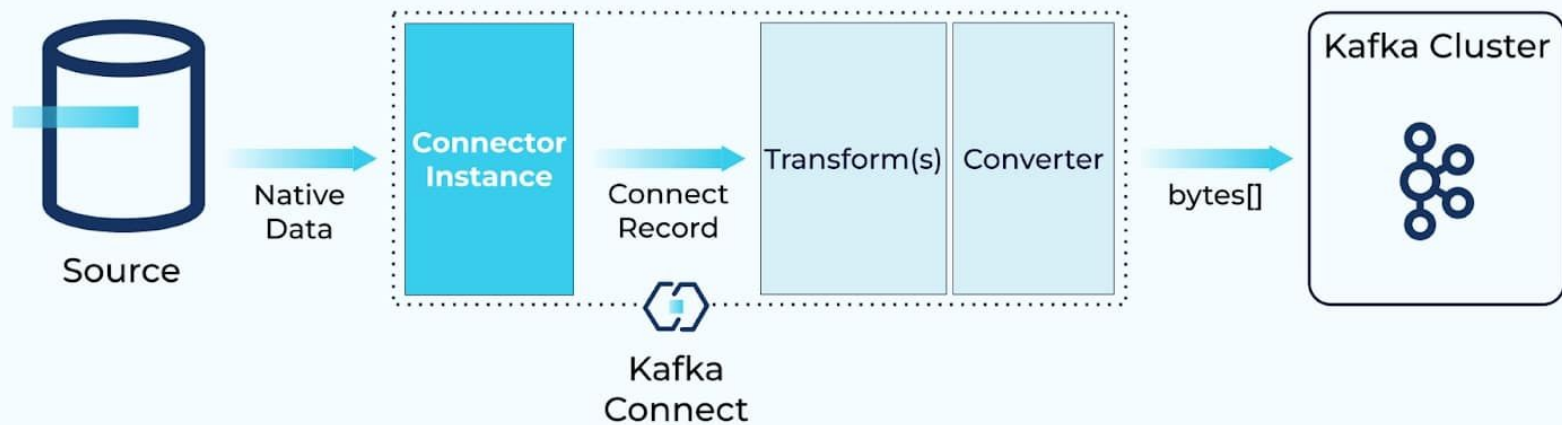- The **Elasticsearch sink connector** takes data from Kafka Connect, and using the Elasticsearch APIs, writes the data to Elasticsearch
- The **S3 connector** from Confluent can act as both a source and sink connector, writing data to S3 or reading it back in

# Add a Connector Instance with the REST API

Copy

```
curl -X PUT -H  "Content-Type:application/json" http://localhost:8083/connectors/sink-
elastic-01/config \
    -d '{
    "connector.class": "io.confluent.connect.elasticsearch.ElasticsearchSinkConnector",
    "topics"         : "orders",
    "connection.url" : "http://elasticsearch:9200",
    "type.name"      : "_doc",
    "key.ignore"     : "false",
    "schema.ignore"  : "true"
}'
```

# What is the Role of the Connector?

CONFLUENT
Developer



Source
Native Data
Connector Instance
Connect Record
Kafka Connect
Transform(s)
Converter
bytes[]
Kafka Cluster

# Converters Serialize/Deserialize the Data

CONFLUENT
**Developer**

Source → **Native Data** → Connector Instance → **Connect Record** → Transform(s) | **Converter** → **bytes[]** → Kafka Cluster

Kafka Connect

There are a ton of different converters available, but some common ones include:

- Avro – io.confluent.connect.avro.AvroConverter
- Protobuf – io.confluent.connect.protobuf.ProtobufConverter
- String – org.apache.kafka.connect.storage.StringConverter
- JSON – org.apache.kafka.connect.json.JsonConverter
- JSON Schema – io.confluent.connect.json.JsonSchemaConverter
- ByteArray – org.apache.kafka.connect.converters.ByteArrayConverter

# Serialization and Schemas

# Single Message Transforms



@TheDanicaFine | developer.confluent.io

# Common uses for SMTs include:

- Dropping fields from data at ingest, such as personally identifiable information (PII) if specified by the system requirements
- Adding metadata information such as lineage to data ingested through Kafka Connect
- Changing field data types
- Modifying the topic name to include a timestamp
- Renaming fields

# Deploying Kafka Connect

# Deploying Kafka Connect

CONFLUENT
Developer

JDBC Source

S3 Sink

S3 Task #1

JDBC Task #1

@TheDanicaFine | developer.confluent.io

# Tasks Are the Unit of Parallelism and Scale

CONFLUENT
**Developer**

| JDBC Source | | S3 Sink |

| S3 Task #1 | JDBC Task #1 | JDBC Task #2 |

@TheDanicaFine | developer.confluent.io

# Connect Worker

CONFLUENT
**Developer**

**JDBC Source**

**S3 Sink**

Worker

**S3 Task #1**

**JDBC Task #1**

**JDBC Task #2**

@TheDanicaFine | developer.confluent.io

# Kafka Connect Distributed Mode



@TheDanicaFine | developer.confluent.io

# *Multiple Workers vs Multiple Clusters*



CONFLUENT
Developer

Worker

**S3 Task #1**

Worker

Kafka Connect
Cluster #1

Worker

**JDBC Task #1**

Worker

**JDBC Task #2**

Kafka Connect
Cluster #2

Kafka Cluster

Offsets_1
Configs_1
Status_1

Offsets_2
Configs_2
Status_2

@TheDanicaFine | developer.confluent.io

# Kafka Connect Standalone Mode

CONFLUENT
**Developer**

**Worker**

| S3 Task #1 | JDBC Task #1 | JDBC Task #2 |

Offsets

# Error Handling in Kafka Connect

Kafka Cluster

Source topic messages

Kafka Connect

Sink data

@TheDanicaFine | developer.confluent.io

# Serialization Challenges - Wrong Converter



**Kafka Cluster**

Source topic messages

| JSON | JSON | JSON | JSON | JSON |

**Kafka Connect**

```
"Value.converter":
"AvroConverter"
```

```
org.apache.kafka.common.errors.SerializationException:
Unknown magic byte!
```

Sink data

# Serialization Challenges - Multiple Formats

CONFLUENT
Developer

**Kafka Cluster**

Source topic messages

| Avro | Avro | JSON | Avro | JSON |

**Kafka Connect**

`"Value.converter":`
`"AvroConverter"`

`org.apache.kafka.common.errors.SerializationException`

Sink data

# Error Tolerances - Fail Fast (Default)

# Error Tolerances - Dead Letter Queue

CONFLUENT
**Developer**

**Kafka Cluster**

Source topic messages

| Avro | Avro | Avro | JSON | Avro |

Dead letter queue

| JSON |

**Kafka Connect**

```
"Value.converter":
"AvroConverter"
```

errors.tolerance=all
errors.deadletterqueue.topic.name=my_dlq

**Failed messages written to DLQ**

Sink data

@TheDanicaFine | developer.confluent.io

# Reprocessing the Dead Letter Queue

**CONFLUENT Developer**

**Kafka Cluster**

Source topic messages

| Avro | Avro | Avro | JSON | Avro |

Dead letter queue

| JSON | |

**Kafka Connect**

`"Value.converter":`
`"AvroConverter"`

`"Value.converter":`
`"JsonConverter"`

```
errors.tolerance=all
errors.deadletterqueue.topic.name=my_dlq
```

Sink data

@TheDanicaFine | developer.confluent.io

# Troubleshoot the Dead Letter Queue

## MySqlSinkConnector_0

⇄ See in Stream lineage

**Overview**    Settings    Events

⏸ Pause

**Running**
This connector is running.

**Messages processed**

241

Total in last 7 days

**Messages behind**

0

Max lag in the last minute

**Messages in DLQ**

241

Total in last 7 days

# Dead Letter Queue Message Header

CONFLUENT
Developer

## dlq-lcc-q2nkw2

⇄ See in Stream lineage

Overview | Messages | Schema | Configuration

**Producers**
Bytes in/sec    0

**Consumers**
Bytes out/sec    8.68K

**Message fields**

- topic
- partition
- offset
- timestamp
- timestampType
- > headers
- key
- value

▶ ⏸   🔍 Filter by keyword    Jump to offset ▾    🔍 0 / Partition: 0    ▤ ▥

＋ **Produce a new message to this topic**

⌄   ▒▒▒▒B▒W�▒▒Item_9
       Partition: 0        Off

⌄   ▒▒▒▒ÌX▒▒▒Item_4
       Partition: 0        Off

⌄   ▒▒▒É▒▒▒W▒▒▒Item_9
       Partition: 0        Of

⌃ **Value**   Header   Key

▒▒▒▒B▒W�▒▒Item_948M�▒<E�#@▒City

⌄   ▒▒▒▒ÌX▒▒▒Item_413 ▒[l▒4▒@▒City_3
       Partition: 0    Offset: 168    Times

⌄   ▒▒▒É▒▒▒W▒▒▒Item_983▒▒Cl[▒▒?▒City_
       Partition: 0    Of

⌃ **Value**   Header   Key

```
 1 ▾    [
 2 ▾      {
 3          "key": "task.generation",
 4          "stringValue": "0"
 5        },
42 ▾      {
43          "key": "__connect.errors.exception.class.name",
44          "stringValue": "io.confluent.connect.jdbc.sink.TableAlterOrCreateException"
45        },
46 ▾      {
47          "key": "__connect.errors.exception.message",
48          "stringValue": "Table \"orders\" is missing and auto-creation is disabled"
49        },
```

@TheDanicaFine | developer.confluent.io

# *Troubleshoot a Failed Connector*

## MySqlSinkConnector_0

⇄ See in Stream lineage

**Overview**   Settings   Events

‖ Pause

---

### Failed

Found nested structure in input data. Ensure that your input events are a flat struct of primitive fields.

### Messages processed

254

Total in last 7 days

---

### Messages behind

32

Max lag in the last minute

### Messages in DLQ

254

Total in last 7 days

@TheDanicaFine | developer.confluent.io

CONFLUENT
**Developer**
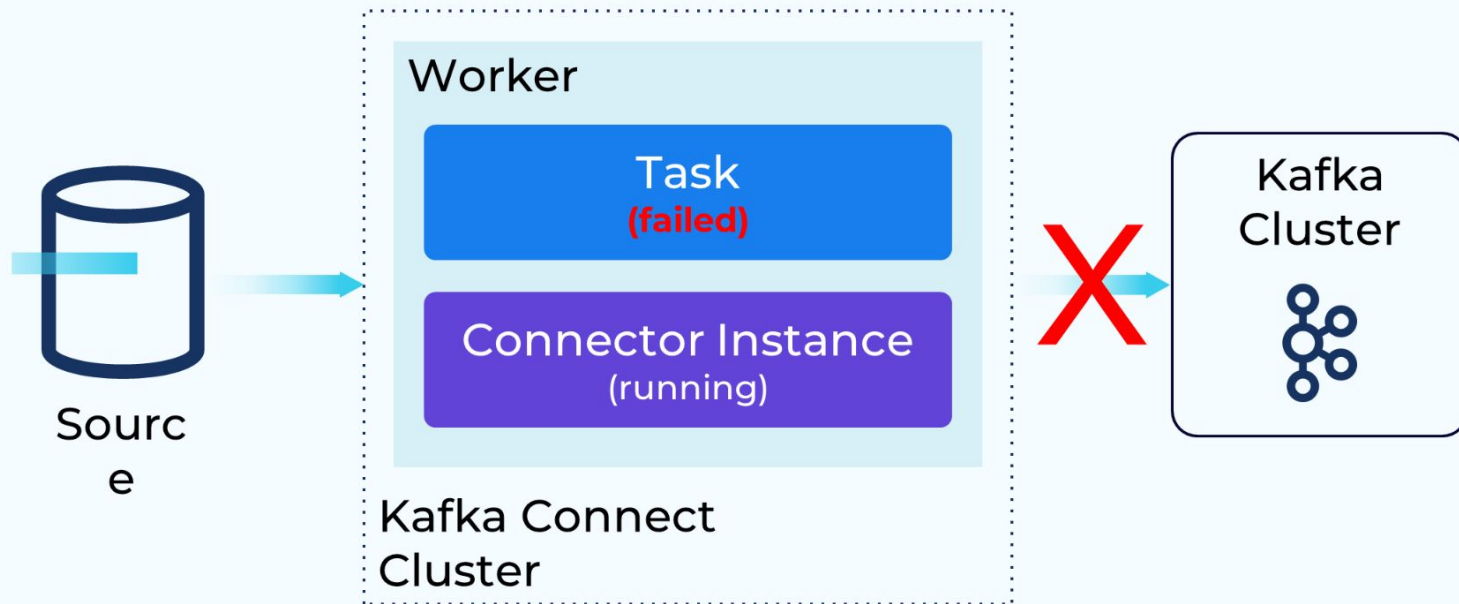
# Confluent Connect API

**CONFLUENT Developer**

```
GET /connect/v1/environments/{environment_id}/clusters/{kafka_cluster_id}/connectors/{connector_name}/status

{
  "name": "MySqlSinkConnector_0",
  "connector": {
    "state": "FAILED",
    "worker_id": "MySqlSinkConnector_0",
    "trace": "Found nested structure in input data. Ensure that your input events are a flat struct of
primitive fields.\n"
  },
  "tasks": [
    {
      "id": 0,
      "state": "USER_ACTIONABLE_ERROR",
      "worker_id": "MySqlSinkConnector_0",
      "msg": ""
    }
  ],
  "type": "sink"
}
```

@TheDanicaFine | developer.confluent.io

# Troubleshooting Scenario

# Getting Connector and Task Status

```
$ curl -s "http://localhost:8083/connectors/jdbc-sink/status" | \
jq '.connector.state'
"RUNNING"
```

```
$ curl -s "http://localhost:8083/connectors/jdbc-sink/status" | \
jq '.tasks[0].state'
"FAILED"
```

# *Getting Task Status*

```
$curl -s "http://localhost:8083/connectors/jdbc-sink/status" | jq '.tasks[0].trace' | sed 's/\\n/\n/g; s/\\t/\t/g'

"org.apache.kafka.connect.errors.ConnectException: Exiting WorkerSinkTask due to unrecoverable exception.
    at org.apache.kafka.connect.runtime.WorkerSinkTask.deliverMessages(WorkerSinkTask.java:618)
    at org.apache.kafka.connect.runtime.WorkerSinkTask.poll(WorkerSinkTask.java:334)
    at org.apache.kafka.connect.runtime.WorkerSinkTask.iteration(WorkerSinkTask.java:235)
    at org.apache.kafka.connect.runtime.WorkerSinkTask.execute(WorkerSinkTask.java:204)
    at org.apache.kafka.connect.runtime.WorkerTask.doRun(WorkerTask.java:200)
    at org.apache.kafka.connect.runtime.WorkerTask.run(WorkerTask.java:255)
    at java.base/java.util.concurrent.Executors$RunnableAdapter.call(Executors.java:515)
    at java.base/java.util.concurrent.FutureTask.run(FutureTask.java:264)
    at java.base/java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1128)
    at java.base/java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:628)
    at java.base/java.lang.Thread.run(Thread.java:829)
Caused by: org.apache.kafka.connect.errors.ConnectException: java.sql.SQLException: No suitable driver found for
jdbc:mysql://localhost/demo
    at io.confluent.connect.jdbc.util.CachedConnectionProvider.getConnection(CachedConnectionProvider.java:59)
    at io.confluent.connect.jdbc.sink.JdbcDbWriter.write(JdbcDbWriter.java:64)
    at io.confluent.connect.jdbc.sink.JdbcSinkTask.put(JdbcSinkTask.java:84)
    at org.apache.kafka.connect.runtime.WorkerSinkTask.deliverMessages(WorkerSinkTask.java:584)
    ... 10 more
Caused by: java.sql.SQLException: No suitable driver found for jdbc:mysql://localhost/demo
    at java.sql/java.sql.DriverManager.getConnection(DriverManager.java:702)
    at java.sql/java.sql.DriverManager.getConnection(DriverManager.java:189)
    at io.confluent.connect.jdbc.dialect.GenericDatabaseDialect.getConnection(GenericDatabaseDialect.java:247)
    at io.confluent.connect.jdbc.util.CachedConnectionProvider.newConnection(CachedConnectionProvider.java:80)
    at io.confluent.connect.jdbc.util.CachedConnectionProvider.getConnection(CachedConnectionProvider.java:52)
    ... 13 more
"
```

# Kafka Connect Log4j Logging

- **The log is the source of truth**

    ```
    $ confluent local services connect log
    ```

    ```
    $ docker-compose logs kafka-connect
    ```

    ```
    $ cat /var/log/kafka/connect.log
    ```

- **The Log4j properties files controls what is logged, the log message layout, and where log files are stored**

    ```
    /etc/kafka/connect-log4j.properties (default location)
    ```

# *Identify the Problem Cause*

CONFLUENT
Developer

```
[2022-07-19 23:57:28,600] ERROR [jdbc-sink|task-0] WorkerSinkTask{id=jdbc-sink-0} Task threw an uncaught and
unrecoverable exception. Task is being killed and will not recover until manually restarted
(org.apache.kafka.connect.runtime.WorkerTask:207)
org.apache.kafka.connect.errors.ConnectException: Exiting WorkerSinkTask due to unrecoverable exception.
     at org.apache.kafka.connect.runtime.WorkerSinkTask.deliverMessages(WorkerSinkTask...)
     at org.apache.kafka.connect.runtime.WorkerSinkTask.poll(WorkerSinkTask.java:334)
     at org.apache.kafka.connect.runtime.WorkerSinkTask.iteration(WorkerSinkTask.java:...)
     at org.apache.kafka.connect.runtime.WorkerSinkTask.execute(WorkerSinkTask.java:204)
     at org.apache.kafka.connect.runtime.WorkerTask.doRun(WorkerTask.java:200)
     at org.apache.kafka.connect.runtime.WorkerTask.run(WorkerTask.java:255)
     at java.base/java.util.concurrent.Executors$RunnableAdapter.call(Executors.java:515)
     at java.base/java.util.concurrent.FutureTask.run(FutureTask.java:264)
     at java.base/java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1128)
     at java.base/java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:628)
     at java.base/java.lang.Thread.run(Thread.java:829)
Caused by: org.apache.kafka.connect.errors.ConnectException: java.sql.SQLException: No suitable driver found for
jdbc:mysql://localhost/demo
     at io.confluent.connect.jdbc.util.CachedConnectionProvider.getConnection(CachedConnectionProvider.java:59)
     at io.confluent.connect.jdbc.sink.JdbcDbWriter.write(JdbcDbWriter.java:64)
     at io.confluent.connect.jdbc.sink.JdbcSinkTask.put(JdbcSinkTask.java:84)
     at org.apache.kafka.connect.runtime.WorkerSinkTask.deliverMessages(WorkerSinkTask.java:...)
     ... 10 more
Caused by: java.sql.SQLException: No suitable driver found for jdbc:mysql://localhost/demo
     at java.sql/java.sql.DriverManager.getConnection(DriverManager.java:702)
     at java.sql/java.sql.DriverManager.getConnection(DriverManager.java:189)
     at io.confluent.connect.jdbc.dialect.GenericDatabaseDialect.getConnection(GenericDatabaseDialect.java:247)
     at io.confluent.connect.jdbc.util.CachedConnectionProvider.newConnection(CachedConnectionProvider.java:80)
     at io.confluent.connect.jdbc.util.CachedConnectionProvider.getConnection(CachedConnectionProvider.java:52)
     ... 13 more
```

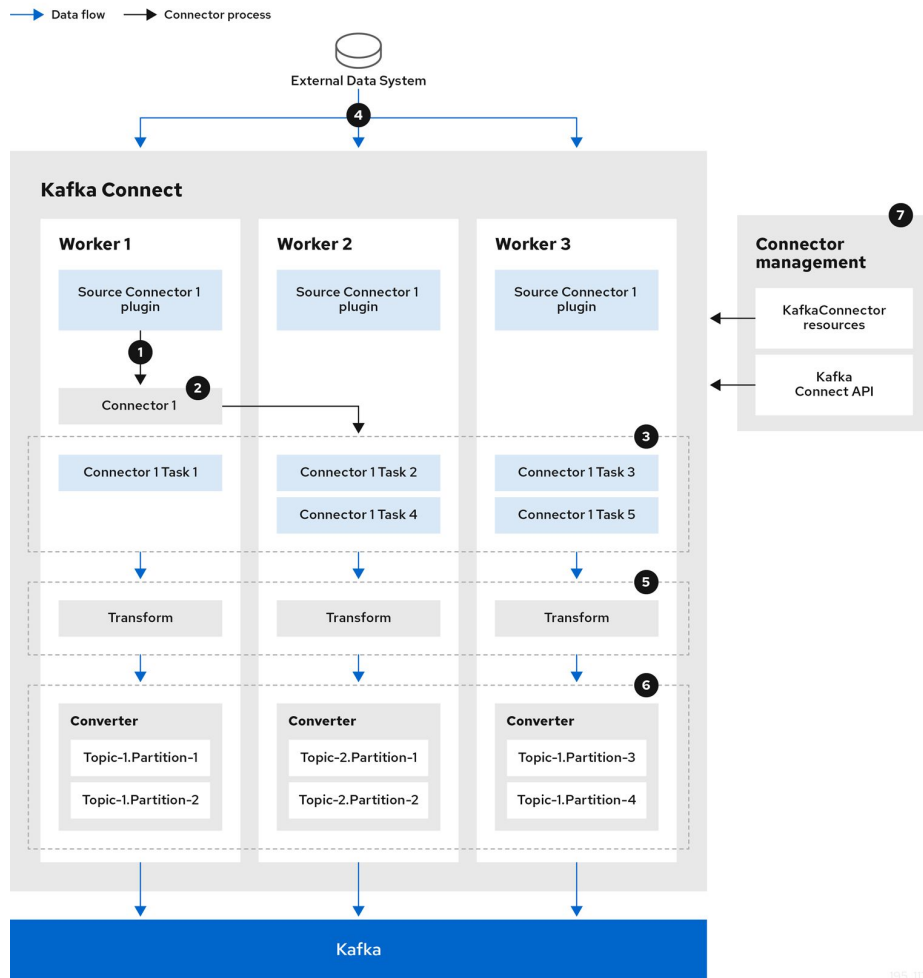**Symptom, not the cause**

**Possible causes**

# Summary

The main Kafka Connect components used in streaming data are as follows:

- Connectors to create tasks
- Tasks to move data
- Workers to run tasks
- Transforms to manipulate data
- Converters to convert data
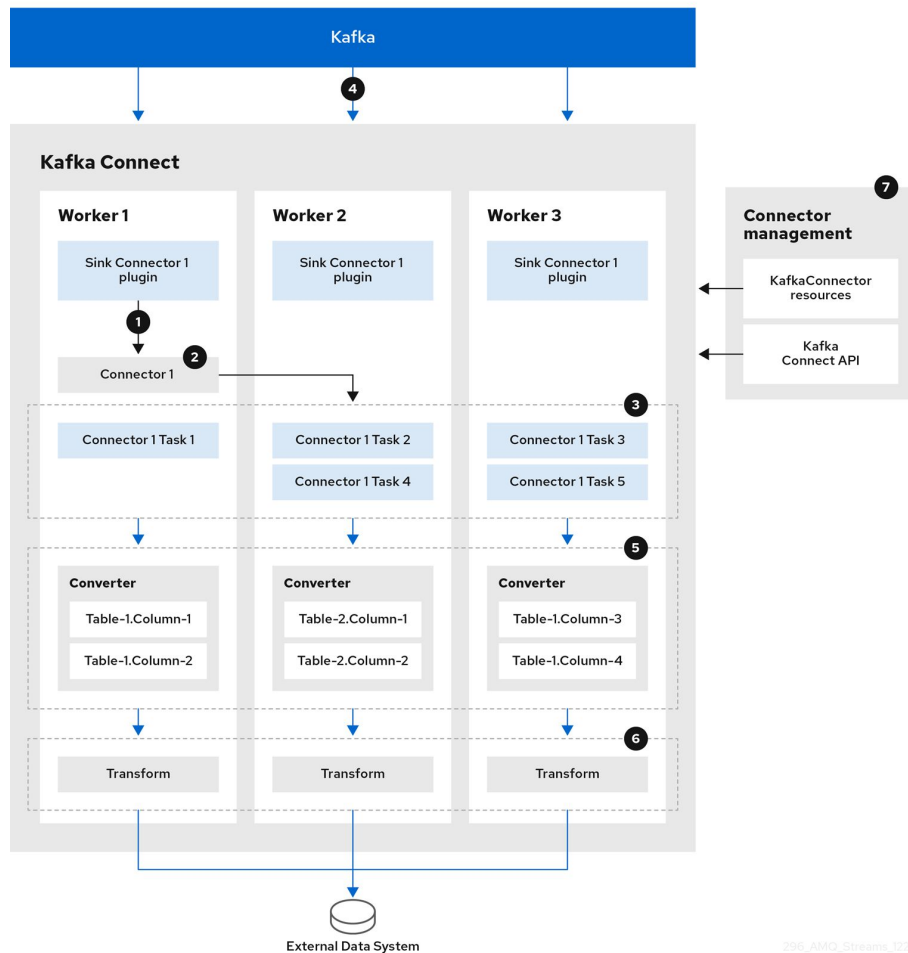
Connectors can be one of the following type:

- Source connectors that push data into Kafka
- Sink connectors that extract data out of Kafka

Source connector streaming data to Kafka

1. A plugin provides the implementation artifacts for the source connector
2. A single worker initiates the source connector instance
3. The source connector creates the tasks to stream data
4. Tasks run in parallel to poll the external data system and return records
5. Transforms adjust the records, such as filtering or relabelling them
6. Converters put the records into a format suitable for Kafka
7. The source connector is managed using KafkaConnectors or the Kafka Connect API

Sink connector
streaming data from
Kafka

1. A plugin provides the implementation artifacts for the sink connector
2. A single worker initiates the sink connector instance
3. The sink connector creates the tasks to stream data
4. Tasks run in parallel to poll Kafka and return records
5. Converters put the records into a format suitable for the external data system
6. Transforms adjust the records, such as filtering or relabelling them
7. The sink connector is managed using KafkaConnectors or the Kafka Connect API