

@rmoff

# From zero to Hero with Kafka Connect

a.k.a.

A practical guide to becoming l33t with Kafka Connect

# \$ whoami

- Robin Moffatt (@rmoff)
- Senior Developer Advocate at Confluent  
(Apache Kafka, not Wikis 😊)
- Working in data & analytics since 2001
- ♠️ Oracle ACE Director (Alumnus)



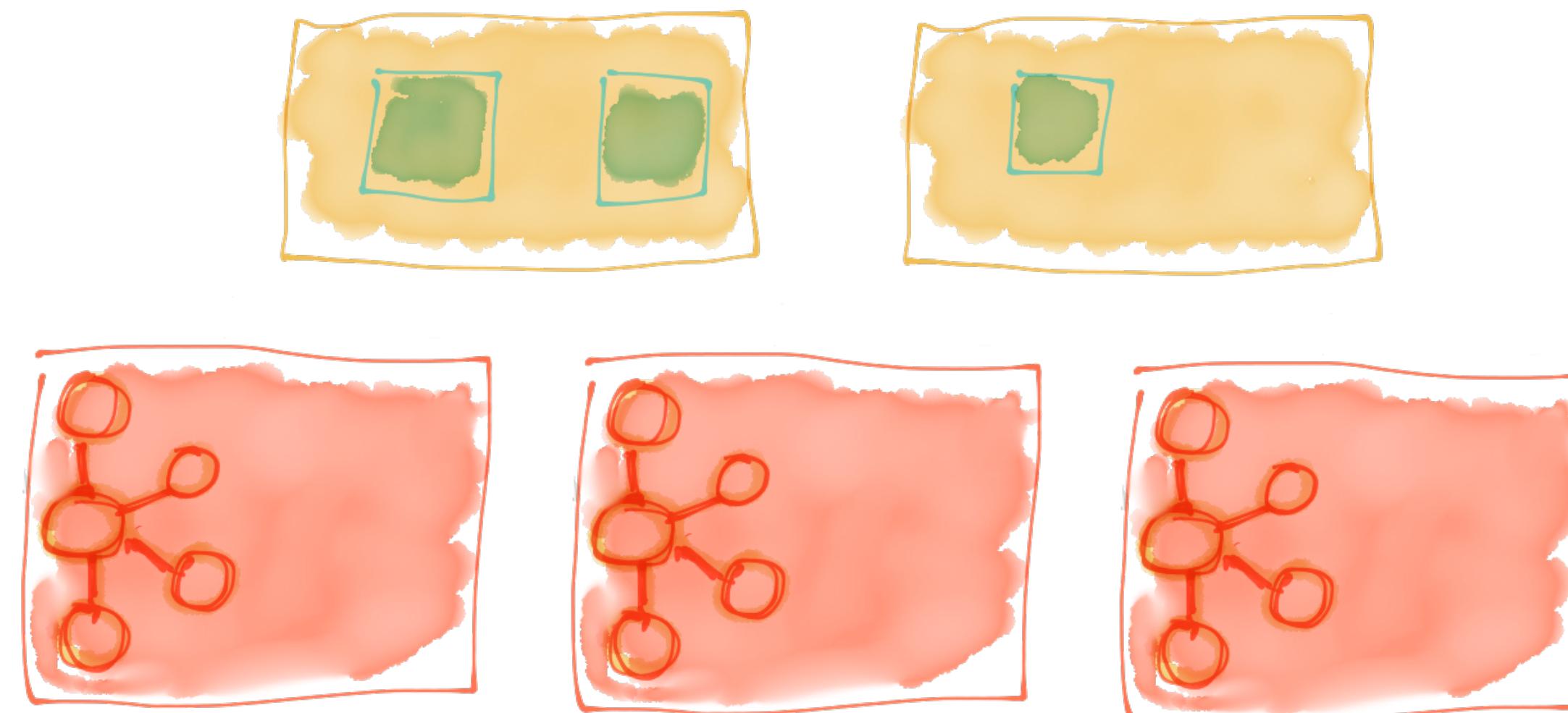
<http://rmoff.dev/talks> • <http://rmoff.dev/blog> • <http://rmoff.dev/youtube>

What is  
Kafka  
Connect?

# Streaming Integration with Kafka Connect



Sources



Kafka Connect

Kafka Brokers



CONFLUENT

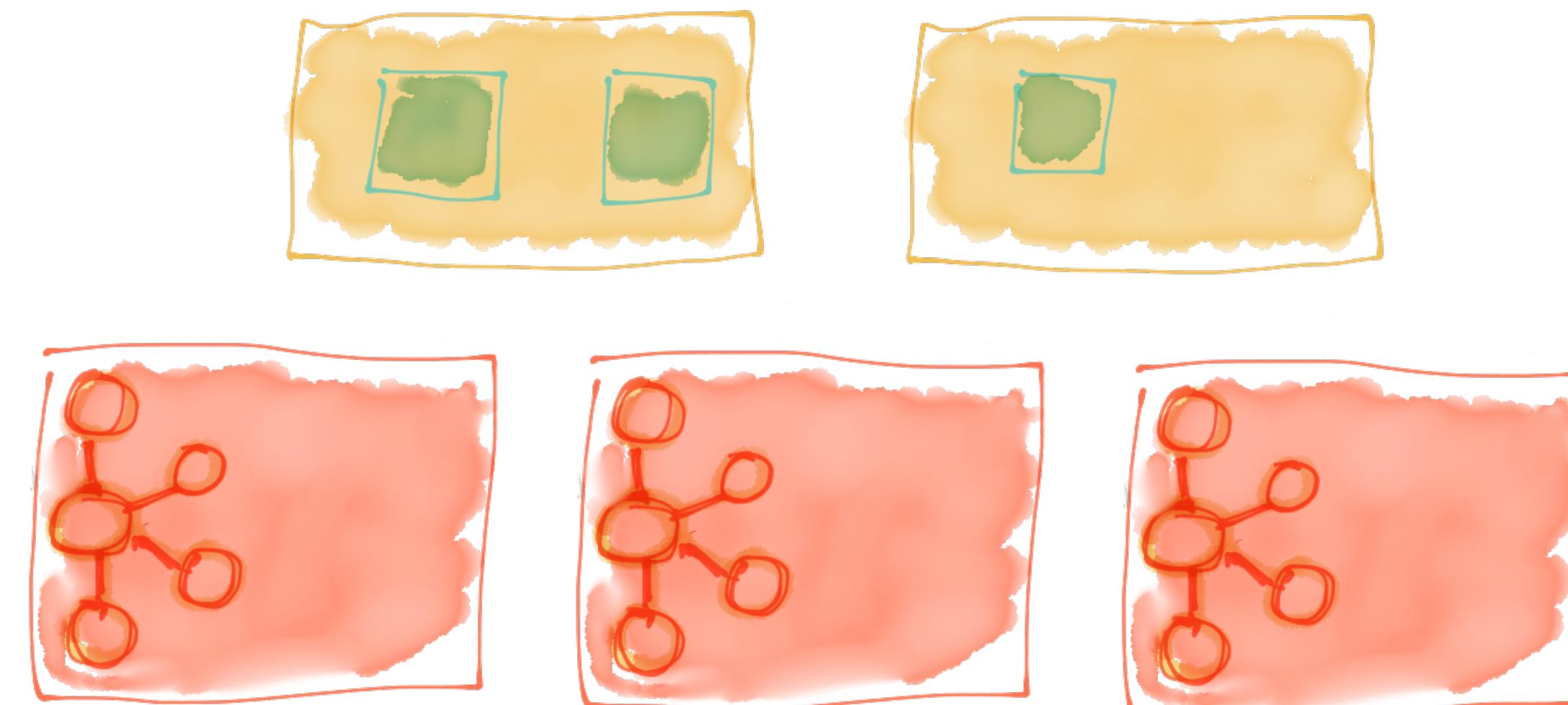
@rmoff

# Streaming Integration with Kafka Connect

Sinks

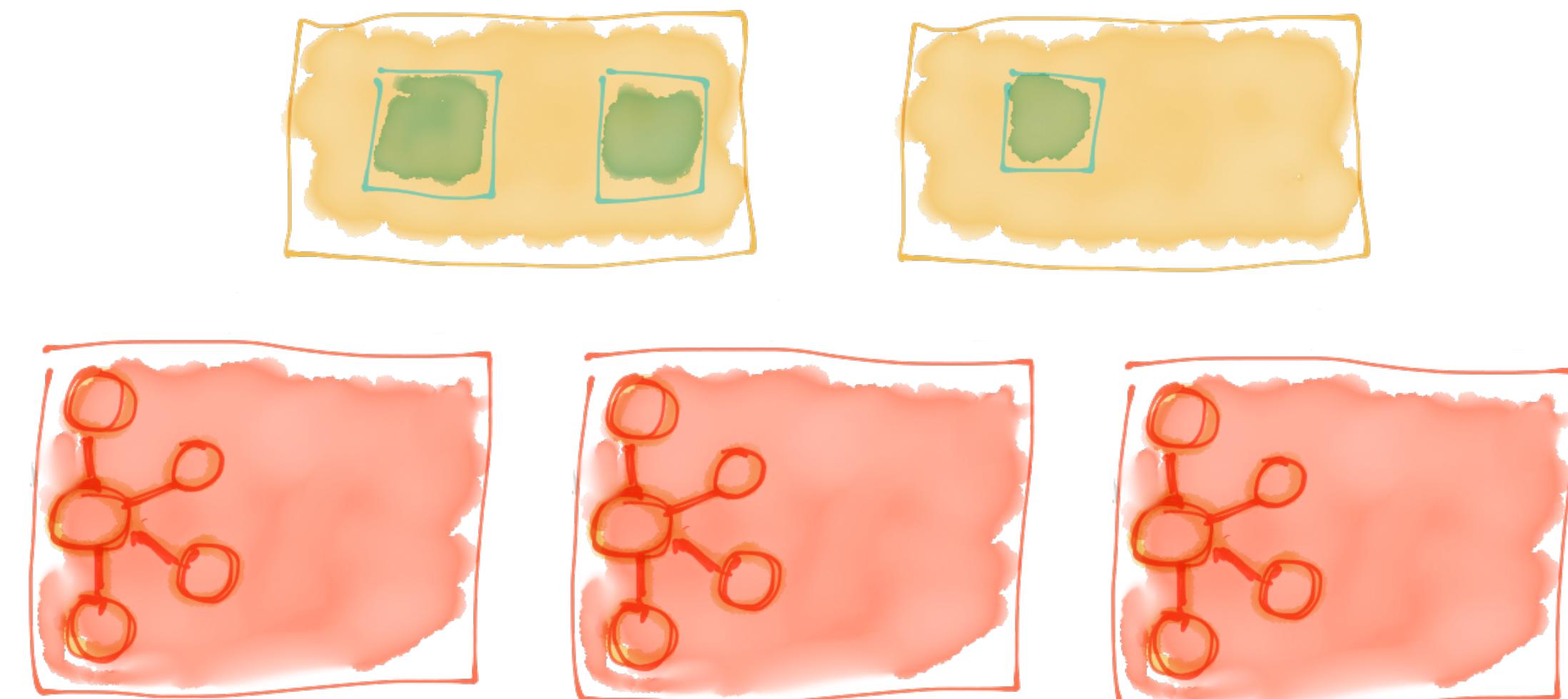
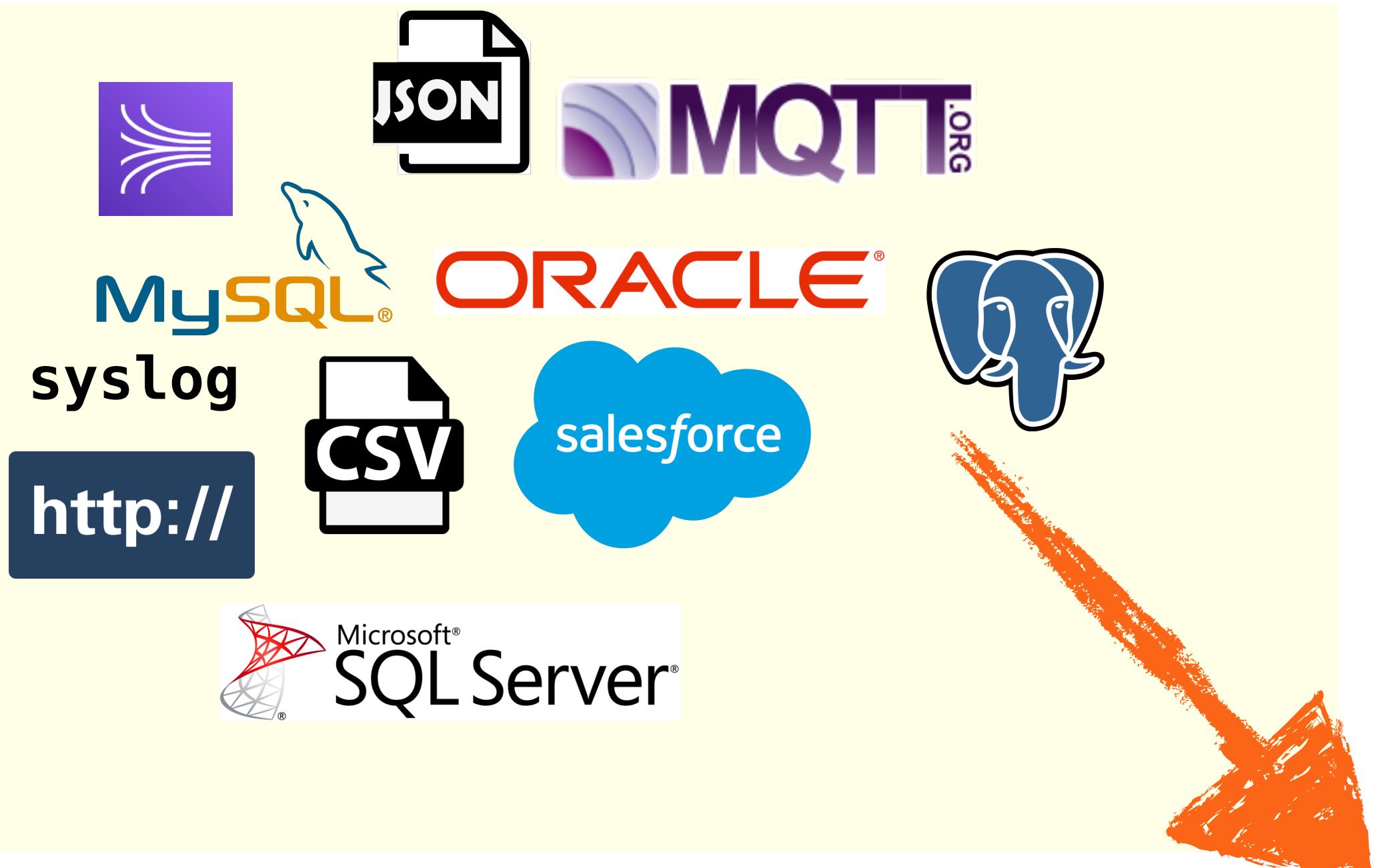


Kafka Connect



Kafka Brokers

# Streaming Integration with Kafka Connect



Kafka Connect

Kafka Brokers

# *Look Ma, No Code!*

```
{
```

```
"connector.class":
```

```
    "io.confluent.connect.jdbc.JdbcSourceConnector",
```

```
"connection.url":
```

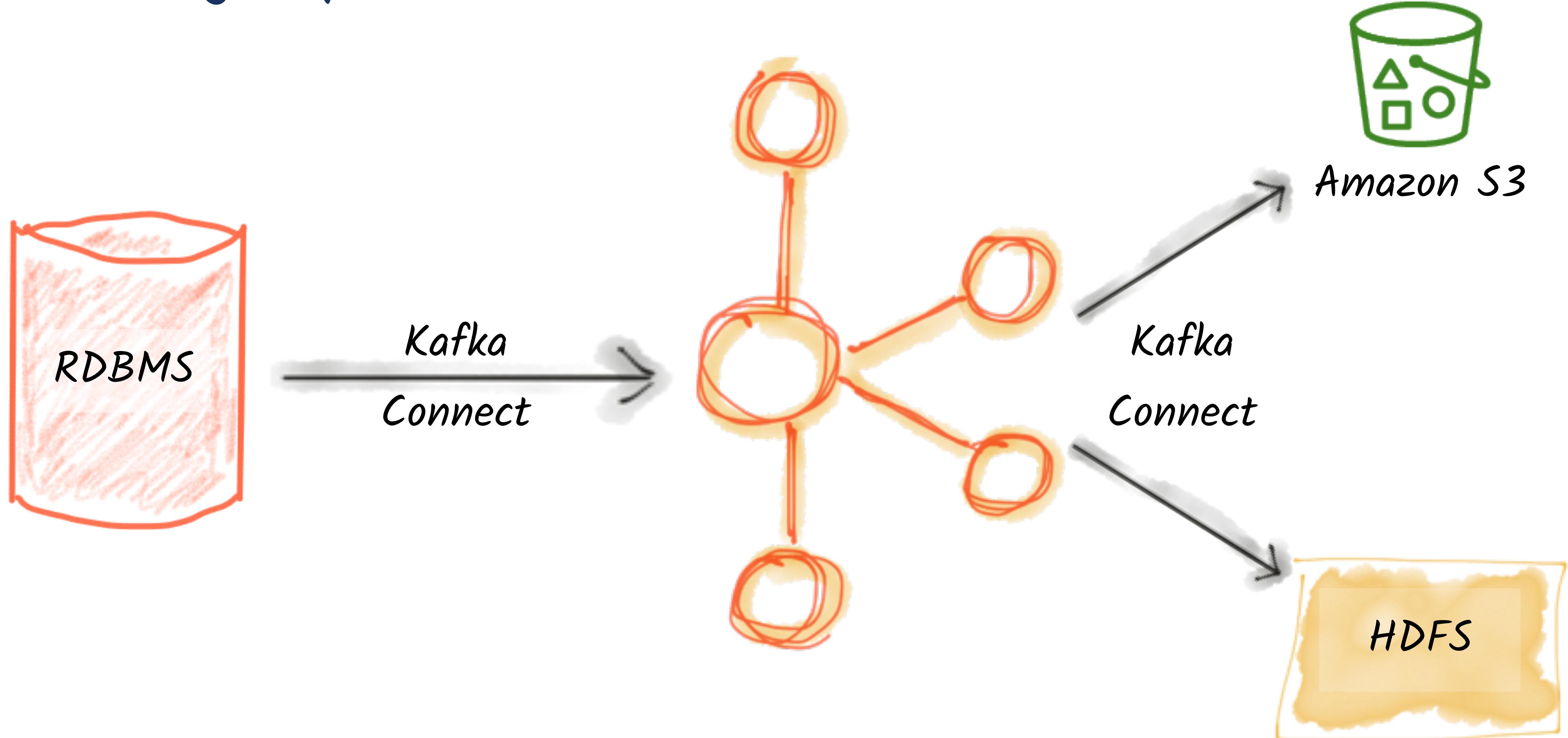
```
    "jdbc:mysql://asgard:3306/demo",
```

```
"table.whitelist":
```

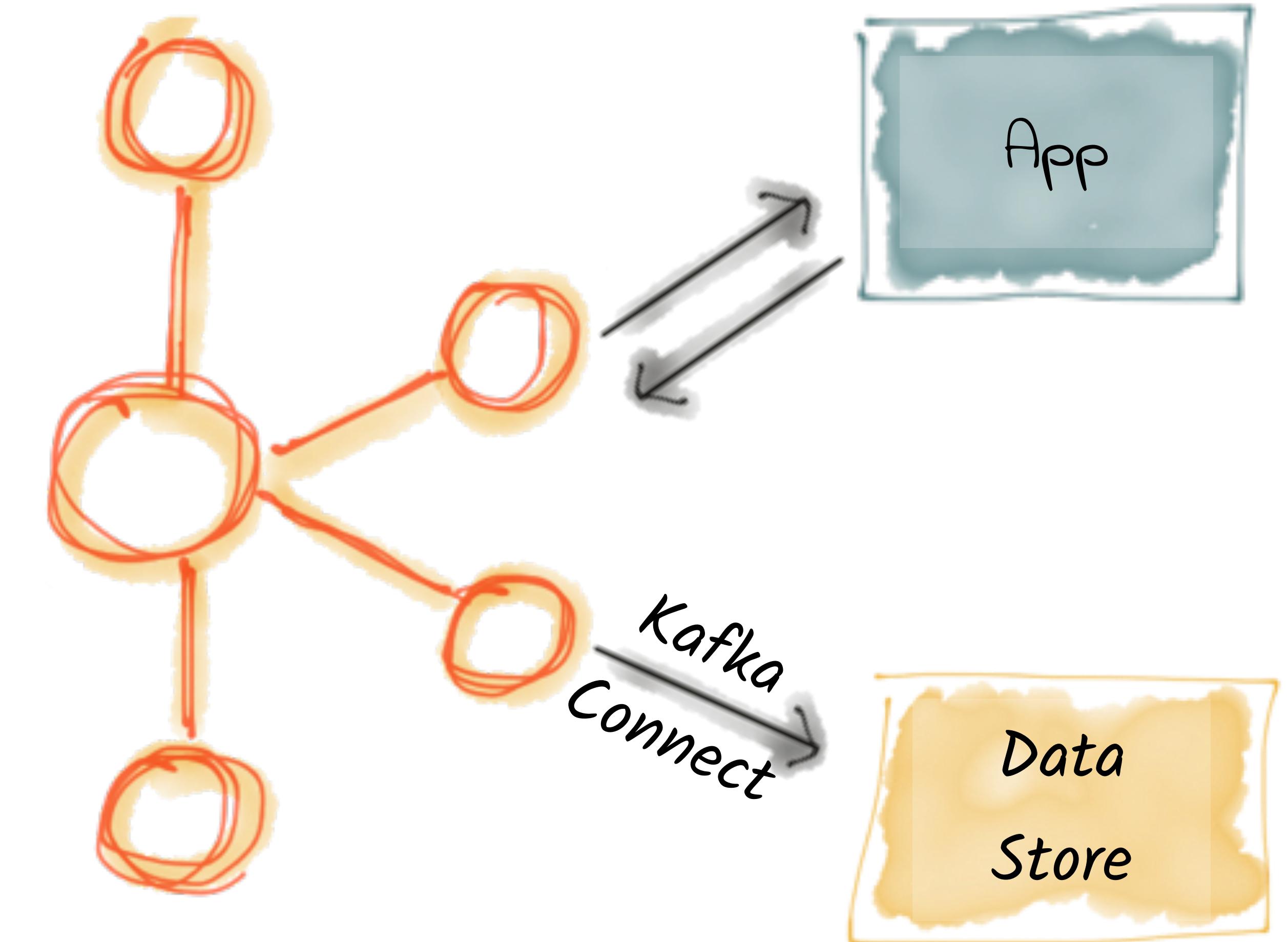
```
    "sales,orders,customers"
```

```
}
```

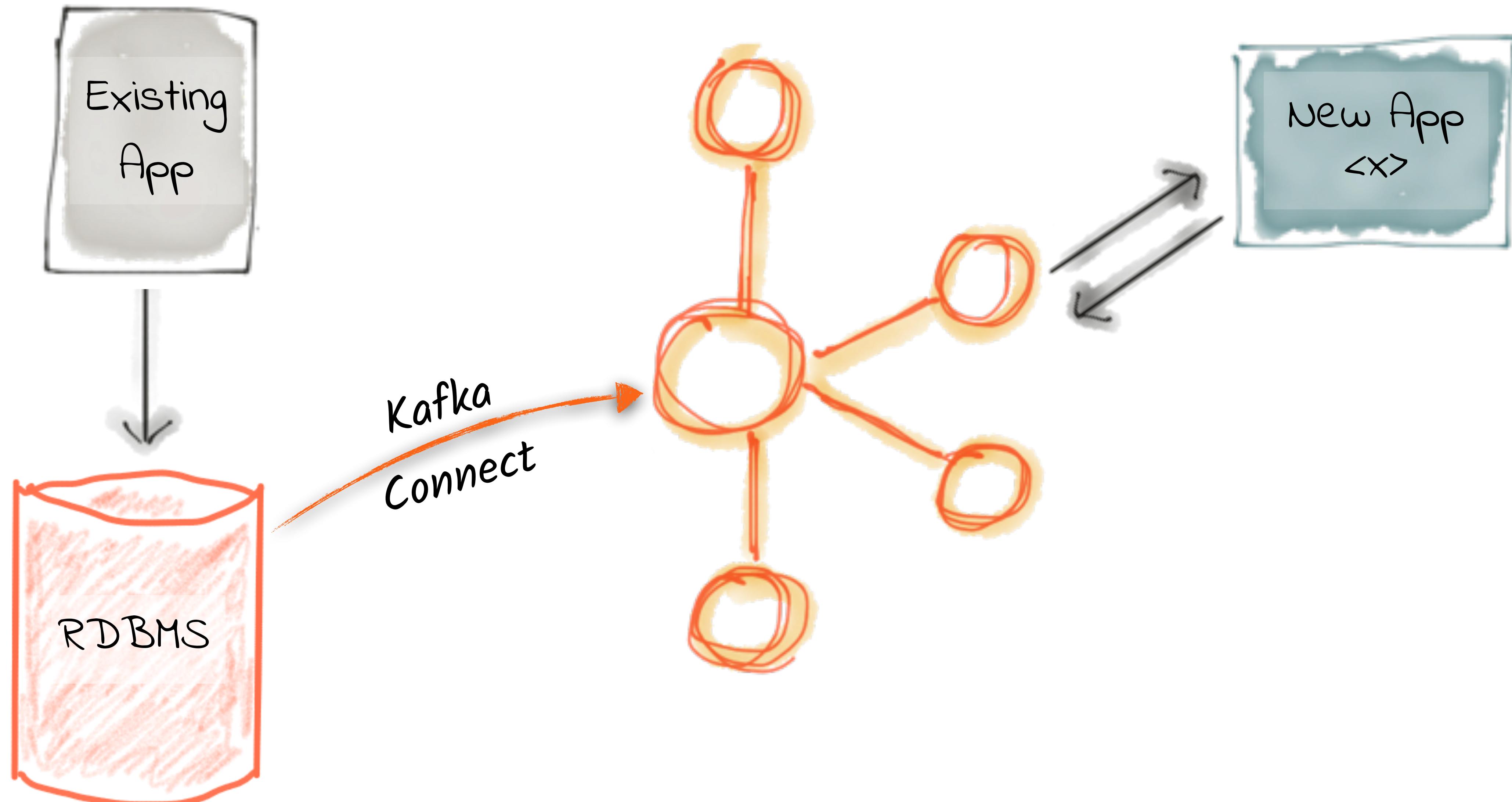
# Streaming Pipelines



# Writing to data stores from Kafka



# Evolve processing from old systems to new



# Demo

<http://rmoff.dev/kafka-connect-code>

# Configuring

## Kafka

## Connect

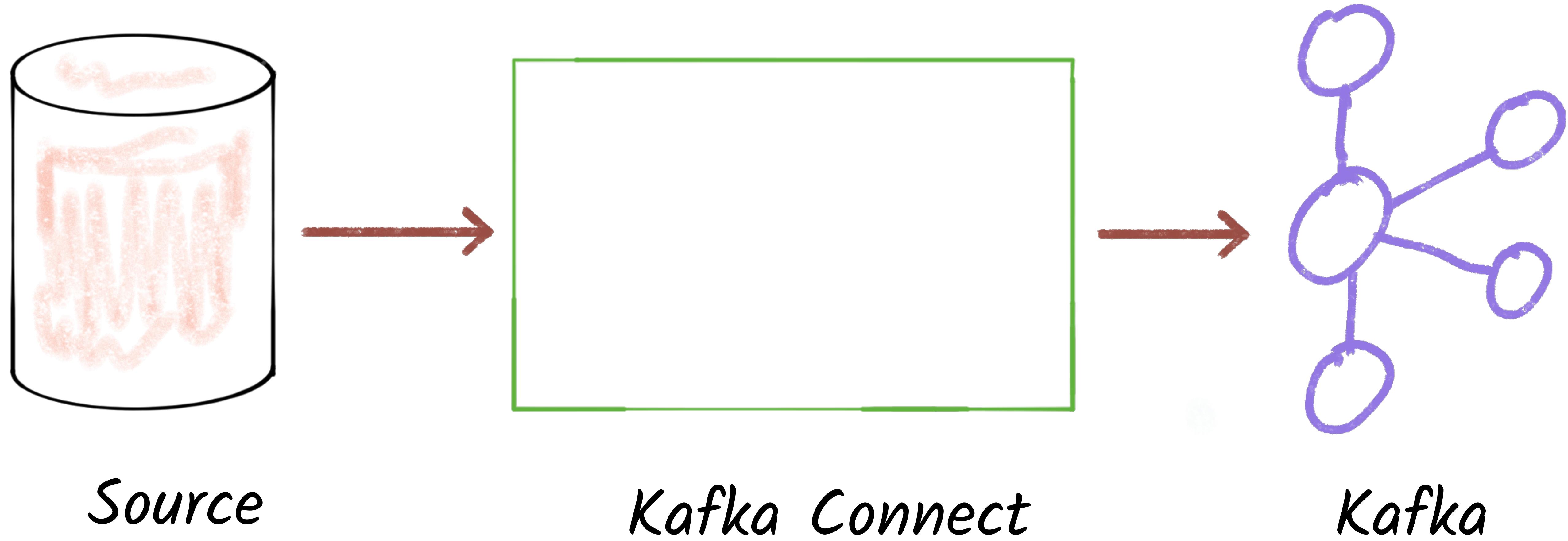
*Inside the API - connectors, transforms, converters*



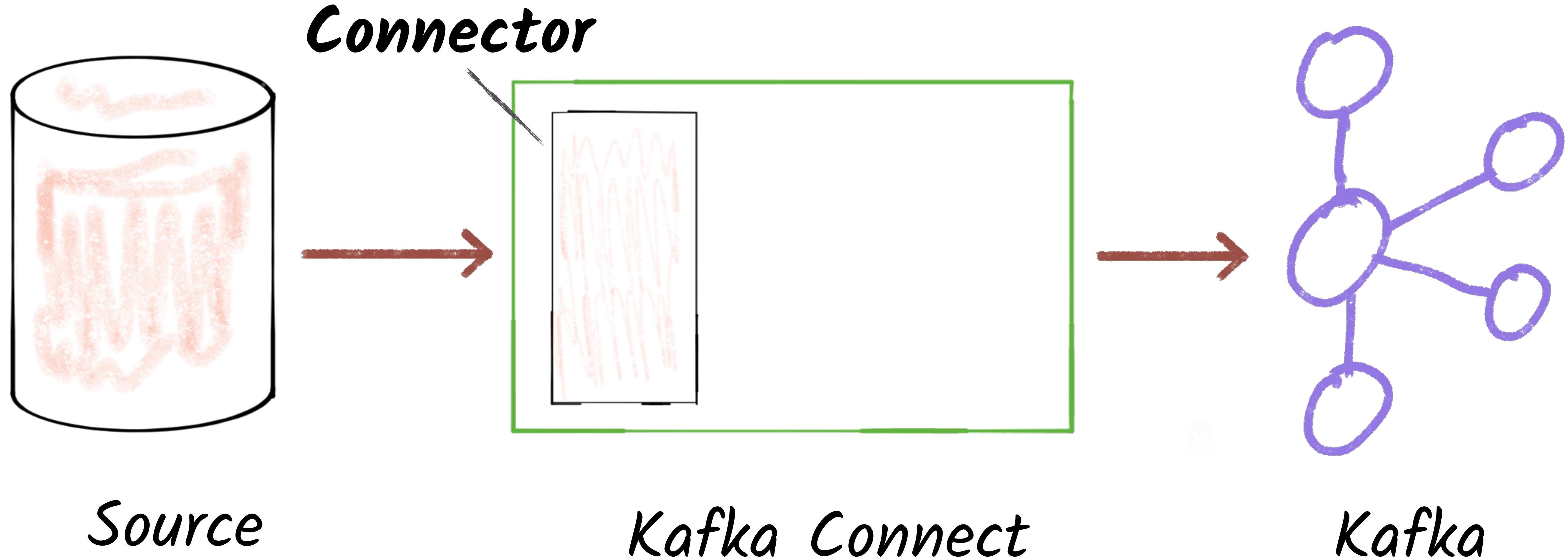
CONFLUENT

@rmoff

# Kafka Connect basics



# Connectors

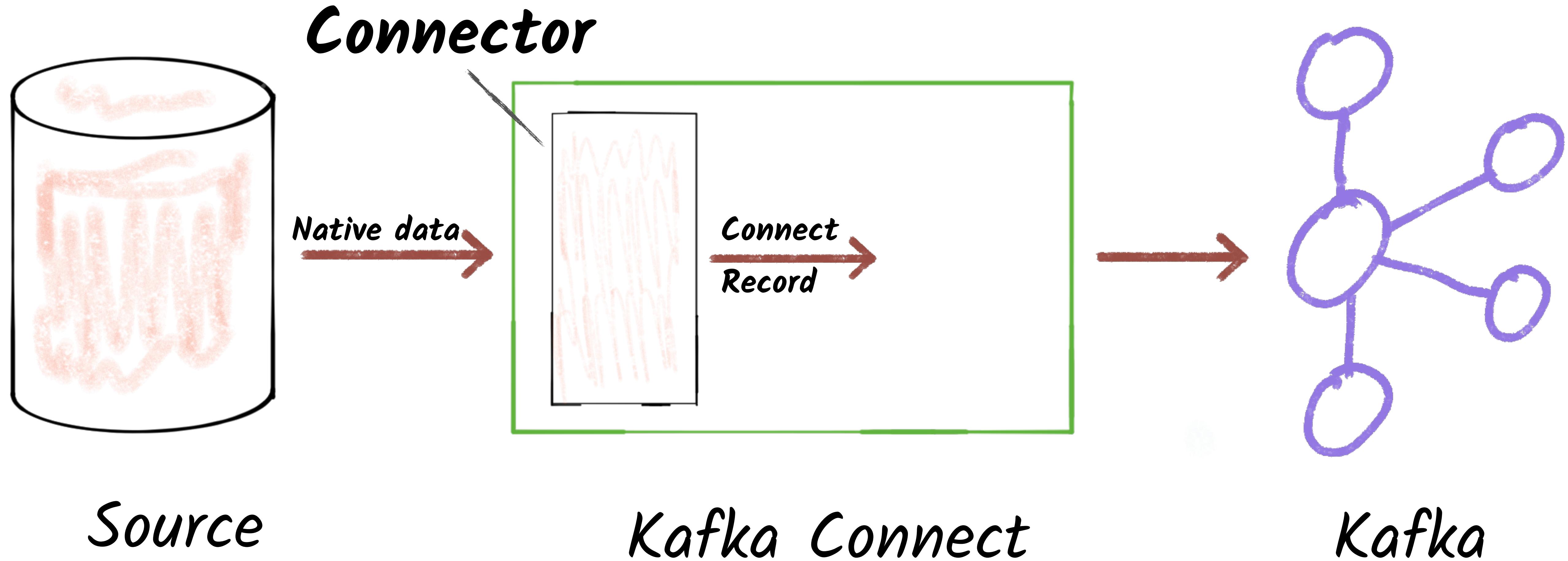


# Connectors

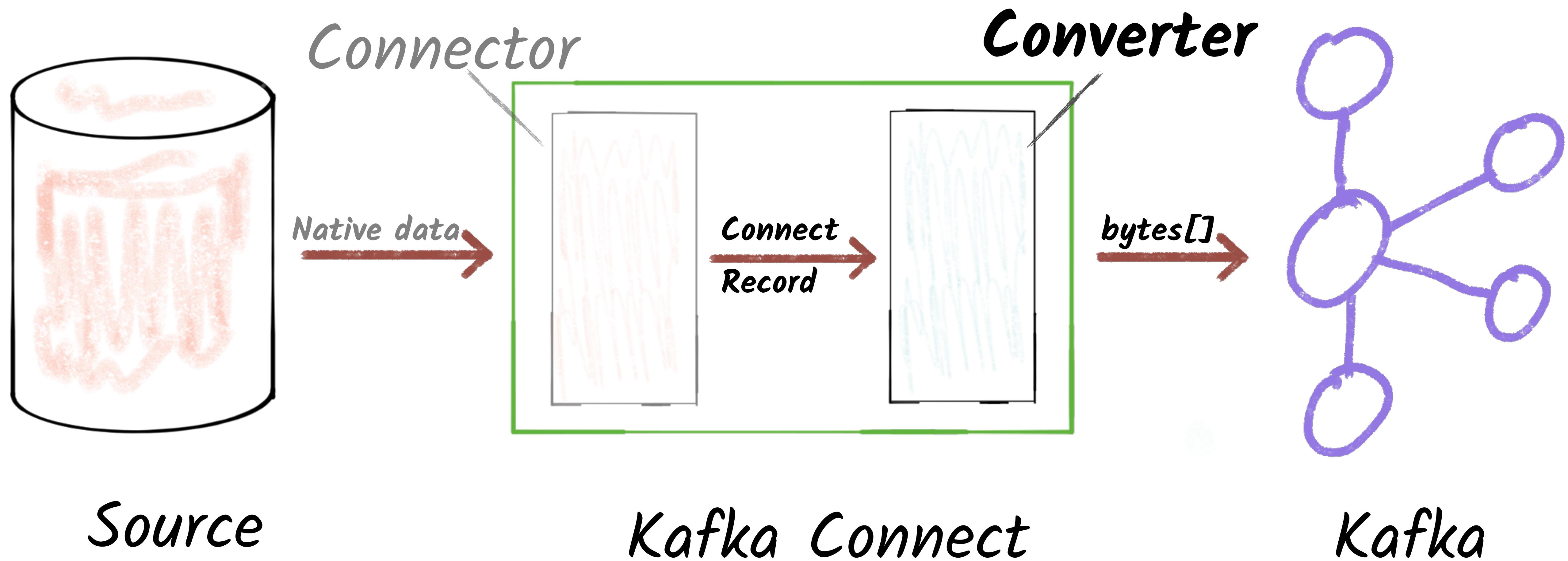
```
"config": {  
    [...]  
    "connector.class": "io.confluent.connect.jdbc.JdbcSinkConnector"  
    "connection.url": "jdbc:postgresql://postgres:5432/",  
    "topics": "asgard.demo.orders",  
}
```



# Connectors



# Converters



# Serialisation & Schemas

Avro



Protobuf



JSON  
Schema



JSON



CSV



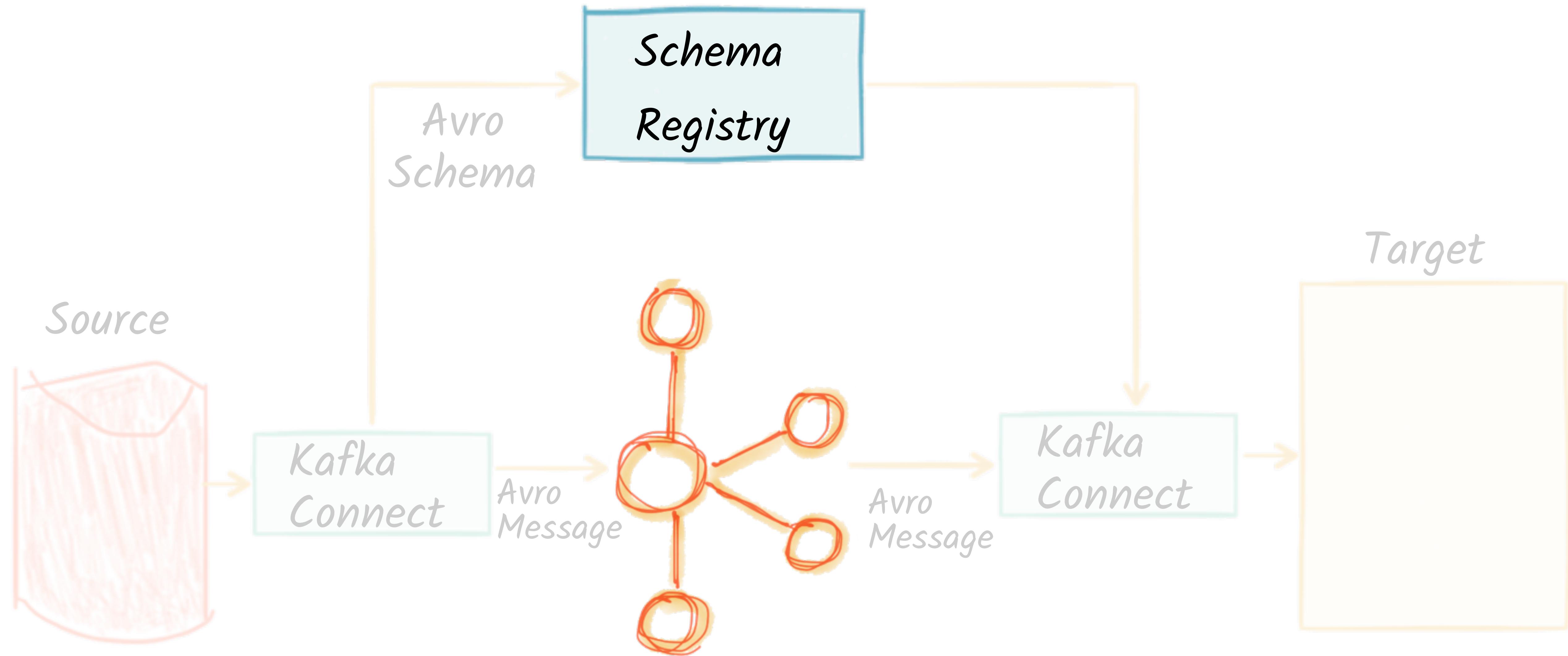
Gwen (Chen) Shapira  
@gwenshap

If your dev process doesn't validate schema compatibility somewhere between your IDE and production - you are screwed and don't know it.

5:50 AM - 5 Apr 2017

<https://rmoff.dev/qcon-schemas>

# The Confluent Schema Registry



# Converters

**key.converter=org.apache.kafka.connect.storage.StringConverter**

**value.converter=io.confluent.connect.avro.AvroConverter**

**value.converter.schema.registry.url=http://localhost:8081**

*Set as a global default per-worker; optionally can be overridden per-connector*

# What about internal converters?

`value.converter=org.apache.kafka.connect.json.JsonConverter`

`internal.value.converter=org.apache.kafka.connect.json.JsonConverter`

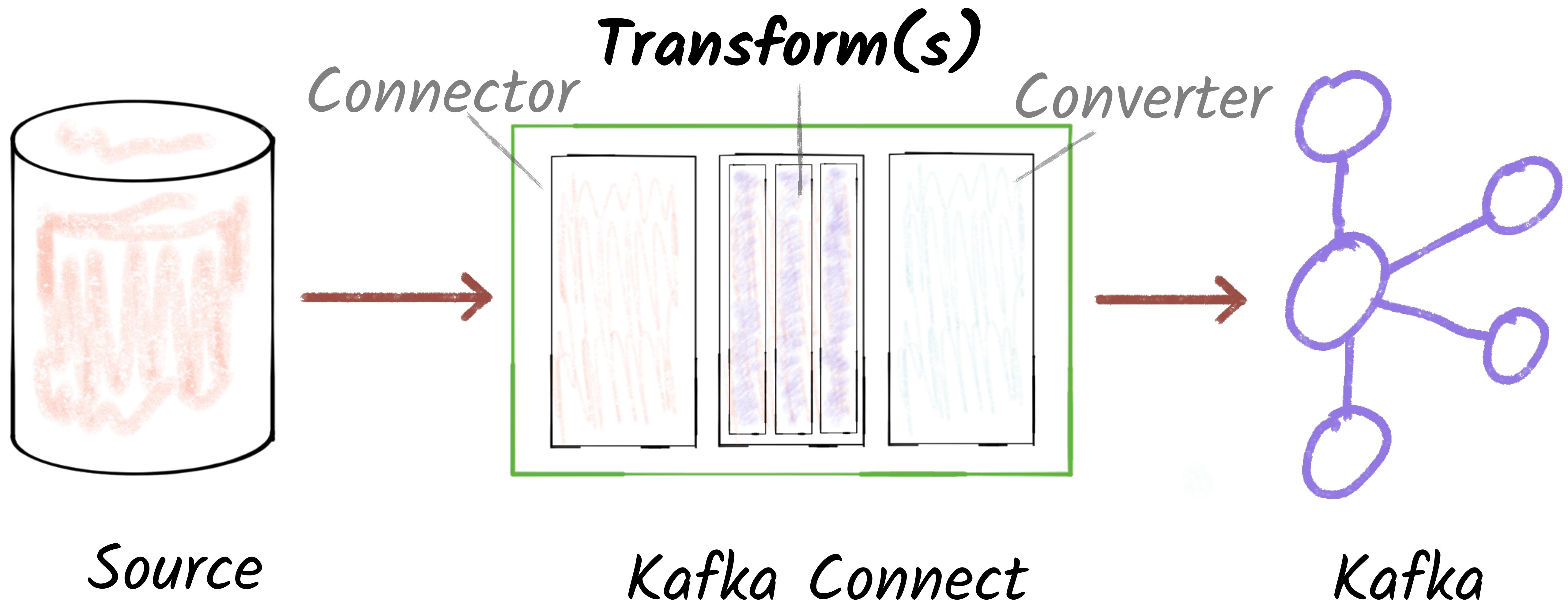
`key.internal.value.converter=org.apache.kafka.connect.json.JsonConverter`

`value.internal.value.converter=org.apache.kafka.connect.json.JsonConverter`

`key.internal.value.converter.bork.bork.bork=org.apache.kafka.connect.json.JsonConverter`

`key.internal.value.please.just.work.converter=org.apache.kafka.connect.json.JsonConverter`

# Single Message Transforms



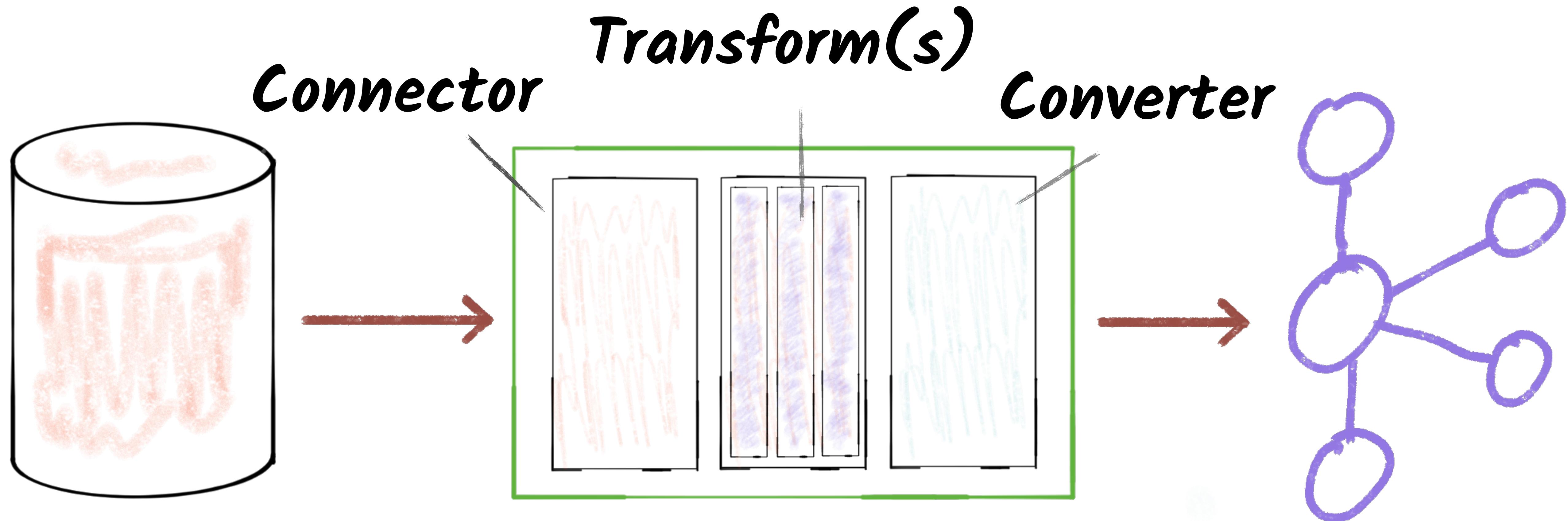
# Single Message Transforms

```
"config": {  
    [...]  
    "transforms": "addDateToTopic, labelFooBar",  
    "transforms.addDateToTopic.type": "org.apache.kafka.connect.transforms.TimestampRouter",  
    "transforms.addDateToTopic.topic.format": "${topic}-${timestamp}",  
    "transforms.addDateToTopic.timestamp.format": "YYYYMM",  
    "transforms.labelFooBar.type": "org.apache.kafka.connect.transforms.ReplaceField$Value",  
    "transforms.labelFooBar.renames": "delivery_address:shipping_address",  
}
```

*Do these transforms*

*Transforms config*    *Config per transform*

# Extensible



# Confluent Hub

— Confluent Hub —

## Discover Kafka® connectors and more

 jdbc

**Filters**

Plugin type ⓘ

- Sink
- Source
- Transform
- Converter

Enterprise support ⓘ

- Confluent supported
- Partner supported
- None

**Results (1)**

[+ Submit a plugin](#)

**Kafka Connect JDBC**

SINK, SOURCE CONNECTOR

The JDBC source and sink connectors allow you to exchange data between relational databases and Kafka. The JDBC source connector allows you to import data from any relational database with a JDBC driver into Kafka topics



Enterprise support:	Verification:	License:
Confluent supported	Confluent built	Free
Installation:	Author:	Version:
Confluent Hub CLI, Download	Confluent, Inc.	5.4.1

hub.confluent.io

# Deploying

# Kafka

# Connect

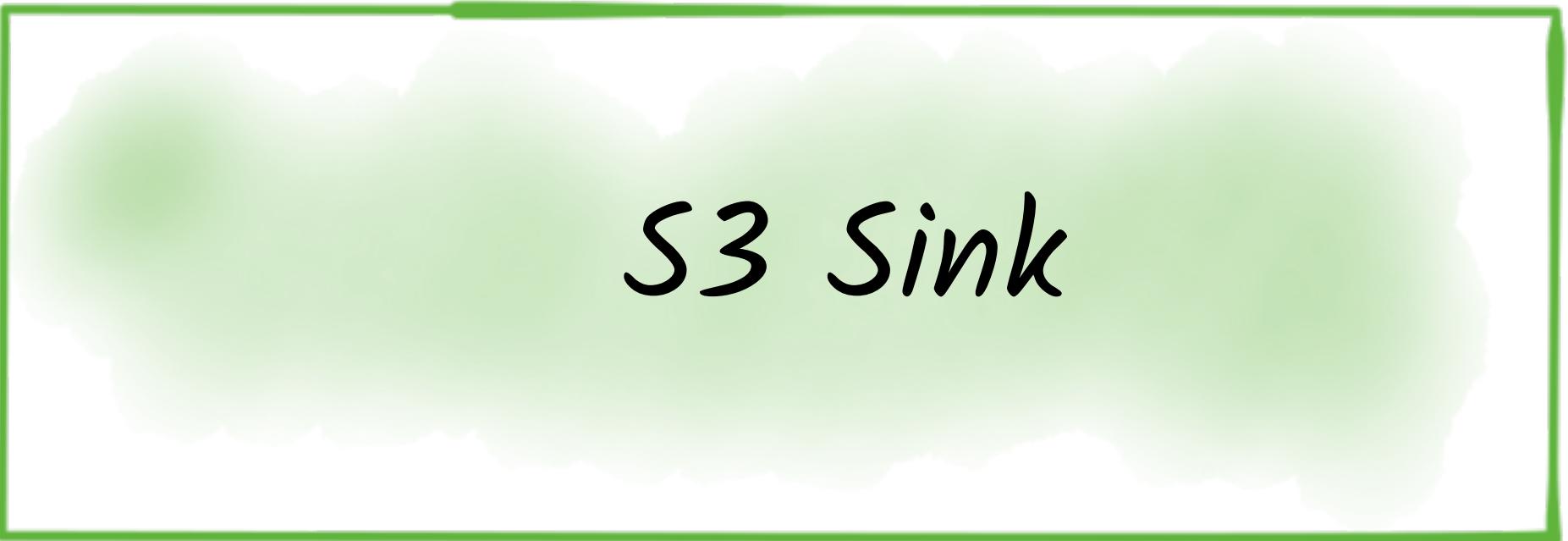
*Connectors, Tasks, and Workers*



CONFLUENT

@rmoff

# *Connectors and Tasks*



# Connectors and Tasks

JDBC Source

S3 Sink

S3 Task #1

JDBC Task #1



# Connectors and Tasks

JDBC Source

S3 Sink

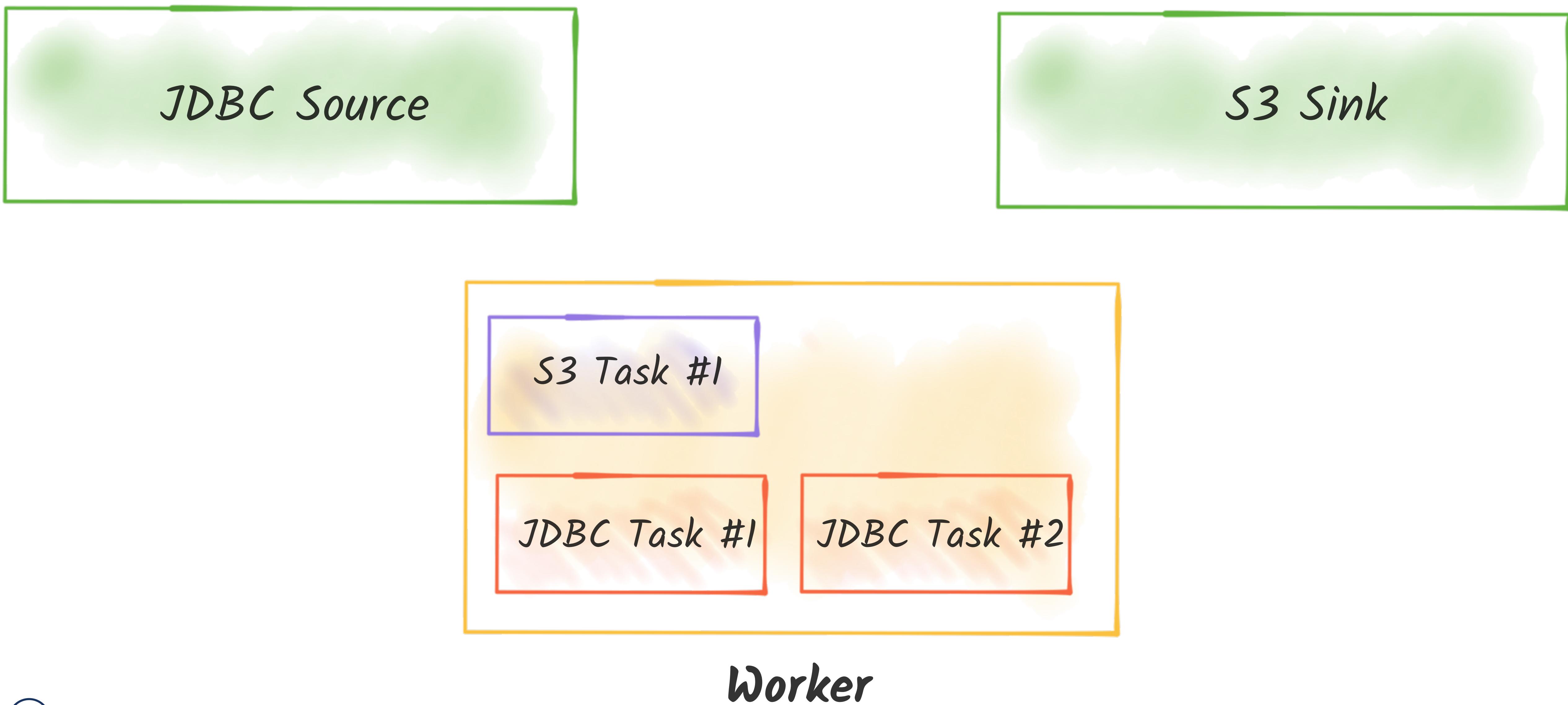
S3 Task #1

JDBC Task #1

JDBC Task #2



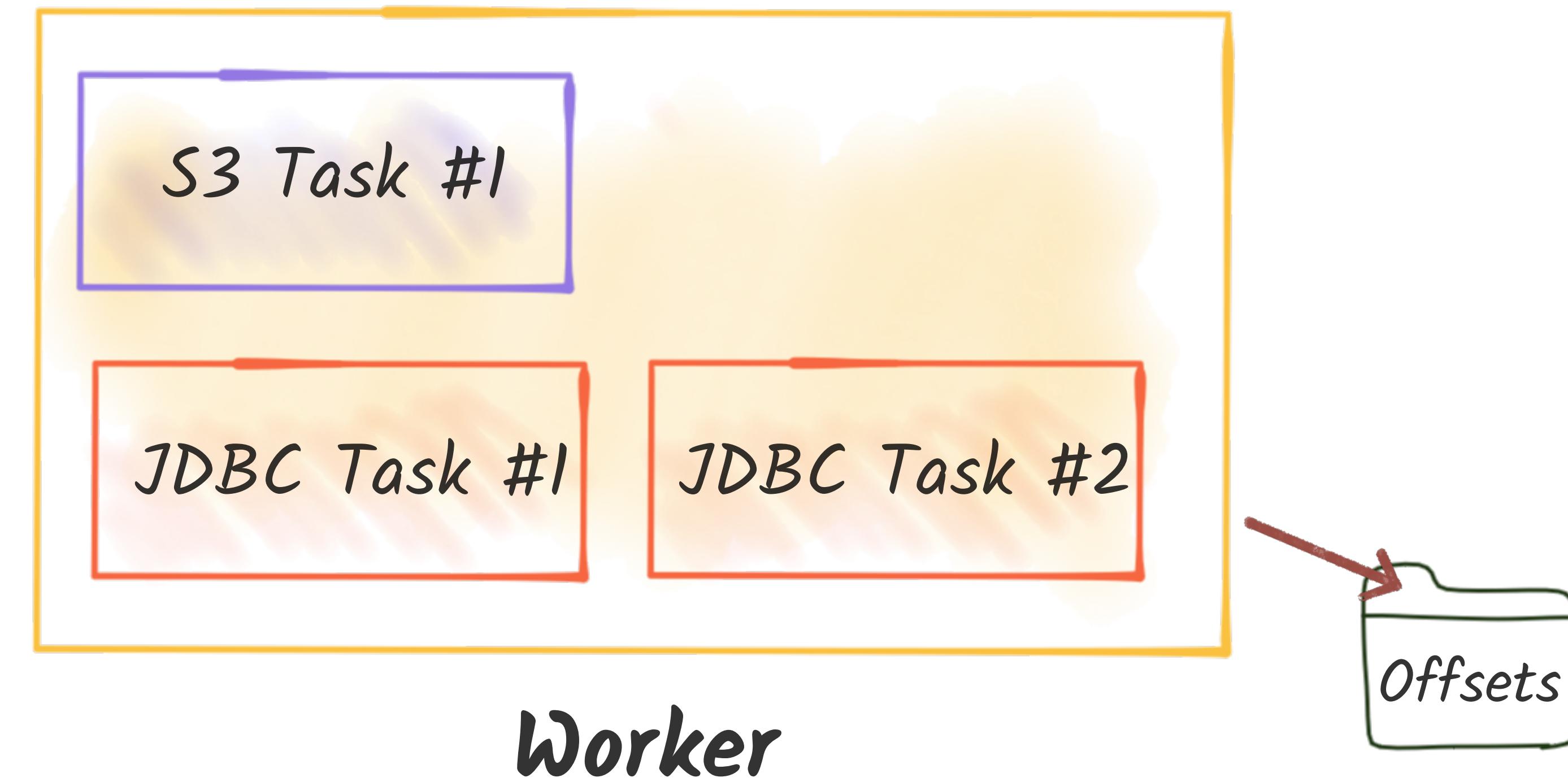
# Tasks and Workers





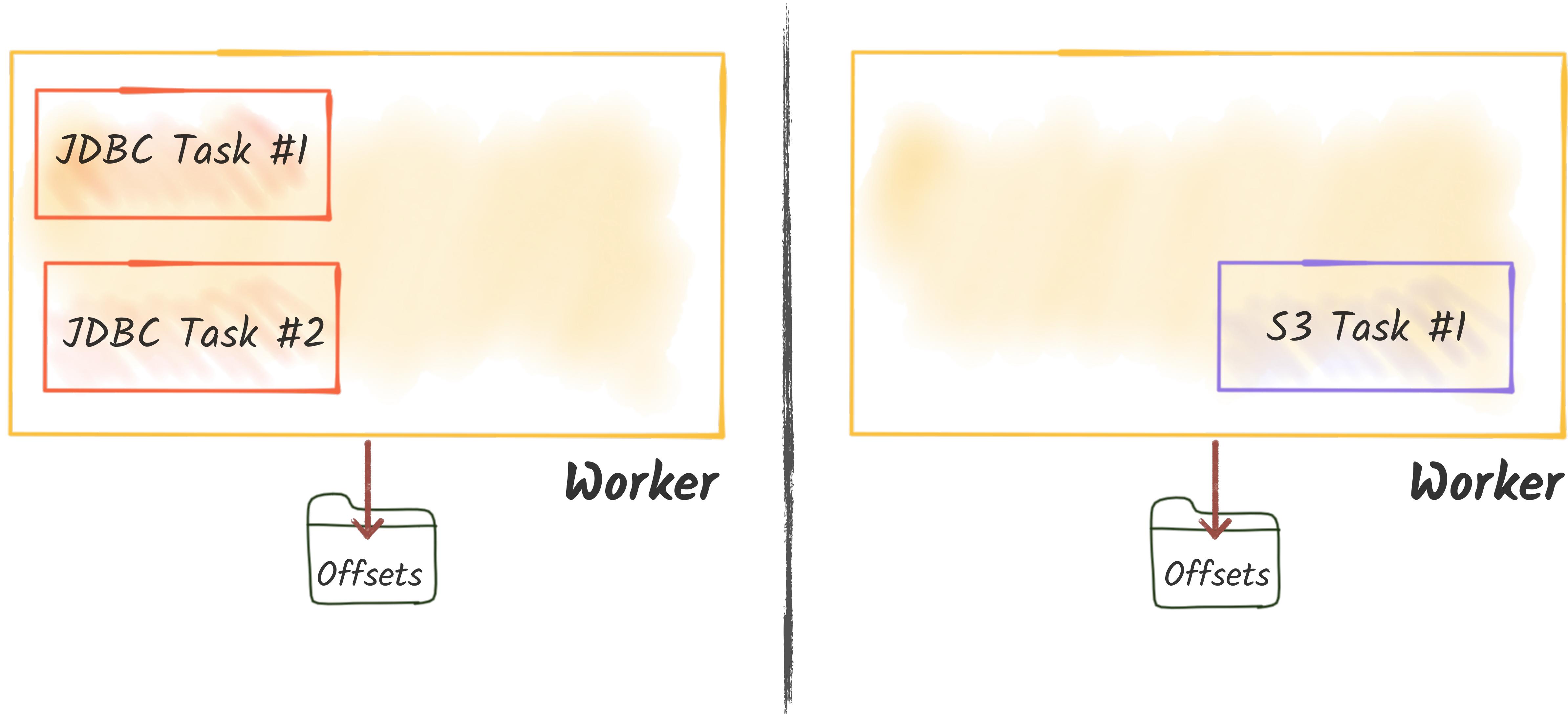
# Kafka Connect Standalone Worker

Fault-tolerant? Nope.



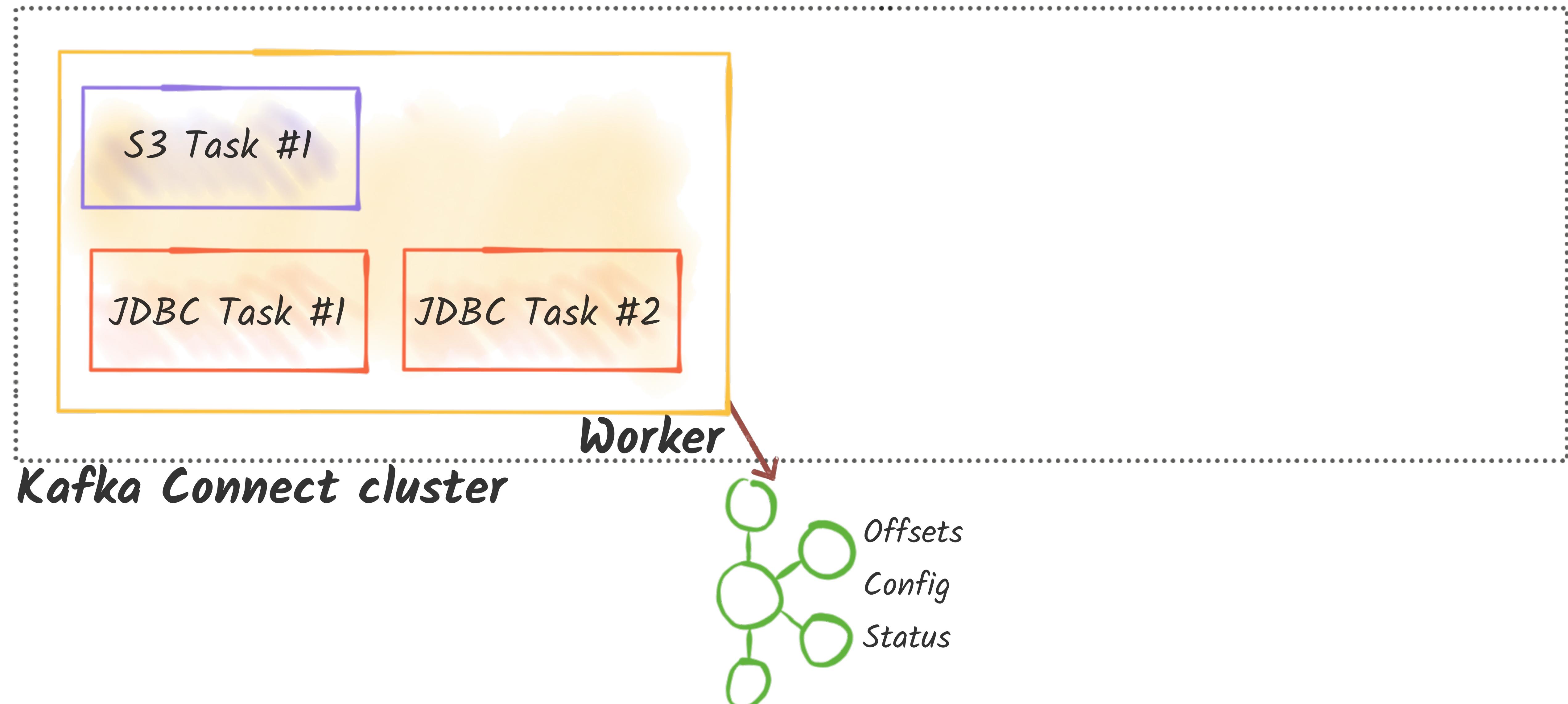
# "Scaling" the Standalone Worker

Fault-tolerant? Nope.



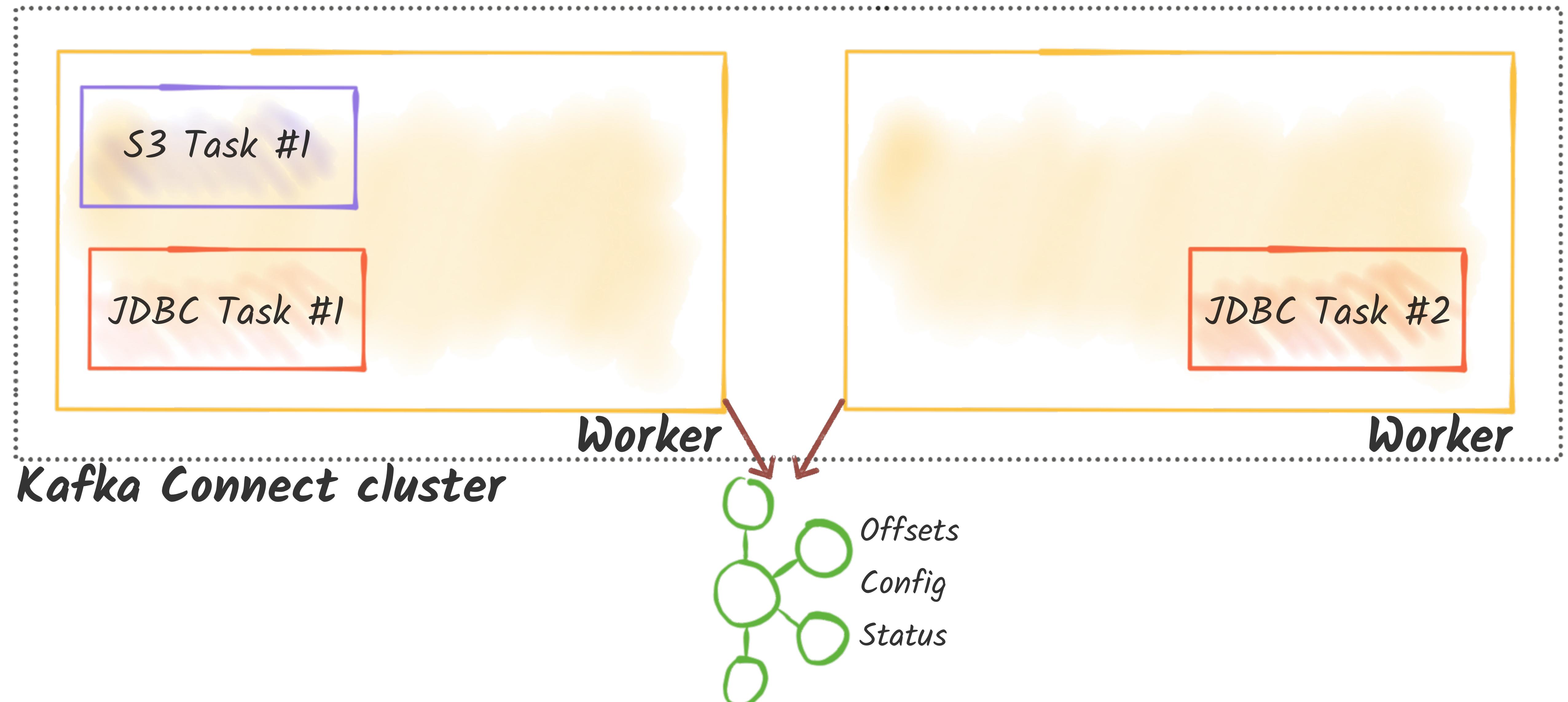
# Kafka Connect Distributed Worker

Fault-tolerant? Yeah!

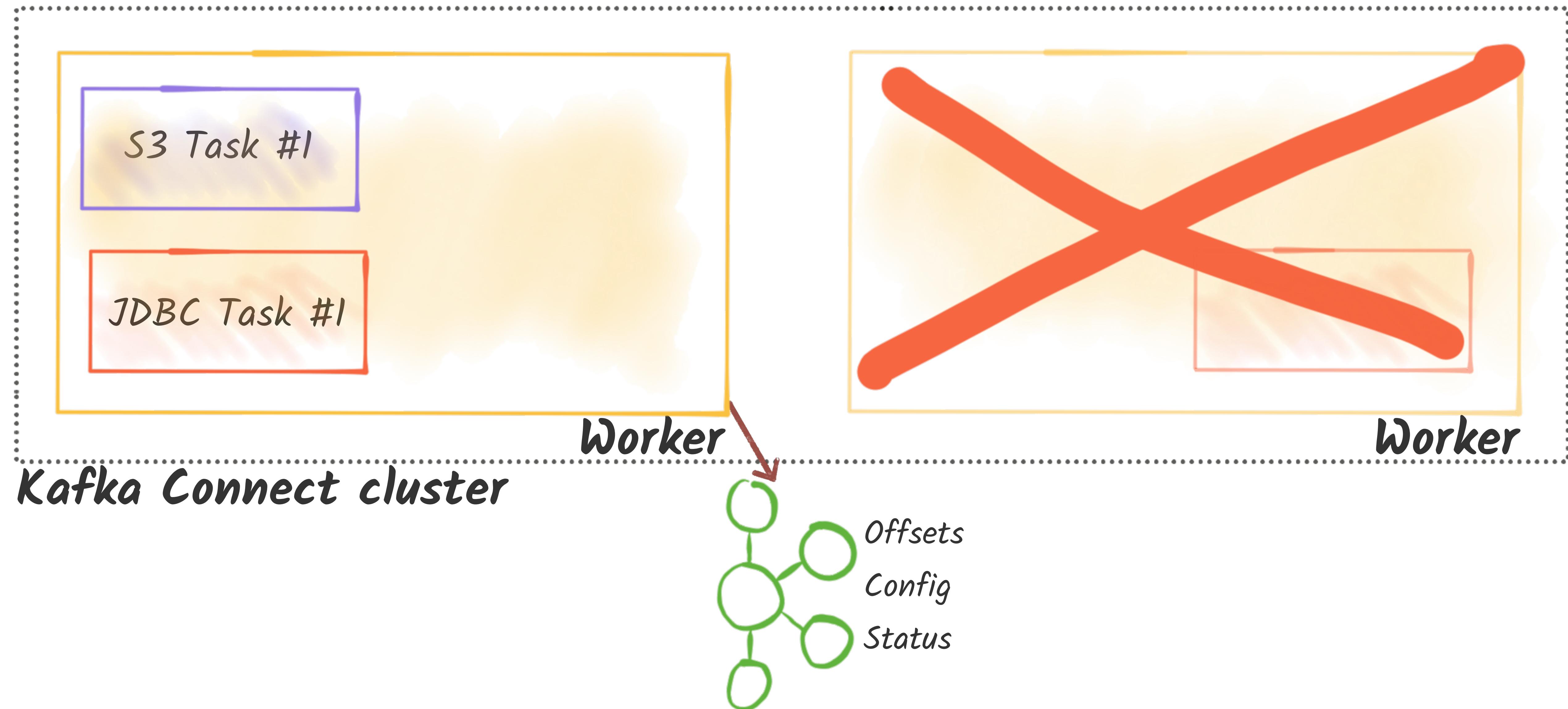


# Scaling the Distributed Worker

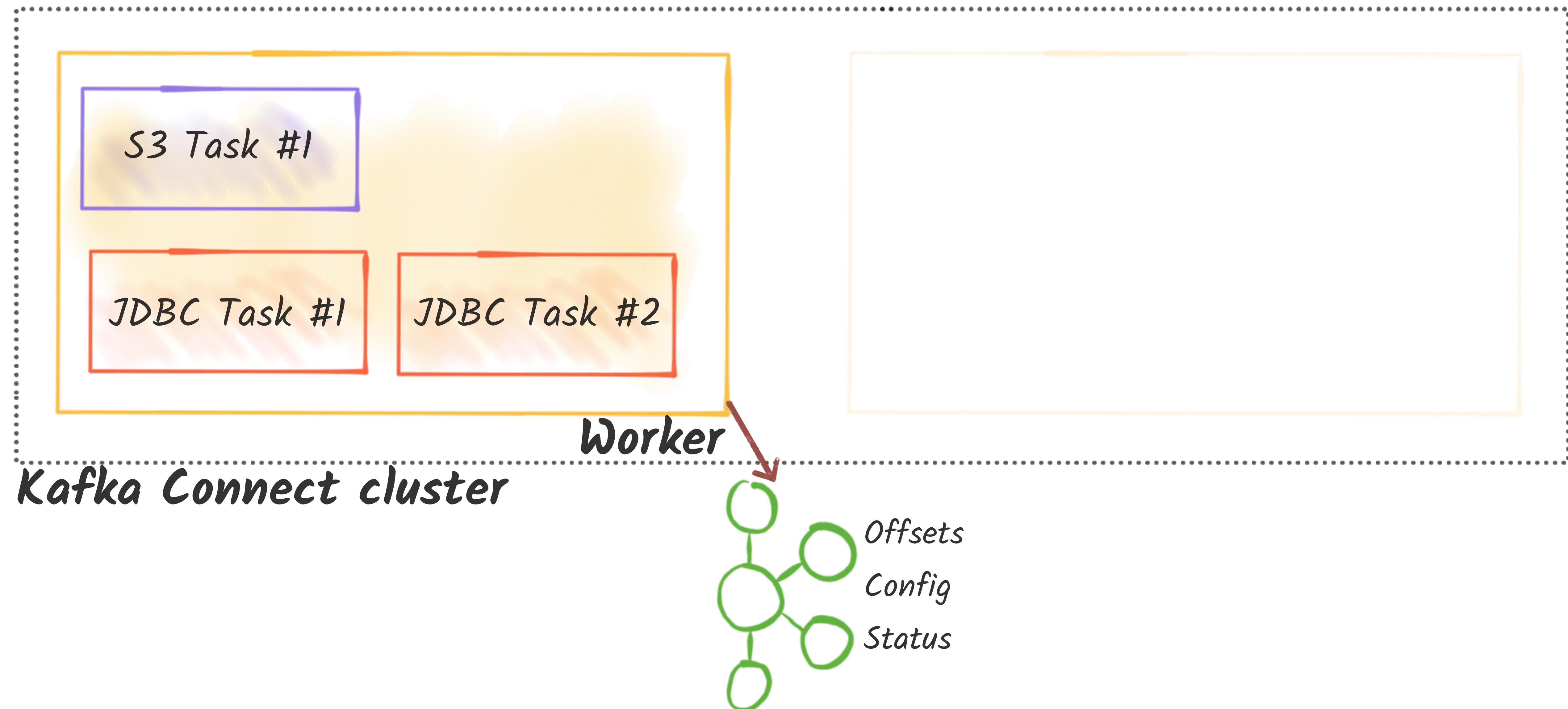
Fault-tolerant? Yeah!



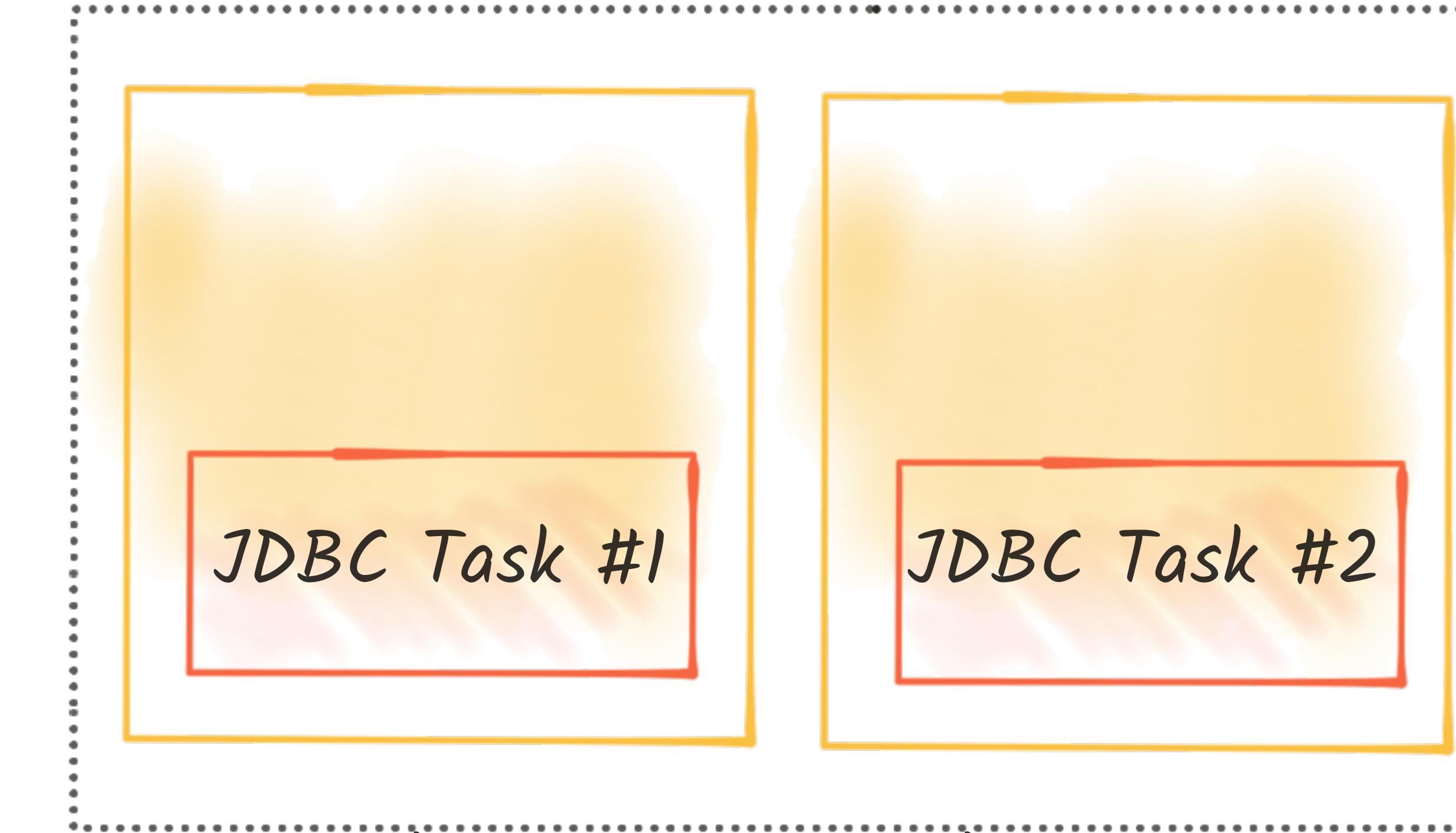
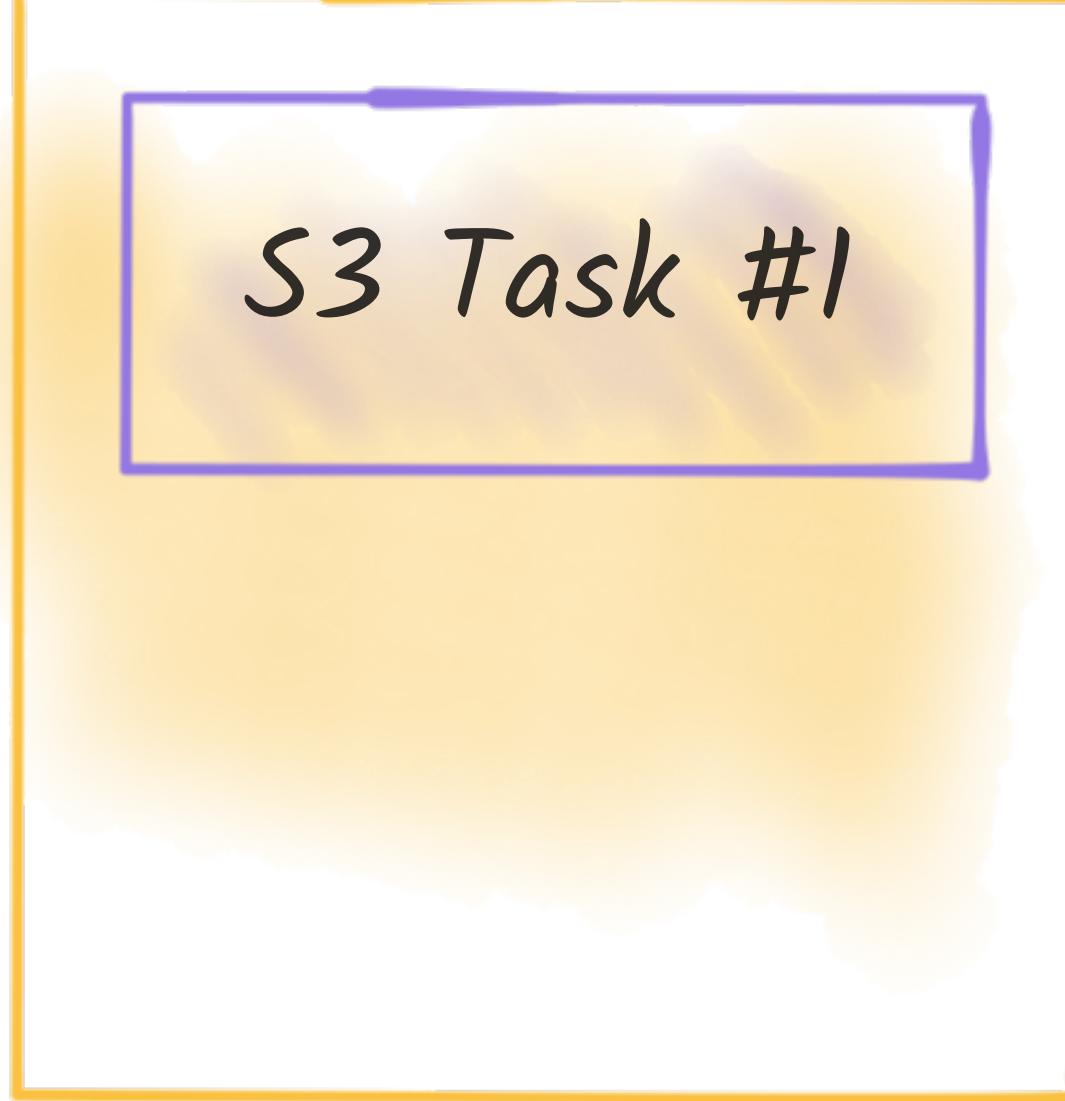
# Distributed Worker - fault tolerance



# Distributed Worker - fault tolerance



# Multiple Distributed Clusters



Kafka Connect cluster #1

Kafka Connect cluster #2



# *Containers*



# Kafka Connect images on Docker Hub

The screenshot shows a Docker Hub search result for the image 'confluentinc/cp-kafka-connect-base'. The page includes the Docker Hub logo, a search bar, and the image details: name, author, update time, description, and type.

**confluentinc/cp-kafka-connect-base** ☆

By [confluentinc](#) • Updated 3 days ago

Confluent Docker Base Image for Kafka Connect

Container

**confluentinc/cp-kafka-connect-base**

# *Adding connectors to a container*

Confluent Hub

JAR



# At runtime

kafka-connect:

```
image: confluentinc/cp-kafka-connect-base:6.0.0
```

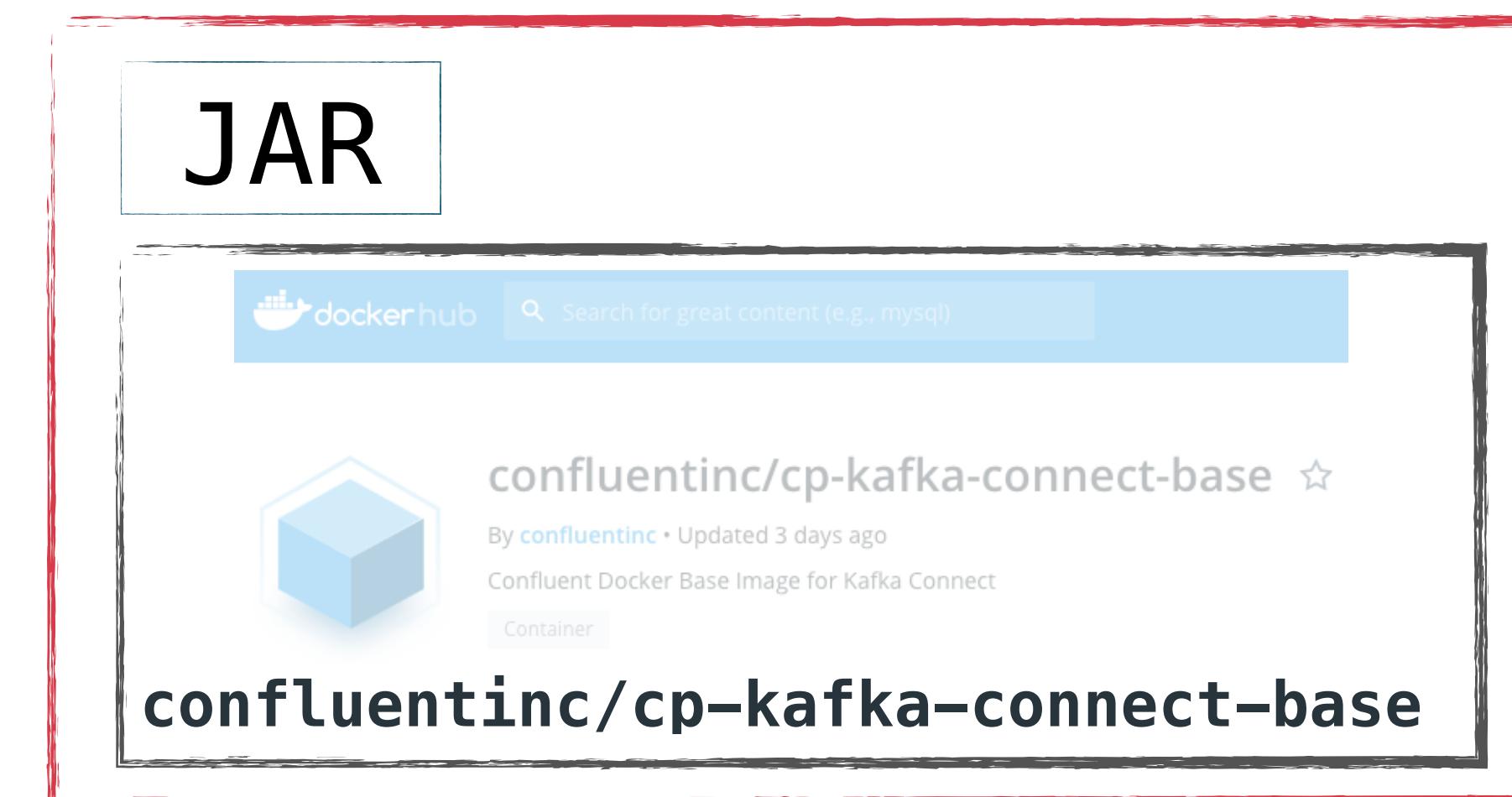
environment:

```
CONNECT_PLUGIN_PATH: '/usr/share/java,/usr/share/confluent-hub-components'
```

command:

- bash
- -c
- |

```
confluent-hub install --no-prompt neo4j/kafka-connect-neo4j:1.0.0  
/etc/confluent/docker/run
```



<http://rmoff.dev/ksln19-connect-docker>

@rmoff

# Build a new image

```
FROM confluentinc/cp-kafka-connect-base:6.0.0
ENV CONNECT_PLUGIN_PATH="/usr/share/java,/usr/share/confluent-hub-components"
RUN confluent-hub install --no-prompt neo4j/kafka-connect-neo4j:1.0.0
```



# Automating connector creation



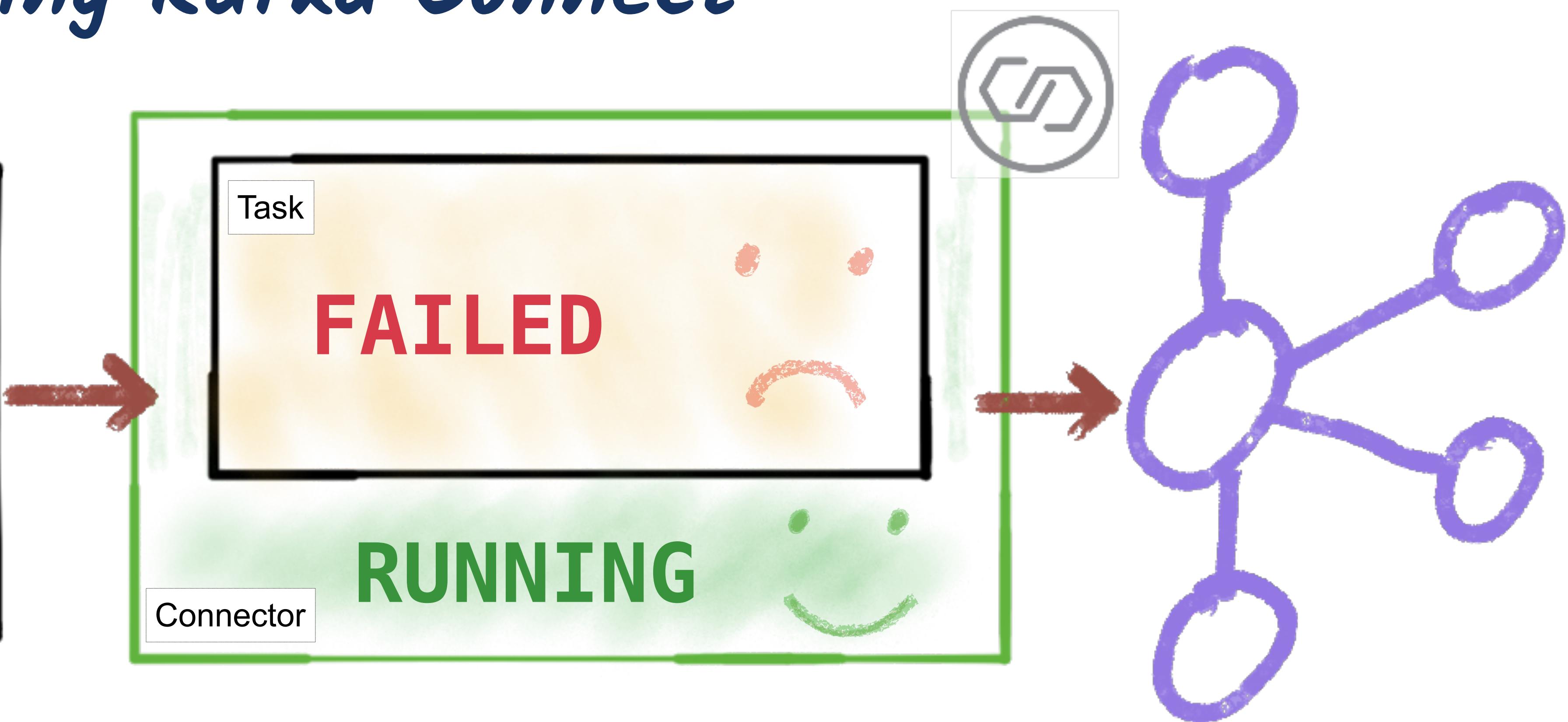
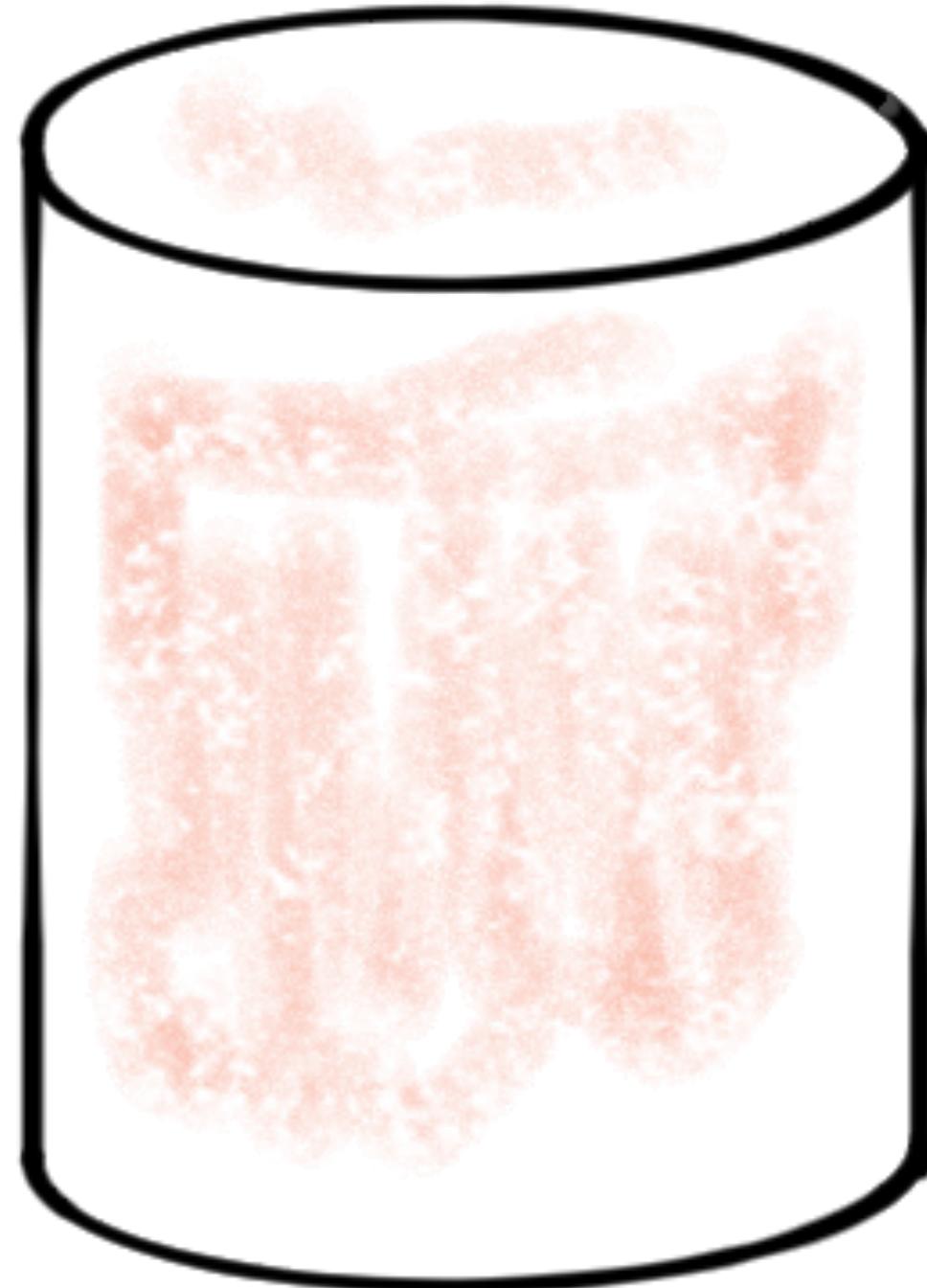
<http://rmoff.dev/ksln19-connect-docker>

```
# Launch Kafka Connect
/etc/confluent/docker/run &
#
# Wait for Kafka Connect listener
while [ $$($ curl -s -o /dev/null -w %{http_code} http://$$CONNECTOR_HOST:8083 ) -ne 200 ]
do
    echo -e $(date) " Kafka Connect listener HTTP state: " $$($ curl -s -o /dev/null -w %{http_code} http://$$CONNECTOR_HOST:8083 )
    sleep 5
done
#
# Create JDBC Source connector
curl -X POST http://localhost:8083/connectors -H "Content-Type: application/json" -d'
{
    "name": "jdbc_source_mysql_00",
    "config": {
        "connector.class": "io.confluent.connect.jdbc.JdbcSourceConnector",
        "connection.url": "jdbc:mysql://mysql:3306/default",
        "connection.user": "root",
        "connection.password": "password",
        "topic.prefix": "mysql_",
        "poll.interval.ms": 10000,
        "transforms": "mysqlRowKey"
    }
}'
```

# Troubleshooting Kafka Connect



# Troubleshooting Kafka Connect



```
$ curl -s "http://localhost:8083/connectors/source-debezium-orders/status" | \
jq '.connector.state'
"RUNNING"
$ curl -s "http://localhost:8083/connectors/source-debezium-orders/status" | \
jq '.tasks[0].state'
"FAILED"
```



# Troubleshooting Kafka Connect

```
curl -s "http://localhost:8083/connectors/source-debezium-orders-00/status"  
| jq '.tasks[0].trace'
```

```
"org.apache.kafka.connect.errors.ConnectException\n\tat  
io.debezium.connector.mysql.AbstractReader.wrap(AbstractReader.java:230)\n\tat  
io.debezium.connector.mysql.AbstractReader.failed(AbstractReader.java:197)\n\tat  
io.debezium.connector.mysql.BinlogReader$ReaderThreadLifecycleListener.onCommunicationFailure(BinlogReader.java:1018)\n\tat  
com.github.shyiko.mysql.binlog.BinaryLogClient.listenForEventPackets(BinaryLogClient.java:950)\n\tat  
com.github.shyiko.mysql.binlog.BinaryLogClient.connect(BinaryLogClient.java:580)\n\tat  
com.github.shyiko.mysql.binlog.BinaryLogClient$7.run(BinaryLogClient.java:825)\n\tat  
java.lang.Thread.run(Thread.java:748)\nCaused by: java.io.EOFException\n\tat  
com.github.shyiko.mysql.binlog.io.ByteArrayInputStream.read(ByteArrayInputStream.java:190)\n\tat  
com.github.shyiko.mysql.binlog.io.ByteArrayInputStream.readInteger(ByteArrayInputStream.java:46)\n\tat  
com.github.shyiko.mysql.binlog.event.serialization.EventHeaderV4Deserializer.deserialize(EventHeaderV4Deserializer.java  
:35)\n\tat  
com.github.shyiko.mysql.binlog.event.serialization.EventHeaderV4Deserializer.deserialize(EventHeaderV4Deserializer.java  
:27)\n\tat  
com.github.shyiko.mysql.binlog.event.serialization.EventDeserializer.nextEvent(EventDeserializer.java:212)\n\tat  
io.debezium.connector.mysql.BinlogReader$1.nextEvent(BinlogReader.java:224)\n\tat  
com.github.shyiko.mysql.binlog.BinaryLogClient.listenForEventPackets(BinaryLogClient.java:922)\n\t... 3 more\n"
```

*The log is the source of truth*

```
$ confluent local log connect
```

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

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

# Dynamic log levels

(Added in Apache Kafka 2.4 / Confluent Platform 5.4)

```
curl -s http://localhost:8083/admin/loggers/ | jq
{
  "org.apache.kafka.connect.runtime.rest": {
    "level": "WARN"
  },
  "org.reflections": {
    "level": "ERROR"
  },
  "root": {
    "level": "INFO"
  }
}
```

```
curl -s -X PUT http://localhost:8083/admin/loggers/io.debezium
-H "Content-Type:application/json" -d '{"level": "TRACE"}'
```

<https://rmoff.dev/kc-dynamic-log-level>



CONFLUENT

@rmoff

# Kafka Connect

```
ache.kafka.connect.runtime.WorkerSourceTask)
[2019-05-07 14:39:13,115] INFO WorkerSourceTask{id=source-debezium-orders-00-0} Finished commitOffsets successfully in 28 ms (org.apache.
kafka.connect.runtime.WorkerSourceTask)
[2019-05-07 14:39:13,116] ERROR WorkerSourceTask{id=source-debezium-orders-00-0} Task threw an uncaught and unrecoverable exception (org.
apache.kafka.connect.runtime.WorkerTask)
org.apache.kafka.connect.errors.ConnectException
    at io.debezium.connector.mysql.AbstractReader.wrap(AbstractReader.java:230)
    at io.debezium.connector.mysql.AbstractReader.failed(AbstractReader.java:197)
    at io.debezium.connector.mysql.BinlogReader$ReaderLifecycleListener.onCommunicationFailure(BinlogReader.java:1018)
    at com.github.shyiko.mysql.binlog.BinaryLogClient.listenForEventPackets(BinaryLogClient.java:950)
    at com.github.shyiko.mysql.binlog.BinaryLogClient.connect(BinaryLogClient.java:580)
    at com.github.shyiko.mysql.binlog.BinaryLogClient$7.run(BinaryLogClient.java:825)
    at java.lang.Thread.run(Thread.java:748)
Caused by: java.io.EOFException
    at com.github.shyiko.mysql.io.ByteArrayInputStream.read(ByteArrayInputStream.java:190)
    at com.github.shyiko.mysql.io.ByteArrayInputStream.readInt(ByteArrayInputStream.java:46)
    at com.github.shyiko.mysql.event.serialization.EventHeaderV4Deserializer.deserialize(EventHeaderV4Deserializer.java:35)
    at com.github.shyiko.mysql.event.serialization.EventHeaderV4Deserializer.deserialize(EventHeaderV4Deserializer.java:27)
    at com.github.shyiko.mysql.event.serialization.EventDeserializer.nextEvent(EventDeserializer.java:212)
    at io.debezium.connector.mysql.BinlogReader$1.nextEvent(BinlogReader.java:224)
    at com.github.shyiko.mysql.binlog.BinaryLogClient.listenForEventPackets(BinaryLogClient.java:922)
    ... 3 more
[2019-05-07 14:39:13,121] ERROR WorkerSourceTask{id=source-debezium-orders-00-0} Task is being killed and will not recover until manually restarted
(org.apache.kafka.connect.runtime.WorkerTask)
```



# Kafka Connect

"Task is being killed and will  
not recover until manually restarted"

*Symptom not Cause*

# Kafka Connect

```
apache.kafka.connect.runtime.WorkerSourceTask)
[2019-05-07 14:39:13,115] INFO WorkerSourceTask{id=source-debezium-orders-00-0} Finished commitOffsets successfully in 28 ms (org.apache.
apache.kafka.connect.runtime.WorkerSourceTask)
[2019-05-07 14:39:13,116] ERROR WorkerSourceTask{id=source-debezium-orders-00-0} Task threw an uncaught and unrecoverable exception (org.
apache.kafka.connect.runtime.WorkerTask)
org.apache.kafka.connect.errors.ConnectException
    at io.debezium.connector.mysql.AbstractReader.wrap(AbstractReader.java:230)
    at io.debezium.connector.mysql.AbstractReader.failed(AbstractReader.java:197)
    at io.debezium.connector.mysql.BinlogReader$ReaderLifecycleListener.onCommunicationFailure(BinlogReader.java:1018)
    at com.github.shyiko.mysql.binlog.BinaryLogClient.listenForEventPackets(BinaryLogClient.java:950)
    at com.github.shyiko.mysql.binlog.BinaryLogClient.connect(BinaryLogClient.java:580)
    at com.github.shyiko.mysql.binlog.BinaryLogClient$7.run(BinaryLogClient.java:825)
    at java.lang.Thread.run(Thread.java:748)
Caused by: java.io.EOFException
    at com.github.shyiko.mysql.binlog.io.ByteArrayInputStream.read(ByteArrayInputStream.java:190)
    at com.github.shyiko.mysql.binlog.io.ByteArrayInputStream.readInt(ByteArrayInputStream.java:46)
    at com.github.shyiko.mysql.event.serialization.EventHeaderV4Deserializer.deserialize(EventHeaderV4Deserializer.java:35)
    at com.github.shyiko.mysql.event.serialization.EventHeaderV4Deserializer.deserialize(EventHeaderV4Deserializer.java:27)
    at com.github.shyiko.mysql.event.serialization.EventDeserializer.nextEvent(EventDeserializer.java:212)
    at io.debezium.connector.mysql.BinlogReader$1.nextEvent(BinlogReader.java:224)
    at com.github.shyiko.mysql.binlog.BinaryLogClient.listenForEventPackets(BinaryLogClient.java:922)
    ... 3 more
[2019-05-07 14:39:13,121] ERROR WorkerSourceTask{id=source-debezium-orders-00-0} Task is being killed and will not recover until manually restarted
(org.apache.kafka.connect.runtime.WorkerTask)
```

①



# Error Handling and Dead Letter Queues



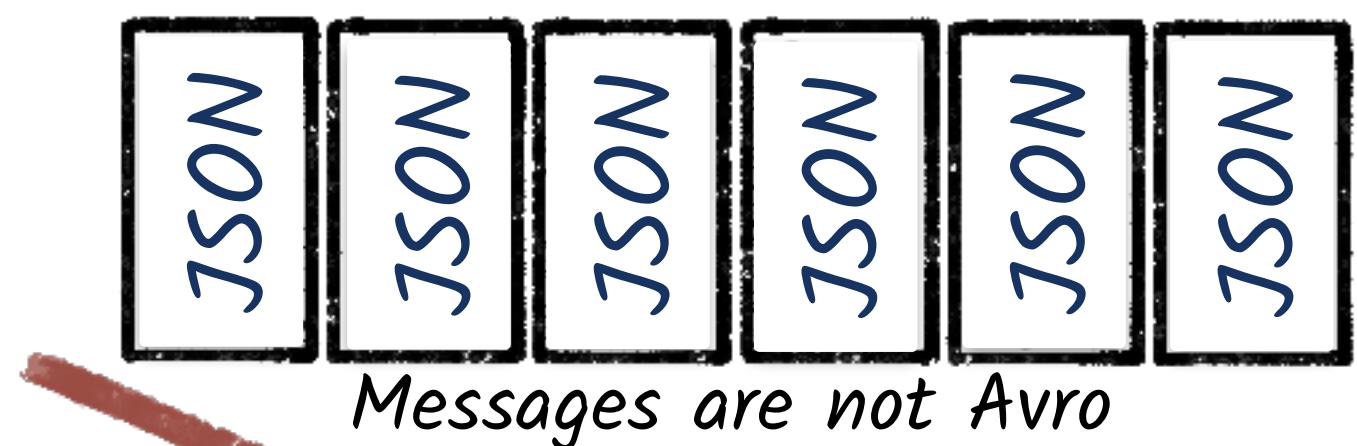
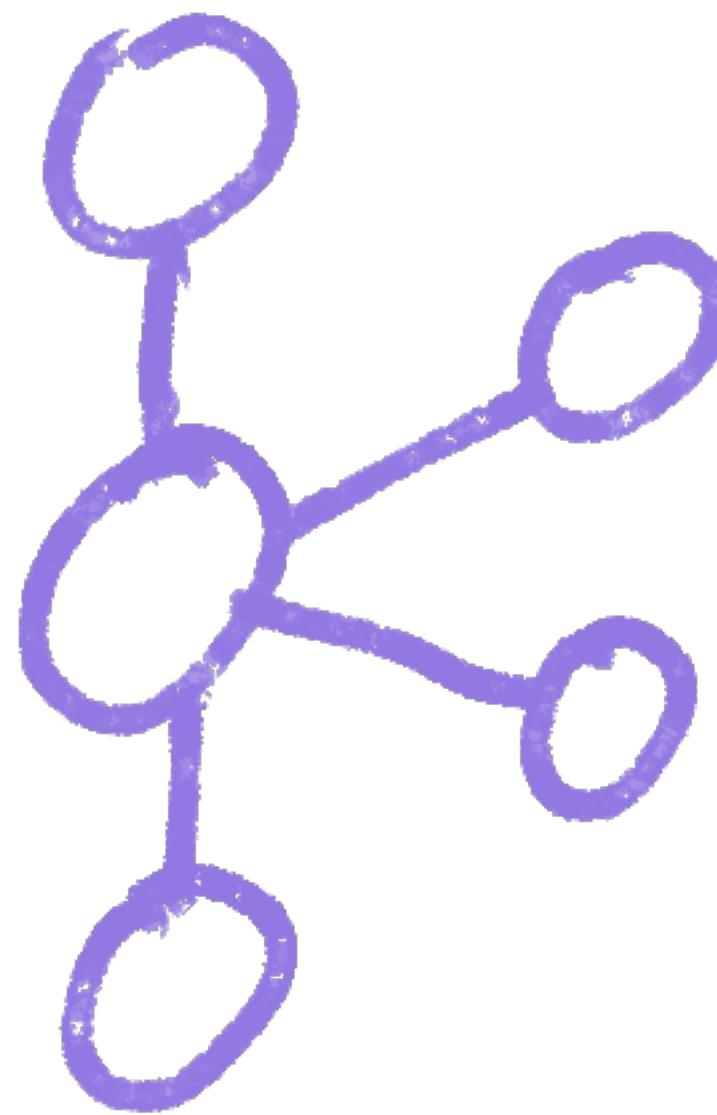
# *Mismatched converters*

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

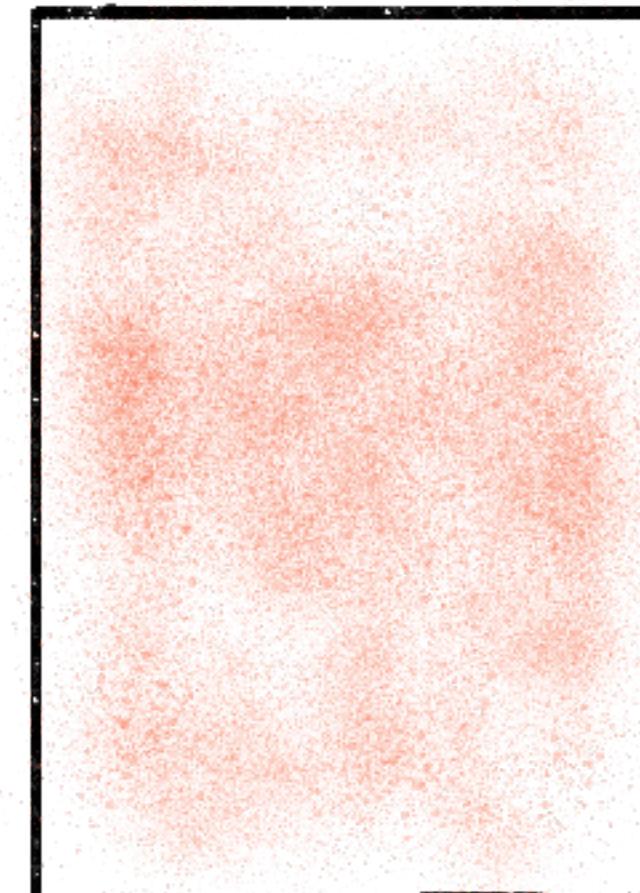
# Mismatched converters

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

*Unknown magic byte!*



**"value.converter":  
"AvroConverter"**

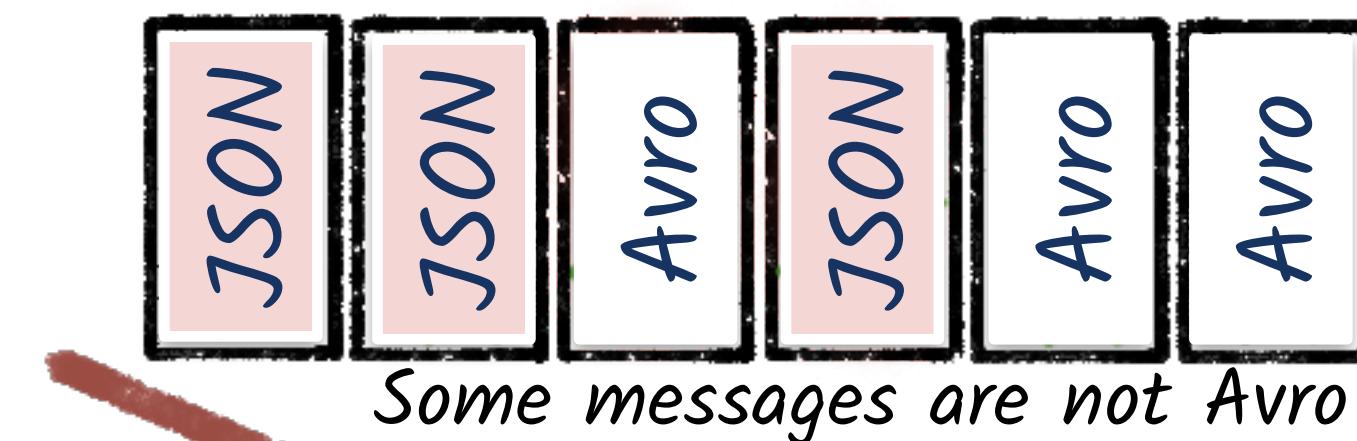
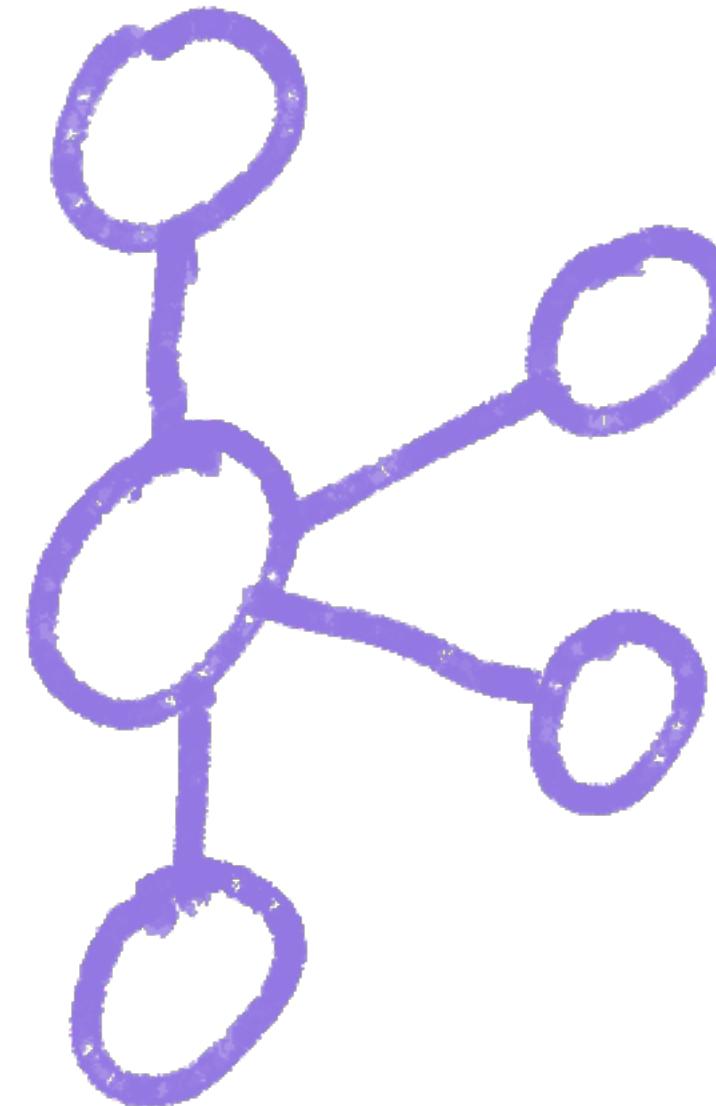


Use the correct Converter for the source data

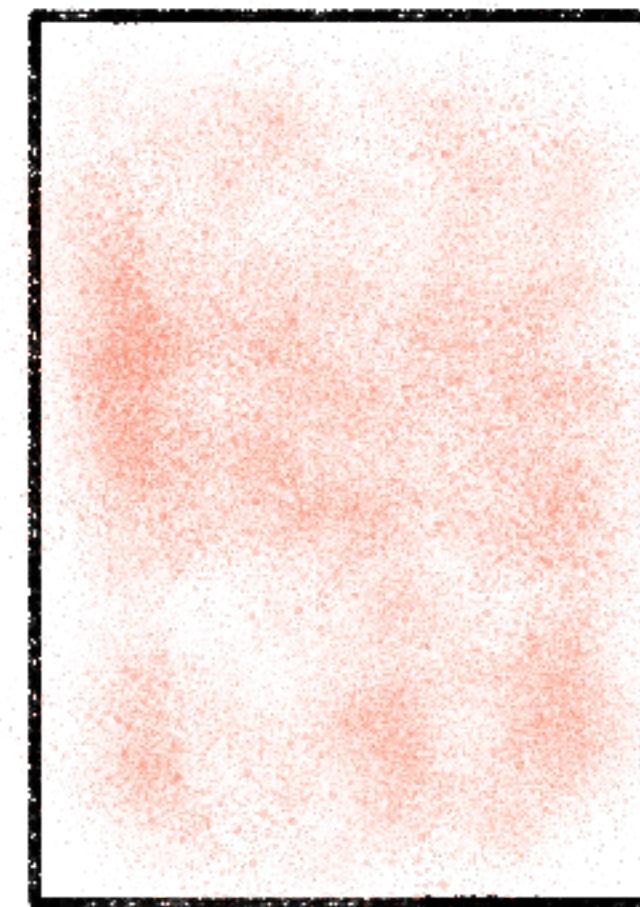
# Mixed serialisation methods

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

*Unknown magic byte!*



`"value.converter":  
"AvroConverter"`



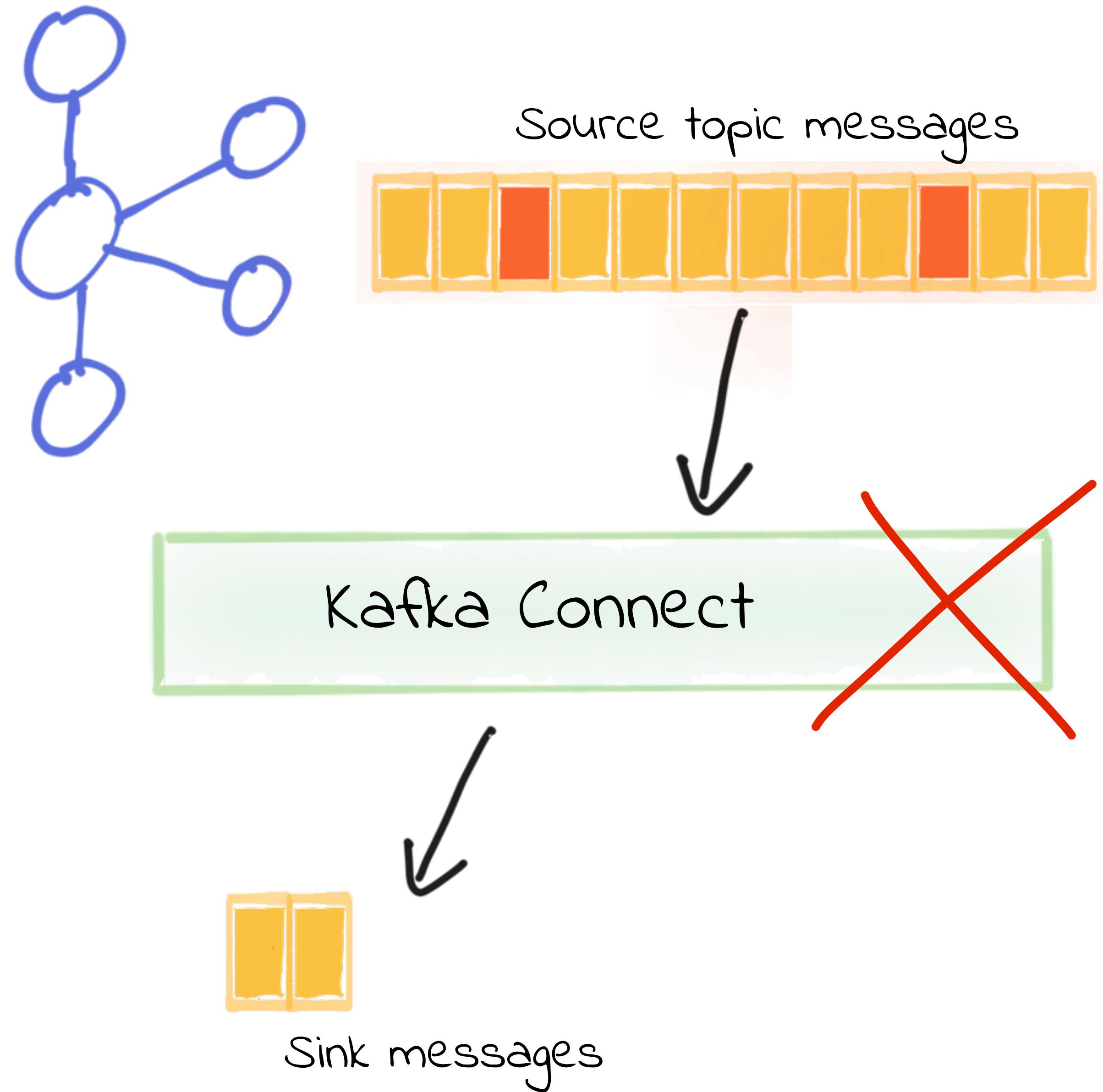
 Use error handling to deal  
with bad messages

# Error Handling and DLQ



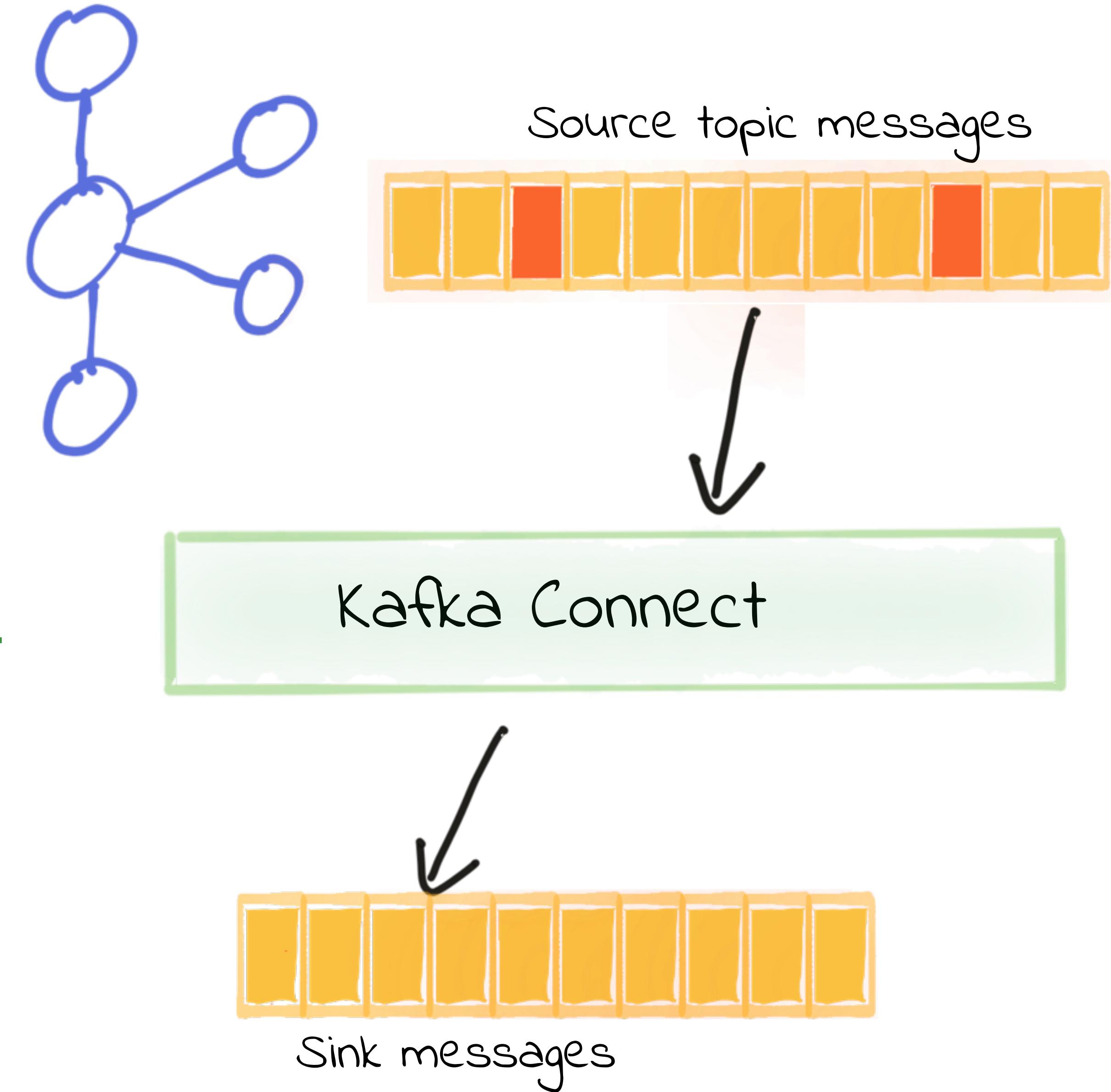
<https://cnfl.io/connect-dlq>

# Fail Fast



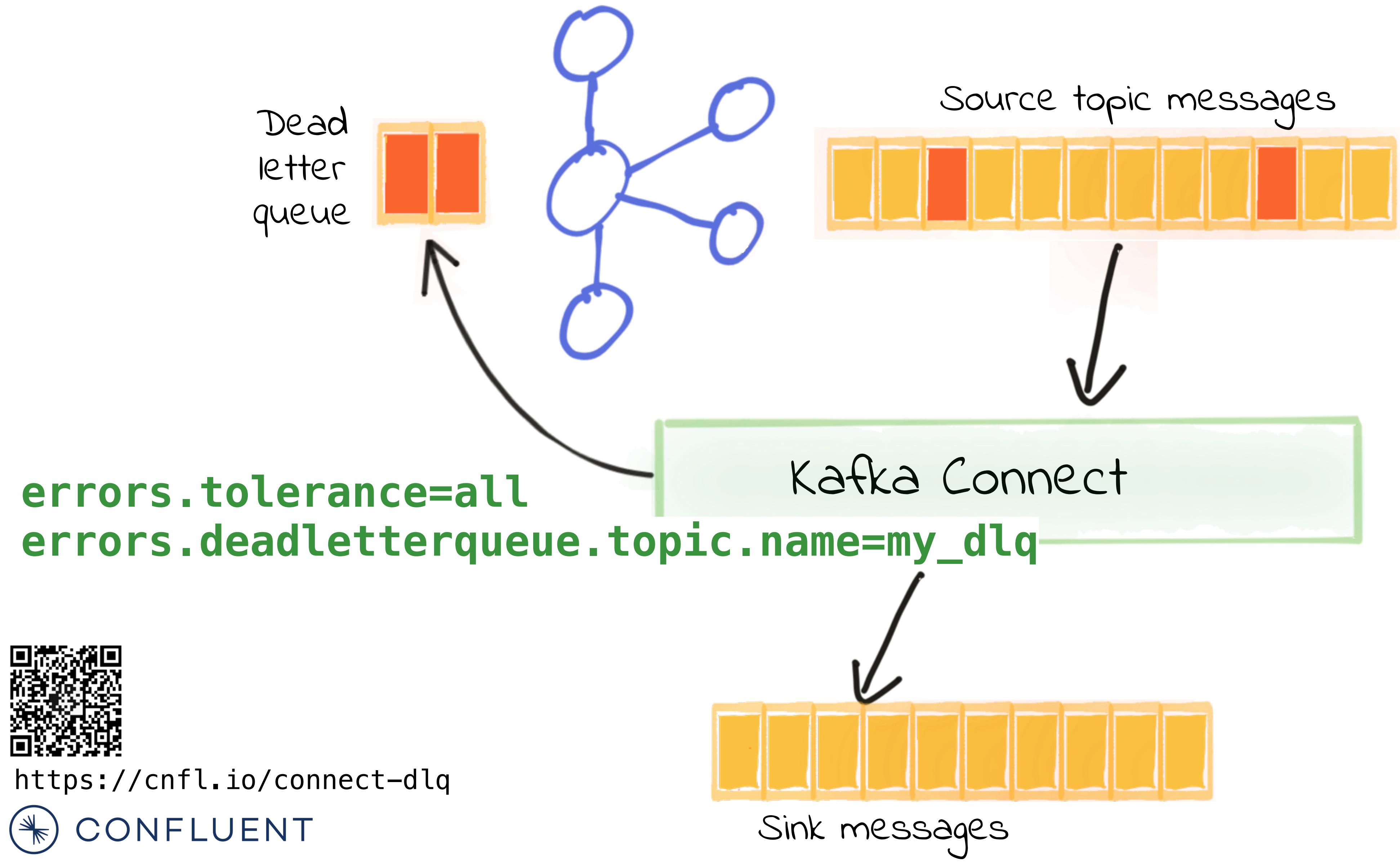
<https://cnfl.io/connect-dlq>

YOLO ^\\_(ツ)\_/^-



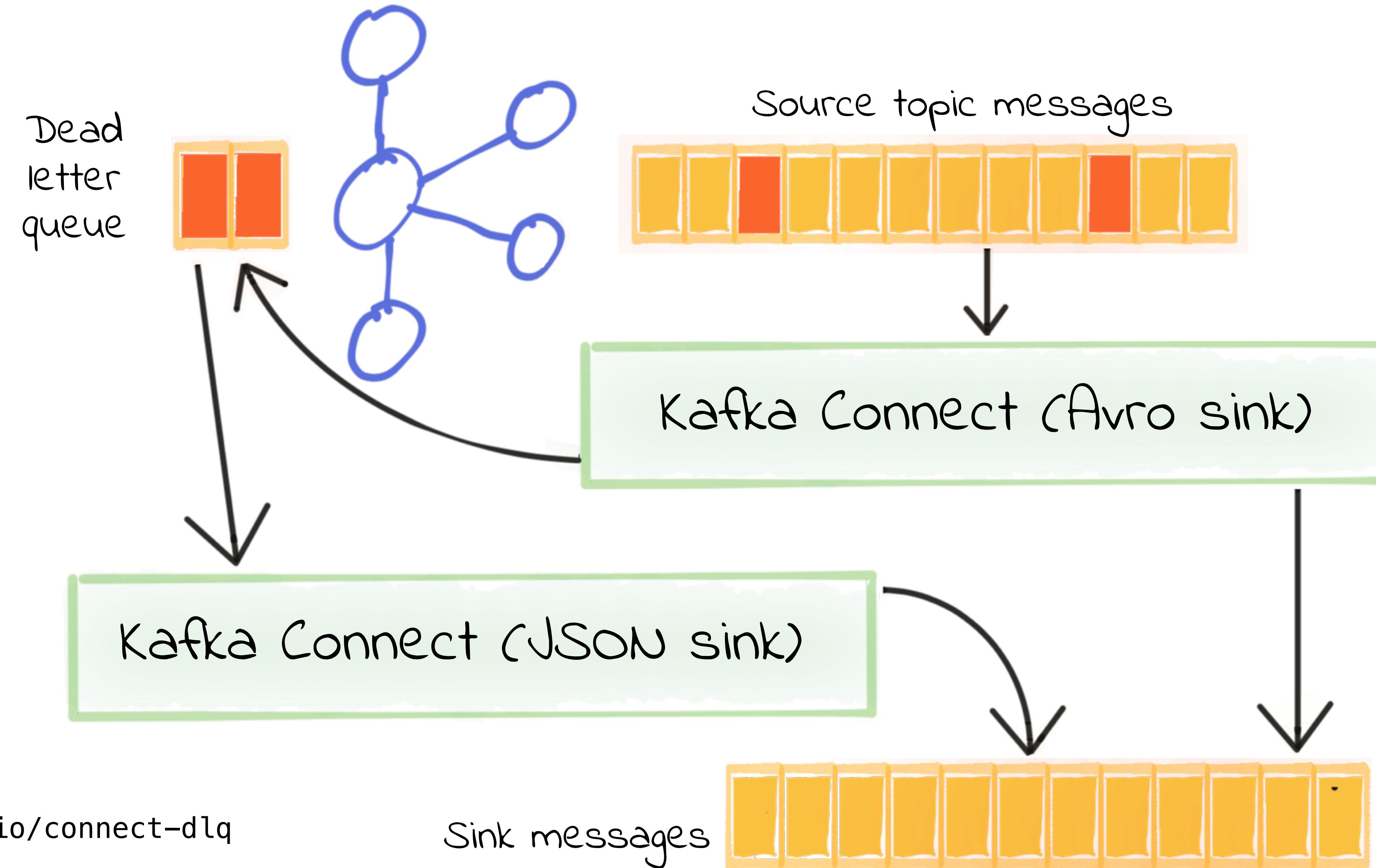
<https://cnfl.io/connect-dlq>

# Dead Letter Queue



<https://cnfl.io/connect-dlq>

# Re-processing the Dead Letter Queue



<https://cnfl.io/connect-dlq>

# Metrics and Monitoring



# REST API

```
[~] Robin@asgard02-2.local
$ curl -s "http://localhost:8083/connectors" | \
jq '.[]' | \
xargs -I{connector_name} curl -s "http://localhost:8083/connectors/{connector_name}/status" | \
jq -c -M '[.name,.connector.state,.tasks[].state]|join(":|:")' | \
column -s : -t | \
sed 's/\"//g' | \
sort

sink-elastic-orders-00      | RUNNING    | RUNNING
sink-elastic-orders-01      | RUNNING    | RUNNING
source-debezium-orders-00  | RUNNING    | RUNNING
[~] Robin@asgard02-2.local
$
```



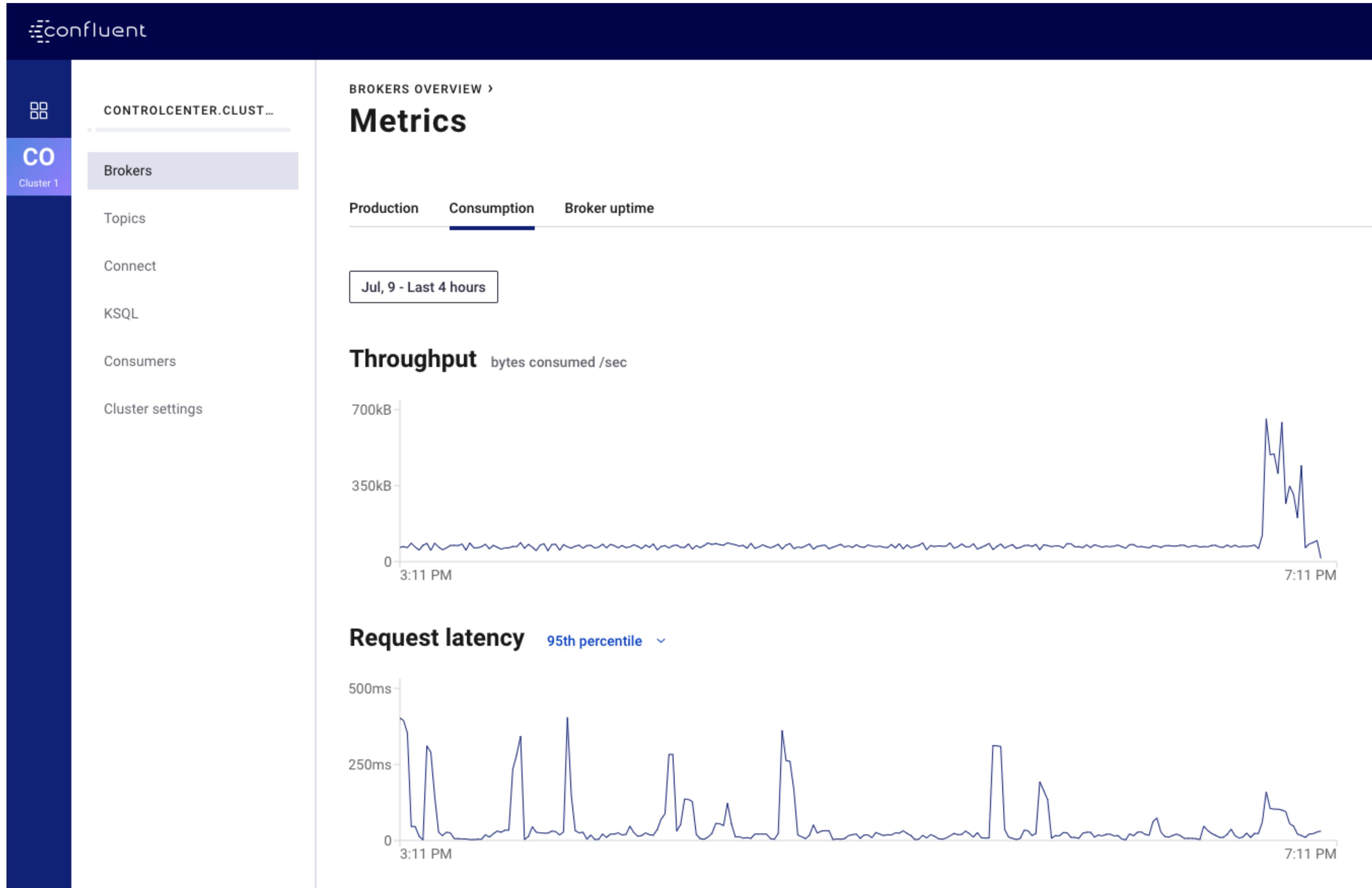
<http://go.rmoff.net/connector-status>



CONFLUENT

@rmoff

# Confluent Control Center

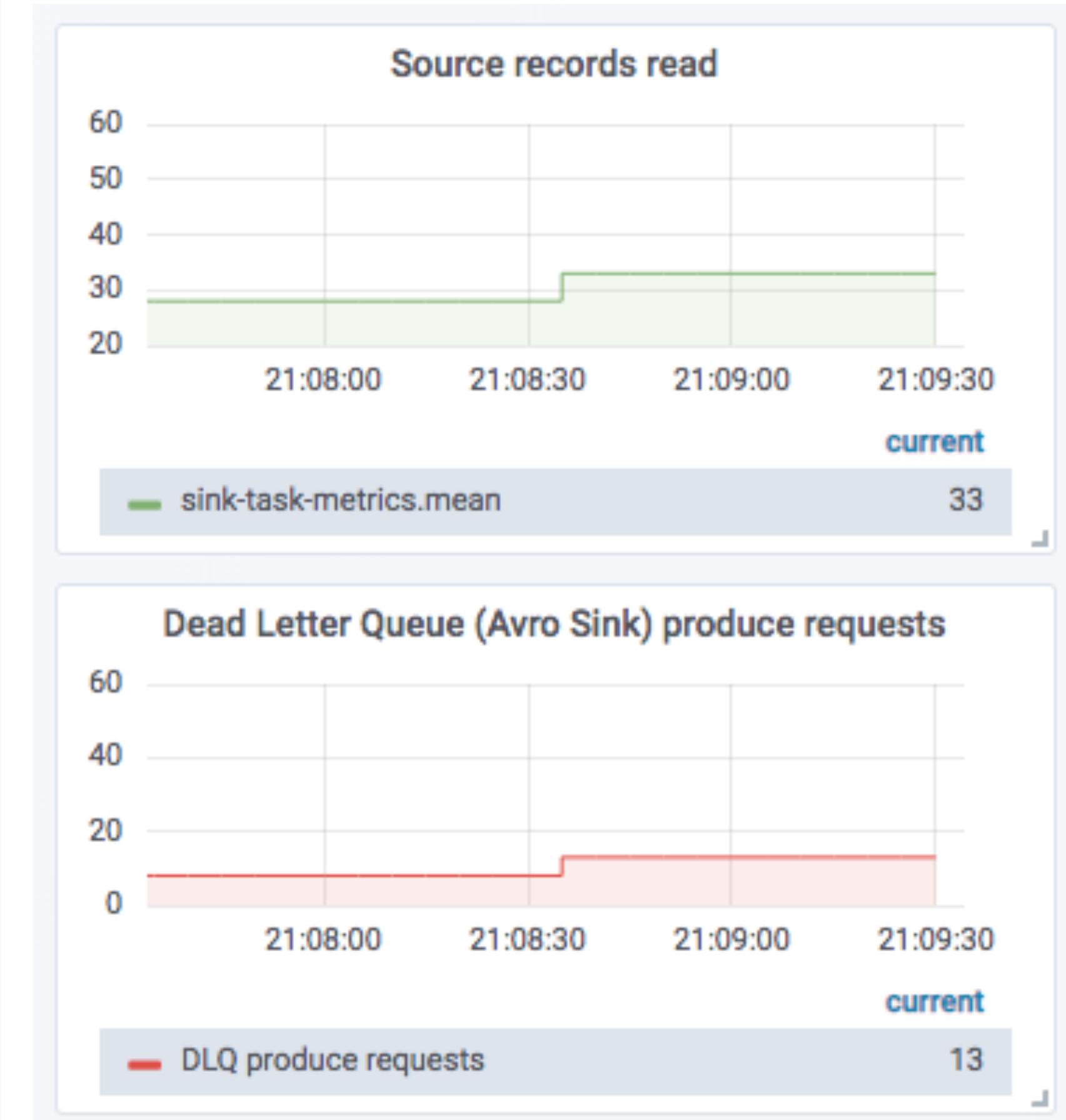
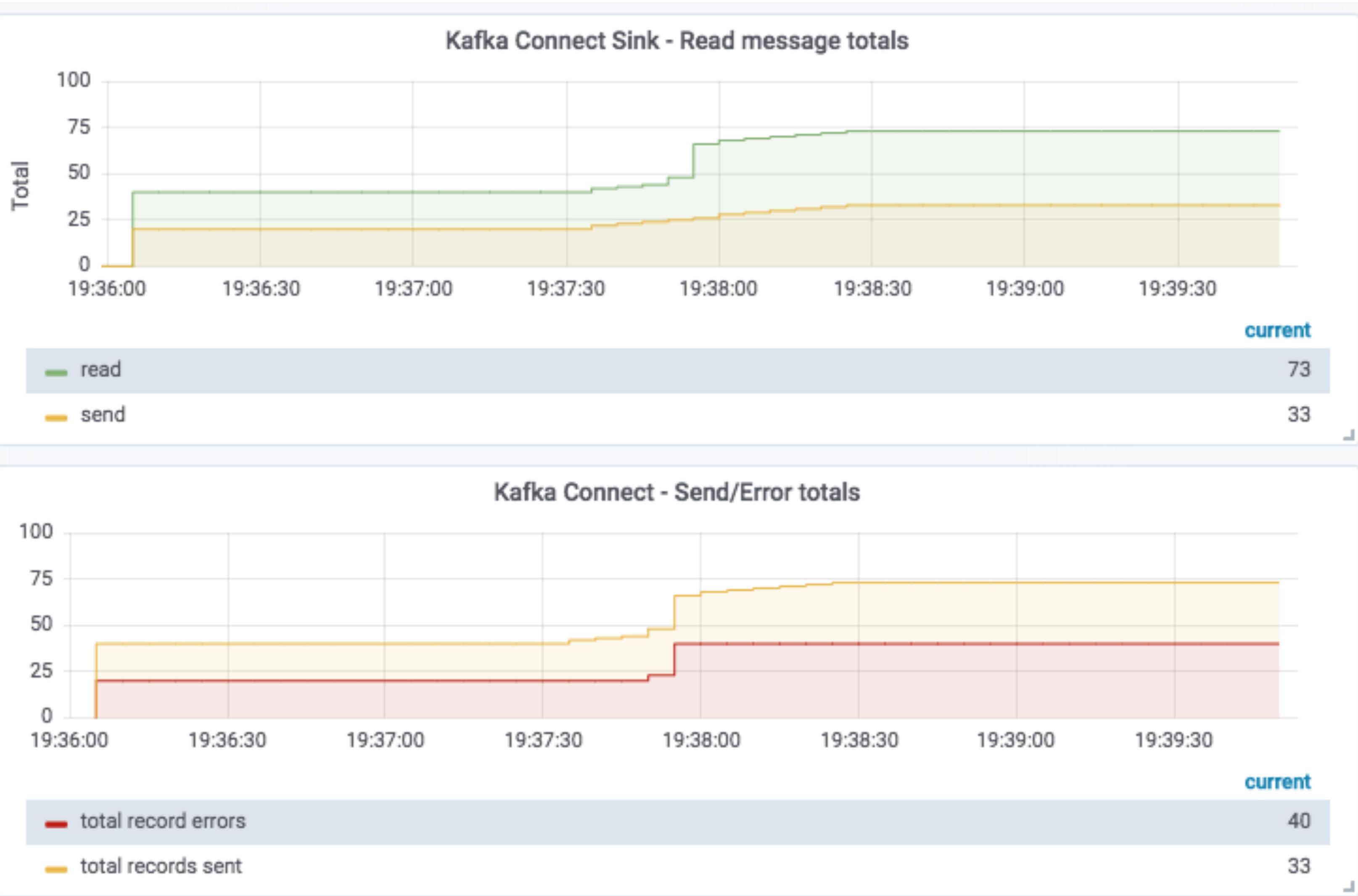


# Consumer lag

The screenshot shows the Confluent Cloud interface for monitoring consumer lag. The left sidebar has sections for MONITORING (System health, Data streams, Consumer lag), MANAGEMENT (Kafka Connect, Clusters, Topics), and ALERTS (Overview, Integration). The 'Consumer lag' section is selected and highlighted in purple. The main content area is titled 'MONITORING > Consumer lag'. A search bar contains the text 'connect-sink'. Below it, a table lists consumer groups and their metrics.

Name	Messages behind	Consumers	Topics
connect-sink-jdbc-TRAIN_MOVEMENTS_ACTIVATIONS_SCHEDULE_00	23	1	1
connect-sink-file-delay-alerts	0	1	1
connect-sink-elasticsearch-movements_01-v01	23	1	1
connect-sink-elasticsearch-movements_activations_schedule_00-v01	23	1	1
connect-sink-s3-tiploc-v00	21283	1	1
connect-sink-s3-train-movements-v00	4613	1	1





Want to  
learn  
more?

CTAs, not CATs  
(sorry, not sorry)



60DEVADV

Free money!  
(additional \$60 towards  
your bill 😊)

\$200 USD off your bill each calendar month  
for the first three months when you sign up

<https://rmoff.dev/ccloud>



confluent cloud

# Fully Managed Kafka as a Service

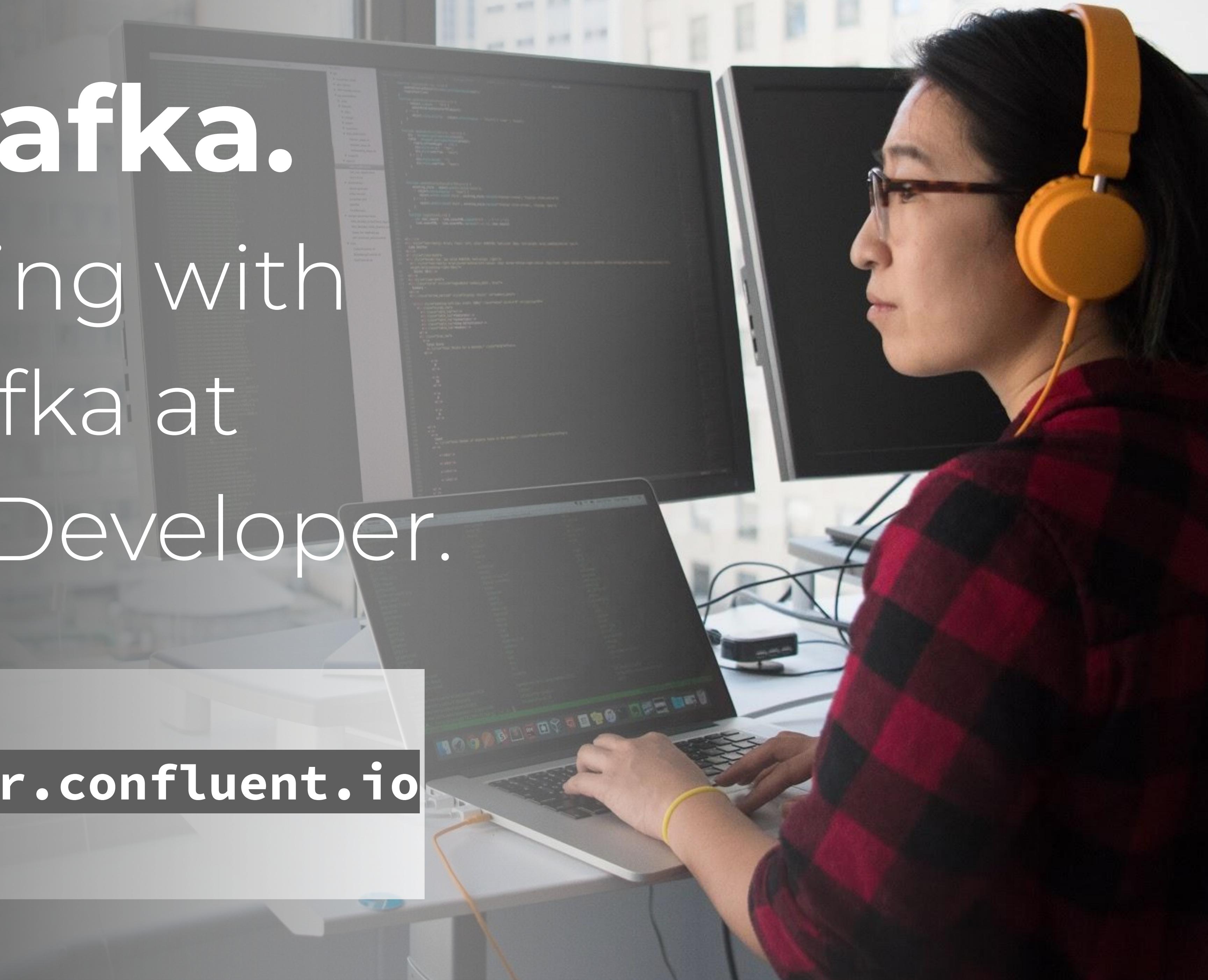


# Learn Kafka.

Start building with  
Apache Kafka at  
Confluent Developer.



[developer.confluent.io](https://developer.confluent.io)



# *Confluent Community Slack group*



**[cnfl.io/slack](https://cnfl.io/slack)**

# Further reading / watching

<https://rmoff.dev/kafka-talks>

## Real-life examples

Here's a nice example using real data to solve a real problem - is my train late now? What are the routes most likely to be delayed?

🚂 [On Track with Apache Kafka: Building a Streaming Platform solution with Rail Data](#)

Moving from 🚂 to 🚗 let's take another real data feed and build some realtime location-based notifications 🚶

🤖 [Building a Telegram bot with Go, Apache Kafka, and ksqlDB](#)

## Integration and data pipelines

Integration between Kafka and other systems? Kafka Connect has you covered ⚡

⌚ [Kafka Connect in 60 seconds](#)

🦸 [From Zero to Hero with Kafka Connect](#)



<https://talks.rmoff.net>

@rmoff