

# BUILDING REALTIME DATA PIPELINES WITH KAFKA CONNECT AND SPARK STREAMING

Ewen Cheslack-Postava  
Confluent



## ABOUT ME

- Engineer @ Confluent
- Kafka Committer
- Kafka Connect Lead





ORACLE  
DATABASE

Spark

OMNITURE

Couchbase

logstash

SPARK  
SUMMIT  
EAST 2017

monetdb



riak



mongoDB

ArangoDB



OrientDB



cassandra

accumulo

neo4j



mongoDB

ArangoDB

VOLTDB



elasticsearch

APACHE  
HBASE



AEROSPIKE



Solr



nuodb



memsql

RethinkDB



Scalaris

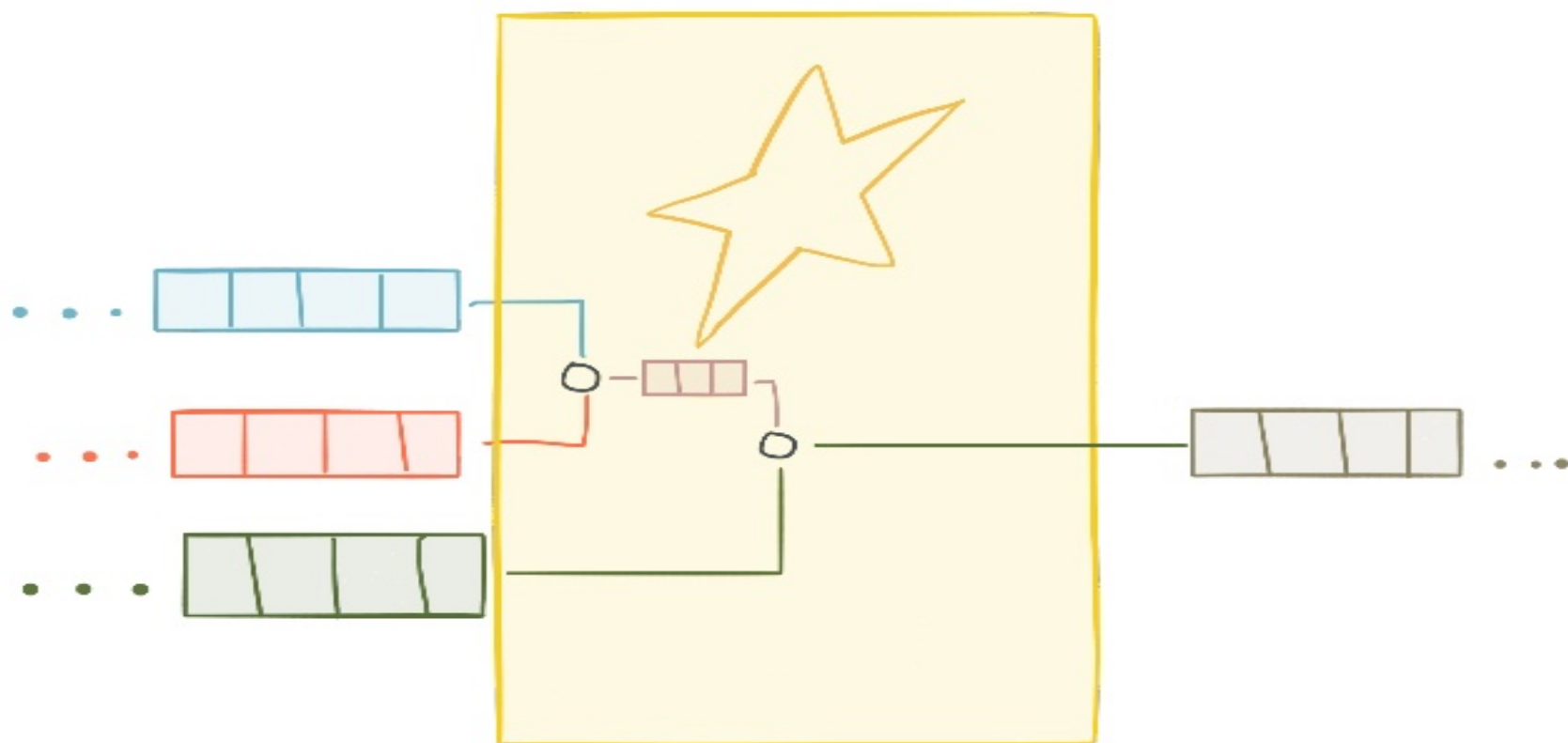


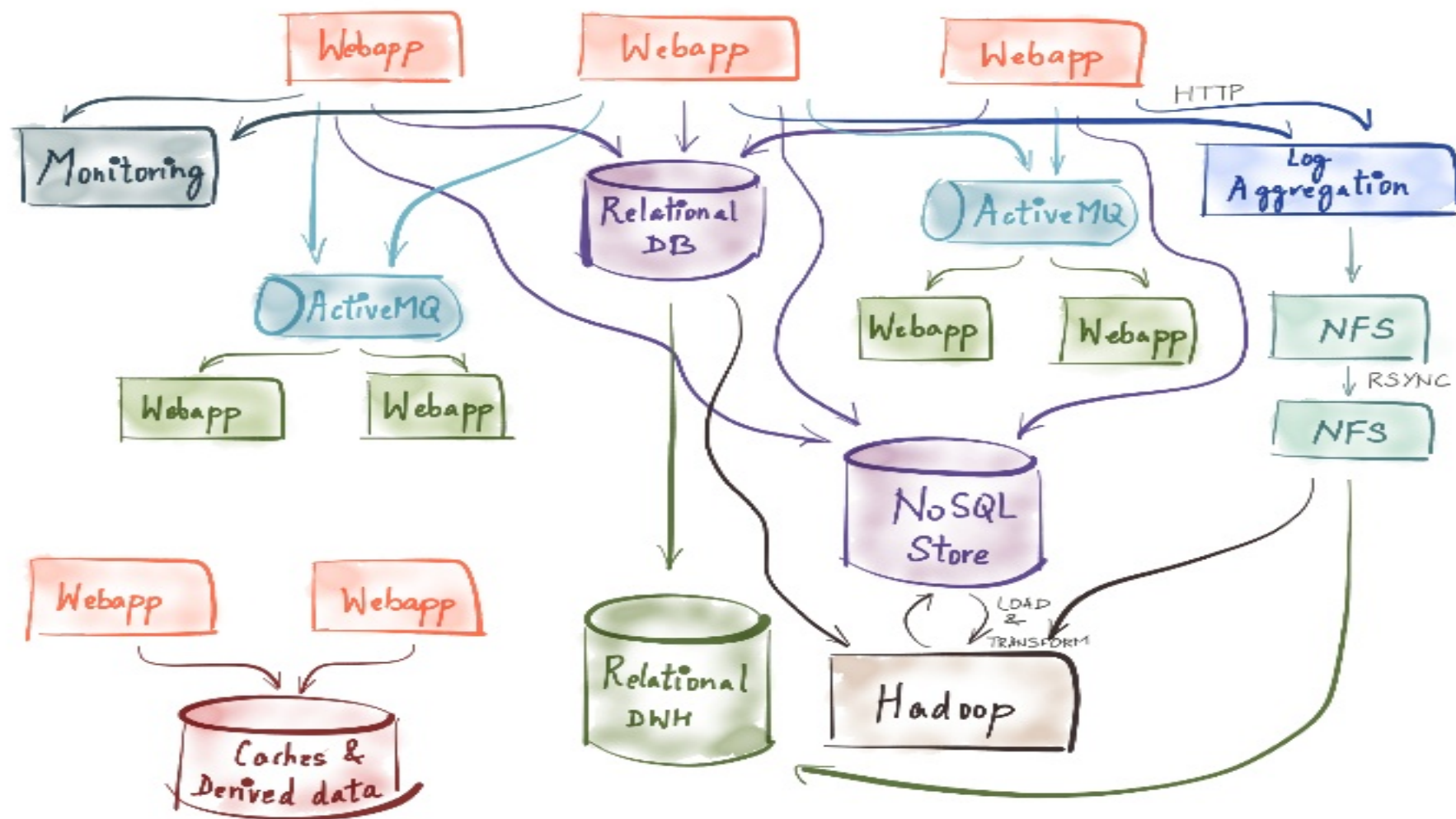
APACHE  
GIRAPH



Flink

# Stream Processing







# ANTI PATTERNS

1. One-off tools
2. Kitchen sink tools
3. Stream processing frameworks

Ad-hoc ← ? → E, T, & L

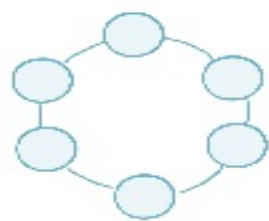
# Introducing Kafka Connect

Large-scale streaming data import/export for Kafka



# GOALS

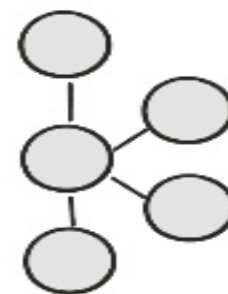
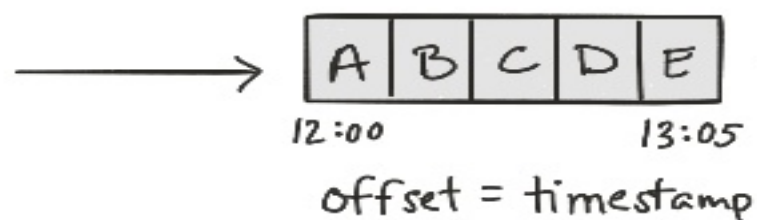
1. Focus on copying
2. Batteries included
3. Standardize
4. Parallelism
5. Scale



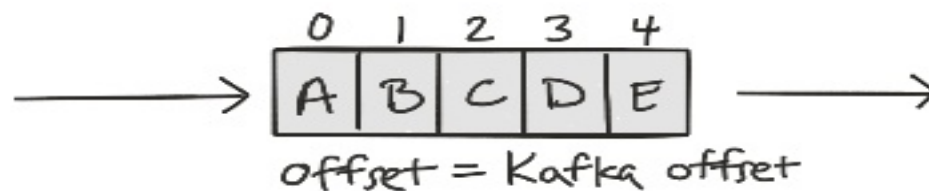
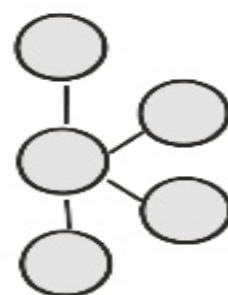
Table

TS	Data
12:00	A
12:20	B
12:30	C
13:00	D
13:05	E

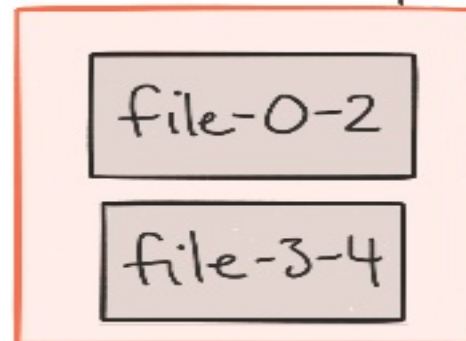
Database Source



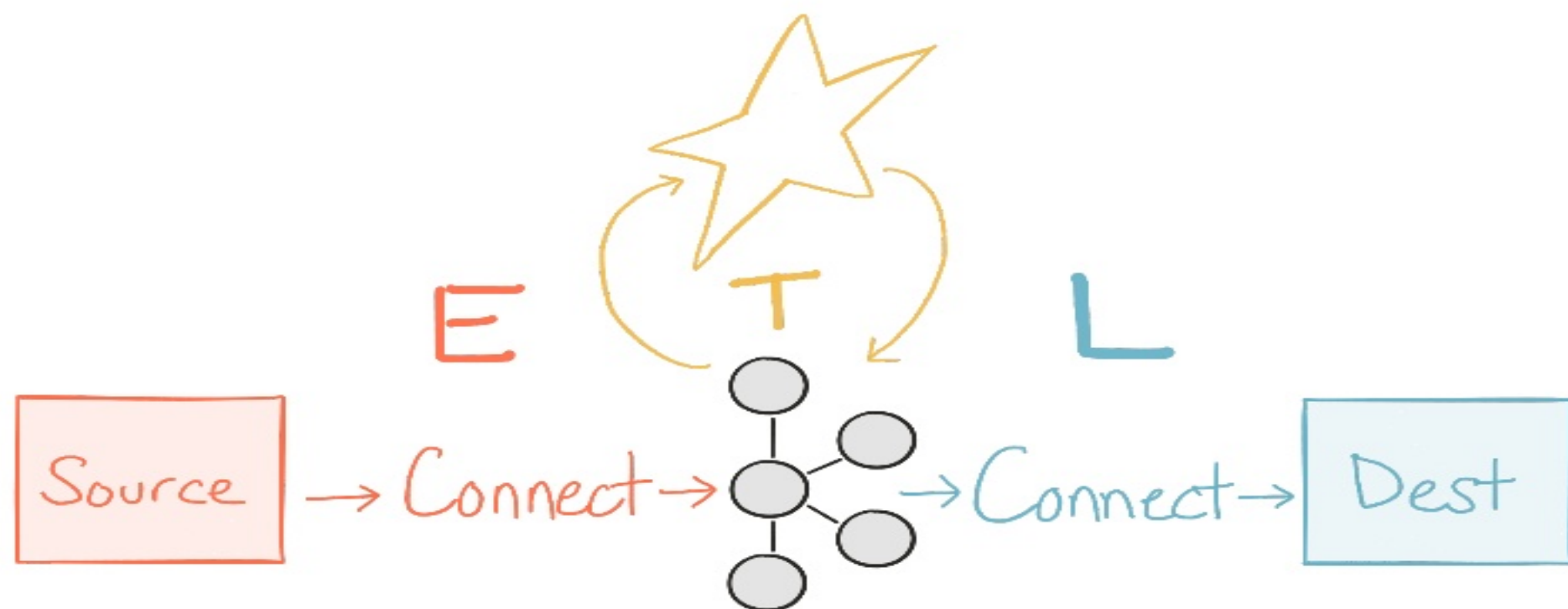
HDFS Sink



HDFS Directory



# Separation of Concerns

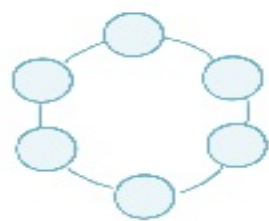


# SPARK STREAMING & KAFKA

INPUT: DIRECT KAFKA STREAMS

- OLD CONSUMER (0.8.2.1+) - SPARK 2.0+
- NEW CONSUMER (0.10.0.0+) - EXPERIMENTAL

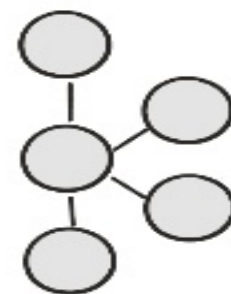
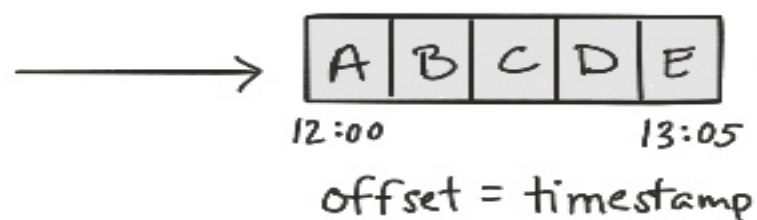
OUTPUT: SPARK KAFKA WRITER



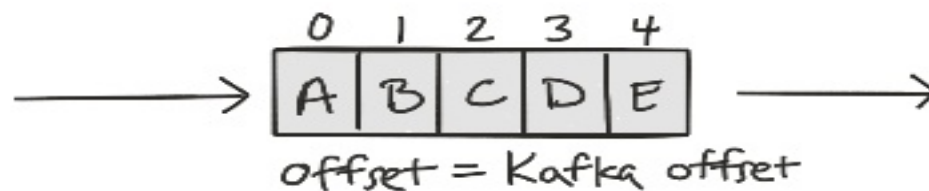
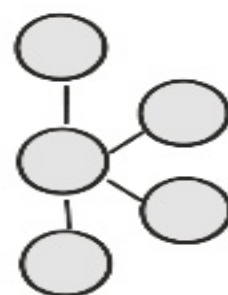
Table

TS	Data
12:00	A
12:20	B
12:30	C
13:00	D
13:05	E

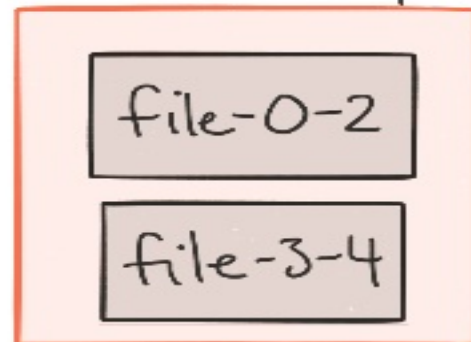
Database Source



HDFS Sink



HDFS Directory

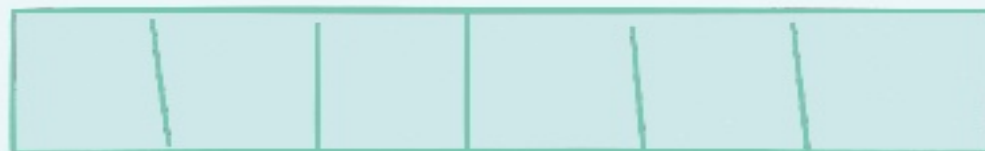


## Partitioned Stream

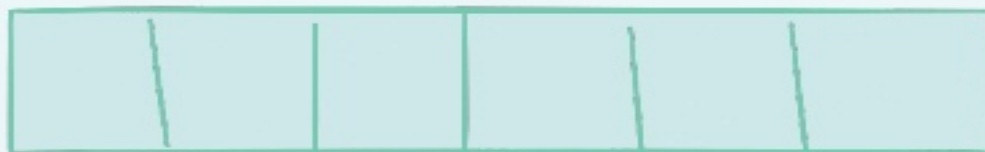
Partition 1



Partition 2



Partition 3



Partition 4



$O_1$   $O_2$   $O_3$   $O_4$   $O_5$   $O_6$

← offsets



## Partitioned Stream - Database

Table 1

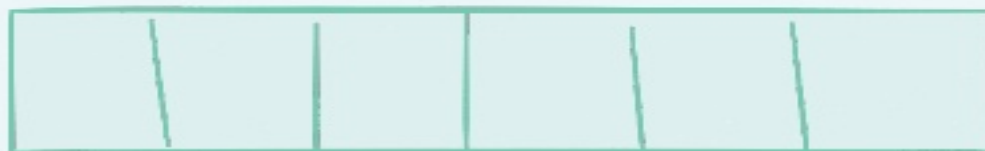


Table 2

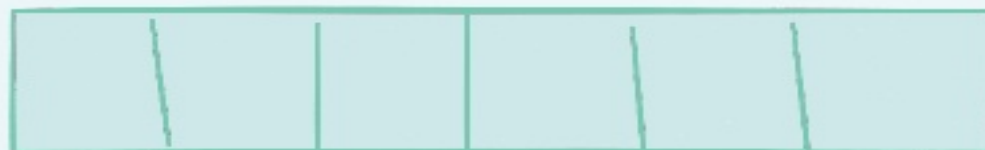


Table 3

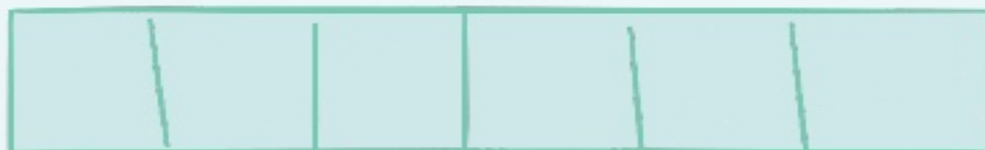
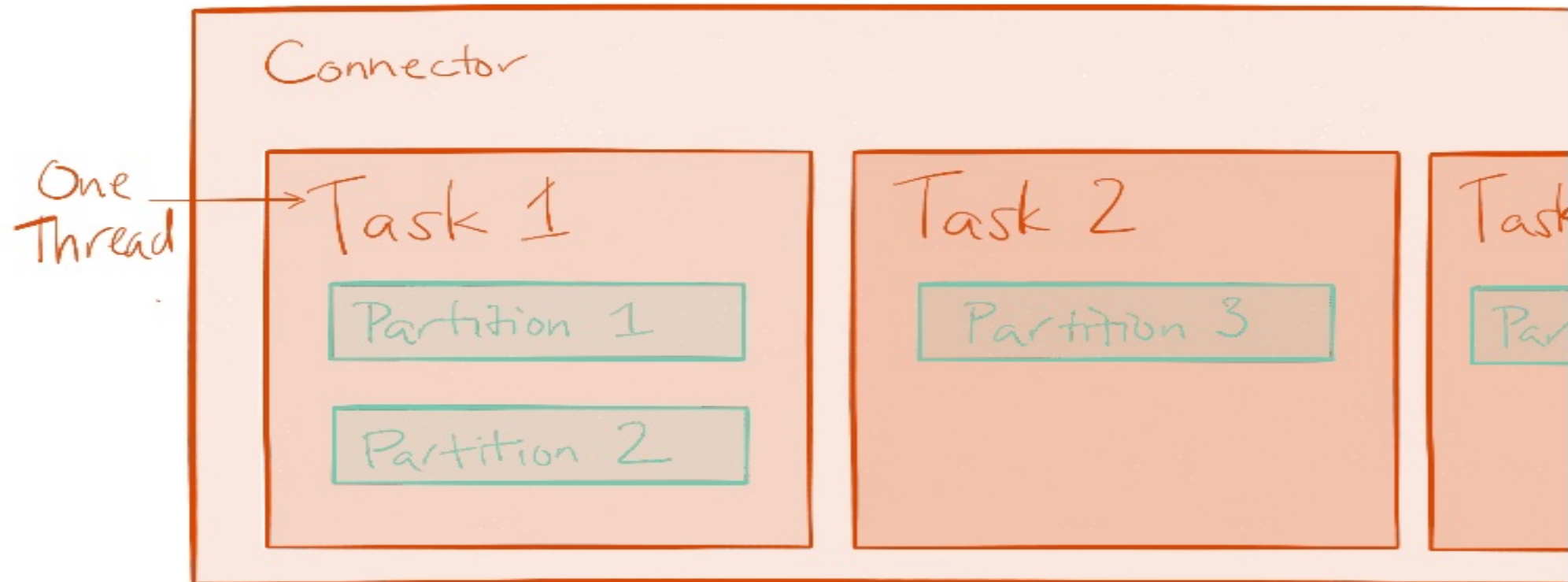
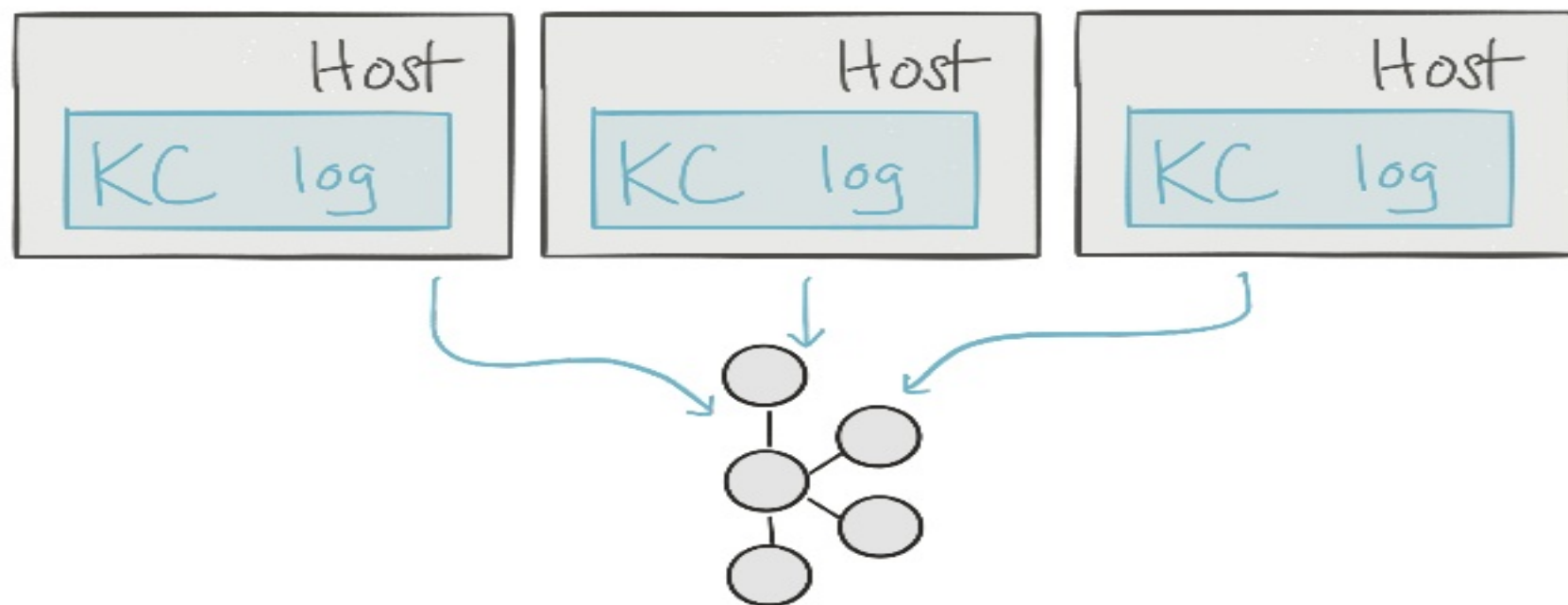


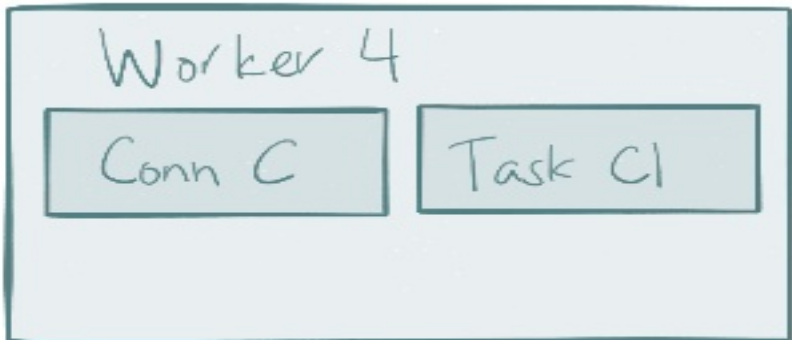
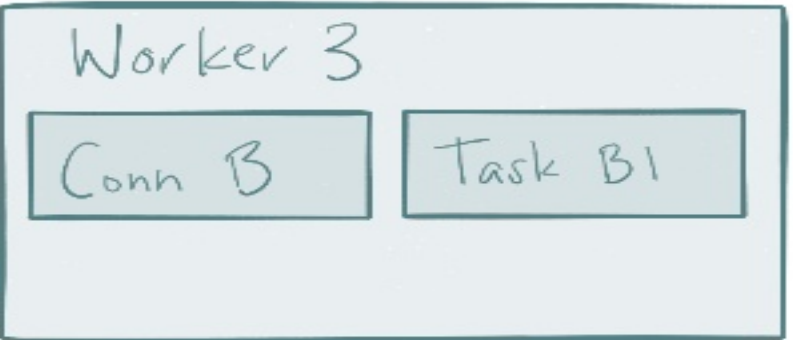
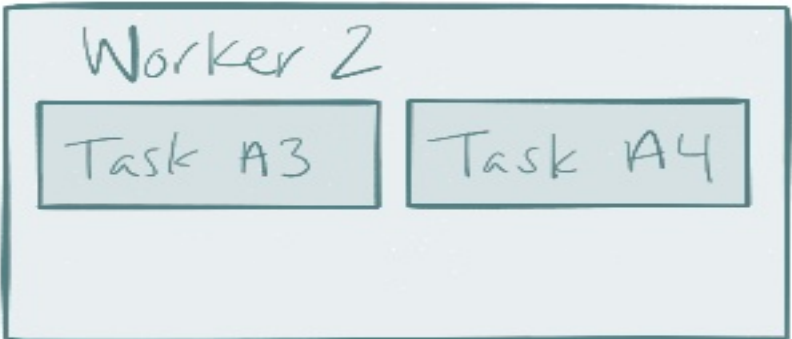
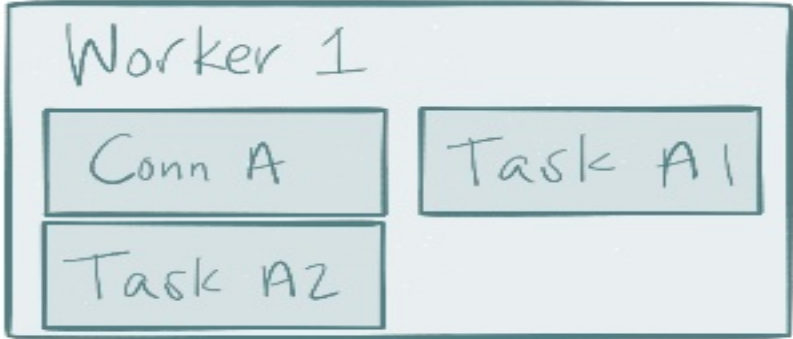
Table 4

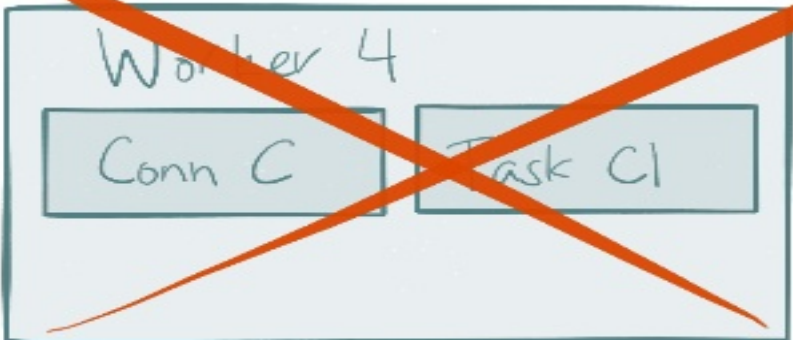
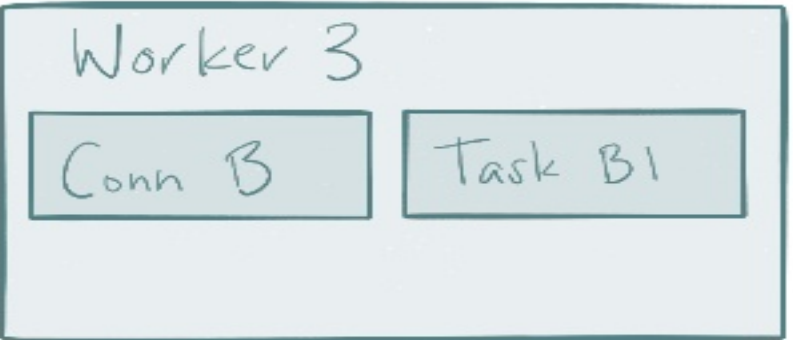
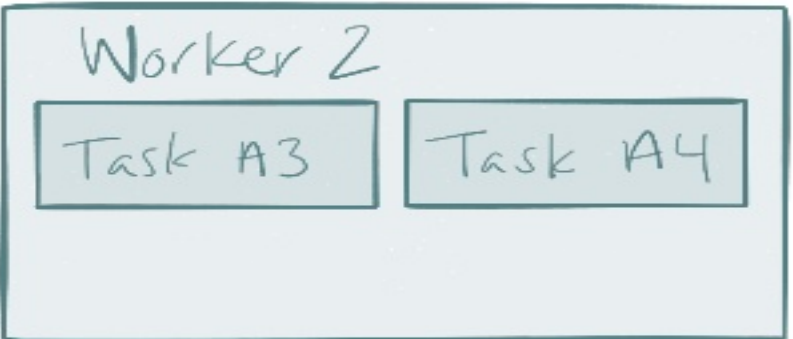
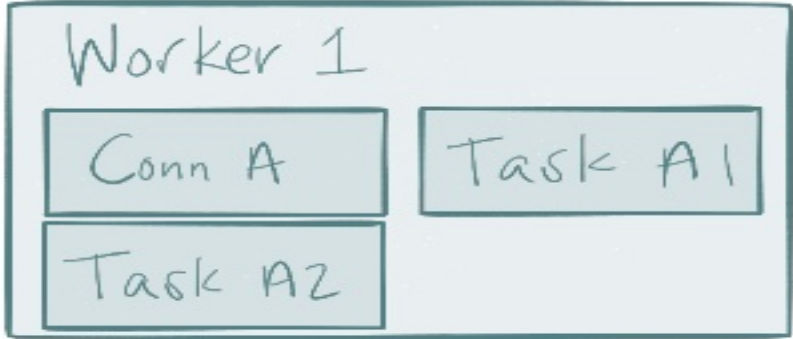


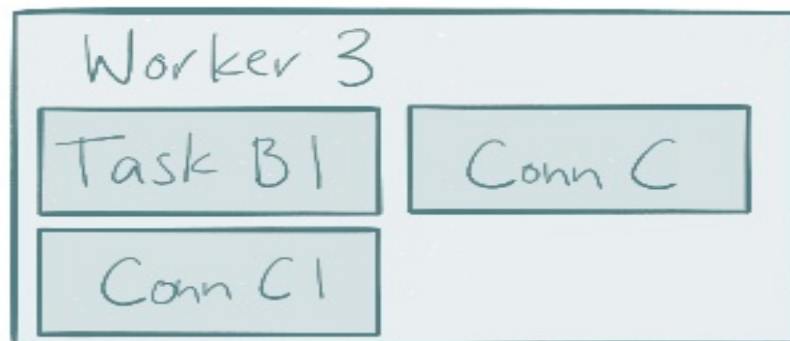
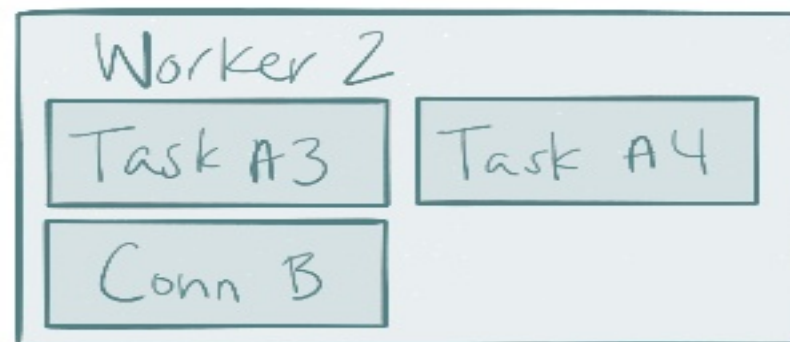
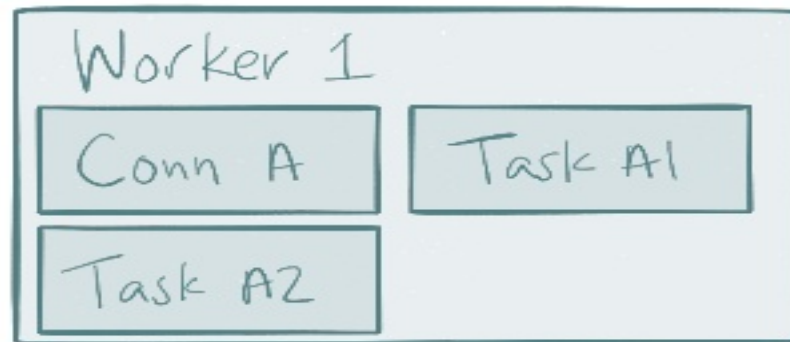
id=1 id=2 id=3 id=4 id=5 id=6





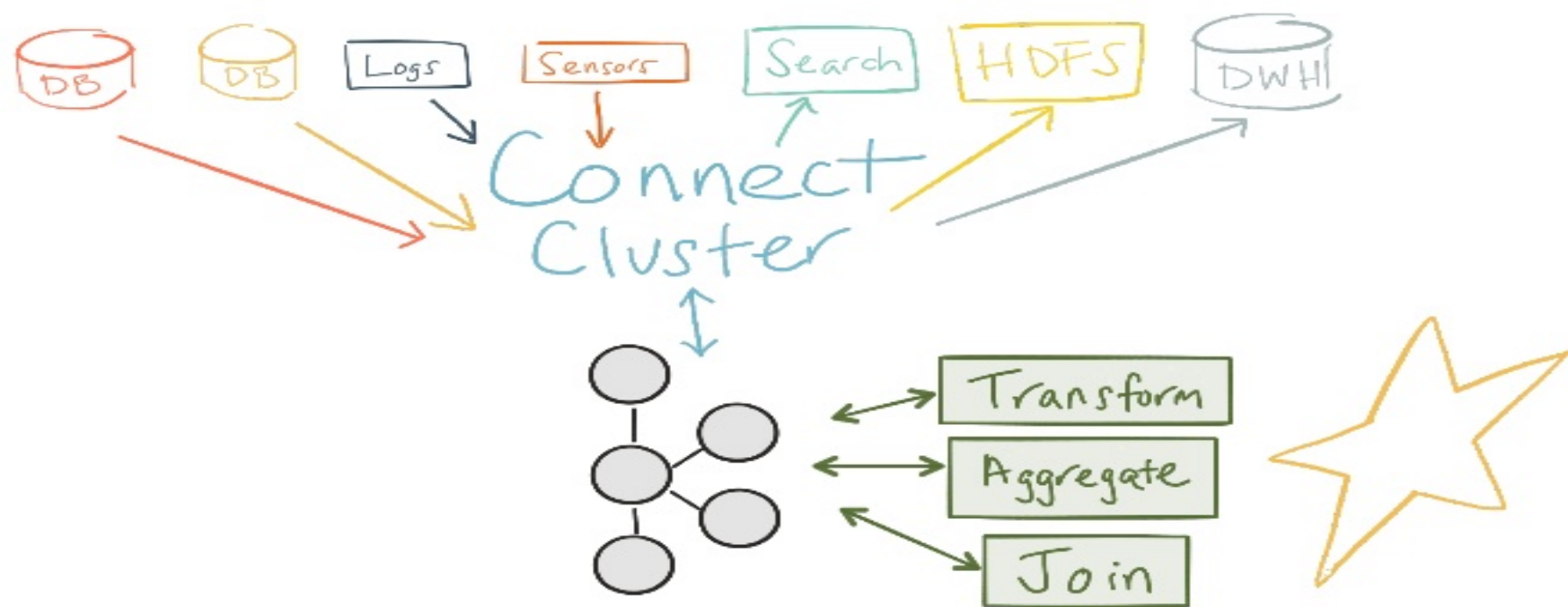








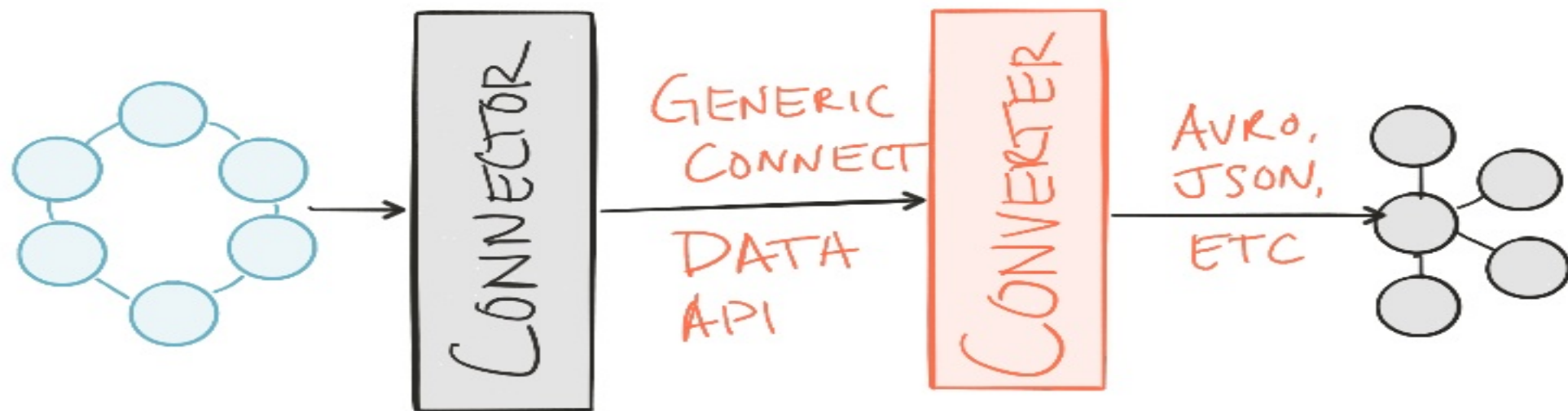
# Data Integration as a Service



# DELIVERY GUARANTEES

- FRAMEWORK MANAGED OFFSETS
- AT LEAST ONCE DEFAULT
- EXACTLY ONCE (w/CONNECTOR SUPPORT)

# CONVERTERS



# SPARK & KAFKA CONNECT

MORE SPARK STREAMING SOURCES + SINKS

REDUCE ADOPTION FRICTION

LEVERAGE KAFKA AS DE FACTO STREAMING  
STORAGE ENGINE

COMING SOON...

IMPROVED KAFKA CLIENT COMPATIBILITY

EXACTLY ONCE SEMANTICS

## Certified Connectors

Certified Connectors have been developed by vendors and/or Confluent utilizing the Kafka Connect framework. These Connectors have met criteria for code development best practices, schema registry integration, security, and documentation.

CONNECTOR	TAGS	DEVELOPER	SUPPORT
HDFS (Sink)	HDFS, Hadoop, Hive	Confluent	<a href="#">Confluent</a>
JDBC (Source)	JDBC, MySQL	Confluent	<a href="#">Confluent</a>
Elastic Search (Sink)	search, Elastic, log, analytics	Confluent	<a href="#">Confluent</a>
DataStax (Sink)	Cassandra, DataStax	Data Mountaineer	<a href="#">Data Mountaineer</a>
Attunity (Source)	CDC	Attunity	<a href="#">Attunity</a>
Couchbase (Source)	Couchbase, NoSQL	Couchbase	<a href="#">Couchbase</a>
GoldenGate (Source)	CDC, Oracle	Oracle	<a href="#">Community</a>
JustOne (Sink)	Postgress	JustOne	<a href="#">JustOne</a>
Striim (Source)	CDC, MS SQLServer	Striim	<a href="#">Striim</a>
Syncsort DMX	DB2, IMS, VSAM, CICS	Syncsort	<a href="#">Syncsort</a>

### Want to build a connector?

Interested Open Source Developers and Vendors can get started with the following resources:

- > [Kafka Connect Overview](#)
- > [Kafka Connect Developers Guide](#)

### Already built a connector?

If you have a connector that you'd like to see added to our list, let us know at:

[confluent-platform@googlegroups.com](mailto:confluent-platform@googlegroups.com)

If your connector is well done, we'll send you a free T-shirt and may even list it here.

### Contact Us

For software vendors seeking to build a connector, our Partner Engineering team is







# THANK YOU

@ewencp

@confluentinc

**<http://confluent.io/download/>**

**<http://connectors.confluent.io>**