

BUILDING REALTIME DATA PIPELINES WITH KAFKA CONNECT AND SPARK STREAMING

Guozhang Wang

Confluent



SPARK SUMMIT 2016
DATA SCIENCE AND ENGINEERING AT SCALE

About Me: Guozhang Wang

- Engineer @ Confluent.
- Apache Kafka Committer, PMC Member.
- Before: Engineer @ LinkedIn, Kafka and Samza.



**What do you REALLY need
for Stream Processing?**



Spark Streaming!



Spark Streaming! Is that All?

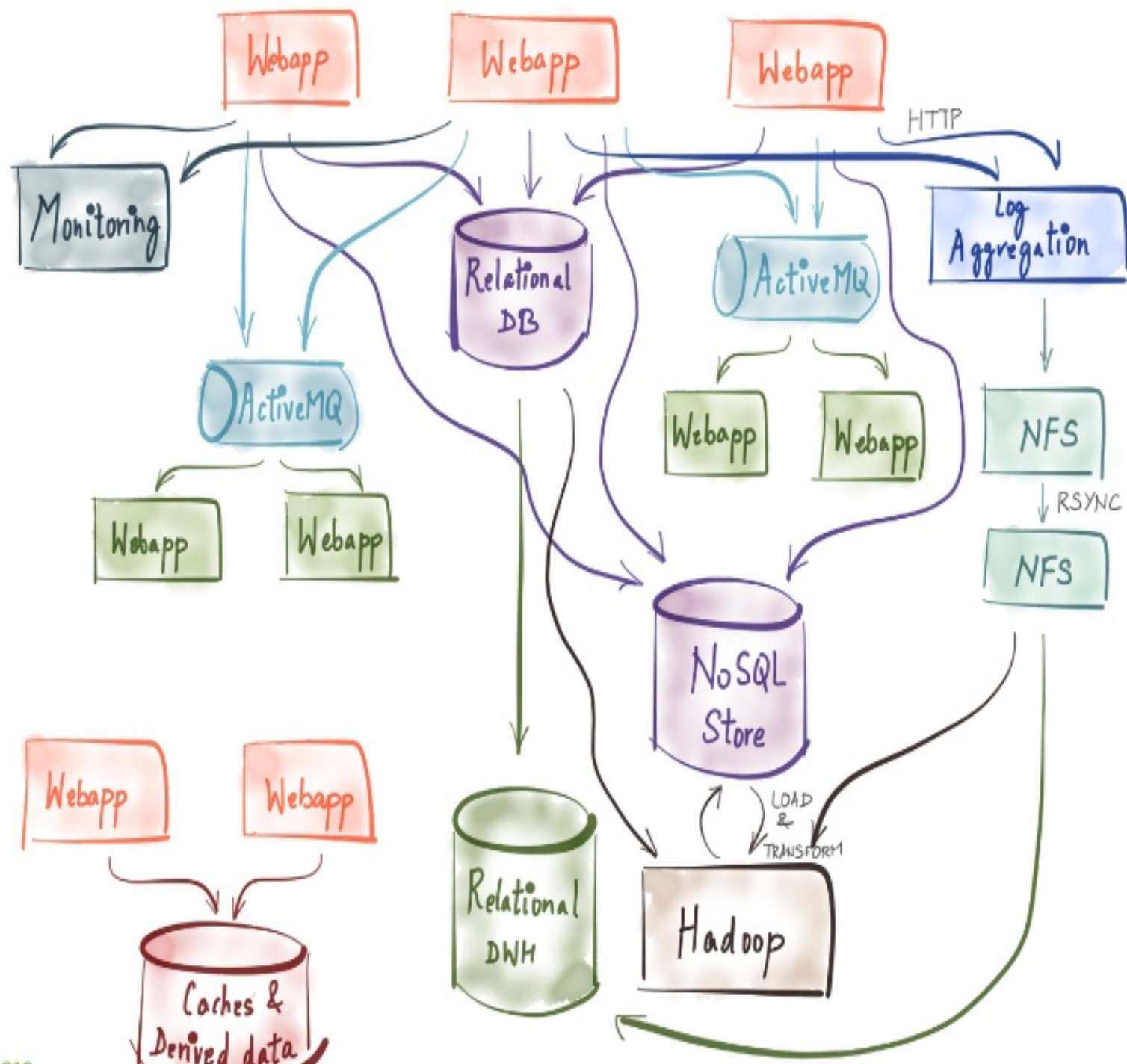


Spark Streaming! Is that All?



Data can Comes from / Goes to..





Real-time Data Integration:

getting data to all the right places



Option #1: One-off Tools

- Tools for each specific data systems
- Examples:
 - jdbcRDD, Cassandra-Spark connector, etc..
 - Sqoop, logstash to Kafka, etc..



Option #2: Kitchen Sink Tools

- Generic point-to-point data copy / ETL tools
- Examples:
 - Enterprise application integration tools



Option #3: Streaming as Copying

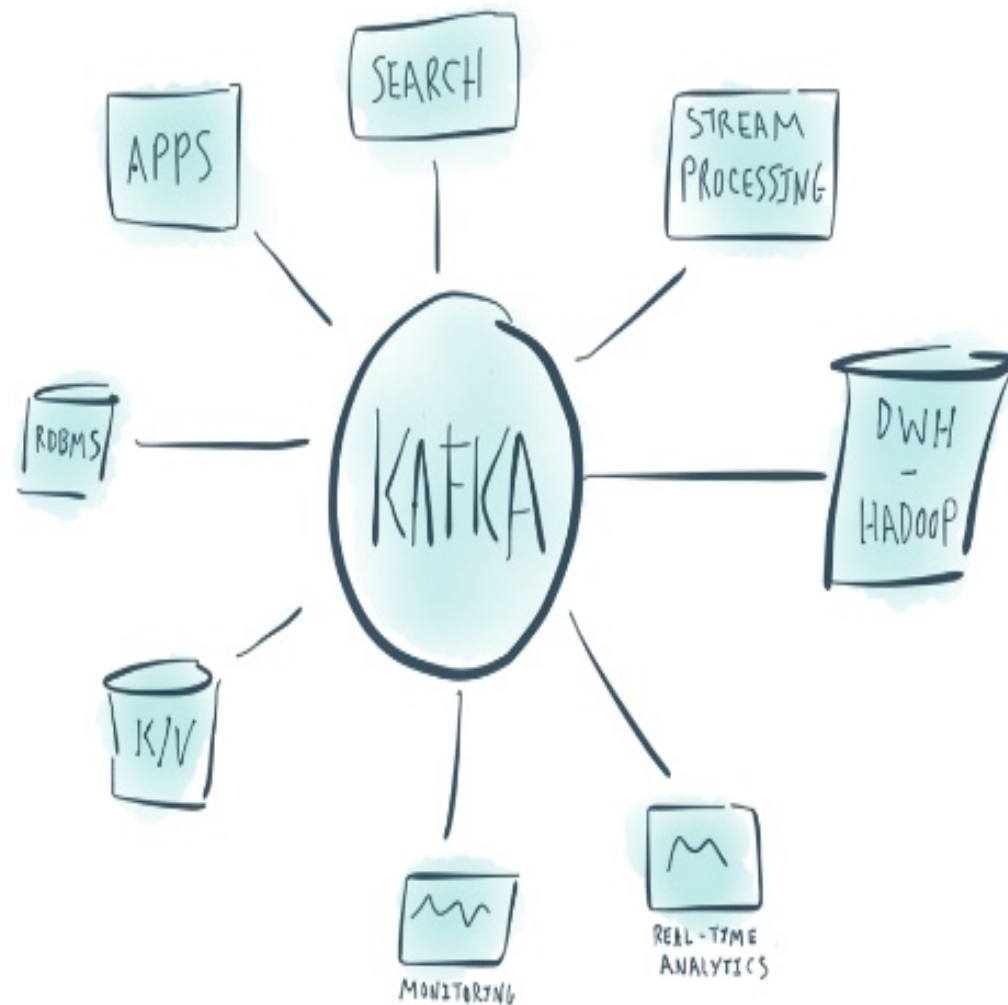
- Use stream processing frameworks to copy data
- Examples:
 - Spark Streaming: `MyRDDWriter (foreachPartition)`
 - Storm, Samza, Flink, etc..



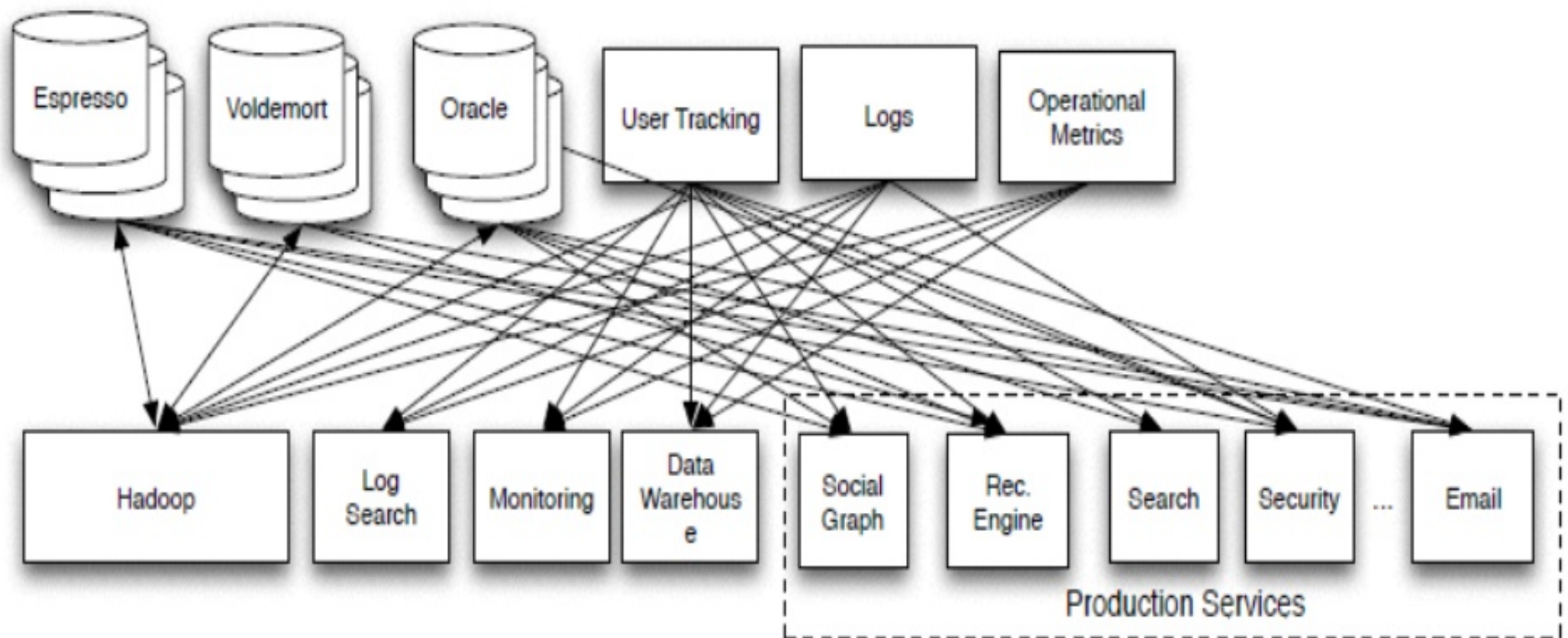
Real-time Integration: E, T & L



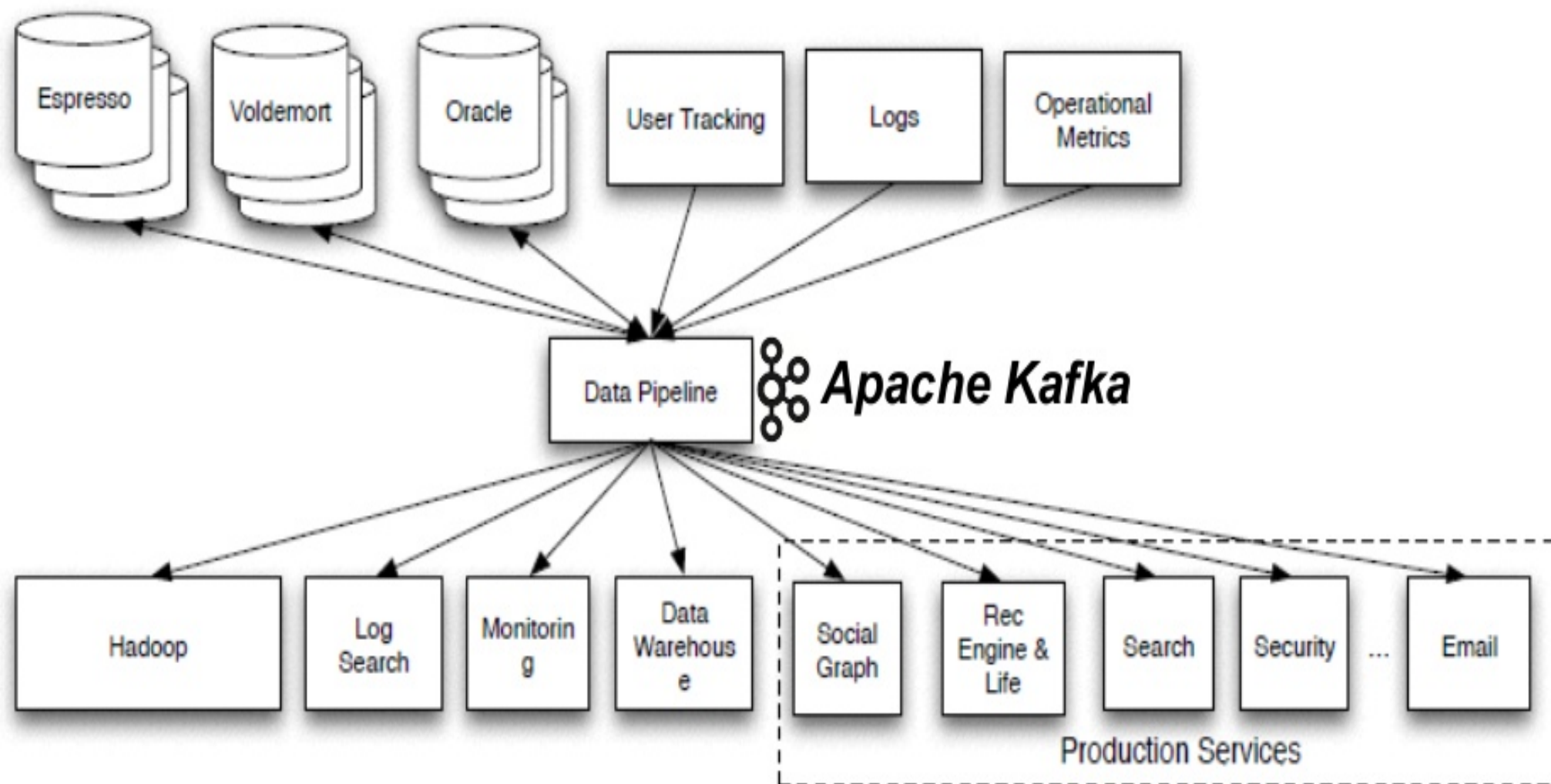
STREAMING PLATFORM



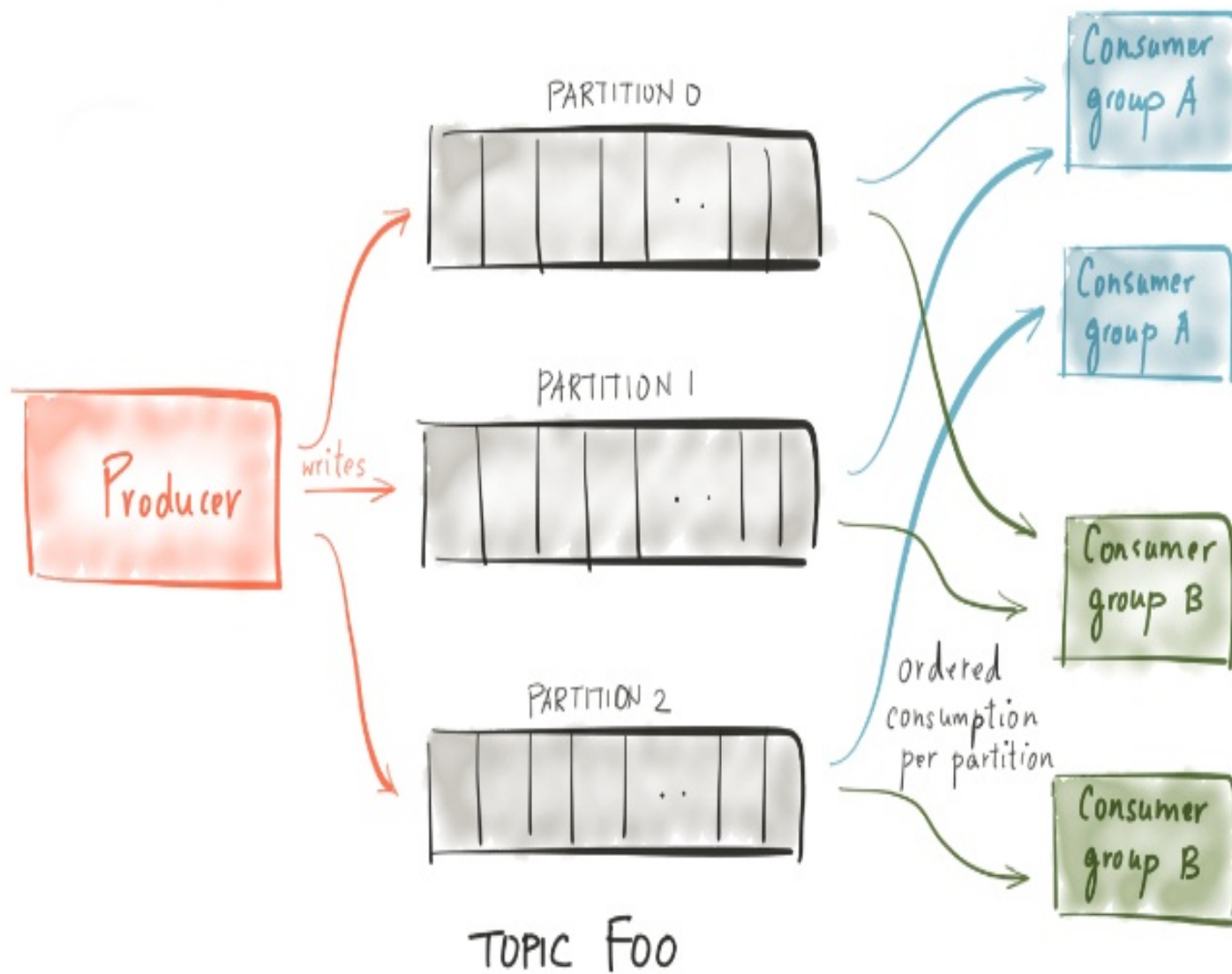
Example: LinkedIn back in 2010



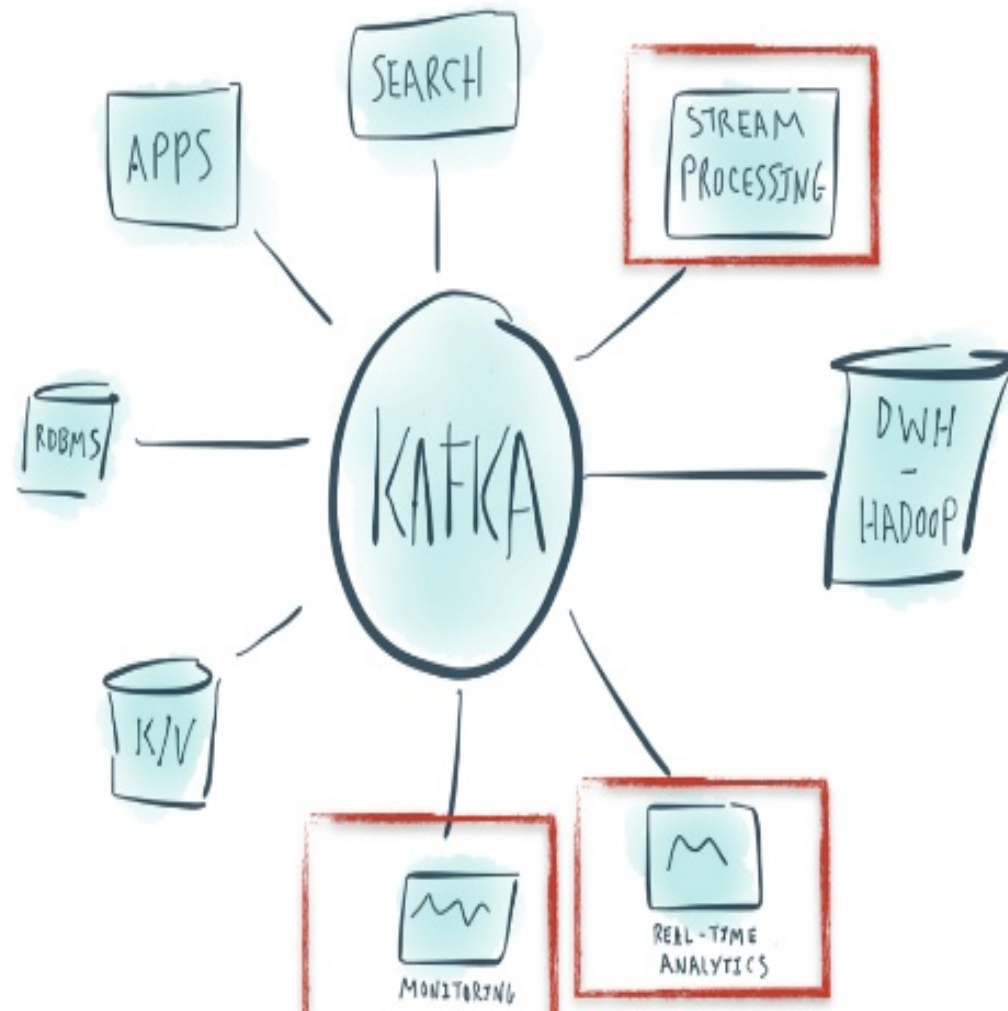
Example: LinkedIn with Kafka

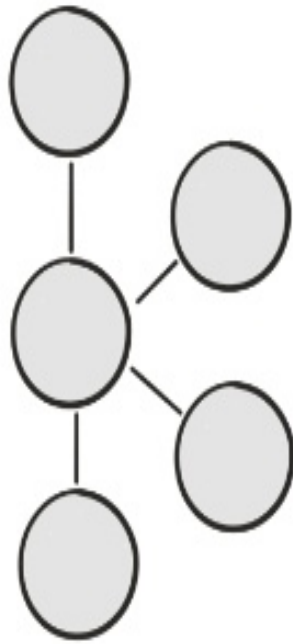


SCALABLE CONSUMPTION



STREAMING PLATFORM





Kafka Connect

Large-scale streaming data import/export for Kafka

GOALS

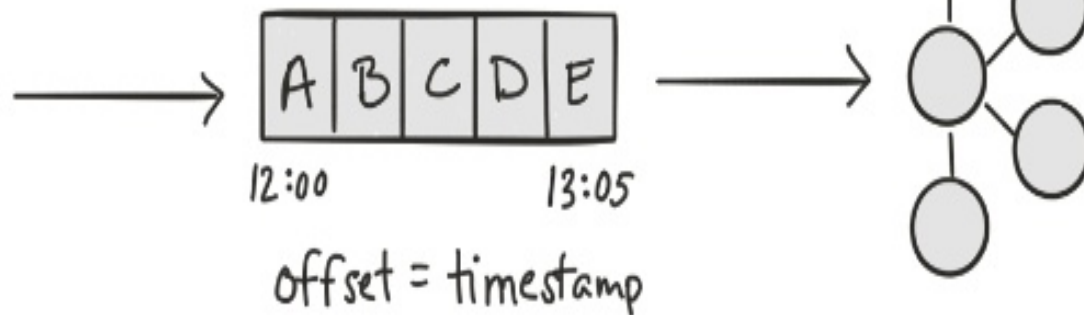
1. Focus on copying
2. Batteries included
3. Standardize
4. Parallelism
5. Scale



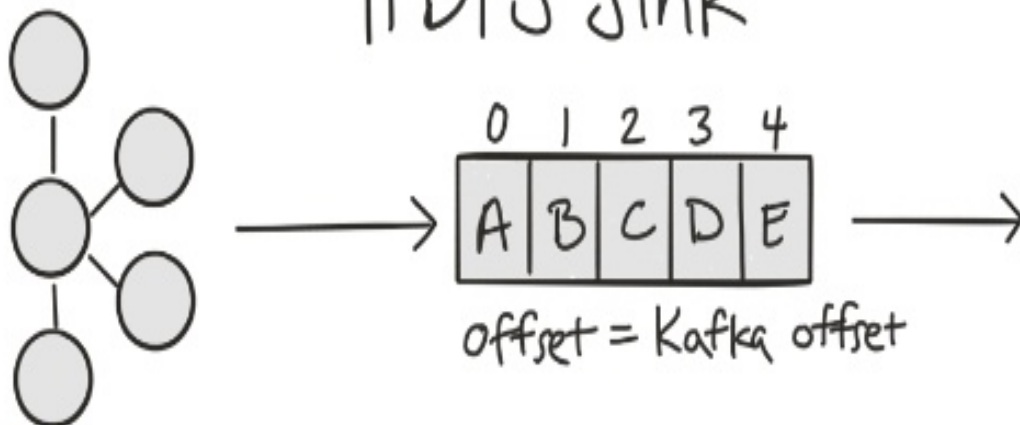
Table

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

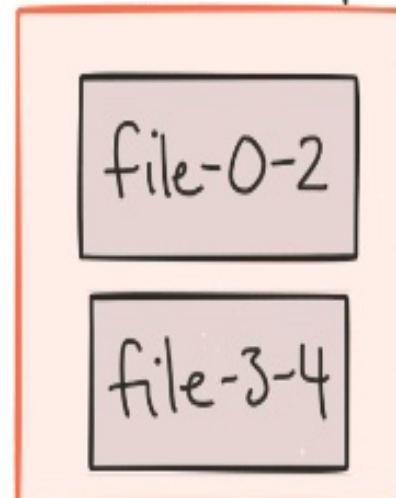
Database Source



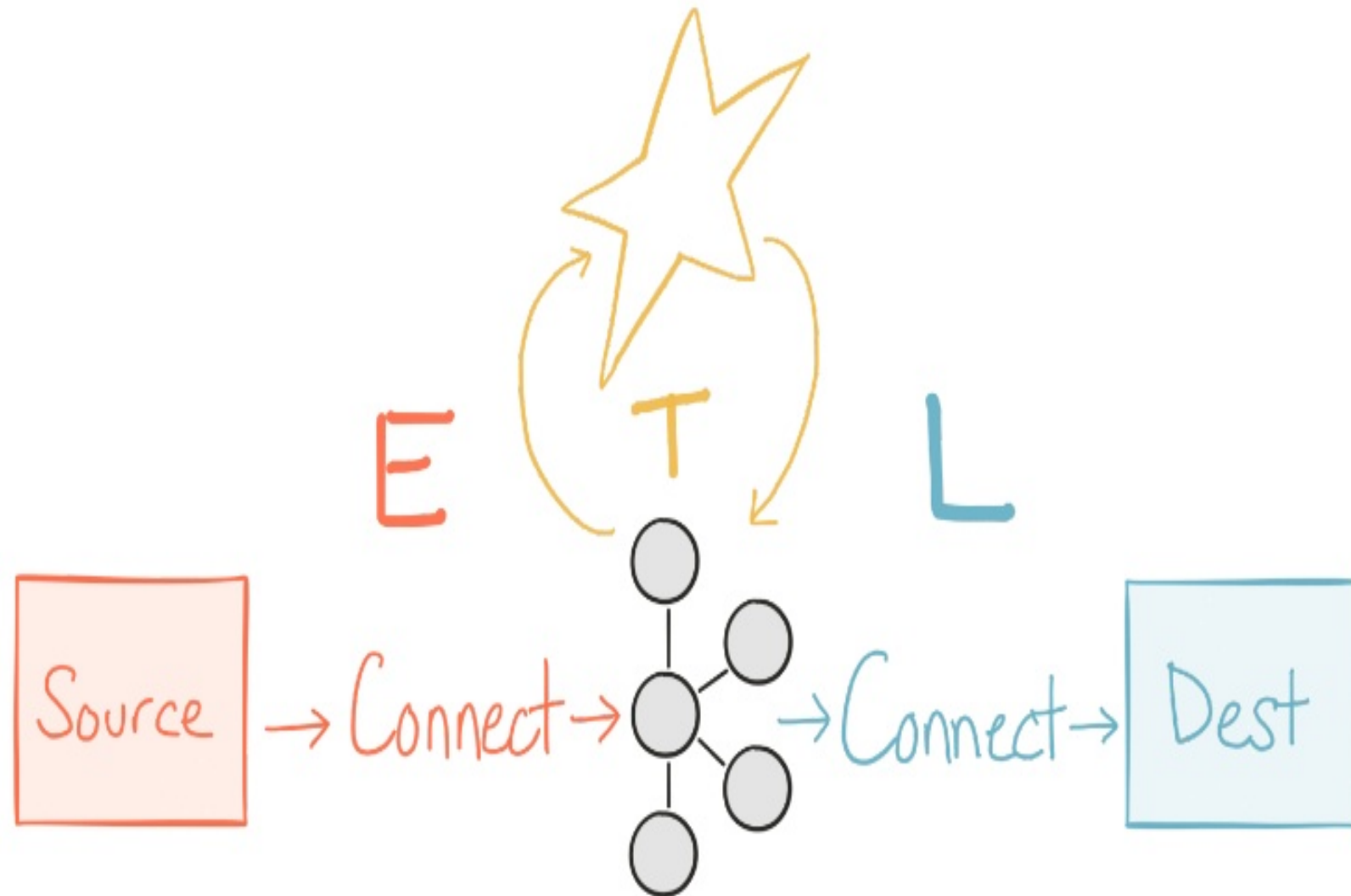
HDFS Sink



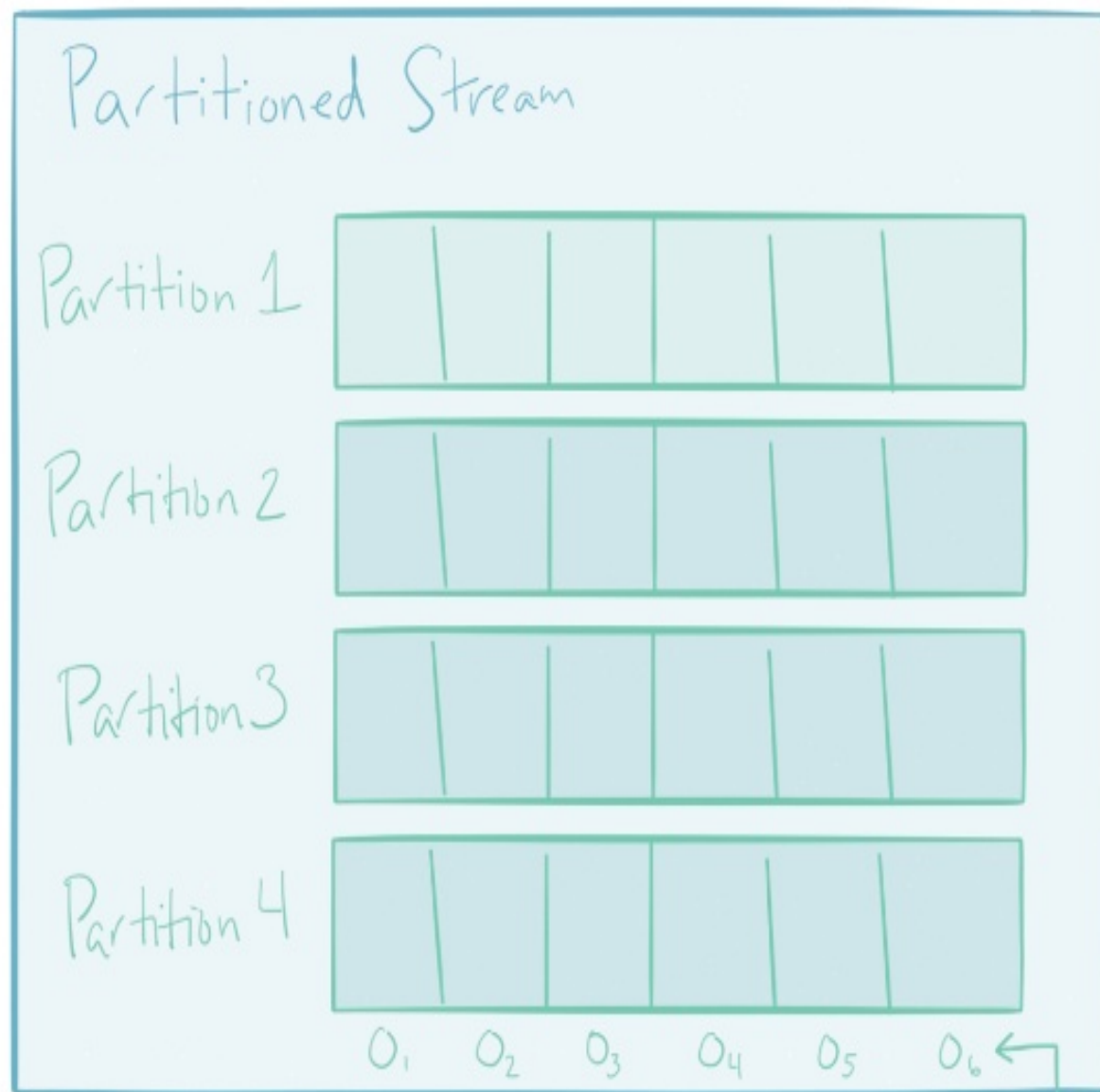
HDFS Directory



Separation of Concerns



Data Model



Data Model

Partitioned Stream - Database

Table 1

--	--	--	--	--	--

Table 2

--	--	--	--	--	--

Table 3

--	--	--	--	--	--

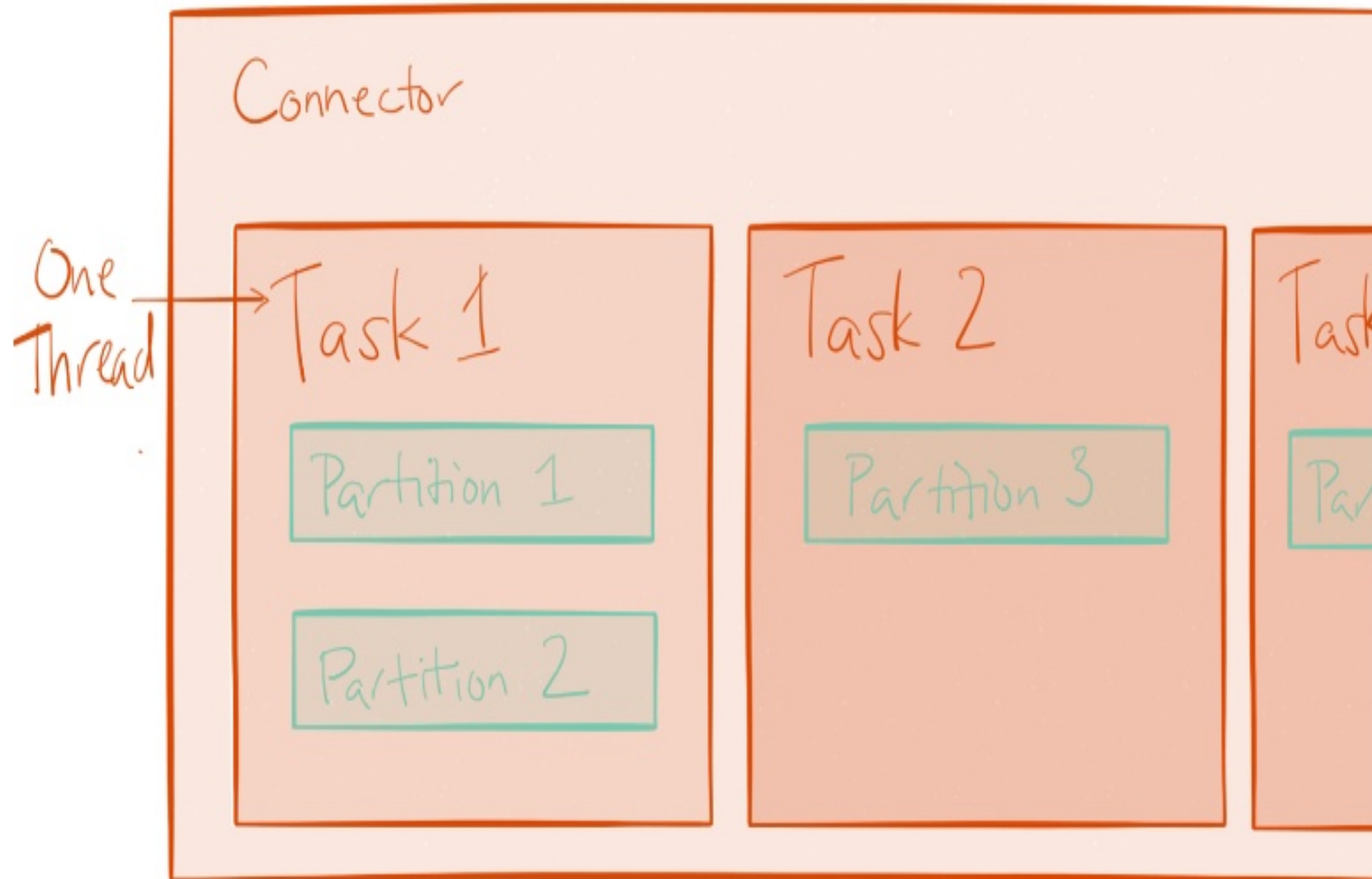
Table 4

--	--	--	--	--	--

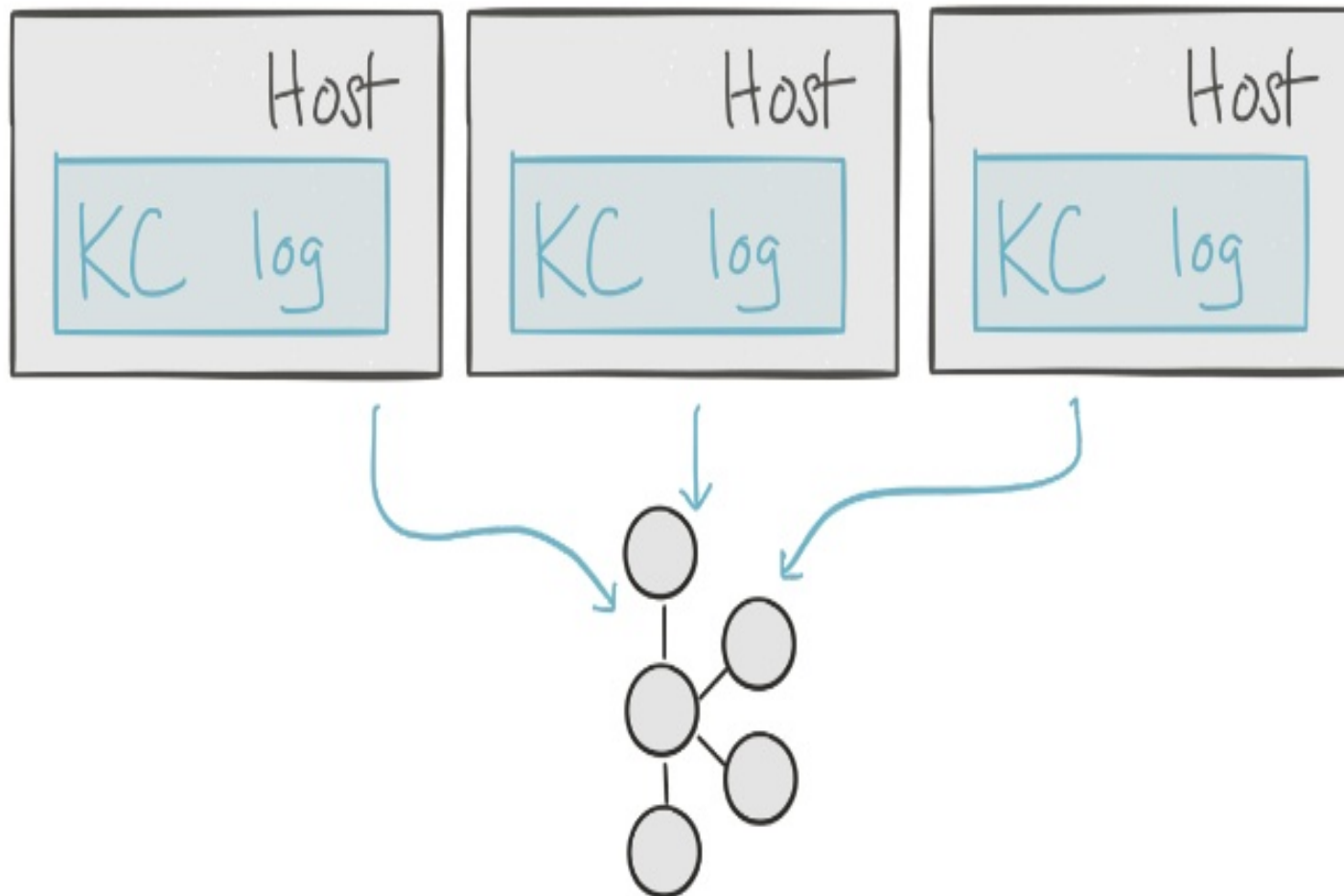
id=1 id=2 id=3 id=4 id=5 id=6



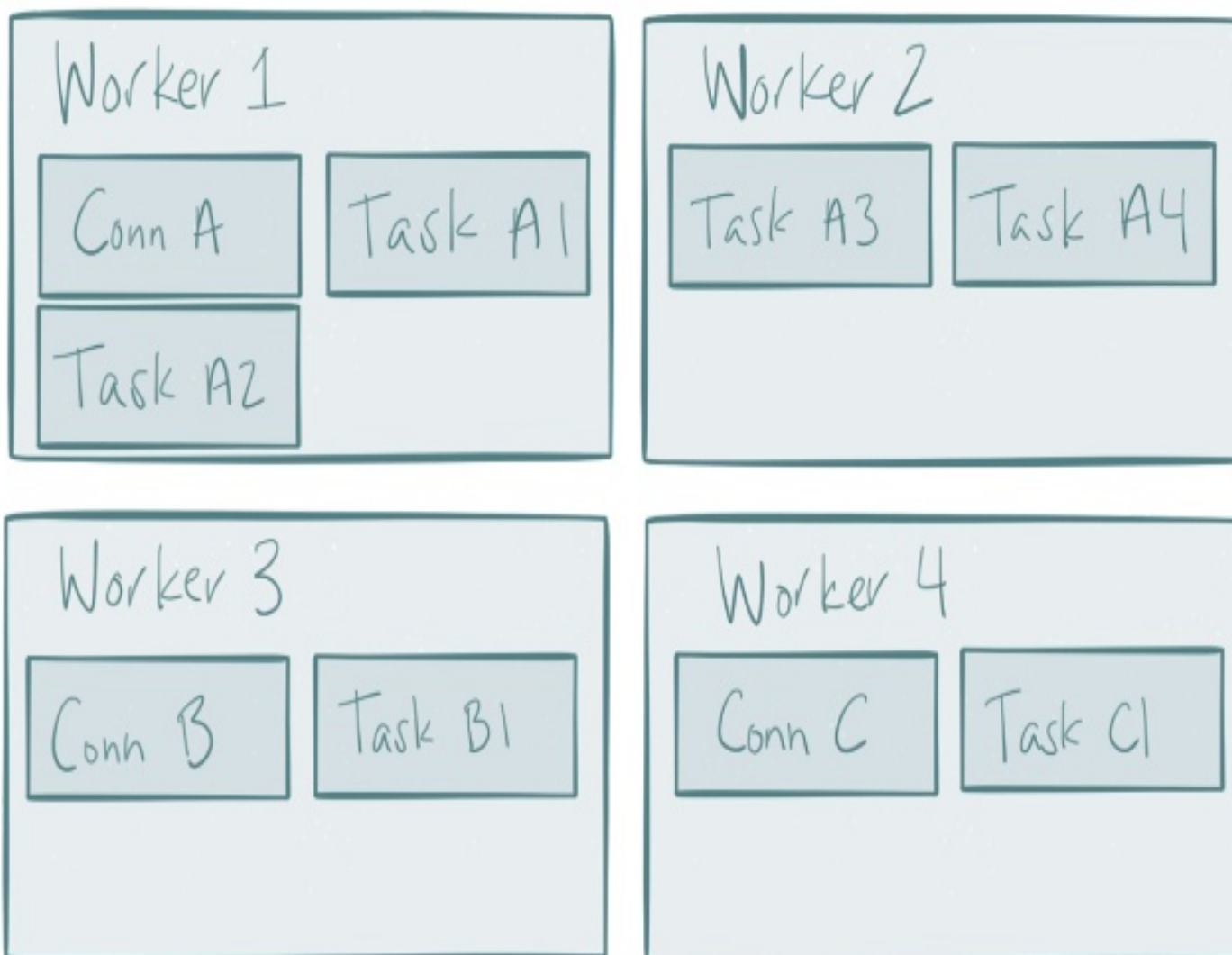
Parallelism Model



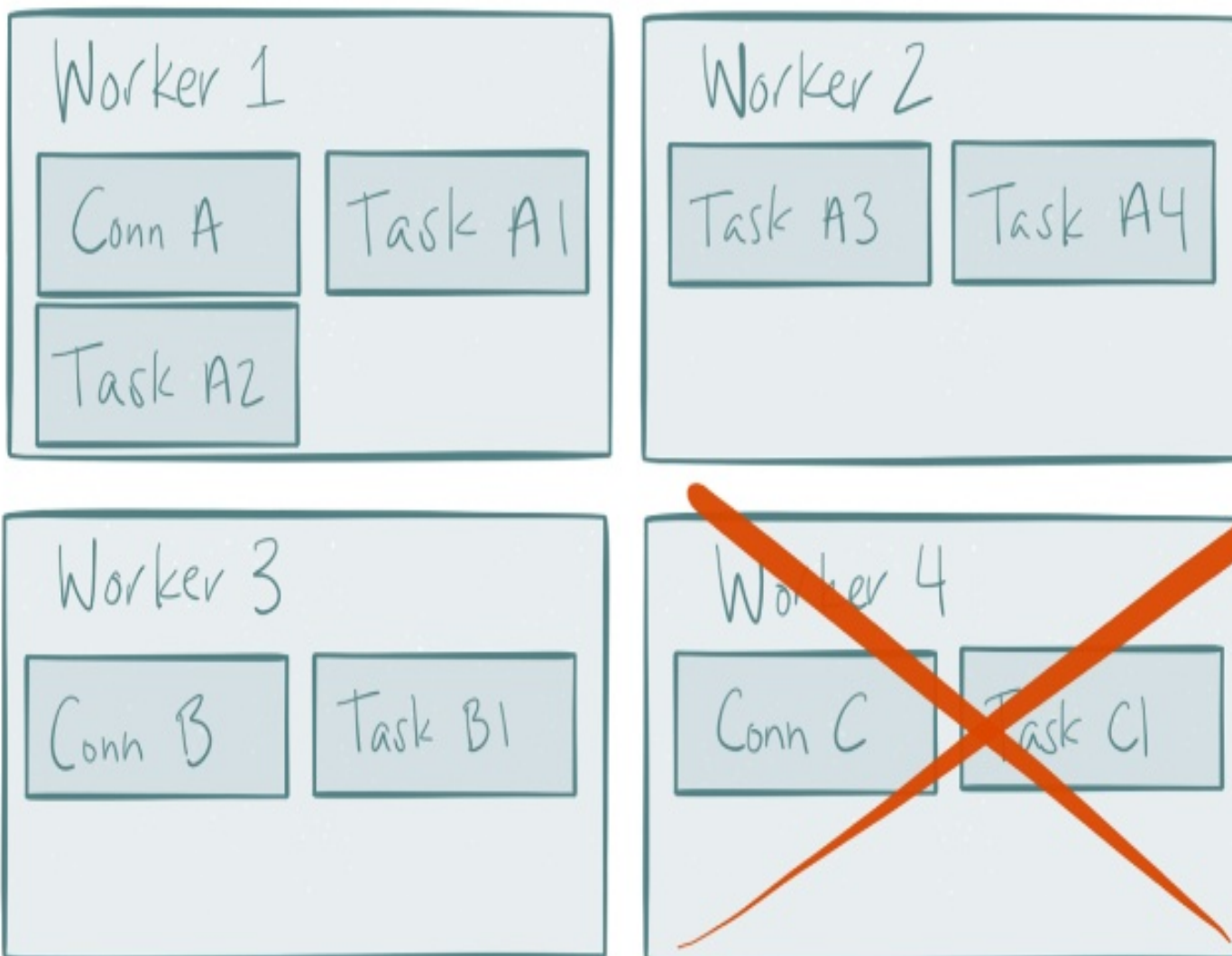
Standalone Execution



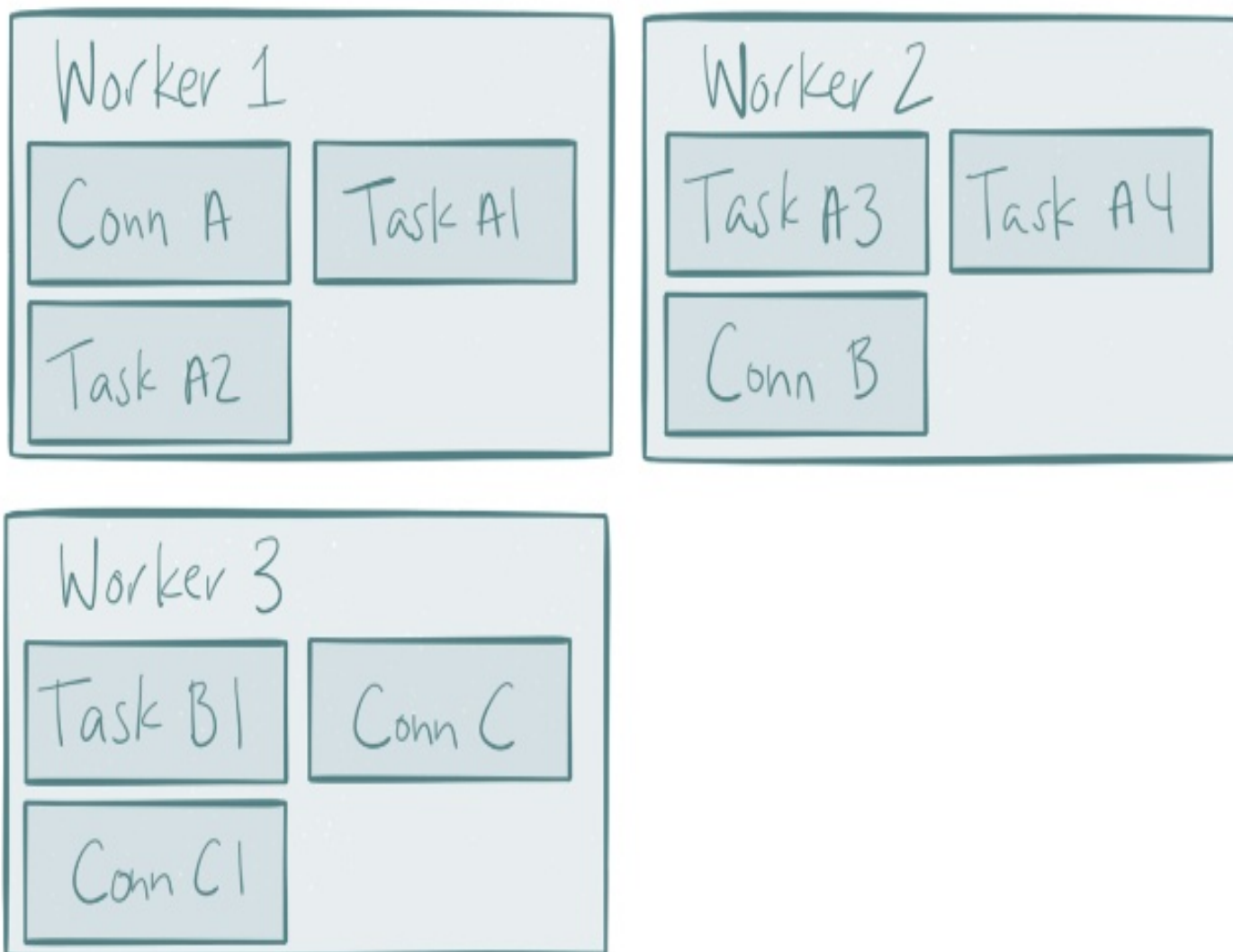
Distributed Execution



Distributed Execution



Distributed Execution



Delivery Guarantees

- Offsets automatically committed and restored
- On restart: task checks offsets & rewinds
- At least once delivery – flush data, then commit
 - Exactly once for connectors that support it (e.g. HDFS)



Format Converters

- Abstract serialization agnostic to connectors
 - Convert between Kafka Connect Data API (Connectors) and serialized bytes
 - JSON and Avro currently supported



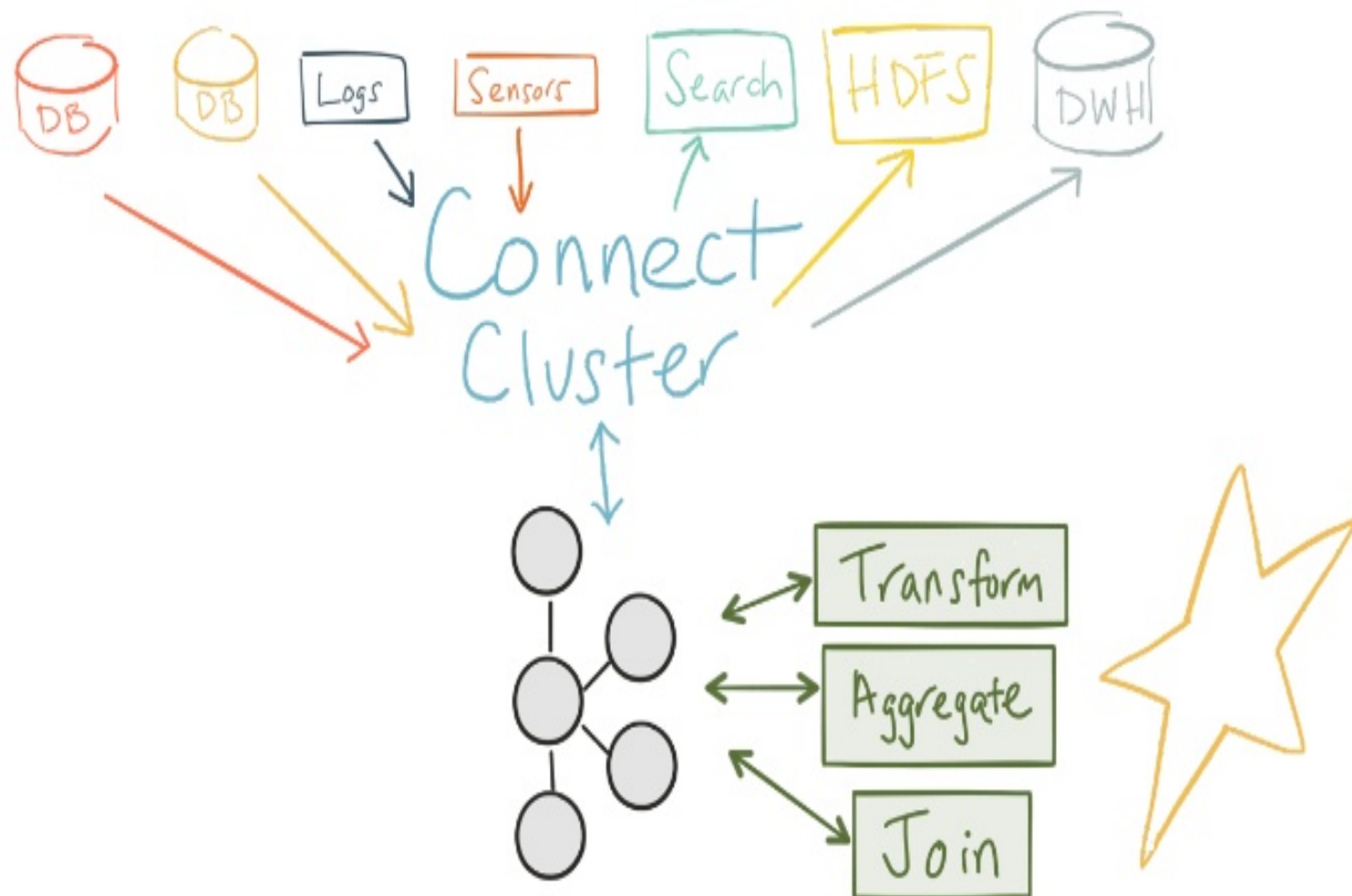
Connector Developer APIs

```
class Connector {  
  
    abstract void start(props);  
  
    abstract void stop();  
  
    abstract Class<? extends Task> taskClass();  
  
    abstract List<Map<...>> taskConfigs(maxTasks);  
  
    ...  
}
```

```
class Source/SinkTask {  
  
    abstract void start(props);  
  
    abstract void stop();  
  
    abstract List<SourceRecord> poll();  
  
    abstract void put(records);  
  
    abstract void commit();  
  
    ...  
}
```



Kafka Connect & Spark Streaming



Kafka Connect Today

- Confluent open source: HDFS, JDBC
- Connector Hub: connectors.confluent.io
 - Examples: MySQL, MongoDB, Twitter, Solr, S3, MQTT, Couchbase, Vertica, Cassandra, Elastic Search, HBase, Kudu, Attunity, JustOne, Striim, Bloomberg ..
- Improved connector control (0.10.0)



THANK YOU!

Guozhang Wang | guozhang@confluent.io | @guozhangwang

Confluent – Afternoon Break Sponsor for Spark Summit

- Jay Kreps – I Heart Logs book signing and giveaway
- 3:45pm – 4:15pm in Golden Gate

Kafka Training with Confluent University

- Kafka Developer and Operations Courses
- Visit www.confluent.io/training

Want more Kafka?

- Download Confluent Platform Enterprise (incl. Kafka Connect) at <http://www.confluent.io/product>
- Apache Kafka 0.10 upgrade documentation at



SPARK SUMMIT 2016
DATA SCIENCE AND ENGINEERING AT SCALE