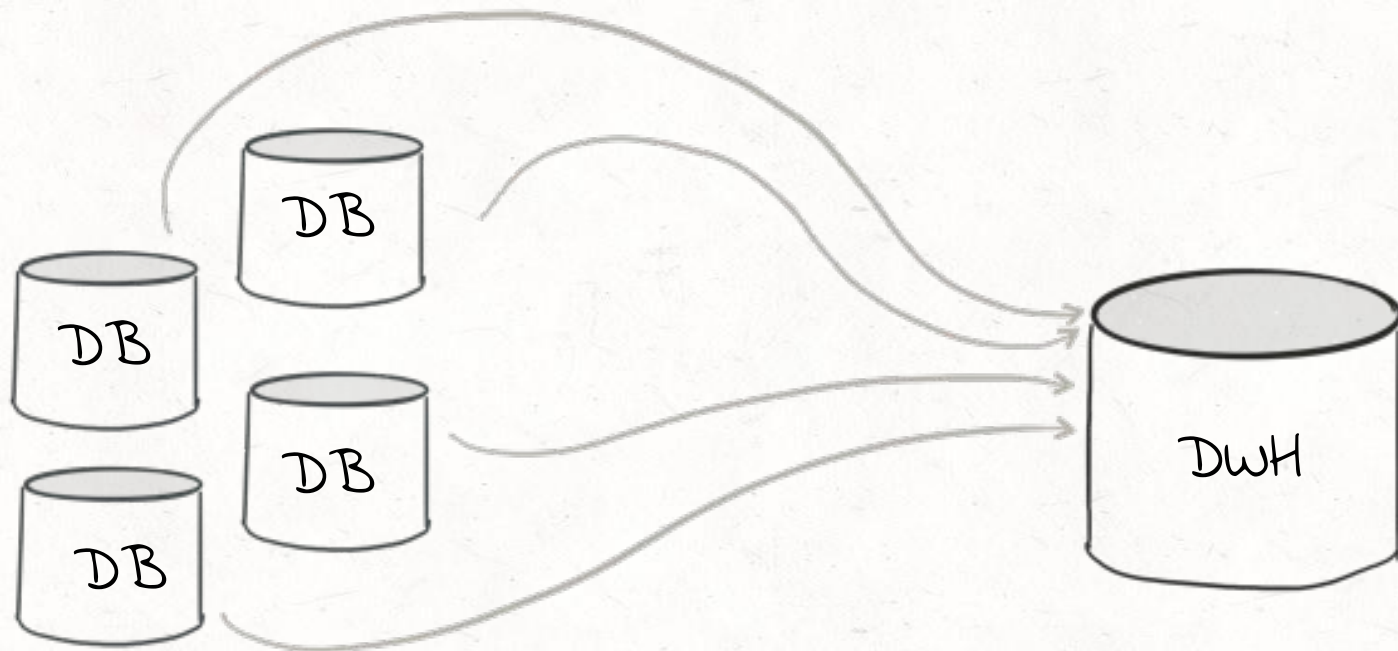# ETL IS DEAD; LONG-LIVE STREAMS

Neha Narkhede,
Co-founder & CTO, Confluent

confluent

"

Data and data systems have really changed in the past decade

confluent

# Old world: Two popular locations for data



operational databases

Relational data warehouse

"

*Several recent data trends are driving a dramatic change in the ETL architecture*

confluent

"#1: Single-server databases are replaced by a myriad of distributed data platforms that operate at company-wide scale

"

#2: Many more types of data sources beyond transactional data - logs, sensors, metrics...
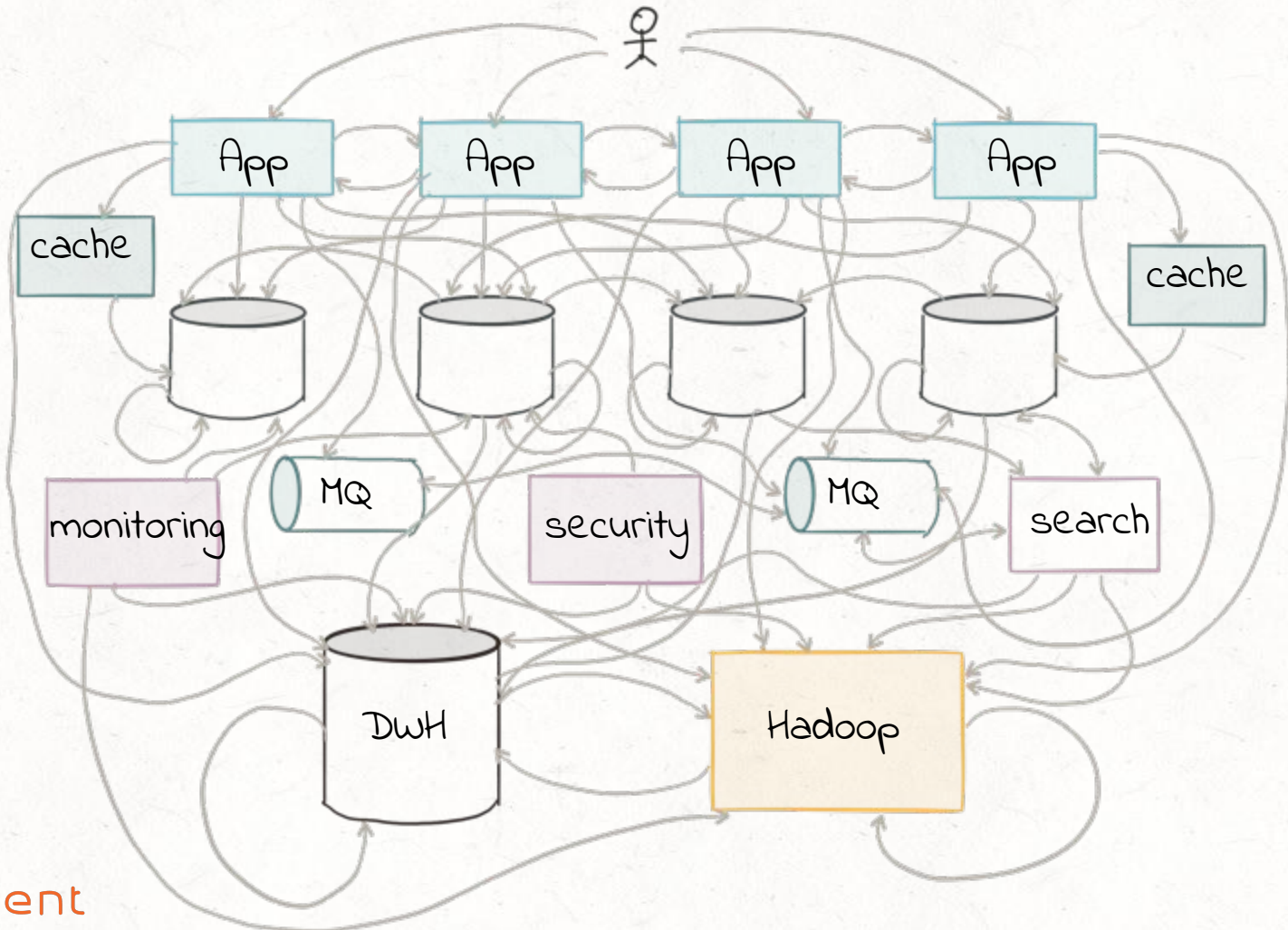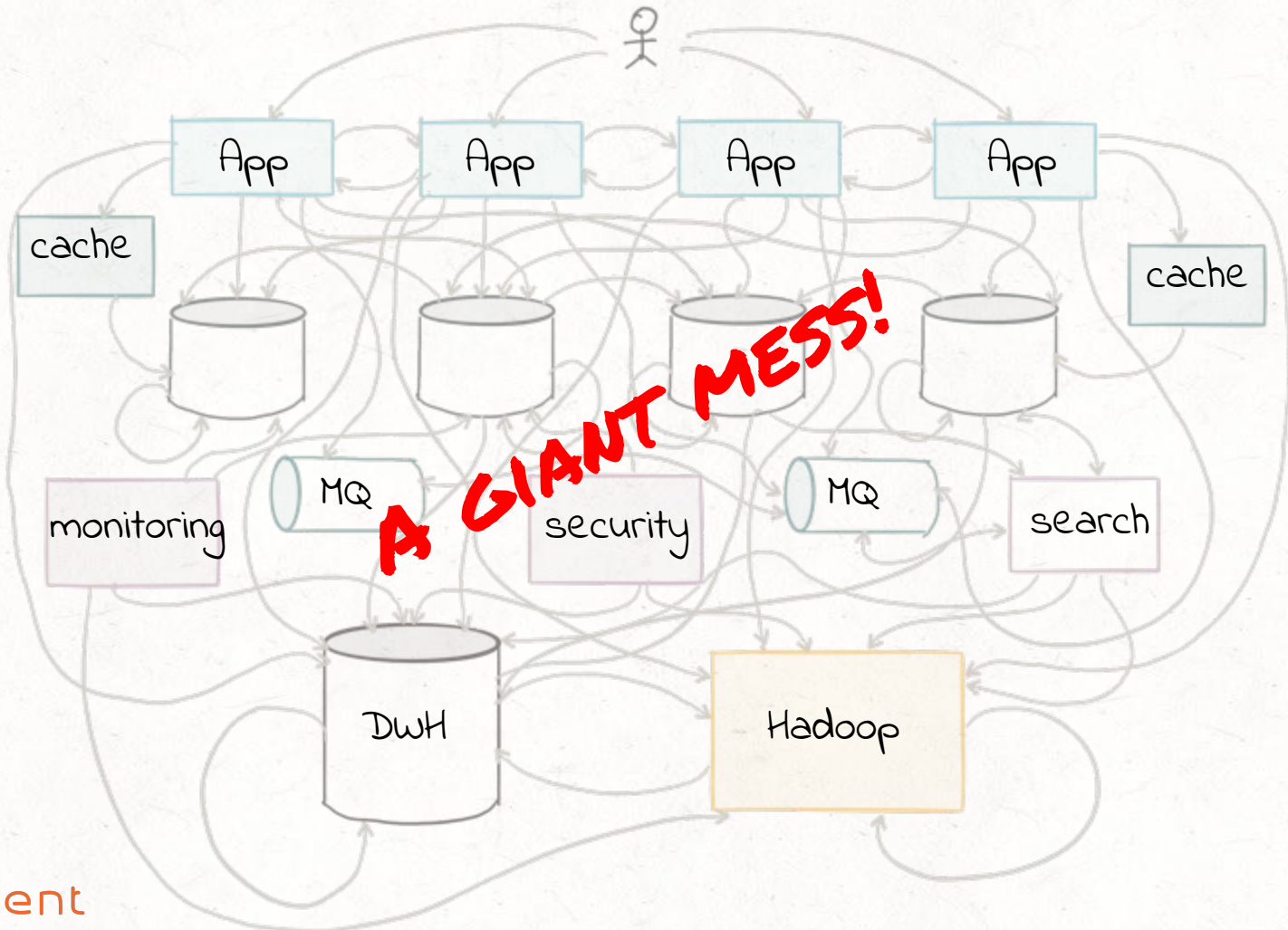
#3: Stream data is increasingly ubiquitous; need for faster processing than daily

"

The end result? This is what data integration ends up looking like in practice
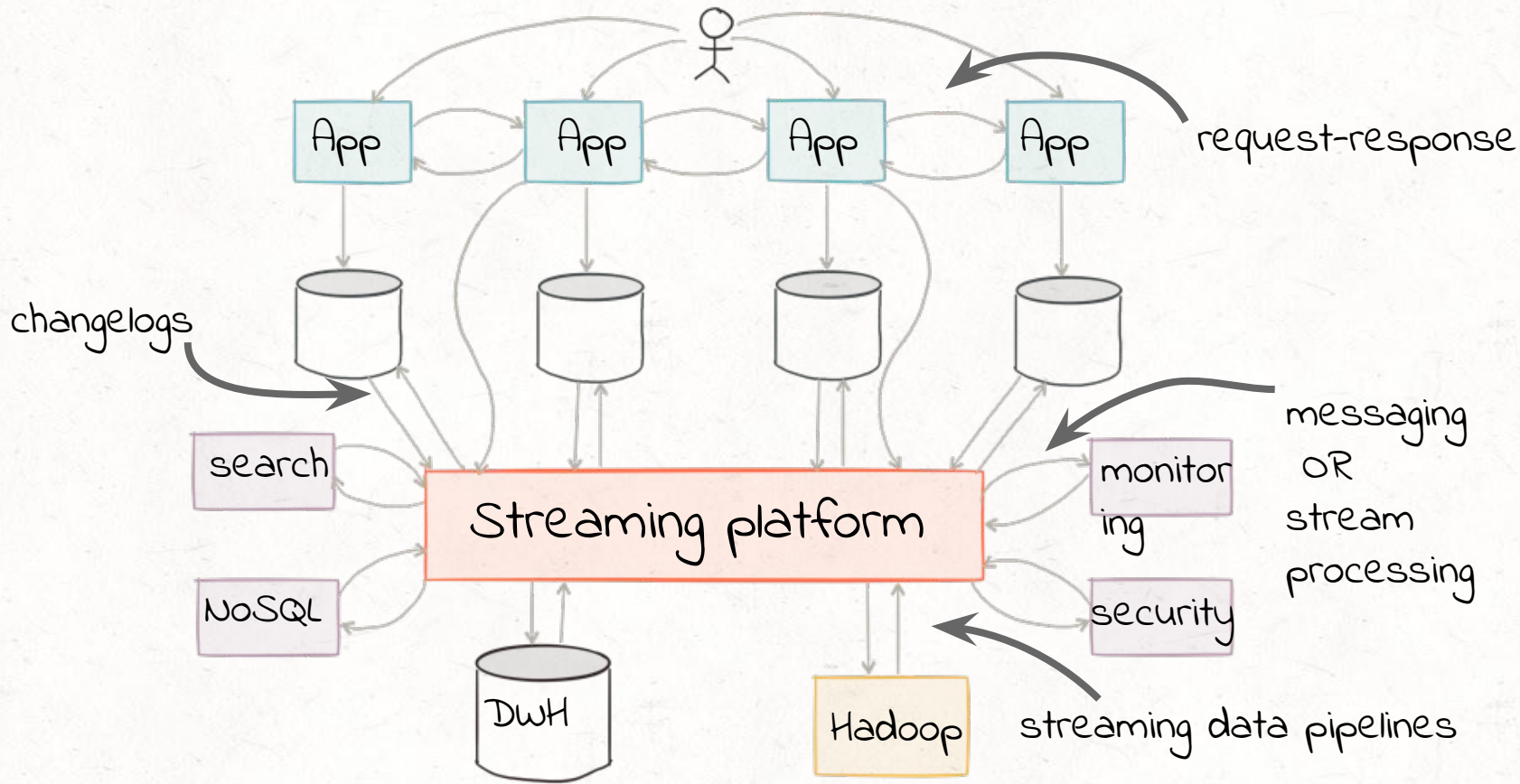
confluent

App · App · App · App

cache · cache

monitoring · MQ · security · MQ · search

DWH · Hadoop

A GIANT MESS!

confluent

"

we will see how transitioning to streams
cleans up this mess and works towards...

request-response

changelogs

messaging
OR
stream
processing

search

Streaming platform

monitor
ing

NoSQL

security

DWH

Hadoop

streaming data pipelines

App App App App

confluent

A SHORT HISTORY OF DATA INTEGRATION

Surfaced in the 1990s in retail organizations for analyzing buyer trends

**Extract** data from databases
**Transform** into destination warehouse schema
**Load** into a central data warehouse

" BUT ... ETL tools have been around for a long time, **data coverage** in data warehouses is still *low!* WHY?

ETL HAS DRAWBACKS

"

#1: The need for a global schema

confluent

**"**

#2: Data cleansing and curation is manual and fundamentally **error-prone**

confluent

" "

#3: Operational cost of ETL is high; it is slow; time and **resource intensive**

"

#4: **ETL** tools were built to narrowly focus on connecting databases and the data warehouse in a **batch** fashion

confluent

**EAI**: A different class of data integration technology for connecting applications in real-time

"

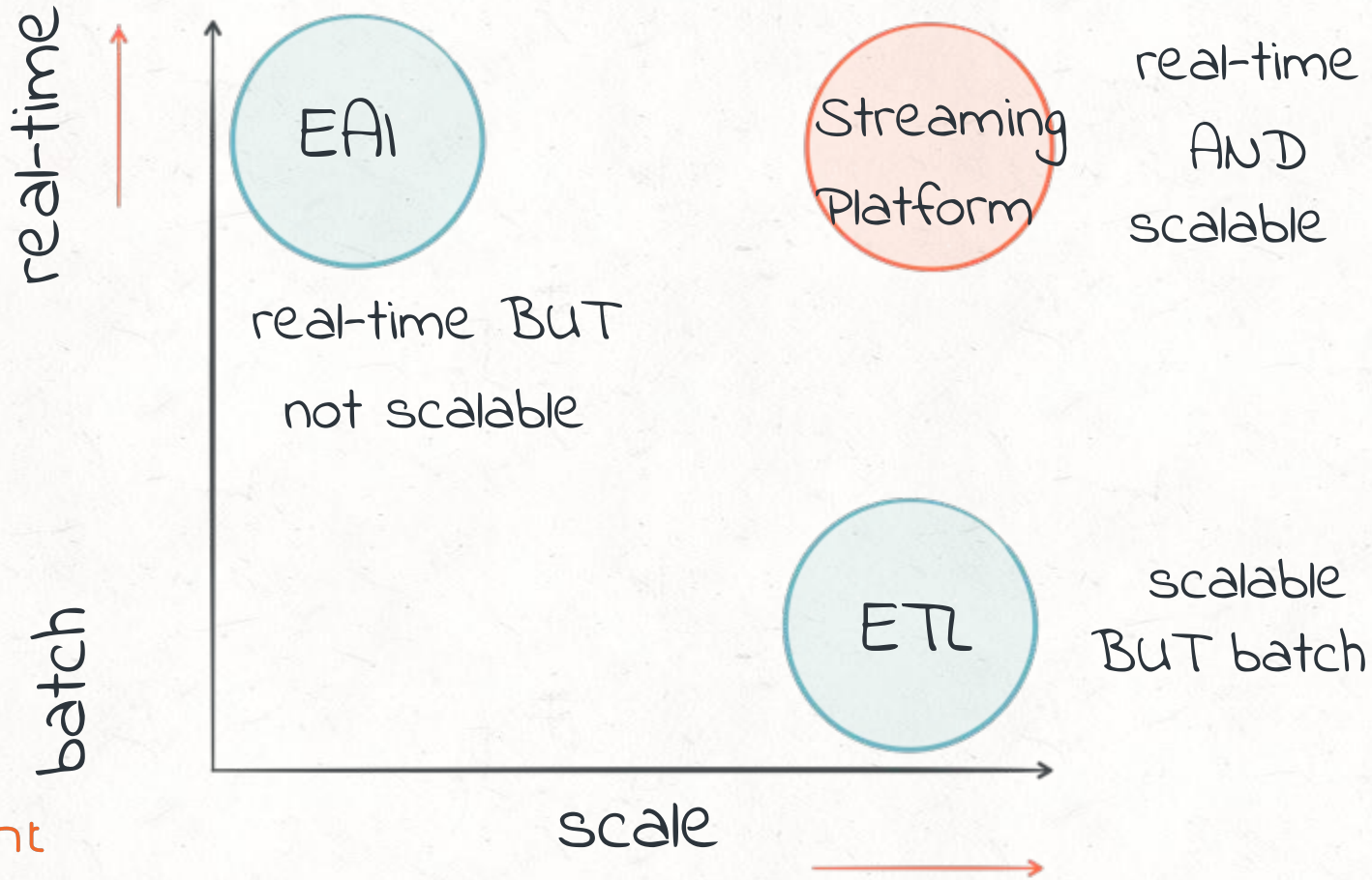EAI employed Enterprise Service Buses and MQs; weren't scalable

ETL AND EAI ARE OUTDATED!

"

Data integration and ETL in the modern world need a **complete revamp**

confluent

"
#1: Ability to process high-volume and high-diversity data

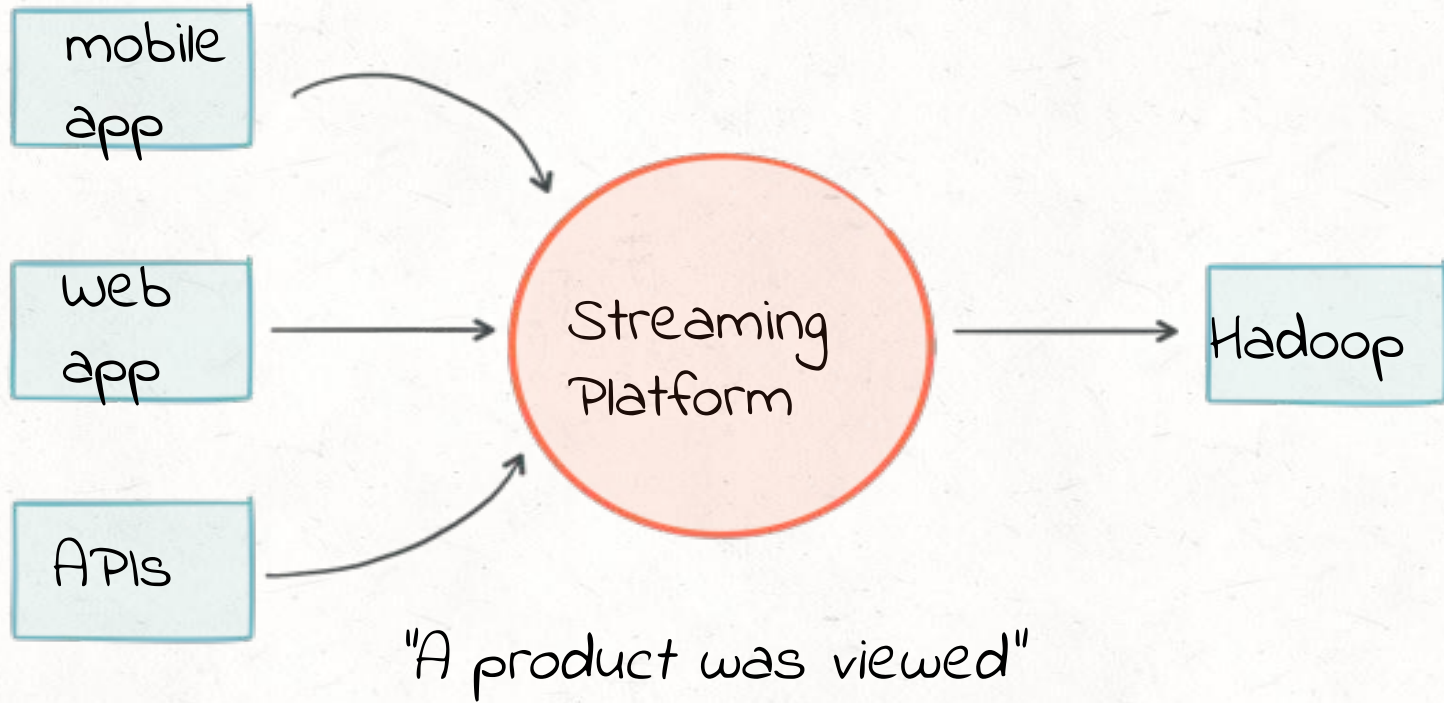#2 Real-time from the grounds up; a fundamental transition to event-centric thinking

Event-Centric Thinking

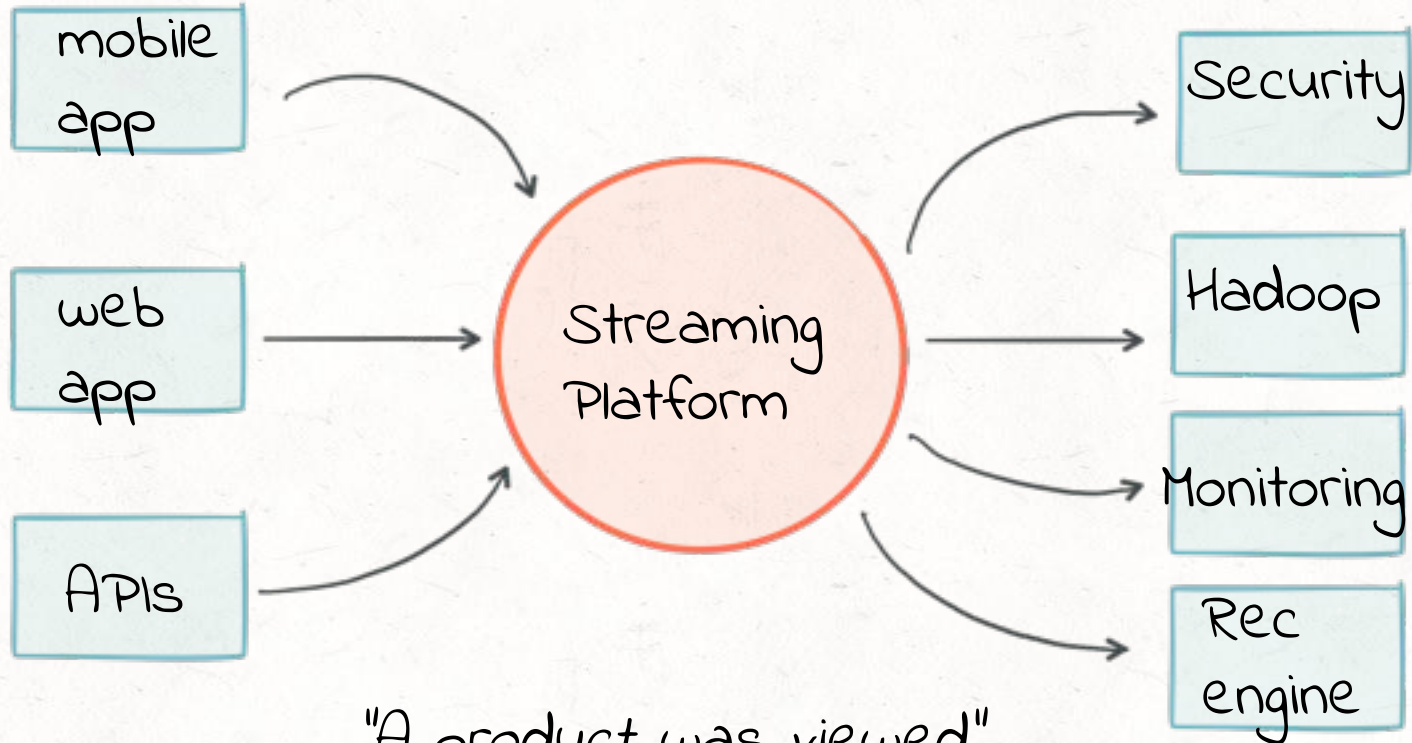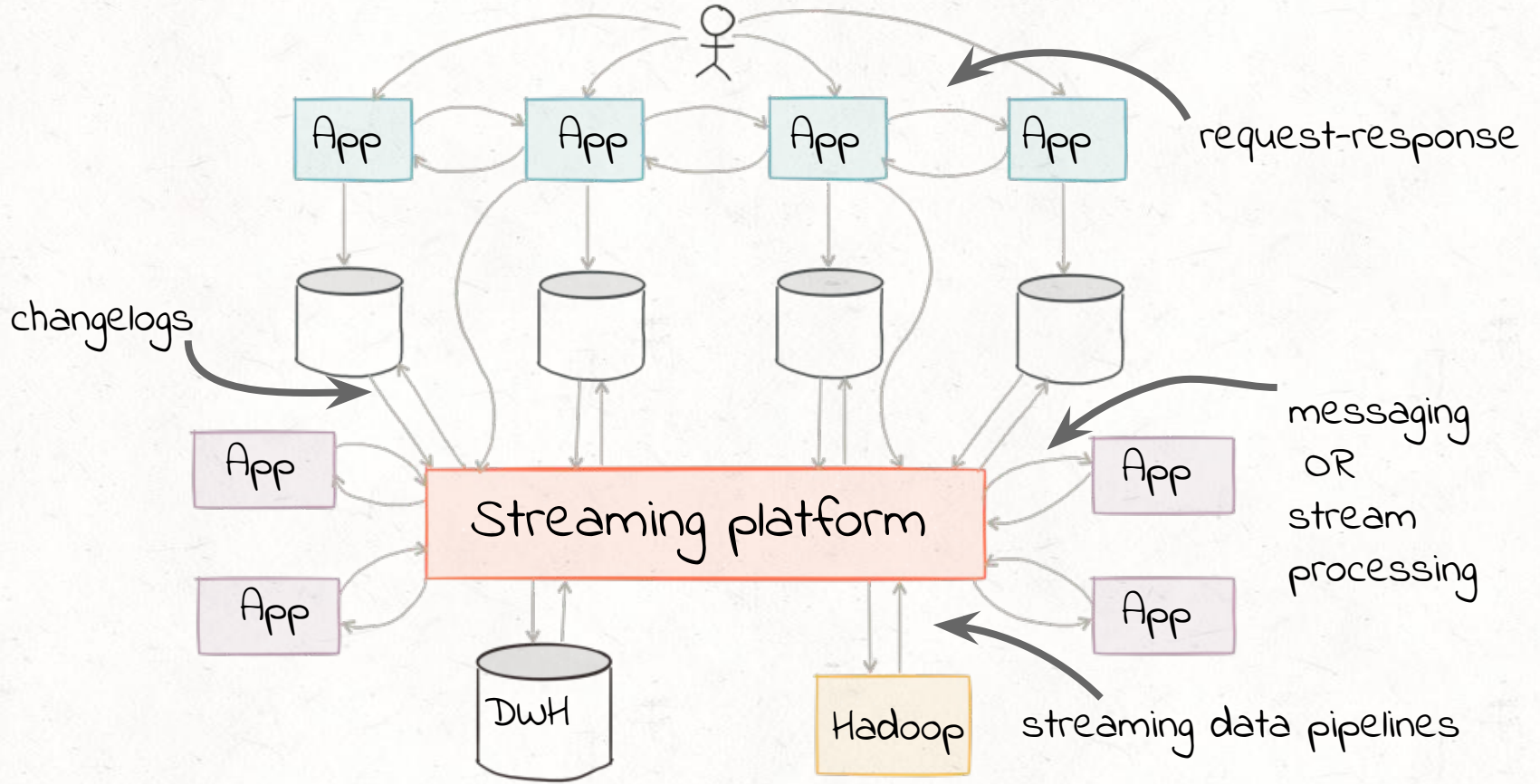# Event-Centric Thinking

"

Event-centric thinking, when applied at a company-wide scale, leads to this simplification ...

confluent

request-response

changelogs

messaging
OR
stream
processing

App

Streaming platform

App

App

DWH

Hadoop

streaming data pipelines

confluent

#3: Enable **forward-compatible** data architecture; the ability to add more applications that need to process the same data ... differently
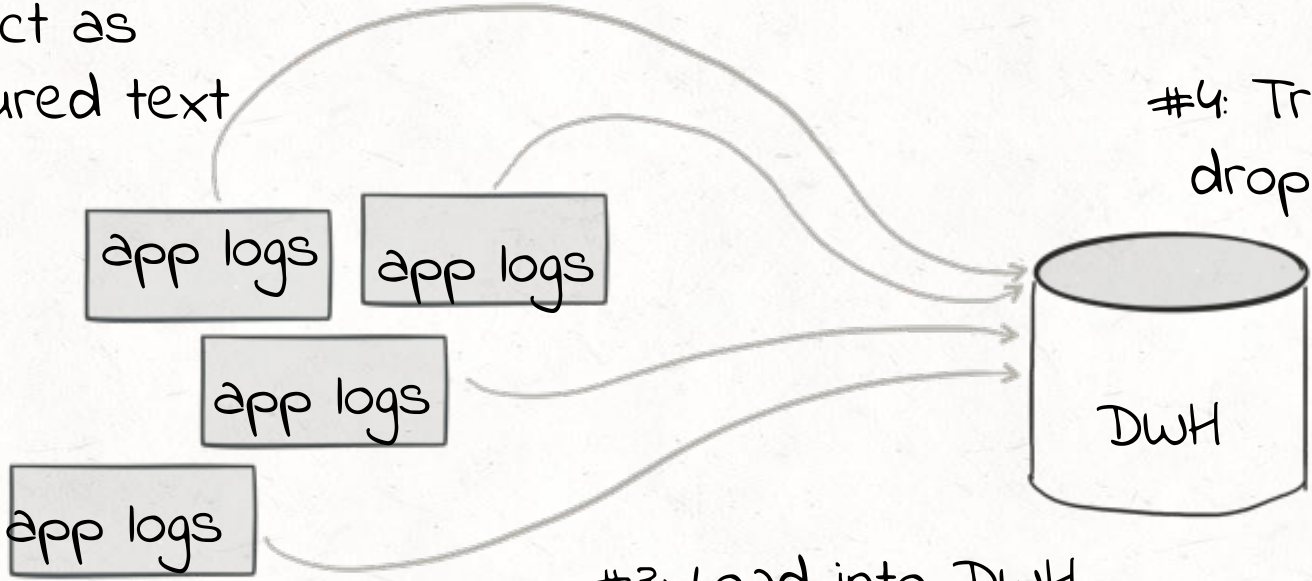
"

To enable forward compatibility, redefine

the T in ETL:

Clean data in; Clean data out

#2: Transform1 = data cleansing
= "what is a product view"

#1: Extract as
unstructured text

#4: Transform2 =
drop PII fields"

app logs

app logs

app logs

app logs

DWH

#3: Load into DWH

confluent

#1: Extract as unstructured text

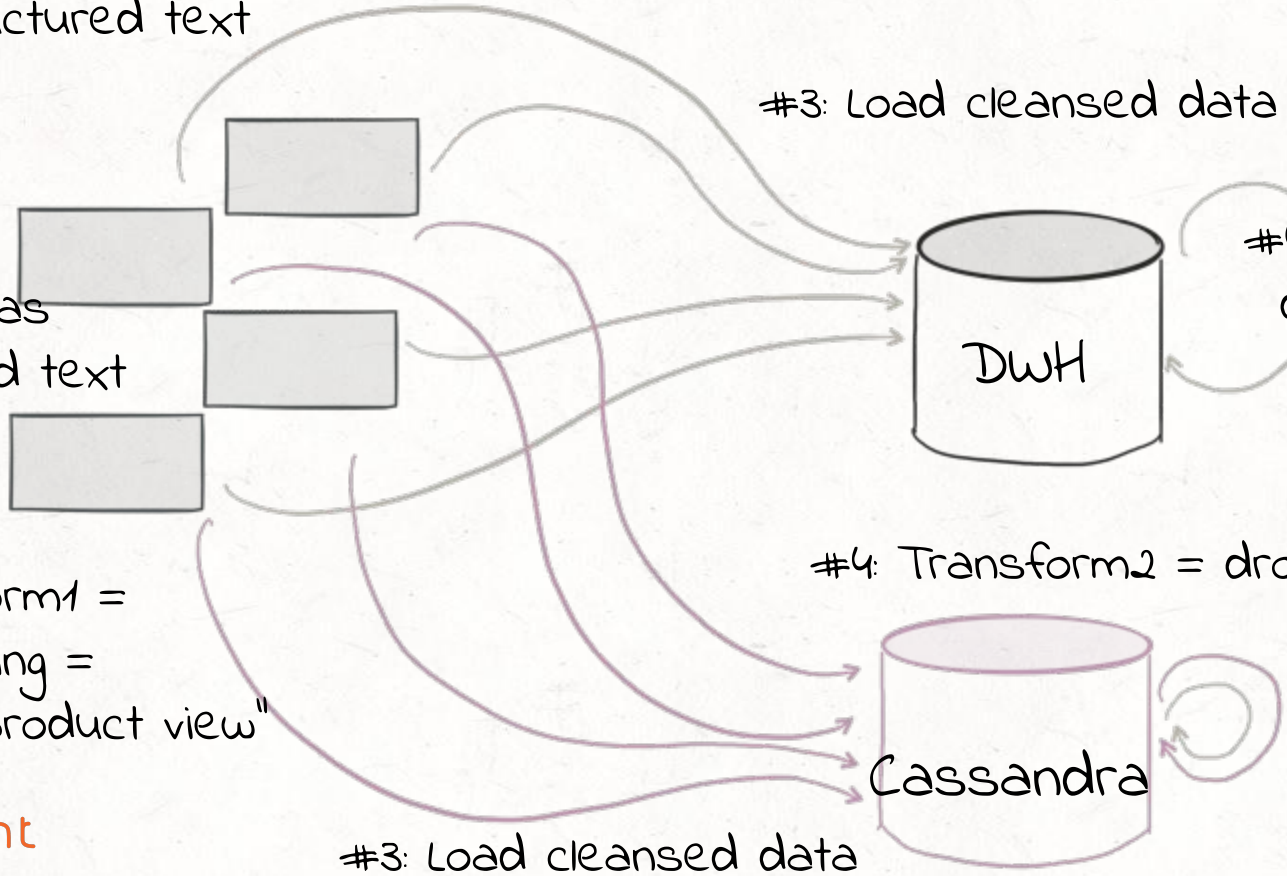#2: Transform1 = data cleansing = "what is a product view"

#3: Load cleansed data

#4: Transform2 = drop PII fields"

#1: Extract as unstructured text again

#2: Transform1 = data cleansing = "what is a product view"

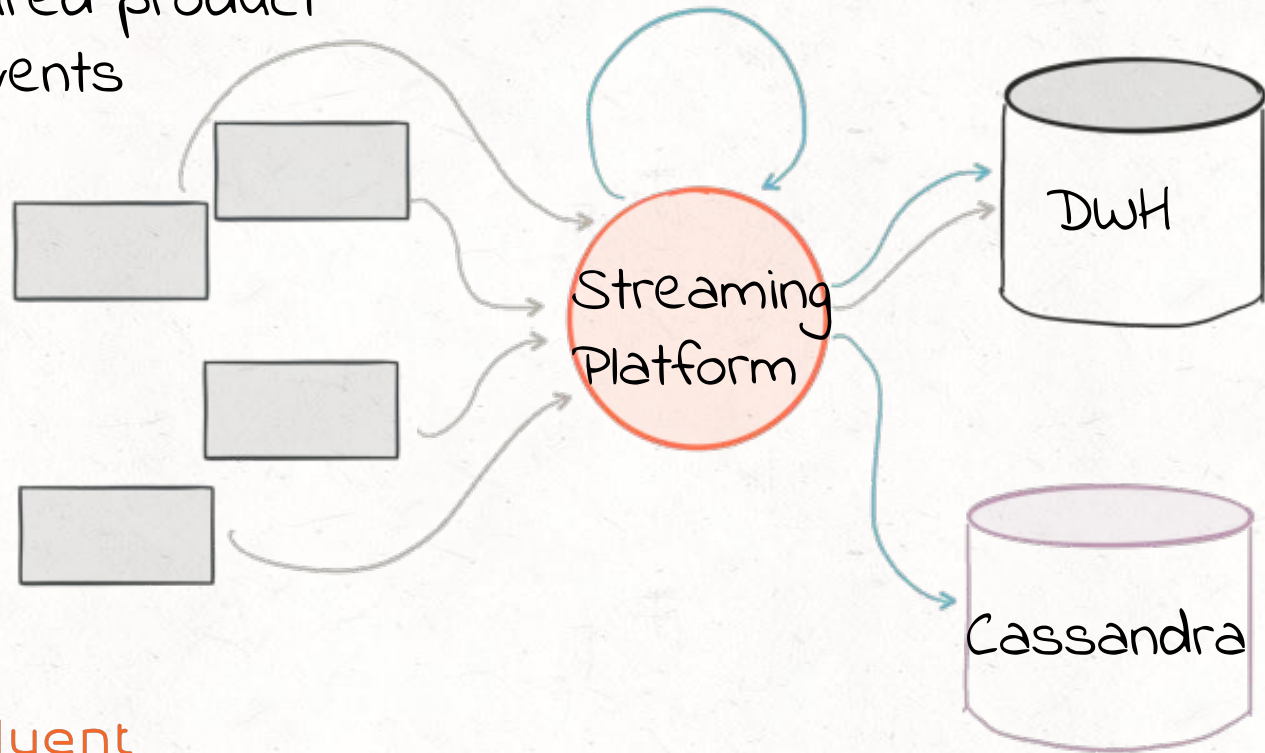#4: Transform2 = drop PII fields"

#3: Load cleansed data

DWH

Cassandra

confluent

#1: Extract as structured product view events

#2: Transforms = drop PII fields"

#4.1 Load product view stream

#4.2 Load filtered product view stream

Streaming Platform

DWH

Cassandra

#4: Load filtered product view stream

confluent

"

To enable forward compatibility, redefine
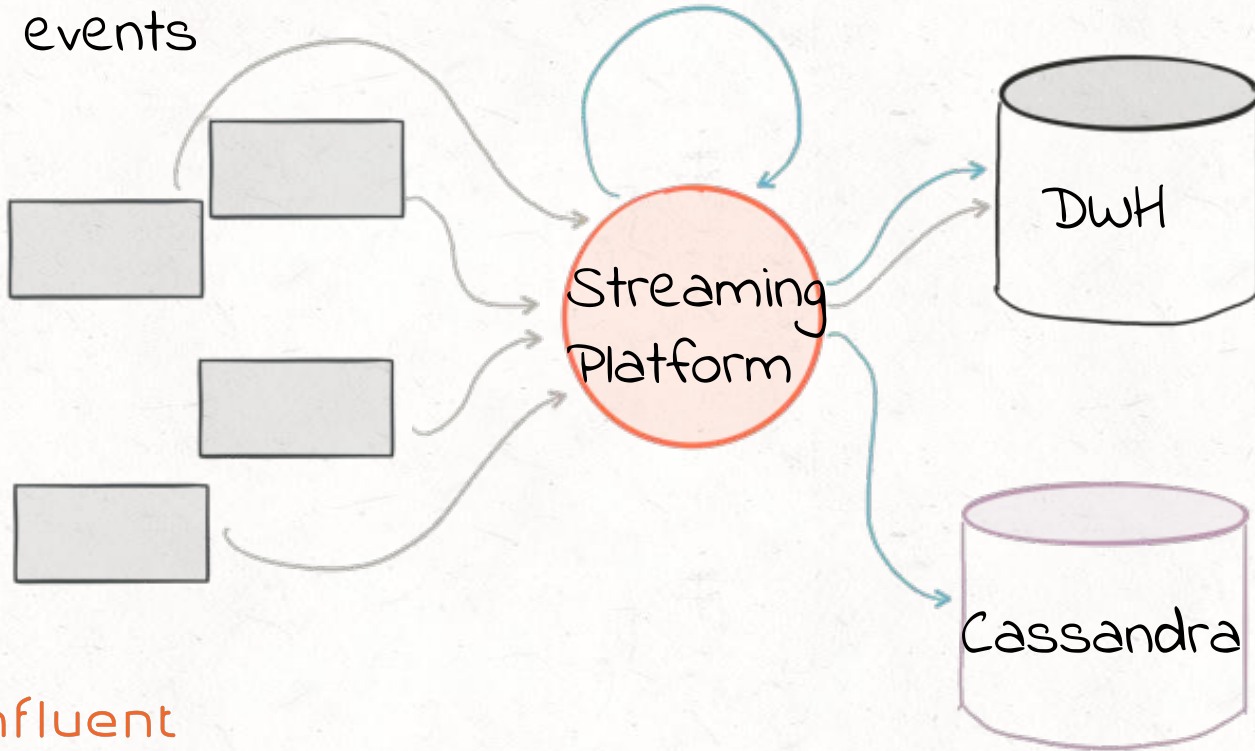
the T in ETL:

Data transformations, not data cleansing!

" Forward compatibility =
Extract clean-data once; Transform many
different ways before Loading into respective
destinations ... as and when required

confluent

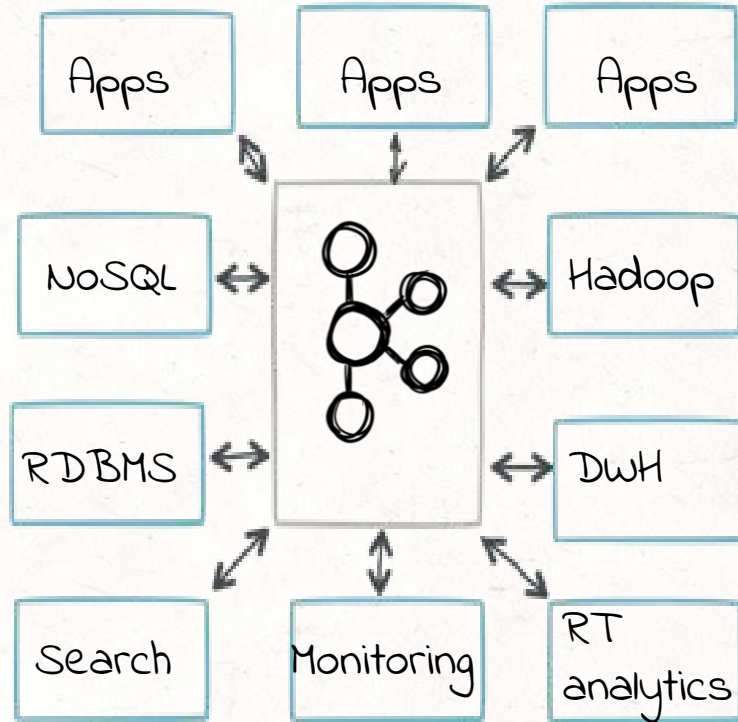# Requirements for a modern streaming data integration solution

- Fault tolerance
- Parallelism
- Latency
- Delivery semantics
- Operations and monitoring
- Schema management

confluent

# Data integration: platform vs tool

Central, reusable infrastructure for many use cases

One-off, non-reusable solution for a particular use case

confluent

# New shiny future of etl: a streaming platform

"
*Streaming platform* serves as the central nervous system for a company's data in the following ways ...

confluent

#1: Serves as the **real-time**, scalable **messaging bus** for applications; no EAI

#2: Serves as the **source-of-truth** pipeline for feeding all data processing destinations; Hadoop, DWH, NoSQL systems and more

confluent

"

#3: Serves as the **building block** for stateful **stream processing** microservices

Streaming

~~Batch~~ data integration

confluent

"

Streaming ~~Batch ETL~~

confluent

a short history of data integration

drawbacks of ETL

needs and requirements for a streaming platform

new, shiny future of ETL: a streaming platform

what does a streaming platform look like and how it enables Streaming ETL?

# Apache kafka: a distributed streaming platform

Apache kafka 6 years ago

> 1,400,000,000,000
messages processed / day

# Now Adopted at 1000s of companies worldwide



NETFLIX

intuit

Pinterest

Goldman Sachs

box

airbnb

ebaY

Adobe

yelp

Square

CISCO

dish NETWORK

salesforce

Cerner

WIKIPEDIA
The Free Encyclopedia

Y!

PayPal

verizon

confluent

"

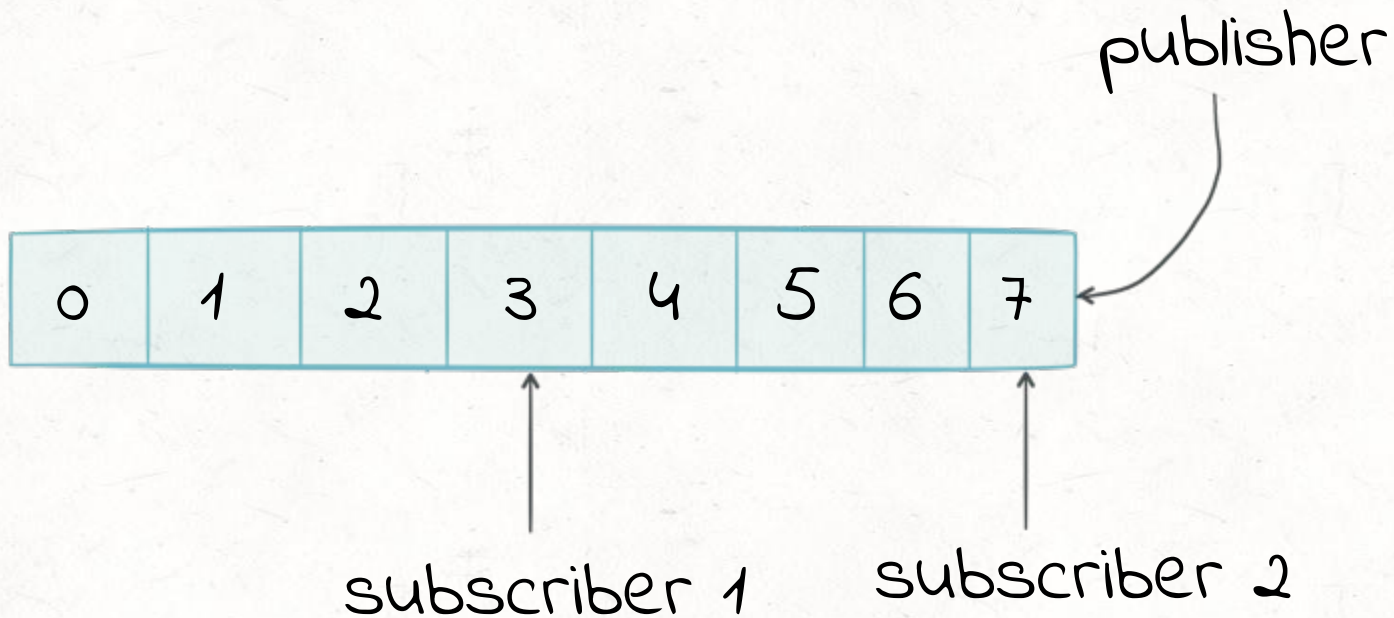what role does Kafka play in the new shiny future for data integration?

**"**

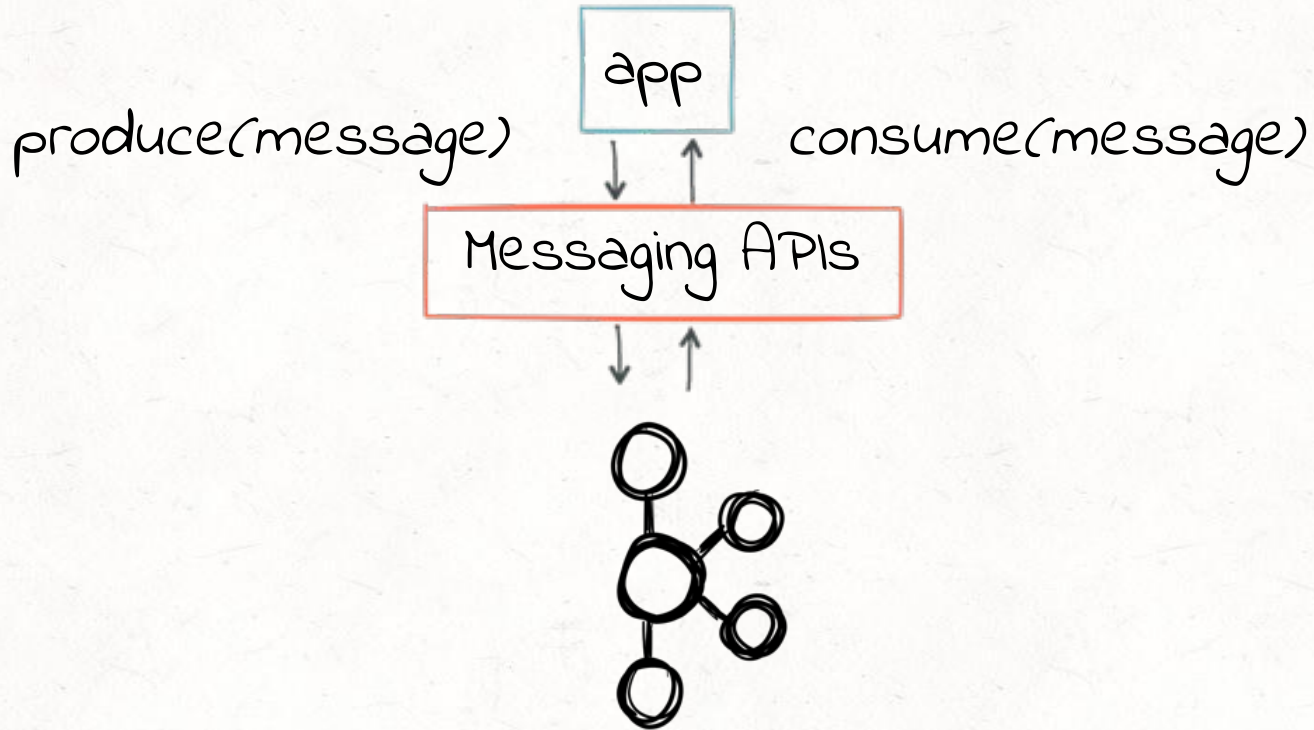#1: Kafka is the de-facto storage of choice for stream data

confluent

# The log

# The log & pub-sub

"
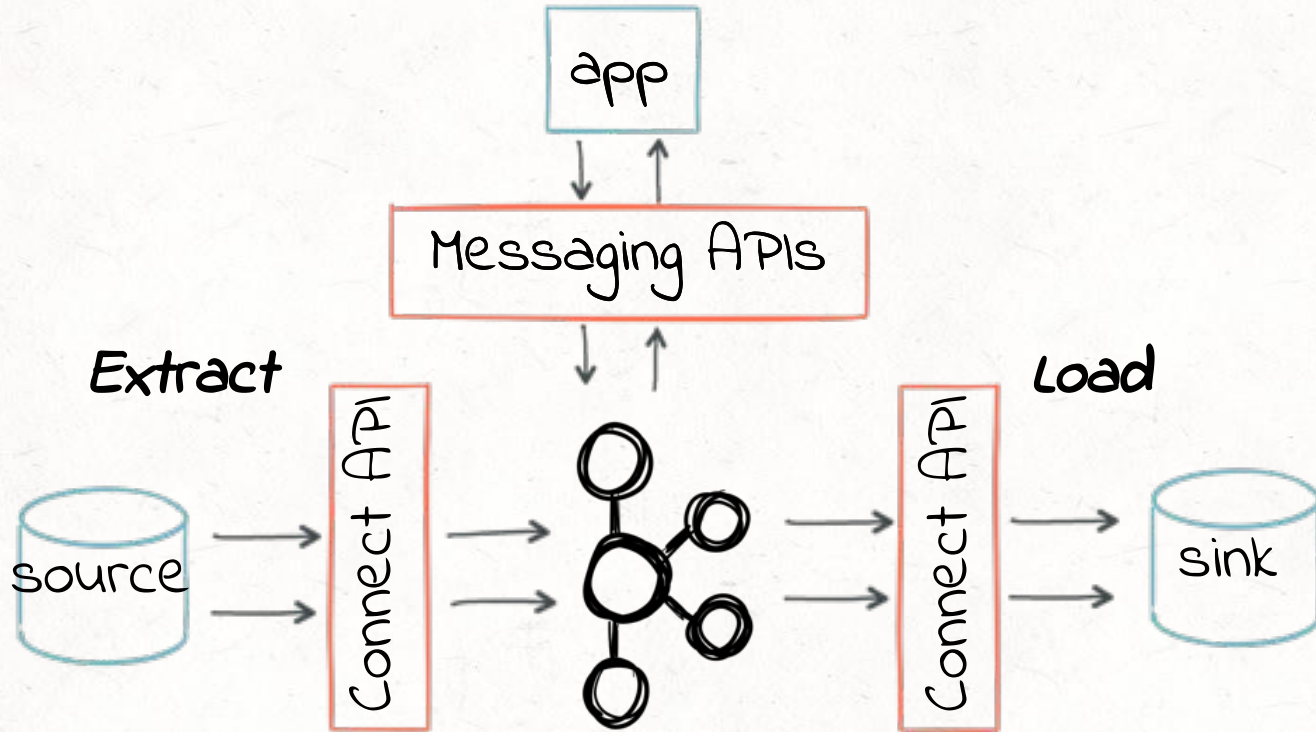#2: Kafka offers a scalable messaging backbone for application integration

confluent

# Kafka messaging APIs: scalable eai

**#3:** Kafka enables building *streaming data pipelines* (E & L in ETL)
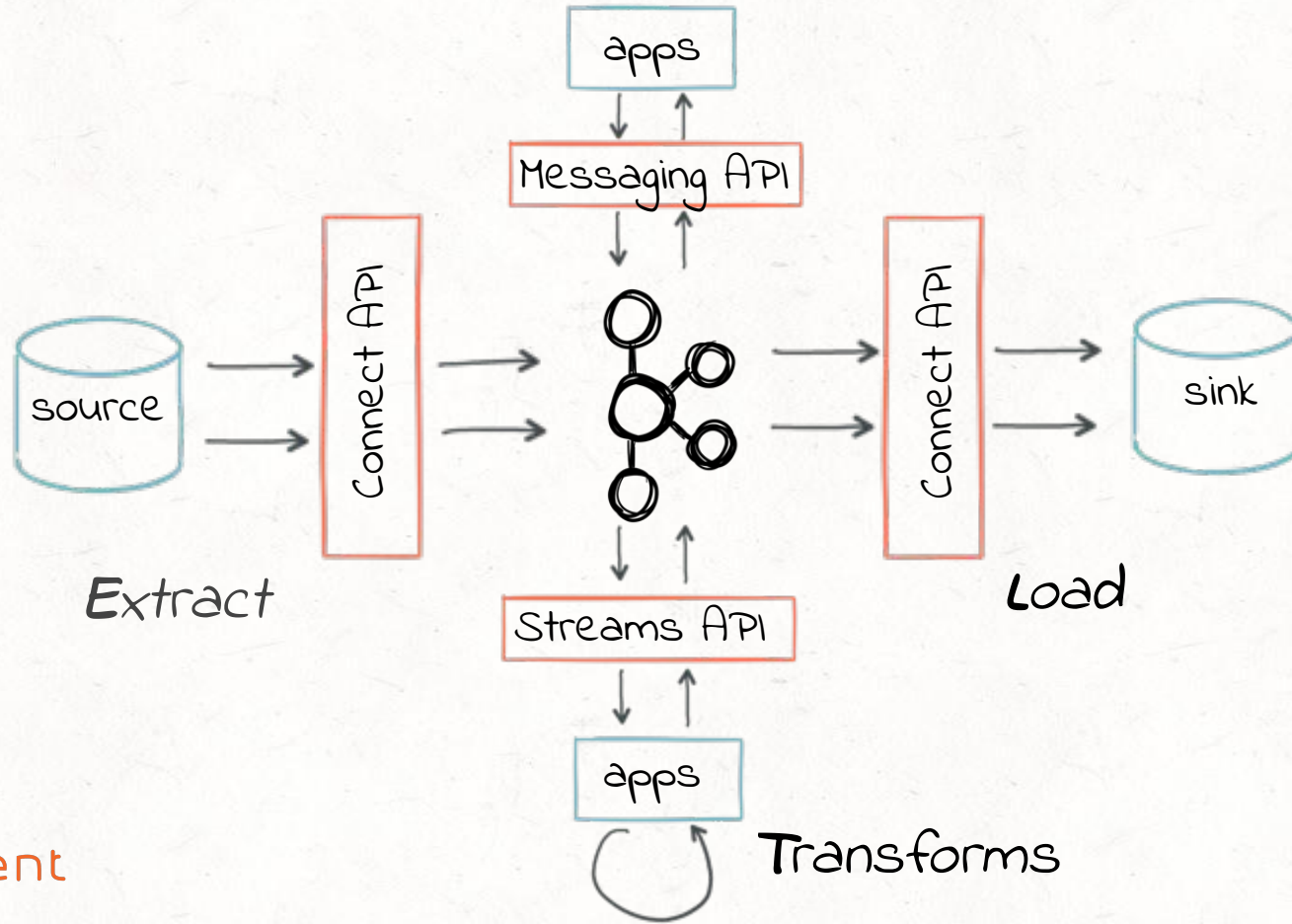
confluent

# Kafka's Connect API: Streaming data ingestion

app

Messaging APIs

Extract

Load

Connect API

source

Connect API

sink

confluent

# Kafka's streams API: stream processing (transforms)



apps

Messaging API

Connect API

source

Extract

Connect API

sink

Load

Streams API

apps

Transforms

confluent

# Kafka's connect API
=
# E and L in Streaming ETL

# Connectors!

Apps    Apps    Apps

NoSQL    Hadoop

RDBMS    DWH

Search    Monitoring    RT analytics

confluent

# How to keep data centers in-sync?

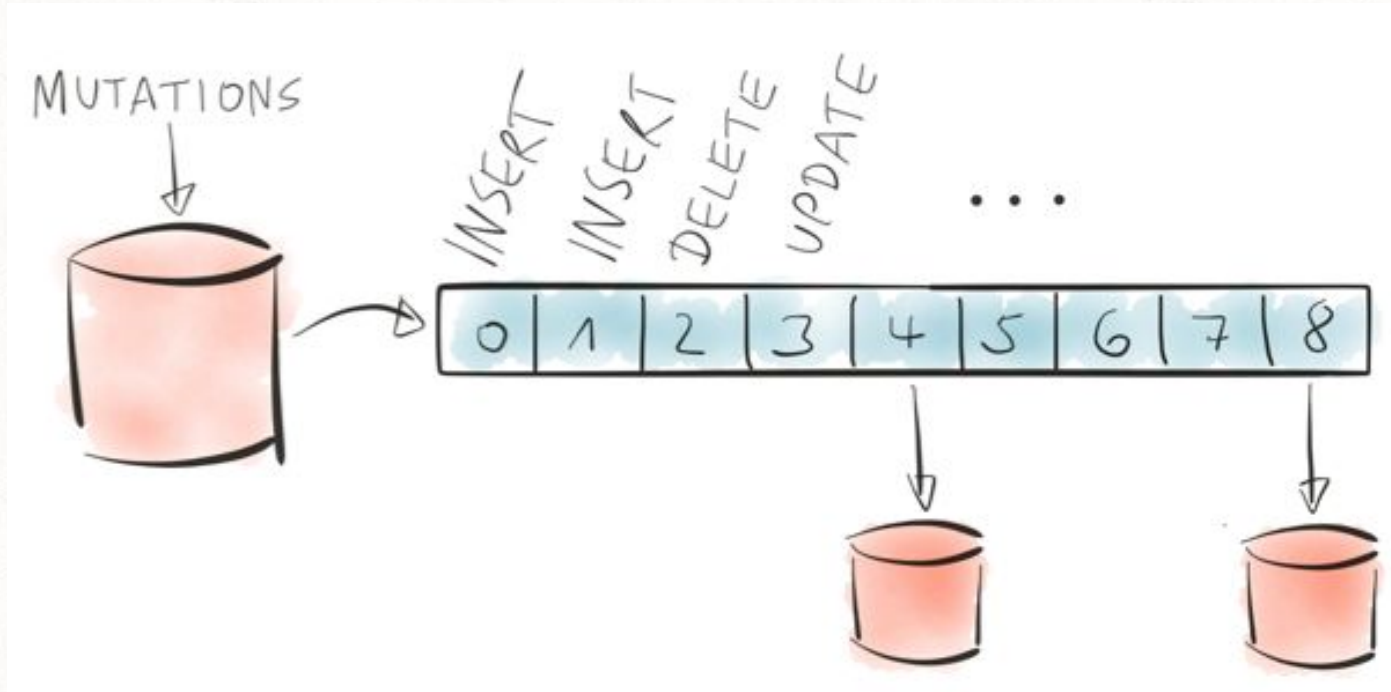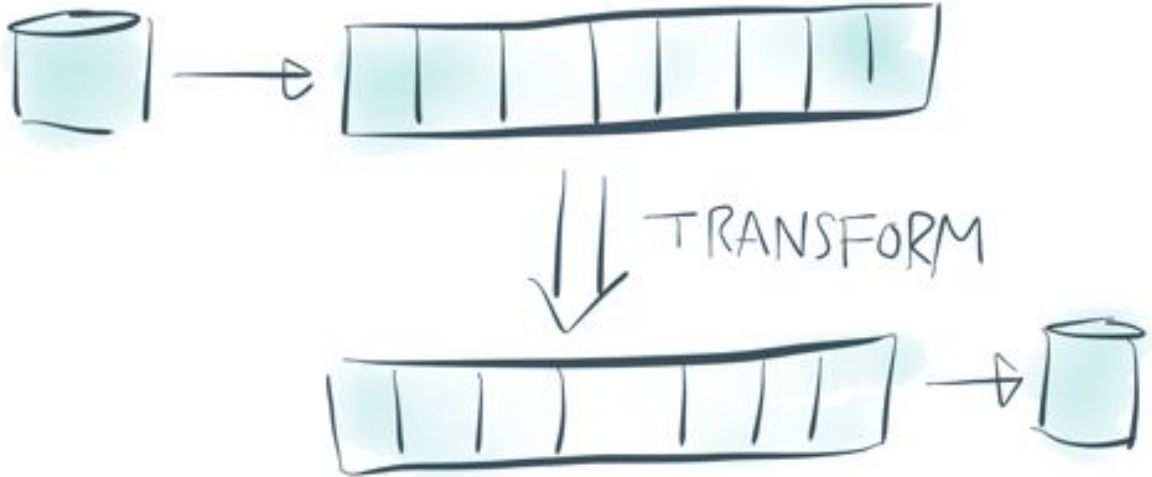

DC1
ON-PREM

DC2
CLOUD

# Sources and sinks

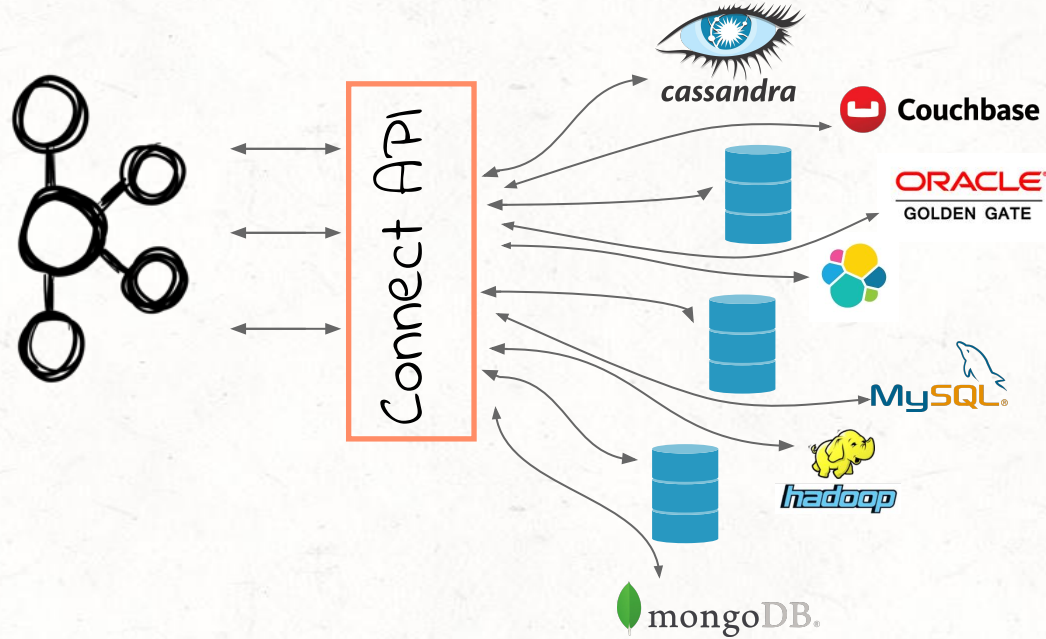# CHANGELOGS

# Transforming changelogs

# Kafka's Connect API = Connectors Made Easy!

- **Scalability**: Leverages Kafka for scalability
- **Fault tolerance**: Builds on Kafka's fault tolerance model
- **Management and monitoring**: One way of monitoring all connectors
- **Schemas**: Offers an option for preserving schemas from source to sink

confluent

# Kafka all the things!



Connect API

cassandra
Couchbase
ORACLE GOLDEN GATE
MySQL
hadoop
mongoDB.

confluent

# Kafka's streams API
# =
# The T in streaming ETL

Stream processing =
**transformations** on stream data

confluent

# 2 VISIONS FOR STREAM PROCESSING

## Real-time Mapreduce VS Event-driven microservices

- Central cluster
- Custom packaging, deployment & monitoring
- Suitable for analytics-type use cases

- Embedded library in any Java app
- Just Kafka and your app
- Makes stream processing accessible to any use case
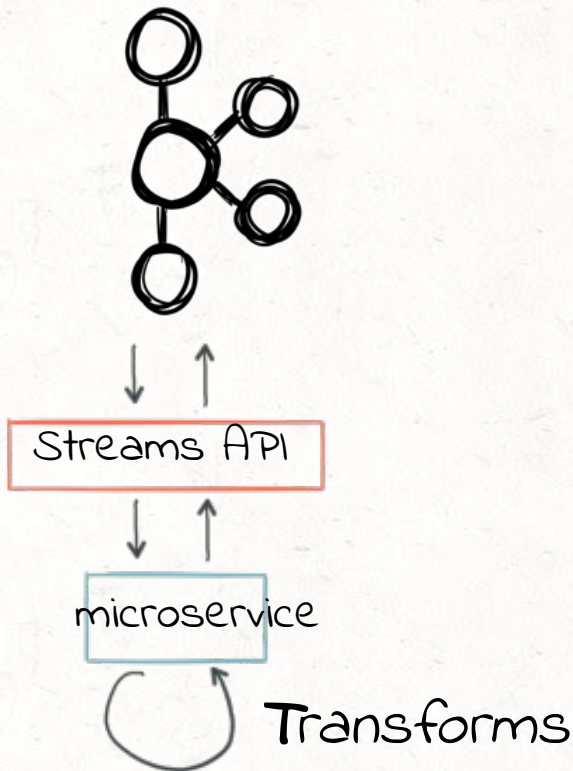
confluent

# Vision 1: real-time mapreduce

# Vision 2: event-driven microservices => Kafka's streams API

"

Kafka's Streams API = Easiest way to do stream processing using Kafka

confluent

**#2: Convenient DSL with all sorts of operators: join(), map(), filter(), windowed aggregates etc**

# Word count program using Kafka's streams API

```java
KStreamBuilder builder = new KStreamBuilder();
KStream<String, String> textLines = builder.stream(stringDeserializer, stringDeserializer, "TextLinesTopic");

KStream<String, Long> wordCounts = textLines
    .flatMapValues(value -> Arrays.asList(value.toLowerCase().split("\\W+")))
    .map((key, value) -> new KeyValue<>(value, value))
    .countByKey(stringSerializer, longSerializer, stringDeserializer, longDeserializer, "Counts")
    .toStream();
wordCounts.to("WordsWithCountsTopic", stringSerializer, longSerializer);

KafkaStreams streams = new KafkaStreams(builder, config);
streams.start();
```
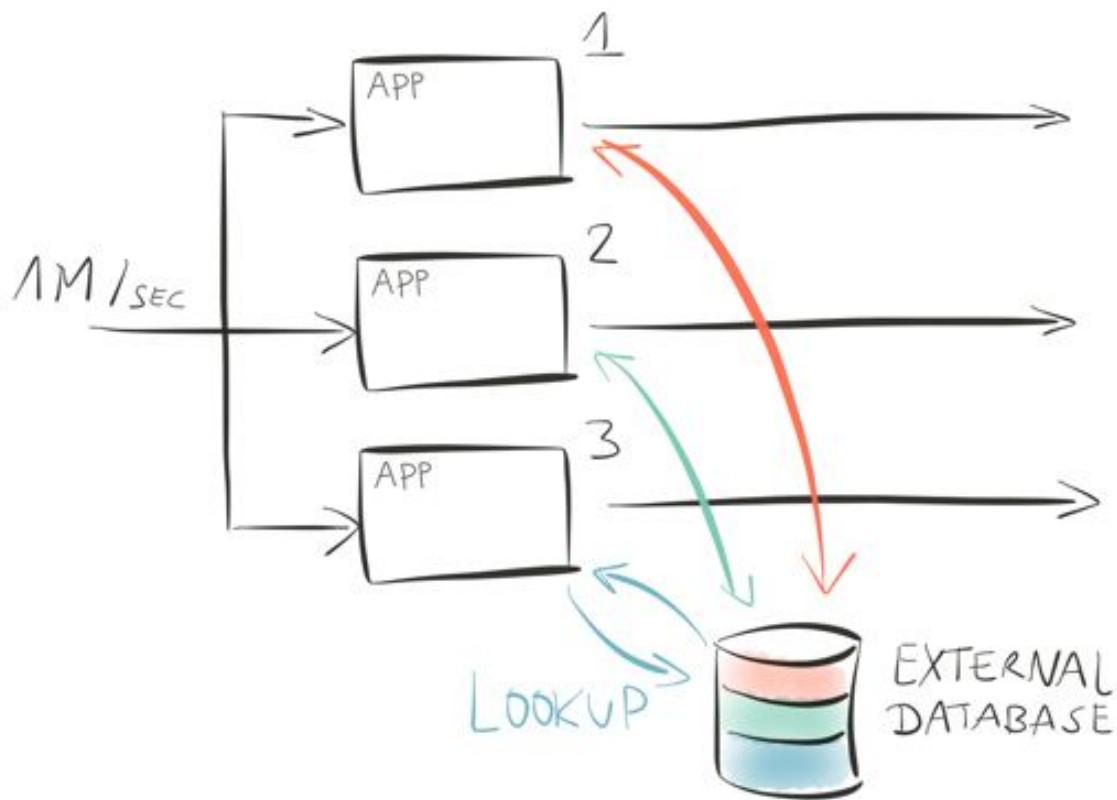
"

#4: Dataflow-style windowing based on event-time; handles late-arriving data
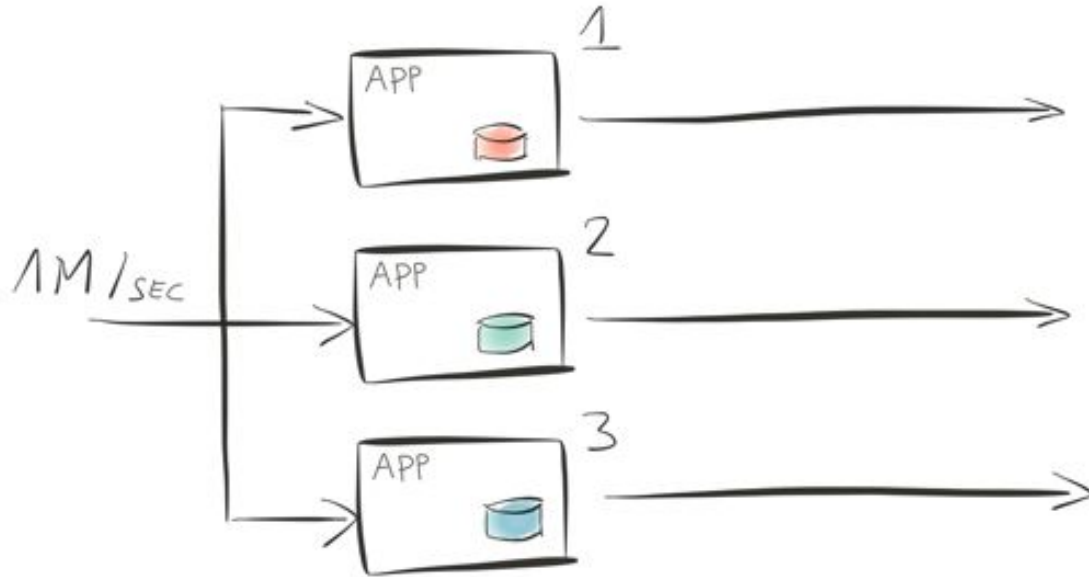
"

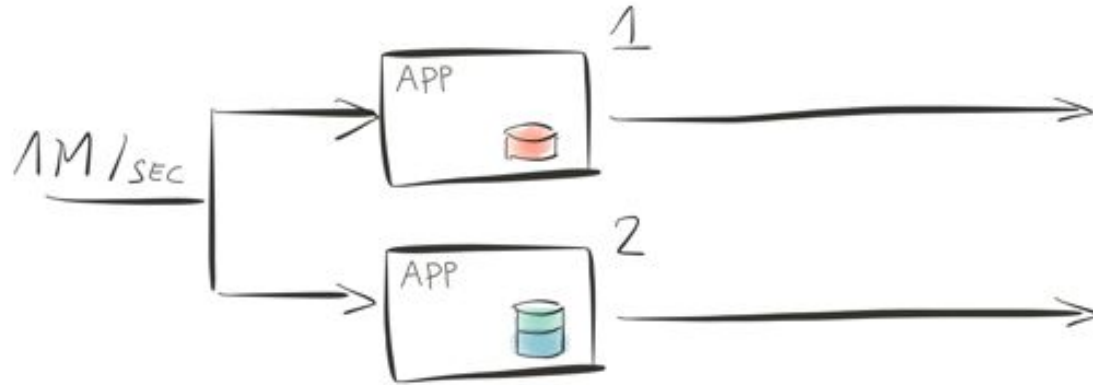#5: out-of-the-box support for local state; supports fast stateful processing

# External state

# LOCAL STATE

# Fault-tolerant local state

**"**

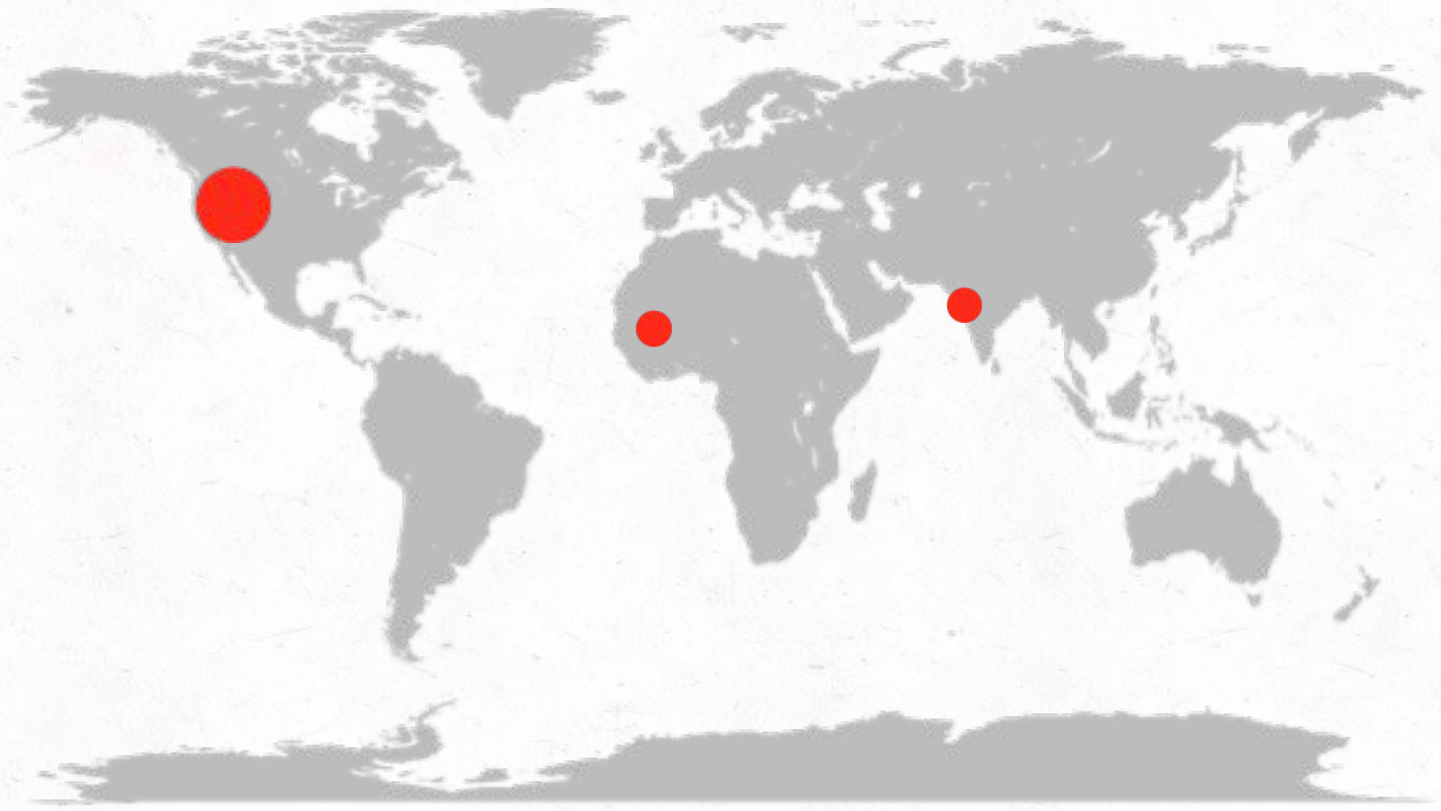#6: Kafka's Streams API allows reprocessing; useful to upgrade apps or do A/B testing
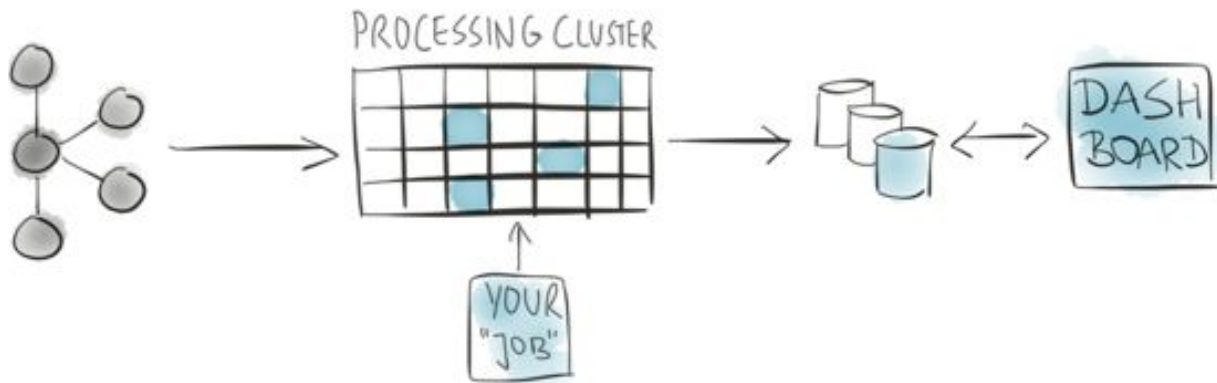
confluent

# REPROCESSING

Real-time dashboard for security monitoring

# Kafka's streams api: simple is beautiful

**Vision 1**



PROCESSING CLUSTER

YOUR "JOB"

DASH BOARD

**Vision 2**



DASHBOARD
KAFKA
STREAMS

confluent

# Logs unify batch and stream processing

ALL YOUR DATA ... EVERYWHERE ... NOW

request-response

changelogs

messaging OR stream processing

Streaming platform

App

App

App

App

DWH

Hadoop

streaming data pipelines

confluent

# Thank you!

@nehanarkhede

confluent