# Embrace the Anarchy :

## Apache Kafka's Role in Modern Data Architectures

Robin Moffatt / Confluent

# $ **whoami**

- Developer Advocate @ Confluent

- Working in data & analytics since 2001

- Oracle Developer Champion 🏆

- Blogging : http://rmoff.net & http://cnfl.io/rmoff

- Twitter: **@rmoff**

  - Geek stuff

  - Beer & Fried Breakfasts

# https://speakerdeck.com/rmoff/



STREAM PROCESSING

We ❤ syslogs: Real-time syslog Processing with Apache Kafka and KSQL—Part 1: Filtering

Robin Moffatt 🍺🏃🥓
@rmoff

#FullEnglish ift.tt/2v7cLfE

# "Apache Kafka is a Streaming Platform

# "Why do we need a streaming platform?

" one of the reasons:

Decoupling
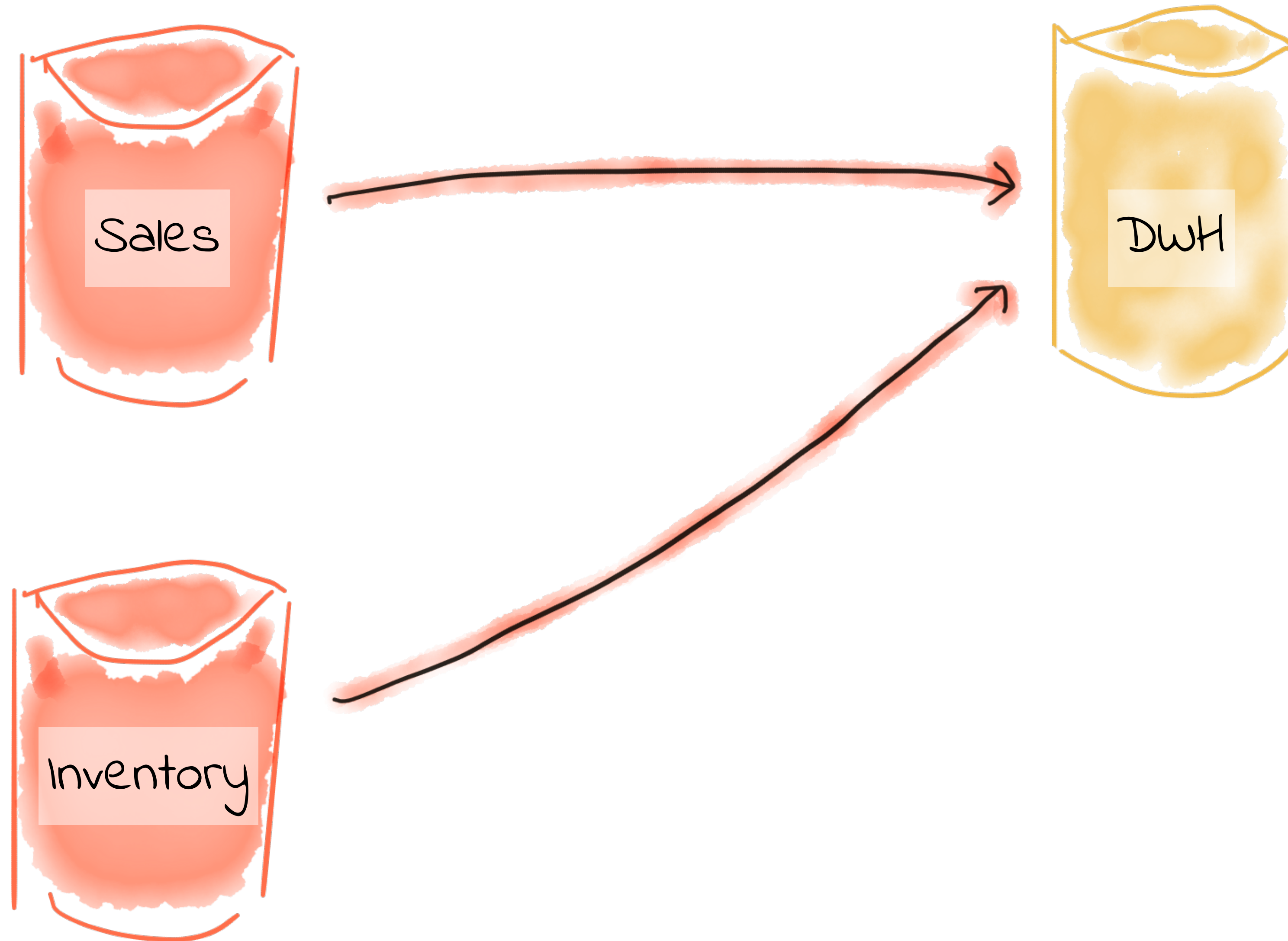
confluent

"

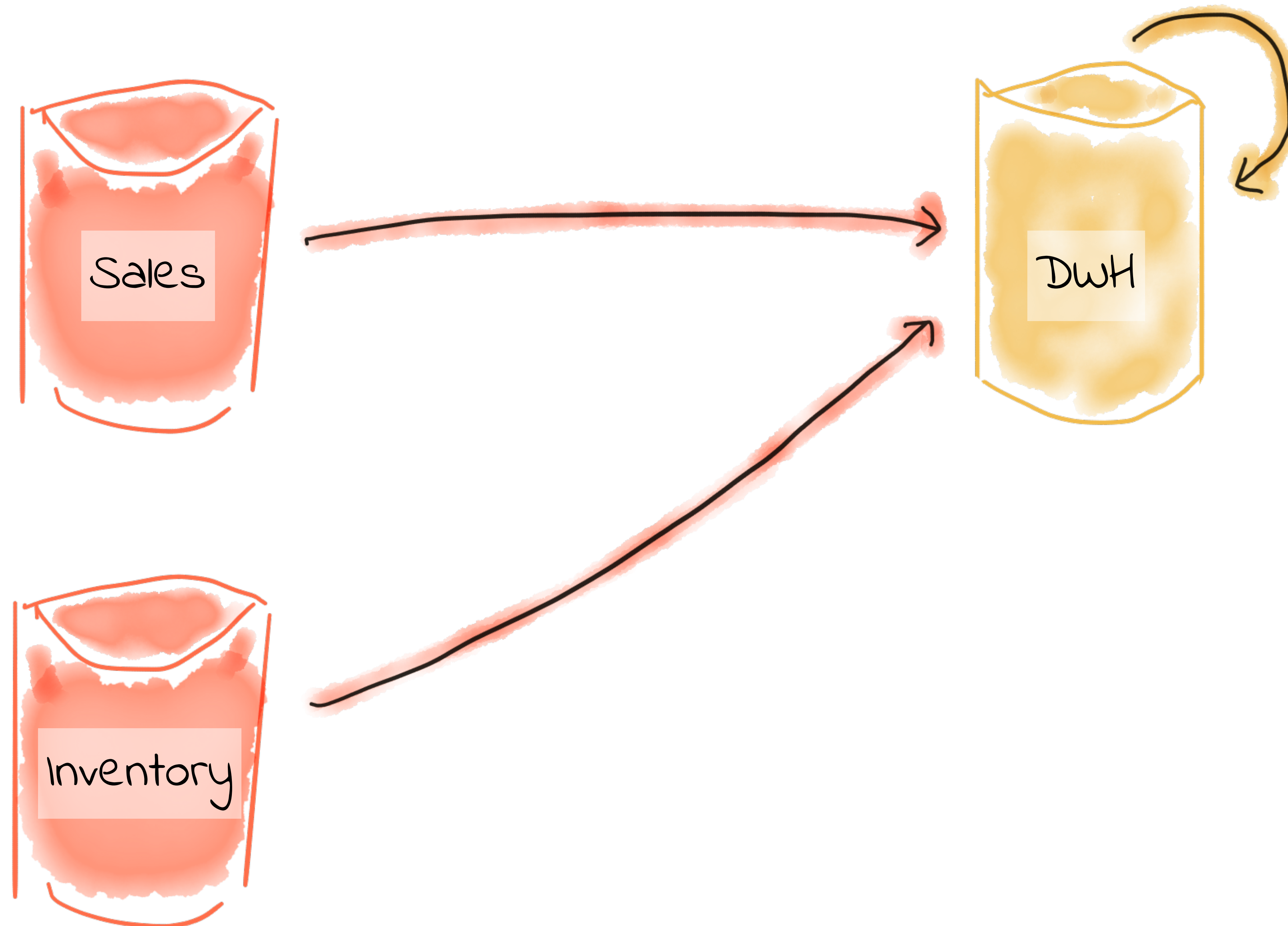# A case in point...Analytics

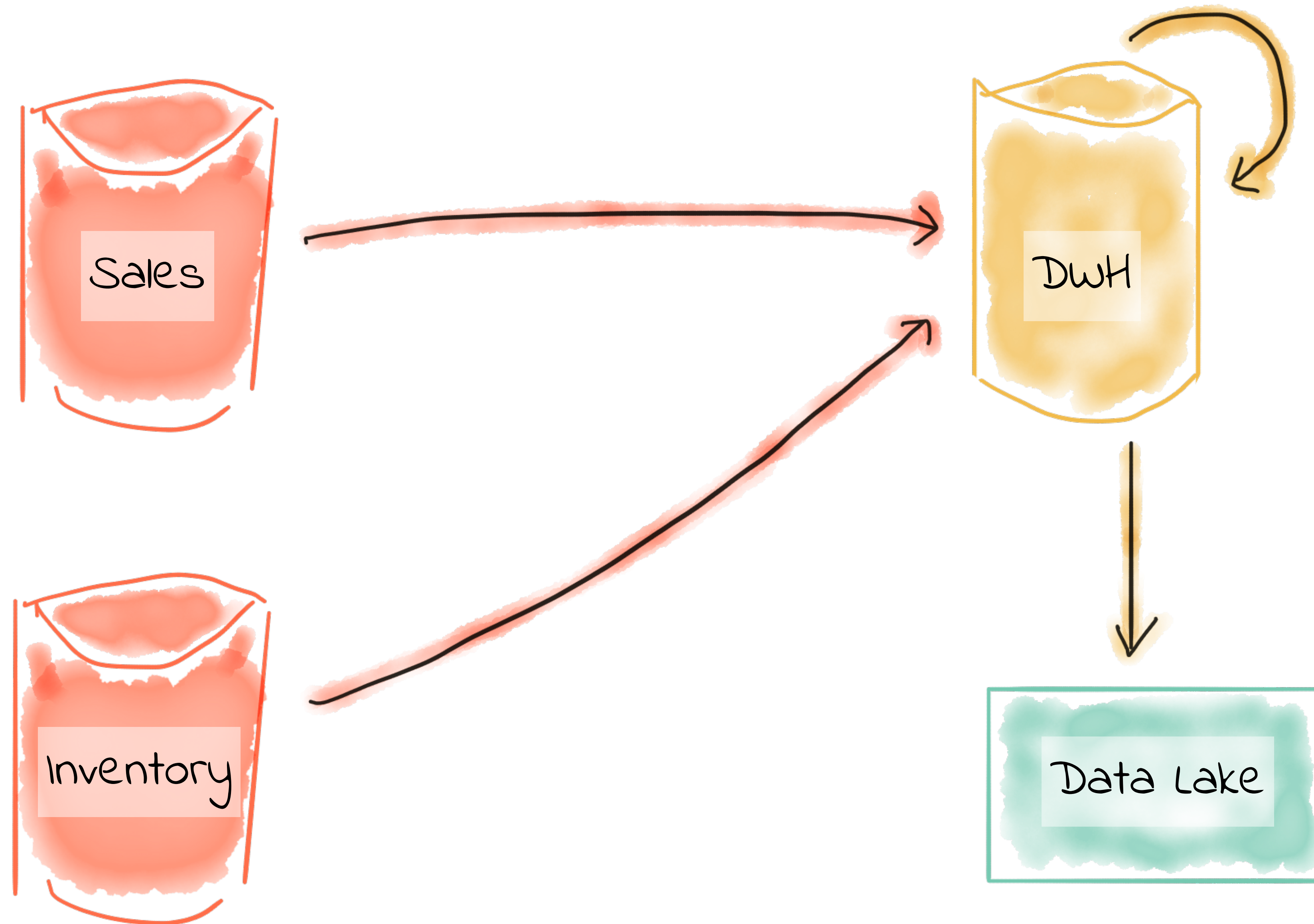# Analytics—In the beginning...

# And then there were more data sources...

# Batch Transformations ... (ETL / ELT)

# Add a Data Lake…

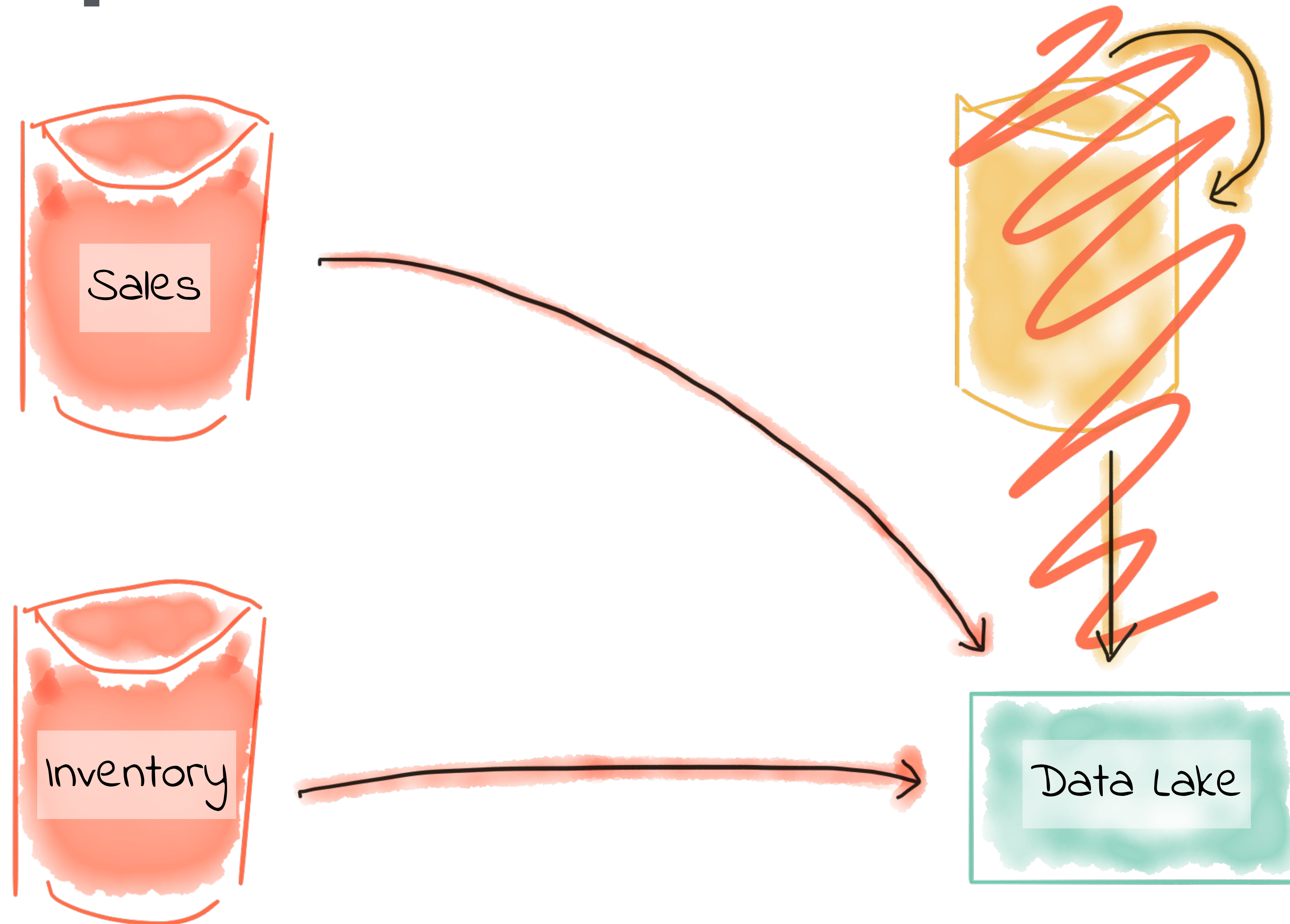# ...or Replace the Data Warehouse



Sales

Inventory

Data Lake

# Still need to do Batch transformations...

Sales

Inventory

Data Lake

**Want your data anytime → SOON ?**

**Batch is Latency built in by Design**

# The World has Changed

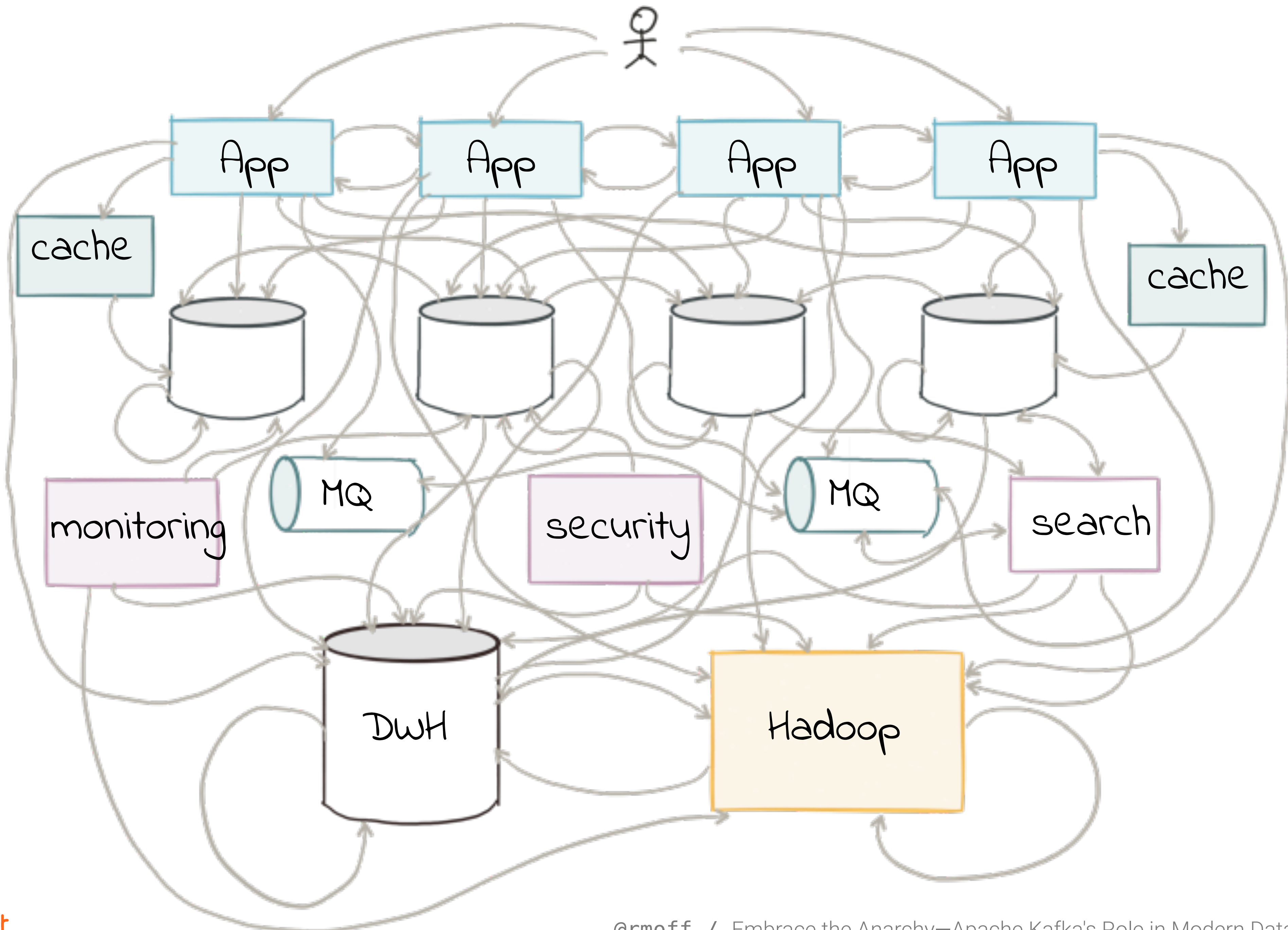Internet of Things

Microservices

Mobile

Machine Learning

Lots of new technologies

(whether you like it or not)

request-response

changelogs

messaging
OR
stream
processing

App

App

KAFKA

App

App

DWH
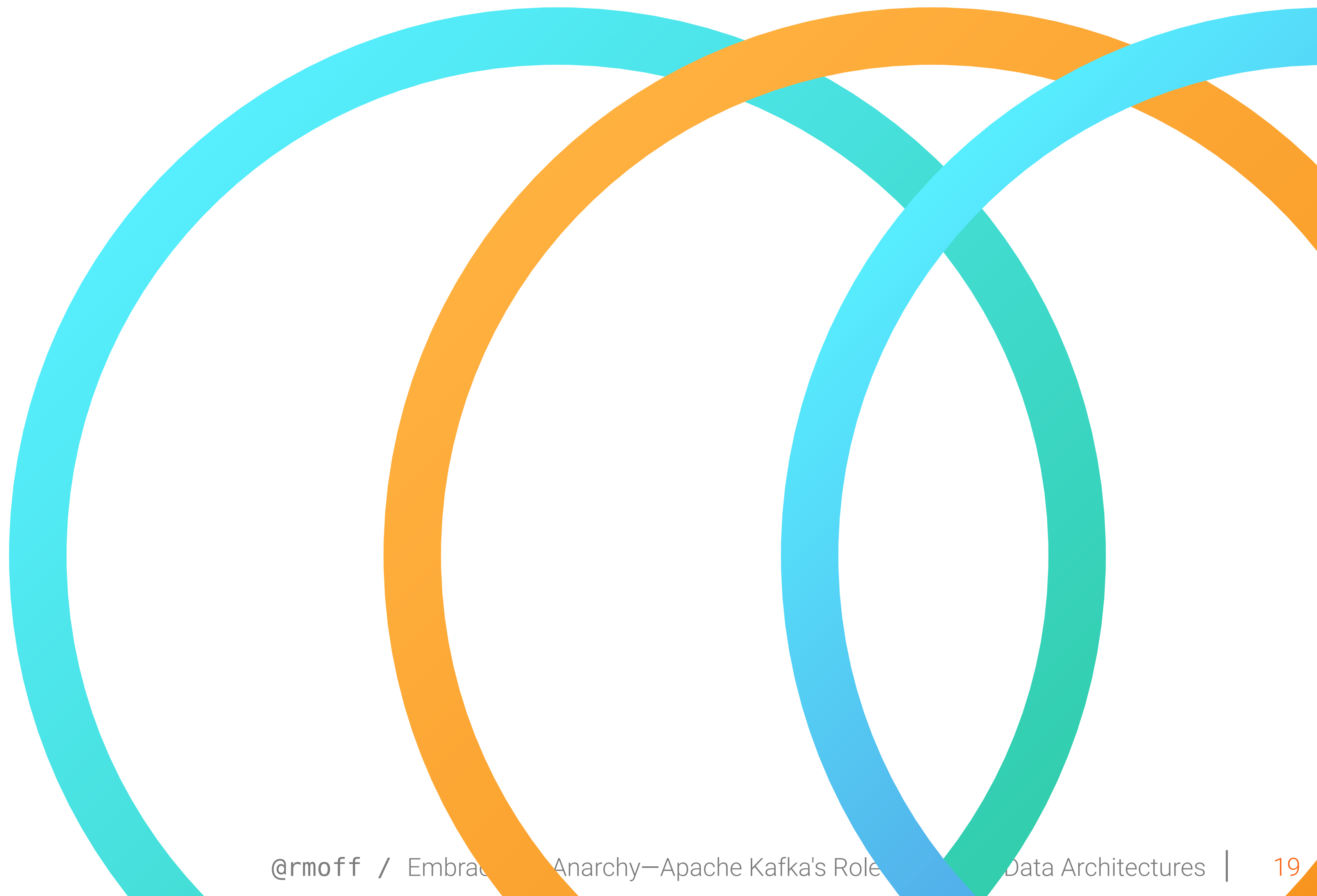
Hadoop

streaming data pipelines

# " Apache Kafka is a Streaming Platform

# Three Lenses

# What is Apache Kafka?

## 01
Messaging
Done Right

## 02
Scalable Streaming
Data Pipelines

## 03
Foundation for
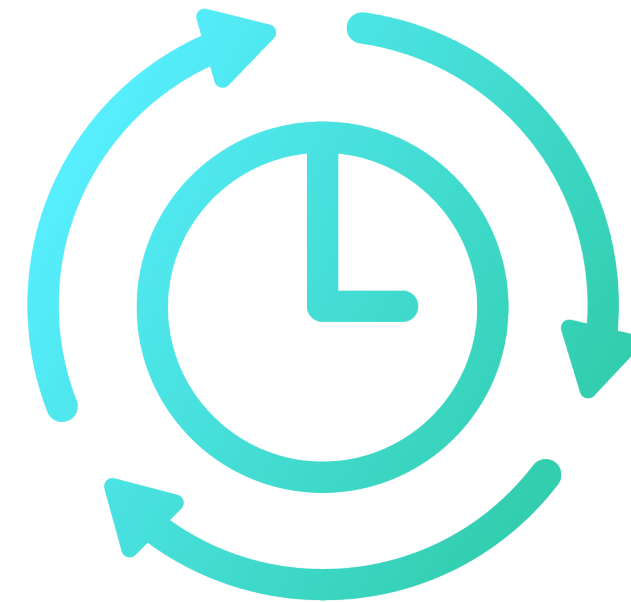Stream Processing
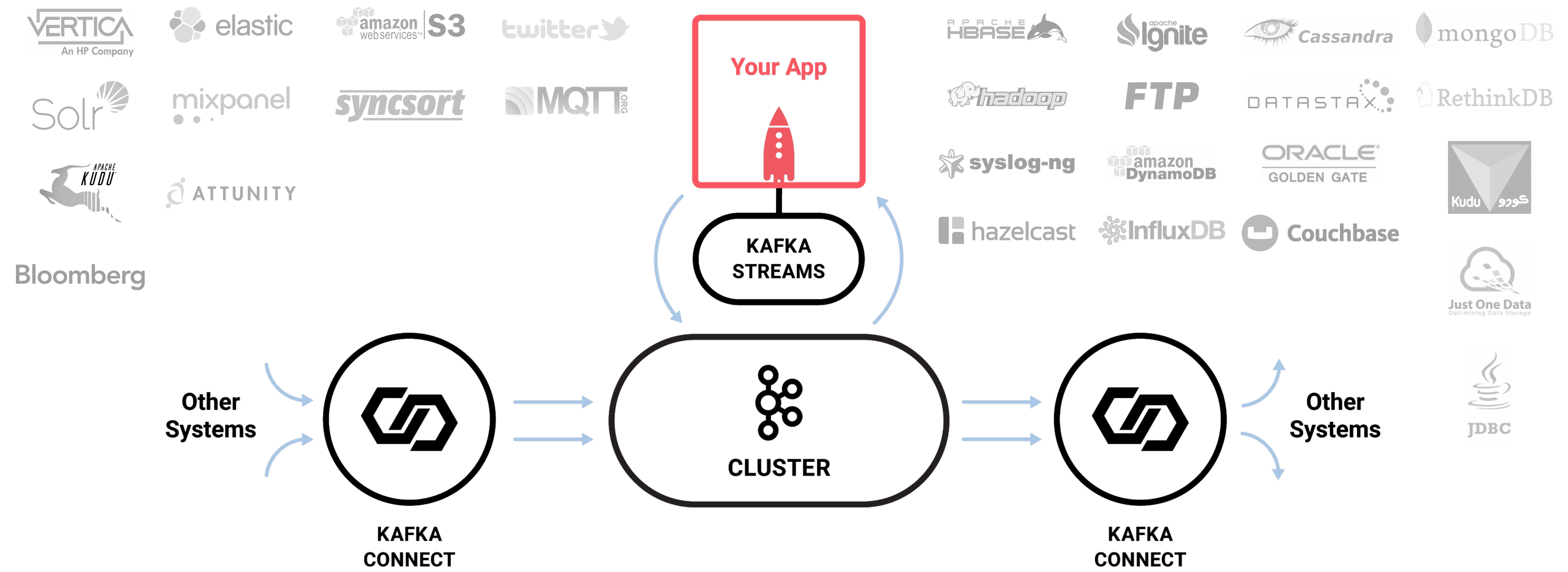
# Lens 1: Messaging Done Right

Scalability

True Storage

Real-Time
Processing

# Lens 2: Scalable Streaming Data Pipelines



VERTICA — An HP Company
elastic
amazon web services™ S3
twitter
Your App
KAFKA STREAMS
APACHE HBASE
apache Ignite
Cassandra
mongoDB

Solr
mixpanel
syncsort
MQTT.ORG
hadoop
FTP
DATASTAX
RethinkDB

APACHE KUDU
ATTUNITY
syslog-ng
amazon DynamoDB
ORACLE GOLDEN GATE
Kudu

Bloomberg
hazelcast
InfluxDB
Couchbase
Just One Data — Optimizing Data Storage

Other Systems
KAFKA CONNECT
CLUSTER
KAFKA CONNECT
Other Systems
JDBC

# Lens 3: Foundation for Stream Processing

# KSQL
**is the**
# Streaming
**SQL Engine**
**for**
# Apache Kafka

# The Streaming Platform



Web · Custom Apps · Microservices · Monitoring · Analytics · ...and more

**APACHE KAFKA®**

any source · NoSQL · Oracle · Bloomberg · Twitter · SFDC · Hadoop · Data Warehouse · any sink

# The Streaming Platform

# Event-Driven
# Scalable
# Decoupled

...and more

any source

NoSQL

Oracle

Bloomberg

SFDC

Data Warehouse

any sink

confluent

"Bold claim: all your data
is event streams

# A Customer Experience

# A Sale

# A Sensor Reading

# An Application Log Entry

# Databases

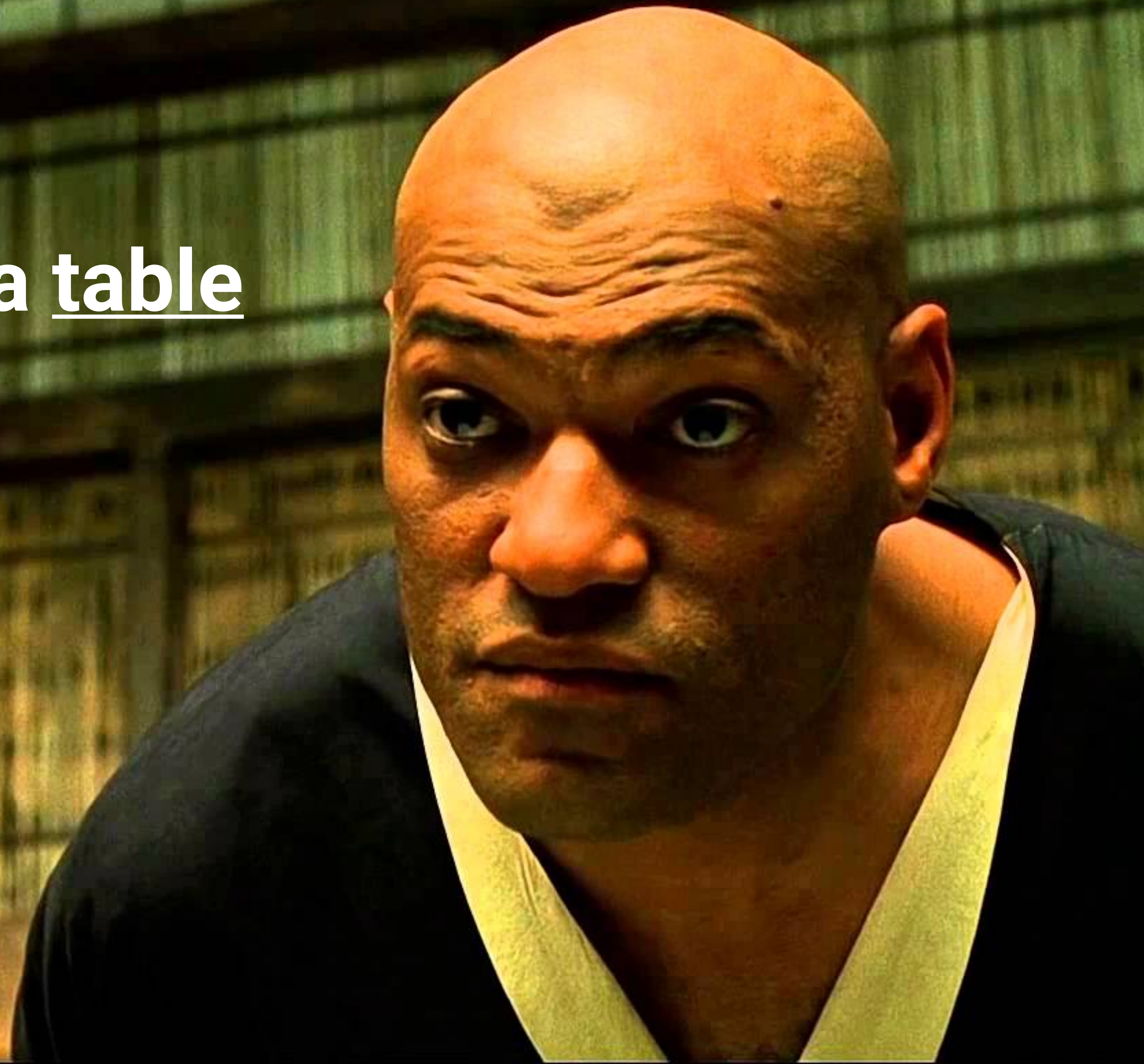Do you think that's a **table**
you are querying?

# The Table Stream Duality

## Stream

| Account ID | Amount |
|------------|--------|
| 12345 | + €50 |
| 12345 | + €25 |
| 12345 | -€60 |

Time →

## Table

| Account ID | Balance |
|------------|---------|
| 12345 | €50 |

| Account ID | Balance |
|------------|---------|
| 12345 | €75 |

| Account ID | Balance |
|------------|---------|
| 12345 | €15 |

The truth is the log.

The database is a cache
of a subset of the log.

—Pat Helland
Immutability Changes Everything
http://cidrdb.org/cidr2015/Papers/CIDR15_Paper16.pdf

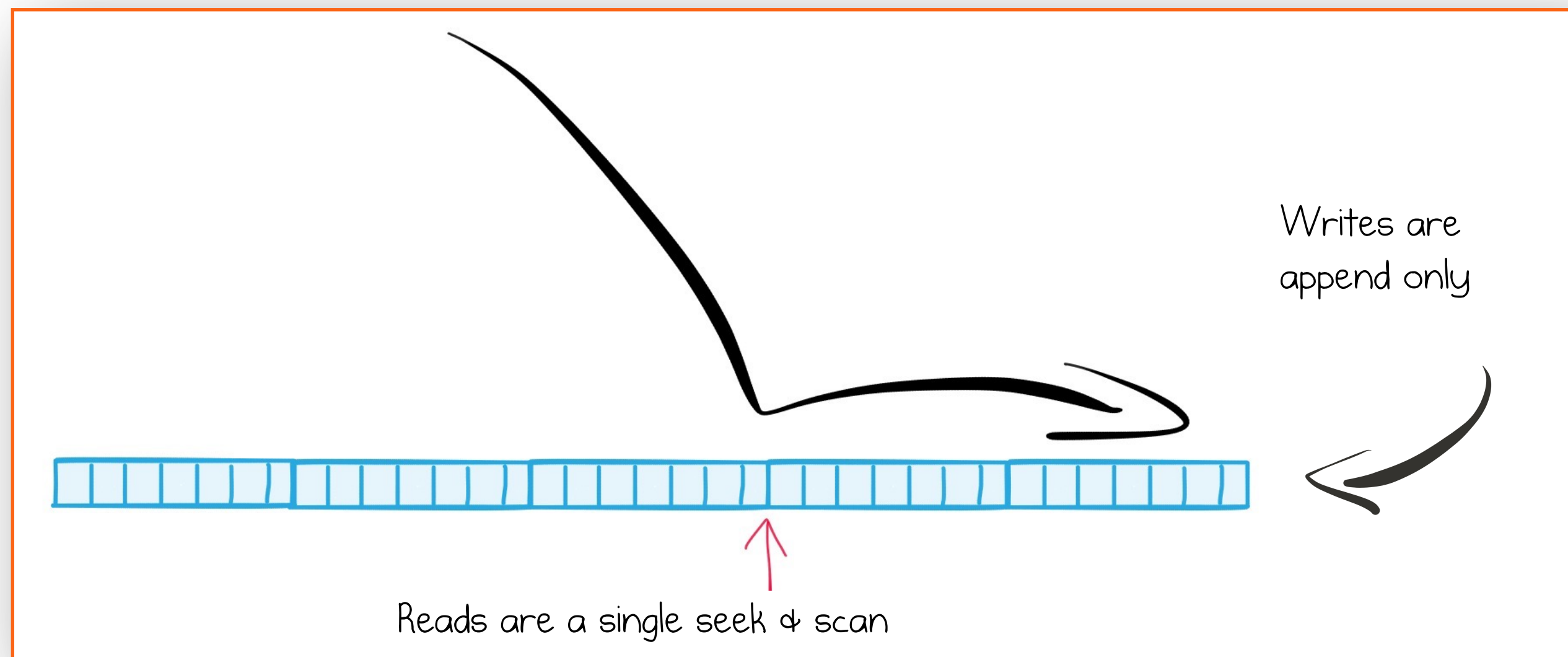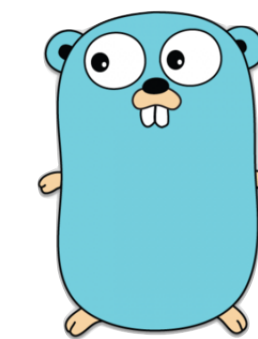# A Brief Look at Kafka's Technology

# Apache Kafka

## Kafka

A **Distributed Commit Log**. Publish and subscribe to streams of records. Highly scalable, high throughput. Supports transactions. Persisted data. Stream processing.

Writes are append only

Reads are a single seek & scan

## Producer & Consumer APIs

Open-source client libraries for numerous languages, to directly integrate with your applications.
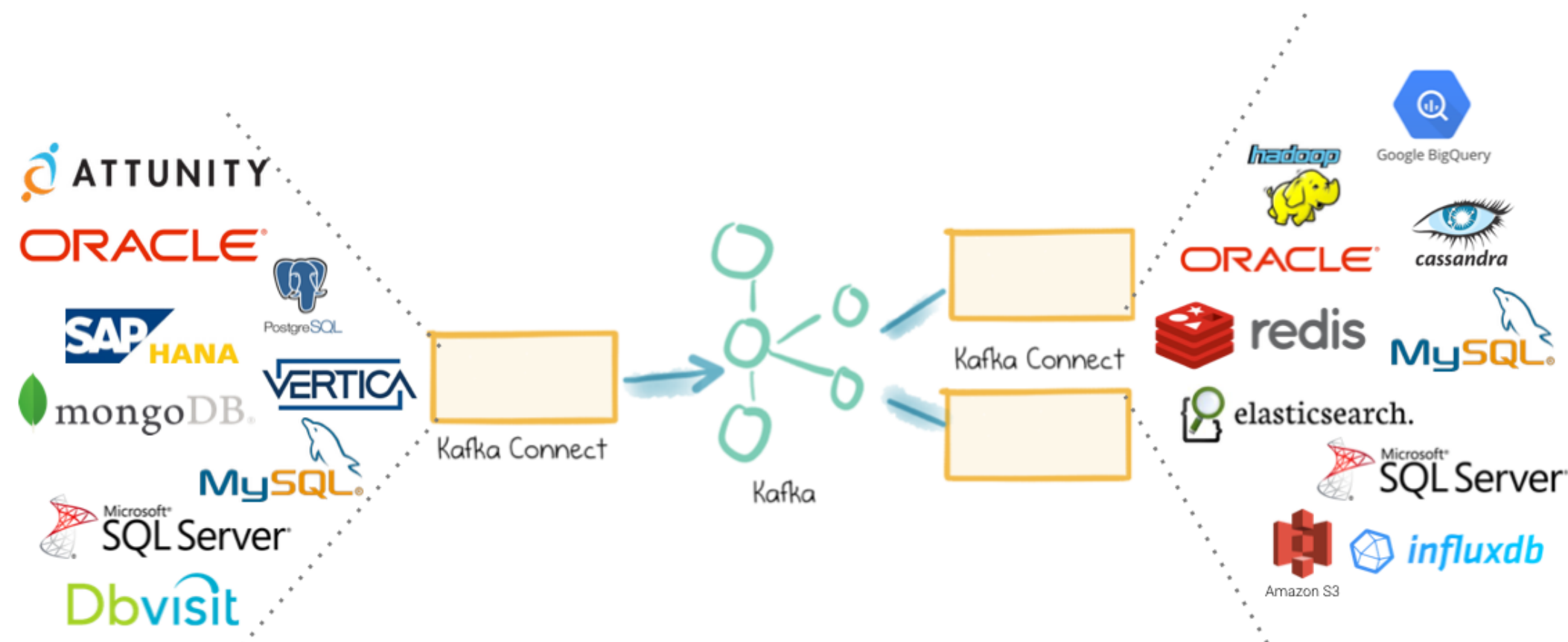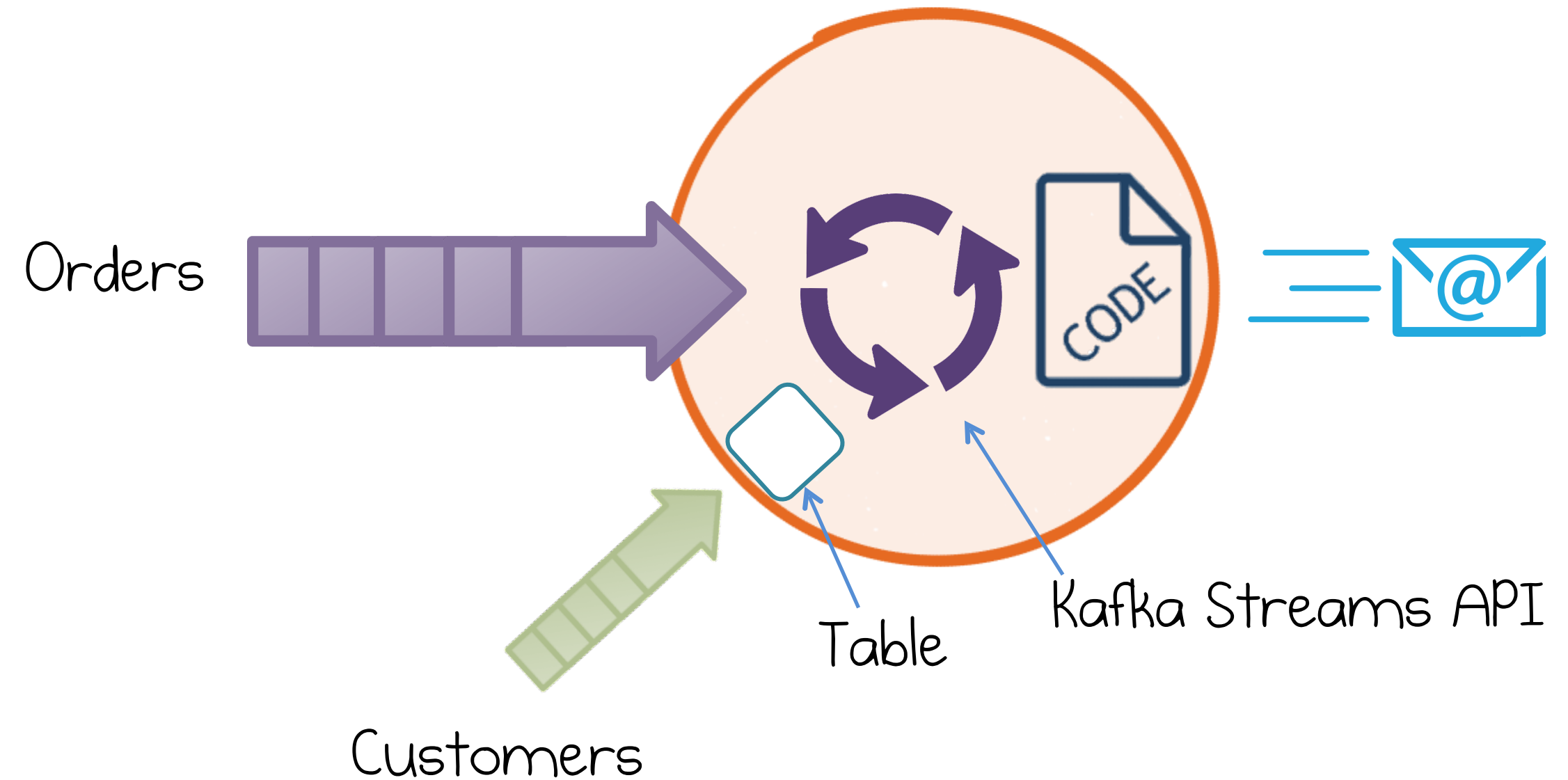
# Apache Kafka

## Kafka Connect API

Reliable and scalable integration of Kafka with other systems – no coding required.



## Kafka Streams API

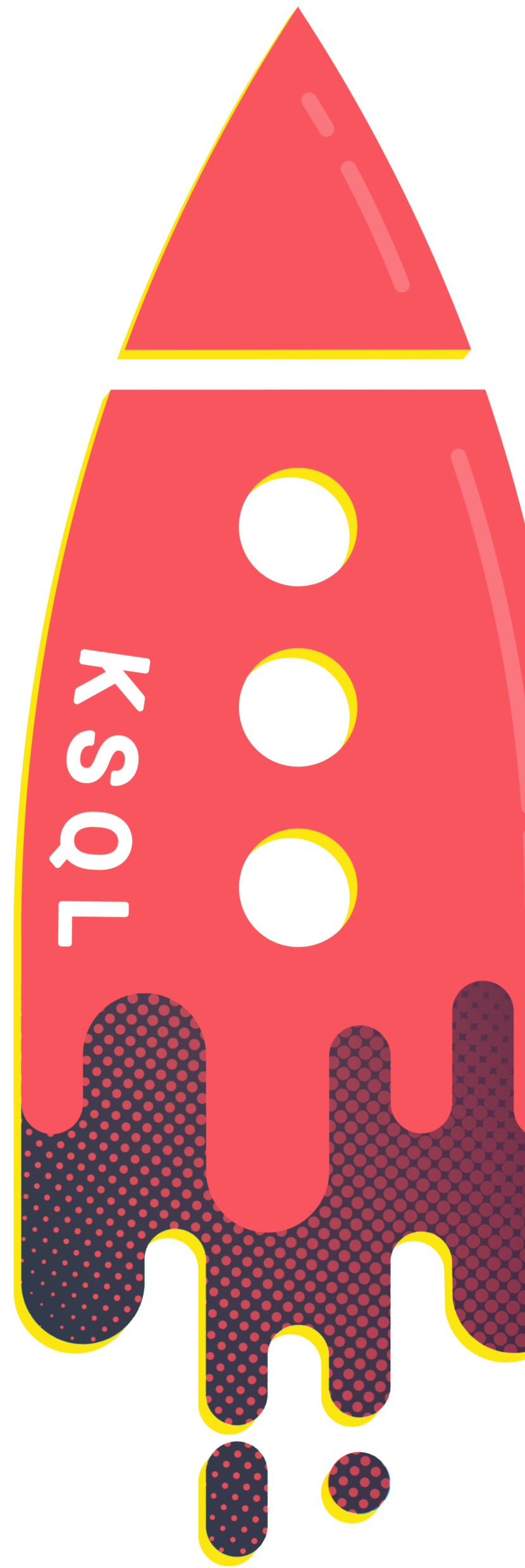Write standard Java applications & microservices to process your data in real-time



Orders

Table

Kafka Streams API

Customers

# KSQL
## is a
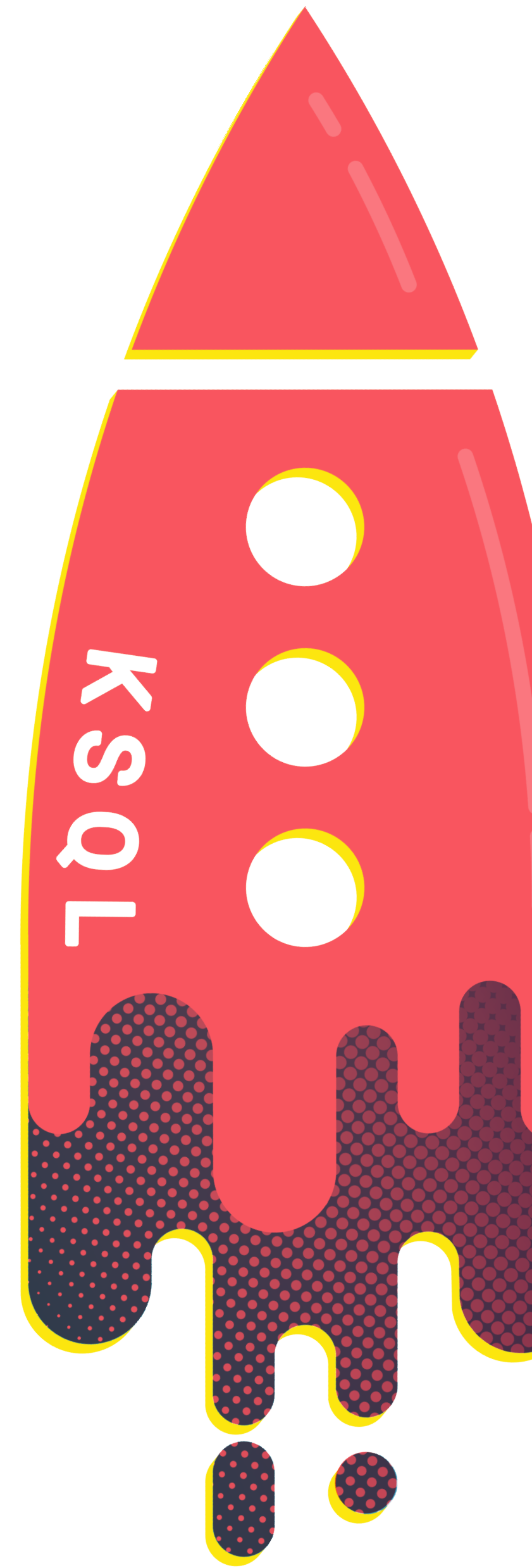## Declarative
## Stream Processing Language

# KSQL
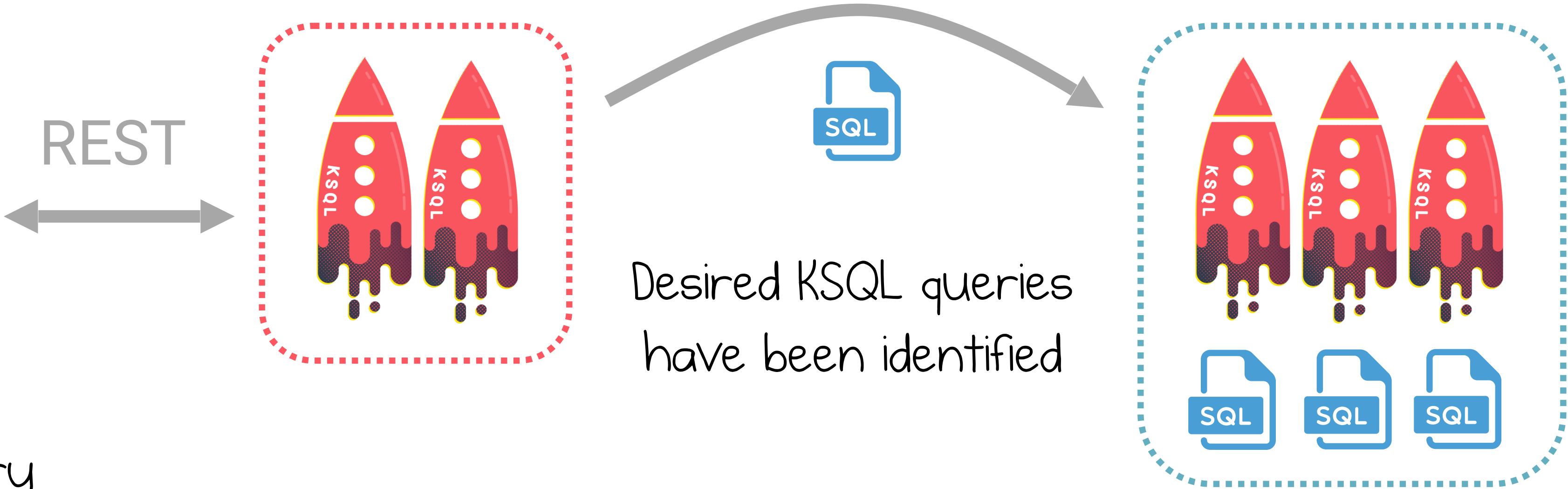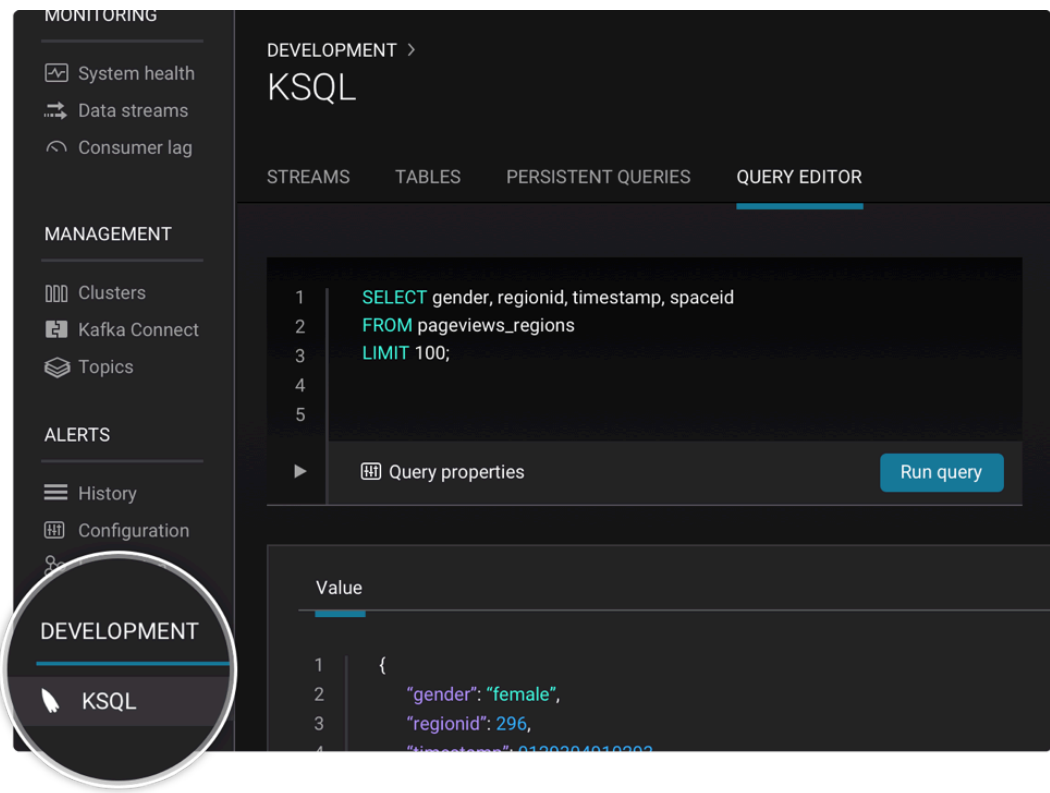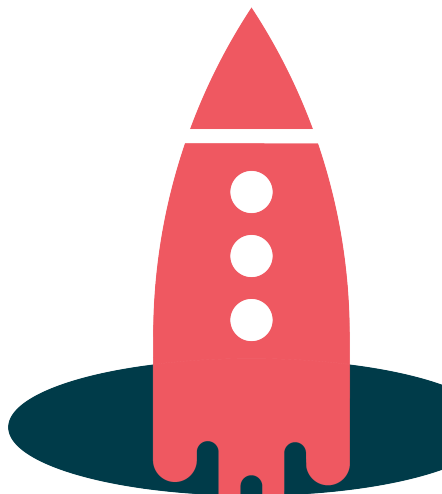## is the
## Streaming
## SQL Engine
## for
## Apache Kafka

KSQL

# KSQL in Development and Production

**Interactive KSQL**
for development and testing

**Headless KSQL**
for Production



REST

Desired KSQL queries
have been identified
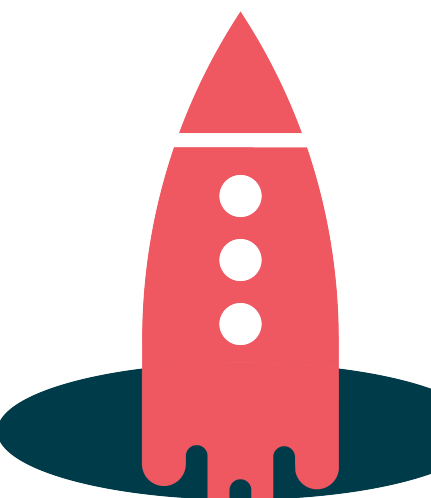
"Hmm, let me try
out this idea..."

# KSQL for Real-Time Monitoring

- Log data monitoring, tracking and alerting
- syslog data
- Sensor / IoT data

```
CREATE STREAM SYSLOG_INVALID_USERS AS
SELECT HOST, MESSAGE
FROM SYSLOG
WHERE MESSAGE LIKE '%Invalid user%';
```
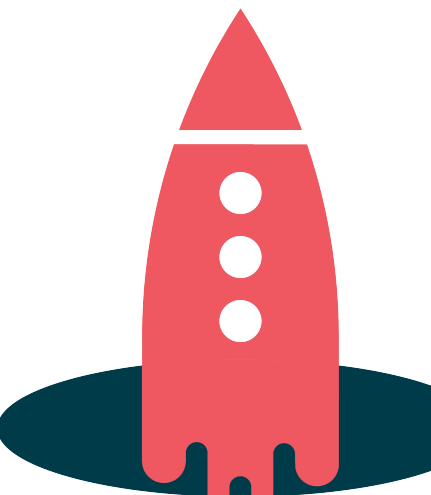
**http://cnfl.io/syslogs-filtering** / **http://cnfl.io/syslog-alerting**

# KSQL for Anomaly Detection

**Identifying patterns or anomalies in real-time data, surfaced in milliseconds**

```
CREATE TABLE possible_fraud AS
  SELECT card_number, count(*)
    FROM authorization_attempts
    WINDOW TUMBLING (SIZE 5 SECONDS)
    GROUP BY card_number
    HAVING count(*) > 3;
```
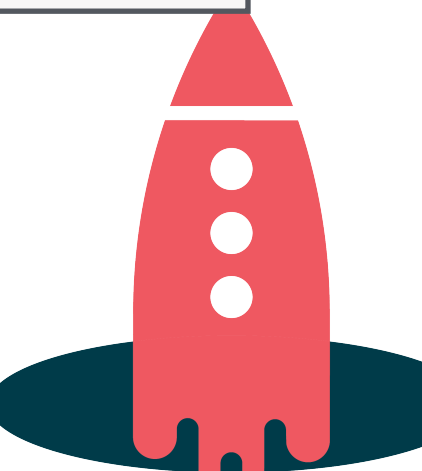
# KSQL for Streaming ETL

**Joining, filtering, and aggregating streams of event data**

```
CREATE STREAM vip_actions AS
  SELECT userid, page, action
  FROM clickstream c
  LEFT JOIN users u
    ON c.userid = u.user_id
  WHERE u.level = 'Platinum';
```
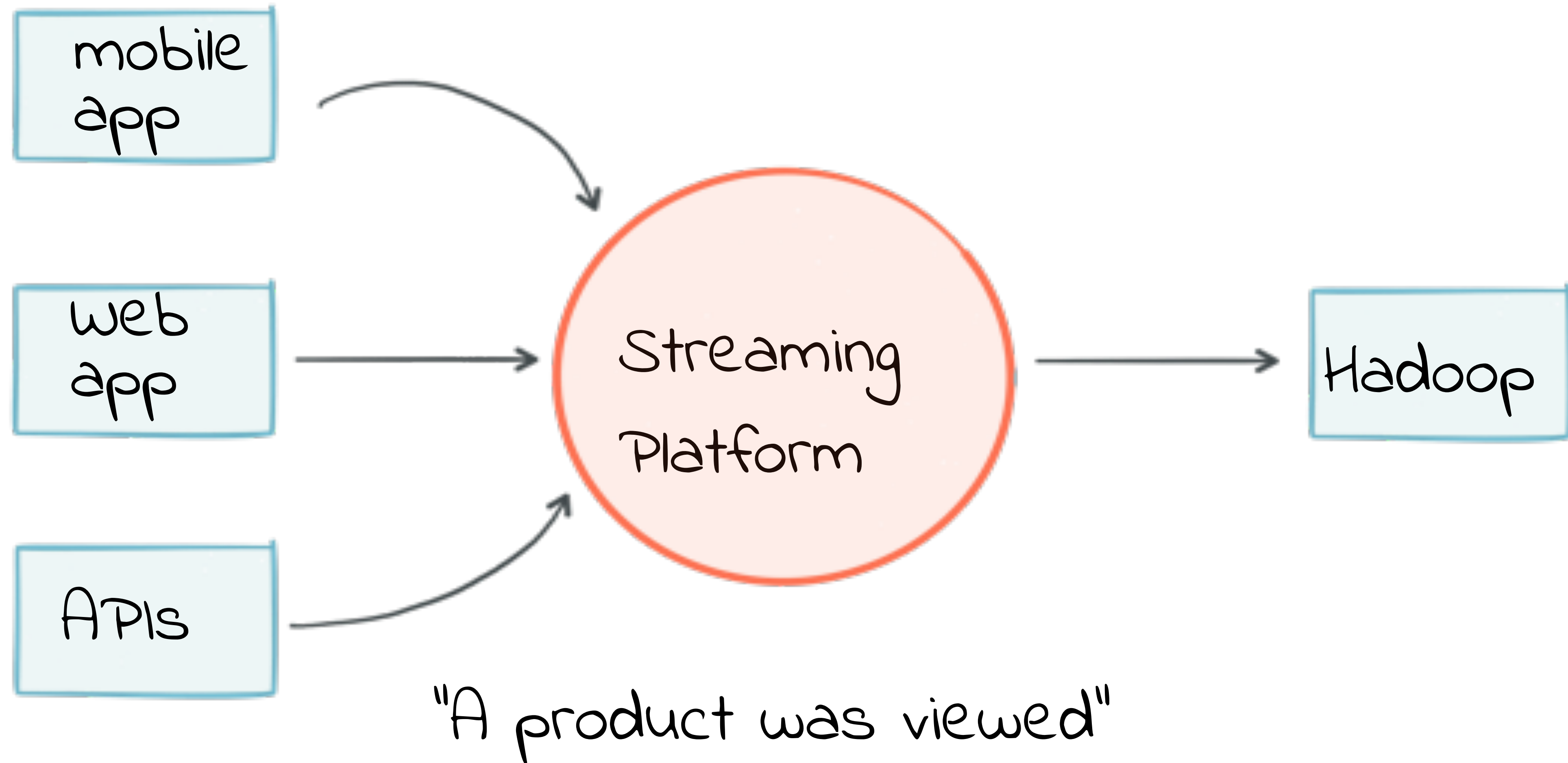
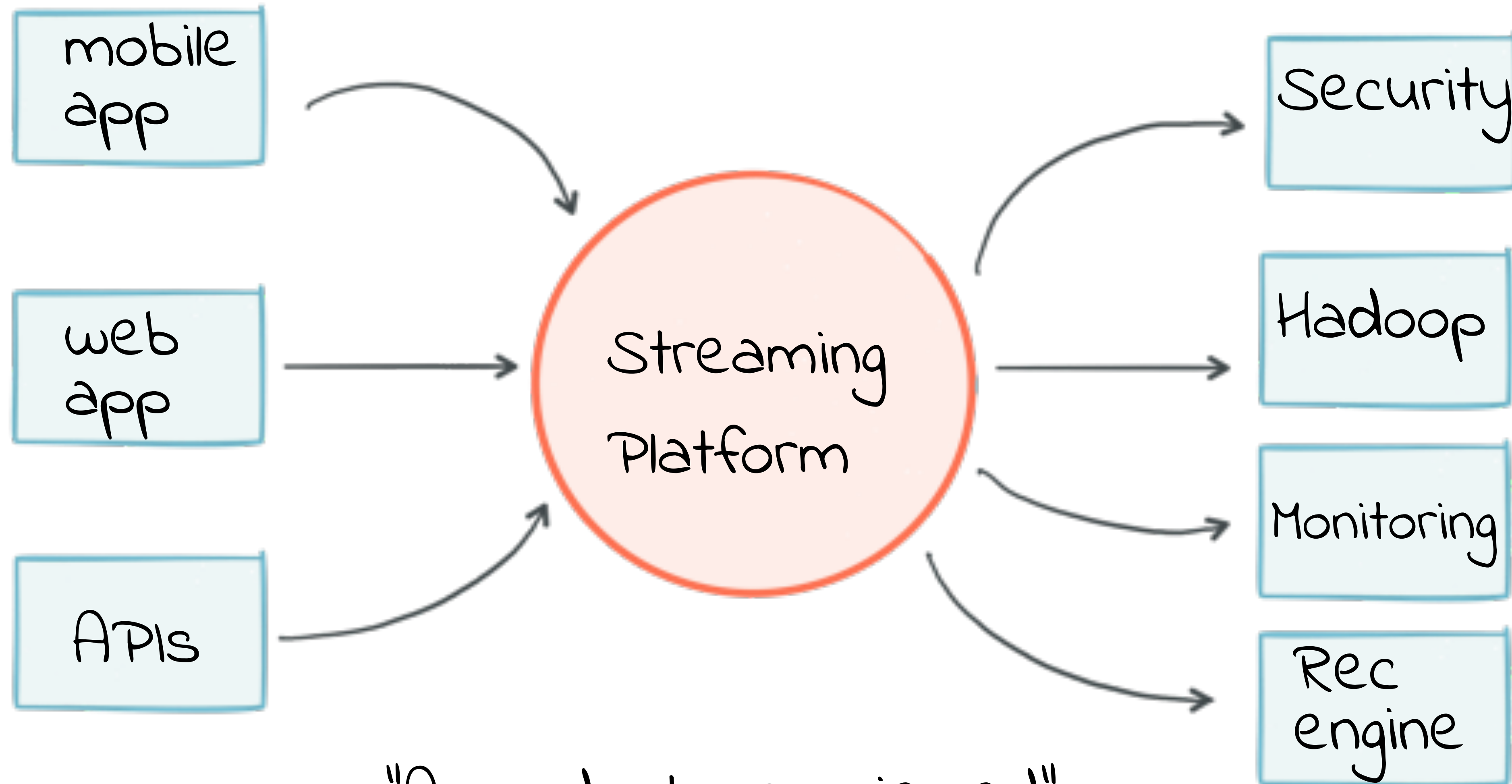# " What Problems does Kafka Solve?

# Event-Centric Thinking

```
┌──────────┐                    ╱‾‾‾‾‾‾‾‾‾‾╲                    ┌──────────┐
│  Web     │         ──────▶   │ Streaming │   ──────▶         │  Hadoop  │
│  app     │                   │ Platform  │                   │          │
└──────────┘                    ╲_____╱                    └──────────┘
```

"A product was viewed"

# Event-Centric Thinking



mobile app

web app

APIs

Streaming Platform

Hadoop

"A product was viewed"

# Event-Centric Thinking



mobile app

web app

APIs

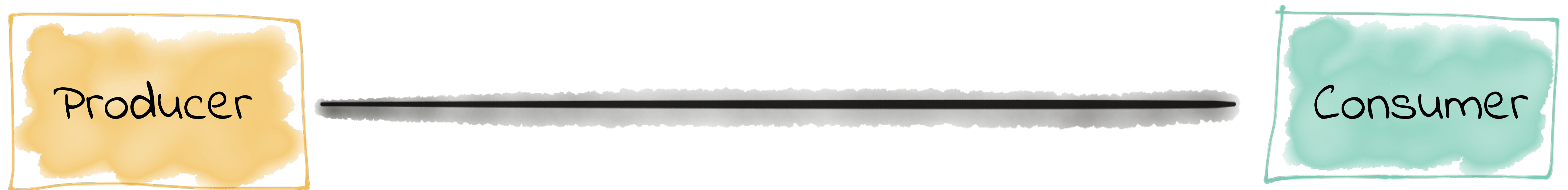Streaming Platform

Security

Hadoop

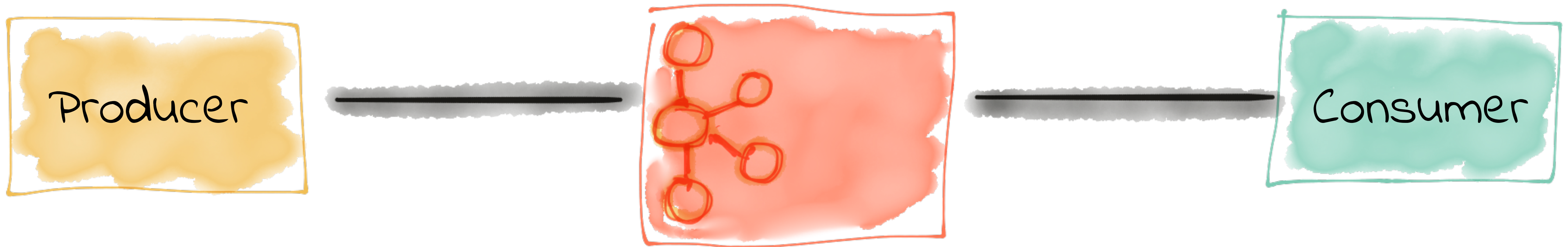Monitoring

Rec engine

"A product was viewed"

# System Availability and Event Buffering

# System Availability and Event Buffering

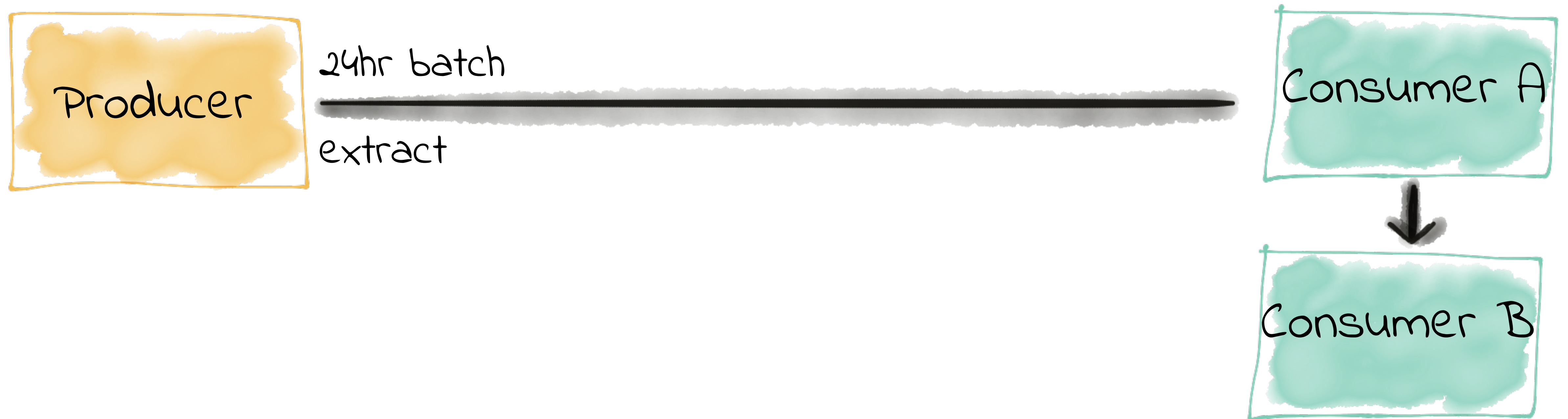Producer

Consumer

# Varying Latency Requirements / Batch vs Stream

Producer

24hr batch

extract

Consumer A

# Varying Latency Requirements / Batch vs Stream



Producer

24hr batch

extract

Consumer A

Consumer B

# Varying Latency Requirements / Batch vs Stream

Producer

24hr batch

extract

Consumer A

Consumer B

# Varying Latency Requirements / Batch vs Stream



Producer

24hr batch extract

Realtime

Consumer A

Consumer B

# Varying Latency Requirements / Batch vs Stream



Producer

Realtime

Consumer A

24hr batch extract

Consumer B

# Varying Latency Requirements / Batch vs Stream



Producer

Realtime

24hr batch extract

Realtime

Consumer A

Consumer B

# Technology & Code/Algo version Changes

Producer → [cluster] → Consumer (v1)

# Technology & Code/Algo version Changes

# Technology & Code/Algo version Changes



Producer → [Kafka] → Consumer (V2)

# "Architectural Patterns with Apache Kafka

confluent

Building for the Future

Tightly-coupled = Inflexible

# Analytics - Database Offload
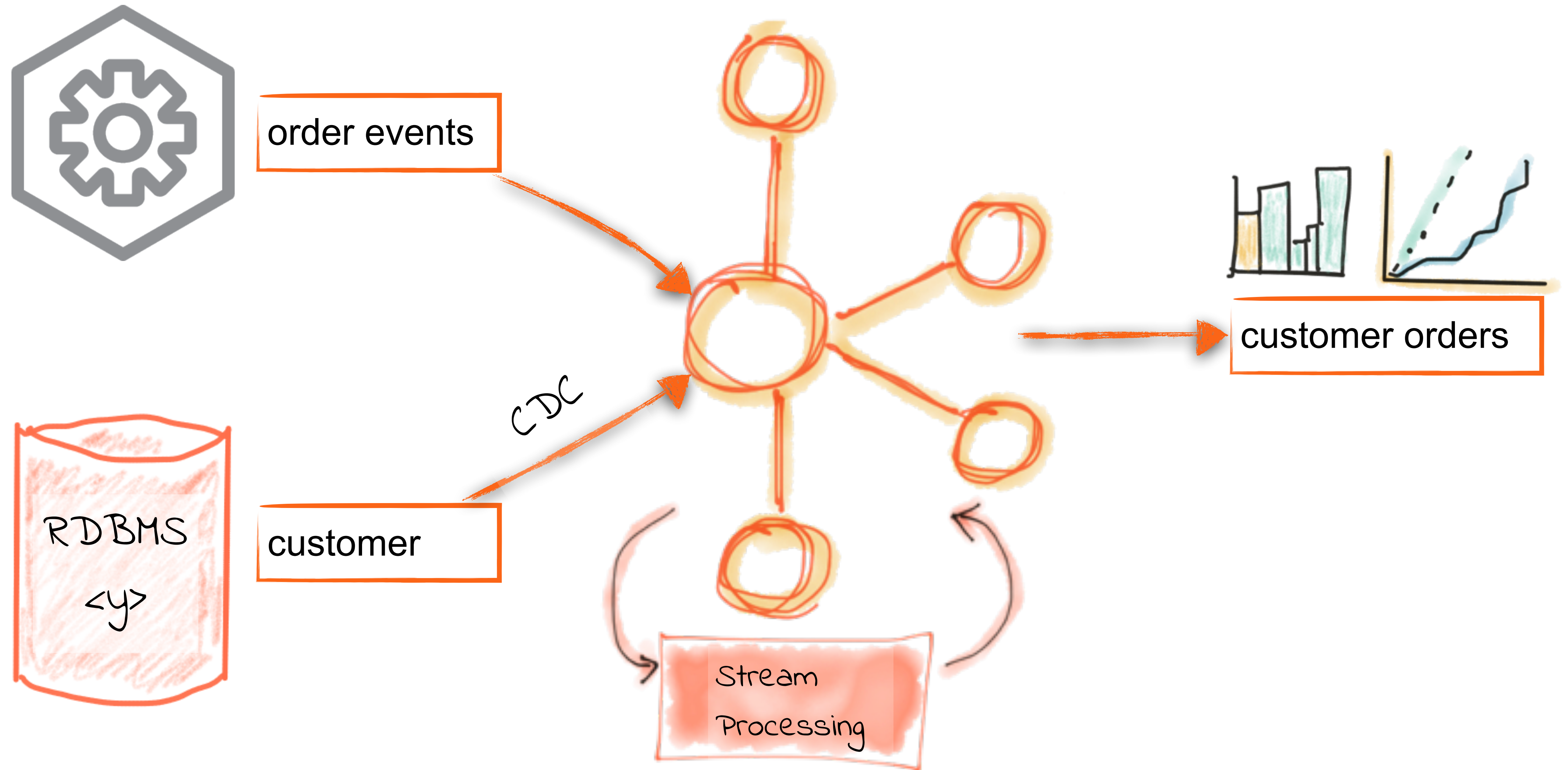


RDBMS

CDC

HDFS / S3 / BigQuery etc

# Stream Processing with Apache Kafka and KSQL



RDBMS

order events

customer

CDC

customer orders

Stream Processing

# Real-time Event Stream Enrichment



order events

customer orders

CDC

RDBMS
<y>

customer

Stream
Processing

# Transform Once, Use Many



order events

customer orders

RDBMS
<y>

customer

CDC

Stream
Processing

New App

# Transform Once, Use Many



order events

RDBMS <y>

customer

CDC

Stream Processing

customer orders

HDFS / S3 / etc
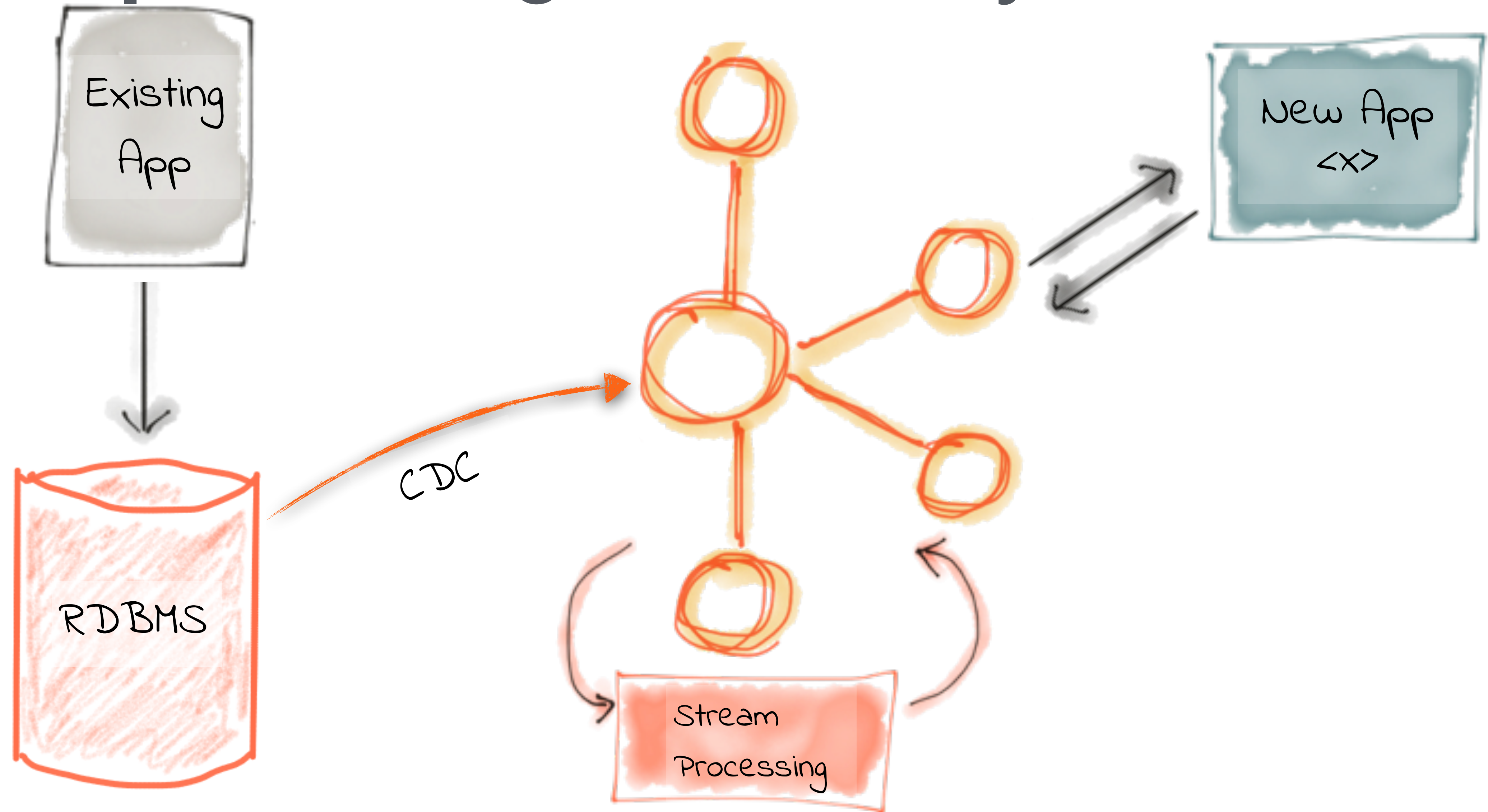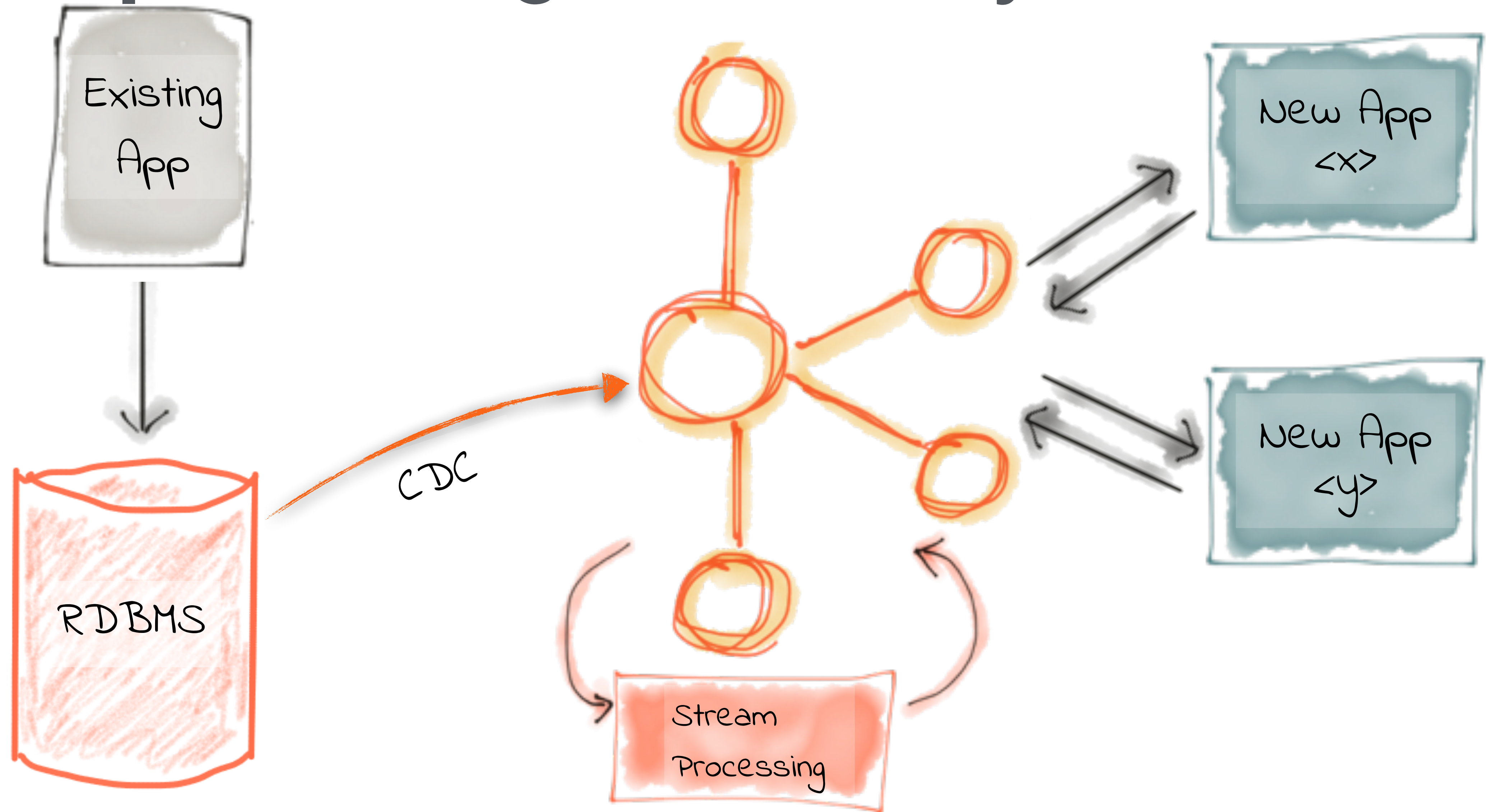
New App <x>

# Evolve processing from old systems to new
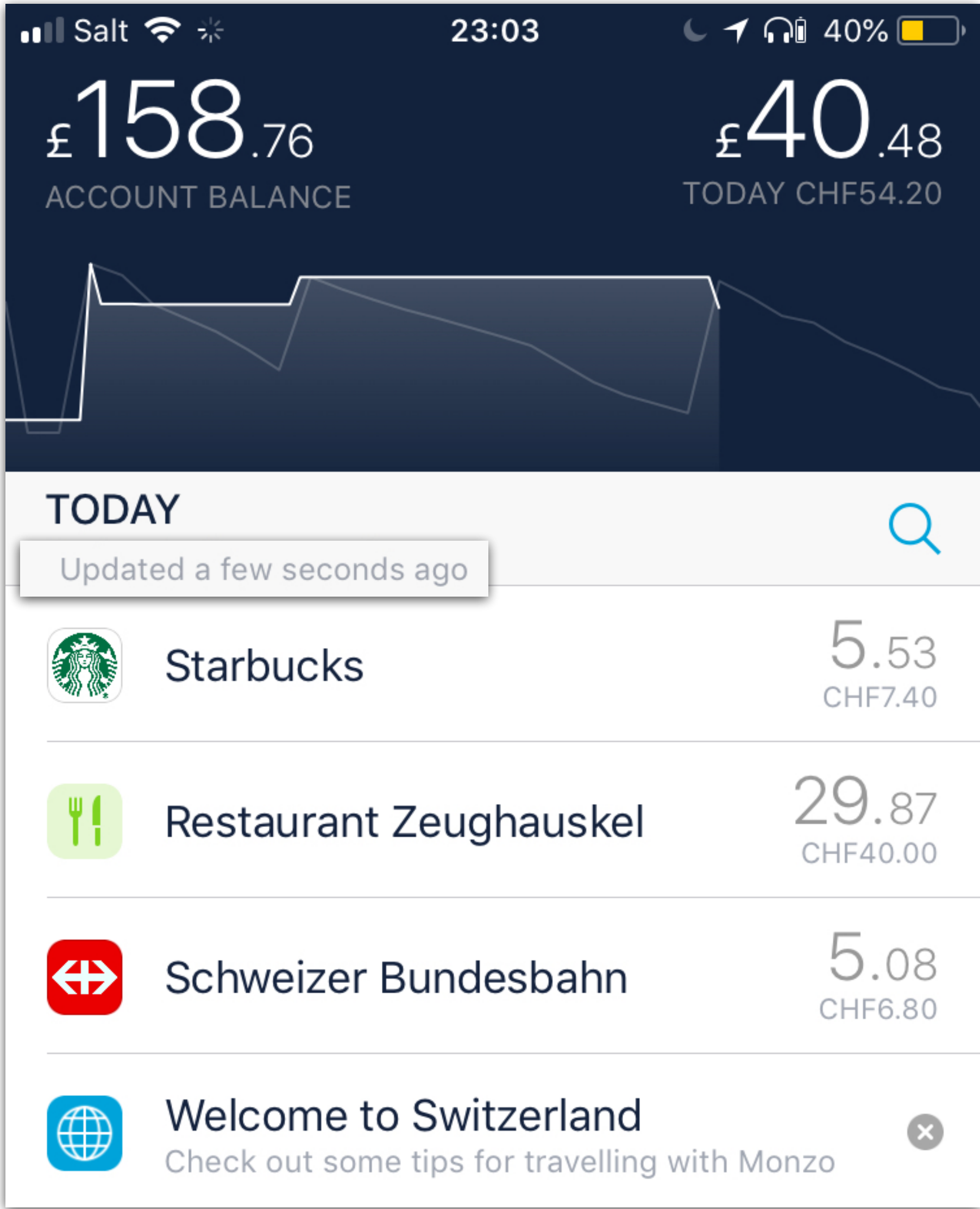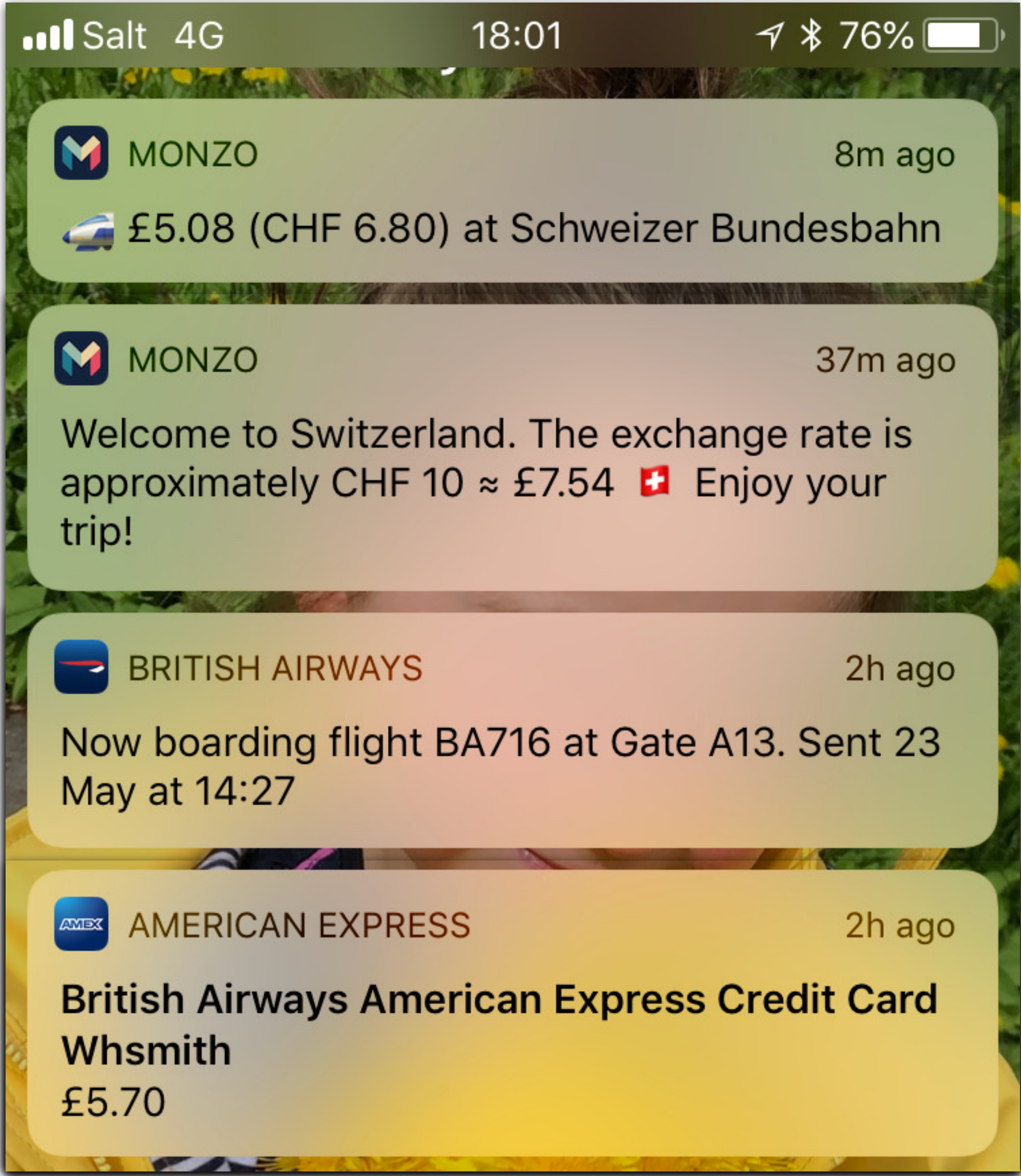
# Evolve processing from old systems to new

Want your data anytime SOON ?

You say that like "latency" is a synonym for "evil"

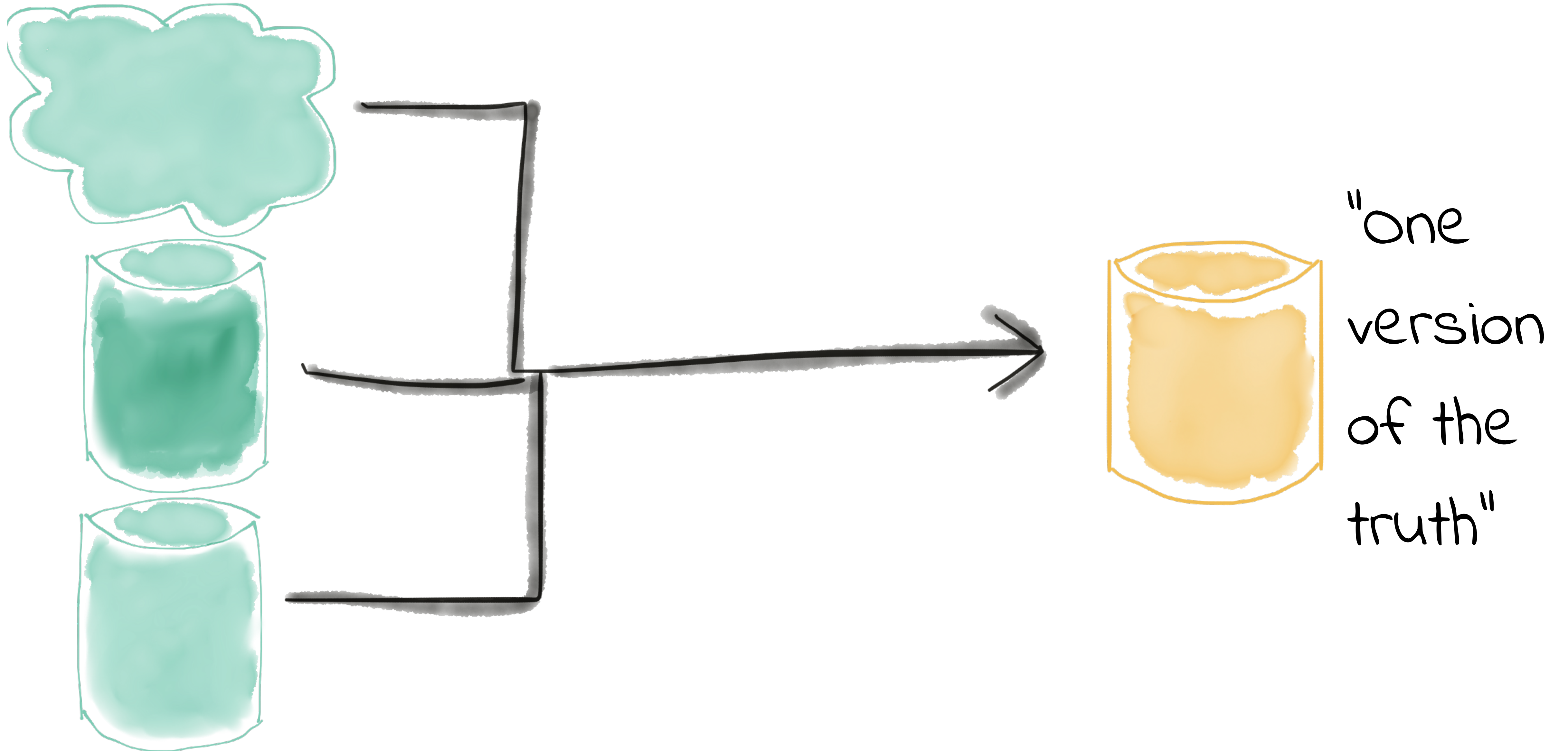Batch is Latency built in by Design

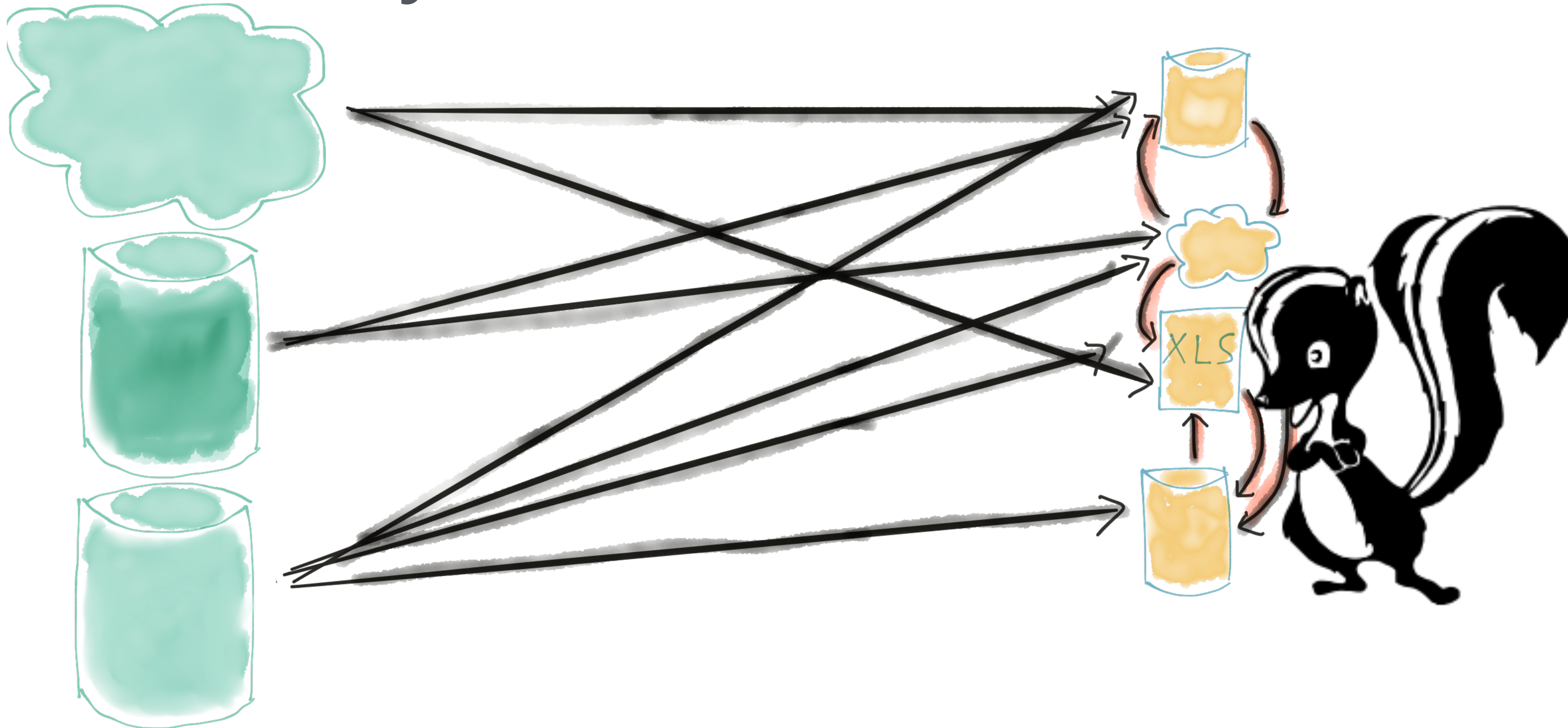# It's all about the Events!

"

# So...Analytics and Kafka
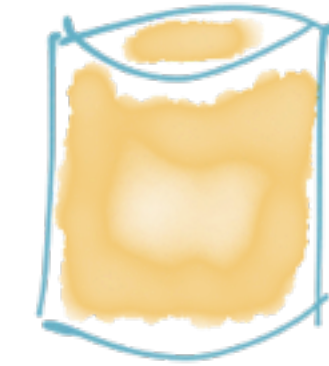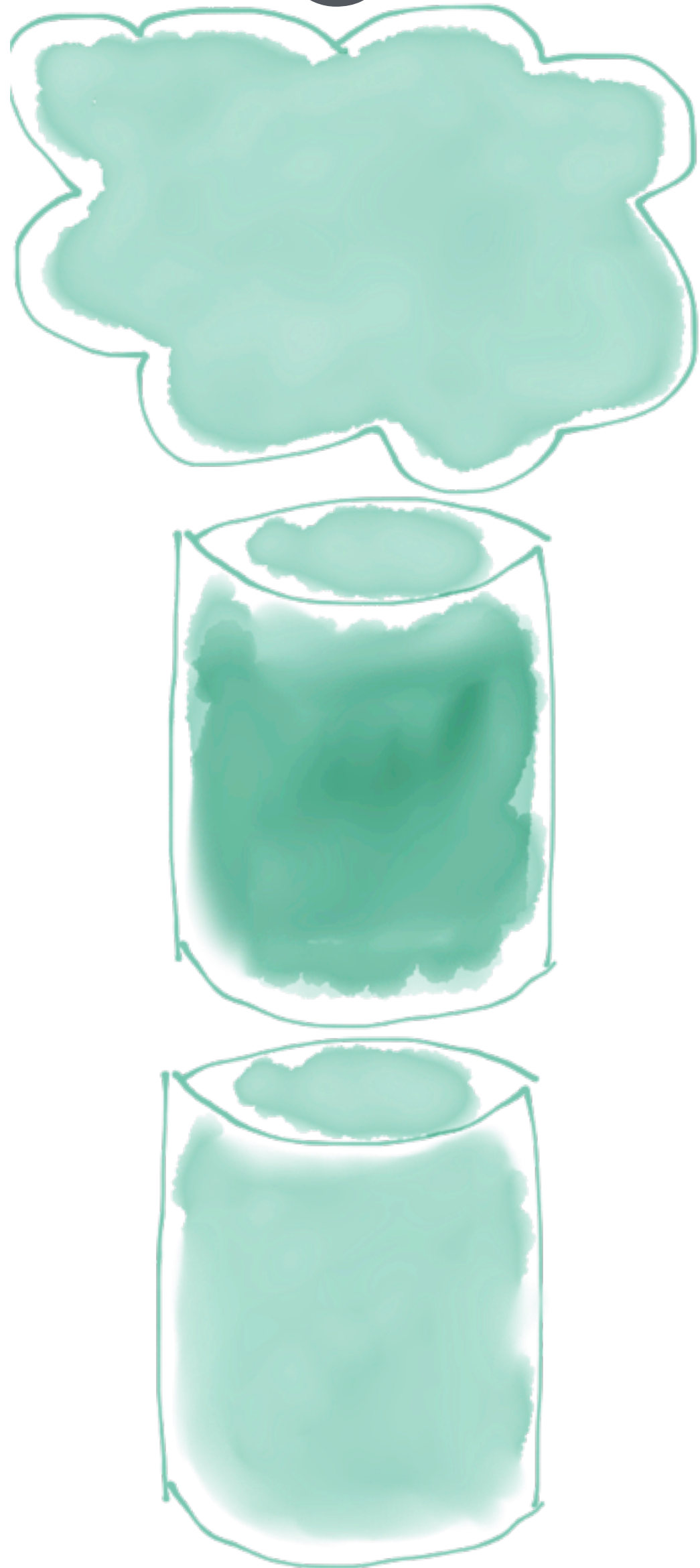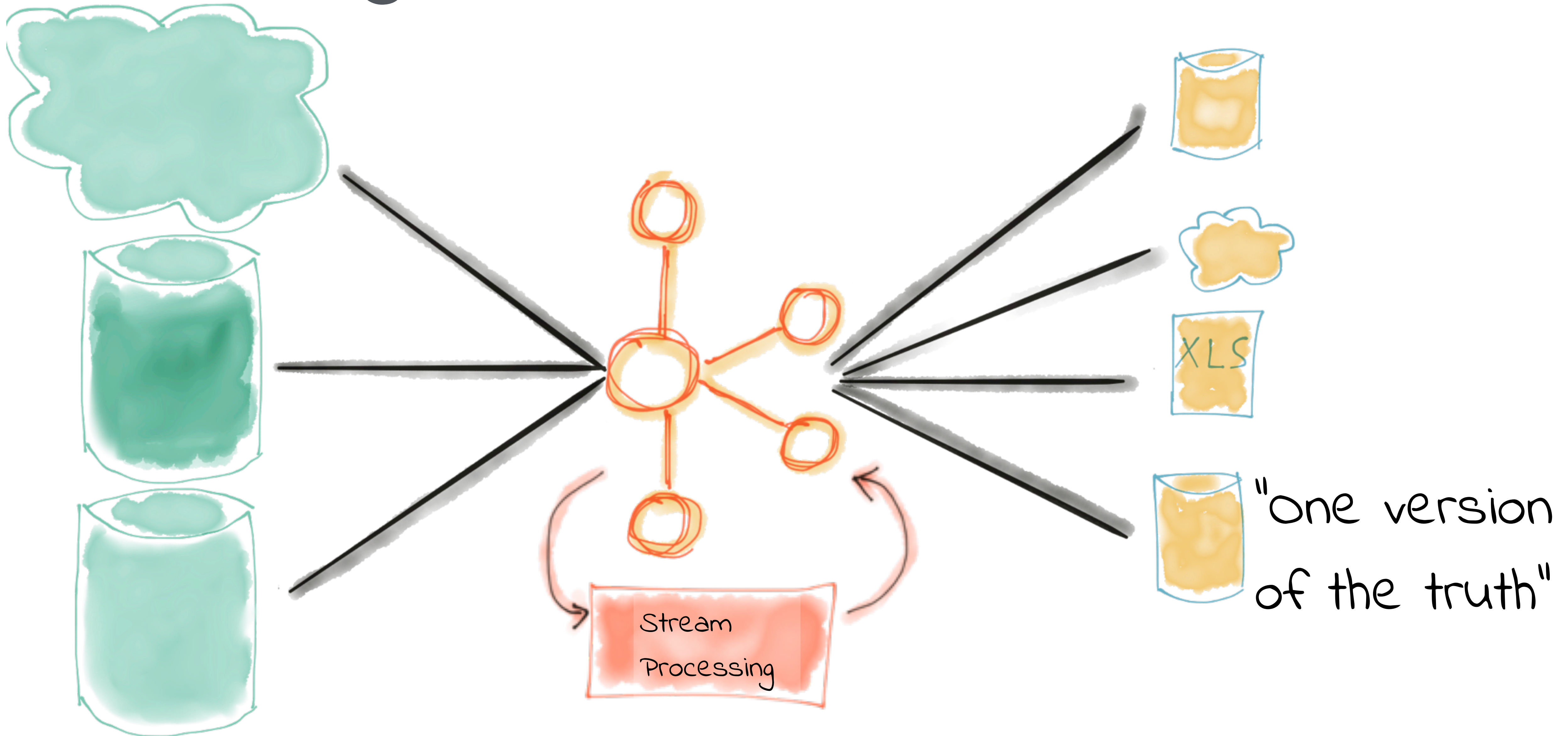
# The Vision!



"one version of the truth"

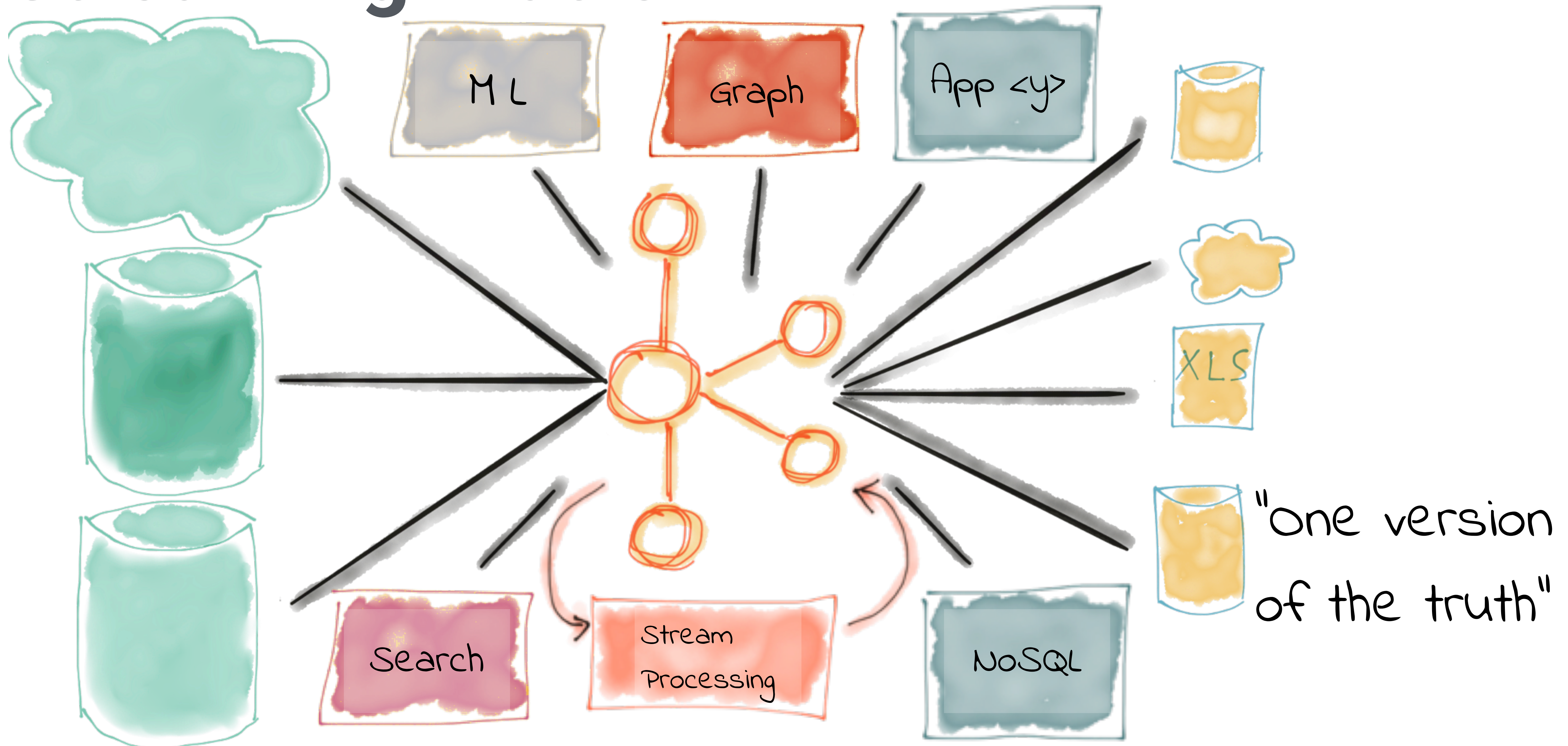# The Reality...

# Pragmatism is...
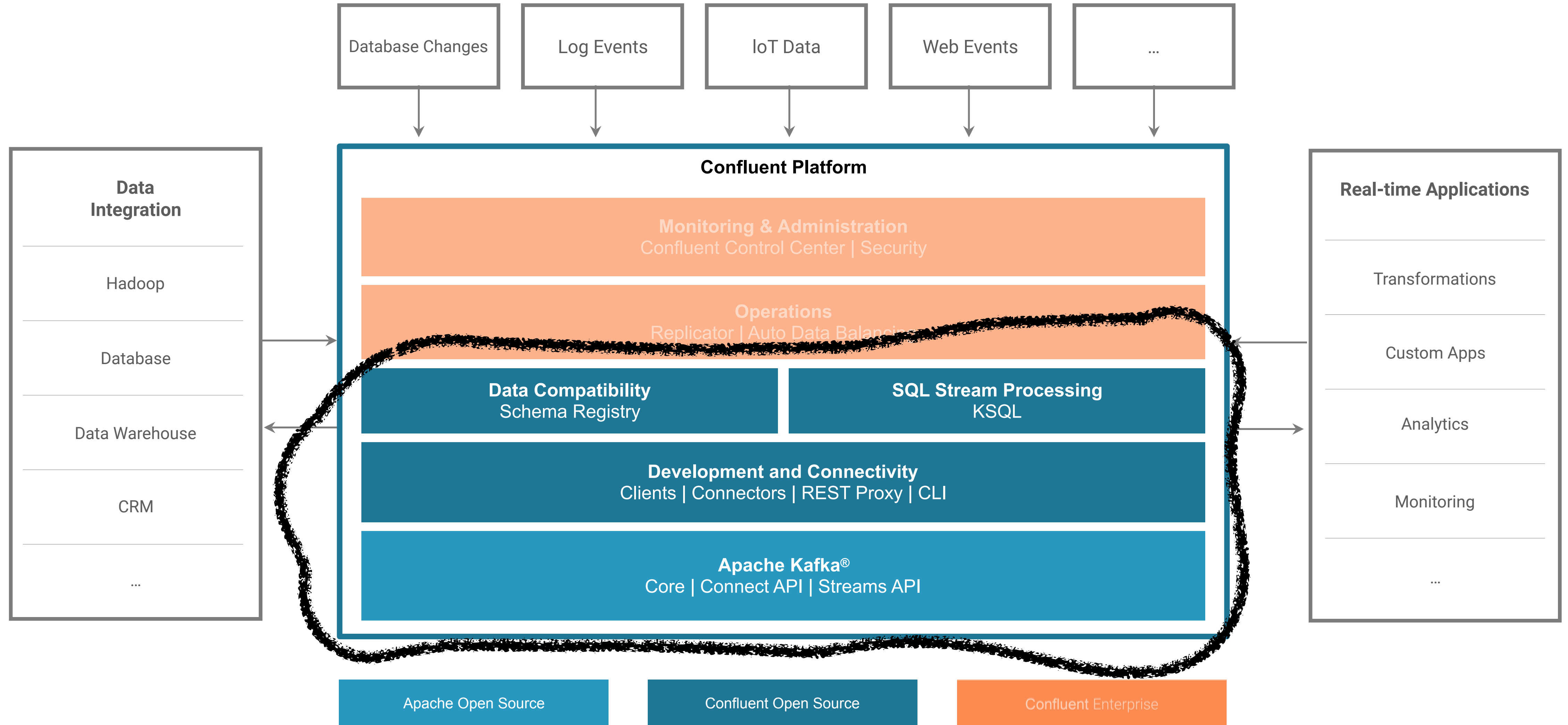
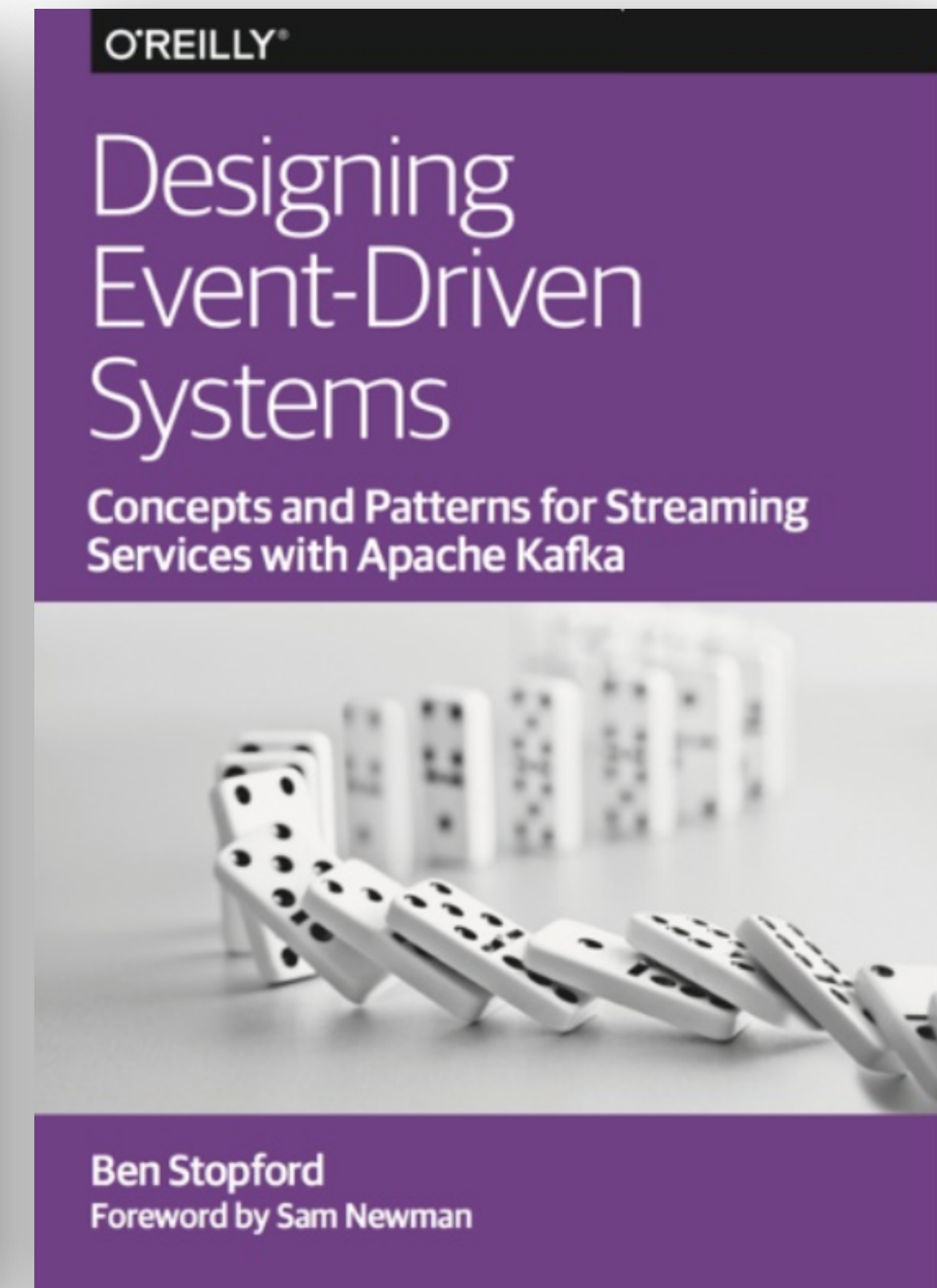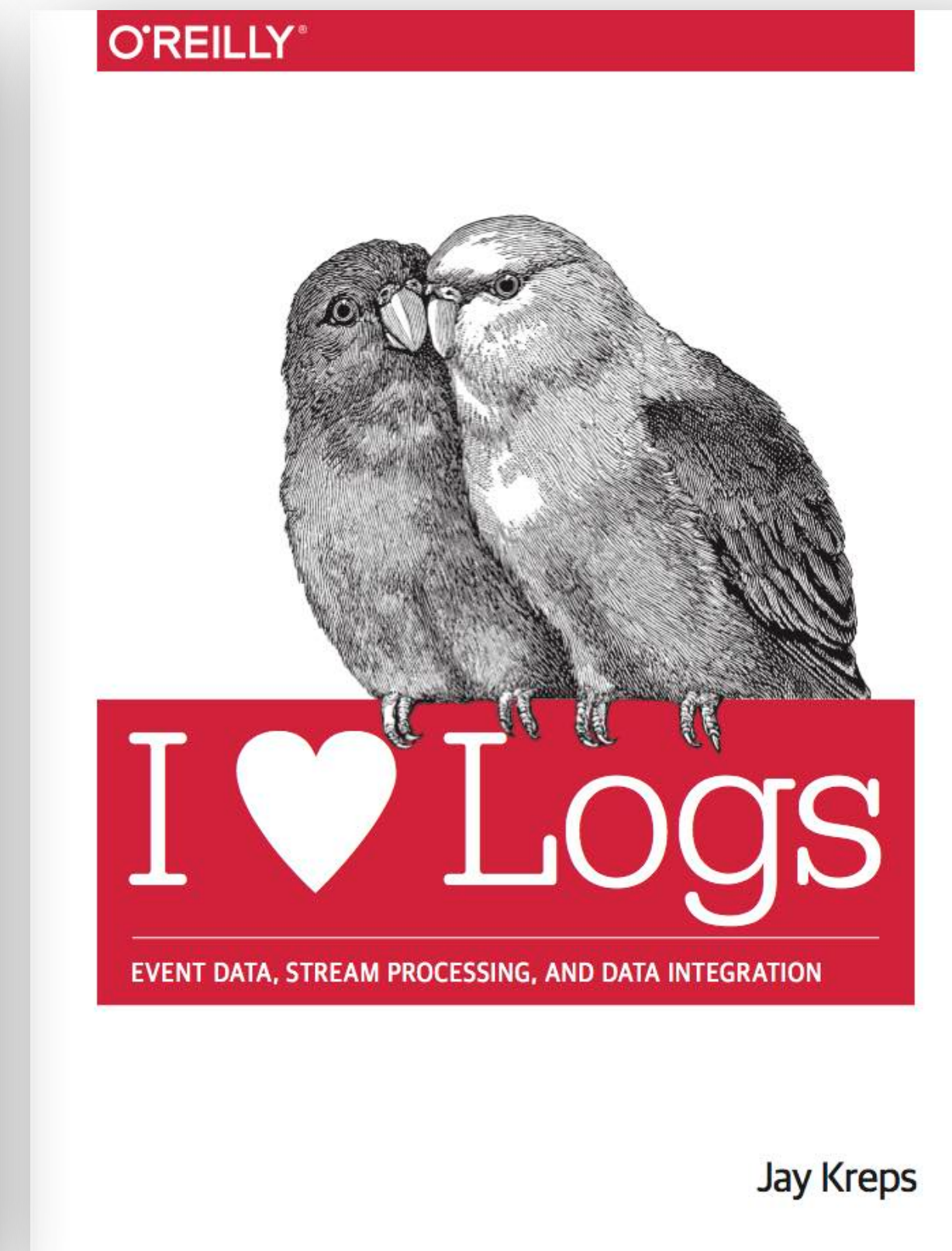"one version of the truth"

# Streaming Platform



Stream
Processing

"one version
of the truth"

XLS

# Streaming Platform



ML

Graph

App <y>

XLS

Search

Stream Processing

NoSQL

"one version of the truth"

# Confluent Open Source :
# Apache Kafka with a bunch of cool stuff! For free!



| Database Changes | Log Events | IoT Data | Web Events | ... |

**Confluent Platform**

**Data Integration**

Hadoop

Database

Data Warehouse

CRM

...

**Monitoring & Administration**
Confluent Control Center | Security

**Operations**
Replicator | Auto Data Balancer

**Data Compatibility**
Schema Registry

**SQL Stream Processing**
KSQL

**Development and Connectivity**
Clients | Connectors | REST Proxy | CLI

**Apache Kafka®**
Core | Connect API | Streams API

**Real-time Applications**

Transformations

Custom Apps

Analytics

Monitoring

...

| Apache Open Source | Confluent Open Source | Confluent Enterprise |

# Free Books!

**https://www.confluent.io/apache-kafka-stream-processing-book-bundle**

# Confluent Streaming Event, Munich



# http://cnfl.io/streaming-event-munich

https://www.confluent.io/download/

http://cnfl.io/slack

@rmoff

robin@confluent.io

#EOF