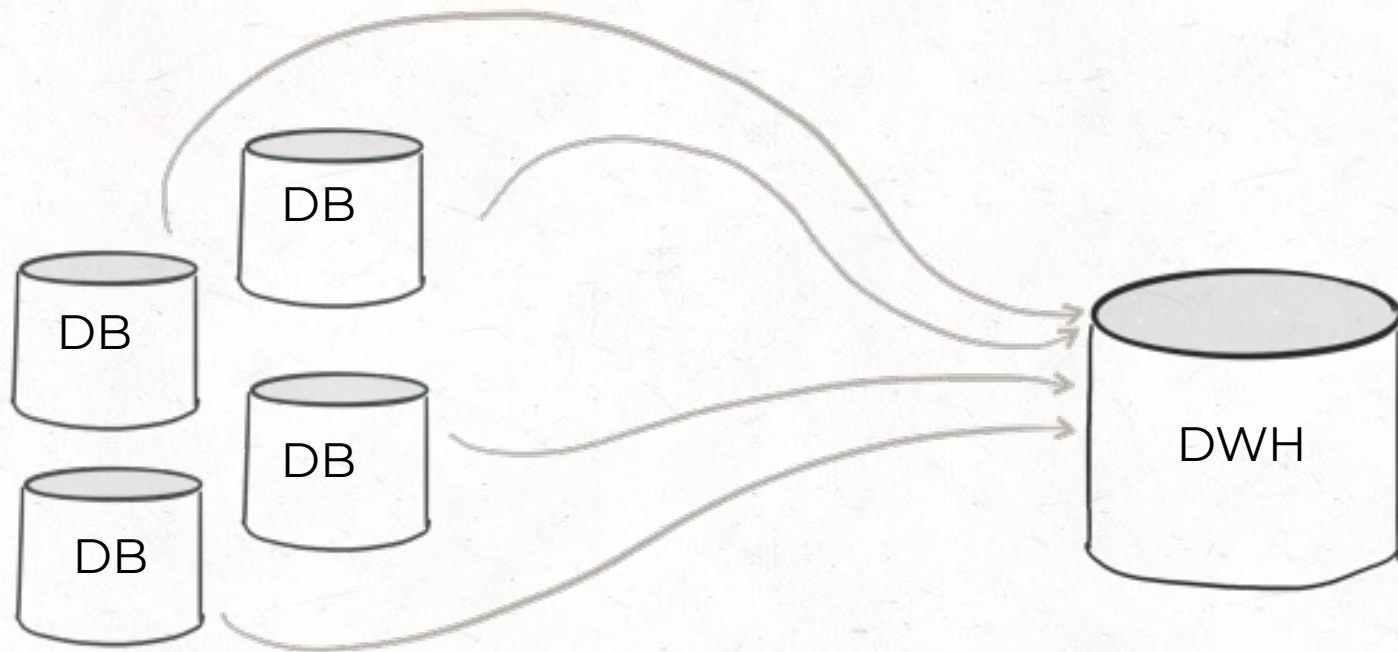# Achieving Streaming ETL

Brian Likosar @liko9
liko@confluent.io

CONFLUENT

" 

*Data and data systems have really changed in the past decade*

CONFLUENT

# Old world: Two popular locations for data



DB

DB

DB

DB

DWH

Operational databases

Relational data warehouse

CONFLUENT

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

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

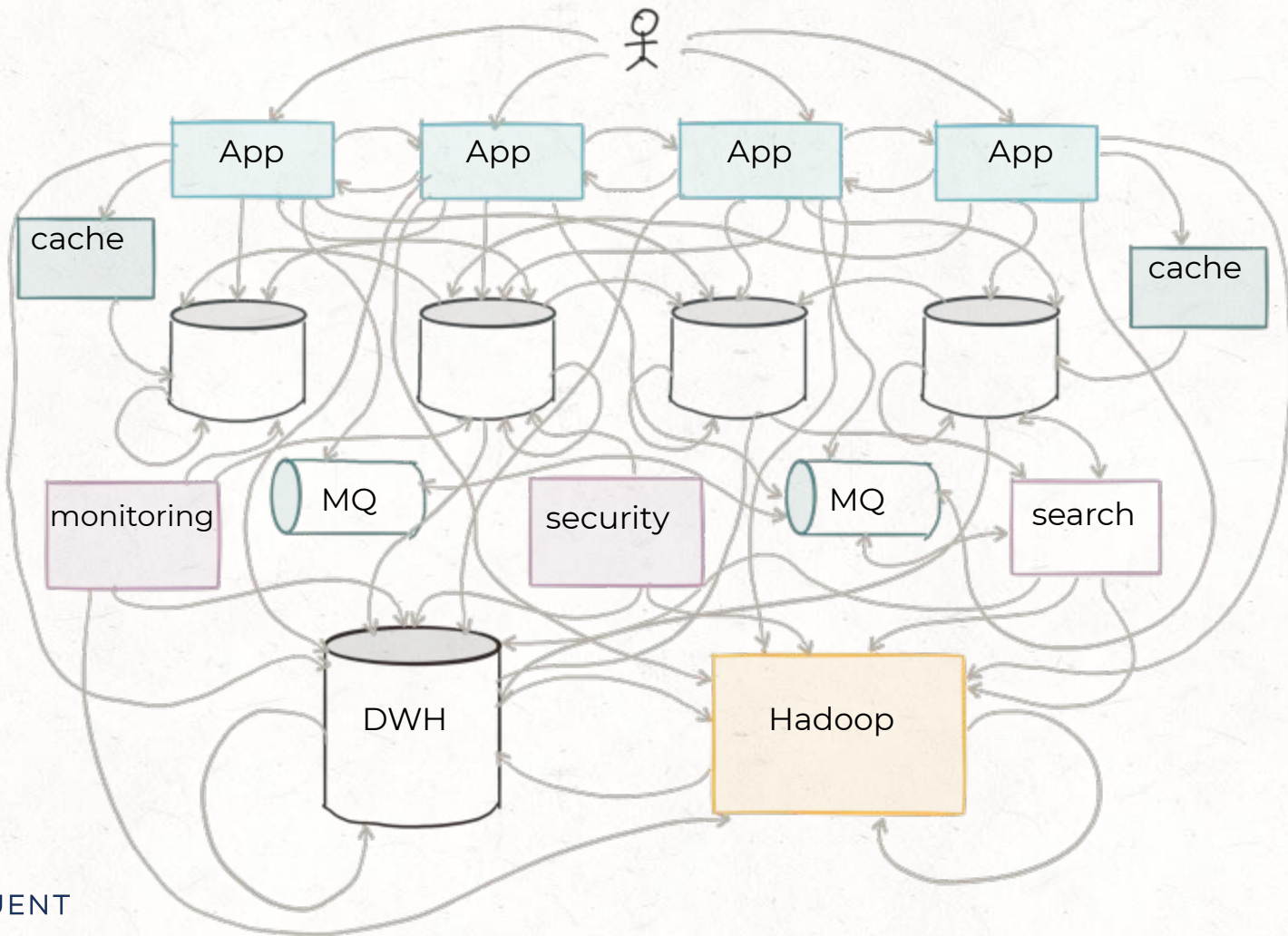*#2: Many more types of data sources beyond transactional data - logs, sensors, metrics...*

"

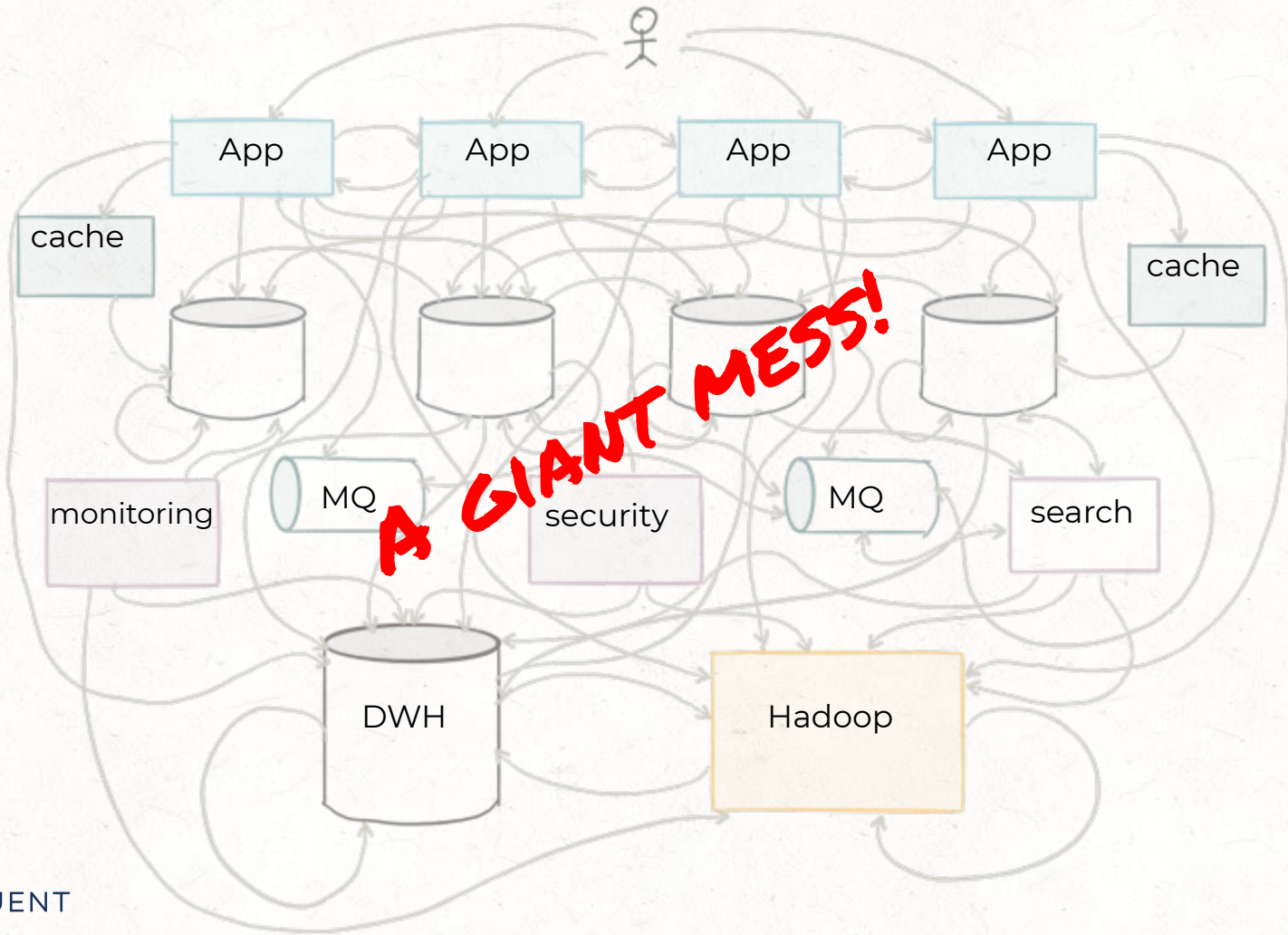*#3: Stream data is increasingly ubiquitous; need for faster processing than daily*

CONFLUENT

App App App App

cache cache

monitoring MQ security MQ search

DWH Hadoop

CONFLUENT

A GIANT MESS!

CONFLUENT

request-response

changelogs

messaging
OR
stream
processing

search

monitoring

Streaming platform

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

CONFLUENT

# #1: The need for a **global** **schema**

"

#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

*Early take on **real-time** ETL*

*=*

*Enterprise Application Integration*

*(EAI)*

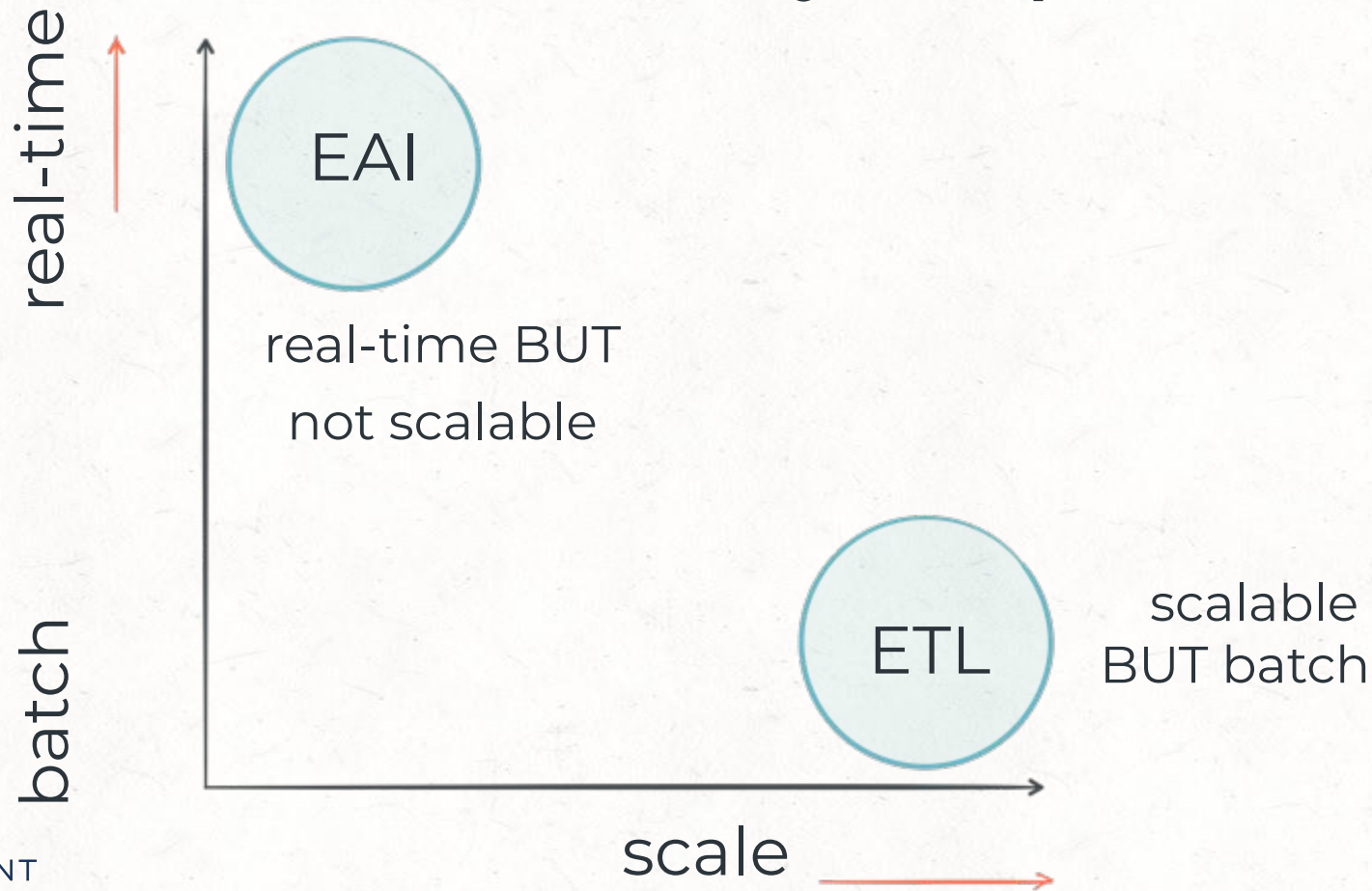**EAI**: *A different class of data integration technology for connecting applications in real-time*

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

CONFLUENT

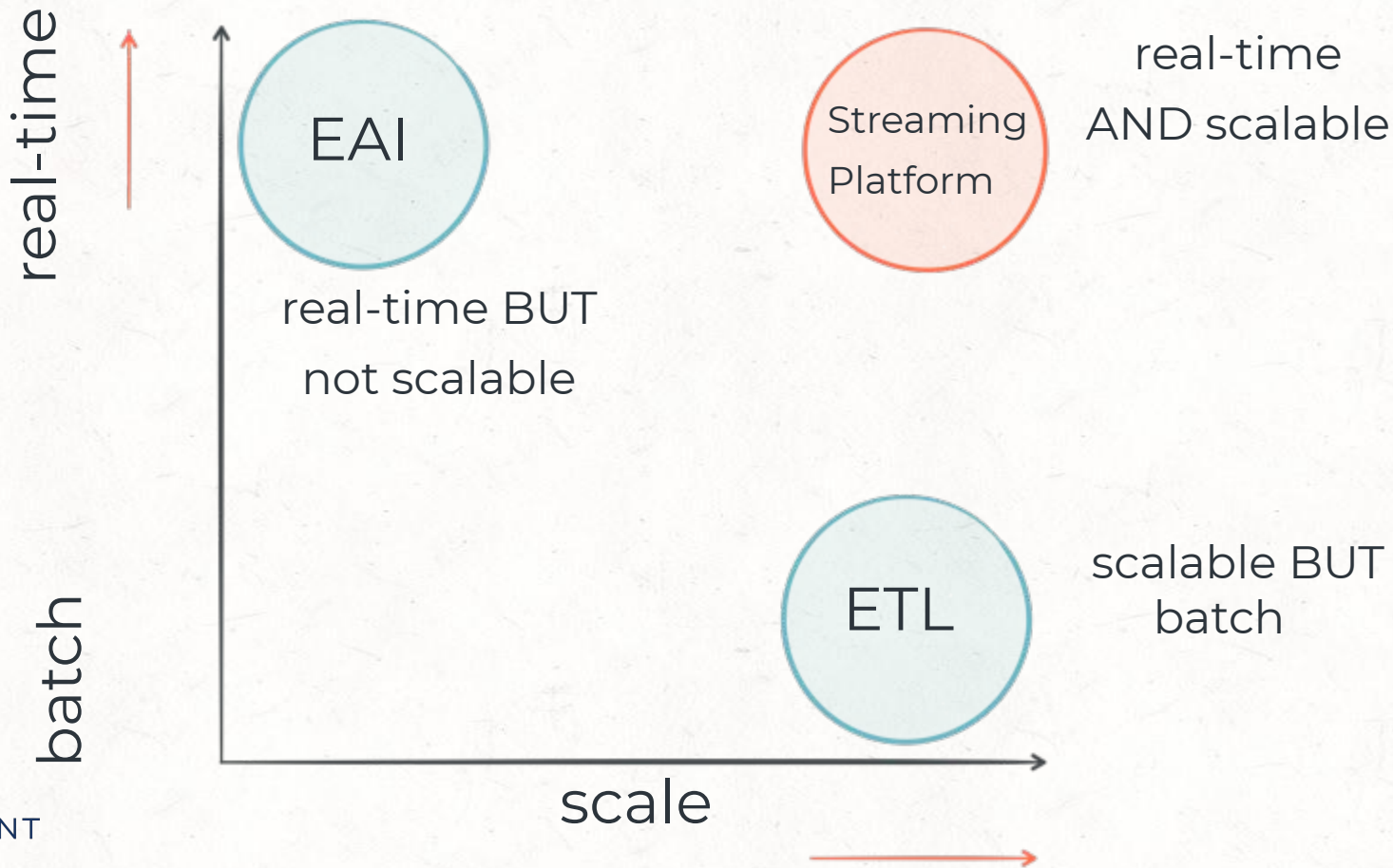# ETL and EAI are outdated!

# Old world: scale or timely data, pick one

real-time

EAI

real-time BUT
not scalable

batch

ETL

scalable
BUT batch

scale

CONFLUENT

"

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

CONFLUENT

# new world: streaming, real-time and scalable

real-time

batch

scale

**EAI**

real-time BUT
not scalable

**Streaming
Platform**

real-time
AND scalable

**ETL**

scalable BUT
batch

CONFLUENT

*Modern streaming world has new set of* **requirements** *for data integration*

CONFLUENT

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

CONFLUENT

#2 A *streaming* platform from the grounds up; a fundamental transition to event-centric thinking
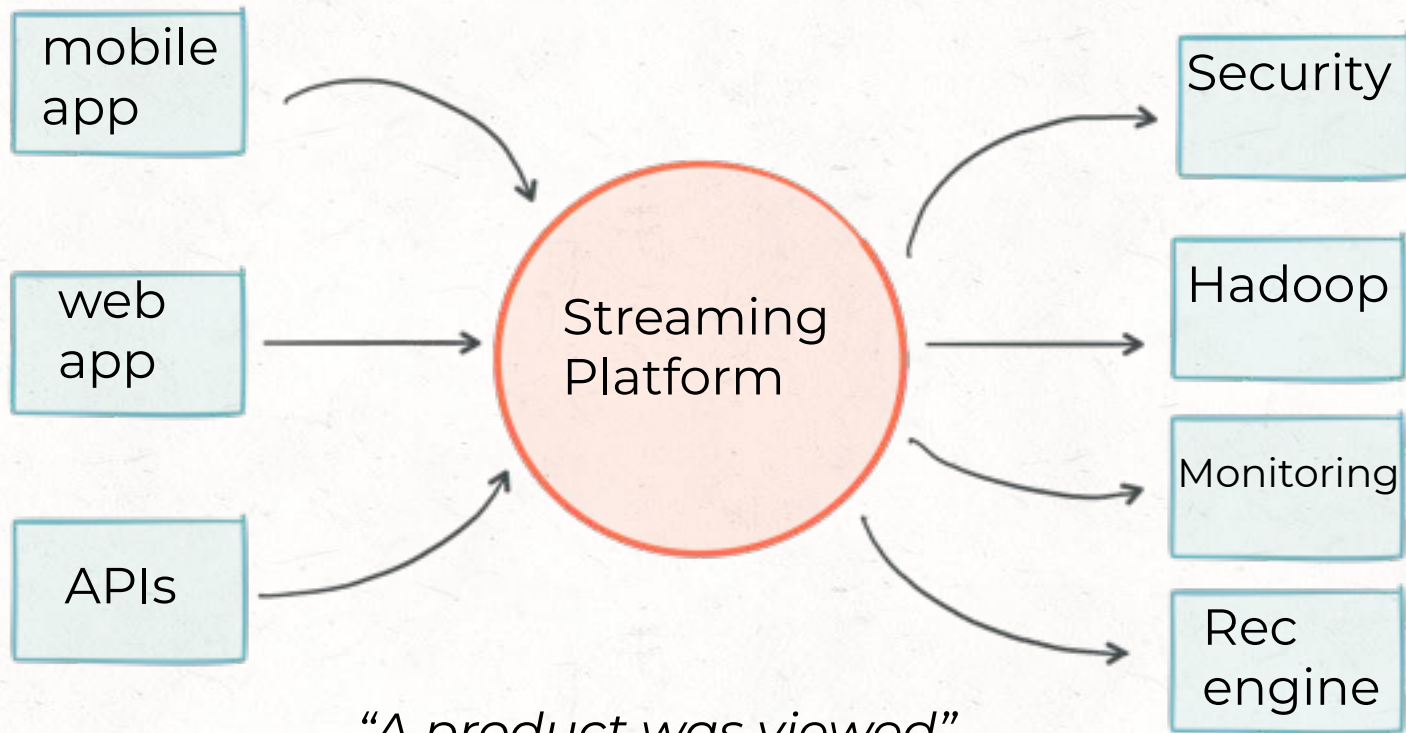
CONFLUENT

# Event-Centric Thinking

Web app → Streaming Platform → Hadoop

*"A product was viewed"*

# Event-Centric Thinking



mobile app

Web app

APIs

Streaming Platform

Hadoop

*"A product was viewed"*

CONFLUENT

# Event-Centric Thinking



*"A product was viewed"*

CONFLUENT

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

CONFLUENT

request-response

changelogs

messaging
 OR
stream
processing

App
App
App
App

Streaming platform

App

App

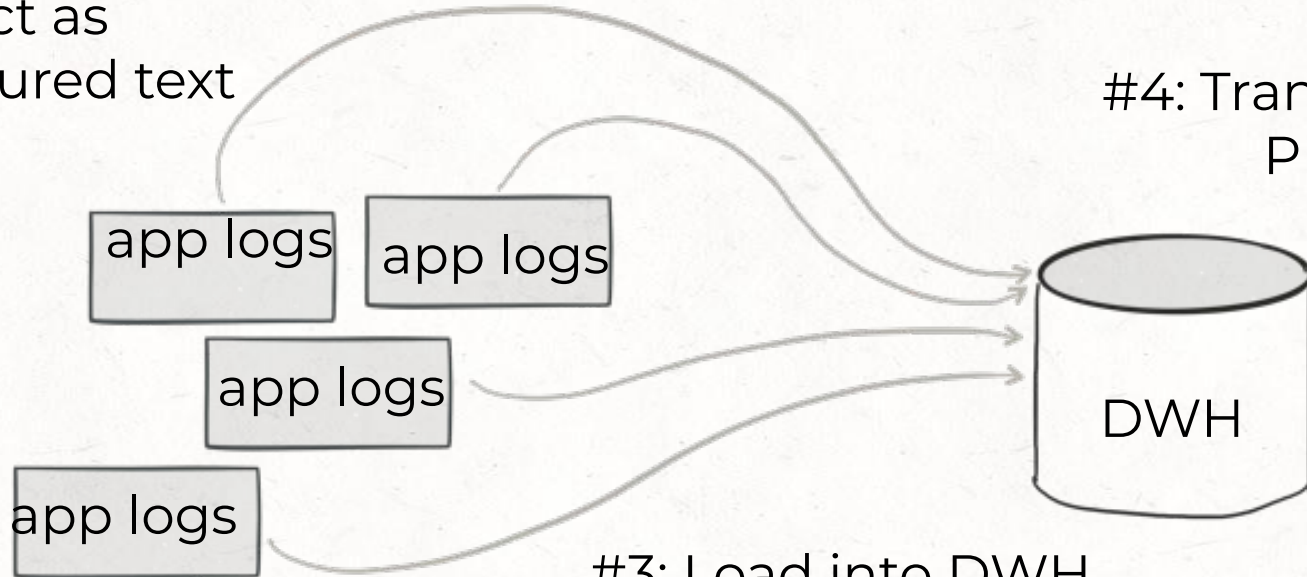App

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 = dro 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"

#1: Extract as unstructured text again

#3: Load cleansed data

#4: Transform2 = drop PII fields"
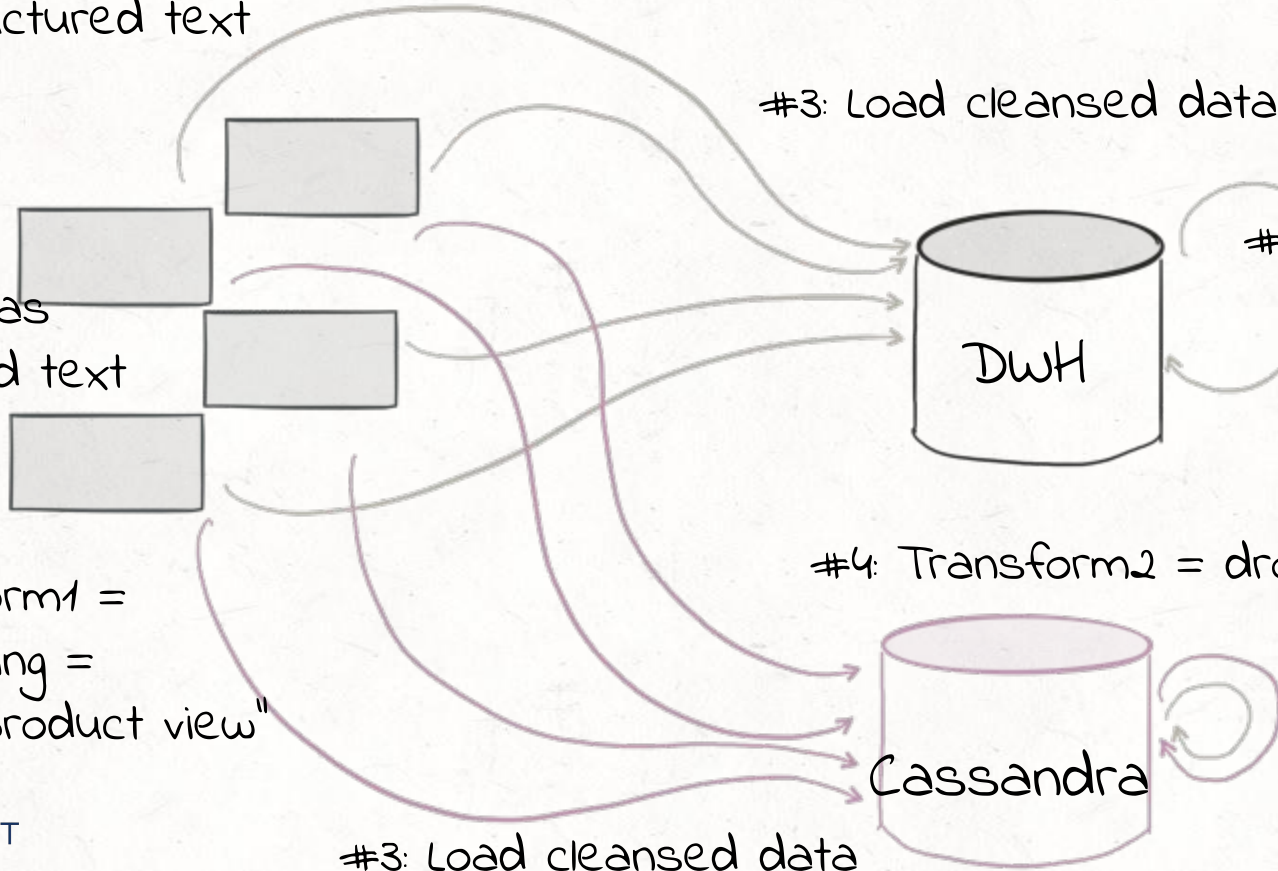
DWH

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

#4: Transform2 = drop PII fields"

Cassandra
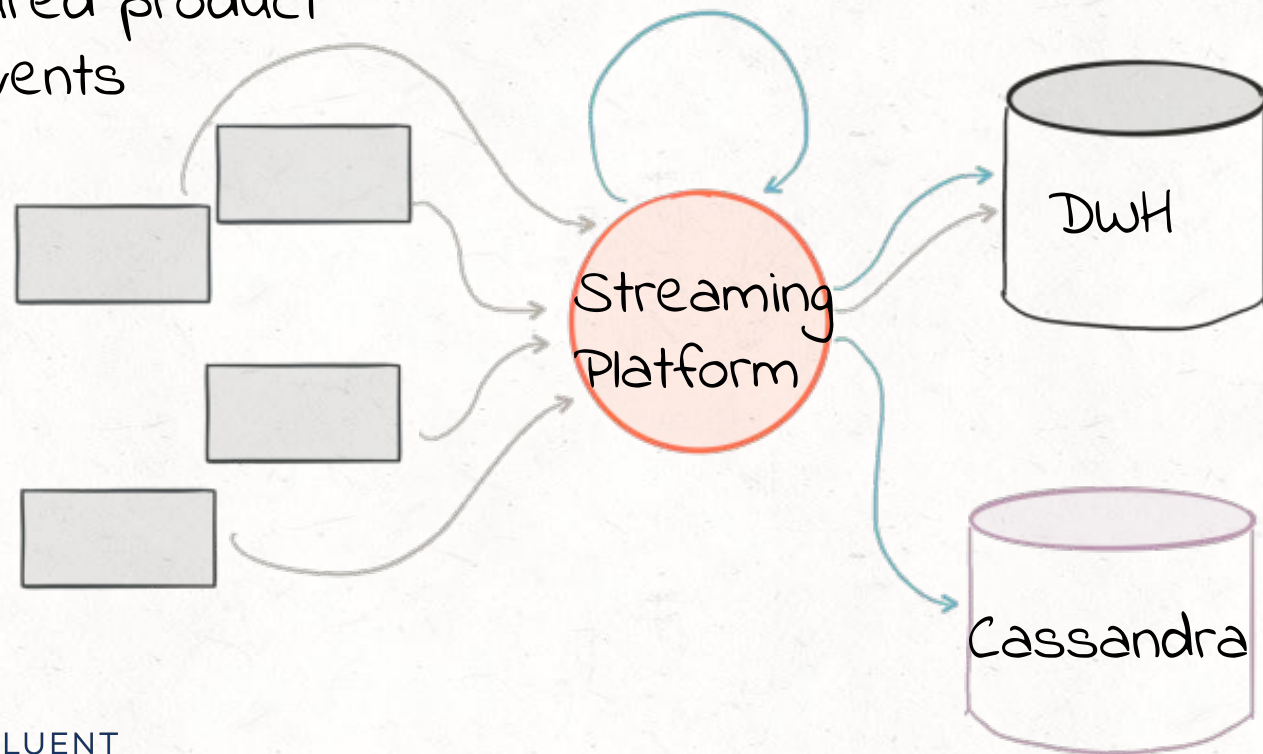
#3: Load cleansed data

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

#4: Load filtered product view stream

Streaming Platform

DWH

Cassandra

CONFLUENT

*In summary, "needs of modern data integration solution? Scale, diversity, latency and forward compatibility*

# Requirements for a modern streaming data integration solution

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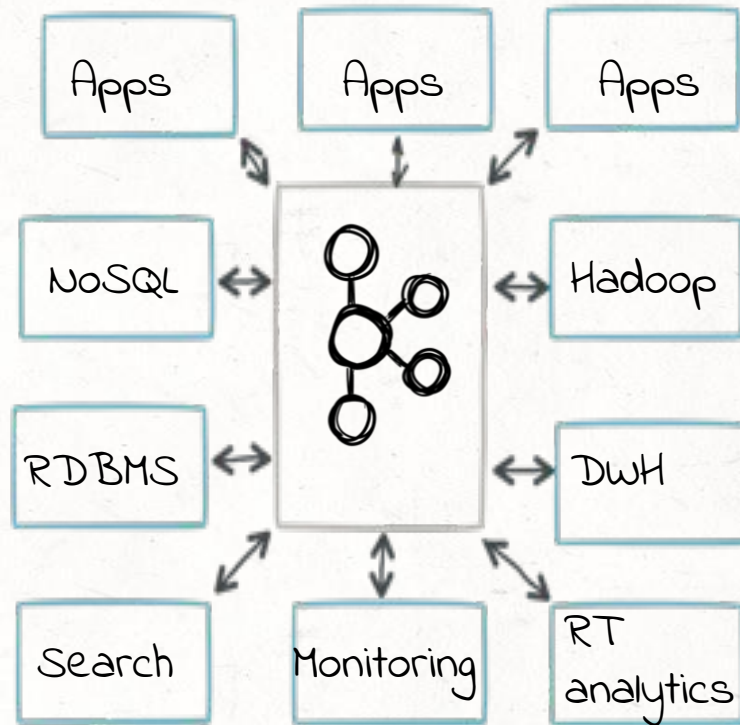- Operations and monitoring
- Schema management

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 ...*"

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

CONFLUENT

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

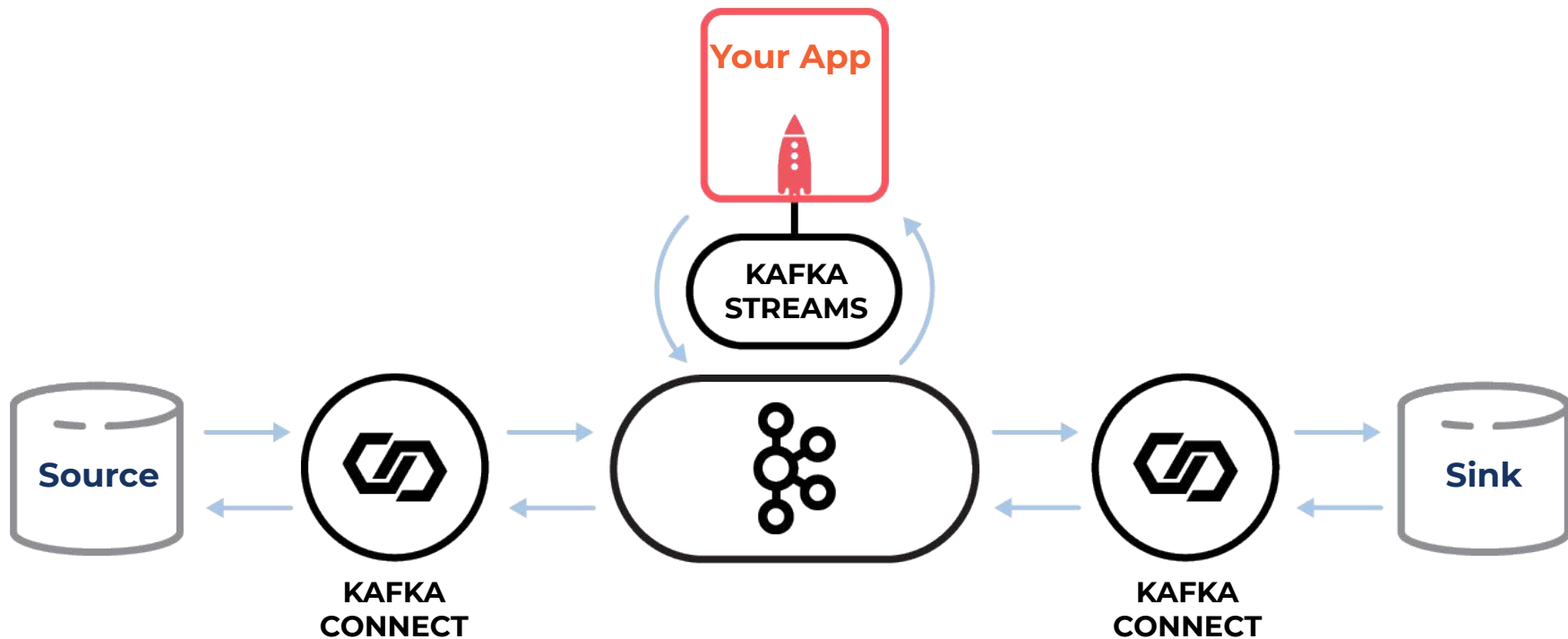#3: Serves as the "building block for stateful stream processing microservices

# What does a streaming platform look like and how does it enable Streaming ETL?

# Kafka Connect and Kafka Streams

# Instantly Connect Popular Data Sources & Sinks

**100+** pre-built connectors

**80+ Confluent Supported**

**20+ Partner Supported, Confluent Verified**

# Confluent Platform

| DEVELOPER | OPERATOR | ARCHITECT |
|---|---|---|
| **Unrestricted Developer Productivity** | **Efficient Operations at Scale** | **Production-stage Prerequisites** |
| **Multi-language Development** Non-Java Clients ● \| REST Proxy ● | **GUI-driven Mgmt & Monitoring** Control Center | **Enterprise-grade Security** RBAC \| Secrets \| Audit Logs |
| **Rich Pre-built Ecosystem** Connectors \| Hub ●\| Schema Registry ● | **Flexible DevOps Automation** Operator \| Ansible ● | **Data Compatibility** Schema Registry ● \| Schema Validation |
| **Event Streaming Database** ksqlDB ● | **Dynamic Performance & Elasticity** Auto Data Balancer \| Tiered Storage | **Global Resilience** Multi-region Clusters \| Replicator |

## ⚙ Apache Kafka

● Open Source | Community licensed

**Self-managed Software**          **Freedom of Choice**          **Fully Managed Cloud Service**

**Enterprise Support**   **Professional Services**   **Committer-driven Expertise**   **Training**   **Partners**

# Two Ways to Deploy Confluent

## SELF-MANAGED SOFTWARE

### Confluent Platform

The Enterprise Distribution of Apache Kafka

Deploy on any platform, on-prem or cloud

## FULLY-MANAGED SOFTWARE

### Confluent Cloud

Apache Kafka Re-Engineered for the Cloud

Available on the leading public clouds