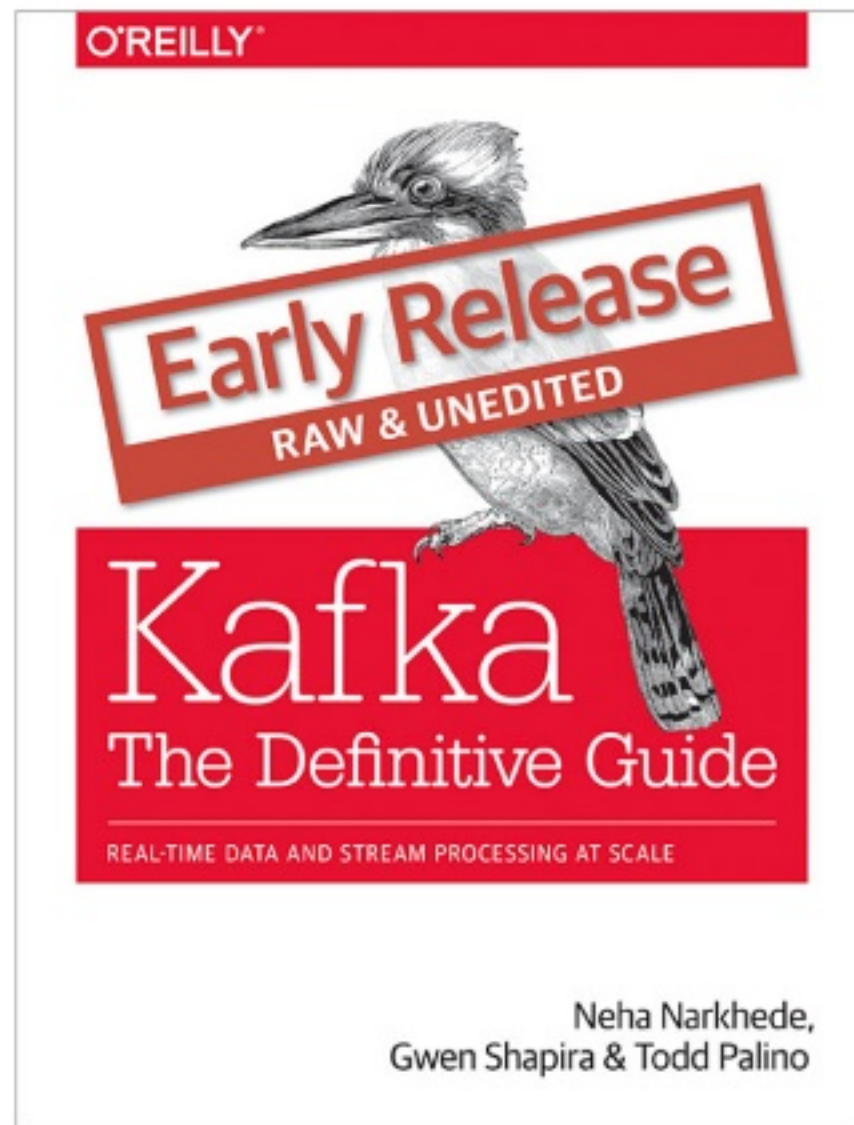# Stream All Things

## Patterns of Modern Data Integration

# About Me

- Working @confluent
- Wrote a book or two
- Product manager
- Apache Kafka PMC member
- Used to be an engineer, consultant, DBA
- @gwenshap
- Github.com/gwenshap



O'REILLY®

Early Release
RAW & UNEDITED

Kafka
The Definitive Guide

REAL-TIME DATA AND STREAM PROCESSING AT SCALE

Neha Narkhede,
Gwen Shapira & Todd Palino

# Evolution of Data Integration

**Data Warehouse**

**Data Lake**

**Data Stream**

- Difficult modeling
- Difficult ingest
- Missing data
- Limited throughput
- Batch only
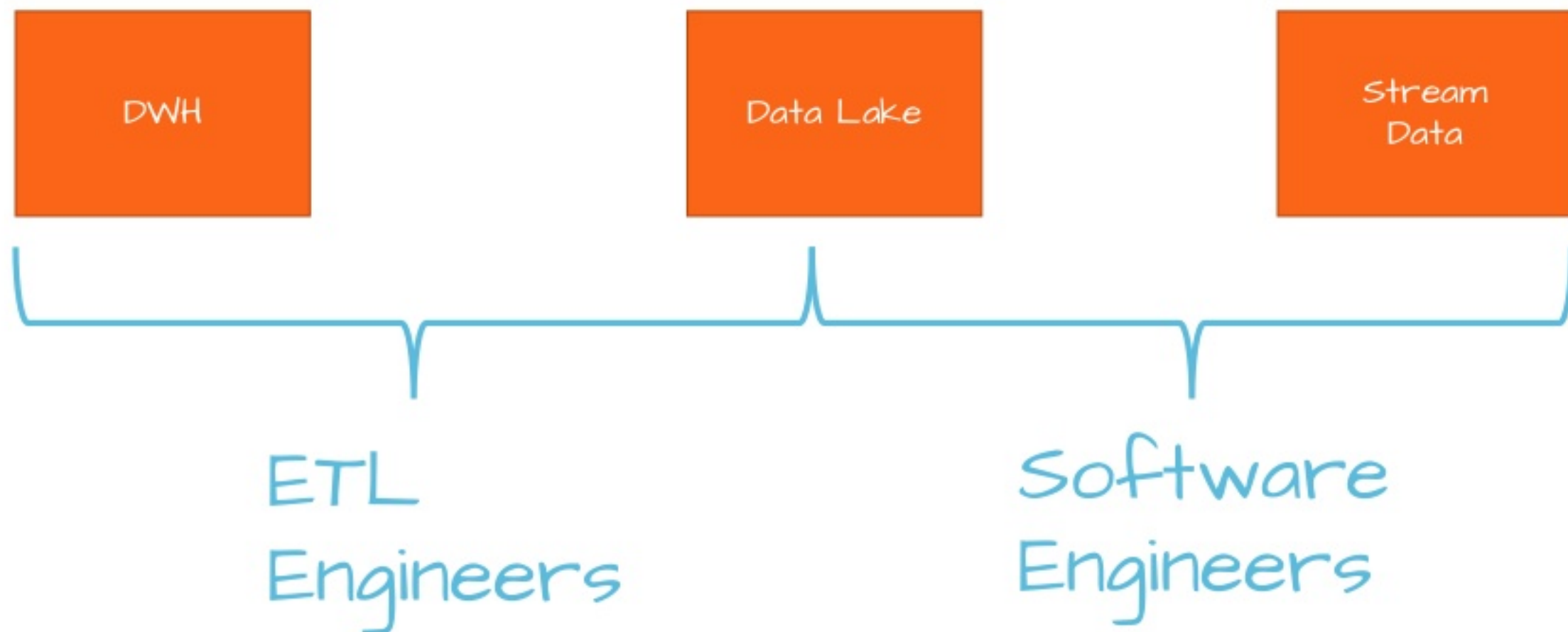
- Still batch only
- Data quality is suspect
- Data model left as exercise for each developer
- Where is my data?

- Kept the scale
- Reduced latency
- Totally new model
- Data producers are 1st class participants
- Microservices? Apps?
- Still figuring the whole thing out

# Evolution of Data Integration



DWH

Data Lake

Stream Data

ETL Engineers

Software Engineers

# Evolution of IT

## 2000

Software engineers
Data modelers
DBAs
Sysadmins
Storage admins
Network admins
ETL engineers
QA

.....

## 2017

Software engineers

# The line between application development and ETL is blurring.

# BOLD CLAIM:
## ALL YOUR DATA
## IS
# EVENT STREAMS

# Few Patterns

1. Stream all things (in one place)

2. Keep Compatible and Process On

3. Ridiculously Parallel Data Integration

4. Streaming Data Enrichment
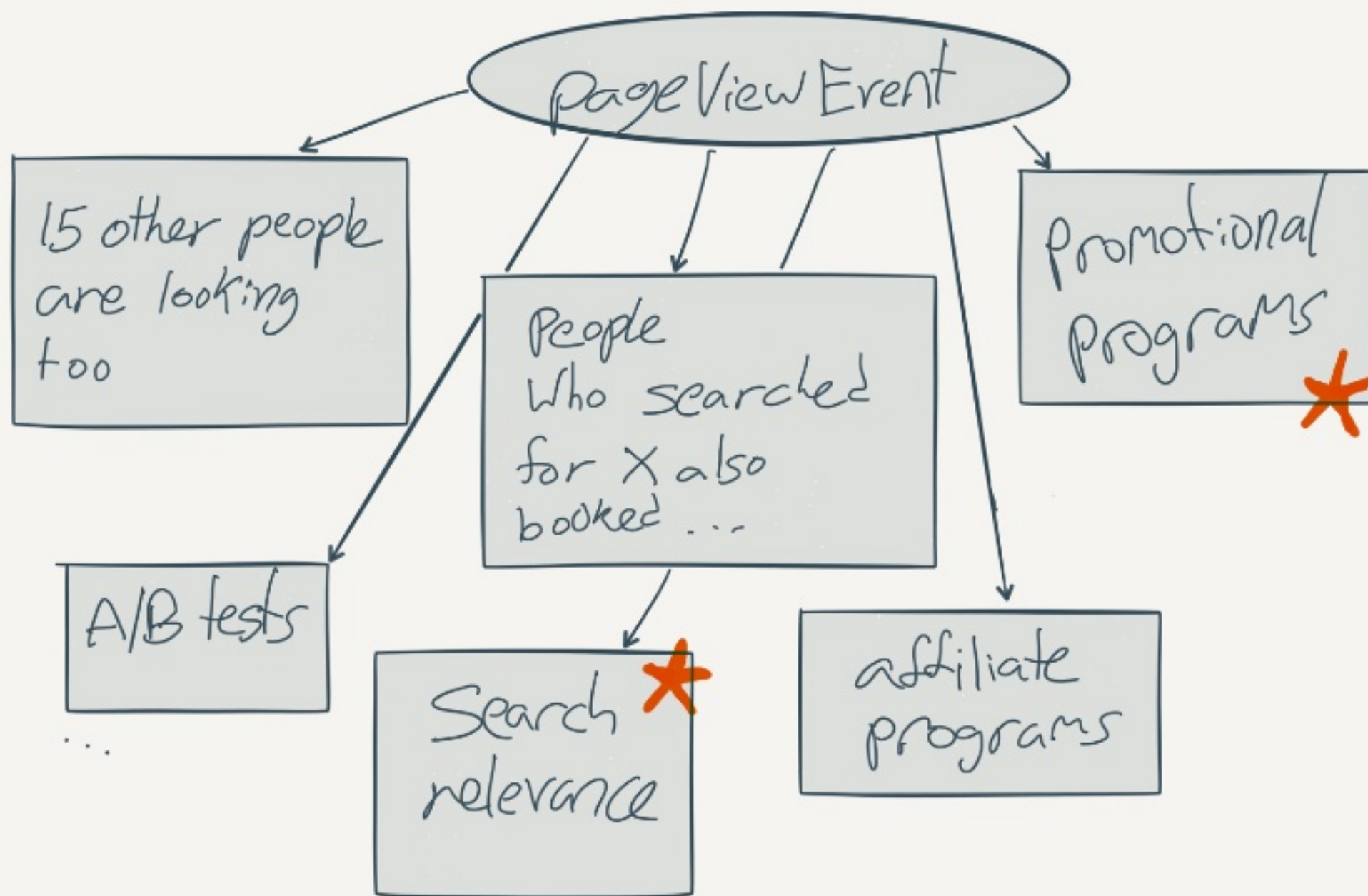
confluent

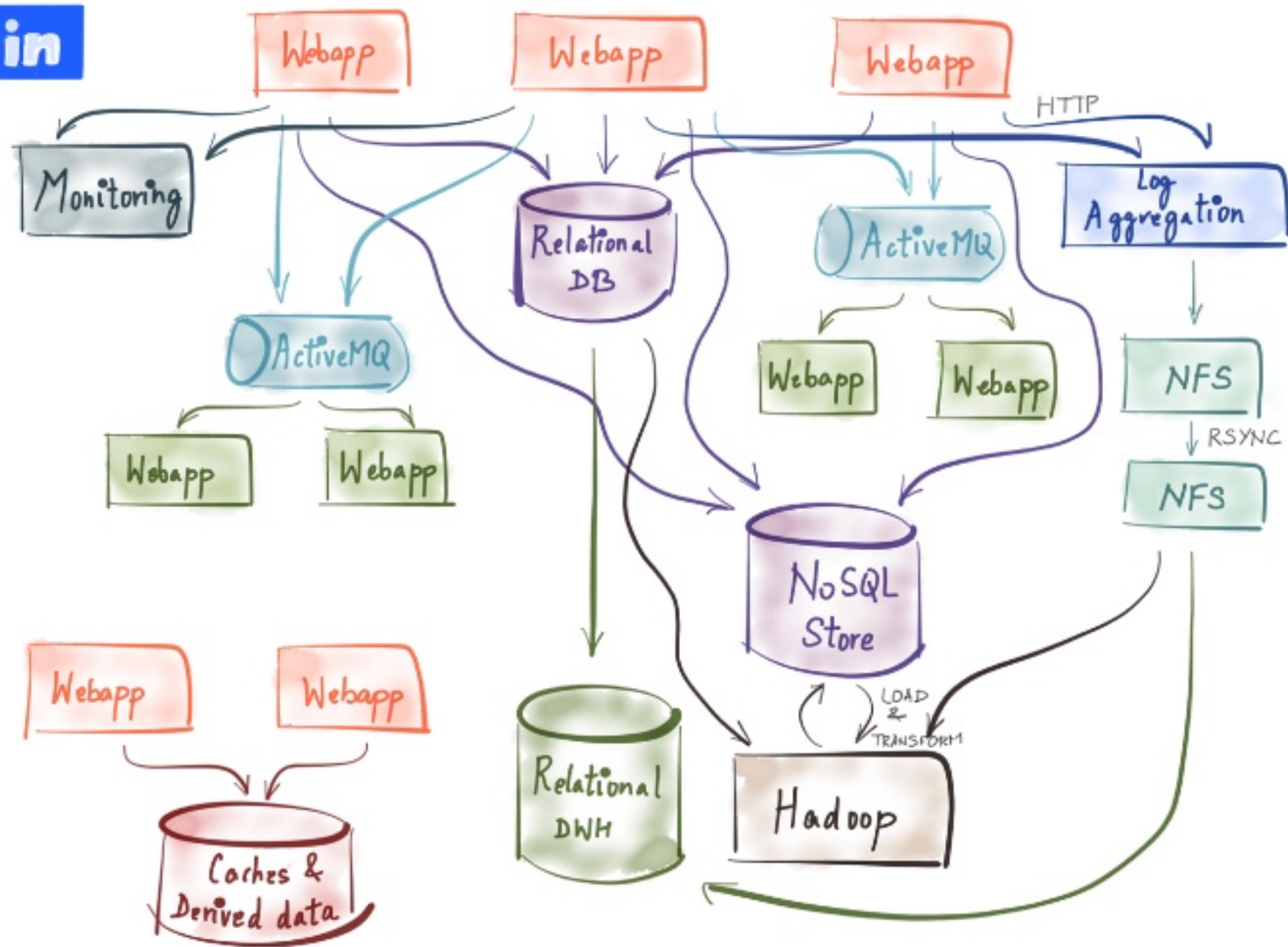# Example: Large Hotel Chain

BOLD CLAIM:
ALL YOUR DATA
IS
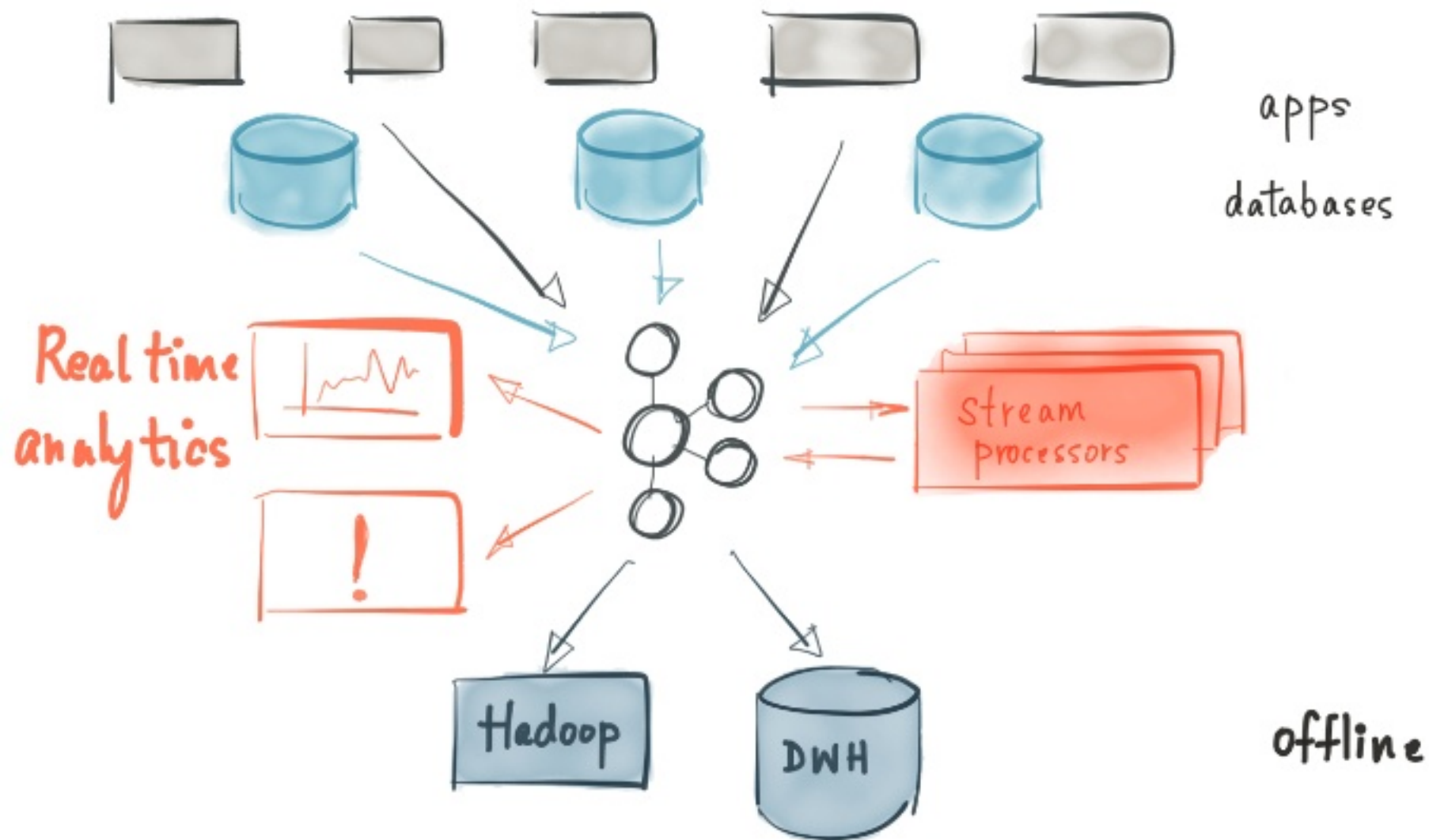EVENT STREAMS

page View Event

# PageViewEvent

```
{
  sessionId: 676fc8983gu563,
  timestamp: 1413215458,
  viewType: "propertyView",
  propertyId: 7879,
  loyaltyId: 6764532
  origin: "promotion",
  ...... lots of metadata....
}
```
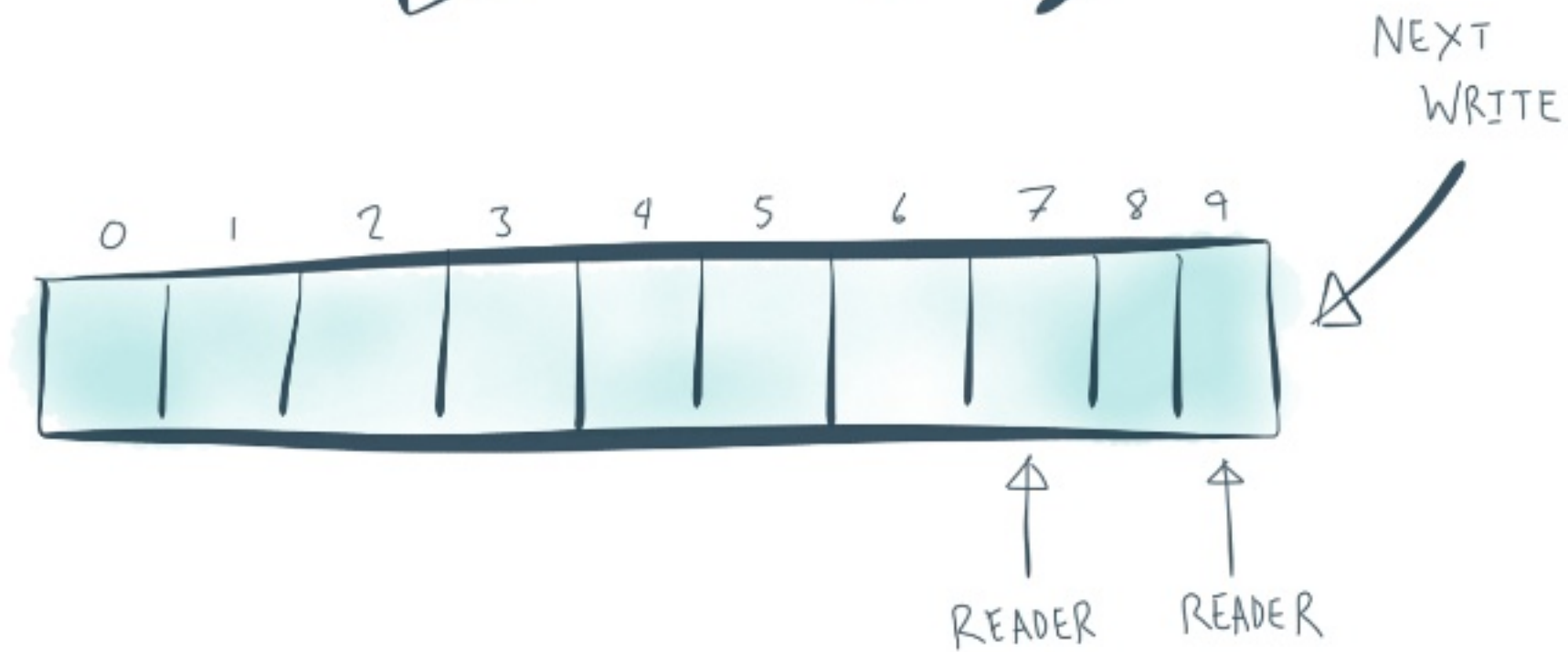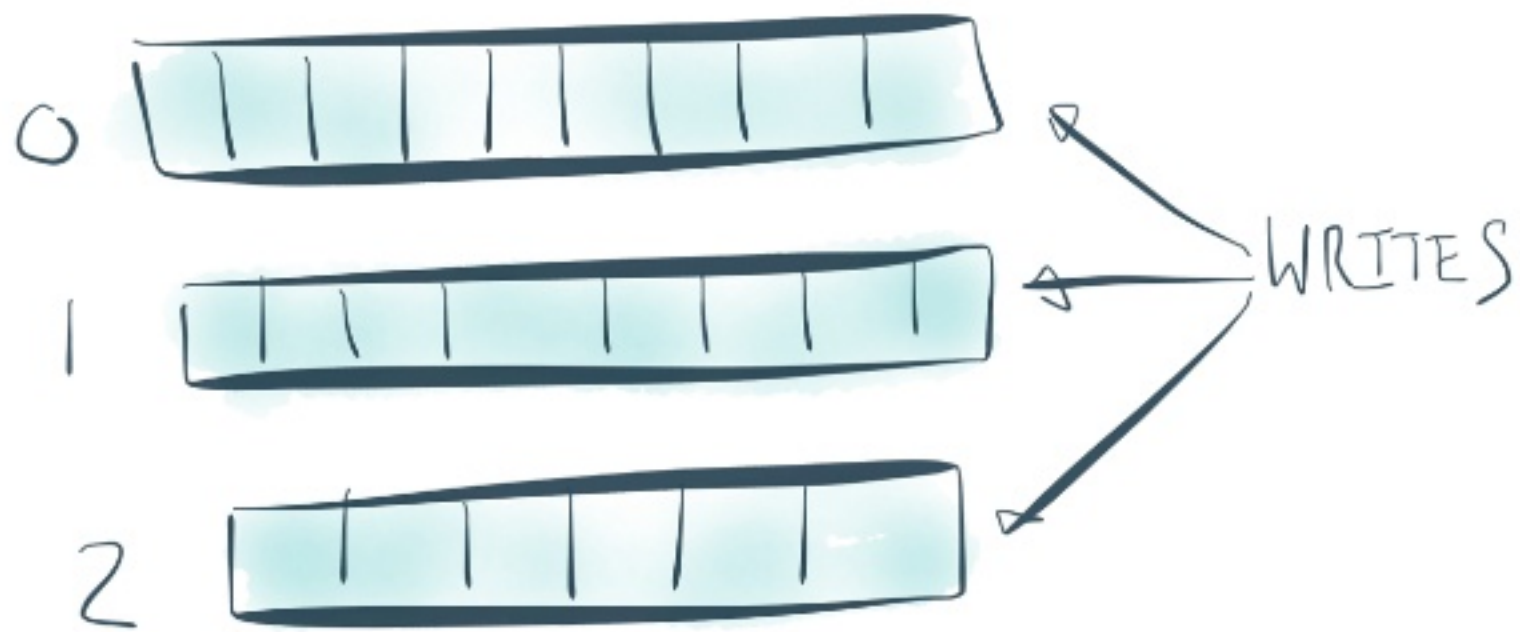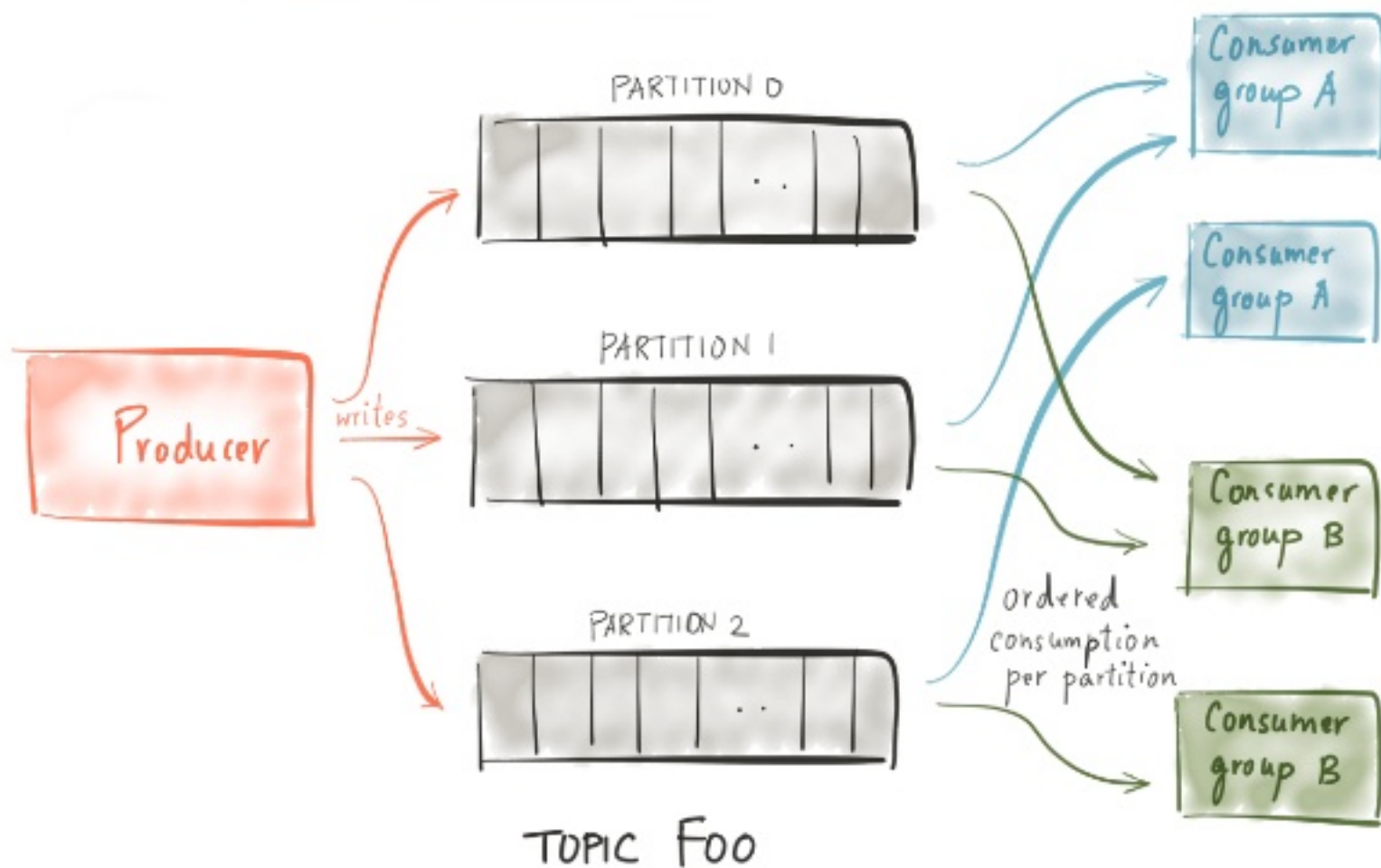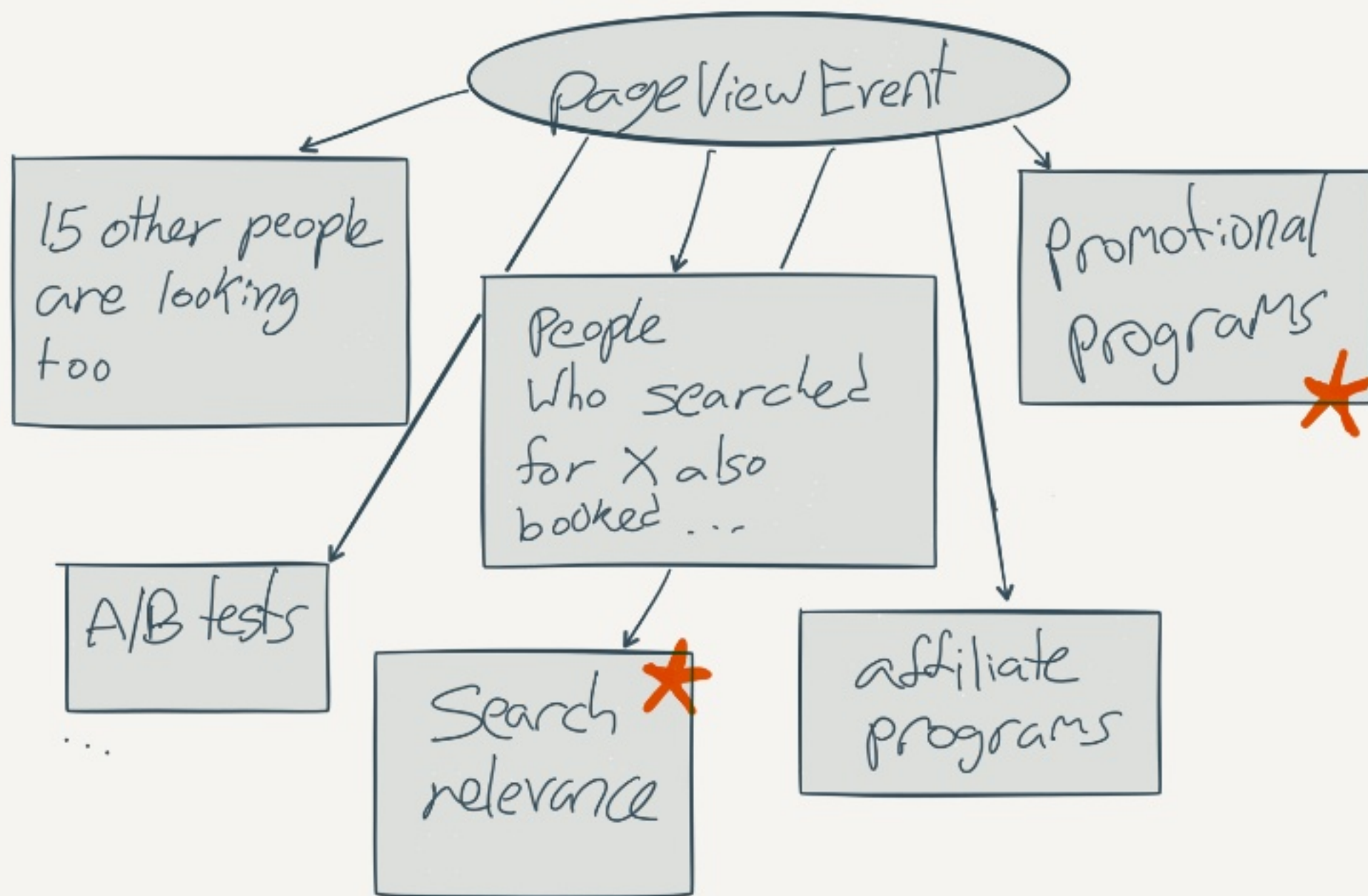
# Pattern #1 – Put all streams in One Kafka



apps
databases

Real time analytics

Stream processors

Hadoop

DWH

offline

# Common Story

PageView

{
  timestamp: 141553216/8

  - - - - }



creatTable
(timstamp number)

Promo

Search

new Date (timestamp)

# Common Story

PageView

{ timestamp: ~~1415532118~~
"Jan 20, 2017"

- - - - }

~~creatTable~~
(timstamp nvmber)

Promo
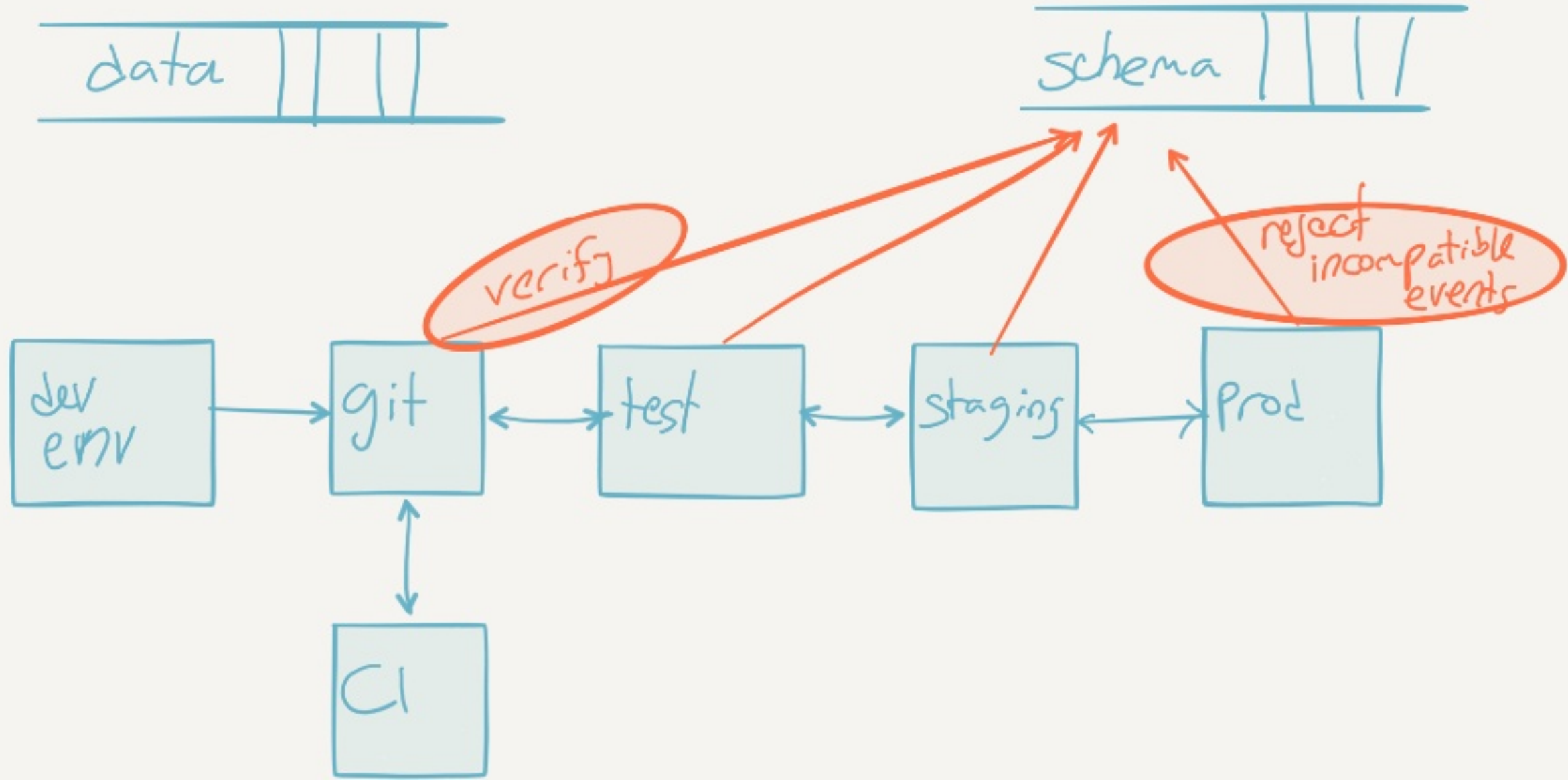
Search

new ~~Date(timestamp)~~

# Pattern #2

APIs between services are Contracts
In Stream Processing World – Event Schemas ARE the API

…except they stick around for a lot longer

# keep events compatible

data | | | |

schema | | | | |

*verify*

*reject incompatible events*

dev env → git ↔ test ↔ staging ↔ prod

git ↕ CI

# Schema Registry

App 1

Serializer

!

App 2

Serializer

!

Kafka Topic

Schema Registry

AVRO

**Example Consumers**

Elastic

Cassandra

Streams App

| **Define** the expected fields for each Kafka topic | **Verify** Schema Compatibility at every stage of dev cycle |
| Automatically **handle** schema changes (e.g. new fields) | **Supports** multi-datacenter environments |

confluent

# Pattern #3 – Ridiculously parallel data integration

```
{
  sessionId: 676fc8983gu563,
  timestamp: 1413215458,
  viewType: "propertyView",
  propertyId: 7879,
  loyaltyId: 6764532,
  origin: "promotion",
  cc: "4444-3333-2222-1111"
}
```
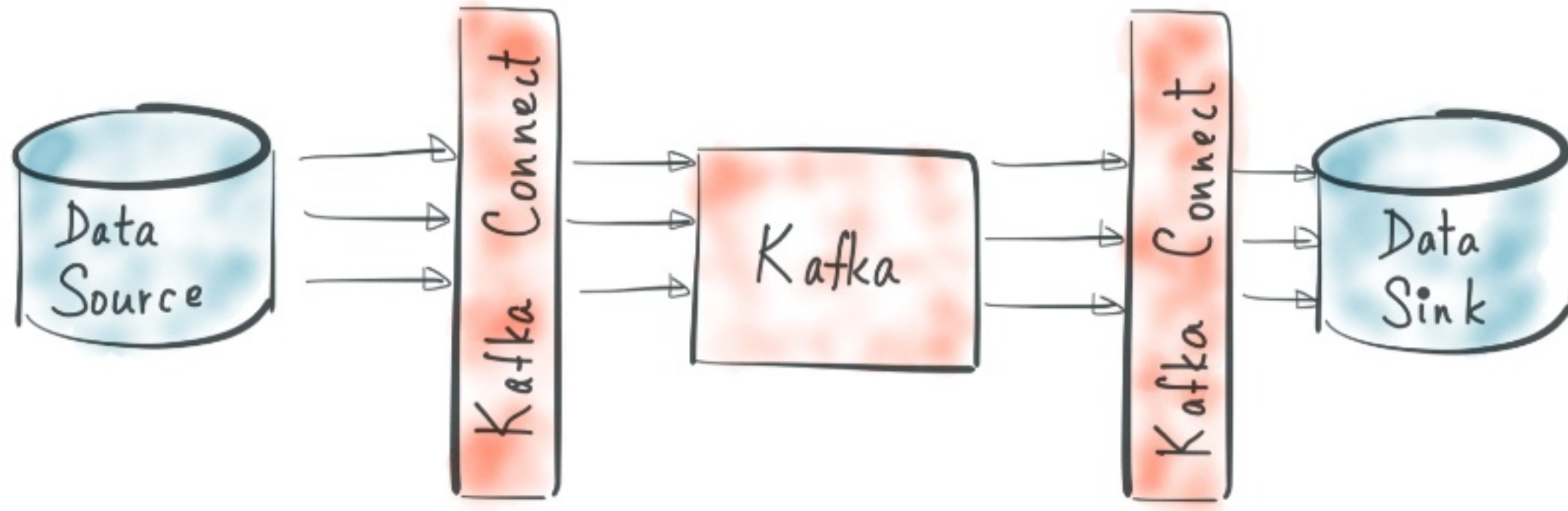
```
{
  sessionId: 676fc8983gu563,
  timestamp: 1413215458,
  viewType: "propertyView",
  propertyId: 7879,
  loyaltyId: 6764532,
  origin: "promotion",
  cc: "xxxx-xxxx-xxxx-1111"
}
```

# Hipster Stream Processing

# Is there a data store involved? KafkaConnect can help
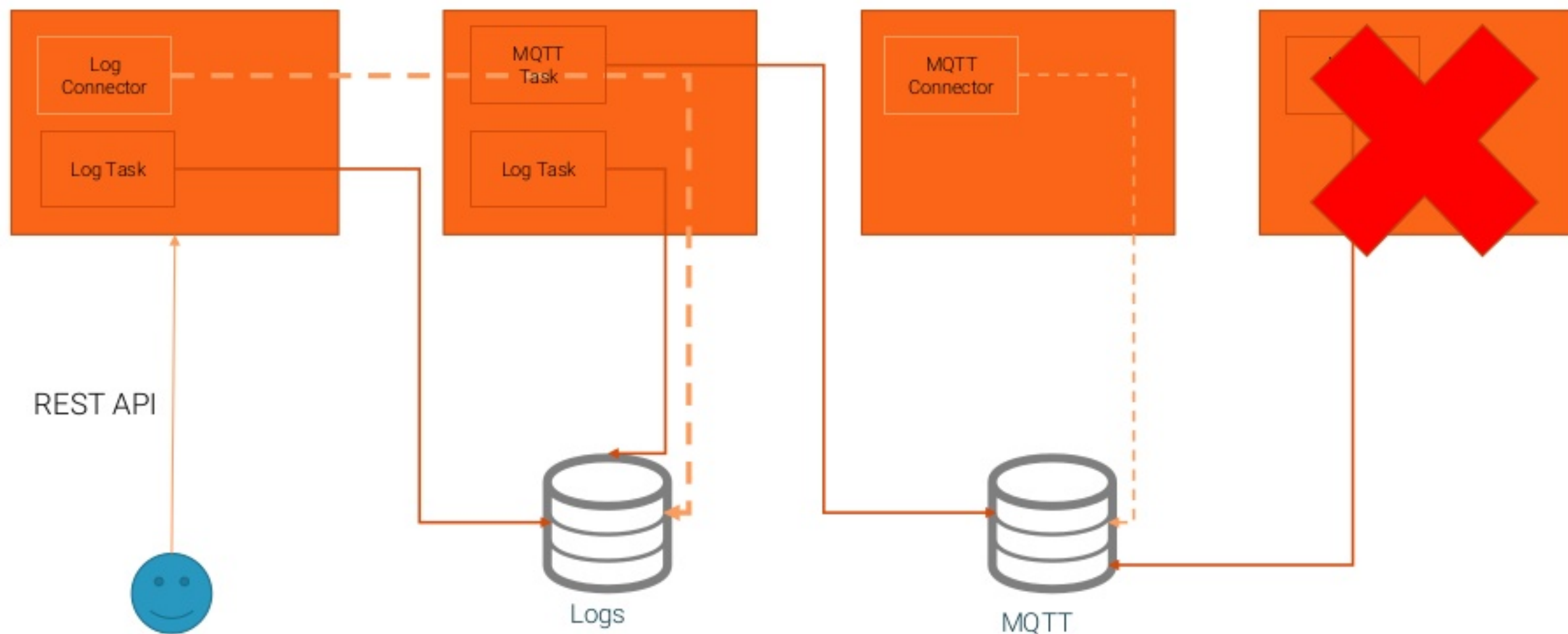
# Ecosystem of Connectors

## Databases

- Java JDBC
- mongoDB
- Cassandra
- DATASTAX
- ORACLE GOLDEN GATE
- Just One Data (Optimizing Data Storage)
- Couchbase
- amazon DynamoDB
- InfluxDB
- Kudu
- APACHE HBASE
- RethinkDB

## Datastore/File Store

- hadoop
- apache Ignite
- FTP
- syslog-ng
- hazelcast

## Analytics

- VERTICA (An HP Company)
- elastic
- mixpanel
- Solr
- APACHE KUDU

## Applications / Other

- amazon web services S3
- Bloomberg
- syncsort
- twitter
- ATTUNITY
- MQTT.ORG

# How Connect Works?



REST API

Logs

MQTT

Mix:

Ridiculously parallel database integration
With
Ridiculously parallel simple transformations

# Promotion!

We want to send Platinum members

Who are looking at beach properties
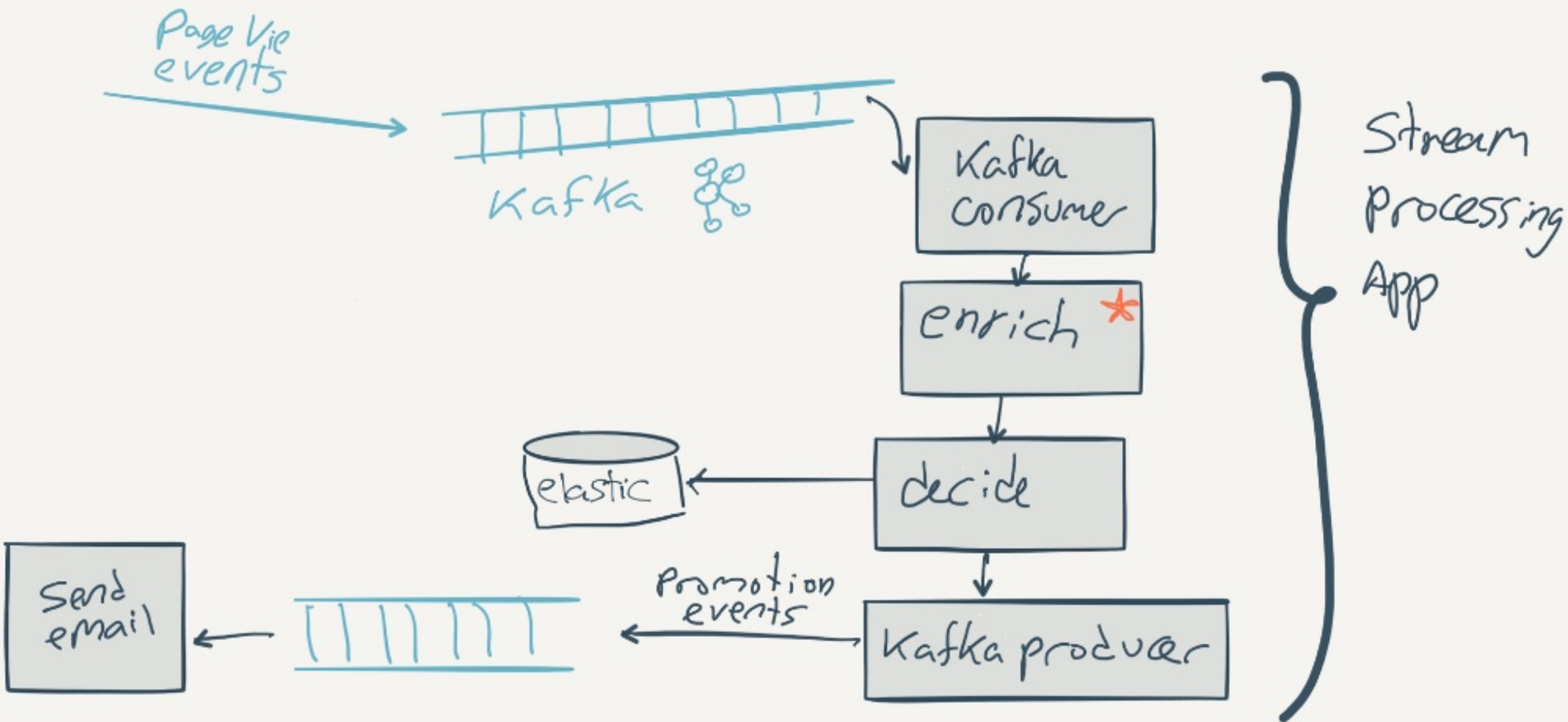
An email about

Discount package in new Hawaii hotel
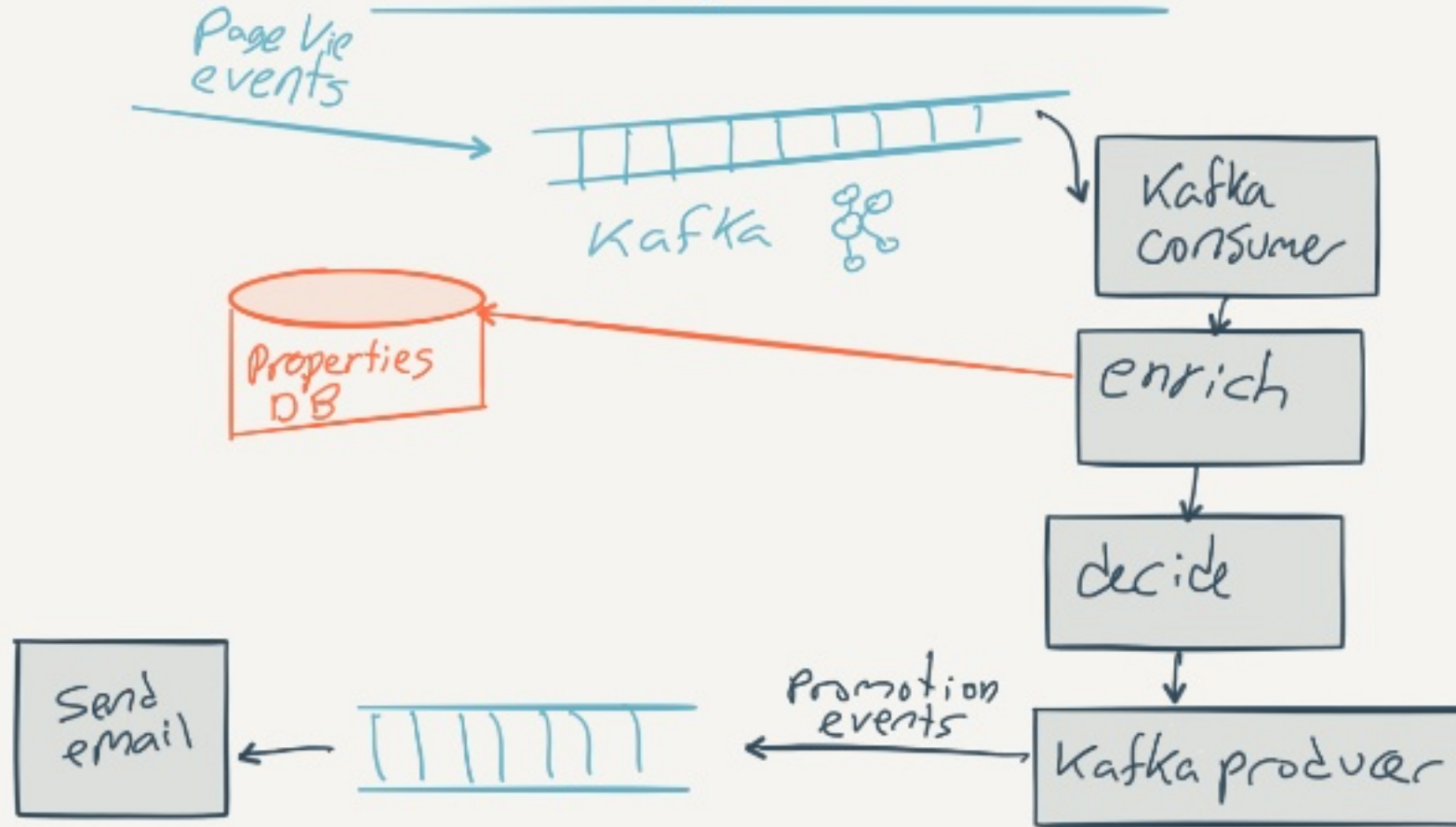
# Pattern #4 – Enriching events

**From Event**

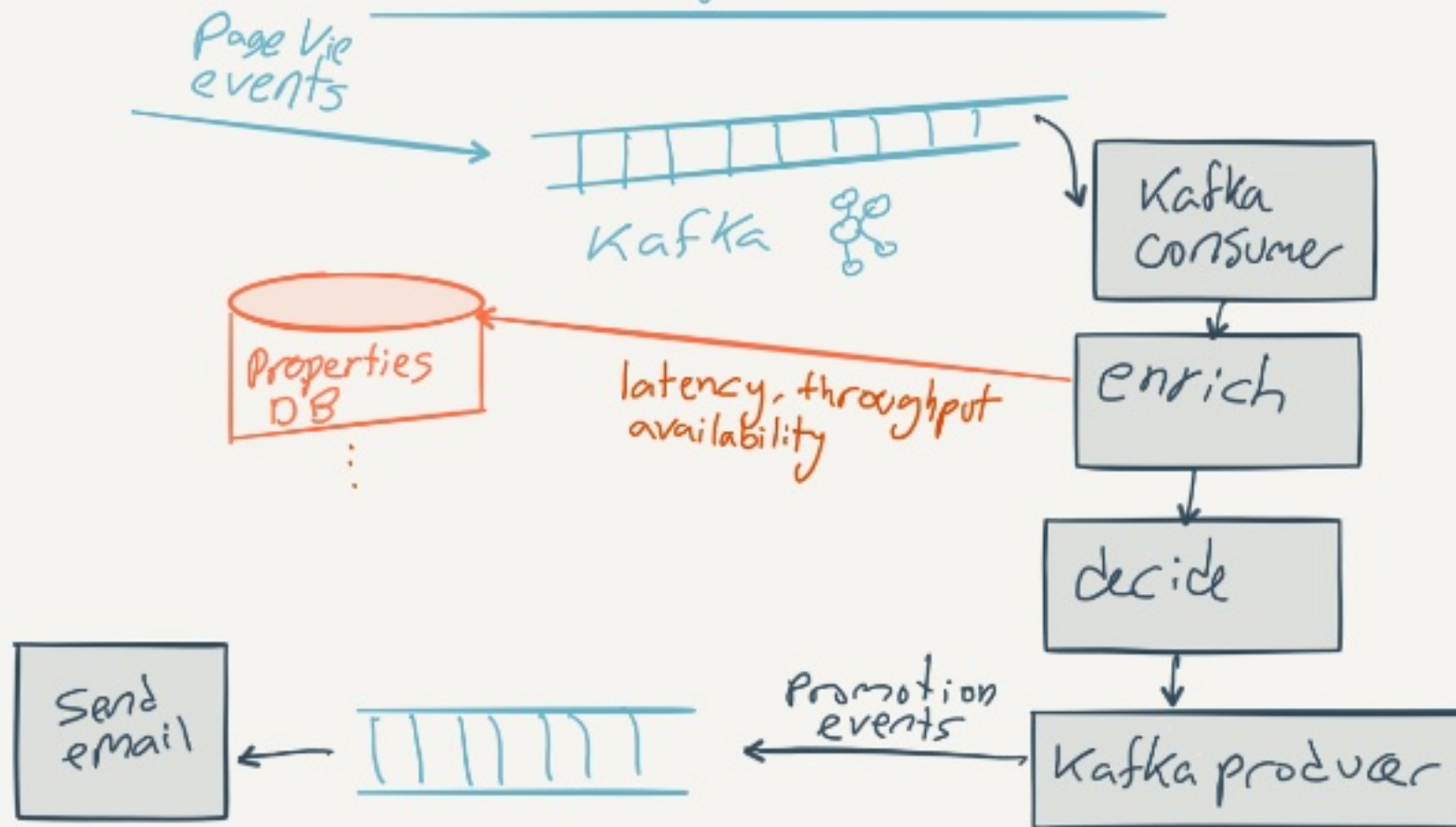sessionID

timestamp

loyaltyID

propertyID

**To Enriched Event**

sessionID

timestamp

loyaltyID

Level

isPlatium

propertyID

Location

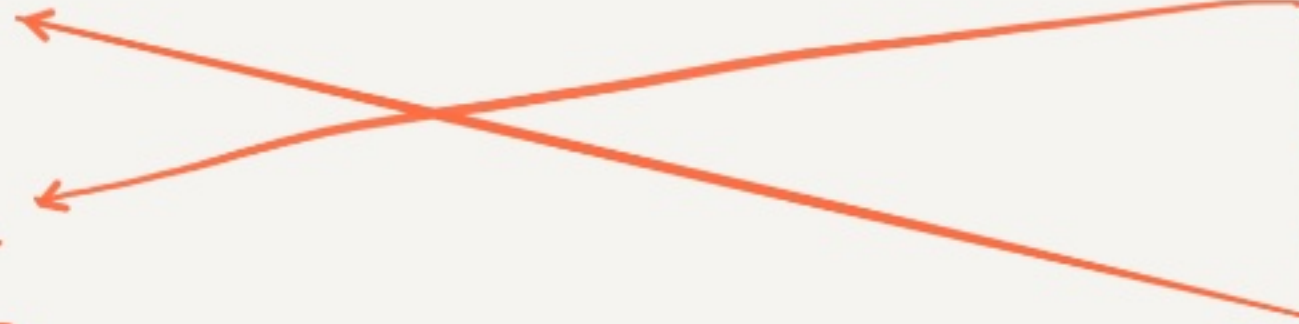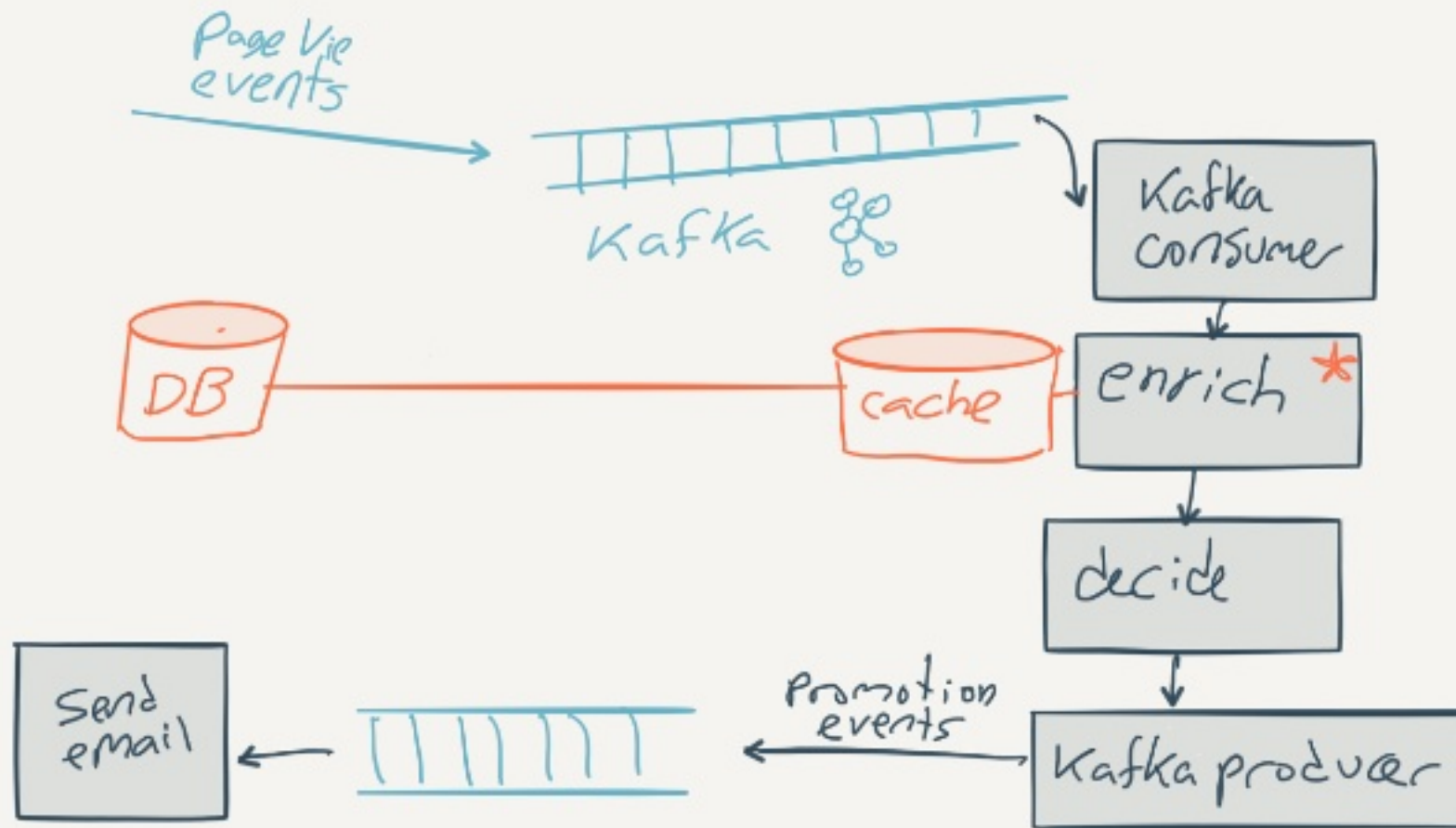isBeach

Page View events

Kafka

Kafka consumer

enrich *

decide

elastic

Send email

Promotion events

Kafka producer

Stream Processing App

# Attempt #1

Page Vie events →

Kafka 🔗

Kafka Consumer

Properties DB

enrich

decide

Promotion events

Kafka producer

Send email ← | | | | | | |

# Attempt #1 — Issues

Page Vie events

Kafka

Properties DB
:

latency, throughput
availability

Kafka consumer

enrich

decide

Send email

Promotion events

Kafka producer

# Stream/table Join

| ID  | level    |
|-----|----------|
| 1   | silver   |
| 2   | platinum |
| 3   | Gold     |
| :   | :        |
| :   | :        |

ID = 2

ID = 1

ID = 53

# Attempt #2

Page Vie
events

Kafka 🐙

Kafka
Consumer
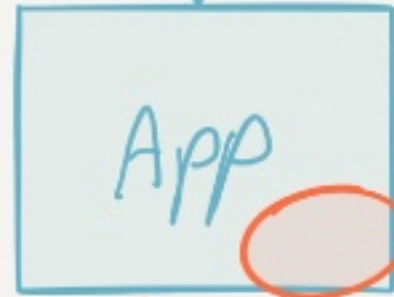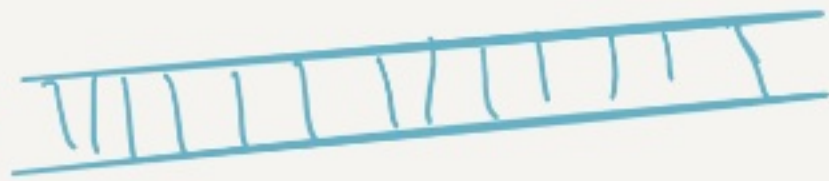
DB

cache

enrich *
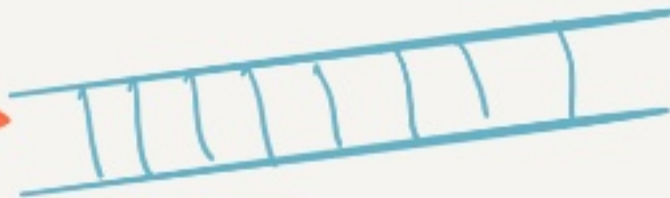
decide

Send
email

Promotion
events

Kafka producer

PageView Event

App

DB

CDC

{loyaltyID, Level}

update profile
set level=platinum
where loyaltyID=538

We use Kafka, Connect and CDC Connectors

To turn state (in DB)

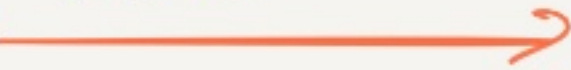Into stream of events (in Kafka)
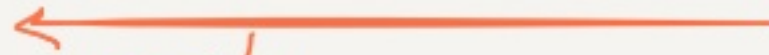
And back into state (in application cache)

# In reality…

- Maintaining the distributed cache isn't trivial
- How do we persist state? Failover? Recover?
- How do we co-partition for joins?

KafkaStreams makes it easy:

```
Ktable loyaltyTbl = builder.table(…, "loyalty-cdc-topic")
PageViewStream.leftJoin(loyaltyTbl, <operation>)
```

Example:
https://www.confluent.io/blog/distributed-real-time-joins-and-aggregations-on-user-activity-events-using-kafka-streams/

# Few Patterns

1. Stream all things (in one place)

2. Keep Compatible and Process On

3. Ridiculously Parallel Single Message Transformations

4. Streaming Data Enrichment