

Scaling Uber Metric System from Elasticsearch to Pinot

Yupeng Fu, Uber
Nan Ding, Uber

Uber

About Us



Nan Ding

Yupeng Fu (yupeng9@github)

- Principal Engineer @ Uber.Inc
- Real-time Data Platform
- Committer: Apache Pinot

- Staff Engineer @ Uber.Inc
- Mobility & Platforms

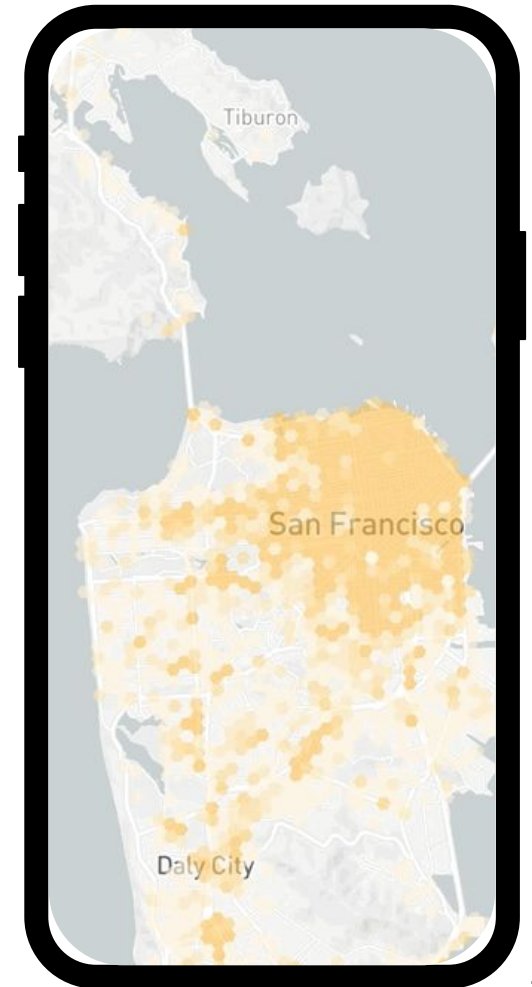
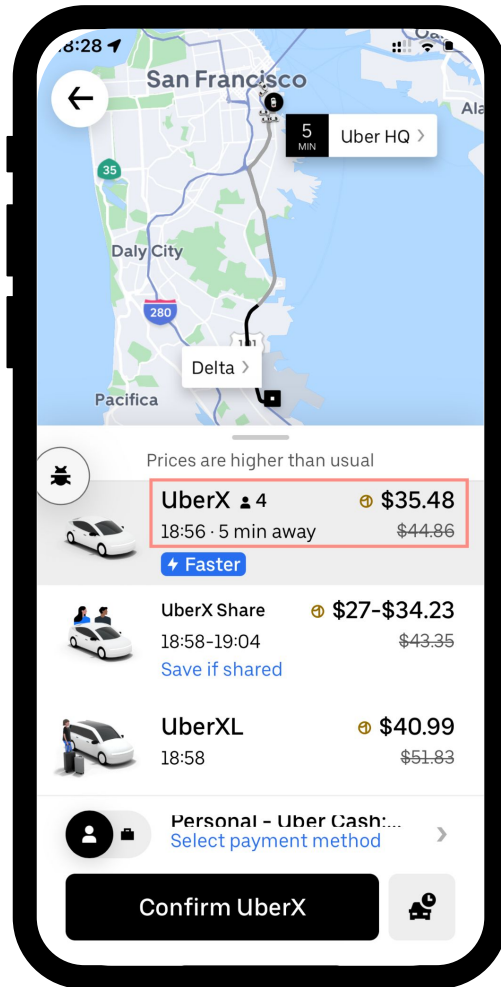
Agenda

- uMetric overview
- Challenges using Elasticsearch
- uMetric over Pinot - Challenge and Architecture
- Pinot at Uber
 - Why Pinot
 - Scale@Uber
 - Features contributed
- Future work

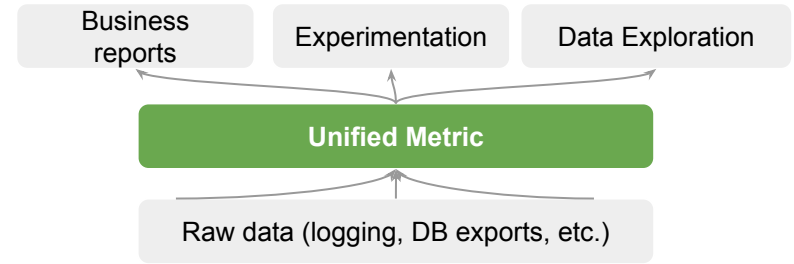
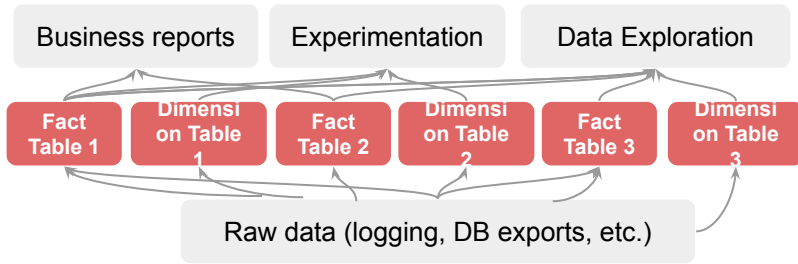
“

Data driven decision making

Capturing a snapshot of the market, and use the data in both real-time and long-retention history for decision-making.



Metric Unification & Standardization



Industry analogs:

Airbnb: Minerva

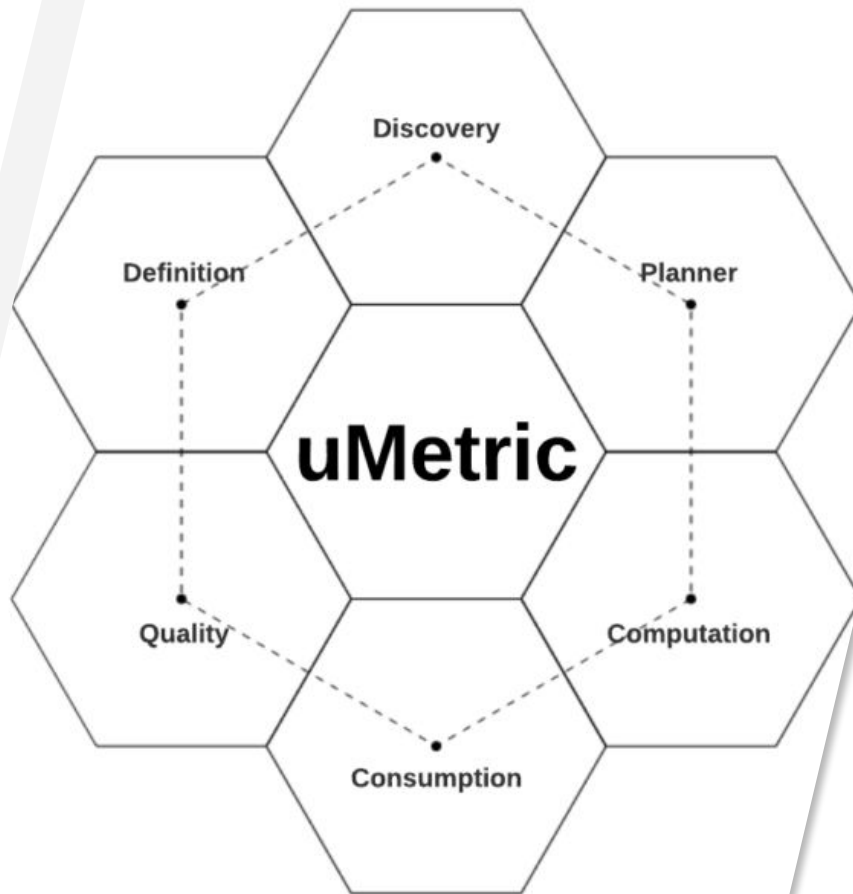
Linkedin: UMP

Transform.co

Liveramp:
Metriql

uMetric

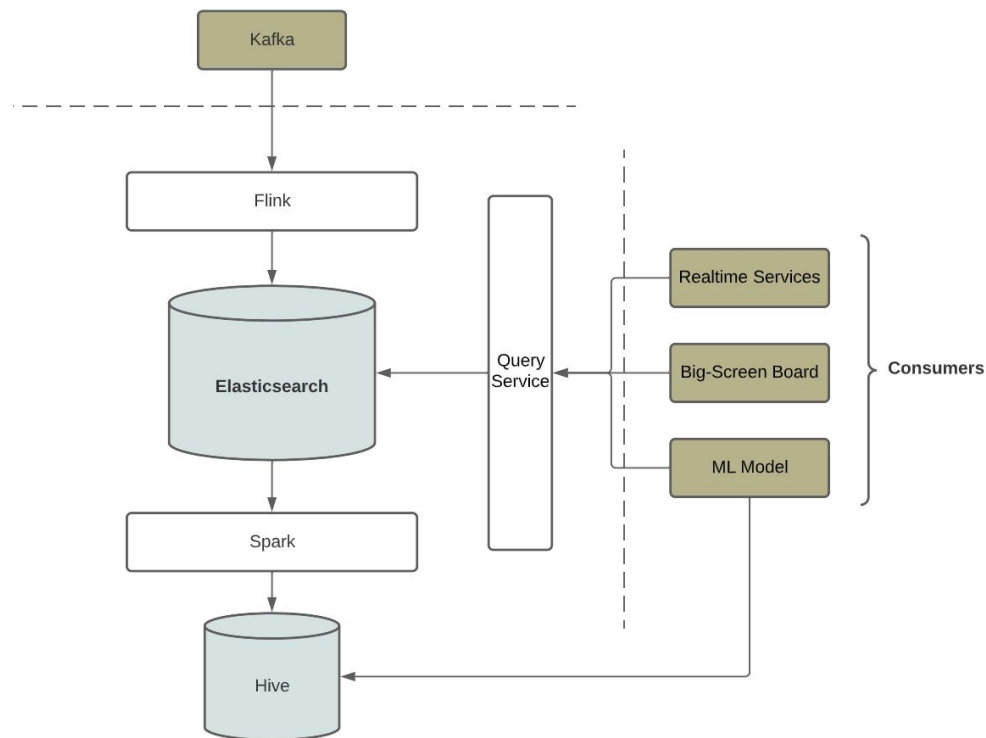
uMetric is a unified metric platform at Uber that manages the full life cycle of a metric: definition, discovery, computation, verification, and serving



Elasticsearch-Oriented Realtime Architecture

Since 2014

- Operable by small team
- Support idempotency insert (upsert)
- Support distributed aggregation
- Linearly scalable
- A matured system



Scale

1.5PB

Dataset Size

4.5T

Documents

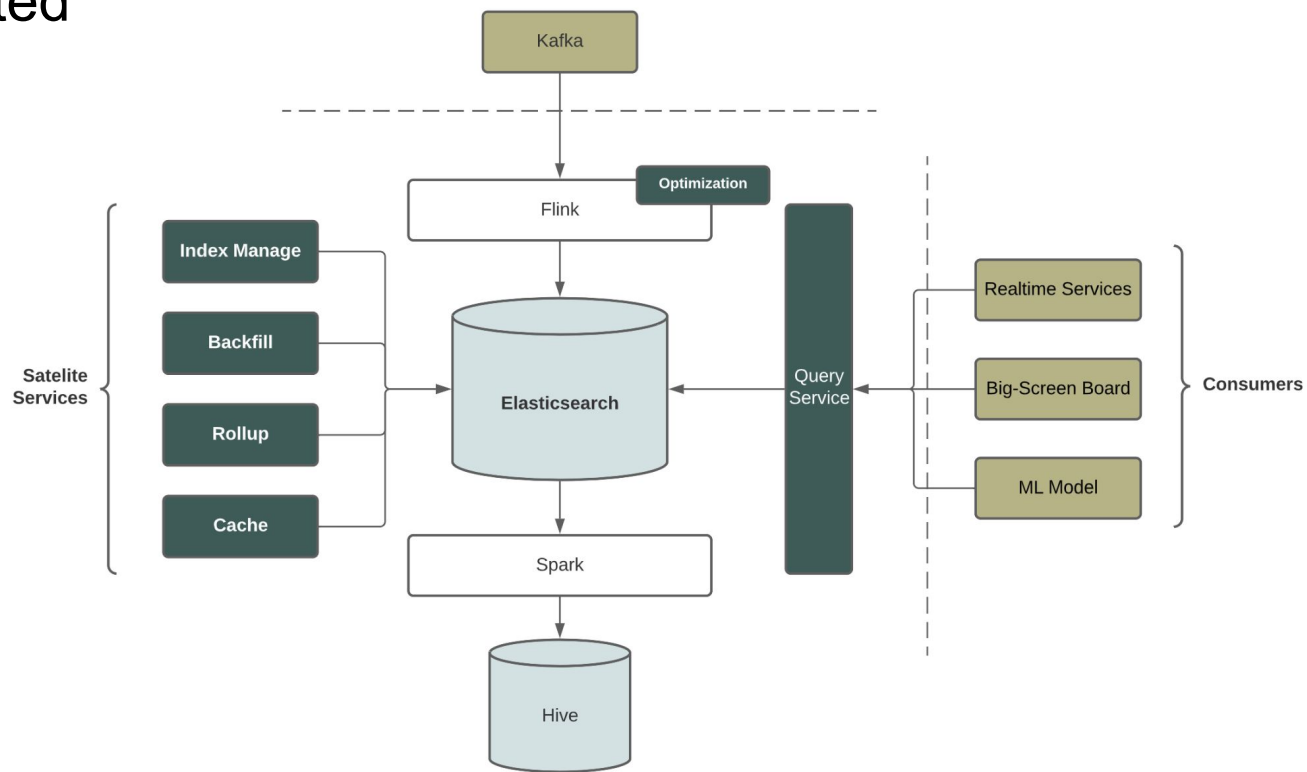
1.3M/s

Write per Sec

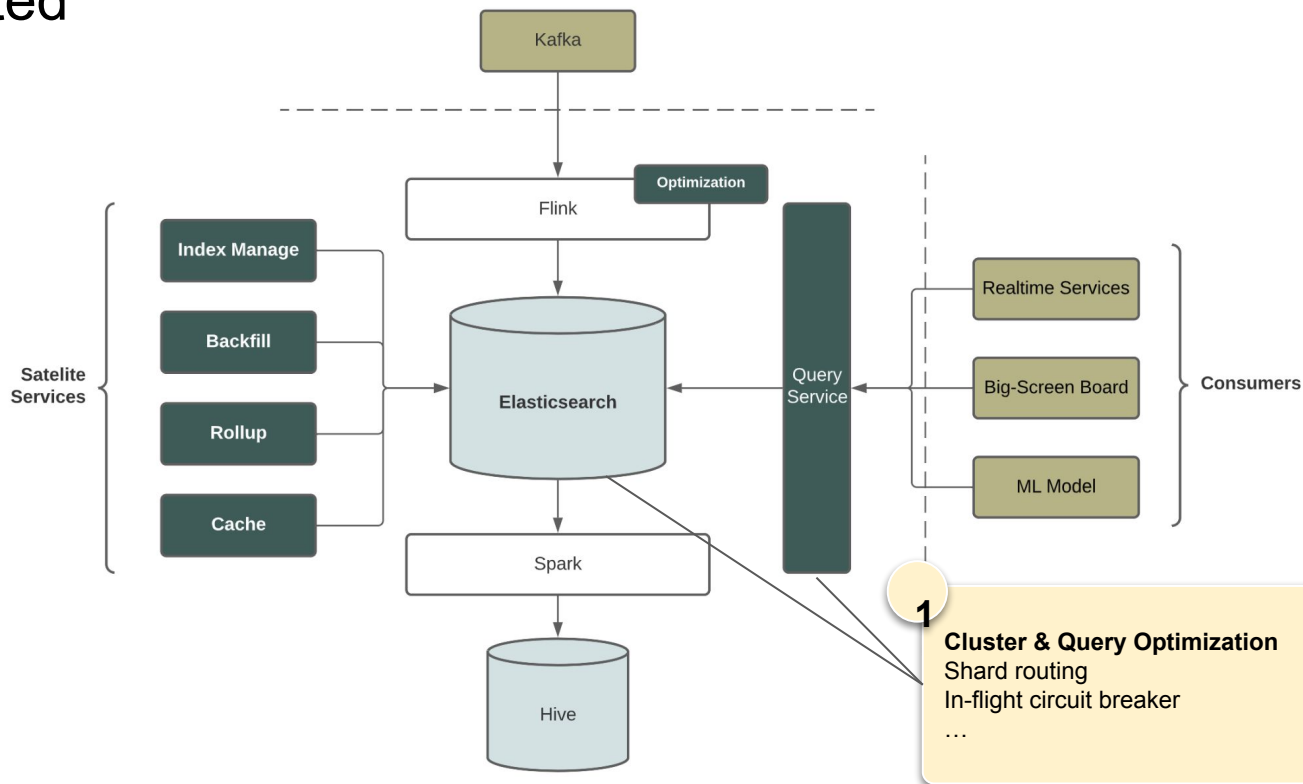
700M/s

Doc Scanned per Sec

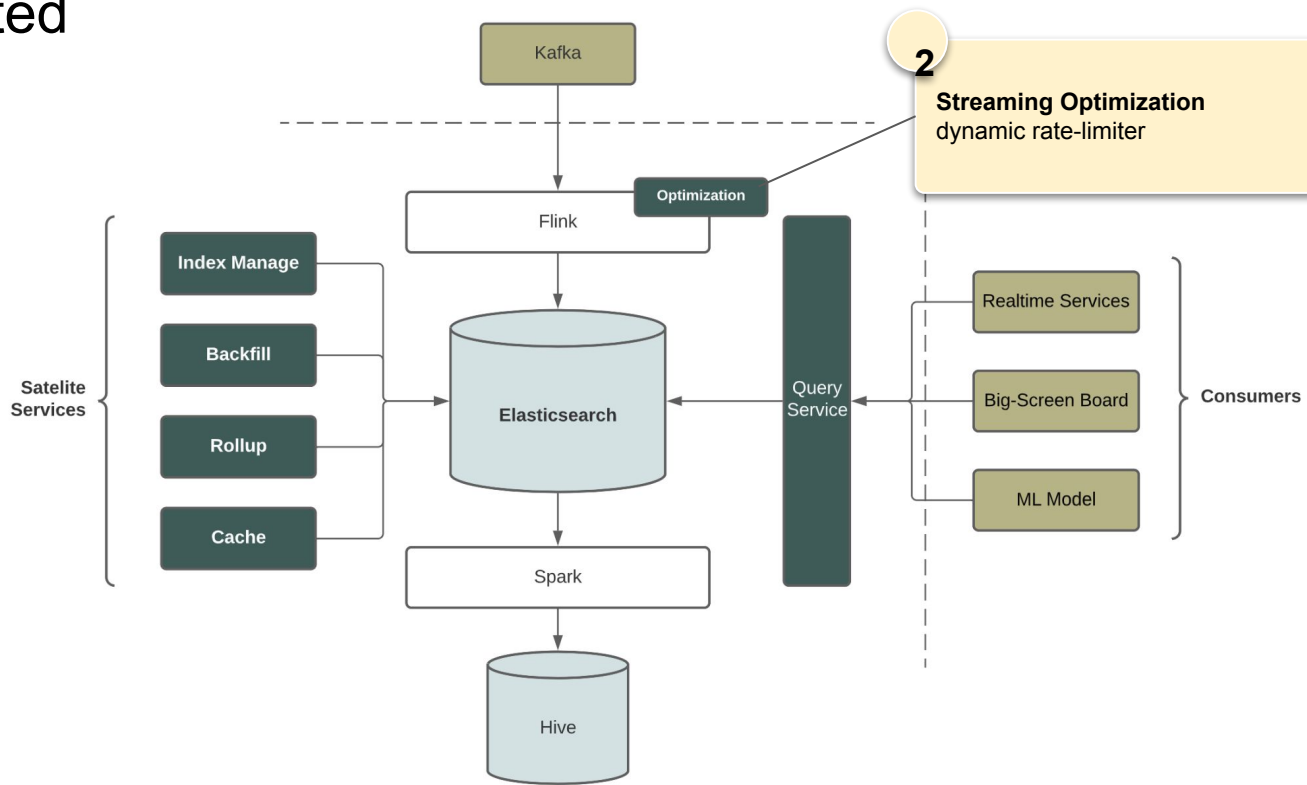
Elasticsearch-Oriented Arch Optimization



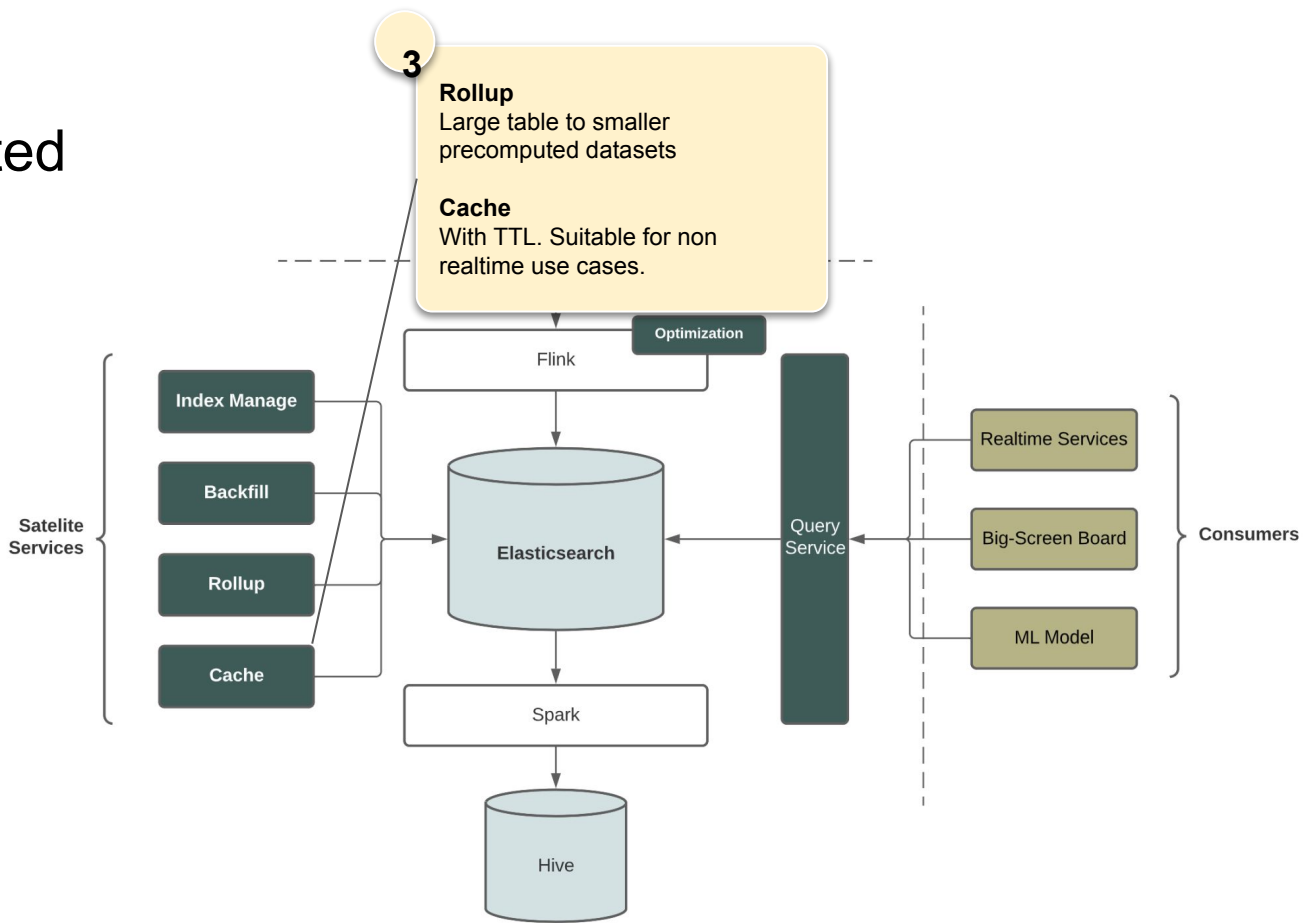
Elasticsearch-Oriented Arch Optimization



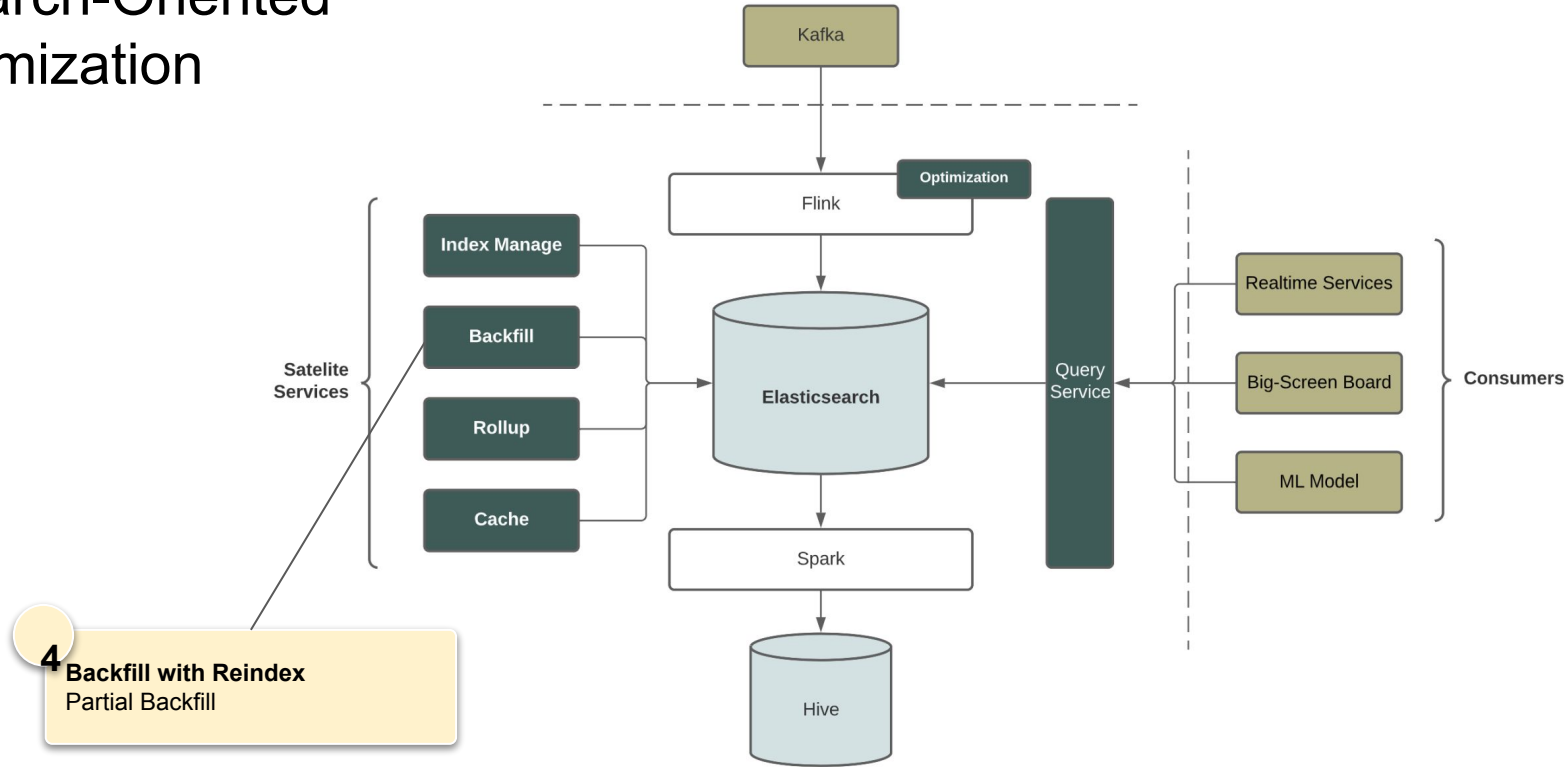
Elasticsearch-Oriented Arch Optimization



Elasticsearch-Oriented Arch Optimization



Elasticsearch-Oriented Arch Optimization



Challenges on Elasticsearch Stack



Reliability

- Cluster constantly in high load
- Frequent SLA breach
- Frequent minor data-loss
- Lack effective mitigation
- Issue cannot be root-caused

Scalability

- Rapidly growing use cases
- Linearly scale cluster / data node has low ROI, and cause derivative issues

Engineering Cost

- Build optimization systems
- High oncall load

Migration toward Pinot



Feature Gaps

- Upsert
- Backfill
- Spark connector
- Nested col support

Safe Migration

- Perf / DQ / Reliability comparison between Pinot and ES
- Drain traffic by table / metric / user / retention / etc.

Performance Tuning

- Multi-cluster by tier
- Dedicated machine spec per cluster
- Fine-tuning per query, with indexes and algorithm optimization

After Migration



Reliability

- Significantly improved overall reliability
- Fine-grained table-level failover

Scalability

- Unblocked new user onboarding
- Scale table up on demand

Engineering Cost

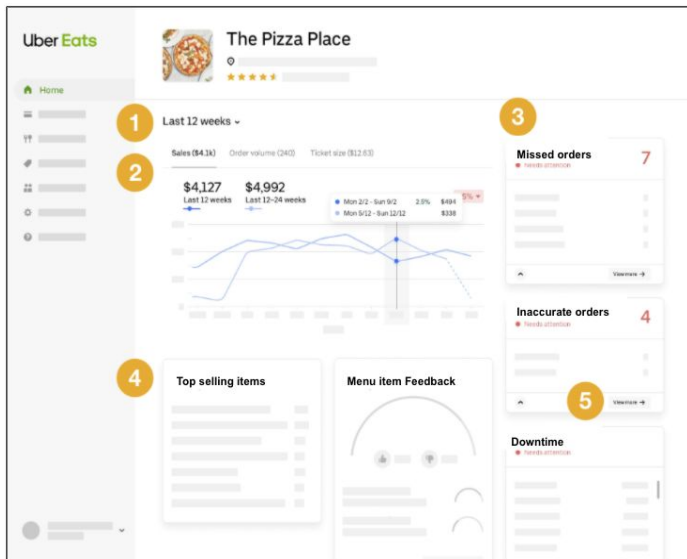
- Deprecated 6 customized system
- Reduced oncall load

Pinot @ Uber

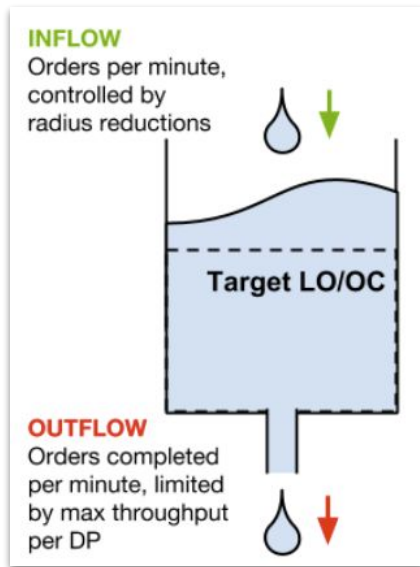
Real-time Analytics at Uber

1. Real-time and actionable insights
2. Time-sensitive decisions
3. User engagement growth

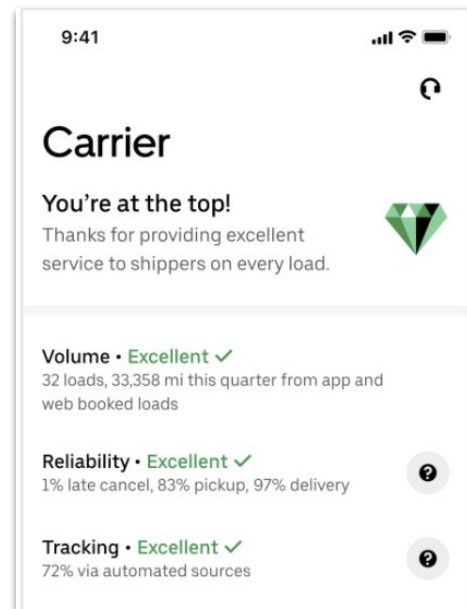
Fast Access to Fresh Data at Scale



Restaurant Performance View



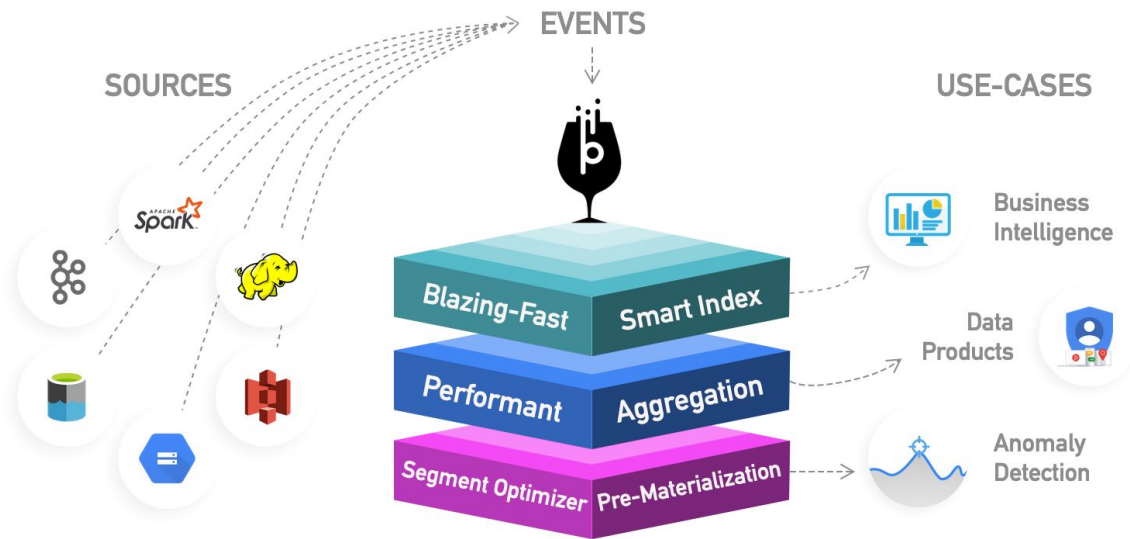
Demand/Supply Management



Freight Carrier Scorecard

Apache Pinot: Fast, Distributed OLAP

- Started by LinkedIn for Metrics System
- Highly available
- Horizontally scalable
- Low latency/High throughput
- Immutable data



1M+

Events/sec

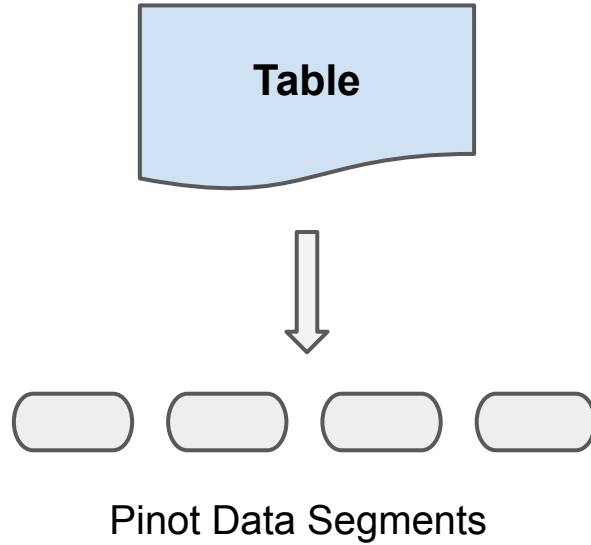
170k+

Peak QPS

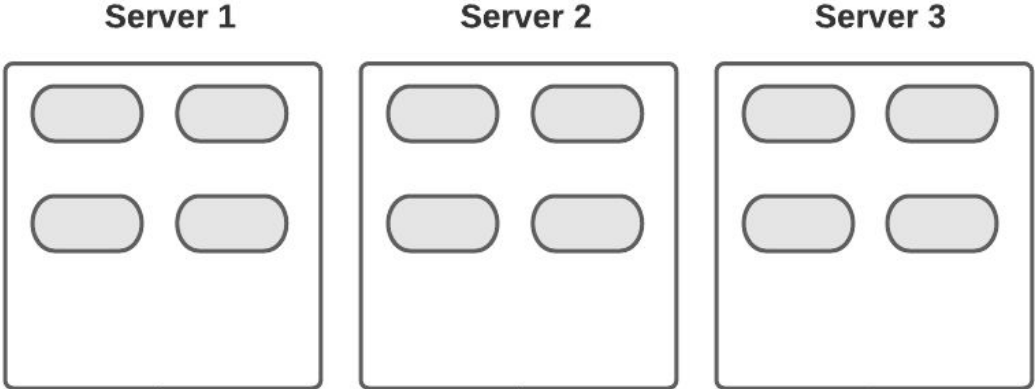
ms

Query Latency

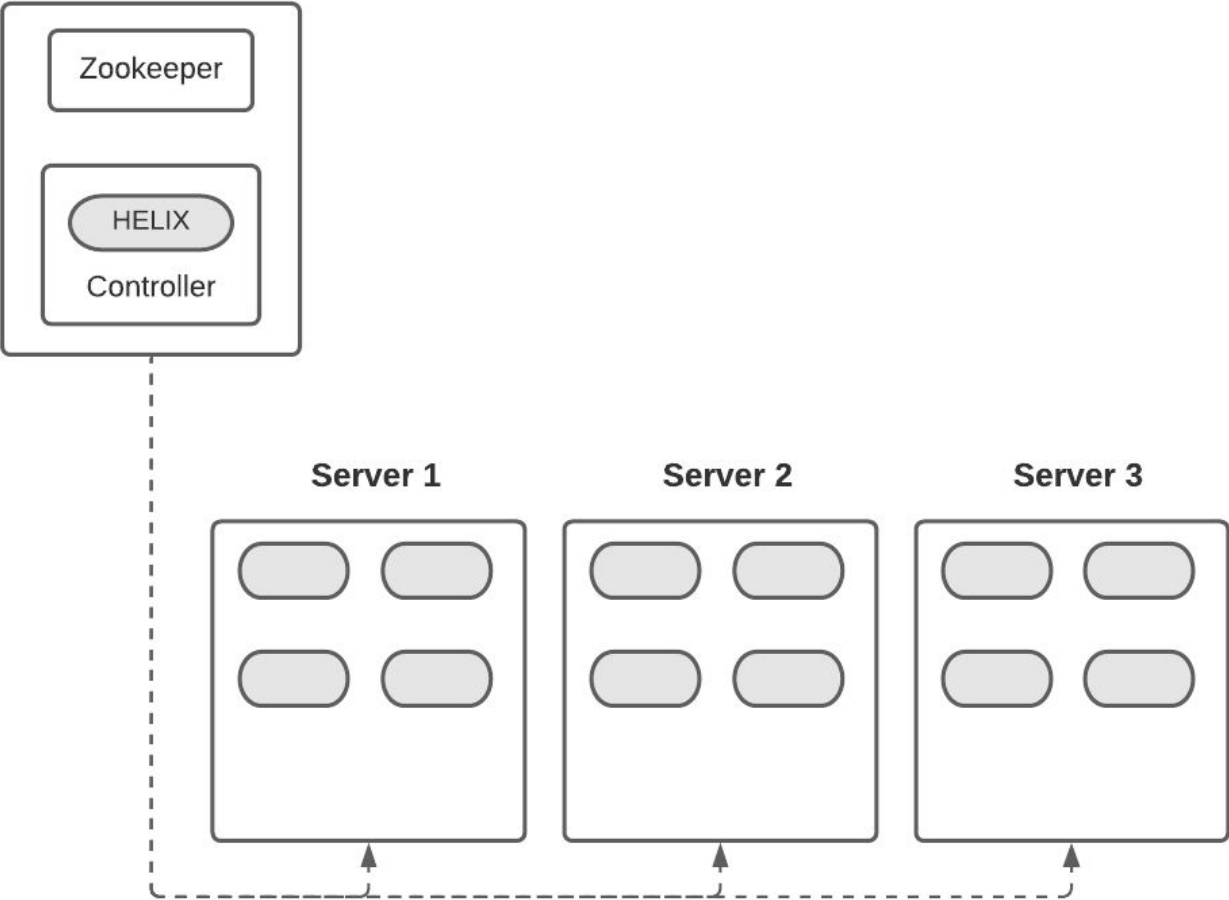
High Level Architecture



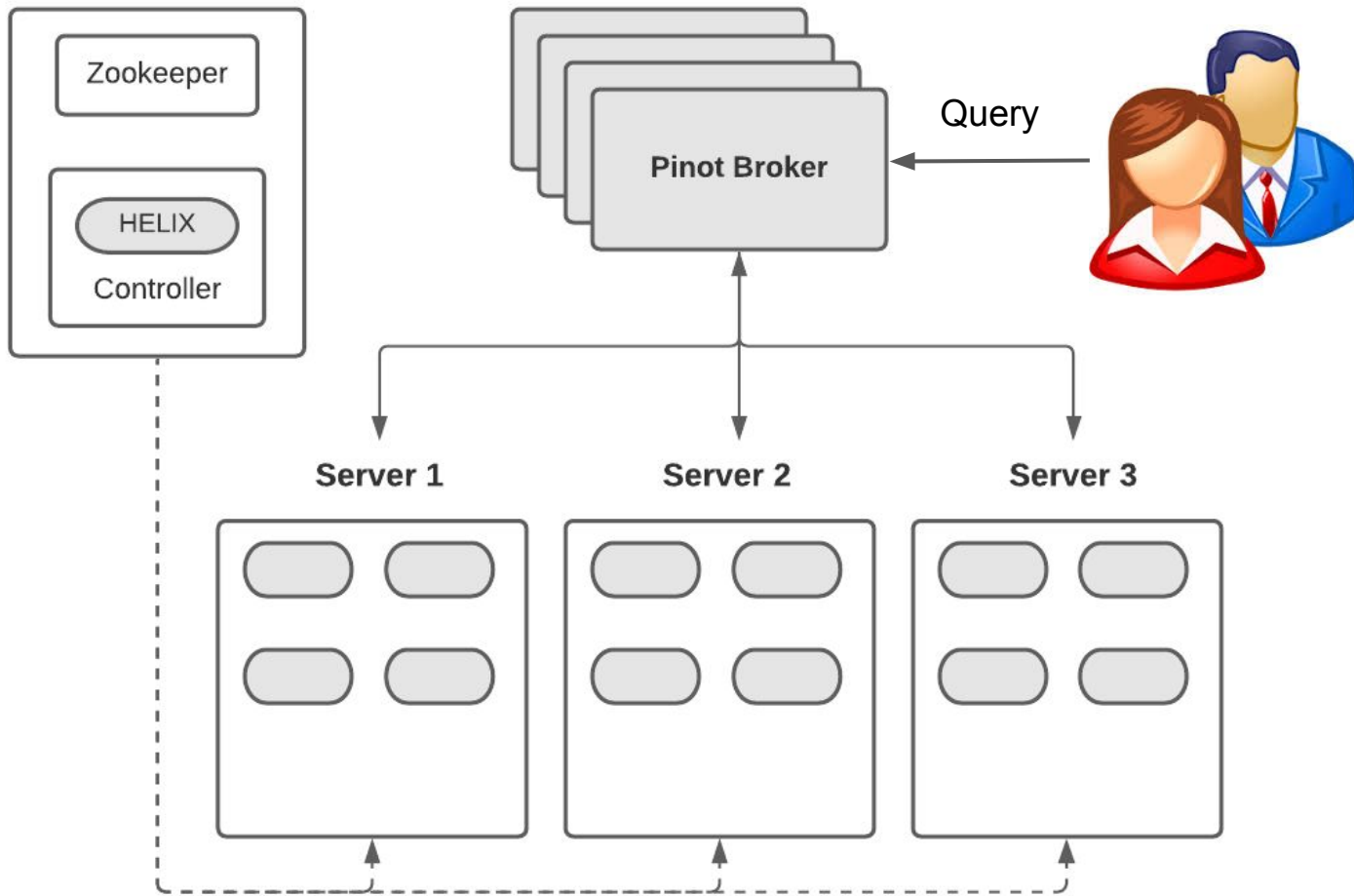
High Level Architecture



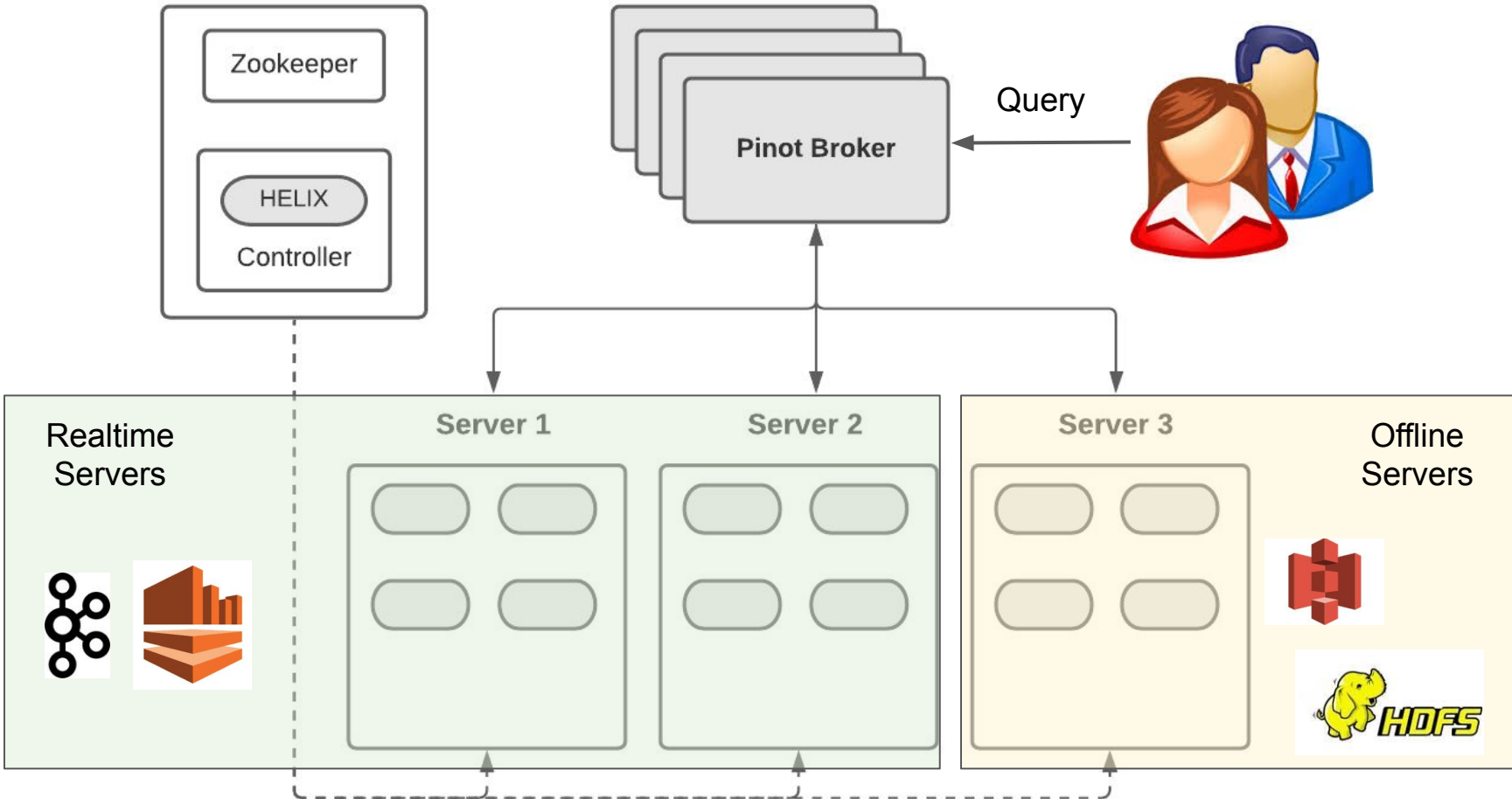
High Level Architecture



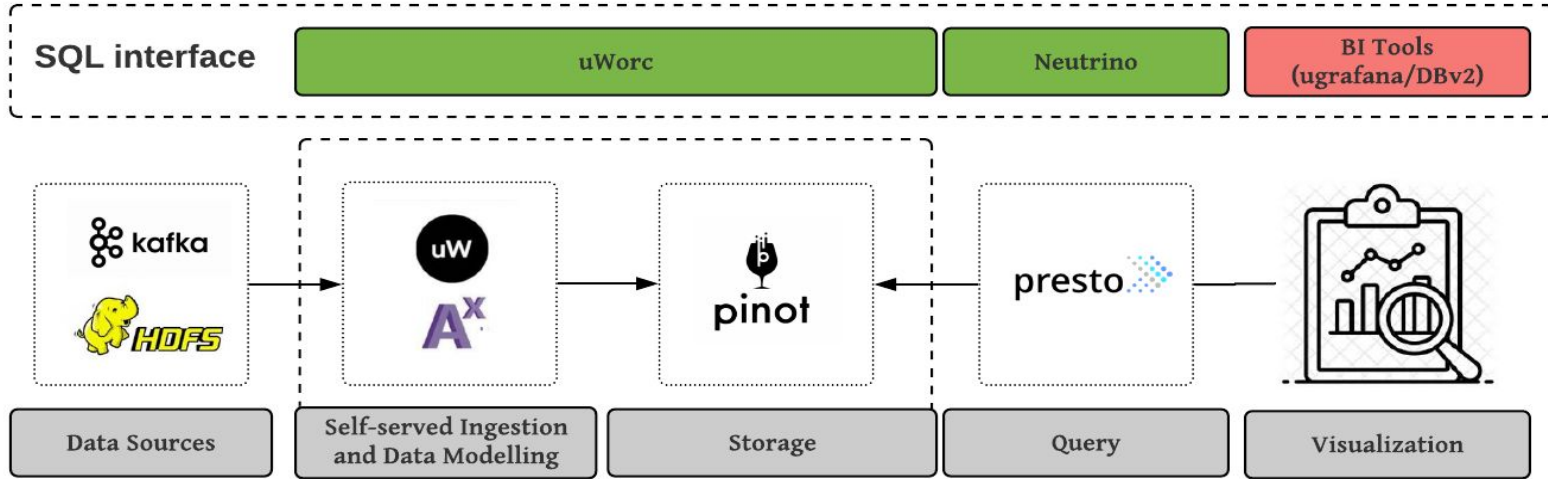
High Level Architecture



High Level Architecture



EVA: Building a World-class RTA platform for Uber



Tier 0 platform, 99.99% uptime

Built on top of **Apache Pinot**

Seconds data freshness

Self service Onboarding (via uWorc)

SQL API (via Presto / Neutrino)

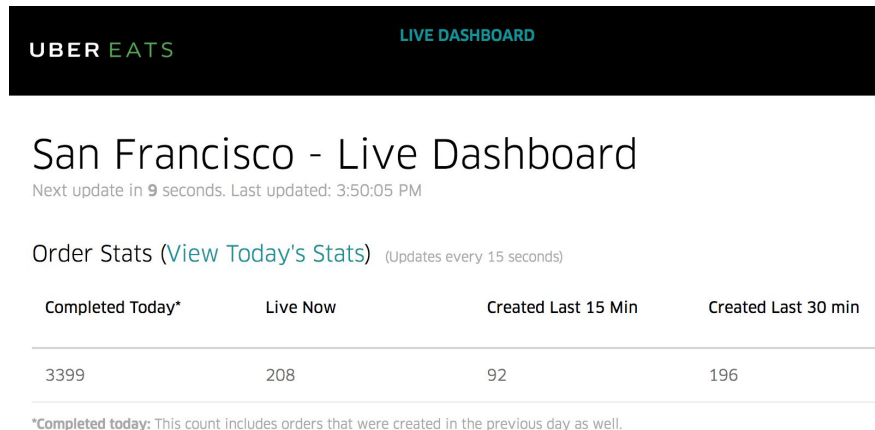
< 100ms @P99 query latency

<https://www.uber.com/blog/operating-apache-pinot/>

Upserts

- First real-time OLAP to support upsert
- Challenge: data stored as immutable segments
- Data partitioning to share nothing
- Released in Pinot 0.6.0
- Partial upsert

```
SELECT current_status,  
       count(*)  
FROM   uberEatsOrders  
WHERE  regionid = 1366  
       AND MinutesSinceEpoch  
       BETWEEN 25432140 AND 25433580  
GROUP BY current_status  
TOP    10000
```



Complex-type (Array, Map) support in Pinot

- Released in Pinot 0.8.0

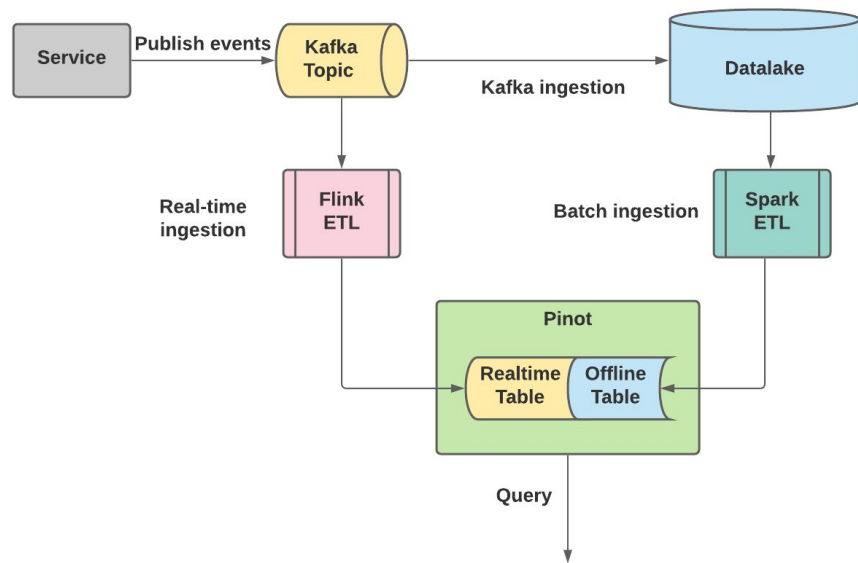
```
▼ object {2}
  rsvp_id : 1869661474
  ▼ group {2}
    ▼ group_topics [2]
      ▼ 0 {2}
        urlkey : paddling
        topic_name : Paddling
      ▼ 1 {2}
        urlkey : hiking
        topic_name : Hiking
    group_id : 28088353
```



```
▼ array [2]
  ▼ 0 {4}
    rsvp_id : 1869661474
    group.group_topics.urlkey : paddling
    group.group_topics.topic_name : Paddling
    group.group_id : 28088353
  ▼ 1 {4}
    rsvp_id : 1869661474
    group.group_topics.urlkey : hiking
    group.group_topics.topic_name : Hiking
    group.group_id : 28088353
```

Spark/Flink connectors with Pinot

- Spark Connector
 - Pinot -> Hive dispersal
 - Perf improvement via GRPC
- Flink Connector
 - Streaming/batch unification for Pinot ingestion
 - Backfill large Pinot Upsert tables



Next

- Leverage StarTree index to replace custom rollup pipelines
- Better cluster isolation and tiering for ease of deployment
- Upsert table compaction / TTL

Uber

Q&A

Apache®, Apache Kafka®, Apache Flink®, Apache Pinot®, Apache Hadoop® and their logo are either registered trademarks or trademarks of the Apache Software Foundation in the United States and/or other countries. No endorsement by The Apache Software Foundation is implied by the use of these marks.

How upsert works

