# Swimming in the data river

## Or, when "streaming analytics" isn't

**Gian Merlino**
gian@imply.io

# Who am I?

Gian Merlino

Committer & PMC member on druid

Cofounder at imply
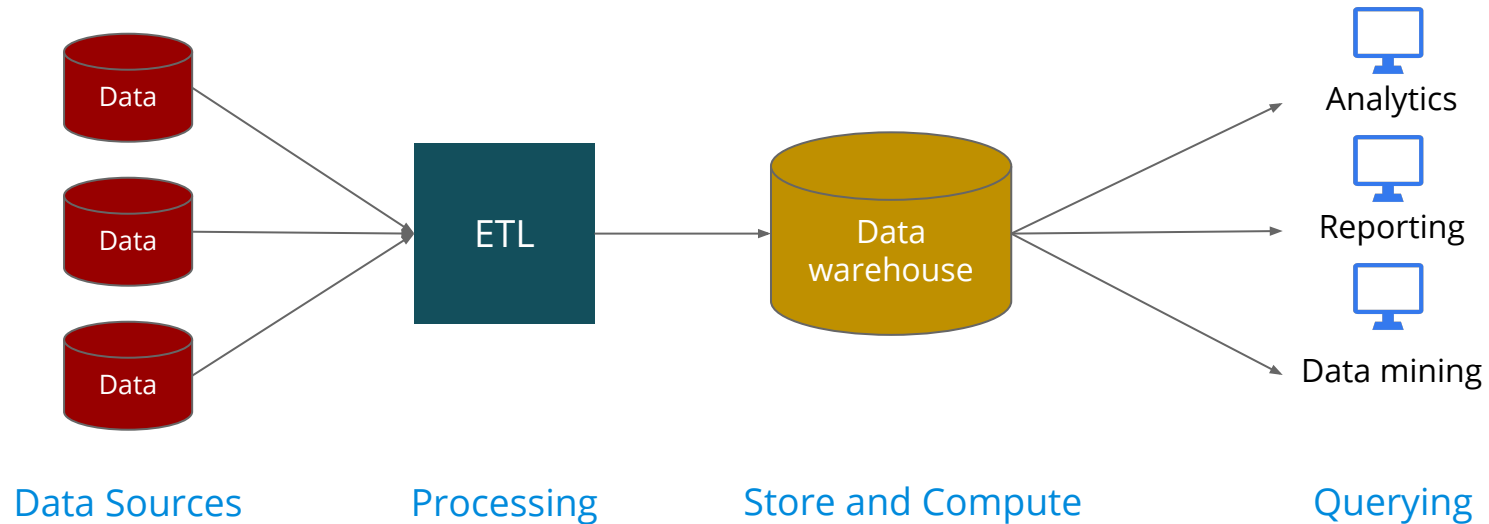
10 years working on scalable systems

# Agenda

- From warehouses to rivers

- What can we do with streaming data?

- Streaming analytics

- Enter the Druid
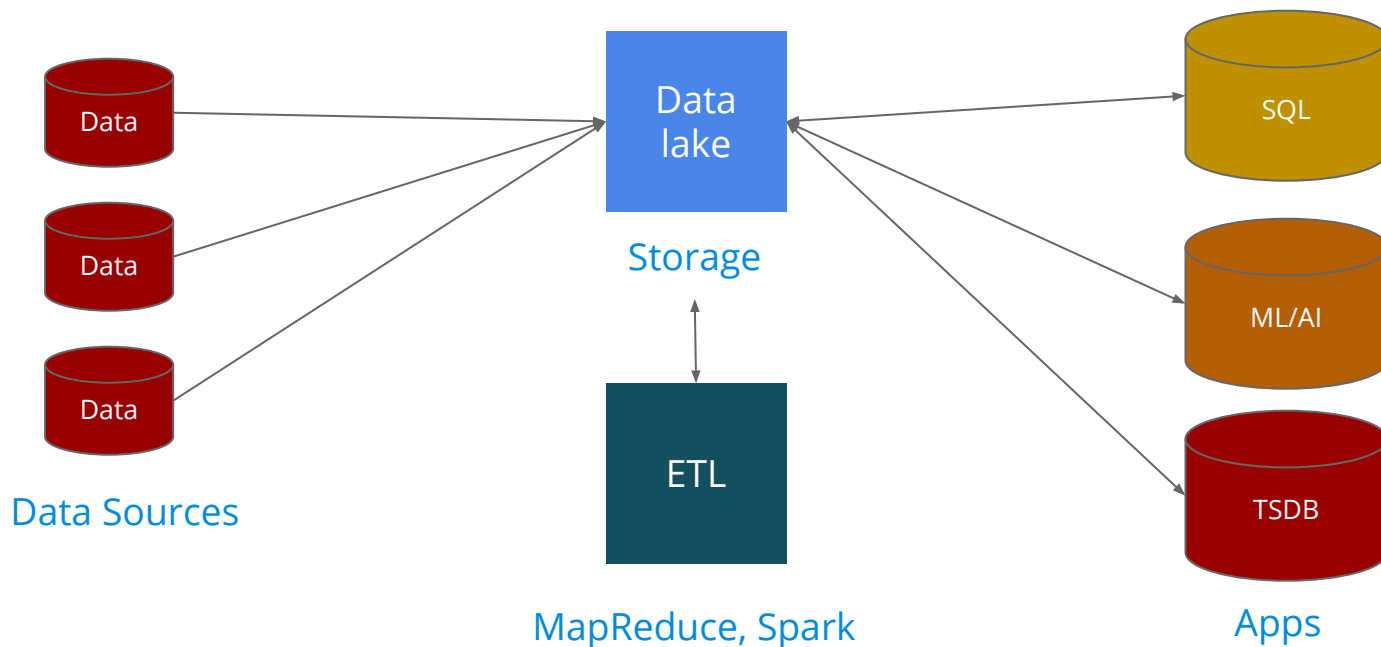
- Do try this at home!

Rolling down the river

# Data warehouses

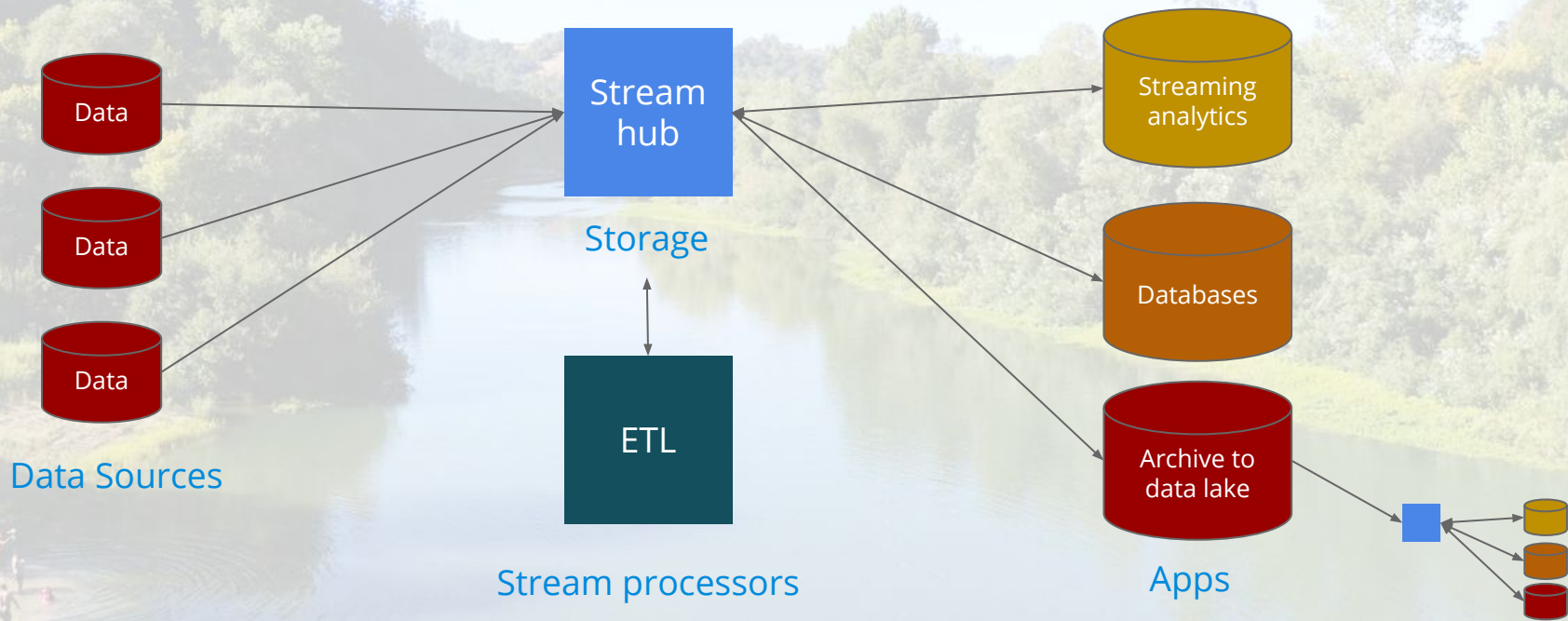Tightly coupled architecture with limited flexibility.



| Data Sources | Processing | Store and Compute | Querying |

# Data lakes

Modern data architectures are more application-centric.



Data Sources

Storage

MapReduce, Spark

Apps

# Data rivers

Streaming architectures are true-to-life and enable faster decision cycles.

Data

Data

Data

**Data Sources**

Stream hub

Storage

ETL

**Stream processors**

Streaming analytics

Databases

Archive to data lake

**Apps**

# Streaming data

8

# Streaming data

# Streaming data

App

Stream hub library

Stream hub

Direct production

# Streaming data



App ← DB connection → Transactional DB → Something DB-specific → Stream hub

Change data capture

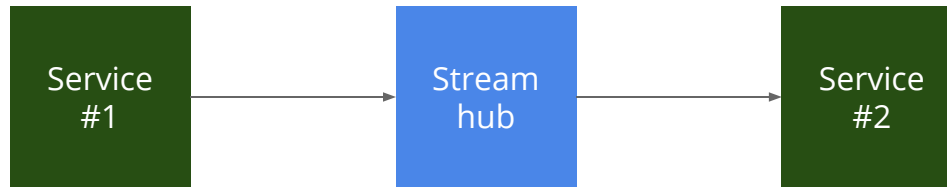# Streaming data

Stream hub → Kafka Connect Kinesis Firehose → EDW or data lake

Streaming data pipeline

# Streaming data



Microservice communication

# Streaming data



Google

top stream analytics products

All    Images    News    Videos    Shopping    More        Settings    Tools

About 214,000,000 results (0.40 seconds)

**Here are the top platforms being used all over the world for Streaming analytics solutions:**

- Apache Flink. Flink is an open-source platform that handles distributed **stream** and batch data processing. ...
- Spark **Streaming**. ...
- IBM **Streams**. ...
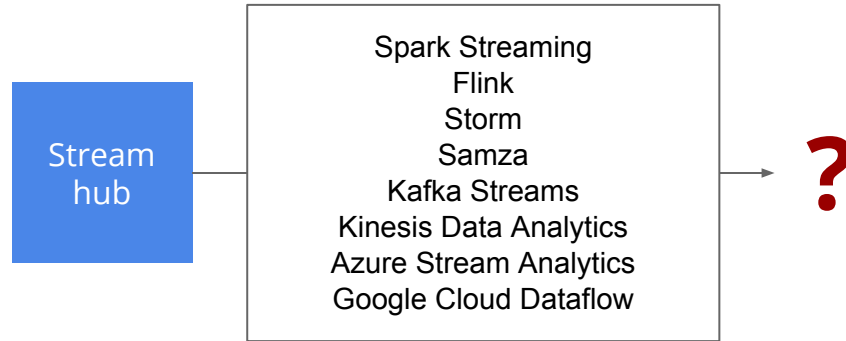- Azure **Stream Analytics**.

Mar 1, 2018

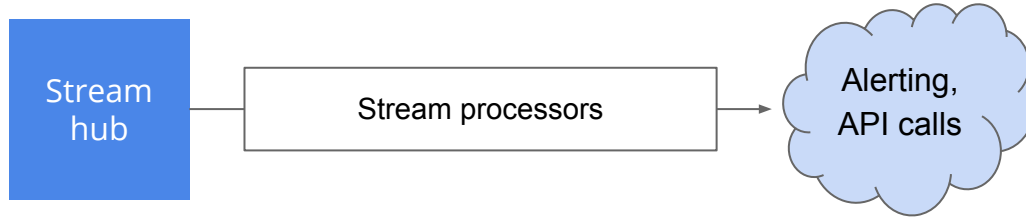5 Streaming Analytics Platforms For All Real-time Applications - Datafloq
https://datafloq.com/read/streaming-analytics-platforms-real-time-apps/4658

About this result    Feedback

# Streaming data

Stream hub → [ Spark Streaming
Flink
Storm
Samza
Kafka Streams
Kinesis Data Analytics
Azure Stream Analytics
Google Cloud Dataflow ] → **?**

15

# Streaming data

Stream hub → Stream processors → Alerting, API calls

Real-time actions

# Streaming data



Stream hub → Stream processors → Storage

HDFS
Cloud storage

Data movement

# Streaming data



Stream hub

Stream processors

Storage

Stream processors

Continuous query

HDFS
Cloud storage

Data movement + enrichment

18

# Streaming data



**Stream hub** → Stream processors → **K / V stores** → **Realtime dashboard**

Stream processors

HBase
Cassandra
Redis

Continuous query + write to serving layer
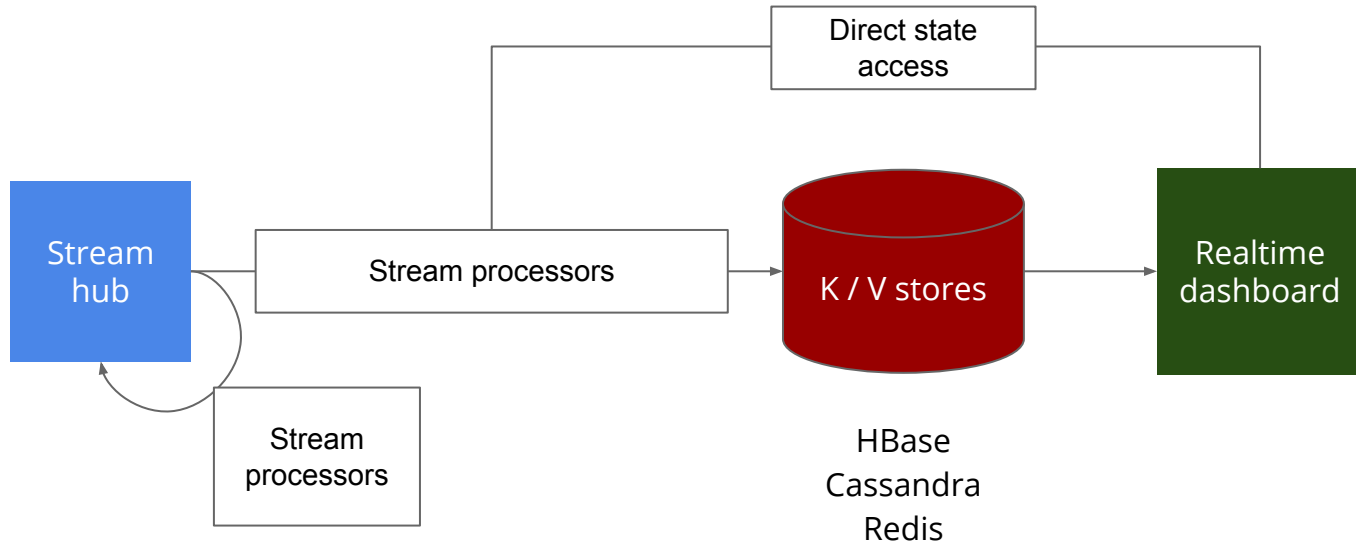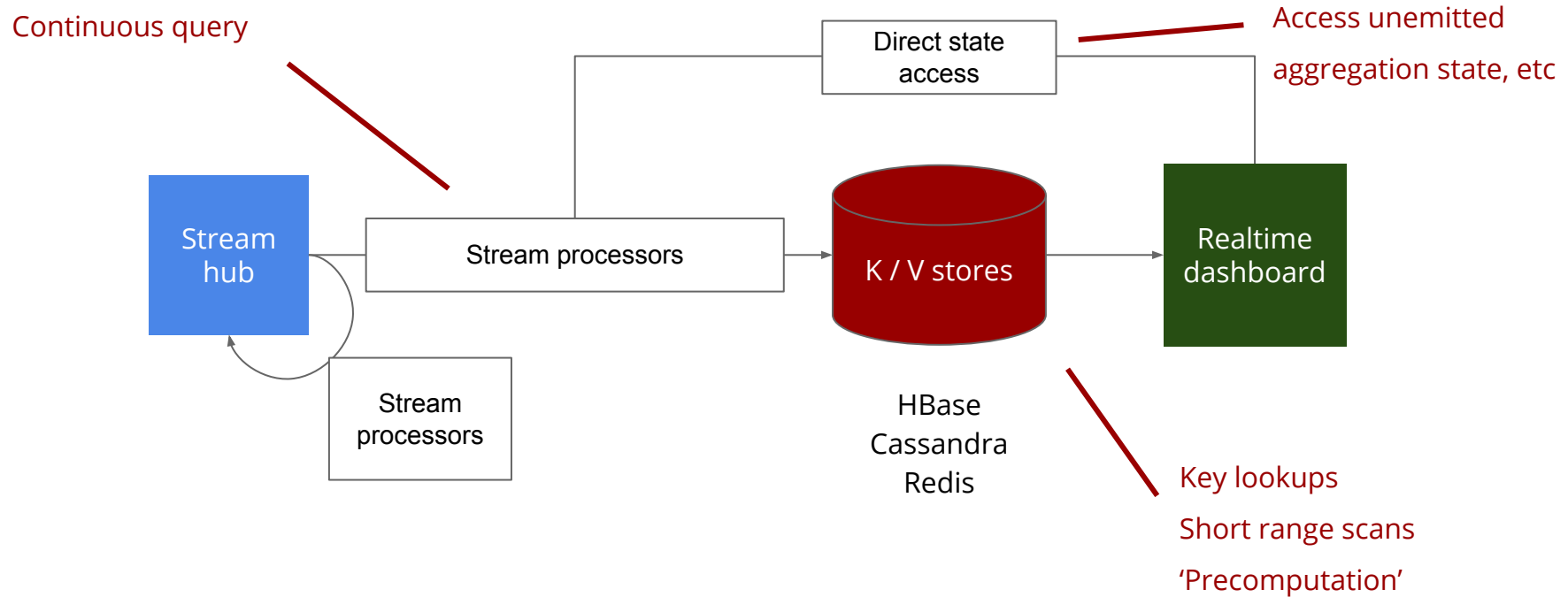
# Streaming data



Continuous query + write to serving layer + unemitted state serving

# Streaming data

Continuous query

Access unemitted
aggregation state, etc

Direct state access

Stream hub

Stream processors

Stream processors

K / V stores

HBase
Cassandra
Redis

Realtime dashboard

Key lookups

Short range scans

'Precomputation'

Continuous query + write to serving layer + unemitted state serving

# The problem



**DevOps Borat** @DEVOPS_BORAT · 23 Mar 2013

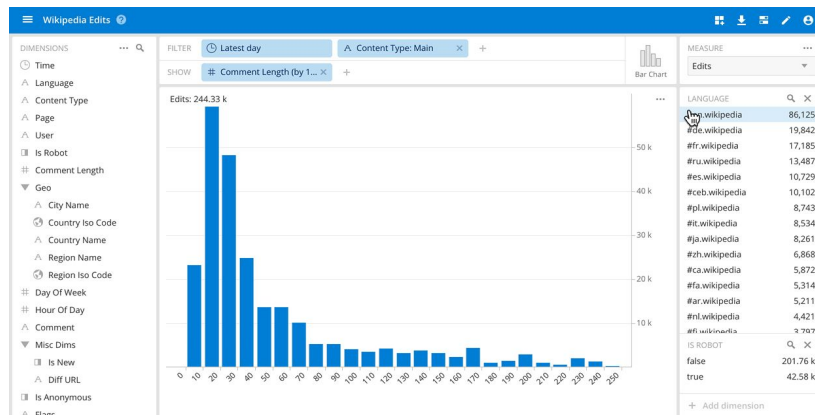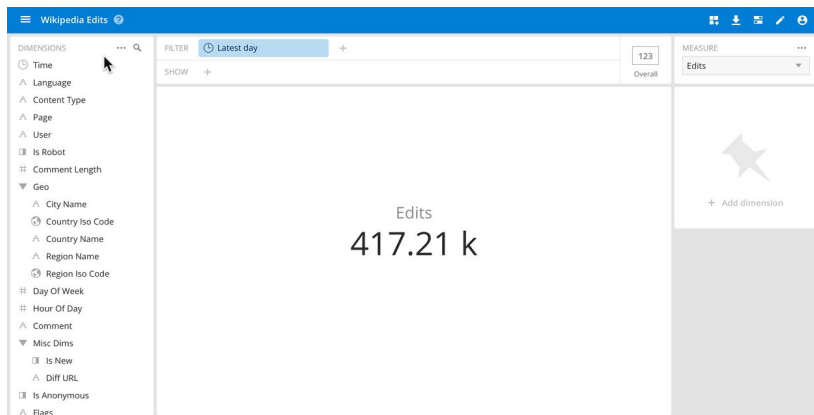In startup we have great of capability for churn out solution. Please send problem, we are pay good money.

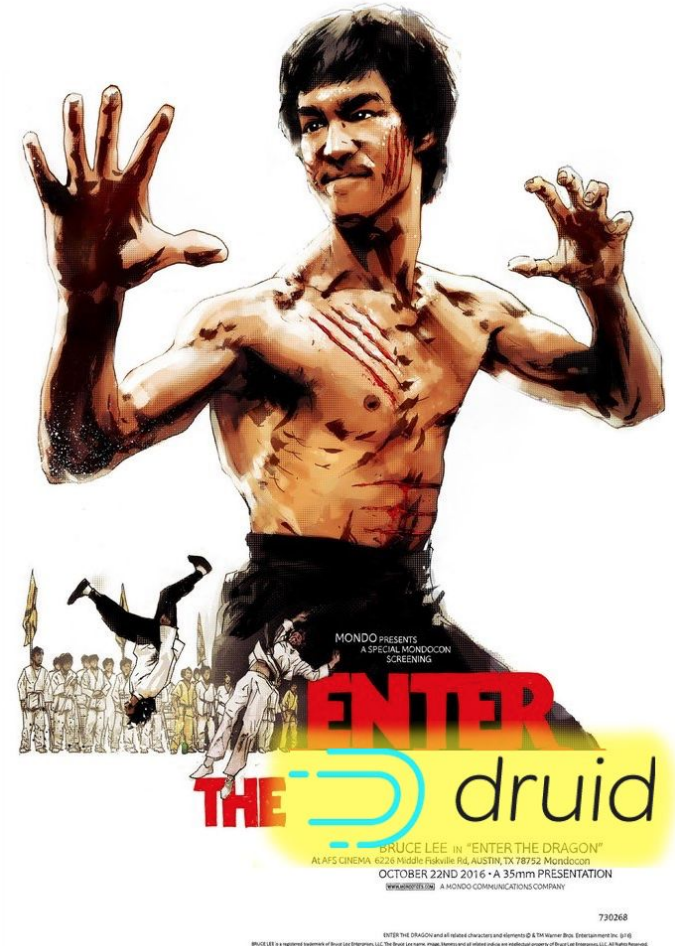↩   🔁 71   ⭐ 28   •••

# The problem

# The problem

- Slice-and-dice for big data streams

- Interactive exploration

- Look under the hood of reports and dashboards

- And we want our data fresh, too

# Challenges

- Scale: when data is large, we need a lot of servers

- Speed: aiming for sub-second response time

- Complexity: too much fine grain to precompute

- High dimensionality: 10s or 100s of dimensions

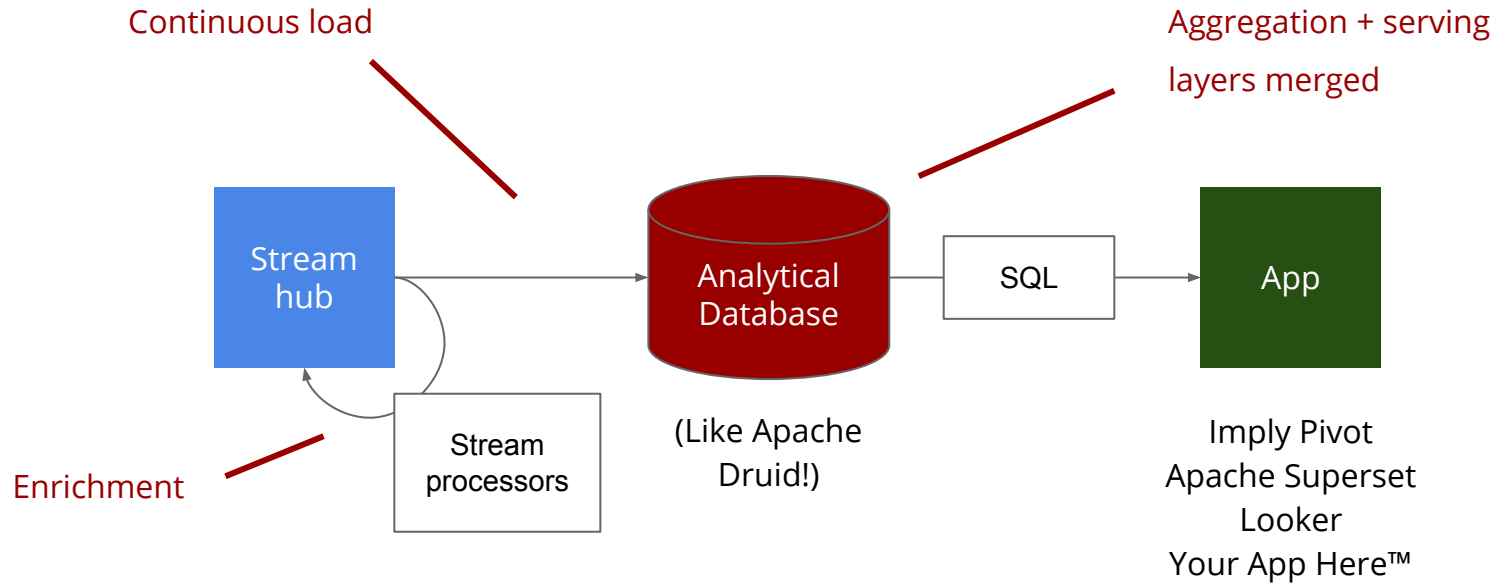- Concurrency: many users and tenants

- Freshness: load from streams

high performance
analytics data store for
event-driven data

# What is Druid?

- **"high performance":** low query latency, high ingest rates

- **"analytics":** counting, ranking, groupBy, time trend

- **"data store":** the cluster stores a copy of your data

- **"event-driven data":** fact data like clickstream, network flows, user behavior, digital marketing, server metrics, IoT

# Streaming data

Continuous load

Aggregation + serving
layers merged

Stream hub

Analytical Database

SQL

App

Stream processors

Enrichment

(Like Apache Druid!)

Imply Pivot
Apache Superset
Looker
Your App Here™

# Key features

- Column oriented

- High concurrency

- Scalable to 100s of servers, millions of messages/sec

- Continuous, real-time ingest

- Indexes on all dimensions by default

- Query through SQL

- Target query latency sub-second to a few seconds

# Use cases

- Clickstreams, user behavior

- Digital advertising

- Application performance management
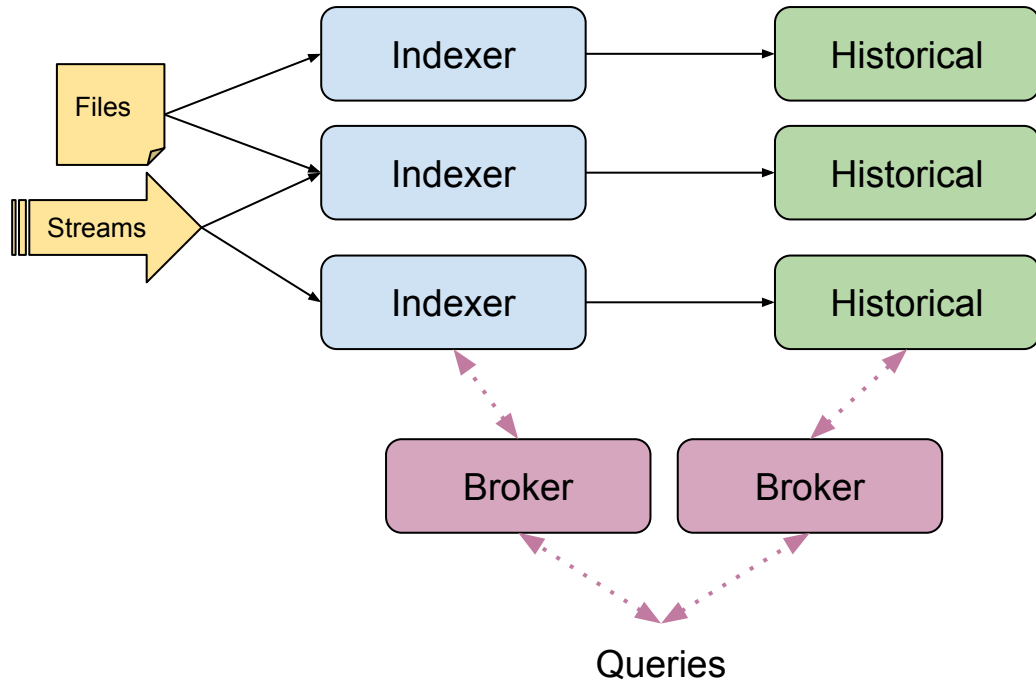
- Network flows

- IoT

# Powered by Apache Druid

# Powered by Apache Druid

From Yahoo:

"The performance is great … some of the tables that we have internally in Druid have **billions and billions of events** in them, and we're scanning them in **under a second**."

*Source: https://www.infoworld.com/article/2949168/hadoop/yahoo-struts-its-hadoop-stuff.html*
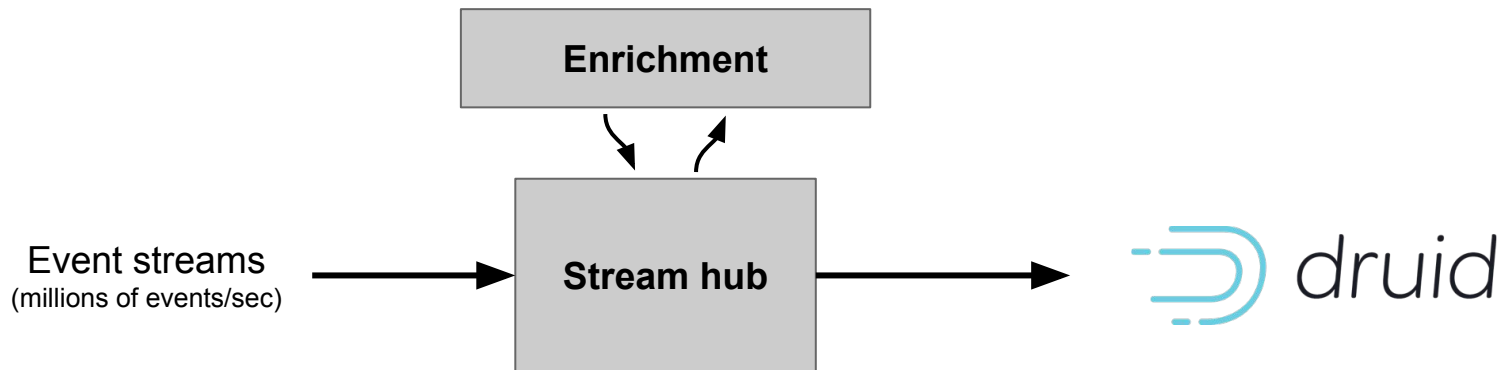
# Architecture

# Why this works

- Computers are fast these days

- Indexes help save work and cost

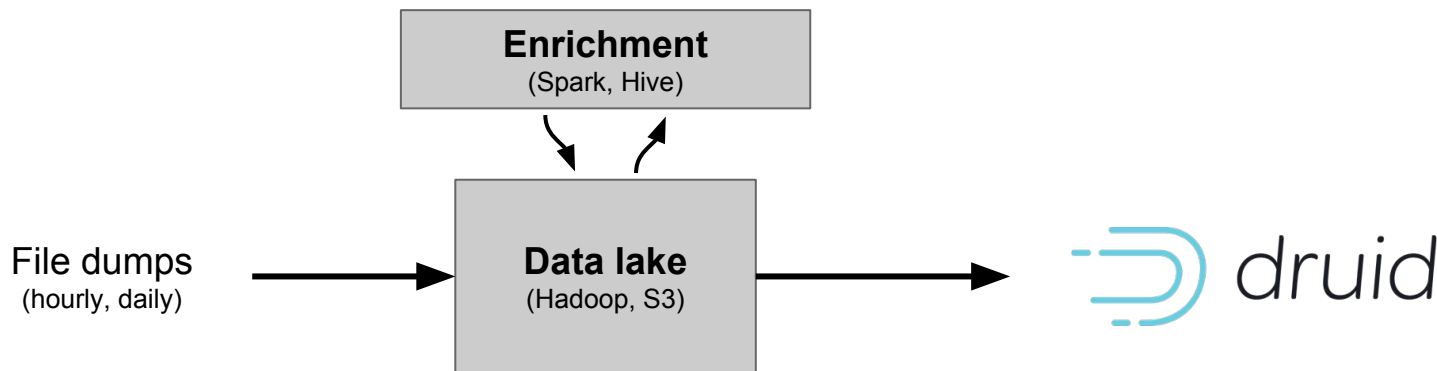- But don't be afraid to scan tables — it can be done efficiently
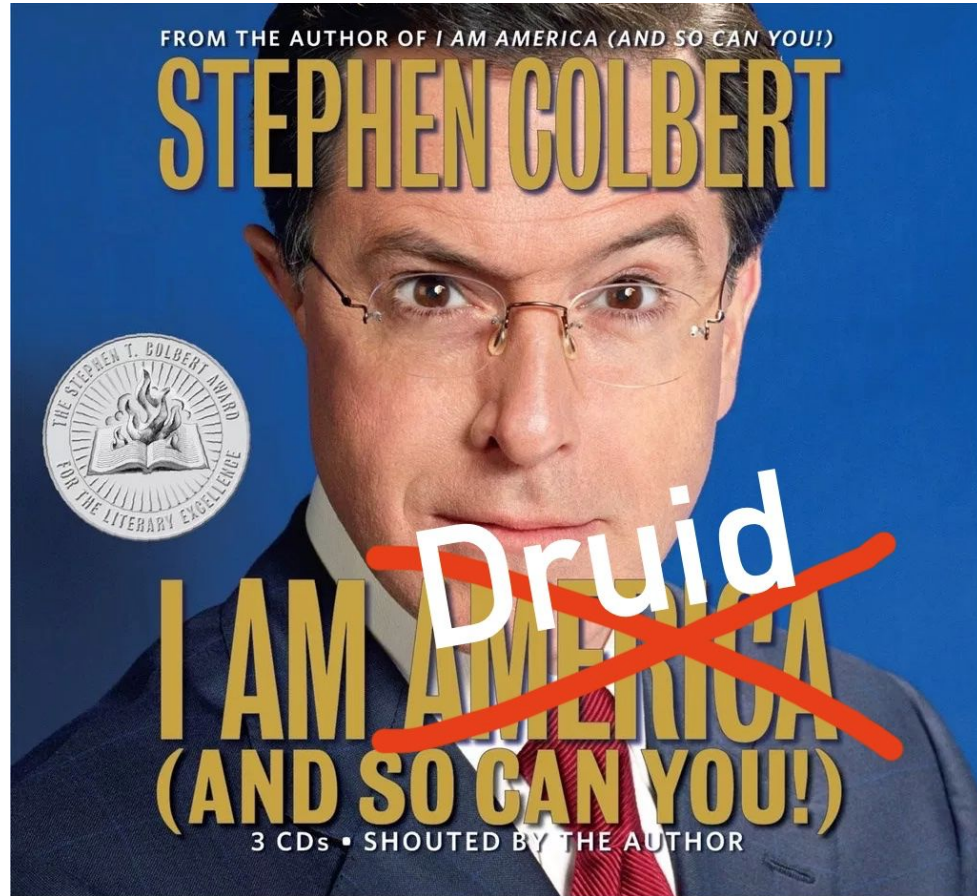
# Integration patterns

# Deployment patterns



- Modern data architecture

- Centered around stream hub

# Deployment patterns



**Enrichment**
(Spark, Hive)

File dumps
(hourly, daily)

**Data lake**
(Hadoop, S3)

druid

- (Slightly less) modern data architecture

- Centered around data lake

# Download

Apache Druid community site (new): https://druid.apache.org/

Apache Druid community site (legacy): http://druid.io/

Imply distribution: https://imply.io/get-started

# Contribute



apache / incubator-druid

Unwatch ▾  550    ★ Unstar  6,848    ⑂ Fork  1,684

<> Code    ⊙ Issues  991    ⑂ Pull requests  137    ▥ Projects  3    ▤ Wiki    ▥ Insights

Apache Druid (Incubating) - Column oriented distributed data store ideal for powering interactive applications
http://druid.io

⊙ 8,622 commits    ⑂ 26 branches    ◌ 409 releases    ⠿ 238 contributors    ⚖ Apache-2.0

Branch: master ▾    New pull request    Create new file    Upload files    Find file    Clone or download ▾

https://github.com/apache/druid

41

# Stay in touch

Follow the Druid project on Twitter!
@druidio

Join the community!
http://druid.apache.org/