# SCALING A DATA PIPELINE:
## MYSTERY TO MASTERY

Dan Goldin

@dangoldin

**triplelift**

# AGENDA

Introduction

AdTech and Data

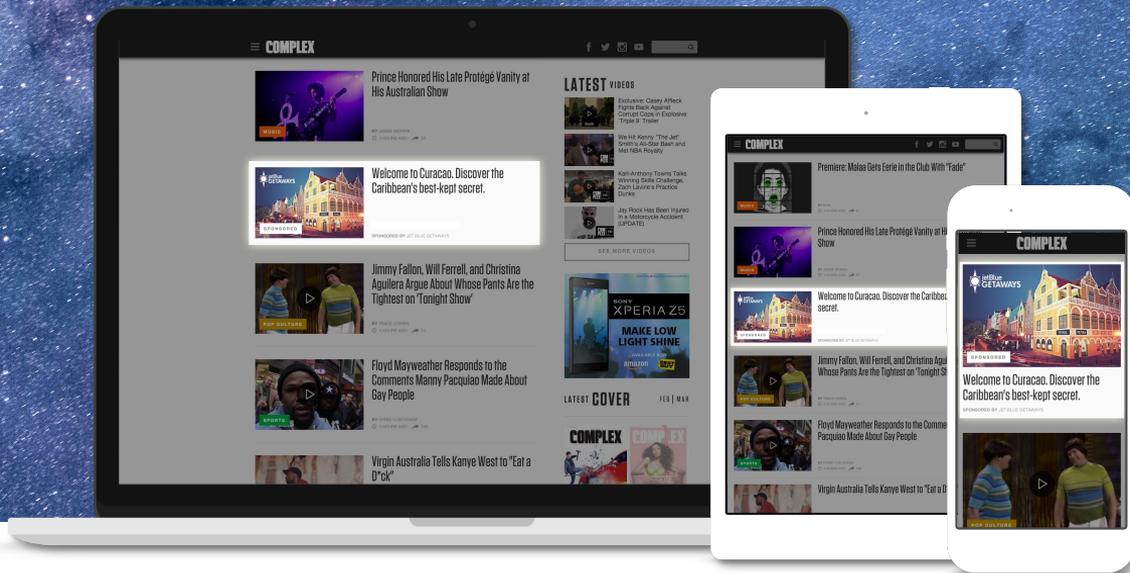The Evolution

Current State

Lessons Learned

Q&A

# INTRODUCTION

## SIMPLE

Render brand's assets to match the unique look and feel of the publisher

## SCALABLE

Bringing scale to high performing consumer friendly formats

## EFFECTIVE

Integrations into the world's largest DSPs – Google, The Trade Desk, Turn, MediaMath, AppNexus and more
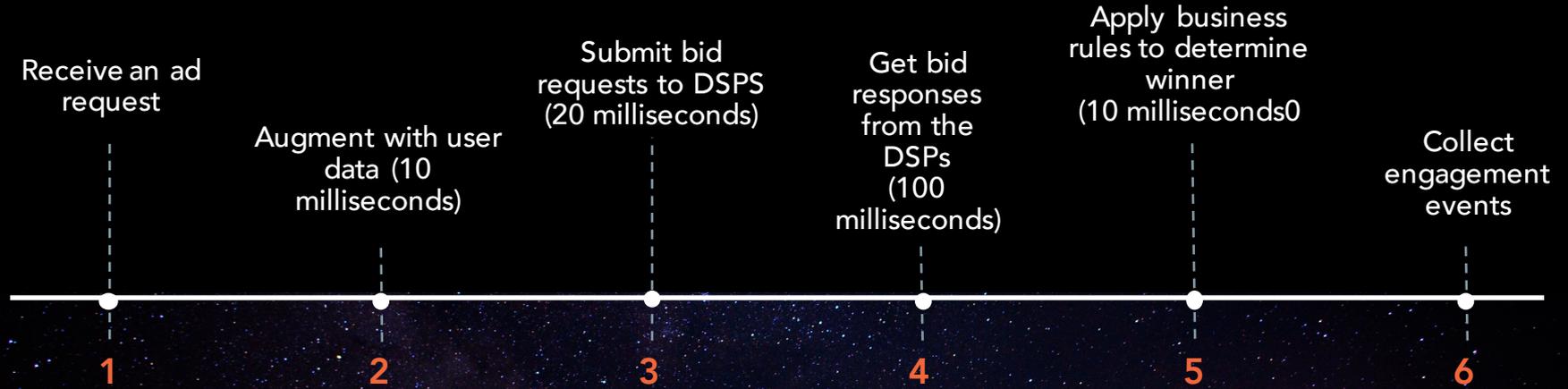
DATA & ADTECH

High **volume** across many **dimensions** that needs to be handled in real time

VOLUME

# REAL TIME BIDDING AUCTION TIMELINE

Receive an ad request

Augment with user data (10 milliseconds)

Submit bid requests to DSPS (20 milliseconds)

Get bid responses from the DSPs (100 milliseconds)

Apply business rules to determine winner (10 milliseconds0

Collect engagement events

1  2  3  4  5  6

# TRIPLELIFT TODAY

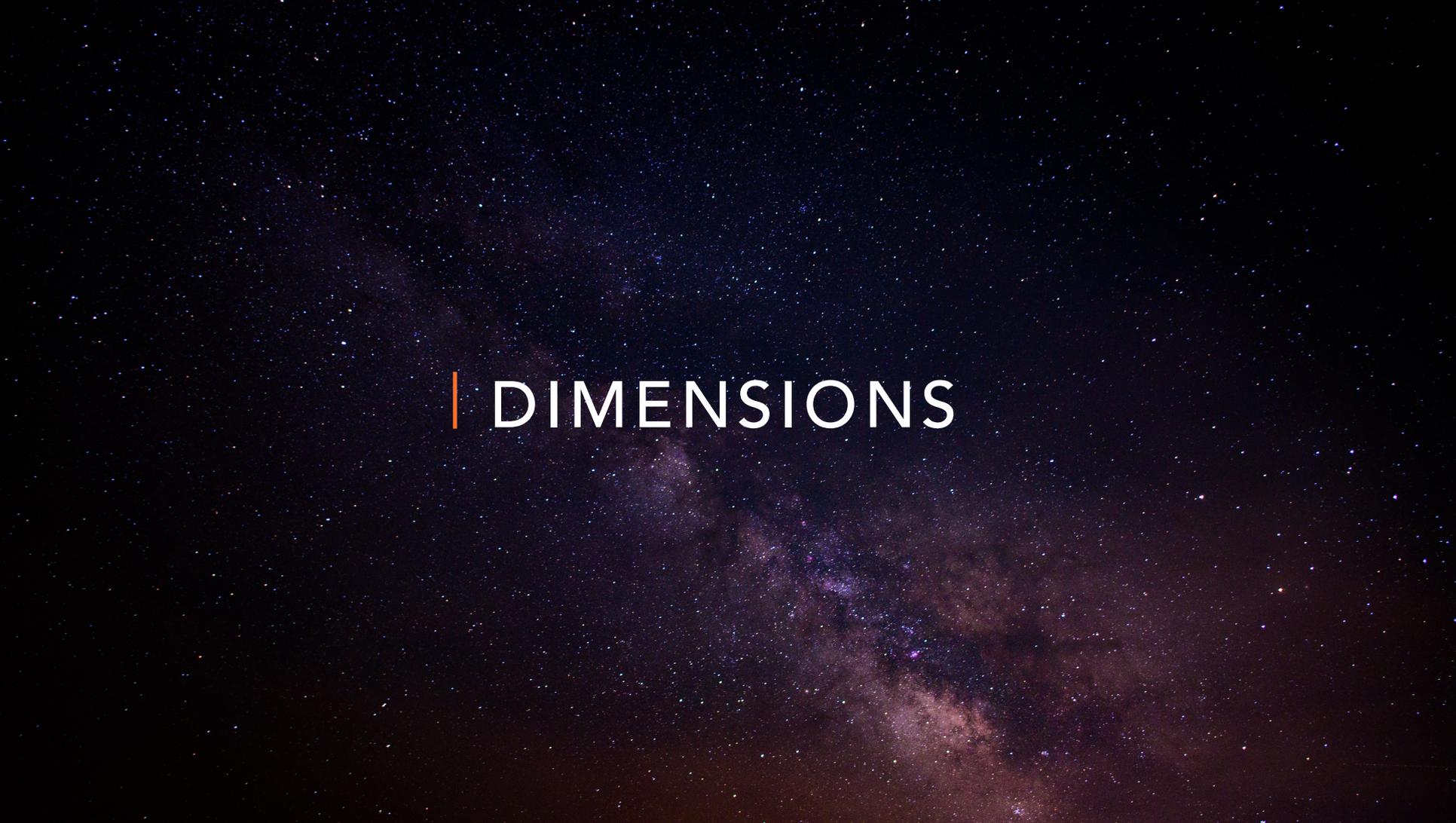**25K+**
Auctions/sec

**2B+**
Auctions/day

**~300**
Bids/auction

**~600B**
Bids/day

**~21B**
Events/day

DIMENSIONS

# LOTS OF DIMENSIONS

**AD REQUEST**

Browser, Device, OS, Geography, Domain, Time

**BID RESPONSES**

Creative, Format, Brand Buyer, Bid Amount

**WINS**

Price

**ENGAGEMENTS**

Click, Duration, Render, View
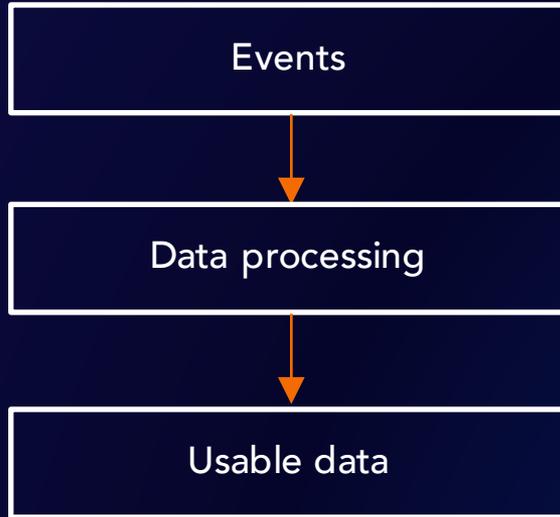
# AUCTION EVENTS

Ad Request

Auction Bid Request

Bid Responses

Win Events

AUCTION

Render

Click

Viewability

Mouseover

Video

More and more and more

| THE EVOLUTION

# SO WHAT'S A DATA PIPELINE?

# SO WHAT'S A DATA PIPELINE?

```
┌─────────────────────────────┐
│          Events             │
└─────────────────────────────┘
               │
               ▼
┌─────────────────────────────┐
│        Aggregation          │
└─────────────────────────────┘
               │
               ▼
┌─────────────────────────────┐
│          Storage            │
└─────────────────────────────┘
               │
               ▼
┌─────────────────────────────┐
│        APIs and UIs         │
└─────────────────────────────┘
```
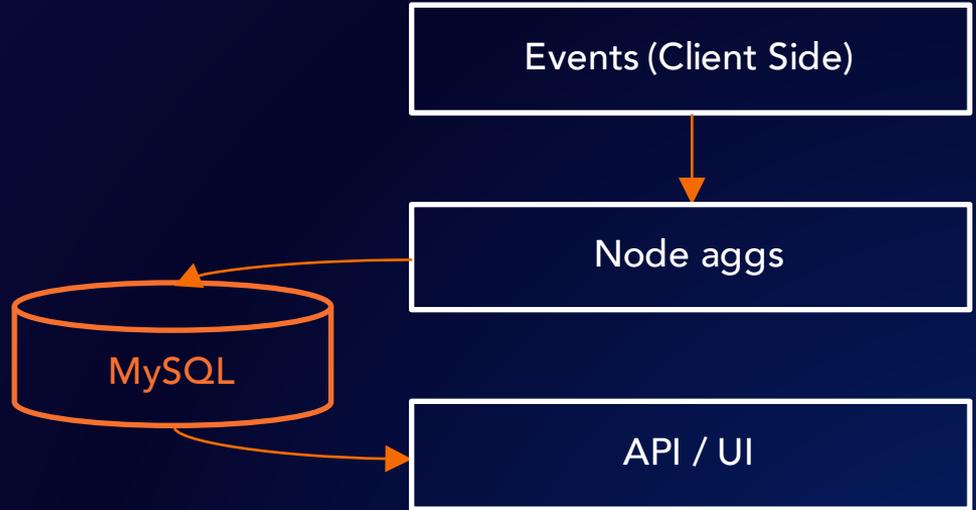
# v1: Sad but true

Implementation highlights
- Variable sample rate
- Keep a running sum in memory and write to MySQL every few minutes

Challenges
- Constant open connection to DB
- Tables became large and unwieldy
- Difficult to slice and dice sampled data
- Easy to lose data
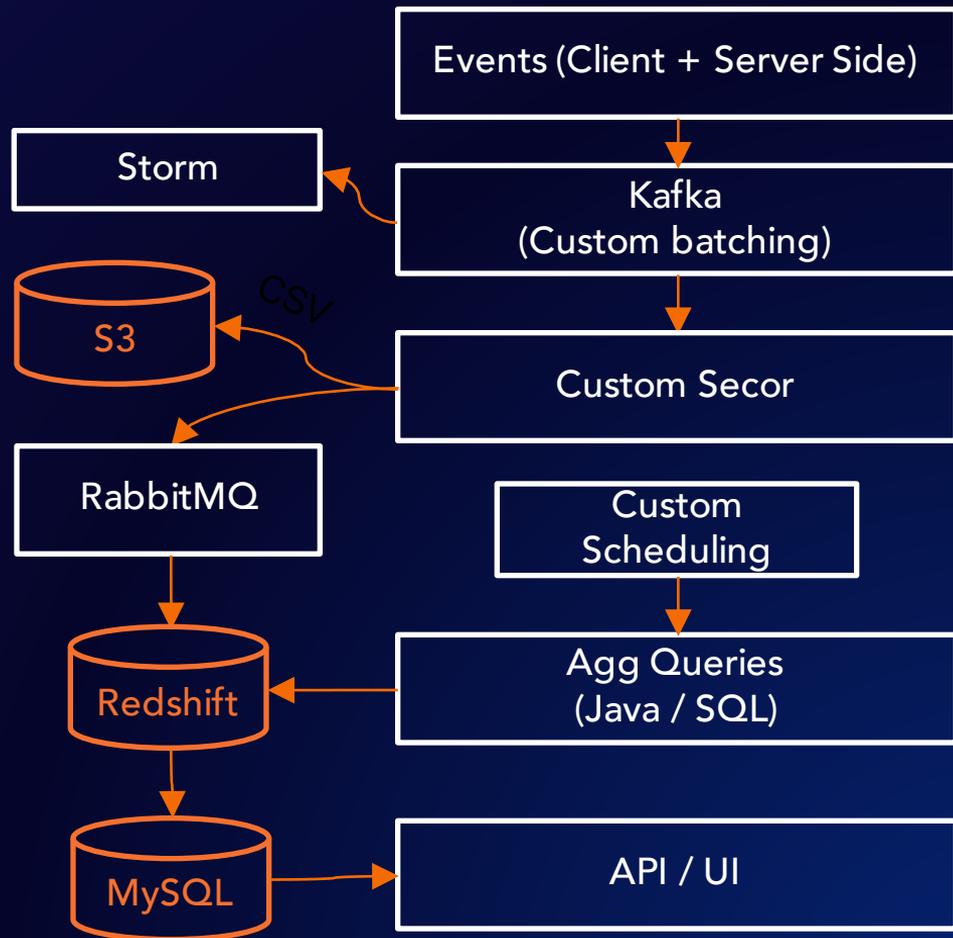
Events (Client Side)

Node aggs

MySQL

API / UI

# v2: Kafka, Secor, and Redshift

Implementation highlights
- Collect every event in Kafka
- Upload to S3 and load into Redshift
- All jobs done through Redshift queries
- Storm to handle real time pacing

Challenges
- Dependencies tough to manage
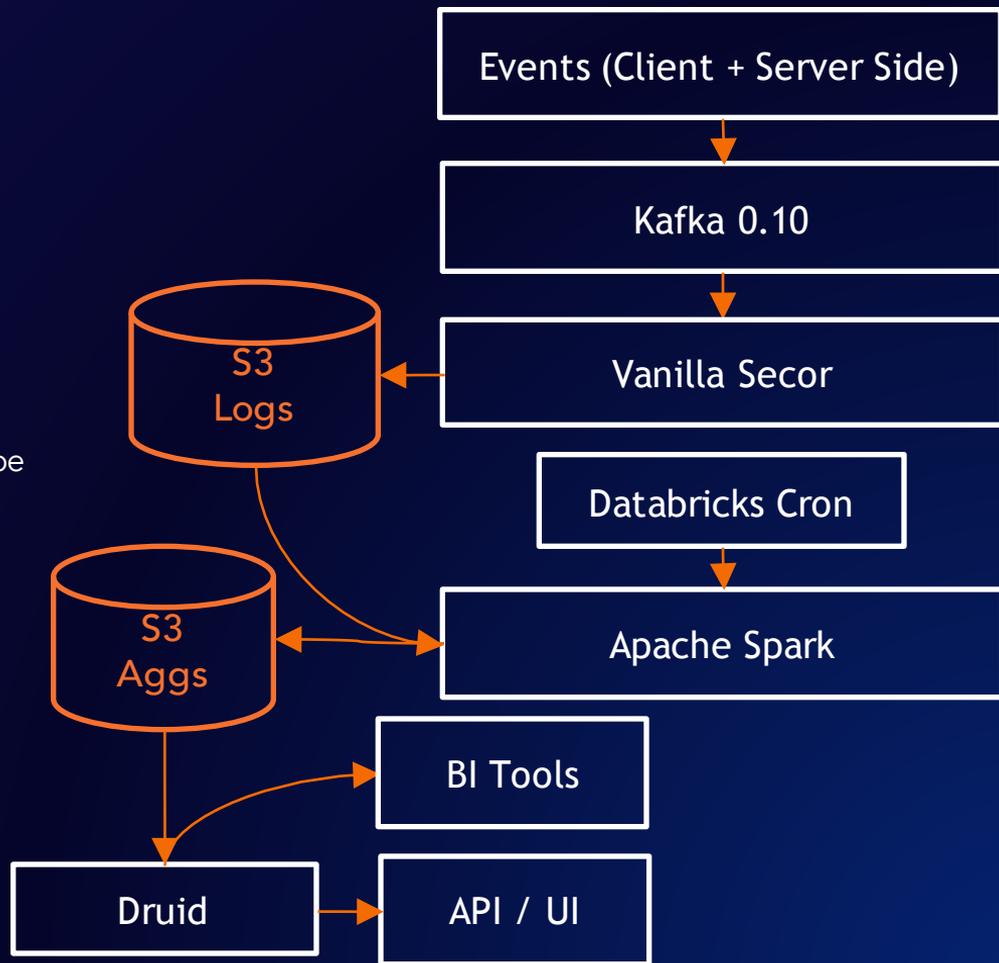- Couldn't do everything via SQL queries
- Redshift became expensive

# v3: Hello Spark; Hello Druid

Implementation highlights
- Kafka 0.10
- Failed attempt at Spark Streaming
- Spark was a big improvement
  - Cheaper & more scalable than Redshift
  - More advanced query logic
- Druid also helped
  - Trivial to scale to 100s of metrics and dimensions
  - Replaced a dozen tables with a single cube
  - Improved query times

Challenges
- More tech to maintain
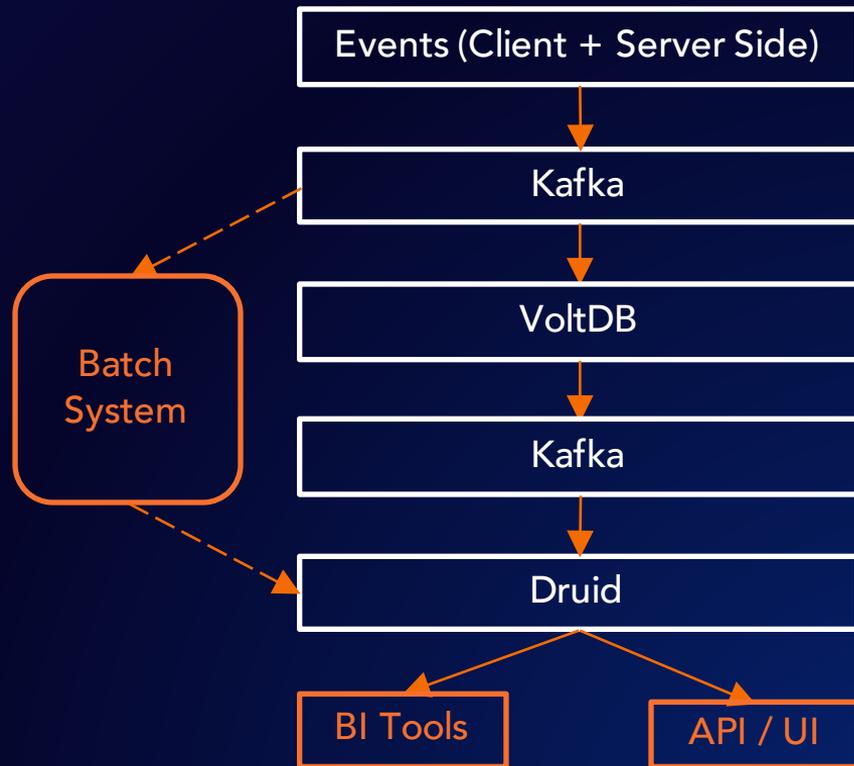- Scheduling still a challenge
- More complex development process

Events (Client + Server Side)

Kafka 0.10

Vanilla Secor

S3 Logs

Databricks Cron

S3 Aggs

Apache Spark

BI Tools

Druid

API / UI

# v4: Lambda, the ultimate?

Implementation highlights
- Introduced VoltDB
- Feeds back into our Druid cluster
- Delays in batch jobs masked

Challenges
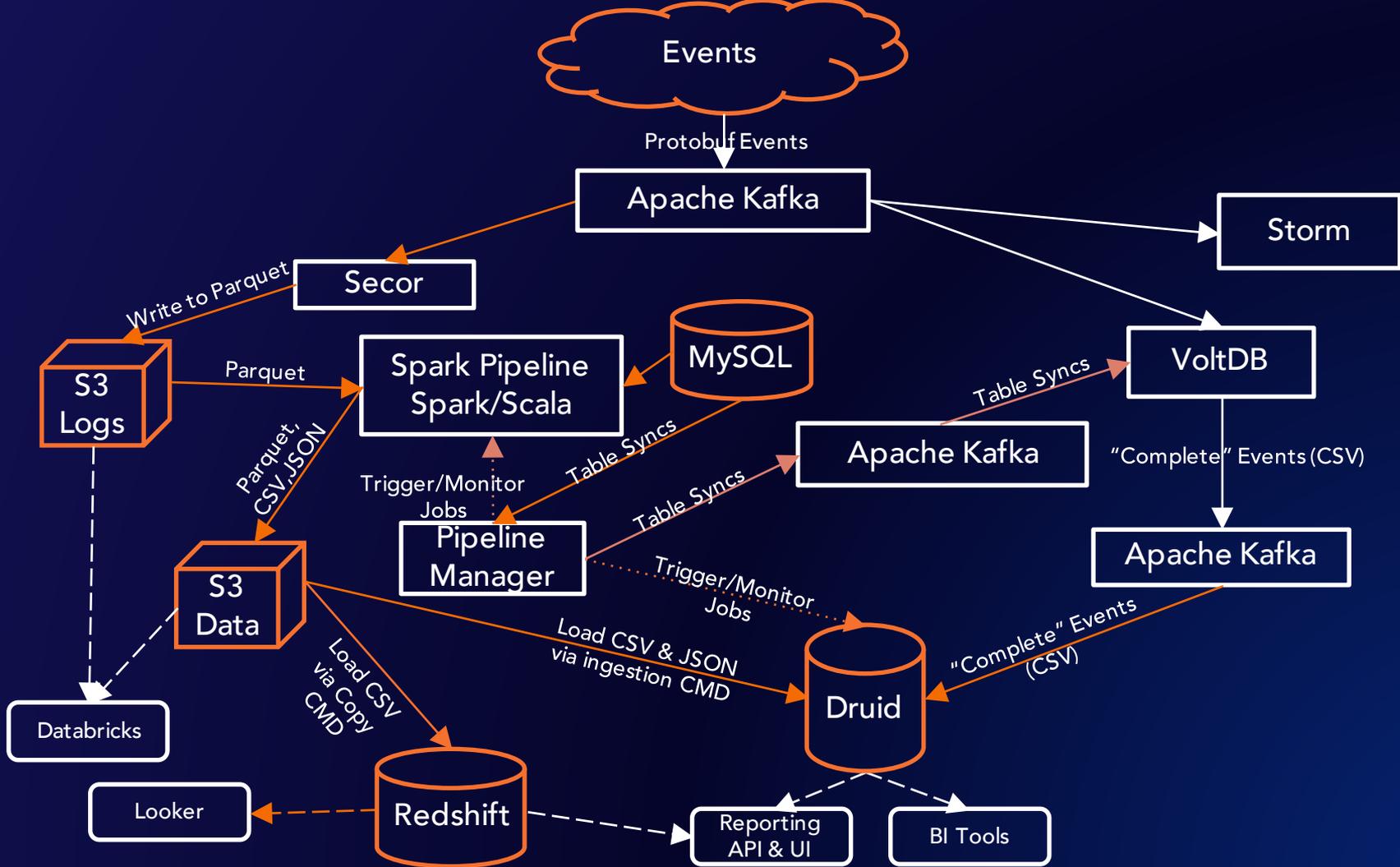- Even more tech to maintain
- Real time can get real expensive

CURRENT STATE

LESSONS LEARNED

- The seemingly simple stuff is difficult
    - Dependencies
    - Scheduling

- Stop hacking open source libraries: Vanilla is an uninspired yet classic and delicious flavor
    - Secor
    - Kafka

- SQL really is everywhere

- Changing code is much easier than changing data

- The big data ecosystem is huge with tons of tools

Q & A

THANK YOU