

adform

# Real-time Business Intelligence with Big Data

Peter Milne  @helipilot50  helipilot50





# Agenda

---

- How big - How fast
- Real-time Bidding overview
- Real-time “ish” Campaign Reporting
- Uncle Pete’s Advice

adform

**Ad Tech**  
**How big – How fast**



# How big – How fast

## How big

1 byte equals 1 second, **how big is a:**

Bytes	Time
1 Byte	1 second
1 Mega Byte	12 days
1 Giga Byte	34 years
1 Terra Byte	34,700 years
1 Peta Byte	34 million years

**Adform has 270-370 million “years” of data**

## How fast



Speed of an **eye blink** is  
200-400 ms.

**In 1 eye-blink, Adform  
process 900,000 –  
1,700,000 bid requests**

# Global Infrastructure

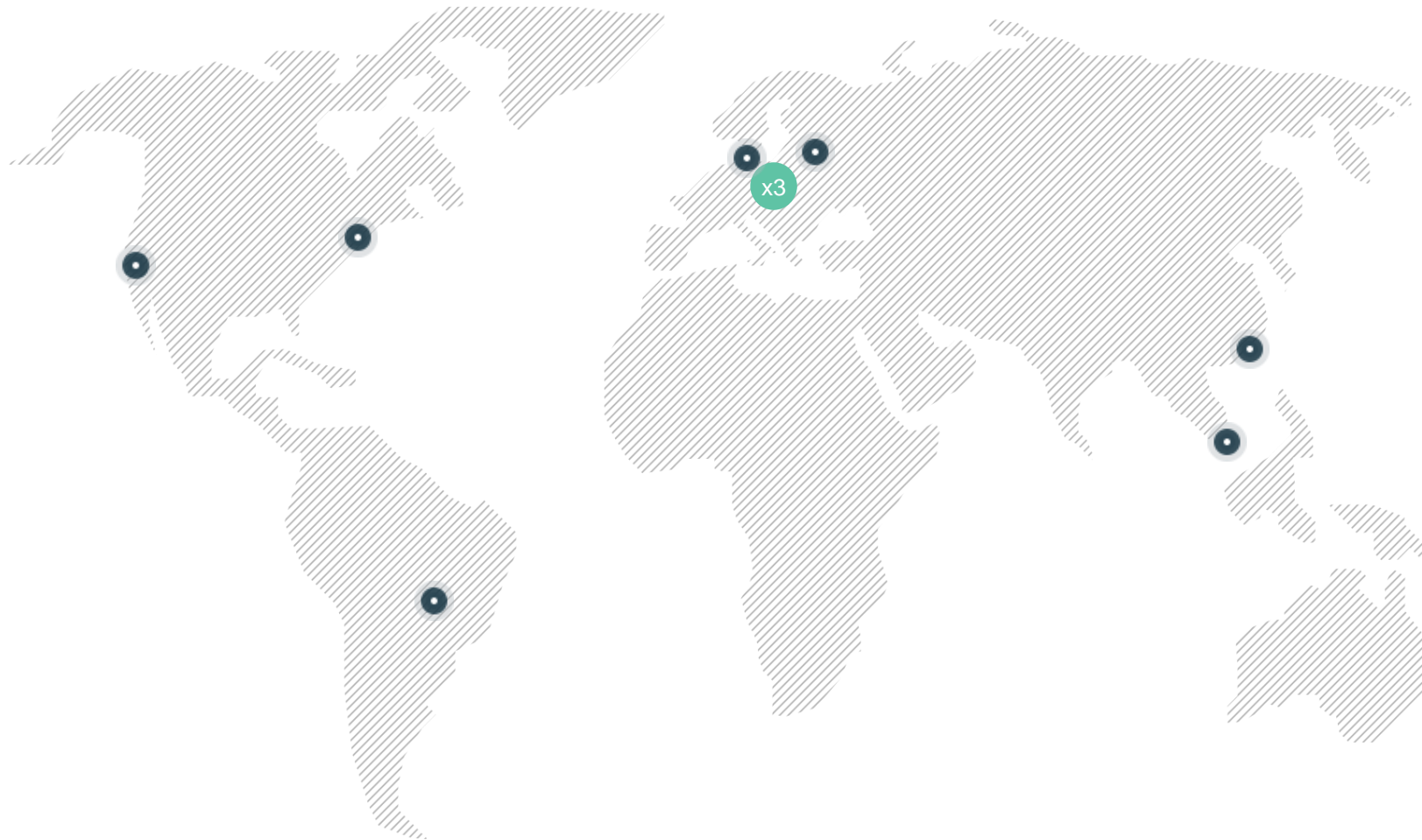
9  
Datacenters

1800+  
Physical Servers

7-11 PB  
HDFS Storage

2+ Mil  
Kafka Messages  
per Second

620+ Gb/s  
Traffic via core



11+ Mil  
Aerospike Ops/S

3200+  
Openstack VMs

1250+  
Hyper-V VMs

130+ Gb/s  
Incoming Network  
Traffic

3,6+ Mil  
Nginx Ops/S

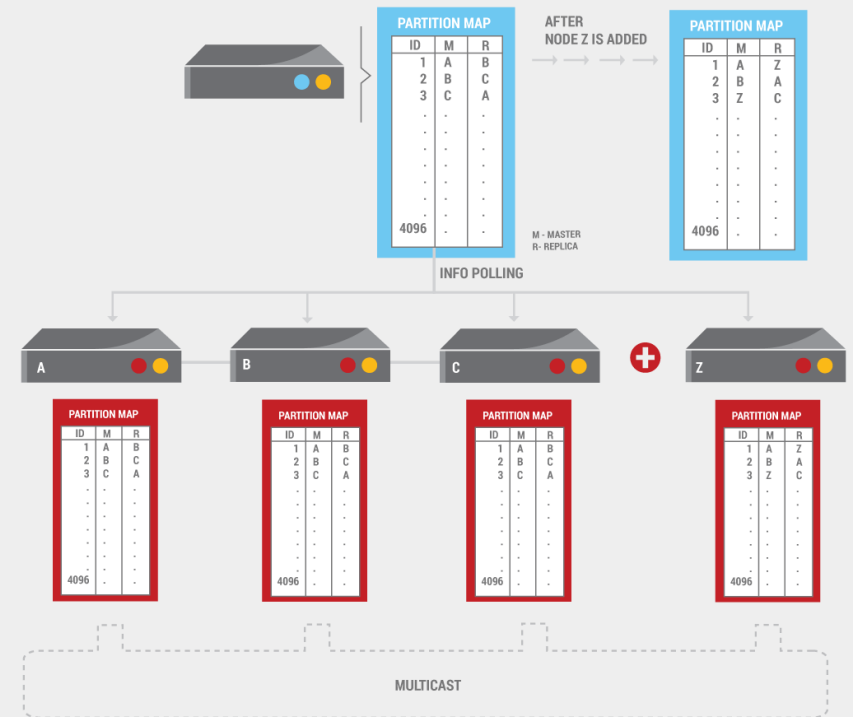
# Fast Data

## Row store – Great for primary key queries

- Total data volume ~**22 TeraBytes**
- Average query latency - **0.9 millisecond**
- Biggest set – **9 billion rows**
- ~**11 million** operations per **second**



`SELECT device, geo, segments FROM cookie_store WHERE cookield = 826701826`



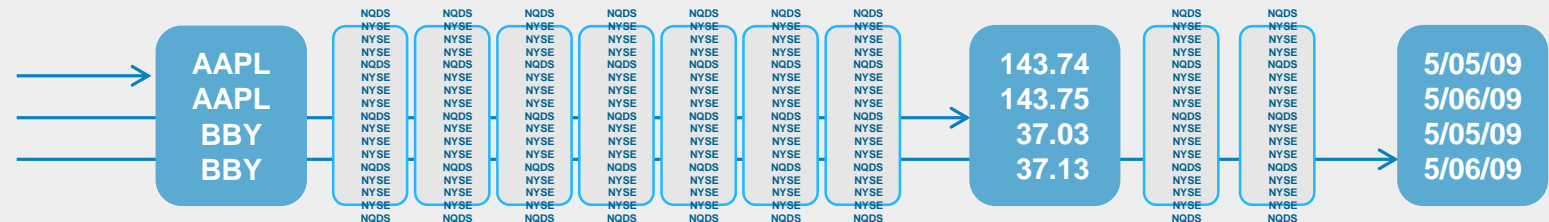
# Analytical Data

## Scalable Column Store - Great for complex multi-filter queries

- Total data volume **~250 TeraBytes**
- Average query latency **0.541 sec**
- Biggest tables **100 billion rows**
- **300 thousand** reporting queries per **day**
- **Ingest 6 billion** rows per **day**



```
SELECT avg(price) FROM tickstore
WHERE symbol = 'AAPL' date =
'5/06/09'
```



Column Store - Reads 3 columns



# Big Data – Data Lake

## Raw Data

For:

- Analytical data
- Extract Transform Load (ETL)
- Large data volumes
- Cold storage / data archiving

Example: **~10 Peta Bytes** client campaign data





adform

# Real-time bidding overview

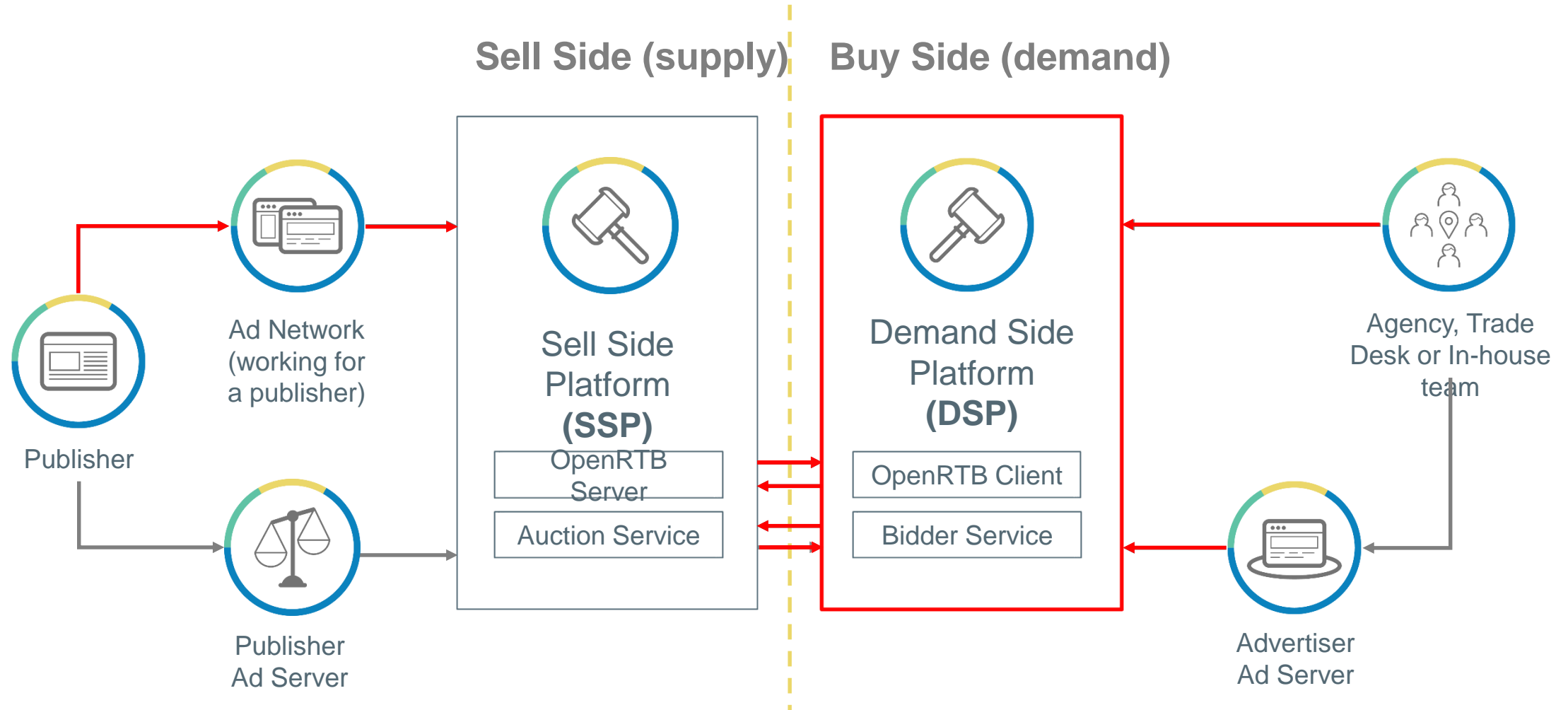


# Digital Advertising

The image shows a screenshot of a Stack Overflow page. The browser address bar shows 'stackoverflow.com/#'. The page features a list of questions on the left and a sidebar on the right. A red rectangular box highlights a Microsoft Azure advertisement. The ad includes the Microsoft Azure logo, the text 'Java, Node.js, PHP, .NET Get started in any language.', and a 'Free account' button. Below the ad, there are sections for 'Jobs near you' and 'Hot Network Questions'.

votes	answers	views	question title	tags	asked	author
0	0	4	leader election. Please refresh	apache-zookeeper, apache-flink, flink-streaming	asked 37 secs ago	kadsank 69
0	0	4	Algorithm improvement for generating unary-binary trees	algorithm, motzkin-numbers, unary-binary-trees	modified 40 secs ago	Guy Coder 9,760
0	0	6	Find max overlap in list of lists	python, list	asked 1 min ago	elcombato 108
0	2	20	Using binding and DataContext in UWP not working	c#, data-binding, uwp, datacontext	answered 1 min ago	Maplestrip 3
0	0	5	Android download file from url	android, get, android-download-manager	asked 1 min ago	Alex RED 53
0	0	2	CommonsXsdSchemaCollection and import statements	spring-ws	asked 1 min ago	gianluca 1
0	0	2	Can someone please explain how this works?fork(),sleep()	linux, fork, wait, zombie-process	asked 1 min ago	Aarthi Vishwanathan 11
0	0	3	Can I use scrapy with Yelp API?	python, api, scrapy, yelp	asked 1 min ago	Billy Jhon 25
0	0	2	My first PIC32MX ISR not firing, code is hanging	mplab, isr, pic32	asked 2 mins ago	Chris 90
0	0	5	Rendering with many and small amount of DOM operations	javascript, performance, google-chrome, layout	asked 2 mins ago	Jaroslav Rewers 153
1	1	9	Display Joomla Intro Article without <p> tag			

# Real Time Biddings - Auction



adform

# Real-time “ish” Campaign Reporting



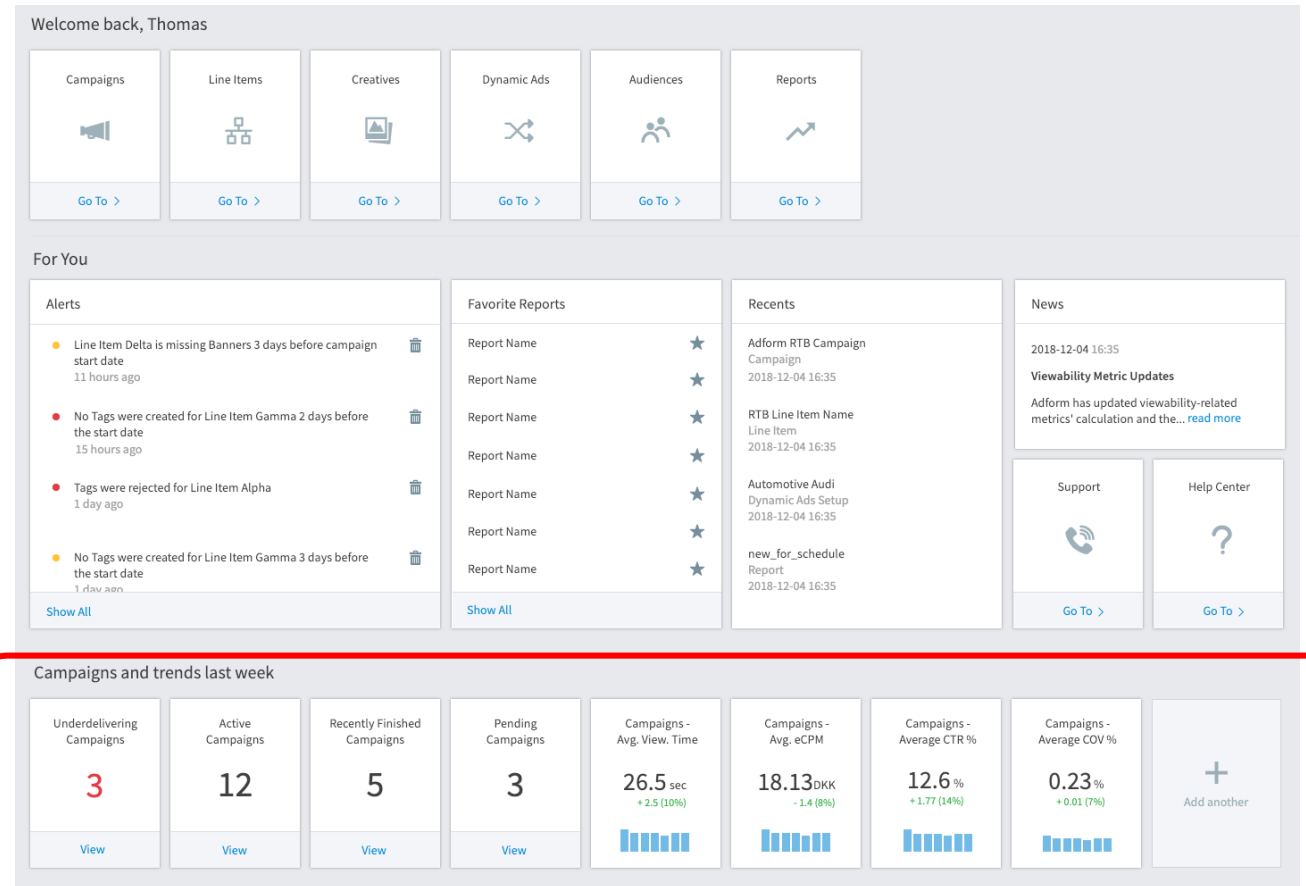
# Campaign reporting

## Fast reports

- **Real enough** Time for the human user
- Graphs and cards in the UI
- Quick visual indicators
- Latency: **< 1 second**

## Slow reports

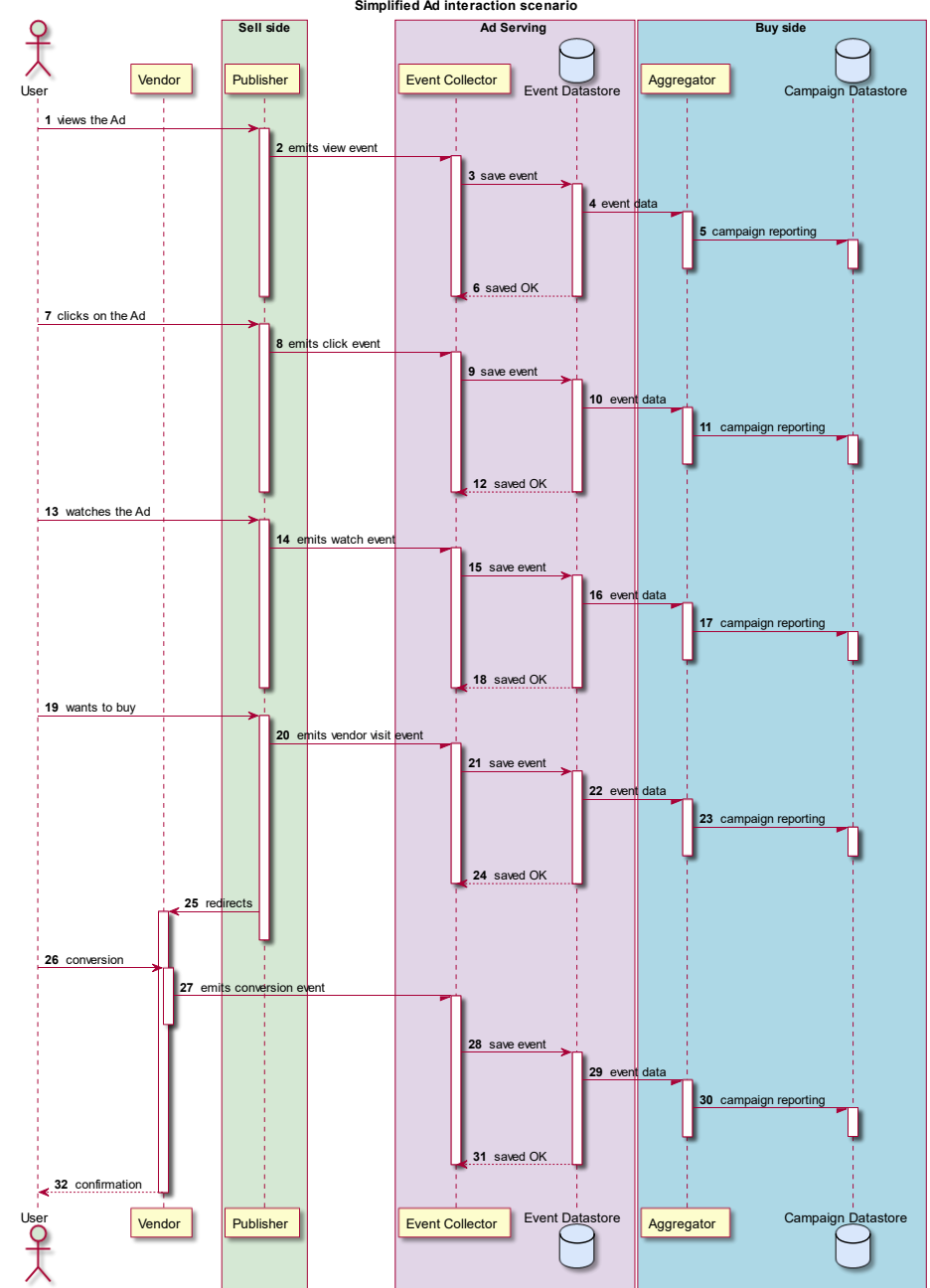
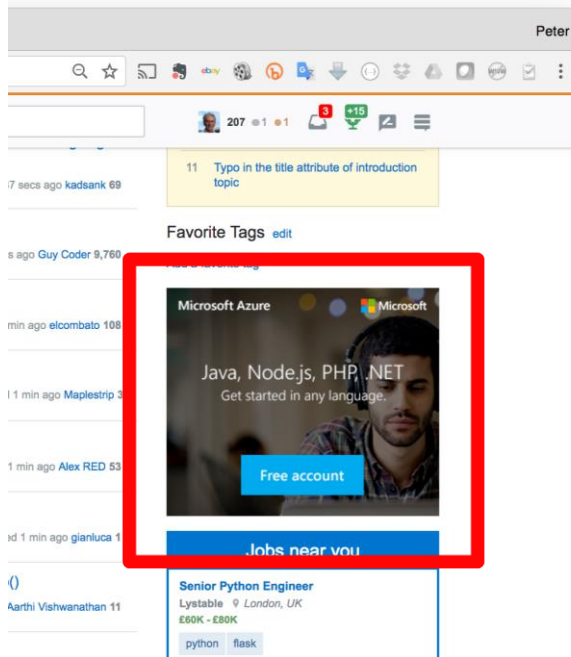
- **Periodic** reporting  
e.g Daily/Weekly/Monthly
- Statistical analysis
- Latency: **0.5 – 10 minutes**



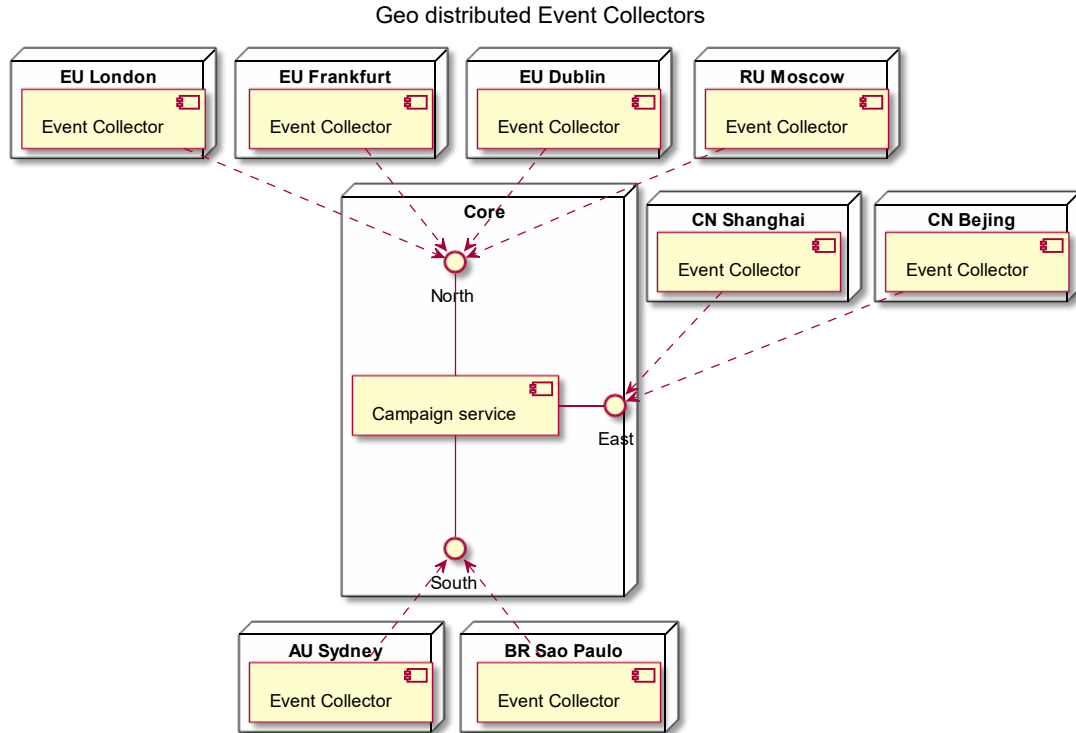
# Ad events

Represent the **value** of the campaign a.k.a **Money**

- Measure the **effectiveness** of the **campaign**.
- **Insights** about the consumer - **Data Science**



# Event Collection

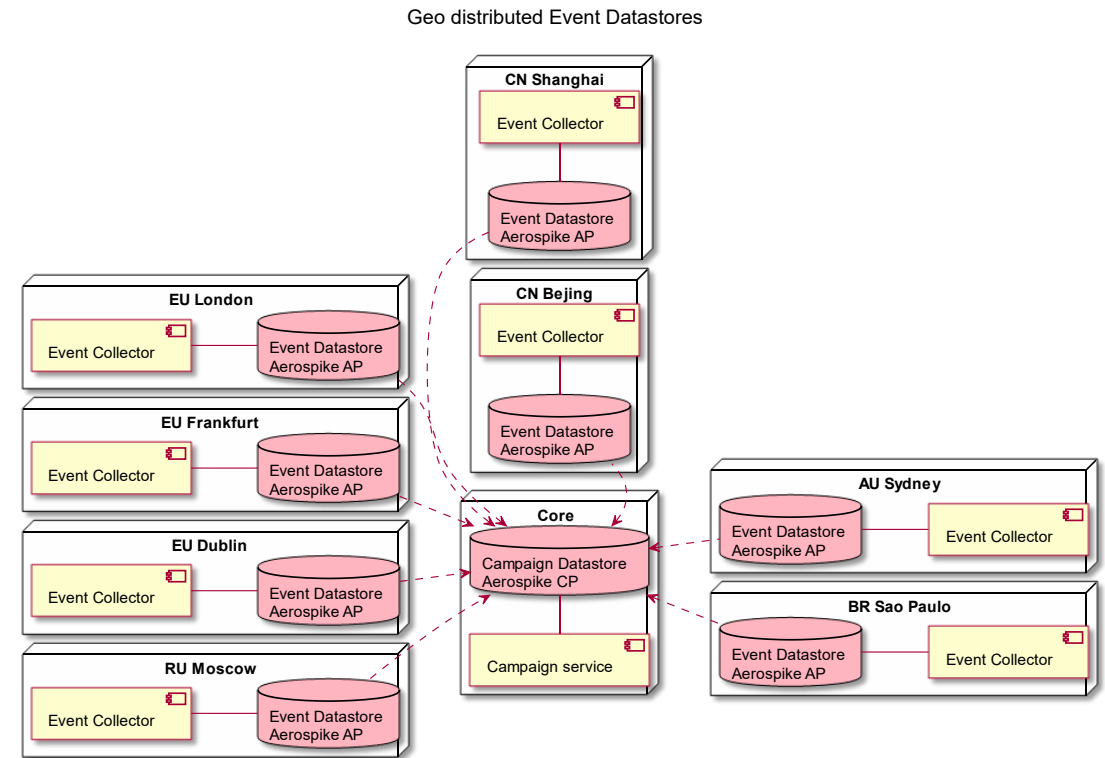


## Event collector

- Super **simple web service** API
- **Collects a portion** of events per Geo location
- **Feeds to** the Campaign service
- Linear and elastic scalable
- No state

# Event Collection with Ad Event datastore

- Datastore acts as buffer
- Lower latency for event collector instance
  - 1-2 ms
- Fewer Event collector instances
- Adds reliability





# Real-time Aggregation and Reduction

## Event collector

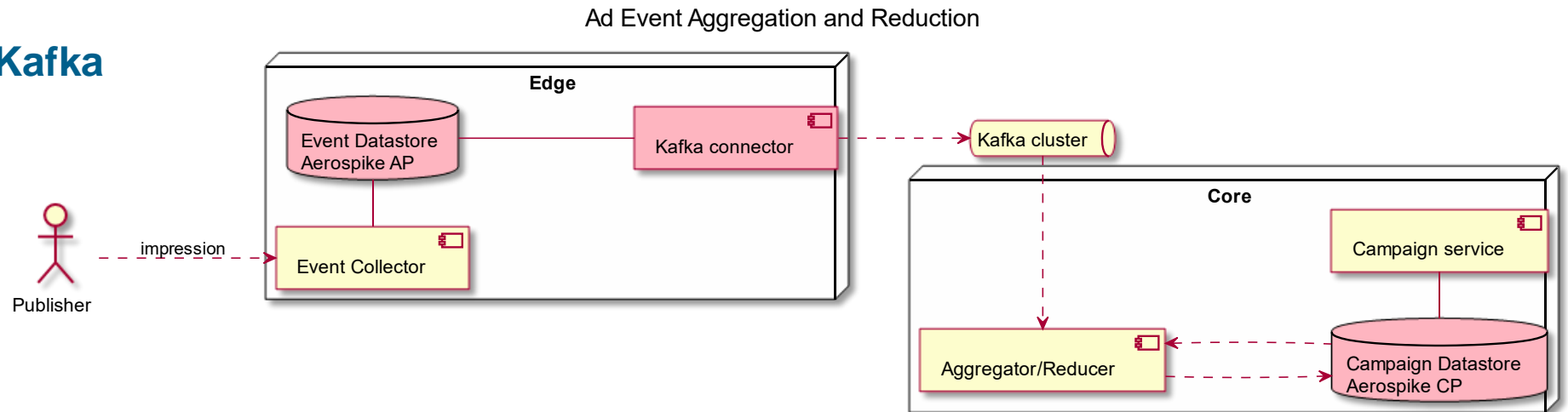
- **Receives** events
- Writes to the **local datastore** with
  - Tag id
  - Event type (click, etc)
  - Timestamp

## Aerospike Kafka Connector

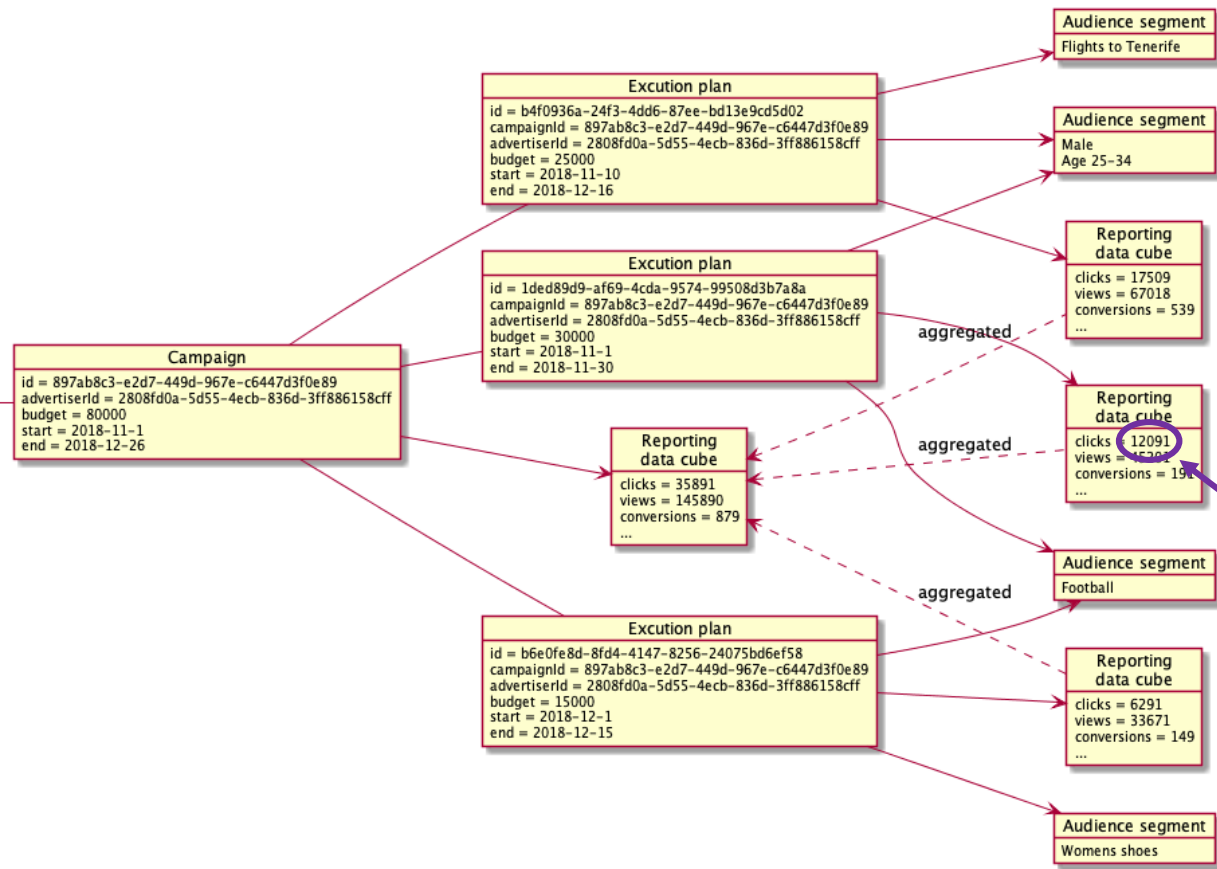
- **Publishes** events to **Kafka**

## Aggregator/Reducer

- Receives **event messages** from the **Kafka** topic
- Maps the **Tag** to the **Campaign/Execution Plan**
- Atomically Increments **counters** (etc) in the **data cube** of the Campaign/Execution plan
- Uses **Aerospike CDT** operations



# Aerospike Complex Data Type Operations



CDT is like a “document”

- Atomic **operation** on a **path** to the **data** element
- Latency < 1ms
- Durable

Example

- **Operation:** Increment (and read)
- **Path:**

`campaign.execution-plan.data-cube.clicks`

# Aerospike datastore summary

---

## User profile datastore

- Size: ~8-12 billion records
- Throughput: ~ 7 million / sec
- Latency: ~0.9 ms
- CAP: Availability, Partitioning
- Uptime: 24/7/365

## Events datastore

- Size: 40-80 million records
- Throughput: ~2-3 hundred thousand / sec
- Latency: ~0.9 ms
- CAP: Availability, Partitioning
- Uptime: 24/7/365

## Campaign datastore

- Size: 40-80 million records
- Throughput: ~2-3 hundred thousand / sec
- Latency: ~0.9 ms
- CAP: Consistency, Partitioning

The logo for Aerospike, featuring the word "AEROSPIKE" in white, uppercase, sans-serif font, with a stylized white arrowhead pointing left, all set against a solid red rectangular background.

AEROSPIKE

*Note: Actual sizes vary seasonally*

adform

# Uncle Pete's Advice

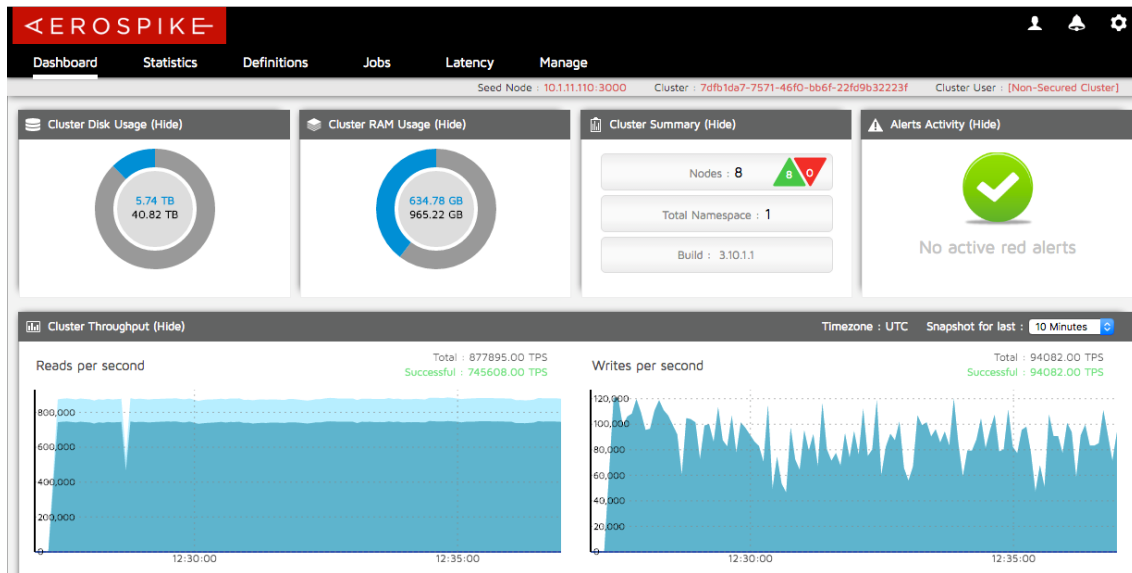
# Why Aerospike



cassandra



AEROSPIKE



## Secret weapon: **Aerospike**

- **8-9 million** reads/s, **1-2 million** writes (Over all clusters)
- **< 1ms latency** to read profile
- **~9 billion** cookies in 12 servers (largest cluster)
- HA with redundancy
- Moved **from Cassandra 5** years ago
  - **Easy** to use
  - **Works** as advertised
  - **Value** for money
- **~150** servers → **12** (largest cluster)

# Open source & Enterprise products

---



## Open source != Free

- Modern **escrow**
- Fork the **source repository**
- **Build** it yourself
- **Survive** on your own



## Enterprise is value for money

- Support
- Training
- Consulting
- Contractual obligation - **A neck to choke**

adform

*Questions?*