# Data LossLess streaming processing
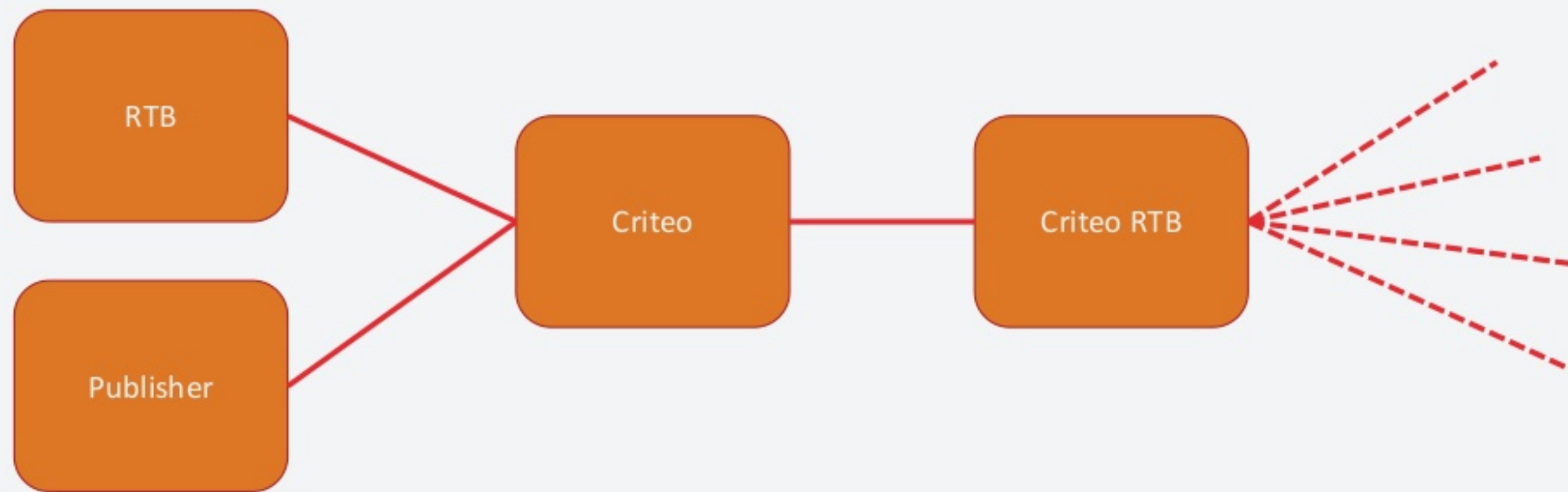
Oleksandr Nitavskyi, Criteo

# Criteo: business model

# Criteo: in numbers

**16 DATA CENTERS**

**MORE THAN**
**25 000 SERVERS**

+15K

PEAK TRAFFIC > **4 mln** QPS

**2 CDH5 CLUSTER OF 4000 NODES** (92TB, 200 GB RAM, 48 cores)

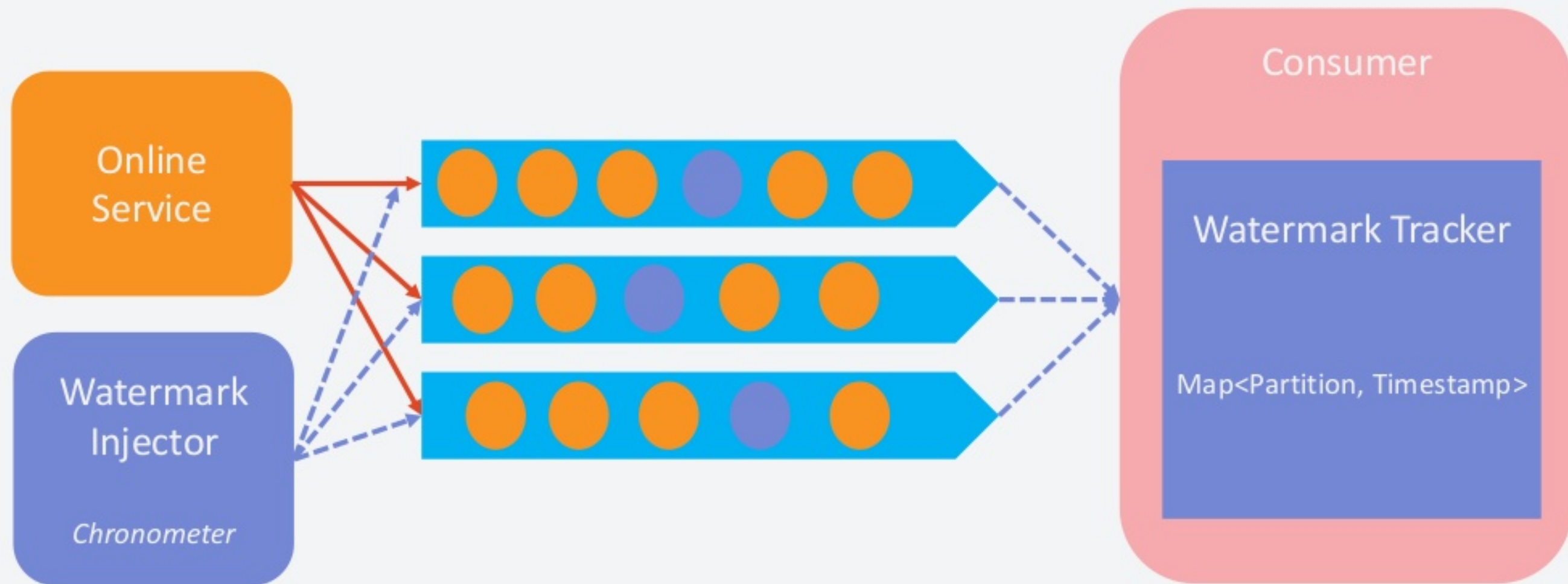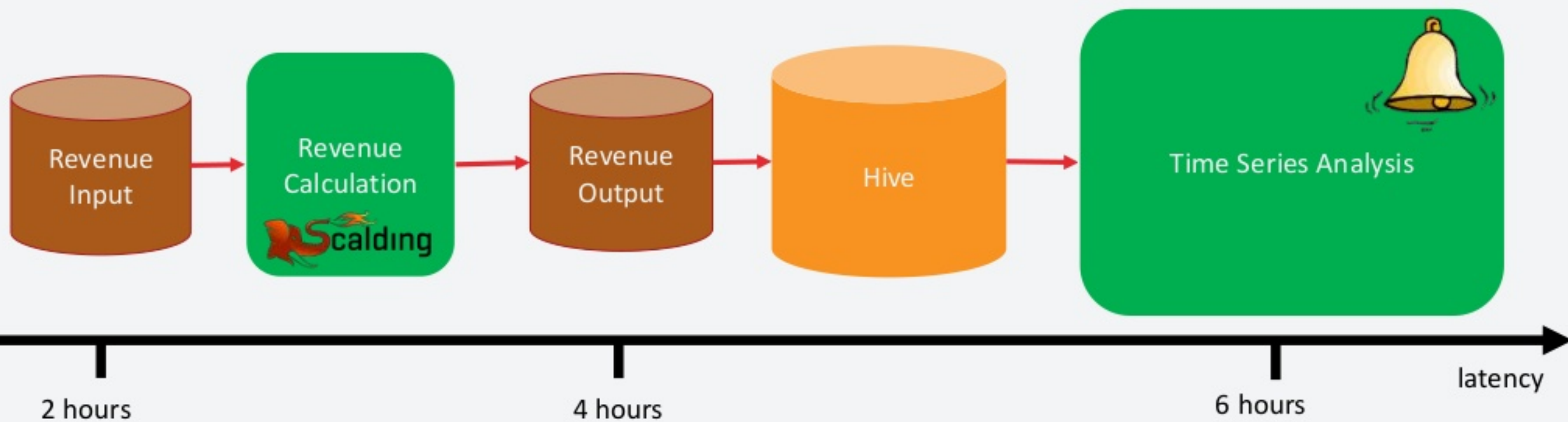# Gathering all data together

# Time series management



- Single Kafka partition is ordered

- Time of partition **max(watermark)**
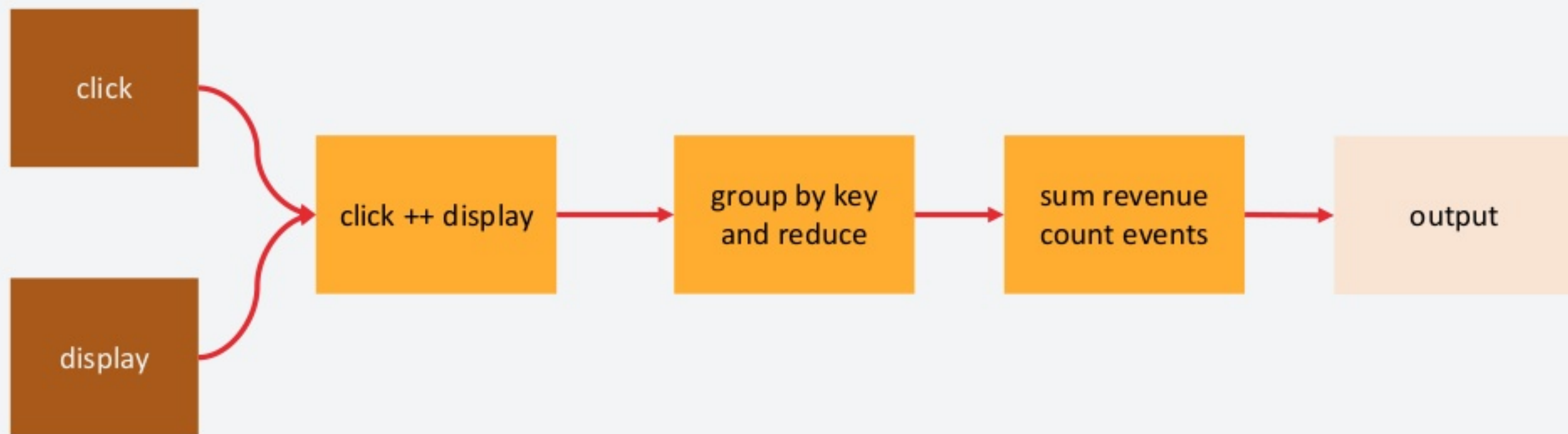
- Global time **min(partitions)**

# Time series management

# Revenue anomalies offline
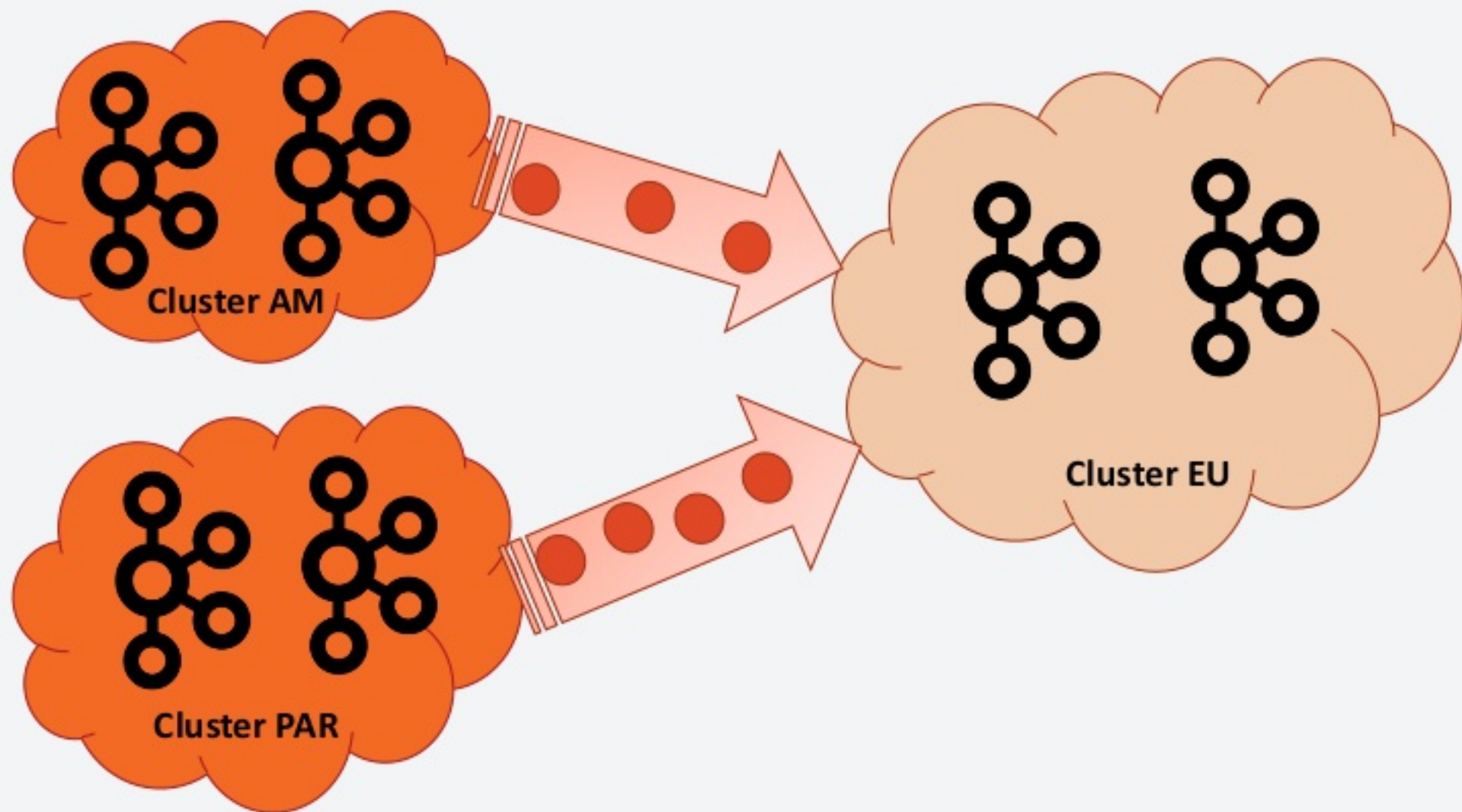


Revenue Input → Revenue Calculation (Scalding) → Revenue Output → Hive → Time Series Analysis

latency

2 hours      4 hours      6 hours

# Business logic

# Stream processing infrastructure



Cluster AM

Cluster PAR

Cluster EU

# Stream processing infrastructure

- A rogue Kafka client can significantly reduce Kafka cluster performance

- We need to shuffle the data uniformly

- A streaming pipeline needs to have platform level (EU, AS, …) data
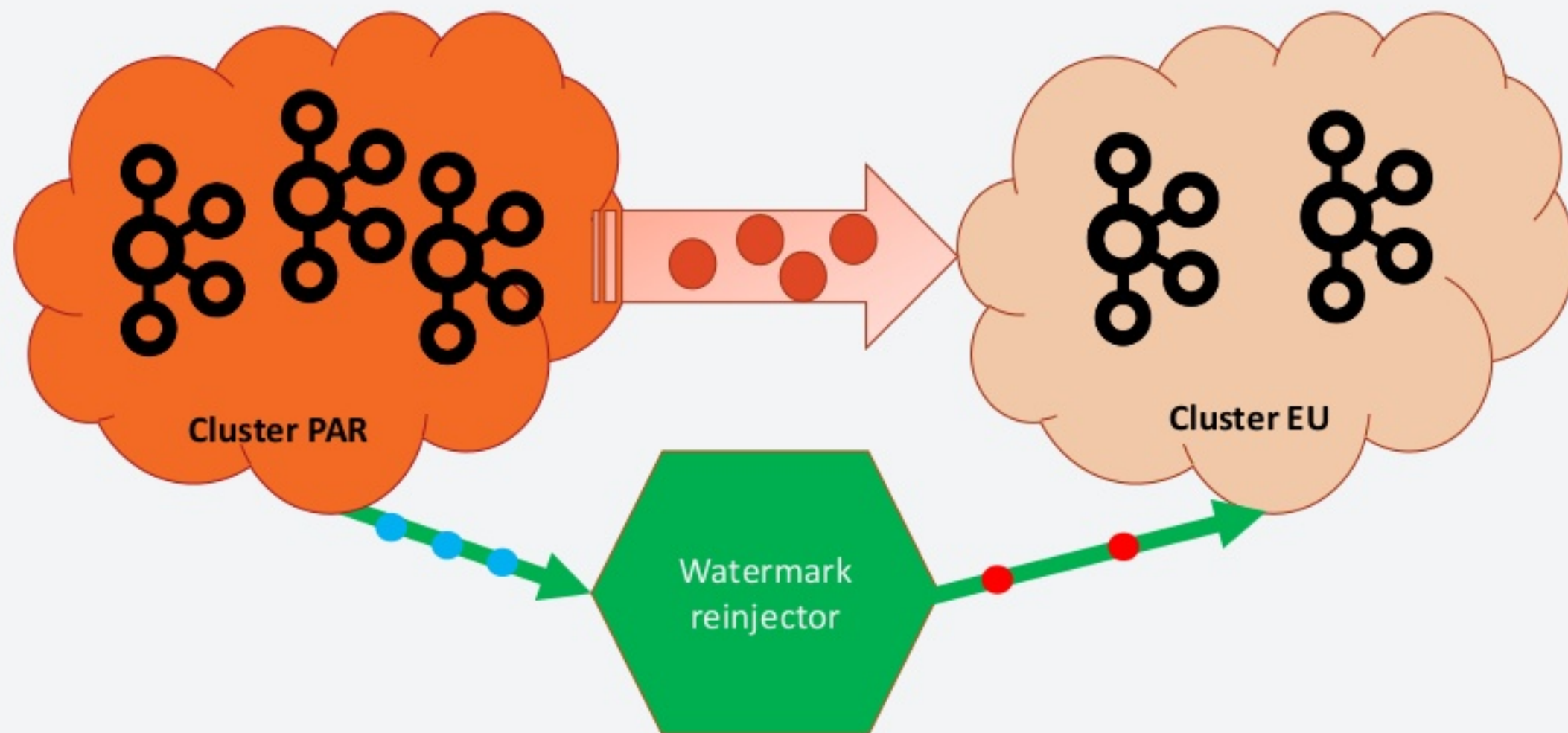
criteo.

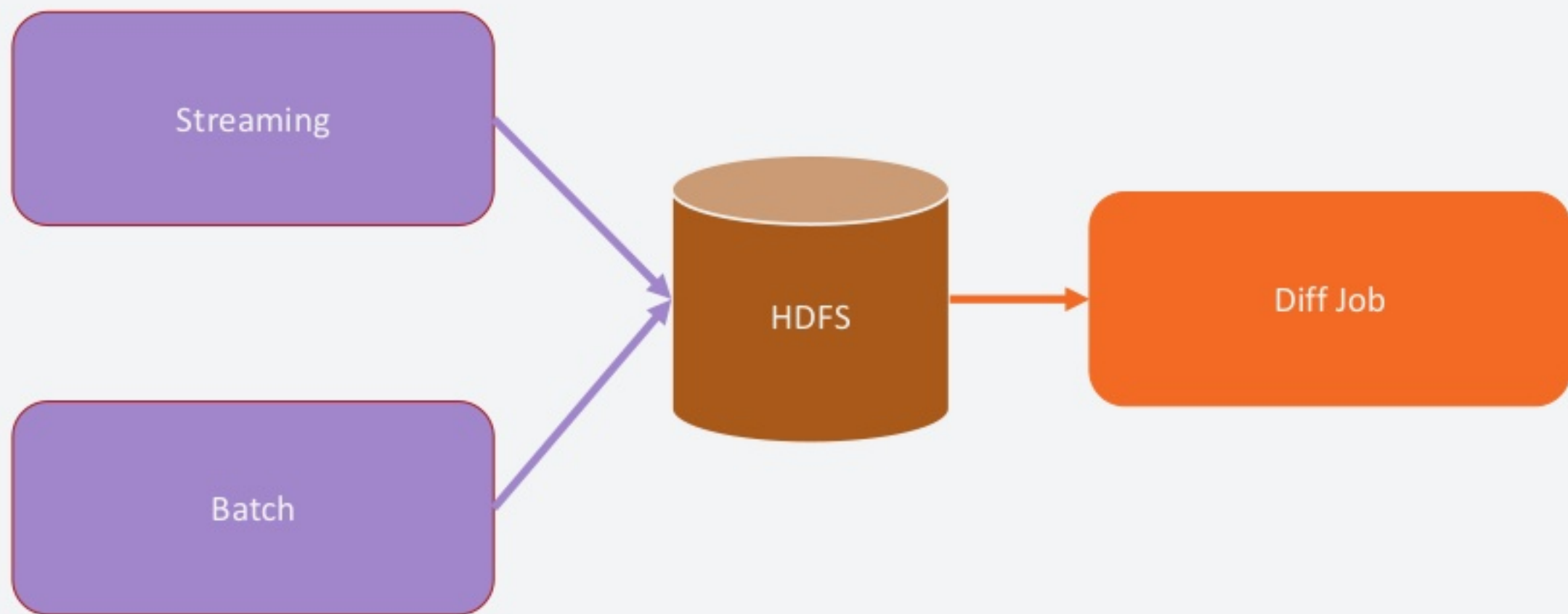# Stream processing first try

# Streaming done right

# Streaming done right

- Flexible event time support

- Robust state management

- Short development cycle

criteo.

# Accurate event time processing



Cluster PAR

Cluster EU

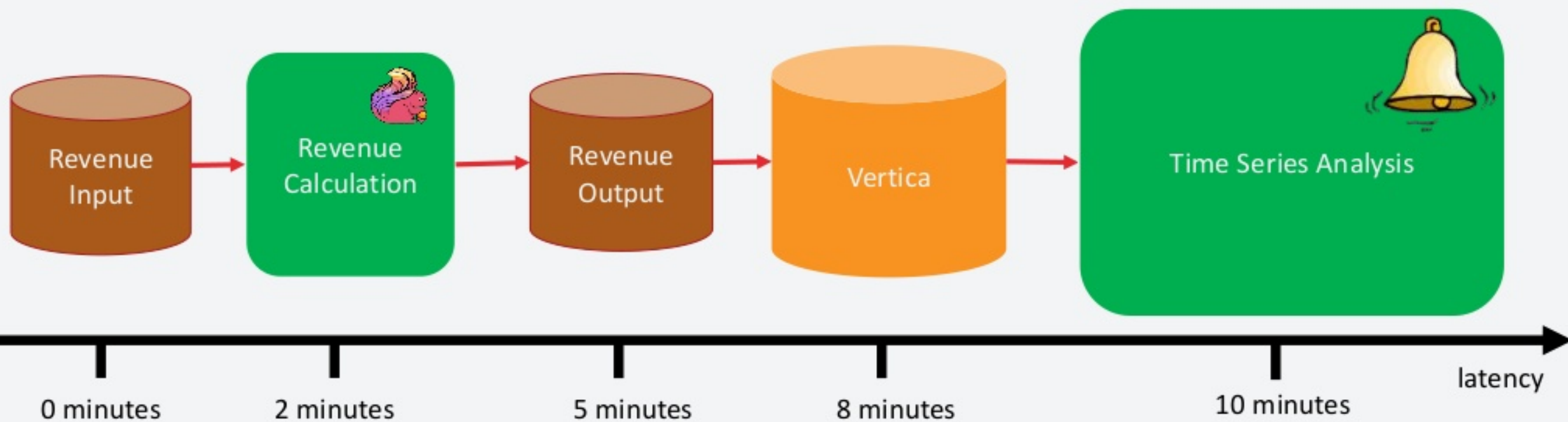Watermark reinjector

# Result verification

# Correctness results

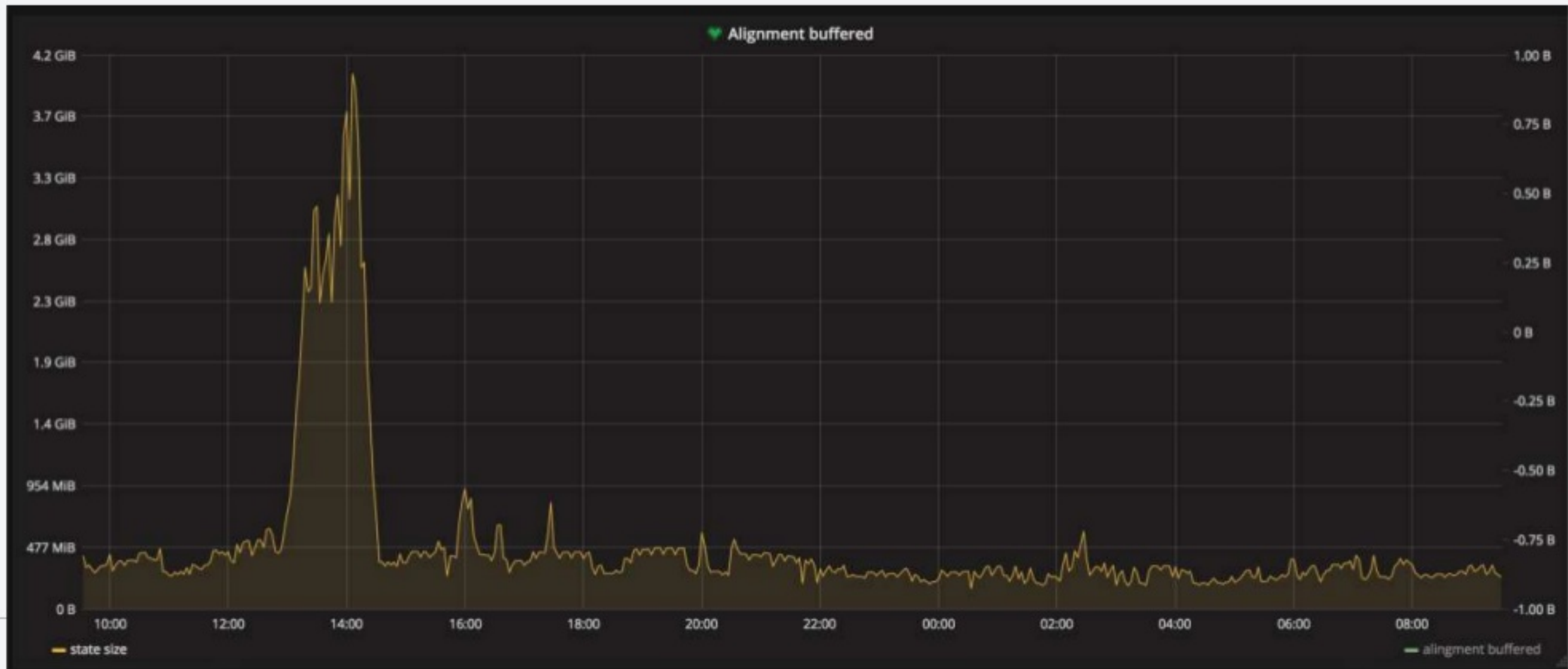## Happy path looked great

- no watermark issues

- small duplication rate

- no lost state of the Flink

- no big latency

- no spike of late events
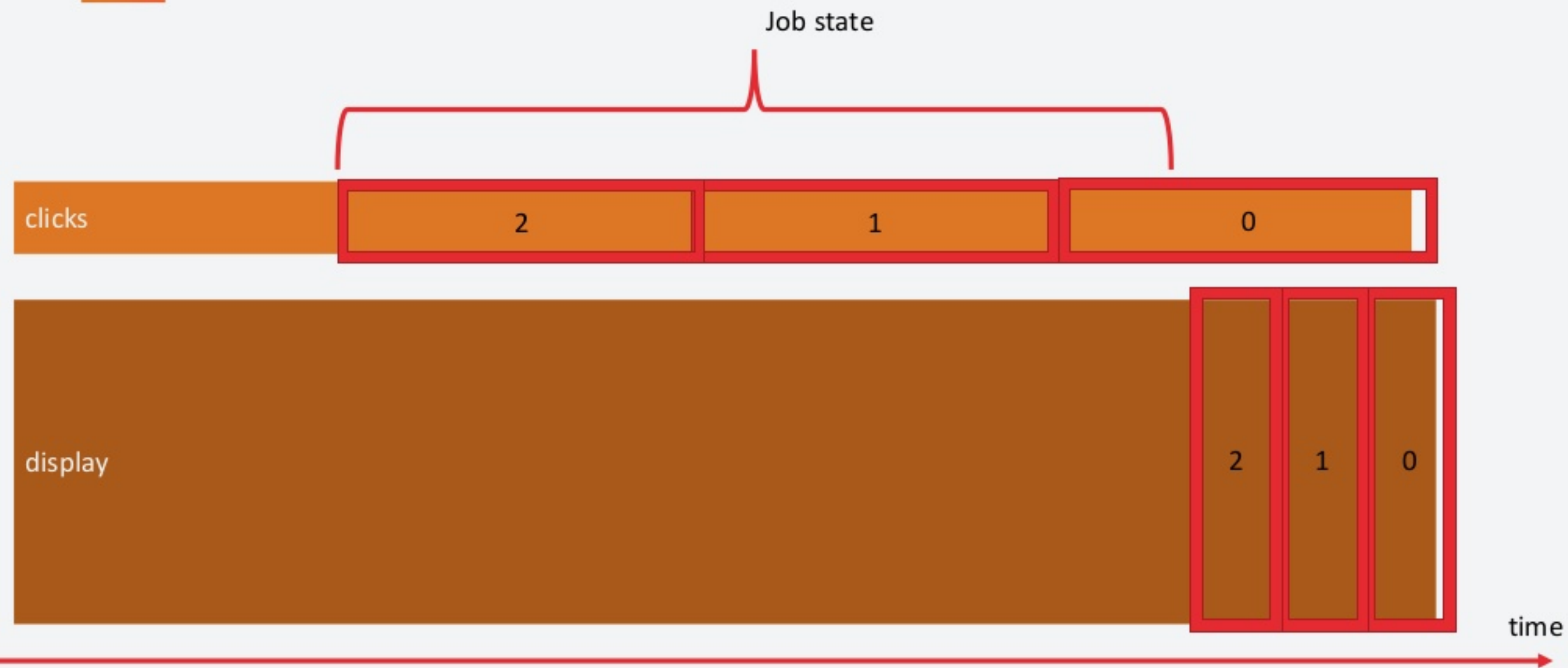
99,9%

# Revenue anomalies online



| Revenue Input | Revenue Calculation | Revenue Output | Vertica | Time Series Analysis |

0 minutes → 2 minutes → 5 minutes → 8 minutes → 10 minutes → latency

# Job state

# Problems with a Kappa I

# More generic problems with a Kappa

partition 0

partition 1

partition 2

partition 3

partition 4

# Streaming enables use-cases

- Our TSE can monitor any ads campaign changes in "realtime"

- We have enabled a short term campaigns which takes only several hours to run

- We efficiently cut campaigns spending when campaign is over

criteo.

# Notes to take

- Simple business use-case can allow to focus on the infrastructure

- Store watermarks within your data for catching up

- Flink is an enabler

"

# Thank you

we are hiring

criteo.