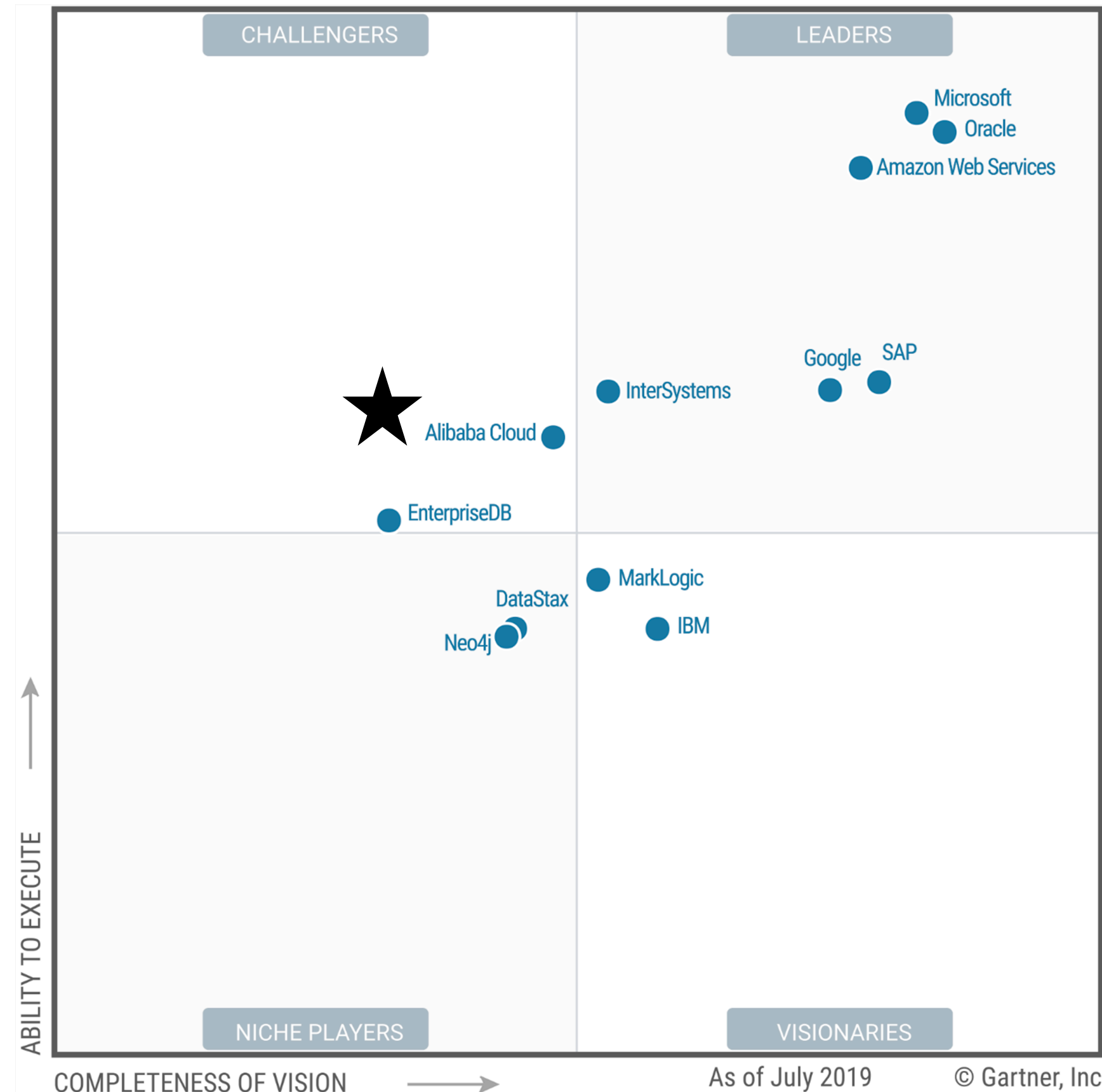# Timon: A Timestamped Event Database for Efficient Telemetry Data Processing and Analytics

Wei Cao, **Yusong Gao**, Feifei Li, Sheng Wang, Bingchen Lin, Ke Xu, Xiaojie Feng, Yucong Wang, Zhenjun Liu, Gejin Zhang

Alibaba Cloud Database Department
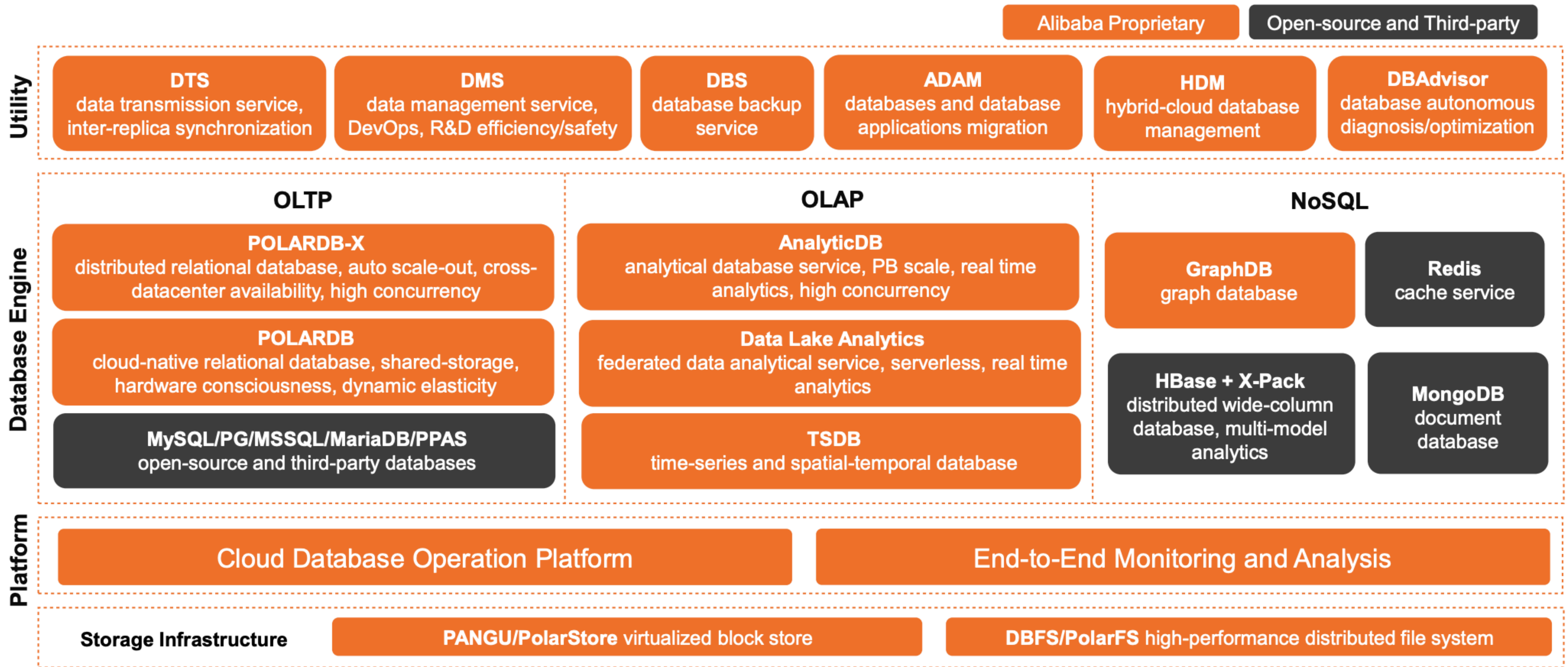
2020.06

# Database services at Alibaba Cloud



As of July 2019    ©Gantner,inc

**Cloud Database Marketing:**

$1^{st}$, Asia Pacific

**Database Products and Services:**

26 Products or Services
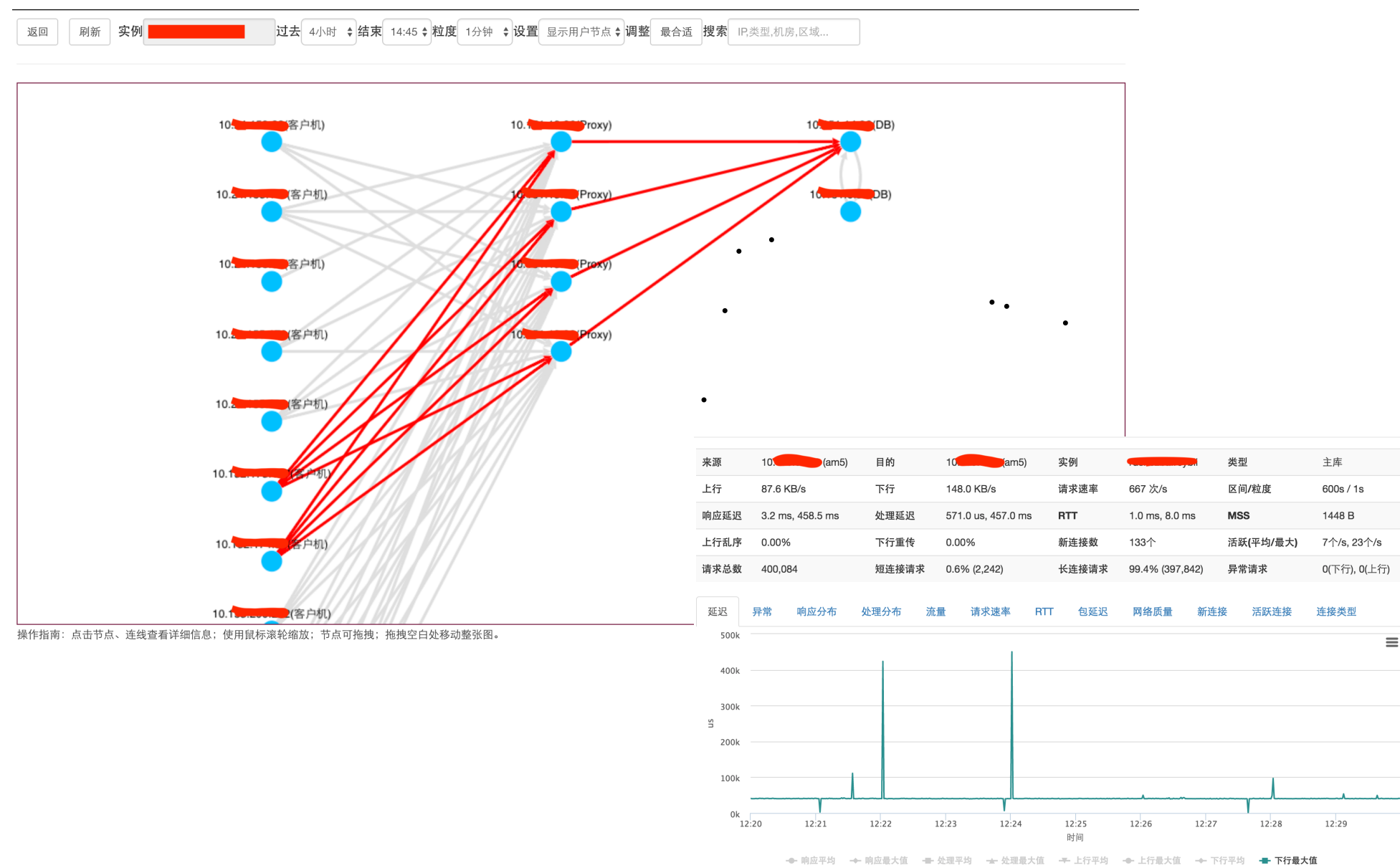
**Enterprise Users:**

100 thousands

**Databases Migrated:**

400 thousands

Expressway for running fully managed databases on Alibaba Cloud.
https://www.youtube.com/watch?v=5VkLDC_uIxM

# Database systems and services at Alibaba and Alibaba Cloud

Alibaba Cloud

Alibaba Proprietary | Open-source and Third-party

## Utility

**DTS**
data transmission service, inter-replica synchronization

**DMS**
data management service, DevOps, R&D efficiency/safety

**DBS**
database backup service

**ADAM**
databases and database applications migration

**HDM**
hybrid-cloud database management

**DBAdvisor**
database autonomous diagnosis/optimization

## Database Engine

### OLTP

**POLARDB-X**
distributed relational database, auto scale-out, cross-datacenter availability, high concurrency

**POLARDB**
cloud-native relational database, shared-storage, hardware consciousness, dynamic elasticity

**MySQL/PG/MSSQL/MariaDB/PPAS**
open-source and third-party databases

### OLAP

**AnalyticDB**
analytical database service, PB scale, real time analytics, high concurrency

**Data Lake Analytics**
federated data analytical service, serverless, real time analytics

**TSDB**
time-series and spatial-temporal database

### NoSQL

**GraphDB**
graph database

**Redis**
cache service

**HBase + X-Pack**
distributed wide-column database, multi-model analytics

**MongoDB**
document database

## Platform

**Cloud Database Operation Platform**

**End-to-End Monitoring and Analysis**

### Storage Infrastructure

**PANGU/PolarStore** virtualized block store

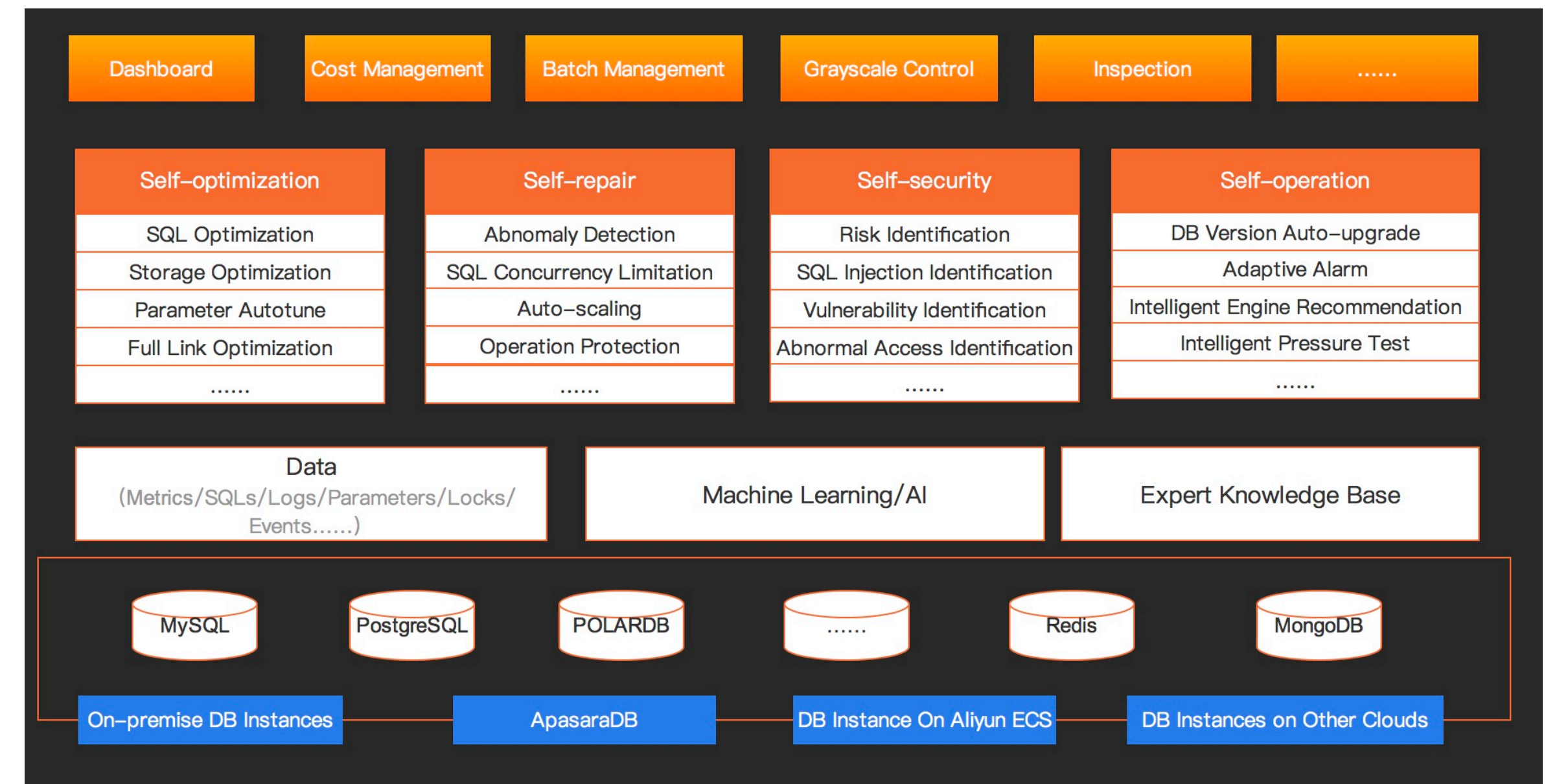**DBFS/PolarFS** high-performance distributed file system

# The scenarios of using telemetry data

## End-to-End Tracking System



SIGMOD'18 TcpRT: Instrument and Diagnostic Analysis System for Service Quality of Cloud Databases at Massive Scale in Real-time.

## Database Autonomy Service (DAS)



https://www.alibabacloud.com/help/doc-detail/64851.htm?spm=a2c63.p38356.b99.2.61d09bb0Xe1MPU

# Challenge of processing telemetry data in cloud database services

- ~10 million objects

  - Cover database engine, network, operating system, and even each individual OS process.

- Support 1 second granularity and hundreds of millions data points per second

  - Find and explain peaks which last a short time.

- Support long-term & multiple granularities queries with low response latency

  - Find trend and periodicity and compare with historical data

- Strict SLA, realtime & accurate even when out-of-order events exist

  - Find and explain anomalies as soon as possible.

  - Not just for monitoring but also for autonomous optimization.

# Sources of out-of-order events

- Distributed Computing Environment

  - Machine failures
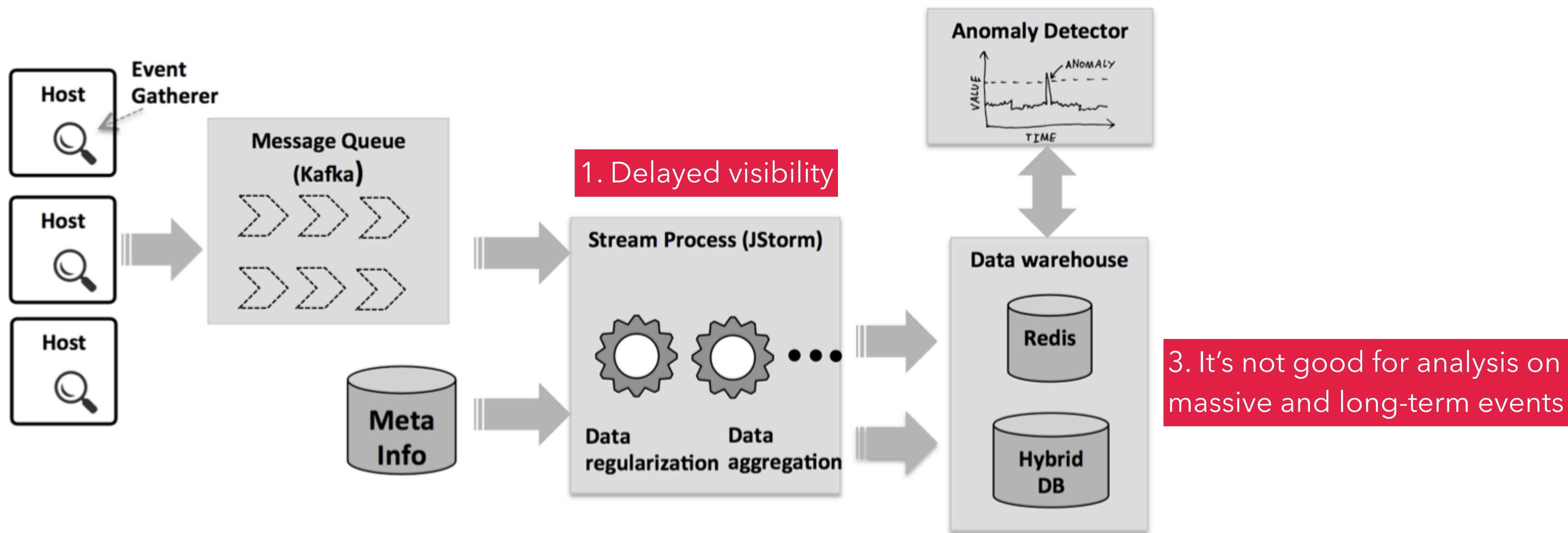
  - Network failures

  - Clock skew

  - etc...

- Application Natures

  - For instance: TcpRT events can't be collected until requests complete. (53.90% out-of-order events)

SQL-1    $T_2$ _____ $T_4$

SQL-2    $T_3$ _____ $T_5$

SQL-3    $T_1$ _____ $T_6$

→ Time

The order of events arrival will be $T_2$, $T_3$, $T_1$

# Our first generation data processing system

# Solution: using blind-write and moving aggregate to storage
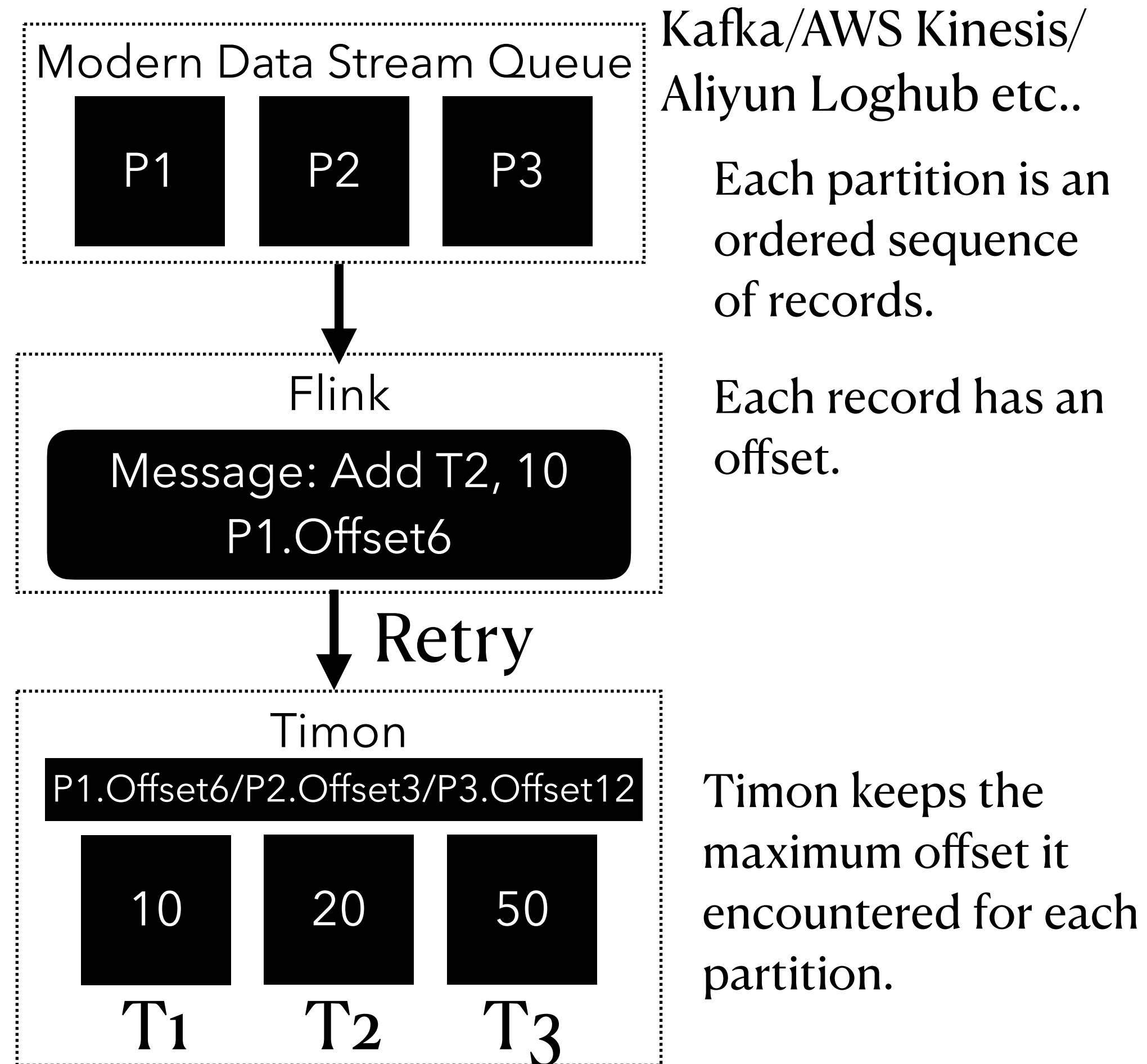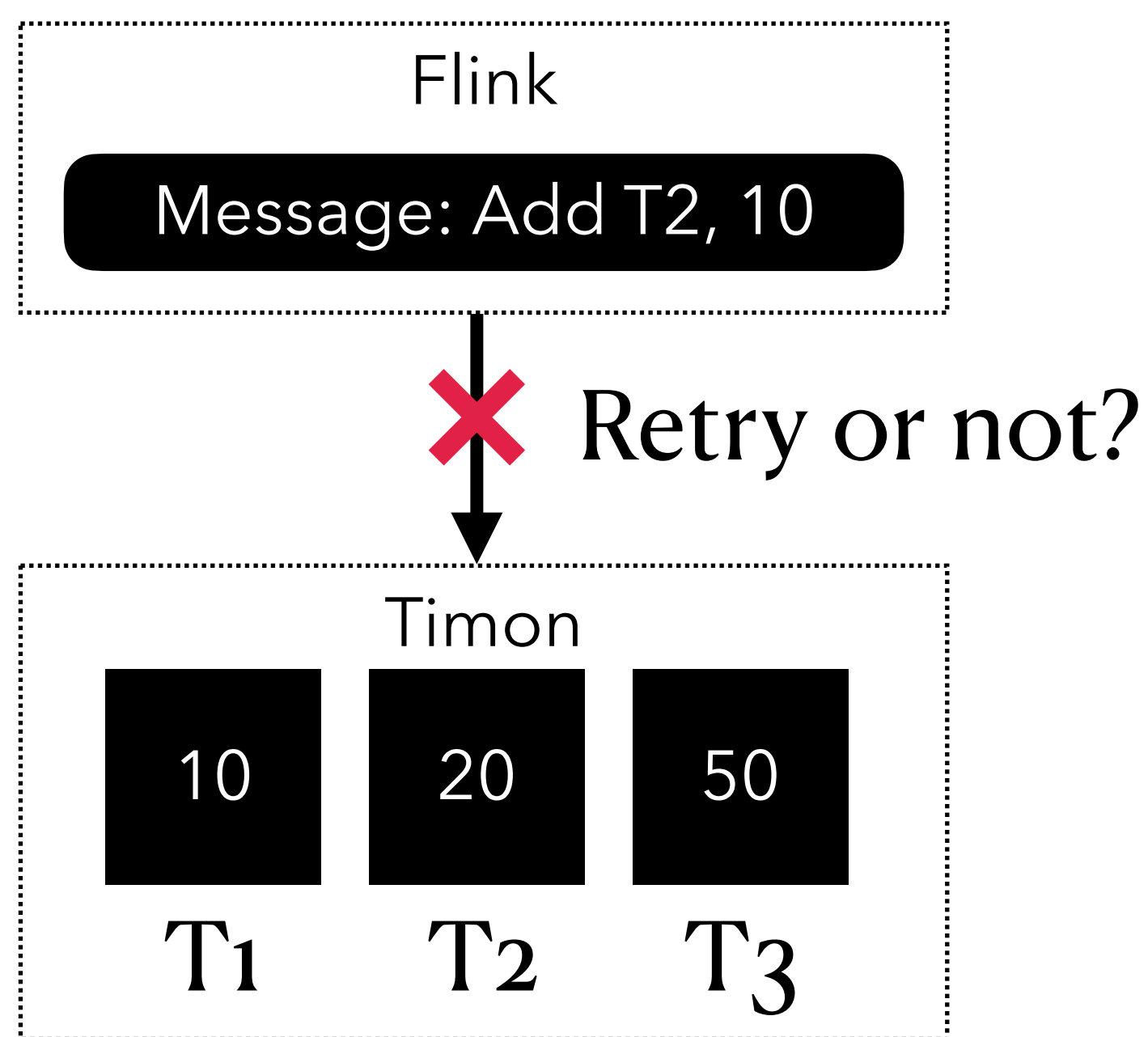
- The most common operators are associative and commutative.
  - **sum**, **max**, **min**, **avg**, **stddev**
  - **quantile**: histogram, t-digest*
  - **distinct**: HyperLogLog

- So it's a good idea to support incremental processing after a data point is written.

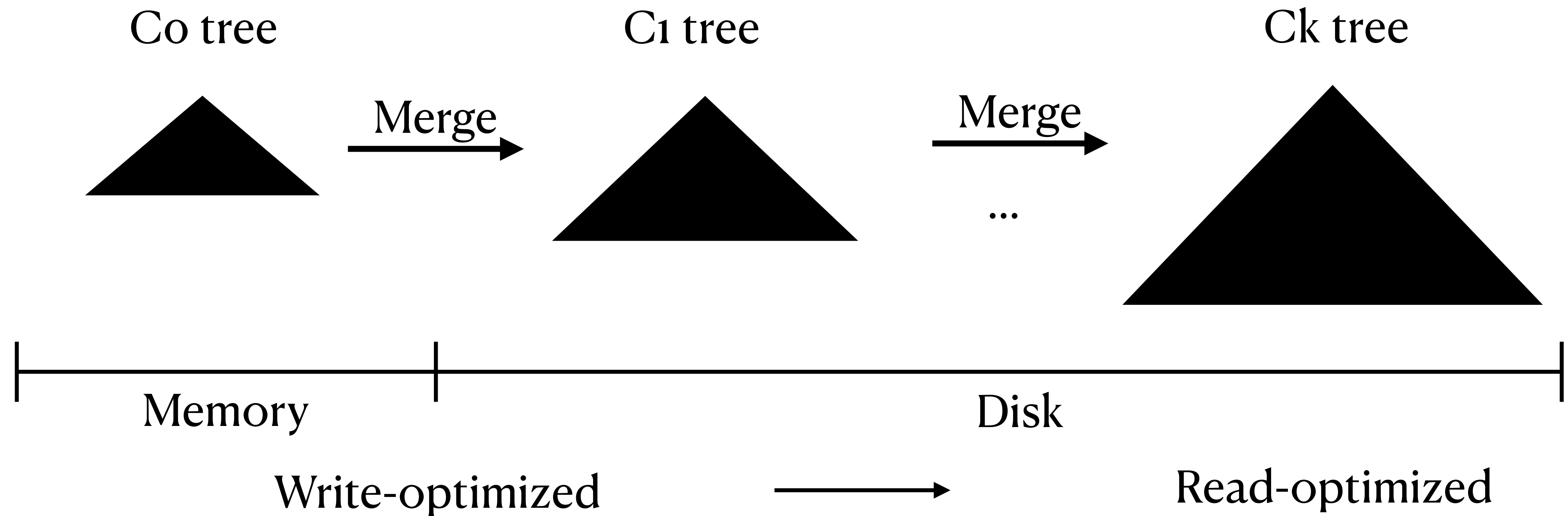\* t-digest: https://github.com/tdunning/t-digest

# Incremental Processing Issue: Idempotence

- If failovers occur while the data points are being written, we have to remove the points that have been successfully written.

# How to support blind-write and incremental processing efficiently
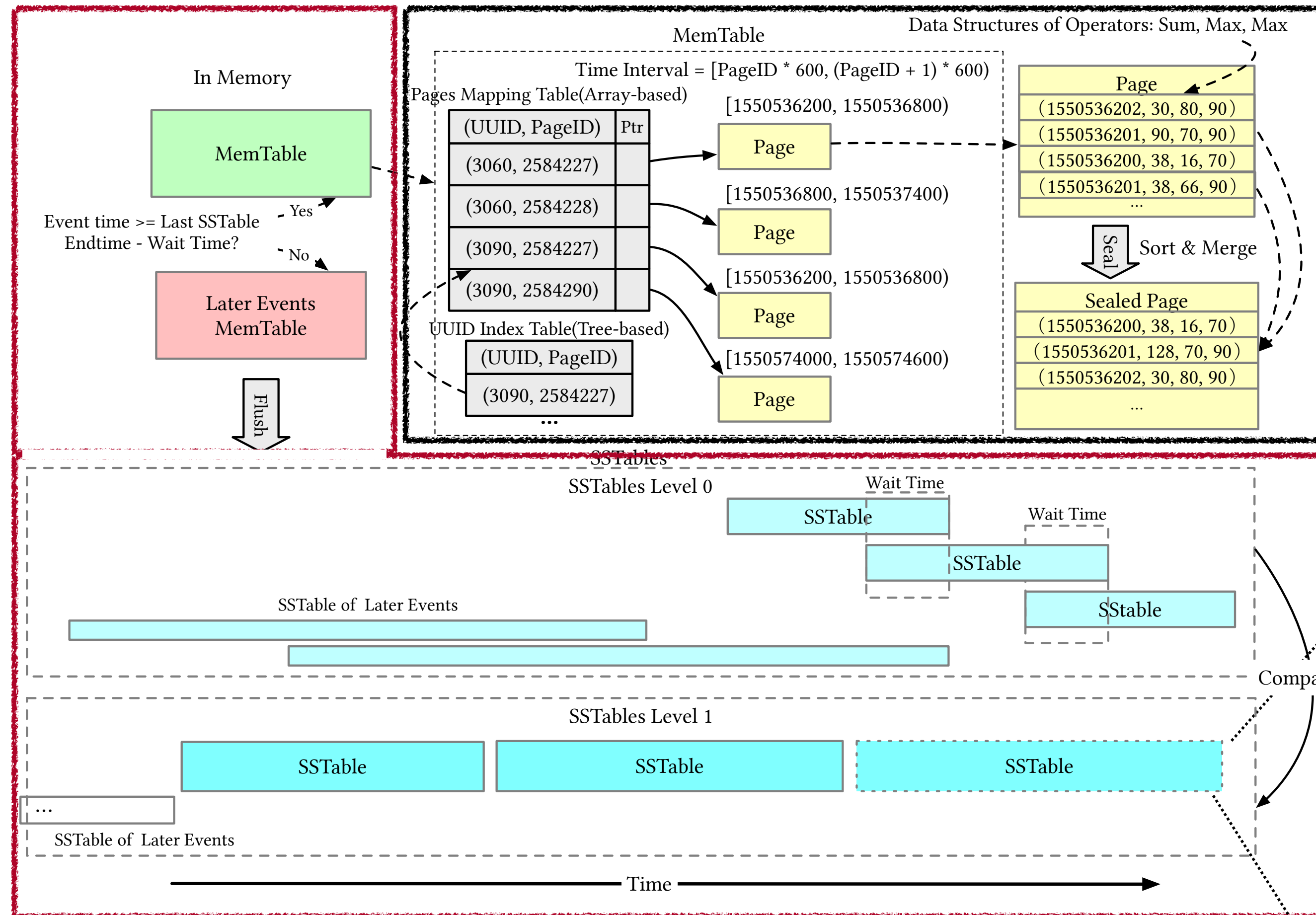
- LSM-Tree by P O'Neil - 1996

Co tree

Merge →

C1 tree

Merge →

...

Ck tree

Memory | Disk

Write-optimized →→→ Read-optimized

Append & merge to tolerate out-of-order events.
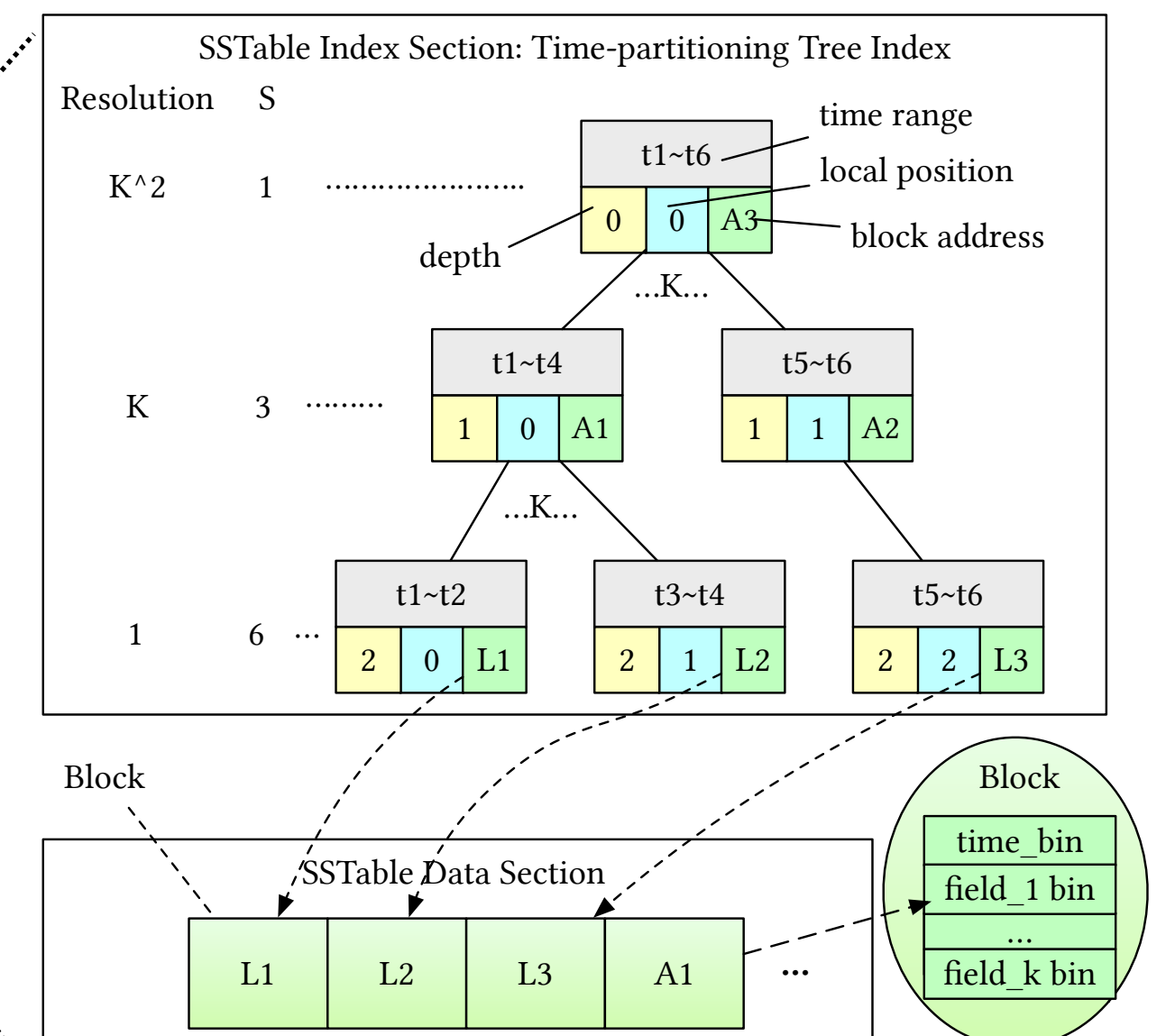
Build time-partitioning tree index for long-term queries.
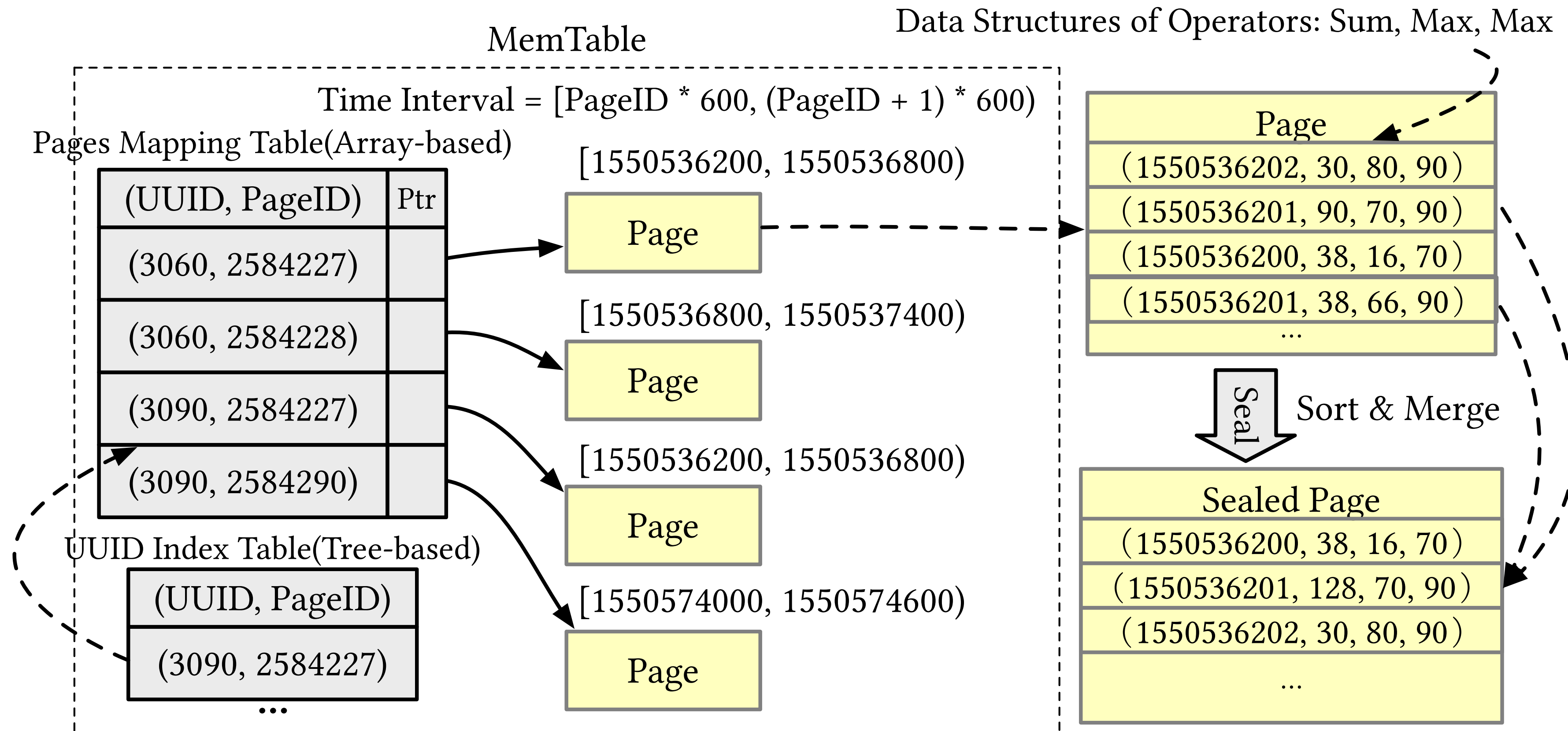
# Time-Segment Log-Structured Merge-Tree



MemTable is optimized for time series data.

Build time-partitioning tree index when compacting.

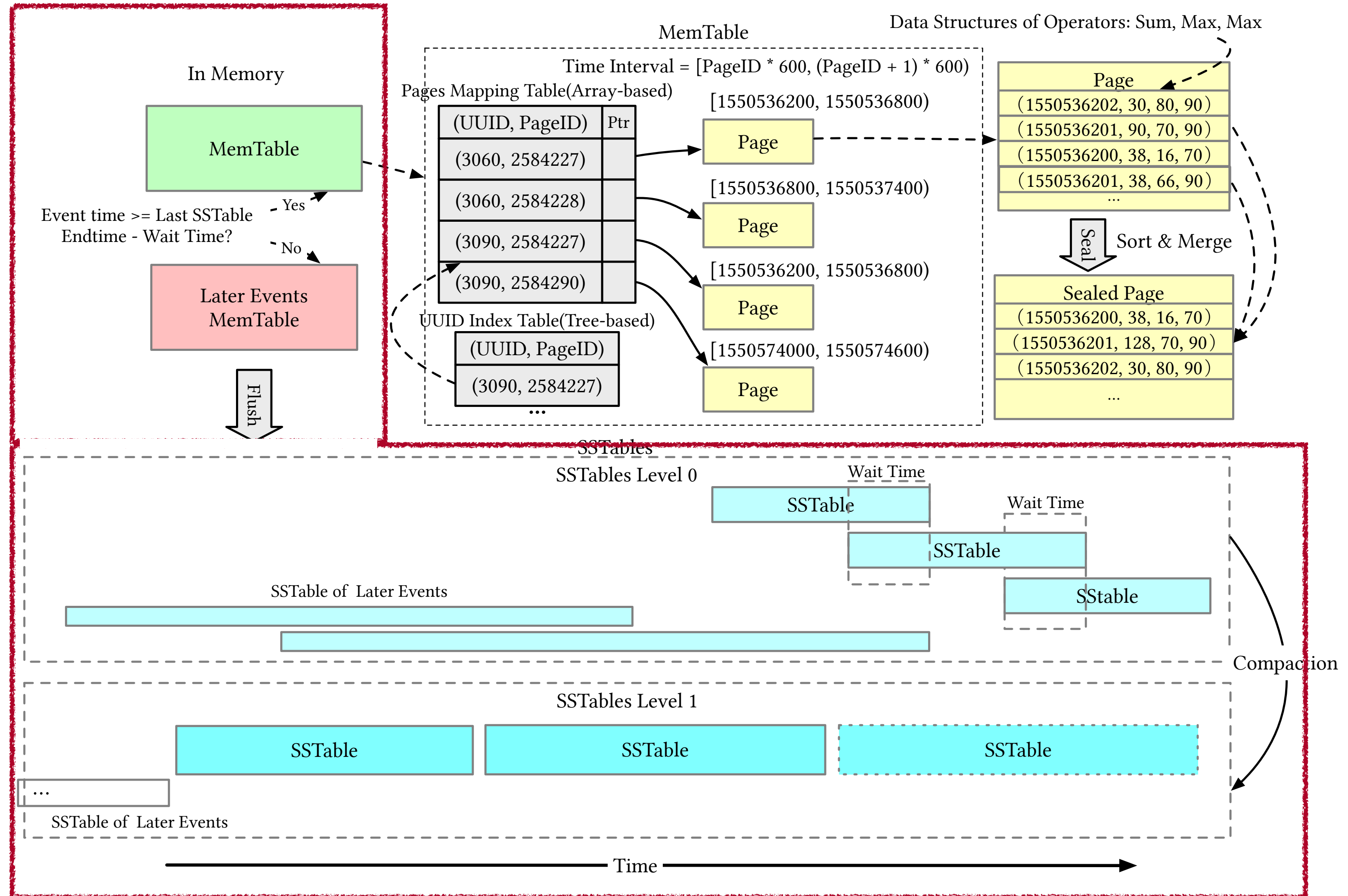Lazy Merge

# Optimized MemTable for Time Series

- Hybrid Tree-based Structure and Array-based Structure

- Pages are used to process sparse events.

MemTable

Data Structures of Operators: Sum, Max, Max

Time Interval = [PageID * 600, (PageID + 1) * 600)

Pages Mapping Table(Array-based)

[1550536200, 1550536800)

| (UUID, PageID) | Ptr |
|---|---|
| (3060, 2584227) | |
| (3060, 2584228) | |
| (3090, 2584227) | |
| (3090, 2584290) | |

Page

[1550536800, 1550537400)

Page

UUID Index Table(Tree-based)

[1550536200, 1550536800)

Page

| (UUID, PageID) |
|---|
| (3090, 2584227) |

[1550574000, 1550574600)

Page

...

**Page**
| | | |
|---|---|---|
| (1550536202, 30, 80, 90) | | |
| (1550536201, 90, 70, 90) | | |
| (1550536200, 38, 16, 70) | | |
| (1550536201, 38, 66, 90) | | |
| ... | | |

Seal → Sort & Merge

**Sealed Page**
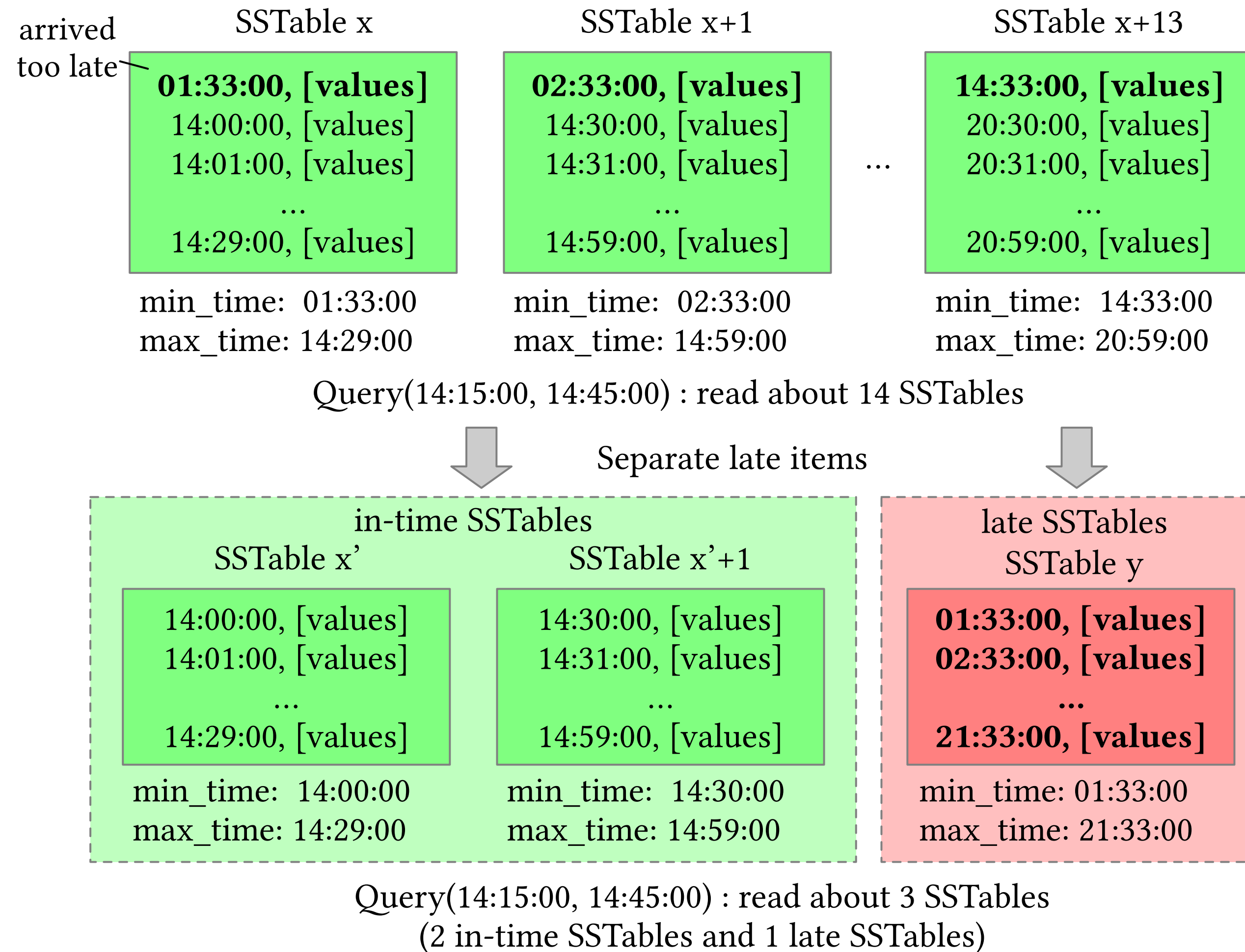| | | |
|---|---|---|
| (1550536200, 38, 16, 70) | | |
| (1550536201, 128, 70, 90) | | |
| (1550536202, 30, 80, 90) | | |
| ... | | |

# Tolerating out-of-order events (1)

- Most of out-of-order events arrive with a delay of less than 5 minutes.

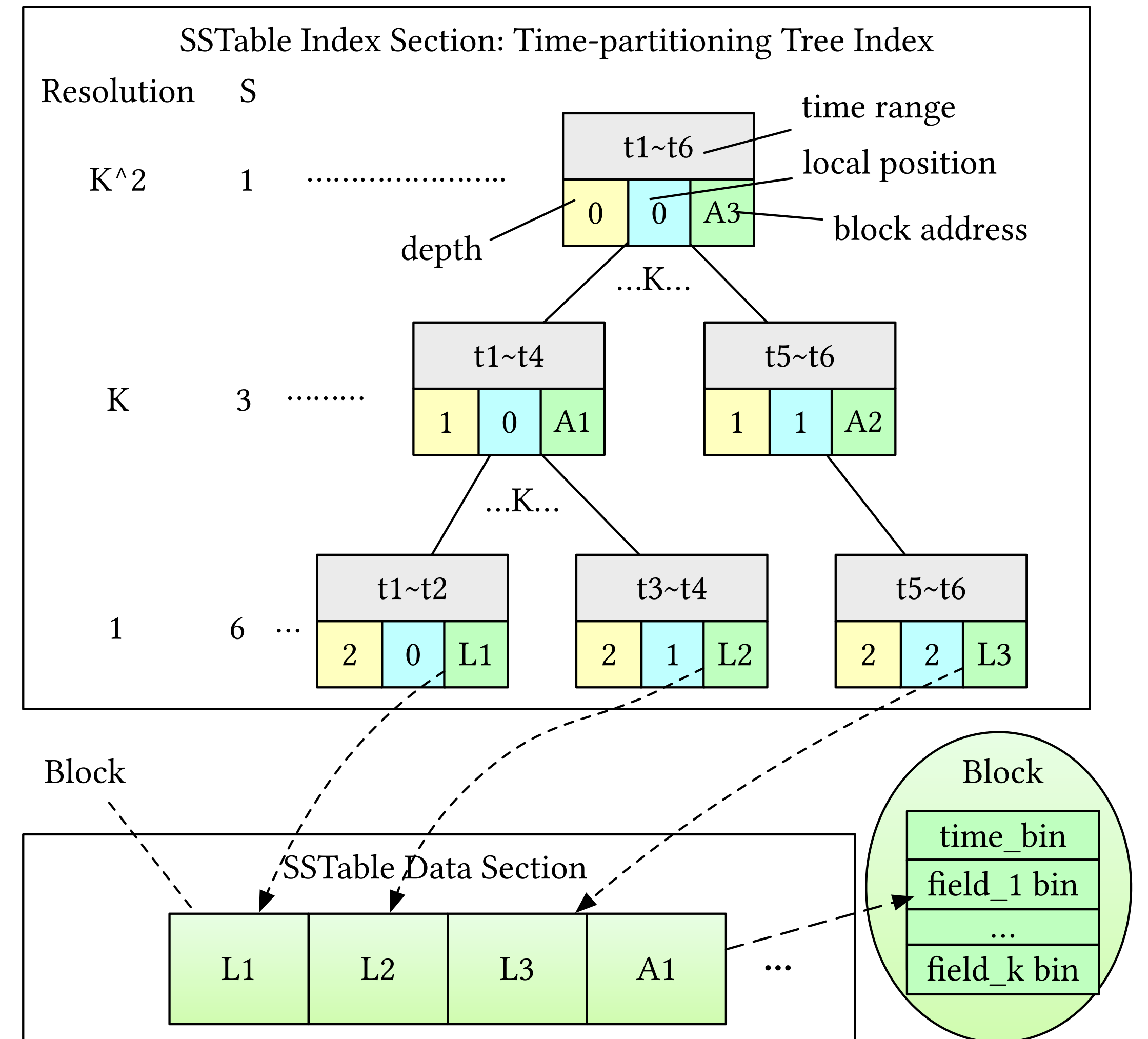- Some latecomers are small in proportion but wide in time range.

# Tolerating out-of-order events (2)

arrived too late

| SSTable x | SSTable x+1 | SSTable x+13 |
|---|---|---|
| **01:33:00, [values]**<br>14:00:00, [values]<br>14:01:00, [values]<br>...<br>14:29:00, [values] | **02:33:00, [values]**<br>14:30:00, [values]<br>14:31:00, [values]<br>...<br>14:59:00, [values] | **14:33:00, [values]**<br>20:30:00, [values]<br>20:31:00, [values]<br>...<br>20:59:00, [values] |
| min_time:  01:33:00<br>max_time: 14:29:00 | min_time:  02:33:00<br>max_time: 14:59:00 | min_time:  14:33:00<br>max_time: 20:59:00 |

...

Query(14:15:00, 14:45:00) : read about 14 SSTables

Separate late items

### in-time SSTables

| SSTable x' | SSTable x'+1 | late SSTables<br>SSTable y |
|---|---|---|
| 14:00:00, [values]<br>14:01:00, [values]<br>...<br>14:29:00, [values] | 14:30:00, [values]<br>14:31:00, [values]<br>...<br>14:59:00, [values] | **01:33:00, [values]**<br>**02:33:00, [values]**<br>...<br>**21:33:00, [values]** |
| min_time:  14:00:00<br>max_time: 14:29:00 | min_time:  14:30:00<br>max_time: 14:59:00 | min_time: 01:33:00<br>max_time: 21:33:00 |

Query(14:15:00, 14:45:00) : read about 3 SSTables
(2 in-time SSTables and 1 late SSTables)

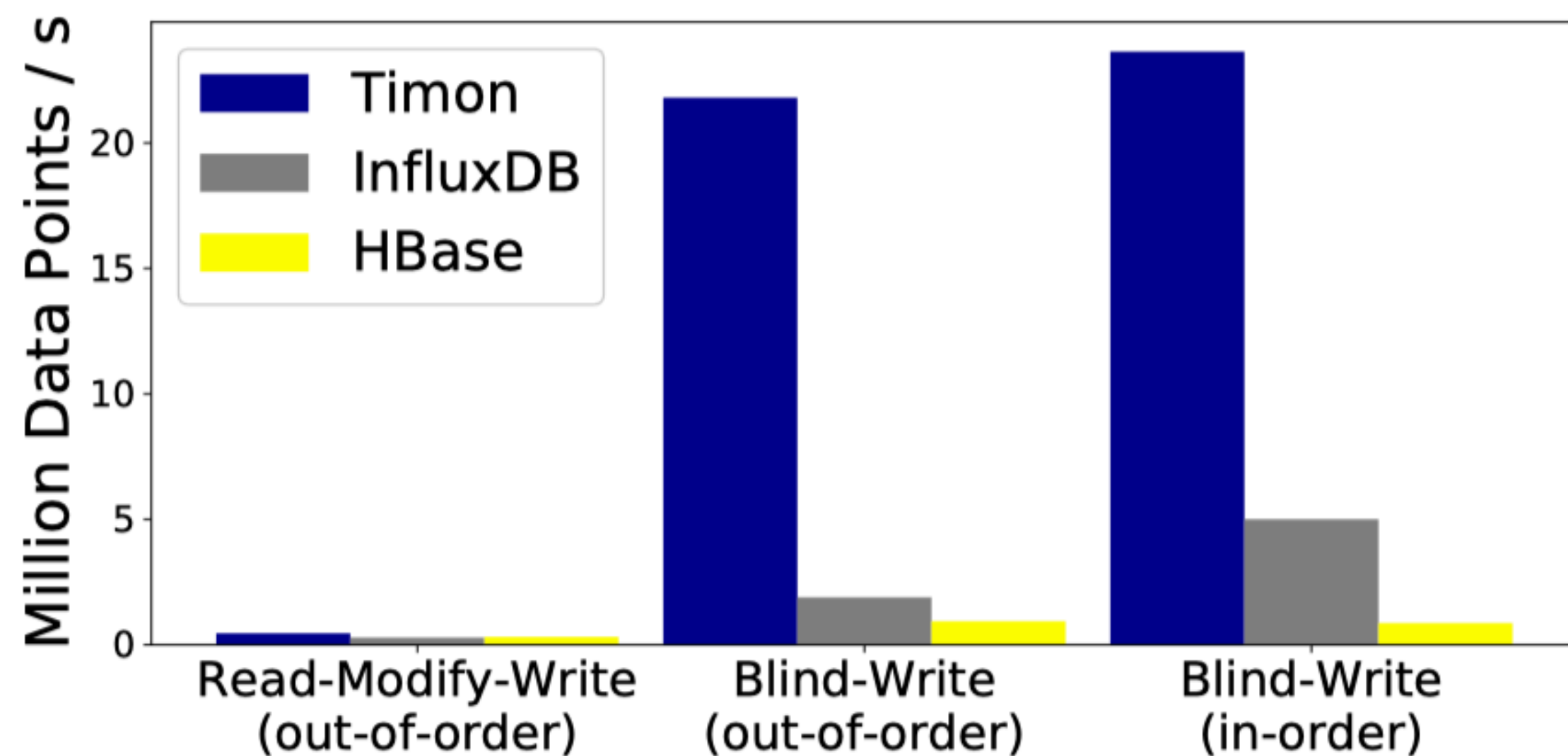# How to process aggregated queries efficiently

- Build time-partitioning tree index by compaction for fast exploration of long-term time-series.

- An aggregated query will scan recent data from MemTable and Lo SSTable, and historical data from time-partitioning tree.
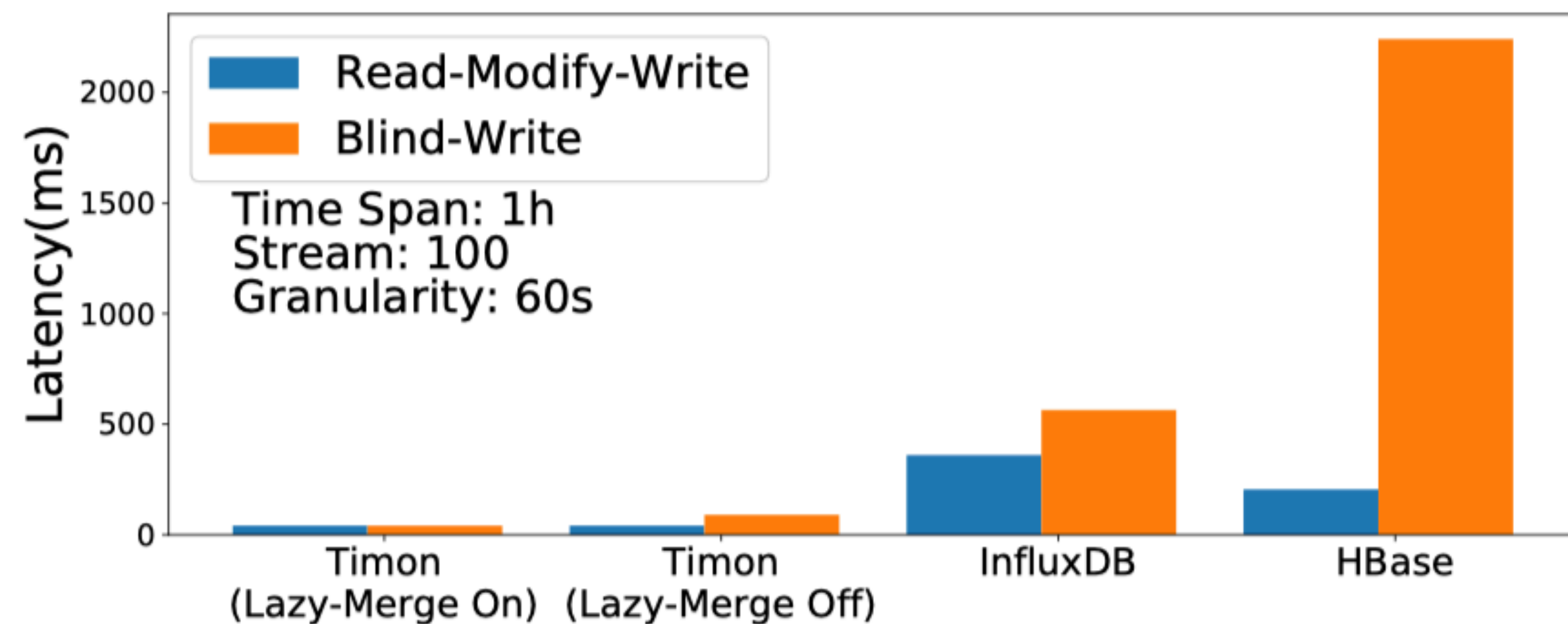
# Timon with user-friendly tools and facilities

- We enhance Timon with user-friendly tools and facilities, such as metric set, materialized view and TQL.

  - Metric set

    - For record which contains dozens of metric values.

  - Materialized view

    - Aggregating data on a higher abstract level, e.g., the region level.

  - TQL

    - SQL-like query languages which allows users to retrieve and analyze the underlying timestamped event data with rich semantics.
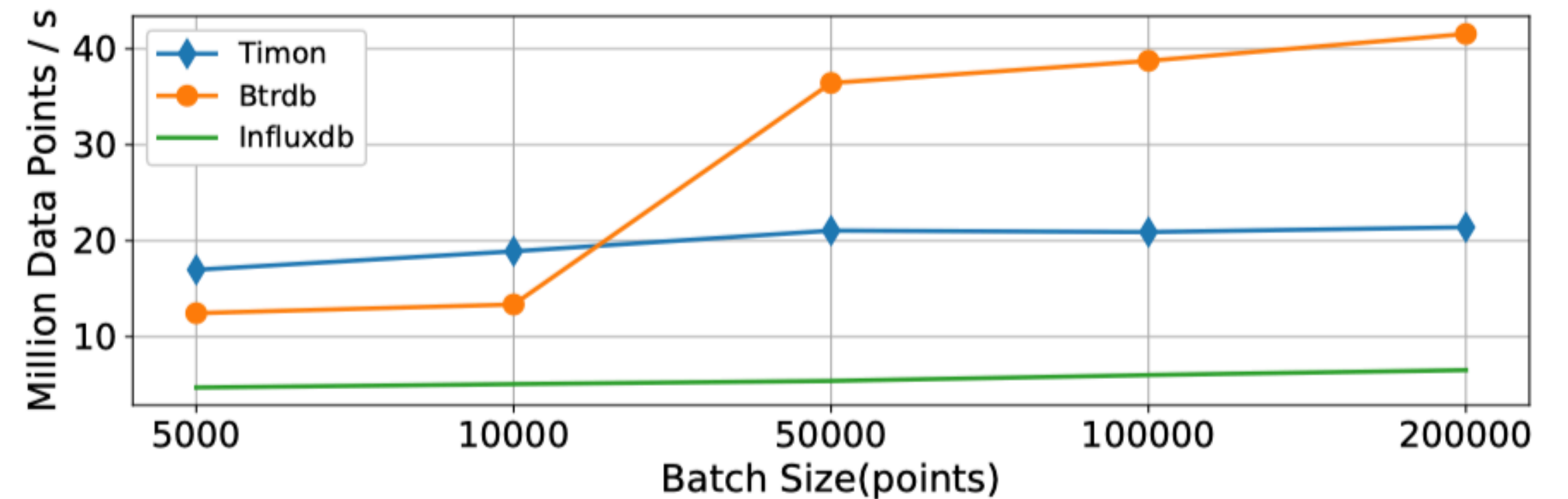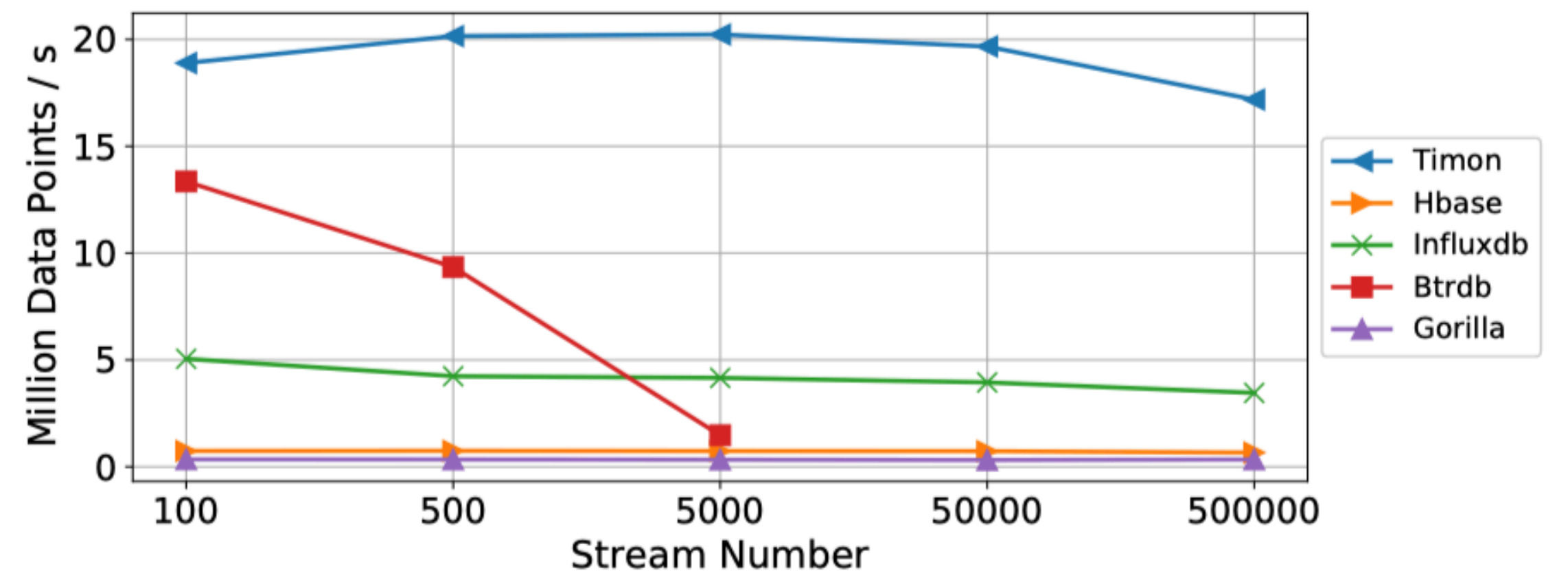
# Benchmark



Throughput with different write mode



Query latency with different write mode.

# Benchmark: Write Throughput as batch size grow or streams grow

- BtrDB (M. P. Andersen at FAST 2016) is a state-of-the-art TSDB with very high performance and long-term time-series exploration support.

- BtrDB is much better for ultra-high frequency data points (i.e., sub-microsecond precision timestamps).

- Timon maintains high write performance stably when the number of streams increases.

- The main difference is that Timon builds Time-partitioning Tree Indexes by a batch and async procedure, but BtrDB inserts records directly into its tree indexes.



(a) Write Throughput as batch size grow



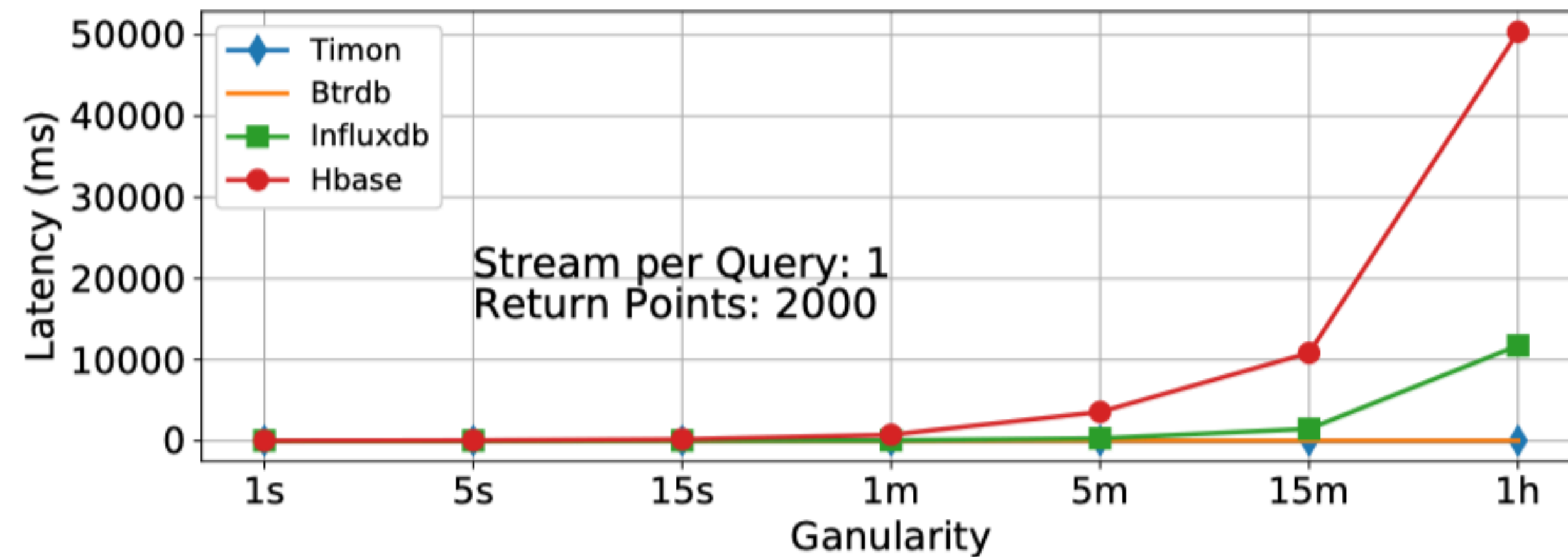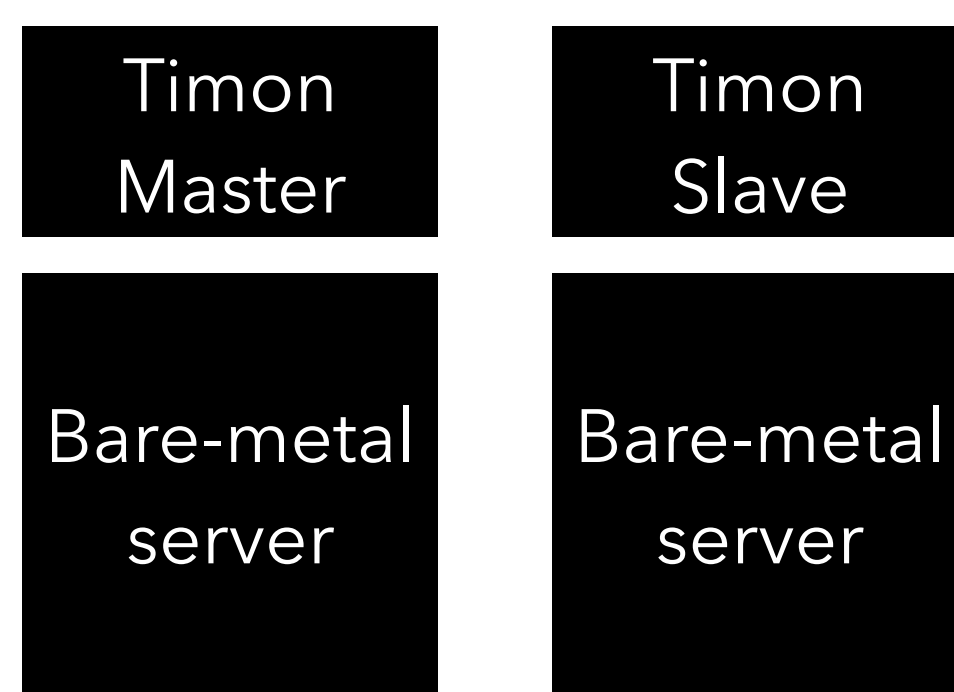(b) Write Throughput as streams grow
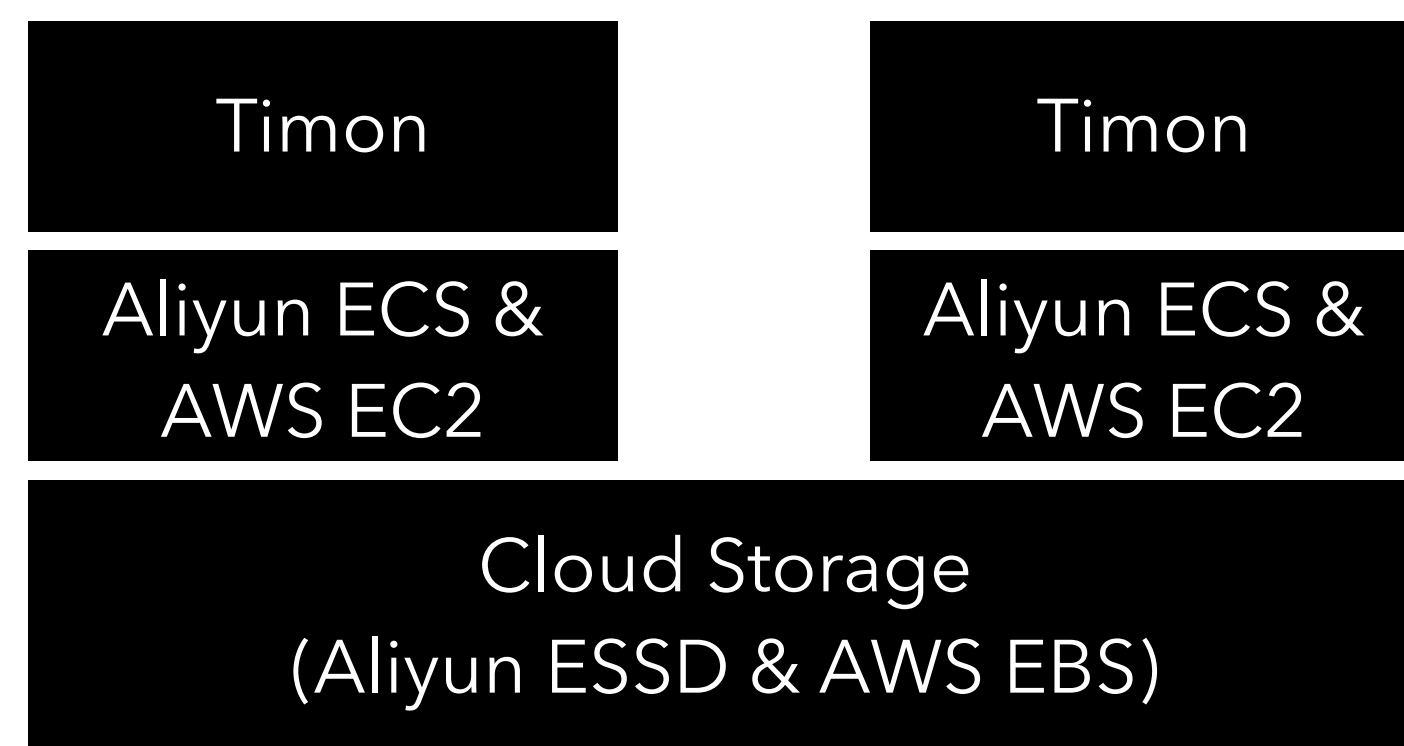
# Benchmark: Query Latency



Figure 17: Query latencies between databases as time granularity changes.

# Deployment

- Deployed in data centers distributed in 21 regions around the world.

- The biggest application cluster: 97 Timon nodes support about 500 million data points writing per second from the system, and the busiest node serves about 18 million data points per second.
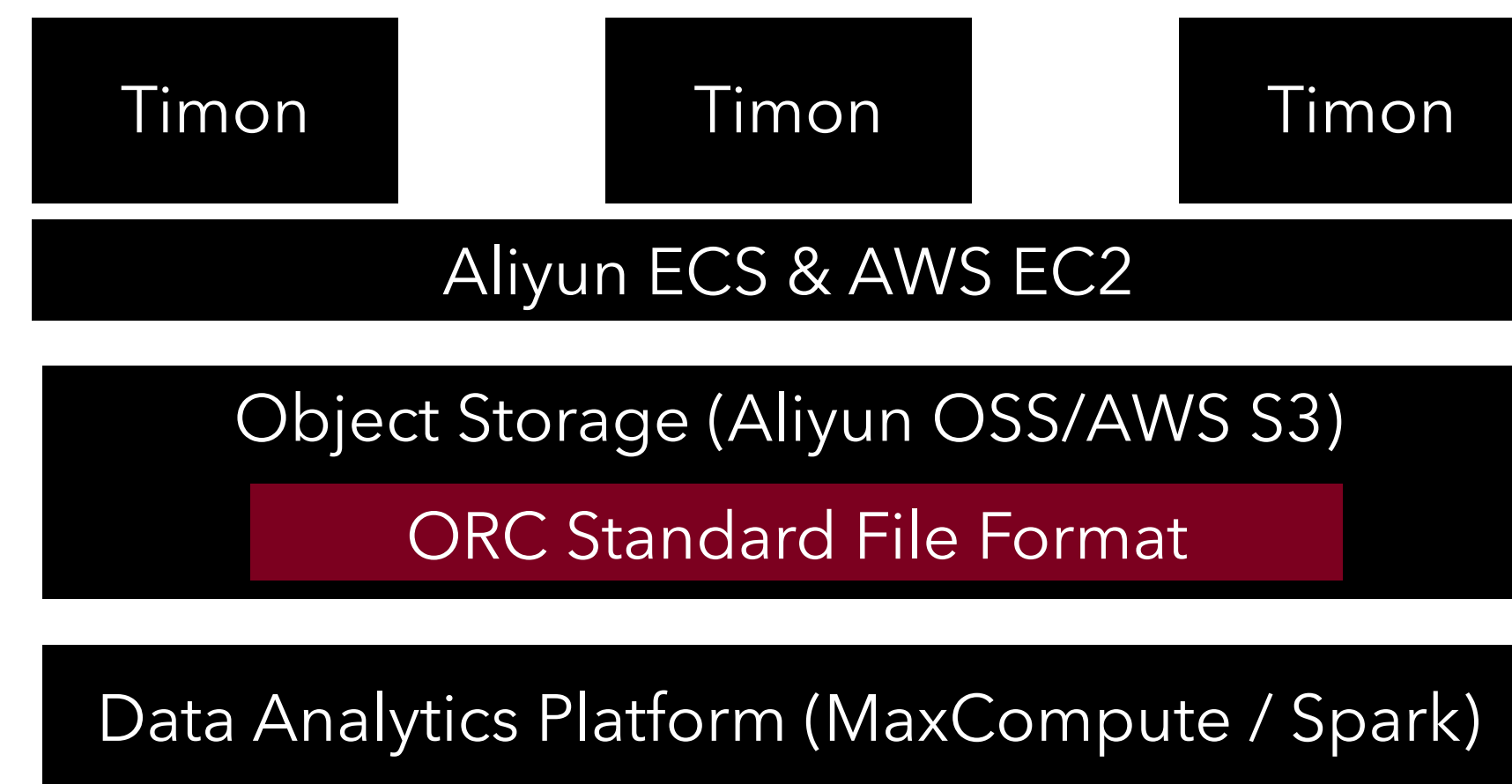
| Timon Master | Timon Slave |
|---|---|
| Bare-metal server | Bare-metal server |

**On Bare-metal Server**

| Timon | Timon |
|---|---|
| Aliyun ECS & AWS EC2 | Aliyun ECS & AWS EC2 |

Cloud Storage
(Aliyun ESSD & AWS EBS)

**On Elastic Compute**

(Each Timon node sees its own data)

| Timon | Timon | Timon |
|---|---|---|

Aliyun ECS & AWS EC2

Object Storage (Aliyun OSS/AWS S3)

ORC Standard File Format

Data Analytics Platform (MaxCompute / Spark)

**On Elastic Compute**

**And Using Object Storage**

# Thanks

Yusong Gao

jianchuan.gys@alibaba-inc.com