# NSF/IUCRC CAC PROJECT

## Experience of Storing and Querying Monitoring Data of Large-scale High Performance Computing Platforms

Jie Li

Doctoral Student, TTU

04/15/2020

Advisors:

Mr. Jon Hass, SW Architect, Dell Inc.

Dr. Alan Sill, Managing Director, HPCC, TTU

Dr. Yong Chen, Associate Professor, CS Dept, TTU

- Time Series, Time Series DBs, InfluxDB

- Problems & Challenges

- Hands-on Experience & Efforts

- Summary & Demo

**Definition:**

Time Series is an <span style="color:red">ordered</span> sequence of values of a variable (e.g. temperature) at <span style="color:red">equally spaced time intervals</span> (e.g. hourly)

**Uses:**

- Time Series Analysis: explore how a given variable changes over time

- Regression Analysis: examine how the changes associated with a specific variable can cause shifts in other variables over the same time period

- Time Series Forecasting: use information regarding historical values and associated patterns to predict future activity

Naqvi, Syeda Noor Zehra, Sofia Yfantidou, and Esteban Zimányi. "Time series databases and influxdb." *Studienarbeit, Université Libre de Bruxelles* (2017).

**Definition:**

A Time Series Database (TSDB) is a database type which is <span style="color:red">optimized</span> for time series or time-stamped data.
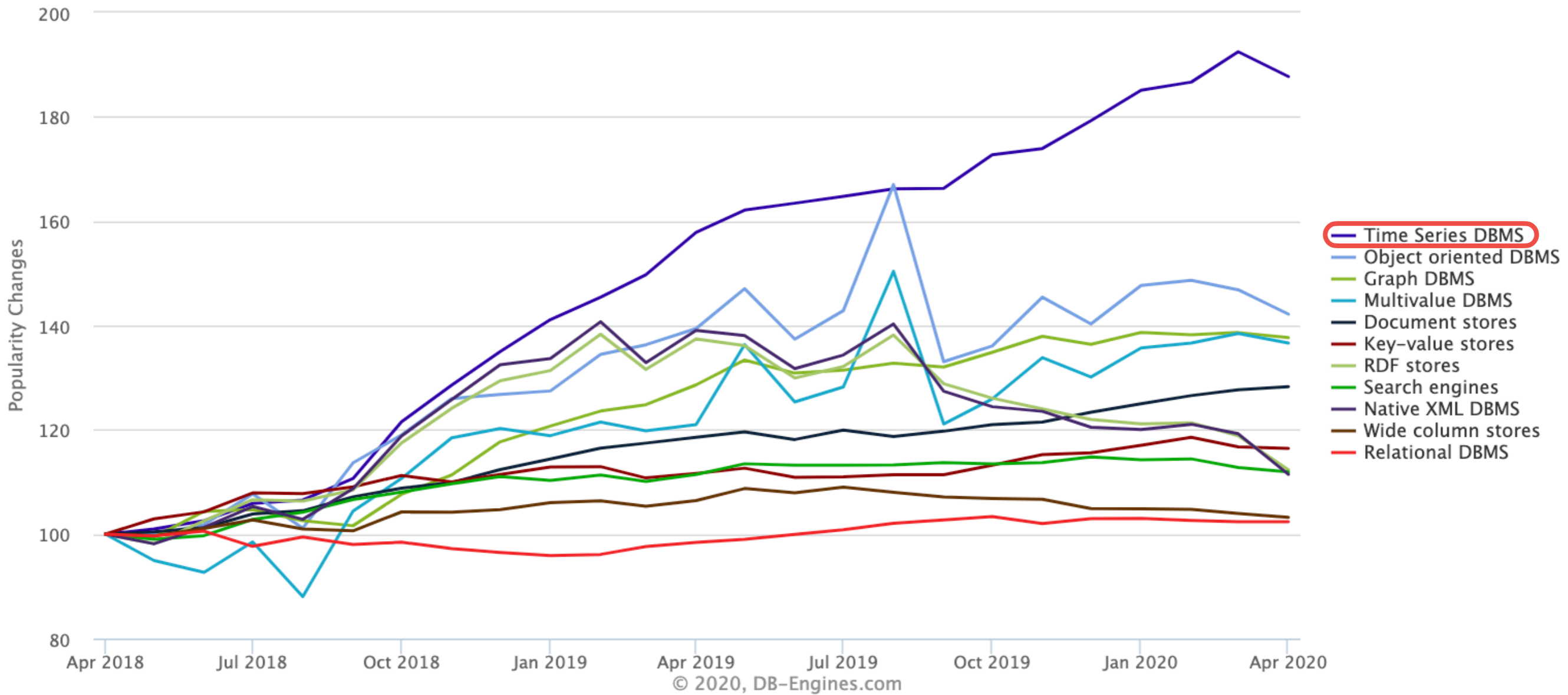
**Properties:**

- Data Location: co-locate chucks of data within the same time range on the same physical part of the database cluster

- Fast, Easy Range Queries: keep the co-related data together to ensure that the range queries are fast

- High Write Performance: ensure high availability and high performance for both read and write operations during peak loads

Naqvi, Syeda Noor Zehra, Sofia Yfantidou, and Esteban Zimányi. "Time series databases and influxdb." *Studienarbeit, Université Libre de Bruxelles* (2017).

**Properties(cont.):**
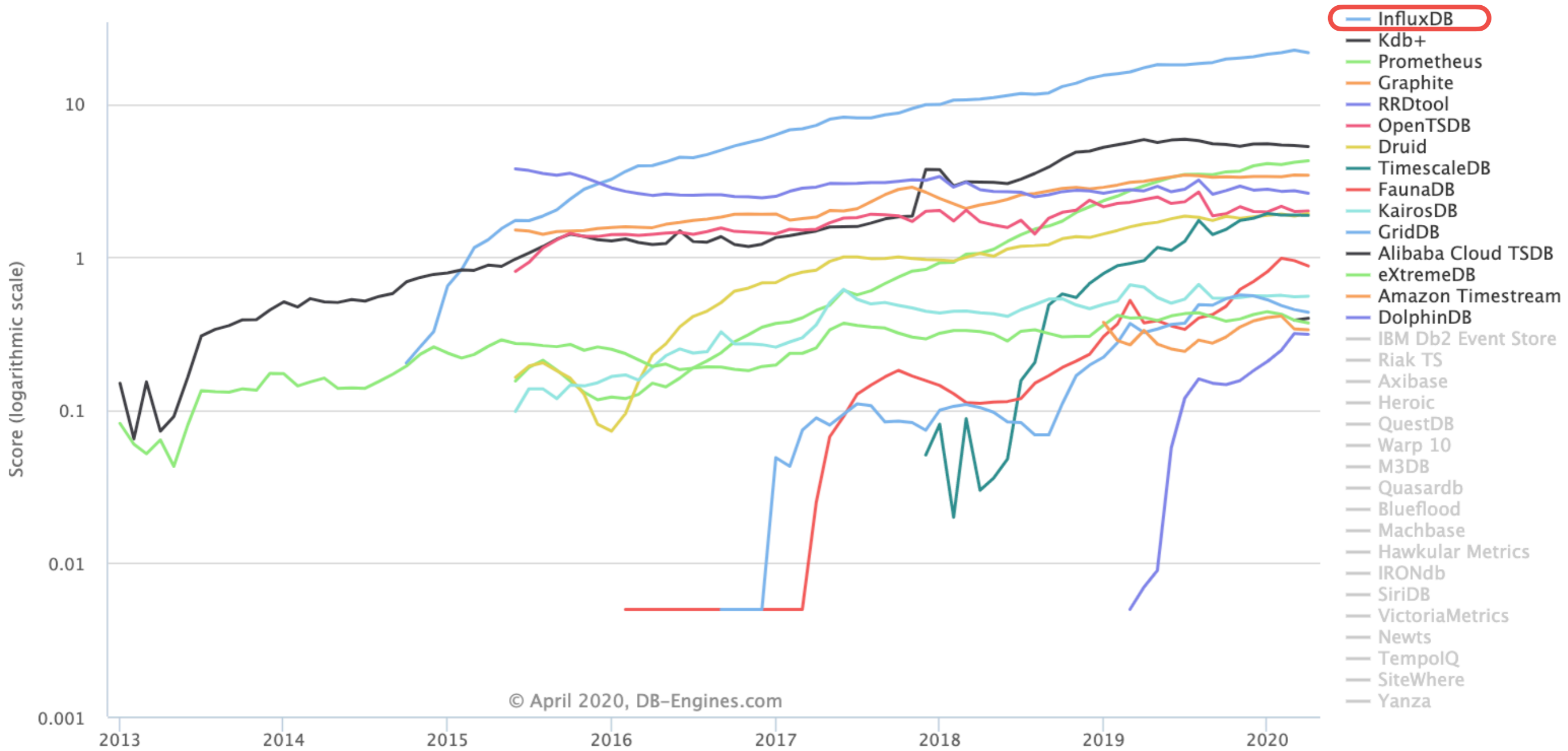
- Data Compression: provide functionality to perform roll-ups in such scenarios for data compaction

- Scalability: take care of scale by introducing functionalities that are only possible when treat time as first concern

- Usability: include functions and operations that are common to time series data analysis
  - data retention policies
  - continuous queries
  - flexible time aggregations
  - range queries etc.

Naqvi, Syeda Noor Zehra, Sofia Yfantidou, and Esteban Zimányi. "Time series databases and influxdb." *Studienarbeit, Université Libre de Bruxelles* (2017).

# TIME SERIES DBS



DBMS (Database Management System) Popularity broken down by database model

# TIME SERIES DBS



Score (logarithmic scale)

10

1

0.1

0.01

0.001

2013    2014    2015    2016    2017    2018    2019    2020

© April 2020, DB-Engines.com

Ranking of Time Series DBMS

Legend:
- InfluxDB
- Kdb+
- Prometheus
- Graphite
- RRDtool
- OpenTSDB
- Druid
- TimescaleDB
- FaunaDB
- KairosDB
- GridDB
- Alibaba Cloud TSDB
- eXtremeDB
- Amazon Timestream
- DolphinDB
- IBM Db2 Event Store
- Riak TS
- Axibase
- Heroic
- QuestDB
- Warp 10
- M3DB
- Quasardb
- Blueflood
- Machbase
- Hawkular Metrics
- IRONdb
- SiriDB
- VictoriaMetrics
- Newts
- TempoIQ
- SiteWhere
- Yanza

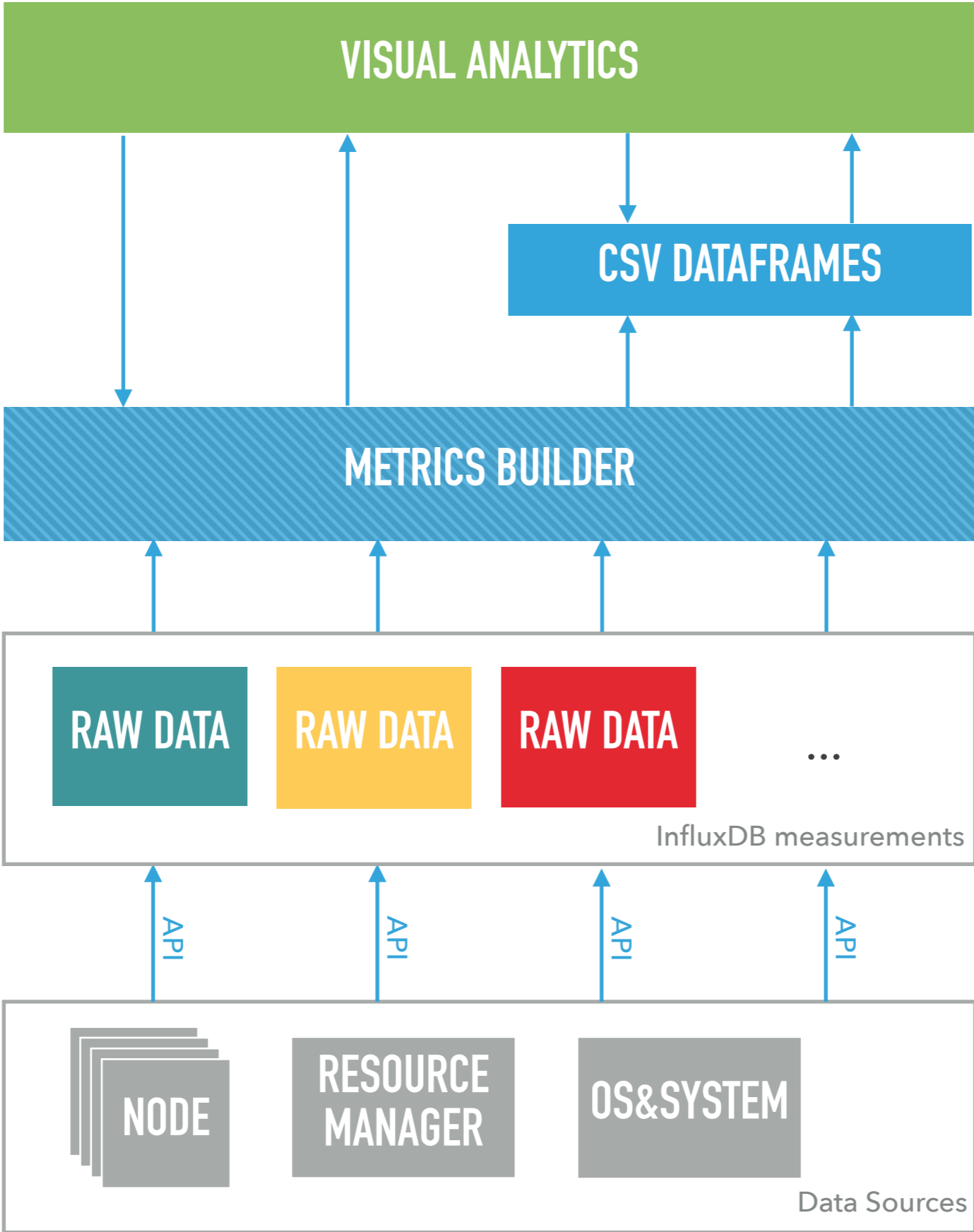https://db-engines.com/en/ranking_trend/time+series+dbms

INFLUXDB

**InfluxDB:**

- Open-source schemaless time series database

- Written in Go and optimized for fast, high-availability storage and retrieval of time series data

- Provides an SQL-like query language

**Data model:**
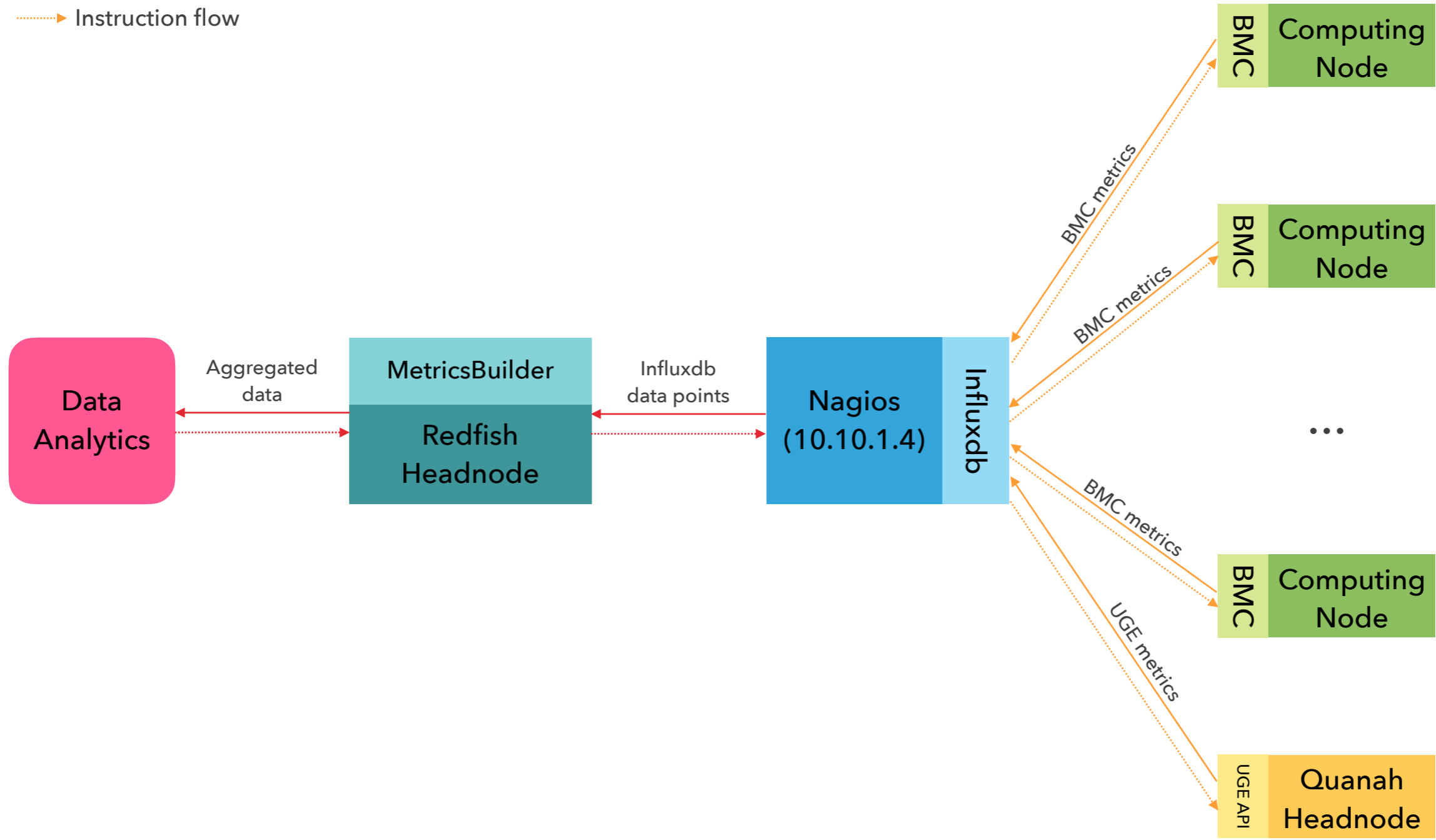
- measurement
- timestamp
- fields
- tags

```
<measurement>[,<tag_key>=<tag_value>[,<tag_key>=<tag_value>]]
<field_key>=<field_value>[,<field_key>=<field_value>] [<timestamp>]
```

# MONITORING FRAMEWORK

# MONITORING FRAMEWORK

data flow

Instruction flow

**Data Analytics**

Aggregated data

**MetricsBuilder**

**Redfish Headnode**

Influxdb data points

**Nagios (10.10.1.4)**

**Influxdb**

BMC metrics

BMC metrics

BMC metrics

UGE metrics

**BMC** | **Computing Node**

**BMC** | **Computing Node**

**BMC** | **Computing Node**

...

**UGE API** | **Quanah Headnode**

# METRICS BUILDER WORKFLOW

**Receive requests from analytics client (HiperViz)**

<span style="color:red">time range:</span>
e.g. 2019-04-20T00:00:00Z, 2019-04-21T00:00:00Z
<span style="color:red">time interval:</span> e.g. 30m
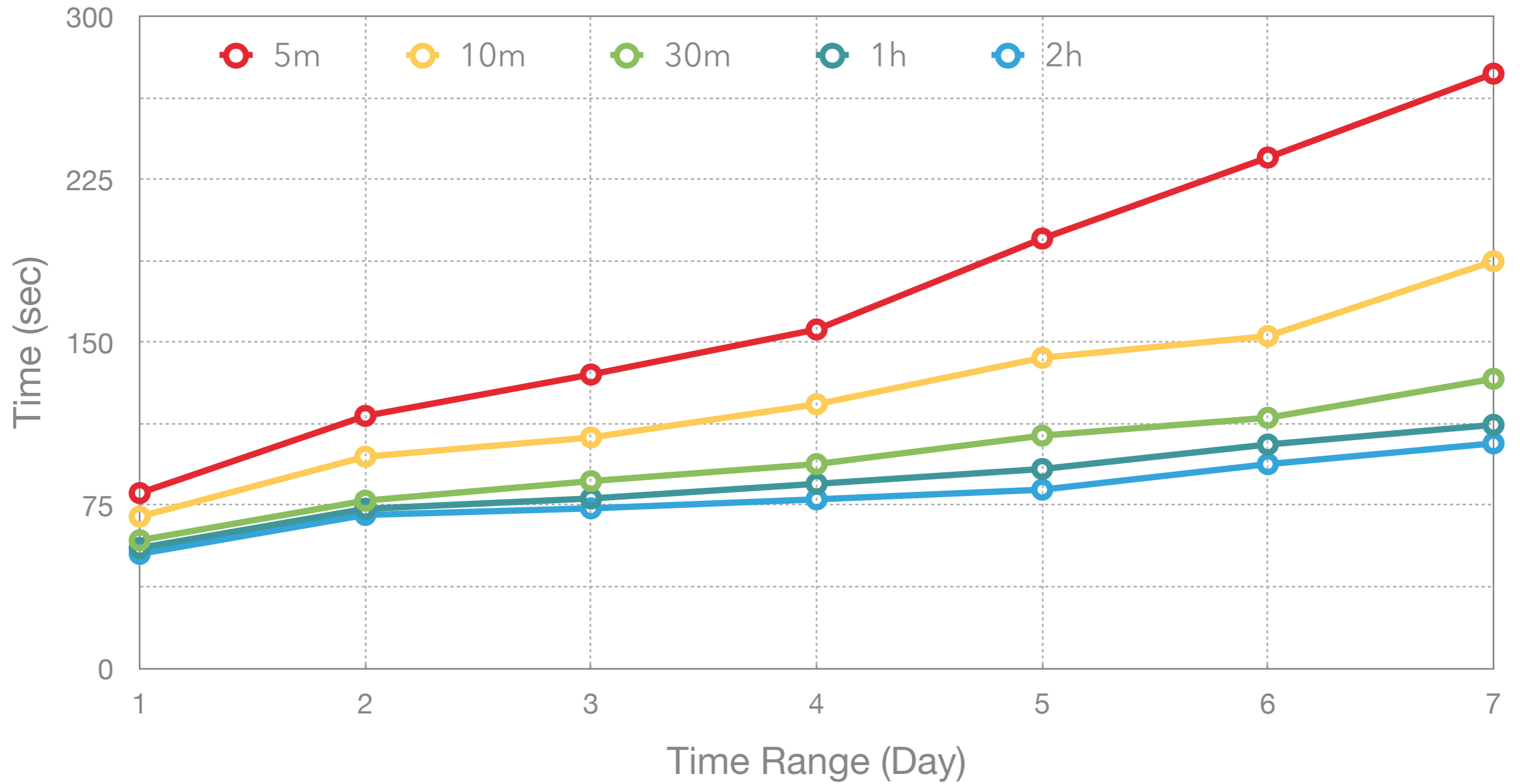<span style="color:red">data type:</span> e.g. Max, Min, Average

**Generate corresponding influxDB query requests**

SELECT max(CPU_Usage) FROM CPU_Usage WHERE host='10.101.1.1' AND time >= '2020-04-10T00:00:00Z' AND time <= '2020-04-11T00:00:00Z' GROUP BY time(30m)

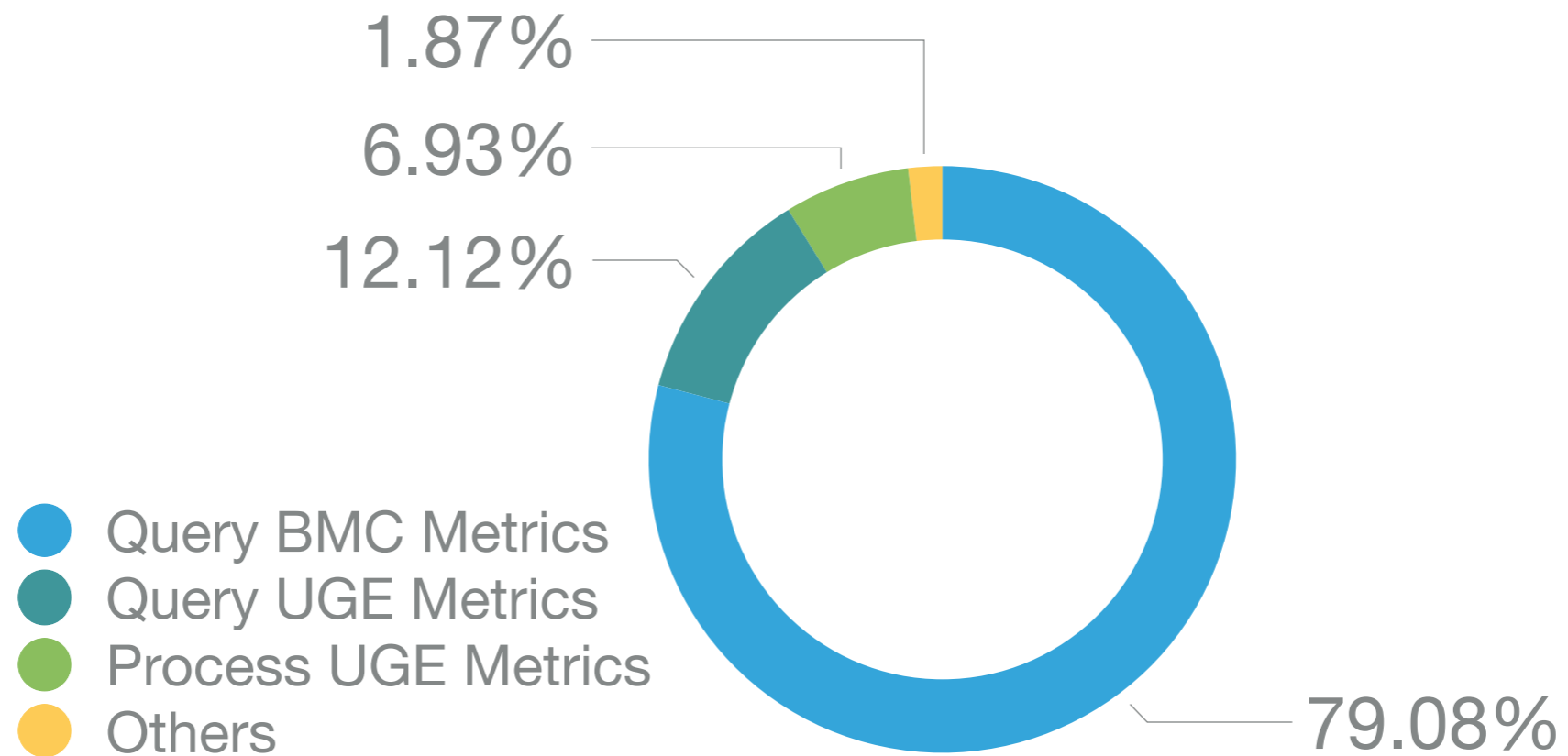**Process data returned from influxDB;
Convert it to csv dataframe;
Return dataframe to analytics client**

- Move data processing from front end to back end
- Provide a uniform API to analytics client
- Act as a middleware that deals with different database design

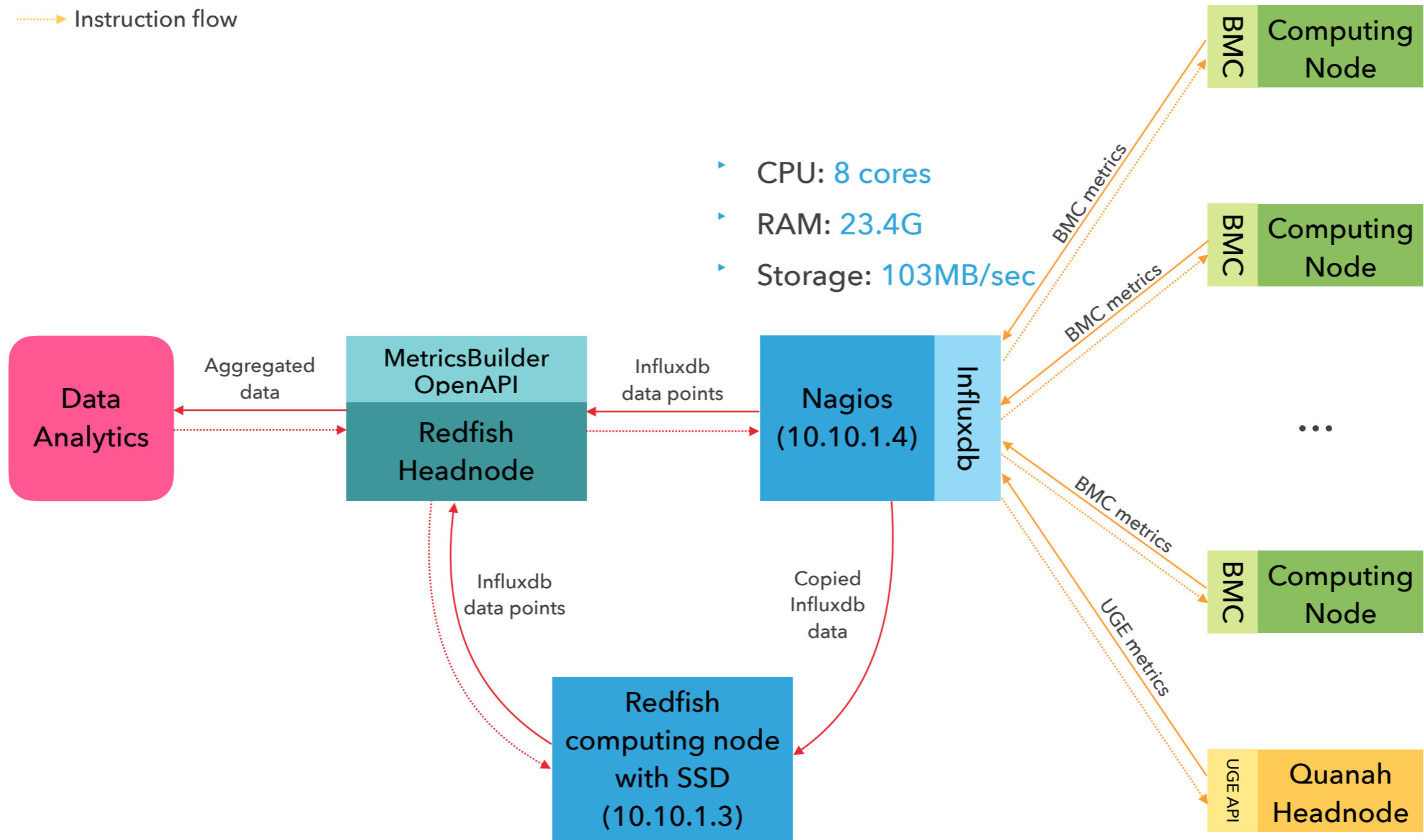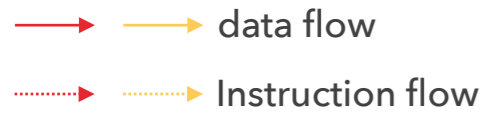# PERFORMANCE



Query and processing time of MetricsBuilder

# ANALYSIS

1.87%

6.93%

12.12%

- ● Query BMC Metrics
- ● Query UGE Metrics
- ● Process UGE Metrics
- ● Others

79.08%

Time Occupation of Query and Processing Metrics

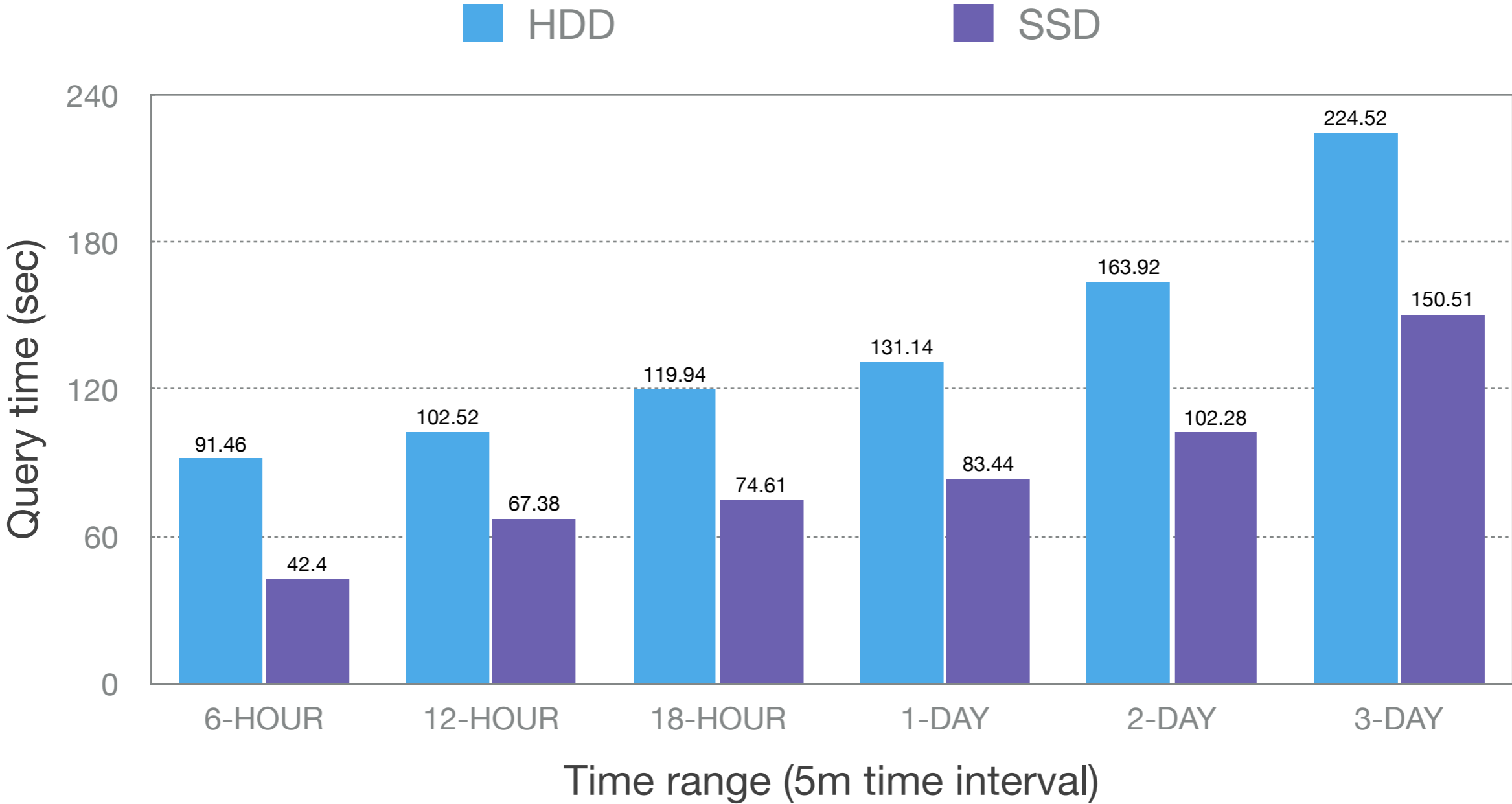**Queuing database occupies ~90% of total running time.**

# Switch from HDD to **SSD**

# EFFORTS - USE SSD

→ → data flow

⋯⋯⋯> ⋯⋯⋯> Instruction flow

- CPU: 8 cores
- RAM: 23.4G
- Storage: 103MB/sec

**Data Analytics**

Aggregated data

**MetricsBuilder OpenAPI**

**Redfish Headnode**

Influxdb data points

**Nagios (10.10.1.4)** **Influxdb**

Influxdb data points

Copied Influxdb data

**Redfish computing node with SSD (10.10.1.3)**

- CPU: 16 cores
- RAM: 94.2G
- Storage: 391MB/sec

BMC | **Computing Node**

BMC | **Computing Node**

BMC metrics

BMC metrics

⋯

BMC metrics

BMC | **Computing Node**

UGE metrics

UGE API | **Quanah Headnode**

# IMPROVEMENT - USE SSD



Query data from SSD performs 1.5x ~ 2.1x faster than HDD

# Redesign Schema

# EFFORTS - REDESIGN SCHEMA

Metrics Saved in different measurements(tables):

- **BMC**:
  - CPU_Temperature
  - Inlet_Temperature
  - Fan_Speed
  - Node_Power_Usage

- **UGE**:
  - Job_Info
  - CPU_Usage
  - Memory_Usage

- **BMC**:
  - cluster_unified_metrics

- **UGE**:
  - Current_Jobs_ID
  - qu_1236124 etc.
  - …

Mar. 14, 2019, 11:44 PM                    Oct. 17, 2019, 9:09 PM

**Previous Schema:**

Measurement: CPU_Temperature
Query: select * from CPU_Temperature WHERE host='10.101.1.1' limit 1

| time | tags | fields |
|---|---|---|
| (epoch time) | cluster<br>host<br>location | CPU1 Temp<br>CPU2 Temp<br>GET_proessing_time<br>Inlet Temp<br>cpuLowerThresholdCritical<br>cpuLowerThresholdNonCritical<br>cpuUpperThresholdCritical<br>cpuUpperThresholdNonCritical<br>error |

# EFFORTS - REDESIGN SCHEMA

**Previous Schema:**

Measurement: qu_1236124
Query: select * from qu_1236124 limit 1

| time | tags | fields |
|---|---|---|
| (epoch time) | cluster location | CPUCores<br>app_name<br>error<br>id<br>nodes_address<br>startTime    "1585883059..04:19 CDT 2020"<br>state<br>submitTime    "1585858122..08:42 CDT 2020"<br>total_nodes<br>user |

# EFFORTS - REDESIGN SCHEMA

```
{
                                    https://10.101.1.1/redfish/v1/Chassis/System.Embedded.1/Thermal
    "@odata.type": "#Thermal.v1_0_2.Thermal",
    "Redundancy": [],
    "Description": "Represents the properties for Temperature and Cooling",
    "Redundancy@odata.count": 0,
    "Fans@odata.count": 4,
    "@odata.id": "/redfish/v1/Chassis/System.Embedded.1/Thermal",
    "@odata.context": "/redfish/v1/$metadata#Thermal.Thermal",
    "Fans": [
        {
            "Status": {
                "State": "Enabled",
                "Health": "OK"
            },
            "UpperThresholdNonCritical": null,
            "MaxReadingRange": 0,
            "Redundancy": [],
            "LowerThresholdCritical": 1050,
            "Redundancy@odata.count": 0,
            "@odata.id": "/redfish/v1/Chassis/System.Embedded.1/Sensors/Fans/0x17%7C%7CFan.Embedded._1",
            "MemberId": "0x17||Fan.Embedded._1",
            "MinReadingRange": 0,
            "UpperThresholdFatal": 17850,
            "ReadingUnits": "RPM",
            "LowerThresholdFatal": 1050,
            "LowerThresholdNonCritical": null,
            "Name": "FAN_1",
            "Reading": 9310,
            "UpperThresholdCritical": 17850,
            "FanName": "FAN_1"
        }
    ],
```

```
"time": 1583792296,

"measurement": "Thermal",

"tags":
    "NodeId": "101.10.1.1"
    "Label": "FAN_1",
"fields":
    "Reading": 9310
```

# EFFORTS - REDESIGN SCHEMA

"time": 1583792296,

"measurement": "UGE",

"tags":
    "NodeId": "101.10.1.1"
    "Label": "CPUUsage",
"fields":
    "Reading": 0.50

"time": 1583792296,

"measurement": "NodeJobs",

"tags":
    "NodeId": "101.10.1.1"
"fields":
    "JobList": ["123456",
                "123457"]

"time": 1583792296,

"measurement": "JobsInfo",

"tags":
    "JobId": "123456"
    "Queue": "quanah"
"fields":
    "StartTime": 1583792200
    "SubmitTime": 1583792200
    "TotalNodes": 1
    "NodeList": ["101.10.1.1"]
    "CPUCores": 10
    "JobName": "test"
    "User": "abc"

Only update when a new job is submitted

# EFFORTS - REDESIGN SCHEMA

Understand measurements

All measurements     : 845,241

  - Numerical measurements : 10

  - Job measurements     : 845,217

  - Other measurements   : 14

(As of Mar. 13, 2020)

Understand sample data points

Convert

All measurements : 5

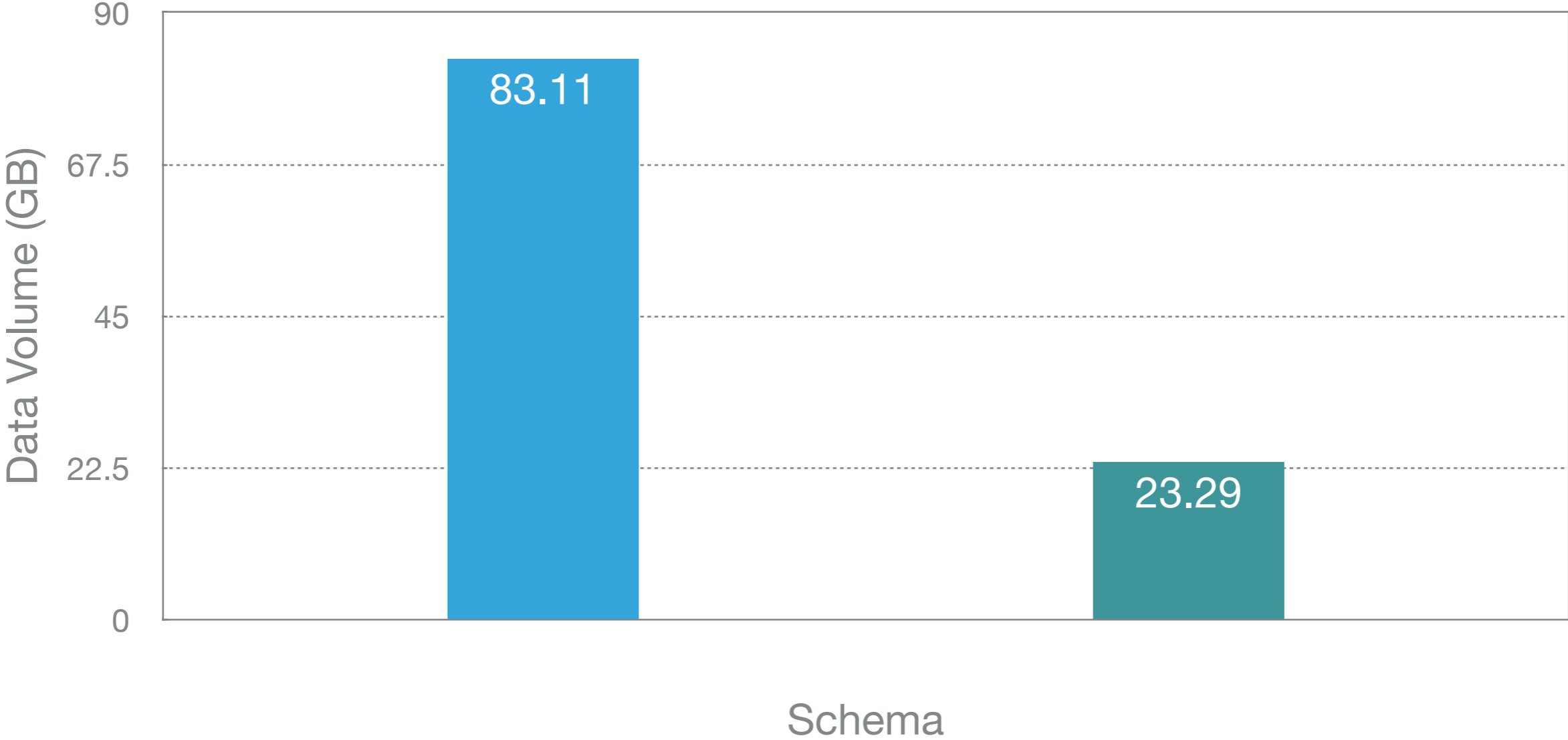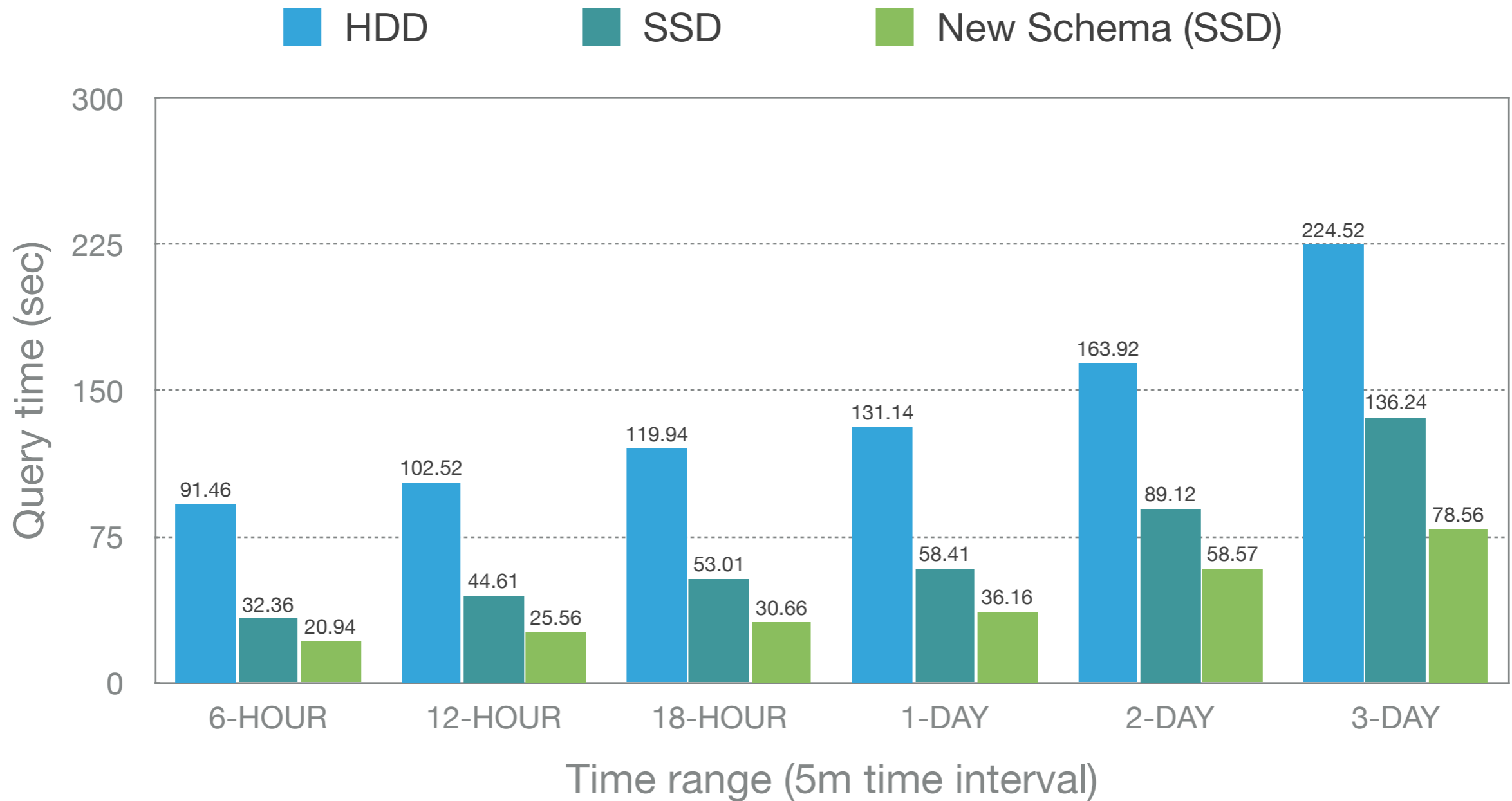| JobsInfo | NodeJobs | Power | Thermal | UGE |
|----------|----------|-------|---------|-----|

# EFFORTS - REDESIGN SCHEMA

March 14, 2019 12:00:00AM – April 10, 2020 12:00:00AM

■ Previsou Schema    ■ Optimized Schema



Data volume in Optimized Schema is 28.02% of the one in previous schema

# IMPROVEMENT - REDESIGN SCHEMA



**HDD**   **SSD**   **New Schema (SSD)**

Query time (sec)

| Time range (5m time interval) | HDD | SSD | New Schema (SSD) |
|---|---|---|---|
| 6-HOUR | 91.46 | 32.36 | 20.94 |
| 12-HOUR | 102.52 | 44.61 | 25.56 |
| 18-HOUR | 119.94 | 53.01 | 30.66 |
| 1-DAY | 131.14 | 58.41 | 36.16 |
| 2-DAY | 163.92 | 89.12 | 58.57 |
| 3-DAY | 224.52 | 136.24 | 78.56 |

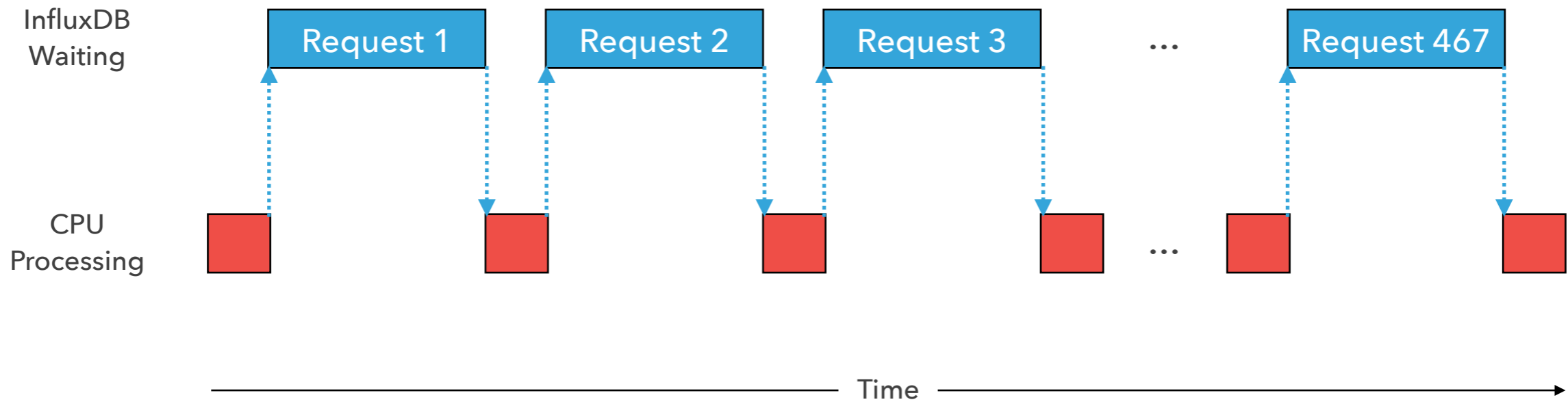Query data from New Schema performs 1.6x ~ 1.76x faster than previous schema

Query data from New Schema performs 2.8x ~ 4.3x faster than previous schema on HDD

**Concurrent** Processing

# EFFORTS - CONCURRENT PROCESSING

Request

SELECT max(Reading) FROM UGE WHERE NodeId='10.101.1.1' AND Label='CPUUsage' AND time >= '2020-04-10T00:00:00Z' AND time <= '2020-04-11T00:00:00Z' GROUP BY time(5m)

InfluxDB Waiting

Request 1     Request 2     Request 3     ...     Request 467

CPU Processing

... ...

Time

Execution Timing Diagram of Previous Implementation

# EFFORTS - CONCURRENT PROCESSING

Request

SELECT max(Reading) FROM UGE WHERE NodeId='10.101.1.1' AND Label='CPUUsage' AND time >= '2020-04-10T00:00:00Z' AND time <= '2020-04-11T00:00:00Z' GROUP BY time(5m)
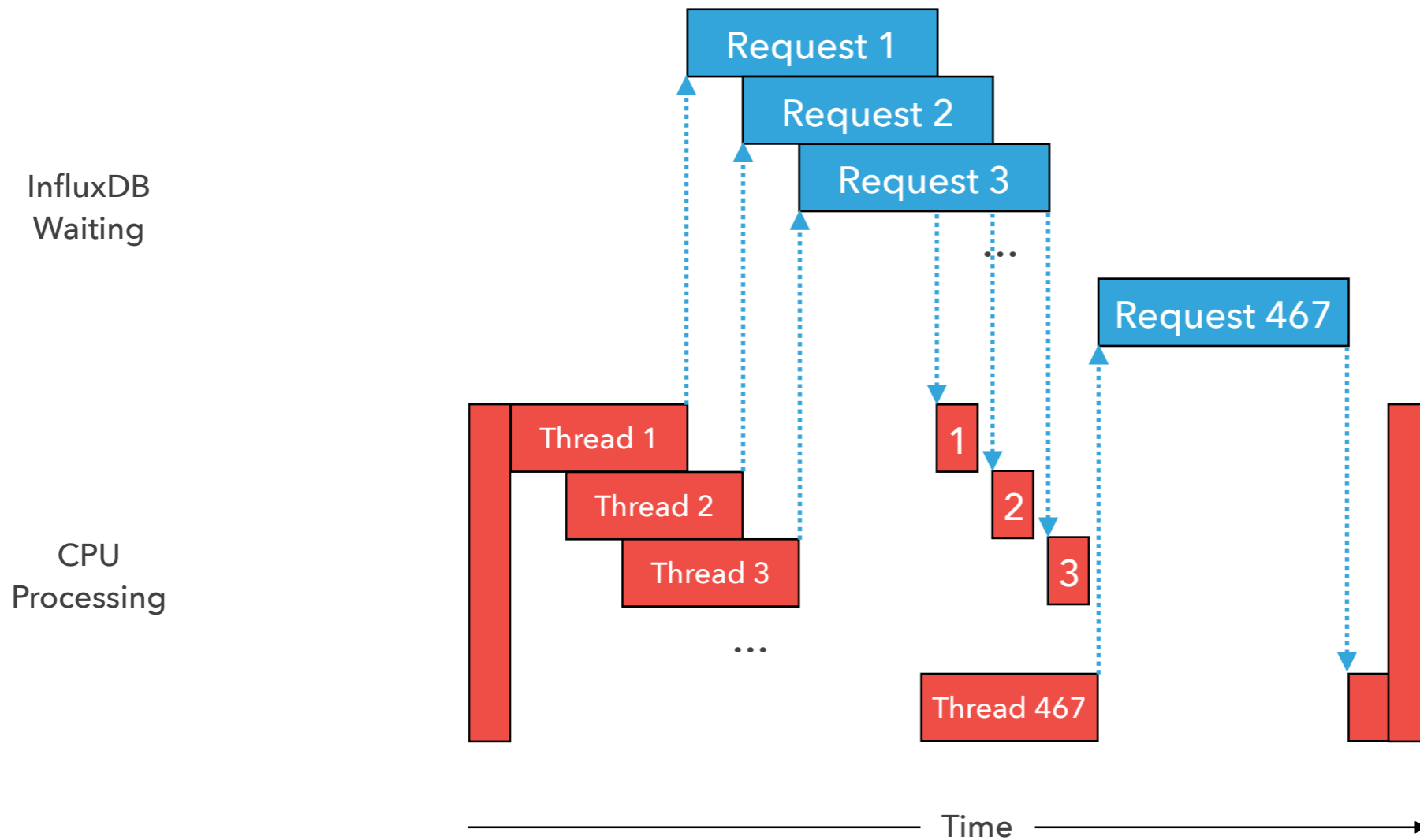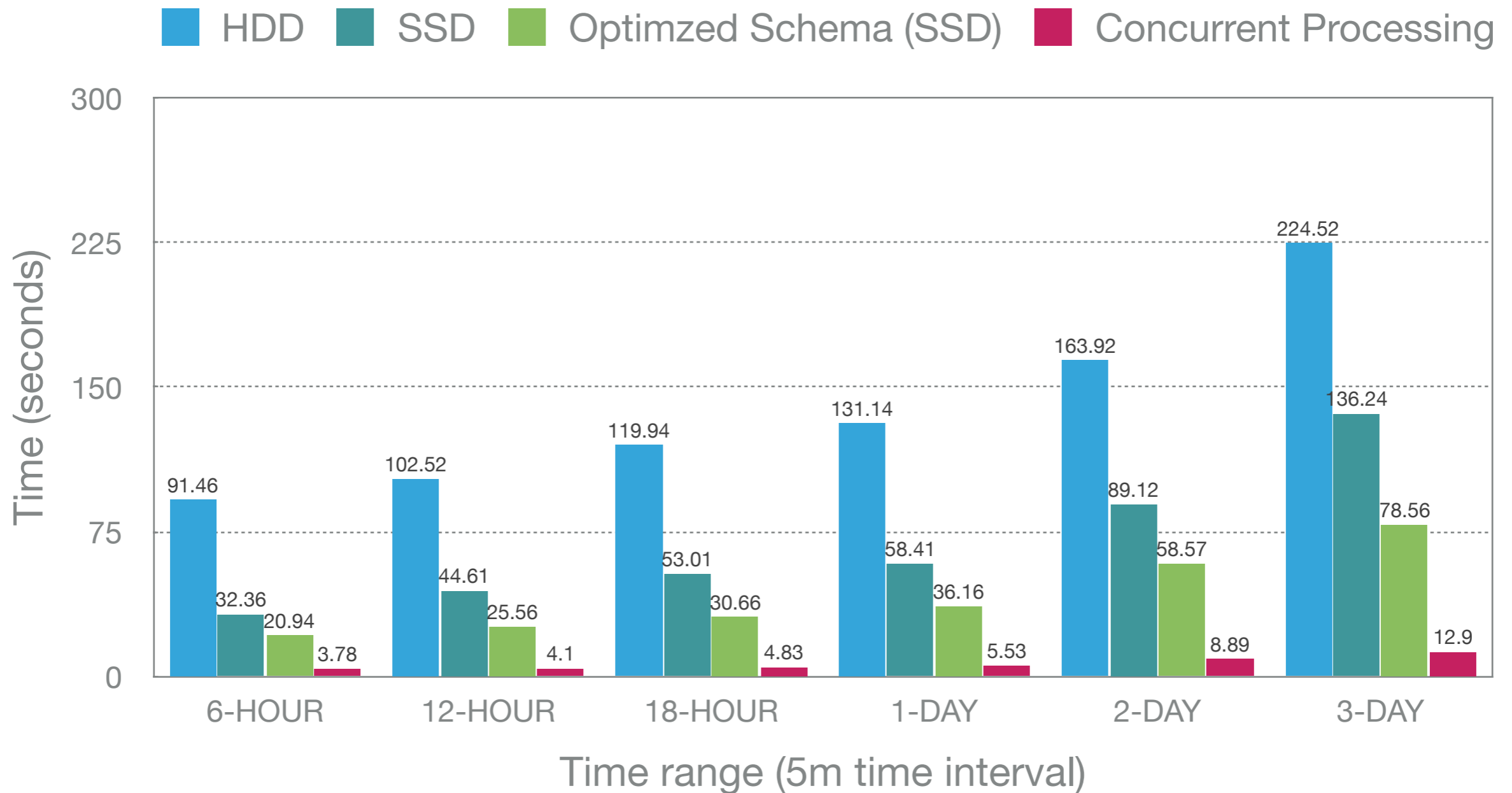
Request 1

Request 2

Request 3

InfluxDB
Waiting

...

Request 467

1

2

Thread 1

3

Thread 2

CPU
Processing

Thread 3

...

Thread 467

Time

**Execution Timing Diagram of Current Implementation**

# IMPROVEMENT - CONCURRENT PROCESSING



**Legend:** HDD, SSD, Optimzed Schema (SSD), Concurrent Processing

| Time range (5m time interval) | HDD | SSD | Optimzed Schema (SSD) | Concurrent Processing |
|---|---|---|---|---|
| 6-HOUR | 91.46 | 32.36 | 20.94 | 3.78 |
| 12-HOUR | 102.52 | 44.61 | 25.56 | 4.1 |
| 18-HOUR | 119.94 | 53.01 | 30.66 | 4.83 |
| 1-DAY | 131.14 | 58.41 | 36.16 | 5.53 |
| 2-DAY | 163.92 | 89.12 | 58.57 | 8.89 |
| 3-DAY | 224.52 | 136.24 | 78.56 | 12.9 |

Time (seconds)
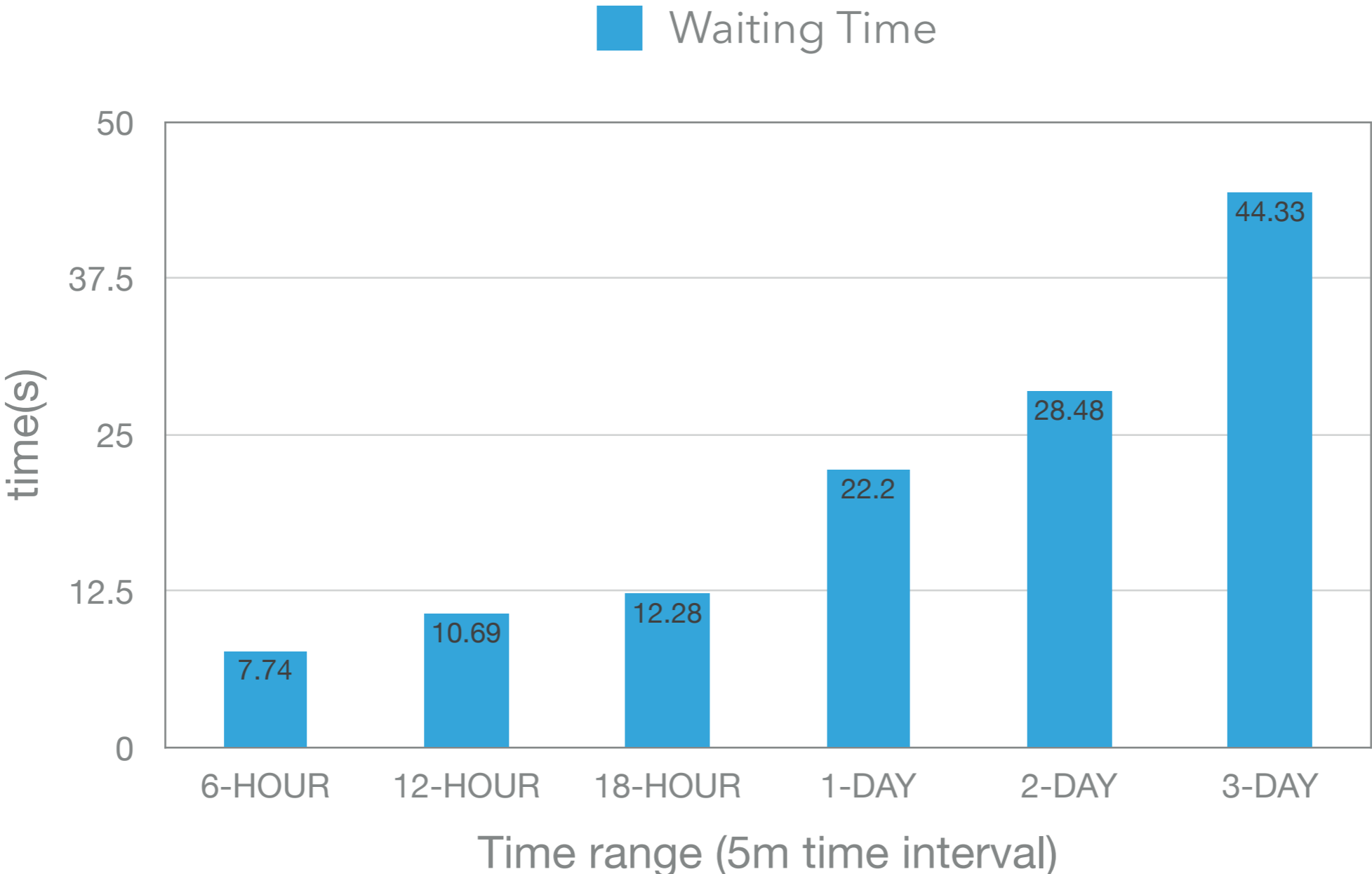
Query data in Concurrent from New Schema performs 5.5x ~ 6.5x faster than in sequence

Query data in Concurrent from New Schema performs 17x ~ 25x faster than previous schema on HDD

# One More Thing…
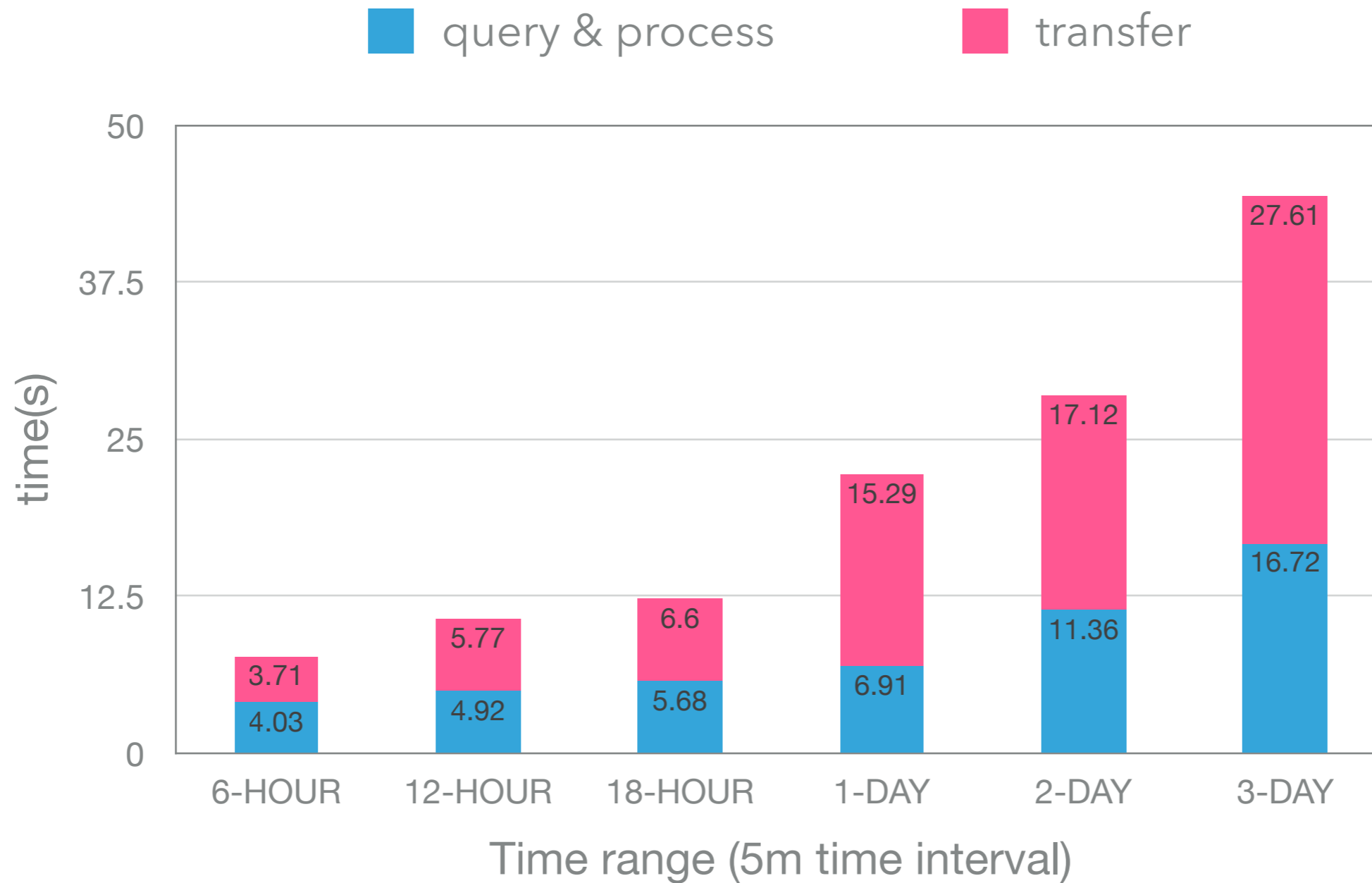
# WAITING TIME



Legend: Waiting Time

| Time range (5m time interval) | Waiting Time |
|---|---|
| 6-HOUR | 7.74 |
| 12-HOUR | 10.69 |
| 18-HOUR | 12.28 |
| 1-DAY | 22.2 |
| 2-DAY | 28.48 |
| 3-DAY | 44.33 |

y-axis: time(s)
x-axis: Time range (5m time interval)
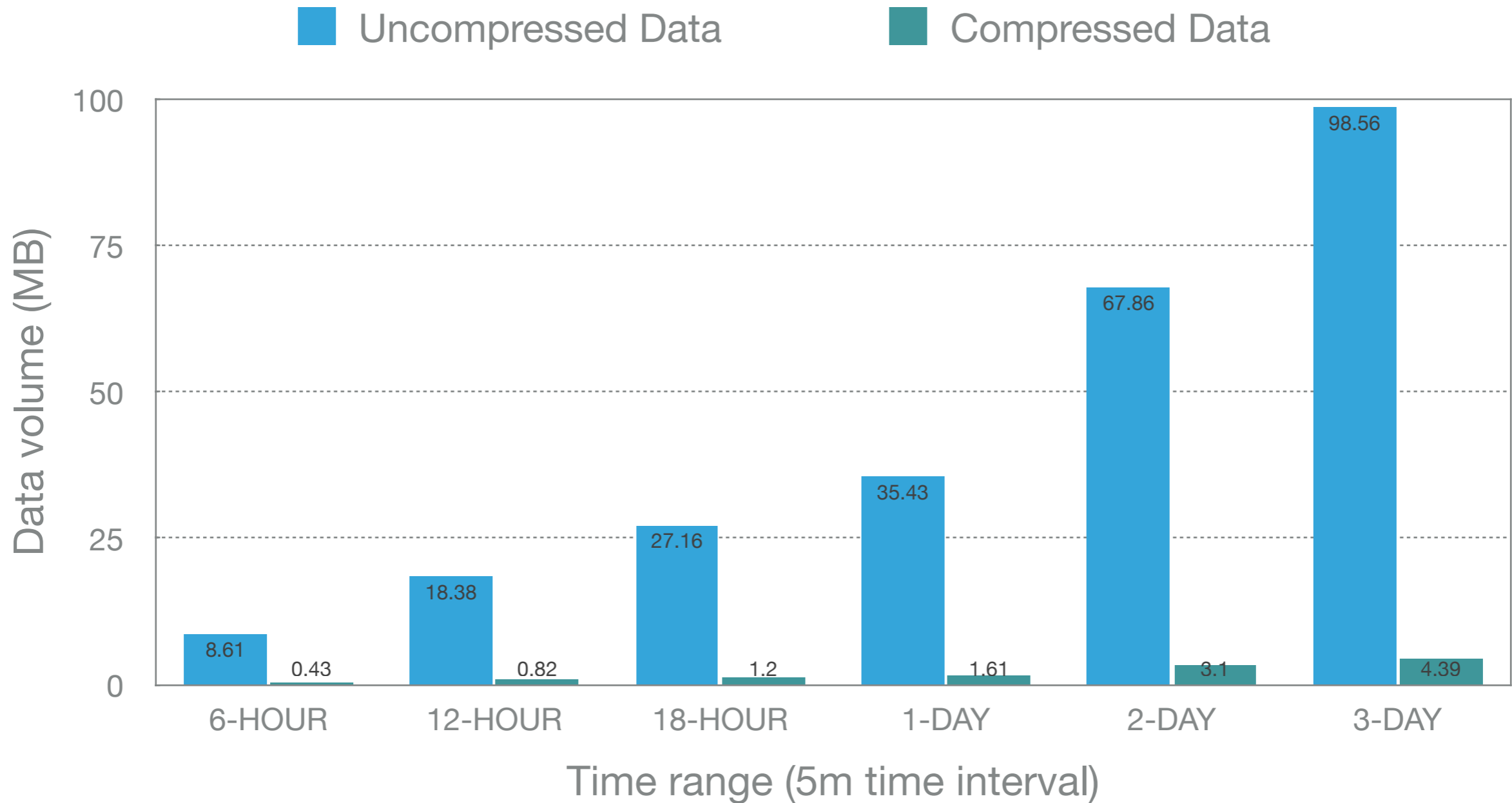
# WAITING TIME DECOMPOSITION



Transfer time is as much as **1.65x** of query & process time
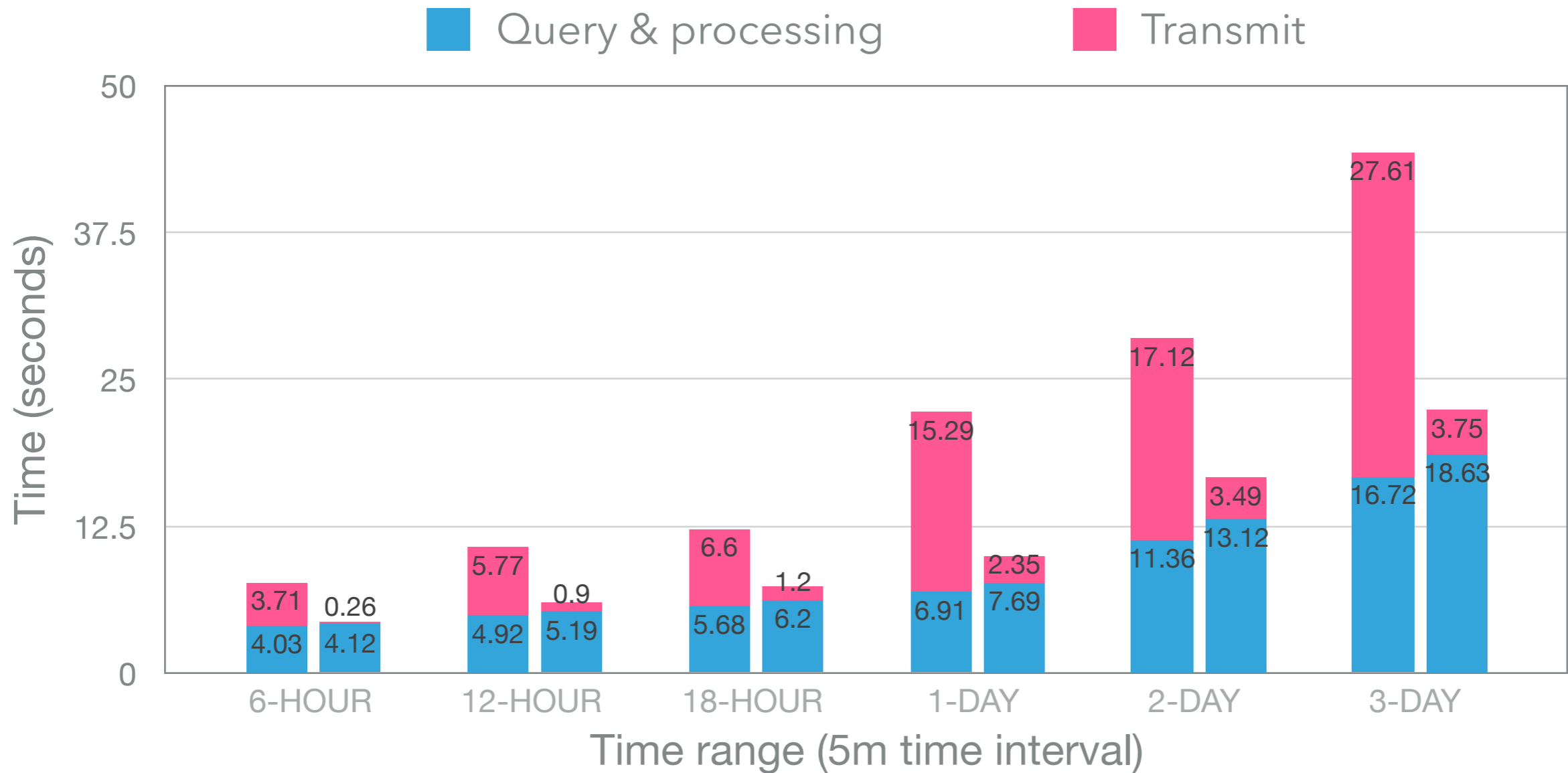
# Transfer **Compressed** Data

# EFFORTS - COMPRESSION

Use **zlib library** for data compression

Legend:
- Uncompressed Data
- Compressed Data

Data volume (MB):

| Time range | Uncompressed Data | Compressed Data |
|---|---|---|
| 6-HOUR | 8.61 | 0.43 |
| 12-HOUR | 18.38 | 0.82 |
| 18-HOUR | 27.16 | 1.2 |
| 1-DAY | 35.43 | 1.61 |
| 2-DAY | 67.86 | 3.1 |
| 3-DAY | 98.56 | 4.39 |

Time range (5m time interval)

Compressed data volume is only **4.45%~5.0%** of uncompressed data

# IMPROVEMENT - COMPRESSION



**Query & processing**   **Transmit**

Time (seconds)

| | |
|---|---|
| 50 | |
| 37.5 | |
| 25 | |
| 12.5 | |
| 0 | |

**6-HOUR**: 3.71 / 4.03, 0.26 / 4.12
**12-HOUR**: 5.77 / 4.92, 0.9 / 5.19
**18-HOUR**: 6.6 / 5.68, 1.2 / 6.2
**1-DAY**: 15.29 / 6.91, 2.35 / 7.69
**2-DAY**: 17.12 / 11.36, 3.49 / 13.12
**3-DAY**: 27.61 / 16.72, 3.75 / 18.63

Time range (5m time interval)

Using compressed data is 1.8x~2.16x faster than using uncompressed data

De-compress 6 hours of data only takes about 0.144 seconds.

1. Switch to **SSD**

2. Redesign **Schema**

3. **Concurrent** Processing

4. Transfer **Compressed** Data

25x Speed Up

2x Speed Up

DEMO