# NSF/IUCRC CAC PROJECT

# MONITORING, VISUALIZING, AND PREDICTING HEALTH STATUS OF HPC CENTERS

Jie Li

PhD Student, TTU

08/30/2019

Advisors:

Mr. Jon Hass, SW Architect, Dell Inc.
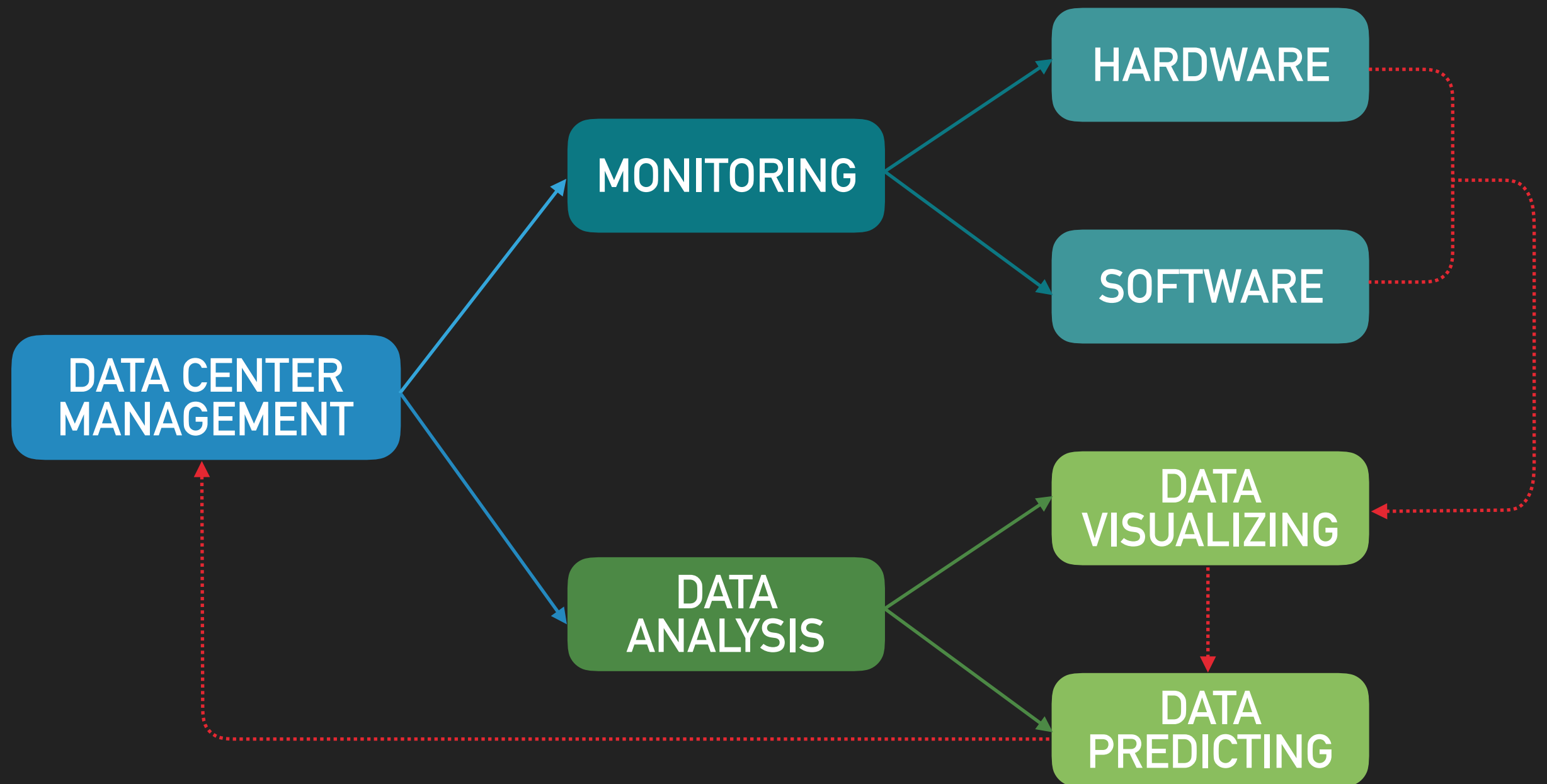
Dr. Alan Sill, Managing Director, HPCC, TTU

Dr. Yong Chen, Associate Professor, CS Dept, TTU
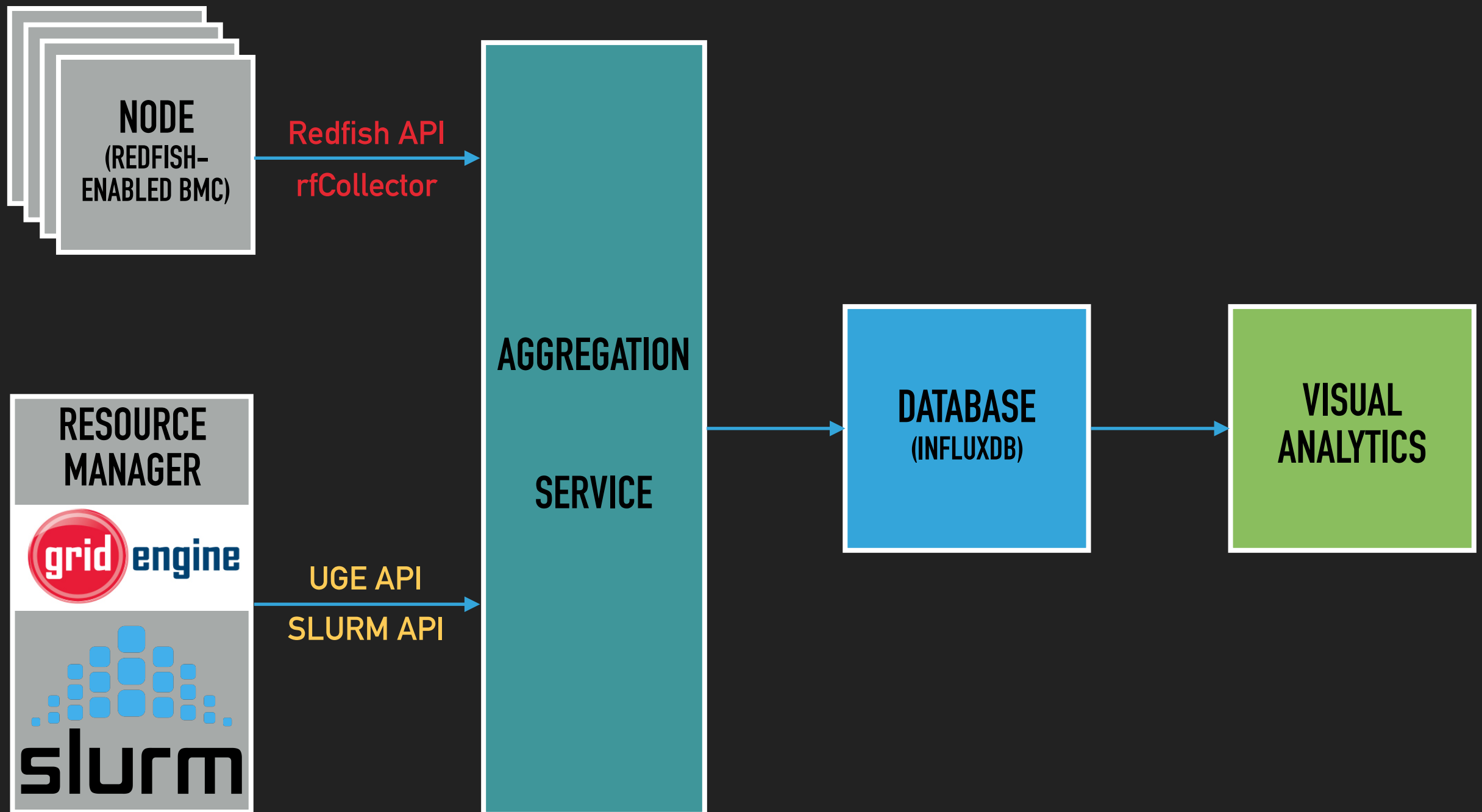
Dr. Tommy Dang, Assistant Professor, CS Dept, TTU

▸ Overview of the project

▸ Interleaving software metrics with hardware metrics

▸ IMPROVE the ways in which systems and applications operate

▸ Measured in terms of performance, reliability, power usage, the ability to meet service level agreements, and similar metrics important to applications and IT infrastructure providers

▸ Operate across multiple time scales, across different size systems, and at multiple level of abstraction(application-centric, infrastructure-centric)

▸ Understand and analyze captured data, in addition to support intelligent problem determination methods, as needed by subsequent management actions. [1]

Ref: [1] Mahendra Kutare et al. Monalytics: Online Monitoring and Analytics for Managing Large Scale Data Centers. ICAC, 2010

# MONITORING FRAMEWORK

**AGGREGATION**

**SERVICE**

▸ Queries and Collects data across the entire cluster

▸ Builds metrics after receiving the monitoring data

▸ Interleaving BMC metrics with Resource Manager metrics

▸ Stores metrics in database(InfluxDB)

**Preparation**

▶ **Builds** the monitoring workload according to a given number of nodes and metrics

**Parallelization**

▶ **Scatters** the monitoring tasks evenly across the available CPU cores as a multi-threaded code and **Gathers** the responses

**Fetch Metric Data**

▶ Each thread **Queries** and **Collects** monitoring data from monitored entities

**Data Processing**

▶ **Interleaves** BMC metrics (hardware data) with Resource Manager metrics (software data), and **Writes** metrics to InfluxDB
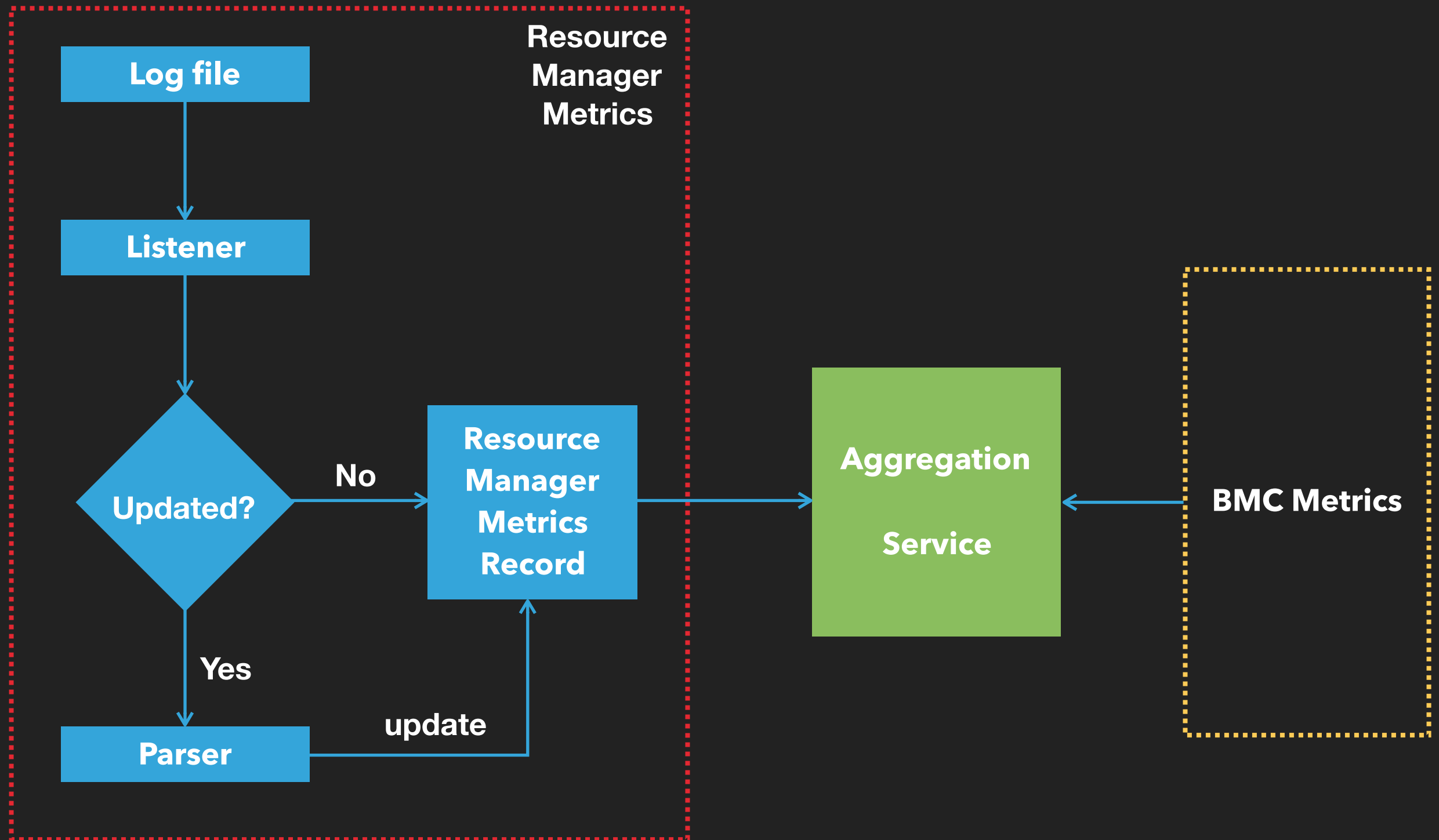
Data Processing

▶ **Interleaves** BMC metrics (hardware data) with Resource Manager metrics (software data), and **Writes** metrics to InfluxDB

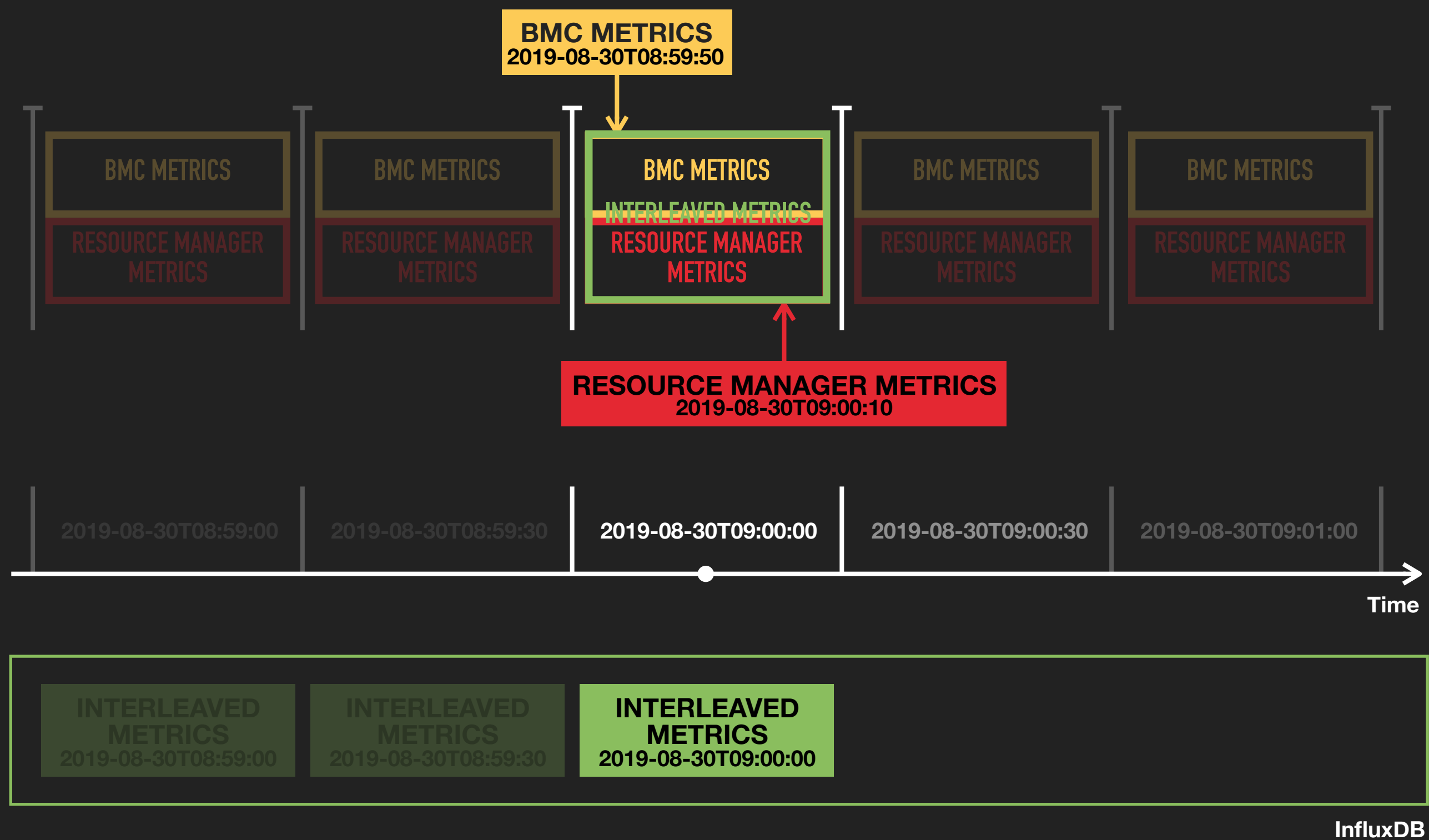At what time interval should we interleave the software data with hardware data?

‣ At the same time interval with one we fetch hardware data ? i.e. about 6 secs
  ‣ Not necessary, the software data may not change as frequently as hardware data
  ‣ May cause extra overhead on the resource manager(UGE/SLURM)
‣ At a relatively large time interval, e.g. 1mins, 5mins
  ‣ Not accurate, some jobs may start and finish during this large time interval
  ‣ The interleaved metrics are affected by how we choose the time interval

‣ Interleaving software data with hardware data as needed
‣ Retrieve data from log file directly instead of from resource manager database
  ‣ Includes resource snapshots of running jobs and other cluster statics
  ‣ All information is time related, whenever there's change, such as a new job submission or a job is finished, these kinds of info are logged into the reporting file

‣ Reporting file (UGE)[1]
‣ Job accounting log file (SLURM)[2]

Ref: [1] http://www.gridengine.eu/mangridengine/htmlman5/reporting.html, [2] https://www.mankier.com/1/sacct

# WORKFLOW

QUESTIONS?/COMMENTS?