

Using Big Data for Anomaly Detection

Javier Cacheiro, CESGA

jlopez@cesga.es



Outline

- Using Big Data to collect Metrics & Logs
- Use case: Anomaly Detection
 - Detecting SSH attacks using log analysis
 - Detecting SLURM job anomalies using metrics

Phases

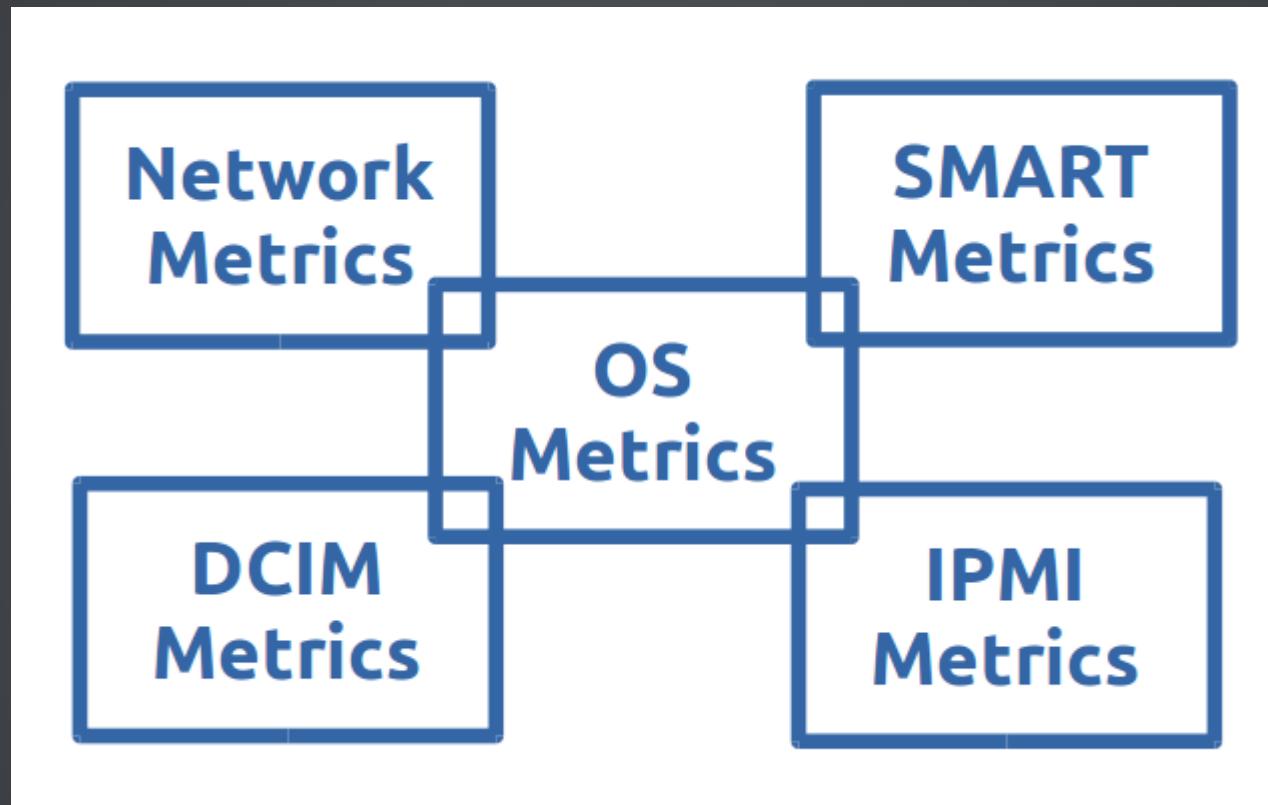
Phase 1: Measure → Collect & Store

Phase 2: Understand → Analyze & Visualize

Phase 3: Control → Monitoring

Phase 4: Improve → Anomaly Detection

Metrics



Logs

OS
Logs

Network
Logs

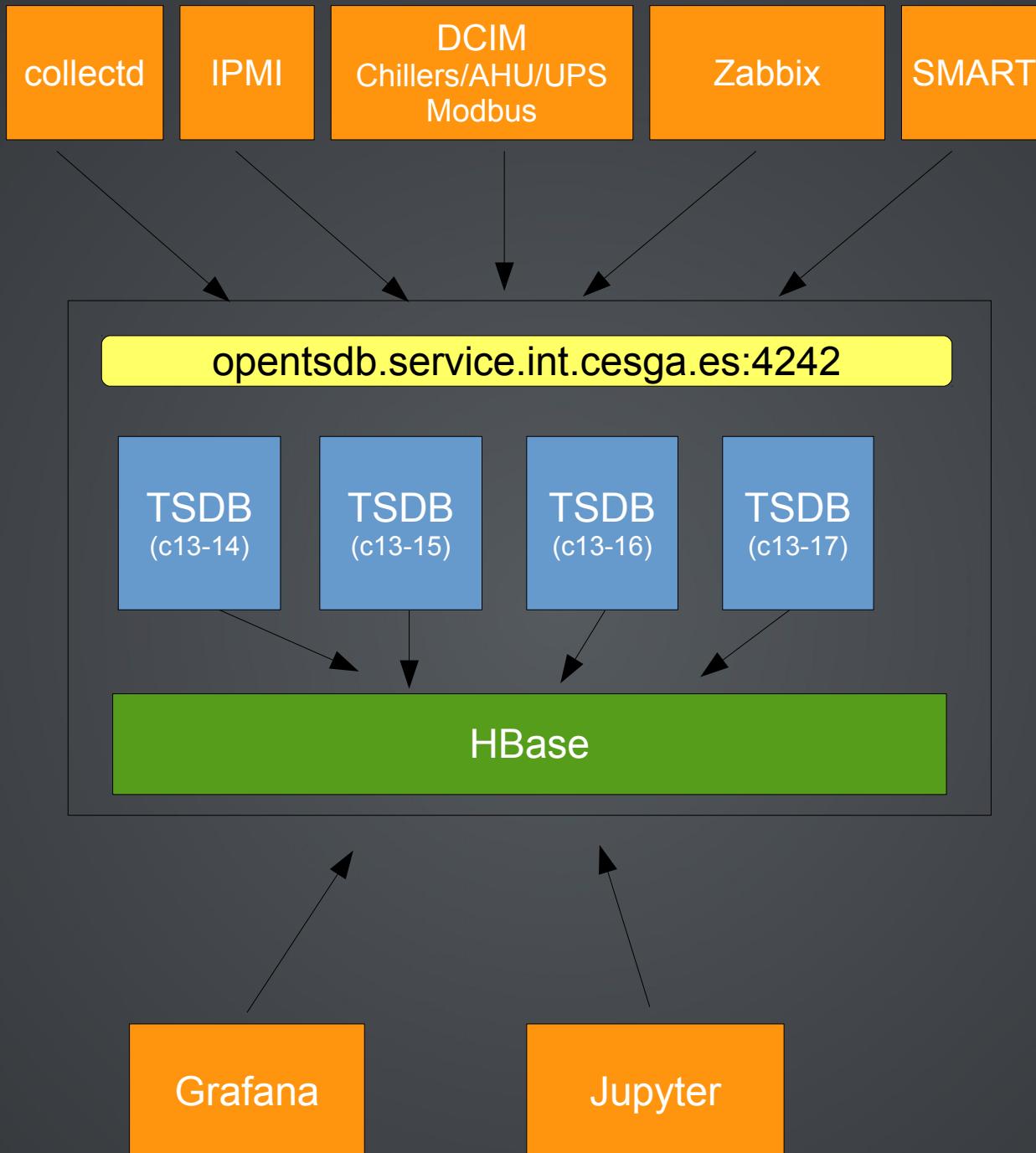
etc

SLURM
Accounting

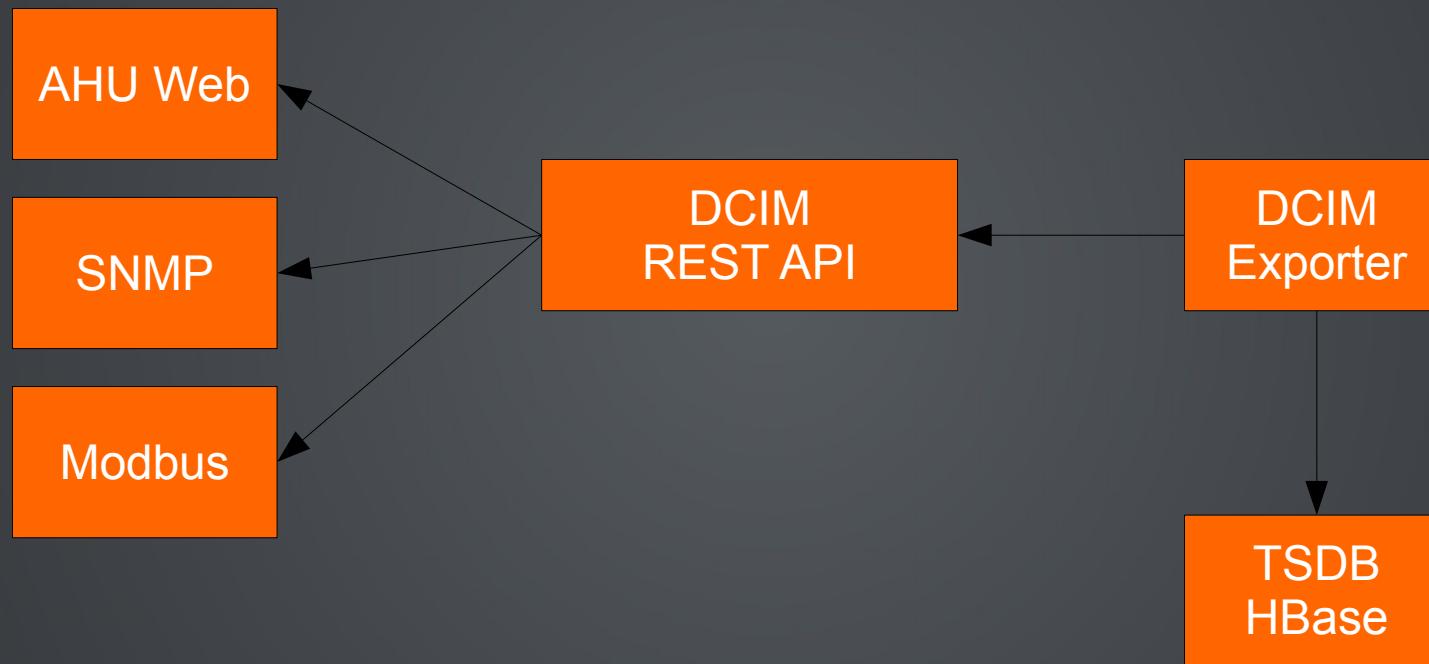
OS
Accounting

33487 metrics
10 Million time series

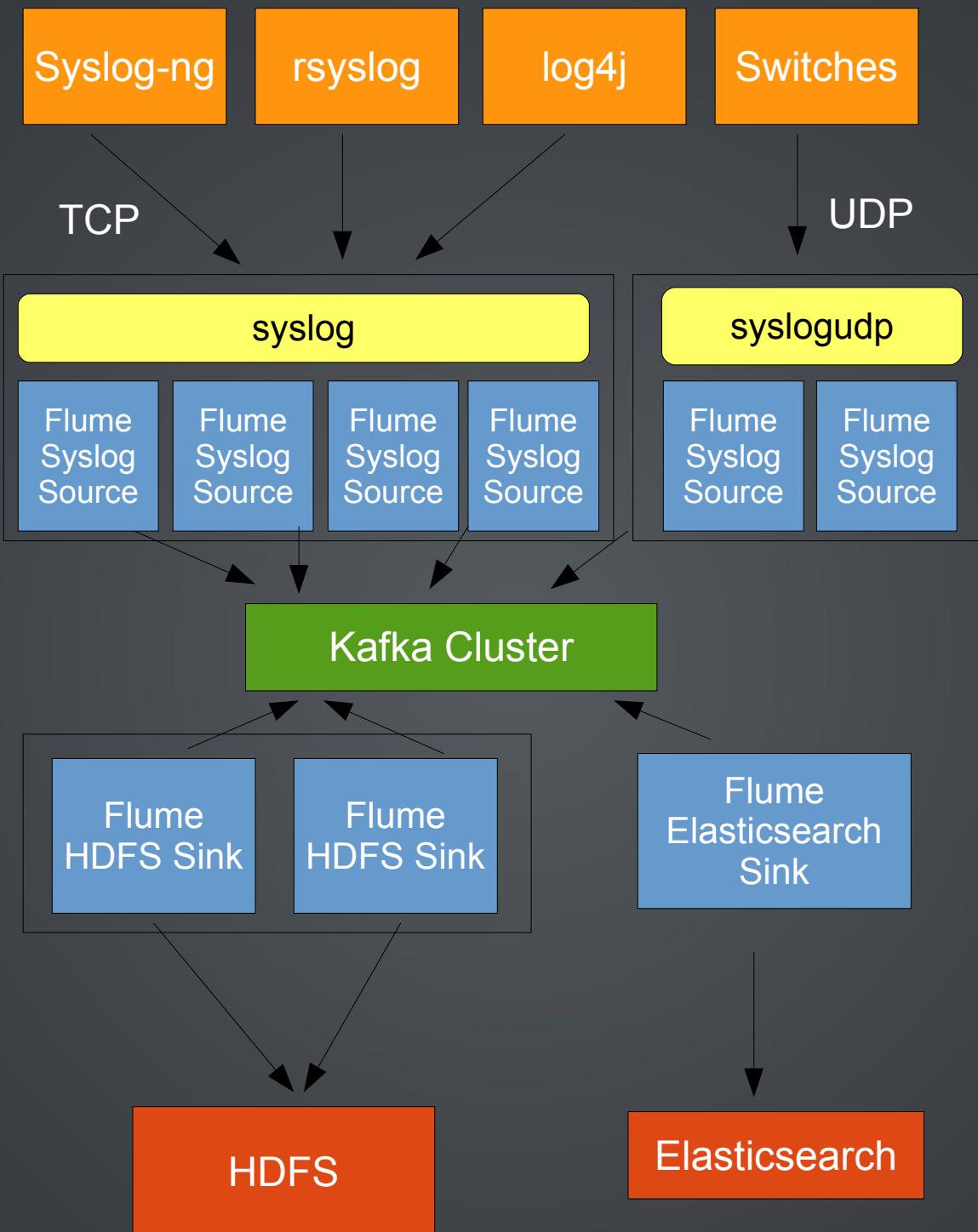
Metrics Collection Architecture



DCIM Collection Architecture

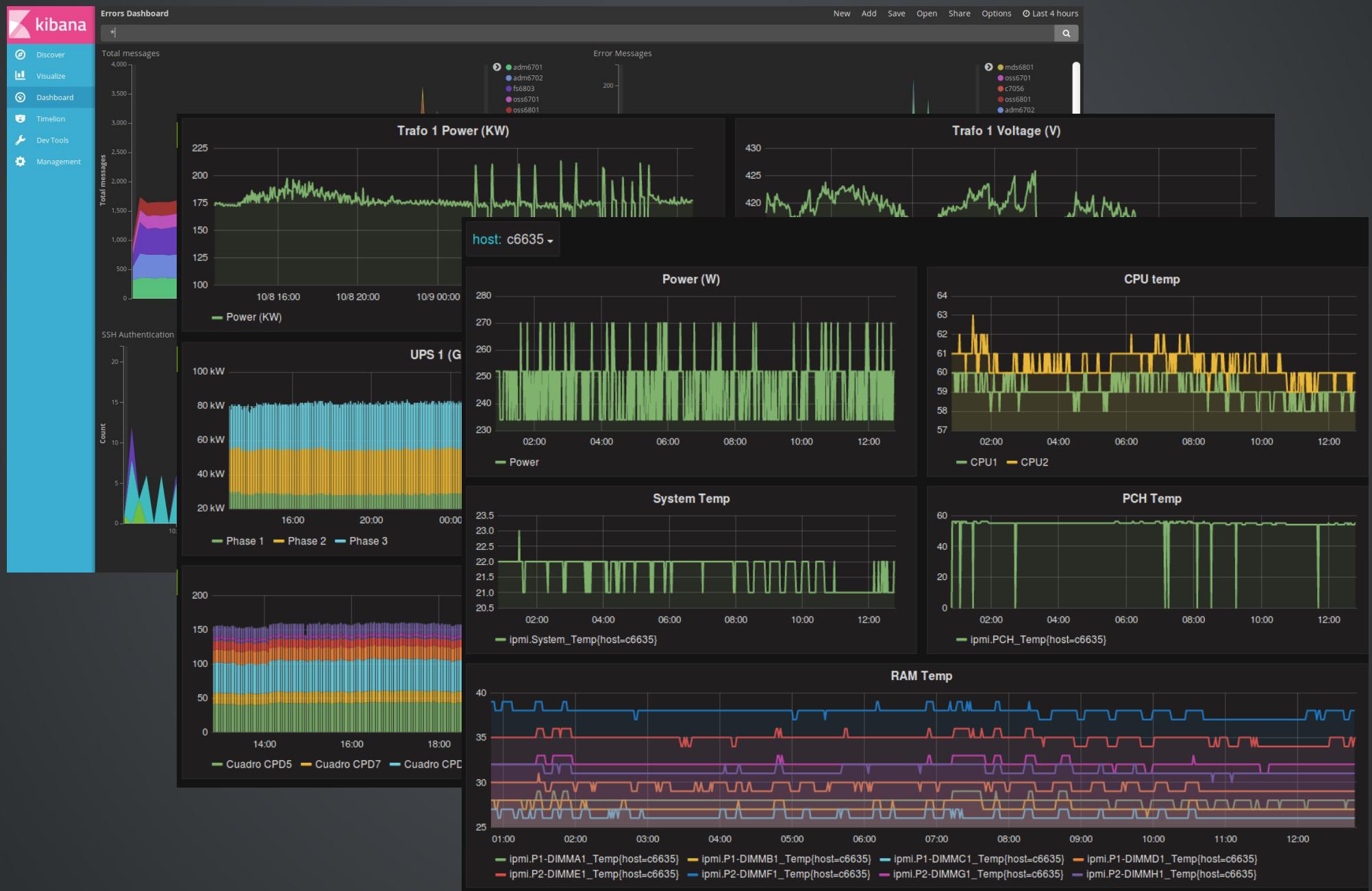


Log Collection Architecture



Phase 2: Understand

Dashboards



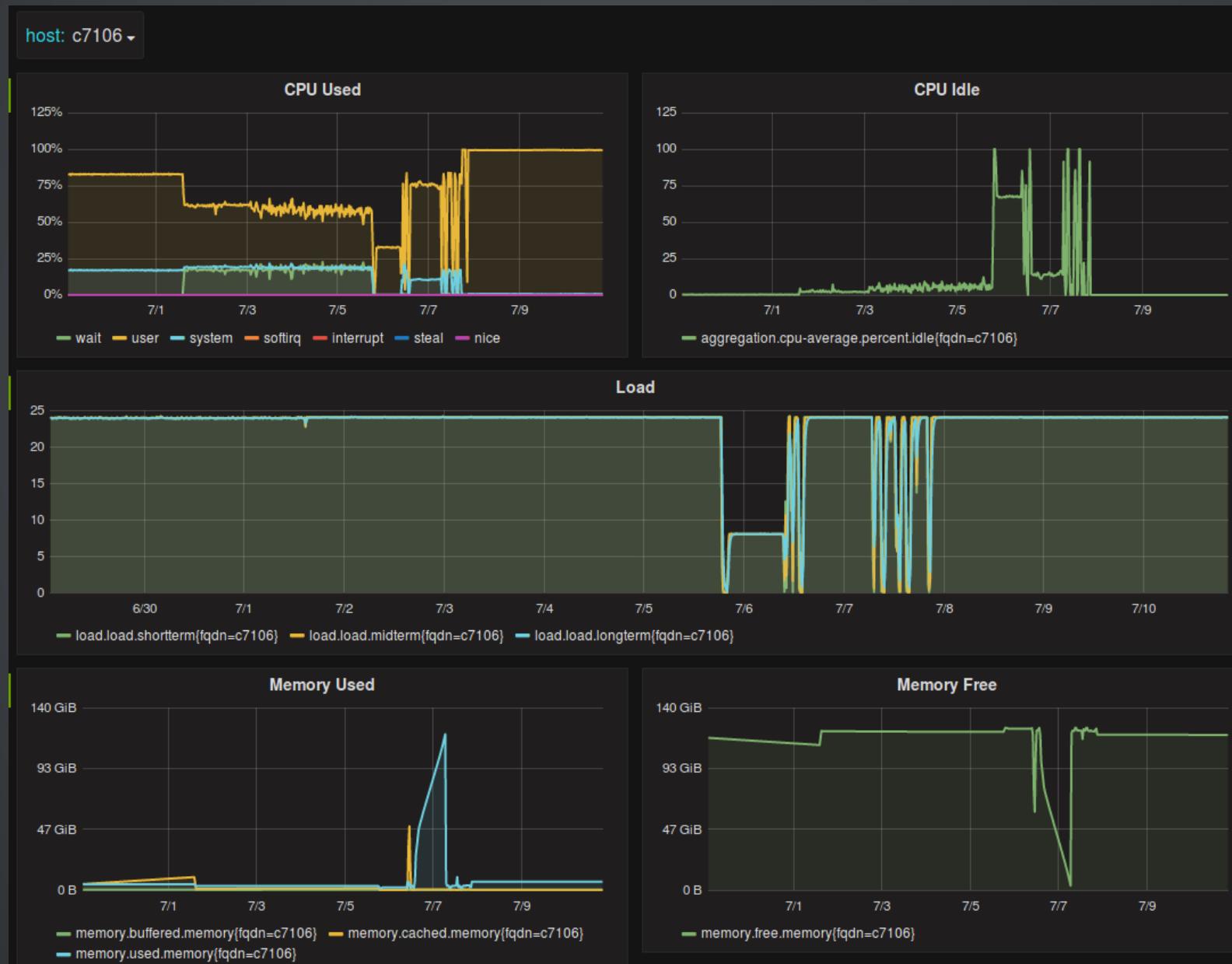
DCIM



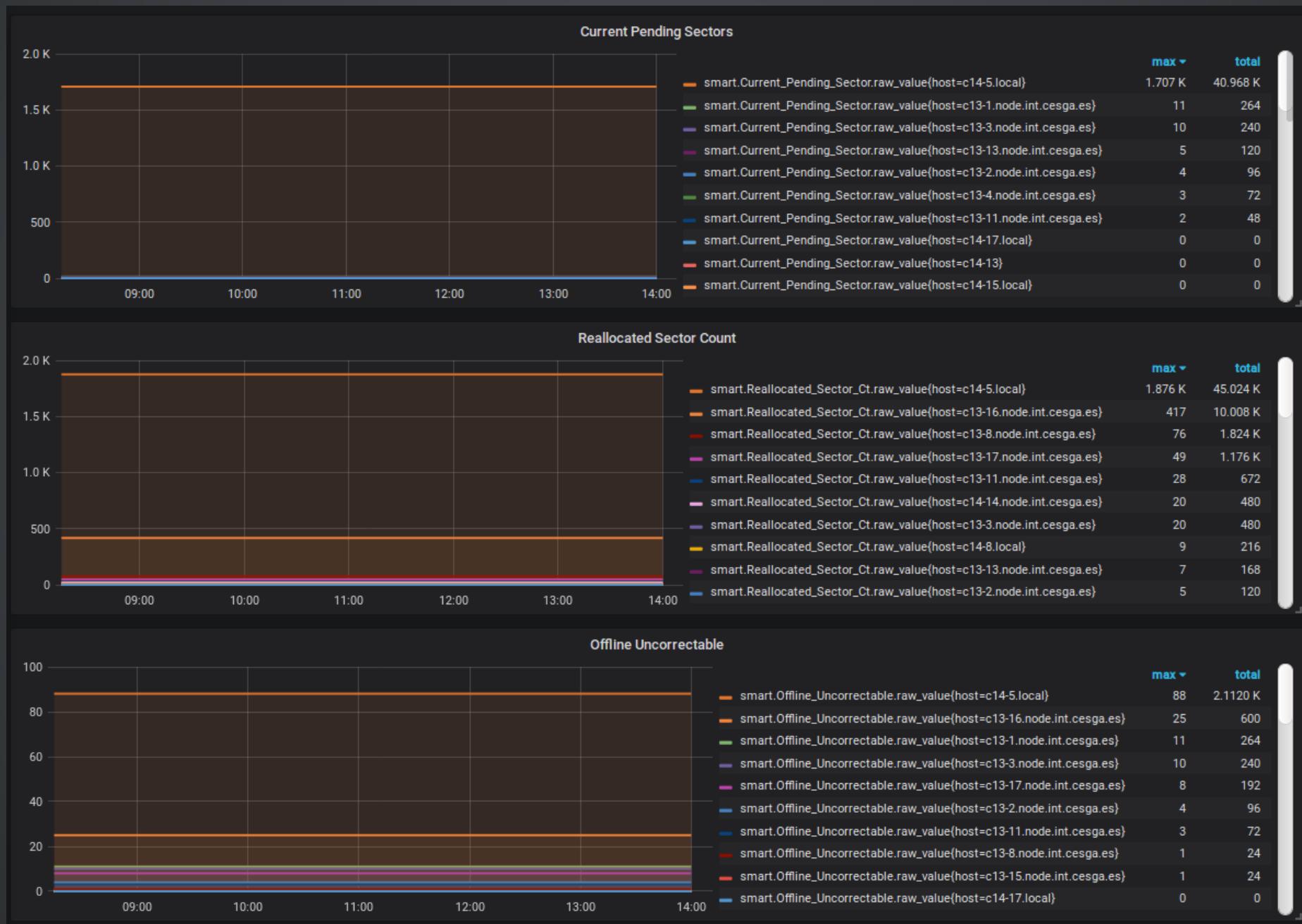
IPMI



System Metrics



SMART



Logs

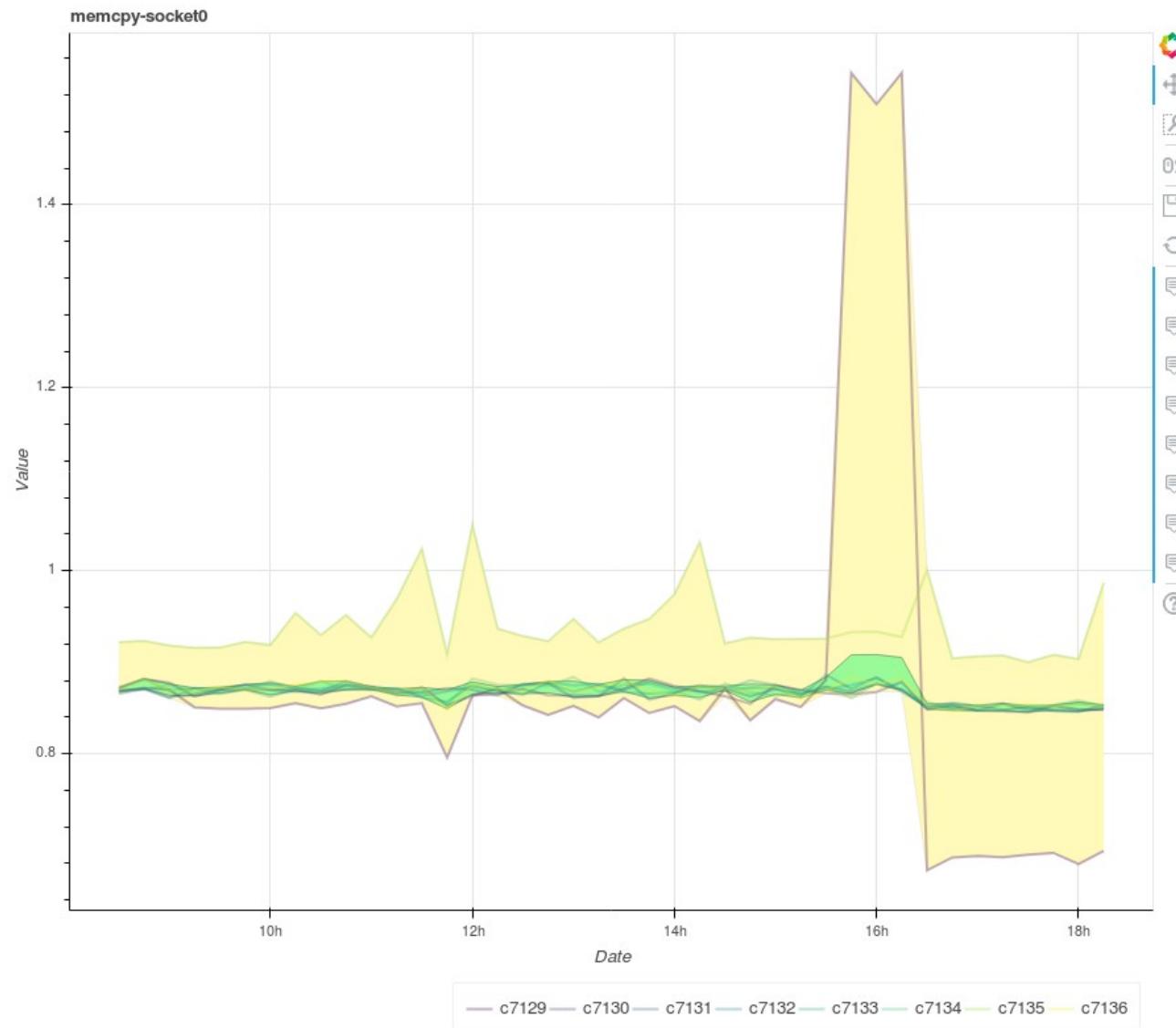




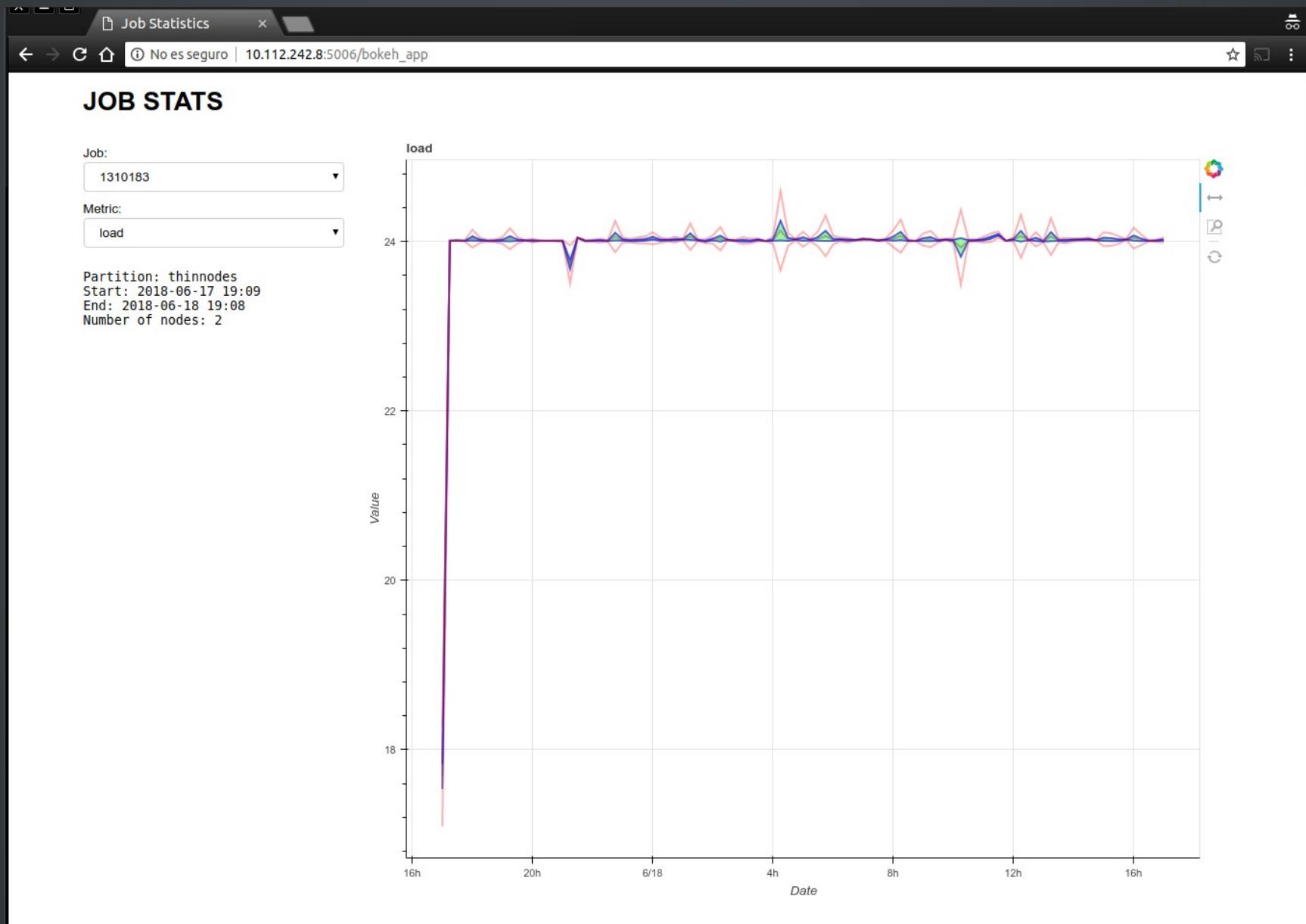
File Edit View Insert Cell Kernel Widgets Help

Trusted Python 2

job_id 1292278
metric memcpy-socket0
frequency 15min



Bokeh App



Phase 3: Control

Log Alerts: Icinga

Metric Alerts: Bosun

Bosun Items Graph Expression Rule Editor Silence Short Link jlopez ?

filter ? 0/125 Errors select all none

Needs Acknowledgement

- ALERT c6614 (critical) memory_performance_FAILING 1 alerts
- ALERT c6607 (warning) load_too_high 1 alerts

Acknowledged

- ALERT c14-8.node.int.cesga.es (critical) collectd_not_publishing 1 alerts
- ALERT c14-8.node.int.cesga.es (critical) smart_not_publishing 1 alerts
- ALERT vpsXX.cesga.es (critical) collectd_not_publishing 1 alerts
- unknown - load_too_high 18 alerts
- unknown - load_too_high 5 alerts

Phase 4: Improve

Anomaly Detection

Types of anomalies in time series

- Outliers
- Change points
- Anomalous time series

Generic Anomaly Detection

- Generic frameworks for Anomaly Detection
- They trigger too many false alerts
- In our experience it is better to trigger specific use cases

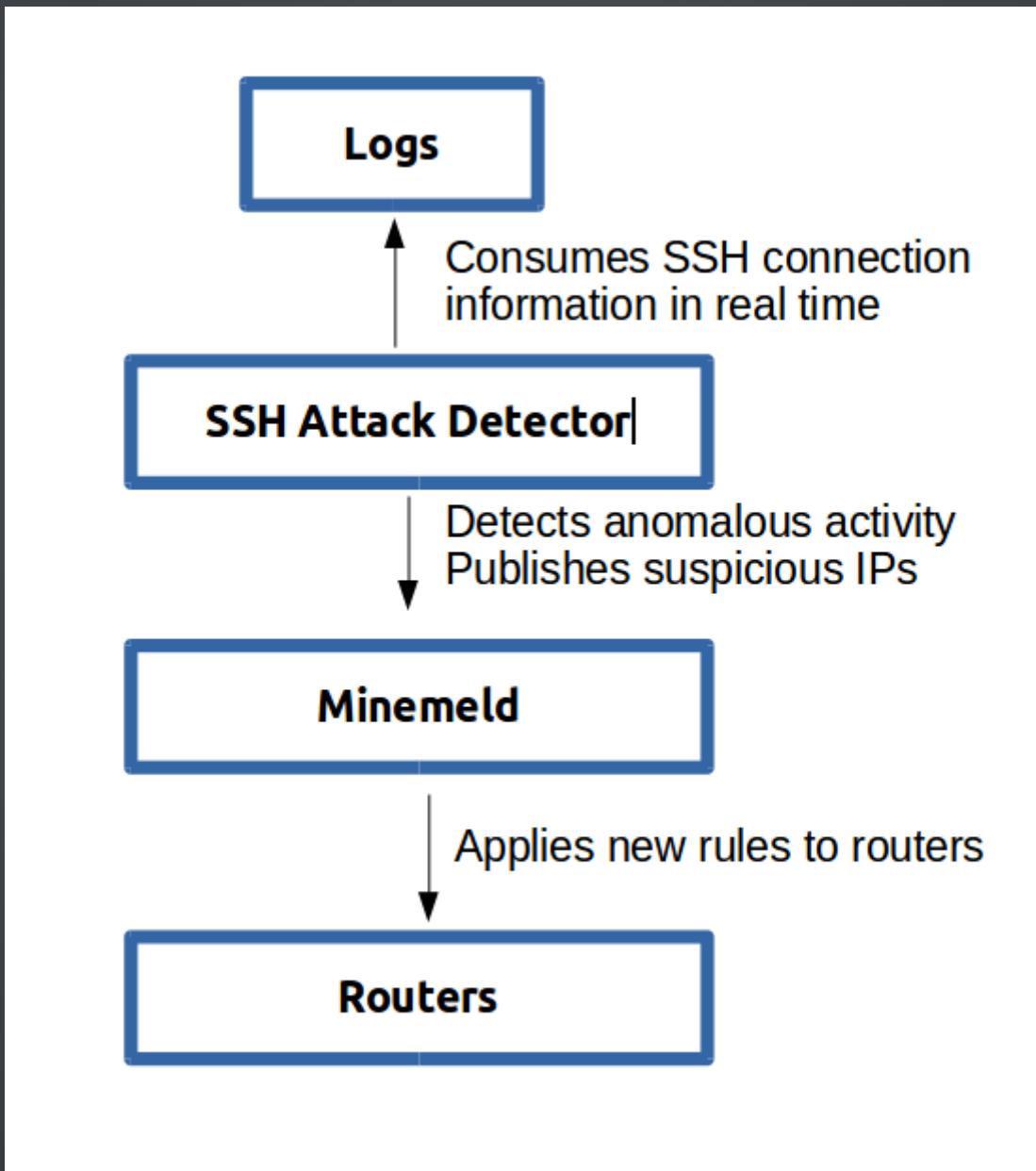
SSH Attack Detection

Problem: Daily our public servers are scanned and attacked

Detection: Correlate real-time SSH connection information to detect attacks

Objective: Automatically update router configuration to stop the attacks

SSH Attack Detection



Results

STATISTICS

7 days

METRIC	CURRENT	HISTORY (LAST 7D)
--------	---------	-------------------

INDICATORS	59	 <p>7 DAYS AGO 534</p>
------------	----	--

METRIC	SINCE ENGINE START	HISTORY (LAST 7D)
--------	--------------------------	-------------------

ADDED	609	
-------	-----	--

AGED_OUT	691	
----------	-----	--

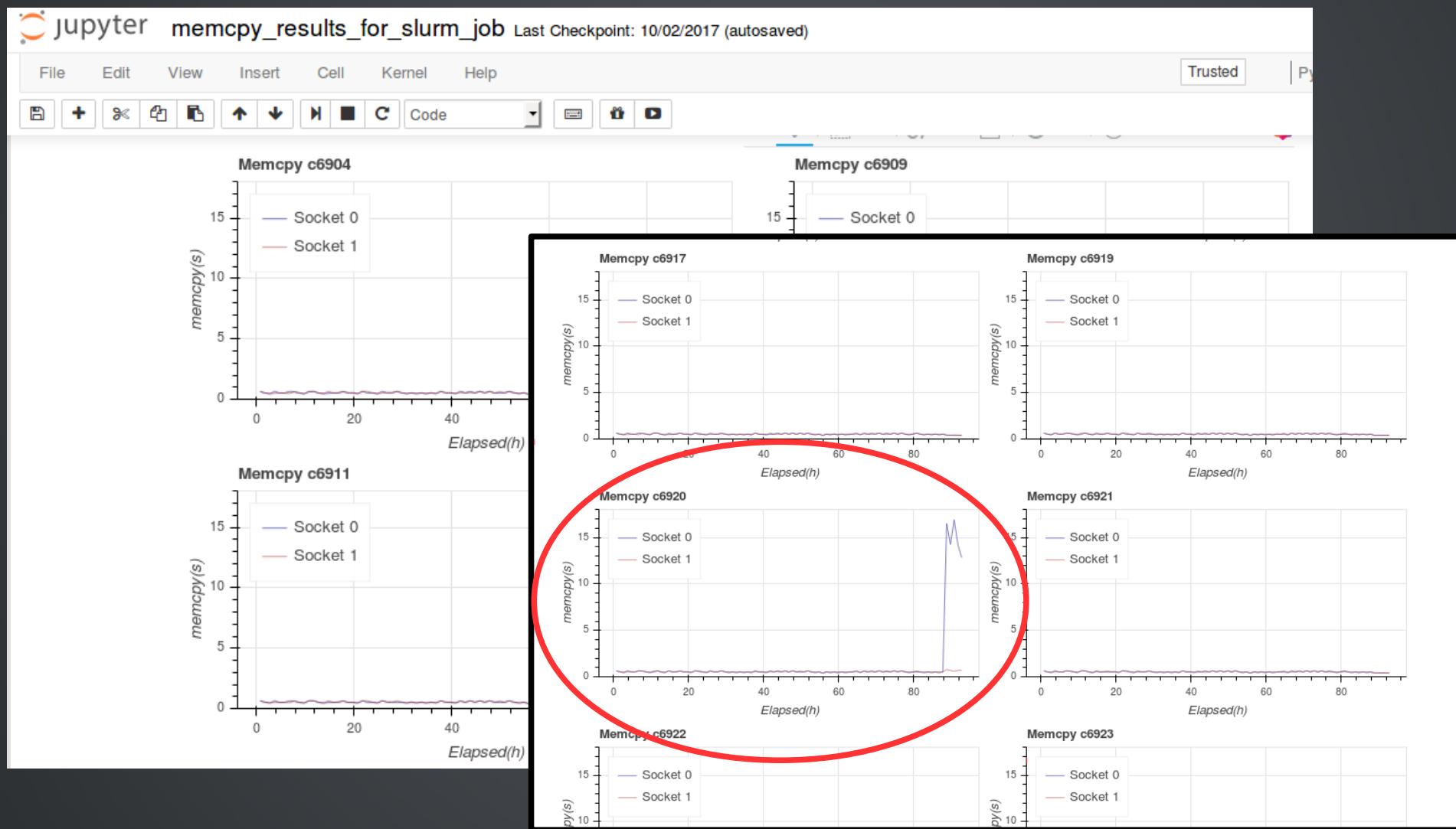
Server Anomalous Performance

Problem: Parallel jobs are cancelled because some of the nodes have poor performance. Computation is lost.

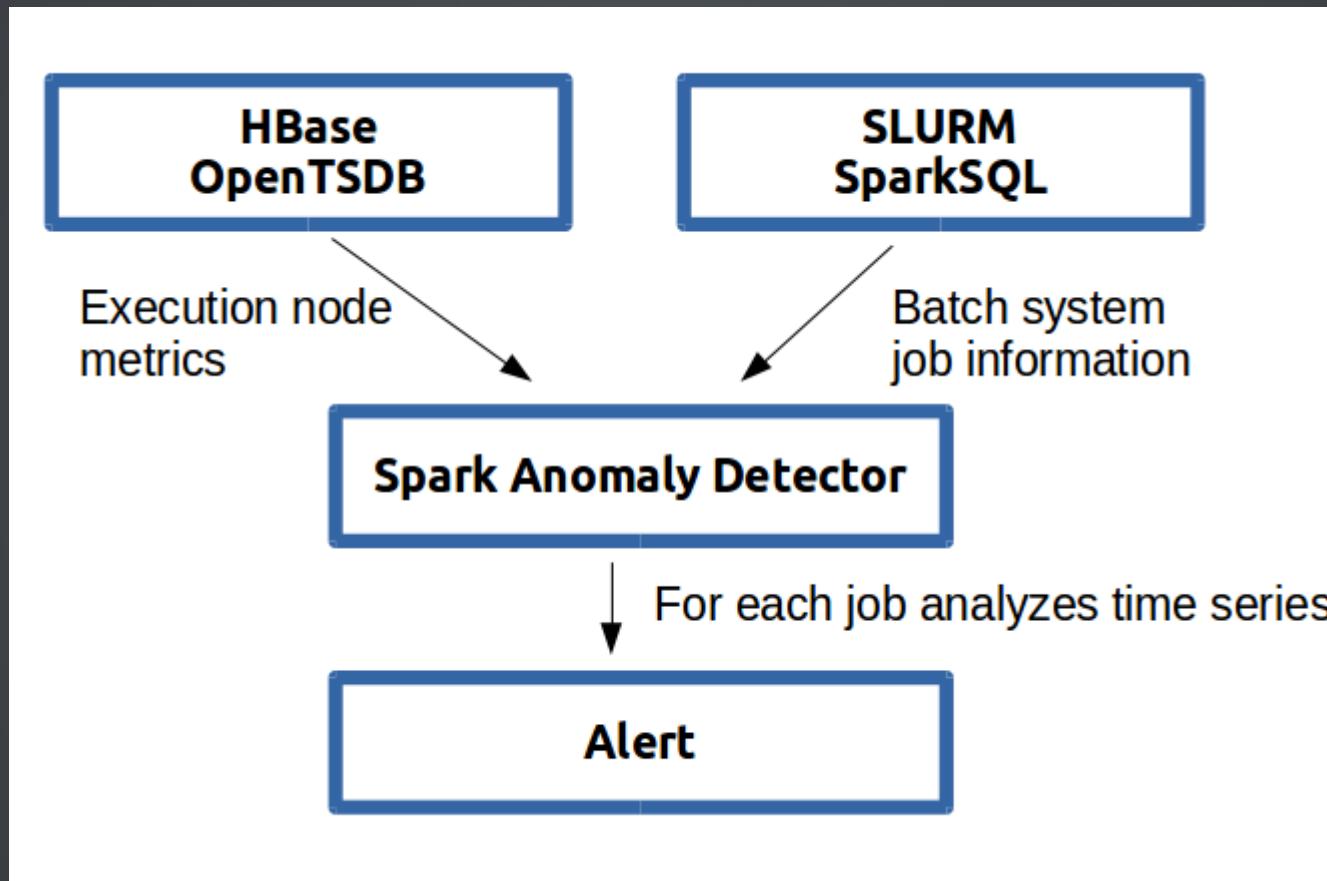
Detection: Analyze & visualize server metrics to spot the anomalous node

Objective: Automatically detect low performance nodes

Analyze & Visualize



Anomalous Performance Detection



Conclusions

- No longer needed to delete old data
- Generic anomaly detection systems generate too many alerts
- Target specific use cases to maintain number of alerts low



XUNTA
DE GALICIA



GOBIERNO
DE ESPAÑA

MINISTERIO
DE ECONOMÍA
Y COMPETITIVIDAD



CSIC
CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS



Unión Europea
Fondo Europeo de
Desarrollo Regional
"Una manera de hacer Europa"



CESGA



THANKS!

[jlopez [at] cesga [dot] es]