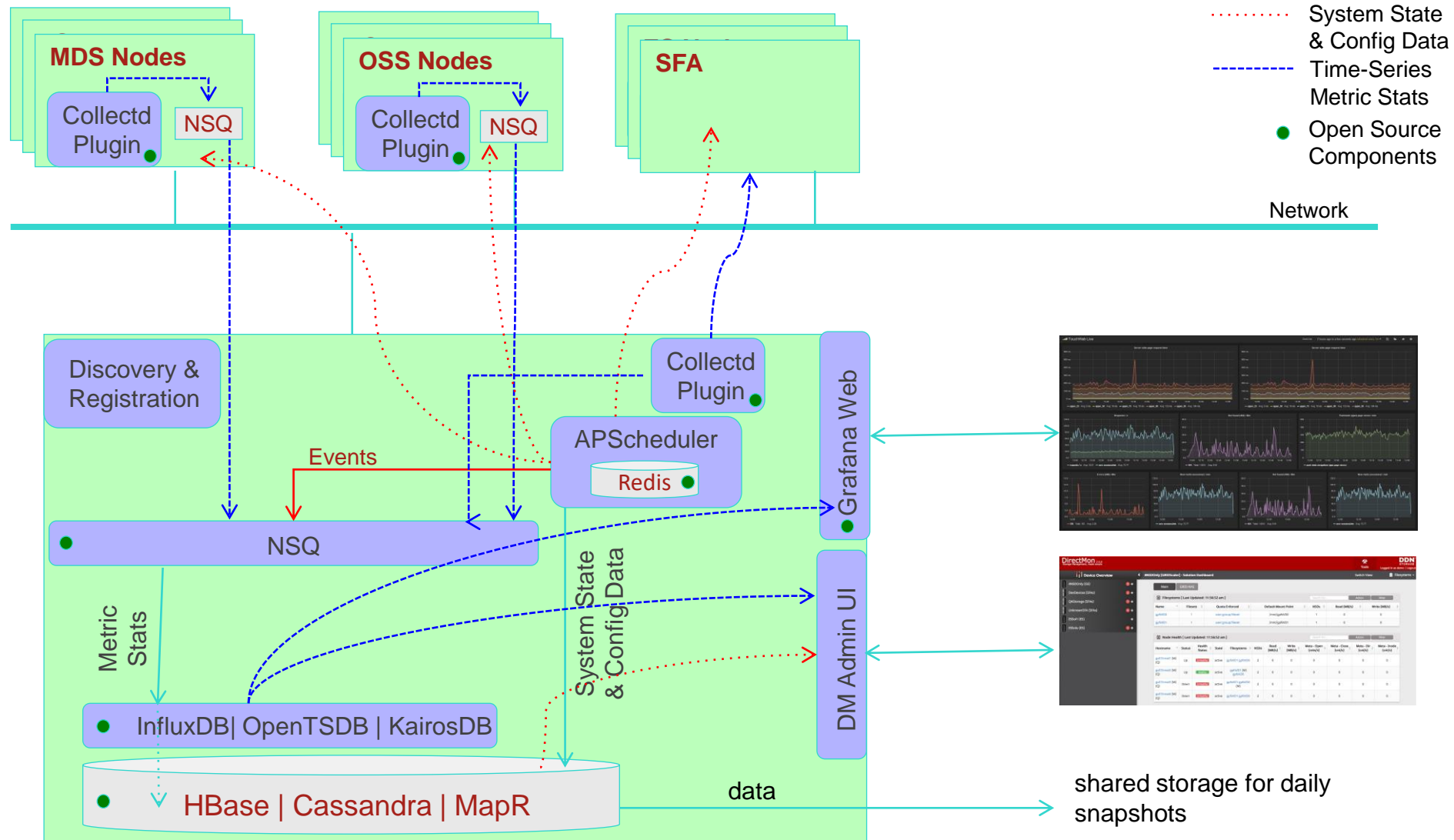# Lustre文件系统的大规模性能监控与IO模式分析

顾政

gzheng@ddn.com

# Background of Lustre Performance monitoring

- ► Activities on the Lustre are black box
  - Users and Administrators want to know "what's going on?"
  - Find "Crazy Jobs" in advance to prevent slow down.
- ► Lustre statistics are valuable big data
  - Not only monitoring and visualization, but also analysis
  - Predictable operations could be possible.
  - It helps optimize applications and data relocation.
- ► Open Source based monitoring tool
  - In general, open source is common in the HPC system and it's straightforward.
  - Various combination is possible and make new use cases.

# C/S monitoring

► Collecting Data from target, usually it could be MDS/MGS, OSS, client.

► Sending collected data to persistent storage.

► Collected data could be reviewed by users friendly.(Time series, Rates etc.)
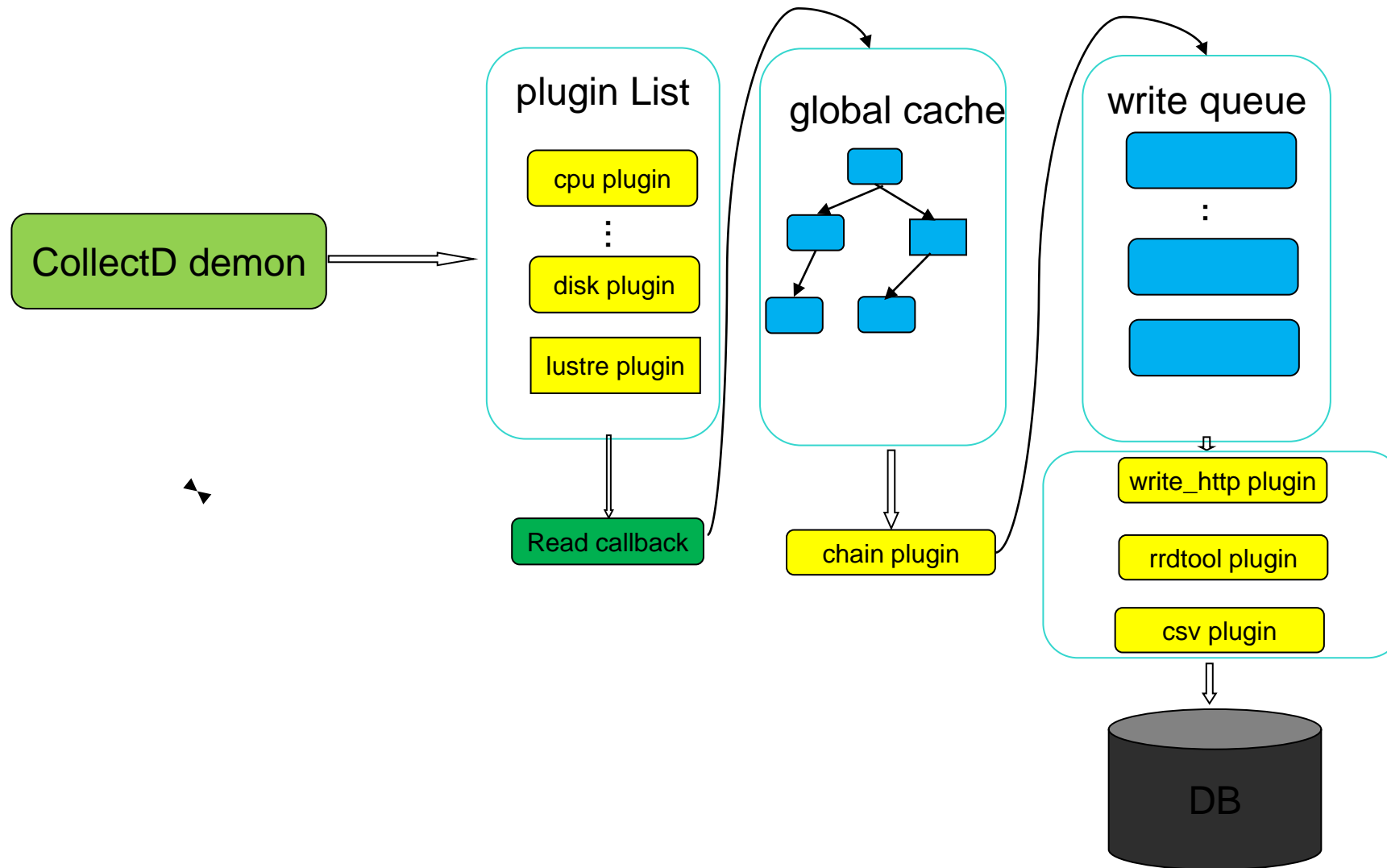
# Standalone Configuration

# Components of Lustre Performance monitoring

► Data collecter
  • Collects statistics from Lustre /proc and sends them to monitoring server over the network.
  • Runs on servers as well as client and routers.

► Backend Storage
  • Receive stats from agents and store them into database.
  • It can be historical and query-able data

► Frontend
  • Collected data is not only visualized, but also analytics.
  • Application I/O analytics

► One-click setup
  • Native one-click script support, no need to assemble components manually, no worries about compatibility

# Flexible data collector

► A lot of agents existed to collect Lustre performance statistics

► Collectd is one of reasonable options

- Actively developed, supported and documented

- Running at many Enterprise/HPC system

- Written in C and over 100 plugins are available.

- Supports many backend database for data store.

- Unfortunately, plugin for lustre is not available, but we made it!

# A glance at Collectd

# Scalable backend data store

► **RDD and SQL based data store dose not scale**

- RDD works well on small system, writing 10M statics into files are very challenging (few million IOPS!)

- SQL is faster than RDD, but still hit next level scalability. And it's complex to make database deign.

► **NoSQL based key-value store shines**

- InfluxDB/Hbase. KairosDB/Cassandra

- key, value and tags are easy adaption for Lustre statics data store. No need complex database schema.

- Need to be aware of managing for statics data archiving. (retention)

# Frontend – Why Grafana ?

**Whamcloud**

▶ **Visualize**

Heatmaps, histograms, graphs to geomaps..

▶ **Alert**

Seamlessly define alerts where it makes sense

▶ **Unify**

Supports InfluxDB, Graphite, Elasticsearch, OpenTSDB and Prometheus.

▶ **Extend**
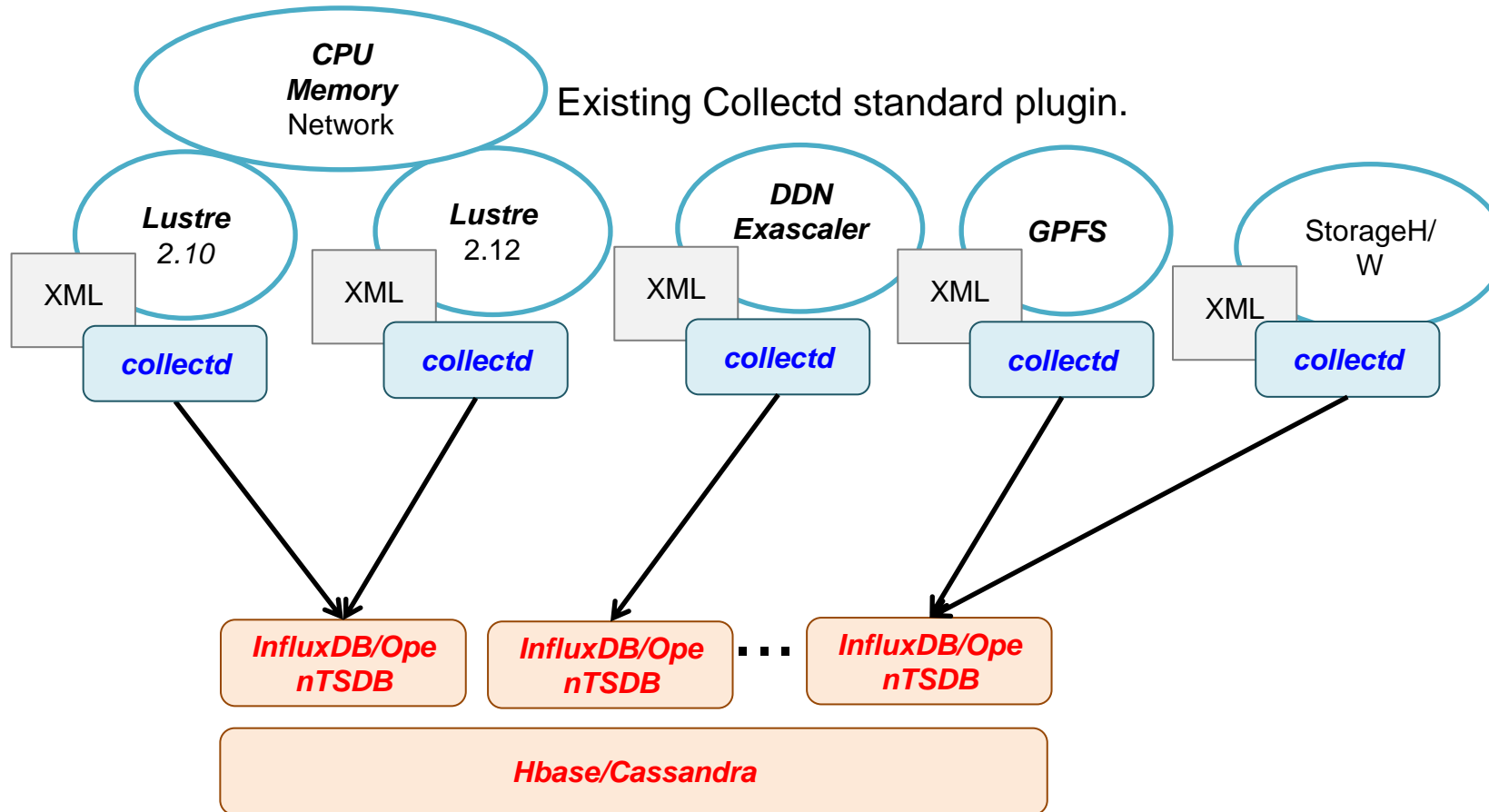
Easy to customize dashboards and plugins

▶ **Open**

Completely open source, and backed by a vibrant community

# Deign of plugin for lustre in collectd

► **A framework consists of two core components**

- Common platform, filedata plugin, collect data by reading and parsing a set of files (not only Lustre)

- Statistics definition layer(XML file and XML parser)

► **Defined XML for Lustre /proc information**

- A single XML file for all definitions of Lustre data collection

- No need to maintain massive error-prone scripts.

- Extendable without core logic layer change.

- Easy to support multiple Lustre version and Lustre distributions in the same cluster.

# Architecture of lustre-plugin

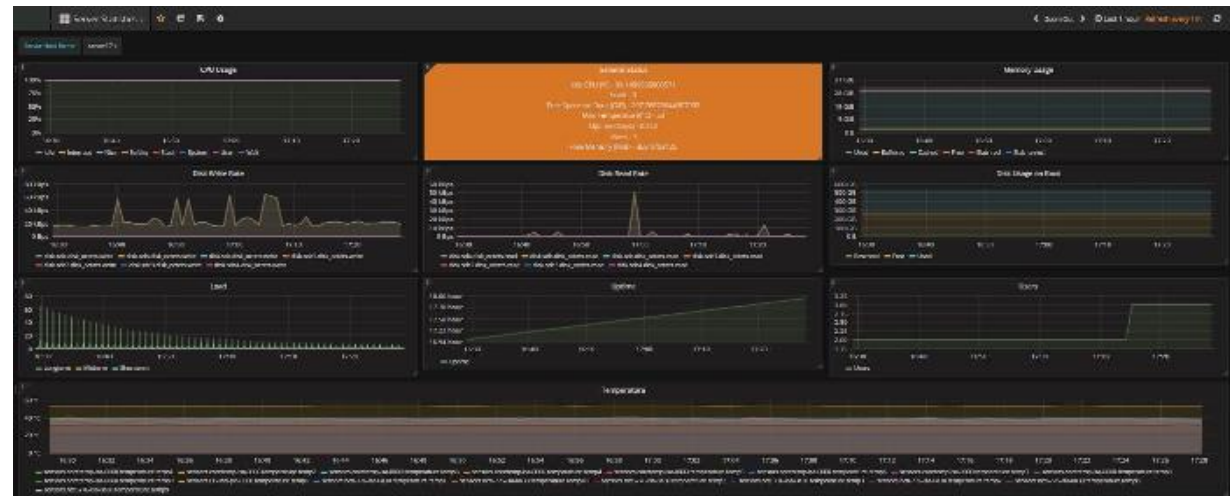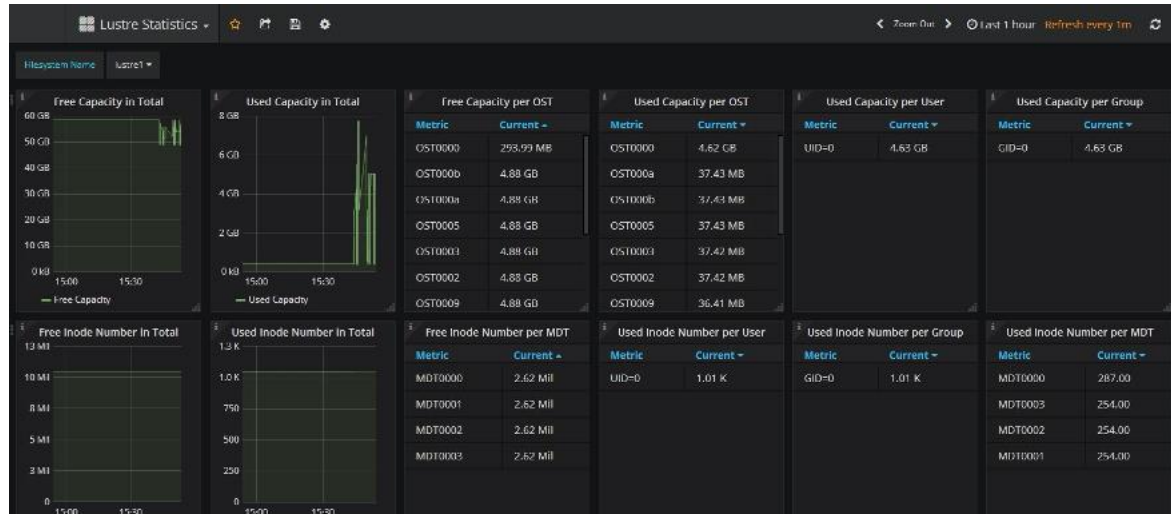# Application aware I/O monitoring

► Scalable backend data store

- Now, we have scalable backend data store InfluxDB/OpenTSDB.
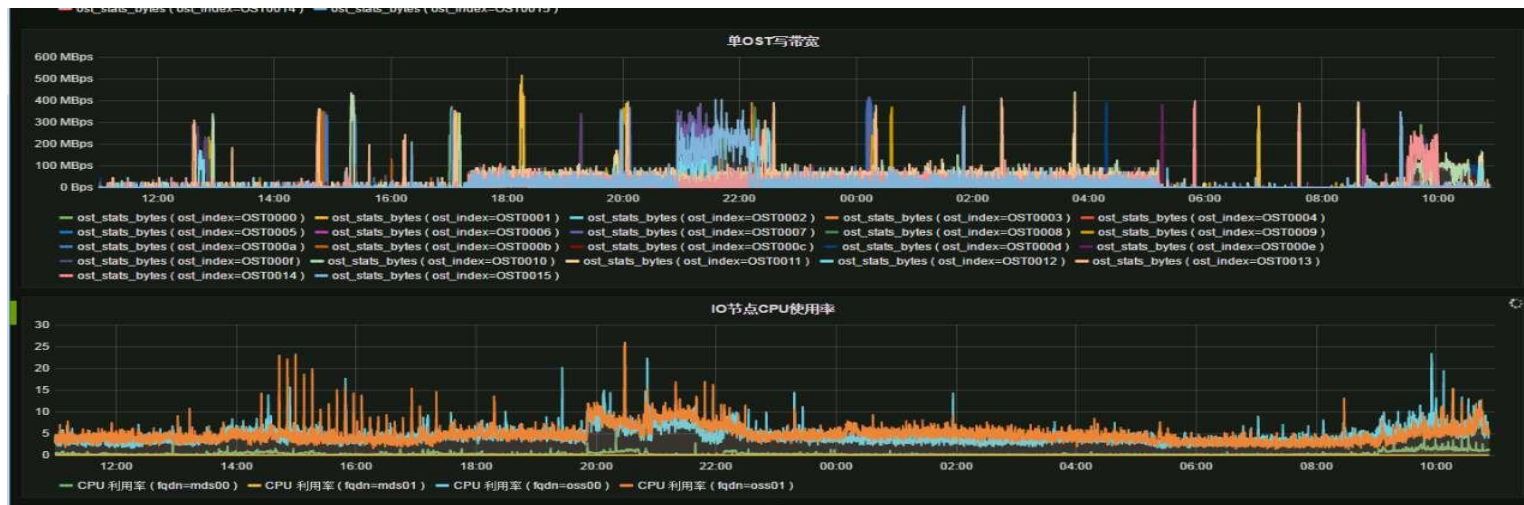
- Store any type of mercies whatever we want to collect.

► Lustre Job stats is awesome, but need to be integration.

- Lustre JOB stats feature is useful, but administrator is not interested in I/O stats just only based on JOBID. (Array jobs. Job associates with another jobs, e.g. Genmic pipeline)

- Lustre performance stats should be associated with all JOBID/GID/UID/NID or custom any IDs.

# Pictures of Lustre PerfMonitoring at customer site

# Pictures of Lustre PerfMonitoring at customer site

# One story about Lustre PerfMonitoring at OIST

Lustre cluster configuration :

► 3PB Lustre filesystem (12 x OSS, 400 x client)

► Lustre jobstats integrated with SLRUM and running on the production system

Lustre PerfMonitoring configuration:

► Unique Lustre Job stats configuration with Collectd Lustre plugin and runs on existing on Jobstats framework.

► Collect jobs stats associated with all UID/GID/JOBID and store them into OpenTSDB.


With the help of Lustre PerfMonitoring, customer found out the root cause so quickly why unexpected burst I/O happened , which they suffered from for a long time.

Whamcloud

**Thank You!**