# Storage System at the supercomputer Fugaku

**Yuichi Tsujita**

**Operation and Computer Technologies Division, RIKEN Center for Computational Science (R-CCS)**

# Outline

- **FEFS: Lustre-based file system enhanced by FUJITSU LIMITED**

- **Storage system at the K computer**

- **Overview of the supercomputer Fugaku**

- **Three-level hierarchical storage system**

- **Monitoring and log collection**

- **Summary**

# FEFS: Lustre-based file system enhanced by FUJITSU LIMITED

# Introduced FEFS in our site

- **FEFS: Fujitsu Exabyte File System**
  - Enhanced Lustre by FUJITSU LIMITED

- **FEFS based on Lustre ver. 1.8**
  - Adopted in the two-level file system of the K computer (hereinafter, "K")
  - High I/O throughput under the huge number of clients
  - Many enhancements to have stable and high performance operations

- **FEFS based on Lustre ver. 2.10  (<- Long Term Support)**
  - Adopted in the 2nd layer storage system of the supercomputer Fugaku (hereinafter, "Fugaku")
  - Cooperative operation with the 1st layer storage system built by SSDs for high throughput I/O in computing and mitigation of load of the 2nd layer storage system
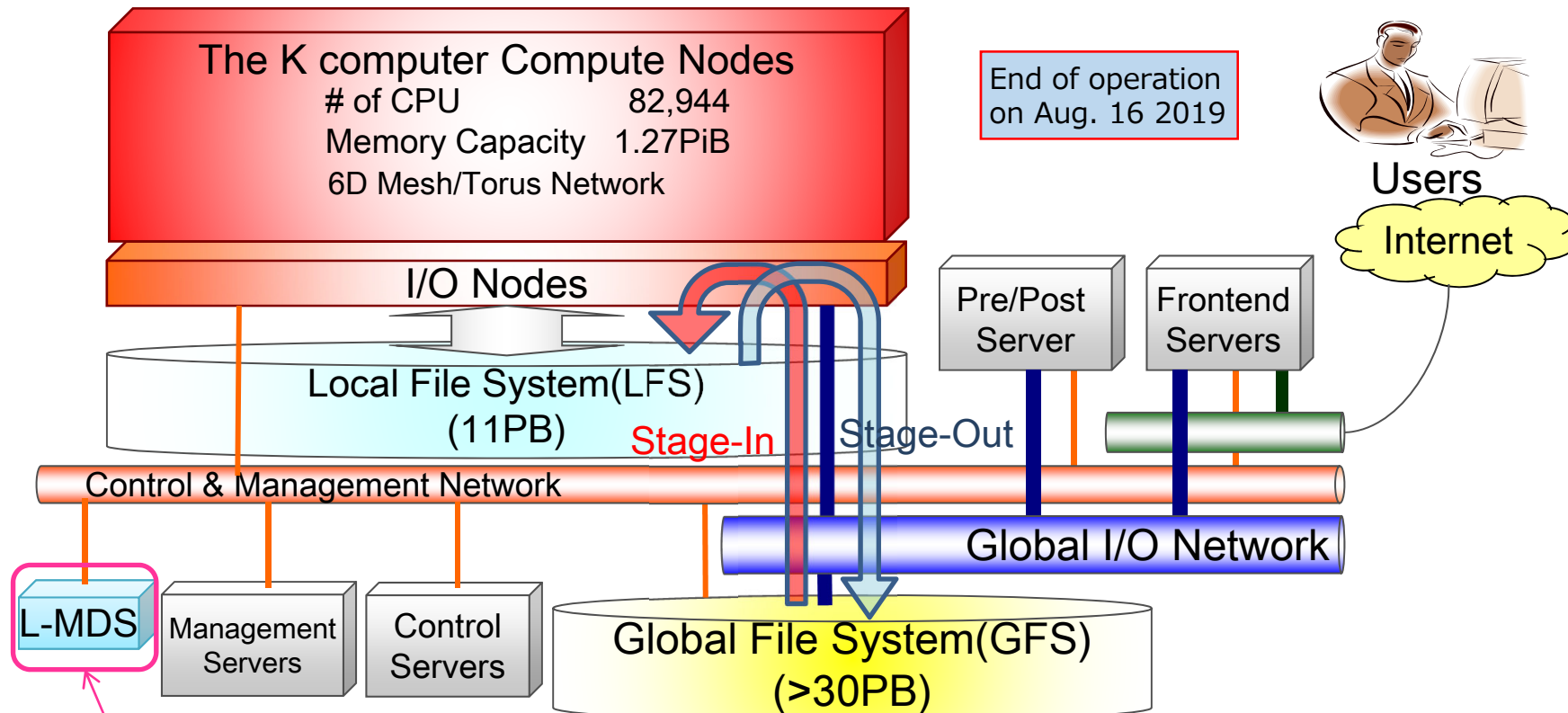  - Full deployment and optimization are still in progress.

# Notable Features of FEFS

- **Enhancements based on Lustre 2.x may contribute to the Lustre community.**
  - FUJITSU LIMITED is a member of the community and they will continue to report bug-fixes and feedbacks to the community with cross relationship.

- **Own enhancements about RAS, system operability, tolerance under high I/O load, and fair-share management among clients are expected to perform well at the 2nd layer storage system.**

# Storage system at the K computer

# Storage system at the K computer

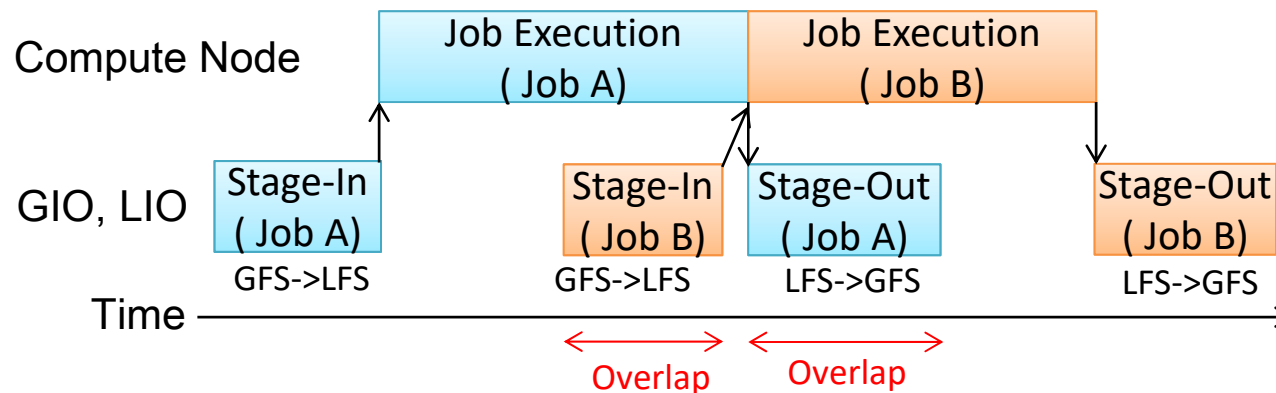- **File staging with two-level local/global file system**



The K computer Compute Nodes
# of CPU 82,944
Memory Capacity 1.27PiB
6D Mesh/Torus Network

End of operation on Aug. 16 2019

Users
Internet

I/O Nodes

Local File System(LFS)
(11PB)

Stage-In   Stage-Out

Pre/Post Server

Frontend Servers

Control & Management Network

Global I/O Network

L-MDS

Management Servers

Control Servers

Global File System(GFS)
(>30PB)

MDS for LFS

FEFS was used for both LFS and GFS.
(FEFS: **F**ujitsu **E**xabyte **F**ile **S**ystem based on Lustre technology)

# File Staging

- **Asynchronous file staging for effective job scheduling and I/O**

| | | | | |
|---|---|---|---|---|
| Compute Node | | Job Execution ( Job A) | Job Execution ( Job B) | |

| GIO, LIO | Stage-In ( Job A) | | Stage-In ( Job B) | Stage-Out ( Job A) | | Stage-Out ( Job B) |
|---|---|---|---|---|---|---|
| | GFS->LFS | | GFS->LFS | LFS->GFS | | LFS->GFS |

Time →

← Overlap → ← Overlap →

- **Pros:**
- ✓ Stable application performance for jobs with the help of overlaps between job executions and file staging
- **Cons:**
- ✓ Pre-defining file name of stage-in/out operation lacks of usability.
- ✓ Data-intensive application which requires a huge storage space affects system utilization because of waiting stage-in/out processing of other jobs.

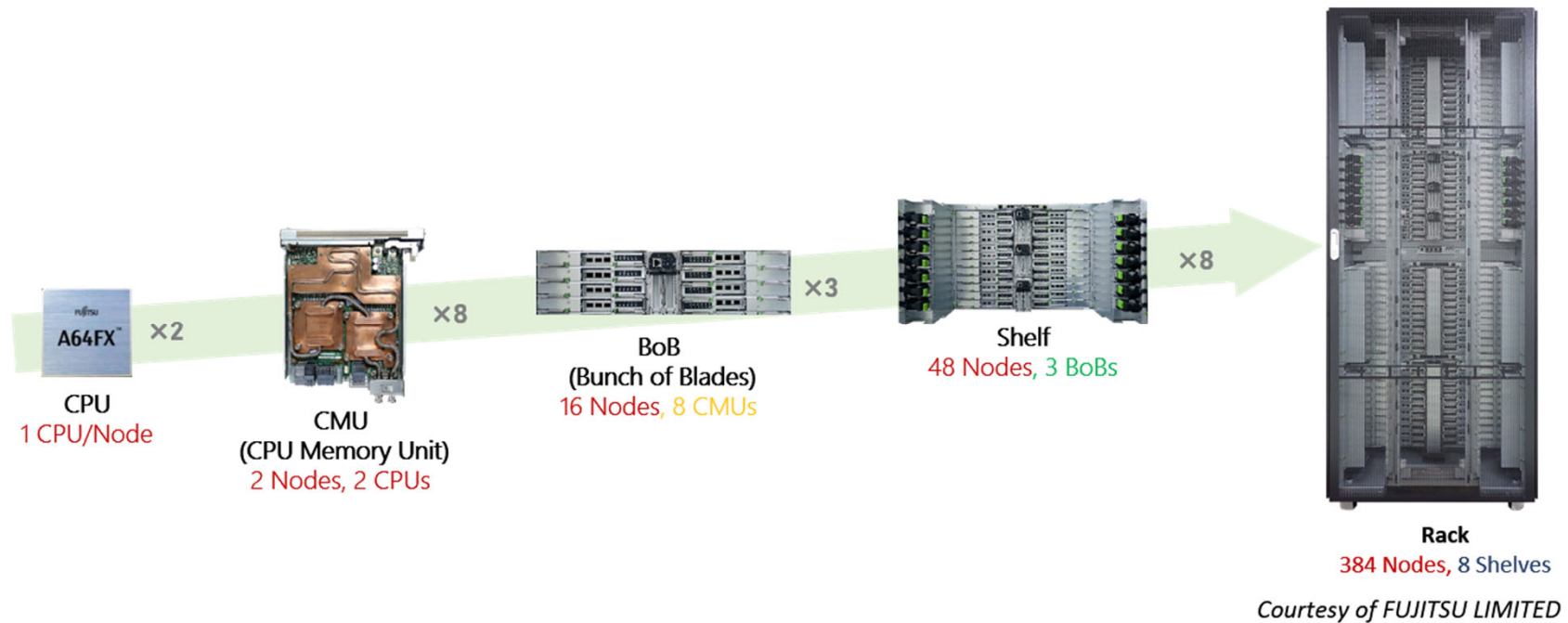# Overview of the supercomputer Fugaku

# From "K" to "Fugaku"

- **Performance of ten racks of "Fugaku" is almost the same performance of "K"(864 racks).**

| | | "Fugaku" | "K" |
|---|---|---|---|
| | CPU Architecture | A64FX<br>Arm v8.2-A SVE (512 bit SIMD) | SPARC64VIIIfx |
| **Node** | Cores | 48 | 8 |
| | Peak DP performance | 2.7+ TF | 0.128 TF |
| | Main Memory | 32 GiB | 16 GiB |
| | Peak Memory Bandwidth | 1,024 GB/s | 64 GB/s |
| | Peak Network Performance | 40.8 GB/s | 20 GB/s |
| **Rack** | Nodes | 384 | 102 |
| | Peak DP Performance | 1+ PF | < 0.013 PF |
| | Process Technology | 7 nm FinFET | 45 nm |

# Hardware Configuration of "Fugaku"

- **From CPU to Rack**



CPU
1 CPU/Node

×2

CMU
(CPU Memory Unit)
2 Nodes, 2 CPUs

×8

BoB
(Bunch of Blades)
16 Nodes, 8 CMUs

×3

Shelf
48 Nodes, 3 BoBs

×8

Rack
384 Nodes, 8 Shelves

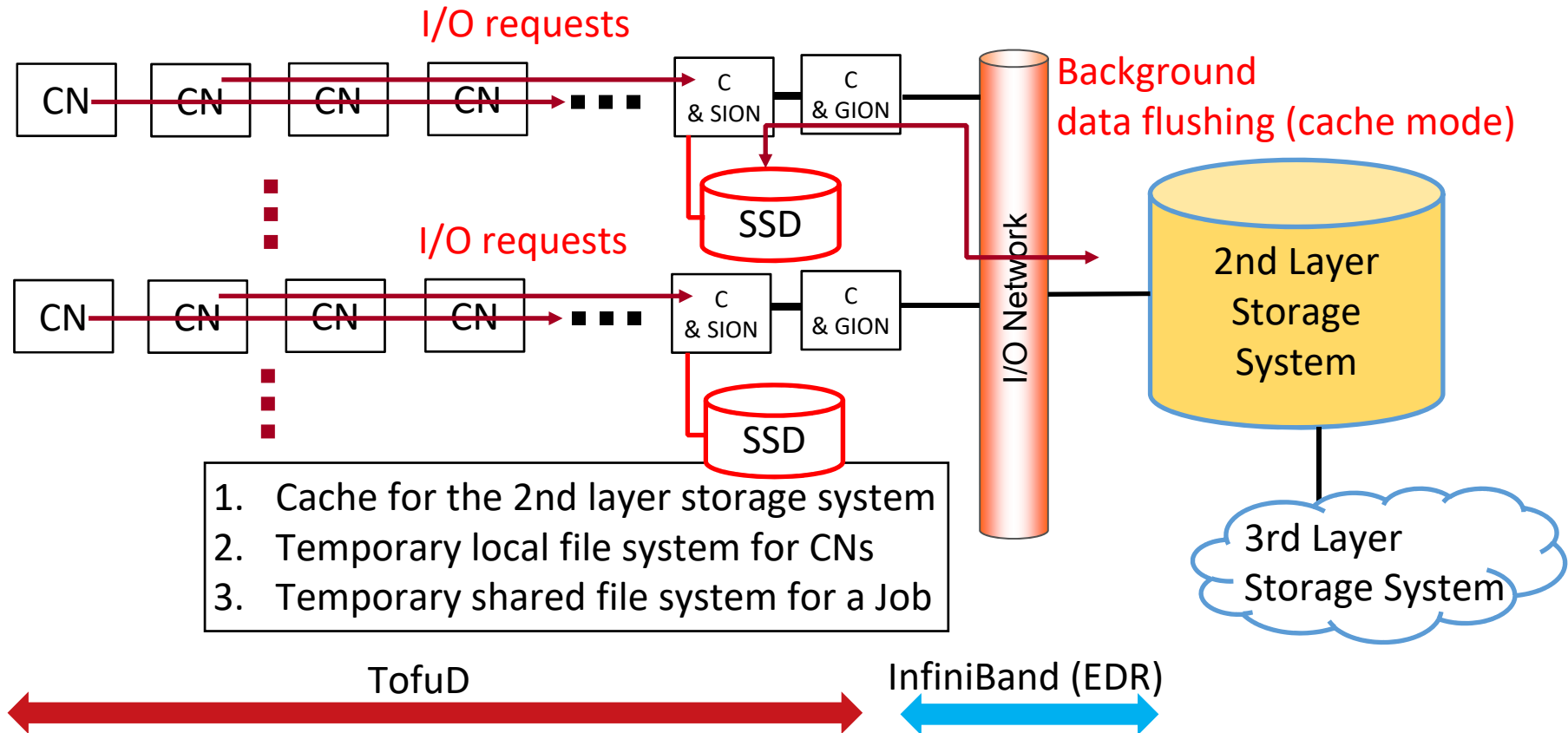*Courtesy of FUJITSU LIMITED*

# System Configuration

- **158,976 nodes**
- **Two types of nodes**
  - Compute node and Compute & I/O node connected by Fujitsu TofuD (6D mesh/torus Interconnects)
- **3-level hierarchical storage system**
  - 1st layer
    - One of 16 compute nodes (CN), called compute & storage I/O node (C & SION), has SSD about 1.6 TB
    - Its services
      - Cache for the 2nd layer file system
      - Temporary file systems
        - ✓ Local file system for CNs
        - ✓ Shared file system for a job
  - 2nd layer (~150 PB, multiple volumes)
    - Fujitsu FEFS: Lustre-based file system
  - 3rd layer
    - Cloud storage service



DDN SFA18KE & SFA18KXE

# Three-level hierarchical storage system

# LLIO: Lightweight Layered I/O Accelerator

- **Cooperative operations with the 2nd layer storage system**

I/O requests

| CN | CN | CN | CN | ∎∎∎ |

C & SION

C & GION

SSD

Background
data flushing (cache mode)

I/O requests

| CN | CN | CN | CN | ∎∎∎ |

C & SION

C & GION

SSD

I/O Network

2nd Layer
Storage
System

1. Cache for the 2nd layer storage system
2. Temporary local file system for CNs
3. Temporary shared file system for a Job

3rd Layer
Storage System

TofuD

InfiniBand (EDR)

# 2nd Layer Storage System

- **Requirements for the 2nd layer storage system of "Fugaku"**
  1. High capacity
  2. High redundancy
  3. High performance

- **FEFS: Lustre-based file system provided from FUJITSU LIMITED**
  - Many experiences and fruitful knowledge through the K computer operation (~8 years) with the FEFS based on Lustre ver. 1.8

- **Installation of FEFS based on Lustre ver.2.10 with enhancements by FUJITSU LIMITED for the 2nd layer storage system of "Fugaku"**
  - RAS (e.g., High availability)
  - QoS
  - Optimized I/O performance
  - Storage management, etc.

Optimizations and parameter setting are in progress.

# I/O nodes and Interconnects

- **I/O nodes and interconnects associated with the 2nd layer storage system**
  - "C & SION", "C & GION"
  - TofuD among "C & SION", "C & GION", and "C & BION"
  - InfiniBand among "C & GION" and the 2nd layer storage system

- **Activities of interconnects and I/O nodes impact performance of the storage system**
  - Monitoring activities of those components with I/O performance/metrics of the storage system would be useful according to our experience at the K computer.
    - Y. Tsujita, "Characterizing I/O Optimization Effect Through Holistic Log Data Analysis of Parallel File Systems and Interconnects," LNCS 12321, pp. 177–190 (2020)

# Monitoring and log collection

# Monitoring and Log Collection

- **Monitoring and log collection of "Fugaku" (in progress)** *
  - Log and metric collection
    - Log collection
      - Logstash/Filebeat
    - Metric collection
      - Prometheus

  Towards stable operation including the storage system

  - Monitoring/alerting and analysis
    - Database
      - Elasticsearch, PostgreSQL
    - Monitoring/alerting
      - Prometheus
    - Visualization
      - kibana, redash, Grafana

  - Node metrics of MDS, MGS, OSS by *node_exporter*
    CPU, memory, disk, network, ⋯
  - Lustre(FEFS) metrics by *lustre_exporter* **
    Bandwidth, IOPS, Stats, ⋯
  - and, more ⋯

* K. Yamamoto, "Operational Data Processing Pipeline," BoF: Operational Data Analytics@SC'19
https://eehpcwg.llnl.gov/conf_sc19.html
** With some enhancements for "Fugaku" configuration

# Monitoring for *dd* write at the 2nd layer storage

- ● *dd*'s write time monitoring of OSTs of each volume (preliminary)
  - ● Periodical monitoring of write times using *dd* on every OST
    - ● Quick investigation of slow OSTs in each volume
    - ● Such approach effectively leads to further investigation about heavy I/O by jobs, system problems, ⋯
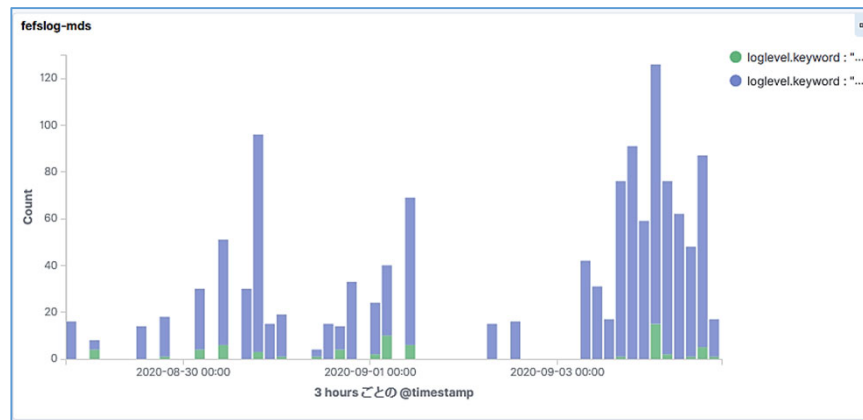
I/O write times over time with indicating by colors
for several performance groups
(Some of OSTs were slow under some I/O
workload stress test.)

I/O BW of each OST over time before/after I/O subsystem
problem indicated by vertical red colored bar
(* Focusing the slowest OST indicated by orange plots)

# Elasticsearch for FEFS log

- **Kibana for Elasticsearch visualization (preliminary)**
  - Quick trouble shooting from a large collection of log data
  - Arrangement in "Fugaku" operation is in progress based on our experiences at the K computer.



Example of evict events generated by MDS (includes both WARN and ERR levels)

# Summary

- **Three-level hierarchical storage system has been introduced at the supercomputer Fugaku.**
- **The 1st layer storage system plays three roles in cooperation with the 2nd layer storage system.**
- **Lustre-based file system (FEFS) developed by FUJITSU LIMITED has been deployed at the 2nd layer storage system based on our experiences at the K computer.**
  - Many enhancements to cope with numerous demands in I/O operations are expected to play important roles at the supercomputer Fugaku.
- **Monitoring activities of I/O nodes and interconnects would be also important aspect at the supercomputer Fugaku based on our experiences at the K computer.**
- **Monitoring/log collection environment is in progress towards stable storage system operation.**
  - Alerting failures and finding root-causes
  - Finding performance bottlenecks and further optimizations, and more⋯