

Status of Lustre-Based Filesystem at the Supercomputer Fugaku

Yuichi Tsujita

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

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

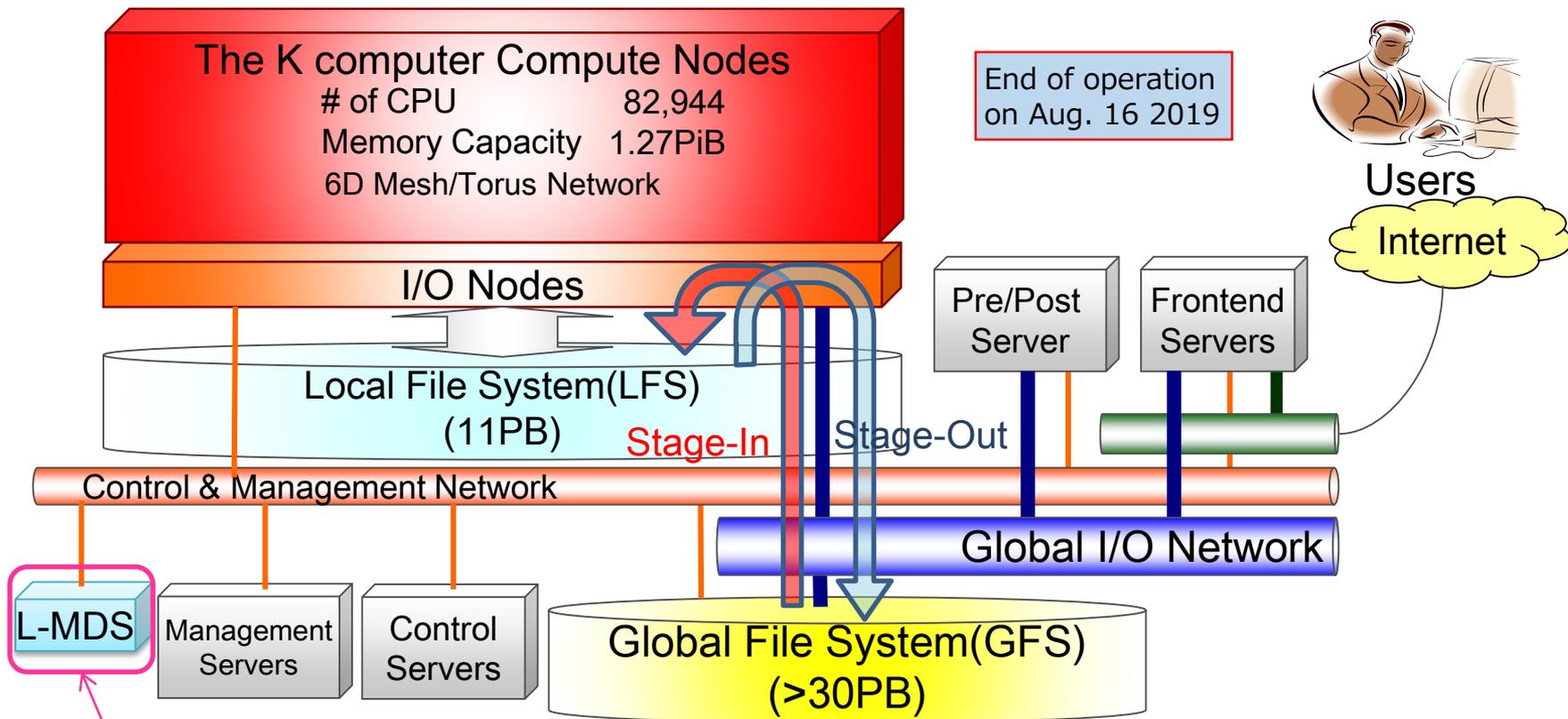
Notable Features of FEFS

- **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.**

Filesystem at the K computer

File system at the K computer

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

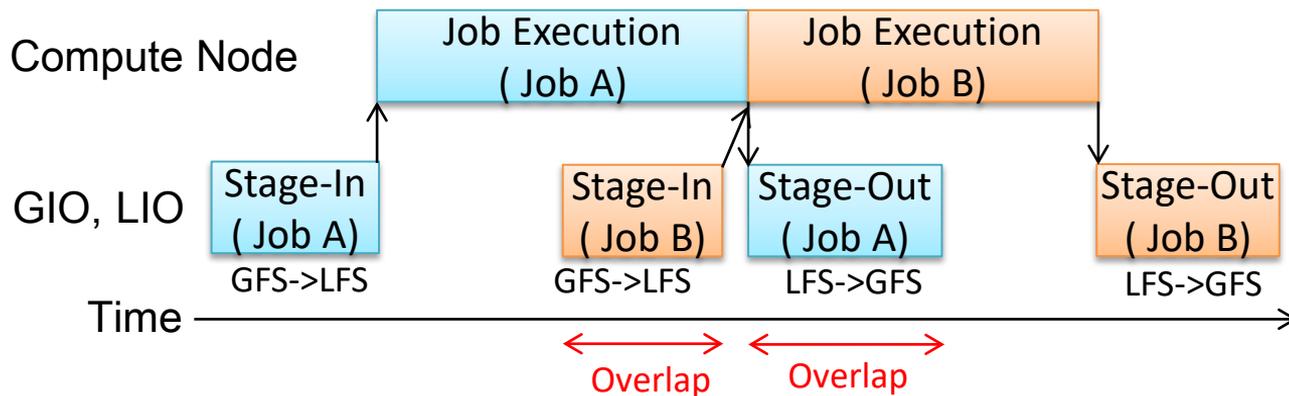


MDS for LFS

FEFS was used for both LFS and GFS.

(FEFS: Fujitsu Exabyte File System based on Lustre technology)

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



- 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

● Performance targets

- 100 times faster than K for some applications (tuning included)
- 30 to 40 MW power consumption

▣ Predicted Performance of 9 Target Applications *As of 2019/05/14*

Area	Priority Issue	Performance Speedup over K	Application	Brief description
Health and Longevity	1. Innovative computing infrastructure for drug discovery	125x +	GENESIS	MD for proteins
	2. Personalized and preventive medicine using big data	8x +	Genomon	Genome processing (Genome alignment)
Disaster Prevention and Environment	3. Integrated simulation systems induced by earthquake and tsunami	45x +	GAMERA	Earthquake simulator (FEM in unstructured & structured grid)
	4. Meteorological and global environmental prediction using big data	120x +	NICAM+ LETKF	Weather prediction system using Big data (structured grid stencil & ensemble Kalman filter)
Energy issue	5. New technologies for energy creation, conversion / storage, and use	40x +	NTChem	Molecular electronic simulation (structure calculation)
	6. Accelerated development of innovative clean energy systems	35x +	Adventure	Computational Mechanics System for Large Scale Analysis and Design (unstructured grid)
Industrial competitiveness enhancement	7. Creation of new functional devices and high-performance materials	30x +	RSDFT	Ab-initio simulation (density functional theory)
	8. Development of innovative design and production processes	25x +	FFB	Large Eddy Simulation (unstructured grid)
Basic science	9. Elucidation of the fundamental laws and evolution of the universe	25x +	LQCD	Lattice QCD simulation (structured grid Monte Carlo)

<https://postk-web.r-ccs.riken.jp/perf.html>

High demands for

- not only computing performance
- but also **storage performance**

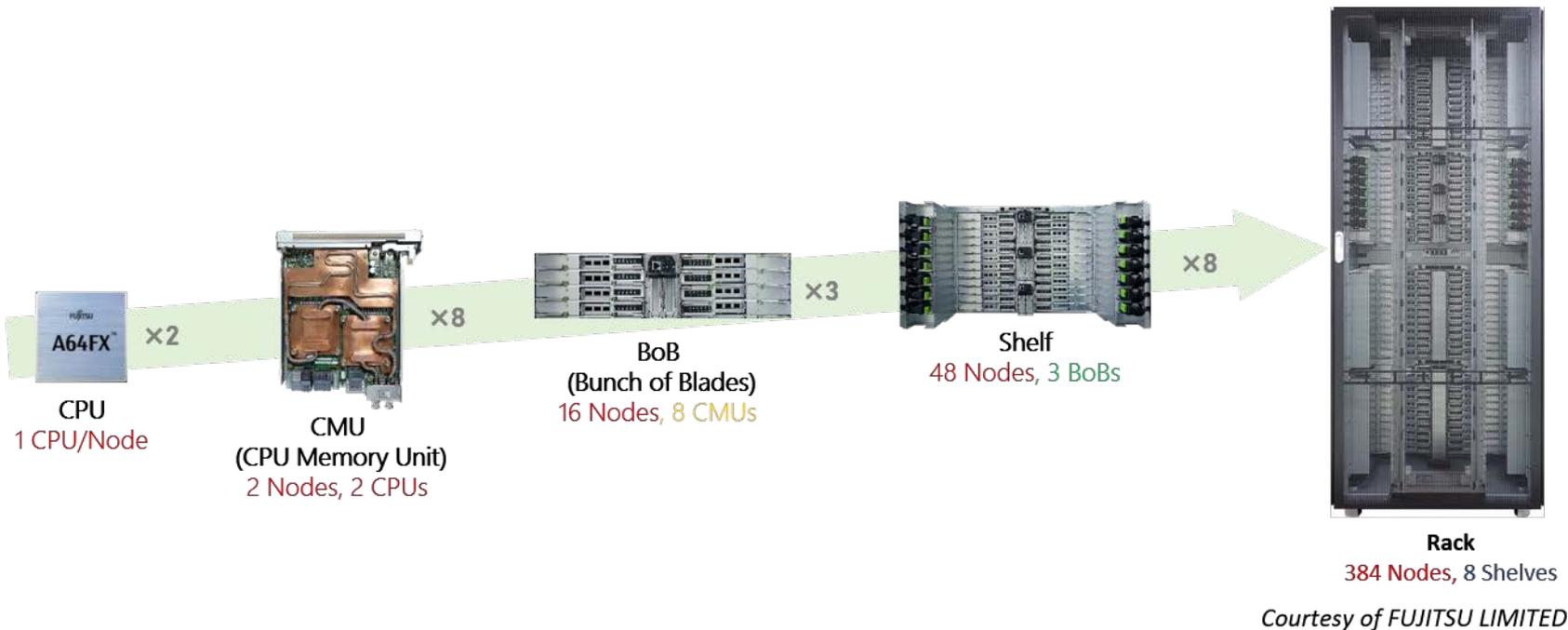
From “K” to “Fugaku”

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

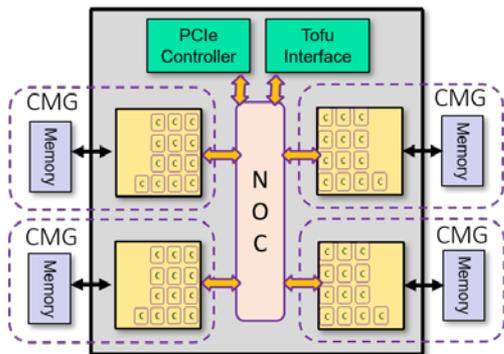
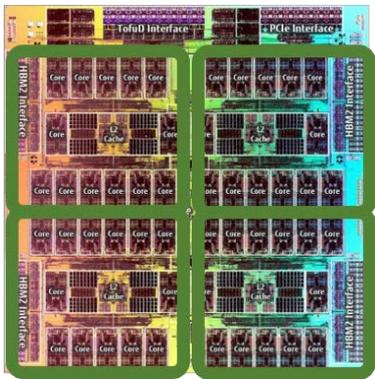
		“Fugaku”	“K”
CPU Architecture		A64FX Arm v8.2-A SVE (512 bit SIMD)	SPARC64VIIIfx
Node	Cores	48	8
	Peak DP performance	2.7+ TF	0.128 TF
	Main Memory	32 GiB	16 GiB
	Peak Memory Bandwidth	1,024 GB/s	64 GB/s
	Peak Network Performance	40.8 GB/s	20 GB/s
Rack	Nodes	384	102
	Peak DP Performance	1+ PF	< 0.013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

Hardware Configuration of "Fugaku"

- From CPU to Rack



● CPU: A64FX Architecture

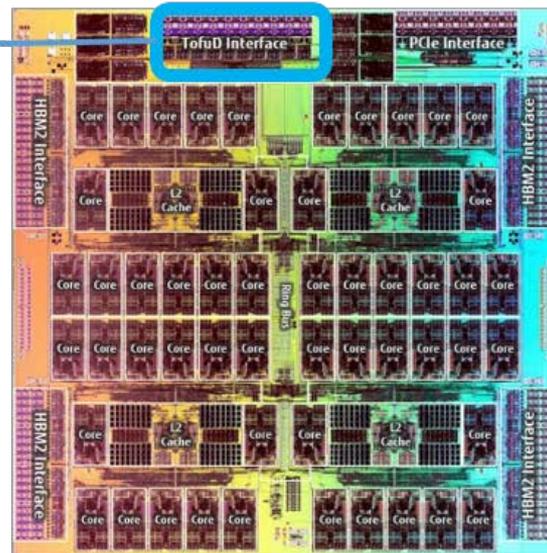
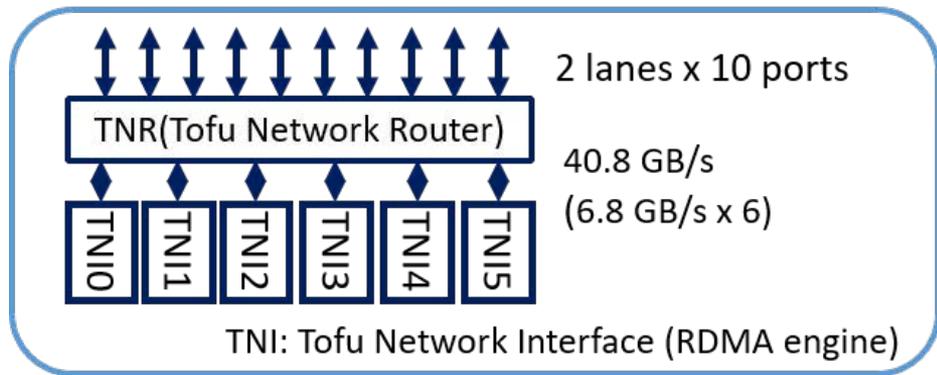


* CMG (Core Memory Group)

<https://github.com/fujitsu/A64FX>

Architecture	Armv8.2-A SVE (512 bit SIMD)	
Core	48 cores for compute and 2/4 for OS activities <i>e.g., I/O</i>	
	Normal: 2.0 GHz	DP: 3.072 TF, SP: 6.144 TF, HP: 12.288 TF
	Boost: 2.2 GHz	DP: 3.3792 TF, SP: 6.7584 TF, HP: 13.5168 TF
Cache L1	64 KiB, 4 way, 230+ GB/s (load), 115+ GB/s (store)	
Cache L2	CMG(NUMA): 8 MiB, 16 way Node: 3.6+ TB/s Core: 115+ GB/s (load), 57+ GB/s (store)	
Memory	HBM2 32 GiB, 1024 GB/s	
Interconnect	TofuD (28 Gbps x 2 lane x 10port)	
I/O	PCIe Gen3 x 16 lane	
Technology	7nm FinFET	

- 6D Torus/Mesh Interconnects among nodes



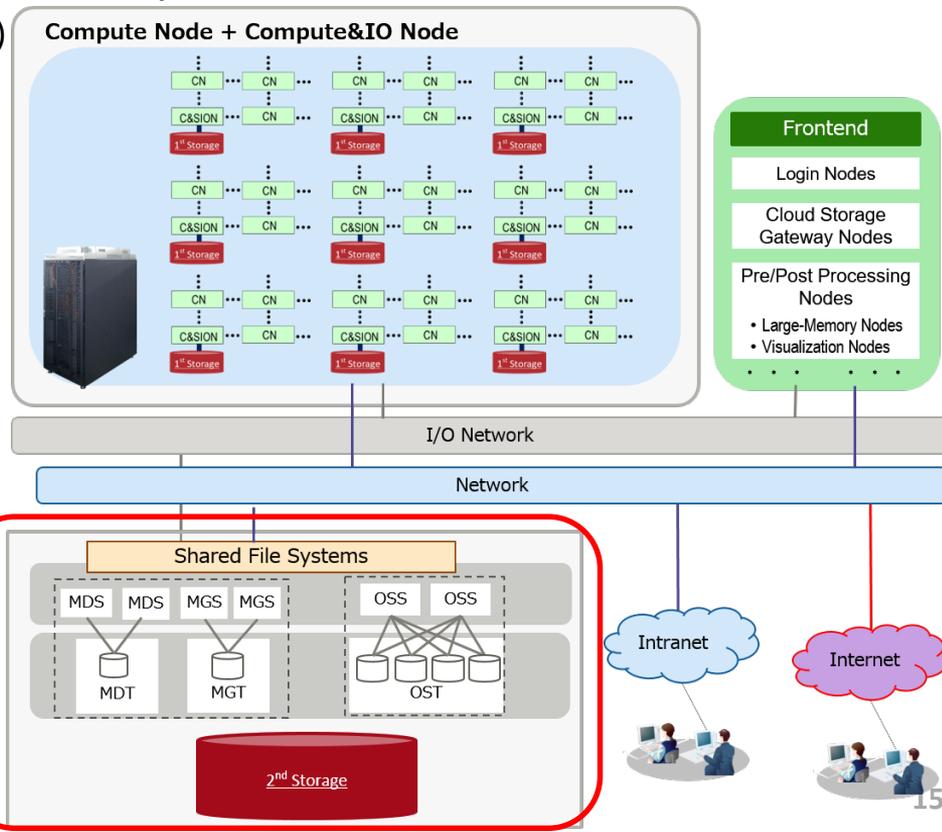
- 6 RDMA engines
- Hardware barrier support
- Network operation offloading capability

Yuichiro Ajima, et al. , “The Tofu Interconnect D,” IEEE Cluster 2018, 2018

- **150k+ 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**

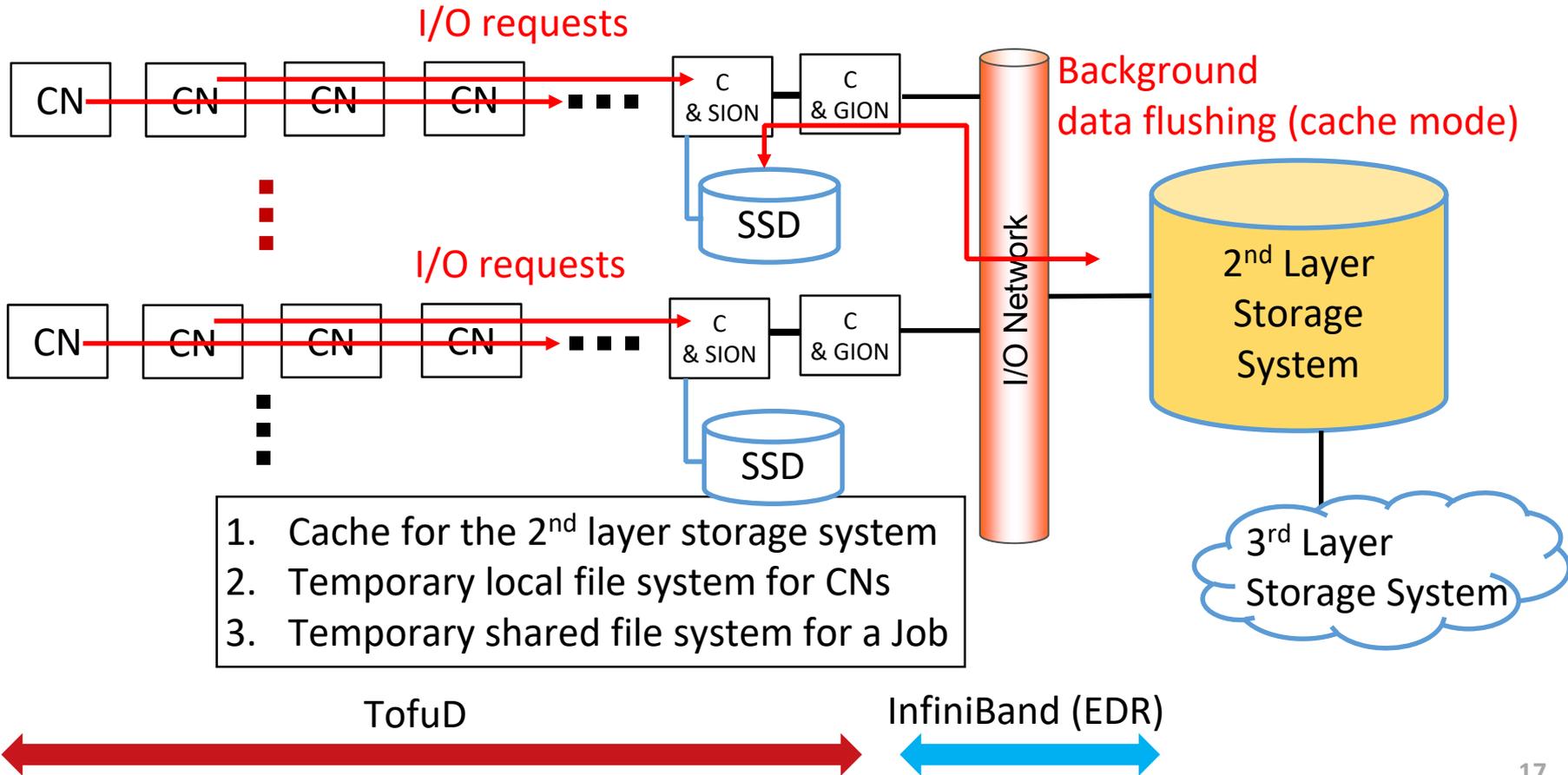
- 1st layer
 - One of 16 compute nodes (CN), called compute & storage I/O node (C & SION), has **SSD** about 1.6 TB
 - Its services
 - Cache for the 2nd layer file system
 - Temporary file systems
 - ✓ Local file system for CNs
 - ✓ Shared file system for a job
- 2nd layer (~150 PB, multiple volumes)
 - **Fujitsu FEFS: Lustre-based file system**
- 3rd layer
 - Cloud storage service



Three-level hierarchical storage system

LLIO: Lightweight Layered I/O Accelerator

- Cooperative operations with the 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 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,” HPC-IODC’20 (<https://hps.vi4io.org/events/2020/iodc>)

Monitoring and log collection

- Monitoring and log collection of “Fugaku” (in progress) *
 - Log and metric collection
 - Log collection
 - Logstash/Filebeat
 - Metric collection
 - Prometheus
 - Monitoring/alerting and analysis
 - Database
 - Elasticsearch, PostgreSQL
 - Monitoring/alerting
 - Prometheus
 - Visualization
 - kibana, redash, Grafana



Towards stable operation including the storage system

- 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 FEFS specific metrics

- ***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 problem, ...

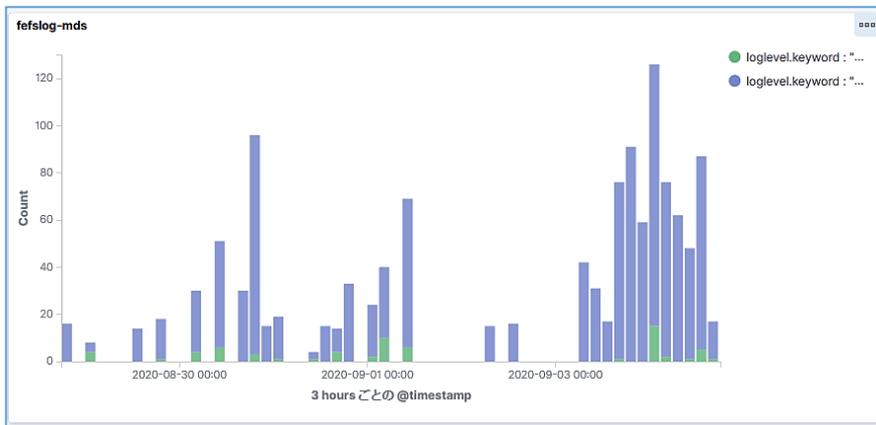


Poor response

Good response

Behavior under some I/O workload stress test
(* Some of OSTs were slow.)

- **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-layer 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.
- **Examinations of 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...