

Key Software Technologies for Future High Performance Computing

Key software issues for future HPC



- Scalability
- Manageability
- Power Efficiency
- Productivity

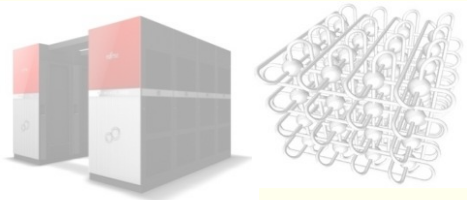
Key software issues for future HPC



- Scalability
- Manageability
- Power Efficiency
- Productivity

Technologies for Scalability

PRIMEHPC Series



Compiler

MPI

Tools

PRIMEHPC OS

Job Manager

System Manager

Mgmt. Nodes

OSS

FEFS

Main Storage

Background:

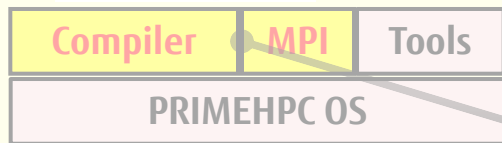
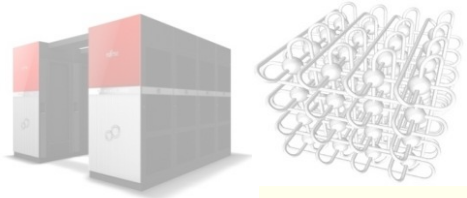
The performance are always top-priority on HPC application. Scalability is the most important for performance on massively parallel computing.

Issues in this area:

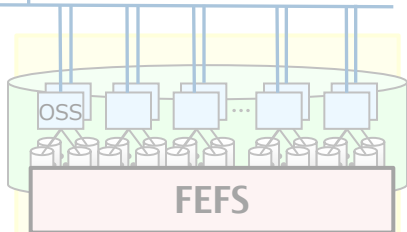
- Improve application parallelization
- Reduce OS Jitter affection
- Keep I/O performance scalability

Technologies for Scalability

PRIMEHPC Series



Mgmt. Nodes

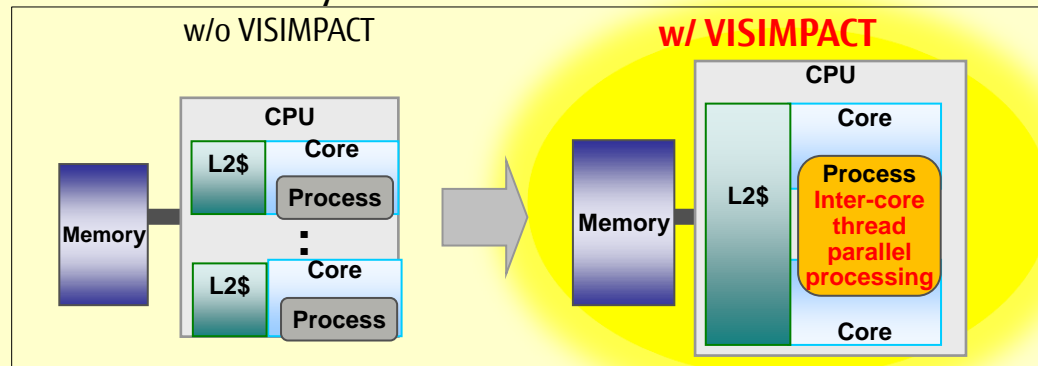


Main Storage

■ Improve application parallelization

Technology: **VISIMPACT + Tofu Optimized MPI**

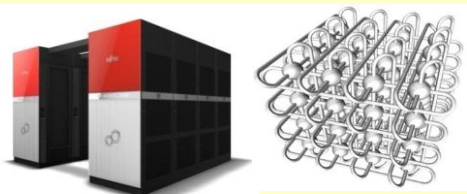
- Easy and efficient inter-core Parallelization automated by VISIMPACT



- Threads-Processes hybrid parallelization with Tofu optimized MPI library.

Technologies for Scalability

PRIMEHPC Series



Compiler

MPI

Tools

PRIMEHPC OS

Job Manager

System Manager

Mgmt. Nodes

FEFS

Main Storage

■ Reduce OS Jitter affection

Technology: Tuned Linux OS for PRIMEHPC

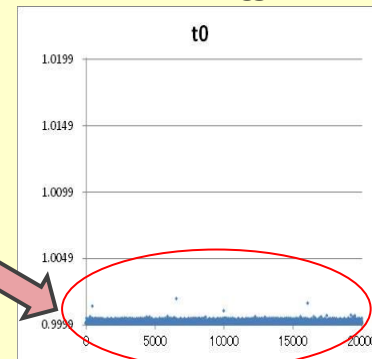
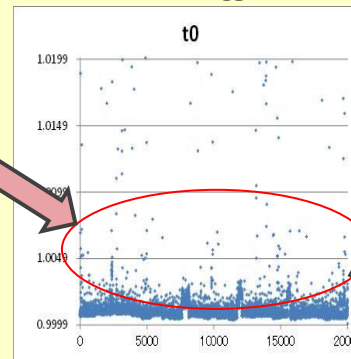
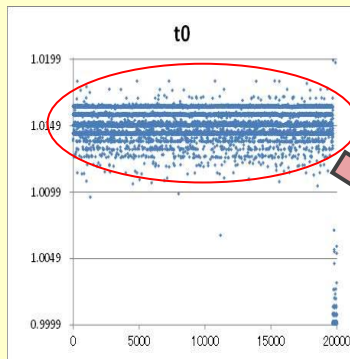
minimized negative effect of OS jitter ultimately by:

- Core-binding technology
- Deliberately selected and tuned system service

x86 cluster
w/o TCS

x86 cluster
w/ TCS

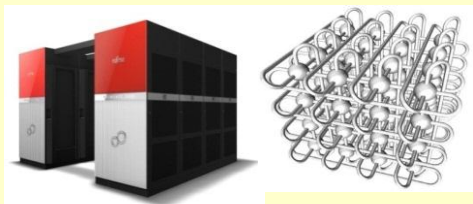
FX10 OS
w/ TCS



TCS: Technical Computing Suite (Fujitsu's System Software Product)

Technologies for Scalability

PRIMEHPC Series



Compiler MPI Tools

PRIMEHPC OS

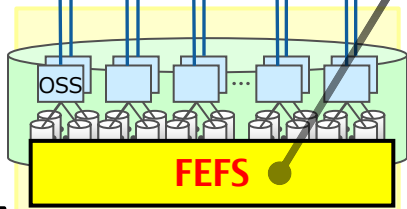


Job Manager



System Manager

Mgmt. Nodes



Main Storage

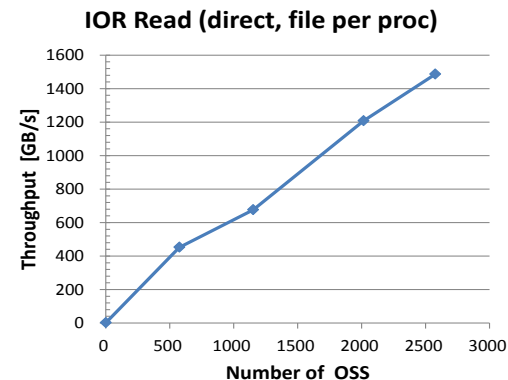
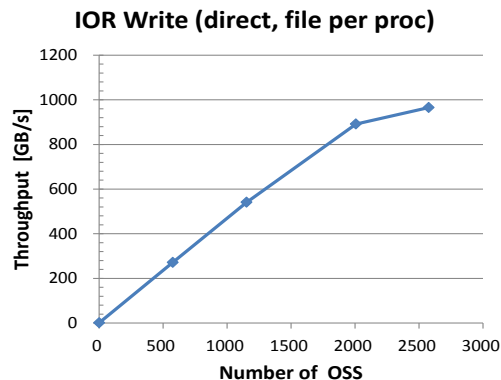
■ Keep I/O performance scalability

Technology: FEFS (Fujitsu Exabyte File System)

Lustre based scalable file system

- Supports up to 8 Exa bytes capacity
- Achieves superb performance on K computer

Write : **965 GB/s** Read : **1,486 GB/s**



Collaborative work with RIKEN

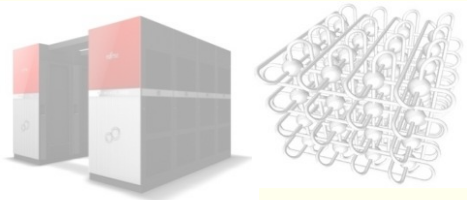
Key software issues for future HPC



- Scalability
- **Manageability**
- Power Efficiency
- Productivity

Technologies for Manageability

PRIMEHPC Series



Compiler

MPI

Tools

PRIMEHPC OS

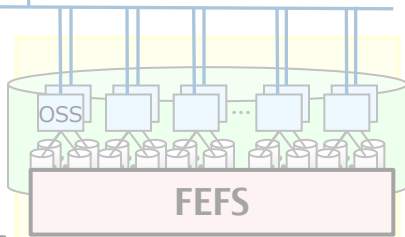


Job Manager



System Manager

Mgmt. Nodes



FEFS

Main Storage

Background:

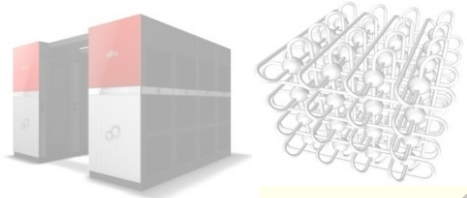
For system administrator, managing system is a key issue. We have already achieved nearly 100,000 nodes system on K computer.

Issues in this area:

- Availability
- Operability

Technologies for Manageability

PRIMEHPC Series



Compiler MPI Tools

PRIMEHPC OS

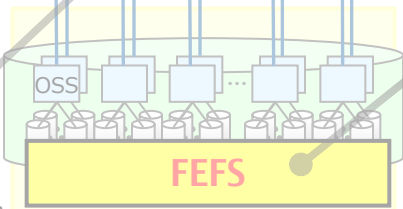


Job Manager



System Manager

Mgmt. Nodes



Main Storage

■ Availability

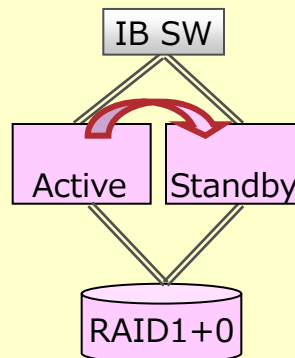
Technology: Automatic Management Node failover

- Immediate failover by Hot Stand-by

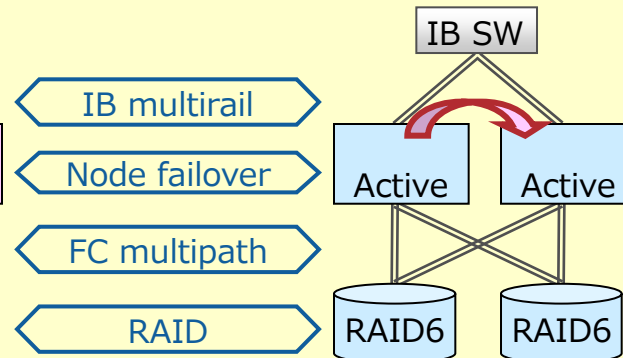
Technology: Automatic MDS/OSS failover

- Monitored by TCS system manager

MDS: Active/Standby

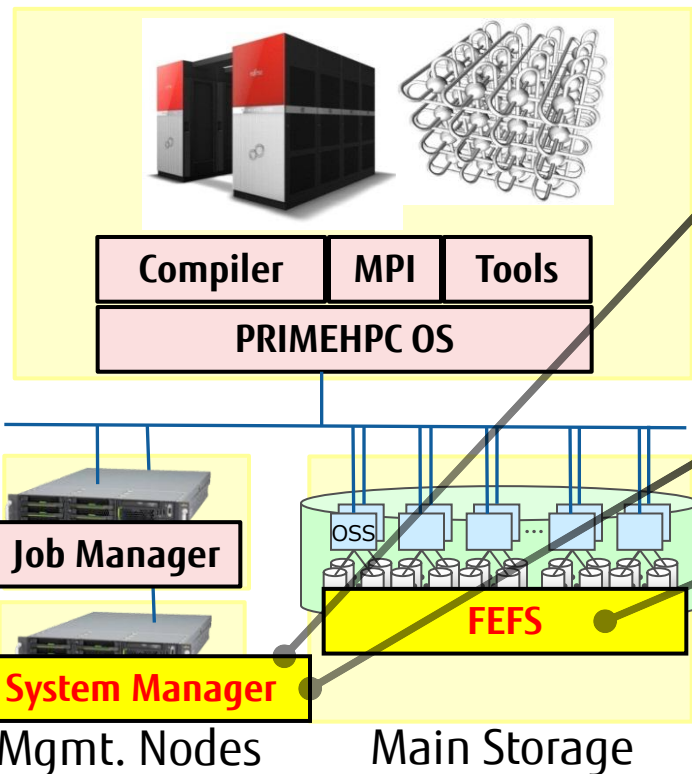


OSS Active/Active



Technologies for Manageability

PRIMEHPC Series



■ Operability

Technology: Centric Management provides single system image for:

- System Installation / Update
- Node status Monitoring (Hardware / Software)
- Power Control / Monitoring
- Support PRIMEHPC / x86 Hybrid Cluster system

Technology: Flexible cluster management provides various physical/logical partitioning.

Technology: QoS/Directory Quota on FEFS facilitates sharing global storage across multi cluster system.

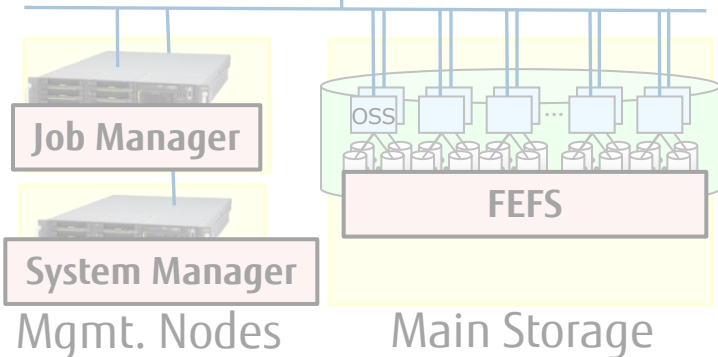
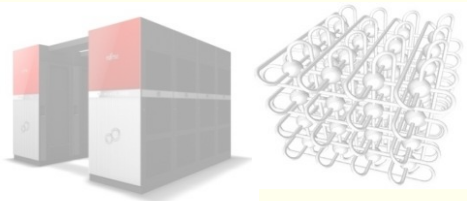
Key software issues for future HPC



- Scalability
- Manageability
- **Power Efficiency**
- Productivity

Technologies for Power Efficiency

PRIMEHPC Series



Background:

$$\text{Power Efficiency} = \frac{\text{"Actual" Throughput}}{\text{Total system power}}$$

Customer Requirement:

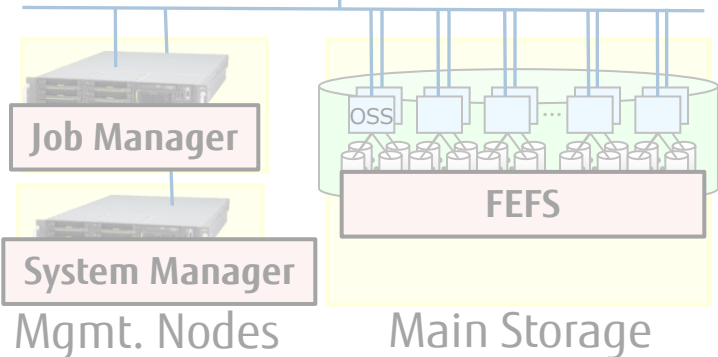
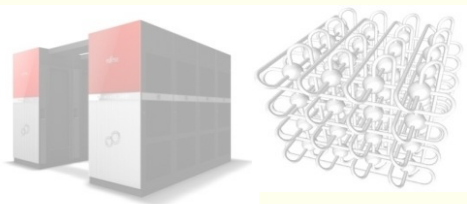
- ✓ Increasing "actual" system throughput
- ✓ Keeping total power at target value

Issues in this area:

- Maximize application efficiency
- Maximize Resource utilization
- System-wide Power Management

Technologies for Power Efficiency

PRIMEHPC Series



■ Maximize application efficiency

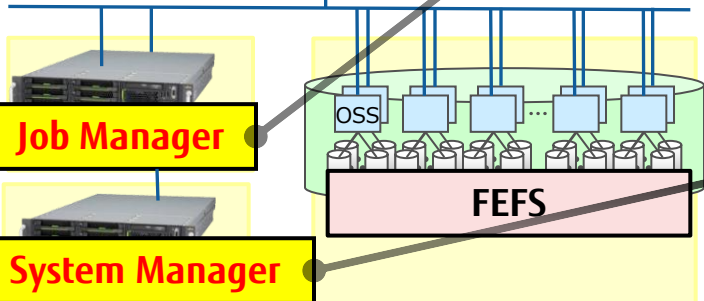
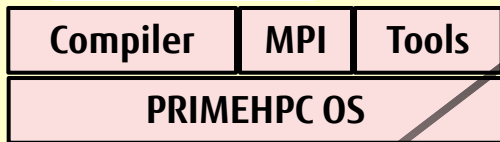
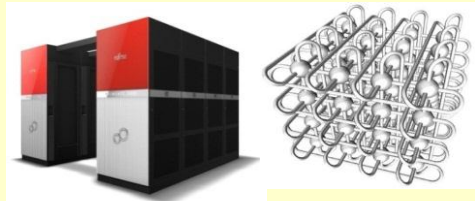
Technology: Optimized OS and languages achieves good efficiency for many applications.

Application	Nodes	Efficiency
LINPACK	88,128	93%
NICAM	81,920	8%
Seism3D	82,944	18%
PHASE	82,944	20%
RSDFT	82,944	52%
FrontFlow/blue	80,000	3%
Lattice QCD	82,944	16%
ZZ-EFSI	82,944	46%

Table is provided by Dr. Minami of RIKEN

Technologies for Power Efficiency

PRIMEHPC Series



Mgmt. Nodes

Main Storage

■ Maximize Resource Utilization

Technology: **Various job allocation** method to increase node/core utilization even on Torus system

- Torus mode/Mesh mode allocation
- Node simplex/share allocation
- Heterogeneous hybrid parallel job allocation

■ System-wide Power Management

Technology: **Centric Power Control** helps to integrate center-wide power capping with:

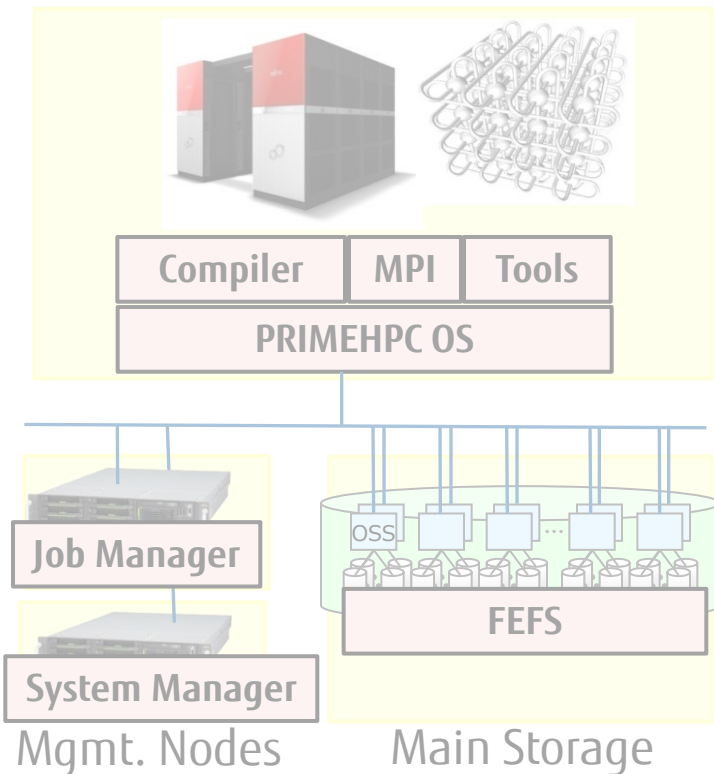
- Interface to control power of nodes or storages
- Power consumption monitoring
- Control power saving mode w/ job manager

Key software issues for future HPC



- Scalability
- Manageability
- Power Efficiency
- **Productivity**

PRIMEHPC Series



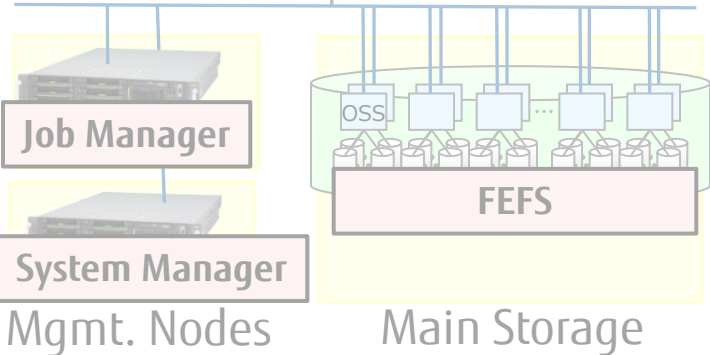
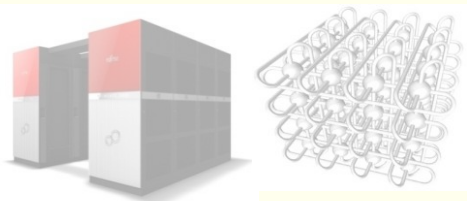
Background:

Because of increasing complexity of node and network architecture, developing applications become more difficult.

Issues in this area:

- **Tuning and Debugging**
- **Portability**

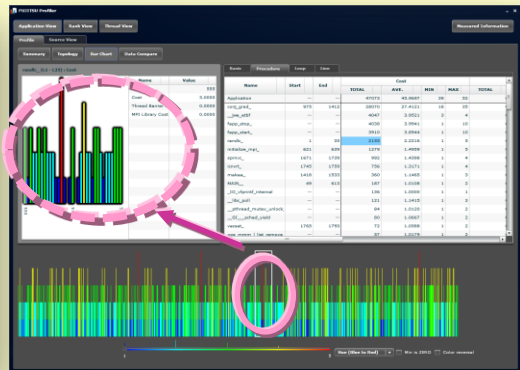
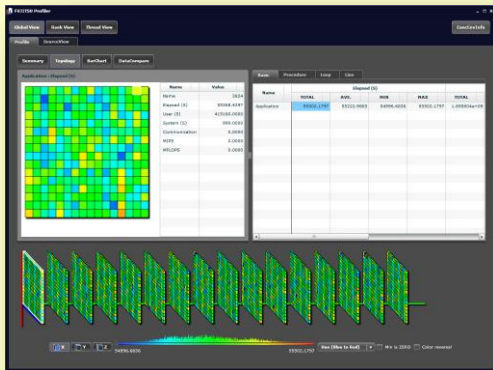
PRIMEHPC Series



Tuning and Debugging

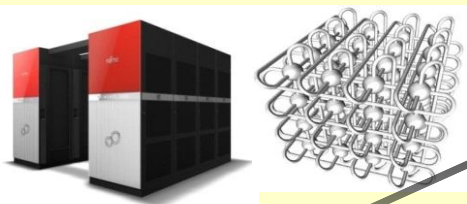
Technology:

- Supports world's standard debugger (DDT)
- Profiler
 - GUI Detailed PA information
 - Optimize communication on Tofu
- Rank Mapping Optimization (RMATT)



Technologies for Productivity

PRIMEHPC Series



Compiler

MPI

Tools

PRIMEHPC OS

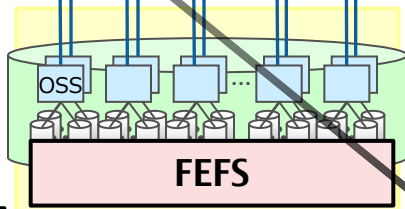


Job Manager



System Manager

Mgmt. Nodes



FEFS

Main Storage

■ Portability

Technology: **Compatibility with K computer**

- Binary compatibility
- Same Architecture applied
1cpu/node, VISIMPACT, Tofu interconnect

Technology: Continues to supports:

- The latest international standards
Fortran 2008, C 11, C++ 11
- De facto standards
GNU C/C++ extensions, OpenMP 4.0, MPI 3.0

Technology: Supports **Generic Linux OS**

- POSIX compliant Linux and generic libraries

Now, we are ready for 100 PFlops!

What's next?

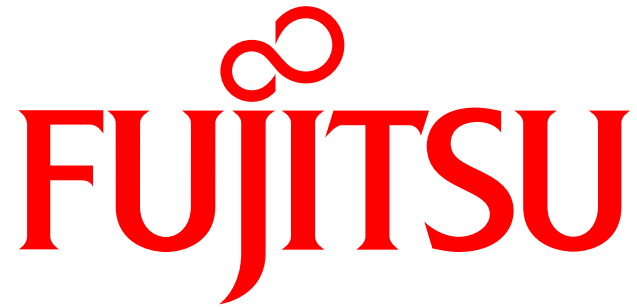
Activities for Exascale Computing

■ Roadmap of Japan's National Project



■ Participate co-design for Exascale System software

- Light-weight Micro kernel next to Linux
- File I/O performance improvement



shaping tomorrow with you