# Advanced Software for the Supercomputer PRIMEHPC FX10

# System Configuration of PRIMEHPC FX10

**Compute nodes**

**6D mesh/torus Interconnect**

• **Data transfer to/from global file system**
• **Data communication for system job operations management**

• **Login**
• **Compilation**
• **Job submission**

**Local file system**
**(Temporary area occupied by jobs)**

IO network (IB), management network (GbE)

Login nodes

**Global file system**
**(Data storage area)**

Job management nodes

File management nodes

Control nodes

System integration node

**Management nodes**

**User**

**Administrator**

• **System operations management**
• **Job operations management**

# System Software Stack

**User/ISV Applications**

**HPC Portal / System Management Portal**

### System operations management
- System configuration management
- System control
- System monitoring
- System installation & operation

### High-performance file system
- Lustre-based distributed file system
- High scalability
- IO bandwidth guarantee
- High reliability & availability

### Compilers
- Hybrid parallel programming
- Sector cache support
- SIMD / Register file extensions

### Job operations management
- Job manager
- Job scheduler
- Resource management
- Parallel execution environment

### VISIMPACT™
- Shared L2 cache on a chip
- Hardware intra-processor synchronization

### Support Tools
- IDE
- Profiler & Tuning tools
- Interactive debugger

### MPI Library
- Scalability of High-Func.
- Barrier Comm.

File system, operations management

Application development environment

**Linux-based enhanced Operating System**
- Enhanced hardware support
- System noise reduction
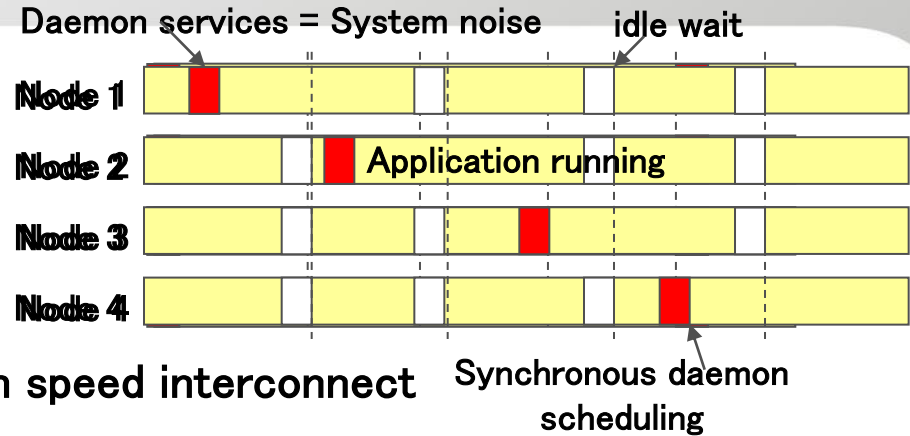- Error detection / Low power

**PRIMEHPC FX10**

# OS (Linux-based enhanced Operating System)

- **Easy existing application porting**
  - POSIX API: Linux kernel 2.6.x, glibc 2.x
- **High performance / High scalability**
  - Enhanced hardware support
    CPU registers, Large memory page, High speed interconnect
  - Reduce system noise in highly parallel program
    Inter-node OS scheduling
- **High availability / low power consumption**
  - Hardware error detection / isolation
    memory patrol, io driver enhance.
  - CPU suspend during system idle state.

Daemon services = System noise    idle wait

Node 1

Node 2   Application running

Node 3

Node 4

Synchronous daemon scheduling

Idle → CPU suspend

Job running

3

# System Software Stack

**FUJITSU**

**User/ISV Applications**

**HPC Portal / System Management Portal**

### System operations management
- System configuration management
- System control
- System monitoring
- System installation & operation

### Job operations management
- Job manager
- Job scheduler
- Resource management
- Parallel execution environment

File system, operations management

### High-performance file system
- Lustre-based distributed file system
- High scalability
- IO bandwidth guarantee
- High reliability & availability

### VISIMPACT™
- Shared L2 cache on a chip
- Hardware intra-processor synchronization

### Compilers
- Hybrid parallel programming
- Sector cache support
- SIMD / Register file extensions

### Support Tools
- IDE
- Profiler & Tuning tools
- Interactive debugger

### MPI Library
- Scalability of High-Func.
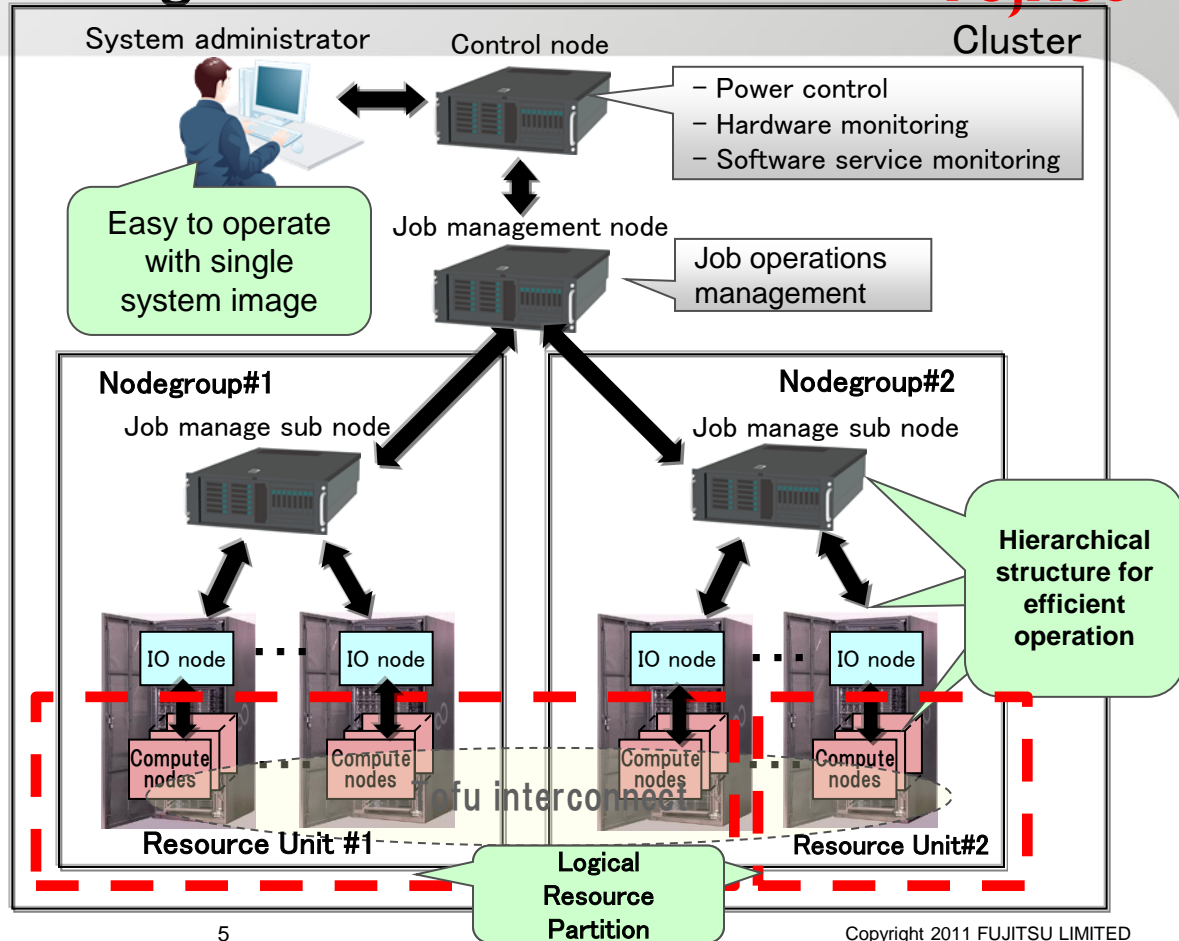- Barrier Comm.

Application development environment

**Linux-based enhanced Operating System**
- Enhanced hardware support
- System noise reduction
- Error detection / Low power

**PRIMEHPC FX10**

# System Operations Management

- Hierarchical structure for efficient system operation and adaptability to large-scale systems
  - The load is distributed by using the job management sub node.
- Easy to operate with a single system image
- The system is efficiently operated by dividing a logical resource partition named "Resource Unit".
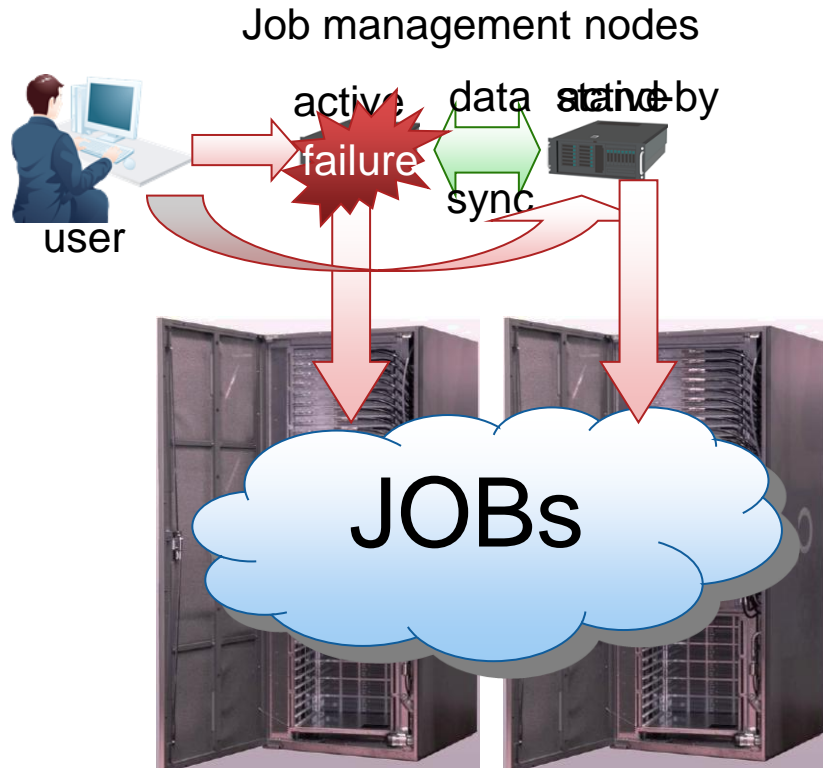
# High Availability System

- **The important nodes have redundancy**
  - Control node
  - Job management node
  - Job management sub node
  - File servers

For example : right figure

- **Continuing job execution even if the job management node is in failed status**
  - The job data always synchronizes between active node and stand-by node.
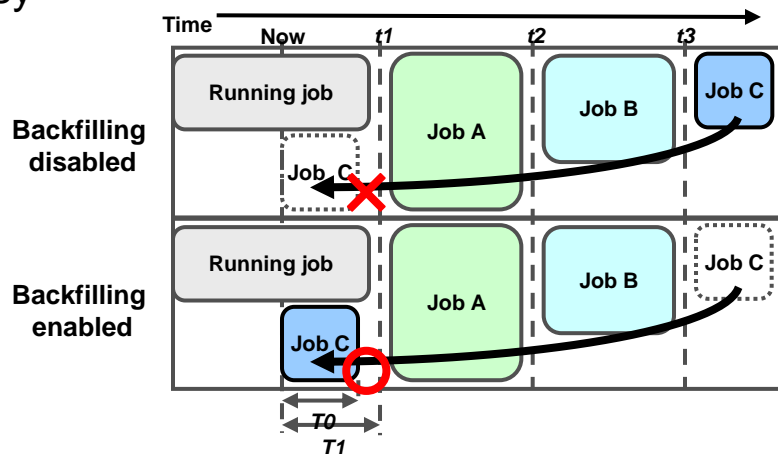  - Alternatively to stand-by node if active node is down.

Job management nodes

active    data    stand-by

failure

sync

user

JOBs

6

# Job Operations Environment

- **Efficient resource usage**
  - Flexible job scheduling based on prioritized resource assignment
  - Interconnect topology-aware resource assignment
  - Backfill scheduling for keeping the resources busy
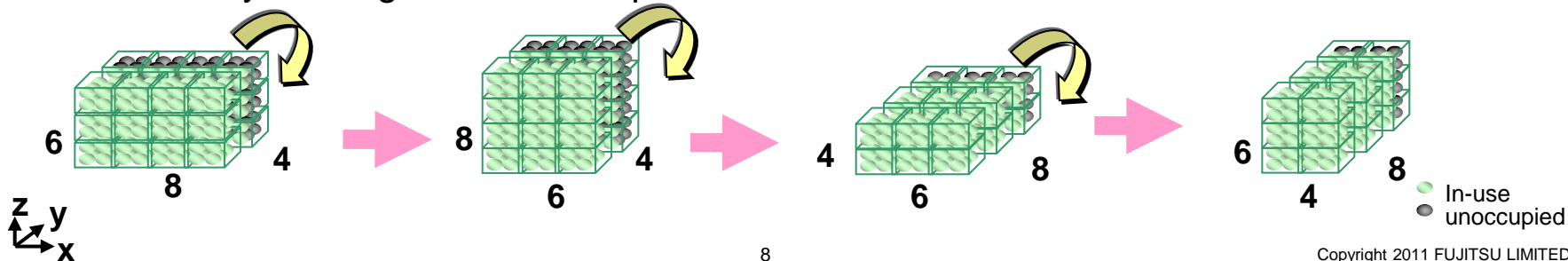  - Asynchronous file staging
- **High availability**
  - Avoids assigning faulty resources to jobs
  - disconnects faulty nodes from job operations
  - Management nodes with failover support

# Resource Assignment

- Interconnect topology-aware resource assignment

    - Treats 12 compute nodes as one interconnect unit

    - Assigns cubic-shaped interconnect unit(s) to a job

    → Using adjacent interconnect unit(s) is suitable for contiguous communication,
      and also avoids interfering with other jobs.

    - Optimizes the alignment of resources

    → Rotating the cubic-shaped interconnect units  This improves total system
      utilization by rotating the cubic shaped interconnect units.



In-use
unoccupied

# System Software Stack

**FUJITSU**

**User/ISV Applications**

**HPC Portal / System Management Portal**

### System operations management
- System configuration management
- System control
- System monitoring
- System installation & operation

### High-performance file system
- Lustre-based distributed file system
- High scalability
- IO bandwidth guarantee
- High reliability & availability

### Compilers
- Hybrid parallel programming
- Sector cache support
- SIMD / Register file extensions

### Support Tools
- IDE
- Profiler & Tuning tools
- Interactive debugger

### Job operations management
- Job manager
- Job scheduler
- Resource management
- Parallel execution environment

### VISIMPACT™
- Shared L2 cache on a chip
- Hardware intra-processor synchronization

### MPI Library
- Scalability of High-Func.
- Barrier Comm.

File system, operations management
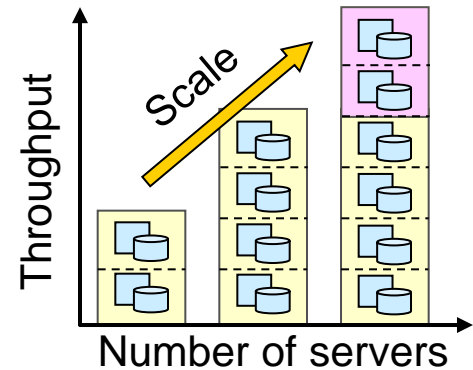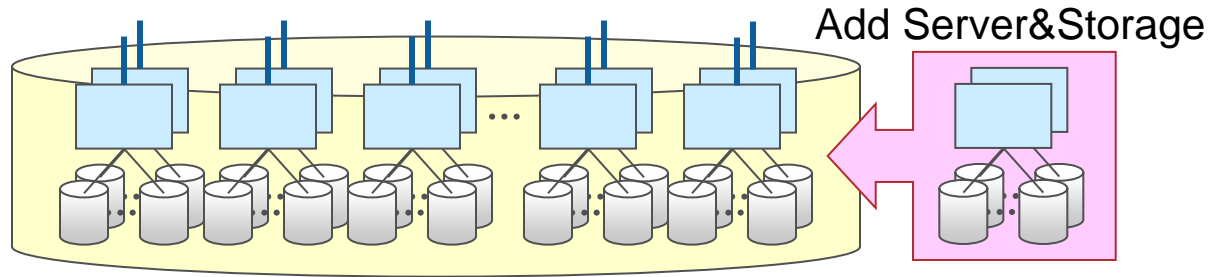
Application development environment

**Linux-based enhanced Operating System**
- Enhanced hardware support
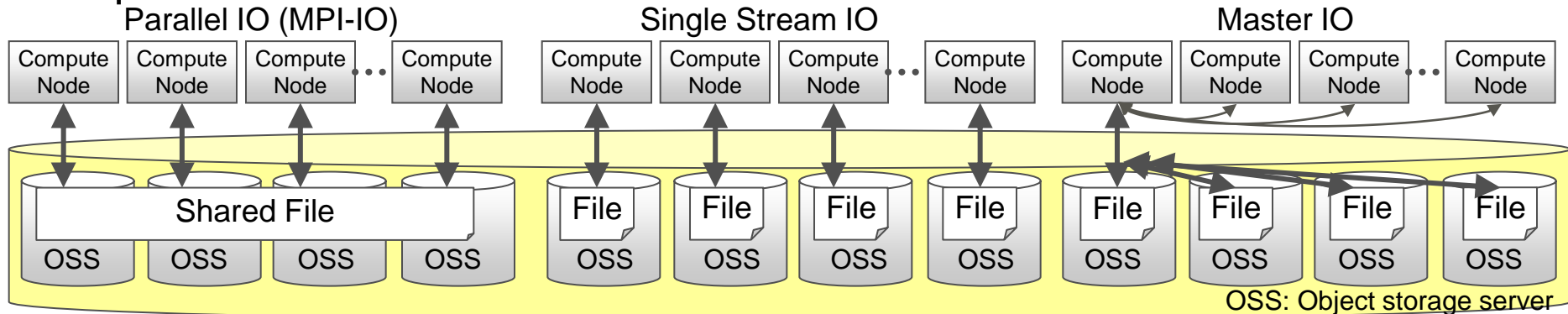- System noise reduction
- Error detection / Low power

**PRIMEHPC FX10**

# High Scalability

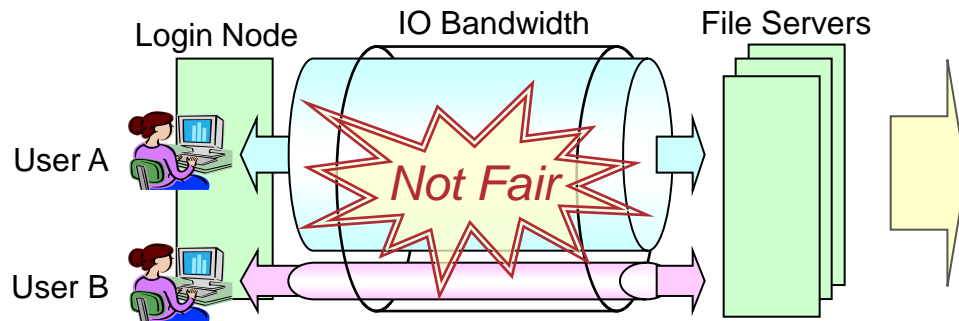■ Achieved high-scalable IO performance with multiple OSSes.

Add Server&Storage



Throughput

Scale

Number of servers

■ Adapted various IO model

Parallel IO (MPI-IO)

| Compute Node | Compute Node | Compute Node | ... | Compute Node |

Shared File

| OSS | OSS | OSS | OSS |

Single Stream IO

| Compute Node | Compute Node | Compute Node | ... | Compute Node |

| File | File | File | File |
| OSS | OSS | OSS | OSS |

Master IO

| Compute Node | Compute Node | Compute Node | ... | Compute Node |

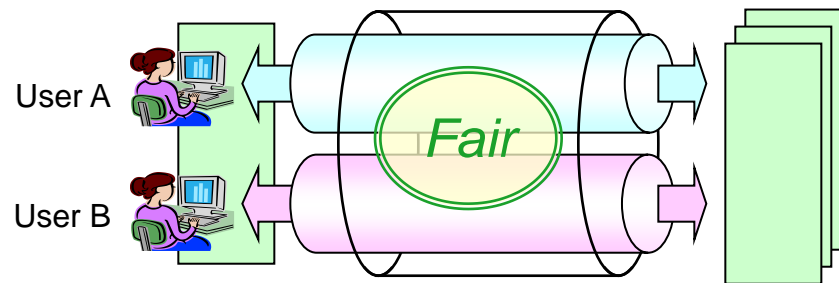| File | File | File | File |
| OSS | OSS | OSS | OSS |

OSS: Object storage server

# IO Bandwidth Guarantee



■ Fair Share QoS: Sharing IO bandwidth with all users.

Without Fair Share QoS

With Fair Share QoS

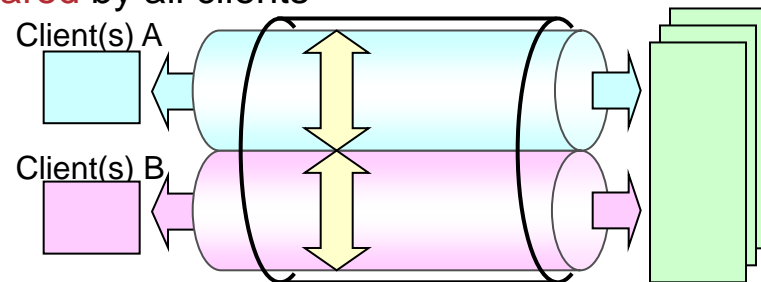Login Node   IO Bandwidth   File Servers

User A

User B

*Not Fair*

*Fair*

User A

User B

■ Best Effort QoS: Utilize all IO bandwidth exhaustively.

Occupied by one client
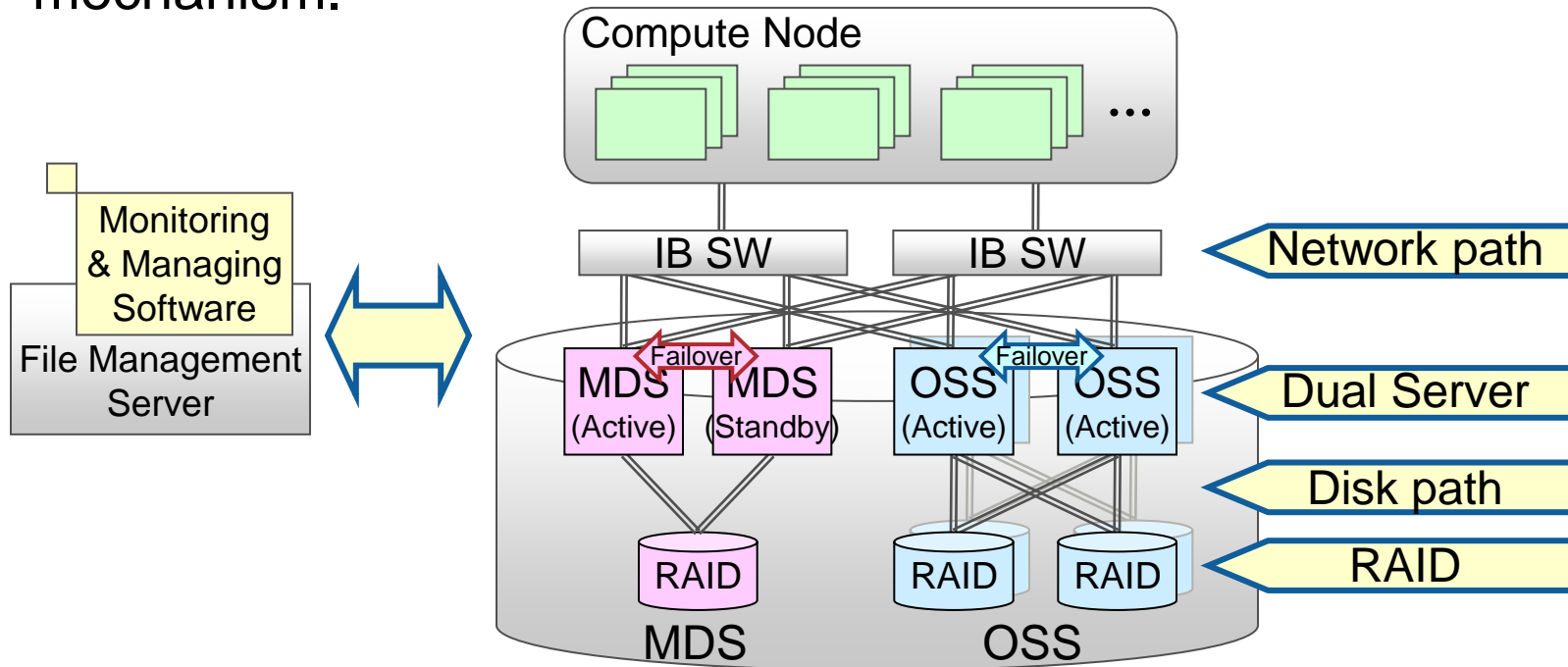
File Servers

Client(s)

Shared by all clients

Client(s) A

Client(s) B

# High Reliability and High Availability

■ Avoiding single point of failure by redundant hardware and failover mechanism.

# System Software Stack

**FUJITSU**

**User/ISV Applications**

**HPC Portal / System Management Portal**

### System operations management
- System configuration management
- System control
- System monitoring
- System installation & operation

### High-performance file system
- Lustre-based distributed file system
- High scalability
- IO bandwidth guarantee
- High reliability & availability

### Compilers
- Hybrid parallel programming
- Sector cache support
- SIMD / Register file extensions

### MPI Library
- Scalability of High-Func.
- Barrier Comm.

### Job operations management
- Job manager
- Job scheduler
- Resource management
- Parallel execution environment

### VISIMPACT™
- Shared L2 cache on a chip
- Hardware intra-processor synchronization

### Support Tools
- IDE
- Profiler & Tuning tools
- Interactive debugger

File system, operations management
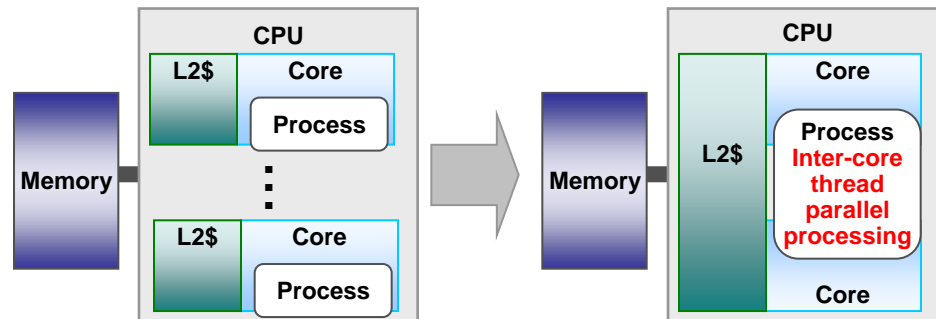
Application development environment

**Linux-based enhanced Operating System**
- Enhanced hardware support
- System noise reduction
- Error detection / Low power

**PRIMEHPC FX10**

# VISIMPACT™
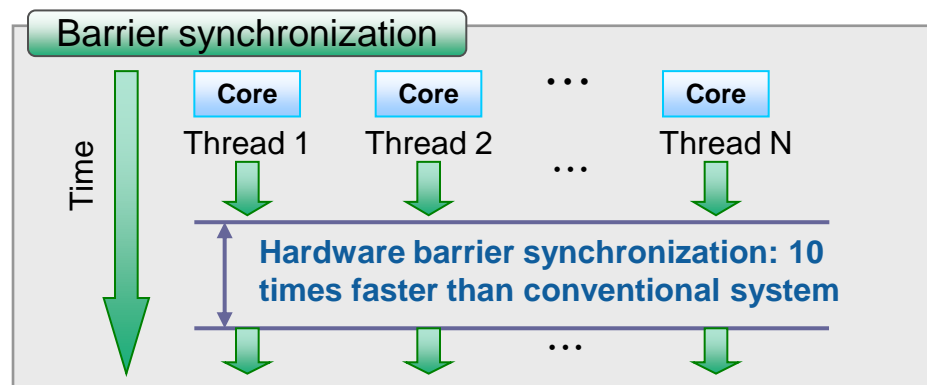(Virtual Single Processor by Integrated Multi-core Parallel Architecture)

- **Mechanism that treats multiple cores as one high-speed CPU**
  - Easy and efficient execution of inter-core thread parallel processing with a multi-core CPU
  - Supports the realization of a highly-efficient Hybrid model (Automatic parallelization + MPI)
- **CPU technologies**
  - Large-capacity shared L2 cache memory decrease in the influence of false sharing
  - Inter-core hardware barrier facilities 6-10 times faster than conventional software barrier

# System Software Stack

**User/ISV Applications**

**HPC Portal / System Management Portal**

## System operations management
- System configuration management
- System control
- System monitoring
- System installation & operation

## Job operations management
- Job manager
- Job scheduler
- Resource management
- Parallel execution environment

## High-performance file system
- Lustre-based distributed file system
- High scalability
- IO bandwidth guarantee
- High reliability & availability

## VISIMPACT™
- Shared L2 cache on a chip
- Hardware intra-processor synchronization

File system, operations management

## Compilers
- Hybrid parallel programming
- Sector cache support
- SIMD / Register file extensions

## MPI Library
- Scalability of High-Func.
- Barrier Comm.

## Support Tools
- IDE
- Profiler & Tuning tools
- Interactive debugger

Application development environment

**Linux-based enhanced Operating System**
- Enhanced hardware support
- System noise reduction
- Error detection / Low power

**PRIMEHPC FX10**

# Programming Model for High Scalability
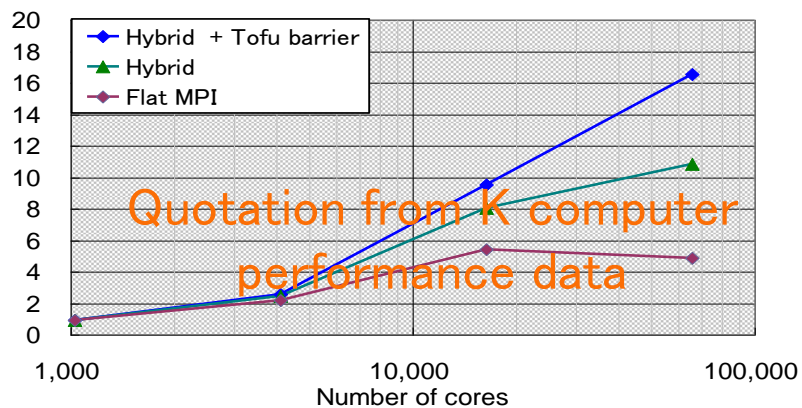
Hybrid parallelism by VISIMPACT and MPI library

- VISIMPACT
  - Automated multi-thread parallelization
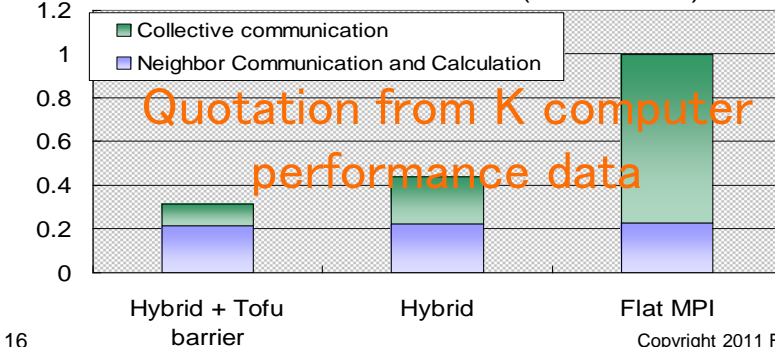  - High performance thread barrier used Inter-core hardware barrier facility
- MPI library
  - High performance collective communications used Tofu barrier facility

### Scalability of Himeno benchmark(XL size)

Performance ratio

- Hybrid + Tofu barrier
- Hybrid
- Flat MPI

Quotation from K computer performance data

Number of cores

### Himeno benchmark detail (65536 Core)

Time Ratio

- Collective communication
- Neighbor Communication and Calculation

Quotation from K computer performance data

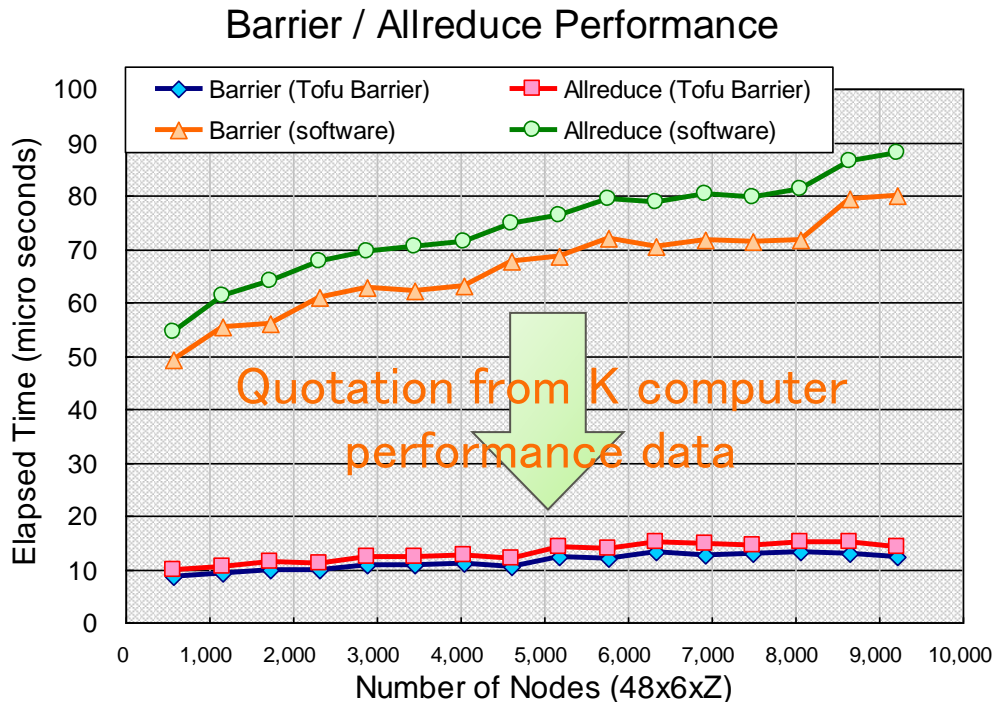Hybrid + Tofu barrier　　Hybrid　　Flat MPI

16

# Customized MPI Library for High Scalability

- ■ **Point-to-Point communication**
  - • Use a special type of low-latency path that bypasses the software layer
  - • The transfer method optimization according to the data length, process location and number of hops
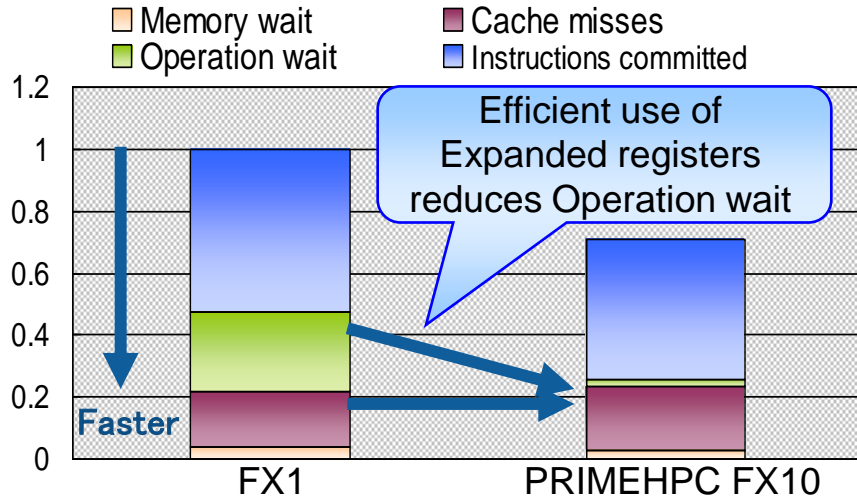- ■ **Collective communication**
  - • High performance Barrier, Allreduce, Bcast and Reduce used Tofu barrier facility
  - • Scalable Bcast, Allgather, Allgatherv, Allreduce and Alltoall algorithm optimized for Tofu network
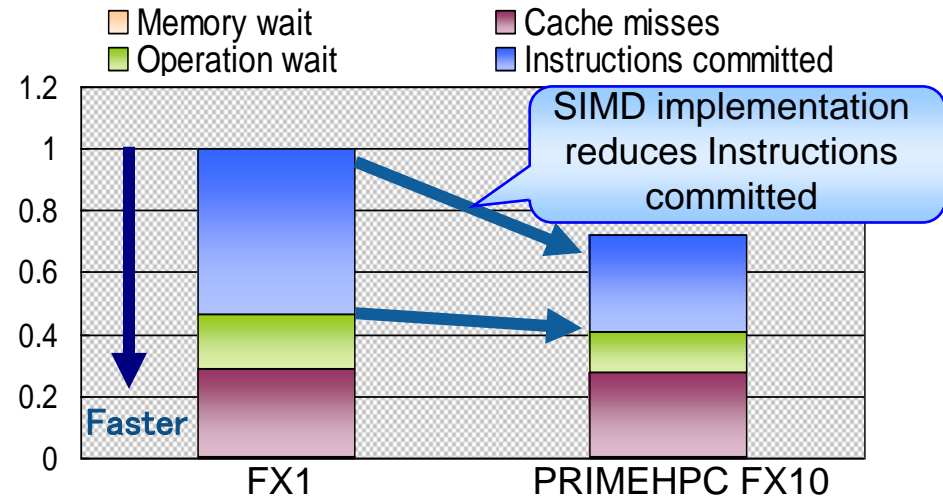
### Barrier / Allreduce Performance



Legend:
- Barrier (Tofu Barrier)
- Allreduce (Tofu Barrier)
- Barrier (software)
- Allreduce (software)

Y-axis: Elapsed Time (micro seconds)
X-axis: Number of Nodes (48x6xZ)

Quotation from K computer performance data

# Compiler Optimization for High Performance

**FUJITSU**

- Instruction-level parallelism with SIMD instructions

- Improvement of computing efficiency used Expanded registers
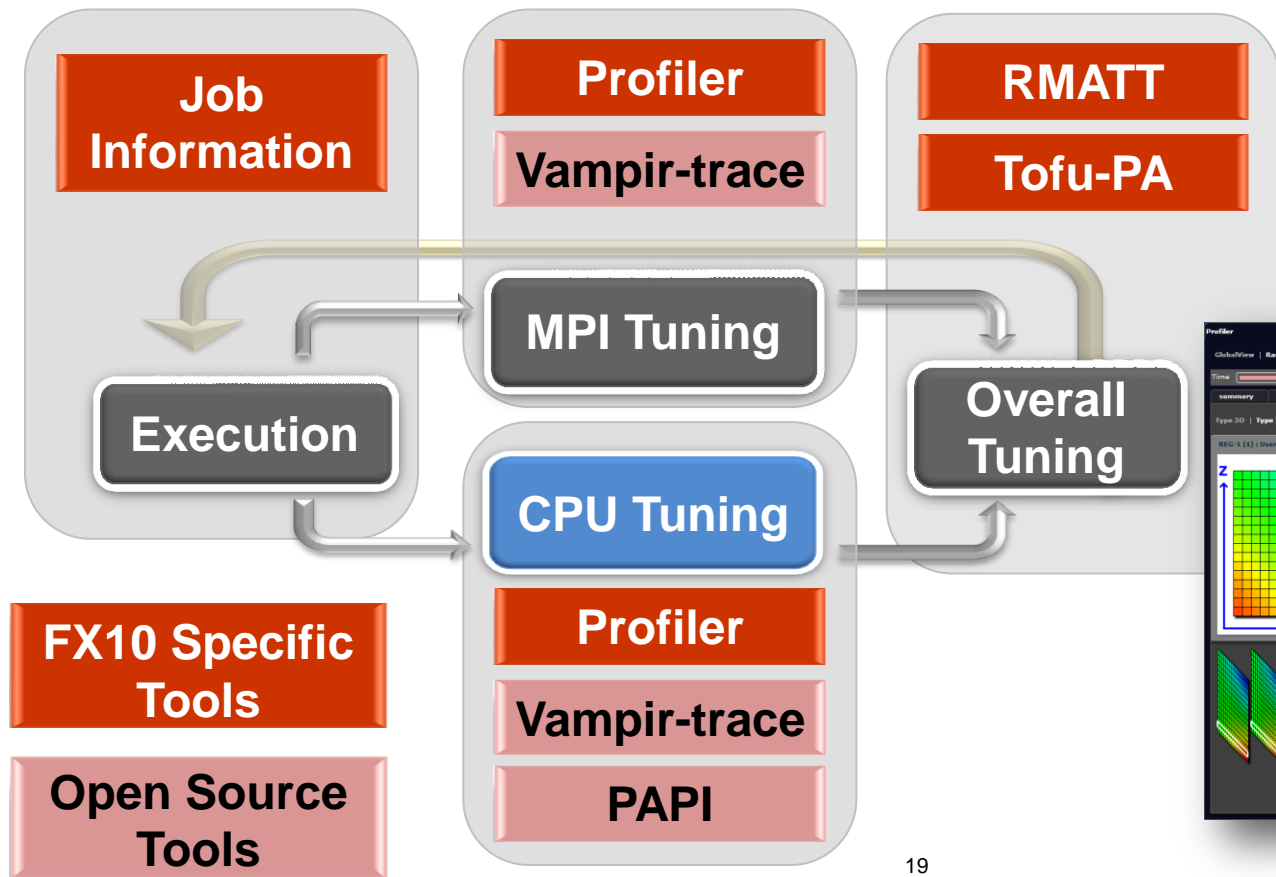
- Improvement of cache efficiency used Sector cache

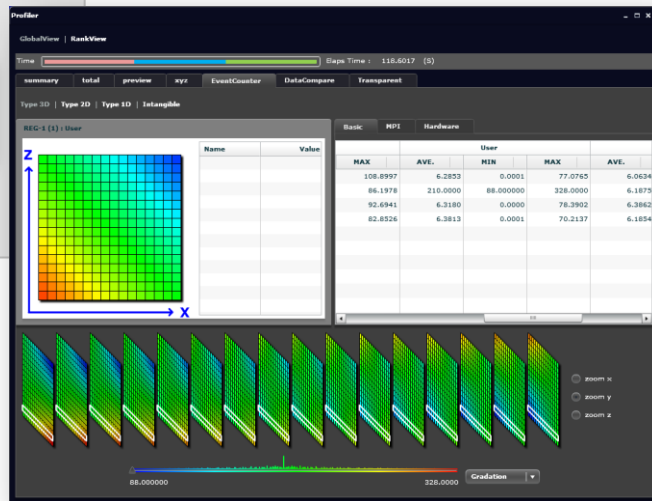

NPB3.3 LU
Execution time comparison (relative values)

Legend: Memory wait, Cache misses, Operation wait, Instructions committed

Efficient use of Expanded registers reduces Operation wait

NPB3.3 MG
Execution time comparison (relative values)

Legend: Memory wait, Cache misses, Operation wait, Instructions committed

SIMD implementation reduces Instructions committed

Faster

FX1    PRIMEHPC FX10

Faster

FX1    PRIMEHPC FX10

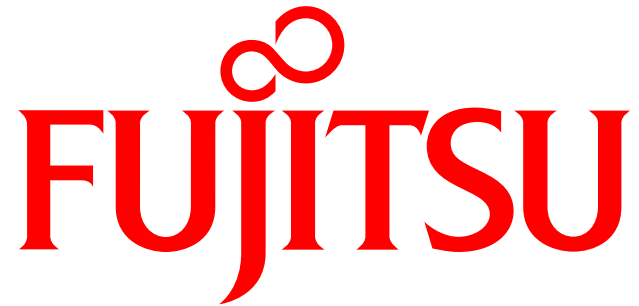# Application Tuning Cycle and Tools



Profiler snapshot

# FUJITSU

shaping tomorrow with you