

## Fujitsu's Technologies Leading to Practical Petascale Computing: *K computer*, *PRIMEHPC FX10* and the Future

November 16th, 2011

Motoi Okuda Technical Computing Solution Unit Fujitsu Limited

## Agenda



### Achievements of *K computer*

- What did K computer achieve?
- Technologies which realize the achievement
- Fujitsu's new Petascale supercomputer PRIMEHPC FX10
  - Fujitsu supercomputer past and present
  - Second generation Petascale supercomputer PRIMEHPC FX10
  - Conclusion

## Latest World Supercomputer Ranking

K computer takes consecutive No.1 on TOP500 List

# 38<sup>th</sup> List: TOP10

SUPERCOMPUTER SITES

TOP 500

<b>*</b>	Rank	Site	Manufacture		Computer	Country	Cores	Rmax [Pflops]	Power [MW]
K compute	1	RIKEN Advanced Institute for Computational Science (AICS)	Fujitsu	K computer SPARC64 VIIIfx 2.0GHz,Tofu interconnect		Japan	705,024	10.510	12.66
	2	National Supercomputing Center in Tianjin	NUDT	Tianhe-1ANUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050JaguarCray XT5-HE Opteron 6-core 2.6 GHz		China	186,368	2.566	4.04
	3	DOE/SC/Oak Ridge National Laboratory	Cray Inc.			USA	224,162	1.759	6.95
	4	National Supercomputing Centre in Shenzhen (NSCS)	Dawning	Dawnin       Image: Construction of the second			58		
	5	GSIC Center, Tokyo Institute of Technology	NEC/HP				10		
	6	DOE/NNSA/LANL/SNL	Cray Inc				ente prico	), 2(000, Ju	)8
	7	NASA/Ames Research Center/NAS	SGI				ters	10	
	8	DOE/SC/LBNL/NERSC	Cray Inc.				nce , November 15, 2011 tors		)1
	9	Commissariat a l'Energie Atomique (CEA)	Bull	Bul	Hans & Week Lill	July Dongam_ Jack Dongarra	Horst Simo	1 <del>//</del> m	59
	10	DOE/NNSA/LANL	IBM	BladeCent Ghz / Oj	University of Mannheim NERSC/Berkeley Lab er QS22/LS21 Cluster, PowerXCell 81 3.2 oteron DC 1.8 GHz, Voltaire Infiniband	University of Tenness	ee NERSC/Be	1.U42	2.35





## The K computer's Performance













### LINPACK performance and its power consumption



Greenness



Program	Measured performance	Ranking
G-HPL(PFlops)	2.118	#1
G-Random Access (Gups)	121.1	#1
G-FFT (Tflops)	34.72*	#1
EP-STREAM (PB/s)	0.812	#1



 $^{\ast}$  : using FFTE5.0ß developed by Prof . Takahashi of Tsukuba Univ.

FUITSU

## Applications Performance, not a paper tiger

### March 2011

- cppmd molecular dynamics
  - 181 million atom simulation
  - Running on 221,184 cores (3.54 PFlops)
  - Sustained performance of 1.316 PFLOPS
  - Efficiency of *37%*



- 2011 Gordon Bell Award finalist
  - RIKEN, University of Tsukuba, University of Tokyo, Fujitsu
    - RSDFT program : first-principles calculations of electron states of a silicon nanowire Hasegawa, Iwata, Tsuji, Takahashi, Oshiyama, Minami, Boku, Shoji, Uno, Kurokawa, Inoue, Miyoshi, Yokokawa
    - 107,292-atom Si nanowire calculation
    - Running on 442,368 cores (7.08PFlops)
    - Sustained performance of 3.08 PFLOPS
    - Efficiency of 43.6 %



Various application programs are in optimization and evaluation stage



## K computer outline and technologies

New Interconnect , *Tofu* 

 6-dimensional Mesh/Torus topology
 High speed, highly scalable, high operability and high availability interconnect for over 100,000 nodes system

Functional interconnect

### ■ SPARC64<sup>™</sup> VIIIfx CPU

- Fujitsu designed and fabricated CPU
- ◆ 8 cores, 128GFlops@2GHz
- HPC-ACE (SPARC V9 Architecture Enhancement for HPC) :128GFlpos
- Main frame CPU level of high reliable design
- Low power consumption : ~58W





LINPACK 10.51 PFlops

RIKEN

17.6 PB

864 racks

88,128 CPUs

705,024 cores



- High memory bandwidth and simple memory hierarchy
- CPU/ICC direct water cooling
  - High reliability, low power consumption and compact packaging

## Agenda



- Achievements of *K computer* 
  - What did K computer achieve?
  - Technologies which realize the achievement
- Fujitsu's new Petascale supercomputer PRIMEHPC FX10
  - Fujitsu supercomputer past and present
  - Second generation Petascale supercomputer PRIMEHPC FX10

Conclusion

### **Fujitsu HPC Servers**





Copyright 2011 FUJITSU LIMITED

## **HPC Platform Solutions**



Full range coverage with choice of HPC platform



Nov. 16<sup>th</sup> 2011, SC11

## **Design targets and features of FX10**



 High Performance
 High peak performance and high effective performance

- Highly parallel application productivity
  - Easy to extract high performance from the highly paralleled programs without inordinate burden to programmers

User requirement of over PFlops computing FX10 design targets

High operability
 Low power consumption
 High reliability and easy to operate

 K computer compatibility
 Binary compatibility
 Same programing environment

## **Design targets and features of FX10**





Nov. 16<sup>th</sup> 2011, SC11

### **PRIMEHPC FX10 System Configuration**



FUÏTSU

## SPARC64™ IXfx

FUjitsu

### High-performance and low-power multi-core CPU

- High performance core by HPC-ACE
  - Register # extension, SIMD operation, software controllable cache, · ·
- VISIMPACT : Support highly efficient hybrid execution model (thread + process)
  - Shared 2<sup>nd</sup> cache, hardware barriers among cores and compiler

#### SPARC64<sup>™</sup> IXfx specifications

Architecture	SPARC64™ V9 + HPC-ACE			
# of FP operations /clock/core	8 (= 4 Multiply and Add )			
No. of cores	16			
Peak performance and clock	236.5 Gflops@1.848GHz			
Memory bandwidth	85 GB/s			
Power consumption	110 W (typical)			

- High performance-per-power ratio and High reliability
  - Water cooling system has lowered the CPU temperature and leak current
  - Wide-ranging error detection/self-recovery functions, instruction retry function



## **Node Configuration**

- Single CPU as a node design
  - ◆ SPARC64<sup>™</sup> IXfx based
  - 32 or 64 GB memory capacity
  - Single CPU per node to maximize memory bandwidth
  - High memory bandwidth of 85 GB/s
- On board InterConnect Controller (ICC)
  - Direct RDMA and global synchronization operations
  - No external switch
- Node type
  - Compute node
    - Consist of CPU, ICC and memory
    - Without I/O capability
    - Four nodes are mounted on a SB (system board)
  - I/O node
    - Same CPU as compute node
    - Includes four PCI Express Gen2 x8 slots
    - 8 GB/s I/O bandwidth per I/O node
    - One node is mounted on an I/O SB (system board)





## New Interconnect : Tofu

- Design targets
  - Scalabilities toward 100K nodes
  - High operability and usability
  - High performance
- Topology and performance
  - User view/Application view : Logical 3D Torus (X, Y, Z)
  - Physical topology : 6D Torus / Mesh addressed by (x, y, z, a, b, c)
    - 10 links / node, 6 links for 3D torus and 4 redundant links
  - Performance : 5GB/sec. /link x 2 (bi-directional)









 $\chi +$ 

#### 3D connection of each node

Copyright 2011 FUJITSU LIMITED

### Tofu ICC : Tofu Interconnect Controller

- ICC integrates Tofu interconnect and PCI Express I/F
- Tofu interconnect
  - Tofu network router (TNR) : 10 x TNRs transfer packets among ICCs,
  - Tofu network interface (TNI): 4 x RMDA communication engines
  - Tofu barrier interface (TBI) : Collective communication capability(Barrier and Allreduce)
- High reliable design

ICC specifications				
# of concurrent connections	4 transmission + 4 reception			
Switching capacity	100 GB/s			
Link speed x number of ports	5 GB/s x Bidirectional x 10 ports			
Operating frequency	312.5 MHz			
PCI Express	5 Gbps x 16 lanes			



## FX10 Software stack



### **Applications**

### **HPC Portal / System Management Portal**

#### **Technical Computing Suite** Automatic parallelization System Management **High Performance** compiler **File System** Fortran **FEFS** System management • C System control • C++ System monitoring Tools and math. libraries System operation support Lustre based high Programming support tools performance Mathematical libraries distributed file **Job Management** (SSL II/BLAS etc.) system • High scalability, high Parallel languages and libraries Job manager reliability and Job scheduler OpenMP availability Resource management MPI Parallel job execution XPFortran

Linux based OS enhanced for FX10

### **PRIMEHPC FX10**

## Language System overview

- Fortran C/C++/Fortran Compiler
- Programming model (OpenMP, MPI, XPFortran)
- Instruction level /Loop level optimization using HPC-ACE
- Debugging and Tuning tools for highly parallel computer



\*1: eXtended Parallel Fortran (Distributed Parallel Fortran) \*2: Rank Map Automatic Tuning Tool

## **FX10 System H/W Specifications**



	FX10 H/W Specifications			
		Name	SPARC64 <sup>TM</sup> IXfx	
	CFU	Performance	236.5GFlops@1.848GHz	
	Node	Configuration	1 CPU / Node	
		Memory capacity	32, 64 GB	
	Rack	Performance/rack	22.7 TFlops	
	System	No. of compute node	384 to 98,304	
		No. of racks	4 to 1,024	
		Performance	90.8TFlops to 23.2PFlops	
		Memory	12 TB to 6 PB	
•	<ul> <li>SPARC64<sup>™</sup> IXfx CPU</li> <li>16 cores/socket</li> <li>236.5 GFlops</li> <li>System rack</li> <li>96 compute nodes</li> <li>6 I/O nodes</li> <li>With optional water cooling exhaust un</li> </ul>		rit	
	Example of the second secon	<ul> <li>System board</li> <li>4 nodes (4 CPUs)</li> </ul>	<ul> <li>System</li> <li>Max. 23.2 PFlops</li> <li>Max. 1,024 racks</li> <li>Max. 98,304 CPUs</li> </ul>	

### K computer and FX10 Comparison of System H/W Specifications



		<i>`</i>	
		K computer	FX10
	Name	SPARC64 <sup>™</sup> VIIIfx	SPARC64 <sup>™</sup> IXfx
	Performance	128GFlops@2GHz	236.5GFlops@1.848GHz
	Architecture	SPARC V9 + HPC-ACE extension	←
CPU	Cache configuration	L1(I) Cache:32KB, L1(D) Cache:32KB	←
		L2 Cache: 6MB	L2 Cache: 12MB
	No. of cores/socket	8	16
	Memory band width	64 GB/s.	85 GB/s.
Nodo	Configuration	1 CPU / Node	←
Node	Memory capacity	16 GB	32, 64 GB
System board	Node/system board	4 Nodes	←
Pack	System board/rack	24 System boards	←
Nauk	Performance/rack	12.3 TFlops	22.7 TFlops

### K computer and FX10 Comparison of System H/W Specifications (cont.)



		湾	
		K computer	FX10
	Topology	6D Mesh/Torus	←
	Performance	5GB/s x2 (bi directional)	←
Interconnect	No. of link per node	10	←
	Additional feature	H/W barrier, reduction	←
	Implementation	No external switch box	←
	CPU, ICC(interconnect chip), DDCON	Direct water cooling	←
Cooling	Other parts	Air cooling	Air cooling + Exhaust air water cooling unit (Optional)

## Agenda



### Achievements of *K computer*

- What did K computer achieve?
- Technologies which realize the achievement
- Fujitsu's new Petascale supercomputer PRIMEHPC FX10
  - Fujitsu supercomputer past and present
  - Second generation Petascale supercomputer PRIMEHPC FX10
  - Conclusion

## "Computing" ideal future

### Achievements of K computer

- Archived LINPACK 10PFlops as scheduled
- ◆ No.1 in 37<sup>th</sup> and 38<sup>th</sup> TOP500 list, #1 in all HPCC class 1 Awards
- Over PFlops performance in real practical applications
- PRIMEHPC FX10 a new Petascale supercomputer for10PFlops era
  - Enhanced technologies applied to K computer
  - Application expertise through K computer project
  - Practical Petascale computing supported by HPC software stack

### Challenges to ExaFlops computing

- 100PFlops class system is *already* in sight
- Our challenges continues to achieve ExaFlops computing
- ExaFlpos computing will more quickly and objectively enable us to discover ways to deliver the future world, that as yet we have been unable to perceive.

### Fujitsu continues to work on realizing an ideal world through supercomputer development and its applications with you

#### Copyright 2011 FUJITSU LIMITED







# shaping tomorrow with you