SC12 Booth presentation



PRIMEHPC FX10 Performance evaluations and approach towards the next -

Toshiyuki Shimizu

Next Generation Technical Computing Unit FUJITSU LIMITED

November, 2012

Outline



- Introduction of Fujitsu petascale supercomputers and key technologies of massively parallel computers
- K computer and FX10 performance evaluations
 - HPC-ACE: CPU architecture extension
 - VISIMPACT: Hybrid parallel execution support
- Challenges towards Exascale computing
- Summary

Fujitsu HPC from workplace to top-end



SPARCE

PRIMERGY x86 Clusters

Celsius workstations





PRIMEHPC FX10 Supercomputers

Customers of large-scale HPC systems



Customer	Туре	Peak
RIKEN (Kobe AICS)	The K computer	11.28 PF
National Computational Infrastructure, Australia (*)	x86 Cluster (CX400), FX10	1.2 PF
The University of Tokyo	FX10	1.13 PF
Central Weather Bureau of Taiwan (*)	FX10	> 1 PF
Kyushu University	x86 Cluster (CX400), FX10	691 TF
HPC Wales, UK	x86 Cluster	>300 TF
Japan Atomic Energy Agency	x86 Cluster, FX1, SMP	214 TF
Institute for Molecular Science	x86 Cluster (RX300), FX10	>168 TF
Institute for Molecular Science (*)	x86 Cluster (CX400)	>136 TF
Japan Aerospace Exploration Agency	FX1, SMP	>135 TF
RIKEN (Wako Lab. RICC)	x86 Cluster (RX200)	108 TF
Institute for Solid State Physics of the Univ. of Tokyo (*)	FX10	90 TF
NAGOYA University	x86 Cluster (HX600), FX1, SMP	60 TF
A*STAR, Singapore	x86 Cluster (BX900)	> 45 TF
A Manufacturer	x86 Cluster	>250 TF

(*) To be operated, Type definitions: FX10=PRIMEHPC FX10, x86 Cluster=Clusters based on PRIMERGY x86 server, SMP= SPARC Enterprise SMP server

Massively parallel computing products



- Single CPU / node architecture for multicore CPU
 - FX1(4 cores) \rightarrow K(8 cores) \rightarrow FX10(16 cores)
 - High memory bandwidth balanced w/ CPU performance
- Key technologies for massively parallel are developed and inherited



Key technologies for massively parallel



Describe two technologies and show evaluation results of their effects by using real applications

HPC-ACE: CPU architecture extension

- Number of floating register extension
- Floating-point reciprocal approximation instructions
- Conditional execution (move, store, and mask generation) instructions

VISIMPACT: Hybrid parallel execution model (process & thread) support

- Automatic parallelization compiler
- Hardware inter-core barrier and shared L2 cache



HPC-ACE extension and its effect

- # of register extension
- Conditional execution instructions
- Floating-point reciprocal approximations

HPC-ACE: # of register extension

- Enhancement of SPARC V9 specification
 - # of integer regs from 32 to 64
 - # of DP FP regs from 32 to 256

The extended registers can be used as same as standard registers with any instructions

- Increasing # of register enables to explorer more parallelism
 - Large loops can be software pipelined



HPC-ACE: Conditional execution instructions-

```
subroutine sub(a, b, c, x, n)
real*8 a(n), b(n), c(n), x(n)
doi = 1, n
  if ( x(i) .gt. 0.0 ) then
    a(i) = b(i) * c(i)
   else
    a(i) = b(i) - c(i)
  endif
enddo
end
```

Conditional execution instructions				
L100: calculation fcmpgted,s fselmovd,s std,s	ns of b(i)*c(i) and b(i)-c(i) %f32,%f34,%f32 %f42,%f40,%f32,%f42 %f42,[%o5+%xg1]			
add bne nt	%xg1,16,%xg1 %iccL100			
nop :	70100, IL 100			

- Conditional branch is eliminated
- SIMDization and software pipelining can be widely applied

HPC-ACE: Floating-point reciprocal approx. Fujirsu



- Limited overlapping of instructions
- Divide and Sqrt are pipelined
- Software pipelining can also be applied

HPC-ACE improves application efficiency





Performance efficiency of real applications Fujirsu

- Measured by petascale K computer and PRIMEHPC FX10
- K computer & FX10 runs broad applications efficiently

# of Cores	Efficiency	IPC	System (Peak)
98,304	30%	1.7	K(1.57PF)
98,304	41%	1.5	
98,304	32%	1.3	
98,304	27%	1.0	
98,304	12%	0.6	
98,304	52%	1.5	
98,304	9%	0.7	
20,480	3%	0.6	K(0.33PF)
32,768	23%	0.9	K(0.52PF)
98,304	40%	1.5	K(1.57PF)
6,144	30%	1.6	FX10(0.1PF)
	<pre># of Cores 98,304 98,304 98,304 98,304 98,304 98,304 98,304 98,304 20,480 32,768 98,304 6,144</pre>	# of CoresEfficiency98,30430%98,30441%98,30432%98,30427%98,30412%98,30452%98,3049%20,4803%32,76823%98,30440%6,14430%	# of CoresEfficiencyIPC98,30430%1.798,30441%1.598,30432%1.398,30427%1.098,30412%0.698,30452%1.598,3049%0.720,4803%0.632,76823%0.998,30440%1.56,14430%1.6

* Measured by K computer: Results are from trial use & not the final.



VISIMPACT and its effect

- Hybrid execution model support
- Scalability and efficiency
- Load imbalance

VISIMPACT (<u>Vi</u>rtual <u>Single Processor by Integrated</u> <u>Multi-core</u> <u>Parallel</u> <u>Archi</u>tecture) FUJITSU

- Hybrid parallel execution is preferable for
 - Good scalability (reduce communications by reducing the # of processes)
 - Larger usable memory per process
 - Efficient use of memory (smaller work area)
- But...
 - Programing is harder due to two level of parallelization



Evaluations with practical meteorological apps.



WRF and COSMO were evaluated

- Operational weather forecasting codes
- Regional models
- Different computational characteristics
 - WRF uses "IO Quilting" mechanism
 - COSMO is memory intensive
 - WRF is thread-parallelized by OpenMP, but COSMO is not

Advantage of hybrid parallelization on FX10 Fujirsu

WRF V3.3.1



(WRF V3.3.1, 1200x700x45 grids, by courtesy of CWB)

Flat-MPI (1x1536) spends much time for communication
The communication includes load imbalance between processes
Load imbalance ratio[†] of Flat-MPI is 21%, while 9.8% for 16x96
* Load imbalance between threads in a process is in the calculation time
Node scalability of hybrid execution (16 threads/proc) is better

 $t:LoadImbalanceRatio = (MaxCalculationTime \div AvgCalculationTime) - 1$

Automatic thread-parallelization on FX10

COSMO-DE(RAPS 5.0)

Sustained Performance & Efficiency



Effective B/F Ratio & Memory BW



(COSMO RAPS 5.0, 421x461x50 grids)

VISIMPACT parallelizes with threads and improves performance when cores are increased

For 384 cores, load imbalance ratio[†] was mitigated from 15% to 7.9%
 Thread parallelization reduces byte/flop and utilizes memory BW

†:LoadImbalanceRatio=(MaxCalculationTime÷AvgCalculationTime)-1

Summary of VISIMPACT (hybrid execution) Fujirsu

- Real applications run on the hybrid parallel mode using VISIMPACT show better performance than Flat MPI
 - By eliminating memory copy
 - By reducing load imbalance between processes
- Hybrid parallel execution model is less memory usage
 - Intrinsic and indispensable features for many core environment
- Reducing a load imbalance between threads in a process should be studied



Challenges towards Exascale computing

Challenges toward exascale computing



- Fujitsu is developing a 100 Petaflops capable system as a midterm goal
- Participates two consecutive national projects for exascale
 - Fujitsu contributed to develop the whitepaper
 - "Report on Strategic Direction/Development of HPC in Japan"
 - Fujitsu has started two-year feasibility study to set the goal & schedule since July 2012



Summary



K computer & FX10 employ massively parallel technologies HPC-ACE

- VISIMPACT
- Evaluation running real applications
 - Good performance efficiency
 - Good scalability and smaller load imbalance

Research and develop massively parallel architecture and approach toward exascale step-by-step

Exascale system



2015

2010

K compute

FUJTSU

shaping tomorrow with you

SC12 Booth presentation for FX10