

Programming for Fujitsu Supercomputers

Koh Hotta

The Next Generation Technical Computing Fujitsu Limited



To Programmers who are busy on their own research, Fujitsu provides environments for

Parallel Programming

Tuning Parallel Programs

System Software Stack



User/ISV Applications								
HPC Portal / System Management Portal								
Technical Computing Suite								
 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 	<u>Compilers</u> Hybrid parallel programming Sector cache support SIMD / Register file extensions Support Tools IDE 						
 Job operations management Job manager Job scheduler Resource management Parallel execution environment 	 <u>VISIMPACT</u>™ Shared L2 cache on a chip Hardware intra-processor synchronization 	 Profiler & Tuning tools Interactive debugger <u>MPI Library</u> Scalability of High-Func. Barrier Comm. 						

Linux-based enhanced Operating System

Red Hat Enterprise Linux

Super Computer: PRIMEHPC FX10

PC cluster: PRIMERGY

Copyright 2012 FUJITSU LIMITED

HPC-ACE architecture & Compiler

Expanded Registers

- reduces the wait-time of instructions
- SIMD instructions
 - reduces the number of instructions
- Sector Cache
 - improves cache efficiency

Compiler Gives Life to HPC-ACE



NPB3.3 LU: operation wait-time dramatically reduced NPB3.3 MG: number of instruction is halved





Large # of Parallelism for Large Scale Systems

- Large # processes need Large Memory & Overhead
 - hybrid thread-process programming to reduce number of processes
- Hybrid parallel programming is annoying for programmers
- Even for multi-threading, the more coarse grain, the better.
 - Procedure level or outer loop parallelism is desired
 - Little opportunity for such coarse grain parallelism
 - System support for "fine grain" parallelism is required
- VISIMPACT solves these problems



Mechanism that treats multiple cores as one high-speed CPU through automatic parallelization

Just program and compile, and enjoy high-speed

You need not think about hybrid



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



- Technologies to make multiple cores to single high speed CPU
 - Shared L2 cache memory to avoid false sharing
 - Inter-core hardware barrier facilities to reduce overhead of thread synchronization
 Barrier synchronization

Thread parallel programs use these technologies



5.000

thread & process Hybrid-Parallel Programming

100

thread parallel in a chip

- Auto-parallel or explicit parallel description using OpenMPTM API
- VISIMPACT enables low overhead thread execution
- process parallel beyond CPU
 - MPI programming
 - Collective communication is tuned using Tofu barrier facility



1,000



MPI Software stack





Customized MPI Library for High Scalability

Point-to-Point communication

- Two methods for inside communication
- •The transfer method selection by
 - •the data length
 - process location
 - •number of hops
- Collective communication
 - •Barrier, Allreduce, Bcast and Reduce use Tofu-barrier facility
 - Bcast, Allgather, Allgatherv, Allreduce and Alltoall use Tofu-optimized algorithm





Application Tuning Cycle and Tools





Copyright 2012 FUJITSU LIMITED

Program Profile (Event Counter Example)

Profiler						- ¤ ×
GlobalView RankView						
Time	Elap:	s Time : 118.60	017 (S)			
summary total preview xyz Evo	entCounter DataCompare	Transparent				
Type 3D Type 2D Type 1D Intangible						
REG-1 (1) : User		Basic MPI	Hardware			
Z	e Value			User		
		MAX 108.8997	AVE.	MIN 0.0001	MAX	AVE.
		86,1978	210.0000	88.000000	328.0000	6,1875
		92.6941	6.3180	0.0000	78.3902	6.3862
		82.8526	6.3813	0.0001	70.2137	6.1854
→ X						•
In the Inc. In the Inc.	h. h. h. h.	b. b.	b. b	. h. l	ln.	
			h. 111.	lla, illa,		
					zo (om x
					🔾 zo	om y
					🔿 🗆 🔿 🗸	om z
				and and	100	
		1 1	1 1	4 4	4	
Δ	<u></u>		<u> </u>	Cradation	-	
88.000000			328.0000	anadation		

3-D job example

- Display 4096 procs in 16 x 16 x 16 cells
- Cells painted in colors according to the proc status (e.g. CPU time)
- Cut a slice of jobs along x-, y-, or z-axis to view

Rank Mapping Optimization (RMATT)

Network Construction Communication Pattern (Communication processing contents between Rank)



addiv

Apply MPI_Allgather Communication Processing Performance

•Rank number : 4096 rank

•Network construction : 16x16x16 node (4096)





FUĨITSU

FUJITSU

shaping tomorrow with you