Oakforest-PACS (OFP) : the largest scale many-core supercomputer in Japan

Taisuke Boku

Center for Computational Sciences, University of Tsukuba

(with courtesy by Toshihiro Hanawa Information Technology Center, the University of Tokyo)



2016/05/12 CCS-LBNL Workshop 2016

Outline

- Supercomputer deployment plan in Japan
- What is JCAHPC ?
- Supercomputer procurement in JCAHPC
- System specification
- Pre-evaluation by KNC (not KNL)
- Conclusions



Deployment plan of 9 supercomputing center (Oct. 2015 + latest)

Fiscal Year		2014 2015 20	016 2017	2018 2019	2020 2021	2022	2023	2024	2025
Hokkaido		HITACHI SR 16000/M1 (Cloud System BS2000 (Data Science Cloud / Storage H/ (10TF, 1.96PB)	172TF, 22TB) 44TF, 14TB) 28000 / WOS7000	3.2 PF (UCC 0.3 PF (C + CFL/M) 0.96 (Cloud) 0.36MW	6MW	30 F CFL	PF(UCC -M) 2N) + MW
Tohoku	W	NEC SX- 9他 (60TF)	-ACE(707TF,160 (31TF), Storage	0TB, 655TB/s) (4PB), 3D Vis, 2MW	~30PF, ~3	0PB/s Men ~3I	n BW (C MW	FL-D/CF	L-M)
Tsukuba		HA-PACS (1166 TF) COMA (PACS-IX) (1001 TI	-) -)	PACS-X 10PF (TP	F) 2MW		100+	PF 4.5MW	/
Tokyo		Fujitsu FX10 (1PFlops, 150TiB, 408 TB/s) Hitachi SR16K/M1 (54.9 TF, 10.9 TiB, 2	Reedbush 1.8	(UCC + TPF) 4.2 ~1.9 PF 0.7 MW	MW 50+ PF (FAC) 3.5MW	(UC	C + TPF)	200+ PF AC) 6.5MW
Tokyo Tech.		TSUBAME 2.5 (5.7 PF, 110+ TB, 1160 TB/s), 1.4MW	TSUBAME 2.5 (3~ TSUE	~4 PF, extended) BAME 3.0 (20 PF, 4~6PE 3.5, 40PF at 2018 if upgra	B/s) 2.0MW adable)	TS	UBAME 4.0 >10PB/s, ~/	(100+ PF, 2.0MW)	
Nagoya		FX10(90TF) Fujitsu CX400(470T Fujitsu F) SGI UV2000 (24TF, 2)	FX100 (2.9PF, X400 (774TF, 7 0TiB) 2MV	81 TiB) 50- V in total	+ PF (FAC/UC	C + CFL up to	-M) 5 4MW	100+F (FAC/U мур to	PF JCC+CFL- 5 4MW
Kyoto		Cray: XE6 + GB8K + XC30 (983TF) Cray XC30 (584TF)	7-8 PF	F(FAC/TPF + UCC 1.5 MW	;) / (FAC,	50-100 TPF + UC	+ PF C) 1.8-2.4	4 MW	
Osaka		NEC SX-ACE NEC (423TF) (22.4TF)	Express5800	3.2 0.7-1PF (I	2PB/s, 5-10Pflop JCC)	/s, 1.0-1.5M	IW (CFL-	M) 25.6 100F 2.0N	9 PB/s, 50- Pflop/s,1.5- /W
Kyushu		HA8000 (712TF, 242 TB) SR16000 (8.2TF, 6TB) FX10 (272.4TF, 36 TB) CX400 (966.2 TF, 183TB)	₩ 15-20 F FX (90.8TF	PF (UCC/TPF) 2.6	MW	(FAG	100-15 C/TPF +	0 PF UCC/TP	F3MW
			Power consu (includes co	mption indicates oling facility)	maximum of po	wer supply			
3		2016/05/12	CCS-I BNI	Workshop 2016					

2016/05/12 CCS-LBNL Workshop 2016

T2K Open Supercomputer Systems



- Same timing of procurement for next generation supercomputers in three universities
- Academic leadership for computational science/engineering in research/education/grid-use on same platform

Kyoto Univ.

416 nodes (61.2TF) / 13TB Linpack Result: Rpeak = 61.2TF (416 nodes) Rmax = 50.5TF



Univ. Tokyo

952 nodes (140.1TF) / 31TB Linpack Result: Rpeak = 113.1TF (512+256 nodes) Rmax = 83.0TF



- Open hardware architecture with commodity devices & technologies.
- Open software stack with opensource middleware & tools.
- Open to user's needs not only in FP & HPC field but also INT world.

Univ. Tsukuba

648 nodes (95.4TF) / 20TB Linpack Result:

Rpeak = 92.0TF (625 nodes) Rmax = 76.5TF



From T2K to Post-T2K

Effect of T2K Alliance

- Three supercomputers are introduced at the same time, sharing wide knowledge for system construction and commodity technology, followed by academic research collaboration among these players
- After T2K, three universities had different time of new system procurement
 - Kyotop U.: four year period of procurement
 - U. Tsukuba: accelerated computing
 - U. Tokyo: T2K + Fujitsu FX10 and other systems
- Post-T2K (with two "Ts")
 - in 2013, U. Tsukuba and U. Tokyo collaborated again for new supercomputer procurement in *much more tight framework*





JCAHPC

- Joint Center for Advanced High Performance Computing (<u>http://jcahpc.jp</u>) (最先端共同HPC基盤施設)
- Very tight collaboration for "post-T2K" with two universities
 - For main supercomputer resources, *uniform specification* to single shared system
 - Each university is financially responsible to introduce the machine and its operation

-> unified procurement toward single system with *largest scale in Japan*

- To manage everything smoothly, a joint organization was established
 - -> JCAHPC



2016/05/12



History of JCAHPC

- March 2013: U. Tsukuba and U. Tokyo exchanged agreement for "JCAHPC establishment and operation"
 - Center for Computational Sciences, University of Tsukuba and Information Technology Center, University of Tokyo
- April 2013: JCAHPC started
 - 1st period director: Mitsuhisa Sato (Tsukuba), vice director: Yutaka Ishikawa (Tokyo)
 - 2nd period (2016~) director: Hiroshi Nakamura (Tokyo), vice director: Masayuki Umemura (Tsukuba)
- July 2013: RFI for procurement
 - at this time, the joint procurement style was not fixed
 -> then a single system procurement was decided
 - to give enough time for very advanced technology for processor, network, memory, etc., more than 1 year of period was taken to fix the specification
- It is the first trial to introduce a shared single supercomputer system by multiple national universities in Japan !



Procurement Policies of JCAHPC

- **based** on the spirit of T2K, introducing open advanced technology
 - massively parallel PC cluster
 - advanced processor for HPC
 - easy to use and efficient interconnection
 - large scale shared file system flatly shared by all nodes
- joint procurement by two universities
 - the largest class of budget as national universities' supercomputer in Japan
 - the largest system scale as PC cluster in Japan
 - no accelerator to support wide variety of users and application fields
 -> not chasing absolute peak performance and inheriting traditional application codes (basically)
- goodness of single system
 - scale-merit by merging budget -> largest in Japan
 - ultra large scale single job execution at special occasion such as "Gordon Bell Prize Challenge"

\Rightarrow Oakforest-PACS (OFP)



CO JCAHPC



Specification of Oakforest-PACS

Total peak performance		e	25 PFLOPS		
Total number of compute nodes		te nodes	8,208		
Compute node	Product		Fujitsu Next-generation PRIMERGY server for HPC (under development)		
	Processor		Next-generation of Intel® Xeon Phi [™] (Code name: Knights Landing), >60 cores		
	Memory	High BW	16 GB, > 400 GB/sec (MCDRAM, effective rate)		
		Low BW	96 GB, 115.2 GB/sec (DDR4-2400 x 6ch, peak rate)		
Inter-	Product		Intel® Omni-Path Architecture		
connect	Link speed		100 Gbps		
	Topology		Fat-tree with (completely) full-bisection bandwidth		
Login	Product		Fujitsu PRIMERGY RX2530 M2 server		
node	# of servers	5	20		
	Processor		Intel Xeon E5-2690v4 (2.6 GHz 14 core x 2 socket)		
	Memory		256 GB, 153 GB/sec (DDR4-2400 x 4ch x 2 socket)		



2016/05/12



Specification of Oakforest-PACS (I/O)

Parallel File	Туре		Lustre File System	
System	Total Cap	acity	26.2 PB	
	Meta data	Product	DataDirect Networks MDS server + SFA7700X	
		# of MDS	4 servers x 3 set	
		MDT	7.7 TB (SAS SSD) x 3 set	
	Object storage	Product	DataDirect Networks SFA14KE	
		# of OSS (Nodes)	10 (20)	
		Aggregate BW	500 GB/sec	
Fast File	Туре		Burst Buffer, Infinite Memory Engine (by DDN)	
Cache System	Total capacity		940 TB (NVMe SSD, including parity data by erasure coding)	
	Product		DataDirect Networks IME14K	
	# of serve	ers (Nodes)	25 (50)	
	Aggregate BW		1,560 GB/sec	



CCS-LBNL Workshop 2016 2016/05/12

10

CO JCAHPC

Full bisection bandwidth Fat-tree by Intel® Omni-Path Architecture





Facility of Oakforest-PACS system

Power consumption			4.2 MW (including cooling)	
# of racks			102	
Cooling system	Compute Node	Туре	Warm-water cooling Direct cooling (CPU) Rear door cooling (except CPU)	
		Facility	Cooling tower & Chiller	
	Others	Туре	Air cooling	
		Facility	PAC	



Software of Oakforest-PACS

	Compute node	Login node				
0S	CentOS 7, McKernel	Red Hat Enterprise Linux 7				
Compiler	gcc, Intel compiler (C, C++, Fortran)					
MPI	Intel MPI, MVAPICH2					
Library	Intel MKL					
	LAPACK, FFTW, SuperLU, PETSc, METIS, Scotch, ScaLAPACK GNU Scientific Library, NetCDF, Parallel netCDF, Xabclib, ppOpen-HPC, ppOpen-AT, MassiveThreads					
Application	mpijava, XcalableMP, OpenFOAM, ABINIT-MP, PHASE system, FrontFlow/blue, FrontISTR, REVOCAP, OpenMX, xTAPP, AkaiKKR, MODYLAS, ALPS, feram, GROMACS, BLAST, R packages, Bioconductor, BioPerl, BioRuby					
Distributed FS	Globus Toolkit, Gfarm					
Job Scheduler	Fujitsu Technical Computing Suite					
Debugger	Allinea DDT					
Profiler	Intel VTune Amplifier, Trace Analyzer & Collector					



CO JCAHPC

(pre) Photo of computation node





Chassis with 4 nodes, 2U size

Computation node (Fujitsu next generation PRIMERGY) with single chip Intel Xeon Phi (Knights Landing, 3+TFLOPS) and Intel Omni-Path Architecture card (100Gbps)



14



Machine location: Kashiwa Campus of U. Tokyo

つくば市 大道語語識 検索 坂東市 **U.** Tsukuba 常総市 茨城空港 日部市 霞ヶ浦 埼玉県 重野町。 秩父 ||越 Balt さいたま 取手市 越谷市 ク井倉 **Kashiwa** 。我孫子市 0 柏市 印西 十九里町 Campus 相認 of U. Tokyo States. 16 神奈川 足立区 海道和線/藤沢の鎌倉 46 大多喜町 6 相模波 御宿町 1千代市 市区 14 众编 RBR 野島崎 伊豆半島 Hongo Campus of U. Tokyo 大島町 4 ΓŒ 三原山 尹豆町 地図を見栄えよく印刷するには、メインメニュー の (印刷) ポタンをご利用ください。 利息村 ライトモード 地図アータ ©2015 Google, ZENRIN 利用規約 プライパシー 15 2016/05/12 Center for Computational Sciences, Univ. of Tsukuba 1/12015/05/20 11:02

Google マップ

https://www.google.com/maps/@?dg=dbrw&newdg=1



Schedule

- 2013/7 RFI
- 2015/1 RFC
- 2016/1 RFP
- 2016/3/30 Proposal deadline
- 2016/4/20 Bid opening
- 2016/10/1 1st step system operation (more than 5% of full system)
- 2016/12/1 2nd step, full system operation
- 2017/4 National open use starts including HPCI
- 2022/3 System shutdown (planned)





System operation outline

Regular operation

- both universities share the CPU time based on the budget ratio
- not split the system hardware, but split the "CPU time" for flexible operation (except several specially dedicated partitions)
- single system entry for HPCI program, and other own program by each university is performed under "CPU time" sharing

Special operation

- massively large scale operation (limited period)
 -> effectively using the largest class resource in Japan for special occasion (ex. Gordon Bell Challenge)
- Power saving operation
 - power capping feature for energy saving scheduling feature reacts to power saving requirement (ex. summer time)



Pre-evaluation environment

- U. Tsukuba and U. Tokyo are running own systems based on current Intel Xeon Phi (KNC) for performance evaluation and tuning
 - Tsukuba: COMA (PACS-IV), 393 nodes, 786 Xeon Phi
 - Tokyo: KNSoaring, 64 nodes, 64 Xeon Phi
- COMA at U. Tsukuba is also dedicated to HPCI and other programs since April 2015.
 - \Rightarrow we will introduce KNL-based small system soon



COMA (PACS-IX) – KNC cluster



- Cray CS300 Cluster
- Intel Xeon Phi (KNC: Knights Corner)
- 393 nodes (2 Xeon E5-2670v2 + 2 Xeon Phi 7110P)
- Mellanox IniniBand FDR, Fat Tree
- Largest Xeon Phi Cluster in Japan, as at Oct. 2015
- File Server: DDN
 1.5PB (RAID6+Lustre)
- 1.001 PFLOPS (HPL: 746 TFLOPS) June '14 TOP500 #51



19

2016/05/12

Xeon Phi tuning on ARTED (with Yuta Hirokawa under collaboration with Prof. Kazuhiro Yabana, CCS)

- ARTED Ab-initio Real-Time Electron Dynamics simulator
- Multi-scale simulator based on RTRSDFT (Real-Time Real-Space Density Functional Theory) developed in CCS, U. Tsukuba to be used for Electron Dynamics Simulation
 - RSDFT : basic status of electron (no movement of electron)
 - RTRSDFT : electron state under external force
- In RTRSDFT, RSDFT is used for ground state

- RSDFT : large scale simulation with 1000~10000 atoms (ex. K-Computer)
- RTRSDFT : calculate a number of unit-cells with 10 ~ 100 atoms



Computation domain and amount

- Parameters for wave function expression
 - k-points (NK), band-number (NB), 3-D lattice points (NL)
 - valuables are in double precision complex with matrix of (NK, NB, NL)
 - for stencil computation, size NL of calculation is performed NKxNB times
- Parameters used in this research (two models)
 - SiO₂ : (4³, 48, 36000 = (20, 36, 50)) -> not enough large
 - Si : (24³, 32, 4096 = (16, 16, 16)) -> larger parallelism on thread
- NK is parallelized by MPI, then NKxNB is parallelized in OpenMP
 - domain of each process: (NK/NP, NB, NL) (NP = number of processes)
 - space domain is not decomposed to minimize MPI communication



CCS-LBNL Workshop 2016 2016/05/12

Stencil code (original)

<pre>integer, intent(in) :: IDX(-4:4,NL),IDY(-4:4,NL),IDZ(-4:4,NL)</pre>				
! NL = NLx*NLy*NLz	indirect index array: keeping nearest neighbor index			
<pre>v(1)=Cx(1)*(E(IDX(1,i))+E(IDX(-1,i)) w(1)=Dx(1)*(E(IDX(1,i))-E(IDX(-1,i))</pre>				
<pre>w(1)=bx(1) (E(1bx(1,1)) E(1bx(1,1)) ! y-computation v(2)=Cy(1)*(E(IDY(1,i))+E(IDY(-1,i)) w(2)=Dy(1)*(E(IDY(1,i))-E(IDY(-1,i)))</pre>))))			
<pre>! z-computation v(3)=Cz(1)*(E(IDZ(1,i))+E(IDZ(-1,i)) w(3)=Dz(1)*(E(IDZ(1,i))-E(IDZ(-1,i))</pre>))))			
<pre>! update F(i) = B(i)*E(i) + A*E(i) - 0.5d0*(v end do</pre>	v(1)+v(2)+v(3)) - zI*(w(1)+w(2)+w(3))			

vector length=4, for DP-complex vector calculation-> 512-bit AVX fittable



2016/05/12

For automatic vectorization





Stencil computation performance

- (NK, NB, NL) = $(8^3, 16, 16^3)$, single process
- 2x performance of Ivy-Bridge Xeon

	Туре	GFLOPS	ratio to peak (%)
Voor Dhi	Original	29.0	2.70
7110D	Compiler Vec.	132.2	12.30
/1101	Explicit Vec.	212.2	19.75
Irre Duidas	Original	53.7	26.85
$E_5 2670x^2$	Compiler Vec.	102.7	51.35
E3-20/0V2	Explicit Vec.	106.9	53.45



CCS-LBNL Workshop 2016 2016/05/12

Relative perf. to CPU on entire code (strong scaling)



2016/05/12

Center for Computational Sciences, Univ. of Tsukuba

Comparison with other systems

Entire code execution performance of ARTED

ratio to peak Vectorization Performance Processor [GFLOPS] [%] 64.4 **KNC** 6.0 compiler 103.4 9.6 explicit 27.4 IvyBridge compiler 54.7 27.9 explicit 55.9 21.7 Haswell 83.2 compiler 70.7 18.4 explicit Sparc64 VIIIfx compiler 13.0 10.1 **K** Computer

exec. efficiency is comparable with K Computer



3x in KNL -> 300 GFLOPS ?

CCS-LBNL Workshop 2016

2016/05/12

Center for Computational Sciences, Univ. of Tsukuba

Summary

- JCAHPC is a joint resource center for advanced HPC by U. Tokyo and U. Tsukuba
- Oakforest-PACS with 25 PFLOPS peak performance with Intel Xeon
 Phi (Knights Landing) and Omni-Path Architecture as interconnect
- small part will start operation on Oct. 2016, and the full system will be available on Dec. 2016
- Under JCAHPC, both universities perform multiple resource sharing programs including HPCI
- JCAHPC is not just an organization to manage the resource but also a basic community for advanced HPC research

2016/05/12

