

JCAHPCの新スーパーコンピュータ Oakforest-PACS

東京大学 情報基盤センター
最先端共同HPC基盤施設(JCAHPC)

中村宏

Oakforest-PACSの全景



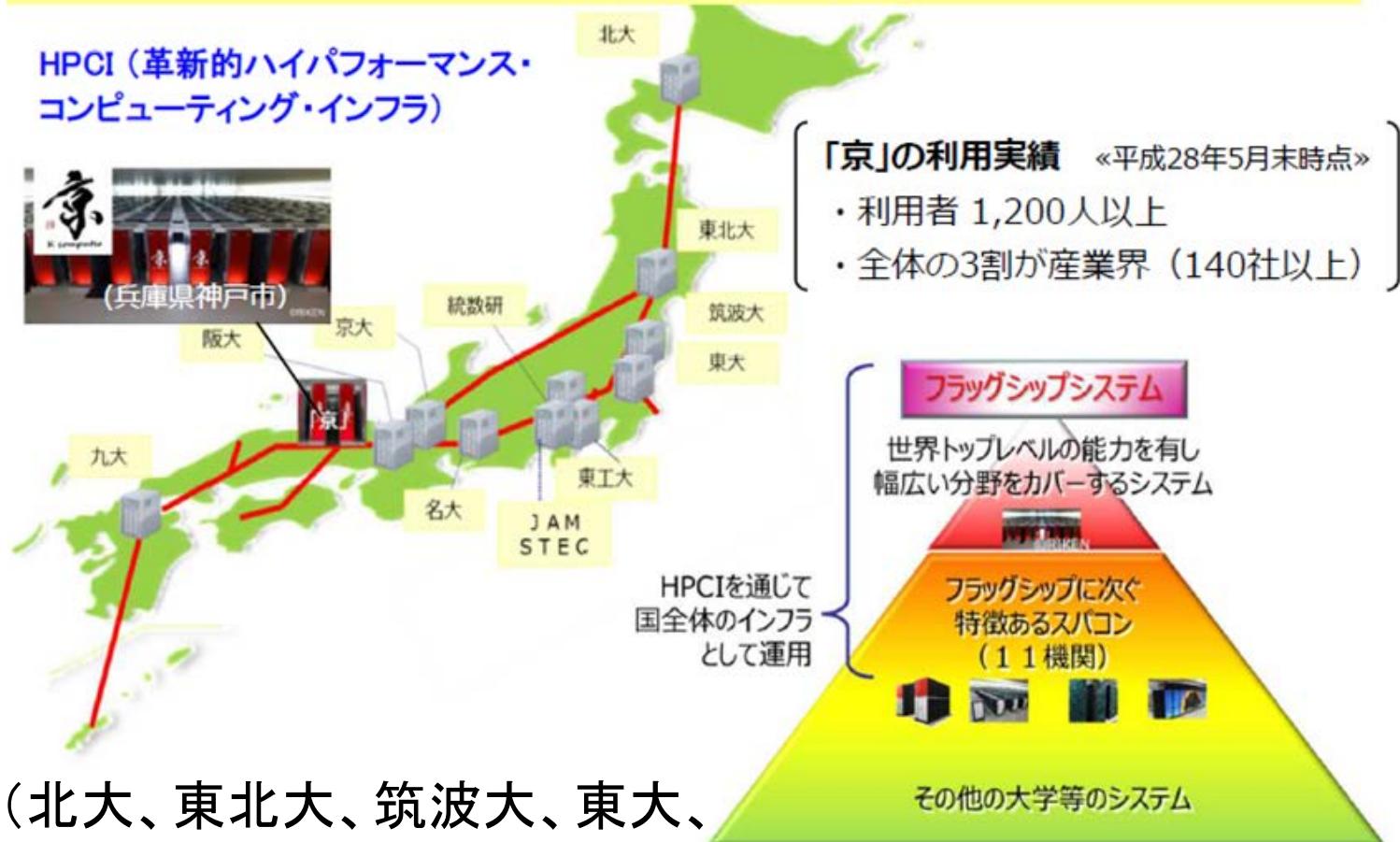
Oakforest-PACS



- 最先端共同HPC 基盤施設(JCAHPC: Joint Center for Advanced High Performance Computing)
 - 東京大学情報基盤センター
 - 筑波大学計算科学研究センター
 - 両センターが共同で、最先端の大規模高性能計算基盤を構築・運営するための組織
 - 東京大学柏キャンパスの東京大学情報基盤センター内
- 2016年12月1日稼働開始
- 8,208 Intel Xeon/Phi (KNL)
- ピーク性能25PFLOPS
- **TOP 500 6位(国内1位), HPCG 3位(国内2位), Green 500 6位(国内2位)(2016年11月)**

HPCI: High Performance Computing Infrastructure

日本全体におけるスパコンインフラ

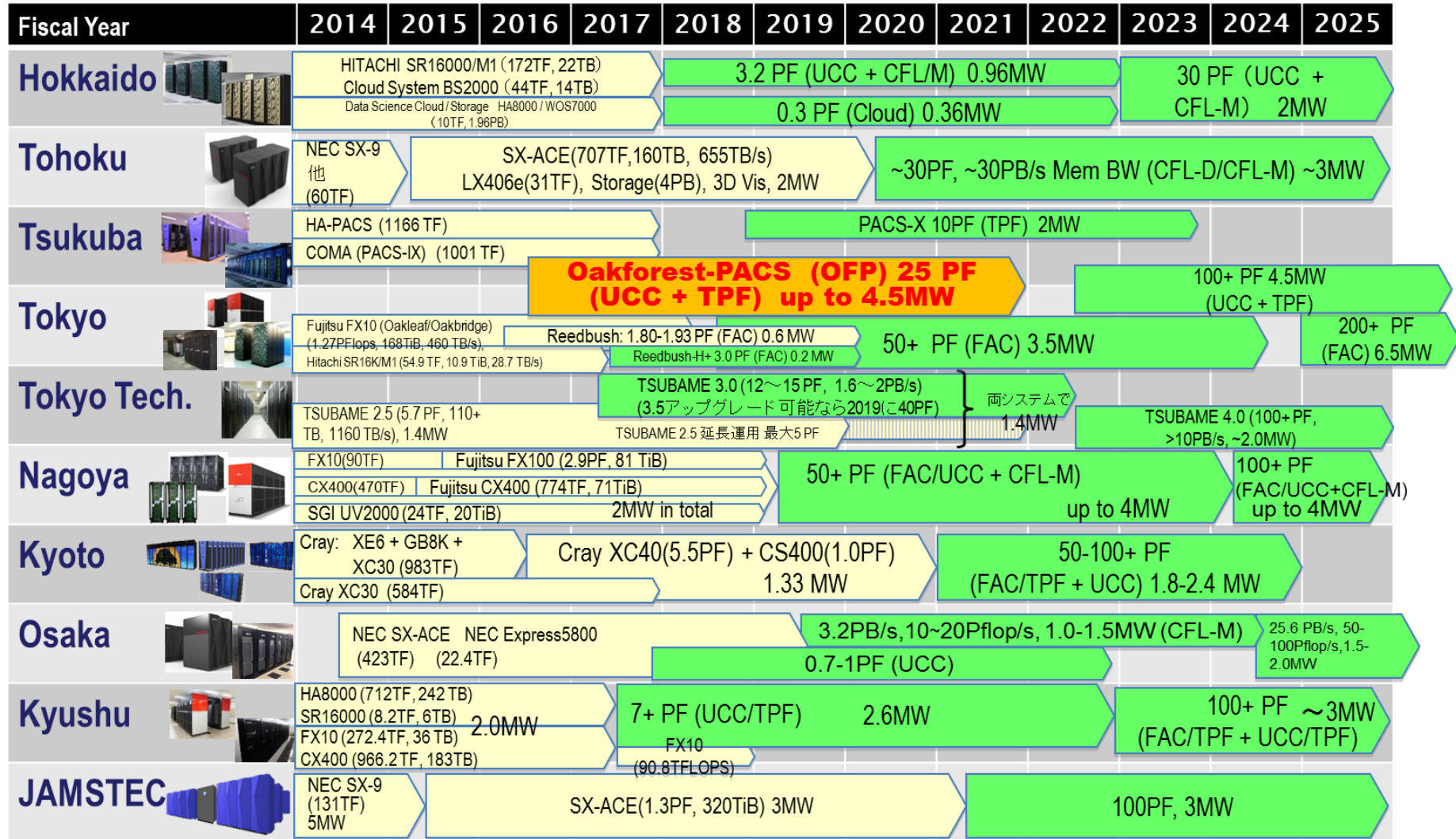


- 9大学(北大、東北大、筑波大、東大、東工大、名大、京大、阪大、九大)の情報基盤センター
- 海洋開発研究機構、統数研 + SINET

HPCI第2階層システムの開発・整備・運用計画

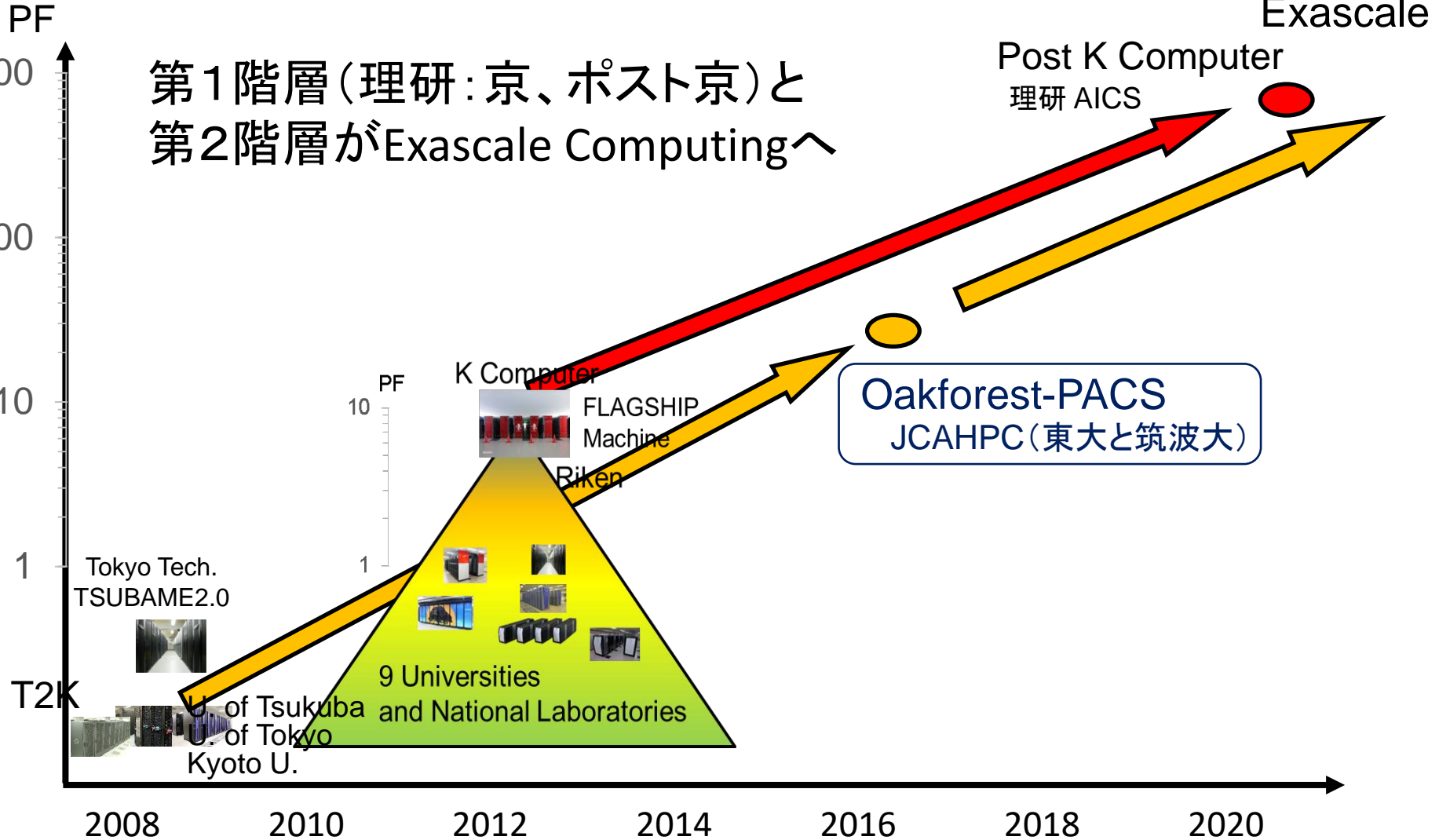
(2016年9月時点)

↓ HPCIコンソーシアムのホームページ掲載 <http://www.hpci-c.jp/>



フラグシップとの両輪として

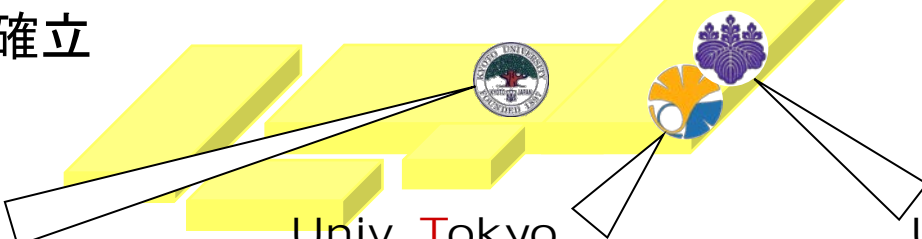
第1階層(理研:京、ポスト京)と
第2階層がExascale Computingへ



OFPまでの道: T2Kオープンスパコンアライアンス

- 筑波大・東大・京大による次世代スパコン技術の推進のための同時調達システム
- 計算科学・計算工学における研究・教育・システム共用に関する大学のリーダーシップの確立

- **オープン**ハードウェアアーキテクチャ
(コモディティ技術による)
- **オープン**ソフトウェアスタック
(オープンソースミドルウェアとツール)
- **オープン**な利用とユーザ知識・アプリケーションの共有



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



2016/12/1
6

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



PCクラスタコンソーシアム

Univ. Tsukuba
648 nodes (95.4TF) / 20TB
Linpack Result:
Rpeak = 92.0TF (625 nodes)
Rmax = 76.5TF



7

T2Kから「ポストT2K」へ

- T2Kオープンスパコンアライアンスの効果
 - 3台のスパコンは同時に調達・運用され、システム構築技術や性能チューニング技術等を共有し、大学間の強いHPC研究コミュニティの発展につながった。
- T2Kの後、異なる次期システム調達スケジュールやシステム開発ポリシーの違いにより、発展的に解消
 - 京大:4年間サイクルの調達
 - 筑波大:演算加速器系システムに傾注
 - 東大:T2KだけでなくFACシステムとしてFX10を導入
- そして「ポストT2K」へ
 - 2013年、筑波大学と東京大学による新たなスパコン導入の枠組み → **JCAHPC**
 - T2Kを越える、より強固な連携によるシステム調達

最先端共同HPC基盤施設 JCAHPC

- Joint Center for Advanced High Performance Computing (<http://jcahpc.jp>)
- 平成25年3月、筑波大学と東京大学は「計算科学・工学及びその推進のための計算機科学・工学の発展に資するための連携・協力推進に関する協定」を締結
- 本協定の下、筑波大学計算科学研究センターと東京大学情報基盤センターが **JCAHPC** を設置
 - 両センターが共同で、最先端の大規模高性能計算基盤を構築・運営するための組織
 - 東京大学柏キャンパスの東京大学情報基盤センター内

JCAHPC: 共同調達への道のり

- 2013活動開始
 - 第1期(2013/4-2015/3):
施設長:佐藤三久(筑波大学)、副施設長:石川裕(東京大学)
 - 第2期(2015/4-):
施設長:中村宏(東京大学)、副施設長:梅村雅之(筑波大学)
- 共同調達・運用へ向けて
 - 2013/7: RFI(request for information)
共同調達は既定路線ではなかった→1システムとして調達へ
 - 複数大学による初めての「1システム」共同調達へ
- どうして共同調達ができたのか? 共同調達は大変・
 - 目標を共有できる、ことに尽きる

2センターのミッション

- 筑波大学計算科学研究センターのミッション：
 - **先端学際科学共同研究拠点**: 最先端の計算科学研究推進
 - 計算機科学と計算科学の協働:**学際的な高性能計算機開発**
→ PACSシリーズの開発: CP-PACS@1996 TOP1
- 東京大学情報基盤センターのミッション：
 - **学際大規模情報基盤共同利用・共同研究拠点**
(8大学の情報基盤センター群からなるネットワーク型) の
中核拠点: 大規模情報基盤を活用し学際研究を発展
 - HPCI資源提供機関: **最先端スパコンの共同設計開発及び運用**、Capability資源および共用ストレージ資源の提供
- 共通部分:
最新のスーパーコンピュータを用いて学際研究を先導

JCAHPC共同調達のポリシー ～2センターで共有したこと～

- T2Kの精神に基づき、オープンな最先端技術を導入
 - T2K: 2008年に始まったTsukuba, Tokyo, Kyoto の3大学でのオープンスパコンアライアンス、3機関の研究者が仕様策定に貢献、システムへの要求事項を共通化
- システムの基本仕様
 - 超並列PCクラスタ
 - HPC用の最先端プロセッサ、アクセラレータは不採用
 - 広範囲なユーザとアプリケーションのため
 - ピーク性能追求より、これまでのコードの継承を優先
 - 使いやすい高効率相互結合網
 - 大規模共用ファイルシステム
- スケールメリットを活かす
 - 超大規模な単一ジョブ実行も可能とする



Oakforest-PACS

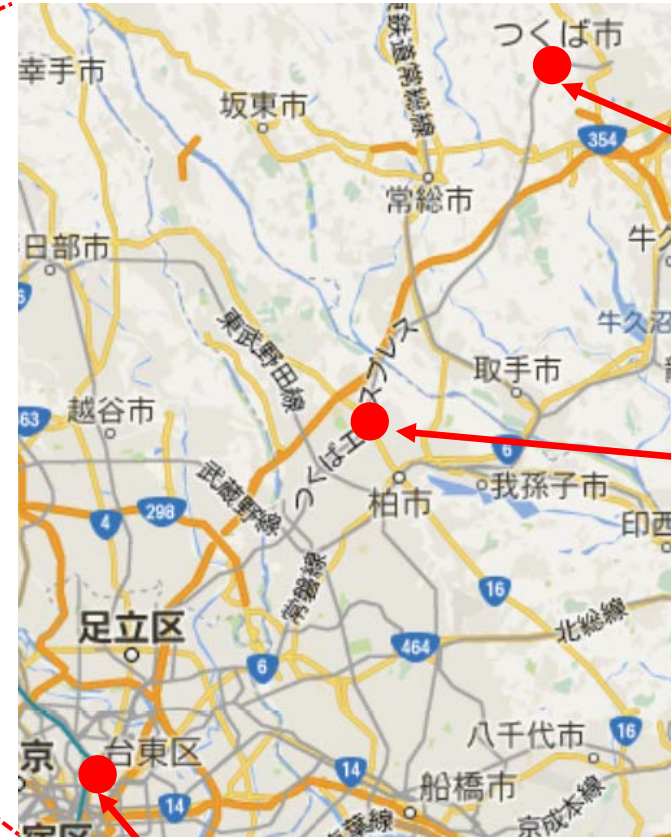
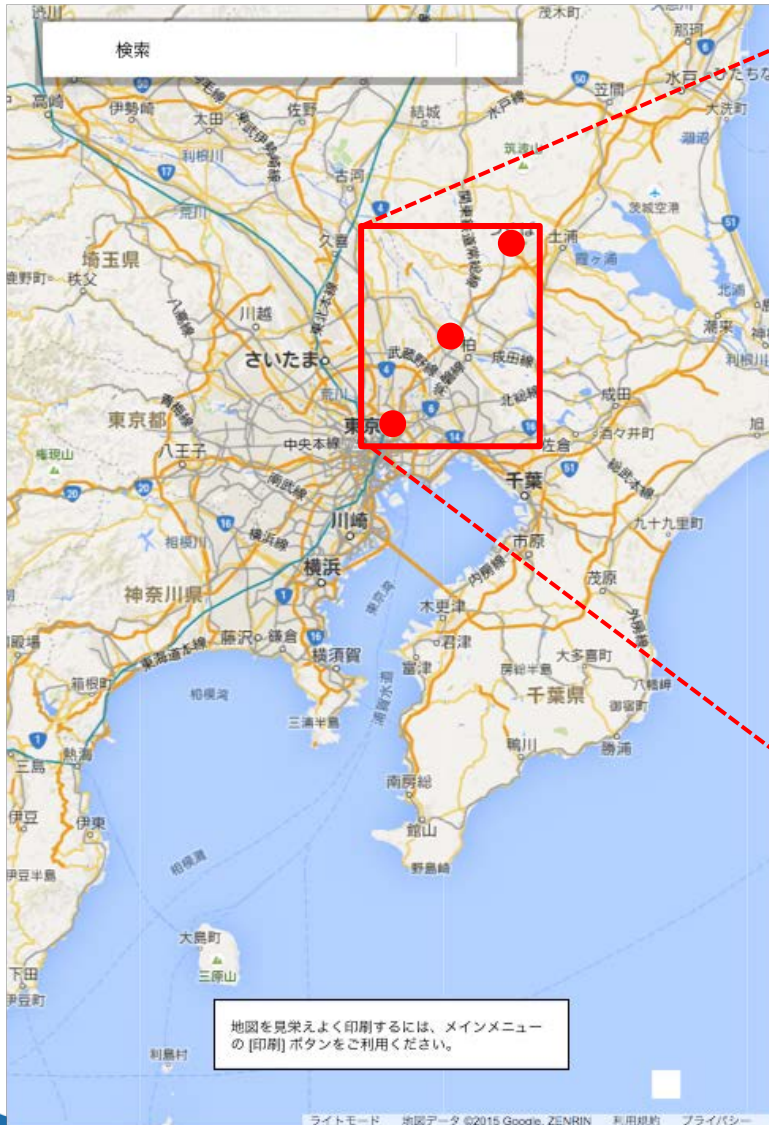
Oakforest-PACSの特徴

- 計算ノード
 - 1ノード 68コア, 3TFLOPS × 8,208ノード = 25 PFLOPS
 - メモリ(MCDRAM(高速, 16GB) + DDR4(低速, 96GB))
- ノード間通信
 - フルバイセクションバンド幅を持つFat-Treeネットワーク
 - 全系運用時のアプリケーション性能に効果, 多ジョブ運用
 - Intel Omni-Path Architecture
- ファイルI/O
 - 並列ファイルシステム
 - 高速ファイルキャッシュシステム(DDN IME) :> 1TB/sec
- 消費電力
 - Green 500でも世界6位
 - Linpack: 4,986 MFLOPS/W(OFP), 830 MFLOPS/W(京)

設置場所：東京大学柏キャンパス

Google マップ

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



筑波大学

東京大学
柏キャンパス

東京大学本郷キャンパス

東大@柏(Oak)
筑波大:PACS

地図を見栄えよく印刷するには、メインメニューの[印刷]ボタンをご利用ください。

ライトモード 地図データ ©2015 Google, ZENRIN 利用規約 プライバシー

Oakforest-PACS の仕様

総ピーク演算性能		25 PFLOPS	
ノード数		8,208	
計算 ノード	Product	富士通 PRIMERGY CX600 M1 (2U) + CX1640 M1 x 8node	
	プロセッサ	Intel® Xeon Phi™ 7250 (開発コード: Knights Landing) 68 コア、1.4 GHz	
	メモリ	高バンド幅	16 GB, MCDRAM, 実効 490 GB/sec
		低バンド幅	96 GB, DDR4-2400, ピーク 115.2 GB/sec
相互結 合網	Product	Intel® Omni-Path Architecture	
	リンク速度	100 Gbps	
	トポロジ	フルバイセクションバンド幅Fat-tree網	

計算ノードとシャーシ



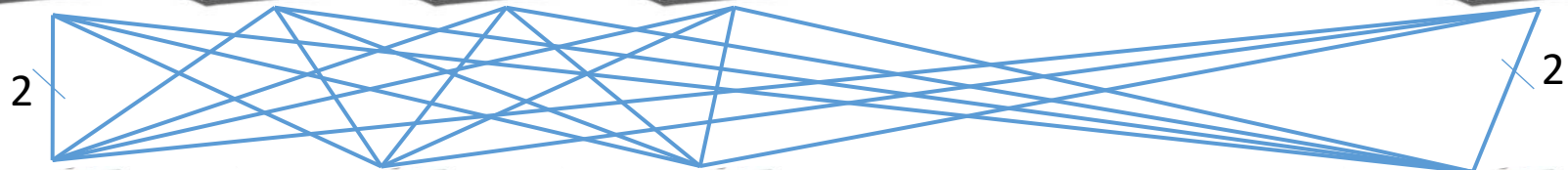
Chassis with 8 nodes, 2U size

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

Intel® Omni-Path Architecture を用いた フルバイセクションバンド幅Fat-tree網



768 port Director Switch
12台
(Source by Intel)



Uplink: 24



48 port Edge Switch
362台



Downlink: 24



コストはかかるがフルバイセクションバンド幅を維持

- システム全系使用時にも高い並列性能を実現
- 柔軟な運用: ジョブに対する計算ノード割り当ての自由度が高い

Oakforest-PACS の仕様 (続き)

並列ファイルシステム	Type	Lustre File System
	総容量	26.2 PB
	Product	DataDirect Networks SFA14KE
	総バンド幅	500 GB/sec
高速ファイルキャッシュシステム	Type	Burst Buffer, Infinite Memory Engine (by DDN)
	総容量	940 TB (NVMe SSD, パリティを含む)
	Product	DataDirect Networks IME14K
	総バンド幅	1,560 GB/sec
総消費電力		4.2MW (冷却を含む)
総ラック数		102

Oakforest-PACS のソフトウェア

- OS: Red Hat Enterprise Linux (ログインノード)、CentOS および McKernel (計算ノード、切替可能)
 - **McKernel**: 理研AICSで開発中のメニーコア向けOS
 - Linuxに比べ軽量、ユーザプログラムに与える影響なし
 - ポスト京コンピュータにも搭載される予定。
- コンパイラ: GCC, Intel Compiler, XcalableMP
 - **XcalableMP**: 理研AICSと筑波大で共同開発中の並列プログラミング言語
 - CやFortranで記述されたコードに指示文を加えることで、性能の高い並列アプリケーションを簡易に開発することができる。
- ライブラリ・アプリケーション: オープンソースソフトウェア
 - **ppOpen-HPC**, OpenFOAM, ABINIT-MP, PHASE system, FrontFlow/blue , LAPACK, ScaLAPACK, PETSc, METIS, SuperLU etc.

各種ベンチマーク

- TOP 500 (Linpack, HPL)
 - 連立一次方程式ソルバー(直接法), 計算速度 (FLOPS値)
 - 規則的な密行列: 連続メモリアクセス
 - 計算性能
- HPCG
 - 連立一次方程式ソルバー(反復法), 計算速度 (FLOPS値)
 - 有限要素法から得られる疎行列 (ゼロが多い)
 - 不連続メモリアクセス
 - 実アプリケーションに近い
 - メモリアクセス性能, 通信性能
- Green 500
 - HPL (TOP500) 実行時の FLOPS/W 値

48th TOP500 List (November, 2016)

	Site	Computer/Year Vendor	Cores	R _{max} (TFLOPS)	R _{peak} (TFLOPS)	Power (kW)
1	National Supercomputing Center in Wuxi, China	Sunway TaihuLight , Sunway MPP, Sunway SW26010 260C 1.45GHz, 2016 NRCPC	10,649,600	93,015 (= 93.0 PF)	125,436	15,371
2	National Supercomputing Center in Tianjin, China	Tianhe-2 , Intel Xeon E5-2692, TH Express-2, Xeon Phi, 2013 NUDT	3,120,000	33,863 (= 33.9 PF)	54,902	17,808
3	Oak Ridge National Laboratory, USA	Titan Cray XK7/NVIDIA K20x, 2012 Cray	560,640	17,590	27,113	8,209
4	Lawrence Livermore National Laboratory, USA	Sequoia BlueGene/Q, 2011 IBM	1,572,864	17,173	20,133	7,890
5	DOE/SC/LBNL/NERSC USA	Cori , Cray XC40, Intel Xeon Phi 7250 68C 1.4GHz, Cray Aries, 2016 Cray	632,400	14,015	27,881	3,939
6	Joint Center for Advanced High Performance Computing, Japan	Oakforest-PACS , PRIMERGY CX600 M1, Intel Xeon Phi Processor 7250 68C 1.4GHz, Intel Omni-Path, 2016 Fujitsu	557,056	13,555	24,914	2,719
7	RIKEN AICS, Japan	K computer , SPARC64 VIIIfx, 2011 Fujitsu	705,024	10,510	11,280	12,660
8	Swiss Natl. Supercomputer Center, Switzerland	Piz Daint Cray XC30/NVIDIA P100, 2013 Cray	206,720	9,779	15,988	1,312
9	Argonne National Laboratory, USA	Mira BlueGene/Q, 2012 IBM	786,432	8,587	10,066	3,945
10	DOE/NNSA/LANL/SNL, USA	Trinity , Cray XC40, Xeon E5-2698v3 16C 2.3GHz, 2016 Cray	301,056	8,101	11,079	4,233

R_{max}: Performance of Linpack (TFLOPS)

R_{peak}: Peak Performance (TFLOPS), Power: kW

HPCG Ranking (SC16, November, 2016)

	Site	Computer	Cores	HPL Rmax (Pflop/s)	TOP500 Rank	HPCG (Pflop/s)	HPCG/ HPL (%)
1	RIKEN AICS, Japan	K computer	705,024	10.510	7	0.6027	5.73
2	NSCC / Guangzhou, China	Tianhe-2	3,120,000	33.863	2	0.5800	1.71
3	JCAHPC, Japan	Oakforest-PACS	557,056	13.555	6	0.3855	2.84
4	National Supercomputing Center in Wuxi, China	Sunway TaihuLight	10,649,600	93.015	1	0.3712	.399
5	DOE/SC/LBNL/NERSC USA	Cori	632,400	13.832	5	0.3554	2.57
6	DOE/NNSA/LLNL, USA	Sequoia	1,572,864	17.173	4	0.3304	1.92
7	DOE/SC/Oak Ridge National Laboratory, USA	Titan	560,640	17.590	3	0.3223	1.83
8	DOE/NNSA/LANL/SNL, USA	Trinity	301,056	8.101	10	0.1826	2.25
9	NASA / Mountain View, USA	Pleiades: SGI ICE X	243,008	5.952	13	0.1752	2.94
10	DOE/SC/Argonne National Laboratory, USA	Mira: IBM BlueGene/Q,	786,432	8.587	9	0.1670	1.94

Green 500 Ranking (SC16, November, 2016)

	Site	Computer	CPU	HPL Rmax (Pflop/s)	TOP500 Rank	Power (MW)	GFLOPS/W
1	NVIDIA Corporation	DGX SATURNV	NVIDIA DGX-1, Xeon E5-2698v4 20C 2.2GHz, Infiniband EDR, NVIDIA Tesla P100	3.307	28	0.350	9.462
2	Swiss National Supercomputing Centre (CSCS)	Piz Daint	Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100	9.779	8	1.312	7.454
3	RIKEN ACCS	Shoubu	ZettaScaler-1.6 etc.	1.001	116	0.150	6.674
4	National SC Center in Wuxi	Sunway TaihuLight	Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway	93.01	1	15.37	6.051
5	SFB/TR55 at Fujitsu Tech. Solutions GmbH	QPACE3	PRIMERGY CX1640 M1, Intel Xeon Phi 7210 64C 1.3GHz , Intel Omni-Path	0.447	375	0.077	5.806
6	JCAHPC	Oakforest-PACS	PRIMERGY CX1640 M1, Intel Xeon Phi 7250 68C 1.4GHz , Intel Omni-Path	1.355	6	2.719	4.986
7	DOE/SC/Argonne National Lab.	Theta	Cray XC40, Intel Xeon Phi 7230 64C 1.3GHz , Aries interconnect	5.096	18	1.087	4.688
8	Stanford Research Computing Center	XStream	Cray CS-Storm, Intel Xeon E5-2680v2 10C 2.8GHz, Infiniband FDR, Nvidia K80	0.781	162	0.190	4.112
9	ACCMS, Kyoto University	Camphor 2	Cray XC40, Intel Xeon Phi 7250 68C 1.4GHz , Aries interconnect	3.057	33	0.748	4.087
10	Jefferson Natl. Accel. Facility	SciPhi XVI	KOI Cluster, Intel Xeon Phi 7230 64C 1.3GHz , Intel Omni-Path	0.426	397	0.111	3.837

運用

- 2017年度3月末までは無料(但し停止期間等あり)
- 計算資源は全系を共用(パーティション分けはしない)
 - 全8,208ノード(25PF)を常に全系で運用できるようにしておき、国内最大の計算資源を有効に活用する

○利用形態

- 各大学独自の利用コース
- HPCI
 - 全資源の20%を「JCAHPC」として拠出, 企業利用可能
- JHPCN(学際大規模情報基盤共同利用共同研究拠点)
 - 全資源の5%程度: 企業共同研究, 国際共同研究も含む(東大のみ)
- 教育(講義, 講習会)
- 大規模HPCチャレンジ: 全ノード占有

HPCIへの資源提供

- 平成29年度課題募集におけるハードウェア資源一覧
 - http://www.hpci-office.jp/pages/h29_boshu_hpci_resource?parent_folder=23

筑波大学 計算科学研究センター	COMA(PACS-IX) ▶ 資源提供元	Xeon +Phi(KNC)	計算ノード：90ノード(約230TFLOPS) 1,800コア (+メニーコアアーキテクチャ 10,980コア) 資源量：756,000ノード時間積 ストレージ：300TB
最先端共同HPC基盤施設 (JCAHPC)※	Oakforest-PACS ▶ 資源提供元	Xeon Phi(KNL)	計算ノード：1,600ノード(4,872 TFLOPS) 108,800コア 資源量：13,824,000ノード時間積 ストレージ：3,000TB
東京大学 情報基盤センター	スーパーコンピュータ FX10 ▶ 資源提供元	SPARC	計算ノード：1,500ノード (354.78TFLOPS) 24,000コア 資源量：12,960,000ノード時間積 ストレージ：500TB
	Reedbush-U ▶ 資源提供元	Xeon	計算ノード：60ノード(72.58TFLOPS) 2,160コア 資源量：518,400ノード時間積 ストレージ：60TB
	Reedbush-H ▶ 資源提供元	Xeon +Tesla	計算ノード：18ノード (193.11-212.73TFLOPS) コア数 648 + GPU 36枚 資源量：155,520ノード時間積 ストレージ：18TB

おわりに

- JCAHPC(最先端共同HPC基盤施設)
- 筑波大学計算科学研究センターと東京大学情報基盤センターが設置
 - 計算科学・工学及びその推進のための計算機科学・工学の発展に資するために連携して設置
- Oakforest-PACS:ピーク性能 25 PFLOPS
 - Intel Xeon Phi (Knights Landing) と Omni-Path Architecture
 - CPU時間を2大学で按分することで柔軟な運用を可能
 - 全系を1システムとして超大規模単一ジョブの実行も可能に
 - 全系システムの稼働は2016/12を予定
 - HPCI資源を含めオープンな資源提供は2017/4を予定
- JCAHPC:最先端HPC研究に寄与する計算資源の提供を目指し、コミュニティに貢献していく予定