

第121回 お試しアカウント付き 並列プログラミング講習会 「Oakforest-PACS実践」

東京大学情報基盤センター

内容に関するご質問は
hanawa @ cc.u-tokyo.ac.jp
まで、お願いします。

講習会概略

- 開催日：
2019年 7月19日（金） 13：30 - 18：00
- 場所：東京大学情報基盤センター 4階 413遠隔会議室
- 講習会プログラム：講師：塙 敏博
 - 13：00 - 13：30 受付
 - 13：30 - 14：30 Oakforest-PACSログイン、Oakforest-PACSシステム紹介
 - (14：30 - 14：45 休憩)
 - 14：45 - 16：15 KNL概要、KNLにおけるOpenMP並列化（講義＋演習）
 - (16：15 - 16：30 休憩)
 - 16：30 - 18：00 Oakforest-PACSでのMPI並列化、Oakforest-PACSでのハイブリッド並列化、性能分析（講義+演習）

Oakforest-PACS利用の準備

配布資料「Oakforest-PACS利用の手引き」参照

ユーザアカウント

- 本講習会でのユーザ名

利用者番号： t00171～

利用グループ： gt00

8/19 9:00まで有効

東大スーパーコンピュータ の概略

東大センターのスパコン

FY 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

2基の大型システム, 6年サイクル(?)

**Yayoi: Hitachi SR16000/M1
IBM Power-7**
54.9 TFLOPS, 11.2 TB

JCAHPC:
Tsukuba, Tokyo

**T2K Tokyo
Hitachi**
140TF, 31.3TB

Oakforest-PACS (OFP)
Fujitsu, Intel Xeon Phi
25PFLOPS, 919.3TB

Big Data & Extreme
Computing

**Oakleaf-FX: Fujitsu PRIMEHPC
FX10, SPARC64 IXfx**
1.13 PFLOPS, 150 TB

BDEC System
60+ PFLOPS

Oakbridge-FX
136.2 TFLOPS, 18.4 TB

Oakbridge-CX
Intel Xeon CLX
6.61 PFLOPS

Integrated Supercomputer
System for Data Analyses &
Scientific Simulations

Reedbush-U/H, HPE
Broadwell + Pascal
1.93 PFLOPS

Exascale System
(2021 Fall ?)

Supercomputer System with
Accelerators for Long-Term
Executions

Reedbush-L
HPE
1.43 PFLOPS



3システム：利用者2,000+, 学外50+%

- Reedbush (SGI, Intel BDW + NVIDIA P100 (Pascal))
 - データ解析・シミュレーション融合スーパーコンピュータ
 - 3.36 PF, 2016年7月～2021年3月末（予定）
 - 東大ITC初のGPUシステム (2017年3月より), DDN IME (Burst Buffer)
- Oakforest-PACS (OFP) (富士通, Intel Xeon Phi (KNL))
 - JCAHPC (筑波大CCS&東大ITC)
 - 25 PF, TOP 500で6位 (2016年11月) (日本1位) (初登場時)
 - Omni-Path アーキテクチャ, DDN IME (Burst Buffer)
- Oakbridge-CX (富士通, Intel Xeon Platinum 8280)
 - 大規模超並列スーパーコンピュータシステム
 - 6.61 PF, 2019年7月～2023年6月
 - 全1,368ノードの内128ノードにSSDを搭載



スーパーコンピュータシステムの詳細

- 以下のページをご参照ください
 - 利用申請方法
 - 運営体系
 - 料金体系
 - 利用の手引などがご覧になれます。

<https://www.cc.u-tokyo.ac.jp/guide>

Oakforest-PACSの紹介

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 16位 (国内2位), HPCG 9位 (国内2位) (2019年6月)

Oakforest-PACSの特徴 (1 / 2)

• 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



Oakforest-PACS の特徴 (2 / 2)

- ファイルI/O
 - 並列ファイルシステム:
Lustre 26PB
 - ファイルキャッシュシステム
(DDN IME) :
1TB/secを超える実効性能,
約1PB
 - 計算科学・ビッグデータ解析・
機械学習にも貢献
- 消費電力
 - Green 500でも世界6位
 - Linpack : 2.72 MW
 - 4,986 MFLOPS/W (OFP)
 - 830 MFLOPS/W (京)



並列ファイル
システム

ファイルキャッシュ
システム



ラック当たり120ノード
の高密度実装

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	

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

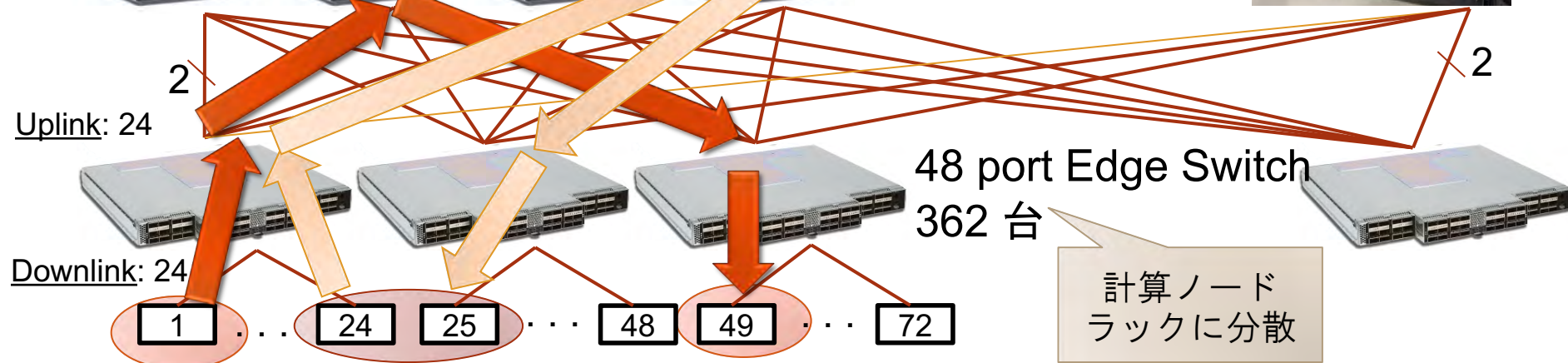
Oakforest-PACS のソフトウェア

- 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, XscalableMP
 - **XscalableMP**: 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.

Oakforest-PACS: Intel Omni-Path Architecture による 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)



コストはかかるがフルバイセクションバンド幅を維持

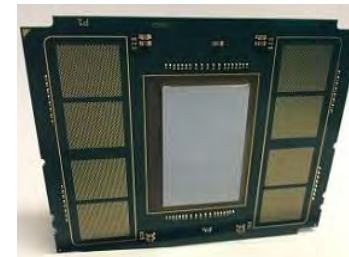
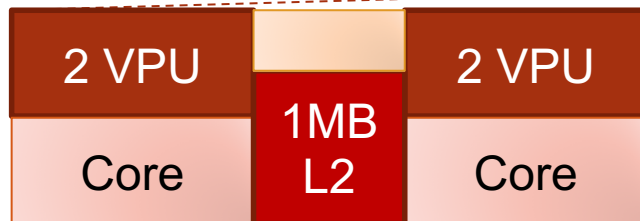
- システム全系使用時にも高い並列性能を実現
- 柔軟な運用: ジョブに対する計算ノード割り当ての自由度が高い

Oakforest-PACS 計算ノード

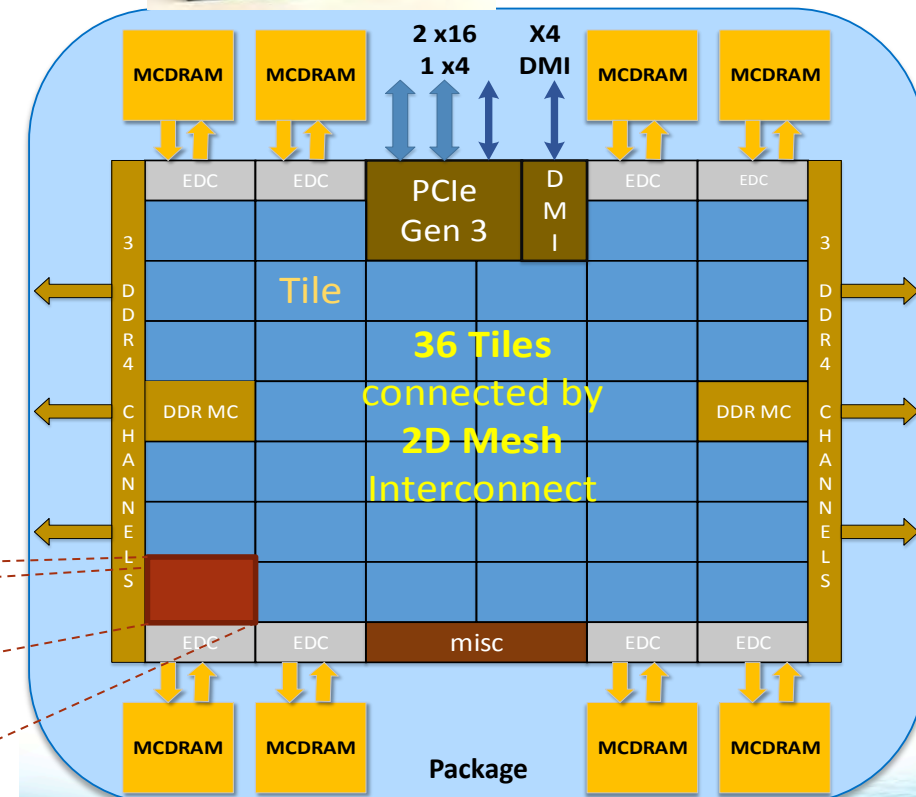
- Intel Xeon Phi (Knights Landing)
 - 1ノード1ソケット
- MCDRAM: オンパッケージの高バンド幅メモリ16GB + DDR4メモリ

ソケット当たりメモリ量:
16GB × 6 = 96GB

MCDRAM: 490GB/秒以上 (実測)
DDR4: 115.2 GB/秒
=(8Byte × 2400MHz × 6 channel)



HotChips27
KNLスライドより



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

<http://www.top500.org/>



R_{\max} : Performance of Linpack (TFLOPS)

 R_{peak} : Peak Performance (TFLOPS),

Power: kW

<http://www.top500.org/>

53rd TOP500 List (June, 2019)

	Site	Computer/Year Vendor	Cores	R_{\max} (TFLOPS)	R_{peak} (TFLOPS)	Power (kW)
1	<u>Summit, 2018, USA</u> DOE/SC/Oak Ridge National Laboratory	IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband	2,414,592	148,600 (= 148.6 PF)	200,795	10,096
2	<u>Sieera, 2018, USA</u> DOE/NNSA/LLNL	IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband	1,572,480	94,640	125,712	7,438
3	<u>Sunway TaihuLight, 2016, China</u> National Supercomputing Center in Wuxi	Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway	10,649,600	93,015	125,436	15,371
4	<u>Tianhe-2A, 2018, China</u> National Super Computer Center in Guangzhou	TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000	4,981,760	61,445	100,679	18,482
5	<u>Frontera, 2019, USA</u> Texas Advanced Computing Center	Dell C6420, Xeon Platinum 8280 28c 2.7GHz, Mellanox Infiniband HDR	448,448	23,516	38,746	
6	<u>Piz Daint, 2017, Switzerland</u> Swiss National Supercomputing Centre (CSCS)	Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100	387,872	21,230	27,154	2,384
7	<u>Trinity, 2017, USA</u> DOE/NNSA/LANL/SNL	Cray XC40, Intel Xeon Phi 7250 68C 1.4GHz, Aries interconnect	979,072	20,159	41,461	7,578
8	<u>ABCi (AI Bridging Cloud Infrastructure), 2018, Japan</u> National Institute of Advanced Industrial Science and Technology (AIST)	PRIMERGY CX2550 M4, Xeon Gold 6148 20C 2.4GHz, NVIDIA Tesla V100 SXM2, Infiniband EDR	391,680	19,880	32,577	1,649
9	<u>SuperMUC-NG, 2018, Germany</u> Leibniz Rechenzentrum	Lenovo, ThinkSystem SD650, Xeon Platinum 8174 24C 3.1GHz, Intel Omni-Path	305,856	19,477	26,874	
10	<u>Lassen, 2019, USA</u> DOE/NNSA/LLNL	IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta V100, Dual-rail Mellanox EDR Infiniband	288,288	18,200	23,047	
16	<u>Oakforest-PACS, 2016, Japan</u> Joint Center for Advanced High Performance Computing	PRIMERGY CX1640 M1, Intel Xeon Phi 7250 68C 1.4GHz, Intel Omni-Path	556,104	13,556	24,913	2,719

(理論) ピーク性能

- OFPに搭載されているXeon Phi 7250のピーク性能を考える。
 - コア数: **68**コア
 - コア当たりのAVX-512ユニット: **2**
 - AVX-512ユニット当たりの同時演算数(倍精度): **8**
 - 積和演算(Fused Multiply Add: FMA): **2**に換算
 - クロック周波数: **1.40** GHz
- ノード当たりピーク性能:
 $68 * 2 * 8 * 2 * 1.40 = \mathbf{3046.4 \text{ GFLOPS}}$
- しかし、AVX-512ユニットは実は**1.40GHz**では動作しない
(より低い周波数)
 - ピークに近い性能が得られるはずのもの(OFPでの実測値)
 - DGEMM(倍精度の行列積): **2200** GFLOPS(ピーク比: 72%)
 - HPL: **2000** GFLOPS(ピーク比: 66%)
 - Top500におけるOFPの登録値はピークの54.4%
- 「**ピーク性能**」の定義を正しく把握しておくことが重要!!
 - **CPUメーカーによって定義の仕方も異なる**

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

HPCG Ranking (June, 2019)

	Computer	Cores	HPL Rmax (Pflop/s)	TOP500 Rank	HPCG (Pflop/s)
1	Summit	2,414,592	148,600	1	2.926
2	Sierra	1,572,480	94.640	2	1.796
3	K computer	705,024	10.510	20	0.603
4	Trinity	979,072	20,159	7	0.546
5	ABCI	391,680	19,880	8	0.509
6	Piz Daint	387,872	21.230	6	0.497
7	Sunway TaihuLight	10,649,600	93.015	3	0.481
8	Nurion (KISTI, Korea)	570,020	13.929	15	0.391
9	Oakforest-PACS	556,104	13.555	16	0.385
10	Cori (NERSC/LBNL, USA)	632,400	14.015	14	0.355

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

Green 500 Ranking (June, 2019)

<http://www.top500.org/>

	TOP 500 Rank	System	Cores	HPL Rmax (Pflop/s)	Power (MW)	GFLOPS/W
1	472	Shoubu system B, Japan	953,280	1,063	60	17.604
2	470	DGX SaturnV Volta, USA	22,440	1,070	97	15.113
3	1	Summit, USA	2,414,592	148,600	10,096	14.719
4	8	ABCI, Japan	391,680	19,880	1,649	14.423
5	394	MareNostrum P9 CTE, Spain	18,360	1,145	81	14.131
6	25	TSUBAME 3.0, Japan	135,828	8,125	792	13.704
7	444	PANGEA III, France	291,024	17,860	1,367	13.065
8	2	Sierra, USA	1,572,480	94,640	7,438	12.723
9	43	Advanced Computing System (PreE), China	163,840	4,325	380	11.382
10	23	Taiwania 2, Taiwan	170,352	900	798	11.285
13	June'18	Reedbush-L, U.Tokyo, Japan	16,640	806	79	10.167
19		Reedbush-H, U.Tokyo, Japan	17,760	802	94	8.576

IO 500 Ranking (June, 2018)

	Site	Computer	File system	Client nodes	IO500 Score	BW (GiB/s)	MD (kIOP/s)
1	JCAHPC, Japan	Oakforest- PACS	DDN IME	2048	137.78	560.10	33.89
2	KAUST, Saudi	Shaheen2	Cray DataWarp	1024	77.37	496.81	12.05
3	KAUST, Saudi	Shaheen2	Lustre	1000	41.00	54.17	31.03
4	JSC, Germany	JURON	BeeGFS	8	35.77	14.24	89.81
5	DKRZ, Germany	Mistral	Lustre2	100	32.15	22.77	45.39
6	IBM, USA	Sonasad	Spectrum Scale	10	24.24	4.57	128.61
7	Fraunhofer, Germany	Seislab	BeeGFS	24	16.96	5.13	56.14
8	DKRZ, Germany	Mistral	Lustre1	100	15.47	12.68	18.88
9	Joint Institute for Nuclear Research	Govorun	Lustre	24	12.08	3.34	43.65
10	PNNL, USA	EMSL Cascade	Lustre	126	11.12	4.88	25.33

<http://www.io500.org/>



IO 500 Ranking (June, 2019)

<http://www.io500.org/>

	Site	Computer	File system	Client nodes/ procs	IO500 Score	BW (GiB/s)	MD (kIOP/s)
1	University of Cambridge, UK	Data Accelerator	Dell EMC Lustre	512 8192	620.69	162.05	2377.44
2	Oak Ridge National Laboratory, USA	Summit	IBM Spectrum Scale	504 1008	330.56	88.20	1238.93
3	JCAHPC, Japan	Oakforest-PACS	DDN IME	2048 2048	275.65	492.06	154.41
4	KISTI, Korea	NURION	DDN IME	2048 4096	156.91	554.23	44.43
5	CSIRO, Australia	bracewell	Dell/ThinkPar Q BeeGFS	26 260	140.58	69.29	285.21
6	DDN	IME140	DDN IME	17 272	112.67	90.34	140.52
7	DDN Colorado	DDN IME140	DDN IME	10 160	109.42	75.79	157.96
8	DDN	AI400	DDN Lustre	10 160	104.34	19.65	553.98
9	KAUST, Saudi	Shaheen2	Cray DataWarp	1024 8192	77.37	496.81	12.05
10	University of Cambridge, UK	Data Accelerator	Dell EMC BeeGFS	184 5888	74.58	58.81	94.57

運用

- 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）
 - 東大分全資源の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，1回24時間程度/月

東大情報基盤センターOakforest-PACSスーパーコンピュータシステムの料金表（2018年4月1日～）

・ パーソナルコース（年間）

- ・ 100,000円 : 1口8ノード(基準)、最大512ノードまで実行可
- ・ 3口まで

・ グループコース

- ・ 400,000円 (企業 480,000円) : 1口8ノード（基準）、最大2048ノードまで

・ 以上は、「トークン制」で運営

- ・ 申し込みノード数×360日×24時間の「トークン」が与えられる
 - ・ パーソナルコースは2ノード相当
- ・ 基準ノードまでは、トークン消費係数が1.0
- ・ 基準ノードを超えると、超えた分は、消費係数が2.0になる
- ・ 大学等のユーザはReedbushとの相互トークン移行も可能

GFLOPS (ピーク性能換算) あたり負担金 (~W)

System	JPY/GFLOPS
Reedbush-U (HPE) (Intel BDW)	61.9
Reedbush-H (HPE) (Intel BDW+NVIDIA P100x2/node)	15.9
Reedbush-L (HPE) (Intel BDW+NVIDIA P100x4/node)	13.4
Oakforest-PACS (Fujitsu) (Intel Xeon Phi/Knights Landing)	16.5
Oakbridge-CX (Fujitsu) (Intel Cascade Lake (CLX))	20.7

トライアルユース制度について

- 安価に当センターのReedbush-U/H, Oakforest-PACSシステムが使える「**無償トライアルユース**」および「**有償トライアルユース**」制度があります。
 - **アカデミック利用**
 - パーソナルコース、グループコースの双方（1ヶ月～3ヶ月）
 - **企業利用**
 - パーソナルコース（1ヶ月～3ヶ月）（RB-U: 最大16ノード, RB-H: 最大2ノード、OFP: 最大16ノード, 最大64ノード）
本講習会の受講が必須、審査無
 - グループコース
 - 無償トライアルユース：（1ヶ月～3ヶ月）：無料（RB-U: 最大128ノード、RB-H: 最大32ノード、OFP: 最大2048ノード）
 - 有償トライアルユース：（1ヶ月～最大通算9ヶ月）、有償（計算資源は無償と同等）
 - **スーパーコンピュータ利用資格者審査委員会の審査が必要（年2回実施）**
 - 双方のコースともに、簡易な利用報告書の提出が必要
- **トライアルユース(無料体験):**利用申請するシステムを初めて使う場合
 - 1ヶ月間、講習会アカウントと同条件
 - 一般利用、有償トライアルユースへの移行も可能

並列プログラミングの基礎

並列プログラミングとは何か？

- 逐次実行のプログラム（実行時間 T ）を、 p 台の計算機を使って、 T/p にすること。



- 素人考えでは自明。
- 実際は、できるかどうかは、対象処理の内容（アルゴリズム）で **大きく** 難しさが違う
 - アルゴリズム上、絶対に並列化できない部分の存在
 - 通信のためのオーバヘッドの存在
 - 通信立ち上がり時間
 - データ転送時間

並列プログラム中の実行主体

- マルチスレッド
 - OpenMP
 - ユーザが並列化指示行を記述
- マルチプロセス
 - MPI (Message Passing Interface)
 - ユーザがデータ分割方法を明示的に記述

共有メモリで動作
→並列化の指示だけでOK

- ノードが物理的に異なる
- 同一ノードでもプロセス毎にメモリは論理的に分離
→明示的な通信が必要

マルチプロセスとマルチスレッドは
共存可能
→ハイブリッドMPI/OpenMP実行

並列プログラミングのモデル

- 多くのMIMD上での並列プログラミングのモデル

1. SPMD (Single Program Multiple Data)

- 1つの共通のプログラムが、並列処理開始時に、**全プロセッサ上で同時に**起動する

• MPI (バージョン1) のモデル



2. Master / Worker (Master / Slave)

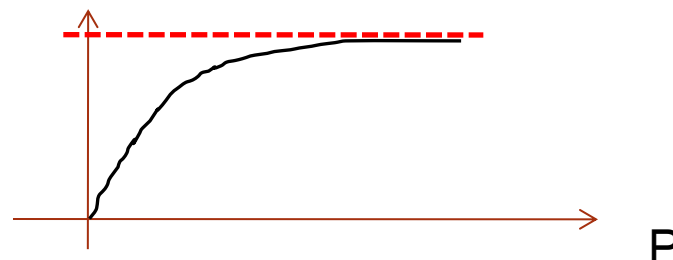
- 1つのプロセス (Master) が、複数のプロセス (Worker) を管理 (生成、消去) する。

性能評価指標

並列化の尺度

性能評価指標—台数効果

- 台数効果 $S_p = T_S / T_p$ ($0 \leq S_p$)
 - 式:
 - T_S : 逐次の実行時間、 T_p : P台での実行時間
 - P台用いて $S_p = P$ のとき、理想的な (ideal) 速度向上
 - P台用いて $S_p > P$ のとき、スーパーニア・スピードアップ
 - 主な原因は、並列化により、データアクセスが局所化されて、キャッシュヒット率が向上することによる高速化
- 並列化効率
 - 式: $E_p = S_p / P \times 100$ ($0 \leq E_p$) [%]
- 飽和性能
 - 速度向上の限界
 - Saturation, 「さちる」



並列性を上げることが鍵

- 京コンピュータ : 88,128ノード x 8コア/ノード
- Oakforest-PACS : 8,208ノード x 68コア/ノード
- Sunway TaihuLight: 40,960ノード x 260コア/ノード
- Summit : 約4,600ノード x 6 GPU/ノード
x (80 or 5120 コア)/GPU
- 1コア当たりではスパコンの方が遅いこともある
 - スパコン向けは、電力効率を狙ってクロック周波数を下げている
=> その分コア数を増やす

Weak ScalingとStrong Scaling

並列処理においてシステム規模を大きくする方法

- Weak Scaling: それぞれの問題サイズは変えず並列度をあげる
 - 全体の問題サイズが（並列数に比例して）大きくなる
 - 通信のオーバーヘッドはあまり変わらないか、やや増加する
- Strong Scaling: 全体の問題サイズを変えずに並列度をあげる
 - 問題サイズが装置数に反比例して小さくなる
 - 通信のオーバーヘッドは相対的に大きくなる

Weak Scaling

それまで解けなかった
規模の問題が解ける



Strong Scaling も重要

同じ問題規模で、短時間に
結果を得る（より難しい）

アムダールの法則

- 逐次実行時間を K とする。
そのうち、並列化ができる割合を α とする。
- このとき、台数効果は以下のようになる。

$$S_p = K / (K\alpha / P + K(1-\alpha))$$
$$= 1 / (\alpha / P + (1-\alpha)) = 1 / (\alpha(1/P - 1) + 1)$$

- 上記の式から、たとえ無限大の数のプロセッサを使っても ($P \rightarrow \infty$)、台数効果は、高々 $1/(1-\alpha)$ である。

(アムダールの法則)

- 全体の90%が並列化できたとしても、無限大の数のプロセッサをつかっても、 $1 / (1 - 0.9) = 10$ 倍 にしかならない！

→高性能を達成するためには、少しでも並列化効率を上げる実装をすることがとても重要である

アムダールの法則の直観例

並列化できない部分(1ブロック) 並列化できる部分(8ブロック)

● 逐次実行



=88.8%が並列化可能

● 並列実行(4並列)



$9/3=3$ 倍

● 並列実行(8並列)



$9/2=4.5$ 倍 \neq 6倍

本講習会の目的

並列アルゴリズムはひとまず置いて...

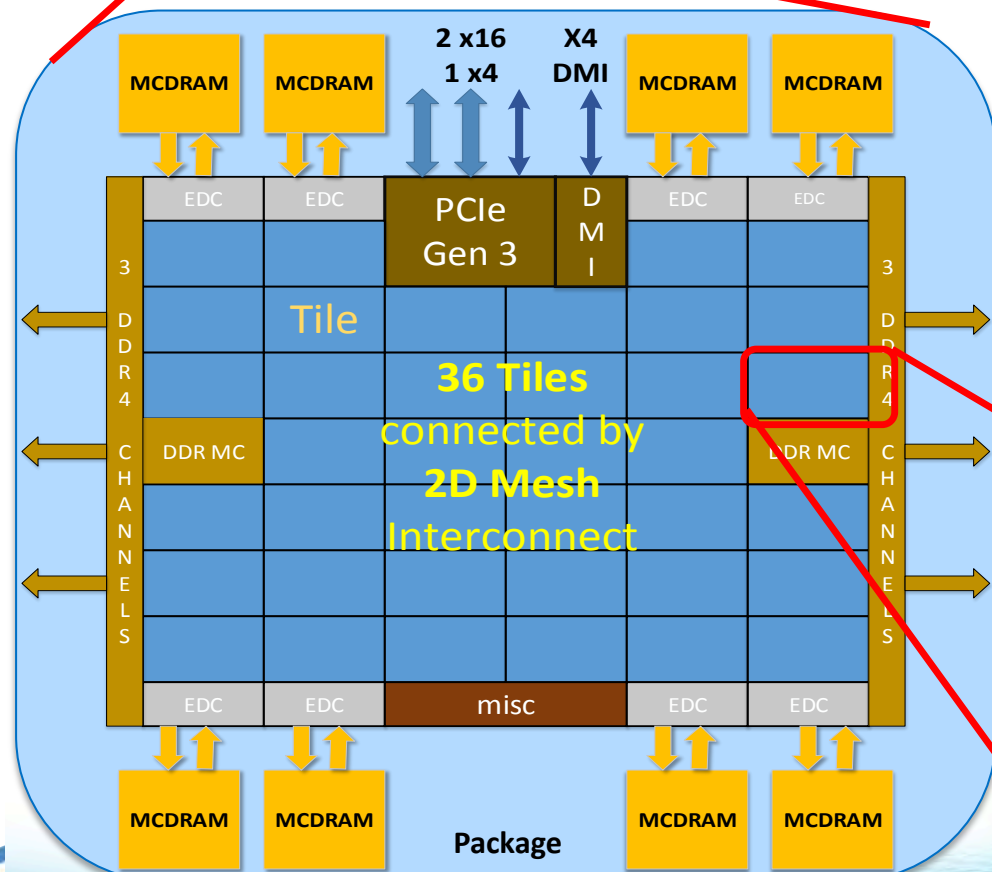
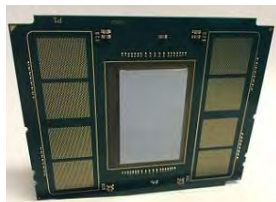
KNLの多数のコアの扱い方、Oakforest-PACSの多数の計算ノードを組み合わせた実行のやり方について学ぼう！

- その他の講習会 (予定) <https://www.cc.u-tokyo.ac.jp/events/lectures/>
 - 並列アルゴリズムを学びたい方
 - OpenMP/OpenACCによるマルチコア・メニィコア並列プログラミング入門
 - その他
 - 「並列有限要素法とハイブリッド並列プログラミング」
 - MPI 基礎
 - MPI 上級編
 - 科学技術計算の効率化入門

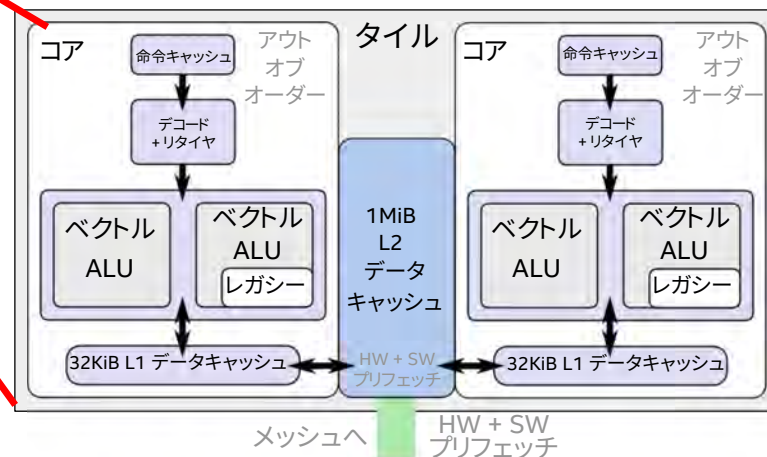
「こんな講習会が欲しい！」
という要望があればぜひアンケートへ！

KNL (Intel Xeon Phi, Knights Landing)の基礎

Intel Xeon Phi (KNL: Knights Landing)



- Atom (Silvermont) コア + AVX512 x2
- 2コア 1タイル
 - 各コアは 4スレッド (HyperThreading)
- 64, 68コア (32, 34タイル)
- **MCDRAM**: オンパッケージの高バンド幅メモリ16GB+ DDR4メモリ



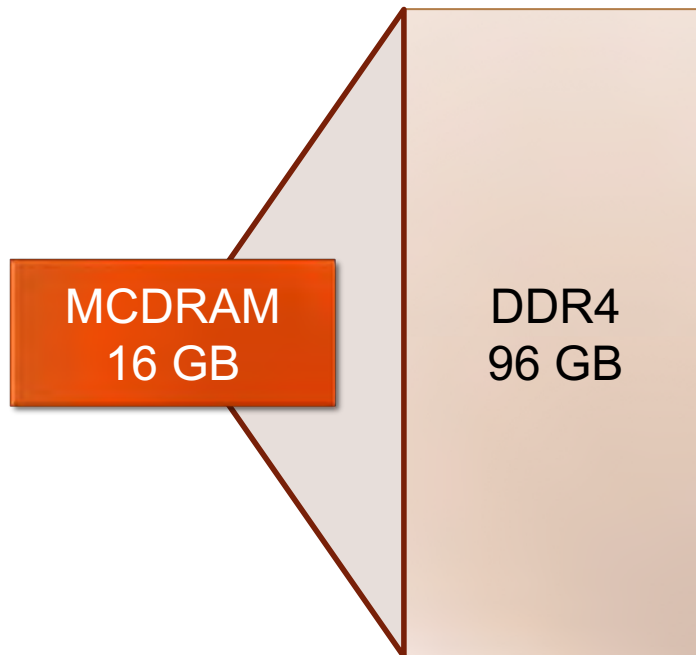
KNLの動作モード

- メモリモード：3種類
 - Flat: MCDRAMとDDR4 が独立したアドレス
 - Cache: MCDRAMはDDR4メモリのキャッシュとして動作
 - Hybrid
- クラスタリングモード: 5種類
 - (All-to-all: アドレス情報が全体に分散... 非推奨)
 - Quadrant, Hemisphere: 内部でアドレス情報が4(または2)に分割 (ユーザからは見えない)
 - SNC-4, SNC-2: NUMAドメインが明示的に4 (or 2)に分割

モードの変更には再起動が必要
=> 現時点では、各モード(Flat, Cache)のジョブ
キューを用意 (regular-flat / regular-cache等)

メモリモード

- Cacheモード
 - MCDRAMはL3キャッシュとして動作



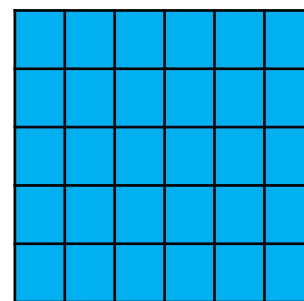
- Flatモード
 - MCDRAMを明示的に使い分け



FlatモードにおけるNUMAドメイン

- `$ numactl -H`
available: 2 nodes (0-1)
node 0 cpus: 0 1 2 3 4 5 6
... 271
node 0 size: 98164 MB
node 0 free: 92027 MB
node 1 cpus: (何もない)
node 1 size: 16384 MB
node 1 free: 15832 MB
node distances:
node 0 1
0: 10 31
1: 31 10

Node 0



CPU



DDR4

Node 1



MCDRAM

クラスタリングモード

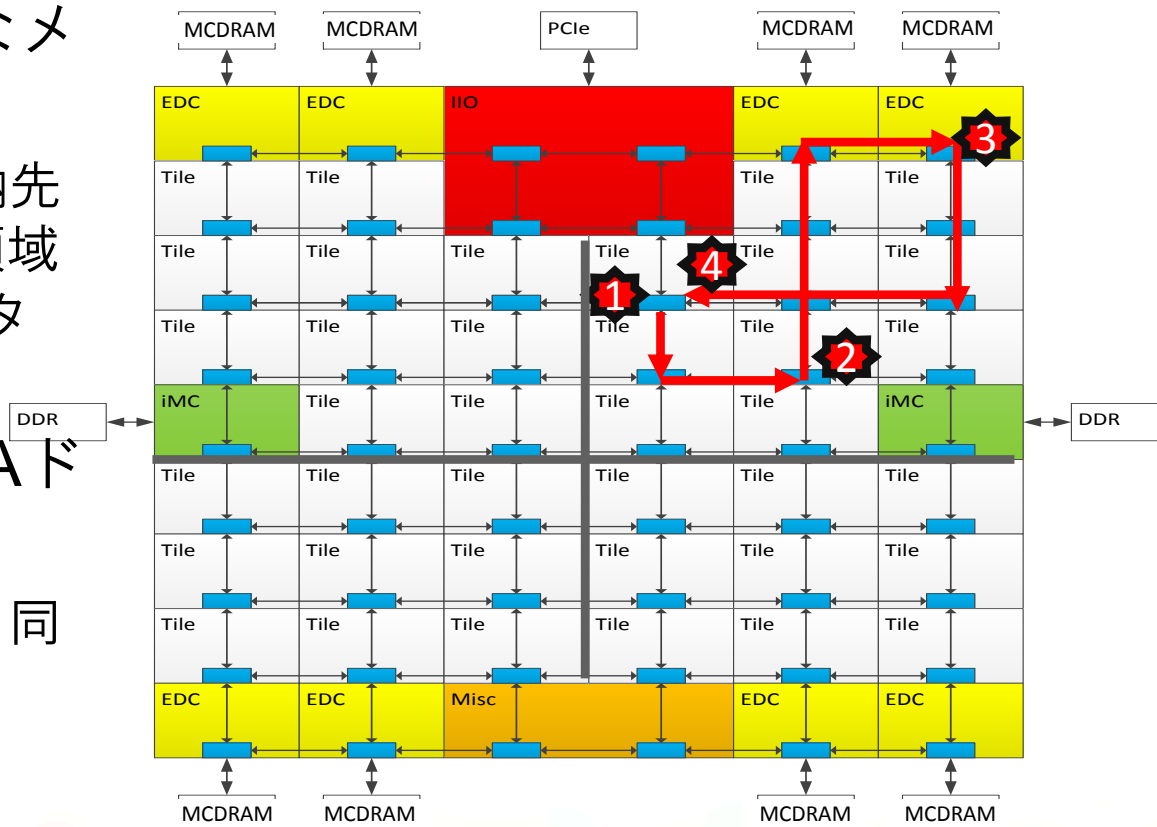
現在OFPではQuadrantのみ

- Quadrant: ソフトウェアからはフラットなメモリ空間に見える

- 内部では4分割、格納先のMCDRAMと同じ領域にキャッシュ情報（タグ）を配置

- SNC-4: 4つのNUMAドメインに見える

- 4ソケットある場合と同じ

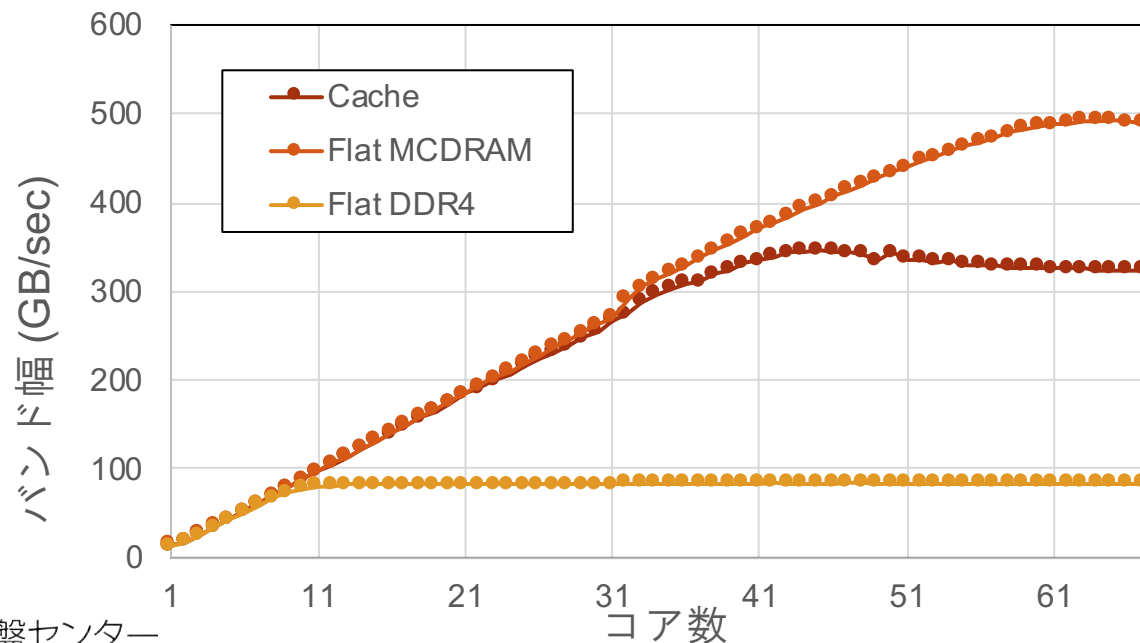


HotChips27
KNLスライドより

Stream Triad性能の比較

- DDR4: micprun -D -k stream
- MCDRAM: micprun -k stream
 - unset KMP_AFFINITY
 - export KMP_HW_SUBSET=66c@2,1t
タイル0を避ける設定（後述）

	DDR4	MCDRAM
Cache	346.5 GB/sec	
Flat	84.8 GB/sec	494.8 GB/sec



FlatモードでのMCDRAMの使い方(1)

(1) numactlによる方法

- MCDRAMはNUMAノード1として認識されている

```
$ numactl -H
available: 2 nodes (0-1)
node 0 cpus: 0 1 2 3 4 ... (省略) ... 268 269 270 271
node 0 size: 98164 MB
node 0 free: 85707 MB
node 1 cpus:
node 1 size: 16384 MB
node 1 free: 15496 MB
node distances:
node  0  1
  0: 10  31
  1: 31  10
```

- コード書き換えを避けたい場合には便利
× 配列毎の細かな使い分けはできない

- プログラム実行時に以下を指定

- \$ numactl --membind 1 ./a.out
(MCDRAM 16GBに入りきらないとエラー)
- \$ numactl --preferred 1 ./a.out
(MCDRAMに収まり切らない場合、DDR4も使用する)
- \$ numactl --interleave 1,0 ./a.out
(MCDRAMとDDR4を交互に使用する)

FlatモードでのMCDRAMの使い方(2)

(2) 環境変数による方法: Intel MPIを使う場合

- I_MPI_HBW_POLICY環境変数
 - export I_MPI_HBW_POLICY=hbw_bind
 - export I_MPI_HBW_POLICY=hbw_preferred
 - export I_MPI_HBW_POLICY=hbw_interleave
- 事実上numactlを呼ぶのと同じ

- コード書き換えを避けたい場合には便利
 - × 配列毎の細かな使い分けはできない
 - × Intel MPIのときしか有効にならない

FlatモードでのMCDRAMの使い方(3)

(3) memkindライブラリによる方法

• C

- `#include <hbwmalloc.h>` を指定
- 以下の通り、メモリ操作関連の関数名に”hbw_”をつける
 - `malloc()` => `hbw_malloc()`
 - `posix_memalign()` => `hbw_posix_memalign()`
 - `free()` => `hbw_free()`
- リンク時に `-lmemkind` を指定

• Fortran

- `Fastmem`の属性を指定
- `!dir$ attributes fastmem :: array`
- リンク時に `-lmemkind` を指定

• `hbw_set_policy()`で切り替え

- `HBW_POLICY_BIND`
- `HBW_POLICY_PREFERRED` (デフォルト)
- `HBW_POLICY_INTERLEAVE`

• 配列毎の細かな使い分けが可能
× 動的配列(allocatable)
にしか適用できない

FlatモードでのMCDRAMの使い方(4) 追加!

(4) memkind+autohbwライブラリによる方法

- ソースコードは変更しなくて良い
- リンク時に `-lautohbw -lmemkind` を指定、または実行時に `LD_PRELOAD=libautohbw.so` を指定
- 環境変数でMCDRAMサイズを指定
 - `export AUTO_HBW_SIZE=4K`
 - 4KB以上の配列はMCDRAM
 - `export AUTO_HBW_SIZE=4K:2M`
 - 4K~2Mの配列をMCDRAMに置く

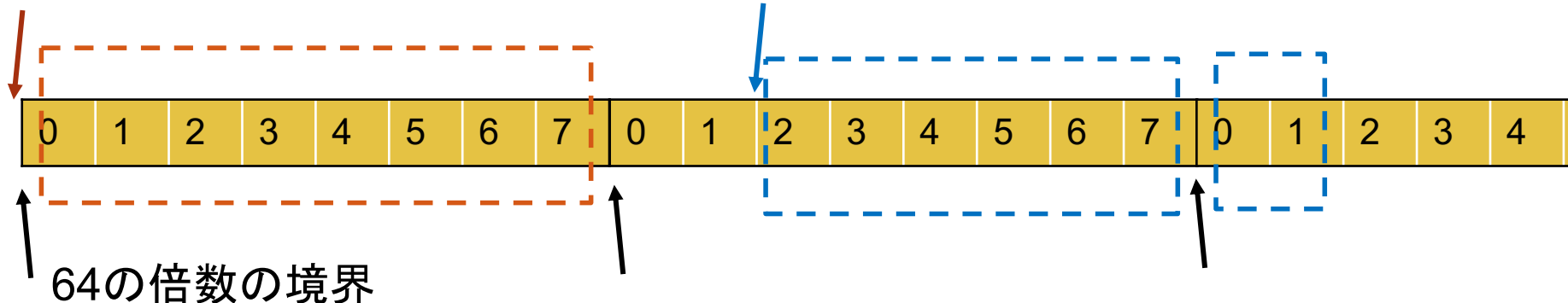
× 動的配列(allocatable)
にしか適用できない
× 配列毎の細かな使い分け
はできない

メモリのアラインメント

- AVX-512命令で512ビットデータを有効に使うためには、データがメモリ境界に揃っていることが必要
 - 512ビット = 64バイトの倍数のアドレスにアクセスするように
- 境界以外では複数回のアクセスになる
 - 例: 倍精度実数(double)を連続8個(=512ビット)読み出す場合

アラインされたアクセス
=効率がいい

境界からずれているため、
実際には複数回アクセスになる



メモリアライメントの指定

• C言語

- `__attribute__((aligned(64)))` を指定 (64 byte = 512 bit)
 - 例: `double A[1024] __attribute__((aligned(64)))`
 - 動的確保の際は `posix_memalign()` を使用
 - 上と同様の例:

```
#include <stdlib.h>
double *A;
posix_memalign(&A, 64, 8192);
```

• Fortran

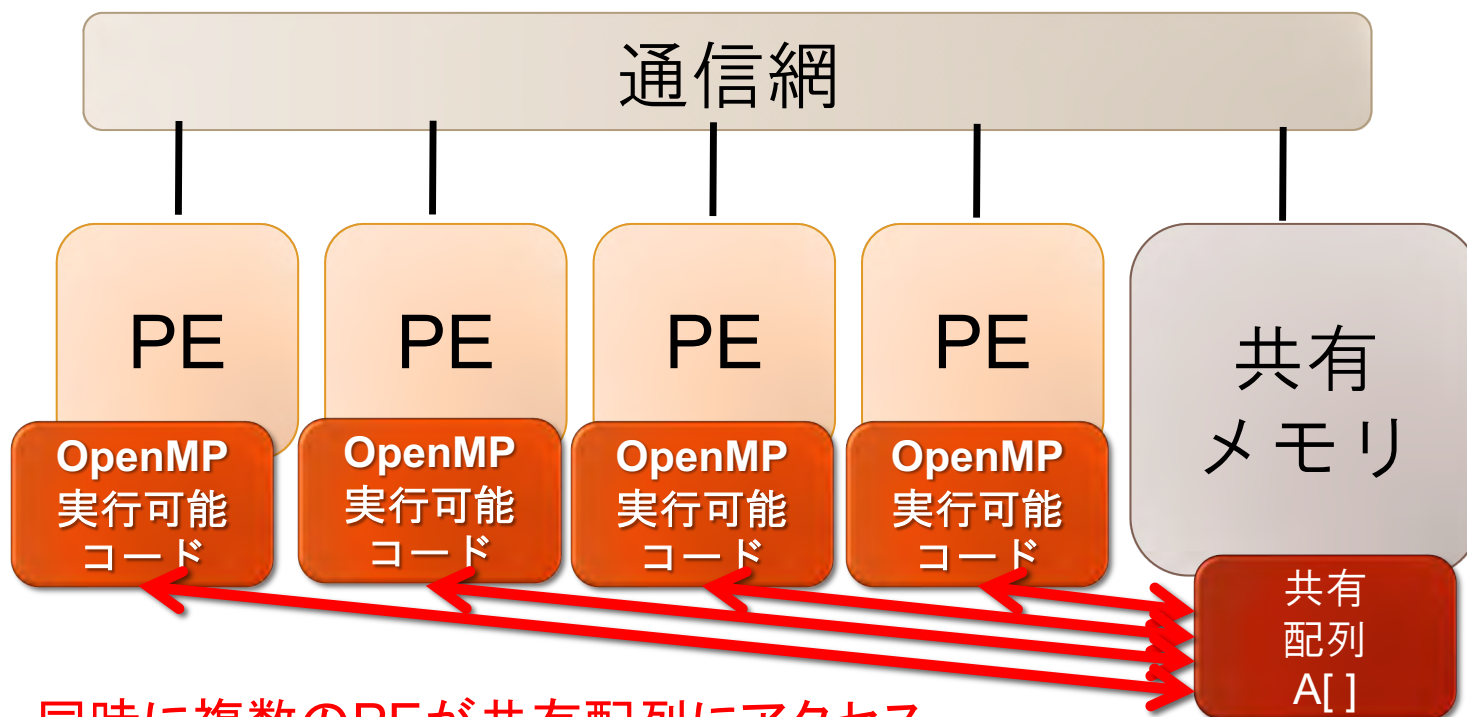
- コンパイル時に `-align array64byte` オプションを指定
- 動的な場合は `attribute` 指定が必要かもしれない(経験上)

```
real*8, dimension(:), allocatable :: A
!dir$ attributes align:64 :: A
```

OpenMPの概要

OpenMPの対象計算機

- OpenMPは共有メモリ計算機のためのプログラム言語



同時に複数のPEが共有配列にアクセス

⇒ 並列処理で適切に制御をしないと、逐次計算の結果と一致しない

OpenMPとは

- **OpenMP (OpenMP Application Program Interface)**とは、共有メモリ型並列計算機用にプログラムを並列化するための：
 1. 指示文
 2. ライブラリ
 3. 環境変数を規格化したものです。
- ユーザが、並列プログラムの実行させるための指示を与えるものです。コンパイラによる自動並列化ではありません。
- 分散メモリ型並列化（MPIなど）に比べて、データ分散の処理の手間が無い分、実装が簡単です。

OpenMPとマルチコア・メニーコアシステム

- スレッド並列化を行うプログラミングモデル
- マルチコア、メニーコア計算機に適合
 - 多数のスレッド実行で高い並列化効率を確保するには、プログラミングの工夫が必要
 1. メインメモリ-キャッシュ間のデータ転送能力が演算性能に比べ低い
 2. OpenMPで並列性を抽出できないプログラムになっている（後述）
- Oakforest-PACS (Intel Xeon Phi 7250, Knights Landing)
 - 68物理コア、272論理コア利用可能
 - プログラム上の工夫が必要 => スレッド数を増やしすぎないためにはMPI+OpenMPハイブリッド実行も解の一つ
- ノード間の並列化はOpenMPではできない
 - ノード間の並列化はMPIを用いる

OpenMPコードの書き方の原則

- C言語の場合
 - `#pragma omp`
で始まるコメント行
- Fortran言語の場合
 - `!$omp`
で始まるコメント行

OpenMPのコンパイルの仕方

- 逐次コンパイラのコンパイルオプションに、OpenMP用のオプションを付ける
 - 例) Intel Fortran90コンパイラ
`ifort -O3 -qopenmp foo.f90`
 - 例) Intel Cコンパイラ
`icc -O3 -qopenmp foo.c`
- 注意
 - **OpenMPの指示がないループは逐次実行**
 - コンパイラにより、自動並列化によるスレッド並列化との併用ができる場合があるが、できない場合もある
 - OpenMPの指示行がある行はOpenMPによるスレッド並列化、指示がないところはコンパイラによる自動並列化
 - 例) Intel Fortran90コンパイラ
`ifort -O3 -qparallel -qopenmp foo.f90`

OpenMPの実行可能ファイルの実行

- OpenMPのプログラムをコンパイルして生成した実行可能ファイルの実行は、そのファイルを指定することで行う
- スレッド数を、環境変数**OMP_NUM_THREADS**で指定
- 例) OpenMPによる実行可能ファイルがa.outの場合

```
$ export OMP_NUM_THREADS=16
```

```
$ ./a.out
```

または

```
$ env OMP_NUM_THREADS=16 ./a.out
```

- **注意**

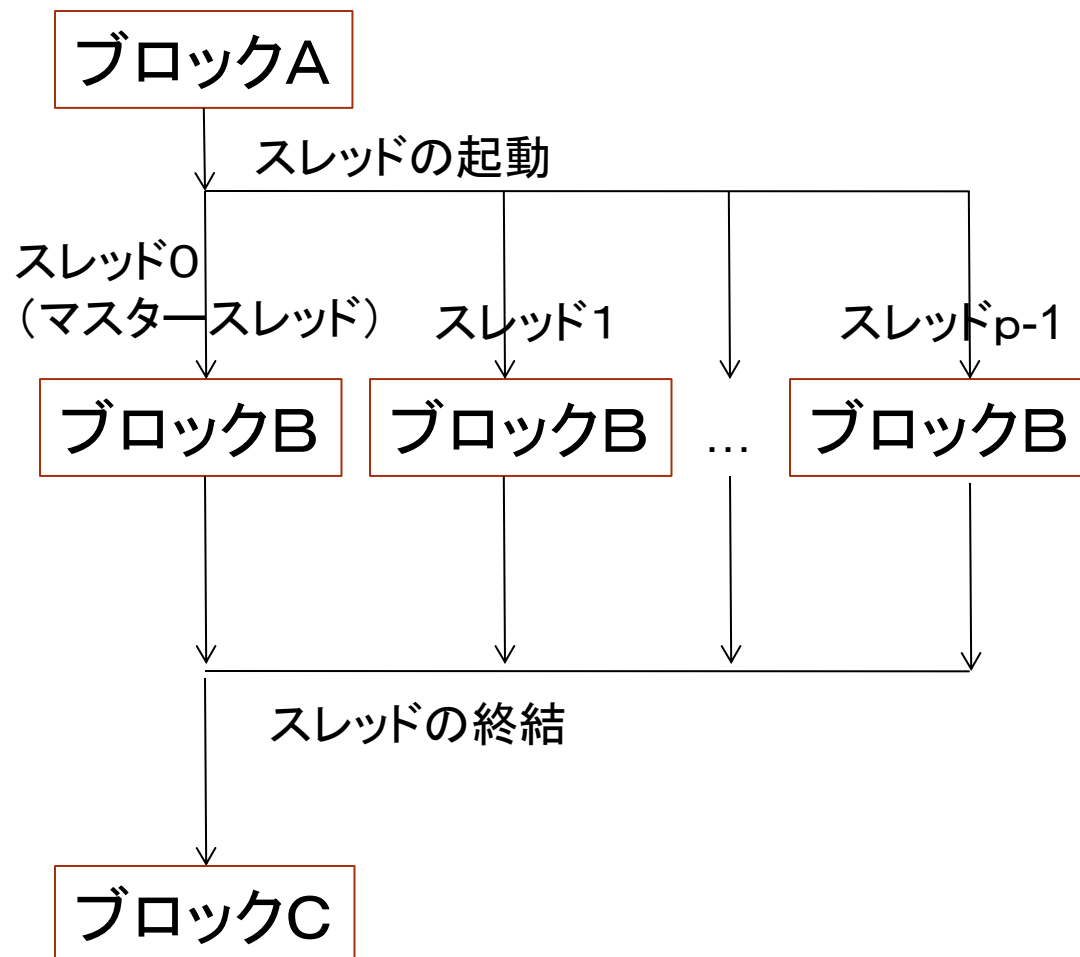
- 逐次コンパイルのプログラムと、OpenMPによるプログラムの実行速度が、OMP_NUM_THREADS=1にしても、異なることがある（後述）
 - この原因は、OpenMP化による処理の増加（オーバーヘッド）
 - 高スレッド実行で、このオーバーヘッドによる速度低下が顕著化
 - プログラミングの工夫で改善可能

OpenMPの実行モデル

OpenMPの実行モデル (C言語)

OpenMP指示文

```
ブロックA
#pragma omp parallel
{
  ブロックB
}
ブロックC
```

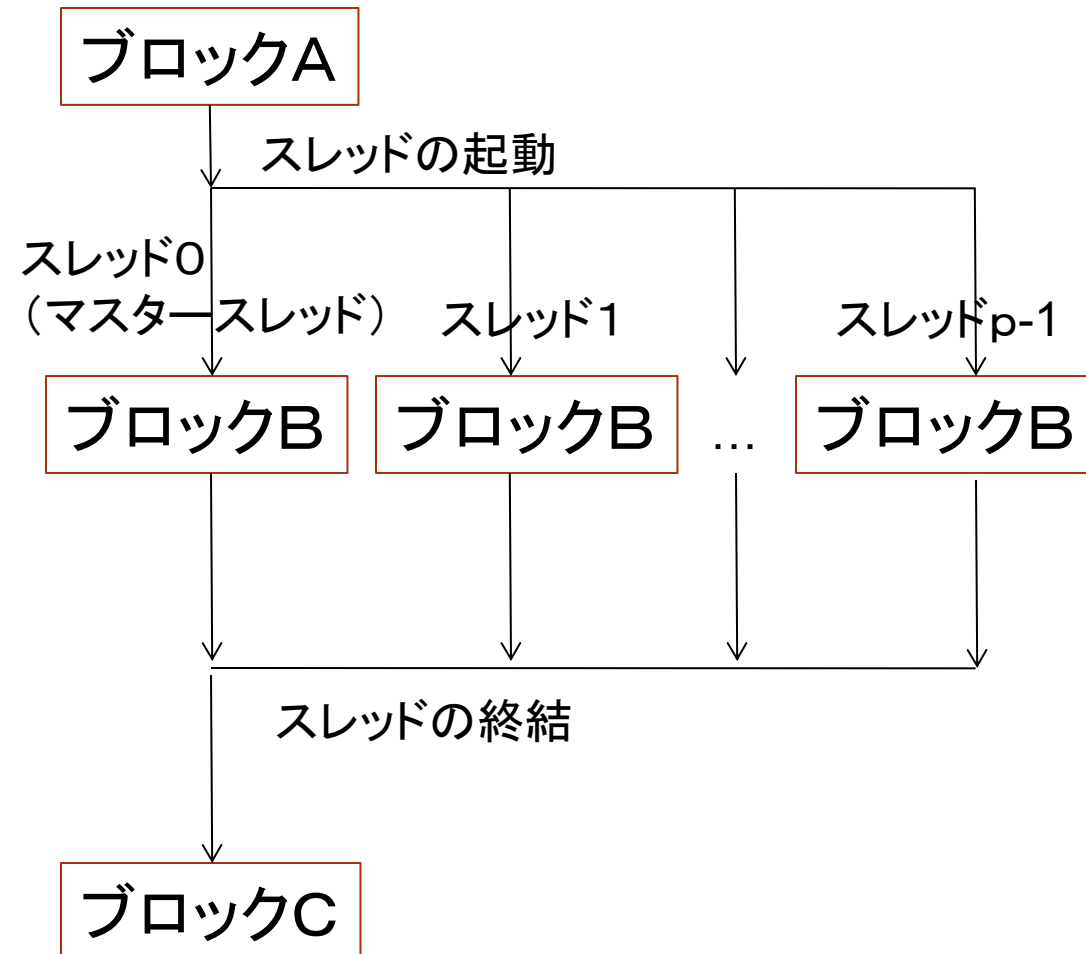


※スレッド数 p は、
環境変数
OMP_NUM_THREADS
で指定する。

OpenMPの実行モデル (Fortran言語)

OpenMP指示文

```
ブロックA  
!$omp parallel  
  ブロックB  
!$omp end parallel  
  ブロックC
```



※スレッド数 p は、
環境変数
OMP_NUM_THREADS
で指定する。

Work sharing構文

- parallel指示文のように、複数のスレッドで実行する場合において、OpenMPで並列を記載する処理（ブロックB）の部分を**並列領域(parallel region)**と呼ぶ。
- 並列領域を指定して、スレッド間で並列実行する処理を記述するOpenMPの構文を**Work sharing構文**と呼ぶ。
- Work sharing構文は、大きく分けて以下の2種がある。
 1. **並列領域内で記載するもの**
 - for構文（do構文）
 - sections構文
 - single構文 (master構文)、など
 2. **parallel指示文と組み合わせるもの**
 - parallel for 構文 (parallel do構文)
 - parallel sections構文、など

代表的な指示文

For構文 (do構文)

```
#pragma omp parallel for
for (i=0; i<100; i++){
  a[i] = a[i] * b[i];
}
```

※Fortran言語の場合は
 !\$omp parallel do
 ~
 !\$omp end parallel do

上位の処理

スレッドの起動

スレッド0

スレッド1

スレッド2

スレッド3

```
for (i=0; i<25; i++){
  a[i] = a[i] * b[i];
}
```

```
for (i=25; i<50; i++){
  a[i] = a[i] * b[i];
}
```

```
for (i=50; i<75; i++){
  a[i] = a[i] * b[i];
}
```

```
for (i=75; i<100; i++){
  a[i] = a[i] * b[i];
}
```

スレッドの終結

下位の処理

※並列化をしても
 正しい結果になることを
 ユーザが保障する。

For構文の指定ができない例

```
for (i=0; i<100; i++) {  
    a[i] = a[i] + 1;  
    b[i] = a[i-1] + a[i+1];  
}
```

- ループ並列化指示すると、逐次と結果が異なる
(a[i-1]が更新されていない場合がある)

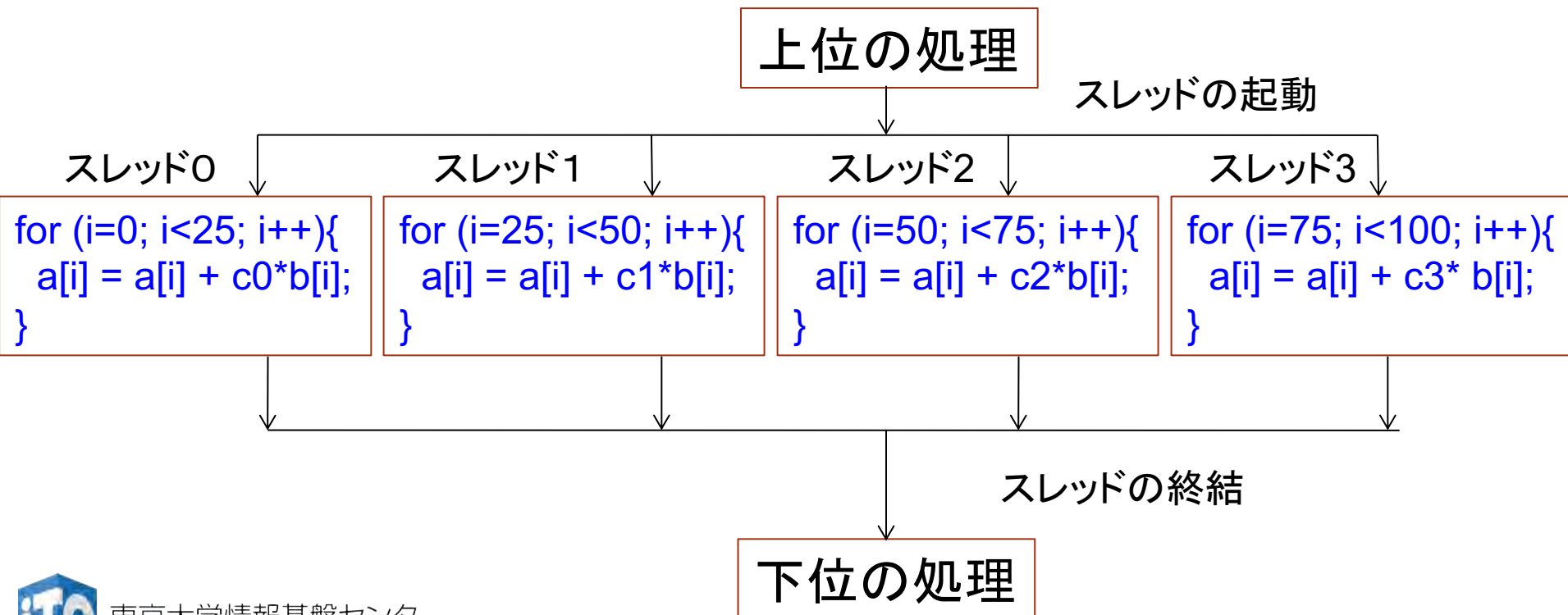
```
for (i=0; i<100; i++) {  
    a[i] = a[ ind[i] ];  
}
```

- ind[i]の内容により、ループ並列化できるかどうか決まる
- a[ind[i]]が参照時点で更新される可能性がない場合にループ並列化できる

Private補助指示文

```
#pragma omp parallel for private(c)
for (i=0; i<100; i++){
  a[i] = a[i] + c * b[i];
}
```

※変数cが各スレッドで別の変数を確保して実行
→高速化される



Private補助指示文の注意（C言語）

```
#pragma omp parallel for private( j )  
for (i=0; i<100; i++) {  
    for (j=0; j<100; j++) {  
        a[ i ] = a[ i ] + amat[ i ][ j ]* b[ j ];  
    }  
}
```

- ループ変数 j が、各スレッドで別の変数を確保して実行される。
- `private(j)` がない場合、各スレッドで共有変数 j のカウントを勝手に行ってしまうため、100回のループ実行にならない。
→ 演算結果が逐次と異なり、エラーとなる。

Private補助指示文の注意 (Fortran言語)

```
!$omp parallel do private( j )
do i=1, 100
  do j=1, 100
    a( i ) = a( i ) + amat( i , j ) * b( j )
  enddo
enddo
!$omp end parallel do
```

- ループ変数 j が、各スレッドで別の変数を確保して実行される。
- `private(j)` が無い場合、各スレッドで共有変数 j のカウントを勝手に行ってしまうため、100回のループ実行にならない。
→ 演算結果が逐次と異なり、エラーとなる。

リダクション補助指示文 (C言語)

- 内積値など、スレッド並列の結果を足しこみ、1つの結果を得たい場合に利用する
 - 上記の足しこみはスレッド毎に非同期になされる
 - **reduction**補助指示文が無いと、**ddot**は単なる共有変数になるため、**並列実行で逐次の結果と合わなくなくなる**

```
#pragma omp parallel for reduction(+:ddot )  
for (i=1; i<=100; i++) {  
    ddot += a[ i ] * b[ i ]  
}
```

ddotの場所はスカラー変数のみ記載可能（配列は記載できません）

リダクション補助指示文 (Fortran言語)

- 内積値など、スレッド並列の結果を足しこみ、1つの結果を得たい場合に利用する
 - 上記の足しこみはスレッド毎に非同期になされる
 - **reduction**補助指示文が無いと、**ddot**は共有変数になるため、**並列実行で逐次の結果と合わなくなくなる**

```
!$omp parallel do reduction(+:ddot )  
do i=1, 100  
    ddot = ddot + a(i) * b(i)  
enddo  
!$omp end parallel do
```

ddotの場所はスカラ変数のみ記載可能（配列は記載できません）

リダクション補助指示文の注意

- **reduction**補助指示文は、排他的に加算が行われるので、**性能が悪い**
 - 多くのスレッドを扱う場合、性能劣化が激しい
- 以下のように、**ddot**用の配列を確保して逐次で加算する方が高速な場合もある（ただし、問題サイズ、ハードウェア依存）

```
!$omp parallel do private ( i )
```

```
do j=0, p-1
```

```
do i=istart( j ), iend( j )
```

```
ddot_t( j ) = ddot_t( j ) + a( i ) * b( i )
```

```
enddo
```

```
enddo
```

```
!$omp end parallel do
```

```
ddot = 0.0d0
```

```
do j=0, p-1
```

```
ddot = ddot + ddot_t( j )
```

```
enddo
```

スレッド数分のループを作成：最大 p スレッド利用

各スレッドでアクセスするインデックス範囲を事前に

各スレッドで用いる、ローカルなddot用の配列ddot_t()を確保し、0に初期化しておく

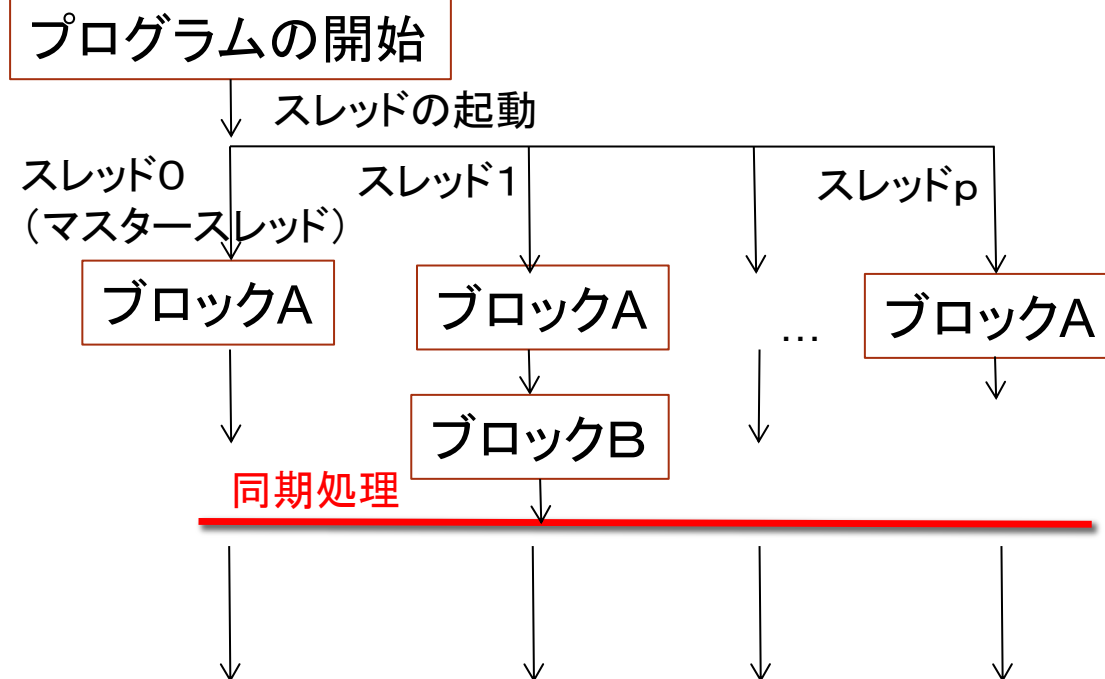
逐次で足しこみ

Single構文

- Single補助指示文で指定されたブロック
どれか1つのスレッドに割り当てる
- **どのスレッドに割り当てられるかは予測できない**
- `nowait`補助指示文を入れない限り、同期が入る

※Fortran言語の場合は
`!$omp single`
 ~
`!$omp end single`

```
#pragma omp parallel for
{
  ブロックA
  #pragma omp single
  { ブロックB }
  ...
}
```



Master構文

- 使い方は、single補助指示文と同じ
- ただし、master補助指示文で指定した処理（先ほどの例の「ブロックB」の処理）は、
必ずマスタースレッドに割り当てる
- 終了後の同期処理が入らない
 - そのため、場合により高速化される

よく使う OpenMP の関数

最大スレッド数取得関数

- 最大スレッド数取得には、`omp_get_num_threads()`関数を利用する
- 型はinteger (Fortran言語)、int (C言語)

● Fortran90言語の例

```
use omp_lib
Integer :: nthreads

nthreads = omp_get_num_threads()
```

● C言語の例

```
#include <omp.h>
int nthreads;

nthreads = omp_get_num_threads();
```

自スレッド番号取得関数

- 自スレッド番号取得には、`omp_get_thread_num()`関数を利用する
- 型はinteger (Fortran言語)、int (C言語)

● Fortran90言語の例

```
use omp_lib
Integer :: myid

myid = omp_get_thread_num()
```

● C言語の例

```
#include <omp.h>
int myid;

myid = omp_get_thread_num();
```


時間計測関数

- 時間計測には、`omp_get_wtime()`関数を利用する
- 型はdouble precision (Fortran言語)、double (C言語)

● Fortran90言語の例

```
use omp_lib
real(8) :: dts, dte

dts = omp_get_wtime()
  対象の処理
dte = omp_get_wtime()
print *, "Elapse time [sec.] =",dte-dts
```

● C言語の例

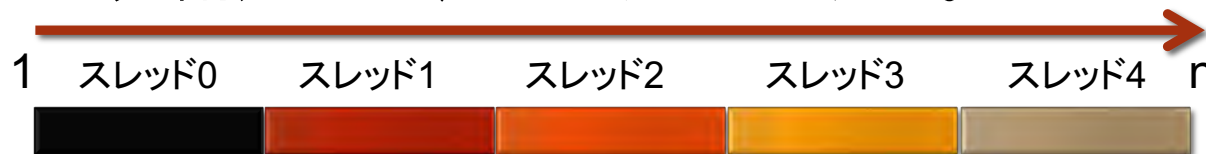
```
#include <omp.h>
double dts, dte;

dts = omp_get_wtime();
  対象の処理
dte = omp_get_wtime();
printf("Elapse time [sec.] = %lf ¥n",
dte-dts);
```

OpenMPにおけるスケ ジューリング

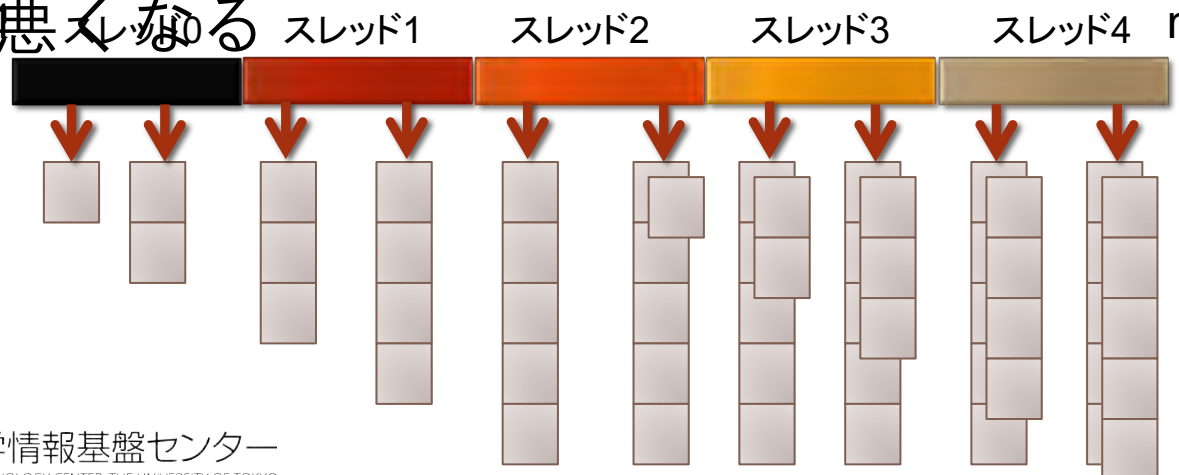
スケジューリングとは（その1）

- Parallel for/do構文では、対象ループの範囲（例えば1～nの長さ）を、単純にスレッド個数分に分割（連続するように分割）して、並列処理をする。



ループ変数の流れ
(反復空間)

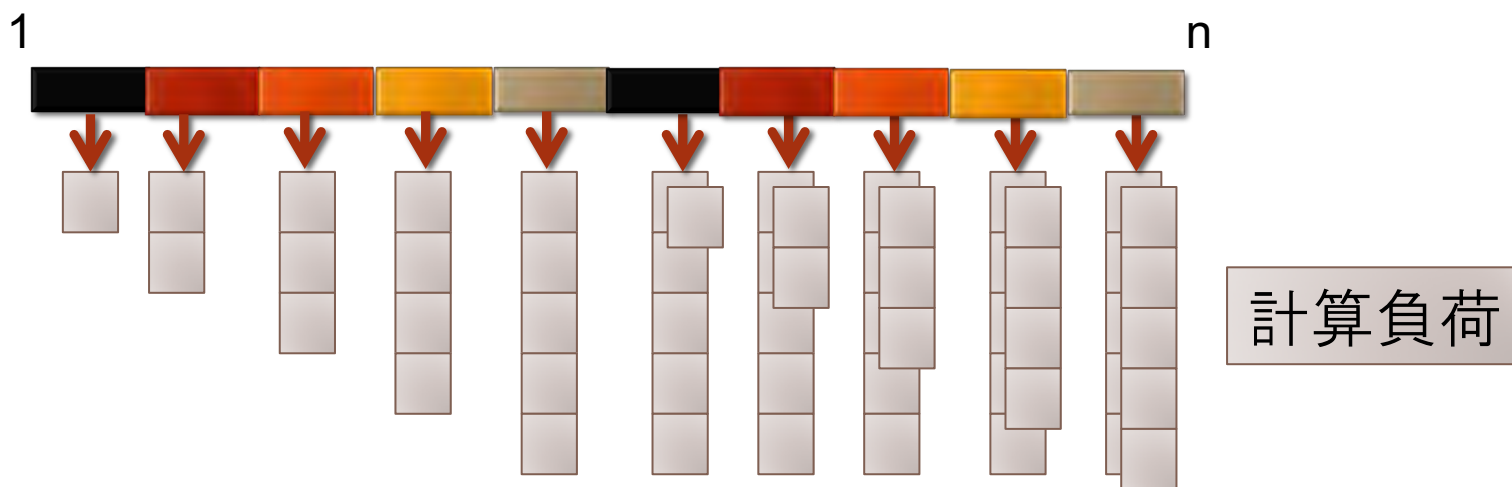
- このとき、各スレッドで担当したループに対する計算負荷が均等でないと、スレッド実行時の台数効果が悪くなる



計算負荷

スケジューリングとは（その2）

- ▶ 負荷分散を改善するには、割り当て間隔を短くし、かつ、循環するように割り当てればよい。



- ▶ 最適な、割り当て間隔（**チャンクサイズ**とよぶ）は、計算機ハードウェアと、対象となる処理に依存する。
- ▶ 以上の割り当てを行う補助指示文が用意されている。

ループスケジューリングの補助指定文 (その1)

- `schedule (static, n)`
 - ループ長をチャンクサイズで分割し、スレッド0番から順番に（スレッド0、スレッド1、・・・というように、ラウンドロビン方式と呼ぶ）、循環するように割り当てる。nにチャンクサイズを指定できる。
 - Schedule補助指定文を記載しないときのデフォルトは、`static`で、かつチャンクサイズは、ループ長/スレッド数。

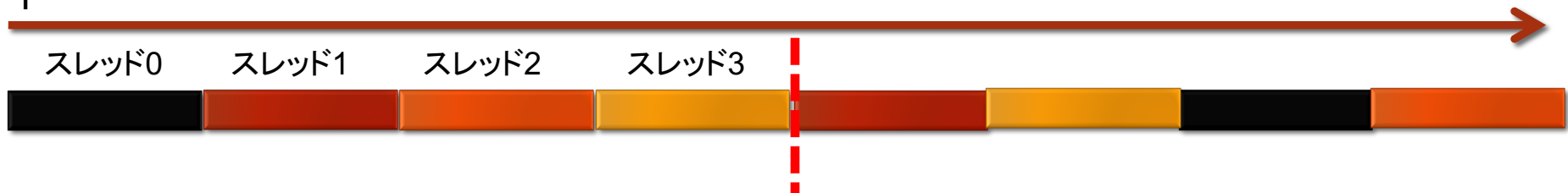
1



ループスケジューリングの補助指定文 (その2)

- `schedule(dynamic, n)`
 - ループ長をチャンクサイズで分割し、**処理が終了したスレッドから早い者勝ち**で、処理を割り当てる。nにチャンクサイズを指定できる。

1



ループスケジューリングの補助指定文 (その3)

• `schedule(guided, n)`

- ループ長をチャンクサイズで分割し、徐々にチャンクサイズを小さくしながら、処理が終了したスレッドから早い者勝ちで、処理を割り当てる。nにチャンクサイズを指定できる。
- チャンクサイズの指定が1の場合、残りの反復処理をスレッド数で割ったおおよその値が各チャンクのサイズになる。
- チャンクサイズは1に向かって指数的に小さくなる。
- チャンクサイズに1より大きいkを指定した場合、チャンクサイズは指数的にkまで小さくなるが、最後のチャンクはkより小さくなる場合がある。

1

スレッド0 チャンクサイズが指定されていない場合、デフォルトは1になる。



ループスケジューリングの補助指示文の使い方

● Fortran90言語の例

j-ループの反復回数が間接参照により決まるので、i-ループの計算負荷が均等であるか不明。実行時にしか、計算負荷の状況がわからないため、dynamicスケジューリングを適用

```
!$omp parallel do private( j, k ) schedule(dynamic,10)
do i=1, n
  do j=indj(i), indj (i+1)-1
    y( i ) = amat( j ) * x( indx( j ) )
  enddo
enddo
!$omp end parallel do
```

● C言語の例

```
#pragma omp parallel for private( j, k )
schedule(dynamic,10)
for (i=0; i<n; i++) {
  for ( j=indj(i); j<indj (i+1); j++) {
    y[ i ] = amat[ j ] * x[ indx[ j ] ];
  }
}
```


ループスケジューリングにおける プログラミング上の注意

- **dynamic、guidedのチャンクサイズは性能に大きく影響**
 - チャンクサイズが小さすぎると負荷バランスは良くなるが反面、処理待ちのオーバヘッドが大きくなる。
 - 一方、チャンクサイズが大きすぎると負荷バランスが悪くなる半面、処理待ちのオーバヘッドが小さくなる。
 - **上記の両者のトレードオフがある。**
 - 実行時のチャンクサイズのチューニングが必須で、チューニングコストが増える。
- **staticのみで高速実装ができる**（場合がある）
 - dynamicなどの実行時スケジューリングは、システムのオーバーヘッドが入るが、staticはオーバーヘッドは（ほとんど）無い。
 - **事前に負荷分散が均衡となるループ範囲を調べた上で、staticスケジューリングを使うと、最も効率が良い可能性がある。**
 - ただし、プログラミングのコストは増大する

OpenMPのプログラミング上の注意 (全般)

OpenMPによるプログラミング上の注意 点

- OpenMP並列化は、

parallel構文を用いた単純なforループ並列化

が主になることが多い。

- 複雑なOpenMP並列化はプログラミングコストがかかるので、OpenMPのプログラミング上の利点が失われる
- parallel構文による並列化は

private補助指示文の正しい使い方

を理解しないと、バグが生じる！

Private補助指示文に関する注意

- OpenMPでは、宣言せずに利用する変数は、すべて共有変数 (shared variable) になる
- C言語の大域変数、Fortran90言語のmodule変数は、そのままでは共有変数になる
 - プライベート変数にしたい場合は、Threadprivate宣言が必要
- parallel構文で関数呼び出ししている場合、その関数内でローカルに宣言している変数も、共有変数になる
 - そのままでは、並列処理で正常動作しない
 - これを防ぐには、以下のコードの変更が必要
 - 上記のローカル変数を引数にした関数呼び出しを作る
 - 上記のローカル変数を大域変数にして、Threadprivate宣言する

Parallel構文の入れ子に関する注意（その1）

- Parallel構文は、do補助指示文で分離して記載できる

```
!$omp parallel
!$omp do private(j,tmp)
do i=1, 100
  do j=1, 100
    tmp = b(j) + c(j)
    a(i) = a(i) + tmp
  enddo
enddo
!$omp end do
!$omp end parallel
```

Parallel構文の
対象が1ループ
なら parallel do
で指定するのが
良い

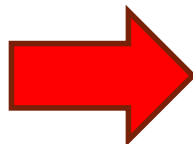


```
!$omp parallel do private(j,tmp)
do i=1, 100
  do j=1, 100
    tmp = b(j) + c(j)
    a(i) = a(i) + tmp
  enddo
enddo
!$omp end parallel do
```

Parallel構文の入れ子に関する注意（その2）

- Parallel構文は、do補助指示文で分離して記載できる
- 複数ループの内側を並列化したい場合は、分離した方が高速になる
 - ただし、外側ループを並列化できる時はその方が性能が良い
 - 外側ループにデータ依存があり、並列化できない場合

```
do i=1, n
!$omp parallel do
  do j=1, n
    <並列化できる式>
  enddo
!$omp end parallel do
enddo
```



```
!$omp parallel
do i=1, n
!$omp do
  do j=1, n
    <並列化できる式>
  enddo
!$omp end do
enddo
!$omp end parallel
```

OpenMPによるSIMD化

OpenMP+SIMD

- **SIMD (Single Instruction - Multiple Data)**
 - 1命令で複数データを処理
 - AVX512命令：1命令で倍精度浮動小数点数8個を処理
 - KNLではAVX512命令を最大限活用しなければ高い性能が得られない
- SIMD節はOpenMP 4.0からサポート
 - コンパイラが自動ベクトル化もしてくれるが、指示をした方が確実
 - コンパイラが誤った答えを出すこともあるので要確認（経験済み）
- コードの書き換えを必要とする場合も（しばしば）ある
- メモリのアラインメントが重要（前述）

OMP For ▪ Do+SIMD

- C言語 :

```
#pragma omp parallel for simd  
for (i=0; i < num; i++) {  
    sum = sum + a[i];  
}
```

- Fortran:

```
!$omp parallel do simd  
do i=1, num  
    sum = sum + a(i)  
!$omp end parallel
```

OMP SIMD単独

- C言語 :

```
#pragma omp parallel for
for (i=0; i < num; i++) {
  #pragma omp simd
  for (j=0; j < num; j++) {
    sum[i] = sum[i] + a[i][j];
  }
}
```

- Fortran:

```
!$omp parallel do
do i=1, num
!$omp simd
do j=1, num
sum(i) = sum(i) + a(j,i)
enddo
!$omp end simd
enddo
!$omp end parallel
```

Declare SIMD

- 関数をまるごとSIMD化したい場合
 - その関数はOMP SIMDの中で呼ばれる場合にSIMD化される

- **#pragma omp declare simd notinbranch**

```
float min (float a, float b)
```

```
{ return a < b ? a : b; }
```

```
void minner (float *a, float *b, float *c) {
```

```
#pragma omp parallel for simd
```

```
for (i=0; i<N; i++)
```

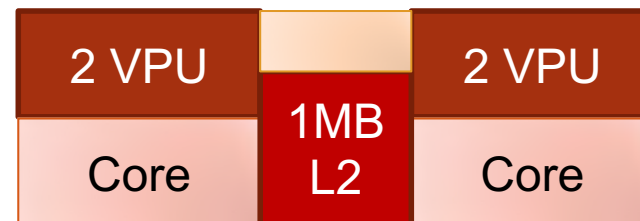
```
    c[i] = min(a[i], b[i]);
```

```
}
```

OFPにおけるOpenMP 実行の注意点

KNL(OFP)におけるHyperThreading

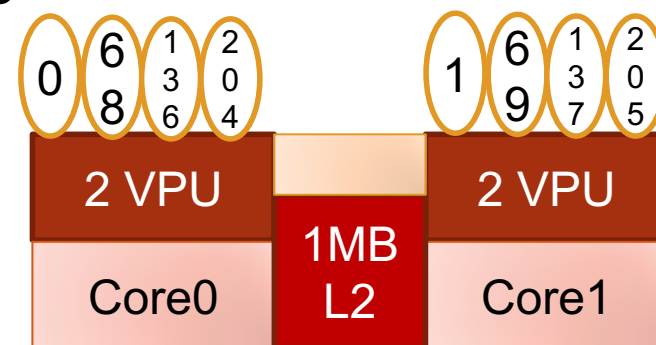
- 物理コア: **68コア**
- HyperThreading(HT)有効: 1コアあたり **4 HT**
=> 合計 **272スレッド**実行可能、ただしスレッドを増やせば速くなる、わけではない
- AVX-512ユニットは**物理コアあたり2個しかない(=136個)**
 - 2つのAVX-512ユニット(VPU)を使い尽くせれば **1 HT**でいいはず
 - AVX-512ユニット1基しか使えていなければ **2 HT**使った方がいい
 - メモリアクセスを隠蔽する必要がある場合は **3~4 HT**がいい場合もある
- 最適化が進むと**1 HT(=68スレッド)**が一番速かったりする



1タイル=2コア

OFPにおけるスレッド割付け

- メニーコア => OS割り込み処理などによるジッタの影響を受けやすい
 - 遅いコアが一人いると全員が影響を受ける
 - 特にタイマ割り込みは 1ms毎にやってくる
- OFPでは、**コア0のみ**タイマ割り込みを受け取るように設定
 - コア0と、タイル内でL2キャッシュを直接共有している**コア1も影響を受ける**
 - HTコアも当然影響を受ける
- コア番号 0,1, 68,69, 136,137, 204, 205 は実行から外した方がいい(場合が多い)
 - 数%以上性能が改善する
 - 具体的な指定方法は後述



1タイル=2コア

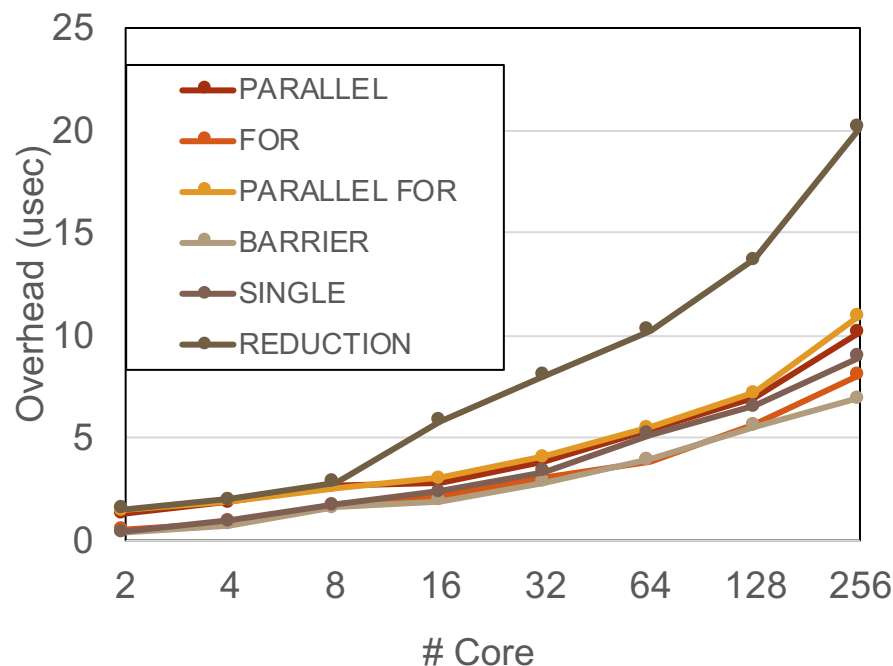
KNLでのスレッド並列数と性能への影響

Xeon Phiは、Xeonなど通常のマルチコアCPUに比べて

- コアのクロックが遅い
 - ReverbushのBroadwell (2.1GHz)に比べて 1.6倍程度
- メモリアクセスレイテンシが大きい
 - チップ内ネットワークが複雑

⇒ OpenMPの指示文に対するオーバヘッドが大きくなる

- 右図：各指示文のオーバヘッド (EPCC syncbench)
- 8コアから16コアで大きな影響があるように見える
 - 2コアは同一タイル、
 - 4~64コアは物理コアのみ、
 - 128~256コアはHT使用



OpenMPにおける対策

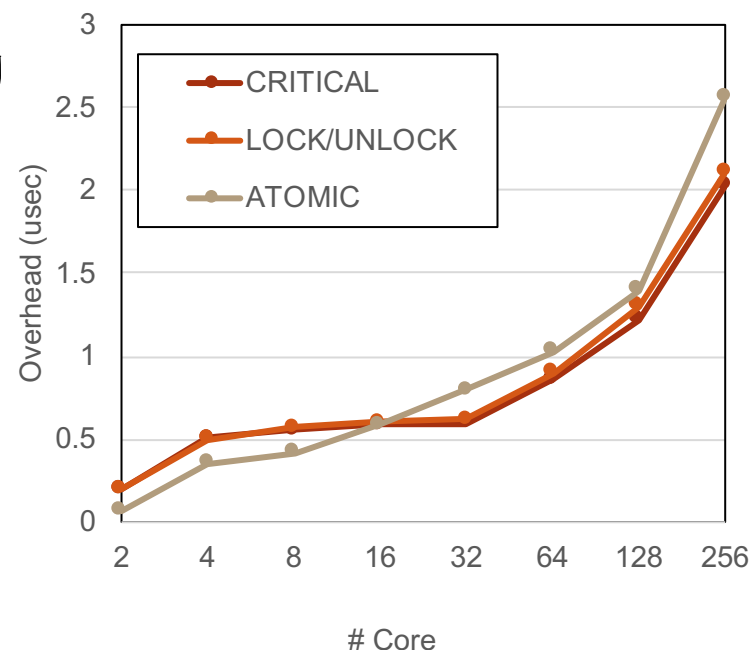
- なるべく Atomic処理やバリアを減らす
 - Parallel for/doなど、様々な指示文にも暗にバリア処理は含まれている
- 最大限大きく Parallel region を取る
 - 毎回Parallel for/doせず、大きな Parallel region の中でomp for/doを繰り返す
- MPI+OpenMPハイブリッド実行により、スレッド並列数を下げる
 - アプリケーションに依存する
 - あとで詳しく

各指示文(Atomic系)のオーバヘッド
(EPCC syncbench)

Xeon Broadwell 32コアでは
Atomic: 0.1us以下

Critical, Lock/Unlock: 0.4us以下

=> メモリアクセスレイテンシの
影響



KNL単体(1ノード)に おける実行

コンパイラ

Intel Compiler 2018 update 1を想定
(OFPのデフォルト, intel/2018.1.163, 2019年7月現在)

- Cコンパイラ: icc
- Fortranコンパイラ: ifort
- 2018 update2/3, 2019 update 1/3/4 もインストール済み
 - メンテナンスのタイミングで新たなバージョンがあれば追加
- 切り替え方法 (例: 2019 update4)
 - コマンドライン, ジョブスクリプト共通
module switch intel intel/2019.4.243
- 新しい方が性能が改善されている可能性が高い
(でも新たなバグもあるかも…)

コンパイルオプション

必須オプション

- `-xMIC-AVX512`
 - KNL向けコード生成
 - Xeonでも実行するなら `-axMIC-AVX512`
- `-O3`
 - 最適化レベル
- `-qopenmp`
 - OpenMP有効
- `-align array64byte`
 - (Fortranのみ) 配列を64byteにアライン (前述)

試した方がいい

- `-qopt-streaming-stores=`
always / **never** / auto
 - **always**: キャッシュに載せない命令を使う、キャッシュの生存時間が短い場合に有効
 - **never**はその逆
 - デフォルトはauto
 - ループの各配列ごとにも制御可能 (以下は**always**相当)
 - **#pragma vector nontemporal** (in C/C++)
 - **!DIR\$ vector nontemporal** (in F)
- `-qopt-prefetch={0~5}`
 - プリフェッチをどのくらい頑張るか
- `-no-vec`
 - SIMD化をあきらめる、ちょっとだけよくなる場合がある

実行時の指定

- 68コアのうち、64コアを使うことを想定
 - ジッタを防ぐための設定、コア 0,1, 68,69, 136,137, 204, 205を排除
 - メモリアクセスより計算の比重が高ければ68コア全て使うべき
- 例：
 1. 1タイルあたり1コアのみ、32コア使用
 - キャッシュ利用効率を重視する場合
 2. 1コア1HT、64コア使用
 3. 1コア2HT、128コア使用
 4. 1コア3HT、192コア使用
 5. 1コア4HT、256コア使用

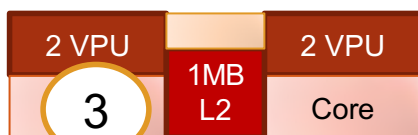
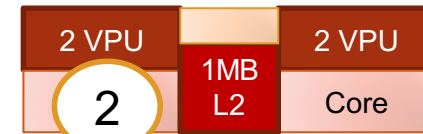
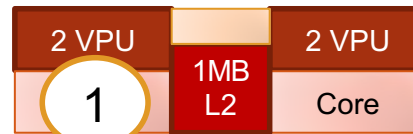
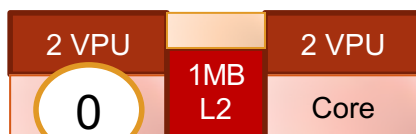
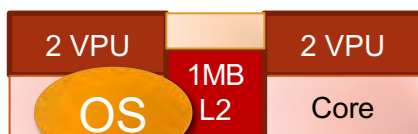
32コアを使用する場合

- タイルあたり1スレッド
 - キャッシュを十分に使えるが実際バンド幅を出すには足りないジョブスクリプト中で

```
#PJM --omp thread=32 または  
export OMP_NUM_THREADS=32
```

(以下同様にスレッド数=
使うコア数)

```
export KMP_AFFINITY=proclist=[2,4,6,8,10,12,14,16,18,20,  
22,24,26,28,30,32,34,36,38,40,42,44,46,48,50,52,54,  
56,58,60,62,64,66],explicit,verbose
```

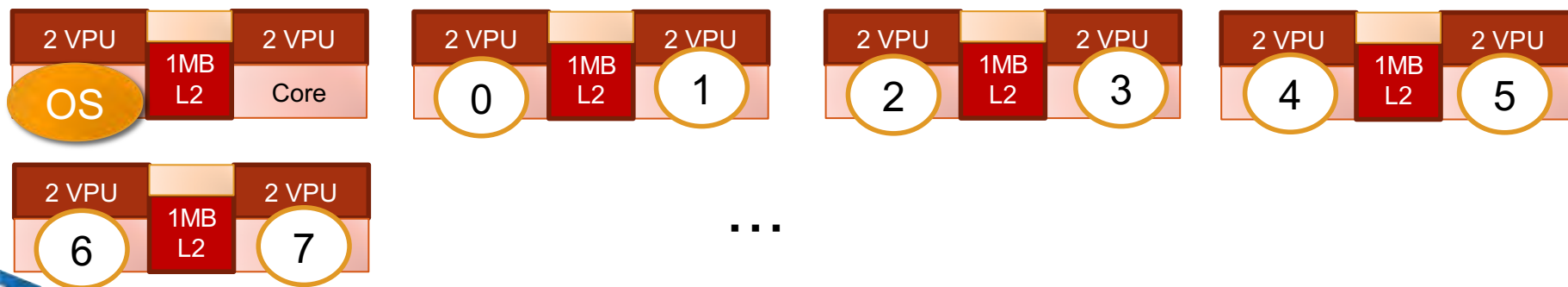


...

64コアを使用する場合

- コアあたり1スレッド
- KMP_HW_SUBSET環境変数を使える
 - 物理コア数c@オフセット, コアあたりスレッド数t
 - 64c@2, 1t: 64コア、コア番号2から使用、コアあたり1スレッド

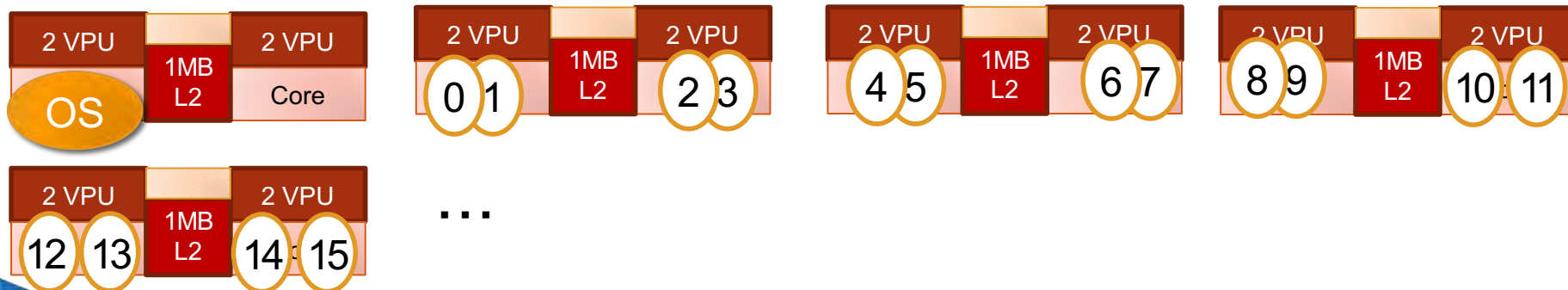
```
export KMP_HW_SUBSET=64c@2, 1t
export KMP_AFFINITY=scatter, verbose
```



128コアを使用する場合

- コアあたり2スレッド
- KMP_HW_SUBSET環境変数: 64c@2,2t
- KMP_AFFINITYにcompactを指定
 - scatterにすると、同一コアにスレッド0と64、スレッド1と65...のようになり、キャッシュが非効率

```
export KMP_HW_SUBSET=64c@2,2t
export KMP_AFFINITY=compact,verbose
```



192, 256コアを使用する場合

192コア

- コアあたり3スレッド

```
export KMP_HW_SUBSET=64c@2,  
      3t  
export KMP_AFFINITY=compact,  
      verbose
```

256コア

- コアあたり4スレッド

```
export KMP_HW_SUBSET=64c@2,  
      4t  
export KMP_AFFINITY=compact,  
      verbose
```


コア割り当ての確認

最初は必ず確認すること

(想定通りになっていなくて悲しい思いをする)

- KMP_AFFINITY環境変数にverboseを指定
- 標準エラー出力 (e+数字のファイル) にレポート出力
- 見るべきは以下のような出力 (最後の方)
 - 32コア実行の例: スレッド0がコア2, スレッド1がコア4, スレッド2がコア6...に割り当て

```
OMP: Info #242: KMP_AFFINITY: pid X tid XX thread 0 bound to OS proc set {2}
OMP: Info #242: KMP_AFFINITY: pid X tid XY thread 1 bound to OS proc set {4}
OMP: Info #242: KMP_AFFINITY: pid X tid XZ thread 2 bound to OS proc set {6}
```

...

ベクトル化、スレッド化、実行効率の確認

- コンパイル時に最適化レポートを確認
 - -qopt-report=5
 - ファイル名.optrpt を確認
- ツールを利用する（後述）
 - Intel Advisor XE
 - ベクトル化、スレッド並列化のアドバイス
 - Intel VTune XE
 - プロファイリング
 - Scalasca, SCORE-P, PAPIなど

KNL/OFPにおける MPI+OpenMPハイブリッド

MPI+OpenMPハイブリッド実行の意義

- MPIランク数の増加に伴い、消費メモリもオーバヘッドも増加
- 通常は以下のような組み合わせが良いとされて来た
 - ノード内: OpenMP
 - ノード間: MPI
- またReedbushのように2ソケットでNUMAを構成していれば
 - ソケット内: OpenMP
 - ソケット間+ノード間: MPI
- フラットMPIも可能だがよろしくない

メニーコアクラスタではそう単純でもない

MPIの特徴

- **メッセージパッシング用のライブラリ規格の1つ**
 - メッセージパッシングのモデルである
 - コンパイラの規格、特定のソフトウェアやライブラリを指すものではない!
- 分散メモリ型並列計算機で並列実行に向く
- 大規模計算が可能
 - 1プロセッサにおけるメモリサイズやファイルサイズの制約を打破可能
 - プロセッサ台数の多い並列システム（Massively Parallel Processing (MPP)システム）を用いる実行に向く
 - 1プロセッサ換算で膨大な実行時間の計算を、短時間で処理可能
 - 移植が容易
 - **API（Application Programming Interface）の標準化**
- スケーラビリティ、性能が高い
 - 通信処理をユーザが記述することによるアルゴリズムの最適化が可能
 - プログラミングが難しい（敷居が高い）

MPIの経緯（これまで）

- MPIフォーラム (<http://www.mpi-forum.org/>) が仕様策定
 - 1994年5月 1.0版 (MPI-1)
 - 1995年6月 1.1版
 - 1997年7月 1.2版、 および 2.0版 (MPI-2)
 - 2008年5月 1.3版、 2008年6月 2.1版
 - 2009年9月 2.2版
 - 日本語版 <http://www.pccluster.org/ja/mpi.html>
- MPI-2 では、以下を強化：
 - 並列I/O
 - C++、Fortran 90用インターフェース
 - 動的プロセス生成/消滅
 - 主に、並列探索処理などの用途

MPIの経緯 MPI-3.1

- MPI-3.0 2012年9月
- MPI-3.1 2015年6月
- 以下のページで現状・ドキュメントを公開中
 - <http://mpi-forum.org/docs/docs.html>
 - <http://meetings.mpi-forum.org>
 - <http://meetings.mpi-forum.org/mpi31-impl-status-Nov15.pdf>
- 注目すべき機能
 - ノン・ブロッキング集団通信機能
(MPI_IALLREDUCE、など)
 - 高性能な片方向通信 (RMA、Remote Memory Access)
 - Fortran2008 対応、など

MPIの経緯 MPI-4.0策定中

- 以下のページで経緯・ドキュメントを公開
 - <https://www.mpi-forum.org/mpi-40/>
- 検討されている機能
 - ハイブリッドプログラミングへの対応
 - MPIアプリケーションの耐故障性 (Fault Tolerance, FT)
 - 永続的(Persistent)コレクティブ
 - 性能アサーションとヒント
 - RMA / One sided通信
 - いくつかのアイデアを検討中
 - Active Messages (メッセージ通信のプロトコル)
 - 計算と通信のオーバーラップ
 - 最低限の同期を用いた非同期通信
 - 低いオーバーヘッド、パイプライン転送
 - バッファリングなしで、インタラプトハンドラで動く
 - Stream Messaging
 - 新プロファイル・インターフェース

MPIの実装

- **MPICH (エム・ピッチ)**
 - 米国アルゴンヌ国立研究所が開発
- **MVAPICH (エムヴァピッチ)**
 - 米国オハイオ州立大学で開発、MPICHをベース
 - InfiniBand向けの優れた実装
- **OpenMPI**
 - オープンソース
- **ベンダMPI**
 - 大抵、上のどれかがベースになっている
例: 富士通「京」、FX10用のMPI: Open-MPIベース
Intel MPI: MPICH、MVAPICHベース
 - 注意点: メーカー独自機能拡張がなされていることがある

MPIによる通信

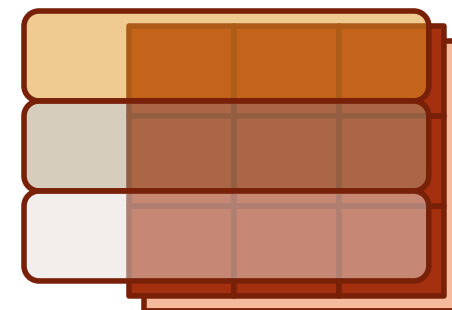
- 郵便物の郵送と同じ
- 郵送に必要な情報：
 1. 自分の住所、送り先の住所
 2. 中に入っているものはどこにあるか
 3. 中に入っているものの分類
 4. 中に入っているものの量
 5. (荷物を複数同時に送る場合の) 認識方法 (タグ)
- MPIでは：
 1. 自分の認識ID、および、送り先の認識ID
 2. データ格納先のアドレス
 3. データ型
 4. データ量
 5. タグ番号

略語とMPI用語

- MPIは「プロセス」間の通信を行います。プロセスは（普通は）「プロセッサ」（もしくは、コア）に一対一で割り当てられます。
- ランク (Rank)
 - 各「MPIプロセス」の「識別番号」のこと。
 - 通常MPIでは、MPI_Comm_rank関数で設定される変数（サンプルプログラムではmyid）に、0～全PE数-1の数値が入る
 - 世の中の全MPIプロセス数を知るには、MPI_Comm_size関数を使う。
（サンプルプログラムでは、numprocs に、この数値が入る）

コミュニケータ

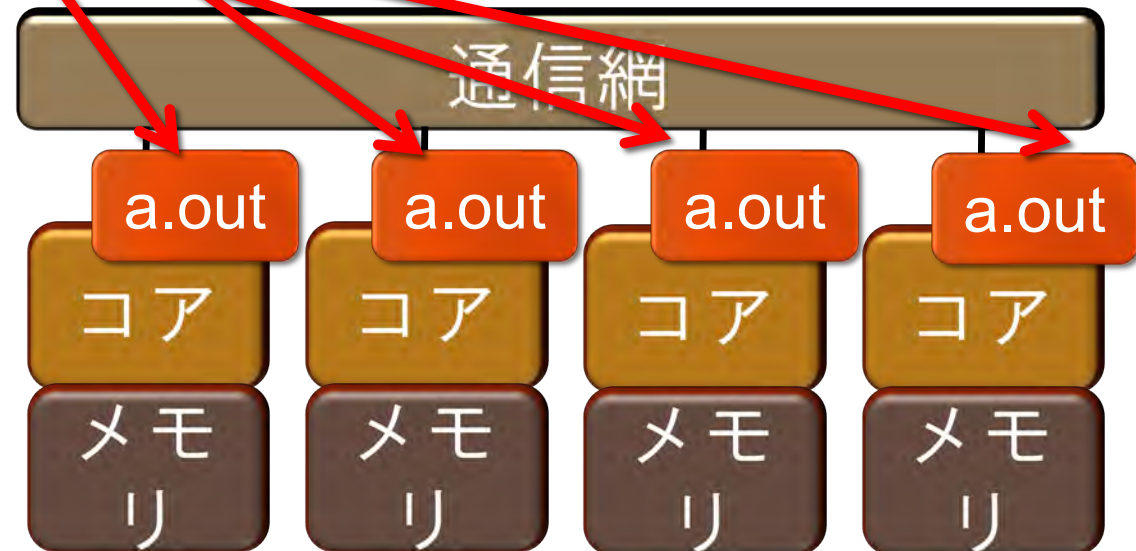
- `MPI_COMM_WORLD`は、**コミュニケータ**とよばれる概念を保存する変数
- コミュニケータは、操作を行う対象のプロセッサ群を定める
- 初期状態では、**0番～`numprocs - 1`番**までのプロセッサが、1つのコミュニケータに割り当てられる
 - この名前が、“**`MPI_COMM_WORLD`**”
- プロセッサ群を分割したい場合、**`MPI_Comm_split`**関数を利用
 - メッセージを、一部のプロセッサ群に放送するとき利用
 - “マルチキャスト”で利用



MPIの起動

- 各ランクでは、同じプログラムが同時に複数起動し開始される = SPMD (Single Program Multiple Data)

```
mpiexec.hydra -n 4 ./a.out
```



基本的なMPI関数

送信、受信のためのインタフェース

MPI関数

- システム関数
 - MPI_Init ; MPI_Comm_rank ; MPI_Comm_size ; MPI_Finalize ;
- 1対1通信関数
 - ブロッキング型
 - MPI_Send ; MPI_Recv ;
 - ノンブロッキング型
 - MPI_Isend ; MPI_Irecv ;
- 1対全通信関数
 - MPI_Bcast
- 集団通信関数
 - MPI_Reduce ; MPI_Allreduce ; MPI_Barrier ;
- 時間計測関数
 - MPI_Wtime

C言語インターフェースと Fortranインターフェースの違い

- C版は、 整数変数*ierr* が戻り値
`ierr = MPI_Xxxx(...);`
- Fortran版は、最後に整数変数*ierr*が引数
`call MPI_XXXX(..., ierr)`
- システム用配列の確保の仕方
 - C言語
`MPI_Status istatus;`
 - Fortran言語
`integer istatus(MPI_STATUS_SIZE)`

C言語インターフェースと Fortranインターフェースの違い

- MPIにおける、データ型の指定
 - C言語
 - MPI_CHAR (文字型)、 MPI_INT (整数型)、
MPI_FLOAT (実数型)、 MPI_DOUBLE (倍精度実
数型)
 - Fortran言語
 - MPI_CHARACTER (文字型)、 MPI_INTEGER
(整数型)、 MPI_REAL (実数型)、
MPI_DOUBLE_PRECISION (倍精度実数型)、
MPI_COMPLEX (複素数型)
- 以降は、C言語インターフェースで説明する

基礎的なMPI関数—MPI_Recv (1/2)

- `ierr = MPI_Recv(recvbuf, icount, idatatype, isource, itag, icommm, istatus);`
- `recvbuf`: 受信領域の先頭番地を指定する。
- `icount`: 整数型。受信領域のデータ要素数を指定する。
- `idatatype`: 整数型。受信領域のデータの型を指定する。
 - `MPI_CHAR` (文字型)、`MPI_INT` (整数型)、`MPI_FLOAT` (実数型)、`MPI_DOUBLE` (倍精度実数型)
- `isource`: 整数型。受信したいメッセージを送信するPEのランクを指定する。
 - 任意のPEから受信したいときは、`MPI_ANY_SOURCE` を指定する。

基礎的なMPI関数—MPI_Recv (2/2)

- **itag**: 整数型。受信したいメッセージに付いているタグの値を指定。
 - 任意のタグ値のメッセージを受信したいときは、**MPI_ANY_TAG** を指定。
- **icomm**: 整数型。PE集団を認識する番号であるコミュニケータを指定。
 - 通常では**MPI_COMM_WORLD** を指定すればよい。
- **istatus**: MPI_Status型（整数型の配列）。受信状況に関する情報が入る。 **かならず専用の型宣言をした配列を確保すること。**
 - 要素数が**MPI_STATUS_SIZE**の整数配列が宣言される。
 - 受信したメッセージの送信元のランクが **istatus[MPI_SOURCE]**、タグが **istatus[MPI_TAG]** に代入される。
 - **C言語**: **MPI_Status istatus;**
 - **Fortran言語**: **integer istatus(MPI_STATUS_SIZE)**
- **ierr(戻り値)**: 整数型。エラーコードが入る。

基礎的なMPI関数—MPI_Send

```
• ierr = MPI_Send(sendbuf, icount, idatatype, idest,  
                 itag,  icomm);
```

- **sendbuf** : 送信領域の先頭番地を指定
- **icount** : 整数型。送信領域のデータ要素数を指定
- **idatatype** : 整数型。送信領域のデータの型を指定
- **idest** : 整数型。送信したいPEのicomm内でのランクを指定
- **itag** : 整数型。受信したいメッセージに付けられたタグの値を指定
- **icomm** : 整数型。プロセッサ集団を認識する番号である
 コミュニケータを指定
- **ierr (戻り値)** : 整数型。エラーコードが入る。

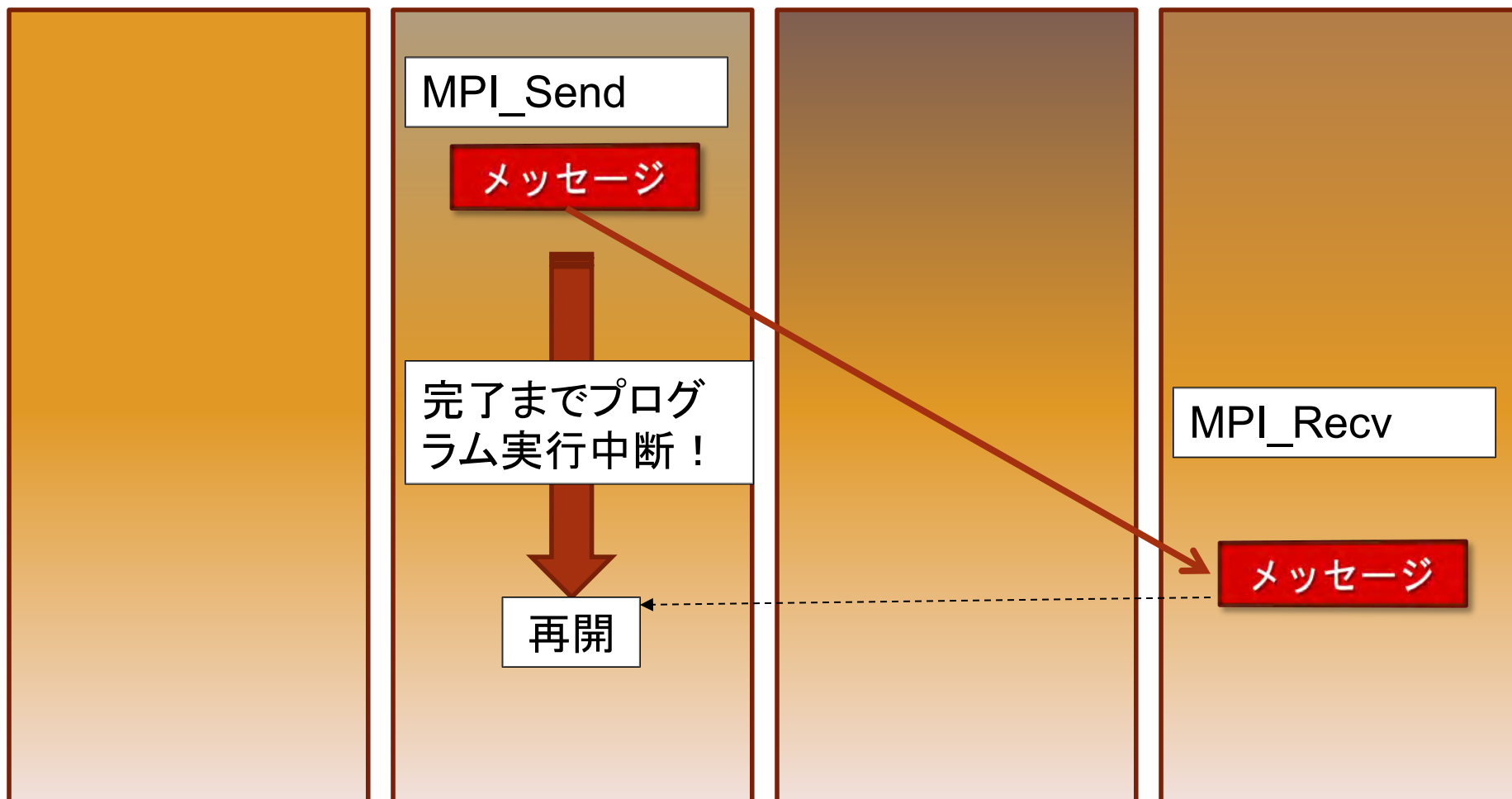
Send-Recvの概念 (1対1通信)

PE0

PE1

PE2

PE3

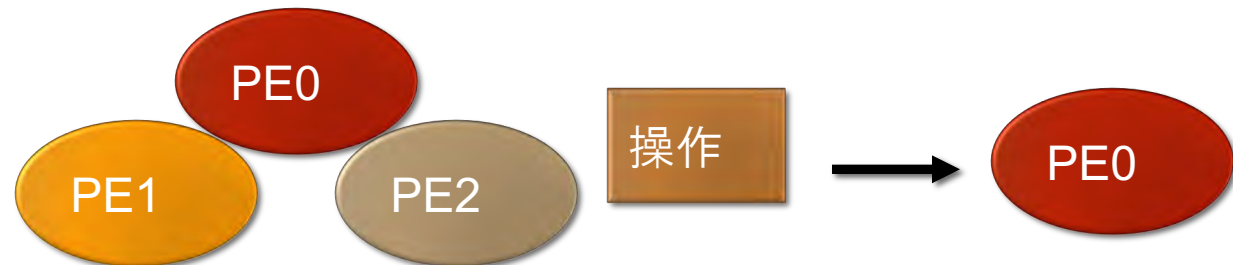


リダクション演算

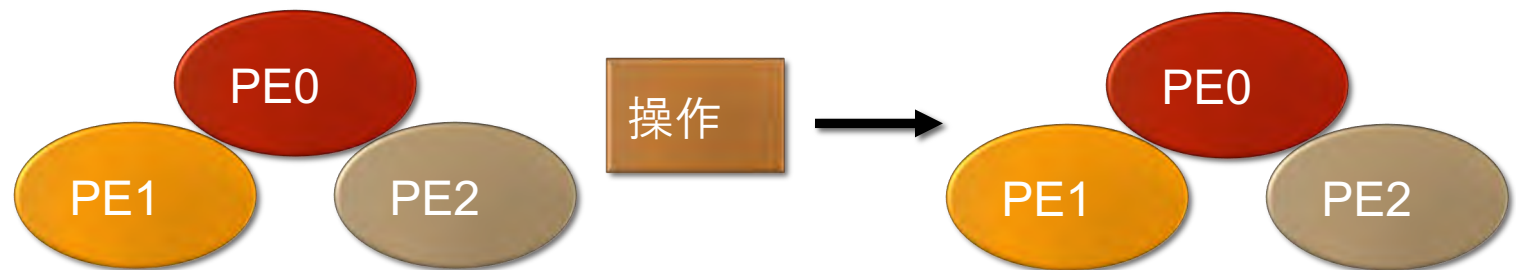
- <操作>によって<次元>を減少
(リダクション) させる処理
 - 例: 内積演算
ベクトル (n 次元空間) \rightarrow スカラ (1次元空間)
- リダクション演算は、通信と計算を必要とする
 - 集団通信演算 (collective communication operation)
と呼ばれる
- 演算結果の持ち方の違いで、2種の
インターフェースが存在する

リダクション演算

- 演算結果に対する所有PEの違い
 - **MPI_Reduce**関数
 - リダクション演算の結果を、ある一つのPEに所有させる



- **MPI_Allreduce**関数
 - リダクション演算の結果を、全てのPEに所有させる



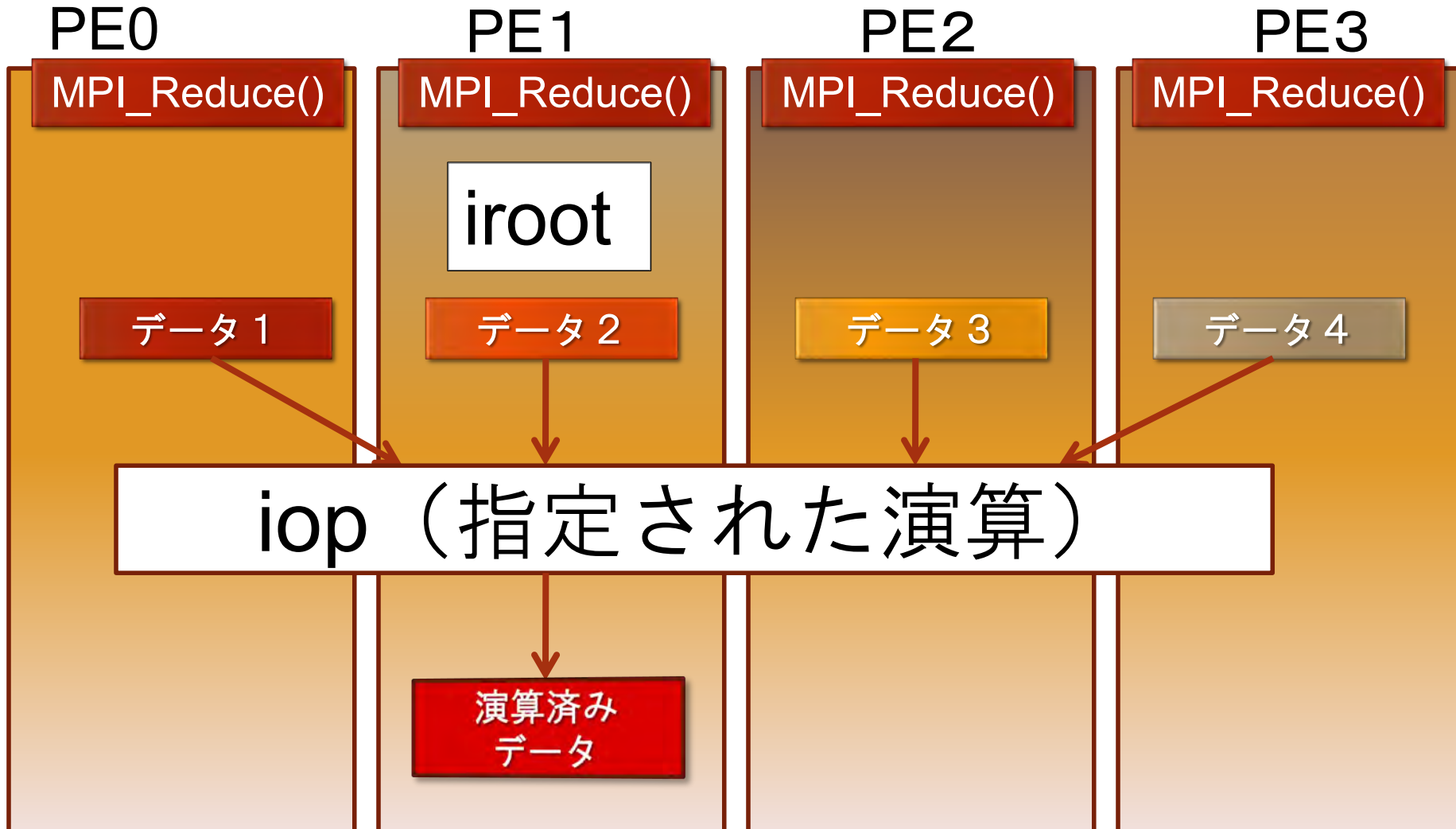
基礎的なMPI関数—MPI_Reduce

- `ierr = MPI_Reduce(sendbuf, recvbuf, icount, idatatype, iop, iroot, icommm);`
- `sendbuf`: 送信領域の先頭番地を指定する。
- `recvbuf`: 受信領域の先頭番地を指定する。 `iroot` で指定したPEのみで書き込みがなされる。
送信領域と受信領域は、同一であってはならない。
すなわち、異なる配列を確保しなくてはならない。
- `icount`: 整数型。送信領域のデータ要素数を指定する。
- `idatatype`: 整数型。送信領域のデータの型を指定する。
 - (Fortran) <最小/最大値と位置>を返す演算を指定する場合は、`MPI_2INTEGER`(整数型)、`MPI_2REAL`(単精度型)、`MPI_2DOUBLE_PRECISION`(倍精度型)、を指定する。

基礎的なMPI関数—MPI_Reduce

- **iop** : 整数型。演算の種類を指定する。
 - **MPI_SUM** (総和)、**MPI_PROD** (積)、**MPI_MAX** (最大)、**MPI_MIN** (最小)、**MPI_MAXLOC** (最大とその位置)、**MPI_MINLOC** (最小とその位置) など。
- **iroot** : 整数型。結果を受け取るPEのicomm内でのランクを指定する。全てのicomm内のPEで同じ値を指定する必要がある。
- **icomm** : 整数型。PE集団を認識する番号であるコミュニケータを指定する。
- **ierr** : 整数型。エラーコードが入る。

MPI_Reduceの概念 (集団通信)



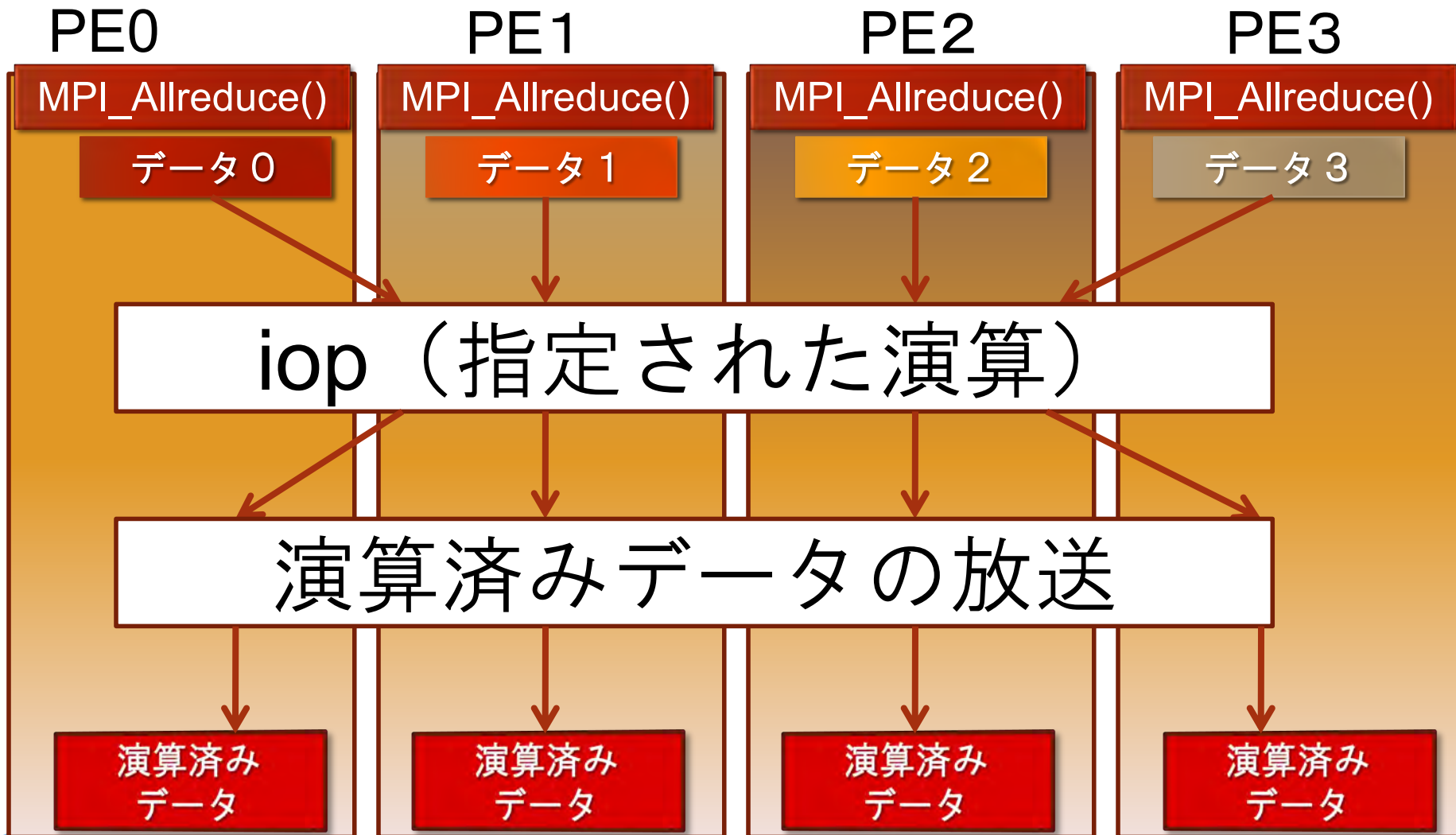
基礎的なMPI関数—MPI_Allreduce

- `ierr = MPI_Allreduce(sendbuf, recvbuf, icount, idatatype, iop, icommm);`
 - `sendbuf`: 送信領域の先頭番地を指定する。
 - `recvbuf`: 受信領域の先頭番地を指定する。 `iroot` で指定したPEのみで書き込みがなされる。
送信領域と受信領域は、同一であってはならない。
すなわち、異なる配列を確保しなくてはならない。
 - `icount`: 整数型。送信領域のデータ要素数を指定する。
 - `idatatype`: 整数型。送信領域のデータの型を指定する。
 - 最小値や最大値と位置を返す演算を指定する場合は、`MPI_2INT`(整数型)、`MPI_2FLOAT`(単精度型)、`MPI_2DOUBLE`(倍精度型) を指定する。

基礎的なMPI関数—MPI_Allreduce

- **iop** : 整数型。演算の種類を指定する。
 - **MPI_SUM** (総和)、 **MPI_PROD** (積)、
MPI_MAX (最大)、 **MPI_MIN** (最小)、
MPI_MAXLOC (最大と位置)、 **MPI_MINLOC**
(最小と位置) など。
- **icomm** : 整数型。PE集団を認識する番号であるコミュニケータを指定する。
- **ierr** : 整数型。 エラーコードが入る。

MPI_Allreduceの概念（集団通信）



ブロッキング、ノンブロッキング

1. ブロッキング

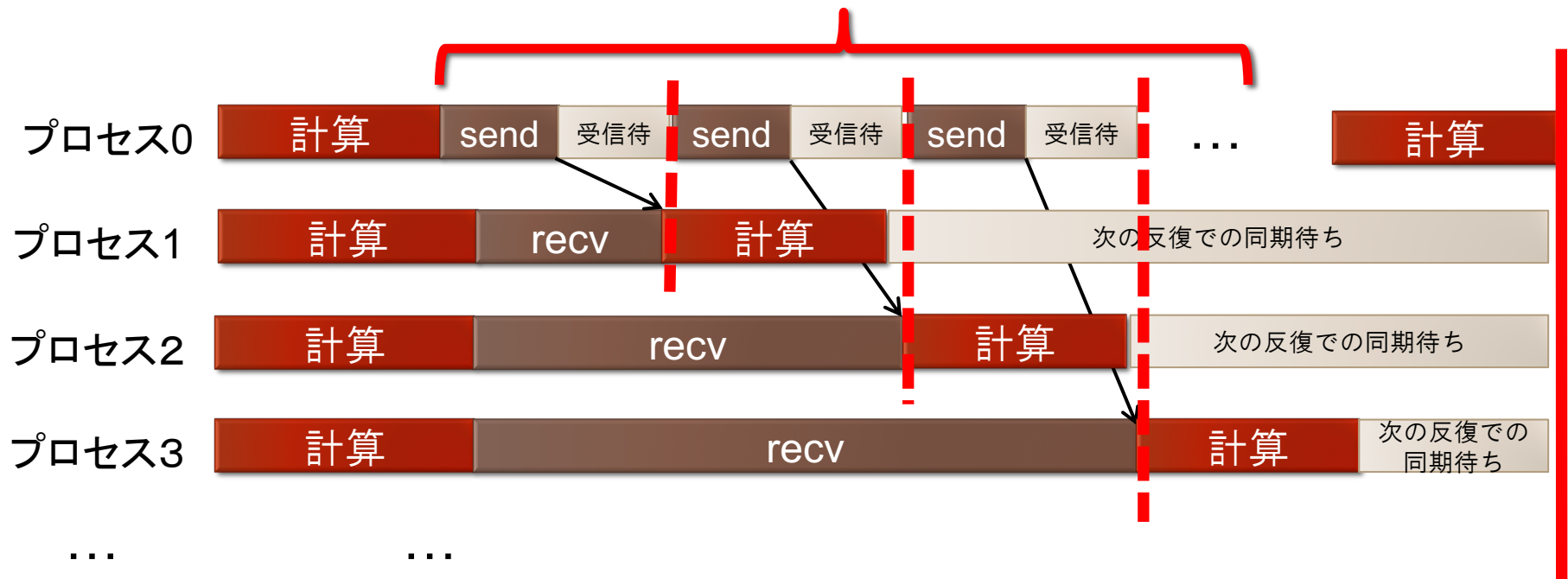
- 送信／受信側のバッファ領域にメッセージが格納され、受信／送信側のバッファ領域が自由にアクセス・上書きできるまで、呼び出しが戻らない
- バッファ領域上のデータの一貫性を保障
- MPI_Send, MPI_Bcastなど

2. ノンブロッキング

- 送信／受信側のバッファ領域のデータを保障せずすぐに呼び出しが戻る
- バッファ領域上のデータの一貫性を保障せず
 - 一貫性の保証はユーザの責任

ブロッキング通信で効率の悪い例

- プロセス0が必要なデータを持っている場合
連続するsendで、効率の悪い受信待ち時間が多発



次の
反復での
同期点

ノンブロッキング通信関数

```
• ierr = MPI_Isend(sendbuf, icount, datatype, idest, itag, icommm, irequest);
```

- `sendbuf` : 送信領域の先頭番地を指定する
- `icount` : 整数型。送信領域のデータ要素数を指定する
- `datatype` : 整数型。送信領域のデータの型を指定する
- `idest` : 整数型。送信したいPEの`icommm` 内でのランクを指定する
- `itag` : 整数型。受信したいメッセージに付けられたタグの値を指定する

ノンブロッキング通信関数

- **icommm** : 整数型。PE集団を認識する番号であるコミュニケータを指定する。
 - 通常ではMPI_COMM_WORLD を指定すればよい。
- **irequest** : MPI_Request型（整数型の配列）。送信を要求したメッセージにつけられた識別子が戻る。
- **ierr** : 整数型。エラーコードが入る。

同期待ち関数

- `ierr = MPI_Wait(irequest, istatus);`

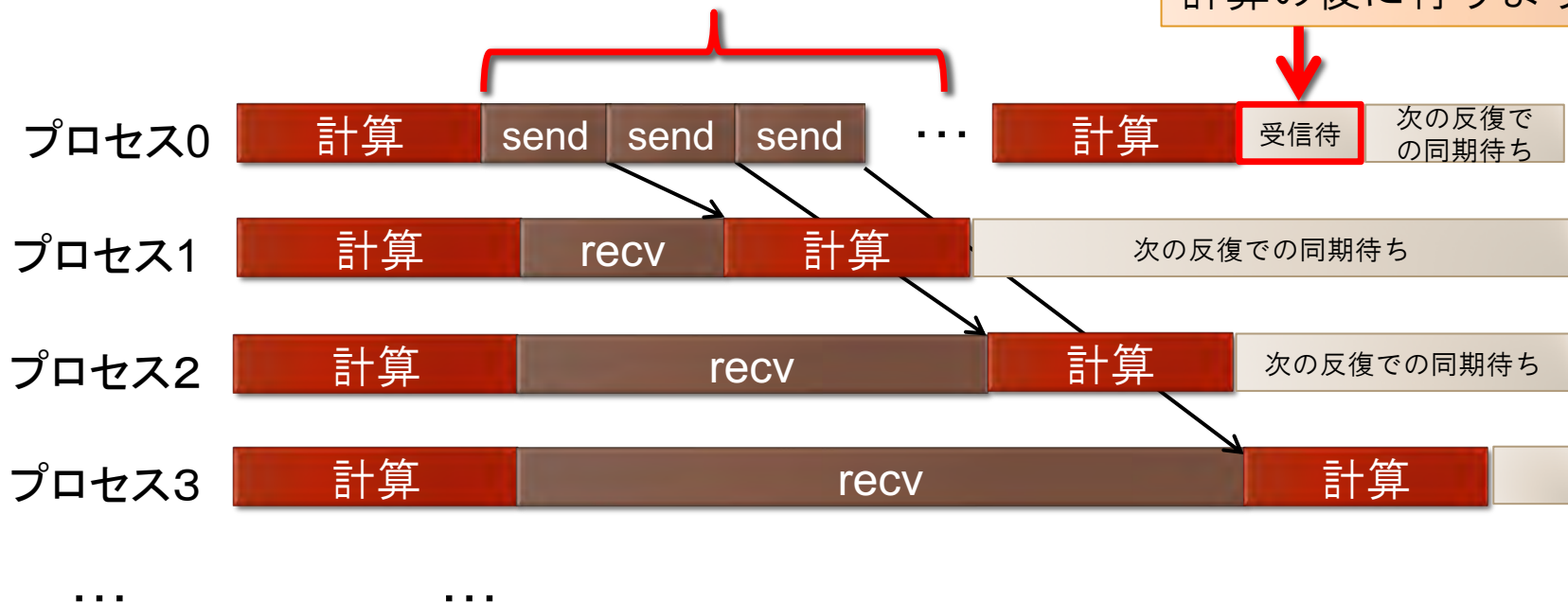
- `irequest` : `MPI_Request`型（整数型配列）。送信を要求したメッセージにつけられた識別子。
- `istatus` : `MPI_Status`型（整数型配列）。受信状況に関する情報が入る。
 - 要素数が`MPI_STATUS_SIZE`の整数配列を宣言して指定する。
 - 受信したメッセージの送信元のランクが`istatus[MPI_SOURCE]`、タグが`istatus[MPI_TAG]`に代入される。
- 送信データを変更する前・受信データを読み出す前には必ず呼ぶこと

ノン・ブロッキング通信による改善

- プロセス0が必要なデータを持っている場合

連続するsendにおける受信待ち時間を
ノン・ブロッキング通信で削減

受信待ちを、MPI_Waitで
計算の後に行うように変更



次の
反復での
同期点

注意点

- 以下のように解釈してください：
 - **MPI_Send**関数
 - 関数中に**MPI_Wait**関数が入っている；
 - **MPI_Isend**関数
 - 関数中に**MPI_Wait**関数が入っていない；
 - かつ、すぐにユーザプログラムに戻る；

参考文献

1. MPI並列プログラミング、P.パチェコ 著 / 秋葉博 訳
2. 並列プログラミング虎の巻MPI版、青山幸也 著、
理化学研究所情報基盤センタ
(<http://acc.riken.jp/HPC/training/text.html>)
3. Message Passing Interface Forum
(<http://www.mpi-forum.org/>)
4. MPI-Jメーリングリスト
(<http://phase.hpcc.jp/phase/mpi-j/ml/>)
5. 並列コンピュータ工学、富田真治著、昭晃堂 (1 9 9 6)

時間計測方法 (C言語)

```
double t0, t1, t2, t_w;  
..  
ierr = MPI_Barrier(MPI_COMM_WORLD);  
t1 = MPI_Wtime();
```

<ここに測定したいプログラムを書く>

```
t2 = MPI_Wtime();
```

```
t0 = t2 - t1;
```

```
ierr = MPI_Reduce(&t0, &t_w, 1,  
MPI_DOUBLE, MPI_MAX, 0,  
MPI_COMM_WORLD);
```

バリア同期後
時間を習得し保存

各プロセッサで、 t_0 の値は異なる。
この場合は、最も遅いものの値をプロセッサ0番が受け取る

時間計測方法 (Fortran言語)

```
real(8) :: t0, t1, t2, t_w
```

```
..  
call MPI_BARRIER(MPI_COMM_WORLD, ierr)  
t1 = MPI_WTIME()
```

<ここに測定したいプログラムを書く>

```
t2 = MPI_WTIME()
```

```
t0 = t2 - t1
```

```
call MPI_REDUCE(t0, t_w, 1, MPI_REAL8, &  
MPI_MAX, 0, MPI_COMM_WORLD, ierr)
```

バリア同期後
時間を習得し保存

各プロセッサで、t0の値
は異なる。

この場合は、最も遅いもの
の値をプロセッサ0番
が受け取る

OFFPにおけるMPI実行

実行上の注意

- プロセスピンニング（アフィニティ）を最適な設定にすることが重要
- ノード数，プロセス数を適切に設定
 - #PJM -L node=使用ノード数**
 - #PJM --mpi proc=全プロセス数**
- 設定スクリプト
 - MPIのみ（フラットMPI）
 - **/usr/local/bin/mpi_core_setting.sh**
 - MPI+ OpenMPハイブリッド
 - **/usr/local/bin/hybrid_core_setting.sh**
 - スクリプトの詳細な使い方は「利用手引書」参照

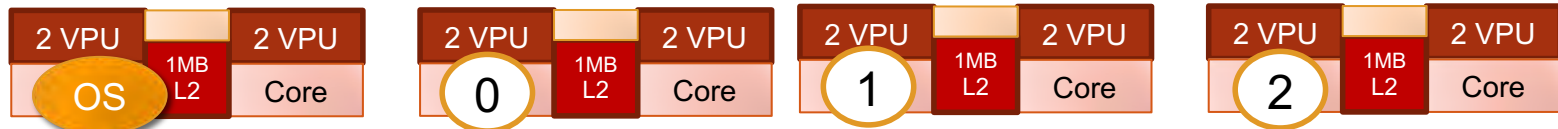
フラットMPI

- ジョブスクリプトに以下を設定

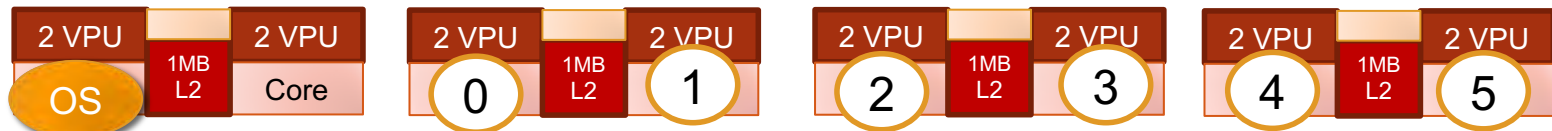
#PJM --mpi proc=全プロセス数

source /usr/local/bin/mpi_core_setting.sh

- 1~33プロセス：先頭1タイルを避けて各タイルに配置



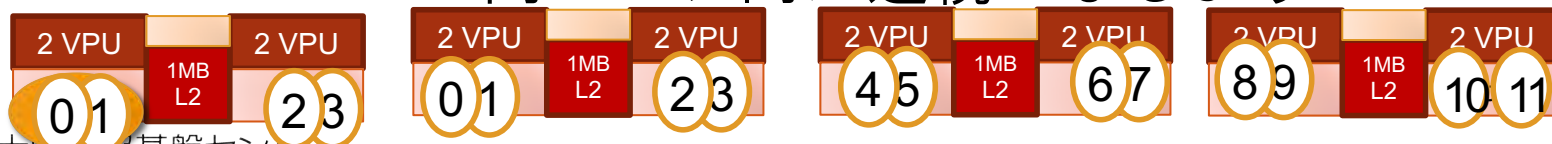
- 34~66プロセス：先頭1タイルを避けて各コアに配置



- 67プロセス：先頭コアを避けて各コアに配置

- 68プロセス：全コアを使うように配置

- 69~272プロセス：同一コア内が連続になるように



フラットMPIでの注意

- 69コア以上でかつ先頭タイルを避けたい場合は、以下の環境変数を追加する必要がある

export

```
I_MPI_PIN_PROCESSOR_EXCLUDE_LIST=0,1,68,69,136,137,204,205
```

- 先頭コアのみ避けるには

export

```
I_MPI_PIN_PROCESSOR_EXCLUDE_LIST=0,68,136,204
```

MPI+OpenMPハイブリッド

- 各プロセスで1~68スレッドまでに対応
- ジョブスクリプトに以下を追加
 - #PJM --mpi proc=全プロセス数**
 - #PJM --omp thread=プロセスあたりスレッド数**
 - source /usr/local/bin/hybrid_core_setting.sh モード**
- モード
 - 1: 先頭の1タイルを避けて**各タイル**に配置, 33スレッドまで
 - 2: 先頭の1タイルを避けて**各コア**に配置, 66スレッドまで
 - 3: 先頭の1コアを避けて**各コア**に配置, 67スレッドまで
 - 4: **各コア**に配置, 68スレッドまで (先頭は避けない)

MPI+OpenMPでの注意

- 69スレッド以上使いたい場合は多少自力で設定が必要
- 例: 128スレッド/プロセス(1ノード1プロセス)使いたい場合

```
#PJM -L node=16
```

```
#PJM --mpi process=16
```

```
#PJM --omp thread=64 わざと68スレッド以下にする
```

```
source /usr/local/bin/hybrid_core_setting.sh 2
```

```
export OMP_NUM_THREADS=128 本来使いたい数
```

```
export KMP_HW_SUBSET=2t コアあたりのスレッド数
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} 実行ファイル名
```

MPI(+OpenMP)実行での注意

- MCDRAMの設定もお忘れなく！！
 - I_MPI_HBW_POLICY環境変数
export I_MPI_HBW_POLICY=hbw_bind
export I_MPI_HBW_POLICY=hbw_preferred
export I_MPI_HBW_POLICY=hbw_interleave
- プロセスへの割り当てが意図通りか確認すること
export I_MPI_DEBUG=5
 - さらにスレッド割り当ても確認
export KMP_AFFINITY+=,verbose

```
[40] MPI startup(): shm and tmi data transfer modes
```

```
...
```

```
[0] MPI startup(): 192          269565      c3850.ofp {2,70,138,206}
```

MPI+OpenMPハイブリッド実行に際して

- MPI関数を、マルチスレッド環境で実行する際は注意が必要
 - 3つのモード: Master only, Funneled, Multiple
 - Multipleでは、集団通信は単純には実行できない、あまりスケールしない
 - Master onlyではParallelリージョンから出なければいけない
 - 特に意識しないのであればFunneledで使うのがよい

Master only

```
#pragma omp parallel
{
  ....
}
MPI_Send( ... );
#pragma omp parallel
{
  ....
}
```

Funneled

```
#pragma omp parallel
{
  ....
  #pragma omp master
  MPI_Send( ... );
  ....
}
```

Master onlyと似ているが
parallel節を閉じなくていい

Multiple

```
#pragma omp parallel
{
  ....
  MPI_Send( ... );
  ....
}
```

本当にマルチスレッド
だが、1対1通信しか
できない

マルチスレッドMPIの初期化

- MPI_Init()の代わりに、MPI_Init_thread()を使用

- C言語:

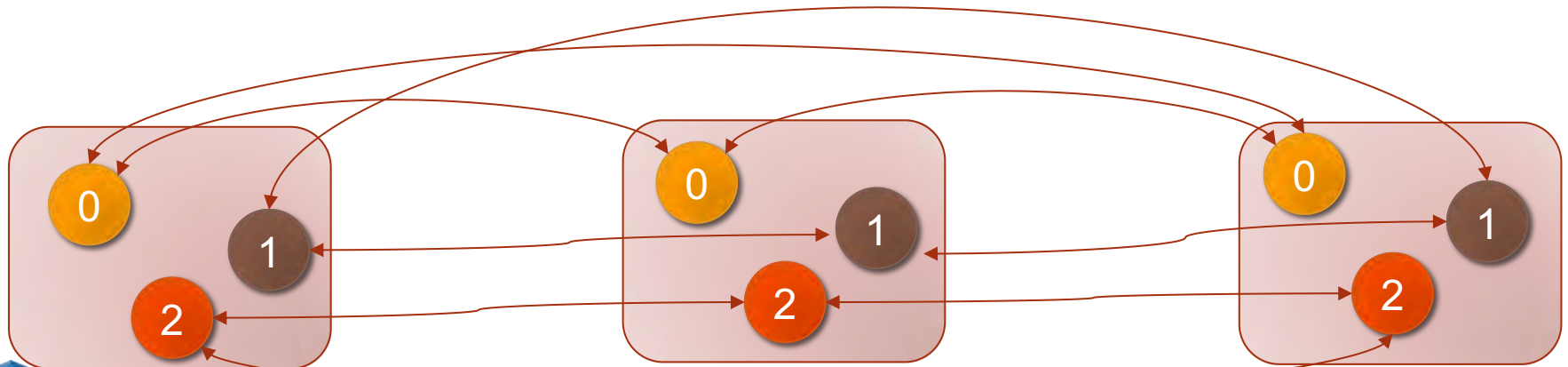
```
int provided;  
MPI_Init_thread(&argc,&argv,MPI_THREAD_FUNNELED,  
&provided);
```

- Fortran

```
integer provided, required  
required=MPI_THREAD_FUNNELED  
call MPI_Init_thread(required, provided, ierr)
```


マルチスレッドMPIの高速化 (IMPI 2019.1以降)

- マルチスレッドMPIのうち、特定の使い方に限って高速化されている
- 各ノード間で同じスレッドID間でだけ通信する場合
- `export I_MPI_THREAD_SPLIT=1`
- `export I_MPI_THREAD_RUNTIME=openmp`



通信と計算のオーバラップ

- "dynamic/runtime"
- "!\$omp master~!\$omp end master"

コア同士の動作モードのばらつきも大きいので、動的スケジューリングが効くと考えられる

```
!$omp parallel private (neib,j,k,i,X1,X2,X3,WVAL1,WVAL2,WVAL3)
!$omp&          private (istart,inum,ii,ierr)

!$omp master          Communication is done by the master thread (#0)
!C
!C- Send & Recv.
(...)
    call MPI_WAITALL (2*NEIBPETOT, req1, stal, ierr)
!$omp end master

!C
!C-- Pure Inner Nodes      The master thread can join computing of internal
                          nodes after the completion of communication

!$omp do schedule (runtime)  export OMP_SCHEDULE="dynamic, [chunksize]"
do j= 1, Ninn
    (...)
enddo

!C
!C-- Boundary Nodes      Computing for boundary nodes are by all threads

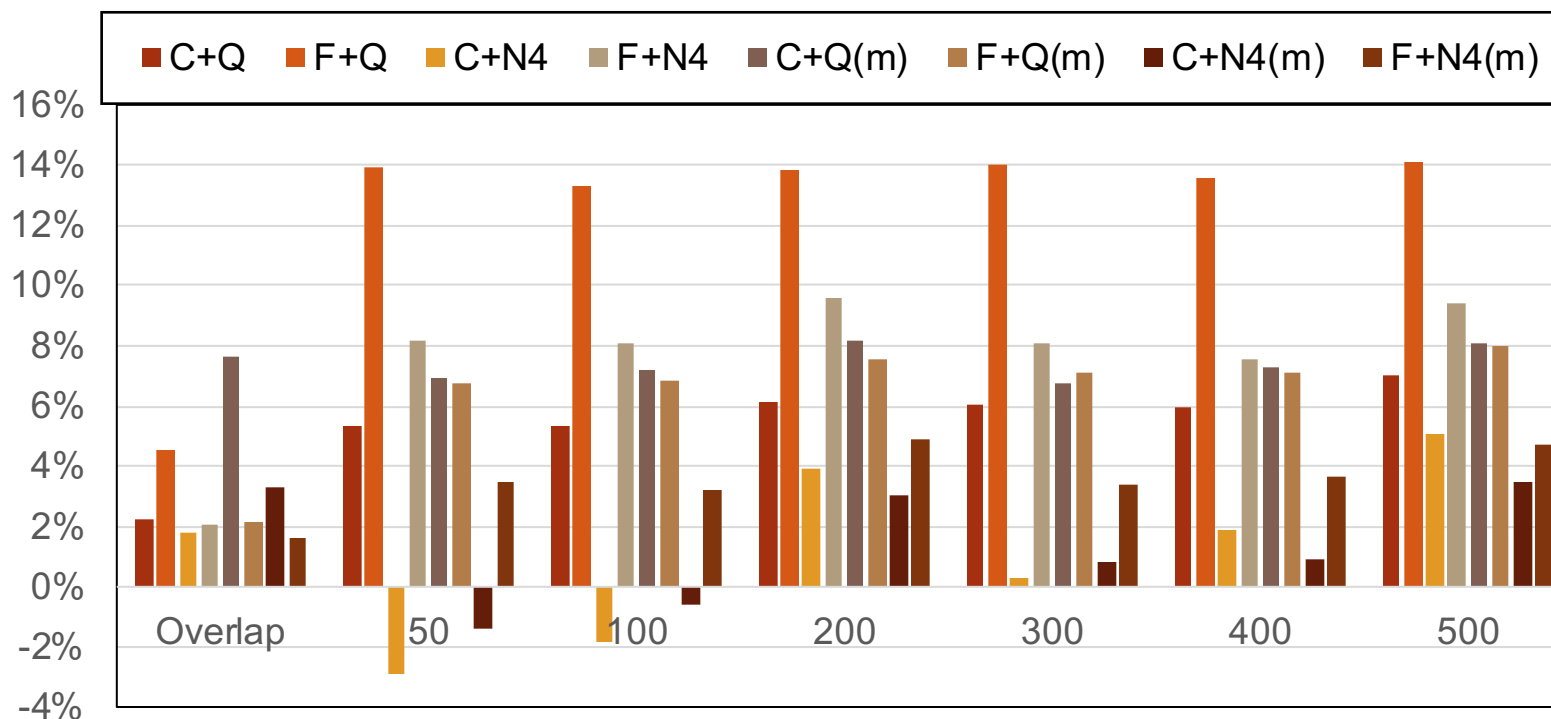
!$omp do          default: !$omp do schedule (static)
do j= Ninn+1, N
    (...)
enddo

!$omp end parallel
```

Idomura, Y. et al., Communication-overlap techniques for improved strong scaling of gyrokinetic Eulerian code beyond 100k cores on the K-computer, Int. J. HPC Appl. 28, 73-86, 2014

Speedup by communication overlapping (OFP: 32 nodes)

- Memory model: C: Cache, F: Flat
- Sub NUMA: Q: Quadrant, N4: SNC-4
- (m) : core binding to avoid “busy” core



通信性能安定化

- MPI通信で触る配列を、2Mバイトにアラインする

Fortranでは、(注: Fortranではバージョン2018から利用可能)

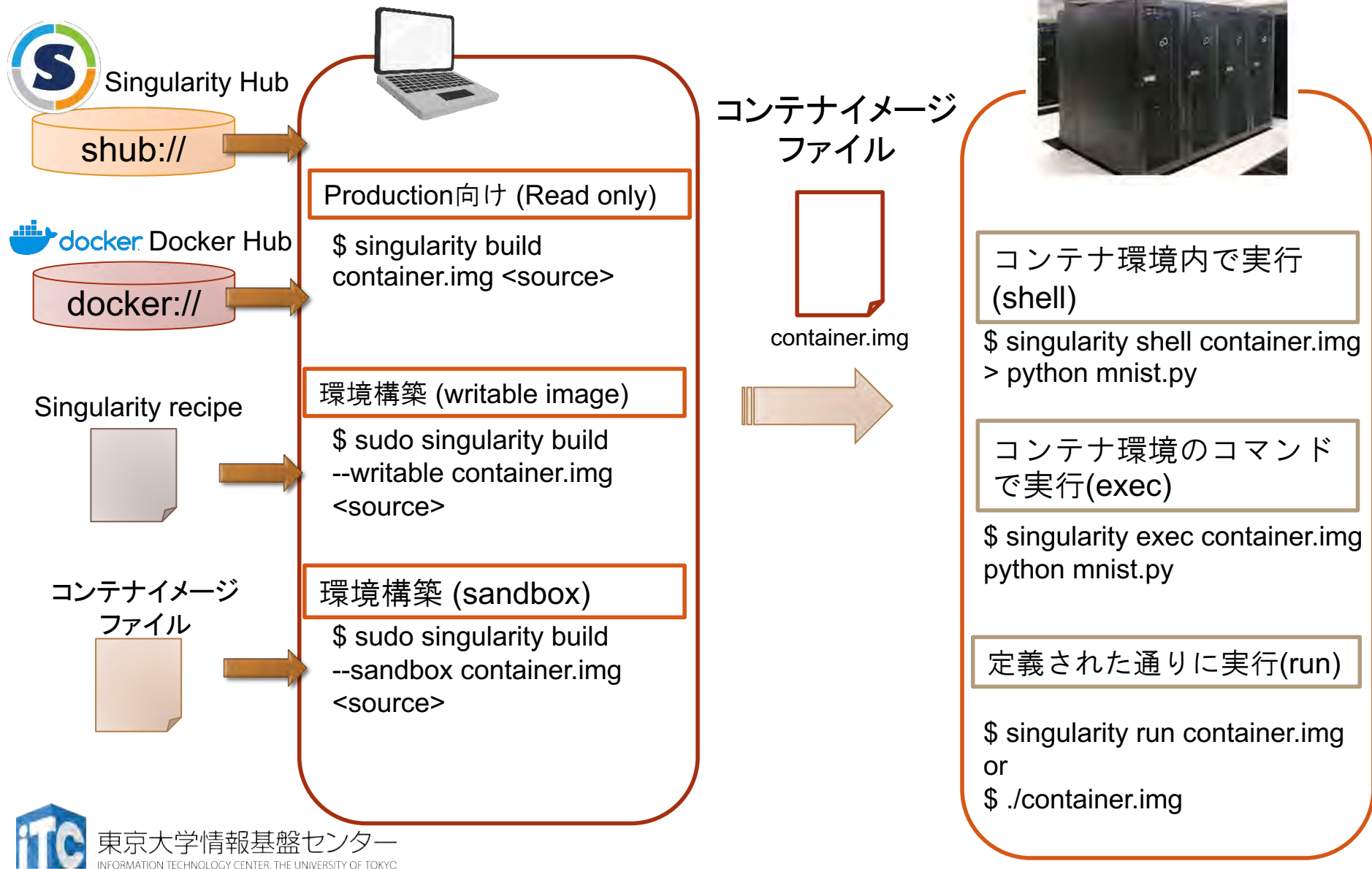
```
!DIR$ ATTRIBUTES ALIGN : 2097152 :: arrayMPI  
real(8), allocatable :: arrayMPI(:, :, :, :)
```

Cでは、(コンパイラのバージョン依存なし)

```
arrayMPI = (double *)__mm_malloc(array_size,  
2*1024*1024);
```

その他の話題

Singularity: コンテナ仮想化



バーストバッファ (DDN: Infinite Memory Engine)

- SSD搭載のサーバが25台, 計50ノード
 - 容量: 960 TB, 理論ピーク性能 1.5TB/sec
 - IO500でTopの性能
- クライアントでパリティを計算して書き込む
 - erasure coding
- Lustreファイルシステムのキャッシュとしても動作
 - /cache でキャッシュへのアクセス (計算ノードのみ)
 - ステージングが必要 (ステージイン, ステージアウト操作)
 - #PJM -x STGIN_LIST=<一覧をリストしたファイル名>
 - #PJM -x STGOUT_LIST=<一覧をリストしたファイル名>
 - ステージインを待ちたいとき
 - `prestige-wait` コマンド

IMEの使い方

1. POSIX:

/work/gt00/t002
xx の代わりに
/cache/gt00/t00
2xx を使う

Intel MPIを使いたい
場合にも有効！！

2. MPI IO:

- a. MPIをIntel MPIから
MVAPICH2に切り替え

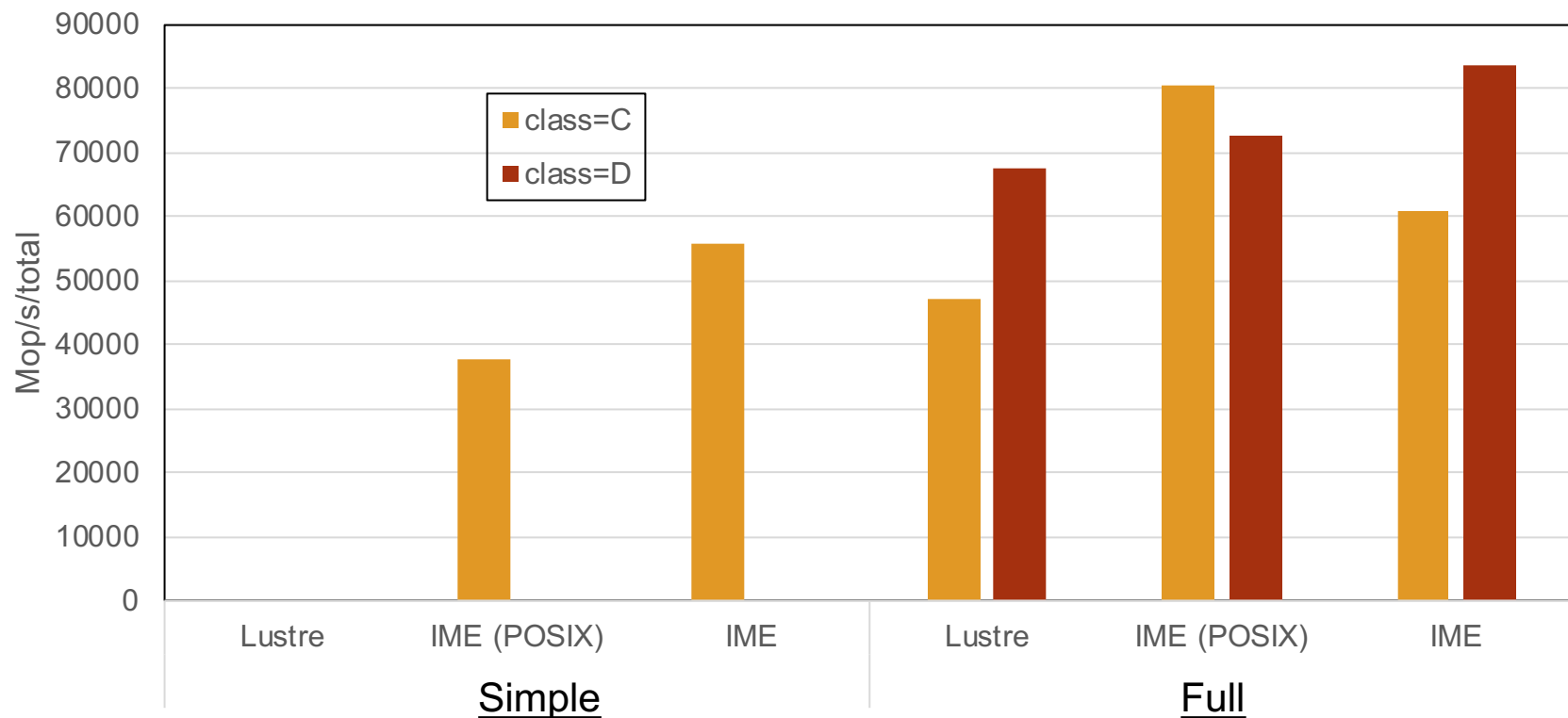
```
module switch impi mvapich2
```

- b. ファイルパスに ime:// を追
加

```
MPI_File_open(MPI_COMM_  
WORLD, "ime:///work/gt00/t002xx...",  
... );
```

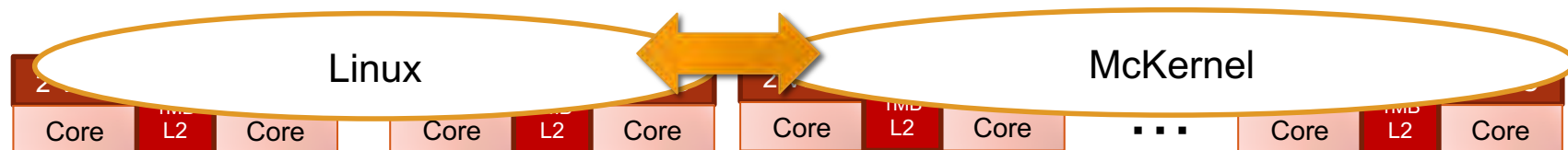

IMEの性能

- NAS Parallel benchmark, BT-IOの例
 - class=C, 64プロセス (16ノード x 4)
 - LustreでのSimpleは時間内に終了せず
 - class=D, 256プロセス(16ノードx16)
 - LustreでのSimpleは時間内に終了せず, IMEでのSimpleはエラー



McKernel

- メニーコアプロセッサ向けのOS
<http://www.pccluster.org/ja/mckernel/index.html>
 - 軽量カーネル
 - Linux 100%互換：Linuxで動作するプログラムは再コンパイル不要でMcKernelでも動作
 - OSによるユーザプログラムへの擾乱(OS Jitter)ゼロ
 - 「ポスト京」スパコンにも搭載予定
- OFPでは，まもなく利用可能になる予定
- OFPでの動作イメージ： 4コア Linux(OSサービス)
+ 64コア McKernel（ユーザーアプリ）



非同期処理の高速化 (IMPI 2019.1以降)

- MPI 3.Xの非同期集団通信、MPI 1.Xの1対1通信
 - **MPI_Ixxx()** の形: MPI_Isend, MPI_Irecv, MPI_Ibcast, ...
 - 事実上MPI_Wait()を呼ぶまでサボる処理系が多い
 - これまでは、途中でMPI_TestやMPI_Probeを呼んで働かせるしかなかった
- 空いているコアを有効に使って通信処理を進めることが可能に
- ジョブスクリプト例:

```
module switch intel intel/2019.1.144
source /work/opt/local/cores/intel/impi/2019.1.144/intel64/bin/mpivars.sh release_mt
# 上の2行はコンパイル時にも必要

export I_MPI_ASYNC_PROGRESS=on
export I_MPI_ASYNC_PROGRESS_PIN=0,1,68,69,136,137,204,205
# ノード内で起動するプロセス数分必要(ここでは8プロセス)
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out # いつも通り
```

演習： サンプルプログラム集

サンプルファイル

- C言語版・Fortran90版共通ファイル：
Samples-knl.tar.gz
- tarで展開後、各サンプルとともにC言語とFortran90言語のディレクトリが作られる
 - [C/](#) : C言語用
 - [F/](#) : Fortran90言語用
- 上記のファイルが置いてある場所
[/work/gt00/z30105](#) ([/home](#)でないので注意)

プログラムをコンパイルしよう(1/2)

1. `cd` コマンドを実行して Lustreファイルシステムに移動する
`$ cd /work/gt00/t001xx`
1. `/work/gt00/z30105` にある `Samples-knl.tar.gz` を
自分のディレクトリにコピーする
`$ cp /work/gt00/z30105/Samples-knl.tar.gz ./`
2. `Samples-knl.tar.gz` を展開する
`$ tar xvfz Samples-knl.tar.gz`
3. `Samples` フォルダに入る
`$ cd Samples-knl`
4. C言語 : `$ cd C`
Fortran90言語 : `$ cd F`
5. アプリ名のフォルダに入る
`$ cd stream`

プログラムをコンパイルしよう(2/2)

6. make する

```
$ make
```

7. 実行ファイル(stream)ができていることを確認する

```
$ ls
```

プログラムを実行しよう

1. streamフォルダ中で以下を実行する
`$ pjsub stream.bash`
2. 自分の導入されたジョブを確認する
`$ pjstat`
3. 実行が終了すると、以下のファイルが生成される
`stream.bash.eXXXXXX`
`stream.bash.oXXXXXX` (XXXXXXは数字)
4. 上記の標準出力ファイルの中身を見してみる
`$ cat stream.bash.oXXXXXX`
5. エラー出力は以下のファイルにあるので念のため確認
`$ cat stream.bash.eXXXXXX`

matmul-omp

- `cd Samples-knl/C/matmul-omp`
- or
- `cd Samples-knl/F/matmul-omp`

- `make`
- `cp /work/gt00/z30105/mat-mat-openmp.bash ./`
- `mat-mat-openmp.bash`を編集
- `pjsub mat-mat-openmp.bash`

SCORE-P & Scalasca

性能解析ツール: 「京」や「富岳」、FX100等でも使える (ようになる)

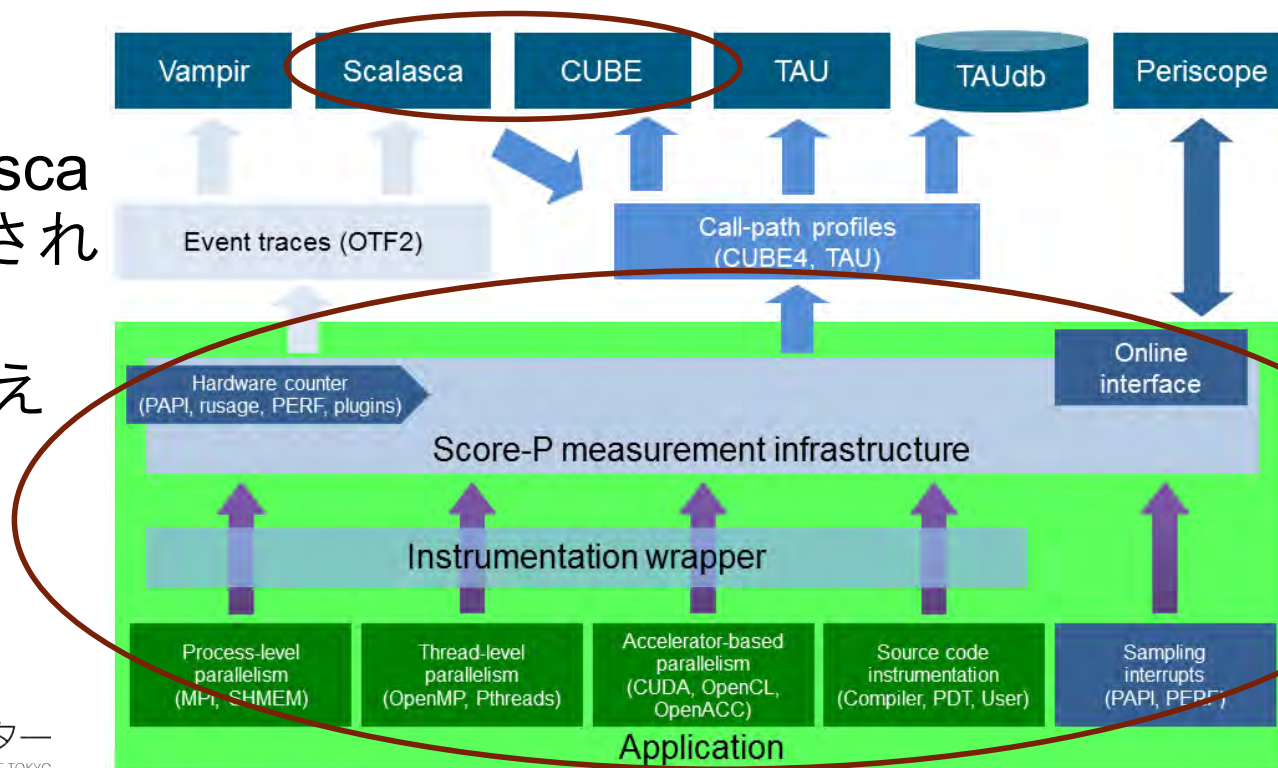
Performance Measurement

- <https://www.vi-hps.org/projects/score-p/>

Performance Analysis

- <https://www.scalasca.org>

- CUBEはScalascaと一緒に配布されている
- PAPIなども使えるようにbuild



実習 : SCORE-P & Scalasca

- 例 : NAS Parallel Benchmarks BT-MZ (NPB3.4-MZ)
- <https://www.nas.nasa.gov/publications/npb.html>
 - 今回はOFP+Scalasca用に少し修正したものを使う
- Scalasca向け環境設定 (ジョブスクリプトでも忘れずに)

```
$ source /work/gt00/share/scalasca.bash
```

```
$ cp /work/gt00/z30105/NPB3.4-MZ.tar.gz
```

```
$ tar xvfz NPB3.4-MZ.tar.gz
```

```
$ cd NPB3.4-MZ/NPB3.4-MZ-MPI
```

- コンパイル時、コンパイラの前に **scorep** をつける

```
$ make bt-mz CLASS=B NPROCS=8 PREP="scorep"
```

```
$ ls bin
```

```
bt-mz.B.x    bt-mz.B16.sh  bt-mz.B16.anlys.sh
```

解析の例：基本的なプロファイリング

```
$ pjsub bt-mz.B16.sh
```

終了してから

```
$ ls
```

```
scorep_bt-mz_B.16x64.2607028
```

```
scorep_bt-mz_B.16x64.2607028-2
```

のようなものが増えているはず

テキスト表示

```
$ cube_stat -r MAIN__ scorep_bt-  
mz_B.16x64.2607028/profile.cubex
```

```
Routine,time
```

```
INCL(MAIN__),2037.46
```

```
EXCL(MAIN__),0.0405123
```

```
setup_mpi_,8.23516
```

```
...
```

解析結果の表示 (GUI)

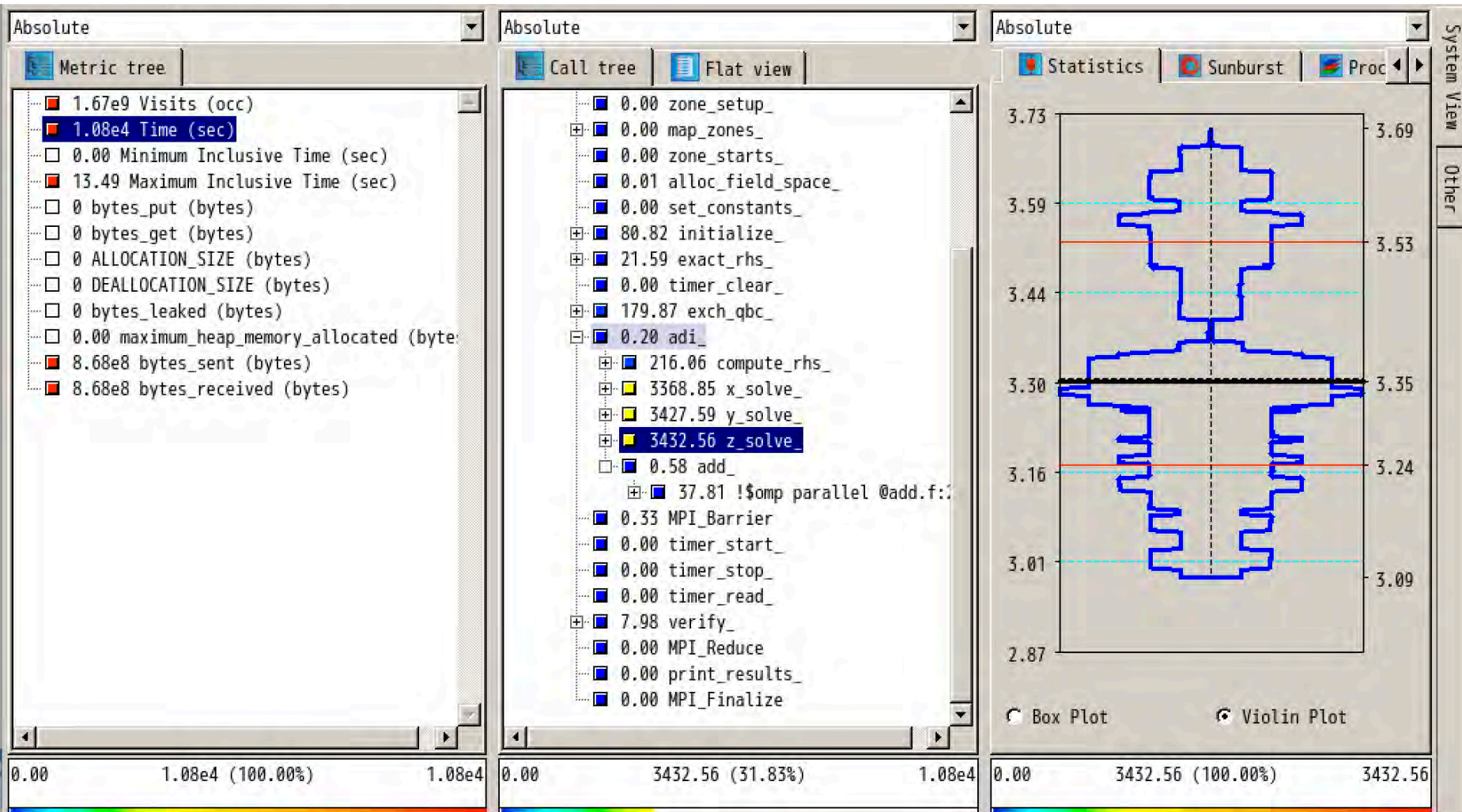
```
$ cube scorep_bt-  
mz_B.16x64.2607028/profile.cubex &  
$ cube scorep_bt-mz_B.16x64.2607028-  
2/profile.cubex &
```

- X Windowのウィンドウが起動する
(sshするときに `-Y` オプション、ノートPC側で表示できる環境が必要)
- cubexファイルをノートPCに転送し (ソースもあった方が便利)、ビューアで確認することもできる

解析結果の例

\$ cube profile.cubex

- スレッド毎の実行時間の分布がわかる



解析の例：より高度な解析 (Scalasca)

```
$ pjsub bt-mz.B16.anlys.sh
```

終了してから

```
$ scalasca -examine scorep_bt-mz_B.16x64.2607029
```

```
$ scalasca -examine scorep_bt-mz_B.16x64.2607029-2
```

```
$ ls scorep_bt-mz_B.16x64.2607029
```

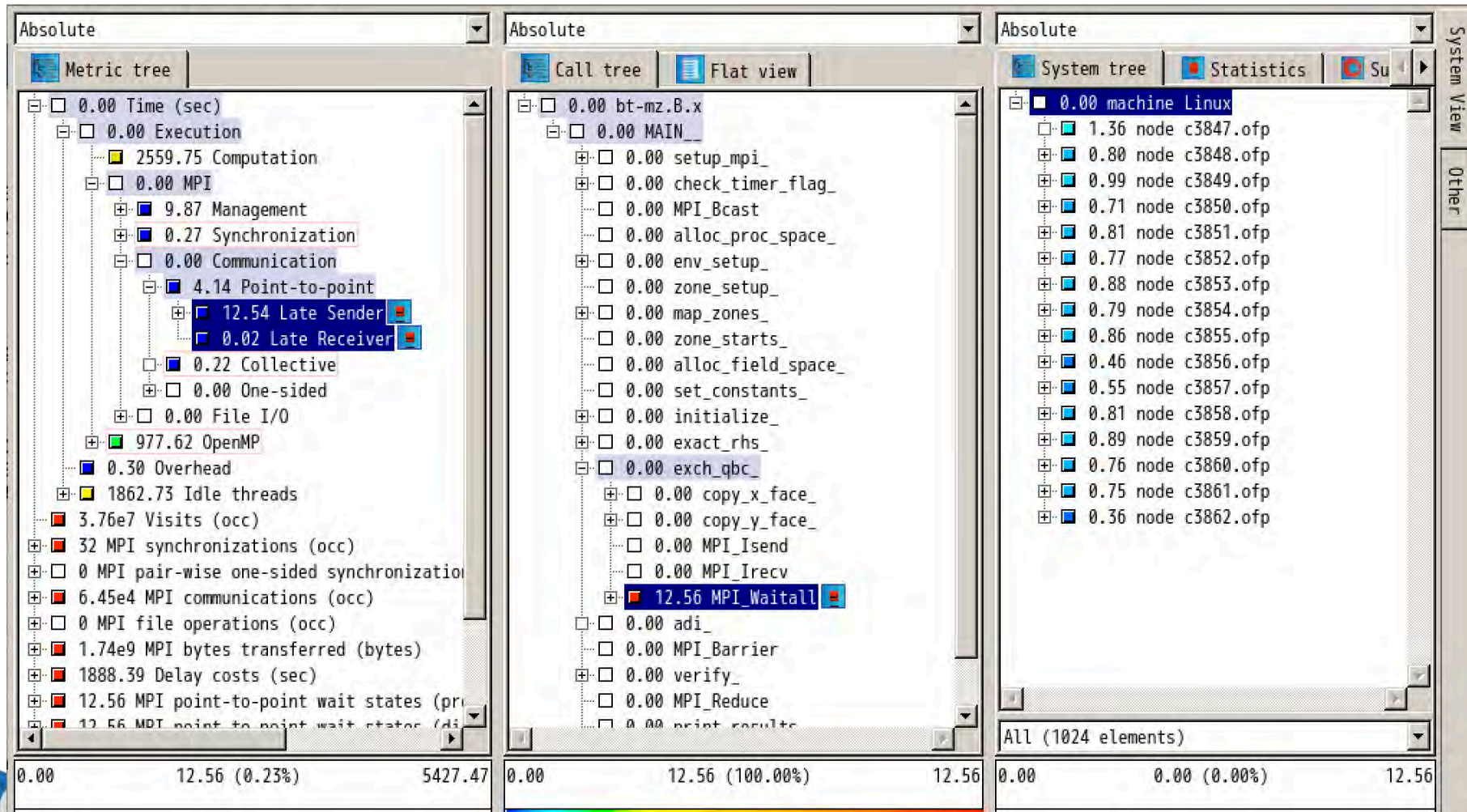
```
MANIFEST.md  profile.cubex  scorep.cfg  
scorep.filter  scorep.log  scout.cubex  scout.log  
summary.cubex  trace.cubex  trace.stat  traces  
traces.def  traces.otf2
```

```
$ cube scorep_bt-mz_B.16x64.2607029/trace.cubex &
```

```
$ cube scorep_bt-mz_B.16x64.2607029-2/trace.cubex &
```

Scalascaによる解析例

\$ cube trace.cubex



参考 : Scalasca & SCORE-P

A full workflow example of Scalasca

- http://apps.fz-juelich.de/scalasca/releases/scalasca/2.5/docs/manual/start_workflow.html

GUIツール: CUBE

- 手元のノートPCにビューアを入れておけば確認が楽
- <http://www.scalasca.org/software/cube-4.x/download.html>