

TSUBAME 2.0における 大規模GPUアプリケーション

東京工業大学 学術国際情報センター
教授

松岡 聡

筑波大学CCS

「先端学際計算科学共同研究拠点」シンポジウム第2回

2011年9月12日



Tsubame2.0 (2010-14)

x30 speedup c.f. Tsubame 1
(2006-2010)

2.4 Petaflops, 1408 nodes

~50 compute racks + 6 switch racks

Two Rooms, Total 160m²



1.4MW (Max, Linpack), 0.48MW (Idle)



TSUBAME2.0 Nov. 1, 2010

World's Smallest Petaflops Supercomputer



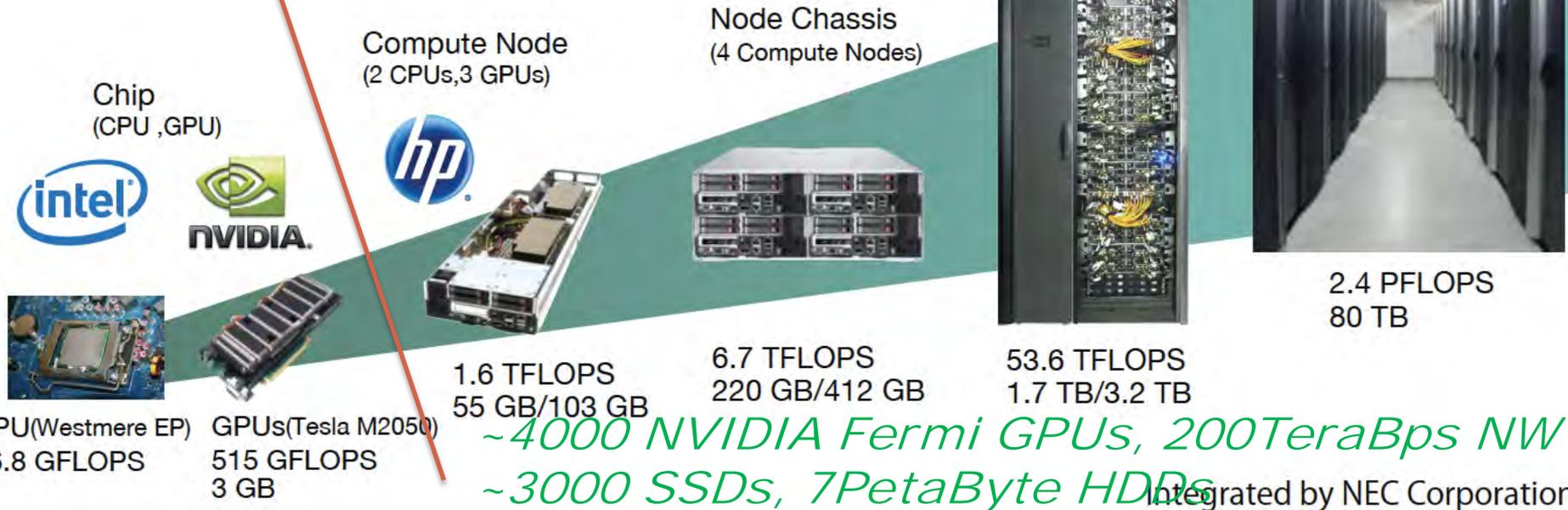
TSUBAME2.0: A GPU-centric Green 2.4 Petaflops Supercomputer

Tsubame 2.0: "Tiny" footprint, very power efficient

- Floorspace less than 200m² (2,100 ft²)
- Top-class power efficient machine on the Green 500

System
(42 Racks)
1408 GPU Compute Nodes,
34 Nehalem "Fat Memory" Nodes

TSUBAME 2.0 New Development





SL390 Compute Node

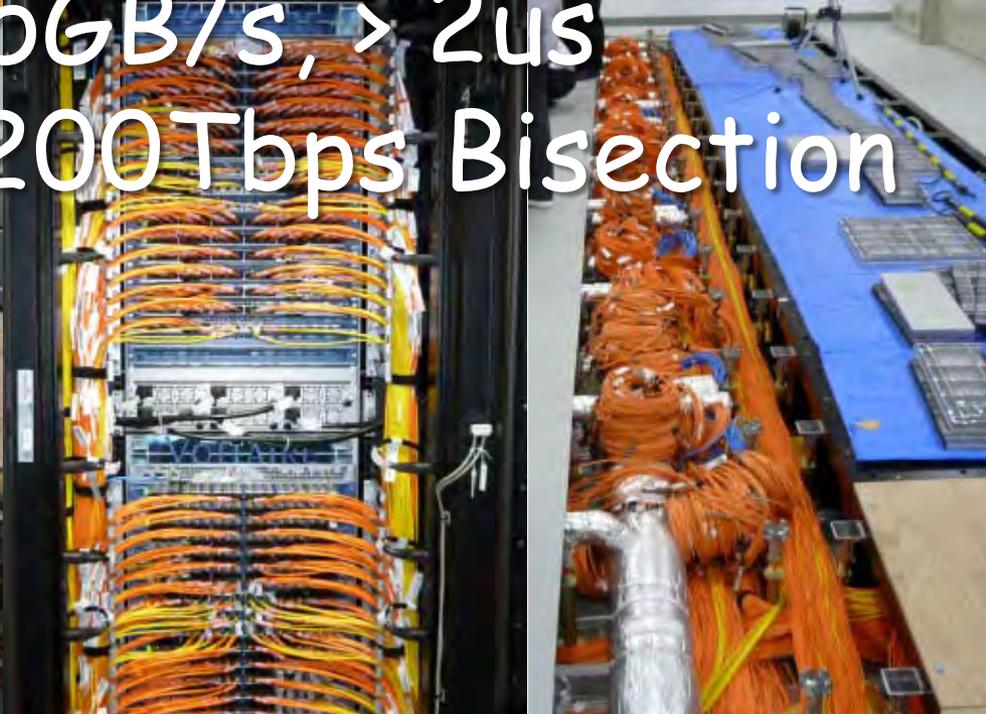
Collaborative Development w/HP

3 GPUs, 2 CPUs, 50-100GB Mem
120-240GB SSD, QDR-IB x 2

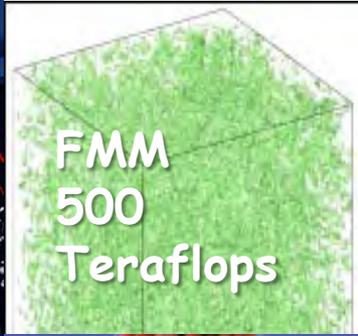
MADE IN T



3500 Fiber Cables > 100Km
w/DFB Silicon Photonics
End-to-End 6.5GB/s, > 2us
Non-Blocking 200Tbps Bisection



Tsubame 2.0's Achievements



4th Fastest Supercomputer in the World (Nov. 2010 Top500)

| Rank | Site | Computer/Year Vendor | Cores | R _{max} | R _{peak} | Power |
|------|---------------------------------------------------------|--------------------------------------------------------------------------------------------|--------|------------------|-------------------|---------|
| 1 | National Supercomputing Center in Tianjin China | Tianhe-1A - NUDT TH MPP, X5670 2.93Ghz 6C, NVIDIA GPU, FT-1000 8C / 2010 NUDT | 186368 | 2566.00 | 4701.00 | 4040.00 |
| 2 | DOE/SC/Oak Ridge National Laboratory United States | Jaguar - Cray XT5-HE Opteron 6-core 2.6 GHz / 2009 Cray Inc. | 224162 | 1759.00 | 2331.00 | 6950.60 |
| 3 | National Supercomputing Centre in Shenzhen (NSCS) China | Nebulae - Dawning TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU / 2010 Dawning | 120640 | 1271.00 | 2984.30 | 2580.00 |
| 4 | GSIC Center, Tokyo Institute of Technology Japan | TSUBAME 2.0 - HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows / 2010 NEC/HP | 73278 | 1192.00 | 2287.63 | 1398.61 |
| 5 | DOE/SC/LBNL/NERSC United States | Hopper - Cray XE6 12-core 2.1 GHz / 2010 Cray Inc. | 153408 | 1054.00 | 1288.63 | 2910.00 |



Fruit of Years of Collaborative Research - Info-Plosion, JST CREST Ultra Low Power HPC...

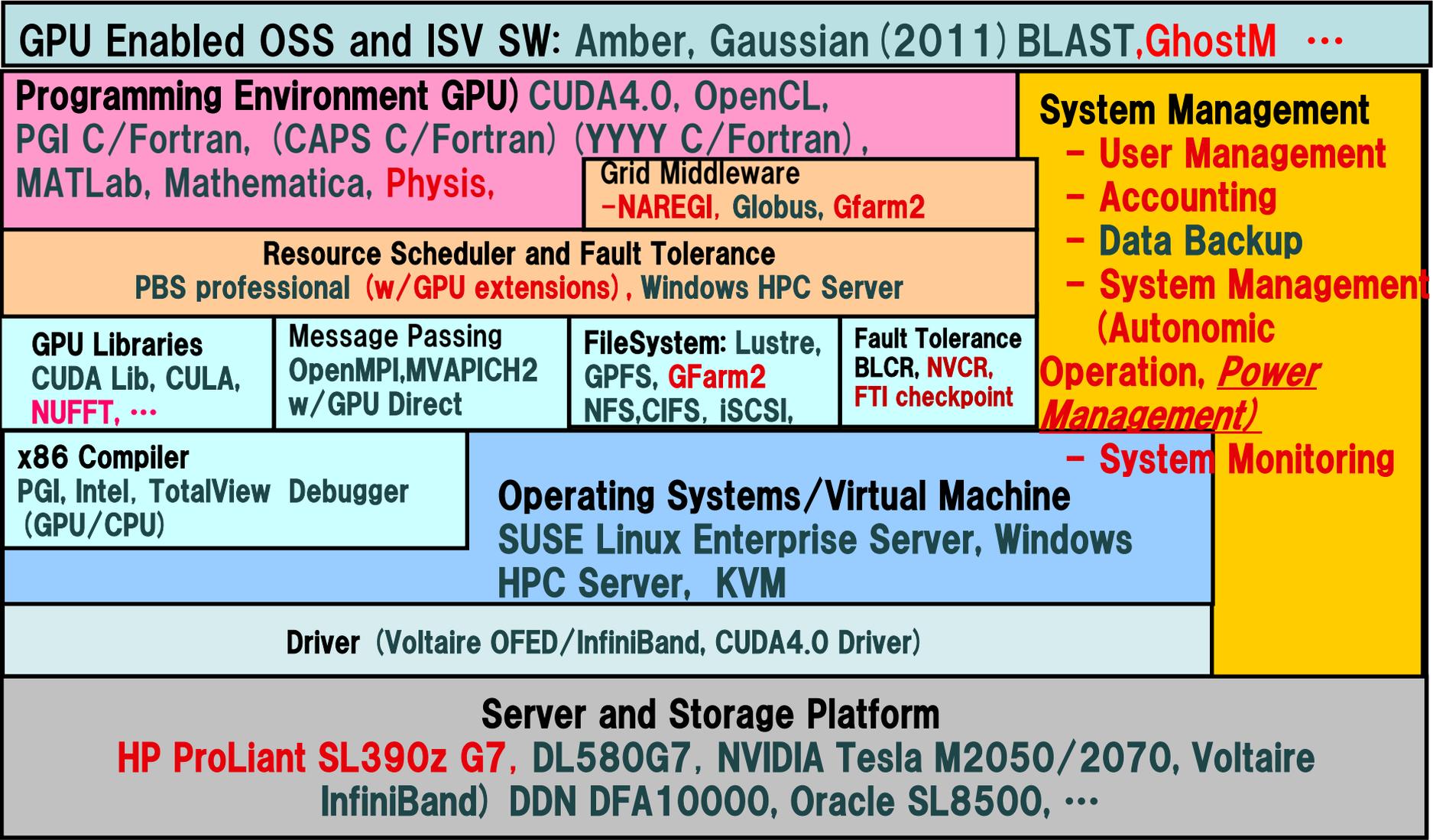


Over 10 Petascale Applications

x66,000 faster
x3 power efficient



TSUBAME2.0 Software Stack (Red: R&D at Tokyo Tech)

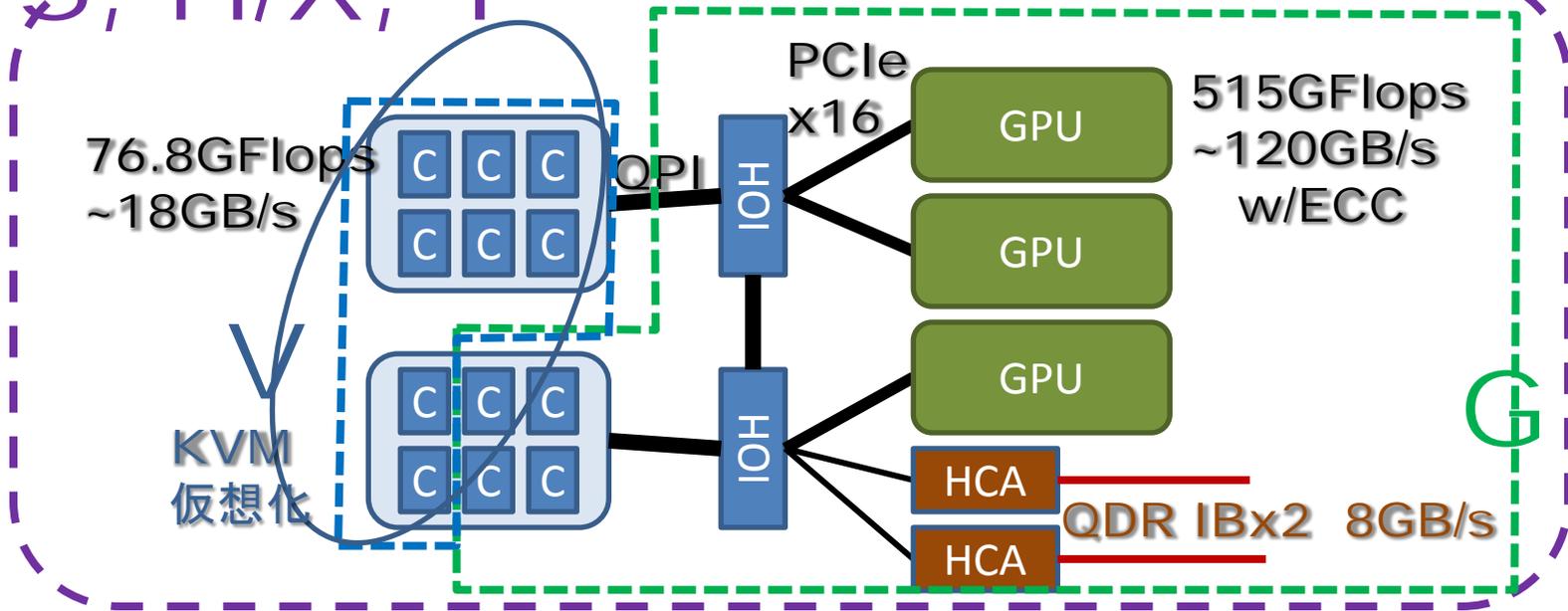


Service List

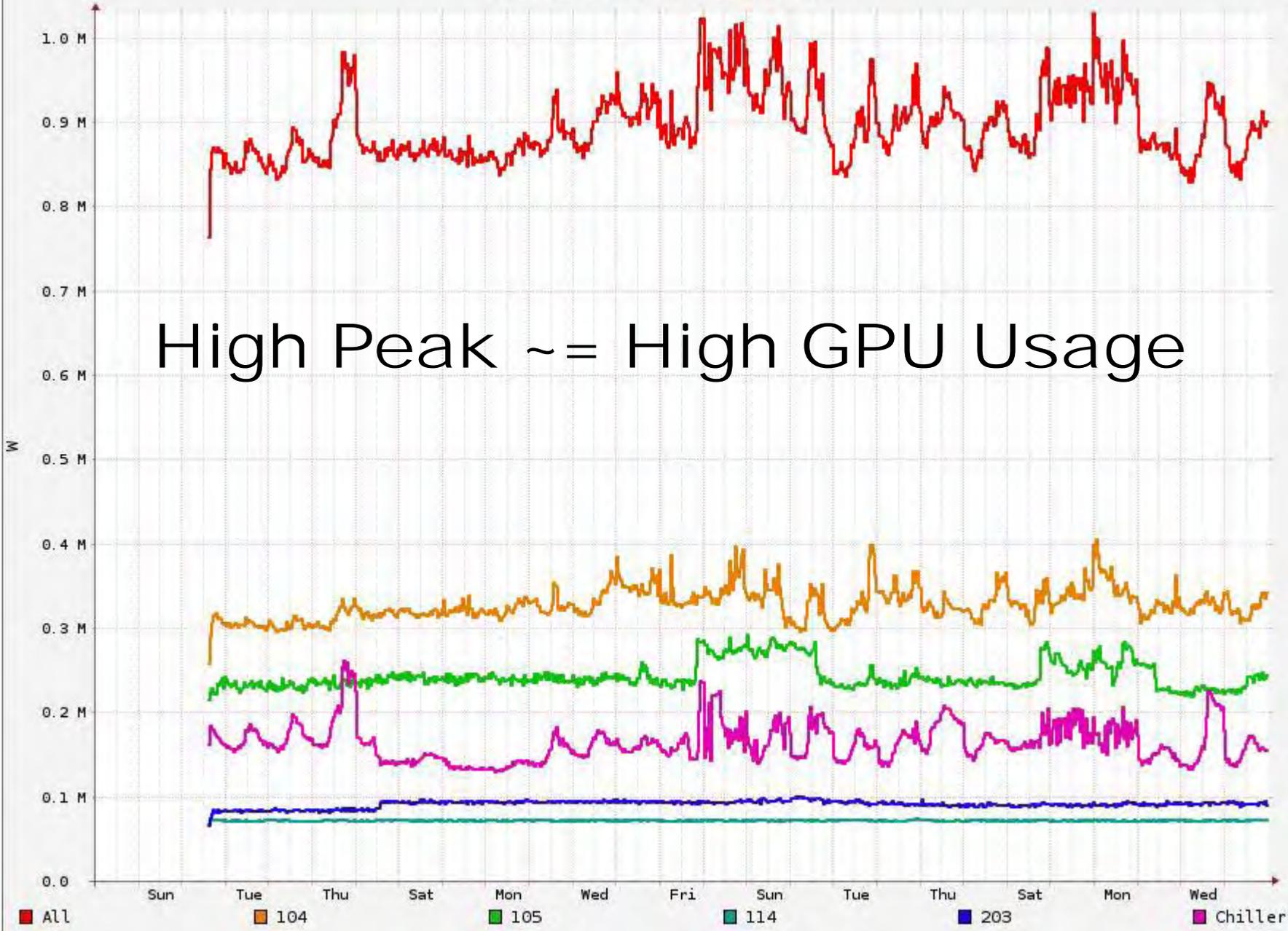
| service | assigned nodes | | | running jobs | users | |
|--------------|----------------|--------------------------|--|--------------|----------------|----|
| <u>S</u> | 85% | 255 / 300 nodes | | 95% | 88 / 92 jobs | 19 |
| <u>S96</u> | 15% | 6 / 40 nodes | | 100% | 6 / 6 jobs | 3 |
| <u>G</u> | 78% | 214 / 274 nodes | | 18% | 7 / 38 jobs | 3 |
| <u>V</u> | 90% | 198 / 220 nodes | | 86% | 397 / 459 jobs | 19 |
| <u>L128</u> | 30% | 3 / 10 nodes | | 100% | 3 / 3 jobs | 2 |
| <u>L128F</u> | 0% | 0 / 8 nodes | | 0% | 0 / 0 jobs | 0 |
| <u>L256</u> | 0% | 0 / 6 nodes | | 0% | 0 / 0 jobs | 0 |
| <u>L512</u> | 0% | 0 / 2 nodes | | 0% | 0 / 0 jobs | 0 |
| <u>H/X</u> | 85% | 32 (+ 321) / 410 nodes | | 94% | 16 / 17 jobs | 5 |
| <u>Y</u> | 98% | 192 / 194 nodes | | 40% | 2 / 5 jobs | 1 |

アカウント人数~5000人

S, H/X, Y



TSUBAME2 Grid Power last month



2010年11月Top500, Green500において TSUBAME2.0上位入賞

- 省エネ性能958MFlops/W ⇒ 世界2位!!
 - Greenest Production Supercomputer in the World賞獲得!!
- 演算性能1.192PFlops ⇒ 世界4位!!
 - 日本のスパコンで10位以内は4年ぶり, 5位以内は5年ぶり



| Rank | Site | Computer/Year | Vendor | TOP 500 [®] SUPERCOMPUTER SITES | | | |
|------|----------------------------------------------------------|----------------------------------------------------------------------------------------------|--------------------------------------------------------|---------------------------------------------|---------|---------|---------|
| 1 | National Supercomputing Center in Tianjin, China | Tianhe-1A - NUDT T1E | MPP, X5670 2.83GHz, NVIDIA GPU, FT-100 8C / 2010, NUDT | | | | |
| 2 | DOE/SC/Oak Ridge National Laboratory, United States | Jaguar - Cray XT5-HE | Opteron 8-core 2.6GHz / 2009, Cray Inc. | 224162 | 1759.00 | 2331.00 | 6950.60 |
| 3 | National Supercomputing Centre in Shenzhen (NSCS), China | Nebulae - Dawning TC3600 Blade, Intel X5850, Nvidia Tesla Q2050 GPU / 2010, Dawning | | 120640 | 1271.00 | 2984.30 | 2580.00 |
| 4 | GSIC Center, Tokyo Institute of Technology, Japan | TSUBAME 2.0 - HP ProLiant SL390s G7, Xeon 8C X5670, Nvidia GPU, Linux/Windows / 2010, NEC/HP | | 73278 | 1192.00 | 2287.63 | 1398.61 |
| 5 | DOE/SC/BNL/NERSC, United States | Hopper - Cray XE6, 12-core 2.1 GHz / 2010, Cray Inc. | | 153408 | 1054.00 | 1288.63 | 2910.00 |
| 6 | Commissariat à l'Énergie Atomique (CEA), France | Tera-100 - Bull bullx super-node S6010/S6030 / 2010, Bull SA | | 138368 | 1050.00 | 1254.55 | 4590.00 |



Heterogeneous Petascale Linpack

Results on TSUBAME 2.0

7:06 – 9:33, October 18, 2010

| T/V | N | NB | P | Q | Time | Gflops |
|-------------------------------------------------|---------|------|----|----|-----------|-----------|
| WR15R2R16 | 2490368 | 1024 | 59 | 69 | 8639.84 | 1.192e+06 |
| Ax-b _oo/(eps*(A _oo* x _oo+ b _oo)*N)= | | | | | 0.0008911 | PASSED |

1.192PFlops with 4071 GPUs

958MFlops/Watt (1244kW)

No.4 in Top500

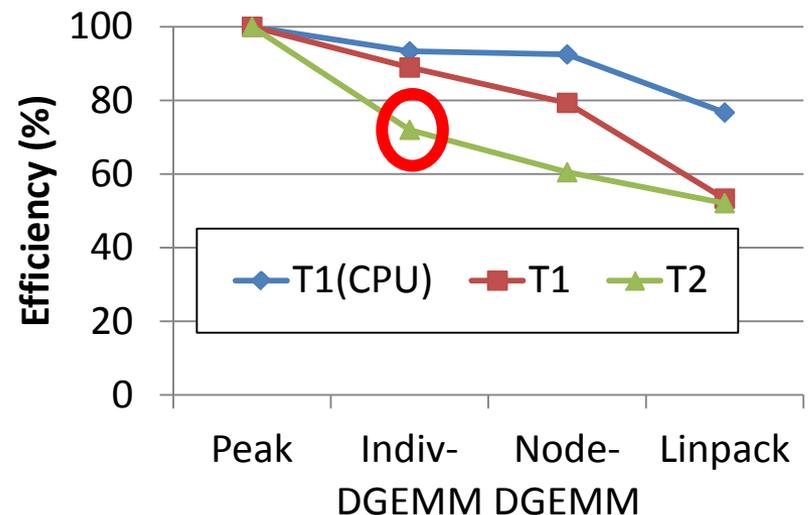
No. 2 in Green500 (Nov 2010)

Overhead analysis

Relative performance is 52%; is this overhead essential on GPU systems?



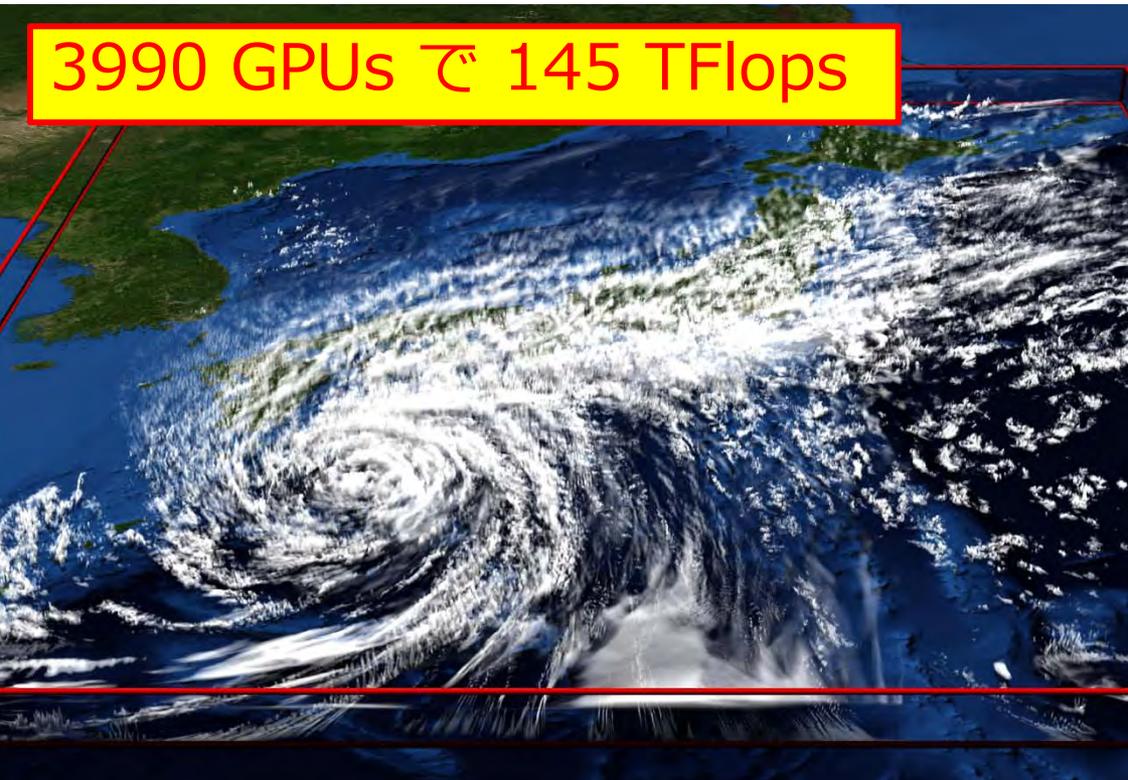
DGEMM kernel on Fermi GPUs exhibits only ~70% performance of peak, which was >90% on older gen.



TSUBAME 2.0 による大規模計算

- 樹枝状凝固成長 (SC11 Tech Paper/Gordon Bell Award Finalist)
- 気象計算 (SC10 Technical Paper/Best Student Paper Finalist)

3990 GPUs で 145 TFlops



4000GPU+16000CPU
cores で 1.017 Pflops



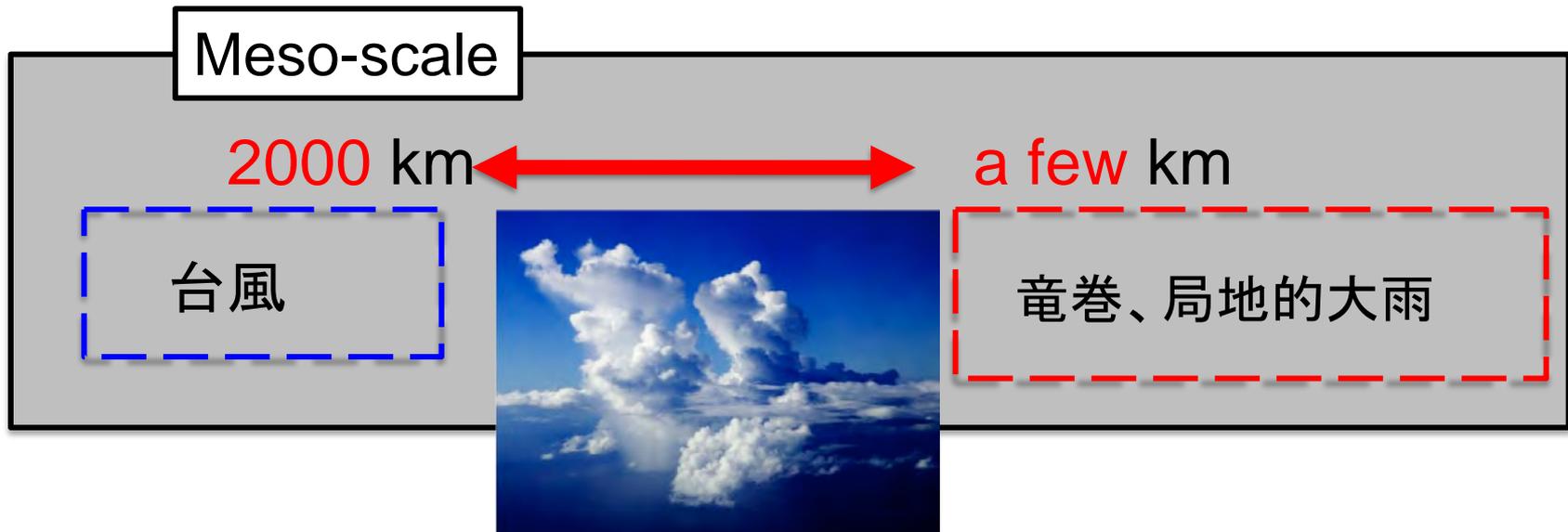
気象シミュレーションコードASUCA

■ ASUCA Production Code

- ✓ 気象庁によって開発が進められる次世代高解像度気象シミュレーションコード
- ✓ Compressible nonhydrostatic equations
Flux form, Generalized coordinate

■ なぜ GPU コンピューティング？

- ✓ 速い、安い、低消費電力→高解像度計算へ



高解像度気象計算の必要性

■ ゲリラ豪雨

- ✓ 突発的で局地的な豪雨（～km - ～10km）
- ✓ 参考：梅雨前線などによる集中豪雨（～100km）

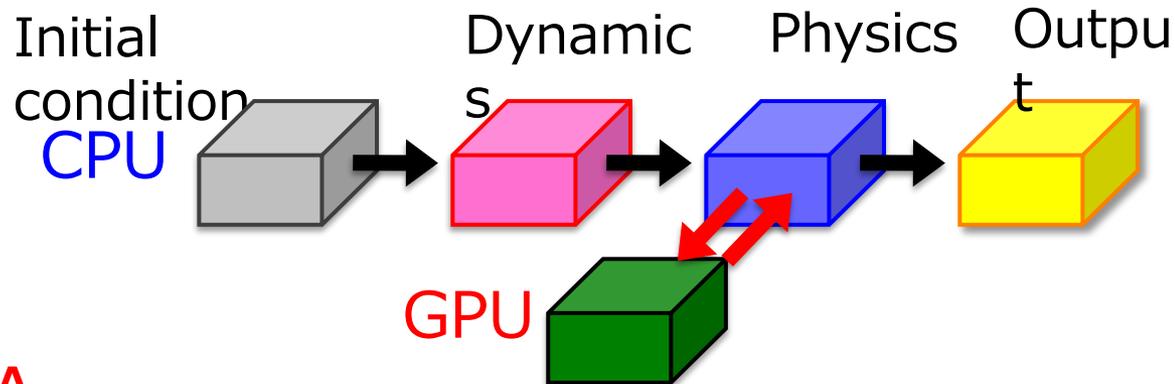
→ 現業の5km 以下の格子で計算することが必須



WRFとASUCAのGPUによる高速化の違い

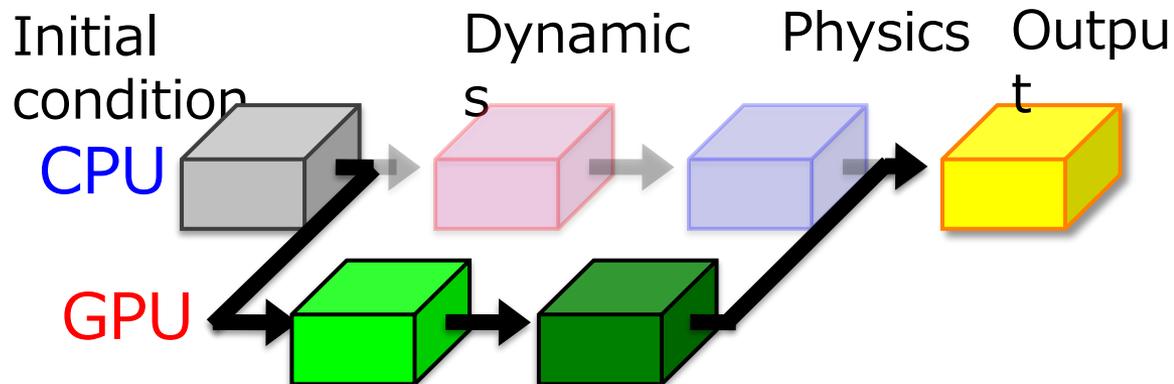
- ✓ WRF (世界的に使われている気象コード)

Accelerator Approach



- ✓ ASUCA

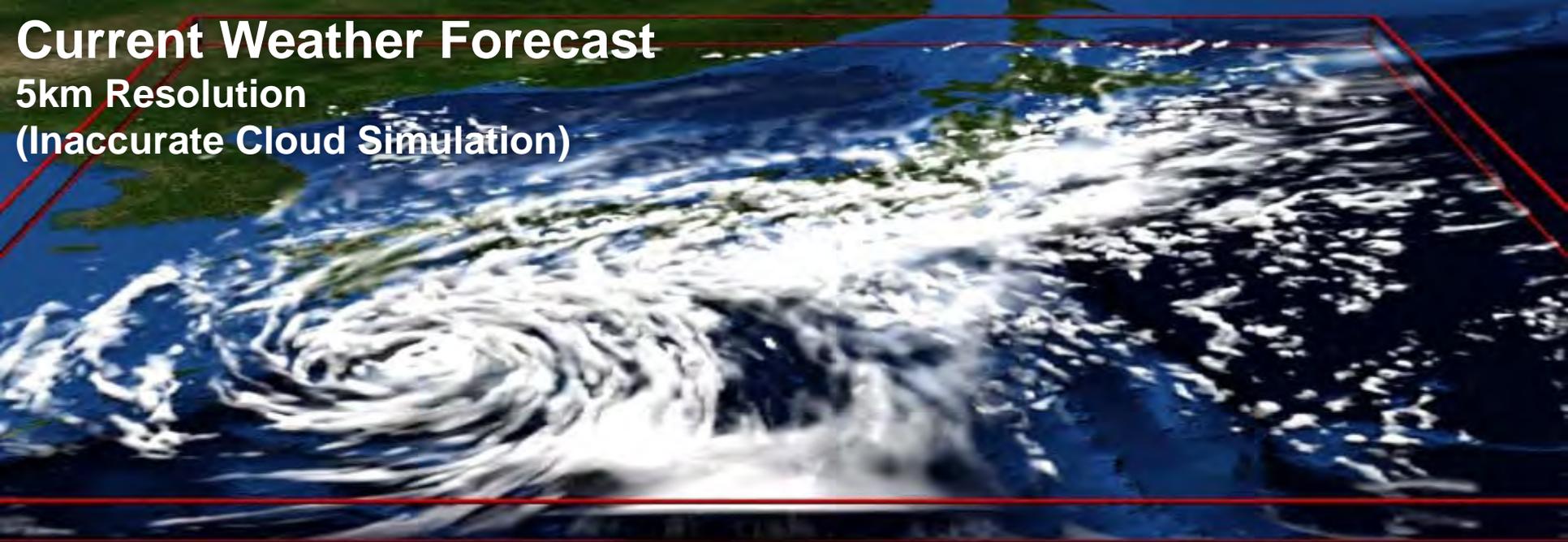
Full GPU Approach → 数十倍の高速化の実現



Current Weather Forecast

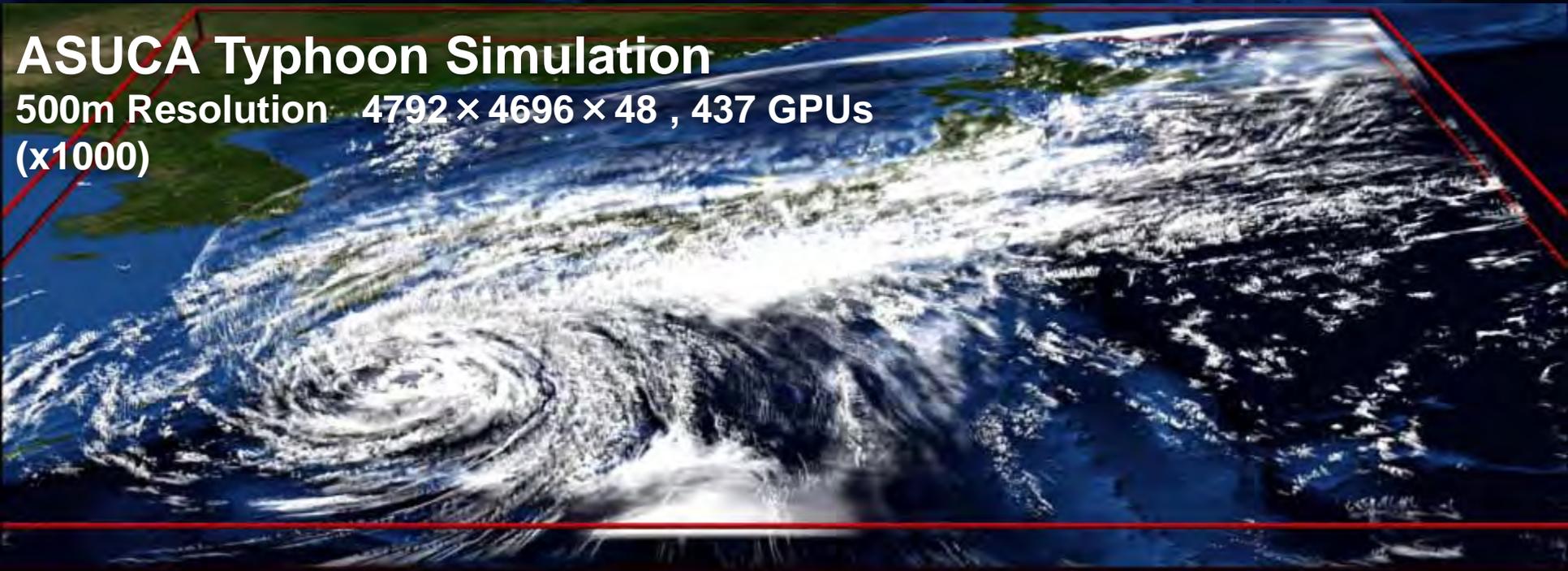
5km Resolution

(Inaccurate Cloud Simulation)

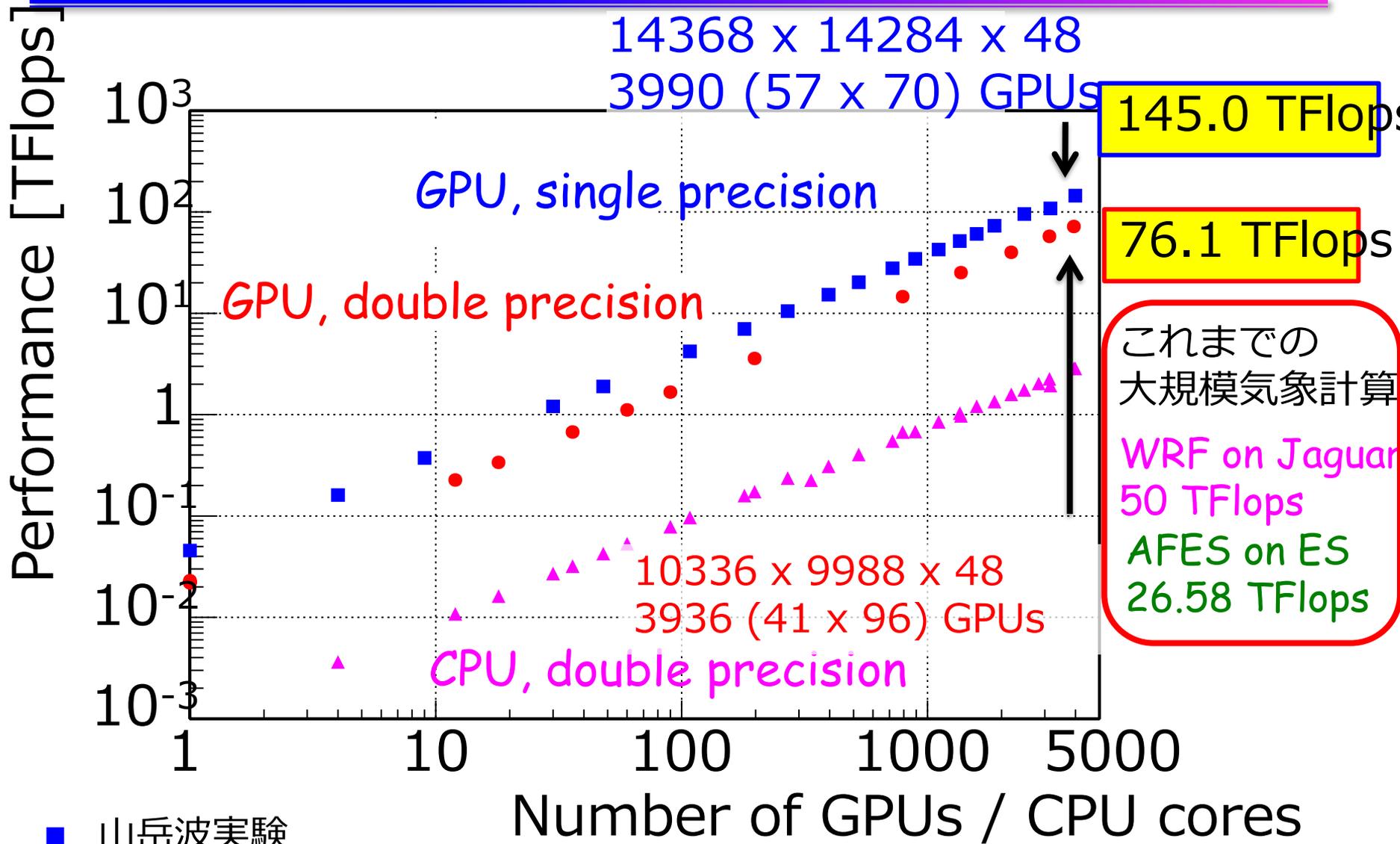


ASUCA Typhoon Simulation

500m Resolution $4792 \times 4696 \times 48$, 437 GPUs
(x1000)



TSUBAME 2.0を用いたASUCAの計算性能



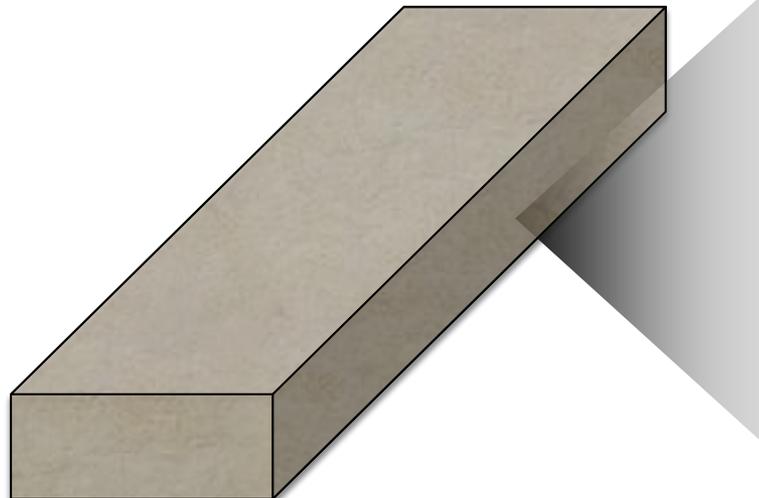
- 山岳波実験
- NVIDIA Tesla M2050 card / Intel Xeon X5670 2.93 GHz on TSUBAME 2.0

金属材料の機械的強度

鑄造

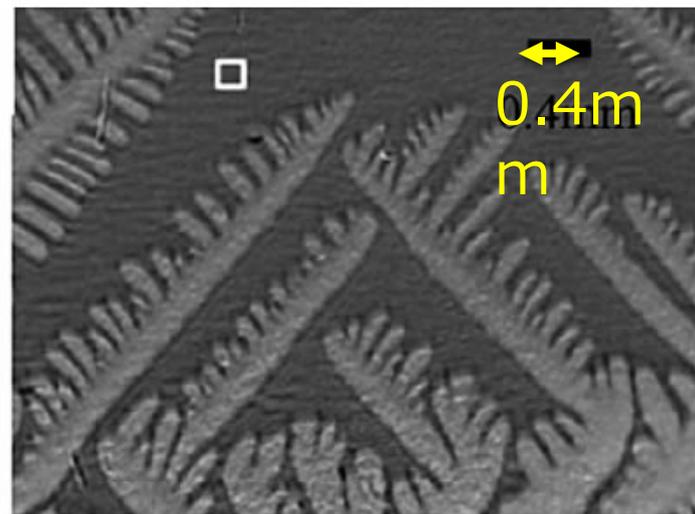


マクロ



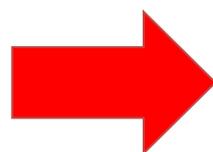
マクロな特性は
ミクロ構造に依存

ミクロ ($\sim\mu\text{m}$)



Spring-8 による実際の観察写真

金属材料の機械的強度や
特性の予測

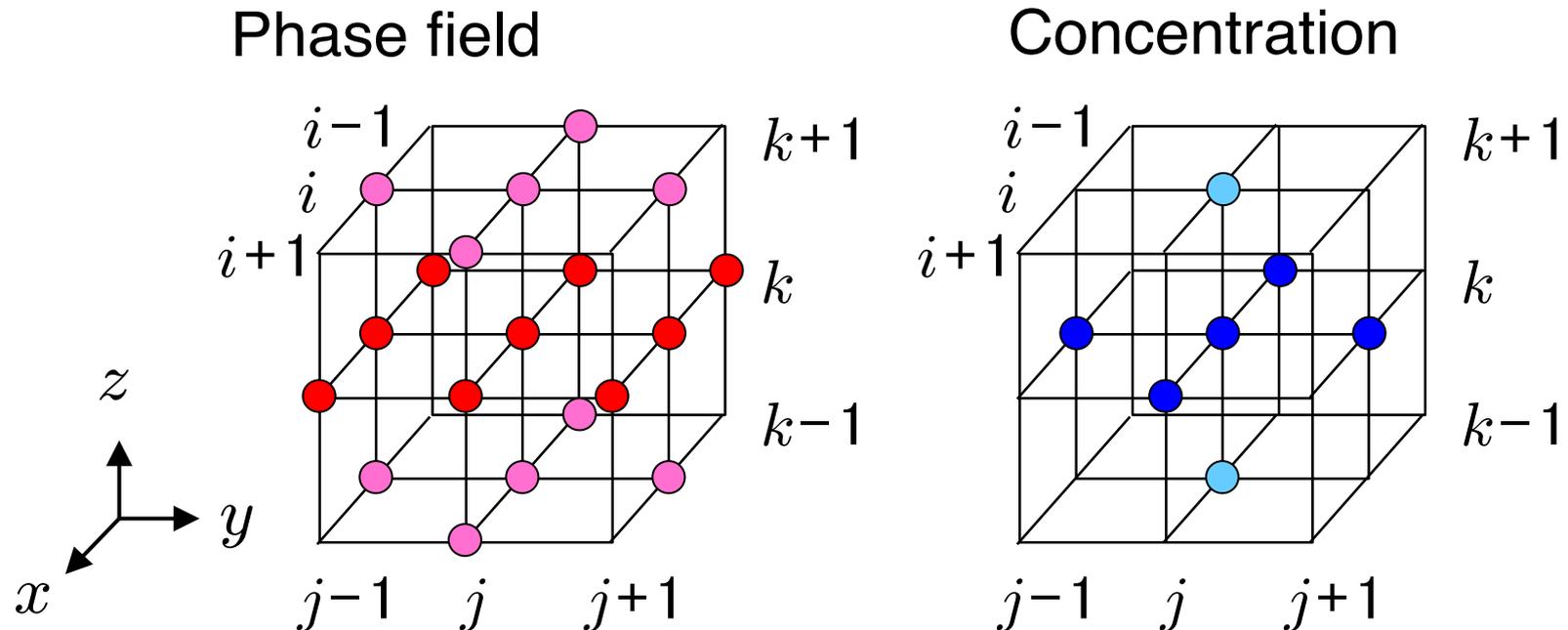


ミクロな組織構造に基づく
大規模シミュレーションが必要

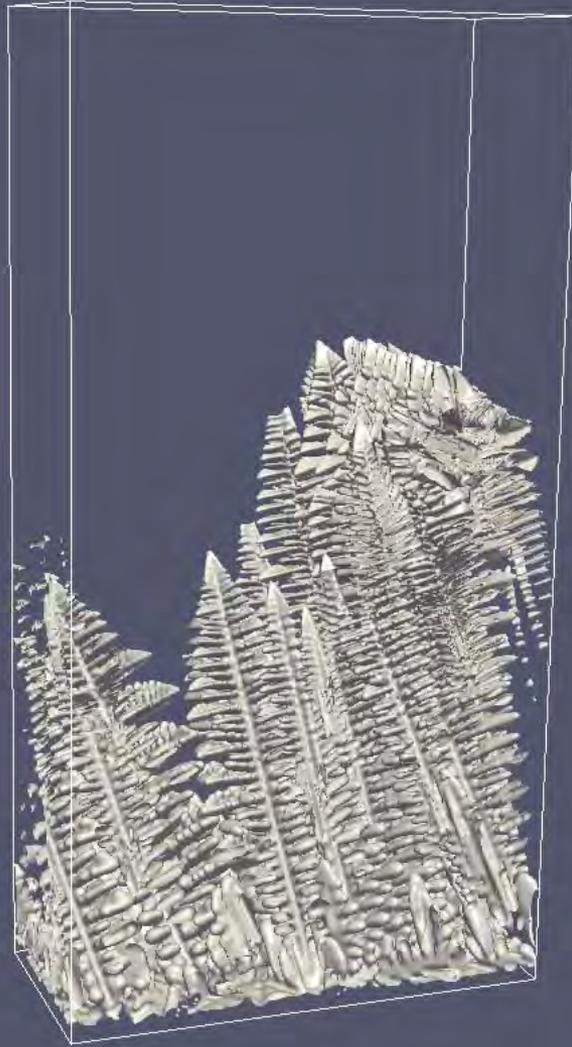
大規模GPU計算を行うことで、
材料組織予測し期待通りの材料を開発へ

フェーズフィールドのデータアクセス

- 空間を格子に分割して計算（気象計算と同じ）
- 1タイムステップ更新するために、
Phase-field変数は**19点**、濃度変数は**7点**
のデータ読み込みが必要。計算量も比較的大。
(気象計算よりもデータアクセスと計算量が多い)



樹枝状凝固成長の計算 (Al-Si 合金)



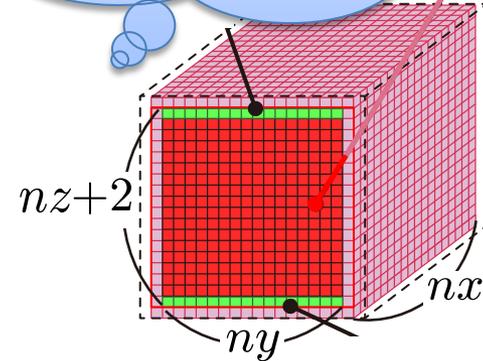
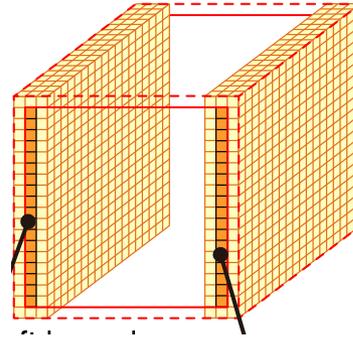
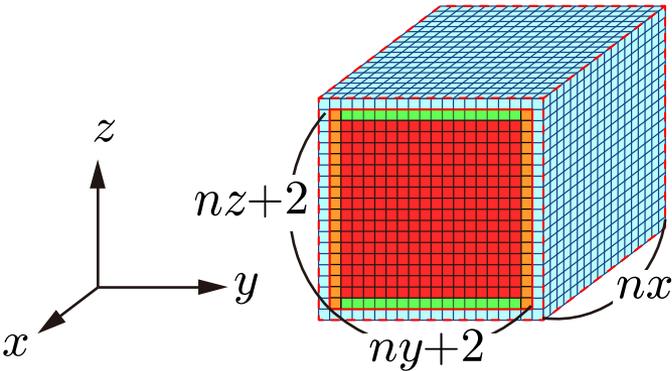
768 x 1632 x 3264 格子でTSUBAME 2.0の1156 GPU_sを利用

計算による通信の隠蔽

CPUとGPUを併用した
ハイブリッド計算

Whole subdomain

Divided domains



CPU計算

GPU計算

MPI通信/
GPU-CPU通信

CPU計算

GPU計算

y方向境界

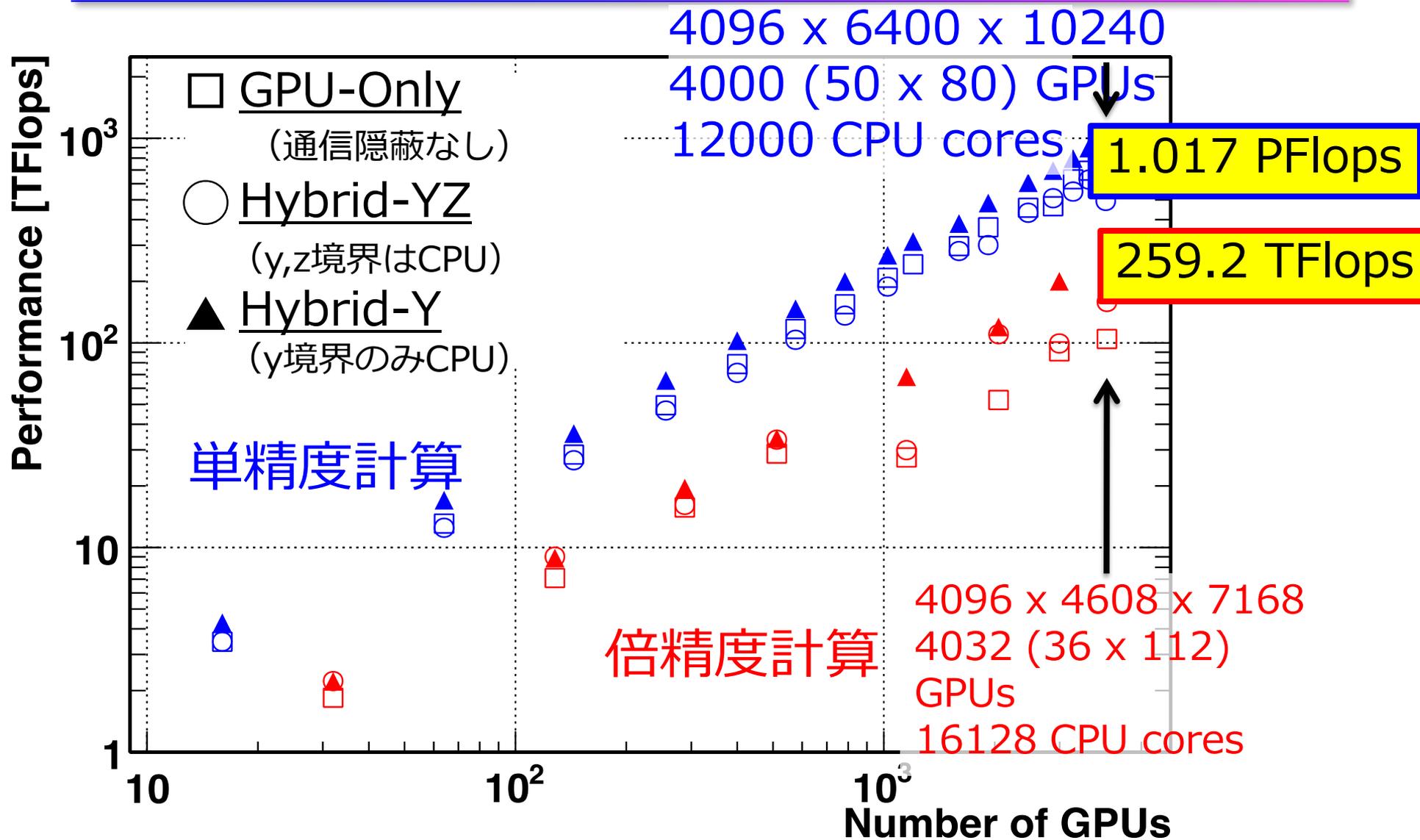
z方向境界

中心領域

time

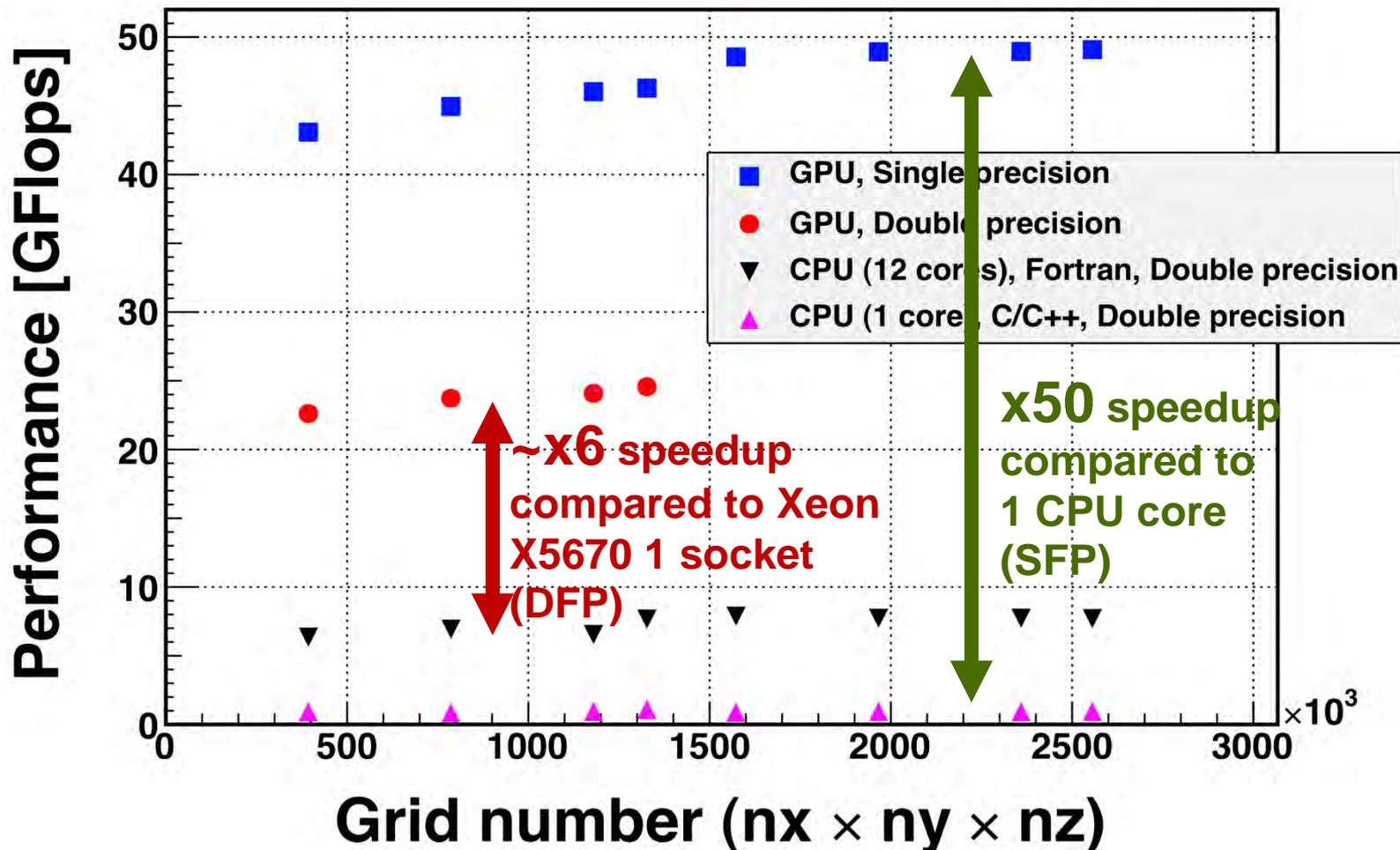
1サブドメインを1GPU + 4CPU cores (OpenMP)で計算
全計算領域を MPI による並列計算

TSUBAME 2.0での計算性能 (Weak scaling)

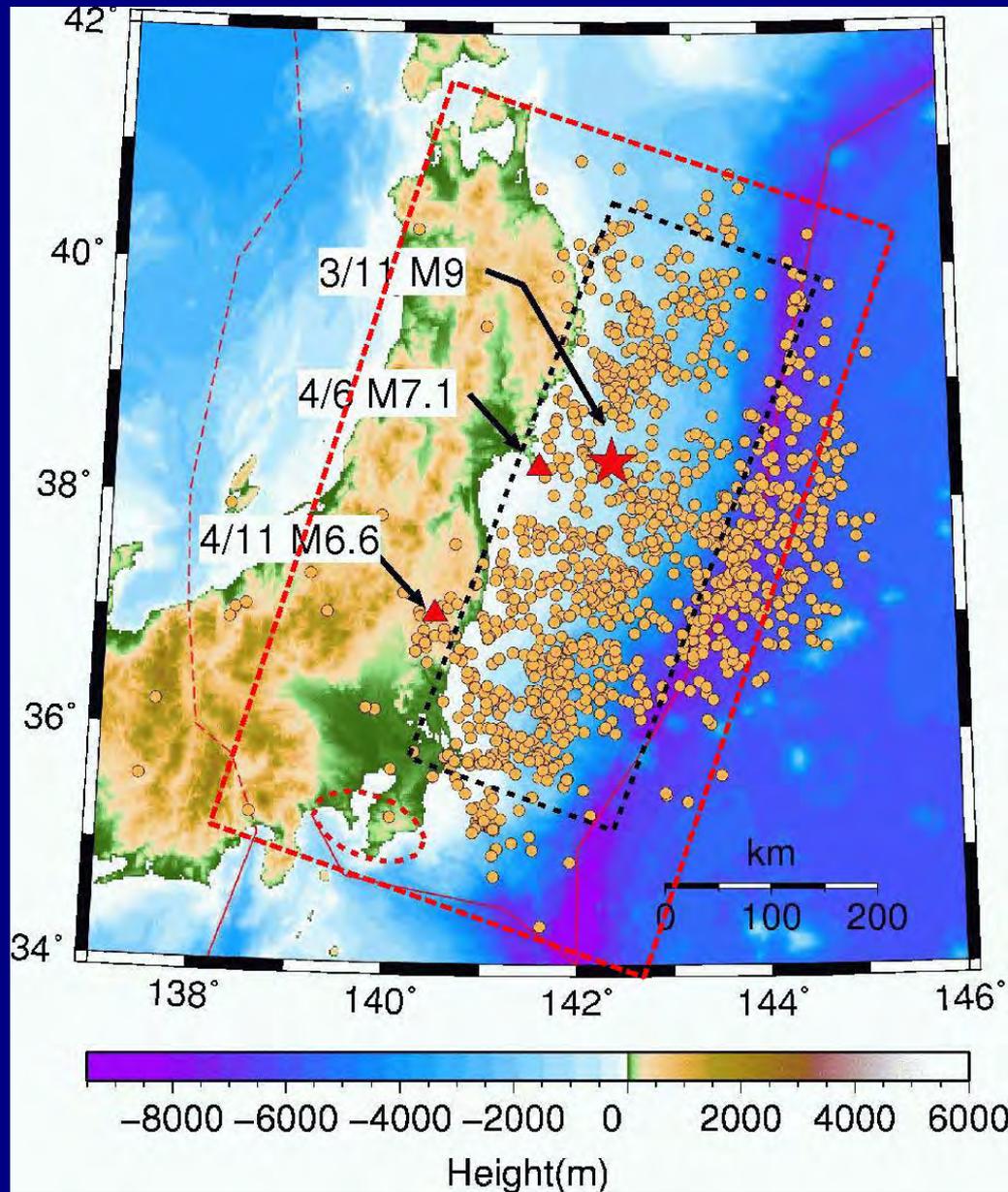


Weak Scaling: GPUあたりの問題サイズ一定で全体サイズを変える

TSUBAME 2.0 (1 GPU)



研究の背景



東北地方 太平洋沖地震 [東工大・岡本ら]

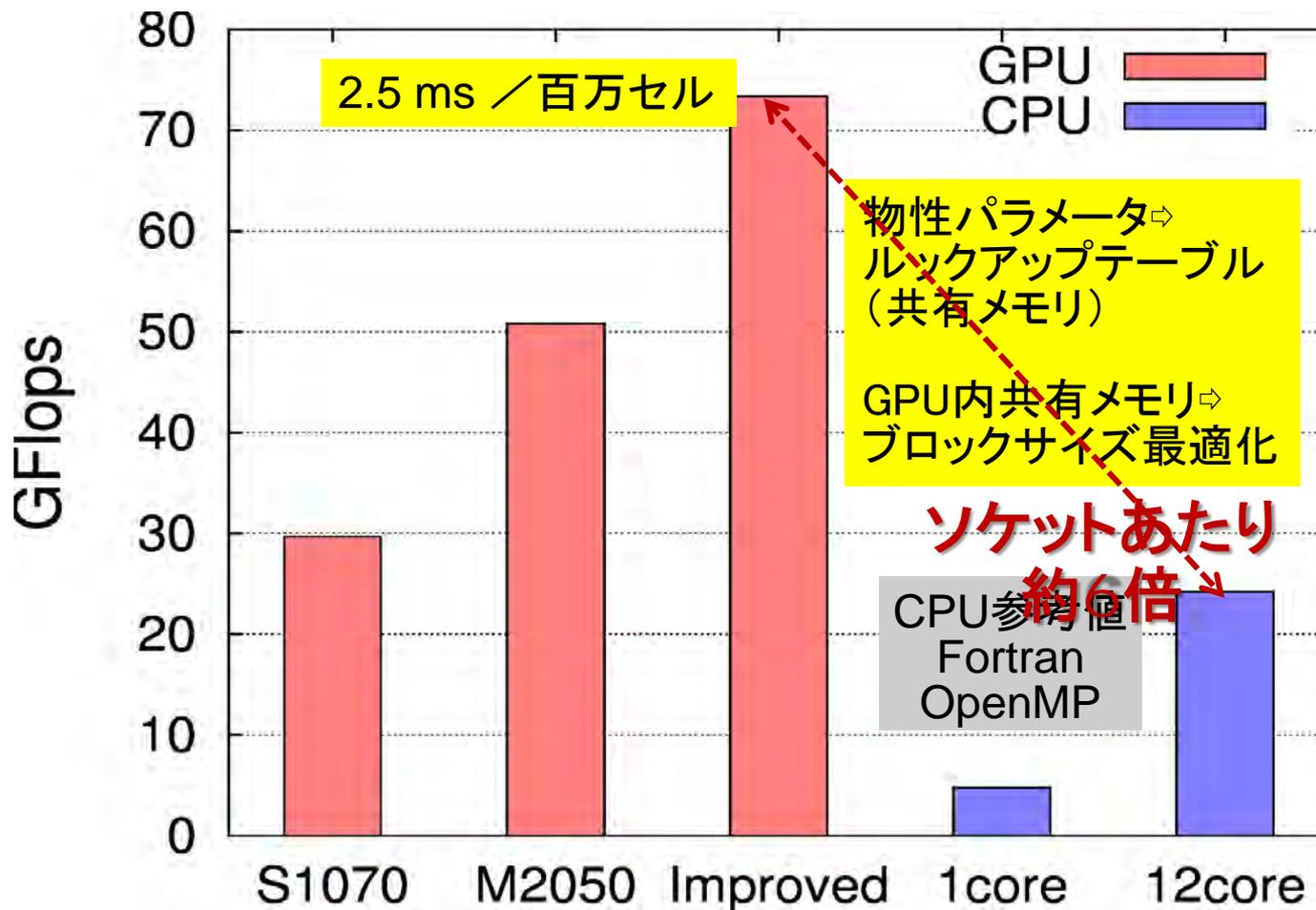
大規模計算が必要
850 x 500 x 200 km

弾性体
 $\Delta x = 125$ m
6800 x 4000 x 1600
1300 GPU

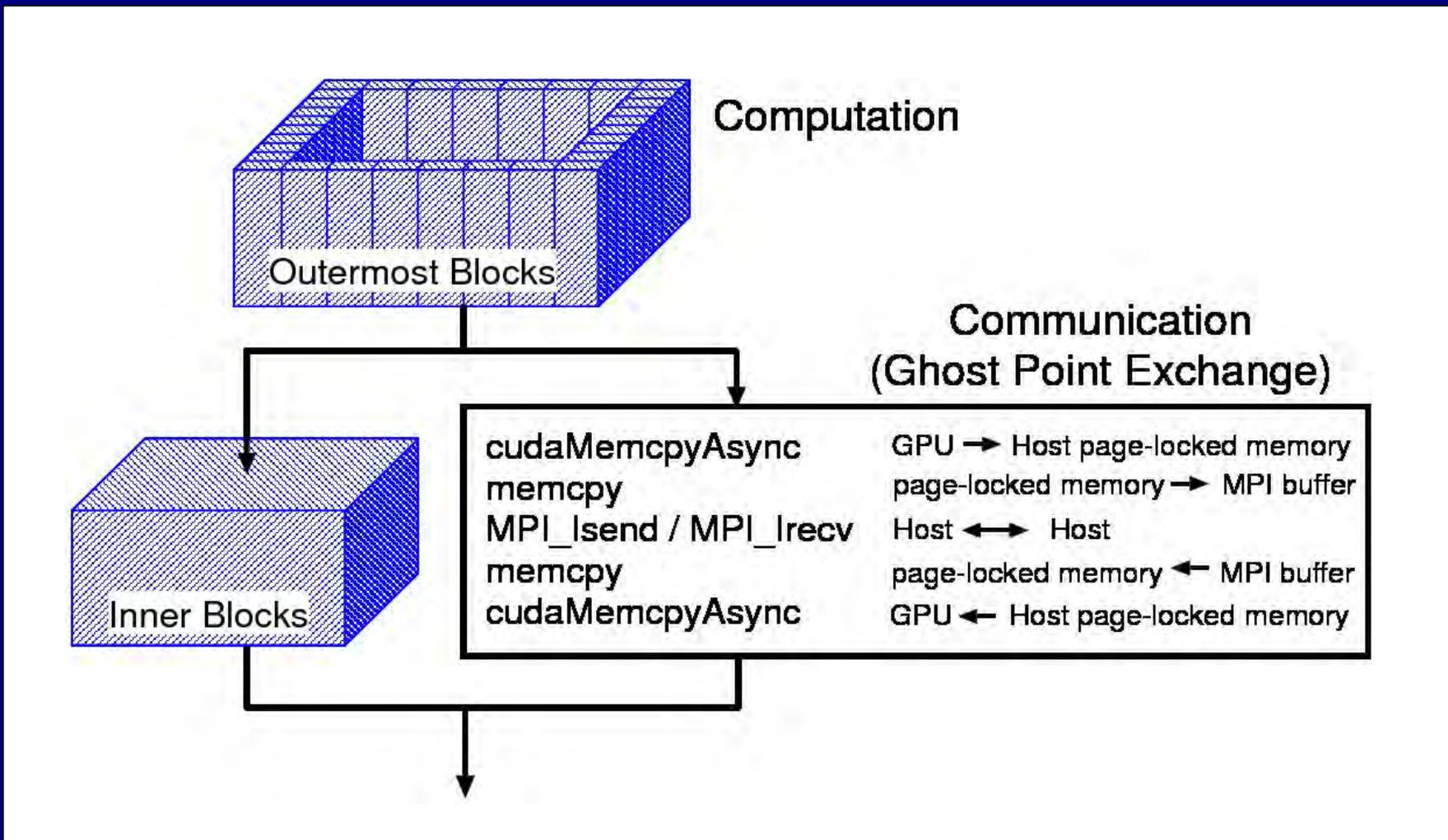
非弾性(5要素)
 $\Delta x = 200$ m
4250 x 2500 x 1000
1300 GPU

GPU1個での計算性能

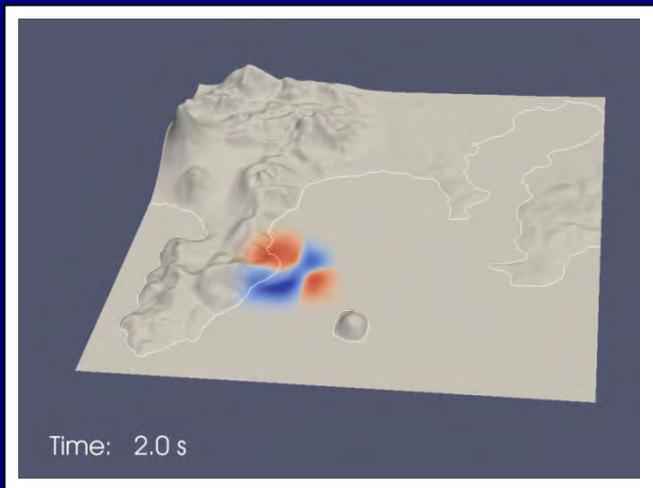
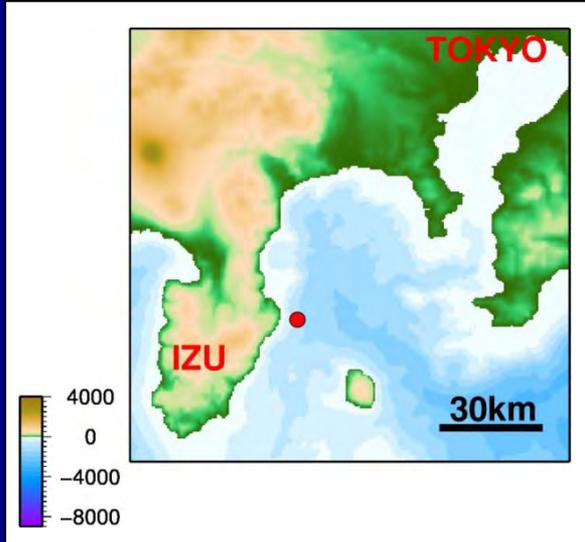
(GPU-MPI 並列コード: 単精度、弾性体)



計算と通信のオーバーラップ

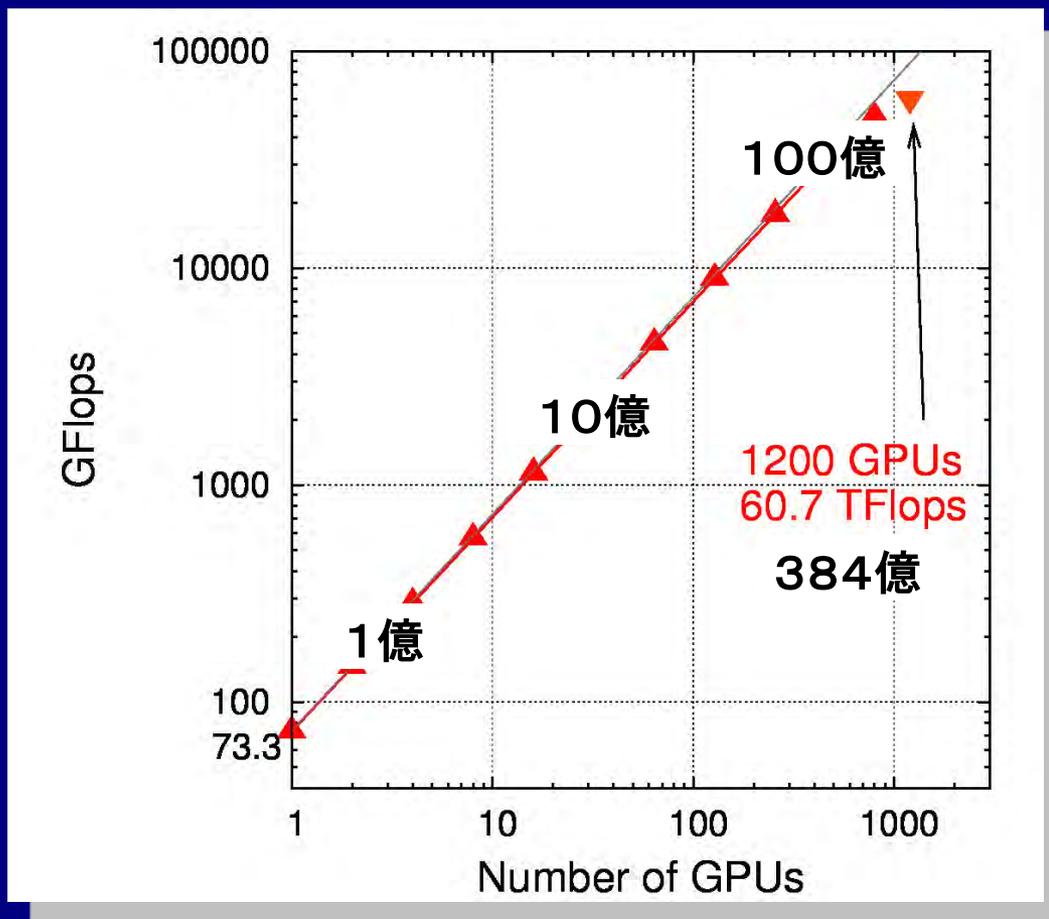


海陸地形を含む計算例



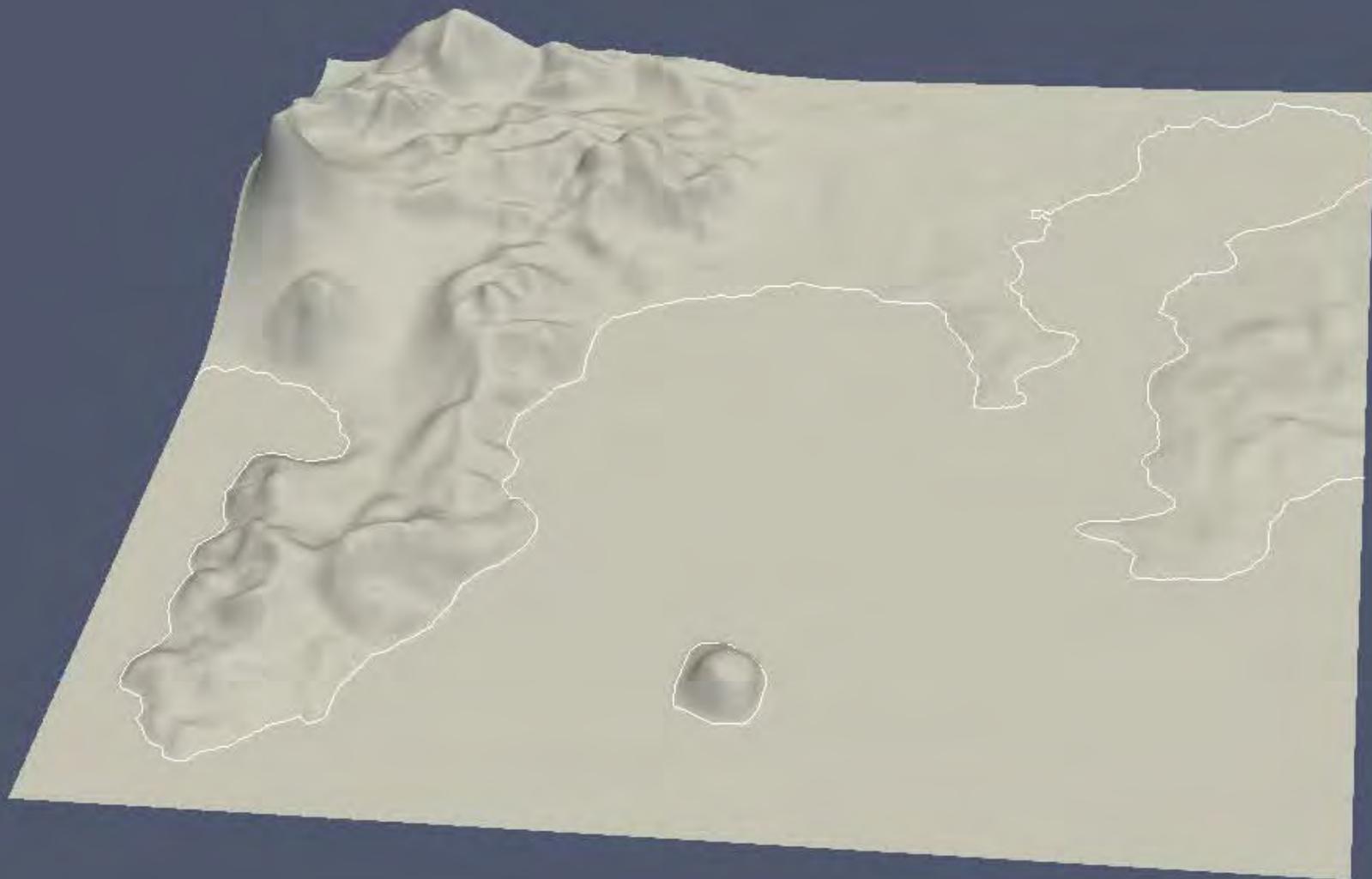
- 地形、海水、堆積物、地殻
(弾性体 3層モデル)
- 2004年 Mw 5.6 深さ 12 km
- 3200 x 3200 x 1280
- 格子間隔 50 m
- **131 億単位セル**
- **400 GPU**
- 20000 ステップ (40 s)
- 計算時間 **2068 秒**
- 非弾性は含まない

スケーラビリティ



領域サイズ: $320 \times 320 \times 320 \times \text{GPU数}$
(弱スケーリング)

- GPUプログラム
 - 3次元領域分割
 - 弾性体版
 - Up to 1200 GPUs
 - 弱スケーリング
 - ほぼGPU数に比例
- 800 GPU まで
 - 2 GPU / ノード
- 1200 GPU
 - 2 GPU / ノード

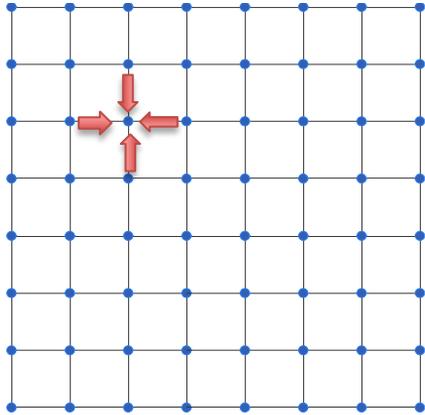


Time: 0.0 s

Physis (Φύσις) DSL Framework

[Maruyama, Matsuoka, etc. SC11]

Physis (φύσις) is a Greek theological, philosophical, and scientific term usually translated into English as "nature." (Wikipedia:Physis)



Stencil DSL

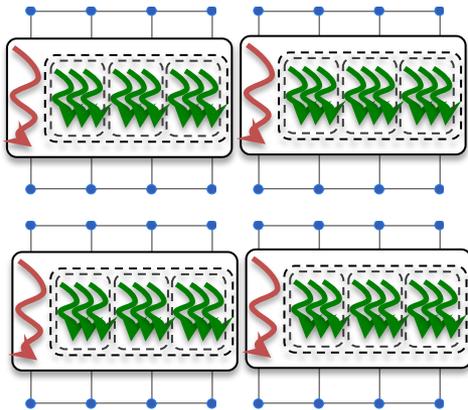
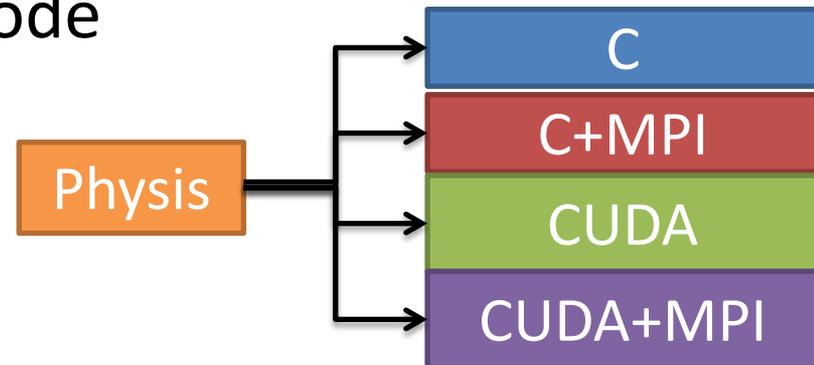
- Declarative
- Portable
- Global-view
- C-based

```
void diffusion(int x, int y, int z,
              PSGrid3DFloat g1, PSGrid3DFloat g2) {
    float v = PSGridGet(g1,x,y,z)
    +PSGridGet(g1,x-1,y,z)+PSGridGet(g1,x+1,y,z)
    +PSGridGet(g1,x,y-1,z)+PSGridGet(g1,x,y+1,z)
    +PSGridGet(g1,x,y,z-1)+PSGridGet(g1,x,y,z+1);
    PSGridEmit(g2,v/7.0);
}
```



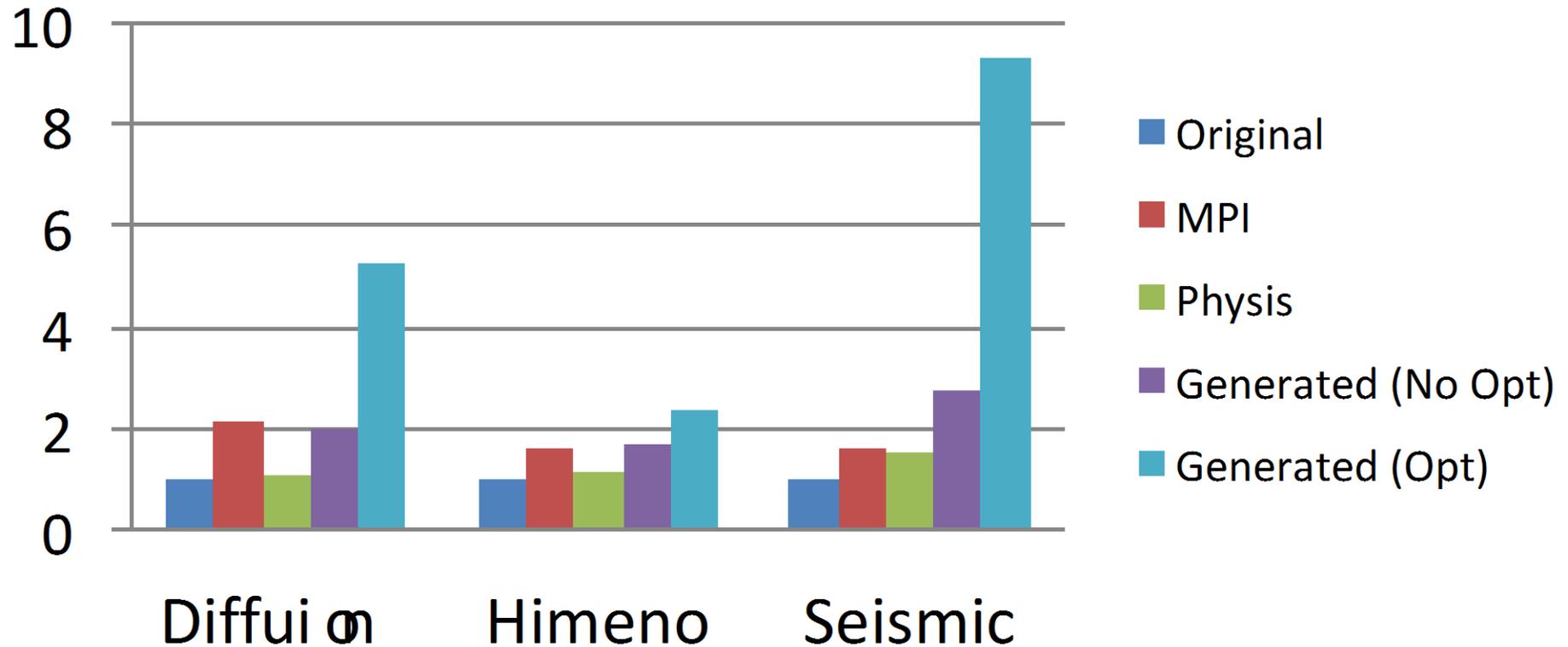
DSL Compiler

- Target-specific code generation and optimizations
- Automatic parallelization



Physis Productivity

Increase of Lines of Code



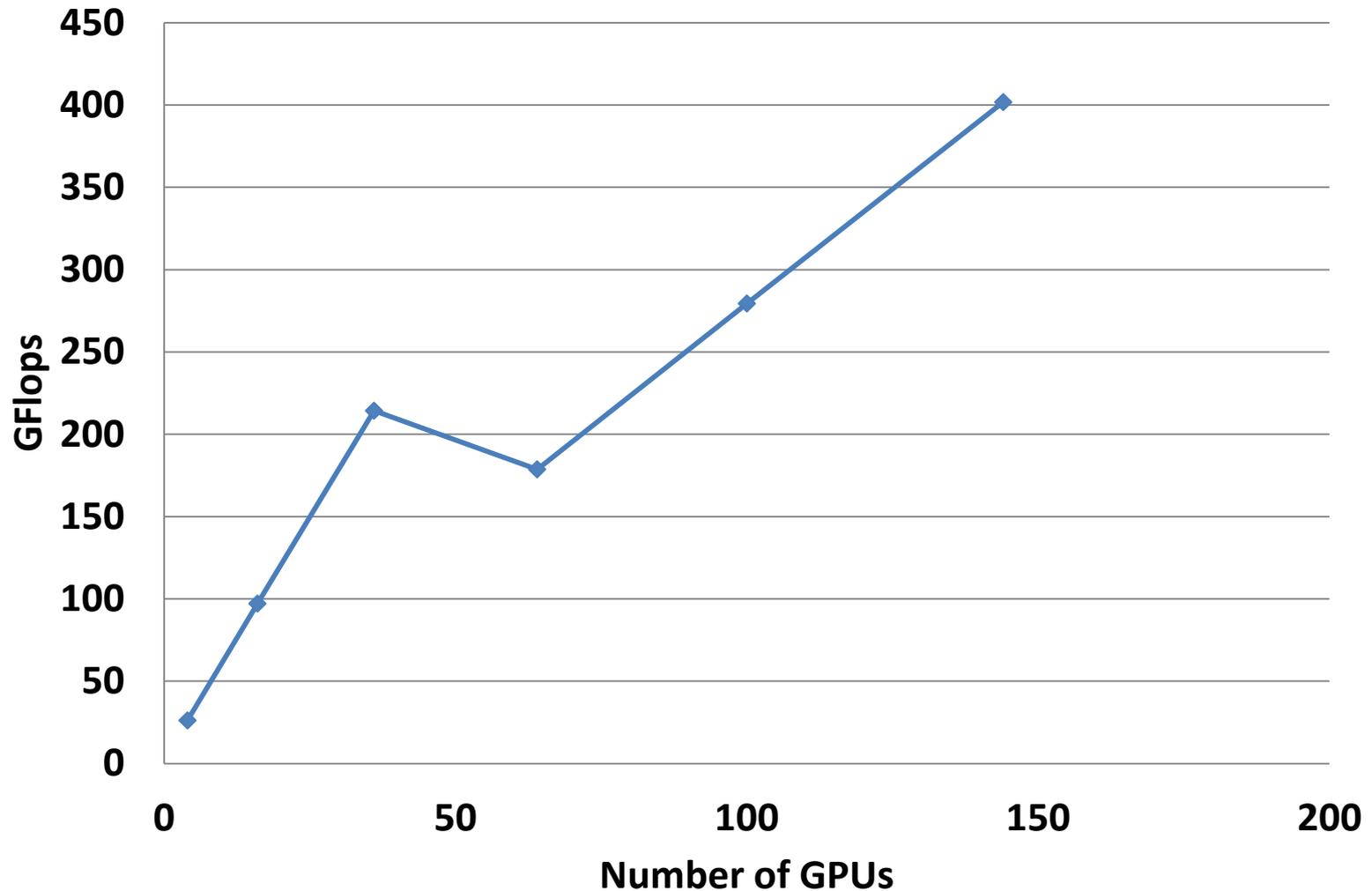
Similar size as sequential code in C

Communication incl. overlap auto-generated

Physis Seismic Weak Scaling

(Original code courtesy Prof. Imamura @ U-Tokyo)

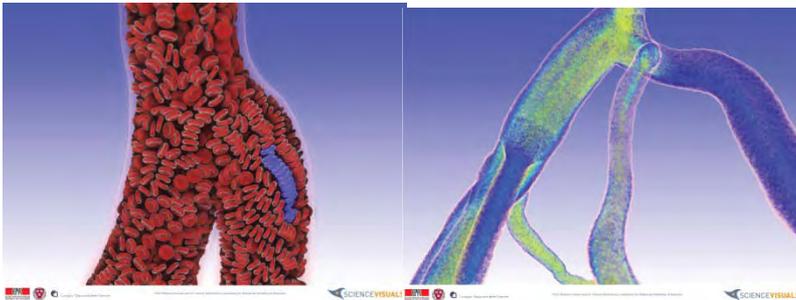
Problem size: 256x256x256 per GPU



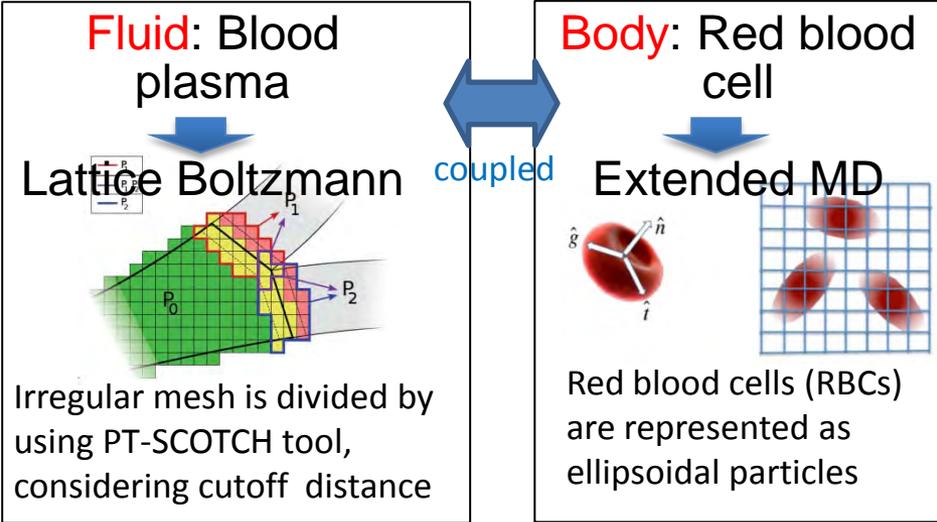
Multiphysics Biofluidics Simulation

[Bernaschi et. al., IAC-CNR Italy]

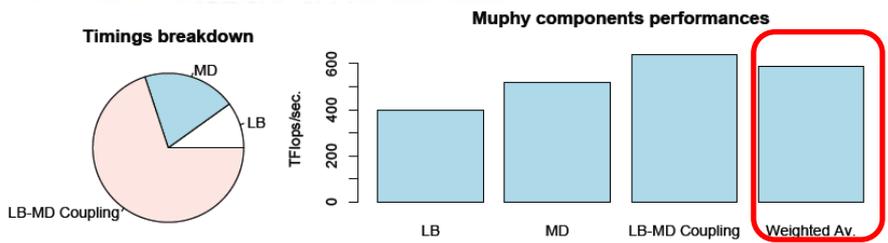
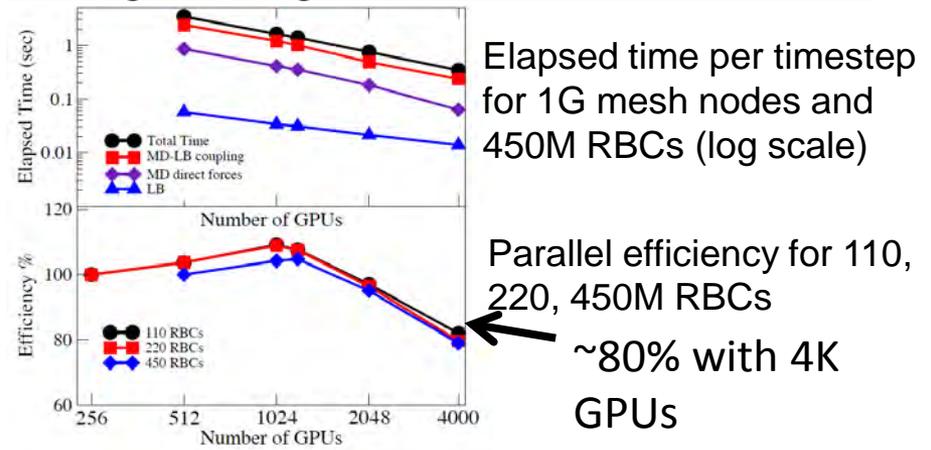
To understand real-life biofluidics problems, simulations of blood flows that accounts for from red blood cells to endothelial stress are conducted



Multiphysics simulation with *MUPHY* software



Strong Scaling Results on TSUBAME2.0



0.6PFlops with 4,000GPUs

for 1G mesh nodes, 450M RBCs
 A complete heartbeat at microsecond resolution can be simulated in 48hours

Large-Scale Metagenomics

[Akiyama et. al. Tokyo Tech.]

Combined effective use of GPUs and SSDs on TSUBAME2.0.

Metagenome analysis: study of the genomes of uncultured microbes obtained from microbial communities in their natural habitats



Collecting bacteria in soil

Two homology search tools are available:

- 1) **BLASTX**, standard software on CPUs
- 2) **GHOSTM**, our GPU-based fast software compatible with BLASTX

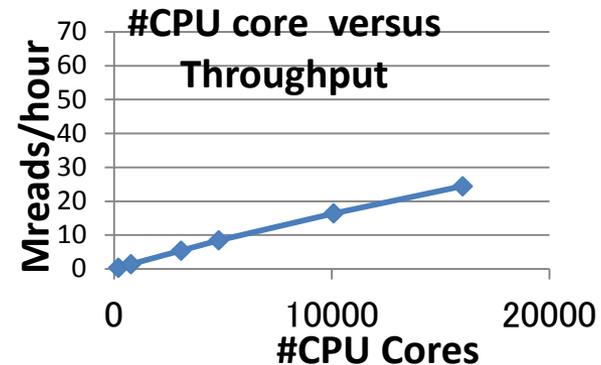
Data: 224million DNA reads(75b) /set

Pre-filtering: reduces to 71M reads

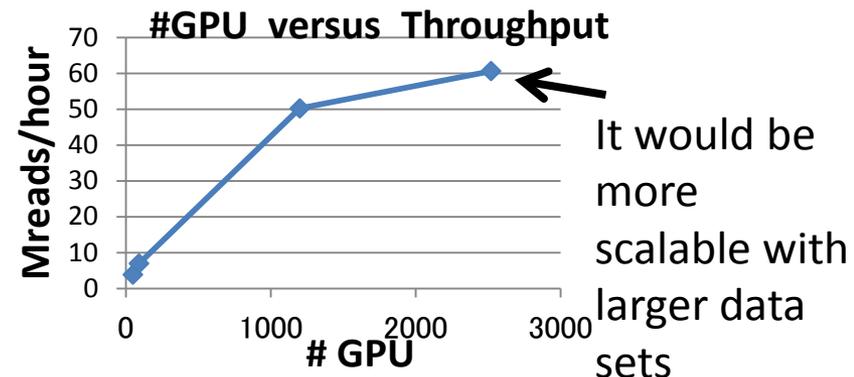
Search: 71M DNA vs. NCBI nr-aa DB (4.2GB)

Results on TSUBAME2.0

BLASTX: 24.4M/hour with 16K cores



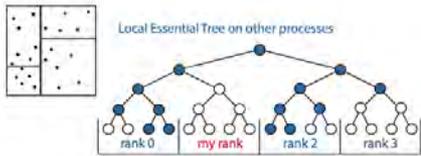
GHOSTM: 60.6M/hour with 2520 GPUs



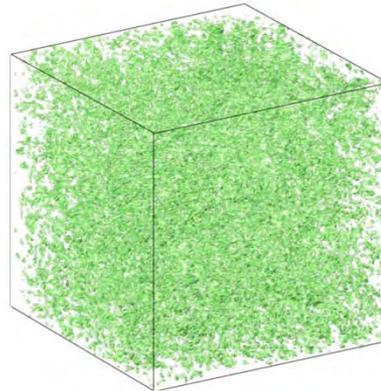
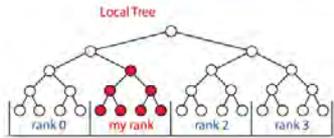
Turbulence Simulation using FMM

[Yasuoka et. al. Keio Univ.]

We simulate turbulent flows using the fast multipole method (FMM), which is more scalable than the traditional FFT-based approach



Domain decomposition and tree construction in FMM



Q criteria in isotropic turbulence

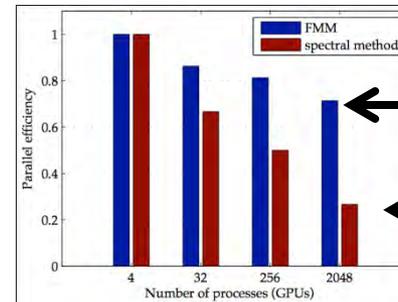
Experimental conditions:

Mesh: 2048^3

Re_λ : 500

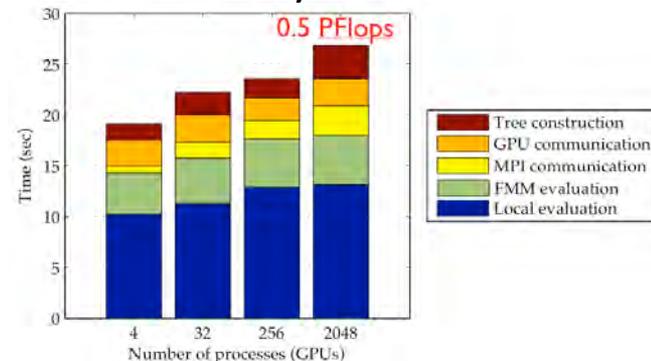
Used library: **ExaFMM** (ours on GPUs) and **FFTW** (on CPUs)

Weak scaling results on TSUBAME2.0
Comparing parallel efficiency
(4×10^6 particles per proc)



ExaFMM: 72%
with 2048GPUs
FFTW: 27%
with 2048CPUs

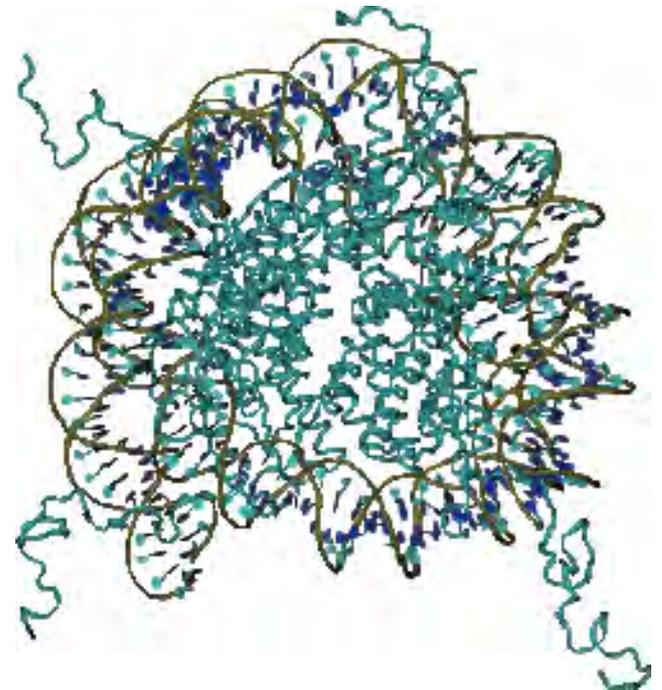
Weak scalability of ExaFMM



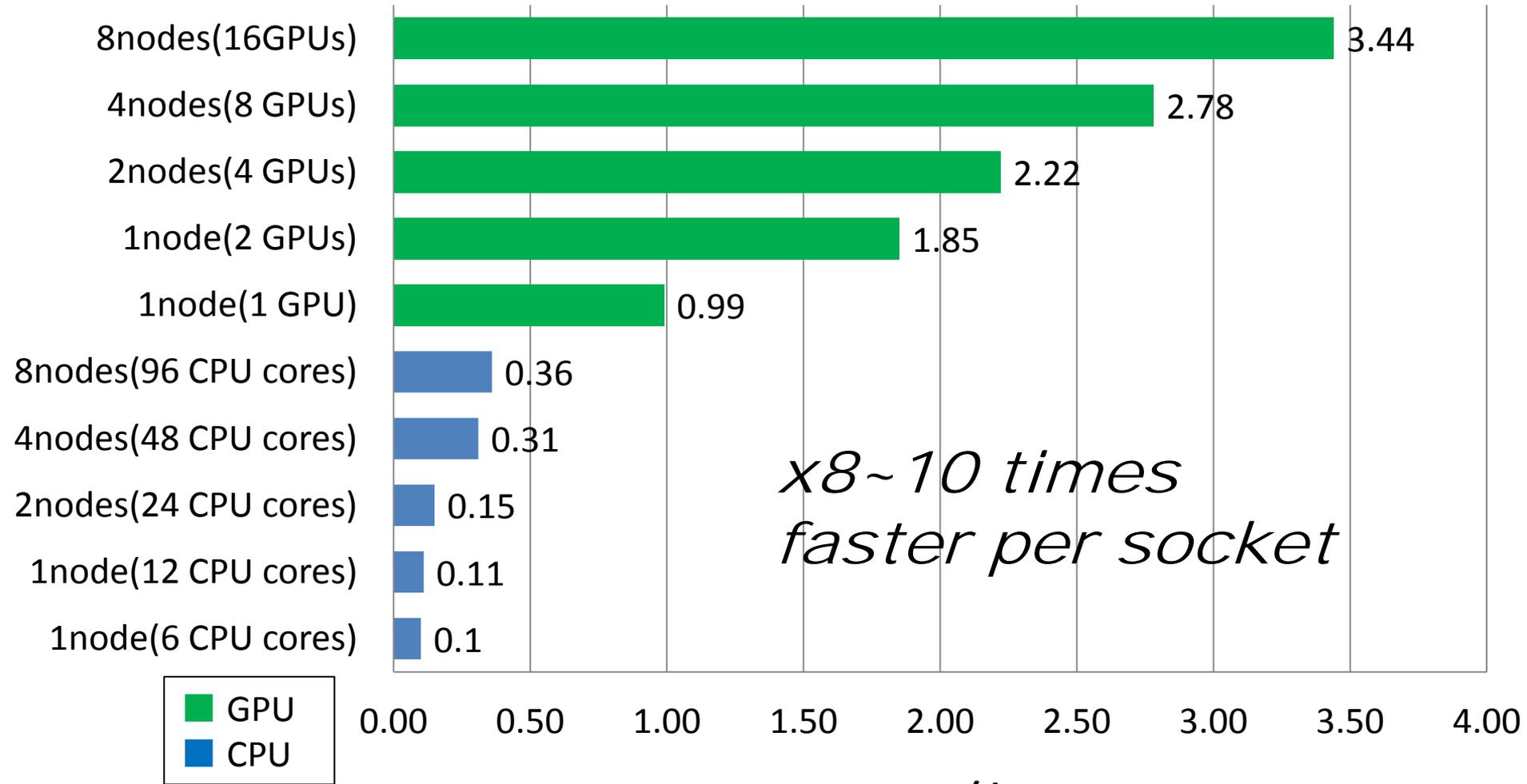
0.5PFlops with 2,048GPUs

AMBERとは？

- UCSDのP.A. Kallman教授のグループにより開発された、主に生体分子向けのモデリング、分子力学および動力学シミュレーションパッケージ
- 現在のバージョンの11より、本格的なGPU対応を開始
- ターゲット
 - Nucleosome (25,095原子)
 - NVEアンサンブル
 - Generalized Born model



Calculation speed

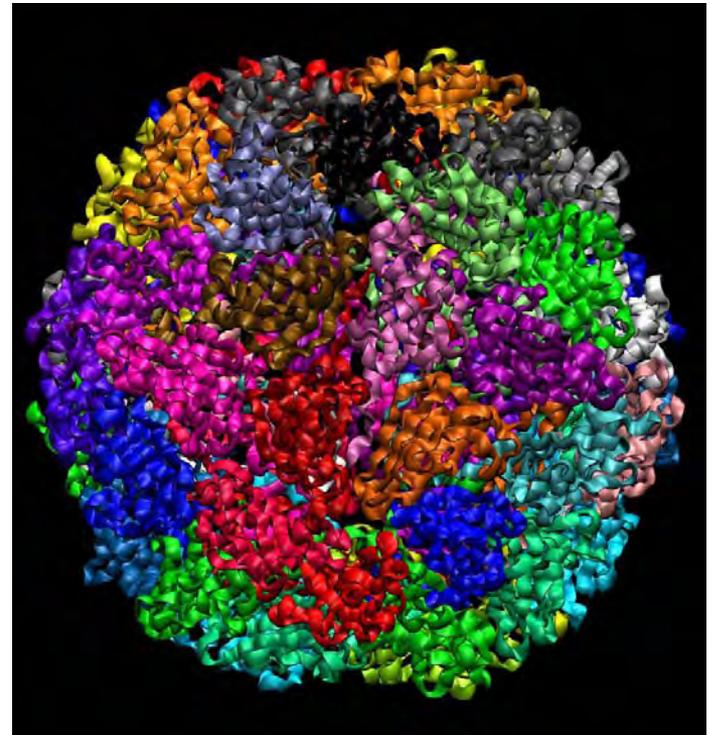


*x8~10 times
faster per socket*



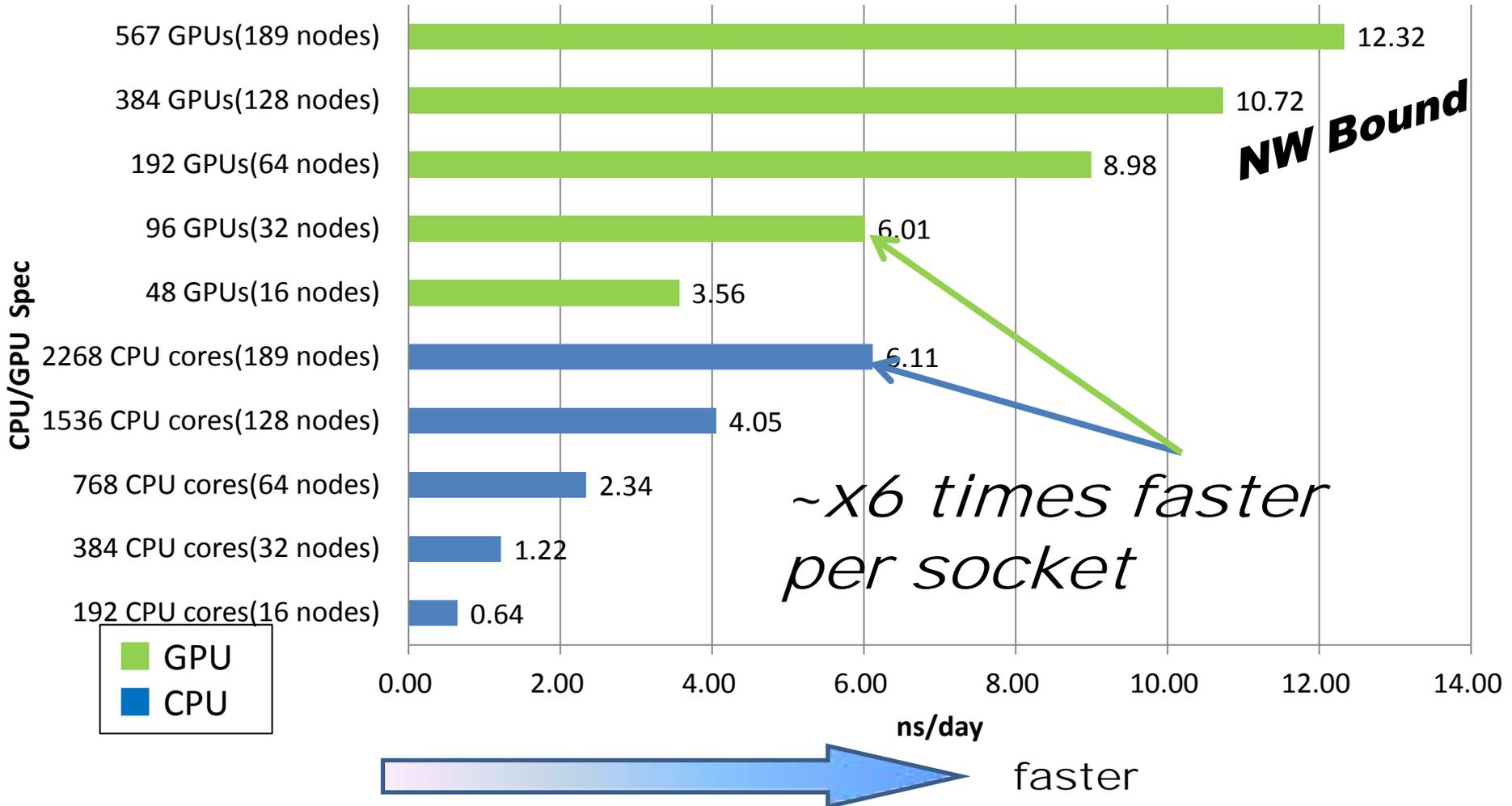
NAMDとは？

- イリノイ大学のKlaus Schulten教授のグループにより開発された、特に並列計算機における高速計算に強い分子動力学シミュレーションパッケージ
- Blue Gene/L, Jaguarを始め、多くの巨大なシステム上で動作
- GPGPU上でも高速実行
- ターゲット
 - STMV (1,066,628)
 - NPT アンサンブル
 - PME (Particle Mesh Ewald)



Calculation speed

stmv (1,066,628 atoms), NPT



TSUBAME の共同利用

平成19年度～ 民間企業へ TSUBAMEを提供

- 文部科学省 先端研究施設共用促進補助事業

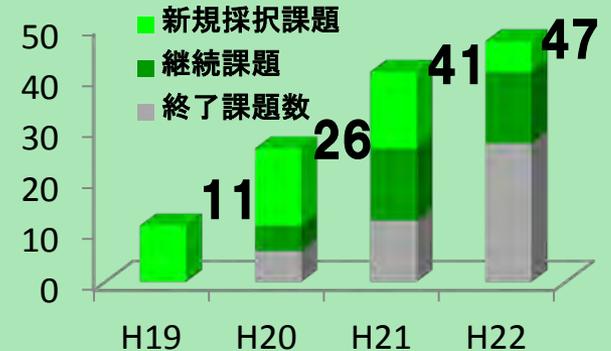
平成21年度～ TSUBAME 共同利用

- 東工大独自事業として、外部利用制度の確立



「産業利用」にて大きな成果

4年間で民間企業の47課題を採択・実施



成果報告会 10/19(水) 13:00～ 蔵前会館
東工大・TSUBAME共用促進シンポ

文部科学省より最高の評価

(22機関中、3機関のみ)

平成22年度～ 学際大規模情報基盤共同利用・共同研究拠点

- 「ネットワーク型」の共同利用・共同研究拠点

日本最先端のスパコン環境を提供し、学術・産業・社会へ貢献
最新GPGPUクラスター環境を提供し、GPGPUの普及へ貢献

産業界でのGPU活用例 on TSUBAME2.0

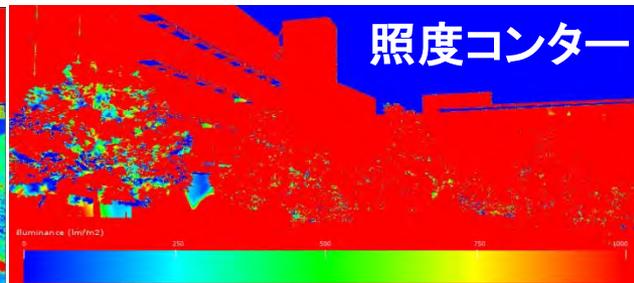
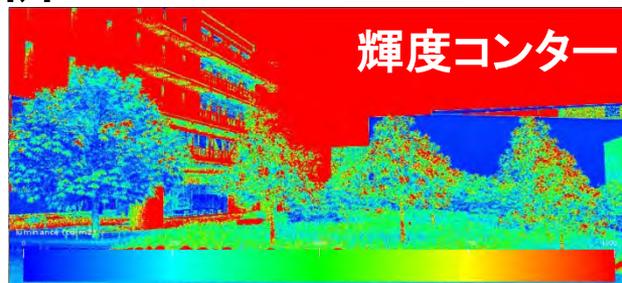
建築物の室内外環境の連成解析とその高速化技術の開発

清水建設株式会社 技術研究所

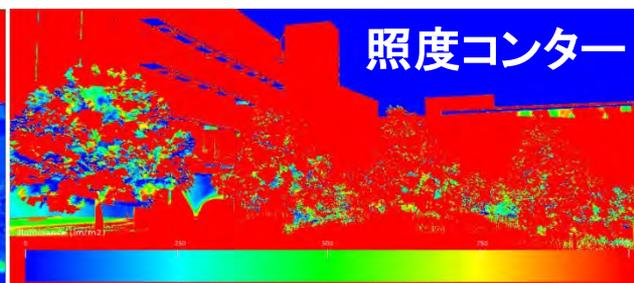
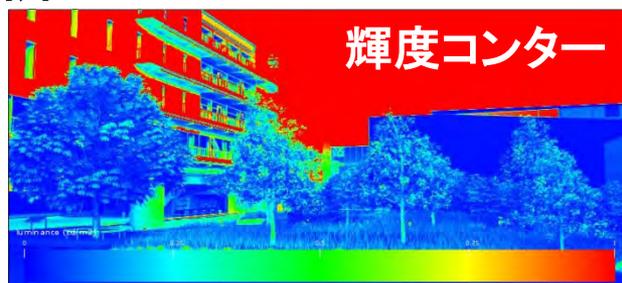
室内外環境の連成解析により建物内外の光環境を素早く可視化
～GPUを活用した超高速可視化プレゼンテーションシステムの開発～

数十～数百のGPUを利用すれば複雑な大規模モデルでも数秒で写真画質の描画が可能

日中の光環境の解析



夕方の光環境の解析



産業界でのGPU活用例 on TSUBAME2.0

| 採択年度 | 課題名 / 実施企業名 | 概要 |
|--------------------------|------------------------------------------------------|------------------------------------------------------------------------|
| H23 現在も 利用中 | 日本ゼオン株式会社 / メソ構造を持つ高分子材料のマルチスケール・シミュレーション | ソフトマテリアルの統合的なシミュレータであるOCTAシステムの、高分子の高精度・大規模・高速な動的平均場法の並列化コードの開発を行う。 |
| H23 現在も 利用中 | 住友電気工業株式会社 / 個別要素法を用いた粉末充填シミュレーションプログラムの並列化とその評価 | 大規模並列化およびGPUコンピューティング技術を駆使した、実プロセス規模の解析の可能性の検討を行う。 |
| H22 ～23 現在も 利用中 | 清水建設株式会社 / 建築物の室内外環境の連成解析とその高速化技術の開発 | 建築物の室内外環境の解析モデル及び連成解析システムとGPUによる数値解法を開発し、建築物の室内外環境の評価を可能にする。 |
| H22 | 株式会社フィアラックス / 分子動力学計算ソフトウェアNAMDのGPGPU大規模並列環境における性能評価 | 分子動力学計算ソフトNAMDを用いて、複数ノードでのGPU並列実行での性能評価と、CPUとの性能比較を行なった。 |
| H21 | ソニー株式会社 新概念による大規模並列電磁界解析技術研究 | 複雑な電子機器の数値モデル化と、環境電磁ノイズ発生メカニズムのGPUによる可視化により、電子機器のノイズ低減の最適設計を行った。 |
| H21 | 株式会社クロスアビリティ CUDAを用いたGPUによるフラグメント分子軌道法の高速度化 | DFTおよびMP2フラグメント分子軌道法をGPUで高速化するプログラム開発を目的とし、電子状態計算のGPGPU高速化モジュール開発を行った。 |

ExaScale Computing Study: Technology Challenges in Achieving Exascale Systems



Peter Kogge, Editor & Study Lead
 Keren Bergman
 Shekhar Borkar
 Dan Campbell
 William Carlson
 William Dally
 Monty Denneau
 Paul Franzon
 William Harrod
 Kerry Hill
 Jon Hiller
 Sherman Karp
 Stephen Keckler
 Dean Klein
 Robert Lucas
 Mark Richards
 Al Scarpelli
 Steven Scott
 Allan Snaveley
 Thomas Sterling
 R. Stanley Williams
 Katherine Yelick

Petaを達成したが中国に抜かれた米国は2018-2020 Exa(10¹⁸)flopへ驀進を開始

Peter Koggeらによる300ページのDoD Exascaleシステムのレポート

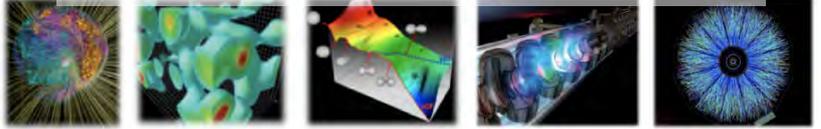
September 28, 2008

This work was sponsored by DARPA IPTO in the ExaScale Computing Study with Dr. William Harrod

Exa-scale Computational Resources

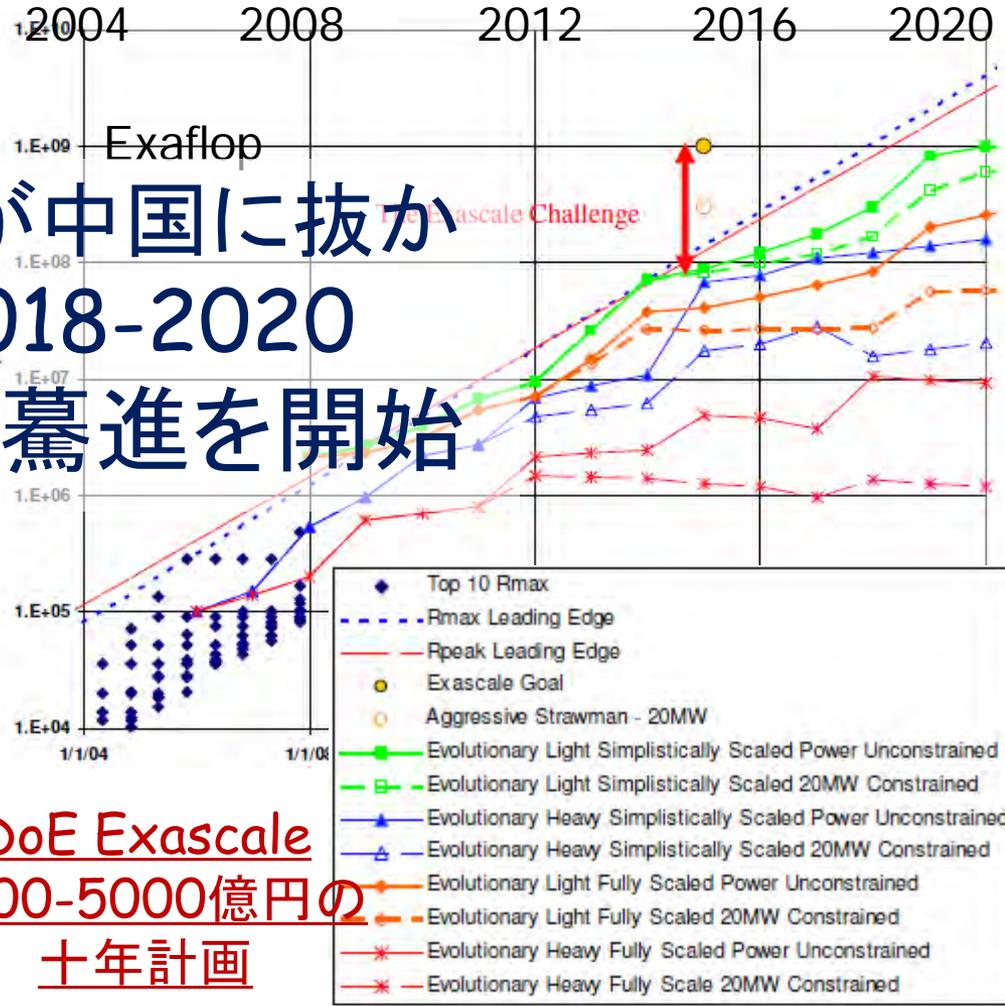
(slide courtesy Martin Savage)

Meeting structured around three main areas of effort
6アプリ分野のExascale Workshop(2008-2009)



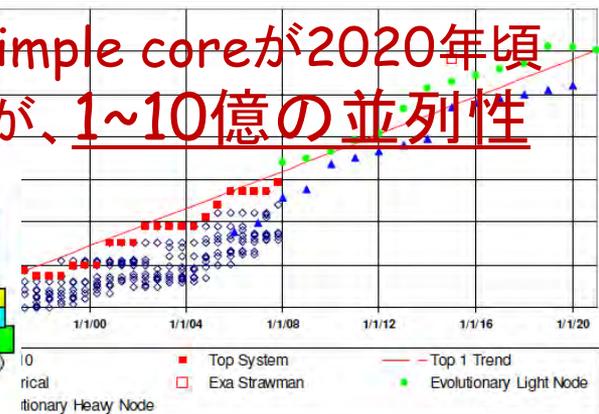
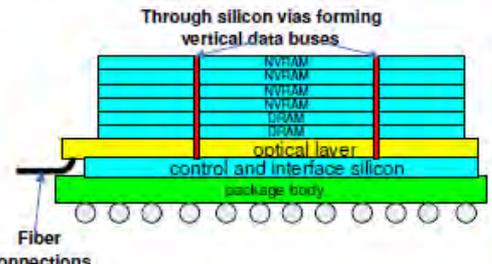
- Exa-scale computing is **REQUIRED** to accomplish the Nuclear Physics mission in each area
- Staging to Exa-flops is crucial :
 - 1 Pflop-yr to 10 Pflop-yrs to 100 Pflop-yrs to 1 Exa-flop-yr (sustained)

Paul Messina June 28, 2009



DoE Exascale 2000-5000億円の十年計画

軽量なsimple coreが2020年頃有望だが、1~10億の並列性



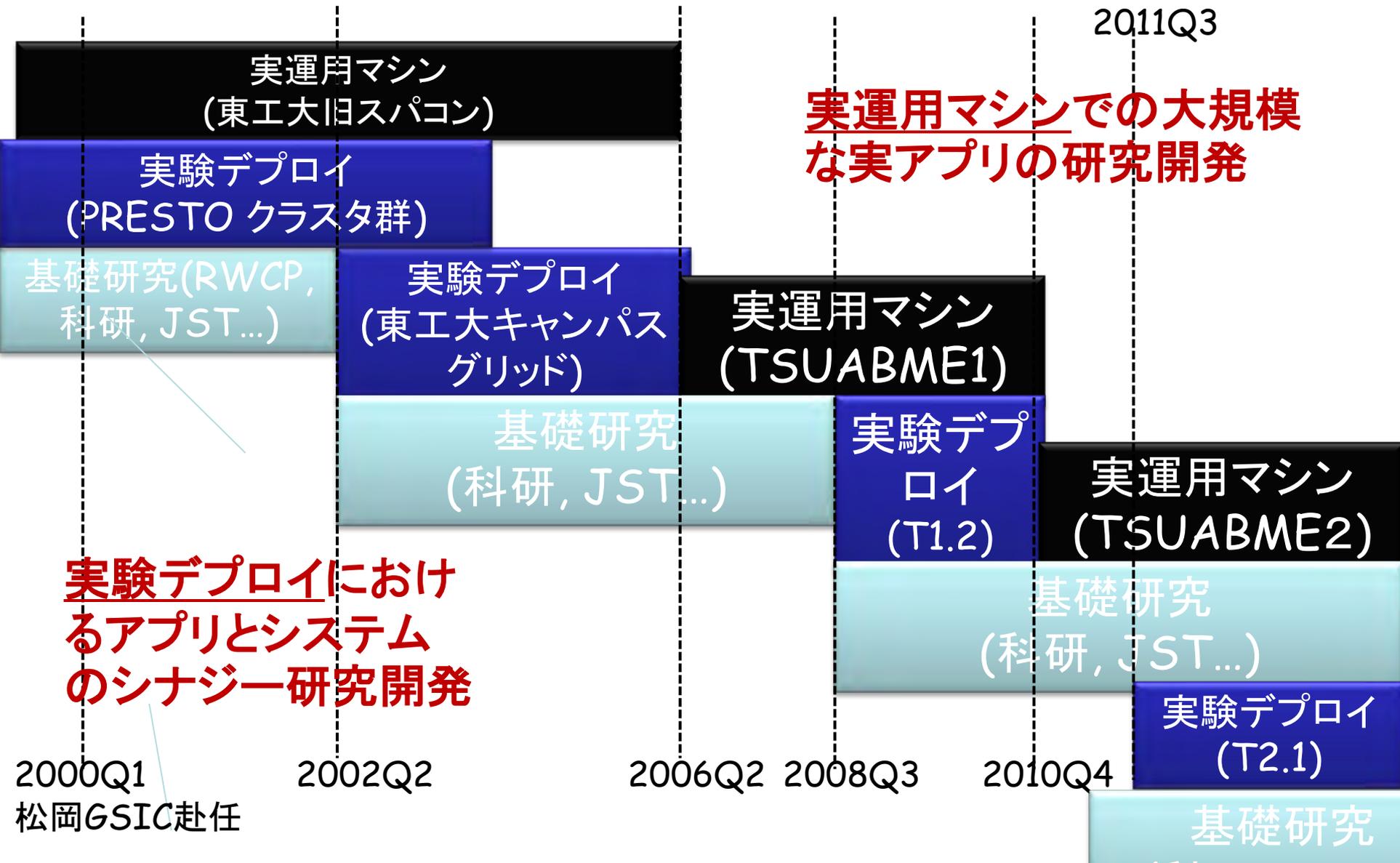
「エクサ10億並列へ」は勇ましいが。。。。

- 電力・エネルギー (松岡 JST CREST ULPHPC + MEXT Green HPC)
- (強)スケーリングの欠落 (Co-Design Center)
- N^2 vs. N 問題により深まるメモリ階層
 - (ネットワークやI/O含む) (藤澤 JST Post Petascale)
 - (レイテンシとバンド幅) (Device & Architecture)
- 極端に低まる信頼性と実行不可能性 (松岡 科研基盤S)
- プログラミングや実行モデル (丸山 JST CREST Post Petascale)

Power Efficiency Compared

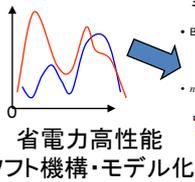
| Machine | Power (incl. cooling) | Linpack Perf (PF) | Linpack Mflops/W | Factor | |
|----------------------|-----------------------|-------------------|------------------|--------|------------------|
| Tsubame1.0 (2006Q1) | 1.8MW | 0.038 | 21 | 2368 | |
| ORNL Jaguar (2009Q4) | ~9MW | 1.76 | 196 | 256 | |
| Tsubame2.0 (2010Q4) | 1.8MW | 1.2 | 667 | 75 | $\times 31.6$ |
| K Computer (2011Q2) | ~16MW | 10 | 625 | 80 | |
| BlueGene/Q (2012Q1) | ~12MW? | 17 | 1417 | 35.3 | |
| Tsubame3.0 (2015Q1) | 1.8MW | 20 | 11000 | 4.6 | $\sim \times 16$ |
| EXA (2018Q4)? | 20MW | 1000 | 50000 | 1 | |

TSUBAME2.0に至る15年間:連続的な multi waterfall model研究開発

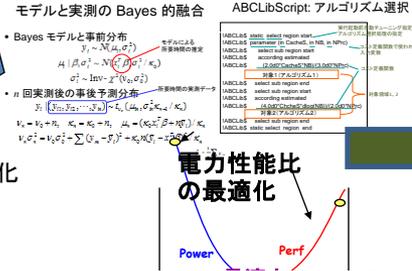


JST-CREST ULP-HPCの成果適用による スパコンのウルトラグリーン化(H23~概算要求)

JST-CREST Ultra Low Power HPC (2006-2013)
スパコンの1000倍の性能電力向上を目指す基礎研究



(提案3) 自動チューニング共通基盤



(3) JST ULPHPC 基礎研究
の適用によるスパコン用
超省電力ドルウェア・ア
プリ等

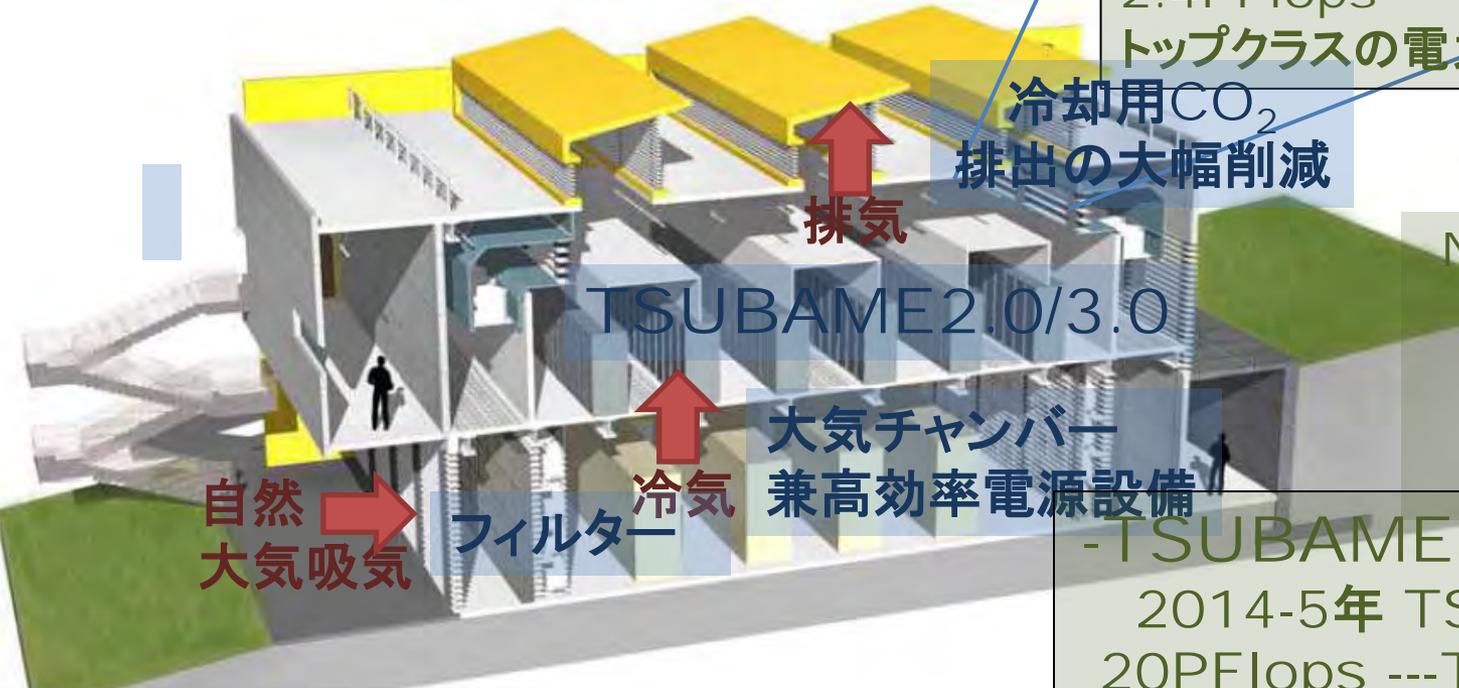


(提案1) ULP-HPC用次世代SW/HW要素利用
技術およびモデル化

(2) TSUBAME2.0 (2010)
2.4PFlops
トップクラスの電力性能



冷却用CO₂
排出の大幅削減



Nvidia Maxwell
(2013)追加
TSUBAME2.5
~10PF
1MWの維持

-TSUBAME3.0への成果-
2014-5年 TSUBAME3.0
20PFlops ---TSUBAME1.0
比で

(1) 自然大気冷却等 年平均PUE ~