JST ULP-HPC 次世代テクノロジのモデル化・ 最適化による超低消費電力 ハイパフォーマンスコンピューティング

東京工業大学・学術国際情報センター 研究基盤部門 教授

代表者 松岡 聡

2012/11/30 ULP最終最終シンポジウム発表資料

過去のスパコンのイメージ... 10年で1000倍の向上→電力を犠牲に









今や電力がスパコンの速度向上の主律束に

2016年ULP-HPC技術により デスクサイド・ペタスケールコンピューティングへ

(Mid-size to Deskside Scaling)



> 10,000 CPU, 1.2 MegaWatt, 350m²





現代版Mooreの法則では「2年で2倍」



Intel, 2009.

10年で電力性能比1000倍の目標



(グラフは冷却電力向上分を除く)

省電力化の手法と有効性

情報基盤/手法	エンタープライズ・ビジネス・ク ラウド	HPC
仮想化による統合化 (Server Consolidation)	0	×
DVFS (Dynamic Voltage/Frequency Scaling)	0	Δ
新デバイス (と、それを生かした アーキテクチャ)	▲(コストや継続性)	0
新アーキテクチャ (と、それを生かす ソフトウェア)	▲(コストや継続性)	Ο
冷却技術 (排熱利用や デバイス省電力化含)	Δ	O (ただし高熱密度)

我々のCREST ULP-HPC全体スキーム



Total System FLOPS/W, MemBW/W considered significant

- "Power" (and cost) is THE constraining factor in large scale supercomputers
 - Determines the size of the system
- Large scale apps performance limited by <u>TOTAL system</u> <u>FLOPS</u>, or <u>TOTAL Memory Bandwidth</u>
 - (Traditional Byte/Flop argument not very useful)
- FLOPS/W = Total Sys FLOPS / Total Sys W
 ~ "Green 500"
- MemBW/W = Total Sys MemBW / Total Sys W
 - ~"Green Graph 500"
- Exascale 50GFlops/W, 5GByte/s / W (B/F=0.1)
 - x20 improvement in 8 years?

ULP-HPC 遂行計画・マイルストーン





ULP-HPC 査読論文•講演等

· 査読論文: 122件

- 国内(和文)誌33件、国際(欧文)誌89件
- 基調講演・招待講演:合計152件
 - 国内会議 67件、国際会議 85件
 - Satoshi Matsuoka, the First Petaflop/s System in the World and Its Impact on Supercomputing, Opening keynote talk, ISC 08 (International Supercomputing Conference) @ Dresden, Germany, June 2008.



- Satoshi Matsuoka. Petascaling Commodity onto Exascale: GPUs as Multithreaded Massively-Parallel Vector Processors - the Only Road to Exascale. <u>Keynote Talk</u>, IEEE Cluster Computing Conference 2009, New Orleans, USA, Sep.3, 2009.
- Takayuki Aoki: Large-scale Stencil Applications on GPU-rich Supercomputer TSUBAME2.0, The annual IEEE International Conference on High Performance Computing (HiPC 2011), Bangalore, India, December 21, 2011

・ 口頭発表・ポスター: 262件

- 国内会議 185件、国際会議 77件
- ・派生プロジェクト: 丸山ポストペタCREST,遠藤ポストペタ CREST、科研基盤(S),日仏FP3C, G8気候シミュレーション、 グリーンスパコン概算要求、TSUBAME2.5...







- 松岡 聡,2010年11月発表のGreen500グリーンスパコンランキングとおいて世界2位、およびthe Greenest Production Supercomputer in the World賞を獲得
- 松岡 聡, 2010年11月発表のTop500スパコンランキングにおいて世界4位
- 額田 彰, IEEE Japan Chapter Young Author Award 2010, 2010年12月
- ・ 下川辺・青木・額田・遠藤・丸山・松岡ら: 2011年11月SC11において, ACM Gordon Bell Prizes: Special Achievement Award in Scalability and Time-to-Solution受賞
- イタリアCNR Massimoらと共同で、遠藤・松岡ら: SC11において、ACM Gordon Bell Prizes: Honorable Mention
- 松岡ら: HPCwire Annual Award:
 - Reader's Choice Award Best application of "green computing" in HPC
 - Editor's Choice Award Best application of "green computing" in HPC
 - Reader's Choice Award Best HPC collaboration between government and industry
- 鯉渕: Best Paper Award, the Second International Conference on Networking and Computing(ICNC), 2011年12月
- 遠藤:情報処理学会山下記念研究賞, 2012年3月
- 松岡ら: HPC Wire "Number Crunching, Data Crunching and Energy Efficiency: the HPC Hat Trick" (TSUBAME2.0のGraph500, Top500, Green500上位に対して)
- 松岡・遠藤・青木: 文部科学大臣表彰科学技術賞(開発部門)「運用世界ーグリーンペタス パコンの開発」

報道等

- TSUBAME2.0を用いた成果・受賞・グリーンスパコン等に関する 多数報道
 - 読売新聞(11/20),朝日新聞(11/10),日経産業新聞(2/24)
 - 日経コンピュータ, WIRED
 - クラウドWatch (Gordon Bell賞受賞に関して)
 - その他asahi.com, HPCwire, IT media等





TSUBAME2.0 2010年11月1日稼働開始 世界最小のペタフロップス・省電カスパコン



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



TSUBAME2.0 Compute Node

Thin Node 1.6 Tflops Peak 400GB/s Mem BW 80GBps NW BW ~1KW max

Infiniband QDR x2 (80Gbps)

HP SL390G7 (Developed for TSUBAME 2.0)

GPU: NVIDIA Fermi M2050 x 3 515GFlops, 3GByte memory /GPU CPU: Intel Westmere-EP 2.93GHz x2 (12cores/node) Multi I/O chips, 72 PCI-e (16 x 4 + 4 x 2) lanes --- 3GPUs + 2 IB QDR Memory: 54, 96 GB DDR3-1333

TSUBAME2.5 4PF Peak, 800GB/s Mem BW

SSD:60GBx2, 120GBx2



SSD: ~200TB

TSUBAME2.0へのULP-HPCの成果の反映

- GPU中心の2.4 PF我が国最速・世界4位のスパコン
 - 1432 nodes, Intel Westmere/Nehalem EX
 - 4224 NVIDIA Tesla (Fermi) M2050 GPUs
 - ~76,600 total CPU and GPU "cores", High Bandwidth
 - 活用するソフト・アプリ・最適化の研究⇒電力効率世界ー?(後述)
- 省電力運用のための仕組みが随所
 - I/O・ネットワークの新デバイス採用(マルチDFB発光素子一体型ケー ブル・SSDデバイス等)による省電力化
 - ノード・ラック・配電版など随所の電力センサーネットワーク
 - 大量の温度センサー(ノード18個=>全体で2万個以、ファンセンサー...)
 - ノード単位の電力キャップ・高効率冷却ノード設計
 - 密閉型水冷ラック/チラー (35KW => PUE=1.28以下)
- <u>今後さらなるULP-CRESTの成果の適用</u>
 - 開発された高性能・省電力GPUライブラリ等の利用
 - 省電力自動チューニング
 - ¹⁶ 省電カスケジューリング(温度感知=>マイグレーション)

TSUBAME2.0の4000超GPUを用いた Linpack実行 [IPDPS10, SACSIS11]

数百~数千アクセラレータを効率利用 するLinpackソフトウェアを継続開発

- TSUBAME1.0→1.1→1.2→2.0
 で連続的に性能向上を実現
- アクセラレータやノードのヘテロ性・ PCI通信の影響・アクセラレータメモ リの小ささに対応する最適化

TSUBAME2.0全体を用いて 1.192PFlops (実行時間2.4h) ⇒ Top500世界4位 958MFlops/Watt ⇒ Green500世界2位 運用スパコン1位





TSUBAME2.0スパコンの電力危機対応



超低消費電力・数値計算アルゴリズム 開発とHPCアプリケーション 東エ大 青木グループ

HPCアプリの関数毎の詳細電力測定

- GPUアプリのメモリ・アクセスの消費電力測定
- GPUアプリの浮動小数点演算の消費電力測定



低速メモリへのアクセス低減が有効





Elaosed time [s]

マルチグリッド Vサイクル

- GPUの Shared Memory を Software Managed Cache として利用
- 通信と計算のオーバーラップ

1m格子解像度による10km×10kmエリアの 都市気流のLESシミュレーション

格子ボルツマン法: D3Q19 モデル

 ■ LES (Large-Eddy Simulation) モデル: Coherent-Structure Smagorinsky model
 ■ メモリアクセス律速な計算にも関わらず 4000 GPU で 600 TFlops,高電力比性能 545 Mflops/W

TBSテレビ,朝日新聞,読売新聞,日経新聞, PC Watch等で 2012年10月,11月に報道







TSUBAME2.0による高電力比効率 合金材料の樹枝状凝固シミュレーション







~1.36 MW 高電力比性能: 1468 MFlops/W (c.f. Green500: 957 MFlops/W)

High Performance FFT on GPU

<u>自動チューニング機能搭載 NukadaFFT</u> <u>library [</u>SC08,SC09] チューニング項目: (1) 基底の分解、(2) スレッド数、 (3) バンクコンフリクト回避のためパディング自動挿入



Performance of 1-D FFT. (Double Precision, batch=32,768, GPU=GeForce GTX 480.) http://matsu-www.is.titech.ac.jp/~nukada/nufft/

NVIDIA純正CUFFTライブラリより高性能 AMD RADEON GPUでもOpenCL版で同程度の性能



GPU間のAll-to-all通信⇒スケールしにくい

- (1) Ibverbs APIを用いてオーバヘッド削減
- (2) 複数の相手ノードと同時通信で競合発 生時のペナルティ軽減
- (3) 複数のリンクに適切に割り当てることに よって競合発生率を低減



TSUBAME2.0 電力消費 (Jaguar比 4-5倍)

	Compute nodes & Interconnect (kW)		Storage (kW)		Storage Cooli (kW) (kW)		ng	Total	(kW)
アイドル時	530		70			70 200		8	00
平均的運用時	680				0	23	30	98	80
Graph500 (CPU)	902				5	34	16	13	23
Earthquake sim. (700nodes)	550/903 利田方法によって消費電				15 ታ	13	20		
ASUCA Weather	9	9 に大きな開き						13	808
NAMD (700nodes)	706,	制	ー 御(の	必要	生		15	527
Turblence Sim.	1190				2	24	10	15	502
Phase-field	1362				3	29	94	17	29
Linpack	1417		~		-	-	-		-
GPU DGEMM	1538		72		410		20)20	

GPU の消費電力測定

GPU 電力の測定装置

Auxiliary Power Lines

A,

PSU

A_{m0}

V.

3

0.8

0.7

0.6

0.5

0.4

0.3

0.2

0.1

0

0.976

0.977

0.978

0.979

0.98

CUDA カーネルの精密な電力効率測定

 $w \approx 72.0 + 1.02 \times 10^{-10} \rho f$





統計的GPU電カモデリング [IEEE IGCC10]



- ・リッジリグレッションによる過学習の防止
- クロスフィッティングによる最適パラメー
 タの決定



- GPUの消費電力を統計的に推定
- ・ 性能プロファイル(パフォーマンスカウンタ)を
 説明変数とした線形回帰モデル



DVFSを導入しても高精度を確認 今後:電力モデルによる電力最適化 線型モデルでも+分な精度

億単位のプロセッサからなるエクサスケール における最適化の実現可能性



OMPCUDA: GPU向けOpenMP処理系(電通大G) CPU向け並列化プログラミング環境でGPUを利用可能に



性能評価

同一ソースコードによる実行性能比較

アプリケーションソースコード (行列演算2重ループ)

#define N 1024 float a[N*N], b[N*N], c[N*N];



提案手法による実行は行列積において
 CPU(逐次実行)の176倍、CPU(4並列実行)の62倍の性能を達成

CUDAによる共有メモリを用いない実装では 22.2GFLOPS、

共有メモリを用いる実装では 83.4GFLOPS を達成可能

並列GPUコード自動生成による ステンシル計算のスケーラビリティ [SIAM PP 11]



ppOpen-ATの省電力拡張

● AT言語ppOpen-ATに省電力化AT機能を実装。任意箇所のエネルギーを最適化可能。



自動チューニングのための数理手法の研究と 自動チューニング数理ライブラリ ATMathCoreLib

自動チューニング ソフトウェアに組み込まれた可変性 を自動的に調整して、できるだけ 高性能・省電力を実現する仕組み



提案手法

- オンライン自動チューニング
- ワンステップ近似
- ワンステップ近似の高速化 ランダムサブセット法
 - ▶ 確率的候補選択法
- 並列処理の自動チューニング

ライブラリ

ATMathCoreLib

として公開中

▶ 並列実験 ▶ 並列試行

AT メーター

4DAC

◆ オンライン自動チューニング $w_i^{(k)} = \mu_i^{(k)} + (k-1)E\left(\min\left\{\mu_i^{(k-1)}, \mu_{\min'}^{(k-1)}\right\}\right)$ ◆ オフライン自動チューニング $w_t = \mu_t + \kappa E_v \left(\min \left\{ \mu_t^{post}(y), \mu_{\min} \right\} \right)$ ◆ 並列自動チューニング $W_n = \max\{t_{P+1}, \dots, t_{P+n}\} + (K-2)\min\{t_{P+1}, \dots, t_{P+n}, t_{ont}\}$ る性能比 Experiment only 0.9 Model only 0.8 Our method 1 0.7 Our method 2 0.6 衣い 0.5 0.4 理論限界(0.3 0.2 0.1 n Xeon2.0 Xeon3.0 Xeon3.8 Core2Duo2.3 Xeon2.4 Crusoe0.86 Core2Duo2.4 PentiumM1.8 Opteron2.6 Opteron2.2 tanium2-1.6 JSparc3-0.9 Power5-1.9 USparc3i-1.0 Power5-1.65 Core2Duo2.3 [>]entiumM2.0 **MobPen4-1.8** ²entium4-3.4

VobPen3-1.1

の外・中アンローリング: n = 1 から 256/512

ltanium2-1.3

国内外の自動チューニング研究をリード

iWAPT (international Workshop on Automatic Performance Tuning) を主催

	開催地		開催地
iWAPT 2006	東京	iWAPT 2010	Berkeley
iWAPT 2007	東京	iWAPT 2011	Singapore
iWAPT 2008	筑波	iWAPT 2012	神戸
iWAPT 2009	東京	iWAPT 2013	Barcelona

編集した書籍, 学会誌特集号

- 「**情報処理**」2009年6月
- Springer "Software Automatic Tuning" 2010 年 9 月
- 「**応用数理**」2010 年 9 月 12 月



様々な学会でセッション実施

- 情報処理学会 HPC 研究会
- 日本応用数理学会 年会
- ・ 日本計算工学会 講演会
- SIAM Parallel Processing (PP)
- SIAM Comp. Sci. and Eng. (CSE)

招待講演

- > 2009 ITBL シンポ(須田)
- ▶ 2009 数値解析シンポ(片桐)
- > 2009 Workshop on Libraries and Autotuning for Petascale Applications (片桐)
- ▶ 2011 HLPP 2011(片桐)

その他関連するもの

- 科学技術分野の文部科学大臣 表彰 若手科学者賞(片桐)
- SIAM News 2008/06 号で紹介
- 文科省科学技術政策研究所 「科学技術動向」2009/11紹介

ppOpen-ATとOMPCUDAの連携





On/Off リンクレギュレーション

マルチパスイーサネット、InfiniBand において On/Off&マルチスピードリンクレギュレーション法を実現

- トラフィック負荷によらず、スイッチの消費電力はほぼ一定
- ・ スイッチ消費電力の中でポートの消費電力が支配的
 - 並列アプリケーションの通信パターンの解析結果より、 利用率の低いリンクをオフ/減速することで電力減
 - 外部操作で電力削減が可能

スイッチの 型番	All except ports	GbE port	Peak(ポー トの割合)
SFS7k(IB)	21.4	0.9	43(50%)
PC62(GbE)	56.8	2.1	155 (63%)
SF-420(GbE)	32.6	1.0	55.4(41%)





NWの電力削減効果 (NAS並列ベンチマーク, Class C)



11月最終報告会デモについて

- 本チームの数値目標がfeasible 10
 であり,
- グループ間の技術統合により 可能である
- ことを示す



- 世代の異なる複数の計算機を 並べ、アプリケーションを実行
 - 各電力性能比が, 10年約1000 倍ラインに乗っていることを示す



計測用システム (デモ展示)



2012年マシン Tesla K10(Kepler) x4 油浸冷却





GT200世代でのGPU化によって60倍以上の電力効率向上。

先進的冷却技術の適用による HPCの更なる省電力化 ・計算システムの省電力化だけでなく冷却システムの 省電力化が必須

- 油浸冷却技術を導入・評価
 - 東工大GSIC概算要求「スパコン・クラウド情報基盤におけ るウルトラグリーン化技術の研究推進」と共同
 - PUE<1.05目標



米国Green Revolution Cooling社製の油槽



自作油槽に4 Kepler GPU マシンを液浸した様子



チップ温度低下・ファン除去







- 現在の世界トップ(2012年BG/Q, 2.1GFlops/Watt)以上の電力性能比
 - 年間のほとんどで、事実上PUE<1.0

TSUBAME2.0から2.5への中間アップグレード (2013年7月)

 TSUBAME2.0のGPU(Fermi 2050)を全部または部分的に最新のアク セラレータに交換



Future Work:より合理的な 電力最適化に向けて システム全体レベルの電力キャップに向けて フィードバック制御システムが必要



観測における課題:

- 電力測定の空間・時間的粒度問題 判断における課題:
- 電力バジェットの割り当て手法
 - 多数ノード間,多数アプリ間
 - CPUジョブとGPUジョブ間の公平性とは何か

2020年 エクサフロップススパコンへ向けて

Scalability of Future Computing



Copyright (c) 2012 Hiroshige Goto All rights reserved.

Machine	Power (incl. cooling)	Linpack Perf (PF)	Linpack MFLOPs/ W	Factor	Iotal Mem BW TB/s (STREAM)	Mem BW MByte/S / W	Factor
Earth Simulator 1	10MW	0.036	3.6	13,400	160	16	312
Tsubame1.0 (2006Q1)	1.8MW	0.038	21	2,368	13	7.2	694
ORNL Jaguar (XT5. 2009Q4)	~9MW	1.76	196	256	432	48	104
Tsubame2.0 (2010Q4)	1.8MW	1.2	667	75 x3	.6 ⁰	305	2 16
K Computer (2011Q2)	~16MW	10	625	80	3800	236	21
BlueGene/Q (2012Q1)	~12MW?	17	~1417	~35.3	~3000	~245	20
(TSUBAME2.5 (2013Q3))	1.8MW	~4	~2000	~22.5	~1000	~580	8.6
Tsubame3.0 (2015Q3)	1.8MW	~20	~11,000	~4.6 ~ x1	~4000 6	~2200 ~XC	2.3 5.6
EXA (2020)	20MW	1000	50,000	1	100K	5000	1



次世代気象予報コード ASUCA の

気象庁数値予報課との共同研究



