PCクラスタワークショップ in 神戸

日立のテクニカルコンピューティングへの 取り組み

2011/2/18

株式会社 日立製作所 中央研究所 清水 正明

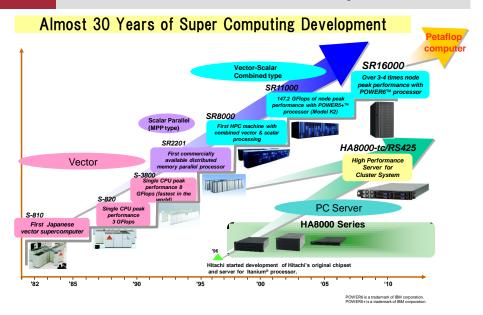


 $\hbox{Copyright} \circledcirc \hbox{Hitachi,Ltd.2011 All rights reserved} \qquad 1$

目 次

- **1** 日立テクニカルサーバラインナップ
- 2 日立サーバラインナップ
- 3 GPUコンピューティングへの取り組み
- 4 SC'10 日立展示

1-1 日立テクニカルサーバ: History & Future



Copyright © Hitachi,Ltd.2011 All rights reserved 3

日立テクニカルサーバ ラインアップ 1-2



SR16000XM1の紹介

電力性能比・価格性能比に優れた POWER7搭載のSR16000シリーズ次世代機 第一弾



日立開発·製造機 SR16000/XM1

日立開発・製造:

IBM社との戦略的アライアンスに基づく、 日立開発・製造のPOWER7搭載サーバ。 EP8000/750と共通プラットフォーム

電力性能比の向上:

SR16000モデルL1/L2と同等の32way SMP構成をPOWER7 4ソケットで実現。 ノード消費電力は約1/3に大幅に削減

抜群の価格性能比:

POWER7の圧倒的な性能と、戦略的な 価格付けにより、価格競合力を強化!

中規模SMPノード・クラスタ:

SRシリーズのコンセプトを受継ぐデザイン。 中規模SMPノードのクラスタシステムにより スループット指向の中規模システムに最適

Copyright © Hitachi,Ltd.2011 All rights reserved 5

日立サーバラインアップ

- ・ブレードサーバ
- ・ラックマウントサーバ

BladeSymphony ラインアップ

ブレードサーバを核に、ストレージ、ネットワーク、管理ソフトウェアを一体化した 統合サービスプラットフォーム BladeSymphony

- 各製品、充実のラインナップで、用途に応じた製品を提供
- 仮想化環境やソリューションを含めたシステム提供も可能



2-2

ハイエンドモデル BS2000/BS2000fx

高性能・高信頼志向のシステム向け



■ 仮想統合を実現する高信頼スケーラブル・ブレードサーバ

仮想化による集約、高速処理に適応した性能・拡張性

- サーバブレード間SMP接続(64cores MAX、メモリ 1TBMAX)

- I/Oスロット拡張装置(64スロット MAX)

日立サーバ仮想化機構Virtage標準搭載(*1)

メインフレームの高信頼・高可用化技術を継承

業界最高レベルの高効率電源

- CSCI Gold基準適合, 80 PLUS® GOLD認証取得(*2)

基幹システムの長期安定稼働を支援 - Eタイプ

ハードウェア長期保守対応

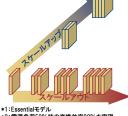
- ロングライフサポートサービス7年/10年(*3)











*1:Essentialモデル *2:電源負荷50%時の変換効率92%を実現 *8:BS2000 Eタイプにて提供 Copyright © Hitachi,Ltd.2011 All rights reserved

Q

小型高集積モデル BS320

より軽く、より小さく 高密度実装を追求 日5320

- 幅広い用途に対応する高集積・省電力ブレードサーバ
 - 高さ6U(約27cm)に最大10ブレード搭載可能
 - 最大重量約98kg/シャーシの軽量設計
 - 用途に応じた多彩なサーバブレードをラインアップ
 - ・日立サーバ仮想化機構Virtageに対応®
 - 高効率電源 (CSCI Silver基準適合, 80 PLUS® SILVER認証取得(+2))
 - ハードウェア長期保守対応 (ロングライフサポートサービス:7年)







出荷開始時期:2010/11/30















*1:PCI拡張サーバブレード**Virtage**モデルで提供 *2:負荷50%時の変換効率89%以上を実現

Copyright © Hitachi,Ltd.2011 All rights reserved 9

HA8000ラインアップ 2-4 HA8000/RS440 Xeon(X7560/X7550/E7540 /E7520) (4Processor) RAID追加機能 出荷開始時期:2010/11/30 ENERGY STARET'N Xeon(X5680/X5670/E5640/E5560/L5630/E5 HA8000/RS220 RAID追加機能 (2Processor) 3.5型 2TB SATA HDD追加 HA8000/TS20 出荷開始時期:2010/11/30 Xeon(X5670/E5640/E5620/E5503) SSD AC200V 低電圧メモリ Xeon(X5670/E5640/E5620/L5630/E55 RAID追加機能 HA8000/RS210 RAID追加機能 3.5型 2TB SATA HDD追加 3.5型 2TB SATA HDD追加 出荷開始時期:2010/11/30 出荷開始時期:2010/11/30 AC200V 低電圧メモリ (1Processor) HA8000/RS110 Xeon(X3480/X3470/X3460/X34 HA8000/SS10 Core i3-540/Pentium G6950 出荷開始時期:2010/11/30 RAID追加機能 RAID追加機能 2TB SATA HDD追加 3.5型 2TB SATA HDD追加 Xeon(X3480/X3470/X3460/X34

Copyright © Hitachi,Ltd.2011 All rights reserved 10

出荷開始時期:2010/11/30

RAID追加機能 3.5型 2TB SATA HDD追加

GPUコンピューティングへの取り組み

- 日立のGPGPUへの取り組み
- ・HPCシステムとアプリケーションの性能

Copyright © Hitachi,Ltd.2011 All rights reserved 11

3-1

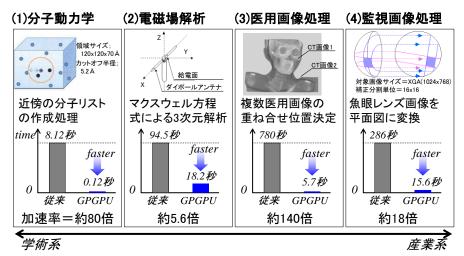
日立のGPGPUへの取組み(1)

計算科学を用いた研究開発分野でGPU利用が拡大中 研究所を中心に技術交流会を定期的に開催

- <利用分野(検討中含む)>
- •原子炉炉心解析
- ・火力・原子力発電の蒸気タービン流れ解析
- ・ボイラ燃焼効率解析
- ・粒子線治療シミュレーション
- ・材料物性・ナノシミュレーション
- ・機械(熱流体,構造,振動)
- •電磁場
- ・ライフサイエンス
- ・金融(実効金利計算)
- •他

日立のGPGPUへの取組み(2)

■GPGPU技術に関し、学術系~産業系アプリの先行評価・提案中

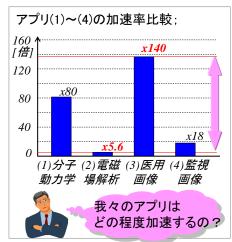


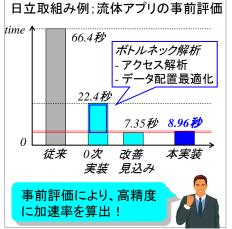
Copyright © Hitachi,Ltd.2011 All rights reserved 13

3-3

GPGPUの性質と日立の取り組み

- ■GPGPU性質:アプリによって加速率に大差。投資判断が難しい。
- ■日立取組:業務アプリを解析し、投資前に加速率を評価可能に。





- ◆社内にはGPUユーザ多数
- ◆利用技術·最適化技術も蓄積中
- ◆ソリューションメニューも整備 (事前評価からサポート)
- ◆GPU対応製品 (PCle x8,x16 搭載) HA8000 他 販売中
- ◆GPU搭載した大規模クラスタ(HPCシステム) 検討中

 $\hbox{Copyright} \circledcirc \hbox{Hitachi,Ltd.2011 All rights reserved} \quad 15 \\$

3 GPUコンピューティングへの取り組み

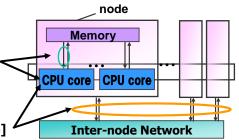
- 日立のGPGPUへの取り組み
- ・HPCシステムとアプリケーションの性能

3-5 システムバランスとアプリケーションの性能

アプリケーションの実効性能(効率)を以下の2点から定量評価

(1) ピーク演算性能に対する メモリバンド幅 [Byte/flop]

(2) ピーク演算性能に対する ネットワークバンド幅 [Byte/flop]



Example of high performance sever

(1),(2)の数値を変化させて実効性能への影響を見る(シミュレーション) ⇒ アプリケーションが求めるシステムバランスを求める

 $\hbox{Copyright} \circledcirc \hbox{Hitachi,Ltd.2011 All rights reserved} \quad 17 \\$

3-6

評価アプリケーション

4種類の並列アプリについて評価を実施 並列化スキームとプロセス間通信パターンは以下の通り

No.	Application	Calculation method	Partition	type
1	Ab initio MD	FFT, DGEM	Band Energy	1
2	Structural Calculation	Finite Element Method	3-Dim. space	2
3	Atmosphere	Difference Method	2-Dim.	3
4	Ocean	Difference Method	2-Dim.	3

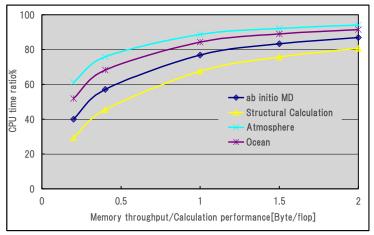
Type	1	2	3
Partition	x	z y	Z y X
Communication Pattern		→ 1	→
MPI function	MPI_allreduce (MPI_sum)	MP_send, MPI_recv MPI_allreduce(MPI_sum)	MP_send, MPI_recv

3-7 メモリバンド幅と演算性能のバランス

<u>(メモリバンド幅 GB/s) / (演算性能 GFlop/s) > 0.4[Byte/flop]</u>

- CPU time ratio becomes lower.

 0.2 ~ 0.4[Byte/flop]
- Better to keep more than 1.0 [Byte/flop]



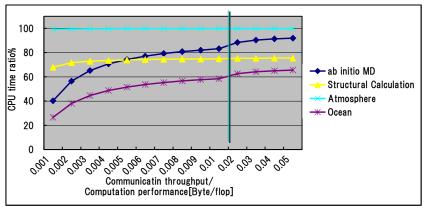
Copyright © Hitachi,Ltd.2011 All rights reserved 19

3-8

ネットワーク性能と演算性能のバランス

<u>(ネットワークバンド幅 GB/s) / (演算性能 GFlop/s)</u> > 0.02 [Byte/flop]

- The ratio of the communication time depends on the application.
- Better to keep more than 0.02 [Byte/flop]



グラフは Memory throughput/Calculation performance[Byte/flop] = 0.4 の場合

◆アプリケーションの要請

(メモリバンド幅 GB/s) / (演算性能 GFlop/s) > 0.4[Byte/flop] (ネットワークバンド幅 GB/s) / (演算性能 GFlop/s)

> 0.02 [Byte/flop]

◆マルチGPUシステムのバランス

<u>(メモリバンド幅 GB/s) / (演算性能 GFlop/s) = 0.25[Byte/flop]</u> (ネットワークバンド幅 GB/s) / (演算性能 GFlop/s)

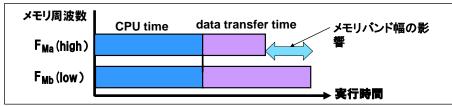
= 0.004 [Byte/flop]

◆実際にアプリケーションの性能はどうなるか?

Copyright © Hitachi, Ltd. 2011 All rights reserved 21

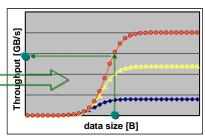
3-10 **GPU**システムでのアプリの性能推定

演算時間: GPUシステムのB/Fより実効効率を計算 サーバのメモリ周波数を変化させて実行時間を測定 実行時間をCPU時間とデータ転送時間に分解



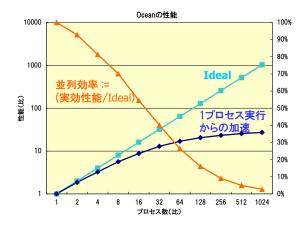
<u>通信時間:</u> PCクラスタで並列実行して 通信プロファイルを取得

> 個々の通信に対して 通信量から通信時間を「 グラフから求める



Ocean の並列性能「推定]

- ・同じ規模の問題をx一方向、y-方向の順で分割を繰り返す ⇒ strong scaling ・1プロセスのメモリ使用量がGPUに収まる最小の並列数を基準(グラフのプロセス数(比)=1) プロセス数(比)=1 のメモリ使用量 2.6GB ⇒ S2050 で利用可能な最大値
- ・プロセス数(比)=1のときの通信時間 ⇒ 全実行時間の9.4%
- ・演算効率は B/Fから推測 ⇒ 3.3%

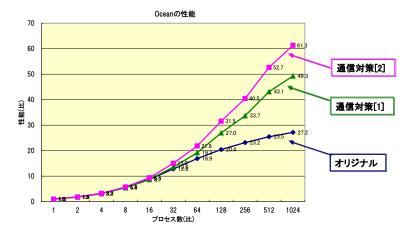


Copyright © Hitachi,Ltd.2011 All rights reserved 23

3-12

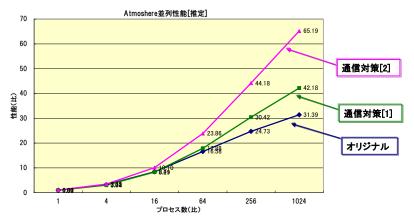
Ocean の並列性能[推定]の改善

- [2] 資源増強による対策 さらに、InfiniBannd を追加して2方向の隣接間通信を同時実行



Atmospher の並列性能[推定]

- ・プロセス数(比)=1 のメモリ使用量 1.8GB
- 演算効率は B/Fから推測 ⇒ 3.8%、プロセス数(比)=1 のとき通信時間は16%



- [1] 通信アルゴリズムによる対策 隣接通信する境界面を多層化して通信回数を削減
- [2] 資源増強による対策 さらに、InfiniBannd を追加して2方向の隣接間通信を同時実行

Copyright © Hitachi, Ltd. 2011 All rights reserved 25

3-14 まとめ:マルチGPUシステムとアプリの性能

◆マルチGPUシステムの特徴

<u>(メモリバンド幅 GB/s) / (演算性能 GFlop/s) = 0.25[Byte/flop]</u> メモリ性能バランスはPCサーバよりやや低め。効率は良い(80%) (ネットワークバンド幅 GB/s) / (演算性能 GFlop/s) = 0.004 [Byte/flop] ネットワーク性能が相対的に低く見える。 レイテンシ > 20 μs GPU Direct は データ長 > 16KB で効果大

◆アプリケーションの並列実行性能

- -GPUのメモリを最大に使用した weak scaling ではネットワークの弱さは 目立たない
- *strong scaling でのスケーラビリティ劣化は早い 今回の評価では 16GPUで約10倍加速、以後急速に劣化
- strong scaling でのスケーラビリティを保つには努力が必要 演算に隠蔽できれば良い

転送データ長が大きい場合はパイプライン化

転送データ長が小さい場合は通信回数の削減

演算数が増えても通信回数削減を検討(shadow領域の多層化など)

4 SC'10 日立展示

 ${\it Copyright} @ {\it Hitachi,Ltd.2011} \ {\it All rights reserved} \quad 27$

4-1

日立ブース

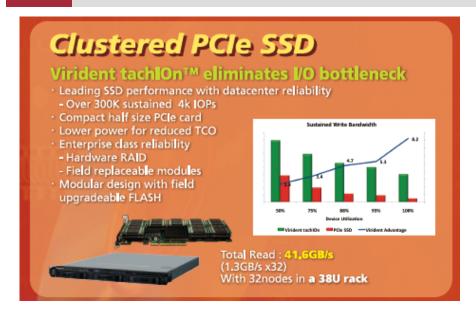


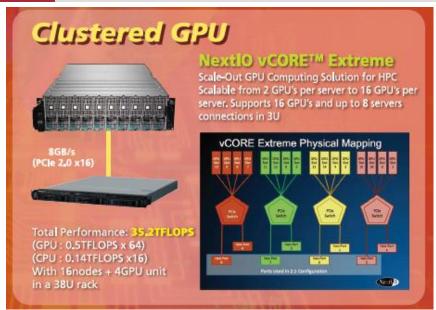
次世代サーバボード



Copyright © Hitachi,Ltd.2011 All rights reserved 29

4-3





Copyright © Hitachi,Ltd.2011 All rights reserved 31



Copyright © Hitachi,Ltd.2011 All rights reserved 32