

## ヘテロ型スーパーコンピュータ TSUBAME の Linpack による性能評価

遠藤 敏夫<sup>†1</sup> 松岡 聡<sup>†1,†2</sup> 橋爪 信明<sup>†3</sup>  
長坂 真路<sup>†4</sup> 後藤 和茂<sup>†5</sup>

TSUBAME スーパーコンピュータは、655 ノード 10480 Opteron core と 360 枚の ClearSpeed SIMD アクセラレータボードを備えるヘテロ型の大規模クラスタシステムである。TSUBAME は Linpack ベンチマークで 38.18TFlops を記録し、この結果により 2006 年 6 月の Top500 ランキングに 7 位としてランクされた。しかしその測定には Opteron のみが使われ、アクセラレータは用いられていない。本論文は汎用 CPU とアクセラレータによる不均一な環境において Linpack 性能を得る上での課題について論じ、解決法を述べる。16 ノード 256 CPU core を用いた予備実験では、アクセラレータ 8 枚を加えたときに 8.2%、16 枚を加えたときに 19% の性能向上が観測された。

### Performance Evaluation of TSUBAME Heterogeneous Supercomputer with Linpack

TOSHIO ENDO,<sup>†1</sup> SATOSHI MATSUOKA,<sup>†1,†2</sup> NOBUAKI HASHIZUME,<sup>†3</sup>  
MASAMICHI NAGASAKA<sup>†4</sup> and KAZUSHIGE GOTO<sup>†5</sup>

The TSUBAME supercomputer is a heterogeneous large-scale cluster system, which is equipped with 10480 Opteron CPU cores on 655 nodes and 360 ClearSpeed SIMD accelerator boards. The TSUBAME system has achieved 38.18TFlops with Linpack benchmark and is ranked 7th in the Top500 supercomputer ranking in June 2006, even though the measurement is done without any accelerator boards. This paper discusses issues to obtain high Linpack performance on heterogeneous systems with general purpose processors and accelerators, and describes solutions. Through preliminary experiments with 256 CPU cores on sixteen nodes, we observed +8.2% speed-up when eight accelerators are added, and +19% with sixteen accelerators.

#### 1. はじめに

東京工業大学学術国際情報センターは 2006 年 4 月に国内最速のスーパーコンピュータ TSUBAME を導入した。このシステムは東工大キャンパスグリッドの中心として、グランドチャレンジシミュレーションを始めとする科学技術計算に加え、全学ストレージサービスや事務系サーバを含めたホスティングサービスにも活用されることが期待されている。このシステム

は Intel 互換 CPU である AMD Opteron と Linux OS を採用することにより、汎用性とプログラミングの容易さを確保する一方で、ClearSpeed 社の SIMD アクセラレータボードを備えることにより、科学技術計算におけるスペース性能比、電力性能比を向上させることをねらっている。システムが備える 10480 の Opteron CPU core の理論性能は 50TFlops、360 枚のアクセラレータの理論性能は 35TFlops であり、合計で 85TFlops である。

Linpack ベンチマークはスーパーコンピュータの科学技術計算性能の指標として、世界ランキングである Top500<sup>2)</sup> でも用いられている。TSUBAME は 2006 年 6 月のランキングにおいて 38.18TFlops を記録し、世界 7 位、国内トップとしてランクされた。この測定に用いられたのは Opteron のみであるため、アクセラレータを併用することによりさらなる性能向上が期待できる。

本論文では、汎用 CPU とアクセラレータの双方を

<sup>†1</sup> 東京工業大学  
Tokyo Institute of Technology  
<sup>†2</sup> 国立情報学研究所  
National Institute of Informatics  
<sup>†3</sup> サン・マイクロシステムズ (株)  
Sun Microsystems K.K.  
<sup>†4</sup> 日本電気 (株)  
NEC Corp.  
<sup>†5</sup> テキサス大学オースティン校  
The University of Texas at Austin

用いた不均一（ヘテロ）な環境における Linpack ベンチマークの結果を報告する。実験には Linpack の並列実装である High-Performance Linpack(HPL)<sup>8)</sup> を用いたが、これは本来均一な環境のために設計されている。不均一な環境に適合させるため、以下の技法を提案する：アクセラレータの有無による起動プロセス数の調整、HPL の look-ahead 手法の改良、MPI 実装の spin wait の変更。TSUBAME の 16 ノード 256 CPU core を用いた実験を通して、アクセラレータ併用による性能向上のためにはこれらの技法が必須であることを示す。

## 2. TSUBAME システム

TSUBAME では、655 ノードの計算サーバ SunFire X4600 と、合計 1.1PBytes のストレージサーバが InfiniBand により接続されている (図 1)。以下、本論文に関連の深い部分について概要を示す。

**ファットノード型計算サーバ:** 各 SunFire ノードは、dual core 2.4GHz Opteron CPU \* を 8 個を持ち、16 CPU core が 32GB のメモリを共有する。またノードは 10Gbps のバンド幅の InfiniBand host channel adapter (HCA) を 2 つ持つ。オペレーティングシステムは 64bit 対応 SuSE Linux Enterprise Server 9 であり、Linux カーネルバージョンは 2.6.5 である。また 655 ノードのうち 360 ノードが、PCI-X バスによって接続される ClearSpeed アクセラレータボードを持つ。

**高速インターコネクト:** 各ノードは InfiniBand により 288 ポートの Voltaire ISR9288 スイッチに接続される。スイッチ間は、InfiniBand 24 本により接続される\*\*。

アプリケーションプログラムは、Message passing interface(MPI) の一実装である Voltaire MPI により並列プログラムを動作させることができる。Voltaire MPI は InfiniBand を直接アクセスすることにより高性能を実現する。

**アクセラレータ:** ClearSpeed Advance Accelerator Board<sup>1)</sup> は、理論性能 96GFlops の演算能力を持つ PCI-X ボードである。アクセラレータボードは 2 つの ClearSpeed CSX600 SIMD プロセッサと 1GB の DRAM を持つ。各プロセッサには、0.5GFlops の演算能力を持つ 96 個の PE が含まれる。なお、アクセラレータ上の演算の入出力データは、1.06GBytes/s

\* 16 ノードのみ、2.6GHz CPU

\*\* 部屋をまたがって設置されている部分については光ファイバーを利用する

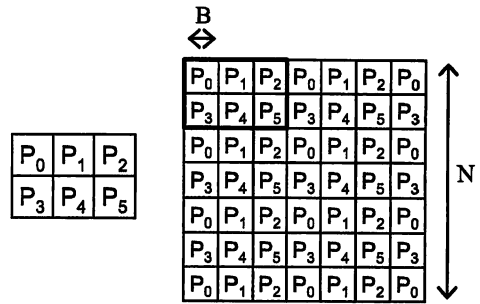


図 2 (左)  $P \times Q = 2 \times 3$  のプロセス構成例、(右) 二次元ブロックサイクリック分割

の PCI-X バスを介してホストと通信する必要がある。

アクセラレータを利用するために、SIMD 並列プログラミング言語  $C^*$  と、BLAS(基本線形演算) ライブラリが提供されている。本論文ではこれらのうち BLAS ライブラリを利用する。

## 3. HPL

### 3.1 HPL の概要

HPL は正方密行列を係数とする連立一次方程式をブロック化ガウス消去法で解く、MPI 並列ソフトウェアである。指定された行列サイズ  $N$  に対して乱数行列を生成し、方程式を解き、その速度を FLOPS 値で評価する。

計算に参加するプロセス群は概念的にサイズ  $P \times Q$  のプロセス格子を形成し、行列はプロセス格子に従って二次元ブロックサイクリック方式で分散される (図 2)。以下、行列サイズを  $N$ 、ブロックサイズを  $B$  とする。計算のほとんどの部分をガウス消去法が占め、その各ステップ (ステップ番号  $k$  とする) は、以下のような処理からなる。

**パネル分解:** 第  $k$  ブロック列はパネル列と呼ばれ、その箇所の LU 分解を部分ピボット選択を用いて行う。

**パネルブロードキャスト:** パネル列の各ブロックの内容を他プロセスへブロードキャストする。ここではプロセス格子の各行内での通信が発生する。

**行交換通信:** 部分ピボット選択の結果に基づき、行交換を行う。ここではプロセス格子の各列内での通信が発生する。

**更新計算:** パネル列と、行交換後の第  $k$  ブロック行の内容を用い、行列の未分解部分の更新計算を行う。行列積演算が発生する。

以上のうち、パネル分解の計算量総計は  $O(N^2 B)$ 、パネルブロードキャストと行交換通信の通信量総計は

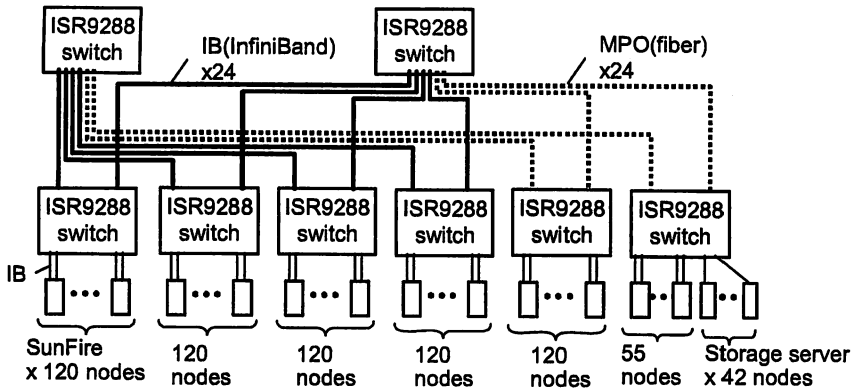


図 1 TSUBAME のシステム構成図

$O(N^2(P+Q))$ , 更新計算の計算量総計は  $O(N^3)$  である。このことから、最も時間がかかるのは更新計算であり、その傾向は  $N$  が大きい程強いと分かる。そのため、並列 Linpack ベンチマークにおいては、 $N$  をメモリ量の限界に近づけるように大きくとり、高速な行列積を行う数値演算ライブラリを用いることにより、性能を上げることができる。

### 3.2 10,368CPU を用いた性能評価

TSUBAME システムの 648 ノード 10,368CPU core を用いた HPL の実行結果を示す。MPI ライブラリとしては Voltaire MPI を、数値演算ライブラリとしては、GOTO BLAS<sup>5)</sup> を用いた。GOTO BLAS ライブラリはユーザが指定した数のスレッドを用い、行列演算をノード内で並列化して行うことができる。本実行では GOTO BLAS が用いるスレッド数を 2 とし、各ノードに 8 プロセスずつ起動することにより、ノード内の 16 CPU core を利用する。

実験に用いた HPL のパラメータを表 1 に示す。パネルブロードキャストのアルゴリズムとしては、占有バンド幅を下げることを優先し、リング型トポロジー (1ring) とした。Look-ahead は HPL が備える最適化機構の一つであり、各ステップの更新計算の最中に以降のステップのパネルブロードキャストを行うことにより、通信待ちの時間を削減するものである。実験では次のステップ (depth=1) のブロードキャストとオーバーラップさせることとした。

行列サイズ 1,334,160, プロセス数  $36 \times 144 = 5184$  の実験において、実行時間約 11.5 時間、速度 38.18TFlops を達成した。システムの Opteron の合計理論性能は 49.87TFlops のため、実効性能は 76.6% となる。

Matrix size	1334160
Block size	240
Process mapping	Row-major
# of processes ( $P \times Q$ )	$36 \times 144$
Panel factorization	Right-looking
NBMIN, NDIV	4, 2
Panel broadcast	1ring
Look-ahead depth	1
Swap	Mix (threshold=240)
Matrix form	L1 trans, U trans
Equilibration	yes
Alignment	8 double words

表 1 10,368CPU を用いた評価における HPL のパラメータ

### 3.3 不均一環境における問題点

前節では TSUBAME の汎用 CPU を用いたが、これに加えて SIMD アクセラレータも併用することにより、さらなる性能向上が期待される。しかしながら、HPL では各プロセスにほぼ均一の行列データが分散され、計算量もほぼ均一となる。そのため、異なる速度を持つ汎用 CPU と SIMD アクセラレータを有効利用するための手法が必要となる。

HPL の性能に大きく影響するのは行列積の速度であり、 $M$  を  $B \leq M \leq N$  とするとき、 $M \times B$  行列と  $B \times M$  行列の積を特に多用する。GOTO BLAS の場合 (2 スレッド利用)、本論文の実験で用いるサイズの行列について、行列積性能は安定的に 8.4 ~ 8.8GFlops である。一方、ClearSpeed により提供される BLAS ライブラリ (以下 CS BLAS) の性能を図 3 に示す。なお、SIMD アクセラレータの PE 数から  $M, B$  を 192 の倍数としている。性能は行列サイズによって大きく変動し、このグラフの範囲では 18GFlops ( $M = 1920, B = 384$ ) から 38GFlops ( $M = 11520, B = 960$ ) となることが分かる。不均一環境では、GOTO BLAS と CS BLAS の双方を有効利用するために、各プロセスの計

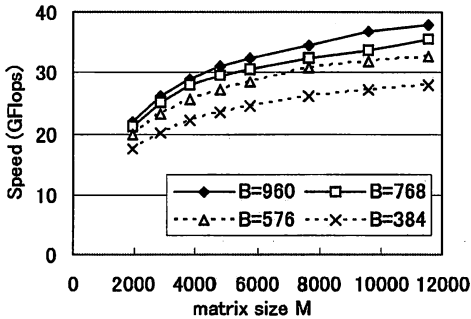


図3 ClearSpeed BLAS によるサイズ  $(M \times B) \times (B \times M)$  の行列積性能

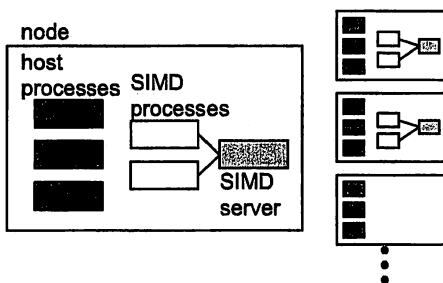


図4 CPU とアクセラレータを利用するためのプロセス構成

算性能をなるべく均等化する必要がある。

なお、図3の性能は理論値96GFlopsに比べると効率40%以下である。この原因はPCI-Xバスのバンド幅の低さに加え、現在のライブラリの計算スケジューリングと通信スケジューリングの問題であり、将来の版で改善されることが期待される<sup>6)</sup>。

#### 4. 提案手法

TSUBAMEシステムにおいては、汎用CPUとアクセラレータが混在するという不均一性に加え、現在は半数強のノードのみがアクセラレータを持つ。このような環境で、均一環境向けに設計されているHPLを、大きく変更せずに効率的に動作させるため、以下の技法を提案する。

##### 4.1 プロセス構成

これまでと同様にGOTO BLASを用いるプロセスをホストプロセスと呼ぶ。それに加えて図4のように、CS BLASを用いるプロセス(SIMDプロセス)を導入する。ホストプロセスとSIMDプロセスの起動数を調節することにより、各プロセスの計算能力を均等に近付けることを目的とする。

ただし、現在のCS BLASは同時に単一のプロセ

スによってしかアクセスできない。これではプロセス数調整の際の制限が厳しかったため、CS BLASを直接アクセスするためのSIMDサーバを設け、アクセラレータを持つ各ノードに1つずつ起動する。複数のSIMDプロセスは、mmapによる共有メモリ上を通して、SIMDサーバへ行列積計算を依頼する。SIMDプロセスは、更新演算以外はホストプロセスと同様に動作する。

次章の実験では、アクセラレータを持つ各ノードに7つのホストプロセスを起動し、それらが用いるGOTO BLASのスレッド数は2つとする。また、SIMDプロセス一つあたりの性能がホストプロセスよりも遅くならないよう、SIMDプロセスの数は3とする。なお、ホストプロセス数を8とはしなかったのは、SIMDサーバのCPU利用率が無視できず、そのための余裕が必要なためである。アクセラレータを持たないノードには、ホストプロセスのみを8つずつ起動する。

##### 4.2 Look-aheadの改良

3.2節で述べたように、HPLは更新計算とパネル通信をオーバラップさせる最適化を採用しているが、以下に述べるような問題がある。現在の実装では、更新対象部分のうち一部について行列積を行い、その度にパネル通信メッセージが到着しているかチェックするという処理を繰り返す。パネル通信が終わった後は残りの更新対象部分を一度の行列積関数で計算する。このアルゴリズムでは、パネル通信以前の行列積が細切れに行われることになる。図3に示したように、CS BLASでは行列積対象が小さいときに大きく性能が低下してしまうため、望ましくない。

そこで行列積の細切れ化を防ぐためにスレッド化を行い、更新計算の行列積関数呼び出しを、通信とは別スレッドで行うこととする。

##### 4.3 MPIのspin waitの変更

本論文の実験に用いているVoltaire MPIの現在の実装は、MPI\_Send, MPI\_Recvなどのブロック通信がCPUを消費し続けるという問題がある。これはMPI関数の内部でspin waitをしているためと考えられる。本章で述べる手法ではCPU数を越える数のスレッドを起動するため、CPUの浪費は望ましくない。

回避策として、ブロック通信関数の代わりにMPI\_Isendなどのノンブロック通信関数をusleepでスリープしながら断続的に呼び出すように置き換えた。これによりCPUの浪費は抑えられるが、本来よりもメッセージの到着に気づくのが遅れる問題が起こる。将来のVoltaire MPIの改良によりこの回避策は必要なくなるだろう。

## 5. 不均一環境での性能評価

CPU と SIMD アクセラレータの両方を用いた場合の Linpack 性能を報告する。TSUBAME システムのうち、ClearSpeed アクセラレータを持つ 16 ノードを実験に用いた。HPL パラメータは、明示的に示すもの以外は表 1 に示すものを用いた。

図 5 に結果を示す。プロセス構成は 4.1 節で述べた通りとし、4.2、4.3 節で述べた技法を適用した。また、ホストプロセスについて、Linux 2.6 の sched.setaffinity システムコールを用い CPU へバインドした。グラフ中の SIMD16 はアクセラレータを全て併用した場合、SIMD8 は 16 枚のうち 8 枚のみを併用した場合、CPU は汎用 CPU のみを利用した場合である。プロセス数は SIMD16 で  $8 \times 20$ 、SIMD8 で  $8 \times 18$ 、CPU で  $8 \times 16$  である。

ブロックサイズ  $B$  は、試行したうちで最も性能の良いものを選択した。SIMD16 と SIMD8 で  $B = 576$ 、CPU で  $B = 240$  である。図 3 によれば  $B$  がより大きい方が CS BLAS の行列積性能は上がるが、HPL の性能は下がった。この原因については調査中であるが、更新計算以外によるボトルネックが  $B$  に伴い大きくなるためと推測している。

アクセラレータを併用することにより、サイズが約 90,000 以上の行列において性能向上が得られていることが判る。 $N = 184,320$  において、SIMD16 は 1169GFlops、SIMD8 は 1062GFlops を達成している。CPU のみの 981.2GFlops に比較すると、SIMD16 では 19%、SIMD8 では 8.2% の向上が得られた\*。

次に 4.2、4.3 節で述べた技法の影響を図 6 に示す。グラフは全てアクセラレータ 16 枚併用の場合であり、二つの提案技法を用いる場合 (Opt)、用いない場合 (NoOpt) を示す。また、 $B = 576$ 、 $B = 960$  の場合を示す。NoOpt に比べると Opt の性能は高く、提案技法の効果は観測された。特に NoOpt( $B=576$ ) の性能の頭打ちが顕著であるが、これは元々の HPL が行列積の細切れ化を行った場合に、その悪影響は  $B$  が小さい場合に強く及ぶためと考えられる。

これまでの実験ではノードあたり 7 個のホストプロセス、3 個の SIMD プロセスという比較的細粒度なプロセスを用いた。プロセス粒度を変更した場合の性能を図 7 に示す。Fine はこれまでと同じ設定の場合で、Coarse ではノードあたり 2 個のホストプロセス

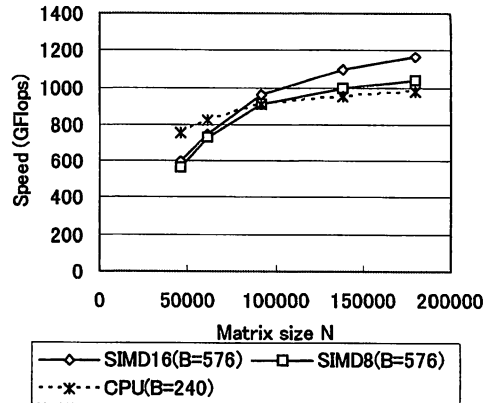


図 5 アクセラレータ併用した場合の 16 ノードの Linpack 性能

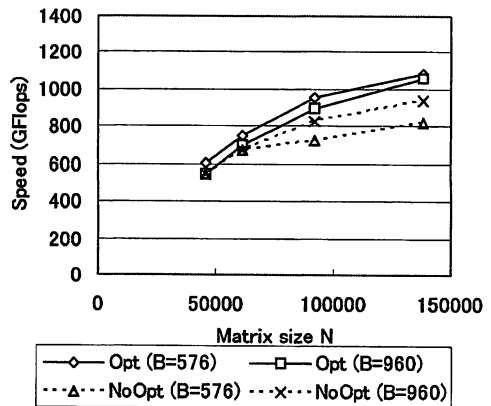


図 6 提案技法の影響の評価。アクセラレータ 16 枚使用時

(GOTO BLAS スレッド数は 7)、1 個の SIMD プロセスを用いた。このグラフにおいては、プロセス粒度はあまり影響しないことが判る。特に、Fine よりも Coarse の方が各プロセスの持つ部分行列が大きくなり、行列積性能は有利と予想されたが、この実験においてはその影響は見られていない。一方、将来 ClearSpeed BLAS の性能が向上した場合には、プロセス数を細かく調整できる Fine の方が有利になると考えられる。

## 6. 関連研究

不均一なクラスタ環境において行列演算を効率的に行うアプローチの一つは、各プロセスにノード性能に比例したサイズの部分行列を割り当てることである<sup>3),4)</sup>。このために二次元再帰分割をはじめとする分割法が提案、評価されてきた。しかし既存の多くの並

\* なお、より大きい  $N$  を用い CPU のみで 1042GFlops を達成している

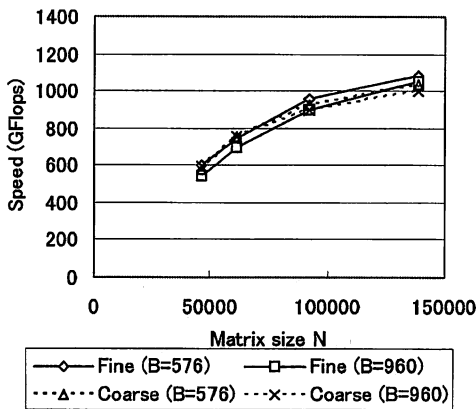


図7 プロセス粒度を変更した場合の性能。アクセラレータ 16 枚使用時

列プログラムは均一な分割を行っており、それを不均一な分割に変更するにはプログラム全般に渡る変更が必要である。

本研究のアプローチはより笹生ら<sup>9)</sup>のものに近い。その研究では HPL の一部を改変し、高性能ノードにおけるプロセスが他プロセスの整数倍の部分行列を持てるようにした。本研究は SIMD アクセラレータを持つファットノードクラスターで大規模評価を行った点、細粒度プロセスが実用的であることを示した点が異なる。

岸本らは各 CPU に異なる数の HPL プロセスを割り当てた場合の性能モデルを提案し、少ない数の実験試行から、ふさわしいプロセス数を推定する手法を提案した<sup>7)</sup>。このモデルにおいては、本研究で指摘したようなプログラムの実装上の問題点については触れていない。

## 7. おわりに

Linpack ベンチマークの実装の一つである HPL を基に改変したプログラムを用い、汎用 CPU である Opteron と ClearSpeed SIMD アクセラレータを備える TSUBAME システムの評価を行った。アクセラレータ併用による性能向上を得るために、パネル通信の lookahead 機構の改変など、HPL への部分的な変更が必要であることを示した。今回の実験ではアクセラレータ側の演算ライブラリとして ClearSpeed BLAS を用いたが、将来の版で性能向上が期待されている。その場合にも、本論文の手法はプロセス数調整により容易に適応可能である。また、今回は 16 ノードというシステムの一部を用いた実験を行ったが、今後シス

テム全体の評価を行っていききたい。

今回 HPL に対して行った変更は HPL 特有のものであるが、複数プロセスでアクセラレータを共有する手法や、通信と計算を非同期に行うアプローチは、広範囲のアプリケーションに適用可能と考えられる。今後 TSUBAME 上で、アクセラレータを併用する、より一般的な高性能計算を実現するための実行モデルを研究する予定である。

謝辞本研究にあたって ClearSpeed 社の Chris Piercy 氏、John Gustafson 氏、Sun microsystems 社の Antoine Petitet 氏に協力、助言頂いた。本研究の一部は科学研究費補助金 (特定領域研究 課題番号 18049028, 若手研究 (B) 課題番号 17700050) の援助による。

## 参考文献

- 1) ClearSpeed Technology Inc. <http://www.clearspeed.com/>.
- 2) TOP500 supercomputer sites. <http://www.top500.org/>.
- 3) Phyllis E. Crandall and Michael J. Quinn. Block data decomposition for data-parallel programming on a heterogeneous workstation network. In *Proceedings of IEEE International Symposium on High Performance Distributed Computing (HPDC)*, pages 42–49, 1993.
- 4) Toshio Endo, Kenji Kaneda, Kenjiro Taura, and Akinori Yonezawa. High performance LU factorization for non-dedicated clusters. In *Proc. of CCGrid*, 2004.
- 5) Kazushige Goto. High-performance BLAS. <http://www.cs.utexas.edu/users/flame/goto/>.
- 6) John Gustafson. ClearSpeed Technology Inc., private communication.
- 7) Yoshinori Kishimoto and Shuichi Ichikawa. An execution-time estimation model for heterogeneous clusters. In *Proceedings of IEEE International Parallel & Distributed Processing Symposium (IPDPS)*, 2004.
- 8) A. Petitet, R. C. Whaley, J. Dongarra, and A. Cleary. HPL - a portable implementation of the high-performance Linpack benchmark for distributed-memory computers. <http://www.netlib.org/benchmark/hpl/>.
- 9) 笹生 健, 松岡 聡, and 建部 修見. ヘテロなクラスター環境における並列 LINPACK アルゴリズム. In *並列処理シンポジウム JSP2002 論文集*, pages 71–78, 2002.