

クラスタ型ベクトル並列スーパーコンピュータ S-3000 クラスタシステムのアーキテクチャと特性評価

田中 輝雄¹, 深川 正一¹, 井村 淳一², 後 保範²,
玉置 由子³, 榊原 忠幸¹, 稲上 泰弘³

tetanaka@kanagawa.hitachi.co.jp

¹(株)日立製作所 汎用コンピュータ事業部

²(株)日立製作所 ソフトウェア開発センタ

³(株)日立製作所 中央研究所

概要

最大演算性能 96 GFLOPS の S-3000 クラスタシステムを開発した。主記憶共有型ベクトルプロセッサ S-3800 相当のシステムを高速接続機構を介して最大 4 組接続するクラスタ型ベクトル並列スーパーコンピュータである。2 階層の並列処理機構として、クラスタ内は S-3800 で開発した主記憶共有並列処理機構を利用し、クラスタ間はメッセージパッシング型の並列処理方式を基本に、高速データ通信機構を開発した。実機による特性評価を行ない、2 階層の並列処理機構が効率良く実行できることを確認、また、連立一次方程式の LU 分解法プログラムを用いて実効性能 78.2 GFLOPS を実現し、80% 以上の実効効率を得た。

The cluster-type vector parallel supercomputer "S-3000 Cluster System" - its architecture and characteristic evaluations -

Teruo TANAKA¹, Masakazu FUKAGAWA¹, Junichi IMURA², Yasunori USHIRO²,
Yoshiko TAMAKI³, Tadayuki SAKAKIBARA¹, Yasuhiro INAGAMI³

tetanaka@kanagawa.hitachi.co.jp

¹ General Purpose Computer Division, Hitachi Ltd.

² Software Development Center, Hitachi Ltd.

³ Central Research Laboratory, Hitachi Ltd.

Abstract

We have developed the *S-3000 Cluster System*, a cluster-type vector parallel supercomputer which achieves 96 GFLOPS as a peak performance. The *S-3000 Cluster System* is composed of a maximum of four machines similar to the S-3800 TCMPs (Tightly Coupled Multi Processors) linked by a high speed interconnection network. The *S-3000 Cluster System* uses two levels of parallelism. The first is a TCMP system with a shared memory. The second is a distributed memory scheme composed of four TCMP systems which supports low-overhead message-passing. We have shown effectiveness of this approach by performance evaluation on a real machine. We have obtained sustained performance of 78.2 GFLOPS and over 80% of the theoretical peak performance on a LU decomposition algorithm for large scale linear equations.

1. はじめに

並列処理方式の実用化によりスーパーコンピュータの性能トレンドは従来の半導体テクノロジーの進歩に比例したトレンド曲線から大幅に上方修正された。現在4台から16台程度のベクトルプロセッサが主記憶を共有する主記憶共有型システムや、富士通VPP500[1]のような完全分散記憶型のシステムが実用化されている。日立では平成4年12月に主記憶共有型スーパーコンピュータS-3800を開発、出荷している[2-5]。

また最近、NECやCRAYからは次世代スーパーコンピュータとして、主記憶共有型システムをネットワークあるいはファイルシステムで複数台接続するクラスタ構成のシステムが発表されている。このようなクラスタ型システムとしては、従来からイリノイ大学のCedar[6]やCRAY社[7]などの例がある。

クラスタ型システムは、

- (1) 現状の主記憶共有型システムの上位互換性を持つため、いままでの運用を継続でき、計算センタ向きである、
- (2) クラスタを構成する主記憶共有型システムは比較的大容量の主記憶を共有するため、並列化できない応用に対しても1台のプロセッサにより大規模計算を行なうことができる、
- (3) ある程度の並列処理までは、主記憶共有型で対応することができるため、従来の自動ベクトル化/並列化機構を用いることができる、

などの特徴を持つ。さらに、

- (4) クラスタ間の結合網を用いてデータを授受し、クラスタ間で協調処理を行ない、超大規模の数値シミュレーションを実行することができる。

しかし、この協調処理はクラスタ内/クラスタ間の2階層の異なる並列処理で実行するために、並列処理オーバヘッドの増大、プログラミングの複雑さなどの問題が発生する恐れがある。

このようなクラスタ型のシステムとして、S-3800相当のスーパーコンピュータシステムをひとつのクラスタとし、このクラスタを高速接続機構（以下では、HSICと略す）を用いて接続し、クラスタ間で協調して大規模数値シミュレーションの実行を可能とするクラスタ型ベクトル並列スーパーコンピュータS-3000クラスタシステムを開発した。

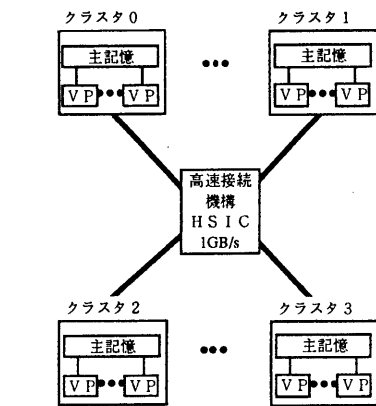


図1 スター型結合によるクラスタ型並列システム
Fig. 1 Star connected cluster type vector-parallel system.

本論文では、S-3000クラスタシステムのアーキテクチャおよび大規模数値計算処理実行時の特性解析結果を示す。

2. 2階層並列処理

S-3000クラスタシステムは、クラスタ内/クラスタ間の2階層の並列処理を行なう。

(1) クラスタ内並列処理

主記憶共有型の並列処理機構を用いる。そのためS-3800において開発したコンパイラによる自動ベクトル化機構、自動並列化機構をそのまま利用することができる。

(2) クラスタ間並列処理

クラスタ間並列処理では、クラスタ間通信処理の扱いが重要となる。多くの数値シミュレーションでは、次の特徴を持つ。

- (a) 転送バタンの静的に決定することができる。
- (b) したがって、プログラム上でデータを送信側プロセッサ（またはクラスタ）で送信する順番に受信側プロセッサ（またはクラスタ）で受け取ることを保障することができる。
- (c) 大容量データ転送が主でありスループット重視である。このため、通信形態としてはメッセージパッシング型が効率が良い。ここでは通信処理の高速化であり、(a)ソフトウェアオーバヘッドの削減、(b)不要な同期処理の排除が重要課題となる。

3. S-3000クラスタシステムの実現

3.1 構成

高速接続機構（HSIC）は、最大1GB/sのデータ転送能力を持つGBオーダの大規模メモリシステムである。複数の計算機システム（本論文ではクラスタ）から共有することが可能であり、HSICと個々のクラスタ間の最大データ転送スループットは最大500MB/sである。HSICは複数のクラスタと接続することが可能であることから、種々の接続形態が考えられる。

たとえば、図1に示すように、HSICを中央に配置しハブとして位置付け、すべてのクラスタから接続するスター型の接続である。これは、すべてのクラスタ間で直接データのやり取りを行なうことができるので汎用性が高い。しかし、複数のクラスタか

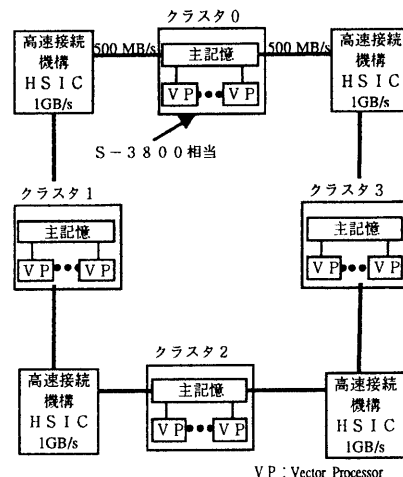


図2 リング結合によるクラスタ型並列システム
Fig. 2 Ring connected cluster type vector-parallel system.

ら同時にアクセスが行われたときはHSICにおいて競合が発生し、クラスタ間のデータ転送処理の性能低下をもたらす。大規模シミュレーションでは、SPMD(Single Program Multiple Data stream)的に記述できる処理が多く、クラスタ間データ転送処理は同時に発生する可能性が高い。そこで、すべての隣接するクラスタ間にHSICを配置し、リング状に接続することにより、HSICの台数分の通信性能向上をはかる(図2)。

今回の開発では、大規模シミュレーションの高速処理を行なうためにリング型の接続形態を取ることにした。また以下では、並列処理のプログラミングスタイルはSPMDを対象とする。

3.2 クラスタ間高速データ転送方式

クラスタ間の並列処理方式の概要を近接作用問題を例に説明する。基本的な近接作用問題は、領域を分割し、それぞれを分散記憶に配置し、分割された境界面のデータをクラスタ間で転送しあう。このデータ転送にメッセージパッシング方式を用い、送り手側のsend処理と、受け手側のreceive処理により実現する。

メッセージパッシング型転送で送信側クラスタのSEND命令と受信側クラスタのRECEIVE命令を独立に実行させるためには、中間にバッファ領域を必要とする。ふつうは、このバッファ領域を受手側の主記憶上に確保するが、S-3000クラスタシステムでは、この役割をHSICが行なう。クラスタ間データ通信の原理を図3に示す。それぞれのデータ転送パスごとにHSIC上にFIFO型のバッファを構築する。そして、データ転送命令(SENDおよびRECEIVE)でこのFIFOを指定する。データの送受信を監視する制御変数をそれぞれの論理的なパスごとに準備することにより、制御変数に対する排他処理を必要としない。

HSIC上のFIFOの実現方法とそれを用いたクラスタ間の転送プロトコルを図4に示す。ここでは、それぞれのクラスタごとにHSICを分割して割り付け、各クラスタの管理領域に対して、他方のクラスタからの書き込みを許さない方式とした。これは双方からの書き込みを許すことにより必要となるシリアル処理を不要にするためである。一方のクラスタの管理対象とするHSIC上の領域をそれぞれ他方のクラスタから読出しのみできるようにし、一方通行のデータ転送パスを構築する。

この領域を対称形に構成することにより、双方向のデータ転送パスの実現を可能とする。このとき、FIFO領域の構成(分割数、サイズ)はプログラム開始時に、そのプログラムに最適な割り付けをプログラム中で行なう。分割数を多くすることにより、送信側SEND命令の受信側RECEIVE命令からの先行発行数を大きく許すことができる。もし、並行して2つ以上の異なるデータを転送したいときは、このデータ転送パスを複数設定する。なお、HSIC上へのFIFO領域の確保は、各クラスタのHSIC割り付け時にOSが設定する。

FIFOはHSICのメモリ上に構成したFIFOバッファとその制御変数NrNsからなる。このプロトコルでは、1回のデータ転送処理に、FIFOへのデータ書き込み、FIFOからのデータ読出しのほか、制御変数NrNsの読出し、書き込みを含めて6回のHSICアクセス処理が必要となる。しかし、事前の制御変数NrNsの読出し時に、各クラスタ内にそのコピーを確保することにより、FIFOバッファがあふれそうになるまでは、新たに制御変数NrNsの読出し処理は必要なくなる。

高速化のためには制御変数NrNsをHSICから独立させ、専用制御論理を実現することも考えられるが、今回は大量データ転送時のスループット重視としHSICのアクセスバスを共用した。専用制御論理を実現したときの性能は4.1節で考察する。

3.3 ユーザインタフェース

プログラムを実行するためのジョブの起動、終結方法および並列化プログラミング方法を示す。

3.3.1 ジョブの起動、終結

PVM(Parallel Virtual Machine)の環境を利用する。ユーザはあるひとつのクラスタのプログラム(SPMD型)を実行すること

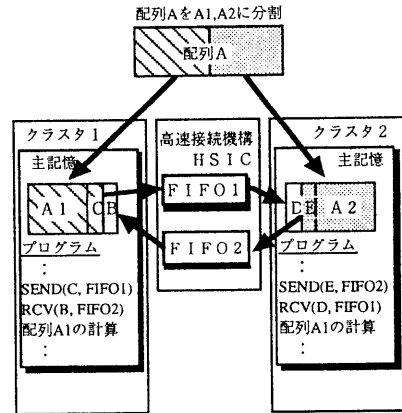
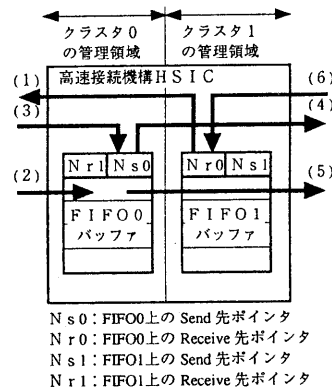


図3 S-3000クラスタシステムによる並列処理の概要
処理 $A_{new}(i)=f(A(i-1), A(i), A(i+1))$ の並列実行

Fig.3 Outline of parallel processing in the S-3000 Cluster system
example: $A_{new}(i)=f(A(i-1), A(i), A(i+1))$



Ns0: FIFO0上の Send 先ポイント
Nr0: FIFO0上の Receive 先ポイント
Ns1: FIFO1上の Send 先ポイント
Nr1: FIFO1上の Receive 先ポイント

データ転送手順(例: クラスタ0からクラスタ1へ)
<送信側(クラスタ0)>
(1) Nr0を読出し、チェック。[HSIC Read命令]
(2) FIFOバッファへデータ書き込み。[HSIC Write命令]
(3) Ns0の更新。[HSIC Read命令]
<受信側(クラスタ1)>
(4) Nr0を読出し、チェック。[HSIC Read命令]
(5) FIFOバッファからデータ読出し。[HSIC Read命令]
(6) Nr1の更新。[HSIC Write命令]
以上を繰り返し、実行。

注: クラスタ内のコピーで判断できる時は(1)(4)は不要

図4 クラスタ間転送プロトコルの概要
Fig.4 Protocol for data transfer between clusters.

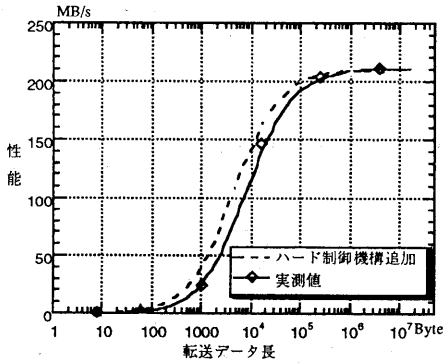


図5 クラスタ間データ転送基本性能 (パターン1)
Fig.5 Performance of data transfer (test pattern 1).

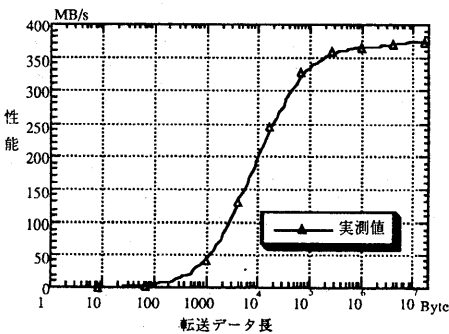


図6 クラスタ間データ転送基本性能 (パターン2)
Fig.6 Performance of data transfer (test pattern 2).

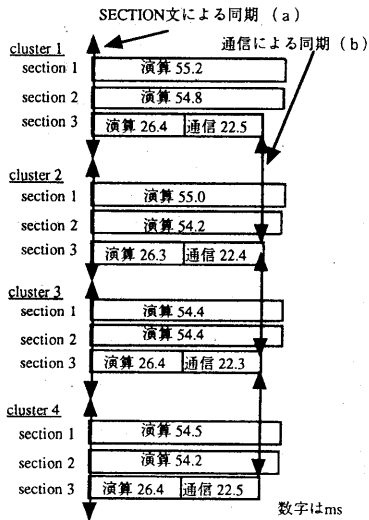


図7 クラスタ並列処理における同期の振舞い
Fig.7 Behavior of synchronization of cluster parallel execution.

により、PVMライブラリが自動的に他のクラスタのプログラムを起動する。それぞれのクラスタごとに実行されるプロセスが、自分がどのクラスタで実行されているか、あるいは何台のクラスタで実行されているかを、実行時にPVMライブラリから受け取ることができる。このため、使用クラスタ数をパラメータにすることが可能となり、クラスタ数を変化させても基本的にプログラムを一種類で作成することができる。

プログラム中で、クラスタ間のデータ転送バスを設定/解放処理を行なうが、この設定/解放処理については基本的なテンプレートが準備してあり、特に、特殊な使い方を必要としなければ、そのままそのテンプレートを用いて実行することができる。

3. 3. 2 ユーザプログラム

ユーザプログラムでは、2階層の並列処理のためのプログラミングを行なう。

クラスタ内は従来の自動ベクトル化および並列化コンパイラを適用する。今回の開発では、各クラスタにはクラスタ間データ転送のための専用の通信プロセッサは組込んでいない。そのためクラスタ間データ転送はクラスタ内のプロセッサのひとつが対応する。高速化のために重要な[8]データ転送処理と他プロセッサによるベクトル演算処理とのオーバラップを実現するために、基本的にコンパイラによるsection型の並列処理方式を利用する。ただし、基本はSPMD型のプログラミングスタイルであり、これをクラスタ内にも適用することにより、クラスタ内のk台のプロセッサに対して、1台のプロセッサが通信処理を含む演算処理をk-1台のプロセッサが同一の演算処理のみを実行するため、sectionごとのプログラムは高々2通りを作成すればよい。

クラスタ間については、メッセージパッシング型のデータ通信ライブラリが用意されており、明示的に通信命令(SEND, RECV)を記述する。送信命令(SEND)は、転送すべきデータの先頭アドレス(転送すべき配列の先頭位置)、データ数および転送バス名を指定する。受信命令(RECV)は、受信先のアドレス(受信すべき配列領域の先頭位置)、転送データ数および転送バス名を指定する。クラスタ間の同期処理は通信命令で行ない、データの受信により送信側クラスタとの順序関係を保障する。転送データの到着/未到着の制御はHSIC上のFIFOの管理で行なわれる。

4. 特性評価

S-3000クラスタシステム上でプログラムを実行し、特性評価を行なった。ここでは、基本通信性能、クラスタ並列処理での同期処理の挙動、およびLU分解法プログラムを用いたアプリケーションレベルでの評価結果を示す。

4. 1 転送基本性能評価

2つのクラスタ間のデータ通信能力を評価した。

パターン1: ひとつのデータ(D1)を2つのクラスタ間で送りあう処理である。SEND命令とRECEIVE命令を交互に実行する。そのため、HSIC上でのデータ転送による競合は発生しない。

パターン2: パターン1に対して、さらに同時にもうひとつのデータ(D2)を送りあう。HSIC上でつねに競合が発生する。これはまた、一方のクラスタから他方のクラスタに対して、ひとつのデータを分割してバイブライン的に転送することと性能的に等価になる。パターン1の結果を図5に示す。実線が実測値である。転送デー

タ長がほぼ数100KBのオーダーになると、200MB/sのデータ転送能力を発揮する。

破線のハード制御機構追加とは、3.2節に示したFIFOバッファの制御変数 $N_r N_s$ をHSIC上でなく、専用ハードウェアで実現した場合の予測値である。ハード制御機構の追加により、転送データ長が数KB~数100KBまでの場合は同じ効果を得るデータ長を約1/2に短縮することができる。

不完全LU分解による前処理などのアルゴリズムの高速化を行なう場合は、このような短い転送データ長が重要となる。

ボタン2に対する結果を図6に示す。データD1、D2合わせて、ボタン1と比較し、ほぼ2倍のデータスループットとなる。性能が完全に2倍にならないのは、HSIC上でのハードウェアによるデータ転送競合が発生しているためである。

4.2 クラスタ並列処理の挙動

S-3000クラスタシステムでは、クラスタ内の並列処理とクラスタ間の並列処理を同時に実行する。並列処理において重要な項目のひとつは同期処理である。クラスタ内はコンパイラによるsection型の並列処理を行ない、SECTION文によってクラスタ内プロセッサ間の同期をとる。一方、クラスタ間の並列処理では、HSIC上のFIFOにおけるデータの授受によって実質的に同期をとる。したがって、全体の同期はこれら2つの同期により、緩やかに行なわれる。図7に、この同期処理のようすを近接作用問題で代表的なRed/Black SORプログラムの1反復ステップでの状況をもとに説明する。各クラスタではsection型の並列処理を行ない、section3を担当するプロセッサが通信を行なう。データの依存関係を保つために、クラスタ間通信処理で得られたデータをもとに計算を行なう境界領域の計算は、section3を担当するプロセッサが確実に実行する。図からわかるように、section文による同期(a)は各クラスタ内のプロセッサ(各sectionを実行するプロセッサ)間で行なわれ、通信による同期(b)は各クラスタのsection3間で行なわれる。これらの同期は全体での同期でないため自由度が存在するが、各クラスタごとの処理の分割が均一であること、クラスタ内の(演算のみを実行する)プロセッサに対する処理の割付が均等であることが満たされると、図のように、整然と実行されているようすがわかる。図7は特別な反復ステップを取り出したのではなく、他の反復ステップも同様である。したがって、緩やかな同期機構で接続されているにもかかわらず、クラスタにおける2階層の並列機構におけるオーバーヘッドはほとんど無視することができることがわかる。

4.3 LU分解法を用いた性能評価

応用レベルの評価として、ここでは、LU分解法に基づく連立一次方程式の直接解法を扱う。この解法を評価することにより、文献[9]から多くの他のスーパーコンピュータや超並列計算機と比較検討することができる。

4.3.1 LU分解法のベクトル並列アルゴリズム

まず、プログラムをS-3000クラスタシステム向きに並列化する。ここでは、各クラスタへのデータの分割方法、2階層の並列化手法について示す。簡単のために、クラスタ数を3、クラスタ内のベクトルプロセッサ数を2とする。もちろん、他の構成でも一般性は失われない。

(1) クラスタ間のデータ分割

行列を各クラスタに対しインタリーブ的に分割する。これはLU分解法では反復処理ごとに行列内の計算対象となる領域が小さ

く変化するためである。インタリーブした各々のブロックの列数は、各クラスタ内のすべてのベクトルプロセッサが効率よく動作可能な幅を確保する。

(2) 並列化アルゴリズム

図8に並列化アルゴリズムのフローチャートを示す。インタリーブ的に分割したために、部分軸選択処理(pivoting)は各反復ステップ毎に担当するクラスタが入れ変わる。この処理自体はベクトル化はできるが並列化は困難である。ベクトル化についてもかなり効率が悪い。したがって、K+1回目の部分軸選択処理をK回目の行列消去処理(decomposition)と並列に実行する。

4.3.2 LU分解法の性能評価

実行結果を、表1、図9および図10に示す。まず、表1はプロセッサ数およびクラスタ数を変化させた各S-3000クラスタシステムでの実行結果を示す。

最大構成のモデル412(4クラスタ×3ベクトルプロセッサ、最大演算性能96GFLOPS)では78.2GFLOPSを実現した。これは81%の効率、半性能長N1/2は4880である。

比較のために、主記憶共有型スーパーコンピュータS-3800の値を表1に示す。同一のプロセッサ数であれば、S-3000クラスタシステムでも遜色ない性能を実現できる。

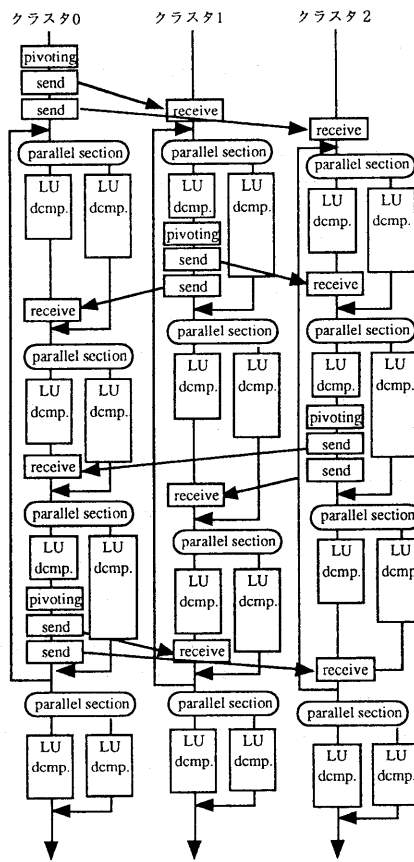


図8 LU分解法のアルゴリズムフローチャート
Fig.8 Algorithm flowchart for LU decomposition.

LU分解法のアルゴリズムは、1回の反復ステップにおいて、演算量が $O(n^2)$ であるのに対して、データ通信量が $O(n)$ であり少なくて済む。これが、効率80%以上を実現した理由であるが、さらに、効率を上げるためには、転送処理の強化として、通信専用プロセサの実現が重要となる。また、今回のS-3000クラスタシステムは転送能力を確保するために、リング型の結合を採用した。LU分解法のアルゴリズムで発生するデータ通信形態は放送であり、今回、開発したFIFO型のデータ転送路に、放送の機能を追加することも重要である。

5. おわりに

スーパーコンピュータの高速化のひとつのアプローチであるクラスタ型ベクトル並列スーパーコンピュータとして、S-3000クラスタシステムを開発した。重要課題であるクラスタ間のデータ転送手段として、排他制御が不要な低オーバーヘッドのクラスタ間高速通信機構を実現した。

実測により2階層の並列化によるクラスタ間の協調処理が整然と動作することを確認。また、LU分解法のプログラムを並列化し、最大性能の80%以上という高い効率を実現することができた。半性能長 $n \approx 1/2$ も比較的小さく適用範囲が広いといえる。

今後は、通信専用プロセサの実現、通信オーバーヘッドのさらなる低減により、適用応用範囲の拡大をはかっていきたい。

参考文献

- [1] T.Utumi, M.Ikeda and M.Takamura : Architecture of the VPP500 Parallel Supercomputer, Supercomputing'94 pp.478-487(1994.11).
- [2] Ishii, K., H.Abe, S.Kawabe and M.Hirai : An Overview of the HITACHI S-3800 Series Supercomputer, Proc. of Supercomputer'92, Springer-Verlag (1992).
- [3] K.Kitai, T.Isobe, Y.Tanaka, Y.Tamaki, M.Fukagawa, T.Tanaka and Y.Inagami : Parallel Processing Architecture for the Hitachi S-3800 Shared-Memory Vector Multiprocessor, ICS'93 (1993.7).
- [4] T.Sakakibara, K.Kitai, T. Isobe, S.Yazawa, T.Tanaka, Y. Inagami and Y. Tamaki : Scalable Parallel Memory Architecture with a Skew Scheme, ICS'93, pp.157-166 (1993.7).
- [5] T.Sakakibara, K.Kitai, T. Isobe, S.Yazawa, T.Tanaka, Y. Tamaki and Y. Inagami : An Interprocessor Memory Access Arbitrating Scheme for the S-3800 Vector Supercomputer, ISPAN(1994.12).
- [6] Gajski, D., Kuck, D., Lawrie, D. and Sameh, A. : Cedar -A Large Scale Multiprocessor, ICPP, pp.524-529(1983)
- [7] Hwang Kai : Advanced Computer Architecture: Parallelism, Scalability, Programmability, McGraw-Hill, Inc., pp.422-423 (1993).
- [8] 田中, 浜中, 村松 : 識別子を用いたデータ転送方式を基本とするMIMD型並列計算機アーキテクチャ, JSPP'89, pp.115-112 (1989.2).
- [9] Dongarra J.J. : Performance of Various Computers Using Standard Linear Equations Software, Comp. Science Dept., Univ. of Tennessee (1994.12).
- [10] M.Nakanishi, H.Ina and K.Miura: A High Performance Linear Equation Solver on VPP500 Parallel Super computer, Supercomputing'94, pp.803-810(1994.11).
- [11] 上村, 清水, 石畑, 堀江 : LINPACKベンチマークの並列ベクトル処理-並列計算機API1000用数値計算アクセラレータ

による実現一, JSPP'93, pp.185-192 (1994.5).

- [12] Hockney R.W. and C.R.Jesshope: Parallel Computers-2, Architectures, Programming and Algorithms, Adam Hilger Ltd., Bristol, UK and Philadelphia (1988).

表1 LU分解法プログラム実行結果
Table.1 Execution results for LU decomposition.

S-3000クラスタシステム						
モデル名	構成	実効性能	問題	N1/2	最大性能	効率
	VP数	GFLOPS	規模		GFLOPS	%
mode412	12(3×4)	78.2	31200	4880	96	81
mode408	8(2×4)	54.1	31200	3760	64	85
mode404	4(1×4)	27.2	31200	2680	32	85
mode309	9(3×3)	59.0	26940	3180	72	82
mode306	6(2×3)	40.9	27000	2400	48	85
mode303	3(1×3)	21.5	27000	1560	24	89
mode206	6(3×2)	40.6	21600	2160	48	85
mode204	4(2×2)	27.9	21600	1640	32	87
mode202	2(1×2)	14.5	21600	1100	16	90
S-3800TCMPシステム						
モデル名	構成	実効性能	問題	N1/2	最大性能	効率
	VP数	GFLOPS	規模		GFLOPS	%
mode480	4	28.4	15500	830	32	89
mode380	3	21.6	15680	760	24	90
mode280	2	14.6	15680	570	16	91
mode180	1	7.4	15680	470	8	92

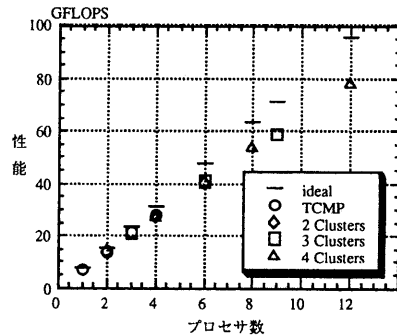


図9 LU分解法による性能評価
Fig.9 Performance for LU decomposition.

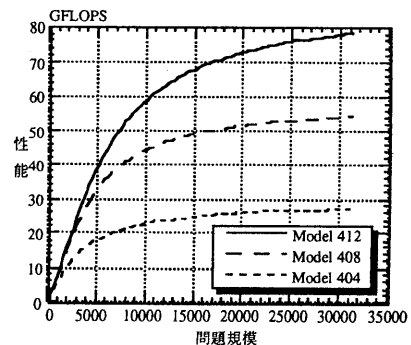


図10 4クラスタシステムでのLU分解法の性能
Fig.10 Performance for LU decomposition on the 4 cluster system.