

A-11

## 高速通信インタフェース DIMMnet-1 での

## 並列アプリケーションによる評価

## Evaluation with a parallel application for a fast communication interface that is the DIMMnet-1

三橋 彰 浩<sup>†</sup> 濱田 芳 博<sup>†</sup> 中 條 拓 伯<sup>†</sup>田 邊 昇<sup>††</sup> 工 藤 知 宏<sup>†††,\*</sup>AKIHIRO MITSUHASHI,<sup>†</sup> YOSHIHIRO HAMADA,<sup>†</sup>  
HIRONORI NAKAJO,<sup>††</sup> NOBORU TANABE<sup>††</sup> and TOMOHIRO KUDOH<sup>†††,\*</sup>

## 1. はじめに

近年、パーソナルコンピュータ (PC) やワークステーション (WS) の性能が著しく向上し、また価格も安価になってきているため、PC/WSを相互接続したネットワークで構成されるクラスタによって、スーパーコンピュータに匹敵する計算性能を持った計算機を安価に導入できるようになった。このため、数十台から数百台の PC/WS クラスタによる並列分散処理が広く行われるようになってきている。

我々は現在、PCIバスでなく DIMM スロットに接続するタイプの NIC である DIMMnet-1<sup>1)</sup> を用いたクラスタシステムを開発している。NIC を DIMM スロットに接続することでバンド幅の拡大とレイテンシの大幅な削減が可能となり、ボトルネックを解消することができる。また、クラスタシステムを構築し、実際に利用するためには Message Passing Interface (MPI) のようなメッセージパッシングライブラリが不可欠であり、本稿ではクラスタシステムへの MPI ライブラリ実装の足掛かりとして、RHINET2 スイッチによって接続された DIMMnet-1 を用いたバリア同期に要する時間を計測し、DIMMnet-1 の低レベルな関数を用いた簡単な並列アプリケーションによる評価を行う。

## 2. DIMMnet-1 の概要

## 2.1 ユーザレベル通信

ユーザレベル通信は、OS をバイパスしユーザプロセスが直接 NIC にアクセスして通信を行う方法で、ユーザ空間とカーネル空間の切替えのオーバーヘッドを削減することができる。ユーザレベル通信では、通信にかかわる情報がユーザプロセスから直接 NIC に渡された情報がどのプロセスからのもので、そのプロセスは NIC 上のリソースにアクセスする権利を持っているかどうかを NIC 側で判断する必要がある。NIC は、内部にユーザプロセスの仮想アドレスから物理アドレスへの変換を行うための Translation Look-aside Buffer (TLB) を持っている。この TLB に格納される属性によりプロセス間の保護を実現している。

## 2.2 Atomic on-the-fly (AOTF)

AOTF は、ホストプロセッサからの書き込みによってパケットを送信する機構である。ユーザプロセスにより AOTF 用の領域にデータの書き込みが行われると、あらかじめ設定しておいたヘッダシートを元にパケットが生成され、パケットとして送信される。このため、送信データが CPU

レジスタ上に存在していれば CPU がそのデータを所定のアドレスに書き込むという一命令だけでパケット送信が可能である。AOTF では一度に最大 8bytes のデータを送信することができ、低レイテンシな通信を実現している。

## 2.3 Block on-the-fly (BOTF)

BOTF はユーザプロセスがパケットを作成し、送信する機構である。ユーザプロセスは、送信をするための領域である Window 上にヘッダを含めたパケットのすべてを書き込んだ後、Window 上のキックアドレスに書き込みを行うことで DIMMnet-1 の NIC コントローラである Martini<sup>2)</sup> に送信開始を指示する。BOTF ではヘッダを除いて一度に最大 464bytes のデータが送信可能であり、高バンド幅な通信を実現している。

## 2.4 OTF 受信機構

OTF 受信機構とは、アドレス変換や DMA コントローラの起動をすること無しに、パケットヘッダの情報から所定の長さのデータ部を直接メモリに書き込む機構である。

DIMMnet-1 では AOTF 送信に限り、リモートアドレスを物理アドレスで登録することができ、受信時のリモートにおけるアドレス変換のオーバーヘッドを削除することが可能である。AOTF 送信に限って立てることができるヘッダ中のフラグを受信部が判定し、アドレス部と 1~8bytes のデータ部を書込みバッファに押し込んでいく。書き込みバッファは Martini 上のオンチップメモリである低遅延共有メモリ (LLCM) に、書き込めるタイミングで書き込む。

このように DIMMnet-1 では送信側の AOTF と受信側の OTFR が共同して極めて低遅延な通信を実現している。

## 2.5 評価環境

評価環境を表 1 に示す

測定環境	(a)	(b)
CPU/FSB	PentiumIII 850MHz/100MHz	
マザーボード	D6VAA	
OS	PentiumIII Linux(Kernel 2.4.2) SCORE 5.0/MPICH 1.2.3	
NIC	DIMMnet-1(Martini2nd) 光モジュール	Myrinet M2M-PC132C-10450
スイッチ	RHINET2/SW2 <sup>3)</sup> 通過遅延 240ns	Myrinet M2M-DUAL-SW8

環境 (b) は MPLBarrier 測定時のみ使用。

## 3. バリア同期

## 3.1 バリア同期の実装

DIMMnet-1 におけるバリア同期の実現方法<sup>4)</sup>としては、N-1 進木のツリーバリア方式を用いた。実装においては到着木、励起木共に同じ木構造を用いる。到着木においてリーフノードはルートノードへ 1bytes の AOTF 送信にて同期フラグを転送する。ルートノードはリーフノードからの同期フラグを 8bytes の変数にまとめ、全てのリーフノードに対するポーリングを 1 回の read で行う。励起に関しては

<sup>†</sup> 東京農工大学

Tokyo University of Agriculture and Technology

<sup>††</sup> (株) 東芝 研究開発センター

TOSHIBA Corporate Research & Development Center

<sup>†††</sup> 新情報処理開発機構

Real World Computing Partnership

\* 現在、独立行政法人 産業技術研究所 グリッド研究センター

Presently with National Institute of Advance Industrial Science and Technology

ルートノードはN-1回のAOTF送信を行い、リーフノードへ同期フラグを送信する。ここでSW2ではマルチキャストをサポートしているため、バリア同期マルチキャスト版では、このリーフノードへのAOTF送信を1回にまとめて実装を行った。

### 3.2 バリア同期評価

図1にDIMMnet-1を用いて実装したバリア同期と、Myrinetを使用したSCore/MPIでのバリア同期の測定結果を示す。DIMMnet-1を用いたバリア同期のレイテンシはMyrinetを用いたバリア同期に比べてほぼ10倍の性能を示した。また、マルチキャストを用いたDIMMnet-1のバリア同期は、用いない場合のバリア同期に比べ、レイテンシの増加が少ないことがわかる。それはノード数増加に伴う、ツリー筋起時のAOTF送信の増加をマルチキャストで一度に送信するからである。

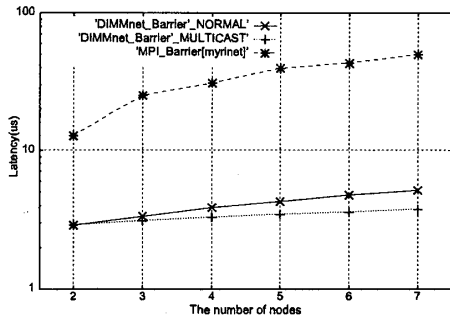


図1 バリア同期実行時間

## 4. 評価

### 4.1 評価に用いた並列アプリケーション

アプリケーションとしてRadixSortとLU分解を用いた。RadixSortはデータが6,000,000個で基数が10であり、LU分解に使用した問題サイズは1024×1024の実数密行列である。ノード間の同期には、マルチキャストを備えたDIMMnet-1バリア同期を用い、データ通信にはAOTFとBOTFを用いた。第2版のMartiniではBOTF送信時に連続ワードを立て続けに送信すると正常に通信できないという不具合があるため、送信間隔をあけることで回避している。第3版のMartiniではこの問題は改善される予定である。

### 4.2 演算時間・通信時間・同期待ち時間

表2に、各アプリケーションの4ノードでの実行時間の内訳、すなわち演算時間、通信時間、および同期待ち時間(バリア同期に要した時間+受信待ち時間)をそれぞれ示す。各実行時間において、台数効果を理想値に近づけるために、RadixSortでは通信時間の短縮。LU分解では主に同期待ち時間の短縮が考えられる。

	演算時間	通信時間	同期待ち時間
RadixSort	80.8%	18.3%	0.9%
LU分解	84.1%	0.4%	15.5%

### 4.3 評価結果

図2に、各アプリケーションにおいて1ノード時の実行時間を1とした場合の台数効果を示す。

RadixSort、LU分解共に、NノードでN倍という理想の性能向上に近い台数効果は得ることはできなかったが、台数を増やすごとに安定した台数効果を得ることができた。

### 4.4 考察

4ノードでのRadixSortの実行時間において、通信時間は18.3%を占める。現状のBOTF送信では送信間隔をあけてデータを送信しているため、送信バンド幅は11MB/s、uncashedで通常命令を用いたメモリへの読み込みバンド幅は26MB/sに抑えられている。しかし、第3版のMartiniを用いることで190Mbytesの通信を行うことが可能とな

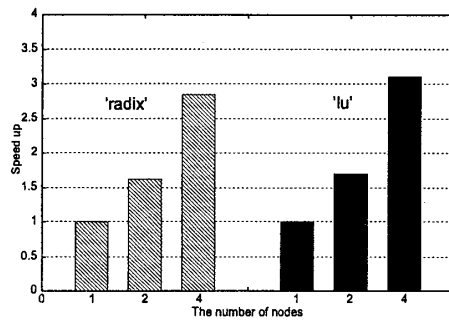


図2 台数効果

り、読み込みバンド幅はuncashedの属性で35Mbytesに向上する。このときの通信時間の割合は14%になり、台数効果は3倍に向上する。このとき、メモリへの読み込みにSIMD命令を用いた場合、読み込みバンド幅を74MB/sにすることが可能であり、通信時間の割合を6%に減少させ、台数効果は3.4倍に向上させることができる。また、Pentium4において使用できるキャッシュラインをフラッシュする、読み込みバンド幅466MB/sのCLFLUSH命令を使用した場合、通信時間の割合は0.01%に減少する。実行時間における通信の割合がSIMD命令の時点で6%まで減少しているため、台数効果は3.4倍と変わらないが、より大きな問題サイズを扱った場合の通信時間による性能の低下は解消することができる。

LU分解においては、実行時間における通信時間の割合は少ない。これはBOTF送信、通常読み込み命令という組合せでほぼ理想の通信が行われていることを示している。また、性能向上のためには同期待ち時間の減少が重要であり、アプリケーションの並列化について再考する必要がある。

## 5. 結論

マルチキャストを用いたDIMMnet-1のバリア同期に関しては、各ノードへの同期フラグの転送をAOTF送信を用いることによって、Myrinetを使用したSCore/MPIでのバリア同期に比べ10倍ほど高性能なバリア同期を実装することができた。LU分解などの計算と通信が同時に行われ、頻繁にバリア同期が用いられるようなアプリケーションの場合、DIMMnet-1で実装したバリア同期は性能向上に貢献すると考えられる。

4ノードシステムでの評価の結果、BOTFはLU分解のような通信と演算を別々に行うことができ、低粒度な通信が多発するアプリケーションに対しては非常に有効であるとわかった。また、RadixSortのような大きいデータをまとめて送信するアプリケーションに対しては、現在の不完全な実装状態にあるMartiniを用いていたのでは、DIMMスロットやBOTFが潜在的に持っている性能を生かすことは難しいといえる。

## 参考文献

- 1) 田邊, 山本, 工藤, “メモリスロットに搭載されるネットワークインタフェースMEMnet”, 情報処理学会計算機アーキテクチャ研究会, Vol99, No67, pp.73-78(1999)
- 2) 山本, 田邊, 西, 他, “高速性と柔軟性を併せ持つネットワークインタフェース用チップMartini”, 情報処理学会研究報告2000-ARC-140, pp.19-24(2000)
- 3) 西, 多昌, 西村, 山本, 工藤, 天野, “LASN用8Gbps/port 8x8 One-chipスイッチ:RHiNET-2/SW”, JSP2000 pp173-180, (May 2000)
- 4) 田邊, 濱田, 須田, 山本, 今城, 中條, 工藤, 天野, “DIMMスロット搭載型ネットワークインタフェースDIMMnet-1の通信性能評価”, 情報処理学会計算機アーキテクチャ研究会2001-ARC-145, pp.51-56(Dec. 2001)