

GigaE PM: Gigabit Ethernet を用いた 高速通信機構の設計と評価

住元 真 司[†] 堀 敦 史[†] 手塚 宏 史[†]
原田 浩[†] 高橋 俊 行[†] 石川 裕[†]

我々は、Gigabit Ethernet を用いたクラスタシステム上で並列アプリケーションを実行することを目的とした高速通信機構 (GigaE PM) を開発している。GigaE PM は、信頼性を確保した高バンド幅、低遅延通信を実現し、かつ、同時に TCP/IP などの既存の標準プロトコルもサポートしている。GigaE PM では、クラスタシステム向けプロトコル処理は NIC 上のファームウェアで、既存の通信プロトコル処理は従来どおりカーネルで処理する。プロトタイプシステムを Essential 社 NIC に実装し評価した結果、Pentium II 400 MHz PC 上で 48.3 μ s のラウンドトリップ遅延 (4 バイトメッセージ) と、56.7 MB のバンド幅 (1,468 バイトメッセージ) を実現している。さらに GigaE PM 上に MPI 通信ライブラリを構築し、GigaE PM と TCP/IP 上の MPI 通信ライブラリを用いたベンチマーク評価を行っている。NAS 並列ベンチマーク IS (クラス S) の評価では、GigaE PM の結果は TCP/IP の結果に比べ、1.8 倍の性能向上を得ている。

GigaE PM: The Design and Evaluation of High Performance Communication Facility Using a Gigabit Ethernet

SHINJI SUMIMOTO,[†] ATSUSHI HORI,[†] HIROSHI TEZUKA,[†]
HIROSHI HARADA,[†] TOSHIYUKI TAKAHASHI[†] and YUTAKA ISHIKAWA[†]

A high performance communication facility, called the *GigaE PM*, has been designed and implemented for parallel applications on cluster systems using a Gigabit Ethernet. The GigaE PM provides not only a reliable high bandwidth and low latency communication, but also supports existing network protocols such as TCP/IP. In the design of the GigaE PM, a reliable communication protocol for a parallel application is implemented on the firmware of network interface card while existing network protocols are handled by an operating system kernel. A prototype system has been implemented using an Essential Communications Gigabit Ethernet card, and evaluated its performance. The performance results show that a 48.3 μ s round trip time for a four byte message, and 56.7 MB/sec bandwidth for a 1,468 byte message have been achieved on Intel Pentium II 400 MHz PCs. We have implemented MPICH-PM on the top of GigaE PM, and evaluated the performance using NAS parallel benchmarks. The results show that the IS (class S) performance on GigaE PM is 1.8 times faster than that of on TCP/IP.

1. はじめに

多数の PC/WS をクラスタ結合した価格対性能比が高い並列システムが広まりつつある。このようなシステムを実現するには 2 つのアプローチがある。1 つは、ギガビットクラスのシステム間ネットワークを用いた専用の高速通信機構を利用する方法があり、既存の並列計算機と比較して安価に同等性能が得られる³⁾。この通信機構としてたとえば、Myrinet¹⁰⁾ を用

いた AM⁶⁾、FM⁷⁾、PM^{1),2)} などがある。

もう 1 つのアプローチは、MPI、TCP/IP などのコモディティソフトウェアと 100BASE-TX、ATM-LAN などのコモディティ LAN を利用する方法がある。現状のコモディティ LAN のバンド幅はたかだか 100 Mbps くらいであり、並列計算機のネットワークと比べてはるかに遅いが、既存の LAN 環境を利用し、簡単、かつ、安価に構築できる。

近年、Gigabit Ethernet が普及し始め、次世代のコモディティ LAN となりつつある。Gigabit Ethernet を利用すれば、高性能なクラスタシステムが LAN 環境上で安価に構築可能になる。しかし、Gigabit Ethernet

[†] 新情報処理開発機構つくば研究センター
Tsukuba Research Center, Real World Computing
Partnership

上で既存の通信プロトコルを用いて、並列アプリケーションに必要なバンド幅と通信遅延を実現することは難しい。たとえば、Gigabit Ethernet 上での TCP/IP 性能は物理層のバンド幅が 125 MB/s にもかかわらず、DEC 社製 Alpha プロセッサ (533 MHz) の計算機 (OS は Windows NT) と Packet Engine 社製の Gigabit Ethernet を用いた評価結果では、29.6 MB/s のバンド幅しか出ないと報告されている⁴⁾。TCP/IP 処理が性能上のボトルネックとなっており、Gigabit Ethernet の物理層性能を最大限に引き出せる通信機構が必要である。

我々は、LAN 環境で既存の通信プロトコルをサポートし、同時に、並列アプリケーションのための高バンド幅、低遅延通信を実現する GigaE PM 通信機構を設計実装している。GigaE PM を使えば、並列アプリケーションと既存の分散アプリケーションが共存できる環境を構築できる。

Gigabit Ethernet は、ハードウェアレベルでメッセージ転送を保証しないため、これを保証するプロトコル処理が必須である。このため、GigaE PM の設計では、プログラム可能なプロセッサを持つネットワークインタフェースカード (以降、NIC) を想定し、ホスト計算機 (以降、ホスト) と NIC 間のデータ交換のオーバーヘッドを最小にするため、GigaE PM の通信プロトコルを NIC 上で実現している。この結果、GigaE PM のプロトタイプでは信頼性のある通信を実現しながら、48.3 μ s のラウンドトリップ遅延 (4 バイトメッセージ) と、56.7 MB/s のバンド幅 (1,468 バイトメッセージ) 性能を実現している。

本論文では、コモディティ LAN である Gigabit Ethernet を用いた GigaE PM の設計と評価について述べる。2 章でクラスタ通信における要件について述べ、Gigabit Ethernet を用いたクラスタ通信はどうあるべきかについて述べる。3 章で GigaE PM の設計課題、4 章で GigaE PM の設計と実装について述べる。5 章では GigaE PM の遅延、バンド幅性能、NAS 並列ベンチマークについて TCP/IP と比較評価する。6 章に高速通信機構に関する関連研究をまとめ、7 章に結論を述べる。

2. クラスタ通信としての Gigabit Ethernet

本章では、並列計算向けの通信ネットワークの要件について述べ、Gigabit Ethernet 上でクラスタ通信をどう実現すべきかについて述べる。

2.1 並列計算向け通信ネットワークの要件

本論文で想定している並列計算を実行するアプリ

ケーションは、複数のプロセスで実現されており、それぞれのプロセスは 1 つのホスト上で実行され、他の多数のプロセスと通信する。

並列計算向け通信ネットワークの要件について以下に述べる。

- 信頼性のある通信：並列計算向け通信ネットワークは、メッセージ到着とメッセージの順序が必須である。たとえば、UDP/IP は、メッセージ転送を保証しないので、UDP/IP だけでは使えず、別途メッセージ転送を保証する機構が必要になる。
- 少ないメモリ使用量：並列計算向け通信ネットワークは、多数の計算機が結合されるため、通信ネットワーク維持に必要なメモリ使用量は最小であるべきである。たとえば、TCP/IP を用いた 1000 ノードのクラスタを考えた場合、他ノードとの通信のためノードあたり 999 ポートのコネクションを張る必要があり、かつ、それぞれのコネクションについて送受信バッファが必要となる。送受信バッファがそれぞれ 128 KB の場合、約 256 MB の送受信バッファが必要になり、ノードの搭載メモリ量で結合可能なノード数が制限される。
- 高い通信性能：上記 2 つの要件を満たしたうえで高い通信性能 (高バンド幅、低遅延) を実現する必要がある。
- 低コスト：並列計算向け通信ネットワークは、多数のノードを結合するため、低コストで実現されるべきである。

2.2 Gigabit Ethernet のクラスタ通信への適用

Gigabit Ethernet をクラスタ通信へ適用しようとした場合、Gigabit Ethernet はハードウェアレベルでメッセージ転送を保証しないため、これを保証するプロトコル処理が必須である。このため、GigaE PM の設計では、プログラム可能なプロセッサを持つ NIC を採用し、GigaE PM の通信プロトコルを NIC 上で実現している。NIC 上にプロトコルを実現することで、プロトコル処理に必要な制御フレームの送受信を出入力バス (たとえば、PCI バス) を経由する必要がなくなり、ホストと NIC 間のデータ交換のオーバーヘッドをより少なくできると考えたからである。

また、実装すべきプロトコルとして、メッセージ転送を保証したデータグラム転送プロトコルを採用することにより、メモリ使用量の少ない信頼性のある通信を実現することとした。

3. 設計課題

本章では、GigaE PM の設計課題について述べる。

- (1) シンプルな並列計算向き通信プロトコルの実現：前章の議論より、通信プロトコルは、多数のプロセスと通信するため、巨大な送受信バッファを必要としてはならない。また、通信プロトコルは、NIC 上の限られたハード資源内で実現可能なくらいシンプルでなければならない。
- (2) メッセージの到着保証と順序保証の実現：Gigabit Ethernet はメッセージ転送を保証していないので、通信プロトコルはメッセージの到着保証と順序保証をサポートする必要がある。
- (3) NIC とホスト間の情報交換コストの最小化：高い通信性能を実現するため、NIC とホスト間の情報交換コストを最小化する必要がある。
- (4) 並列計算向け通信プロトコルと TCP/IP など他プロトコルとの共存の実現：通信機構は TCP/IP など他のプロトコルが同時に利用される環境で利用されるため、通信機構は高性能な通信とともに既存の通信プロトコルも同時にサポートする必要がある。

4. 設計と実装

4.1 仮想ネットワークと API

GigaE PM 通信機構は並列計算アプリケーションに適用できるように設計されている。GigaE PM は、PM¹⁾ で採用されている仮想ネットワークをサポートするためのチャンネルを提供している。このチャンネルはメッセージ転送を保証したデータグラム通信を提供する。

並列アプリケーションのそれぞれのプロセスは、仮想ネットワーク上の同じ番号のチャンネルを排他的に利用する。チャンネルの数は NIC のハードウェア資源に依存する。現状の実装では、4 つのチャンネルをサポートしている。チャンネル数よりも多くのチャンネルを必要とする並列プロセスを実現するために、我々は PM¹⁾ のチャンネルの多重化を実現する SCore-D システムを開発している。

GigaE PM の API は PM¹⁾ の API に準拠しており、PM 上に開発された MPI を含む SCore システムとアプリケーションは、再コンパイルにより実行することができる。

4.2 信頼性のある通信

メッセージの到着と順序性を保証するために STOP and GO 方式による受信バッファフロー制御を備えた GO back N 方式を採用した。

4.2.1 GO back N プロトコル

GO back N プロトコルでは、送信側は i 番目から $i+N$ 番目までのデータを受信側からの ACK メッセージを待つことなく送信できる。受信側は i 番目のメッセージの受信時に送信側に ACK メッセージを送信する。送信側は i 番目の ACK メッセージを受信後に $i+1$ 番目から $i+1+N$ 番目のデータを受信側に送信できる。

もし、送信側が一定時間内に i 番目の ACK メッセージを受信しない場合は、送信側は i 番目とそれに続くメッセージを再度送信する。これは、 i 番目のデータが失われたことを意味する。送信側は i 番目から $i+N$ 番目までのメッセージのバッファ領域を確保する必要がある。TCP/IP と違い、受信側は N 個のメッセージ分バッファ領域を確保する必要がない。

もし、受信側がシーケンス番号 i より大きな番号のメッセージを受信した場合、受信データは廃棄され、送信側に LOSE メッセージを送信する。この点も TCP/IP の sliding window プロトコルと異なる。

もし、受信側のメッセージバッファがフルになった場合は、STOP メッセージが送信側に送信される。受信側は十分なバッファ領域が確保できた時点で GO メッセージを送信側に送信する。これは STOP and GO フロー制御と呼ばれている。

通信プロトコルの詳細は付録参照のこと。

4.3 ホストと NIC 間の情報交換

ホストと NIC 間の情報交換コストを定量的に見積もることは高バンド幅、低遅延の通信機構を設計するうえできわめて重要である。本節では、この情報交換のコストを解析し、このコストを基に通信機構を設計するうえで重要なディスクリプタの配置、アクセス方法、および、ホストと NIC 間の通知方法について議論する。

ホストと NIC 間の情報交換コストを Essential 社の PCI Gigabit Ethernet NIC を用い、Linux 2.1.131 を搭載した Pentium II 400 MHz の PC 上で解析した。表 1 に、ホストと NIC 間のデータ転送、割込み、システムコール (ioctl) にかかるコストを示す。これらのコストは以下のようにして測定した。

プロセッサによるホストとホスト間のデータ転送：

10 MB のユーザメモリ空間に対して、ホストプロセッサ上でデータ転送プログラムを作成して測定し、得られた測定値を基に式を作成した。

プロセッサによるホストと NIC 間のデータ転送：

ユーザメモリ空間に NIC メモリを mmap() し、ホストプロセッサでこの領域に対するデータ転送プ

表1 ホスト NIC 間データ転送, 割り込み, システムコールコスト
Table 1 Data transfer cost between the host and NIC memories and interrupt/system call processing.

	コスト (μs)
N ワード転送コスト (1 ワードは 4 バイト)	
プロセッサコピー: ホストからホスト	$0.03 \times N$
プロセッサコピー: ホストから NIC	$0.25 \times N$
プロセッサコピー: NIC からホスト	$0.57 \times N$
NIC DMA 機能: ホストから NIC	$2.30 + 0.03 \times N$
NIC DMA 機能: NIC からホスト	$2.00 + 0.03 \times N$
割り込み処理コスト	6.5
システムコールコスト	1.0

プログラムを作成して測定し, 得られた測定値を基に式を作成した.

NIC DMA によるホストと NIC 間のデータ転送: NIC 上でファームウェアを作成して測定を行った. 測定は DMA 起動のためのレジスタ設定と DMA 完了までの時間を NIC の持つタイマを用いて測定し, 得られた測定値を基に式を作成した.

割り込みのコスト:

NIC 上のファームウェア, および, ホストプロセッサ上で割り込み処理を行うデバイスドライバと測定プログラムを作成し, NIC ファームウェアがプロセッサへの割り込みレジスタに書き込んだ後, ホストプロセッサ上のデバイスドライバの割り込み処理関数の先頭にたどり着くまでの時間を PCI Gigabit Ethernet NIC の持つタイマを用いて測定した.

システムコールコスト:

デバイスドライバと測定用プログラムを作成し, このデバイスドライバへの 4 引数の `ioctl()` が完了する時間をプロセッサのタイマを用いて測定した.

GigaE PM における情報交換はメッセージディスクリプタを使って行われる. それぞれのエントリはメッセージバッファのアドレス, メッセージサイズとメッセージタイプに依存した送信側あるいは受信側の識別子から構成される.

4.3.1 ディスクリプタの配置

メッセージディスクリプタの置き場所は, ホストメモリ上か NIC メモリ上に置くことができる. このディスクリプタはホストと NIC の両方がアクセスするため, 両方のアクセスコストの総和が最小になるような場所に置くべきである. それぞれに必要なコストは, 表 1 のデータを基に計算すると以下となる.

- (1) ホストメモリ上: NIC DMA によるホストから NIC への 3 ワード転送コスト = $2.39 \mu\text{s}$

- (2) NIC メモリ上: プロセッサコピーによるホストから NIC への 3 ワード転送コスト = $0.75 \mu\text{s}$

Essential 社の NIC を用いた我々の実装では, アクセスコストを最小にするためメッセージディスクリプタは NIC 上に置いている.

4.3.2 ディスクリプタへのアクセス

もし, ユーザプロセスが, 直接送信メッセージディスクリプタを更新し, NIC に新しいディスクリプタが準備できたことを知らせるのであれば, 送信時においてシステムコールは必要ない. この場合安全な通信を保証するためには, 以下の 2 つを行う必要がある.

- (1) メッセージディスクリプタへのアクセスは, ユーザプロセスが適切な場所のみを参照可能とし, かつ, 他のユーザプロセスからのアクセスを禁止する.
- (2) ユーザプロセスが間違ったメッセージアドレスを書かないかどうか, NIC が送信メッセージディスクリプタをチェックする.

Essential 社の NIC では, NIC 上メモリへのアクセスに制御レジスタの更新が必要である. このため, もし, ユーザプロセスが制御レジスタの更新が可能な場合, ユーザプロセスは他の制御レジスタも更新できるため危険である.

GigaE PM では安全な通信を保証するため, ディスクリプタの更新はシステムコールを用いてカーネルで行う. なお, カーネルでディスクリプタ更新を行うことによるオーバーヘッドは, 表 1 より $1.0 \mu\text{s}$ と見積もることができる.

4.3.3 通 知

ホストから NIC, および, NIC からホストへの通知方式とそのコストを表 1 のデータを基に計算した結果を以下に示す.

ホストから NIC: 2 つの方式が選択可能

- (1) ホストが NIC メモリのフラグを更新: プロセッサコピーによる 1 ワード転送コスト = $0.25 \mu\text{s}$
- (2) NIC がホストメモリのフラグをポーリング: NIC DMA による 1 ワード転送コスト = $2.33 \mu\text{s}$

以上の計算結果より, ホストが NIC メモリに書く方式を採用する.

NIC からホスト: 3 つの方式が選択可能

- (1) NIC がホストメモリのフラグを更新: NIC DMA による 1 ワード転送コスト = $2.03 \mu\text{s}$
- (2) ホストが NIC メモリのフラグをポーリング: プロセッサコピーによる 1 ワード転送コスト =

0.57 μ s

- (3) NICがホストに割り込み: 割り込み+プロセッサコピーによる1ワード転送コスト= 7.05 μ s

以上の計算結果より, ホストがNICメモリのフラグをポーリングする方式を採用する. しかし, 安全な通信を保証するためポーリングをシステムコールで行うのは, さらに 1.0 μ s のコストがかかりオーバーヘッドが大きい.

そこで, Essential社のNICでは, ホストが制御レジスタを更新しないで直接アクセスできる少量のメモリ領域があるため, この領域にユーザプロセス用の受信メッセージディスクリプタを置くこととした. ユーザプロセスが, このディスクリプタをRead-Onlyで直接参照することにより, システムコールを使わなくてもディスクリプタを参照できる.

4.4 GigaE PM 実装の概要

本節では, GigaE PM 実装の概要について述べる. GigaE PM は, NIC上のファームウェア, デバイスドライバ, および, ユーザライブラリから構成される.

- APIは, $PM^{1,2)}$ と同じである. 現状, Zero-Copy機能は未実装.
- 送受信バッファは, ホストメモリ上にあらかじめ割り当てられ, pin-downされる. ユーザ領域へはmmap()でマップして利用される.
- NICのレジスタはRead-Onlyでユーザ空間にマップされる.
- GigaE PMのMTUは, 1,468バイトである.

4.5 送受信処理の概要

本節では, 図1を利用してGigaE PMの送受信処理の概要を説明し, ホストとNIC間の情報交換がどのように行われるかを示す.

送信側:

- (1) ユーザプログラムは`_PM_getsendBuf()`を発行してメッセージバッファを獲得する(図1中sc1). 引数はチャンネル番号とバッファサイズである. このバッファ上でメッセージを作成する.
- (2) ユーザプログラムが`_PM_sendmsg()`を実行すると, GigaE PMドライバが起動(図1中sc2)され, このドライバは送信メッセージの情報(メッセージアドレス, サイズ, 送信先識別子)をNICのディスクリプタに書く(図1中sc3).
- (3) NICはユーザ空間からNICメモリ領域にメッセージを転送する(図1中sd1).

受信側:

- (1) NICが新しいメッセージを受信した場合, ユーザプロセス上にマップされた受信メッセージ領

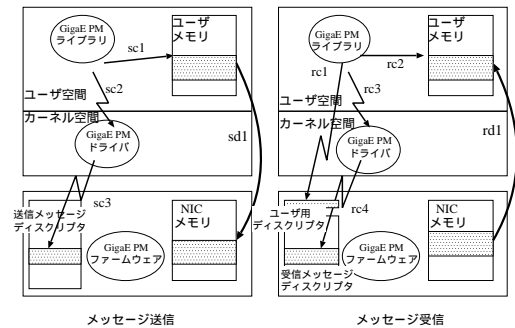


図1 GigaE PMでのディスクリプタ, ホスト, NIC
Fig.1 Descriptors, host, and NIC in the GigaE PM.

域に受信メッセージを転送する(図1中rd1). そして, NICは到着フラグを更新する.

- (2) ユーザプログラムが`_PM_receive()`を実行すると, `_PM_receive()`はシステムコールを発行することなくメッセージ到着フラグを参照する(図1中rc1). もし, このフラグが立っていたら受信メッセージディスクリプタより, メッセージバッファアドレス, サイズ, 送信者IDを獲得し, ユーザプログラムはメッセージ処理を行う(図1中rc2).
- (3) ユーザプログラムが`_PM_putreceiveBuf()`を実行すると, GigaE PMドライバが起動され(図1中rc3), このメッセージバッファが解放されたことを知らせる(図1中rc4).

4.6 GigaE PM ファームウェア

GigaE PM ファームウェアは, 他の多くのPCI Ethernet NICがハードウェアの機能として備えているPCI DMAとネットワークインタフェースのDMAを制御し, メッセージ転送を行う. クラスタ通信についてはGigaE PM通信プロトコル処理を行う. ホストプロセッサとの情報交換は, NICメモリを介して行う.

他の多くのNICの持つ機能を用いて実装しているため, NIC上にプログラム可能なプロセッサとホストプロセッサがアクセス可能なメモリを持つNICであれば, GigaE PMを実装することが可能である.

4.7 GigaE PMと他プロトコル

GigaE PMの通信プロトコルはNICで処理されるが, 他の通信プロトコルは, 他のEthernet NICと同様ホストで処理される. このプロトコルの切替えはEthernetフレームのタイプで行っている. GigaE PMの場合には特殊なタイプを用い, タイプがGigaE PMのものでないパケットが届いた場合, そのパケットはそのままホストに転送され, ホストのデバイスドライバの割り込みハンドラが起動される.

表 2 測定環境

Table 2 Machine environment.

ハードウェア	Pentium II 400 MHz, 440 BX チップセット, 256 MB SDRAM メモリ
NIC	Essential 社 PCI Gigabit Ethernet NIC model EC-440-SF (33 MHz clock, PCI DMA, MAC SEEQ 8100 MAC), 1 MB MEM
スイッチ	Extreme 社 Summit 1
ホスト OS	Redhat 5.1 Linux (2.1.131 kernel)

5. 評価

この章では、基本性能として遅延とバンド幅の測定結果、および、NAS 並列ベンチマークの性能評価を TCP/IP と比較する。表 2 に測定環境を示す。

GigaE PM のバンド幅と遅延の測定では、次の 2 つの場合について評価を行った。なお、TCP/IP の結果はすべてスイッチを用いて接続した場合である。

- (1) 2 台のホストを Extreme 社スイッチを用いて接続した場合
- (2) 2 台のホストをスイッチなしで直接接続した場合

5.1 バンド幅

図 2 に GigaE PM, GigaE PM 上の TCP/IP, Essential 社ファームウェア上での TCP/IP のバンド幅の測定結果を示す。GigaE PM は、1,468 バイトメッセージ時に 56.7 MB/s のバンド幅性能である。これに対して TCP/IP の場合は、Essential 社製ファームウェア利用、1,280 バイトメッセージのときに 33.0 MB/s のバンド幅性能である。

GigaE PM 上の TCP/IP と Essential 社ファームウェア上の TCP/IP 性能の比較では、4 バイトから 256 バイトメッセージまでは、GigaE PM を用いた場合が性能が良く、512 バイトから 1,280 バイトメッセージまでは、Essential 社ファームウェアを用いた場合が性能が良い。

5.2 遅延

図 3 にラウンドトリップ遅延の結果を示す。GigaE PM は、4 バイトメッセージ時にスイッチを介さない場合は 48.3 μ s のラウンドトリップ遅延であり、Summit 1 スイッチ経由の場合に 58.3 μ s となる。これは、スイッチの遅延が単方向で 5.0 μ s であることを意味する。

一方、TCP/IP のラウンドトリップ遅延は Essential 社製ファームウェア利用、4 バイトメッセージ時に 414.0 μ s であるのに対して、GigaE PM を用いた場合には、131.4 μ s であった。GigaE PM は、高速通信プロトコルとともに、既存のネットワークプロトコルにおいても高い性能を実現している。

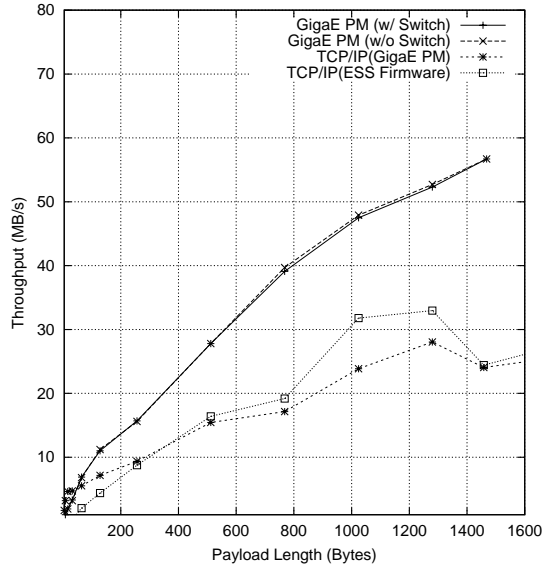


図 2 バンド幅性能
Fig. 2 Bandwidth.

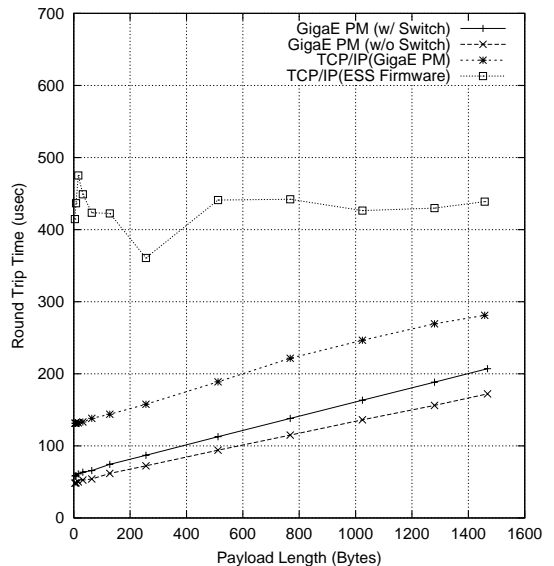


図 3 ラウンドトリップ遅延
Fig. 3 Round trip time.

Essential 社製ファームウェアを利用した場合のラウンドトリップ遅延がメッセージサイズに関係なく 400 μ s 付近の値を示しているのは、ホストプロセッサへの割り込み負荷を抑えて通信バンド幅性能を向上させるため、割り込みを NIC の持つタイマにより一定間隔遅らせているからである。図 2 における GigaE PM 上の TCP/IP よりも Essential 社製ファームウェアを用いた場合の方が通信バンド幅性能が良いのはこの理

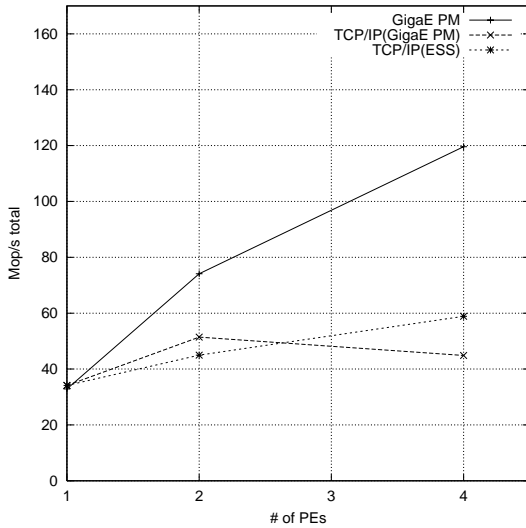


図4 CG クラス S
Fig. 4 CG Class S.

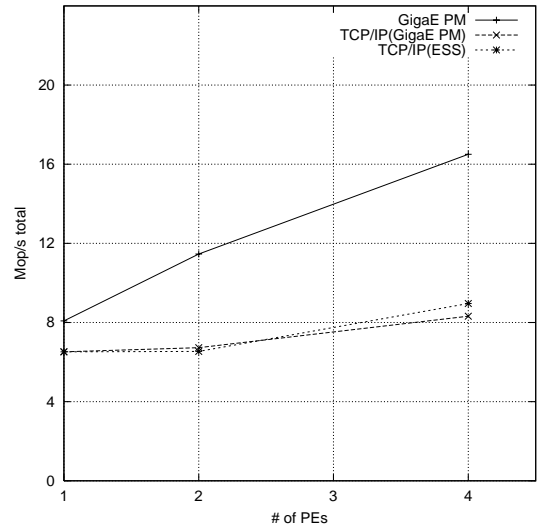


図5 IS クラス S
Fig. 5 IS Class S.

由による。

5.3 NAS 並列ベンチマーク

GigaE PM 上に MPICH-PM(MPICH 1.0.11 ベース)を移植し, NAS 並列ベンチマーク 2.3⁵⁾の性能を測定した. NAS 並列ベンチマーク 2.3 には 8 つのベンチマークがあるが, この中で低い通信遅延と高いバンド幅が必要な CG (Conjugate Gradient), IS (Integer Sort) の性能を測定した. 結果を図 4, 図 5 に示す. クラスは S である. 比較のため, GigaE PM 上の TCP/IP, Essential 社製ファームウェアを用いた結果を載せる. TCP/IP に用いた MPI は MPICH 1.0.11 の ch_p4 である.

結果より, GigaE PM では 4 ノード時に 3.63 倍 (CG), 2.0 倍 (IS) の性能向上が見られるのに対し, TCP/IP では, 4 ノード時に最大, 1.72 倍 (CG), 1.37 倍 (IS) の性能向上であった. 4 ノード時の GigaE PM の結果は, TCP/IP の結果より, 2.0 倍 (CG), 1.84 倍 (IS) 高速であり, GigaE PM の通信性能の高さを反映した結果である.

5.4 メモリ使用量比較

表 3 に N ノード, および, 1000 ノードのクラスタでの TCP/IP と GigaE PM における送受信バッファと制御に必要なメモリ使用量を示す. TCP/IP の送受信バッファのメモリ量は Linux 2.1.131 のデフォルト値, 制御構造体は tcp_opt 構造体のサイズである.

表 3 の結果より, GigaE PM の送受信バッファに必要なメモリ量はノード数に依存しないこと, 制御構造体も 1000 ノードのクラスタにおいても 23 KB と,

表 3 N ノードクラスタにおけるメモリ使用量の比較

Table 3 Comparison of memory usage in an N node cluster.

	TCP/IP	GigaE PM
N ノードクラスタ		
送信バッファ (KB)	$128 \times (N - 1)$	120
受信バッファ (KB)	$128 \times (N - 1)$	128
制御構造体 (KB)	$0.274 \times (N - 1)$	$0.023 \times N$
1000 ノードクラスタ		
送信バッファ (KB)	127,987	120
受信バッファ (KB)	127,987	128
制御構造体 (KB)	274	23

NIC 上メモリに収まることが分かる. 反面, コネクションベースの TCP/IP はノード数に比例した送受信バッファが必要になるため, ノードの搭載メモリ量で結合可能なノード数が制限される.

6. 関連研究

数多くの高性能通信機構が開発されている. AM⁶⁾, AM-II⁸⁾, FM⁷⁾, VMMC-2⁹⁾, と PM¹⁾ は Myrinet¹⁰⁾ をベースとしている. Myrinet はギガビットクラスのネットワークでハードウェアレベルでメッセージ転送を保証している. それゆえにネットワーク上でメッセージが失われることは考える必要はなく, 受信バッファのフロー制御を行えばよい.

システムコールのオーバーヘッドを減らすために, PM と U-Net¹¹⁾ ではシステムコールを使っていない. PM は Myrinet 上で $15 \mu\text{s}$ のラウンドトリップ遅延を実現しているため, システムコールのオーバーヘッドは大き

い。しかしながら，GigaE PM の場合，システムコールのオーバーヘッドは， $48.3 \mu\text{s}$ のラウンドトリップ遅延の中で $2.0 \mu\text{s}$ (4.1%) であるため，それほど大きな影響はない。

GigaE PM と同様，AM-II⁸⁾ も，信頼性のある通信を NIC 上で実現し，かつ，安全な通信をサポートしている。AM-II はシステムコールを使わずに安全な通信を実現しているが，GigaE PM はシステムコールを使い，安全，かつ，低遅延，高バンド幅通信を実現している。

VIA¹²⁾ (Virtual Interface Architecture) はマイクロソフトの Windows 上のギガビットクラスのネットワークに広く実装されている。VIA はその上に socket や MPI のような通信ライブラリが実装されることを想定している。VIA はコネクションベースの通信をサポートしており，信頼性のある通信は VIA specification Version 1.0 ではオプションである。

もし，コネクションベースの通信を Gigabit Ethernet のような信頼性のないネットワーク上で実現する場合，通信プロトコル処理は，より大きな通信バッファを必要とする。よって，我々は，コネクションベースで並列アプリケーション（特にデータ並列アプリケーション）向き的高速通信のサポートは難しいと考える。GigaE PM のプロトコルはシンプルで他のコネクションベースの通信機構と比べ大きなメッセージバッファを必要としないため，NIC 上で実現でき，かつ Gigabit Ethernet 上で信頼性のある通信を実現することができる。

7. おわりに

本論文では，Gigabit Ethernet 上で並列計算向き的高速通信を実現するためには，NIC とホスト間の情報交換を考慮した設計が重要であることを示した。高速通信機構 GigaE PM はこれを満たすよう設計されている。

プロトタイプを Essential 社の Gigabit Ethernet NIC 上に実装，評価した結果，4 バイトメッセージ時に $48.3 \mu\text{s}$ のラウンドトリップタイム，1,468 バイトメッセージ時に 56.7MB/s のバンド幅が得られている。同時に，GigaE PM は，性能を落とすことなく，TCP/IP プロトコルもサポートしている。また，MPICH-PM 上の NAS 並列ベンチマーク（クラス S）で評価した結果，4 ノード時の GigaE PM の結果は，TCP/IP の結果より，2.0 倍（CG），1.84 倍（IS）高速である。

本論文の寄与するところは，Gigabit Ethernet を用

いた高速通信機構の設計である。GigaE PM は，他の多くの NIC が持つハードウェア機能を用いて実装しているため，NIC 上にプログラム可能なプロセッサとホストプロセッサがアクセス可能なメモリを持つ NIC であれば実装することができる。

GigaE PM は，信頼性のある，高バンド幅，低遅延通信を提供するだけでなく，同時に TCP/IP などの既存の通信プロトコルもサポートしている。この機能を用いることにより，並列アプリケーションと分散アプリケーションが共存した高性能なクラスタシステムが構築できる。

参 考 文 献

- 1) Tezuka, H., Hori, A., Ishikawa, Y. and Sato, M.: PM: An Operating System Coordinated High Performance Communication Library, Sloot, P. and Hertzberger, B. (Eds.), *High-Performance Computing and Networking*, Vol.1225, Lecture Notes in Computer Science, pp.708–717, Springer-Verlag (Apr. 1997).
- 2) Tezuka, H., O'Carroll, F., Hori, A. and Ishikawa, Y.: Pin-down Cache: A Virtual Memory Management Technique for Zero-copy Communication, *IPPS/SPDP'98*, pp.308–314, IEEE (Apr. 1998).
- 3) <http://www.rwcp.or.jp/lab/pdslab/benchmarks/>.
- 4) <http://www.packetengines.com/products/performance/gnicntperf.htm>.
- 5) <http://www.nas.nasa.gov/Software/NPB/>
- 6) http://now.cs.berkeley.edu/AM/lam_release.html.
- 7) Pakin, S., Lauria, M. and Chein, A.: High Performance Messaging on Workstations: Illinois Fast Messages (FM) for Myrinet, *Proc. Supercomputing '95*, San Diego, California (1995).
- 8) Chun, B.N., Mainwaring, A.M. and Culler, D.E.: Virtual Network Transport Protocols for Myrinet, *Hot Interconnect'97* (Aug. 1997).
- 9) Dubnicki, C., Bilas, A., Chen, Y., Damianakis, S. and Li, K.: VMMC-2: Efficient Support for Reliable, Connection-Oriented Communication, *Hot Interconnect'97* (Aug. 1997).
- 10) Boden, N.J., Cohen, D., Felderman, R.E., Kulawik, A.E., Seitz, C.L., Seizovic, J.N. and Su, W.-K.: Myrinet – A Gigabit-per-Second Local-Area Network, *IEEE MICRO*, Vol.15, No.1, pp.29–36 (Feb. 1995).
- 11) Basu, A., Buch, V., Vogels, W. and von Eicken, T.: U-Net: A User-Level Network Interface for Parallel and Distributed Computing, *Proc. 3rd International Symposium on High*

Performance Computer Architecture HPCA
(Feb. 1997).

12) <http://www.viarch.org/>.

付 録

GigaE PM プロトコル

ここでは GigaE PM の通信プロトコルの詳細について述べる。それぞれのチャンネルに送信バッファと受信バッファがあり、すべての送信メッセージは送信バッファ、すべての受信メッセージは受信バッファに格納される。バッファは、TCP/IP と違い peer to peer 通信単位ではなくチャンネル単位に割り当てられる。

これ以降利用する用語を定義する。メッセージは、 Msg (送信側 ID , 受信側 ID , メッセージ順序番号) で表現する。 N は送信側が ACK メッセージを待つことなく非同期に送信できるメッセージ数を示す。 T はタイムアウト時間、 $SBuf(r, i)$ は i 番目の送信メッセージを格納する送信バッファ、 $STime(r, i)$ は i 番目のメッセージが受信側 r に送信された時刻、 $MsgIdSent(r)$ は受信側 r に送信されたメッセージ順序番号の最大値、 $MsgIdAcked(r)$ は受信側 r から受信した ACK メッセージの順序番号の最大値、 $MsgIdRecv(s)$ は受信側が送信側 s から受信したメッセージ順序番号の最大値とする。初期化時点において $MsgIdSent(r)$, $MsgIdAcked(r)$ と $MsgIdRecv(s)$ は 0 である。

送信側と受信側における GigaE PM プロトコルを以下に示す：

送信側において：

S1 送信側 s は、次の条件を満たすメッセージ $Msg(s, r, i)$ を受信側 r に送信する。送信時には、以下の処理を行う。

$$\forall Msg(s, r, j) \text{ where } MsgIdSent(r) < j < \\ MsgIdAcked(r) + N.$$

- (1) $Msg(s, r, i)$ を受信側に送信する。
- (2) 送信メッセージ $Msg(s, r, i)$ のバッファ $SBuf(r, i)$ を作成し、維持する。
- (3) 現在の時刻を $STime(r, i)$ に格納する。
- (4) $MsgIdSent(r) \leftarrow MsgIdSent(r) + 1$.

S2 現在時刻と $STime(s, r, i)$ の差が T より大きい場合、以下を実行する。

- (1) $MsgIdSent(r) \leftarrow i$.
- (2) S1 を実行する。

- もし次の条件を満たすLOSEメッセージ $LOSE(r, k)$ を受信した場合、以下を実行する。

$$LOSE(r, k) \text{ where } MsgIdAcked(r) < k < \\ MsgIdAcked(r) + N.$$

(1) 次の条件を満たす送信メッセージバッファ $SBuf(r, i)$ を解放する。

$$\forall SBuf(r, i) \text{ where } MsgIdAcked(r) < k.$$

(2) $MsgIdAcked(r) \leftarrow k$.

(3) S1 を実行する。

- もし次の条件を満たすSTOPメッセージ $STOP(r, k)$ を受信した場合、以下を実行する。

$$STOP(r, k) \text{ where } MsgIdAcked(r) < k < \\ MsgIdAcked(r) + N.$$

(1) 次の条件を満たす送信メッセージバッファ $SBuf(r, i)$ を解放する。

$$\forall SBuf(r, i) \text{ where } MsgIdAcked(r) < k.$$

(2) $MsgIdAcked(r) \leftarrow k$.

(3) 送信を停止する。

- もし次の条件を満たすGOメッセージ $GO(r, k)$ を受信した場合、S1 を実行する。

$$GO(r, k) \text{ where } MsgIdAcked(r) \leq k < \\ MsgIdAcked(r) + N.$$

- もし次の条件を満たすACKメッセージ $ACK(r, k)$ を受信した場合、以下を実行する。

$$ACK(r, k) \text{ where } MsgIdAcked(r) \leq k < \\ MsgIdAcked(r) + N.$$

(1) 次の条件を満たす送信メッセージバッファ $SBuf(r, i)$ を解放する。

$$\forall SBuf(r, i) \text{ where } MsgIdAcked(r) \\ \leq k.$$

(2) $MsgIdAcked(r) \leftarrow k$.

(3) S1 を実行する。

受信側において：

- 受信側 r がメッセージ $Msg(s, r, i)$ を受信したとき、

- $MsgIdRecv(s) + 1 = i$ を満たし、かつ、受信側 r の受信バッファがフルでない場合、以下の処理を行う。

(1) ACK メッセージ $ACK(r, i)$ を送信側 s に送信する。

(2) 受信メッセージをホストに転送する。

(3) $MsgIdRecv(s)$

$$\leftarrow MsgIdRecv(s) + 1.$$

- $MsgIdRecv(s) + 1 = i$ を満たし、かつ、受信側 r の受信バッファがフルの場合、以下の処理を行う。

(1) STOP メッセージ $STOP(r, i)$ を送信側 s に送信する。

(2) これ以降、到着するすべての受信メッセージを廃棄する。

- $MsgIdRecv(s) + 1 < i$ を満たす場合、以下
を処理を行う。これは、メッセージが失われ
たことを意味する。

(1) LOSE メッセージ

$LOSE(r, MsgIdRecv(s))$

を送信側 s に送信する。

- もし、ホストの受信バッファに十分な空きがあり、
かつ、STOP メッセージが送信済みの場合、以下
を実行する。

(1) GO メッセージ $GO(r, MsgIdRecv(s)+1)$
を送信側 s に送信する。

(2) 受信側 r は、受信処理を再開する。

(平成 11 年 8 月 31 日受付)

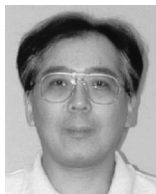
(平成 11 年 12 月 2 日採録)



住元 真司 (正会員)

1986 年同志社大学工学部電子工学
科卒業。同年(株)富士通入社(株)
富士通研究所にて並列オペレー
ティングシステム、並列分散システムソ
フトウェアの研究開発に従事。1997

年より新情報処理開発機構に出向。コモディティネッ
トワークを用いた高速通信機構の研究開発に従事。並
列分散システムのアーキテクチャ、システムソフト
ウェア等に興味を持つ。



堀 敦史 (正会員)

1979 年早稲田大学理工学部電気
工学科卒業。1981 年同大学院理工
学研究科計測制御工学専攻修士課程
修了。同年(株)三菱総合研究所入
社。1992 年より技術研究組合新情

報処理開発機構に出向。JSPF'98 最優秀論文賞受賞。
並列オペレーティングシステムの研究に従事。並列プ
ログラミング言語、並列アーキテクチャ等に興味を持
つ。工学博士(東京大学工学部)。



手塚 宏史 (正会員)

1980 年東京大学教養課程中退。
1981 年(株)生活構造研究所入社。
1985 年ソニー(株)入社。1988 年
(株)ソニーコンピュータサイエンス
研究所入社。1990 年ソニー(株)入
社。1993 年北陸先端科学技術大学院大学研究生。1995
年より技術研究組合新情報処理開発機構研究員。現在
に至る。オペレーティングシステム、リアルタイム処
理、マルチメディア処理等に興味を持つ。日本ソフト
ウェア科学会会員。



原田 浩 (正会員)

1988 年東京理科大学理学部物理
学科卒業。同年(株)ソフトウェア・
リサーチ・アソシエイツ入社。1997
年より技術研究組合新情報処理開発
機構研究員。現在に至る。オペレ
ーティングシステム、並列・分散システム等に興味を持
つ。ACM 会員。



高橋 俊行 (正会員)

1993 年東京理科大学理工学部情
報科学科卒業。1995 年同大学院修
士課程修了。1995 年~1998 年東京
大学理学系研究科情報科学科博士課
程。1998 年より新情報処理開発機
構研究員。現在に至る。プログラミング言語における
メタレベルアーキテクチャと並列計算ソフトウェア技
術に興味を持つ。理学修士。



石川 裕 (正会員)

1987 年慶應義塾大学院理工学
科電気工学科博士課程修了。同年電
子技術総合研究所入所。1988~1989
年カーネギー・メロン大学客員研究
員。1990 年日本ソフトウェア科学会
高橋奨励賞を受賞。1993 年から新情報処理開発機構
に出向。並列・分散システム、適応可能並列プログラ
ミング言語/環境/処理系、リアルタイム処理等に興味
を持つ。日本ソフトウェア科学会、ACM、IEEE 各会
員。工学博士。