

# メモリスロット搭載型ネットワークインタフェース DIMMnet-1 における細粒度通信機構

田邊 昇      山本 淳二      工藤 知宏

新情報処理開発機構つくば研究センタ

筆者らはネットワークインタフェース(NIC)をメモリスロットに搭載するというアプローチ(MEMOnetアーキテクチャ)に基づき、高性能なPCクラス用NICのプロトタイプDIMMnet-1を開発中である。DIMMスロットを有する安価なパソコン上で、片方向あたり8Gbpsという通信バンド幅と、8バイトデータ送信時にソフトウェアオーバーヘッド1命令サイクル、ハードウェア遅延45nsという極めて高速なメッセージ生成が可能である。本報告では従来型NICの問題点をバンド幅と遅延時間の観点から整理し、DIMMnet-1においてはどのようにこれらを克服しているかを解説する。次にDIMMnet-1における細粒度通信アーキテクチャとその通信遅延性能に関する予測を述べる。

## The Finegrain Communication Facilities on DIMMnet-1 Network Interface Inserted in a Memory Slot

Noboru Tanabe      Junji Yamamoto      Tomohiro Kudoh

Real World Computing Partnership

The finegrain communication architecture on DIMMnet-1 network interface based on MEMOnet is described. MEMOnet is the network interface (NIC) architecture that the NIC inserted in a memory slot. The DIMMnet-1's link bandwidth is 8Gbps per direction. The overhead to generate a message is only 1 CPU cycle and the hardware delay is only 45ns. The excellent performance is available for cheap personal computers with DIMM slots. We also present the many problems of the old fashioned NICs. These problems become obvious if the 8Gbps optical link is introduced to the NIC with the same old fashioned architecture. Most of all these problems are newly solved on the DIMMnet-1.

### 1 はじめに

近年、高性能PCを多数用いて並列処理を行なういわゆるクラスタコンピューティングが注目されている。これらには、Myrinet[1]に代表されるシステムエリアの低遅延ネットワークを用いたものと、100BaseTやGigabit Ethernetなどの汎用の高速ネットワークを用いたもの[2]がある。これらはいずれもPCIバスに接続されるタイプのネットワークインタフェース(NIC)を用いている。

筆者らはRHINET-1[3]という1Gbpsの光インタコネクションを用いた、PCIバスに接続されるタイプのNICを開発し、RHINET-2[4]という8Gbpsの光インタコネクション[5]を用いた、64bit66MHzPCIバスに接続されるタイプのNICを開発中である。

しかし、光インタコネクションの持つ大きなバンド幅を有効に活用するには従来のPCIバスではバンド幅および遅延ともに力不足である。一方、PCIX[6]、InfiniBand[7]といった次世代のサーバー向け入出力の規格が提案されつつある。よってサーバー機に関してはこのボトルネックは改善されつつあるが、最も価格性能比においてメリットのあるエンドユーザー用PC向けに、これらの入出力バスが普及するかどうか不透明である。

このような問題意識にたち我々は、従来のようにPCIバス等の入出力バスではなく、メモリスロットに搭載されるタイプのネットワークインタフェースを検討している。このようなクラスのネットワークインタフェースをMEMOnet[8]と名付けた。MEMOnetは安価なPC上で、PCIバスのバンド幅や遅延時間の限界を超越したネットワークインタフェースを実現可能と思われる。

そのようなクラスに属するNICの研究としては、1990年代中期、MINI[9]という複数のSIMMソケットに搭載されるエクステンダを介したボード上に画像用のデュアルポートDRAMを用いたATM用NICがあった。一方、DIMMを拡張した新しいメモリモジュールが近年提案され、規格化が進められた。それが日本の工業規格であるEIAJや国際

規格であるJEDECで認可されたPEMM(Processor Enhanced Memory Module)[10]である。我々はこのPEMM規格を基盤に、MEMOnetのプロトタイプとしてDIMMスロットに搭載されるDIMMnet-1を開発中である。

光インタコネクションの有効活用から出発した経緯から、筆者らMEMOnetを提唱した際に、バンド幅に関する評価[8]を示した。しかし、MEMOnetの利点はバンド幅向上はもとより、通信遅延の低下に貢献すると考えている。MEMOnetにより、極めて低遅延かつ高バンド幅のネットワークを実現できれば、Shasta[11]のようなアプリケーションの逐次バイナリコードからプライベートでないロード・ストアに対して一貫性管理コードを挿入するような細粒度で高頻度なメッセージ通信を発生するアプローチを、より効果的にする可能性を秘めている。

本報告ではまず、従来型NICの問題点をバンド幅と遅延時間の観点から整理し、DIMMnet-1のアーキテクチャにおいてはどのようにこれらを克服しているかを解説する。次にDIMMnet-1の遅延時間に関する改善に焦点をあて、細粒度通信機構の構成と動作を紹介し、通信遅延時間がどの程度になるかという見通しを報告する。

### 2 従来型NICの問題点

本章では、8Gbpsの光通信リンクを活用するDIMMnet-1で解決する課題を明らかにすべく、従来のNICの問題点を明らかにする。それらは大別して、PCIバス利用に起因するもの(P:PCI)、資源の不足に起因するもの(C:Capacity)、アーキテクチャ上の工夫の不足によるもの(A:Architecture)の3つの要因があり、それらが絡み合っているものもある。以下、各節のタイトルに付した記号(P,C,A)は関連する上記の要因を示す。

#### 2.1 バンド幅の不足

##### 2.1.1 ピークバンド幅不足(P)

光モジュールを使えば8Gbps(1GB/s)のバンド幅を実装することは可能である。一方、PCIバスは133MB/s、

64bit66MHzPCIでも532MB/sにすぎない。安価なコンシューマ用PCには64bit66MHzPCIやPCI-X等のサーバー用のI/Oバスが現状では搭載されておらず、将来的にも搭載されるか不透明である。よって当面は安価なコンシューマ用PCでは133MB/sという制約を受ける。

### 2.1.2 主記憶アクセスの競合 (P)

PCIバスからのDMAアクセスも、CPUからの主記憶アクセスも同じ主記憶へのアクセスとなるので、競合が発生する。通常のPCの主記憶のバンド幅は800MB/s~1GB/s程度で、そこに光通信による1GB/sを流し込むとそれだけでバンド幅が足らなくなる。

### 2.1.3 各種PCIカードとの競合 (P)

イーサカード、SCSIカード、ビデオカード、ビデオキャプチャカード等のPCIカードともPCIバスバンド幅や、主記憶バンド幅を分け合わなければならない。例えば、ビデオキャプチャカードとMyrinetでフルサイズ、フルレートの動画転送(約35MB/s)を主記憶経由でやろうとすると、PCIのバンド幅を使い切る。

### 2.1.4 NICメモリへのアクセス集中 (P,C,A)

Myrinetで言えば送信と受信を並行して行わせるならばNIC上のメモリをアクセスするDMA転送が4種類(主記憶⇄NICメモリ、NICメモリ⇄通信リンク)と、NIC上のプロセッサからの命令フェッチ、データフェッチが全て一箇所のNICメモリ(小容量高速SRAM)に集中するので、通信リンクを8倍のバンド幅に引き上げたら、NICメモリでバンド幅を確保することが困難である。

特に通信リンクがI/Oバスのバンド幅より大きい場合は、通信リンクから入力されるパケットをNICメモリを介さず主記憶上に受けようとするNICあるいはネットワークに存在する小容量バッファを溢れさせてかえってパフォーマンスダウンにつながる事が予想される。このため、通信は必ずNICメモリ経由とならざるを得ず、このことはNICメモリのバンド幅不足に拍車をかける。さらなるアーキテクチャ上の工夫により送信時のDMAを排除できるが、そこまでやっているNICは見当たらない。

### 2.1.5 高遅延ゆえの通信時期の集中 (P,C,A)

従来は通信起動遅延時間が短くても数 $\mu$ s程度かかるので、計算しながら即通信するのではなく、メッセージのサイズを大きくして通信回数を減らすためにある程度同じ行き先のデータが溜まるまで計算を行ってから通信を行わざるを得ない。よってSPMDタイプのアプリケーションでは通信の集中するフェーズと、通信が発生しないフェーズがくっきり分かれがちとなる。このため実質的なバンド幅を低めてしまう。さらに、メモリに溜め込んでから送るので不必要な主記憶バンド幅消費を引き起こすこともある。

## 2.2 大きな遅延時間

### 2.2.1 PCIバスの遅い周波数 (P)

PCIバスはFSBや主記憶のバスと比べて通常1/4程度の周波数で動いている。よってPCIバス上で1クロックで済む処理でも、FSBや主記憶のバスにおける1クロックの4倍程度の時間がかかる。

### 2.2.2 バス調停の遅延 (P)

PCIバスには複数のバスマスタが接続される可能性があるためバスの調停にかかる遅延がかかる。この遅延は主記憶バスの1/4程度の周波数で行われる事象である。

### 2.2.3 NICメモリの容量不足 (C)

従来のNIC上のメモリは外付けSRAMによる高速性と引き換えに小容量であったため、主記憶上に通信に必要な情報を一時的または永続的に置かざるを得ないことがある。このため、DMAのセットアップ時間やバス調停時間のオーバーヘッドが通信のたびににかかるようになる。

### 2.2.4 ディスクリプタの動きの激しさ (P,C)

ディスクリプタを主記憶上に置かざるを得ない構成のNICを用いる場合、NICと主記憶とCPUの間をディスクリプタが行き交い、その度に調停フェーズも入る。

また、SRAMで実装されているためにNICメモリの容量が少ないMyrinetなどの場合、NICメモリ上に配置されるVIを保持する領域の容量には制約が大きいため、ノード数の大きなクラスターにVIA[12]を載せるにはNICメモリ上ではなく主記憶上にVIを配置せざるを得なくなる。その場合も同様の遅延が生じる。

### 2.2.5 パケット廃棄対応の複雑なプロトコル (A)

パケットがネットワーク中で廃棄されることが前提になっていることが多いために、プロトコルが複雑になり、それらを大方ソフトで行っているため遅延時間がかかる。

### 2.2.6 低速プロセッサ上での複雑な処理 (A,P)

ホストのCPUへの割り込みやNIC上資源へのCPUからのアクセスに伴うオーバーヘッド(ハードによってはこれがカーネルモードへの移行が強要される)を避けるために、近年のホストCPUと比較して10分の1以下の周波数で動くNIC上のプロセッサで複雑なプロトコル処理を行わせているために、大きな遅延時間がかかる。

### 2.2.7 イベント通知オーバーヘッド (P,A)

PCにおける割込みオーバーヘッドは極めて大きい。それを避けるためにPCIバス上のNICへのポーリングをする場合があるが、メモリを1ワード読むための時間より、PCIバスを1ワード読むための時間は大きく、バスの調停も引き起こす。つまり、PCIバス上の資源は安易にポーリングを行える場所ではない。

さらにMyrinetの場合はNIC上のレジスタが同一ページにマップされるため、マルチユーザー環境では事実上ユーザーモードでのNICレジスタへのポーリングはできない。よってNICメモリ上に設けたフラグをポーリングするか、DMAにより1ワードのDMAで主記憶上のフラグを更新しなければならない。PCIバス上のNICメモリへのポーリングがバンド幅に悪影響を与えるため、例えばMyrinet上のPM2.0ではDMAで主記憶上のフラグを更新するようになった。これでバンド幅は向上するが、遅延時間はDMAの増加で1 $\mu$ s程度増加させてしまう。

### 2.2.8 ヘッダー情報再利用性の乏しさ (A)

パケットとなるべきデータの格納アドレスや、宛先関連のパケットのヘッダーの元になる情報を含むディスクリプタは部分的にも再利用されることはなく、別のパケットを生成するたびにNICと主記憶とCPUの間を行き交う。

### 2.2.9 ディスクリプタの冗長ビット (A)

ディスクリプタは通常ドライバのソフトやNICのDMAが扱いやすいようなデータ構造を持つ場合があり、冗長なビットを多数含んでいると共に、パケットのヘッダーを生成するにもNICが数クロックを費やして毎回変換を行っている。例えば1ワードや1ダブルワードといった粒度の細かい通信を行おうとする場合はディスクリプタの方が本来送るべき情報より多くの遅延を生む。

### 2.2.10 冗長なDMA転送 (P,C)

高バンド幅のNICは前述の理由で、NIC上のメモリに一度データを受け、改めて主記憶にDMA転送する必要があった。よって、送受信で1回ずつの冗長なDMA転送が存在した。アーキテクチャ上の工夫により送信時のDMAを排除できるが、排除していないのが実状である。

## 3 MEMOnet と DIMMnet-1

### 3.1 MEMOnet とは

MEMOnetとは筆者らによって1999年8月より提唱されているNICのクラスである。従来のMyrinet [1]や

RHiNET [3] [4] のように PCI バス等の I/O バスに接続されるのではなく、「主記憶が搭載されるメモリスロットに接続される NIC」と定義されている。

MEMOnet には基盤となる汎用パソコンの主記憶の実装仕様で大きく分けて三つ (SIMM, DIMM, RIMM) が考えられる。SIMM 形式のものを SIMMnet、DIMM 形式のものを DIMMnet、RIMM (Rambus 仕様) 形式のものを RIMMnet と呼んでいる。MINI[9] は SIMMnet の一例ということになる。

我々は MEMOnet の有効性を実証すべく、DIMMnet のプロトタイプ DIMMnet-1 を現在開発中である。

### 3.2 DIMMnet-1 の概要

DIMMnet-1 は、PC100 または PC133 仕様の DIMM スロットにささるネットワークインタフェースである。その基本構造を図 1 に、主な仕様を表 1 に示す。

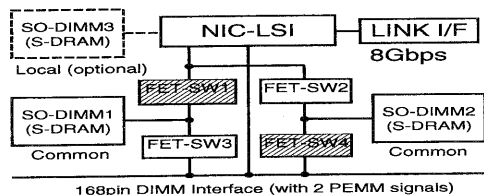


図 1: DIMMnet-1 の基本構造

NIC-LSI は低遅延の FET バススイッチにより 2 バンクの SO-DIMM (ノート型 PC で用いられる汎用部品) 間を切り替えてつ、RHiNET-2 互換の光リンクインタフェースとデータの送受信をする。DIMM スロットの信号はじかに NIC-LSI に入力される NIC 制御ポートを有する。

SO-DIMM のシリアルプレゼンスディテクトデータを読み、それに対応する適切なシリアルプレゼンスディテクトデータを DIMM スロットに出力する。さらにオプションでローカルな SO-DIMM を 1 バンク持つ。

メモリバス側のインタフェースは日本電子機械工業会規格の「プロセッサ搭載メモリ・モジュール (PEMM) 動作仕様標準」[10] に準拠させる。PEMM 規格準拠のチップセットやマザーボードは現状では存在しないので、PEMM 準拠モード以外にも、PEMM で追加された 2 つの信号 (バンクメモリへのアクセスを待たせる信号と割込み信号) が無くても動作するモードの二つのモードを有する。

表 1: DIMMnet-1 の主な仕様

インタフェース	PCI133 ベースの PEMM
通信リンク	12ch 光接続モジュール MDS4212A, MDR4212A
NIC メモリ (共有)	PC133 SO-DIMM × 2
NIC メモリ (ローカル)	PC133 SO-DIMM × 1
NIC-LSI のテクノロジー	0.18 $\mu$ m CMOS
通信リンクバンド幅	各方向 8Gbps (全二重)
NIC メモリバンド幅	2128MB/s (ローカル メモリ使用時 3192MB/s)
送信時 NIC 遅延時間	45ns
受信時 NIC 遅延時間	45ns
メモリアクセス遅延増	0 クロック

## 4 DIMMnet-1 における改善点

本章では、2 章に列挙した問題点に対する DIMMnet-1 における改善を 2 章に列挙した順に述べる。

### 4.1 バンド幅の不足への対応

#### 4.1.1 メモリスロットで 1GB/s

メモリスロットを CPU とのデータ受け渡し口にしており、そのバンド幅は光通信リンクの片方向分 (1GB/s) を

有する。これは PCI バスの 8 倍、64bit66.6MHz PCI バスの約 2 倍の能力を有し、1Gbps 程度の通信リンクを PCI バスで支えている Myrinet や RHiNET-1 と比較して、8 倍の能力になった通信リンクを有する DIMMnet-1 はバランスを維持している。

#### 4.1.2 ダブルバッファでバンド幅倍増

DIMMnet-1 では主記憶のバスに NIC が搭載される MEMOnet アーキテクチャをとる。NIC 上に大容量の SO-DIMM が最低 2 枚搭載され、ある瞬間ではその片方が CPU の方を向いており、もう片方が NIC の方を向いているように低遅延 FET バススイッチが制御される。つまり NIC メモリがダブルバッファの構成を取ることで、見かけ上合計で 2GB/s のメモリバンド幅があるように見え、CPU と NIC の同時並列アクセスが実現される。

#### 4.1.3 PCI カードとはバンド幅を非共有

DIMMnet-1 は PCI バスとは十分に隔離された場所をアクセスするので、他のイーサカード、SCSI カード、ビデオカード、ビデオキャプチャカード等の PCI カードとは NIC がバンド幅を共有しない。

#### 4.1.4 内蔵メモリ、バイパス等によるアクセス分散

DIMMnet-1 は、RHiNET/SW に用いられた ASIC またはその後継シリーズを使用して実装される。このため内蔵メモリも比較的大きなものを実装でき、内蔵 CPU のキャッシュや、通信リンクと CPU とのインタフェースなどに使用される FIFO も内蔵される。それらのバッファが溢れるか、NIC メモリへのアクセスは通信のソースまたはデスティネーションが NIC メモリではない限り発生しないので、NIC メモリへのアクセス要求は軽減される。さらに、On the fly 送信においては通信リンクと CPU とのインタフェースによる FIFO バッファ間がバイパスされ、NIC メモリは経由しないので NIC メモリへのアクセス要求は発生しない。

#### 4.1.5 低遅延化による通信時期の分散

DIMMnet-1 では後述する様々な対策により、通信遅延時間を大幅に短縮する。従来のメッセージ発生が数マイクロ秒以上かかっていたのに対し、DIMMnet-1 では最短で  $45\text{ns} + \alpha$  ( $\alpha$  はチップセット遅延) でメッセージを発生させることができる。その結果、通信粒度を粗くするために通信時期が集中してしまうことを防止・軽減でき、実効的な通信バンド幅は高くなったように見える。

## 4.2 大きな遅延時間への対応

### 4.2.1 メモリバスの高い周波数

DIMMnet-1 では主記憶のバスに NIC が搭載される MEMOnet アーキテクチャをとるので、安価な PC 上でも PCI バスと比べて 4 倍程度の周波数で動作可能な NIC を構築できる。DIMM スロットでは 100MHz ~ 133MHz 程度が現在の主流だが、その程度の周波数は FPGA では実装困難ながら ASIC によれば実現可能な領域である。

### 4.2.2 ダブルバッファによる調停の削減

DIMMnet-1 では主記憶のバスに FET バススイッチを介して 2 つのメモリバンクが接続されており、NIC はある瞬間は 2 つあるメモリバンクのうち少なくとも一つを排他的にアクセス可能な状態にあり、別バンクの NIC メモリや、別スロットにある通常の DIMM への CPU からのアクセスと競合しない。NIC の DMA による送受信が必ずしも CPU との調停を必要としないことを意味する。さらに他の PCI デバイスとの主記憶の競合に関しては一切ない。よって従来の PCI バスによる NIC の場合はあらゆる DMA アクセスが CPU や他の PCI バス上のマスタデバイスとの調停を必要としていたのに対して、バス調停の遅延の総和を縮めることができる。しかもこの遅延は主記憶バスという PCI バスの 4 倍程度の周波数で行われる事象によるものである。

#### 4.2.3 内蔵メモリと大容量 NIC メモリ

DIMMnet-1 では比較的潤沢な内蔵メモリを持った ASIC により実装され、基本的には内蔵メモリが溢れたり、デスティネーションが外部メモリである場合に限り外部メモリへのアクセスが発生する。このため、外部メモリへの速度的要求は Myrinet 等の従来型 NIC と比較して小さく、ノート PC 用の安価な汎用部品である SO-DIMM による外部メモリによって NIC メモリを構成することができる。高々 4MB 程度しか無い Myrinet 等と比較して、安価に大容量の NIC メモリを持つことができ、NIC メモリの容量不足に起因する望ましくない使い方に伴う DMA やバス調停時間のオーバーヘッドが回避できる。

#### 4.2.4 NIC 上に配置されるディスクリプタリスト

DIMMnet-1 では主記憶と同等の遅延時間でアクセス可能な大容量の NIC メモリを搭載するために、従来の NIC で用いられていた DMA 転送モードにおけるディスクリプタを主記憶上に保持したり、通信前に CPU が作成し NIC に転送する必要は無く、NIC メモリ上に配置することができる。このため CPU はディスクリプタリストへのトップポインタのみを NIC にセットしてやるだけで良く、NIC と主記憶と CPU の間をディスクリプタが行き交い、その度に調停時間や DMA 設定時間、転送遅延時間がかかるといった従来の NIC で発生していた遅延時間を排除できる。またノード数の大きなクラスタに VIA を載せる際に NIC メモリ上に VI を配置できるので、従来の NIC で発生していた遅延時間を排除できる。

#### 4.2.5 選択的信頼性維持ポリシー

DIMMnet-1 ではパケットそのものが消失することは bit 付け等のハード的なエラーが起こった時のみで、ethernet や ATM のように正常な動作をしつつパケットが廃棄されることはない。DIMMnet-1 で用いられる光モジュールの bit error rate は  $10^{-20}$  であり高い信頼性を持っている。ただし台数や設置状況によってはハードの信頼性を十分高いと見られる場合と見られない場合があると考えられる。

さらに DIMMnet-1 の場合はある瞬間では NIC が確実に占有できる容量もバンド幅も大きな領域 (NIC メモリの片方のバンク) を持っており、受信の際のメモリアクセスの優先度を高くすることにより受信パケットの欠落を従来の NIC より低く抑えさせることが可能である。

一方、上位ソフトウェアレイヤによる対応 (チェックポイントリングなど) を仮定する動向も一部で存在する。例えば VIA1.0 [12] のように信頼性のある通信は必須事項ではない。RWCP で開発された SCORE-D においてもコンシステントチェックポイントリング [13] が実現されている。

よって DIMMnet-1 ではハードで通信の信頼性を保証するモードと、保証しない代わりに高速なモードを切り替え可能とする。信頼性を保証しないモードでは、ハード的にエラーは検出可能とすることで上位ソフトによる対応を促す。一種の同期手段としても用いられるプログラマから見える ACK 付きの通信と、ACK なしの通信を明示的に指示可能とし、適宜に ACK 付きの通信を追加したり間引くことで信頼性と性能の妥協点を調整可能とする。

信頼性を保証するモードでは NIC がパケットにシーケンス番号を付け、通信ログを NIC メモリ上に取り、NACK が返ってきた場合、NACK を引き起こしたシーケンス番号までさかのぼって再送を行う。受信バッファオーバーフローを起こした場合は NACK が送信元に返され、受信側は同一ノードからのそのシーケンス番号以降のパケットを廃棄する。通信ログを残す際に余計なメモリアクセスが増加するが、ローカルメモリを設置することにより、通信ログ保存に伴うオーバーヘッドを軽減可能とする。

#### 4.2.6 高速プロセッサ上での簡素な処理

DIMMnet-1 では、Myrinet の LANai プロセッサのように PCI バスのクロックまたはその半分の周波数 (33.3MHz) で動作する低速な CPU ではなく、メモリバスのクロック

(100~133MHz) で動作する NIC-LSI に内蔵される CPU によって処理される。その処理内容は、パケットそのものが消失することは十分頻度が低くチェックポイントリング等の上位ソフトウェアで扱うべき事象として扱う比較の簡素なプロトコル処理である。

#### 4.2.7 制御ポート等による軽いポーリング

DIMMnet-1 では、DIMM スロット上に FET バススイッチを介して接続される NIC メモリと、DIMM スロットに直に接続される NIC 制御ポートを有する。このポートは CPU からしかアクセスされないため調停も、バンク切り替えの煩わしさも無く、一部の特権レジスタを除く内部のレジスタや内部メモリをユーザーモードでアンキャプチャ可能な主記憶と同様の比較的小さなオーバーヘッドでアクセス可能とする。ポーリングによるアクセスと NIC の送受信が全て PCI バスを通過する従来の NIC と異なり、独立の経路でアクセスがされるのでポーリング NIC の送受信のバンド幅低下といった問題を誘発しない。よって、頻繁にポーリングすることによって割込みオーバーヘッドを排除することが可能である。

NIC メモリ上にあるワードをポーリングするケースでは、その時点で CPU 側から見えているバンクへのポーリングは上記のレジスタへのポーリングと同様にアンキャプチャ可能な主記憶へのポーリングコストと同等である。その時点で CPU 側から見えていないバンクへのポーリングは、NIC 制御ポートを介したバンク切り替え要求後、NIC 制御レジスタのポーリングを行ってバンクが切り替わったことを確認するというオーバーヘッドが追加される。これと数マイクロ秒を要する割込みオーバーヘッドに比べればはるかに低オーバーヘッドである。

#### 4.2.8 ヘッダーバッファとディスクリプタリスト

DIMMnet-1 では、後述するシングルモード On the fly 送信においては、あらかじめヘッダーバッファと NIC メモリ上にパケットヘッダーの種となるデータ (ヘッダーシード) を登録しておき、そこに保存されている同一ページ内へのワードまたはダブルワード送信は、既に登録済みのパケットヘッダー情報が再利用される。

DMA 転送モードにおけるディスクリプタは NIC メモリ上にリスト形式で配置する。このため CPU はディスクリプタリストへのトップポインタのみを NIC にセットすれば良く、その再利用が可能である。例えば、ループの繰返しごとに同じ通信を行うのであれば、ループ内に存在する通信のディスクリプタのリストを NIC メモリ上に保持しておくことにより、一連の通信のディスクリプタは再利用され、通信起動オーバーヘッドは大幅に短縮される。

#### 4.2.9 冗長ビットを排除したヘッダーシード

DIMMnet-1 では、後述するシングルモード On the fly 送信においては、あらかじめヘッダーバッファと NIC メモリ上にパケットヘッダーの種となるデータ (ヘッダーシード) を登録する。これは最初の登録を行う API の回数の中で生成され、プロセスグループ ID とオフセット値を追加するだけでヘッダーになる冗長ビットをそぎ落としたヘッダーに限りなく近い状態のデータであり、実際の通信発生時に未圧縮の状態では NIC に渡されない。

#### 4.2.10 主記憶的 NIC メモリと On the fly 送信

DIMMnet-1 では、従来の NIC メモリよりも低オーバーヘッドかつ高バンド幅で、桁違いに大容量である性質を生かし、NIC メモリをメッセージの終点とし、ユーザープロセスから直接、NIC 上のメモリに届いたデータをアクセスする。DMA 転送の回数は DMA 送信モードでは送信時に NIC メモリから通信リンクへの DMA1 回と、受信時に通信リンクから NIC メモリへの DMA1 回のみであり、送受信で 1 回ずつの DMA 転送を削減している。

さらに、On the fly 送信においては NIC メモリは経由しないので NIC メモリからの DMA は完全に排除される。

## 5 DIMMnet-1 の細粒度通信機構

### 5.1 細粒度通信 API の概要

DIMMnet-1 では RWCP で開発され MPICH-PM や SCore などのベースになっている PM API に、表 2 に一部を列挙した DIMMnet-1 の細粒度通信の拡張を施した PM-D (PM DIMMnet extension) API を当面提供する予定である。

表 2: PM-D API の概要

NIC メモリ領域	<code>_pmAlloc</code>	割当て
	<code>_pmFree</code>	解放
NIC メモリバンク	<code>_pmBank</code>	状態を得る
	<code>_pmBBlock</code>	ロック
	<code>_pmBRelease</code>	ロックの解除
Block mode On the fly 通信	<code>_pmOpenB</code>	書き込み口生成
	<code>_pmCreateH</code>	ヘッダシードの生成
	<code>_pmKickB</code>	起動アドレスの生成
	<code>_pmCloseB</code>	書き込み口の閉鎖
Single mode On the fly 通信	<code>_pmOpenS</code>	書き込み口生成
	<code>_pmMapS</code>	ヘッダシードのマッピングとヘッダバッファへの登録
	<code>_pmUnmapS</code>	ヘッダシードのアンマッピングとヘッダバッファからの削除
	<code>_pmCloseS</code>	書き込み口の閉鎖

### 5.2 シングルモード On the fly 通信

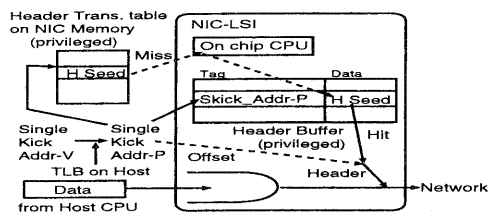


図 2: シングルモード On the fly 通信

シングルモード On the fly 通信におけるパケット生成メカニズムを図 2 に示す。シングルモード On the fly 通信に際しては、あらかじめ `_pmOpenS()` により起動用の仮想アドレス `skick.addr` を確保し、`_pmMapS()` によりヘッダバッファに `skick.addr` をタグとして連想されるヘッダシードをセットしておく。送信すべきデータがレジスタ上に存在すれば、CPU がレジスタ上にあるデータを仮想アドレス `skick.addr` に書き込むというわずか 1 命令でメッセージ送信を起動できる。

ヘッダシードは送信すべきパケットのヘッダから、リモートアドレス部の下位が削除されたものである。これが登録されるヘッダバッファやヘッダ変換テーブルはユーザーモードからは直接は触れることのできない場所に配置され、`_pmMapS()` の内部でカーネルモードに移行して、カーネル空間にマップされる NIC レジスタ経由でセットされるので、プロテクションをつかさどるプロセスグループ ID (PGID) をユーザーが勝手に書き換えたりできない。よって、この送信モードに限り、`_pmMapS()` の内部でカーネルモードからリモートアドレスを物理アドレスで登録でき、受信時のリモートにおけるアドレス変換オーバーヘッドを削除可能である。

なお、このモードで発生するパケットは、単純なリモートライトを行うだけではない。リモートノード `remote.node` のアドレス `remote.addr` (`pv.flag` で仮想アドレスか物理アドレスかを指定) に `command` を転送し、ステータスを仮想アドレス `stat.addr` に、結果を仮想アドレス `result.addr` に返す。`command` には ACK なしリモートライト、ACK

付きリモートライト、ACK なしリモート間接ライト、ACK 付きリモート間接ライト、リモートリード、返り値なしホットメッセージ、返り値付きホットメッセージがある。ホットメッセージとは、リモートノードの NIC プロセッサでのハンドラーの仮想アドレスを指定した一種のリモートプロセスジャコールである。

### 5.3 ブロックモード On the fly 通信

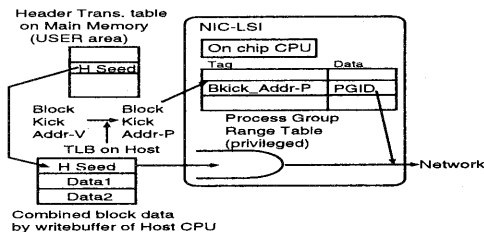


図 3: ブロックモード On the fly 通信

ブロックモード On the fly 通信におけるパケット生成メカニズムを図 3 に示す。シングルモード On the fly 通信があらかじめ NIC 側にヘッダシードを設定しておくのに対し、ブロックモード On the fly 通信は `_pmCreateH()` で生成したヘッダシードをユーザー空間に保持し、これをブロックモード On the fly 送信用の FIFO 入口 `bpush.addr` (write combining 属性の領域) に書き込み、送りたいデータをそれに後続して書き込み、最後に `bkick.addr` (uncachable 属性の領域) にダミーを書いて先行する `bpush.addr` に書き込まれたデータをパケットにしてネットワークに出力する。

こうして CPU のライトバッファによって書き込まれた細かいデータはライトバッファの容量 (通常、キャッシュのラインサイズ) までブロック化されてバーストライトが発生し、最後の `bpush.addr` への書き込みがライトバッファに残ったデータを全て NIC のブロックモード On the fly 送信用の FIFO に吐き出す。

その際、ヘッダシードがユーザー空間にあるからといって事前に PGID を書き換えても `bpush.addr` に対応する物理アドレスを頼りに NIC が正しい PGID に付け替えるので他人に成りすますことはできない。

その他、送信データが 8 バイト以上になっても構わないことと、物理アドレスでのリモートアドレス指定ができないことを除けば、概ねシングルモード On the fly 送信のパケットが実行する動作は同じである。

## 6 性能見積もり

### 6.1 最短送信時 NIC 遅延時間

シングルモード On the fly 送信の際に、ヘッダバッファがヒット場合が最短送信遅延時間を達成する。その際の各クロックサイクルでの動作は以下ようになり、NIC が搭載されるメモリスロット上に最初の信号が発生してから 6 クロック (133MHz 動作時に 45ns) で通信リンクインタフェースへの出力が始まる。

1. SDRAM の 8 バイトライト動作第一サイクル目。
2. 上記の第二サイクル目。
3. メモリスロット上のアドレスとデータを入力 FIFO に取り込む。
4. 取り込んだアドレスの上位ビットからヘッダバッファを検索する。
5. ヒットして得られたヘッダシードに取り込んだアドレスの下位をはめ込みヘッダを完成させ、送信 FIFO にヘッダを突っ込む。
6. 送信 FIFO からヘッダを読み出し、リンクに送出。

なお、実際の送信に際しては、CPU がレジスタ上にある `double` のデータを `skick.addr` に書き込む 1 命令を実

行してからメモリスロット上に信号が現れるまで数クロック（チップセットに依存）の遅延が加算される。

CPU からみればあらかじめヘッダーバッファにヘッダーシードをセットしておけばわずか1命令でメッセージ送信が起動できる点で究極的に高速な送信を実現している。

## 6.2 最短受信時 NIC 遅延時間

物理アドレス指定の8バイトのリモートライトの際に、NIC の方に向けているバンクに書き込む場合に、最短送信遅延時間を達成する。その際の各クロックサイクルでの動作は以下のようになり、通信リンクインタフェース上に最初の信号が発生してから6クロック（133MHz 動作時に45ns）でNIC メモリへの出力が始まる。

1. 通信リンクから受信 FIFO にヘッダー1を取り込む。
2. 受信 FIFO からヘッダー1を取り出す。
3. 受信 FIFO からヘッダー2を取り出し、CRC を計算。ヘッダー中のコマンドを検査し、物理アドレス指定8バイトリモートライトであることを知る。
4. 受信 FIFO からデータを取り出し、CRC を計算。
5. 受信 FIFO から CRC を取り出し、チェック後、受信データとアドレスをライト FIFO に突っ込む。
6. アドレスから書き込みバンクをチェックし、アクティブであることを知り、NIC メモリへの書き始める。

なお、仮想アドレス指定の場合は第六クロックサイクルでNIC 内 TLB による物理アドレスへの変換が入るために最短でもう1クロック時間がかかる。アクティブではないバンクに受信する場合は、バンクがCPU からロックされていない場合は、もう1クロック時間がかかる。

## 6.3 最短リモートライト遅延時間

チップセット遅延時間と最短送信遅延時間と最短受信遅延時間の合計が2ノード間直結時の最短リモートライト遅延時間となる。チップセット遅延時間を多目に見積もって10クロック（75ns）とすると、合計で165nsとなる。

表3に各NIC でのプロセス間単方向通信遅延を示す。

表3: 各NIC でのプロセス間単方向通信遅延

DIMMnet-1	165ns	1倍	本論文
Memory channel 2	2.2 $\mu$ s	13倍	文献[14]
PM on Myrinet	7.5 $\mu$ s	45倍	文献[15]
GigaE PM II	44.6 $\mu$ s	270倍	文献[2]

## 7 おわりに

従来型NIC の数多くの問題点をバンド幅と遅延時間の観点から列挙し、現在開発中のDIMMnet-1がいかにしてこれらの問題を克服しているかを示した。I/OバスではなくメモリスロットにNIC を配置するという斬新なアプローチに起因する様々なメリットを享受できるようDIMMnet-1では様々な工夫をしている。

今回示した細粒度通信機構によれば、CPU の命令セットアーキテクチャから根本的に変えない限り実装テクノロジーやソフトが進歩してもなかなか短縮できなかつた遅延時間を、Alpha ベースのシステムで用いられているMemory channel 2の13倍、高性能なクラスター用NIC として各所で活用されているMyrinet 上のPM の45倍、汎用Gigabit Ethernet で構成するGigaE PM の270倍という圧倒的な遅延時間短縮の見通しを得た点で画期的である。

DIMMnet-1 は新たなハード性能のバランス点を与え、従来のバランス点の上に成り立っていた並列ソフト技術の再構築を促すような技術であると思われる。DIMMnet-1 の設計思想に基づくNIC 製品が流布するようになれば、例えば逐次バイナリコードをそのまま並列化してしまうShastaのような魅力的なソフトウェア技術を、性能面で真に使える技術に押し上げ、普及させていく原動力となるのかもしれない。そのような評価は今後の課題である。

今後は、今回触れなかつたホットメッセージ等を実行する主体であるNIC 内CPU 等の一部ハードウェアを共有する形で同一LSI 内に実装されるRHiNET-2 との両立を模索しつつNIC 内CPU 回りの設計を進め、2001年春の完成を目指し、DIMMnet-1 の開発を進める予定である。

## 謝辞

バンクメモリへのコンパイラ対応の経験を御教授いただいた早稲田大学の笠原博徳教授、MEMnet という名称についてコメントをいただいた東大の松本尚助手、従来型NIC の問題点を御議論いただいた新情報処理開発機構の手塚宏史氏、佐元真司氏に感謝致します。

## 参考文献

- [1] Myricom corp. <http://www.myri.com/>
- [2] 住元, 堀, 手塚, 原田, 高橋, 石川 "GigaE PM II: Gigabit Ethernet による高速通信ライブラリの設計", 情報処理学会研究報告 99-ARC-134 (SWoPP'99), pp.61-66 (1999.8)
- [3] 山本, 建部, 横山, 土屋, 宮脇, 清水, 天野, 工藤 "高性能並列計算用ネットワーク RHiNET-1 の実装と評価", 情報処理学会研究報告 HOKKE2000 (2000.3)
- [4] Tomohiro Kudoh, Shinji Nishimura, Junji Yamamoto, Hiroaki Nishi, Osamu Tatebe and Hideharu Amano "RHiNET: A network for high performance parallel processing using locally distributed computers", IWIA'99 (1999.11)
- [5] S. Nishimura, T. Kudoh, H. Nishi, N. Matsudaira, K. Harasawa, S. Akutsu, K. Tasho, H. Amano "A network switch using optical interconnection for high performance parallel computing using PCs", Parallel Interconnect'99 (1999.10)
- [6] "Computer Makers Propose New PCI Design", <http://www.techweb.com/wire/story/TWB19980904S0008>
- [7] InfiniBand Trade Association, <http://www.sysio.org/>
- [8] 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" 情報処理学会研究報告 99-ARC-134 (SWoPP'99), pp.73-78 (1999.8)
- [9] Ron Minnich, Dan Burns and Frank Hady "The Memory Integrated Network Interface" IEEE Micro, Vol. 15, No. 1, (1995.2)
- [10] 日本電子機械工業会 "日本電子機械工業会規格: プロセッサ搭載メモリ・モジュール (PEMM) 動作仕様標準", EIAJ ED-5514 (1998.7)
- [11] D. J. Scales, K. Gharachorloo, and C. A. Thekkath "Shasta: A Low Overhead, Software-Only Approach for Supporting Fine-Grain Shared Memory", ASPLOS'96 (1996.10)
- [12] Intel corp. "Intel VI Architecture Developer's Guide V1.0", <ftp://download.intel.com/design/servers/vi/intel.pdf>
- [13] 西岡, 堀, 手塚, 石川 "クラスターにおけるコンシステントチェックポイントの実現", 並列処理シンポジウム JSP'99, pp.207-214 (1999.6)
- [14] M. Fillo and R. B. Gillett "Architecture and Implementation of MEMORY CHANNEL 2", Digital Technical Journal, Vol.9(1), (1997)
- [15] H. Tezuka, A. Hori, Y. Ishikawa, and M. Sato "PM: An Operating System Coordinated High Performance Communication Library", High Performance Computing and Networking'97 (1997)