

並列コンピュータ Cenju-3 用 Mach における NORMA IPC の実現

菅原 智義

C.Howson

高野 陽介

小長谷 明彦

NEC C&C 研究所

我々は分散メモリ型並列コンピュータ Cenju-3 をターゲットに、Mach マイクロカーネルをベースにした並列 OS DenEn を開発している。DenEnにおいては、Mach-IPC と互換性を持つつ、高トラフィックに対応できるネットワーク IPC が重要である。Mach では Mach-IPC をネットワーク拡張した NORMA IPC が用意されているが、DenEn の要求を満たすものではない。そこで、我々はオリジナルの NORMA IPC をベースに、新たに NORMA IPC/DE を構築した。本稿ではこの NORMA IPC/DE の実装と評価について述べる。NORMA IPC/DE は、前述の要求を満たしつつ、小量メッセージでは低レイテンシ、大量メッセージでは高スループットを実現している。現時点での性能はレイテンシ $210\mu s$ 、スループット 22M バイト / s である。

An Implementation of NORMA IPC on Cenju-3

Tomoyoshi Sugawara

Christopher Howson

Yosuke Takano

Akihiko Konagaya

C&C Research Laboratories NEC Corporation
1-1, Miyazaki 4-Chome Miyamae-ku Kawasaki Kanagawa 216 Japan

In this paper, we describe an implementation and its evaluation of NORMA IPC/DE. We have modified original NORMA IPC to take advantage of Cenju-3's high-speed and reliable network. Our goals is to provide the same local Mach-IPC semantics over the network, while supporting high message traffic with low latency and high throughput. We have achieved a latency of $210\mu s$ for short messages, and throughput of 22Mbytes/s for large messages.

1 はじめに

分散メモリ型、すなわち、多数の要素プロセッサを共有メモリを使わずに高速なネットワークで結合した並列コンピュータが多くベンダーで開発され、商用化され始めている。このような並列コンピュータのためのOSは、科学技術計算などのために高度な計算能力を提供するだけでなく、データベースやマルチメディアなどの非数値を扱うような並列アプリケーションも効率よく実行できることを望まれている。これを実現するために、我々は分散メモリ型の並列コンピュータ Cenju-3[1] をターゲットにして、Machマイクロカーネル[3] をベースにした並列OS DenEn[2]を開発している。

Machではタスク間の通信に Mach-IPC を用いるが、これとほぼ同じセマンティクスを持った、NORMA(NO Remote Memory Access) IPC[4] と呼ばれるネットワーク IPC がカーネルレベルで提供されている。Machでは OS サーバに対するサービス要求は元より、マイクロカーネルに対するシステムコールも大半は Mach-IPC を使用しているため、分散メモリ型の並列コンピュータに Mach を載せた場合、別のノードの OS サーバやマイクロカーネルのサービスを利用することができる。

NORMA IPC は分散メモリ型の並列コンピュータを意識して設計されてはいるものの、同時に複数のノード間通信を確立できない、Mach-IPC のフロー制御のセマンティクスを反映していないなど、Cenju-3 プラットフォームに適合しない点が多い。そこで、我々は NORMA IPC をベースに、分散メモリ型並列コンピュータの通信性能を十分に引き出す NORMA IPC/DE を新たに構築し、Cenju-3 上に実装した。本稿では、NORMA IPC/DE の設計思想、実装方式、評価結果および高速化の検討結果について報告する。

2 ネットワーク IPC に対する要求事項

我々のネットワーク IPC に対する要求事項は次の 3 点である。

• 高スループット、低レイテンシ

スループットを高くすることは、画像処理などのデータ転送量の多いアプリケーションの性能向上に大きく貢献する。また、分散共有メモリやタスクマイグレーションの実現のために、OS レベルでも高スループットが要求される。

また、OS サーバを分散化した場合、リモートのサーバに対する IPC とローカルのサーバに対する IPC との間に大きな性能差があると使い物にならない。そのため、レイテンシも小さくする必要がある。

• 複数のノード間通信の同時確立

分散メモリ型の並列コンピュータ上で並列アプリケーションを実行する場合、ノード間の通信は必然的に多くなる。また、OS サーバを分散化した場合、1つの OS サーバに複数のノードのユーザタスクからサービス要求が届くことは頻繁に起りうる。このような高トラフィックの状況でも、すべての IPC を処理できることが必要である。

• Mach-IPC との互換性

シングルプロセッサシステム用の OS サーバをそのまま分散化し、安全に動作させるためには、インターフェースだけでなく、セマンティクスも Mach-IPC と互換にする必要がある。

3 NORMA IPC

3.1 Mach-IPC との関係

NORMA IPC は Mach-IPC のネットワーク拡張であり、セマンティクスはほぼ同じである。Mach-IPC では、メッセージはポートと呼ばれる通信チャネルを介して受け渡される。実際には、ポートはメッセージのキューであり、一旦、キューにつながれたメッセージはカーネルにより保護される。また、ポートにはポートライトと呼ばれる送受信の権利があり、ポートに対してアクセスできるタスクはポートライトを持つタスクに限定される。これらにより、Mach-IPC はセキュリティの高い通信を実現している。

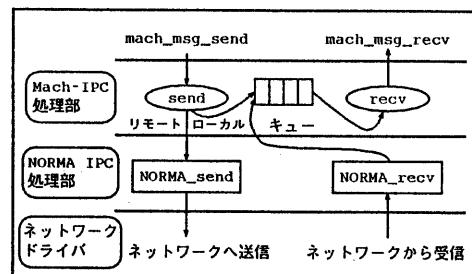


図 1: Mach-IPC と NORMA IPC の関係

NORMA IPC を使う場合、カーネルインターフェースは Mach-IPC のものをそのまま用いる。図 1 に示すように、NORMA IPC の処理は Mach-IPC のキュー操作の代わりに行なわれる。Mach-IPC の処理ルーチンは、ポートがリモートなものである場合は、ポートのキューに入れる代わりに NORMA IPC の送信処理ルーチン (NORMA_send) にメッセージを渡す。リモ-

トノードに送られたメッセージは、NORMA IPC 受信処理ルーチン (NORMA_recv) により Mach-IPC のキューにつなげられる。

また、Mach-IPC で扱われるメッセージは、図 2 に示すようにヘッダとデータタイプ、データから構成される。データはデータ本体を含む inline データと、データへのポインタのみを含む out-of-line (OOL) データに 2 種類に分けられる。OOL データの転送では仮想メモリの操作でコピーが行われるため、大量データの転送を比較的小さいオーバヘッドで行なうことができる。さらに、ポートライトもデータの 1 つとして扱うことができ、ポートライトの生成、複製、移動は IPC として実現されている。

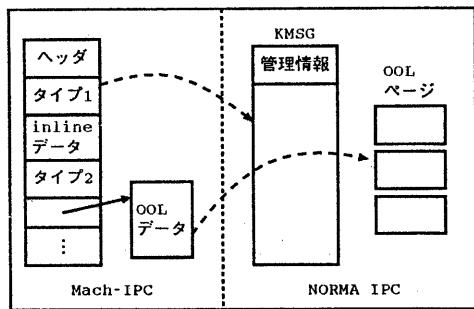


図 2: メッセージの関係

NORMA IPC ではこれらのデータをすべて扱うことができる。ただし、ポインタをそのまま扱うことはできないので、OOL データとそれ以外は別々に扱われる。ヘッダとデータタイプ、inline データはひとまとめにされ KMSG として、OOL データはページ毎に区切られて OOL ページとして扱われる。

3.2 仕様上の問題点

Mach-IPC の機能の中で、NORMA IPC ではサポートされていないものがある。

1 つは Message Accept の通知 (Notification) で、送信側でメッセージの受信がすでに済んだかどうかを知るための機能であり、ノンブロッキングの RPC を実現するために使われる。

もう 1 つは、ポートのキューレベルでのフロー制御で、送信側のポートのキューにつながるメッセージの数を制限することによりフロー制御を行なう機能であり、受信側が処理できなくてメッセージが溜った時に送信側をブロックさせる働きをする。低レベルなフロー制御では、不当なユーザタスクが一方的に OS サーバにメッセージを送りつけた場合に、OS サーバの要求

受け付けのキューが占有されることを防ぐことができないので、OS サーバを安全に動かすためにはこの機能が必要不可欠である。

3.3 プロトコルの問題点

NORMA IPC のプロトコルでは、受信側はメモリが不足した時、メッセージを受けとらずに捨てる。また、送信側はタイムアウトすると該当するメッセージを再送する。このプロトコルでは複数のノードから 1 つのノードに通信要求が集中した場合、メッセージを途中まで受信したところで受信側の空きメモリが不足すると、受信側はメッセージを捨て続ける。一方、送信側はタイムアウトと再送を繰り返す。結果的に、この受信側ノードに関わっているすべての通信が停止する。

この状態を回避するために、NORMA IPC ではある時点での通信できるノードを 1 つだけに限定している。1 つのノードと通信している間に他のノードからのメッセージが届いても、単に棄却するだけである。このため、同時に複数のノード間通信を確立することはできない。

また、信頼性の高いネットワークを仮定しているにも関わらず、プロトコルは、低い信頼性を仮定している点も問題である。

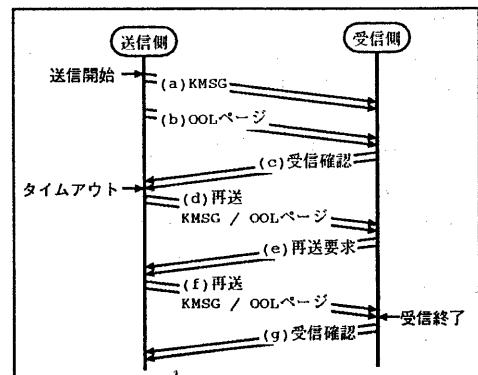


図 3: NORMA IPC のプロトコル

以下に NORMA IPC の具体的なプロトコルを示す（図 3 参照）。以下の説明でメッセージとはすべてネットワークを流れるメッセージを指す。

1. 送信側は、最初に KMSG メッセージを送る (a)。ページサイズを越える場合はページサイズに分割する。KMSG メッセージの送信が終了すると、

- 次に OOL ページメッセージを送る (b)。これらのメッセージの送信時にはタイムアウトを設定する。
2. 受信側は、KMSG または OOL ページメッセージを受けとる度に受信確認メッセージを返す (c)。十分な空きメモリがない場合は、メッセージを捨てる。途中のメッセージが失われてしまった場合は、メッセージを再送してもらうために再送要求メッセージを返送する (e)。すべてのメッセージを受信し、すべての確認メッセージを返送した時、受信は完了する (g)。
 3. 送信側は、受信確認メッセージが届かずについにタイムアウトした場合、該当するメッセージを再送する (d)。再送要求メッセージを受けとった場合も、同様に該当するメッセージを再送する (f)。送信したすべてのメッセージに対する受信確認メッセージを受けとった時、送信は完了する。

3.4 実装上の問題点

オリジナルの NORMA IPC の実装では、ほとんどの処理が割り込みレベルで行なわれるため、Mach 標準のメモリ割り当て機能を使えず、メッセージをそのまま上のレベルに渡すことができない。このため、ネットワークドライバのレベルで受信したメッセージをすべて別のバッファにコピーしている。

コピーのオーバヘッドはメッセージサイズが小さい時は無視できるが、サイズが大きくなると処理時間の大半を占めることになるため、この実装ではスループットを高くすることは難しい。

4 NORMA IPC/DE の設計方針

2 節で列挙したネットワーク IPC に対する要求と、3 節で述べた NORMA IPC の問題点を踏まえて、NORMA IPC/DE を設計する。設計の要点をまとめると次のようになる。

プロトコルの変更

同時に複数のノードからの通信を受け付けられるようにプロトコルを変更する。NORMA IPC のプロトコルの問題点である、通信の途中で受信側のメモリが不足する状態を回避するために、通信の始めに必要なメモリを予約して、予約が成功した時だけ KMSG や OOL ページを送ることにする。このために、メモリ予約用のメッセージを追加する。これにより、通信の途中で受信側がブロックすることなくなる。また、受けとられることのないメッセージが送られなくなるので、ネットワークのトラフィックを無駄に増やすこともなくなる。

これにより、タイムアウトや再送をプロトコルに組み込む必要がなくなるので、これらを取り除き、プロトコルを簡素化する。

また、サーバに対するネットワーク IPC のレイテンシを小さくするために、数百バイトまでの小量の inline データ転送のためのプロトコルを簡単なものにする。サイズを数百バイトに限定した理由は、サーバに対する要求、応答のメッセージの大半は数十～百数十バイト程度に収まるからである。

プロトコルメッセージ受信用バッファの確保

変更したプロトコルでは、最初のメモリ予約のためのメッセージだけは、常に受信できなければならない。しかし、受信用のバッファを動的に割り当てる方法では、メモリが不足した時、メッセージを受信できなくなる。そこで、全ノード数分のプロトコルメッセージ（KMSG と OOL ページ以外のメッセージ）受信用のバッファを固定的に確保する。これにより、空きメモリの有無に関係なくプロトコルメッセージだけは受信できることが保証される。

受信データの直接利用

大量データを効率よく転送するために、ネットワークドライバで扱うバッファをそのまま Mach-IPC のレベルまで渡せるようにする。ドライバレベルでは、カーネルが提供する一般的なメモリ割り当て手段を使うことはできないので、AST(Asynchronous Software Trap) レベルでこれを行なうこととする。AST レベルとは、カーネルレベルからユーザレベルに戻る時のことを指し、すべてのカーネルスレッドが休止した後に AST レベルのハンドラが動作する。

Message Accept 通知機能のサポート

Message Accept 通知機能をサポートするためにプロトコルを追加する。これにより、ポートキューレベルでのフロー制御も可能になる。

5 NORMA IPC/DE の設計および実装

前章の設計方針に基づいて、NORMA IPC/DE を設計し、実装した。ここでは、ターゲットである Cenju-3 のハードウェアについて簡単に説明し、その後で、NORMA IPC/DE の設計、実装について述べる。

5.1 Cenju-3

Cenju-3 は我々が開発した分散メモリ型並列コンピュータである。要素プロセッサには R4400 75MHz を用いており、256 台の最大構成時の性能は 38.4GIPS、12.8 GFLOPS である。

NORMA IPC を実現するにあたり最も関係が深いのは通信機能である。ネットワークのハードウェアは 4 段の多段接続網で、スループットは 40M バイト /s、1 段当たりの遅延時間は 250ns である。Cenju-3 では、ネットワークのハードウェアに見合う十分な性能を実現するための専用のインターフェースハードウェア（以下インターフェースと略す）を持つ。このインターフェースの特徴を以下に示す。

- 通信の種類は、メッセージ通信とリモート DMA が提供されている。
- プロセッサとは独立に、キューイングされたデータを自動的にパケットに分割してネットワークに送り出すことができる。パケットに分割される前のデータの最大長は 128K バイト、パケットの最大長は 508 バイトである。
- インターフェースからメモリへのアクセスは、物理アドレスに限られる。また、インターフェースに渡すアドレスもすべて物理アドレスを指定する必要がある。
- インターフェースが使用するメモリ領域は、ソフトウェアで明示的にキャッシュ制御を行なう必要がある。

5.2 NORMA IPC/DE のプロトコル

オリジナルの NORMA IPC に対して行なったプロトコルの変更を以下に示す。

第 1 に、メモリ予約のための送信要求メッセージと、メモリ割り当てができたことを知らせる送信許可メッセージを加えた。また、不要になった再送機能をプロトコルから取り除いた。

第 2 に、Mach の IO 処理に使われる iodone スレッド（後述）を送信側のページマッピングに流用したかったので、一度に送ることができる OOL ページの数を 8 ページに制限した。そして、8 ページ単位で OOL ページを送信するために、ページ送信許可とページコピー終了のメッセージを加えた。

第 3 に、小量の inline データ、すなわち KMSG のサイズが小さいメッセージに用いる簡単なプロトコルでは、送信要求メッセージに KMSG を含めて送ることにした。これにより、通常のプロトコルに比べ、送信許可と KMSG メッセージの 2 つを減らすことができる。Cenju-3 の 1 パケットのサイズが 508 バイトなので、これから送信要求の情報などを減じた 440 バイトまでの KMSG を持つメッセージをこのプロトコルの対象にした。

第 4 に、大量データ転送時のコピーを減らすために、OOL ページの転送にはリモート DMA を利用することにし、送信許可とページ送信許可メッセージの中には書き込み先を示す物理アドレスのリストを持たせた。

以下に、実装した NORMA IPC/DE のプロトコルを示す。（図 4 参照）

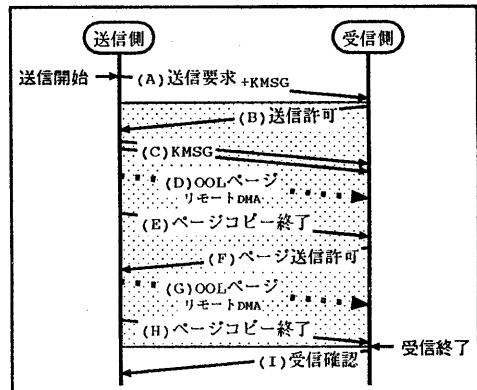


図 4: NORMA IPC/DE のプロトコル

1. 送信側で KMSG のサイズと OOL ページの数を計算し、それらを送信要求メッセージとして、受信側に送る。この時、KMSG のサイズが 440 バイト以下ならば、KMSG も付けて一緒に送る (A)。
2. 受信側は送信要求メッセージを受けとると、KMSG と OOL ページを受信するのに必要なメモリをすべて割り当てる。KMSG が含まれている時は、割り当てたメモリに KMSG の内容をコピーし、6 の処理へ移る。KMSG が含まれていない時は、送信許可メッセージを返送する (B)。
受信すべき OOL ページがある場合は、割り当てたページの先頭物理アドレスをリストにして、送信許可メッセージに含める。この時、OOL ページが 8 ページ以上ならば、8 ページ分だけのアドレスリストを含める。
3. 送信側は送信許可メッセージを受けとると、まず KMSG メッセージを送る (C)。128KB を越える KMSG は分割される。
OOL ページは、KMSG メッセージの送信終了後、受信許可メッセージに含まれるアドレスリストを基に、OOL ページをページ単位でリモート DMA により転送する (D)。OOL ページ（最大 8 ページ）の転送終了後、ページコピー終了メッセージを送る (E)。

4. 受信側は、KMSG メッセージを受けとると、その内容を割り当てたメモリにコピーする。

ページコピー終了メッセージを受けとると、まだ、受けとっていない OOL ページがある場合、次の 8 ページまでの物理アドレスのリストを含んだページ送信許可メッセージを送信側に返す (F)。
5. 送信側は、ページ送信許可メッセージを受けとると、3 と同様に OOL ページを転送し (G)、ページコピー終了メッセージを送る (H)。すべての OOL ページの転送が終るまで、4、5 を繰り返す。
6. 受信側はすべてのデータを受信し終ったら、受信終了メッセージを送信側に送る (I)。

このプロトコルでは、440 バイト以下の KMSG だけを含む場合は、2 メッセージ、440 バイトを越える KMSG だけを含む場合は、最低 4 メッセージが必要である。OOL ページを含む場合は、8 ページ毎に 2 メッセージずつ必要になる。

5.3 送受信バッファの管理

複数のノードからのメッセージを同時に受け付けるために、ノード毎に 256 個 (256PE 構成の場合) のプロトコルメッセージ受信用バッファを固定的に確保した。メッセージ受信バッファの他にもノード毎に必要な情報があるので、1 ノード当たり 716 バイト必要になった。256 台分だと 182K バイト程度になる。無視できるほどに小さいとは言えないが、全体のメモリサイズからすれば、この程度ならば問題ない。

また、KMSG や OOL ページを格納したバッファを Mach-IPC のレベルに直接渡せるように、Mach 標準のメモリ割り当て機能 (kalloc と vm_page_grab) を使った。これらのメモリ割り当て機能を使えば、カーネルレベルでもユーザレベルでもメモリをフリーすることができる。また、OOL ページ用のバッファは、ネットワークインターフェースハードウェアにより直接、書き込まれるため、すべてのページをページアウトされないように固定し、さらに、受信許可メッセージを返す前にすべてのページについてキャッシュフラッシュしている。

送信バッファに関しては、KMSG も OOL ページも Mach-IPC の処理部分すでに割り当てられているので、新たに割り当ては必要ない。ただし、リモート DMA を使うためには OOL ページに物理メモリをマッピングしなければならない。この処理はブロッキングの恐れがあるので、スレッドレベルで行なう必要がある。Mach では、ディスクからの DMA でマッピングの処理を行なう iodone スレッドが用意されているので、こ

れを OOL ページのマッピングの処理にも流用した。iodone スレッドによって行なわれるマッピングの処理は、8 ページに限られている。

さらに、OOL ページをリモート DMA する前にキャッシュインパリディートする必要がある。この処理には、時間がかかることが予想されていたので、1 ページ毎にキャッシュインパリディートとリモート DMA を繰り返し、複数ページの送信の際はこれらをオーバラップできるようにした。

5.4 Message Accept 通知プロトコル

Mach-IPC の Message Accept 通知機能を加えるために、新たに Message Accept 通知プロトコルを加えた。以下にプロトコルを示す。

1. 送信側では、シーケンス番号を KMSG に書き込んで送信する。
2. 受信側で、受信タスクがポートからメッセージを受信する時に、KMSG の中のシーケンス番号と送信側のポート番号を Message Accept リストに登録する。このリストは後で、AST レベルの Message Accept ハンドラが Message Accept メッセージとして送信側に送る。
3. 送信側では、Message Accept メッセージを受けてると、リストのポート番号に対応するポートに対して Message Accept 通知の処理を行ない、ポートのキューから Mach メッセージを取り除く。

このように、このプロトコルは送信側のポートのキューにつながるメッセージの管理も一緒に行なうので、ポートのキューレベルでのフロー制御もできるようになる。

6 評価

ここでは Cenju-3 上に実装した NORMA IPC/DE の定量的評価を示し、その考察を述べる。

6.1 小量データ転送のレイテンシ

小量データ転送のレイテンシを調べるために、4 バイトから 2K バイトまでの inline データについて、メッセージの往復にかかる時間を測定した。メッセージの送信には mach_msg_send、受信には mach_msg_recv を用い、mach_msg_send の前と mach_msg_recv の後で時間を計り、その差の半分を片道の所要時間とした。また、比較のために Mach-IPC でも同様の測定を行なった。これらの結果を inline データのサイズに対する所要時間のグラフ (図 5) に示す。ここで測定結果の精度は $0.1\mu s$ である。

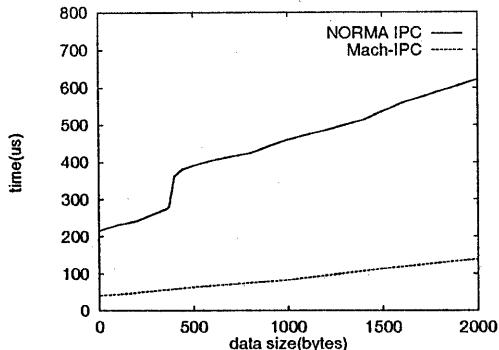


図 5: inline データのレイテンシ

このグラフを見ると、440 バイト付近でのレイテンシの跳ね上がりから、5.2節で述べた簡単なプロトコルと通常のプロトコルの差が読みとれる。その差は $100\mu s$ 程度である。この差は 2 つのメッセージ（送信許可と KMSG）を減らしたことによるものである。したがって、各プロトコル処理にかかる時間に大きな差はない仮定すると、1 つのメッセージの送信から受信までの時間、すなわち、メッセージの生成、キャッシュ操作、送信ハードウェアインターフェースの操作、ネットワークでの遅延、受信側プロセッサの割り込み処理、メッセージ受信後のプロトコル処理にかかる時間の合計は、およそ $50\mu s$ であると考えられる。

ここで、4 バイトデータの場合についてレイテンシの内訳を考えてみると Mach IPC の処理時間はグラフから約 $40\mu s$ であり、やりとりされるネットワークメッセージは送信要求と受信確認の 2 つであるから、計算上のレイテンシは、

$$\begin{aligned} & \text{ネットワークメッセージの送信から受信までの時間} \times 2 \\ & + \text{Mach-IPC の処理時間} \\ = & 50(\mu s) \times 2 + 40(\mu s) \\ = & 140(\mu s) \end{aligned}$$

程度になるはずである。実際にかかっている時間は約 $210\mu s$ であるので、 $70\mu s$ 程余計にかかっていることになる。いかにプロトコル処理に時間がかかったとしても、 $70\mu s$ は大き過ぎるので、何か他にレイテンシを増やす原因があると考えられる。

そこで、4 バイトデータの場合について、プロトコル処理にかかる時間の内訳を調べた。その結果を図 6 に示す。この図は送信側で `mach_msg_send` を呼んでから、受信側で `mach_msg_recv` が戻るまでの時間を、プロトコルの処理毎に示したものである。ここでの測定時間の精度は $1\mu s$ 程度である。

これを見ると前の計算では入っていなかった Message Accept 通知プロトコルの処理がレイテンシに含まれて

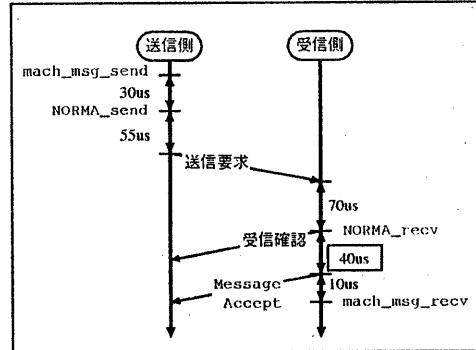


図 6: レイテンシの内訳

いることがわかる（図中四角枠の中の $40\mu s$ ）。Message Accept 通知プロトコルの処理は、受信タスクがポートからメッセージを受けとった後に行なわれるものなので、このレイテンシの計算には関係ないと考えていたが、実際にはこのような形で関わっていた。

6.2 大量データ転送のスループット

大量データ転送のスループットを調べるために、1 ページから 25 ページまでの OOL データについて、メッセージの往復にかかる時間を測定した。ページサイズは 16K バイトである。測定の方法は inline データのレイテンシ測定の場合と同じである。ここでの測定結果の精度は $0.1\mu s$ である。スループットは、測定結果の半分を片道の所要時間とし、送信したデータサイズをその時間で割って求めた。この結果を OOL データのページ数に対するスループットのグラフ（図 7）に示す。ここでは、平均のスループットとともに 1000 回測定した中の最大のスループットも併記する。

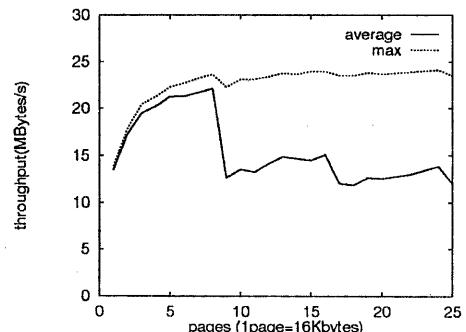


図 7: OOL データのスループット

最初に、平均のスループットを見ると、8 ページまではページ数が増えるほど高くなっていますが、8 ページ

の時で平均 22M バイト /s 程度である。これは、送信要求、受信許可、受信確認のメッセージによるオーバヘッドが、ページ数が増えるほど小さくなるためであると考えられる。

また、8 ページを越えると平均のスループットが激減し、16 ページ、24 ページを越えたところでも減少が見られるが、これには次の原因が考えられる。

1. 1 度のリクエストで送信できる量を 8 ページに制限しているので、8 ページを越えるとページ送信許可メッセージが必要になる。
2. 送信側では、8 ページの送信が終了した後で、次に送る 8 ページに物理メモリをマッピングするために iodone スレッドを呼んでいる。このスレッド内での処理がブロックすることがあり、最悪の場合、数十 ms の間、再スケジューリングされないことがある。

原因 1 は、8 ページ毎に 1 メッセージ分、すなわち、数十～数百 μ s の処理時間を増大させるだけであり、これは全体の所要時間（8 ページで 6ms 程度）の数 % に過ぎず、スループットに与える影響は小さい。一方、原因 2 は最小の処理時間の数倍に相当する処理時間を増大させることになるため、影響は極めて大きい。したがって、9 ページ以上の転送でのスループット減少の主な原因是、この予測できない iodone スレッドの待ち時間にあると考えられる。

また、9 ページ以上のスループットが 8 ページ単位で単調増加せず、上下しているのも、原因 1 によるものと考えられる。

一方、最大のスループットは、iodone スレッドによる処理時間が最小である場合と見なすことができる。このグラフから、8、16、24 ページのところに山があり、24 ページの時には 24M バイト /s 程度のスループットになることがわかる。したがって、iodone スレッドの影響がなければ、ページ数が増えるほどスループットが上がるものと予測できる。

7 高速化の検討

前節の評価を踏まえて、今後の高速化のための検討を行なう。

まず、レイテンシに関しては、Message Accept 通知のプロトコルメッセージを他のメッセージの一部に含めて送ることで、50 μ s 程度減少できると考えられる。Message Accept 通知のための情報は 1 つのポートについたかだか 2 ワードしか必要でないので、他のメッセージに含めて送ることができる。

また、スループットに関しては、iodone スレッドの代わりにページ送信専用のスレッドを用意することに

より、必要以上に長い時間ブロックしないようにすることが可能である。これにより、9 ページ以上の転送の平均のスループットを最大のスループットにかなり近づけることができると考えられる。

さらに、最大のスループットをあげるために、ページマッピングの単位を大きくする方法がある。しかし、これは一時的とはいえユーザタスクが利用できる物理メモリを減少させてるので、メモリサイズとの兼ね合いでマッピングの単位を決めなければならない。

8 まとめ

Mach のネットワーク IPC である NORMA IPC をベースに、新たに NORMA IPC/DE を構築し、分散メモリ型並列コンピュータ Cenju-3 上に実装した。NORMA IPC/DE は Mach IPC と同じセマンティクスを持ち、複数のノード間通信の同時確立をサポートし、さらに、大量メッセージについては高スループット、小量メッセージについては低レイテンシを実現している。

まだ性能面では改良の余地はあるものの、現時点で、小量メッセージ転送ではレイテンシ 210 μ s、大量メッセージ転送では最大スループット 22M バイト /s の性能を出している。

今後はさらに 7 節で述べた高速化を行ない、再評価を行なう予定である。

参考文献

- [1] 広瀬、加納、丸山、中田、浅野、稻村、 “並列コンピュータ Cenju-3 のアーキテクチャ”， 情報処理学会 計算機アーキテクチャ研究会報告 Vol.107-16, pp.121-128, Jul. 1994.
- [2] 高野、小西、C. Howson、荒木、菅原、小長谷、 “Mach マイクロカーネルをベースとした並列 OS DenEn の実現”， 情報処理学会 '95 全国大会 予稿集.
- [3] Joseph S. Berra III, “A fast Mach network IPC implementation”, In Proceedings of the Usenix Mach Symposium, pp.1-12 Nov. 1991.
- [4] David B. Golub, Randall Dean, Alessandro Forin and Richard Rasid, “Unix as an Application Program”, USENIX Summer Symposium, Jun, 1990.