

PUMA-III における 1 Gbps FC ネットワークの実現技術

5 B-1

新家 正総 陣崎 明

(株)富士通研究所

1. はじめに

ワークステーション (WS) をネットワークで結合し分散/並列処理を行う WS クラスタが注目されている。高速標準ネットワークの普及に伴い、WS クラスタへの期待はより高まっている。

我々は標準ネットワークによる高性能 WS クラスタの実現を目指し、1Gbps ファイバーチャネル¹ (FC) を用いた WS クラスタ PUMA-III² を開発している。

高性能な WS クラスタ実現の鍵はプロセス間通信の高速化である。しかし標準ネットワークの利用では独自ネットワークと比較してオーバーヘッドがかかる点が問題となる。PUMA-III の最終的な課題は、標準ネットワークを用い、なおかつ高速なプロセス間通信を実現することである。

今回 PUMA-III 開発の第一ステップとして、FC の機能のうち FC1 以下のレイヤを利用するネットワークアダプタを開発した。本稿ではこのアダプタハードウェアに焦点をあて、検討と開発の結果を述べる。

2. アダプタハードウェアの検討

WS クラスタのプロセス間通信はスループットと遅延の両面で高速性を実現しなければならない。そこでアダプタハードウェアの要件としてはアダプタ内部の処理の高速化とアダプタインタフェースの高速化が重要である。

しかしながらこれらの要件を実現する上で、アダプタを実現する方式には様々な選択肢がありうる。まずアダプタ内部の処理方式については、おおまかには制御 CPU を用いるか専用ハードウェアで制御するかを選択がある。またアダプタインタフェースの方式でも、まず DMA とプログラム I/O の選択肢があり、更に DMA の場合には DMA の終了認識の方式としてポーリングと割り込みの選択肢がある。以下ではこれらの点のうち、DMA とプログラム I/O の選択に焦点を絞って考える。

DMA は CPU の負荷なしに効率良く高速転送が可能である。しかしメッセージが非常に短く DMA 起動処理時間が無視できない場合や、転送データ

が DMA 領域に置けない場合はプログラム I/O の方が性能が優るかもしれない。例えばイリノイ大のメッセージ通信方式である Fast Messages³ (FM) では、短いメッセージの性能に重点を置き、送信にプログラム I/O を用いている。しかしながら、どの程度の長さのメッセージにおいてプログラム I/O より DMA が劣るかは、DMA の起動方式、I/O アクセス性能、DMA 性能、ソフトウェアオーバーヘッドなどに依存するため一般に明らかではない。

我々は今回 DMA 方式で SBus を用いた FC アダプタを開発した。以下では SBus FC アダプタについて概要を述べ、DMA がどの程度有効なのか評価した結果を述べる。

3. SBus FC アダプタ

SBus FC アダプタ (図 1) は、SBus の DVMA 空間上のデータを DMA で送受信する。DMA の起動、終了確認は以下のようにして行う。

送信の場合はバッファアドレスと転送サイズをアダプタにセットし DMA 起動をかける。

受信の場合は複数の DMA アドレスをアダプタが保持できるように受信アドレス FIFO を用意した。受信アドレスは 64 個まで設定できる。また受信済バッファごとに DMA 末尾アドレスがわかるように受信 RESULT FIFO を設けた。受信認識はポーリングで行う。

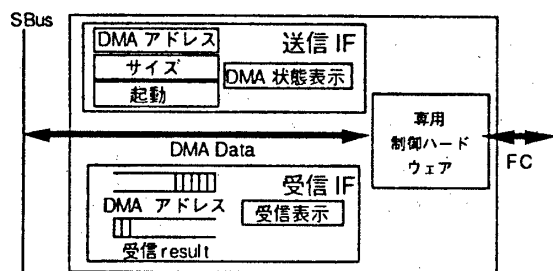


図 1 SBus FC アダプタ

アダプタの内部は専用の制御ハードで構成した。制御ハードはファイバーチャネルのプロトコル階層における FC0 (物理インタフェースとメディア)、FC1 (encode, decode など) の機能をもつ。パケット形式は FC2 と異なる独自方式であり、ホストから受け取ったデータにパケットの先頭、末尾の識別子を付加してネットワークに転送する。

4. SBus FC アダプタの評価

アダプタの評価は、1つの送信バッファから同じデータを送り続ける試験(独自プロトコル⁴による試験)とイリノイ大のFast Messagesを我々のハードウェア上に移植したFC FM⁴を用いて行った。イリノイ大のFMは送信にプログラムI/Oを用いているが、我々は新たにユーザ空間からの直接DMAで実装している。

独自プロトコルではソフトウェアオーバーヘッドがほとんど無い。送信起動のための3wordのI/Oアクセスと、ポーリングだけがオーバーヘッドとなっている。FMではシーケンス番号管理などのソフトウェアオーバーヘッドが若干必要である。

測定に用いたWSはSPARCstation2(CPU Weitek倍速、SBus 58MB/s)とSPARCstation5(CPU microSPARC、SBus 35MB/s)である。比較として用いたMyrinet³によるイリノイ大FMはSPARCstation20を使用している。

図2に各種プロトコルのスループットを示す。独自プロトコルではSS2で最大52MB/s、SS5で28MB/sの性能を得た。SS5の方が遅い理由は2つある。1つはSS5のDMAの最大バーストサイズがSS2の半分の32byteのためである。もう一つはSS2ではDMA転送中はCPUからのポーリングがほとんど入らないのに対して、SS5ではポーリングでDMAがブロックされているためである。

次にFC FMとイリノイ大Myrinet FMを比較する。SS2ではほぼ512byte以上、SS5ではほぼ256byte以上でFC FMがMyrinet FMより高性能であり、2KBでの性能はそれぞれ40MB/s、25MB/sとなる。またMyrinet FMより性能の低い領域でも、SS5ではほぼMyrinet FMに近い性能を出せている。このことから、256byte以下の小さいメッセージサイズでも、DMAでプログラムI/Oに近い性能が出せることがわかる。

図3にFC FMの遅延を示す。SS5でMyrinet FMとほぼ同等の性能が得られた。SS2ではCPU性能が低いために小さなメッセージサイズでMyrinet FMより遅延が大きくなっているが、512byte以上ではDMAの効果により逆にMyrinet FMよりも小さな遅延となっている。

5. まとめ

DMA方式のSBus FCアダプタはSS5で最大25MB/s以上の性能が得られ、256byte以下の小さなメッセージにおいてもイリノイ大FMに近い性能が得られる。またSS2での最大性能は40MB/s以上である。これらの結果からDMA方式のアダプタはWSクラスタのネットワークアダプタとして

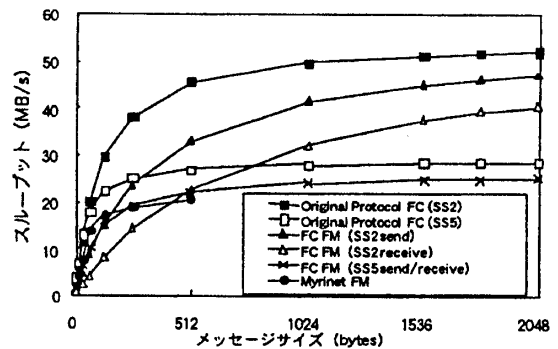


図2 スループット性能

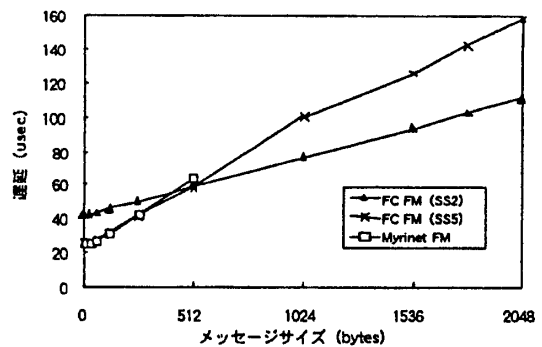


図3 遅延

有望なことがわかった。

6. 今後の予定

今回開発したSBus FCはFC1レイヤまでのサポートになっており、FC2以上の複雑な機能を実現することができない。そこでFC2より上位のプロトコル処理をハードウェアで行う処理エンジンを開発中である。プロトコル処理ハードウェアでは、FCだけではなくTCP/IPなどの処理をアシストすることも考えている。

またFC1以下のアダプタとしても、SBusのボトルネックを解消するためにPCIバスを用いたFCボードを開発中である。PCIの場合理論性能で125MB/s程度が得られ、より高速な通信性能が期待できる。

参考文献

- [1] FIBRE CHANNEL: PHYSICAL AND SIGNALING INTERFACE (FC-PH) REV4.3, ANSI, June 1994
- [2] 新家他: PUMA-III: 1Gbps ファイバーチャネルを用いた高速分散共有メモリシステム, 第50回情報処理学会全国大会, 1H-10, (1995-3)
- [3] Scott Pakin 他: High Performance Messaging on Workstations: Illinois Fast Messages (FM) for Myrinet, Supercomputing, Dec. 1995
- [4] 小林他: PUMA-IIIにおける各種メッセージプロトコルの実装と評価, 本大会予稿