

100Mb/s 光ファイバリングネットワーク を用いたプロセッサ間結合方式

木村 行男、星子 隆幸
NTT 電気通信研究所

センタシステムにおける計算機の処理量増大、システムの段階的成長、高信頼化等の要望に応じて、多数台の疎結合マルチプロセッサシステムを構築するため、伝送速度100Mb/sの光ファイバリングネットワークを開発した。

本ネットワークのノードであるプロセッサ間結合装置の通信効率に関連する特徴的な機能とその適用効果は次の通りである。(1)複数サブチャネルを用いて、個々の通信プロセス間で直接転送を行うことにより、通信処理プログラムのダイナミックステップ数を削減(2)通信待ち合せ機能適用によるアテンション割込み処理時間の削除(3)分割転送方式により、長データをチャネル速度で転送し、プログラムに対するデータ長制限を撤廃

本稿では、上記プロセッサ間結合装置の機能およびPCIの実効データ転送能力を紹介し、上記効果を上げるためにソフトウェアで対処すべき事項を述べる。

Processor to Processor Communication System Using a 100Mb/s Optical Token Ring Network

Yukio Kimura and Takayuki Hoshiko
NTT Electrical Communications Laboratories
1-2356, Take, Yokosuka-shi, Kanagawa-ken, 238-03 Japan

For the purpose of constructing a loosely coupled multiprocessor system which can meet rapid increases in demands for increased processing power, step-by-step growth of systems, and higher reliability in center systems, we have developed a 100Mb/s optical fiber ring network.

Principal functions of this ring network node (processor-to-processor communication interface unit : PCI) and its application effects are as follows;

- (1)Reduction of dynamic steps of a communication control program by communicating directly between any pairs of communication processes using multiple subchannels.
- (2)Deletion of time for processing attention interruption by preissuing a communication synchronizing command.
- (3)Reduction of time for transferring long data and removal of data length limit by using the split transfer method.

This paper describes the above functions, actual data transfer ability, and ways of deriving the effects of the above functions by software.

1. はじめに

データ通信サービスの高度化、システムの段階的成長にともない、センタシステム内のプロセッサ（HOST、FEP等）間の高速な接続が要求されるようになってきている。この要求に応じて多数台の疎結合マルチプロセッサシステムを構築するため、伝送速度100Mb/s光ファイバリングネットワークを開発した。

本稿では、100Mb/s光ファイバリングネットワークのノードであるプロセッサ間結合装置（PCI；Processor-to-processor Communication Interface unit）について、その特徴的な動作の概要および性能評価結果を紹介する。またこのPCIを用いて通信性能向上を図るためにプログラムで対処すべき事項について述べる。

2. 100Mb/s光ファイバリングネットワークの概要

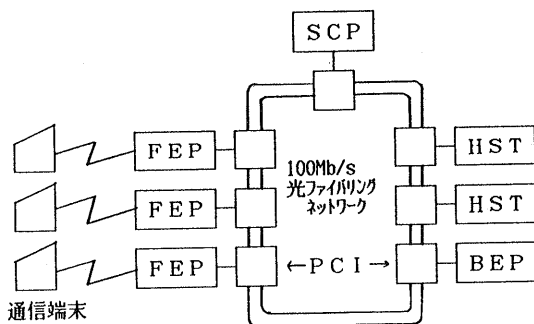
本稿で述べる100Mb/s光ファイバネットワークは、以下のねらいのもとに開発した。

- (1) センタシステムの処理能力拡大に対して柔軟に対処可能なシステム拡張性の提供
- (2) 経済的なセンタ内プロセッサ間通信の実現
- (3) センタシステムの長時間連続運転とシステム障害時の高速リカバリの実現

本ネットワークを用いたセンタシステム構成例を図1に示す。

2.1 主な物理仕様及びデータリンク仕様

本ネットワークのネットワークノード（ここではPCIとよぶ）とプロセッサとの接続については、中大型汎用計算機において多く使われている高速で標準的なチャネルインタフェース（I/Oインタフェース）を用いている。ネットワーク伝送路については、伝送速度100Mb/sで、センタ内接続の他、隣接センタ間接続も考慮してPCI間距離が



PCI：プロセッサ間結合装置

SCP：システム制御用プロセッサ

HST：ホストプロセッサ

FEP：フロントエンドプロセッサ（通信制御用プロセッサ）

BEP：バックエンドプロセッサ

図1 光ファイバリングネットワークによるセンタシステム構成例

最大2kmまで可能な光ファイバ伝送方式を用いている。ネットワークトポロジーは、100Mb/sの高速転送を多数台（数10台以上）のPCI間で光を用いて信頼性良くかつ経済的に実現可能な、リング形式を用いた。ネットワークアクセス方式は、伝送効率の良いトークンパッシング方式を採用している。

リングネットワークは、PCIまたは伝送路障害がシステム全体の障害につながる性質をもっているため、本ネットワークでは伝送路を2重化するとともに代替バス切り替え、ループバック切り替えによる障害箇所自動切り離し、送信権自動回復を行っている。これらの動作は特定のPCIによる集中制御ではなく完全分散制御であり、ハードウェアで自動的におこなわれる。

2.2 プロセッサ間通信方式

高速のプロセッサ間通信実現のためには、伝送路の効率を向上させるだけでなく、伝送路を流れるトラフィックの制御の改善、すなわち通信制御に要するオーバヘッドを削減し、ハードウェア・ソフトウェア処理双方の高速化が必要である。そのため、本方式では通信制御に要するソフトウェア処理高速化に有効

な機能をPCIでサポートしている。

この機能は、複数サブチャネルを用いた多重通信方式による通信プロセス間直接転送機能であり、1対1チャンネル結合のプロセッサ間通信で使用されている機能をN対Nプロセッサ間通信用に拡張したものである。(図2参照)

以下に本機能ならびに本機能をN対N環境で用いるための通信先アドレス指定およびデータ転送方法、分割転送方式等特徴的な機能について述べる。

(1) サブチャネル多重による通信プロセス間直接転送機能

本機能は、通信プロセス毎にサブチャネルを用いて通信バスを特定し(サブチャネル多重通信方式)、送受信バッファを個々の通信プロセス側で用意して、該バッファ間でのデータの直接転送を行うものであり、同一送受信バッファの連続使用が可能のため、バッファエリアのフィックス化/フロート化を最適な契機で実施することが可能である。すなわち、高トラヒック、長データ転送の環境下で送受信バッファ管理や割り込み処理、データ移送等の通信処理に要するダイナミックステップ数を削減することができる。

なお、多重通信の効率化のために、各サブチャネル対応にチャンネルを早期に解放し、アテンション割り込みを伴わずに相手からの通信起動を待ち合わせる通信待ち合わせ機能も併せてサポートしている。

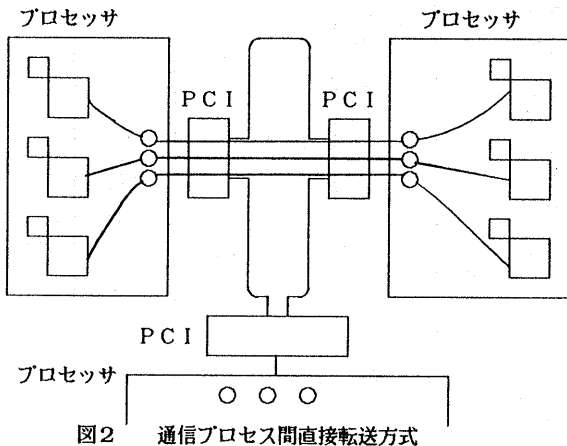


図2 通信プロセス間直接転送方式

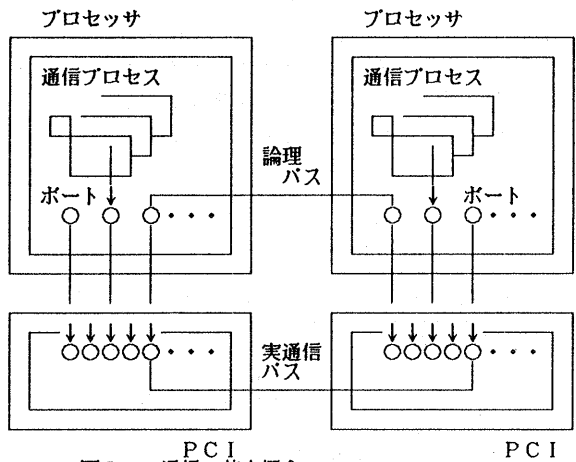


図3 PCI通信の基本概念

(2) 通信先アドレス指定

サブチャネル多重による通信プロセス間直接転送機能をN対N通信環境で用いる場合の対処として通信相手のノードアドレスとサブチャネルアドレスを、チャンネル仕様下で効率的に指定する通信先アドレス指定法を以下に示す。

本プロセス間通信は、図3に示すように通信プロセス間の仮想回線である論理バスと論理バスの窓口であるポート及びPCI間のデータ転送を行う実通信バスを用いて行われるものとする。個々の通信プロセスは論理バスを指定して通信を行うが、通信に先立って論理バスを開設する必要がある。各プロセッサ内に論理バスを開設するための窓口として、ポートを設定する。論理バスはポート名の対で意識される。ポート名はPCIアドレス+ポート番号で指定される。PCI間のデータ転送は実通信バスを指定して行われる。実通信バスは送受信双方のPCIアドレス+PCIサブチャネルアドレスで識別される。実通信バスの指定にあたり、現行仕様では自PCIのサブチャネルをデバイスとして指定できるのみなので、前もって実通信バスを設定し、自PCIサブチャネル指定のみで通信できれば効率的である。このため、実通信バスを記憶し管理するバス管理テーブル(CTBL)をPCI内に設置した。通信時の

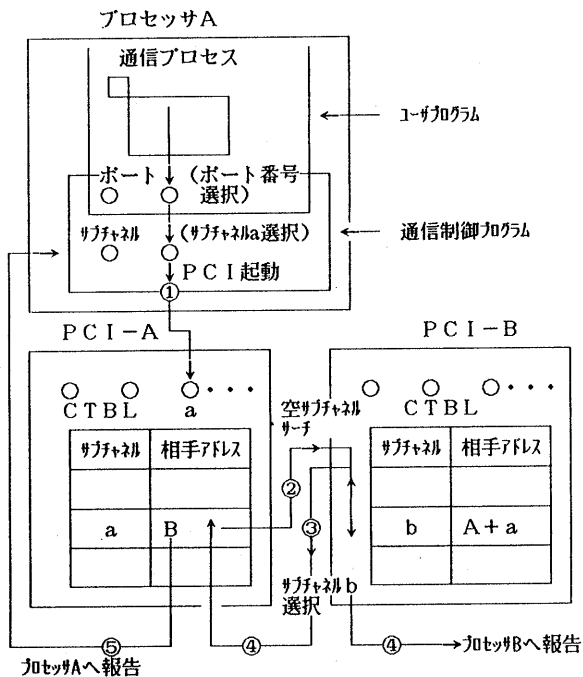
前処理として、専用コマンドによりCTBLを設定し実通信バスを確立する。通信をとまなうコマンド（Writeコマンド、Readコマンド等）発行時、相手アドレスはCTBL検索によりPCIハードが自動設定する。この実通信バス設定の動作例を図4に示す。

(3) データ転送

チャンネル結合されたプロセッサ間通信においては、送受信双方のプログラムが発行する通信用コマンドで同期をとってデータ転送が行われる。このため、PCI-PCI間のデータリンク手順としては、単に上記通信用コマンド、データ、関連応答情報のみを伝達するだけの手順を用いている。PCIで採用している通信用コマンドレベルのデータ転送動作例を図5に示す。

(4) 分割転送方式

PCI-PCI間のデータ転送は、長データ時のメッセージ通過時間を短くするために、分割転送方式を用いている。分割転送方式は、図6に示すようにPCI内の交替バッファを用いて、プロセッサ-PCI間のデータ送受信とPCI-PCI間のデータ送受信を並行動作させるものである。転送データ長が長くなると、同期化に要する時間がデータ転送時間に比べて無視できる程小さくなるため、チ



(実通信バス設定手順)

- (1) 通信制御プログラムがPCI-Bとの実通信バス設定用に選択したPCI-Aのサブチャンネルaに対して通信バス設定コマンド（SPA）により、相手PCIアドレス情報Bを転送する。
- (2) サブチャンネルaに対応するCTBLエントリの相手アドレス域にBを設定し、上記SPAをPCI-Bに転送する。
- (3) SPAを受信したPCI-Bは、CTBL内で相手アドレスが確立していない任意のサブチャンネルbを選択し、bの相手アドレス域にA+aを設定する。
- (4) PCI-Bのサブチャンネルbに対応する相手アドレスが確立したことをPCI-Bの上位プロセッサおよびPCI-Aに報告
- (5) PCI-AはPCI-Bからの上記報告を基に、サブチャンネルaの相手アドレスとしてbを追加設定し、一連のSPAコマンド処理が完了したことをPCI-Aの上位プロセッサに報告する。

図4 実通信バス設定動作例

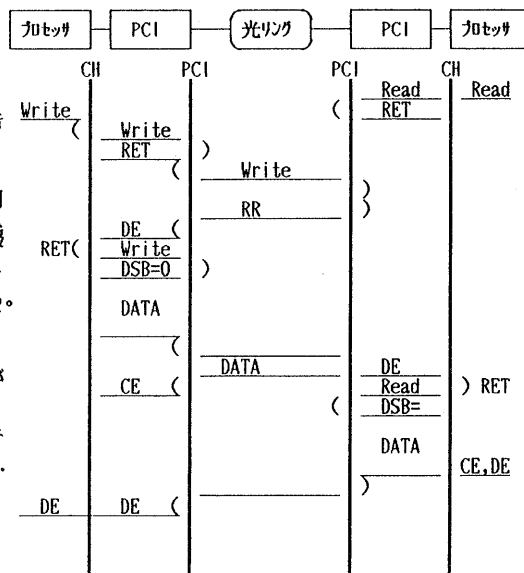
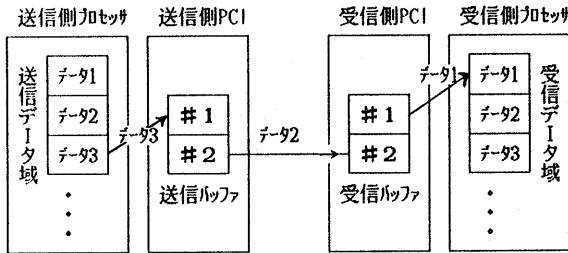


図5 データ転送動作

チャンネル転送速度でデータ転送を行うことが可能になる。この方式は、ハードウェアが一定長にデータを分割するために、プログラムが指定するデータ長に制限をつける必要がない。



(注) PCI内バッファは交替バッファ方式をとっている。

図6 分割転送方式

3. PCI性能

PCIの性能として、送受信一对のPCI間での単位時間内の転送データ長をとると、性能は次式で近似される。

PCI性能 =

$$\frac{\text{転送データ長}}{\text{ソフト走行時間} + \text{PCIハード動作時間}}$$

ここで転送データ長とは、1回のコマンドで転送されるデータの長さである。またPCIハードウェア動作時間(T)は、1コマンドを処理する際にPCIのハードが占有される時間であり、転送データ長の関数として次式で示される。

$$T = A + B + C + D$$

- A: PCIファームウェア処理時間
- B: チャンネル転送時間 (転送データ長に比例)
- C: 伝送路送信待ち時間
- D: リング伝送路遅延時間

また、ソフト走行時間は送受信の際にPCIハードと直列に走行する通信制御プログラムの走行時間を示す。

上式に基づいたPCI性能算出結果を図7、

図8に示す。

図7は転送データ長を変化させた場合の性能との関係である。チャンネルとPCI間のデータ転送時間を除くPCIハード動作時間が変わらないため、転送データ長の増加に伴い性能は向上し、各々チャンネル転送速度で飽和する。

図8は短データ(200B)転送時のソフト走行時間に対するPCI性能との関係である。高速のプロセッサの場合は、ソフト走行時間増加の影響は少ないが、低速プロセッサでは影響が大きい。低速プロセッサにPCIを接続する場合にはソフトウェアダイナミックステップ数の削減が重要であることを示している。

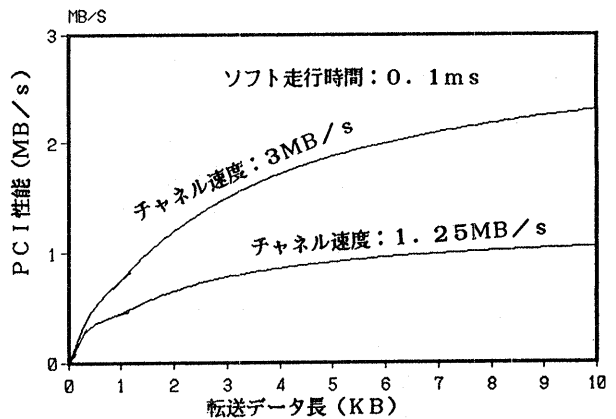


図7. 転送データ長に対するPCI性能

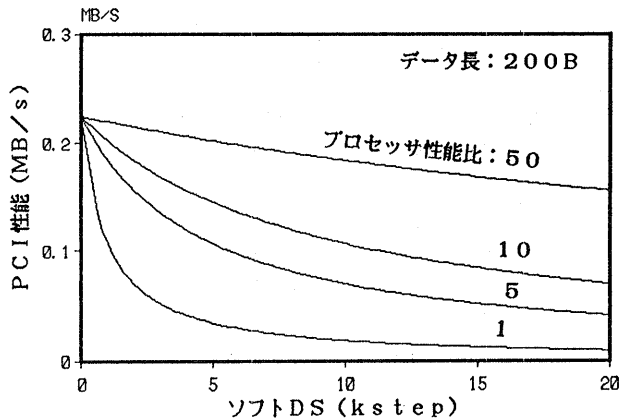


図8. ソフトDSに対するPCI性能

4. 光ファイバリングネットワークを用いたプロセッサ間結合システムの効率的な使用方法

PCIを介した通信の効率化を図るため、以下の点に留意する必要がある。

(1) Writeコマンド発行回数の削減

PCIでは、性能向上の点から、1 Writeコマンドあたりのデータ長は、なるべく長い方が好ましい。分割転送方式の採用により、データ長の増大によるPCIハード動作時間の増大はチャンネル転送時間の増大分だけである。短い電文を多数転送する場合には、短電文を一まとめにして1回のコマンドで転送するデータチェーンを使用したほうがよい。

この多数の短データを一まとめて転送する効果を、転送データ長をパラメータとして図9に示す。

(2) Readコマンドの先行発行

PCIでは、アテンション割り込みをなくし、複数サブチャンネルを用いた効率的な多重動作を行うためにコマンド待ち合わせ機能を用いている。このため、ReadコマンドをWriteコマンドに先行して発行できるプログラムの作りしておくことが有効である。すなわち、各サブチャンネルに対し最初に先行Readをまとめて発行し、各Read動作が終了次第すぐに先行Read発行を行うようにすることで効率化が図れる。

(3) 多数サブチャンネルの使用

2. 2節で述べたように、多数サブチャンネルの使用によって通信処理に要するダイナミックステップ数の削減が可能である。

またPCIハードウェアでも、チャンネル早期解放により、動作中のチャンネル空き時間(コマンド発行時のRET(再試行要求)からRR受信時まで、及びチャンネルデータ転送後のCEからACK受信まで:図5)を他サブチャンネルの動作に用いることにより、多重動作が可能である。

従って、PCIを用いてプロセッサ間通信を行う場合、多数のサブチャンネルを使用することで、通信効率向上が図れる。

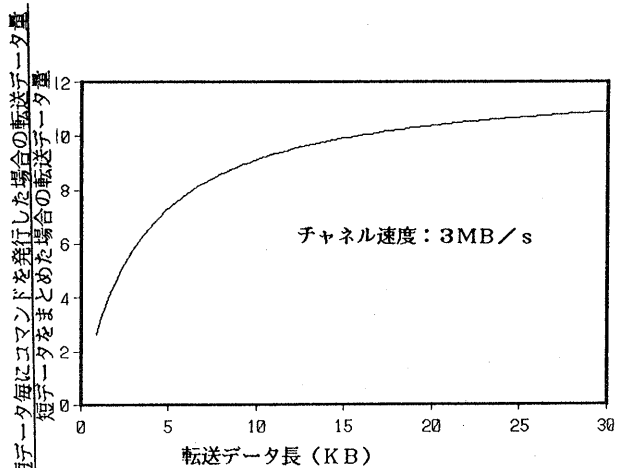


図9. 短データをまとめて転送する効果

5. おわりに

100Mb/s光ファイバリングネットワークを用いたプロセッサ間通信における特徴的な動作(サブチャンネル多重、長データ分割転送等)および性能について述べた。また、これら機能を生かすための使用方法を述べた。

本ネットワークは、既に十数システムのDIPSセンタに導入され、商用に供されている。

参考文献

- (1) Stuck, B.W.: Calculating the Maximum Mean Data Rate in Local Area Networks, Computer, Vol.16, No.5, pp72-76(1983)
- (2) Bux, W.: Local-Area Subnetworks: A Performance Comparison, IEEE Trans. Commun., Vol. COM-29, No.10, pp1465-1473(1981)
- (3) 小柳津他: DIPS複合構成システムプロセッサ間結合装置、通研実報, 35, No. 3, pp283-293(1986)

(4) 星子他：100Mb/s 光トークンリングを用いたプロセッサ間結合システム，情処論文誌、第27巻、第4号、435頁-444頁（1986）

(5) 中野、森：疎結合計算機システムにおける高速計算機間通信方式、情処第32回計算機アーキテクチャ研究会（1981）

(6) 星子他：共通バスを介したプロセッサ間通信効率化の検討、情処第24回全国大会、6H-6（1982）