

PIE64 のネットワーク・インタフェース・プロセッサの概要

5N-5

小池 汎平, 清水 剛, 島田 健太郎, 田中 英彦  
(東京大学 工学部)

1 はじめに

並列推論マシン PIE64q[1] のネットワーク・インタフェース・プロセッサ (Network Interface Processor: NIP) は、各推論ユニット (Inference Unit: IU) 内部と、相互結合網[2]とのインタフェースを行ない、データ転送やプロセス間同期などのプリミティブ・オペレーションをネットワークを介して2台のIU間で実行することにより、PIE64で行なわれる並列推論処理のうち、並列処理機能を支援する。NIPはマスタ部とスレーブ部にわかれる。マスタNIPは、UNIRED及びSPARCから、並列処理用プリミティブ・オペレーション実行の要求を受け付け、相手IUまでのネットワークの経路を接続したのち、接続先IU内のスレーブNIPと協調して、プリミティブ・オペレーションを実行する。PIE64は2系統のネットワークを持ち、各IUには、それぞれのネットワークのために、2系統のNIPが用意されている。本稿では、現在検討を進めているPIE64のNIPの機能について述べる。

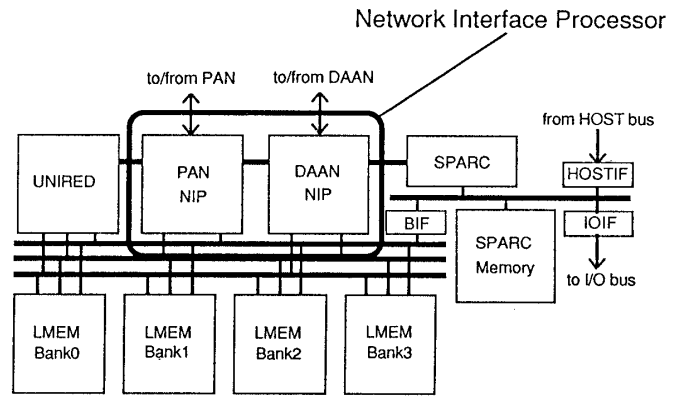


図1: Network Interface Processor

2 NIPの機能

並列推論マシン PIE64 のネットワーク・インタフェース・プロセッサ (Network Interface Processor: NIP) は、各推論ユニット (Inference Unit: IU) 内部と、相互結合網とのインタフェースを行ない、データ転送やプロセス間同期などのプリミティブ・オペレーションを2台のIU間でネットワークを介して実行することにより、PIE64で行なわれる並列推論処理のうち、並列処理機能を支援する。図1に示すように、NIPは相互結合網と、推論ユニットの内部バスに接続され、IUのローカル・メモリを直接アクセスし、ネットワークにデータを送り出す。又、NIPは、コマンドバスを通じてUNIRED/SPARCより、プリミティブ実行のコマンドを受け取る。

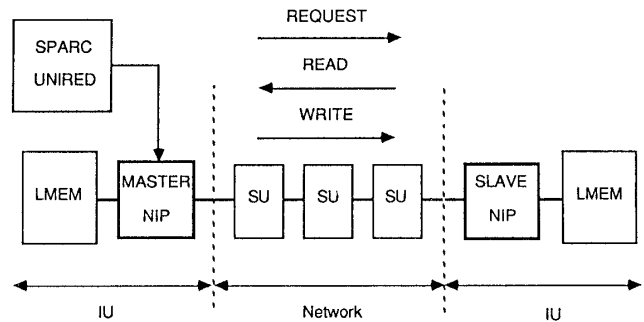


図2: NIPの役割

NIPは、マスタ部とスレーブ部にわかれる。図2に示すように、マスタNIPは、UNIRED/SPARCより様々な並列処理用プリミティブ・オペレーション実行の要求を受け付け、相手IUまでのネットワークの経路を接続したのち、接続先IU内のスレーブNIPと協調して、各種の並列処理用プリミティブ・オペレーションを実行する。マスタNIPとスレーブNIPは、大部分のハードウェアを共通化できるので、同一チップをモード切替えて動作させる。PIE64は2系統のネットワークを持ち、各IUには、それぞれのネットワークのために、2系統のNIPが用意されている。

3 NIPのプリミティブ・オペレーション

NIPのプリミティブ・オペレーションは、基本的に、

<sup>0</sup>The Network Interface Processor of PIE64

Hanpei KOIKE, Tsuyoshi SHIMIZU, Kentaro SHIMADA  
Hidehiko TANAKA, the University of Tokyo

- データ転送用
- プロセス間同期用
- その他

に分けられる。

データ転送用プリミティブ・オペレーションとしては、

- read1 read2 readn readx
- writel write2 writen writex
- writel1 writel2 writeln writelx

がある。read、write、writelは、それぞれ、相手IU内データの読み出し、相手IUローカルメモリへのデータの書き込み、最小負荷IUへのデータの書き込み(負荷分散)、を表す。また、末尾の1、2、n、xは、転送するデータのサイズを表し、それぞれ、1ワード、2ワード、データの先頭からサイズを読み出す、プリミティブ・オペレーションのオペランドとしてサイズを直接指定する、ことを意味する。前3者は、それぞれ、FLENGの、変数型、リスト型、ベクタ型に直接対応する。

基本的には、データ転送は、コピー元、コピー先のメモリアドレスを指定した、2つのIUのローカルメモリ間でのメモ

り内容のブロック転送である。ただし、1ワードの読み出しでは、読み出した値が、コマンドバスを介して、読み出し要求元の UNIREG 又は SPARC に直接返され、又、1ワードの書き込みでは、書き込む値は、書き込みコマンドとともに、コマンドバスを介して、直接 NIP に渡される。

プロセス間同期用のプリミティブ・オペレーションとしては、

- suspend
- bind
- activate

が用意され、FLENG の持つプロセス間同期機構を、ハードウェアレベルで直接サポートする。これらのプリミティブオペレーションは、互いに密接に関連して動作し、以下に示すように、場合によってお互いに他を呼び合うことがある。

suspend は、ある IU 内で、あるコンテキストの処理が、別の IU 内に置かれた論理変数の値が定まっていなかったためにサスペンドを起した時、その論理変数の置かれた IU に対して、サスペンドしたコンテキストのアドレスを送り付け、未束縛論理変数の保持するサスペンションリストにサスペンドしたコンテキストを登録するためのオペレーションである。もし、suspend オペレーションの実行前に、既に別のゴールによって論理変数に値がバインドされていた場合は、スレーブ NIP は、送られてきたコンテキストをサスペンションリストに登録する代わりに、そのコンテキストに対する activate 動作を開始する。

bind は、ある IU 内で Active Unification が行なわれたときに、他 IU 内に置かれた論理変数に対し値をバインドするためのオペレーションである。値をバインドするに先立ち、スレーブ NIP は、その論理変数のサスペンションリストに登録されている個々のコンテキストに対し、activate 動作を行なう。もし、論理変数にバインドする値が、他の未束縛の論理変数であった場合は、2つの論理変数の保持するサスペンションリスト同士をマージする必要があり、このために、スレーブ NIP は、バインドする論理変数に対し、自 IU 内の論理変数が保持するサスペンションリストを用いた suspend オペレーションを起動する。

activate は、論理変数のサスペンションリストに登録されている個々のコンテキストの置かれた IU に対し、論理変数にバインドされた値を送りつけ、ゴールの管理を行なうプロセッサに対し割り込みをかけるオペレーションである。

NIP には、以上で示した基本的なプリミティブオペレーションの他に、モード設定や、タイムアウト時間などのパラメータの設定のためのコマンドが用意される。

#### 4 NIP の通信プロトコル

NIP による1回のプリミティブオペレーションの実行は、

- UNIREG 又は SPARC からのコマンドの到着
- マスタ NIP による前処理
- ネットワークの接続
- ネットワークを介したデータのやりとり
- ネットワークの解除
- マスタ NIP 及びスレーブ NIP による後処理

の順で行なわれる。

ネットワークの接続時に、マスタ NIP は、最初に、ネットワークに対しプリミティブ・オペレーションのオペランドであ

る相手 IU 内のデータを指すポインタを送出する。この最初の1ワードの送りで、ポインタの各フィールドは、次のように用いられ、可能な限りの情報が相手プロセッサに渡される。このポインタのプロセッサ番号フィールドは、ネットワークは、経路を接続するのに用いる。ネットワークが接続されるとともに、ポインタはそのまま相手 IU のスレーブ NIP に転送され、ポインタのデータ型を表すタグ部と GC 時に用いられるマーク部のフィールドは、実行すべきプリミティブ・オペレーションを、ローカルメモリアドレスフィールドは、スレーブ NIP のアドレスレジスタに設定される。以上の操作は、ネットワークに衝突がなければ、ネットワークの接続に1-2クロック、接続の確認に1クロック、合計2-3クロック程度で実行できると考えられる。しかし、ネットワーク上で接続要求の衝突が頻繁に起きる場合は、接続要求は、一定時間でタイムアウト処理する。ネットワーク接続の後、基本的には、1ワード/クロックで、データの転送を行ない、プリミティブ・オペレーションの実行が行なわれる。この間、マスタ/スレーブ両 NIP 内の制御回路は、1クロック遅延のあるハンドシェイクによって、同期がとられる。又、PIE64 のネットワークは双方向なので、1クロックのオーバーヘッドでネットワークの方向を逆転させることも可能である。最後に、ネットワークの解放が行なわれ、これには1クロック要する。

#### 5 おわりに

本稿では、並列推論マシン PIE64 のネットワークインタフェースプロセッサの機能について述べた。現在詳細な仕様の検討を進めている。今後に残された課題としては、

- データ転送エラーに対する信頼性向上のための検討などが挙げられる。

#### 参考文献

- [1] 小池, 田中, “並列推論マシン PIE64 の概要”, 本大会.
- [2] Koike, H., Takahashi, E., Yamauch, T. and Tanaka, H.: *The High Performance Interconnection Network of the Parallel Inference Machine PIE64*, Proc. of Computer Architecture Symposium, IPSJ, May, 1988.