

高並列推論エンジンPIEにおける相互結合網の構成

Interconnection Networks on Highly Parallel
Inference Engine - PIE -坂井 修一^ω 田中 英彦^ω 元岡 達^ω
Shuichi Sakai Hiehiko Tanaka Tohru Moto-oka^ω 東京大学 工学部
Faculty of Engineering, The University of Tokyo

1. まえがき

高度な知識情報処理に対応する能力を持つものとして、論理型プログラムへの期待が高まっている。論理型プログラムは、プログラムの仕様と実行ストラテジを分離して検証することができることから、人間の思考を自然に記述できる利点があり、ソフトウェア生産性の向上をもたらすと考えられている。その反面で、現在の計算機上で、実用規模の知識情報処理システムを構築するのに十分な処理速度を得るのは難しい。一方、論理型プログラムには高い並列性が内在しており、並列処理により計算速度の向上が実現できると考えられる。この目的から、我々は高並列推論エンジンPIE(1)~(3)の開発を進めてきた。

PIEはゴール書換えモデルに基づいて論理型言語の処理を行うマシンである。ゴール書換えモデルの他にも手続き指向型・関係表指向型などの方式のマシンが現在検討されている。どの処理方式を採用した場合も、並列推論マシンにおいては、プロセッサ同士あるいはプロセッサ・メモリ間のデータの授受に携わる通信系の役割は重要であり、相互結合網の転送速度はシステムの全体性能を決定するといえる。

開発中の様々な推論向き計算機(10)(11)(12)(13)などにおいて、相互結合網の実現形態に関して触れられたものが幾つか見受けられる。例えば、ALICE(10)では、データ網(7)(8)とリングバスを用いることが提案されている。しかし、一般に、特定の相互結合網を適用する根拠が十分明らかにされていないとはいえず、また、通信量の正確な見積りが報告されている論文は少ない。すなわち、網

の選択は未定であるか、またはなお流動的であるのが現状である。

本稿では、PIEの相互結合網とそれに関連する制御の方式を提案し、検討を加える。簡単に本稿の特徴を述べれば、以下のようなになる。

① マシンの実行環境を想定して行ったシミュレーションの結果(通信頻度、1回のデータ転送量など)を、結合網設計の根拠として用いる。

② 転送される5種類のデータを、4つのグループに分類し、それぞれに対応する機能をもった結合網の設計を行う。

③ 処理ユニットへの負荷の分散状況に動的に適應する結合網を提案する(4.1, 4.2)。

本稿は、まず高並列推論エンジンPIEの処理方式と構成の概要を述べ(2章)、次いでPIEにおけるデータ転送の特徴を抽出し、必要な結合網の種類と機能分担を示す(3章)。続いて、各網の実現方式について述べ(4章)、これに考察を加え、今後の検討課題を明らかにする(5章)。

2. 高並列推論エンジンPIE

2.1 処理方式と機能モジュール

PIEは、ゴール書換えモデルに基づいて、Prologに代表される論理型のプログラムをOR並列に直接実行するマシンである。

ゴールは基本的に互いに独立であり、ゴールプール中に格納されている。プロセッサは、ゴールプール中のゴールを取り出し、単一化(Unification)を行って新たなゴールを生成し、これに縮退(Reduction)を施して、再びゴールプールに格納する。PIEでは、以上の操作

を繰り返して処理を進め、実行結果を得る。

マシン内でゴールは、AND関係で結ばれたゴールリテラル（ゴールフレーム、Goal Frame, GF）として表現される。マシンは、機能単位である推論ユニット（IU）（図1）が複数台結合された構成をとる。IUは、定義節メモリ（DM）、単一化プロセッサ（UP）、メモリモジュール（MM）、アクティビティコントローラ（AC）から成る。DMには単一化を行う定義節、すなわちプログラムが格納される。UPでは単一化・縮退が行われ、親GFおよび定義節の持っている必要なデータをコピーすることによって新しい子GFが生み出される。（単一化と縮退はパイプライン的に実行される。）MMはゴールプールの役割を果し、UPは自IU内のMMのみからGFを受け取る。新しいGFは、現在の負荷の分散状況に応じて、適当なIUのMMに送り込まれる。各GFは推論木上の各ノードと対応づけられており、ACはそのノード操作を行う。このようにPIEではゴールの処理とその制御が分離されている。

PIEの初期モデル（PIE—I）では、新たなGFはUP内で完全コピーにより生成されるとし、シミュレーション評価などを行った。しかし、この方式は常に大きなGFを処理する必要があり、縮退・転送の両操作に大きなオーバーヘッドを生じる（1）（2）。そこで、未定義変数部を含まない構造データ部分（Ground Instance, GI）を、構造メモリ（SM）に格納してGF間で共有し、必要な場合のみUPによって呼び出す方式の導入を検討した（構造データ共有方式）。シミュレーションによって構造データを共有する効果を調べたところ、従来の完全コピー方式と比較して、GF長が約1/2、縮退時間も約1/2になることが示された（2）。

また、PIE—Iでは、1つ1つのDMがすべての定義節を持つことを仮定していたが、扱う問題が大きくなると、これは現実的ではない。DMのキャッシュ化と、定義節のプールである定義節構造メモリ（DSM）の導入によって、この問題に対処することを検討している。

2.2 PIE—IIの全体構成

PIEにおけるSMの設置方法には2通りが考えられる。全IUが1台のSMを共有する集中型と、各IUごとにSMを設ける分散型である。分散型は、メモリアクセス競合が少なく、各バンクが並列動作を行える利点がある反面、

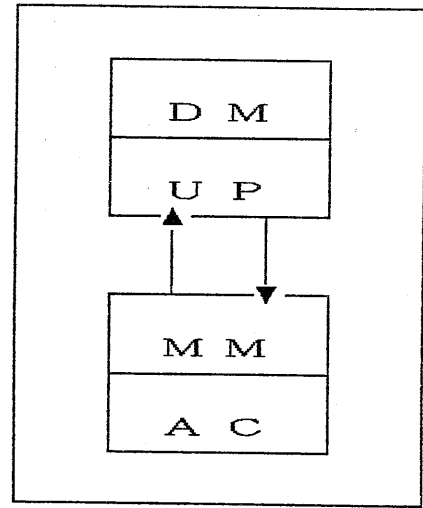


図1 推論ユニット
Fig.1 Inference Unit.

複数バンクにまたがるGIの読み出しに時間がかかり、ガーベジコレクション（GC）の手間も大きくなる欠点を持つ。

PIEでは、集中型SMを採用する。このときSMの処理能力・アクセス競合の点から、1台のSMで賄えるIUの台数には限度があり、16台程度が上限であることが示されている（2）。一方で、マシン全体のIU台数としては数100台から1000台程度を想定している。そこで、PIE第2次モデル（PIE—II）ではマシンを階層的に構成することにし、構造データの共有は低いレベルのシステム内のIU間でしか行わないことにする。

階層構成を採ることには、他にも処理の局所性を活かす利点がある。これによって、

- ① GF転送時間の短縮
- ② ノード情報（コマンド）転送時間の短縮がもたらされる。

図2に、PIE—IIの全体構成を示す。レベル1システム（下位システム）は、IU10数台とSM1台より成り、この中でGFは構造データを共有する。システム内の負荷分散の監視と、他のレベル1システムとのインタフェースは、アクティビティマネージャ（AM）が司る。1つのレベル1システムで処理しきれない仕事の場合、レベル2の結合網を介して、他のレベル1システムにGFが分配される。このGFは、GI部分を完全にコピーして持っており、レベル1システム間での構造データの共有は起らない。レベル2システムを構成するレベル1システムの個数は、目下64程度を想定している。レベル2には、システムマネージャ（System Manager）があり、全システムの実行制御を司

っている。

なお、図中には示さなかったが、各レベル1システム内にはDSMが1個置かれ、定義節プールの役割を果たす。

(今後PIEとはPIE-IIのことを指すものとする。)

2.3 レベル1—レベル2インタフェース
レベル1・レベル2間のインタフェース部では、通常の通信制御の他に、GFと構造データの結合や切り分けなどの作業が必要である。また場合によっては、レベル1システム間で定義節の授受を行うことも考えられる。したがって、インタフェースの管理を行うAMと、SM・DSMは物理的にひとまとまりにして置く方が良く考えられる。

言い換えれば、インタフェース部はAM、SM、DSMを含み、これに作業用メモリであるインタフェースメモリ(IM)と、インタフェースプロセッサ(IP)を付加した構成となる(図3)。詳細は3.2で述べる。

3. PIEにおけるデータ転送と相互結合網

3.1 転送データの分類

PIEの相互結合網の設計に先立ち、マシン

の機能モジュール間を移動するデータの特徴を把握し、要求される転送性能を明らかにする必要がある。

レベル1システム内で転送されるデータは、大きくわけて5種類ある(図4)。ゴールフレーム(GF)、コマンド(ノード情報)、構造データ、定義節、SMアドレスがこれである。レベル2の結合網では、GF、コマンドおよび定義節が転送される。

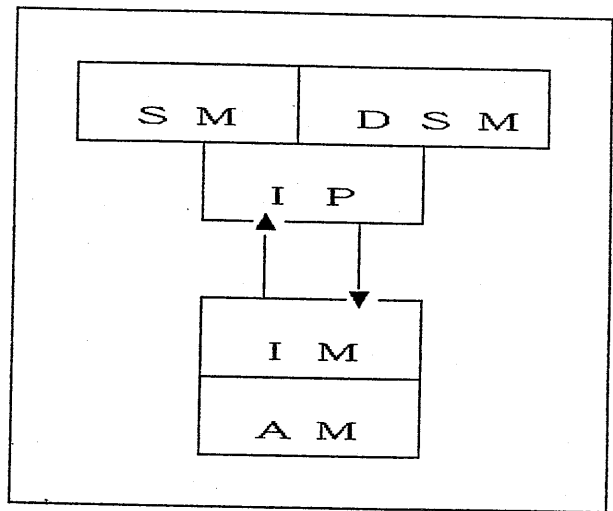


図3 レベル1—レベル2インタフェース部
Fig.3 Level 1—Level 2 Interface Part.

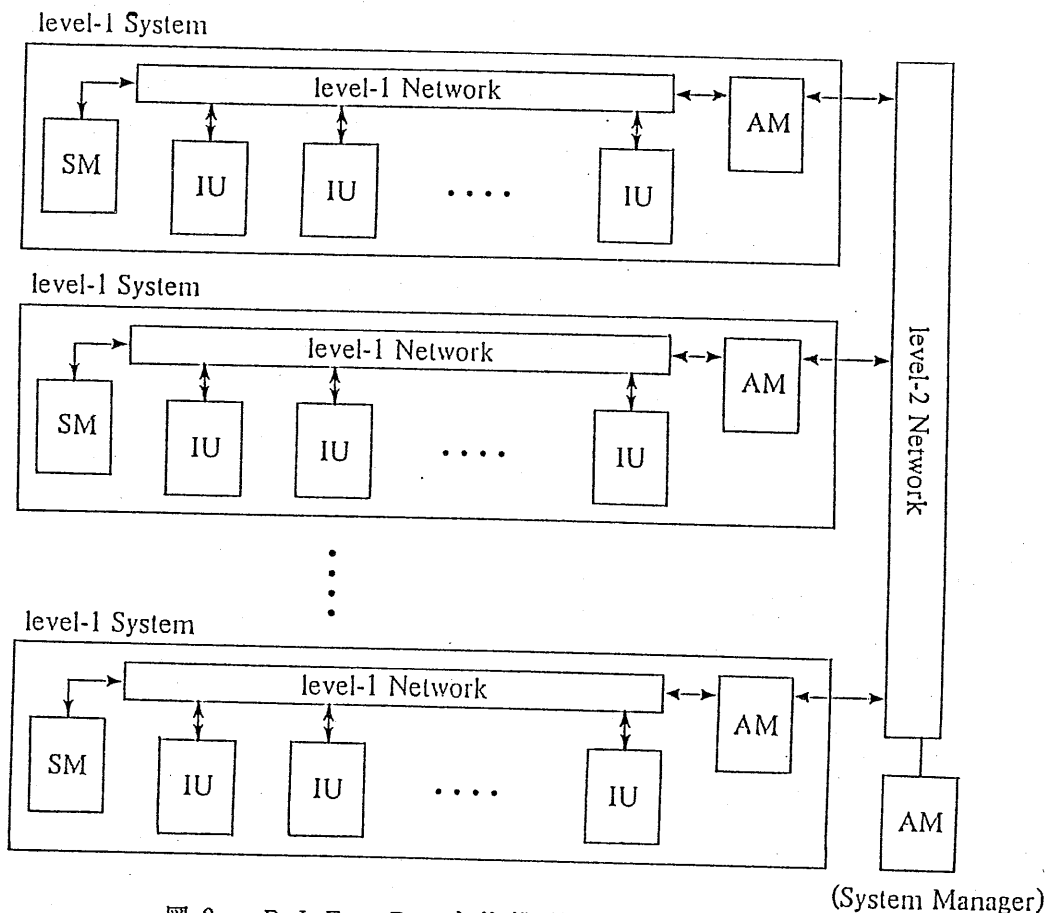


図2 PIE-IIの全体構成
Fig.2 Overview of PIE-II Architecture.

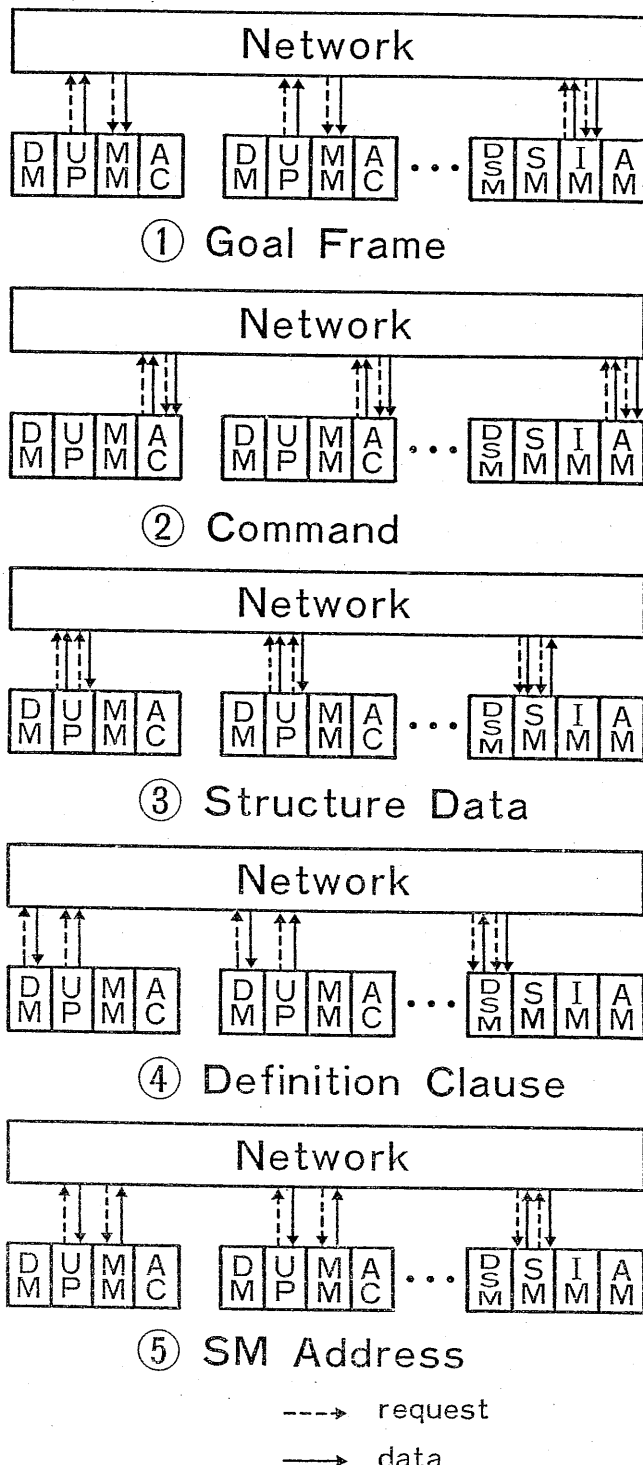


図4 レベル1システムにおけるデータ転送
Fig.4 Data Transfer in Level 1 System.

表1に各データの特徴を示した。大きさ・転送速度は、シミュレーションの結果得られた値である(1)(2)。

① GF

UP内で生成されMMに送られるが、その際、同じレベル1システム内のMMに分配される場合と、IMを経由して他のレベル1へ分配される

表1. レベル1システム内の転送データ
Table 1. Transfer Data in Level 1 System.

データ	FROM	TO	大きさ	転送頻度*(\angle /IU)	
GF	UP UP IM	MM IM MM	数100B	1	
コマンド	AC AC AM	AC AM AC	10~20数B	2~3	
構造データ	Δ G	UP	SM	10~20B	0.5 ~ 1
	LFコマンド	UP	SM	約10B	0.1 ~ 1.4
	LFデータ	SM	UP	10~50B	0.1 ~ 1.4
定義節	UP DSM	DSM DM	数100B 数100B	\ll 1 \ll 1	
SMアドレス	SM MM	UP SM	約4 B 約4 B	0.5 ~ 1	

*UP内での1GFの平均処理時間 T_u を単位時間とした値

る場合がある。

本マシンでは、1GFの転送遅延は問題にならず、処理の流れを乱さないだけの高い転送スループットが要求される。

② コマンド

AC間でやりとりされるノード情報であり、同じレベル1システムのAC間で転送が行われる場合と、IMを経由して他のレベル1へ分配される場合がある。数 μ sの低い転送遅延が要求される。

③ 構造データ

UPで生成され、SMに転送・格納されるもの(Δ G)と、UPでの単一化に用いるためSMから読み出されるもの(追加読み出しデータ、LFデータ)がある。後者の読み出し時には、UPから追加読み出しコマンド(LFコマンド)を発行する必要があるが、本稿においては、便宜上これも構造データの範疇に入れる。

追加読み出しには、数 μ sの早いレスポンスが要求される。

④ 定義節

外部からの入力として与えられる場合と、UPにおいて生成される場合がある。網の立場から見れば、新しい定義節をDSMに格納する転送、DSM内の定義節をDMに読み出す転送に分類されるが、GF転送と比較して、どちらも頻度が小さい。

なお、DSMからDMへの転送は、一般にマルチキャストである(4.1参照)。

⑤ SMアドレス

SMアドレスは、以下の2つの目的で転送される。

(1) 構造データ (ΔG) 生成のとき、新しい GF に SM 内の ΔG を指すポインタを付加する。

(2) ガーベジコレクションのとき、MM 中の生きている GF が参照している構造データに、(SM 内で) 印付けを行う。

前者では、UP で ΔG が生成される速度に追従してアドレスの供給を行うことが必要である。また、表 1 中には (2) に関する転送頻度が示されていないが、これはシミュレーションによるデータが得られていないという意味であり、実際には、(1) の場合と同じくらいの頻度であると予想される。

①から⑤の転送頻度は、表 1 中では、UP 内の 1 GF の平均処理時間 T_u を単位時間とする相対値として与えられている。PIE-I では、この時間は約 700 クロック (試作 UP (1) を用いた場合、約 140 μs) である。PIE-II では、この値は約 1/2 になると予想されている。

3. 2 相互結合網の構成と役割分担

上述の 5 種 (厳密には 7 種) のデータをやりとりする PIE の相互結合網の構成方式を述べる。本節では、結合網の種別と、通信における役割分担を示し、各網の具体的な実現方式については次章で議論することにする。

PIE の結合網としては、5 種のデータそれぞれに、独立の結合網を割り当てる構成が最も自然で、かつ高い転送性能を期待できる方法と考えられる。しかし、網ハードウェア量が大きく必要になり、IU と網のインタフェースが複雑化する難点がある。そこで、大きさなどの点で共通点のある GF と定義節は同一の網で、コマンド・構造データ・SM アドレスはそれぞれ独立な網で転送を行うことにする。このように、ハードウェアの許容範囲内で各種の転送データに個別に対応した網構成を採ることによって、各網の役割分担が明確になり、それぞれの通信機能が単純化され、転送の並列性も高まる。

したがって、レベル 1 システムは 4 種類の網を持つことになる (表 2)。その内訳を以下に列挙して示す。

- (I) 分配網 (Distribution Network, DN)
GF と定義節を IU に分配する。
- (II) コマンド網 (Command Network, CN)
AC 間のコマンド転送を行う。
- (III) 追加読み出し網 (Lazy Fetch Network, LFN)

表 2. PIE の相互結合網
Table 2. Interconnection Networks on PIE

結合網	転送データ	特徴
DN	GF 定義節	高スループット マルチキャスト
CN	コマンド	低遅延
LFN	ΔG LF コマンド LF データ	はやいレスポンス (数 μs 以下)
AN	SM アドレス	ΔG 生成に追従する スループット

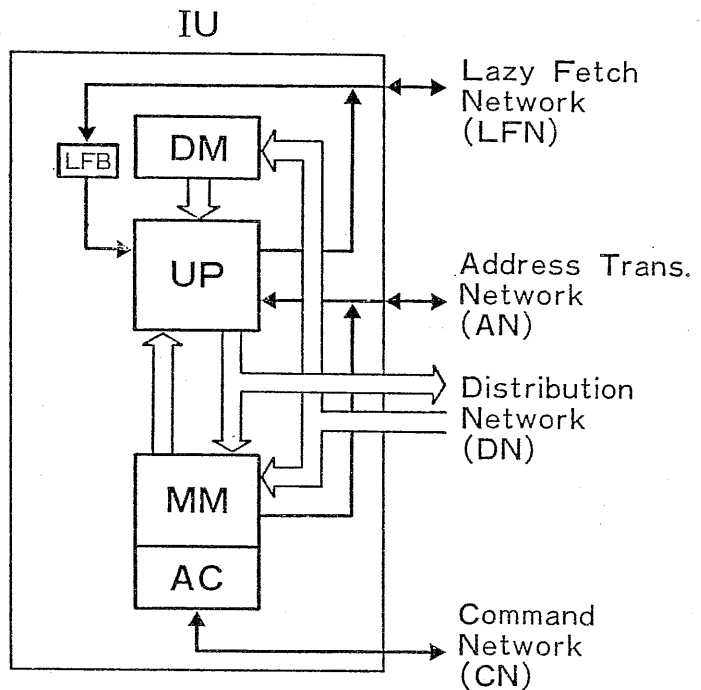


図 5 IU と相互結合網
Fig. 5 IU and the Interconnection Networks.

ΔG , LF コマンド, LF データの転送を行う。

- (IV) アドレス転送網 (Address Transmission Network, AN)

SM アドレスの授受を行う。

同様にレベル 2 には、分配網 (DN) とコマンド網 (CN) がある。

IU と相互結合網の接続を図 5 に、レベル 1 - レベル 2 インタフェース部と結合網の接続を図 6 に示す。

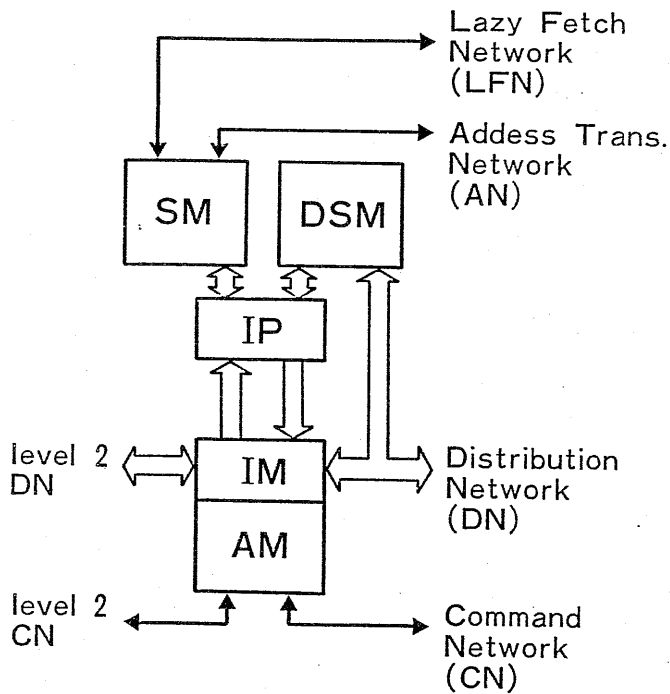


図6 レベル1-レベル2インタフェース部と相互結合網
Fig.6 Level 1 - Level 2 Interface Part and the Interconnection Networks.

4. 各結合網の機能と実現方式

4.1 DN (レベル1)

DNでは、GFと定義節が転送される。両方のデータとも数100Bの大きさであり、本マシンの結合網の中では最も大きな容量が必要となる(表1, 2)。

① GFの転送

PIEでは、GFの分配によって負荷分散を実現する。したがって、

(1) 転送時のスループットが高いこと
の他に、

(2) 負荷の分散が均等であることが要求される。

さらに、GF分配後のコマンドの転送を少なくするようにあらかじめ考慮しなくてはならない場合がある。例えば、notやguardを含むGFを多くのIUに分散させた場合、後に発生するコマンドの送受に要する時間が全体の処理時間に悪影響を及ぼすことが予想される。したがってこのようなGFはできる限り1IUに集中して置くべきである。

すなわち、3番目の要求として、

(3) 後のコマンドの転送量が低く抑えられることが挙げられる。

(1)(2)(3)は、必ずしも同時に満足され得る要求

ではなく、場合によっては相矛盾することもある。3者のトレードオフを考慮したGFの行先制御が必要である。

② 定義節の転送

定義節は、処理の開始時に外部入力として与えられ、あるいは処理中にUPにより生成されてDSM内に格納される。DSMからの読み出しは、DMの要求によって起るが、1つのIUに必要な定義節(群)は、同じレベル1の他のIUにも必要である場合が多く、網はマルチキャスト機能を備えていなくてはならない。

DSMの読み出し要求は、10B程度のコマンドの形で与えられる。この読み出しコマンドの転送には、特殊な網を設けたり他の網を用いたりせず、DNを用いることにする。これは、網インタフェースの複雑化・ハードウェア量の増大を防ぐ意味からも妥当であり、また本コマンドは、転送量・転送頻度がともに小さいことから、GFや定義節の転送に悪影響を及ぼすことはないと考えられる。

①②より、DNでは、1対1のあらかじめ行先の決まっている通信の他に、行先を特定しない負荷の分配とマルチキャストという2種類の特殊な通信が行われることが示された。

最もハードウェアコストの低いDN(レベル1)の実現方法として、バスを用いる方式が挙げられる。しかし、要求されるスループットは1IU(1ポート)あたり数MB/s以上であり、容量の点から適用が不可能である。一方、完全結合網やクロスバススイッチを用いれば高いスループットが得られるが、大きなハードウェア量が必要になる($O(N^2)$)。

格子型網・超立方体網・CCC網などの静的な結合トポロジー(5)の適用も考えられるが、これらは結合の局所性が大きく、問題をうまく分割してIUに割り当てないと、非局所的な通信が全体性能を低下させることになる。

レベル1のDNとしては、オメガ網(図7)(6)を適用する。オメガ網には結合の局所性がなく、したがって、プロセッサ間距離の大きな転送が処理の隘路となる危険がない。また、ハードウェア量も $O(N \log N)$ と妥当である。なお、オメガ網の構成単位であるスイッチング・ユニット(SU)は、ポート数4のものを用いる。

負荷(GF)の均等な分配を実現する方法として、AMによるIUの監視が挙げられる。しかし、完全に負荷の分散状況を監視することは、AMの処理オーバーヘッド・CNの通信オーバー

ットの点から現実的ではない。また、GFの行先MMをあらかじめ決めてしまうと、DN内の閉塞によって待たされるGFが多くなり、全体のスループットが低下する。

逆に、GFの行先を完全にランダムに与えた場合、AM・CNのオーバーヘッドはなくなるが、負荷の不均等から処理効率が著しく低下する。さらにDN内の閉塞の問題があり、このような環境では、理想的な場合の約1/2の転送スループットしか得られないことが解析されている(7)。

本マシンでは、この問題を以下に述べる方法で解決した。

DNを構成するオメガ網には、GFの転送経路とは逆向きに、すなわちMMからUPの向きに、GF転送の許可を与える信号線を設け、負荷を受け入れ可能なMMがあらかじめこの信号をオンにしておく。許可信号が立っている行先のうち、閉塞なしで到達できるものが選ばれ、GFはそこへ転送される。このように簡単な制御機構を付加することで、スループットを高く保ちつつ負荷の均等な分配を実現する。

許可信号のオン・オフを決めるのは、IU内に現存する負荷の量であるが、その閾値は動的にAMが設定する。この作業は厳密に行う必要がなく、負荷の監視を行う周期も、GF転送時間に比較して十分大きく取れるため、AM・CNのオーバーヘッドはともに小さい。

4.2 負荷分散適応型SUの設計

前節で述べたDNは、図8に示すSUで構成される。網は回線交換方式を考え、網制御は各SUが独立に行う分散型である。

本SUは、

- ① 行先を特定するモードと負荷分散を行うモードの2つの制御モードを持ち、ポートごとにモード指定を行う。
 - ② マルチキャスト機能を持つ。
 - ③ 同期式・非同期式の両方で使える。
- などの特徴を持つ。

SUは、データ部(Data Plane)と制御部(Control Plane)より成る。データ部は4入力4出力のクロスバスイッチであり、制御部からの信号(C₀₀~C₃₃)により、ゲートの開閉がなされる。

制御線としては通常見られるreq, ack, stbの他に、rmd(モード指定)とldc(負荷受け入れ可)が設けられている。前者は転送モードの指定(上記)を行い、後者は行先に

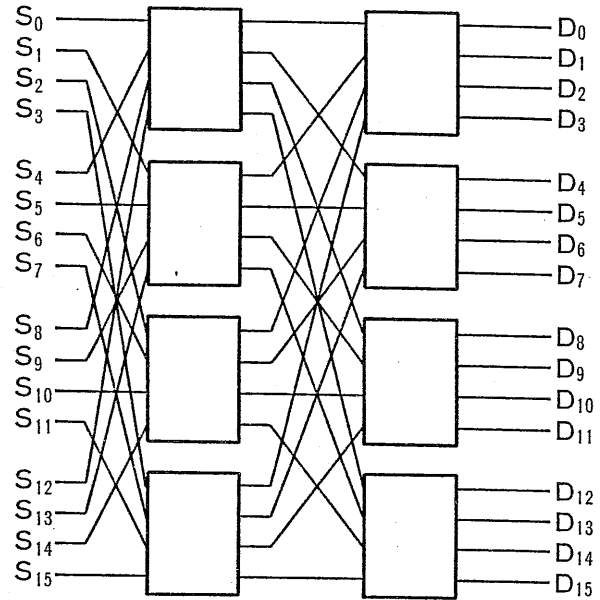


図7 4ポートのSUから構成されるオメガ網
Fig.7 Omega Network Consisting of 4 Port SU's.

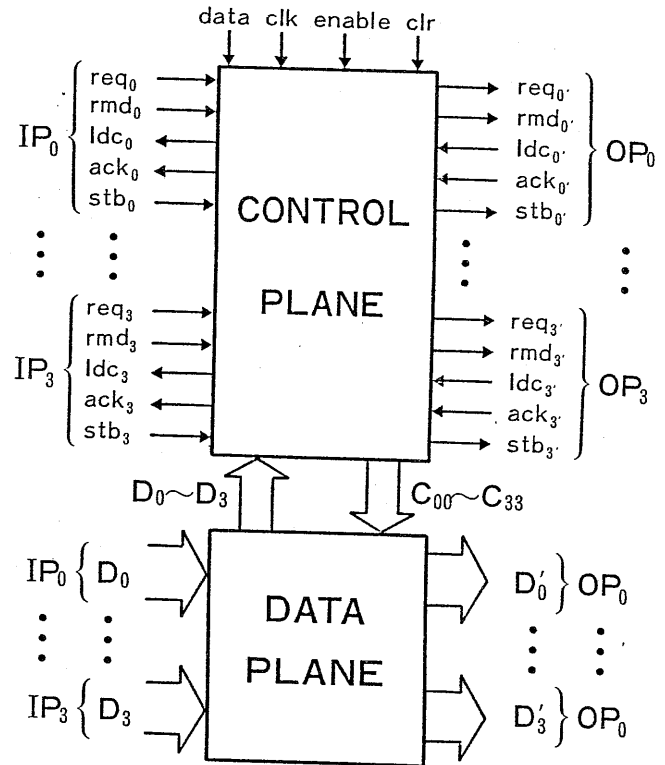


図8 負荷分散適応型SU
Fig.8 SU Adaptive for Load Balancing.

負荷を受入れることができるモジュールが存在するかどうかを示す。

制御部の入力ポート(IP)側には、ルーティング・コントローラがあり、モードに応じて出力ポート(OP)に転送要求を出す。これを受けたOP側は、アービトレーションを行って、

適当なゲート制御信号 (Cij) を立て、次段のSUに転送要求を出す。

行先を特定するモードでは、ソース側はreqを立てると同時にデータ線に行先モジュールの番地を出力し、各段のSUはこの情報をもとにルーティングを行う。負荷分散モードの場合、各段のSUはIdcの立っている任意の空きOPを選択して経路を設定する。

本SUの設計例を付録に示した(DDLによる記述)。4ポートの構成でデータ幅を9bit(1bitはパリティのチェックなどに用いる)とすると、ゲート数は600余り、入出力線の総数は116である。

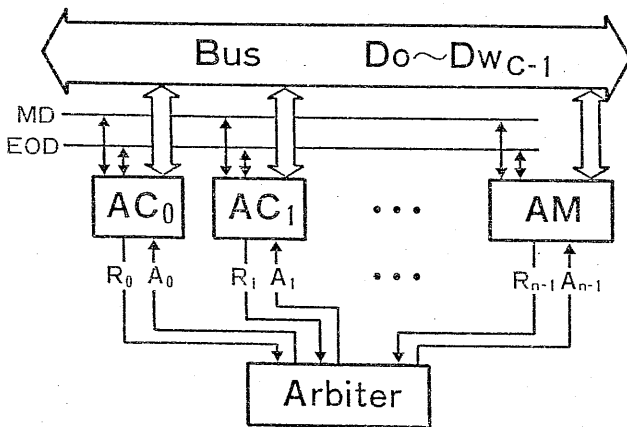
なお、PIEのDN(レベル1)に本SUを適用する場合、規模の点から同期式を採るのが良いと考えられる。

4.3 CN (レベル1)

CNは、AC(またはAM)間でコマンドの転送を行う網である(表1, 2)。

コマンドの転送量は、レベル1全体で0.3~1.2KB/Tuであり、比較的小さい(最大15MB/s程度)。一方で、数μsという低い転送遅延が要求される。

CNは通常特定のAC(AM)間の1対1の通信を行うが、負荷分散の制御に用いられる閾値(4.1参照)を設定する際には、同時に全MMに情報を提供するのが望ましく、マルチキャストの機能が必要である。



MD: Bus Mode EOD: End of Data
Ri: Request i Ai: Acknowledge i

図9 集中管理型個別要求論理バスによるCNの実現
Fig.9 CN Construction by Bus, with Integrated Control and Independent Request Signals.

以上のことから、レベル1のCNとして、総ハードウェア量の小さいバス方式を用いることにする(図9)。バスは、制御の高速性の点から個別要求論理とし、アービトレーションの柔軟性の点から集中管理型を考える。

各AC(AM)は、アービタに転送要求を出し、受理された(ack)のを確認してから、コマンドをバス上に送り出す。その際、線数の節約のため、最初に行先アドレスを送る方式を採り、アドレスバスは設けない。転送の終了はEOD(End of Packet)信号によって知らされる。

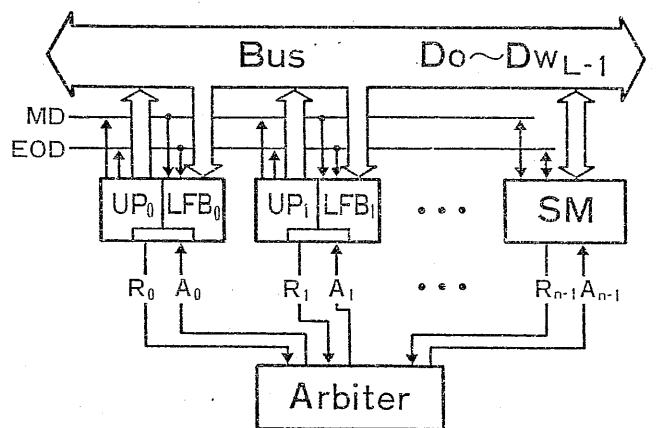
バス幅Wcは、一般に大きいほど高い転送効率が得られると考えられるが、ハードウェア量の増大、データのスキューなどの問題がある。CNの環境では、3B幅(Wc=24)で目標とする転送時間を達成できると考えられる。

4.4 LFN (レベル1)

LFNによって送られるデータは、ΔG, LFコマンド, LFデータの3種である(表1, 2)。

3者の転送量を合計すると、レベル1全体で0.1~2KB/Tu程度である(最大25MB/s)。一方で、追加読み出しの応答(往復とSM内の処理)は、単一化処理を中断して行われるため、十分に短い時間で実現されねばならない(数μs)。

このことと、現在検討中のSMが集中型であ



MD: Bus Mode EOD: End of Data
Ri: Request i Ai: Acknowledge i

図10 集中管理型個別要求論理バスによるLFNの実現
Fig.10 LFN Construction by Bus, with Integrated Control and Independent Request Signals.

ることから、CNと同様に、集中管理型個別要求論理バスによってLFNを構成することにする(図11)。バス幅 W_L は、5B(40本)程度で十分に短いレスポンス時間を得ることができる。なお、転送要求が衝突した際に、後回しにされた追加読み出しの応答が遅くならないように、アービトレーション・転送・SMの内部処理の3つの処理を重畳化する。

4.5 AN (レベル1)

ANは、SMからUPに空きアドレスを提供し、使用中のSMアドレスをMMからSMに知らせるための結合網である(表1, 2)。

各UP内にはSMの空きアドレスを蓄えるバッファがあり、 ΔG 生成時に1つずつこれが使われる。したがって、ANでは転送遅延は問題とならず、 ΔG の生成・消滅に追従するスループットを得ることが問題となる。

これは、レベル1全体で0.1KB/Tu程度と小さい(約1.3MB/s)。したがって、ビットシリアルな時分割多重チャンネル方式のリングバス(図11)をANに適用することにする。

チャンネルには、データの型や有効なデータが入っているか否かを示すヘッダを設ける(この場合、転送単位は固定長でよい)。

4.6 レベル2の結合網

レベル2には、DN・CNの2種類の結合網があり、前者はGFと定義節、後者はコマンドの転送に携わる。

レベル2のDNによって転送されるGFは、構造データを完全にコピーして持つものであり、レベル1のそれと比較して大きい。定義節およびコマンドは、レベル1のそれとほぼ同じ大きさであると考えられる。転送頻度は3者ともにレベル1のものより低いと予想される。

2種類の網に関しては、レベル1システム間でやりとりされるデータのトラヒックを見積り、その結果から実現方式の検討を行う予定である。DNとして、4.1で述べた負荷分散適応型のオメガ網を適用することが考えられ、またCNとして集中管理型個別要求論理バス(4.3)の適用が考えられる。

なお、目下レベル2システムは64個のレベル1システム(全ユニット数1024台程度)より成る場合を想定しており、したがって、2種の網はともに64ポートを持つ。

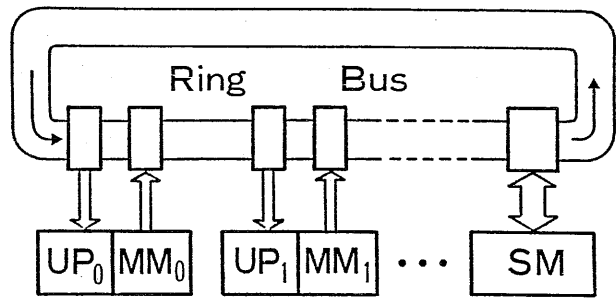


図11. 時分割多重チャンネル型リングバスによるANの実現

Fig.11 AN Construction by Time Slotted Multi-channel Ringbus.

5. 考察

レベル1のDNとしてオメガ網を用いることを考えた。オメガ網は1ルートの多段結合網であり、一般にその欠点として、閉塞率が高いことと信頼性が低いことが指摘されている。本マシンの環境において前者を解決する方法を4.1に示した。

信頼性に関しては、パリティなどによるエラーチェックの他に、次のような対処法が考えられる。

- ① SU内の論理回路の多重化
- ② SUの付加による網の多ルート化(4)(9)

このうち②は、例えば1段分のSUをオメガ網に追加し、故障部を迂回するルーティングを行うことなどが挙げられる。この場合、1組の入出力の組合せに対し、4つの経路が設定できる。

また、本稿では、負荷分散の大まかな監視をAMが行う場合を想定したが、負荷分散に関する制御をDNだけで行う方法が考えられる。また、相互結合網を考慮に入れた負荷分散の評価(解析・シミュレーション)を行い、ここで提案した方式の検証を行う必要がある。

PIEの高機能化(3)とともに、SMに対する要求が増す可能性がある。高頻度のアクセスが起る環境では、SMを多バンク化し、多段結合網でIUと接続する対処法が考えられる。

全体を通じて、PIEの実行環境を想定しての具体的なシミュレーション評価が必要である。

なお、各網のインタフェース部に関しては触れなかったが、データ転送の優先順位決定・エラーチェックなどを含めたインタフェースの設計は重要な課題である。

6. むすび

本稿では、高並列推論エンジンPIEの相互結合網に関して、レベル1システムを中心に、その実現方式の検討を行い、PIE-IIの処理速度に十分に追従する網構成を示した。現在、相互結合網を考慮した負荷分散に関して、より一般的なマシン環境を想定した解析・シミュレーションを行っている。また、本稿で述べた結合網は、PIE-IIの階層的シミュレータ(3)に組み入れられ、全体的な評価が行われる予定である。

文 献

- (1) Moto-oka, T., Tanaka, H., Aida, H., Hirata, K., Maruyama, T. : "The Architecture of a Parallel Inference Engine -PIE-", Proc. Int. Conf. of FGCS 1984 pp.479-488 (Nov. 1984) .
- (2) 平田, 田中, 元岡 : "PIEにおける構造メモリの構成について", アーキテクチャワークショップインジャパン '84シンポジウム (1984-11) .
- (3) 丸山, 平田, 相田, 田中, 元岡 : "高並列推論エンジンPIEの階層的構成とそのシミュレータ", 本研究会技報 (1984-12) .
- (4) 坂井, 計, 田中, 元岡 : "汎用スイッチング・ユニットを用いた高並列計算機の相互結合網", 信学技報, EC 84-18, (1984-07) .
- (5) Feng, T. : "A Survey of Interconnection Networks", IEEE COMPUTER, 14, 12, pp.12-27 (Dec. 1981) .
- (6) Lawrie, D.H. : "Access and Alignment of Data in an Array Processor", IEEE Trans. Comput. C-24, 12, pp. 1145-1155 (Dec.1975)
- (7) Patel, J.H. : "Performance of Processor-Memory Interconnection for Multiprocessors" IEEE Trans. Comput. C-30, 10, pp.771-780 (Oct.1981) .
- (8) Dias, D.M., Jump, J.R. : "Analysis and Simulation of Buffered Delta Networks", IEEE Trans. Comput. C-30, 4, pp.273-282 (Apr.1981) .
- (9) Adams, III, G.B., Siegel, H.J. : "The Extra Stage Cube: A Fault-Tolerant Interconnection Network for Super-systems", IEEE Trans. Comput. C-31, 5, pp.443-454 (May.1982) .
- (10) Darlington, J., Reeve, M. : "ALICE and the Parallel Evaluation of Logic Programs", Reserch Report, Imperial College, London (June 1983) .
- (11) Ciepielewski, A., Haridi, S. : "A Formal Model for OR-Parallel Execution of Logic Programs", IFIP1983, pp.299-305 (Sep.1983) .
- (12) 伊藤, 益田, 清水, 来住, 久野 : "データフロー方式並列推論マシンのアーキテクチャ", Proc. of the Logic Prog. Conf. '84, 7-1, (1984-03) .
- (13) Hasegawa, R., Amamiya, M. : "Parallel Execution of Logic Programs based on Dataflow Concept", Proc. Int. Conf. of FGCS 1984, pp.507-516 (Nov. 1984) .

付録 負荷分散適応型スイッチング・ユニットの構成例 (DDLによる記述) Appendix Construction of a Switching Unit adaptive for Load Balancing (Example Described in DDL) .

```

<SYSTEM> CSU:
<TIME>
  CLK0<(10)>, /* circuit switching SU */
  CLK<(40,0)>. /* arbiter clock */
<ENTRANCE> /* selector clock */
  IPIN(4,12), /* IP input */
  OPIN(4,2), /* OP input */
  DATA,ENABLE,CLEAR.
<TERMANAL>
  IREQ(4), IRMD(4), ILDC(4), IACK(4), ISTB(4), ID(4,9),
  OREQ(4), ORMD(4), OLDLC(4), OACK(4), OSTB(4), OD(4,9),
  C(4,4), IIREQ(4,4), DREQ(4,4), DACK(4,4), ACK(4,4)
  ISTB(1) = IPIN(1,11),
  OLDLC(1) = OPIN(1,0),
  OACK(1) = OPIN(1,1),
  OD(1) = ID(0) & C(0,1) | ID(1) & C(1,1) |
  ID(2) & C(2,1) | ID(3) & C(3,1),
  ILDC(1) = |(OLDLC & OREQ)
/* loop end */
/* boolean end */
<AUTOMATON> IP4 : CLK :
<REGISTER>
  RR(4,4), /* routing register */
  DR(4,4), /* distribution register */
  DLR(4,4) /* dist. latch register */
<TERMINAL>
  DRE(4), DST(4,2)
  
```

```

<DELAY>
IRE(4)<(5)>
<LOGIC>
<<I = 0 : 3>>
TACK = |/ACK(I)

IRE = IREQ & ~IRMD,
DRE = IREQ & IRMD,
<<I = 0 : 3>>
/* address converter */
DST(I,0) = |/(RR(I) & ID(I,0) || ID(I,2) ||
ID(I,4) || ID(I,6)),
DST(I,1) = |/(RR(I) & ID(I,1) || ID(I,3) ||
ID(I,5) || ID(I,7)),
IIREQ(I,0) = ~DST(I,0) & ~DST(I,1) & IRE(I),
IIREQ(I,1) = ~DST(I,0) & DST(I,1) & IRE(I),
IIREQ(I,2) = DST(I,0) & ~DST(I,1) & IRE(I),
IIREQ(I,3) = DST(I,0) & DST(I,1) & IRE(I),
/* routing register update */

|* ENABLE & IREQ(I) *|
RR(I) <- RR(I,1:3) || DATA
/* r.r.update end */
/* address converter end */
/* ditribution selector */
DREQ(I) = DLR(I),
|* ~&/ DACK(I) */
DR(I) <- (~|/DR(I,0:2) || DR(I,0:2),
DLR(I) <- DRE(I) & DR(I)

/* distribution selector end */
/* loop end */
/* logic end */
<END> IP4.

```

```

<AUTOMATON> OP4:CLK0:
<REGISTER>
AR(4,4) /* arbitration register */
<TERMANAL>
IC(4)
<BOOLEAN>
IC = OLDC & (OREQ ~& ~ORMD),
<<J = 0 : 3>>
<<I = 0 : 3>>
ACK(I,J) = C(I,J) & OACK(J)
..
OSTB(J) = ISTB(0) & C(0,J) | ISTB(1) & C(1,J) |
ISTB(2) & C(2,J) | ISTB(3) & C(3,J)
/* loop end */
/* booleazn end */
<LOGIC>
<<J = 0 : 3>>
<<I = 0 : 3>>
C(I,J) = (IC(J) & DREQ(I,J) | IIREQ(I,J))
& AR(J,I), /* AR is special */
DACK(I,J) = DREQ(I,J) & C(I,J)
.. /* I loop end */
OREQ(J) = C(0,J) | C(1,J) | C(2,J) | C(3,J),
ORMD(J) = DACK(0,J) | DACK(1,J) |
DACK(2,J) | DACK(3,J),
|* ~OREQ(J) *|
AR(J) <- (~|/AR(J,0:2) ) || AR(J,0:2)
/* if end */
/* J loop end */
/* logic end */
<END> OP4.
<END> CSU.

```