

並列計算機 A D E N A R T における
データの分散配置方式

森康浩 右田学 若谷彰良 岡本理
松下電器産業(株)半導体研究センター
Email:ymori@vdrl.src.mei.co.jp

並列計算機 A D E N A R T においては、ホスト計算機上の配列データは複数のプロセッサエレメント (P E) に分散配置され、各 P E はこの分散配置されたデータを独自のネットワーク (H X n e t) を用いて交換、編集しながら処理を行ない、処理結果は各 P E から再びホスト計算機上の配列データに戻る。このため、H X n e t の特徴を生かした効率的な P E 間転送を可能にさせる最適なデータの分散配置方式が必要となる。本稿では、本計算機における配列データの分散配置方式について概説し、これを高速に実現するためのハードウェア機構 (I U P L U S) の構成について述べる。最後に、このハードウェア機構を実装した上での A D E N A R T のシステム性能の評価に関して報告する。

The data distribution method
in the parallel computer A D E N A R T

Yasuhiro Mori, Manabu Migita, Akiyoshi Wakatani, Tadashi Okamoto
Matsushita Electric Industrial Co., Ltd. Semiconductor Research Center
3-15, Yakumonakamachi, Moriguchi 579, Japan

In the parallel computer, A D E N A R T, data on the host-computer is distributed to each processor element (P E). Each P E proceeds with their computation while exchanging and editing data via the unique network called H X n e t . The result data on each P E is combined and again sent back to the data-set on the host computer. In order to realize efficient inter-processor communication featuring advantages of H X n e t , appropriate data distribution strategy is crucial. In the report, we summarize the data distribution method in A D E N A R T and show the accelerating mechanism and the hardware structure for it.

1. はじめに

我々は、京都大学工学部（野木達夫助教授）と共同で、並列計算機ADENART（以前はADENAと呼んでいた）を開発した。(1)(2)ADENARTは256個のプロセッサエレメント（以後PEと呼ぶ）を用いたMIMD型の並列計算機であり、PE間通信ネットワーク（以後ハイパースネットワークと呼ぶ）によって2回のPE間データ転送により、すべてのPE間のデータ転送が可能である。このハイパースネットワークの特徴を生かすため、各PEに対して最適な配列データの分散配置が必要となる。また、ADENARTでは、並列動作を”陽に”意識してプログラミングするユーザーのために、FORTRANに並列化構文を追加したADETRAN(3)を用意している。一方、既にFORTRANで開発されたプログラムをそのまま流用したいというユーザーに対しては、並列プログラムに変換する手段が必要となる。この手段の一つとして自動並列化コンパイラAPARC(4)を開発中である。APARCによって生成される並列プログラムは、”陽に”並列動作を意識して記述されたプログラムに比べ、逐次処理部と並列処理部とのデータ転送が頻発することになり、このデータ転送を高速化することはシステム性能を上げる上で重要である。

以上のように、ADENARTのようにホスト計算機と並列処理用の”エンジン”-ADENART本体から構成され、ホスト計算機のデータを複数のPEで分散処理する形式の並列計算機システムにおいては、データの分散配置方式と、ホスト-ADENART本体間のデータ転送によって生じるオーバーヘッドが問題となる。

富士通のAP1000においては、データの分散/収集時のセル（ADENARTのPEにあたる）の切替えによるオーバーヘッドを解消するため、Broadcast networkと呼ぶネットワークと専用のLSIを導入し、データのブロードキャストと同様の手順/速度でデータの分散/収集ができるようにしている。(5)これに対し、

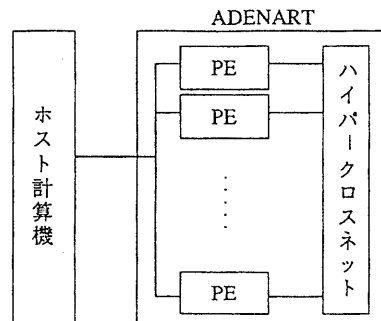
ADENARTでは、後述するようにハイパースネットワークの特徴を生かすため、データの分散配置にともなうオーバーヘッドが大きい。

本稿では、ADENARTにおけるデータの分散配置方式を示した後、これを高速化するためのハードウェア機構（IUPLUSと呼ぶ）について触れ、最後にこれを実装した上でのADENARTのシステム性能の再評価に関して述べる。

2. 並列計算機ADENART

2.1 ADENARTシステム

ADENARTシステムは、第1図に示すようにホスト計算機とADENART本体から構成される。

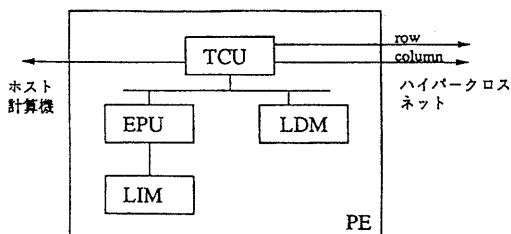


第1図 ADENARTシステム

ホスト計算機はUNIX環境下でユーザーに対するあらゆるサービスを提供するとともに、プログラムの逐次処理部分を実行する。ADENART本体はMIMD型の並列計算エンジンで、プログラムの並列処理部分を実行する。

ADENART本体は、256個のPEとハイパースネットワークから構成される。一つのPEは、ピーク性能20MFLOPSの浮動小数点演算プロセッサ：EPU(Element Processing Unit)、データ転送用コントローラ：TCU(Transfer Control Unit)、命令メモリ：LIM(Local Instruct

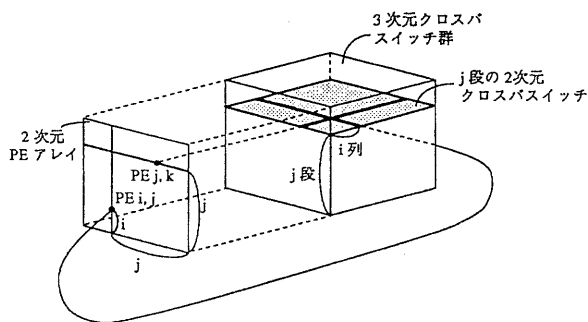
ion Memory, 24ビット*64Kワード)、データメモリ: LDM(Local Data Memory, 72ビット*256Kワード, ECC8ビットを含む)の4つのパーツによって第2図の接続により構成される。



第2図 PEの構成

2. 2 ハイパークロスネット

ハイパークロスネットはADENARTを特徴づけるPE間接続ネットワークであり、自社開発した2種類のLSIによって構成される。第3図にその接続を示すブロック図を示す。



第3図 ハイパークロスネットの構成

ハイパークロスネットは、 16×16 個の2次元のクロスバスイッチを3次的に16段重ねて

配置したクロスバスイッチ群を構成し、 16×16 個のPEのそれぞれに2方向のバスを持たせ、その2方向のバスによって前記のクロスバスイッチに対して2方向から接続する形態になっている。すなわち、 $PE_{i,j}$ はj段i列の16個のクロスバスイッチを介してj行の16個のPE($PE_{j,k}$ $k=0,15$)と接続されることになる。 $PE_{i,j}$ が $PE_{j,k}$ に接続されているので、(式1)に示す2回のPE間転送で任意のPE間のデータ転送を実現できる。

$$PE_{k,1} \leftarrow PE_{j,k} \leftarrow PE_{i,j} \quad (1)$$

3. データ転送に伴う問題点

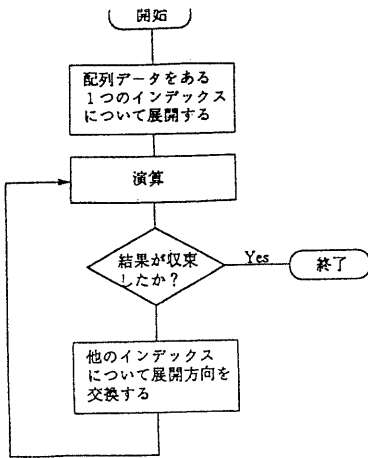
ADENARTでは、データの転送に伴って、次の2点の問題が発生する。

- (1) データの再配置
- (2) データの分散/収集

上記の問題点について、次にその詳細を述べる。

3. 1 配列データの再配置

2章でADENARTのハイパークロスネットは、 $PE_{i,j}$ と $PE_{j,k}$ がクロスバスイッチ群を介して接続された形態のネットワークであることを述べた。ところで、このアーキテクチャーはもとも数値解法アルゴリズムの一つであるADI法(Alternating Direction Implicit method)を高速度に処理することを目的として考案されたものである。ADI法によって問題を処理するとき、第4図に示すようなフローチャートに従って結果が収束するまで演算を繰り返す。ここで、配列データ: $U_{i,j,k}$ をインデックス: i に着目して1次元化することをi方向に展開するという。また、演算ステップ間で行なうデータの展開方向を交換する操作をADE(Alternating Direction Edition)操作といい、これは(式2-1, 2, 3)で示される。



第4図 ADI法による問題処理手順

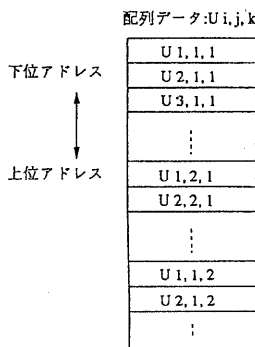
$$U_{i,j,k} \leftarrow U_{i,j,k} \quad (2-1)$$

$$U_{i,j,k} \leftarrow U_{i,j,k} \quad (2-2)$$

$$U_{i,j,k} \leftarrow U_{i,j,k} \quad (2-3)$$

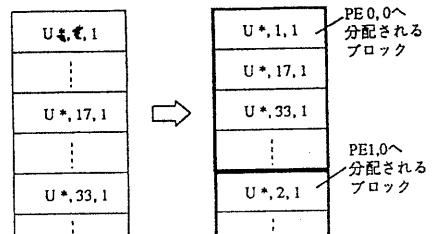
ハイパースネットワークはPE_{i,j}とPE_{j,k}が接続された形態のネットワークであるので、//で囲まれたインデックスを論理的なプロセッサ番号として配列データを各PEに割り当てるとき、(式2-1, 2, 3)は唯1回のPE間転送によって実現されることになり、より効率的で自然なADE操作を実現できる。現実には、PEは有限個であるので、1つの物理PEが複数の論理PEを兼ねることになる。(PEの多重使用)

一般に配列データ: $U_{i,j,k}$ は第5図に示すようにFORTRANコンパイラによってより左側にあるインデックスがより頻繁に変化するようメモリ上に割り付けられる。



第5図 FORTRANコンパイラによる配列データの配置

ADENARTでは配列データのインデックスによって配置先のPEを決定するため、前記したような理由により、あるPEに割り当てるデータはホスト計算機のメモリー上で分散していることになる。このため、このままでは最悪の場合1ワード毎に転送先のPEを切り替えなければならないといった事態が発生し、転送時のオーバーヘッドが大きくなり過ぎる。したがって、ADENARTでは第6図に示すように一つのPEに割り当てるデータ毎にブロック化して再配置する必要がある。このことに起因するオーバーヘッドは、データサイズの増加とともに増大する傾向にあり、転送効率を上げる上で障壁となる。



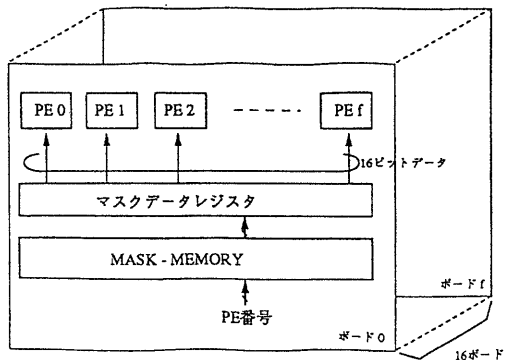
第6図 i方向へ展開する場合のデータのブロック化

3. 2 データの分散/収集

ブロック化されて再配置されたデータを各PEのLDMに分配するためには、1)データを切り分けて、2)転送先のPEを切替えながら転送するといったステップが必要である。

従来、ADENARTにおいては、このオーバーヘッドをおさえるため、第7図に示すハードウェア機構を設けている。(6)

ADENARTでは1枚のボードに16個のPEが実装され、計16枚のボードによって256個のPEが提供される。各ボードにはボード番号が付けられ、固有のアドレスが割り振られる。



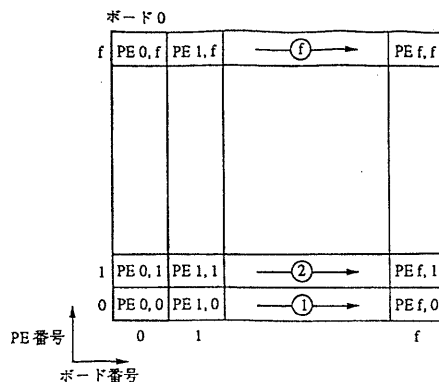
第7図

また、各PEにはボード上に物理的に実装された位置によってPE番号が決められており、このPE番号はボードに実装されているMASK-MEMORYのアドレスに置き換えられる。MASK-MEMORYにはあらかじめ、あるPE番号=アドレスにそのPEだけが選択されるようなデータが記憶される。MASK-MEMORYから出力されるデータの1ビット1ビットは、1つ1つのPEのチップセレクトとして入力される。

一つ一つのPEは、ボード番号とPE番号によって特定され、識別されるが、これらは、実際には、ボードのアドレスとMASK-MEMORYのアドレスに置換されて用いられることになる。

PEへのアクセスは、第8図に示すようにPE 0,0, PE 1,0と続き、PE f,0までいって次の行へ入る。次行へ移行する際に全ボードにPE番号をブロードキャストしておき、1PE/1ボードを選択しておく。その後、ボード番号をインクリメントしながら、データを送る。

したがって、ホストは各ボードに割り当てられたアドレスに対して連続的にデータを書き込んで行くことによって、PEにデータを転送できるが、データの切りわけや転送先のPEのアドレス計算はホスト計算機で実行されているのでより高速化を図る必要がある。

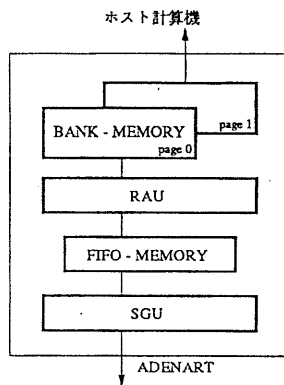


第8図 PEのアクセス順序

4. IUPLUS

IUPLUSは、ホスト計算機とADENAの間であって、3章で述べたデータの再配置、並びに分散/収集によって発生するオーバーヘッドを小さくすることを目的とするハードウェア機構である。IUPLUSのブロック構成図を第9図に示す。

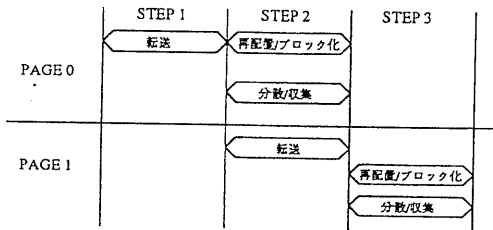
IUPLUSは、BANK-MEMORY, RAU (Re-Arrangement control Unit), FIFO-MEMORY, SGU (Scatter/Gather control Unit)の4つのブロックに分割できる。



第9図 IUPLUSの構成

第10図はIUPLUSの動作をしめすチャート図である。

ホスト計算機のデータは、まず、BANK-MEMORYに1次的に記憶される。(この時点では配列データはFORTRANコンパイラの定めるフォーマットのままである。) つぎに、RAUは配列データの中から一つのPEに転送するデータを連続的に拾いだし、FIFO-MEMORYに詰め込んでいく動作をする。これによって、一つのPEに送るデータがブロック化されフォーマットの変換がなされることになる。(データの再配置) SGUは、ブロック化されて詰め込まれたデータをFIFO-MEMORYから順に読み出して、順次転送先のPEを切替え、分散(収集)させる。(データの分散/収集) ここで、RAUとSGUがFIFO-MEMORYを挟んで非同期に動作できるので、データの最配置と分散/収集といったステップがオーバーラップして実行でき、これらを原因とするオーバーヘッドを削減し、システム性能を向上させることができる。また、BANK-MEMORYとして2変数分=2ページのメモリを用意してあるので、2変数以上のデータをADENARTに転送するとき、ホスト-ADENARTのデータ転送とデータの再配置、分散/収集もオーバーラップさせることができる。



第10図 IUPLUSの動作

5. 性能予測

従来のシステム(ADENART)において、 n ワードのデータを転送する場合、データ転送に要する時間: T は、データ再配置に要する時間: $R(n)$ とデータの分散/収集に要する時間: $Tc \times p$ と、実際のデータ転送時間: $Tt \times n$ の和になり、(式3)で書き表される。

$$T = R(n) + Tc \times p + Tt \times n \quad (3)$$

p : PE数

Tc : PE切替え時間

Tt : 1ワードのデータ転送時間

一方、IUPLUSを実装したシステム(ADENART+)の場合、データ転送時間: Tiu は、データ再配置と分散/収集がほぼオーバーラップされるため、データ再配置時間: $Riu(n)$ と、実際のデータ転送時間: $Tt \times n$ の和となり、次式で書き表される。

$$Tiu = Riu(n) + Tt \times n \quad (4)$$

ここで、 $R(n)$ 、 $Riu(n)$ はデータ再配置の時間関数であり、単純にデータのワード数: n にのみ依存する増加関数である。

第1表はADENARTにおいて、(3)を実測した結果である。

n	$R(n)$	$Tc \times p + Tt \times n$
256k	0.28sec	0.14sec
2048k	1.49sec	1.0 sec

第1表 再配置の実行時間

A DENARTにおいては実際のデータ転送能力として2MB/sec程度あるが、再配置に要する時間を考慮すると、約1/3に落ちる。

A DENART+においては、2回のメモリ間転送:

1) ホスト計算機→BANK-MEMORY,

2) BANK-MEMORY→FIFO-MEMORY

(FIFO-MEMORY→PEは2)にかくれる)によってデータの分散配置が実行される。したがって、1)あるいは2)のデータ転送速度で遅い方の速度が支配的となる。1)の転送速度は10MB/sec程度であるので、A DENART+においてはA DENARTの10倍程度の速度でデータの再配置が実行できると予測される。

APARCはFORTARNで記述されたプログラムを入力としてADENTRANプログラムを出力する並列化コンパイラである。APARCを用いて並列化したプログラムは、初めからADETRANを用いて"陽"に並列記述したプログラムに比べ、並列処理部分の粒度が小さくなりがちである。このため、データ転送に伴うオーバーヘッドが大きいと、十分に満足のいく処理性能を得られないことになり、場合によっては並列化のメリットがないということもありうる。従って、ホスト計算機とA DENART本体とのデータ転送に伴うオーバーヘッドを削減することは、APARCの適用範囲を拡大し、A DENARTの並列プログラム開発環境の向上に貢献する上で重要である。

6. おわりに

並列計算機A DENARTの配列データの分散配置の方式とこれを効率良く行なうためのハードウェア機構について述べた。本ハードウェア機構によって、ホスト計算機とA DENARTとの間でのデータ転送時におけるオーバーヘッドを削減し、APARCの適用範囲を拡大するとともに、その負担を下げることができる。システムにインプリメントした上で、各種のアプリケーションを

用いて評価していく予定である。

参考文献

- (1) T. Nogi: Parallel computation, patterns and waves qualitative analysis of nonlinear differential equations, pp279-318 (1986)
- (2) 谷川他: 並列計算機A DENA, 情報学会計算機アーキテクチャ研究会 報告, CPSY88-11, pp33-40 (1989)
- (3) 若谷他: 並列処理言語ADETRANの実装, 情報学会ソフトウェア基礎論プログラミング言語合同研究会 報告, pp1-10 (1990)
- (4) 材木他: 並列計算機A DENART 用自動並列化コンパイラ: APARC, 並列処理シンポジウムJSP P' 91, pp293-300 (1991)
- (5) 加藤他: 高並列計算機CAP-IIIのロードキャストネットワーク, 情報学会計算機アーキテクチャ報告 83-39, pp229-234, 1991
- (6) 森他: 並列計算機におけるPEへのデータ転送制御方式, 電子通信学会 春期全国大会No. 6, D-318 (1989)