

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号  
特許第4447770号  
(P4447770)

(45) 発行日 平成22年4月7日(2010.4.7)

(24) 登録日 平成22年1月29日(2010.1.29)

(51) Int.Cl.  
G06F 15/80 (2006.01)

F I  
G O 6 F 15/80

請求項の数 39 (全 46 頁)

(21) 出願番号	特願2000-516291 (P2000-516291)	(73) 特許権者	597154922
(86) (22) 出願日	平成10年10月9日 (1998. 10. 9)		アルテラ コーポレーション
(65) 公表番号	特表2001-520418 (P2001-520418A)		Altera Corporation
(43) 公表日	平成13年10月30日 (2001. 10. 30)		アメリカ合衆国 95134 カリフォル
(86) 国際出願番号	PCT/US1998/021478		ニア州 サン ホセ イノベーション ド
(87) 国際公開番号	W01999/019807		ライヴ 101
(87) 国際公開日	平成11年4月22日 (1999. 4. 22)	(74) 代理人	100076428
審査請求日	平成17年10月7日 (2005. 10. 7)		弁理士 大塚 康徳
(31) 優先権主張番号	08/949, 122	(74) 代理人	100112508
(32) 優先日	平成9年10月10日 (1997. 10. 10)		弁理士 高柳 司郎
(33) 優先権主張国	米国 (US)	(74) 代理人	100115071
前置審査			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		最終頁に続く	

(54) 【発明の名称】 相互接続システム及び並列プロセッサとその形成方法

(57) 【特許請求の範囲】

【請求項 1】

N × Mアレイに接続された複数の処理エレメント用相互接続システムであって、NとMは両方とも1より大きく、各処理エレメントはデータ及びコマンドを送受信するための通信ポートを備え、前記複数の処理エレメントが、それぞれM個の処理エレメントからなるNクラスタにグループ化された相互接続システムにおいて、

処理エレメント間接続経路と、

クラスタの間の相互に排他的な処理エレメント間接続経路を組み合わせ、従来型トラス接続処理エレメントアレイの接続性に等価な処理エレメント間接続性を提供するために必要な通信経路の個数を実質的に減少させるように、前記複数の処理エレメントのそれぞれの前記通信ポートに接続され、前記処理エレメントにより制御される制御可能なスイッチエレメントを含むクラスタスイッチとを有し、

前記クラスタスイッチが更に転置処理エレメント間および超立方体補足処理エレメント間に直接通信を提供するための接続部を有し、

データ及びコマンドが、前記制御可能なスイッチエレメントで選択された8つの選択モードのうちの1つで前記通信ポートにおいて送受信可能であり、前記8つの選択モードは、

a) 前記通信ポートを介して東処理エレメントへデータを送信し、同時に、前記通信ポートを介して西処理エレメントからデータを受信するための送信東/受信西モードと、

b) 前記通信ポートを介して北処理エレメントへデータを送信し、同時に、前記通信ポ

ートを介して南処理エレメントからデータを受信するための送信北 / 受信南モードと、  
c) 前記通信ポートを介して南処理エレメントへデータを送信し、同時に、前記通信ポートを介して北処理エレメントからデータを受信するための送信南 / 受信北モードと、  
d) 前記通信ポートを介して西処理エレメントへデータを送信し、同時に、前記通信ポートを介して東処理エレメントからデータを受信するための送信西 / 受信東モードと、  
e) 転置された処理エレメント間の送受信のための転置送信 / 受信モードと、  
f) 距離 1 超立方体処理エレメント間の送受信のための超立方体送信 / 受信モードと、  
g) 選定距離 2 超立方体処理エレメント間の送受信のための超立方体送信 / 受信モードと、

h) 距離 d、d 次元超立方体補足処理エレメント間の送受信のための超立方体補集合送信 / 受信モードとを有することを特徴とする相互接続システム。

10

【請求項 2】

前記モードが前記処理エレメント間に直接接続経路が確立されることを可能にすることを特徴とする請求項 1 に記載の相互接続システム。

【請求項 3】

更に、各処理エレメント制御ポートに制御情報を同時に送り、かつ各処理エレメントにおいてレジスタにロードするために各処理エレメントのデータポートにデータを送るよう  
に接続されたコントローラ及びメモリシステムを有することを特徴とする請求項 2 に記載  
の相互接続システム。

【請求項 4】

20

各通信ポートが B ビット幅の送受信経路を含み、前記 B は 1 に等しいか、1 以上の整数  
であることを特徴とする請求項 2 に記載の相互接続システム。

【請求項 5】

各処理エレメントは、ある通信ポートを経てデータ又はコマンドを選択的に送り、制御  
ポートを介して受信し各々の処理エレメントに存在する制御論理により復号された通信命  
令に基づいて、同時に、別の通信ポートを経てデータ又はコマンドを受信するように接続  
されたことを特徴とする請求項 2 に記載の相互接続システム。

【請求項 6】

前記通信命令は、コントローラから前記制御ポートを経て前記制御論理によって受信さ  
れることを特徴とする請求項 5 に記載の相互接続システム。

30

【請求項 7】

前記クラスタスイッチがオペレーションをサポートし、前記処理エレメントはそれぞれ  
同時にコマンド又はデータを送り、同時に、コマンド又はデータを受け取ることを特徴と  
する請求項 5 に記載の相互接続システム。

【請求項 8】

前記処理エレメントはそれぞれ前記通信ポートの送信部分を経てコマンド又はデータを  
同時に送り、同時に、前記通信ポートの受信部分を経てデータ又はコマンドを受け取るよ  
うに、前記同時オペレーションが選択的にスイッチされることを特徴とする請求項 7 に記  
載の相互接続システム。

【請求項 9】

40

並列プロセッサであって、

各クラスタが、各処理エレメントが合計 B 本のワイヤを経てデータを送受信する通信ポ  
ートを有する M 個の処理エレメントを含む N 個のクラスタと、

前記クラスタの対の間に接続された線幅 (M) (B) に等しいか、それ以下の個数の通  
信経路であって、前記クラスタの対内の各クラスタメンバが当該クラスタの対のもう一方  
のクラスタ内処理エレメントに対してトラス最隣接体である処理エレメントを含み、各  
通信経路が相互に排他的な 2 つのトラス方向である、南と東、又は、南と西、又は、北  
と東、又は、北と西における前記クラスタの対の間の通信を可能にする通信経路と、

各通信経路の処理エレメントにより制御され、前記クラスタの対の間の線幅 2 (M) (B) のワイヤによる通信を、線幅 (M) (B) のワイヤ経路数に等しいか、それ以下に組

50

み合わせるように接続されたマルチプレクサと、  
を有することを特徴とする並列プロセッサ。

【請求項 10】

各クラスタの処理エレメントが北および西トラス方向で一方のクラスタと、南および東トラス方向でもう一方のクラスタと通信することを特徴とする請求項 9 に記載の並列プロセッサ。

【請求項 11】

各クラスタの処理エレメントが北および東トラス方向で一方のクラスタと、南および西トラス方向でもう一方のクラスタと通信することを特徴とする請求項 9 に記載の並列プロセッサ。

10

【請求項 12】

各クラスタの処理エレメントが 2 つの超立方体方向で一方のクラスタと、2 つの超立方体方向でもう一方のクラスタと通信することを特徴とする請求項 9 に記載の並列プロセッサ。

【請求項 13】

少なくとも 1 つのクラスタが超立方体補足対を含むことを特徴とする請求項 9 に記載の並列プロセッサ。

【請求項 14】

クラスタスイッチが前記マルチプレクサを有し、相互に排他的な 2 つのトラス方向から受信した通信を、1 つのクラスタ内の処理エレメントに多重化するように前記クラスタスイッチが接続されることを特徴とする請求項 9 に記載の並列プロセッサ。

20

【請求項 15】

前記クラスタスイッチが 1 つのクラスタ内の処理エレメントからの通信をもう一方のクラスタへ送信するために多重化するように接続されることを特徴とする請求項 14 に記載の並列プロセッサ。

【請求項 16】

前記クラスタスイッチが、1 つのクラスタ内の転置処理エレメント間で通信を多重化するように接続されることを特徴とする請求項 15 に記載の並列プロセッサ。

【請求項 17】

前記 N が M に等しいか、M 以上であることを特徴とする請求項 9 に記載の並列プロセッサ。

30

【請求項 18】

前記 N が M より小さいことを特徴とする請求項 9 に記載の並列プロセッサ。

【請求項 19】

並列プロセッサであって、

各処理エレメントが合計 B 本のワイヤを経てデータを送受信する通信ポートを有し、1 つのクラスタ内の各処理エレメントが 1 つのクラスタ内において前記クラスタの外部の処理エレメントに対するよりも他の処理エレメントに対して物理的に一層近接して形成された、それぞれが M 個の処理エレメントを含む N 個のクラスタと、

前記クラスタの対の間に接続された線幅 (M) (B) のワイヤの数に等しいかそれ以下の通信経路であって、前記クラスタの対内の各クラスタメンバが、当該クラスタの対のもう一方のクラスタ内処理エレメントに対してトラス最隣接体である処理エレメントを含み、各通信経路の処理エレメントによって制御されるマルチプレクサが、南と東、又は、南と西、又は、北と東、又は、北と西、または、2 つの超立方体次元の間である相互に排他的な 2 つのトラス方向における前記クラスタの対の間の通信を可能にした通信経路とを有し、前記マルチプレクサは、前記クラスタの対の間の線幅 2 (M) (B) のワイヤによる通信を線幅 (M) (B) のワイヤ経路数に等しいか、それ以下に組み合わせるように接続されたことを特徴とする並列プロセッサ。

40

【請求項 20】

各クラスタの処理エレメントが北および西トラス方向で一方のクラスタと、南および

50

東トラス方向でもう一方のクラスタと通信することを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 21】

各クラスタの処理エレメントが北および東トラス方向で一方のクラスタと通信することを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 22】

少なくとも 1 つのクラスタが超立方体補足対を含むことを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 23】

クラスタスイッチが前記マルチプレクサを有し、2 つの超立方体方向から受信した通信を 1 つのクラスタ内の処理エレメントに多重化するように前記クラスタスイッチが接続されることを特徴とする請求項 19 に記載の並列プロセッサ。

10

【請求項 24】

前記クラスタスイッチが 1 つのクラスタ内の処理エレメントからの通信を、もう一方のクラスタへ送信するために多重化するように接続されることを特徴とする請求項 23 に記載の並列プロセッサ。

【請求項 25】

前記クラスタスイッチが 1 つのクラスタ内の超立方体補足処理エレメント間での通信を多重化するように接続されることを特徴とする請求項 24 に記載の並列プロセッサ。

【請求項 26】

20

前記 N が M に等しいか、それ以下であることを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 27】

前記 N が M より大きいことを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 28】

前記処理エレメント間通信がビット直列であり、各処理エレメントのクラスタが前記通信経路を経て他の 2 つのクラスタと通信することを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 29】

処理エレメント間の通信経路がデータバスを有することを特徴とする請求項 19 に記載の並列プロセッサ。

30

【請求項 30】

前記通信経路が双方向性であることを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 31】

前記通信経路が単方向性であることを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 32】

P と Q が、前記並列プロセッサと同じ処理エレメントの数を有するトラス接続アレイのそれぞれ行と列の個数であり、P と Q がそれぞれ N と M に等しいことを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 33】

40

P と Q が前記並列プロセッサと同じ処理エレメントの数を有するトラス接続アレイのそれぞれ行と列の個数であり、P と Q がそれぞれ M と N に等しいことを特徴とする請求項 19 に記載の並列プロセッサ。

【請求項 34】

並列プロセッサであって、

次式によって定義された次元に従ったサイズ 4 の d 次元正規トラスに関する処理エレメントのクラスタを有し、

【数 1】

$$\text{reshape}\left(\prod_{q=1}^{d-1}(I_{4^{q-1}} \otimes G \otimes I_{4^{d-q-1}})\text{vec}(T), \underbrace{4, 4, \dots, 4}_d\right)$$

前記処理エレメントにより制御され、前記クラスタ間の相互に排他的な方向の処理エレメント間通信経路を多重化するように接続されることにより、トーラス接続アレイの接続性に等価の処理エレメント間接続性を提供するクラスタスイッチを有し、

ここで、 $d > 1$  であり、上記 reshape は、ベクトルからテンソルを作成する逆演算子で、 $\otimes$  は積の演算子、 $I$  は識別マトリクス、 $G$  は順列行列、

$\otimes$

はクロネッカー積、 $T$  は  $d$  次元の正規トーラスを表すテンソルを定義し、 $\text{vec}(T)$  は、テンソル  $T$  を用いてテンソル  $T$  の次元に沿って  $T$  の要素を積み重ねることによりベクトルに戻す演算子であることを特徴とする並列プロセッサ。

【請求項 3 5】

前記クラスタスイッチは更に、1つのクラスタ内の転置処理エレメント対における処理エレメント間直接通信を提供するように接続されることを特徴とする請求項 3 4 に記載の並列プロセッサ。

【請求項 3 6】

前記クラスタを組み合わせ、同時に、多重化を維持し、或いは前記クラスタスイッチによって、前記クラスタがスケラブルであることを特徴とする請求項 3 5 に記載の並列プロセッサ。

【請求項 3 7】

前記クラスタスイッチは更に 1つのクラスタ内超立方体補足対における処理エレメント間直接通信を提供するように接続されることを特徴とする請求項 3 6 に記載の並列プロセッサ。

【請求項 3 8】

並列プロセッサを形成する方法であって、

各クラスタが次式によって定義され、

$$\text{reshape}(G_N \text{vec}(T), N, N)$$

クラスタが他の少なくとも 1つのクラスタの処理エレメントと相互に排他的な方向においてのみ通信する、各  $N$  個の処理エレメントの  $N$  個のクラスタに処理エレメントを配列するステップと、

前記処理エレメントの制御の下で、前記相互に排他的な方向の通信を多重化するステップとを有し、

ここで、上記 reshape は、ベクトルからテンソルを作成する逆演算子で、 $T$  は正規トーラスを表すテンソル、 $G_N$  は順列行列を定義し、 $\text{vec}(T)$  は、テンソル  $T$  を用いてテンソル  $T$  の次元に沿って  $T$  の要素を積み重ねることによりベクトルに戻す演算子であることを特徴とする並列プロセッサを形成する方法。

【請求項 3 9】

それぞれが  $M$  個の処理エレメントを含む  $N$  クラスタに配列された  $N(M)$  個の処理エレメントを有する並列プロセッサであって、

処理エレメントの対が、クラスタ内、および、1つのクラスタ内の第 1 処理エレメントと前記第 1 処理エレメントを含むクラスタに隣接する 2つのクラスタの 1つのクラスタ内の第 2 処理エレメントとの間の利用可能な通信経路上で、単一ステップで通信するように前記処理エレメントを接続する処理エレメント間通信経路を有し、

1つのクラス内の全ての処理エレメントは完全に接続されており、前記処理エレメント

10

20

30

40

50

が分離された送信および受信ポートを備え、第1クラスタ内の任意の処理エレメントと隣接する第2クラスタ内の任意の処理エレメントとの間の通信が、前記第1および第2クラスタ内の全ての処理エレメントに対して、異なる処理エレメント対の間のM個の通信経路を介して同時に実施可能であることを特徴とする並列プロセッサ。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、データ処理システム及び方法の改良に関し、更に詳細には、改良された並列データの処理アーキテクチャに関するものである。

【0002】

【従来の技術】

多くの計算タスクは、データを並列演算するように開発できる。並列プロセッサの効率は、並列プロセッサのアーキテクチャ、コード化されたアルゴリズム、および、並列エレメント内データ配置に依存する。例えば、イメージ処理、パターン認識、および、コンピュータグラフィックスは、全て、2次元または3次元グリッド内に自然配列されたデータに作用する適用方法である。データは、オーディオ、ビデオ、SONAR信号、または、RADAR信号のような多種多様な信号を表す。離散コサイン変換(DCT)、離散逆コサイン変換(IDCT)、コンボリューション、および、この種データに関して一般に実施される演算等は異なるグリッドセグメント上で同時に実施可能であるので、一時に複数のプロセッサが特定タスクに作用できるようにすることによって、この種の演算を著しく加速できるマルチプロセッサアーレイシステムが開発された。並列処理は、ここに参考として組み込み済みの米国特許第5,065,339号、第5,146,543号、第5,146,420号、第5,148,515号、第5,577,262号、第5,546,336号、及び、第5,542,026号を含む多数の特許の対象とされている。

【0003】

並列処理アーキテクチャに関する従来型の一方法は、最隣接メッシュ接続コンピュータであり、これについては、全て参考としてここに組み込み済みであるR. Cypher、及び、J. L. C. Sanz、「SIMD Architectures and Algorithms for Image Processing and Computer Vision」(イメージ処理およびコンピュータビジョン用SIMDアーキテクチャ及びアルゴリズム)音響効果に関するIEEE議事録、音声および信号処理Vol. 37、No. 12、pp. 2158-2174、1989年12月、及び、K. E. Batcher「Design of a Massively Parallel Processor」(大量並列プロセッサの設計)コンピュータに関するIEEE議事録Vol. C-29、No. 9、pp. 836-840、1980年9月、及び、L. Uhr「Multi-Computer Architectures for Artificial Intelligence」(人工知能用マルチコンピュータアーキテクチャ) New York、N. Y.、John Wiley & Sons, Ch. 8、p. 97、1987年において検討されている。

【0004】

図1Aの最隣接トラス接続コンピュータにおいて、多重処理エレメント(PE)はトラス接続経路MPを介して、それらの北、南、東、西隣接PEへ接続され、全てのPEは同期的単一命令多重データ(SIMD)様式において処理される。トラス接続コンピュータは、循環(ラップアラウンド)接続コンピュータにメッシュ接続コンピュータを加えることによって得られるので、メッシュ接続コンピュータは、トラス接続コンピュータの部分集合を見なすことができる。図1Bに示すように、各経路MPはT送信ワイヤ及びR受信ワイヤを含み、または、図1Cに示すように、各経路MPはB双方向ワイヤを含むことができる。単方向および双方向通信は両方とも本発明の対象であるが、1つの経路において制御信号を除くバスワイヤの全個数は、今後一般的に、Kワイヤと称することとし、ここに、双方向バス設計においては $K = B$ であり、単方向バス設計においては $K = T +$

10

20

30

40

50

Rである。PEはその近傍PEの任意のPEにデータを伝送できるが、一時にただ1つに限られるものと仮定する。例えば、各PEは、1通信サイクルにおいてデータをその東隣に伝送可能である。更に、データ及び命令は1つの同報通信（ブロードキャスト）発送期間内にコントローラから全てのPEへ同時に発送可能であるような同報通信（ブロードキャスト）メカニズムが存在するものと仮定する。

#### 【0005】

通常、ビット直列PE間通信は配線の複雑性を最小限化するために用いられるが、それでもなお、トーラス接続アレイの配線の複雑性は実現上の問題を提起する。図1Aの従来型トーラス接続アレイは、PEの $4 \times 4$ アレイ10に接続される16個の処理エレメントを含む。各処理エレメント $PE_{i,j}$ は、それぞれ、その行番号*i*と列番号*j*がラベル付けされる。各PEは、2点間接続におけるそれぞれ北（N）、南（S）、東（E）、西（W）最隣PEに通信する。例えば、図1Aに示す $PE_{0,0}$ と $PE_{3,0}$ の間の接続は、 $PE_{0,0}$ のNインタフェースと $PE_{3,0}$ のSインタフェースの間の循環部接続であり、アレイをトーラス構成に形成する循環インタフェースの1つを表す。この種の構成において、各行は1組のN相互接続部を含み、N行には $N^2$ の水平接続部がある。同様に、それぞれN垂直相互接続部を有するN列には $N^2$ の垂直相互接続部がある。例えば、図1Aの場合、 $N = 4$ である。従って、例えば、循環接続部を含む $N \times N$ トーラス接続コンピュータの集積回路具体化における金属化ラインのようなワイヤの全個数は $2kN^2$ である。ここに、*k*は各相互接続におけるワイヤの個数である。ビット直列相互接続において数*k*は1に等しくてもよい。例えば、図1Aに示すように、 $4 \times 4$ アレイ10において $k = 1$ の場合、 $2kN^2 = 32$ である。

#### 【0006】

Nが比較的小さい幾つかの用途において、PEアレイ全体が1つの単一集積回路に組み込まれることが望ましい。ただし本発明は、各PEが、例えば、個別のマイクロプロセッサチップであるような実施形態を排除するものではない。トーラス接続されたコンピュータ内のワイヤの全個数は重要な意味をもつので、相互接続部が多量の貴重な集積回路の「不動産」又はチップの有効領域を消費することもあり得る。その上、PE相互接続経路は非常に頻繁に相互に交差し、ICレイアウトプロセスを複雑化し、おそらくは、漏話を介して通信線にノイズを導入する。更に、アレイの北と南および東と西端部におけるPEを接続する循環リンクの長さは、アレイサイズの増大につれて増大する。この長さが増大すると各通信ラインのキャパシタンスを増大させ、それによって、ラインの最大ビットレートを低下させ、当該ラインに更にノイズを導入することになる。

#### 【0007】

トーラスアレイの別の欠点は、転置操作に関連して起きる。処理エレメントとその転置は、通信経路における少なくとも1つの介在エレメントによって分離されるので、転置を用いる演算に待ち時間が導入される。例えば、 $PE_{2,1}$ がその転置 $PE_{1,2}$ からデータを必要とする場合には、当該データは介在する $PE_{1,1}$ 又は $PE_{2,2}$ を経て移動しなければならない。 $PE_{1,1}$ 及び $PE_{2,2}$ が占有されていない場合であっても、これは演算に遅延を自然に導入する。ただし、PEがマイクロプロセッサエレメントとして実現される一般的な場合には、 $PE_{1,1}$ と $PE_{2,2}$ が他の演算を実施し、データ又はコマンドを $PE_{1,2}$ から $PE_{2,1}$ へ転送するために、これらは、整然とした様式において、これらの演算またはコマンドを無視しなければならないという確率が極めて高い。従って、 $PE_{1,2}$ から $PE_{1,1}$ にデータを転送し始めるためであってさえも幾つかの演算が実施され、転置したデータを転送するために演算 $PE_{1,1}$ が強制的に取り消され、これが遅延となる。この種の遅延は、全ての介在PEと共に雪だるま式に増大し、最遠方の転置対に関してかなりの待ち時間が導入される。例えば、図1Aの $PE_{3,1} / PE_{1,3}$ 転置対は、最小限3つの介在PEを持ち、4つの通信ステップに相当する待ち時間を必要とし、なおその上に、一般的な場合には、 $PE_{3,1}$ と $PE_{1,3}$ の間でデータを転送するために、これら全てのPEにおいて取り消されるべき全てのタスクの待ち時間が生じる。

#### 【0008】

10

20

30

40

50

トラス接続アレイのこの種の限界を認識することによるアレイに関する新規技法が、参考としてそれらの全体がここに組み込み済みの「Massively Parallel Diagonal Fold Array Processor」、G. G. Pechanek等、アプリケーション特定アレイプロセッサに関する1993年国際会議、pp. 140 - 143、10月25 - 27、1993年、ベニス、イタリア、及び、「Multiple Fold Clustered Processor Torus Array」、G. G. Pechanek等、VLSI設計に関する第5NASAシンポジウム議事録、pp. 8.4.1 - 11、11月4 - 5、1993年、ニューメキシコ大学、Albuquerque、ニューメキシコに開示されている。これらのトラスアレイ組織の演算技法は、フォールドオーバーエッジとして従来型の最隣接トラスの対角PEを用いるPEアレイのフォールディングである。図2Aのアレイ20に示すように、これらの技法は、循環接続部の個数および長さが減少し、それらの転置PEに密接に近接してPEが位置するようにPE間配線を実質的に減少させるために利用可能である。このプロセッサアレイアーキテクチャは、例えば、それらの全体が参考としてここに組み込まれている米国特許第5,577,262号、第5,612,908号、EP0,726,532、EP0,726,529に開示されている。この種のアレイは、例えば、単フォールド対角フォールドメッシュのようなPE組合わせの不規則性により、従来型トラスアーキテクチャよりも実質的に優れた利益を提供すると同時に、幾らかのPEは2つのグループとしてまとめられ、その他のPEは単独のままである。3フォールド対角フォールドメッシュにおいては、4個のPEおよび8個のPEで構成されるクラスタがある。アレイ全体は三角形であるので、対角フォールド型アレイは、効率的で安価な集積回路の実現にとって本質的な障害となる。なおその上に、対角フォールドメッシュ及び他の従来型メッシュアーキテクチャにおいては、相互接続トポロジは本質的にPE定義の一部分である。この技法は、トポロジにおけるPEの位置を固定し、結果的に、PEのトポロジおよび実現される固定したコンフィギュレーションへのそれらの接続性を限定する。

#### 【0009】

多くの並列データ処理システムは超立方体相互接続トポロジを用いる。超立方体コンピュータは、高度の接続性を供給する方式で相互接続される $P = 2^d$  PEを含む。接続部は幾何学的または算術的にモデル化できる。幾何学モデルにおいて、PEは $d$ 次元超立方体の角に相当し、リンクは超立方体の縁に相当する。 $P = 2^d$  PEの超立方体は、各々が更に小さい超立方体の対応する角の間の接続部をもつ、 $2^{d-1}$  PEの2つの超立方体とみなすことができる。

#### 【0010】

算術モデルにおいて、各PEは、0から $d - 1$ までの一意的2進インデックスを割り当てられる。それらのインデックスの2進表現が厳密にただ1ビット位置だけ異なりさえすれば、任意の2つのPEが接続される。幾何学および算術モデルは、 $d$ 次元の各々を一意的ビット位置と提携させることにより相互に関連付けることができる。従って、1ビット位置だけ異なるインデックスを持つプロパティは2つの $(d - 1)$ 次元超立方体の対応する角を占有することに等価である。例えば、1つのPEには、トポロジ内のその位置を示すラベルを割り当てることができる。このラベル $\{D_0, D_1, \dots, D_{r-1}\}$ は2進表現であり、ここに、各数字は $r - D$ 超立方体上の通信に利用可能な1つの $r$ 次元接続経路を示す。殆どの場合、超立方体における各ノードは、その直接接続されたノードと $D$ だけ異なる1つの数字である。例えば、超立方体における最長経路は、PE $\{D_0, D_1, \dots, D_{r-1}\}$ とその補数 $\{\neg D_0, \neg D_1, \dots, \neg D_{r-1}\}$ の間、例えば、PE101101とPE010010の間の経路である。超立方体トポロジについては、ここに参考として組み込み済みのRobert Cypher、及び、George L. C. Sanz「The SIMD Model of Parallel Computation」(並列コンピュータのSIMDモデル)1994年Springer-Verlag、New York、pp. 61 - 68、及び、F. Thomas Leighton、「Introduction To Parallel Algorithms and Architectures: A

10

20

30

40

50



arrays, Trees, Hypercubes,」(並列アルゴリズム及びアーキテクチャ概論:アレイ、トリ、超立方体)1992年、Morgan Kaufman Publishers、Inc.、San Mateo、CA、pp.389-404に論じられている。超立方体トポロジの1つの欠点は、各プロセッサへの接続部の個数がネットワークのサイズと共に対数的に増大するという点である。その上、超立方体内のPE間通信には、特に、PEが相互に補数である場合、実質的な待ち時間が課される。

#### 【0011】

多次元超立方体は、トーラス、対角フォールドトーラス、又は、他のPE配列構成にマップ可能である。この種のマッピングについて、以下に簡単に検討することとする。この検討に係る図面およびこの出願に含まれる他の全ての図面においては、別途注記しない限り、各PE相互接続を単線として示すが、線(ライン)は双方向トライステートリンク又は2つの単方向リンクである相互接続リンクを表す。双方向トライステートリンクは、当該リンク上におけるデータ衝突を防止する制御スキーム下における、1つのリンク上での多重点における信号源(ソース)の生成をサポートする。単方向リンクは、あらゆるインターフェイス信号用として、二点間単一源および単一受信機対を用いる。更に、ビット直列および多重ビット並列具体化例についても検討する。

#### 【0012】

超立方体は、トーラス上にマッピング可能で、2次元トーラスはプロセッサエレメント(PE)で構成され、図1A及び1Dに示すように、各PEは、頂部PEラベルによって示されるように、トーラスノード(行と列)、及び、各PE内の底部ラベルによって示される超立方体PE番号と連携する。超立方体PE番号またはノードアドレスは、各数字が接続性次元を表すr次元(rD)超立方体に関するr数字表現として与えられる。超立方体内の各PEは、それらのノードアドレスがそれ自体から厳密に1つの数字だけ変化するこれらのPEのみに接続される。この相互接続スキームは、図1A及び1Dに示すように、4D超立方体が4×4トーラスにマッピングされることを可能にする。図1Aは、ただ1つの単一2進数字のみが順次数の間で変化するグレイコード符号化 $PE_{G(i),G(j)}$ を用いて $PE_{i,j}$ ノードをコード化する。例えば、10進数列0、1、2、3は、2進数列では00、01、10、11と表されるが、グレイコード数列は00、01、10、11となるはずである。図1Dは、最隣接トーラスへの代替超立方体マッピングを示す。

#### 【0013】

超立方体マシンの最も初期の具体化例の1つは、Caltech、C. Seitz「The Cosmic Cube」ACM通信、Vol.28、No.1、pp.22-33、1985年記載の6D超立方体であるコスミックキューブであった。コスミックキューブは、複数命令列複数データ(MIMD)モードで実行し、超立方体接続プロセッサ間で通信するためにメッセージ受け渡しを用いるインテル8086プロセッサによって実現された。他の超立方体具体化例NCUBEは、特注プロセッサが超立方体のノードを形式するチップを用いる1つの大型コンフィギュレーションとしての10-D超立方体から成る。NCUBEはMIMD型マシンであるが、単一プログラム複数データ(SPMD)モード作動もサポートする。この場合、各ノードプロセッサは同一プログラムのコピーを持つので、異なる条件付きコードストリームを独立的に処理することができる。Thinking Machines Corporationによって作成されたコネクションマシン(CM)は別の超立方体具体化例であった。初期のCM-1マシンは、ビット直列処理セルの4×4グリッドを含む各ノードを有する12D超立方体であった。

#### 【0014】

これらのような従来型の超立方体具体化例の1つの欠点は、各処理エレメントが、各超立方体次元に関する少なくとも1つの双方向性データポートを所有しなければならないことである。

#### 【0015】

以下に、更に詳細に検討するように、本発明の1つの態様は、我々のPEがネットワークトポロジから結合解除され、ただ1つの入力ポートと1つの出力ポートだけを必要とする

10

20

30

40

50

ことである。

【 0 0 1 6 】

更に、各追加された超立方体次元は、各 P E におけるポートの個数を増大するので、データポート専用 P E の占める割合が過度に大きくなって各 P E の設計が早急に非実際的となる。更にその上、「直径」としての待ち時間が更に一層大きくなることによって補足 P E 間通信の負担が大きくなる。即ち、超立方体の補足 P E 間の通信ステップの個数が拡大する。換言すれば、ノードアドレスとその補集合の間の接続を提供することにより、超立方体 P E ノード間の最長経路は、難しく、かつ高価になり、スケーラブルでなくなるはずである。

【 0 0 1 7 】

従って、プロセッサの並列アレイ内の処理エレメント間に高度の接続性を提供し、同時に、処理エレメントを相互接続するために必要な配線を最小限化し、かつ P E 間通信が遭遇する通信待ち時間を最小限化することが非常に望ましい。多重プロセッサアレイのアーキテクチャ及びプロセッサ相互接続の更なる改良の必要性が存在し、以下に、一層十分に検討するように、本発明は、これら及び他のこの種の必要性を扱う。

【 0 0 1 8 】

【発明の要約】

本発明は、処理エレメント間の接続性を改良する処理エレメントのアレイに向けられ、同時に、従来型トラス又は超立方体処理エレメントアレイの配線必要条件と比較して、アレイの相互接続の配線必要条件を実質的に軽減することに向けられる。好ましい実施形態において、1つのアレイは、転置操作の待ち時間および P E ノードとその超立方体補足(hypercube complement)ノードの間の通信待ち時間の実質的な減少を達成する。その上、このアレイは、循環配線の長さをアレイ全体の次元から切り離し、それによって、最長相互接続配線の長さを減少させる。同様に、通信 P E 間で衝突を引き起こすことのないアレイ通信パターンであるためには、特定のトポロジがその P E ノードから必要とする近傍接続部の個数と関係なく、1つの P E 当たり、ただ1つの送信ポートと1つの受信ポートだけが必要である。そのアレイの好ましい集積回路具体化例は、矩形または方形アウトラインを呈示するように組合わされた類似の処理エレメントクラスタの組合わせを含む。処理エレメントの類似性、処理エレメントクラスタの類似性、および、アレイ全体のアウトラインの規則性は、本アレイを特にコスト的にも効率的な集積回路の製造に適している。

【 0 0 1 9 】

本発明に従ってアレイを形成するには、先ず、処理エレメントを、単一命令多重データ(SIMD)演算の通信必要条件を利用するクラスタに組み合わせる。次に、そのクラスタ内で処理エレメントを完全に接続する。次に、処理エレメントをグループ化し、1つのクラスタのエレメントがクラスタ内において、および、ただ2つの他のクラスタの構成メンバーと通信するようにしてもよい。更に、各クラスタの構成処理エレメントは、ただ2つの相互に排他的な方向のみにおいて、他のクラスタの各々の処理エレメントと通信する。定義により、単方向能力を有する SIMD トラスにおいて、北/南方向は東/西と相互に排他的である。処理エレメントクラスタは、名称が示唆するように、相互に物理的に密接に近接して形成されることが好ましいプロセッサのグループである。集積回路の具体化例において、例えば、クラスタの処理エレメントは相互に出来る限り近接して配置されることが好ましく、アレイ内の他の任意の処理エレメントよりも相互に一層近接していることが好ましい。例えば、処理エレメント従来型 4 × 4 トラスアレイに対応するアレイは、それぞれ4つのエレメントで構成される4つのクラスタを含み、各クラスタは北および東方にのみ相互に通信し、南および西方に他のクラスタと通信し、または、南および東方に他のクラスタと通信し、北および西方に他のクラスタと通信する。このように P E をクラスタ化することにより、多重化を介して、P E クラスタ間の通信経路は共用可能であり、従って、アレイに必要な相互接続配線を実質的に減少させることができる。

【 0 0 2 0 】

好ましい実施の形態において、クラスタを構成する P E は、処理エレメント、それらの転

10

20

30

40

50

置、及び、超立方体補足 P E が同一クラスタ内に所在し、クラスタ間通信経路を介して相互に通信し、それによって、従来型トーラスアレイ上で実行される転置操作および従来型超立方体アレイ上で実施される超立方体補足 P E 間通信と連携した待ち時間を除去するように選定される。その上、従来型循環経路は任意の P E から P E への経路と同様に扱われるので、最長通信経路は、アレイ全体の大きさに関係なく、クラスタ間スペーシングと同じ程度に短い。

#### 【 0 0 2 1 】

各 P E は、仮想 P E アドレス記憶ユニット及びコンフィギュレーションコントローラを含む。仮想 P E 数およびコンフィギュレーション制御情報は、クラスタスイッチの設定を決定し、それによって P E アレイのトポロジを再構成するように組合わされる。この再構成は、例えば、コントローラからディスパッチされた命令に応答してなされても良い。アレイ内 P E はクラスタ化され、P E とその転置がクラスタ内で組合わされ、P E とその超立方体補集合が同一クラスタ内に含まれる。その上、各クラスタ内で完全に P E 間を接続するためのクラスタスイッチと組合わされた動的再構成は、多種多様なトポロジでアレイを再構成する能力を提供する。

#### 【 0 0 2 2 】

他の態様において、クラスタ内 P E は、当該クラスタ内で P E を完全に接続し、当該クラスタ内における各仮想 P E に 2 つの外部直交クラスタへの同一アクセスを可能にするクラスタスイッチに対して同一インタフェースを有利に所有することができる。さて、本発明の教示に従い、適所にクラスタスイッチを備えた 2 つのネットワークが実際に所在する。

1 つはクラスタ内 P E を相互に完全に接続するネットワーク、及び、もう 1 つは P E を他のクラスタ P E へ接続するネットワークであり、これによって、トーラスと超立方体の接続性に必要な接続経路を提供する。クラスタスイッチに対する内部接続経路は、転置と超補足体の接続性を提供する。異なる仮想 P E 配列構成に関して、転置はクラスタを横断して実施される。この種の 4 P E クラスタスイッチ、及び、他の 4 P E クラスタへのその相互接続に関して、任意のクラスタに対して生成されるただ 4 つの出力バスが有っても差し支えない。これら 4 つのバスの各々は、任意のクラスタにおいて、2 つの直交クラスタ接続点を持つ。本発明に係るマニフォールドアレイ処理において、強化された接続性超立方体が提供可能であり、この超立方体内には 4 個のノードの各クラスタはわずか 4 つの出力バスを持ち、各バスのファンアウトは 3 であり、1 つはスイッチに対し、1 つは直交クラスタの各々に対する。1 つの仮想ノード当たり 3 つの受信される信号がある。1 つはスイッチにとって内部信号であり、1 つは直交クラスタの各々からの信号である。

#### 【 0 0 2 3 】

本発明のこれら及び他の特徴、態様、及び、利点は、添付図面と共に次の詳細な記述から、当該技術分野における当業者にとって明白となるであろう。

#### 【 0 0 2 4 】

##### 【発明の実施の形態】

一実施形態において、本発明に係るマニフォールドアレイプロセッサは、一方の 1 つのクラスタのエLEMENT が、他方のただ 2 つのクラスタの部材と直接通信するようにクラスタまたはグループとしての P E を組み合わせ、各クラスタの組成処理ELEMENT は、ただ 2 つの相互に排他的な方向において、もう一方のクラスタの各々の処理ELEMENT と直接通信する。このように P E をクラスタ化することにより、P E クラスタ間の通信経路は共用可能であり、従って、アレイにとって必要な相互接続配線を実質的に減少させることができる。その上、各 P E は単一の送信ポートと単一の受信ポート、又は、2 方向性の場合、順次的、又は、タイムスライス通信実施の場合には、単一の送受信ポートを有する。その結果、個別 P E はアレイアーキテクチャから切り離される。即ち、各 P E が N 通信ポートを有する従来型の N 次元ハイパキューブ接続アレイとは異なる。単一の送信ポートおよび単一の受信ポートを用いる具体化例においては、アレイ内の全ての P E は送信と受信を同時に実施する。従来の 6 D ハイパキューブの場合において、これは、各 P E に対して、6 個の送信ポートと 6 個の受信ポートからなる、合計 12 個のデータポートを必要とする。本

発明の場合には、ハイパキューブ（超立方体）の大きさに関係なく、ただ1つの送信ポートと1つの受信ポートからなる、合計2つのがデータポートだけが必要である。上記のとおり、2方向、順次的、または、タイムスライスデータ通信が用いられる場合には、送受信データポートを1つの送受信データポートに組み合わせることができる。各PEは仮想PE記憶ユニット及びコンフィギュレーション制御ユニットを含む。仮想PE番号およびコンフィギュレーション制御情報は、通信の方向を制御し、PEアレイのトポロジを再構成するために、クラスタスイッチの設定を決定するように組み合わせられる。この再構成は、例えば、コントローラから発送された命令に回答してなされても差し支えない。1つのPEとその転置（トランスポーズ）が1つのクラスタ内で組み合わせられるようにアレイ内のPEがクラスタ化され、PEおよびそのハイパキューブ補集合は同じクラスタに含まれる。

10

#### 【0025】

本実施形態において、クラスタを含むPEは、処理エレメントとそれらの転置が同一クラスタ内に配置され、クラスタ内通信経を介して相互に交信するように選択される。説明のために、処理エレメントは従来のトラスアレイであるとみなされ、例えば、処理エレメント $PE_{0,0}$ は、従来型のトラスアレイの「北西」コーナ、即ち、行0と列0に位置する処理エレメントとみなされる。従って、新規クラスタアレイのレイアウトは、従来型のアレイプロセッサのレイアウトとは実質的に異なるが、従来のトラス及び新規なクラスタアレイの対応する処理エレメントに同一データが供給される。例えば、新規なクラスタアレイの0,0エレメントは、従来型のトラス接続アレイの0,0エレメントを作動させるのと同じデータを受け取るはずである。その上、本記述で用いる方向はトラス接続アレイの方向を意味するものとする。例えば、エレメント間の通信が北から南へ実施される場合に、これらの方向は、従来型のトラス接続アレイ内通信の方向を意味する。

20

#### 【0026】

PEは、単一命令ストリーム・単一データストリーム（SISD）型の単一マイクロプロセッサチップであっても差し支えない。包含されるコンセプトを示すために、以下の記述に限定されることなく、基本的なPEについて記述することとする。図3Aは、本発明の新規PEアレイ用の各PEとして用いられる適当な実施形態を示すPE40の基本構造を示す。説明を簡単にするために、インタフェース論理回路およびバッファは図示されていない。命令バス31は、SIMDコントローラ29からディスパッチされた命令を受け取るように接続され、データバス32は、メモリ33又はPE40にとって別の外部のデータソースからのデータを受け取るように接続される。レジスタファイル記憶媒体34は、実行ユニット36にソースオペランドデータを供給する。命令デコーダ/コントローラ38は、命令バス31を介して命令を受け取るように、かつバス21を経てレジスタファイル34内のレジスタに制御信号を供給するように接続される。ファイル34のレジスタは、経路22を経てそれらの内容をオペランドとして実行ユニット36へ供給する。実行ユニット36は、命令デコーダ/コントローラ38から制御信号23を受け取り、レジスタファイル34に経路24を経て結果を供給する。更に、命令デコーダ/コントローラ38は、Switch Enable（スイッチイネーブル）とラベルを付した出力ライン39にクラスタスイッチイネーブル信号を供給する。クラスタスイッチの機能については、図5及び図6の検討に関連して以下に更に詳細に検討することとする。データ又はコマンドのPE間通信は、Receive（受信）とラベル付けされた入力37において受け取られ、Send（送信）とラベル付けされた送信出力35から送信される。

30

40

#### 【0027】

仮想PE記憶ユニット42は、それぞれのストア43及び読み出し45ラインを介して命令デコーダ/コントローラ38へ接続される。仮想PE番号は、新規仮想PE番号を記憶ユニット42へ送るデコーダ/コントローラ38で受け取られた命令を介して、コントローラ29によってプログラム可能である。仮想PE番号は、接続ネットワークによって課される限界内において、コントローラ29によってトポロジ内の各PEの位置を動的に制御するために使用可能である。

#### 【0028】

50

コンフィギュレーションコントローラ 44 は、それぞれのストア 47 及び読み出し 49 ラインを介して命令デコーダ/コントローラ 38 へ接続される。コンフィギュレーションコントローラ 44 は、例えば現行コンフィギュレーションのようなコンフィギュレーション情報を供給し、制御情報をクラスタスイッチへ供給する。これらクラスタスイッチは、アレイ内の他の P E への P E の接続を制御する。デコーダ/コントローラ 38 は、コンフィギュレーションコントローラ 44 からの現行コンフィギュレーションと、仮想 P E 記憶ユニット 42 からの仮想 P E アドレスと、コントローラ 29 からの命令によって運ばれた、例えば「転置 P E 間通信」のような通信操作情報を組み合わせ、この情報をクラスタスイッチに伝達する。デコーダ/コントローラ 38 は、図 6 に関連してさらに詳細に検討するように、この情報を使用してクラスタスイッチに関する適切な設定を決定し、スイッチイネーブルインタフェース 39 を介してこの情報を伝送するスイッチ制御論理回路を含む。スイッチ制御論理回路、クラスタスイッチ命令デコーダ/コントローラ、および、コンフィギュレーションコントローラは、P E の境界外において、クラスタスイッチに組み込み可能である。新規 P E ノードはトポロジ接続から独立して定義されるので、これらの機能は分離可能である。本実施の形態において、全体の論理と全体の機能性は、制御機能が独立している場合であっても、制御機能を分離しないことによって改良される。

#### 【0029】

図 3 D において、クラスタスイッチ制御の更なる詳細を示すために、例えば、適当なクラスタスイッチ 60 を示す。このクラスタスイッチ 60 は、図示されるように、A、B、C、D の 4 グループに分割され、各グループは 4 入力マルチプレクサと 3 入力マルチプレクサから成る。これらのグループの各々は、P E クラスタ内の特定の P E と連携し、この連携は点線矢印によって示される。例えば、P E<sub>0,0</sub> は「A」グループのマルチプレクサ a<sub>1</sub> および a<sub>2</sub> と連携する。これらグループ内のマルチプレクサは、それらに関連する P E によって制御される。図に示すように、これらのマルチプレクサを制御することにより、正常な S I M D 動作モードが保存される。

#### 【0030】

S I M D モードの動作時において、全ての P E はコントローラ 29 から送信命令を受け取り、同期して作動する。P E の I D に依存する演算を一意的に指定する命令を含む全ての命令が、全ての P E にディスパッチされる。これらの命令は、全ての P E によって受け取られ、復号されてから、コントローラ 29 から送信された命令によるプログラム制御の下で選択可能な P E イネーブル/ディスエイブルフラグを用いて、命令内のオペコード、その命令コードの拡張フィールドオペコードに依存する P E の全て又は P E の部分集合によって実行される。オペコード及びその拡張フィールドは、受け取った命令を実行する 1 つのセットを含む P E の集合を指定する。P E イネーブル/ディスエイブルフラグは、P E が応答し得る活動レベルを決定する。例えば、以下に、適宜使用できるフラグを示す。

- ・ レベル 0 : 完全に不能化された P E 。
- ・ 受け取った全ての命令が N O P として扱われる。
- ・ レベル 1 : 部分的にイネーブルにされた P E : P E が制御情報を受け取る :
- ・ 例えば仮想 P E の I D、飽和 / 未飽和モード、等のような制御情報のローディングを可能にする。
- ・ 例えば読取り状態レジスタのように制御情報の記憶を可能にする。
- ・ 全ての演算および通信命令が N O P として扱われる。
- ・ レベル 2 : 部分的にイネーブルにされた P E ; P E が制御情報を受け取る :
- ・ 例えば、仮想 P E の I D、飽和 / 未飽和モード、等のような制御情報のローディングを可能にする。
- ・ 例えば読取り状態レジスタのように制御情報の記憶を可能にする。
- ・ 全ての演算命令は N O P として扱われる。
- ・ 全ての通信命令が実行される。
- ・ レベル 3 : 完全にイネーブルにされた P E :
- ・ 受取られた全ての命令が実行される。

## 【 0 0 3 1 】

所与サイズのマニフォールドアレイに関しては、例えば 4 D、5 D、または、6 D ハイパキューブの選定により、許容されたコンフィギュレーションが前以て決定されていても差し支えない。この種の一実施形態において、可能なノード識別は「ハードワイヤード(hardwired)」である。即ち、集積回路具体化例の一部として不揮発に固定される。次に、所与のコンフィギュレーションに関する仮想 P E 番号は、コントローラ 2 9 からアレイ内の全ての P E 4 0 に送られる単一命令によって示唆される。この命令は、適当な仮想 P E 番号をそれぞれの P E に割り当てるために、各 P E 内のデコーダ/コントローラ 3 8 によって解釈されることが好ましい。各デコーダ/コントローラ 3 8 は、P E およびコンフィギュレーションに関する仮想 P E 番号を含む各 P E 記憶エリア 4 2 内のそれぞれのロケーションに関して、効果的にテーブルルックアップ動作を実施できる。

10

## 【 0 0 3 2 】

その中で類似エレメントが、図 3 A の P E 4 0 の指定番号を共有する図 3 B の P E 4 0 ' は、命令デコーダ/コントローラ 3 8、および、レジスタファイル 3 4 に接続されたインタフェース制御ユニット 5 0 を含む。この制御ユニット 5 0 は、信号ライン 2 5 を経てデコーダ/コントローラ 3 8 から獲得された制御信号に基づいて、例えば並直列変換、データ暗号化、および、データフォーマット変換のようなデータフォーマット変換を提供する。P E 4 0 " を示す図 3 C の別の実施形態において、送信経路 3 7 は 1 つ又は複数の実行ユニット 3 6 によって生成され、受信経路 3 5 は、直接またはインタフェース制御ユニット 5 0 を介してレジスタファイル 3 4 に接続される。インタフェース制御ユニット 5 0 は、1 つ又は複数のライン 2 5 を経て命令デコーダ/コントローラ 3 8 から受信した制御信号に基づいてデータをフォーマット化する。このインタフェース制御ユニットによって実施されるデータフォーマット化は、例えば、並列から直列変換、直列から並列への変換、データ暗号化、および、データフォーマット変換を含んでもよい。

20

## 【 0 0 3 3 】

更に、代替 P E 4 0 " は、各 P E 4 0 " へのローカルメモリブロック 4 8 および 5 2 の追加を含む。図 3 A 及び 3 B からの、ロード経路バス 2 6 及びストア経路バス 2 6 ' の両方を含むデータバス 3 2 の詳細を図 3 C に示す。これらのバスは両方とも、3 状態(トライステート)技術を用いるか、多重一方向のバスで実現されても良く、また例えば 1 6 ビット、3 2 ビット、6 4 ビットのように、種々のバス幅を持つことができる。例えばアドレスおよびバスプロトコル信号のような様々な制御信号を適宜用いることができる。更に、バス 2 6 及び 2 6 ' の各々は、1 つがコントローラ 2 9 a の DMA ユニットによって直接的に制御される 2 つのバスとして実現可能である。コントローラロードバスは、例えば、内部コントローラレジスタの内容を P E へロードするために使用可能であり、読取りバスは、例えば状態ビットのような内部 P E レジスタの読み取り用に使用できる。コントローラ 2 9 は、コントローラ 2 9 をメモリへ接続するインタフェースライン 6 1 を経てこれらのバスへのアクセスを持つ。DMA ロードバスは、メモリ 3 3 から P E ローカルメモリ 4 8 へメモリのブロックのローディングに用いられ、DMA 読取りバスは、ローカルメモリ 4 8 からデータのブロックをメモリ 3 3 に記憶するために用いられる。DMA 機能はコントローラ 2 9 の一部分であることが好ましい。メモリスイッチ 4 6 は、バス 5 1 と 5 3 およびバス 5 5 と 5 7 を介して P E ローカルメモリ 4 8 を P E レジスタファイル 3 4 へ接続するために用いられる。同様に、メモリスイッチ 4 6 は、バス 2 6 と 2 6 ' 及びバス 5 5 と 5 7 を介してメモリ 3 3 を P E レジスタファイル 3 4 へ接続する。メモリスイッチ制御信号は制御インタフェース 5 9 を経て、ロードおよびストアユニット 2 8 から受け取られる。ロードおよびストアユニット 2 8 は、インタフェース 2 7 を介して命令デコーダ/コントローラ 3 8 からロードおよびストア命令情報を受け取る。この受信した命令情報に基づき、ロードおよびストアユニット 2 8 は、メモリスイッチ 4 6 のためのスイッチ制御を生成する。コントローラ 2 9 から P E のアレイに発送される全ての命令は、アレイの各 P E 内で同じ仕方で解釈される。この発送された P E 命令は個別 P E または P E のグループへ結合されない。

30

40

50

## 【 0 0 3 4 】

P E の  $4 \times 4$  アレイを図 4 に示す。それぞれ 4 つの P E を含む 4 つのクラスタ 5 2、5 4、5 6、5 8 は、図 4 のアレイに組合わされる。クラスタスイッチ 8 6 及び通信経路 8 8 は、1 9 9 7 年 6 月 3 0 日付けで提出され、参考としてその全体がここに組み込み済みである米国係属出願 0 8 / 8 8 5 , 3 1 0 にさらに詳細に説明された仕方においてクラスタを接続する。ただし、この図において、各処理エレメントは、2 つの入力と 2 つの出力ポートを有する好ましい実施形態として示されているが、クラスタスイッチ内での多重化のための他の層は、各 P E に関して入力用として 1 つ、及び出力用として 1 つの所定数の通信ポートを装備する。P E 当たり 4 つの近傍伝送接続部を有する一方向通信の標準トラスにおいて、すなわち、P E 当たりただ 1 つの送信方向がイネーブルにされる場合、各 P E において、4 つの多重化伝送経路が必要とされる。これは、P E の一部分として定義される相互接続トポロジに起因する。最終結果として、標準トラス内に  $4 N^2$  個の多数の伝送経路が所在する。マニフォールドアレイにおいて、等価接続性および無制限通信である場合、わずかに  $2 N^2$  個の多重化伝送通路が必要とされる。マルチプレクサおよび  $2 N^2$  個の伝送経路によって消費される領域は、 $4 N^2$  伝送経路によって消費される領域よりも著しく少ないので、このように伝送経路を減少させることは、集積回路の不動財部分を大幅に節減することを意味する。通信経路は、トラス接続アレイ内の通信方向に対応して N、S、E、W とラベル付けされる。

10

## 【 0 0 3 5 】

完全なクラスタスイッチ 8 6 の一実施の形態を図 5 の構成図に示す。北、南、東、西方向の出力は既に図示したとおりである。クラスタスイッチ 8 6 には他の多重化層 1 1 2 が追加される。この多重化層は A とラベル付けされた東 / 南方向の受信と、B とラベル付けされた北 / 西方向の受信との間で選択し、それによって、各 P E の通信ポートの必要条件を受信ポート 1 つ及び送信ポート 1 つに減少させる。その上、T とラベ付けされたクラスタ間の転置接続部を介して転置 P E  $_{1,3}$  と P E  $_{3,1}$  の間の多重化された接続が実施される。特定のマルチプレクサに関して T マルチプレクサイネーブル信号が出力されると、転置 P E からの通信は、そのマルチプレクサと連携した P E において受信される。

20

## 【 0 0 3 6 】

好ましい実施の形態において、全てのクラスタは、例えば P E とその転置 P E の間の経路のような転置経路を含む。これらの図は、全体の接続スキームを示すものであり、多層集積回路の具体例は、一般的な設計上の定例的な選択問題として実施されるルーチンアレイの相互接続を如何にして完全に達成するかを説明することを意図しない。集積回路のレイアウトに際して、I C 設計者は、本発明に係る実際の I C により実現されるアレイを取り入れるプロセスにおいて、各種のトレードオフを分析するであろう。例えば、多数のインタフェースの配線長さを減少させるために、クラスタスイッチは P E クラスタ 1 内に分散させてもよい。

30

## 【 0 0 3 7 】

多次元アレイをサポートし、かつ実現を簡素化する 4 P E クラスタにおける接続性の拡張に必要な多重化の変化を図 6 のクラスタスイッチ 6 8 6 に示す。説明を簡明にするために、別途注記されない限り、単方向性リンクであるものと仮定する。図 6 におけるマルチプレクサは、各データ経路入力と連携したイネーブル信号を有する。これらのイネーブル信号は、S I M D コントローラからの個別信号ラインによって、個別 P E 内において受信されて復号され、かつディスパッチ済み命令から、または、スイッチクラスタから間接的に生成可能である。個別制御メカニズムは、S I M D コントローラからのディスパッチ済み命令を受け取り、その命令を復号し、多重化イネーブル信号を生成するスイッチクラスタ内に、個別のデコーダ / コントローラメカニズムによって提供され得る。好適な本実施形態においては、クラスタスイッチ多重化イネーブル信号が P E 内で生成される。それぞれ 4 つの 4 から 1 への (4 to 1) 送信マルチプレクサ、及び 4 つの 3 から 1 (3 to 1) のマルチプレクサが、それぞれ 4 / 1、3 / 1 とラベル付けされ、この好ましい実施形態内に用いられる。

40

50

## 【 0 0 3 8 】

図 5 に示すクラスタスイッチ 5 8 6 までの図 6 に示す拡張部は、8 個の 2 入力送信マルチプレクサ {  $x_1, x_3, x_5, x_7, x_2, x_4, x_6, x_8$  } を、4 個の 4 入力マルチプレクサ 4 / 1 により置換えられる。{  $A, B$  } および {  $A, B, T$  } とラベル付けされたイネーブル信号と関連する、図 5 に示す 4 個の 2 入力および 3 個の 3 入力の受信マルチプレクサは、4 個の 3 入力受信マルチプレクサ 3 / 1 によって置き換えられ、送信ラインのゲート / バッファリングがゲート / バッファ B 1 - B 8 によって追加される。4 個の 4 入力送信マルチプレクサは、クラスタ 5 2 内の 4 個の P E の間で完全な接続性を提供する。図 6 のクラスタスイッチ 6 8 6 は、図 5 のクラスタスイッチ 5 8 6 によって提供される接続性のスーパーセットを表す。クラスタスイッチにおいて、内部配線をレイアウトするには多くの方法があり、図 6 の表現は接続点を示すが、多層シリコンにより如何にして接続部を実現するかは示さない。

10

## 【 0 0 3 9 】

P E クラスタ 5 2、5 4、5 6、5 8 は、図 7 の構成図における  $4 \times 4$  マニフォールドアレイに組織される。送信ラインのゲート / バッファリングは、一般的な場合、即ち、P E クラスタ及びそれらのクラスタスイッチが同一シリコン上に配置されることなく、クラスタ間の信号用配線またはケーブル配置に要求される物理的距離だけ分離される場合に、必要とされる。更に、電力およびレイアウトの観点から、ノイズ、及び電力を減少させるために送信信号をゲート / バッファリングすることが重要であることもあり得る。ただし、本実施形態において、マニフォールドアレイ組織は、単一チップまたは集積回路に組み込まれ、4 個の P E クラスタから成るクラスタスイッチは、ゲート / バッファリング回路が除去可能であるように、物理的に近接してまとめて配置される。図 6 のクラスタ組織は、図 9 A のクラスタ 9 8 6 A において、このバッファリングが除去された状態を示す。

20

## 【 0 0 4 0 】

同様に、図 7 の  $4 \times 4$  マニフォールドアレイのバッファリングは、図 8 A に示すマニフォールドアレイ 8 0 0 A の好ましい単一チップ実現のために除去される。 $4 \times 4$  マニフォールドアレイ 8 0 0 A が図 8 A に示すように接続される場合には、クラスタ内の各 P E は、相互に直交する他の 2 つのクラスタへ接続可能である。図 8 B は、 $4 \times 4$  マニフォールドアレイ 8 0 0 B における 1 個の P E、即ち P E<sub>1,3</sub> に関する出力送信接続性を示す。 $4 \times 4$  マニフォールドアレイ 8 0 0 B における各 P E は、そのクラスタ内の別の P E に到達可能であり、他の 2 つのクラスタへ接続可能である。図 8 C は、アレイ 8 0 0 C のための  $4 \times 4$  P E ノード入力 (受け取る) 接続性を示す。本発明に係る  $4 \times 4$  マニフォールドアレイの好ましい一実施形態において、任意の 2 つのノード間の最大通信距離は 2 である。

30

## 【 0 0 4 1 】

クラスタ 9 5 2 内の P E を接続するクラスタスイッチ 9 8 6 B を含む  $2 \times 2$  マニフォールドアレイ 9 0 0 B を図 9 B の構成図に示す。この図における任意の P E は、クラスタ 9 5 2 内のあらゆる他の P E と通信可能である。例えば、P E<sub>00</sub> はデータを自分自身、P E<sub>01</sub>、P E<sub>10</sub>、又は P E<sub>11</sub> に送ることが可能であり、P E<sub>01</sub> は、P E<sub>00</sub>、それ自身、P E<sub>10</sub>、または、P E<sub>11</sub> と通信可能であり、当該クラスタ内の他の 2 個の P E に関しても同様である。転置操作に関しては、P E<sub>00</sub> 及び P E<sub>11</sub> は何もしないか、P E<sub>01</sub> が P E<sub>10</sub> と通信し、せいぜいこれら自身と通信するに過ぎない。超立方体 (ハイパーキューブ) の状況に関して、P E<sub>00</sub> は P E<sub>1,1</sub> と通信し、P E<sub>0,1</sub> は P E<sub>10</sub> と通信する。図 9 B の  $2 \times 2$  マニフォールドアレイ 9 0 0 B から  $4 \times 4$  マニフォールドアレイへの移行は、図 9 A においてクラスタ 5 2 に関して示すように、4 個の  $2 \times 2$  を接続するようにマルチプレクサの追加集合を加えることに関連する。図 9 A のクラスタスイッチ 9 8 6 A は、マルチプレクサ 9 9 0 の追加集合を含む。従って、本発明は、プロセッサ、ノード等の相互接続に、高度に柔軟かつスケラブルな方法を提供することが理解されるはずである。

40

## 【 0 0 4 2 】

図 1 0 ~ 図 1 3 の構成図は、マニフォールドアレイ 1 0 0 0 に関するそれぞれの最隣接東、西、北、南通信のためのそれぞれの経路を示す。各経路は矢印によって示される。各マル

50



チプレクサにおけるただ 1 つの入力経路が、所与の時点においてイネーブルにされる。データ転送を選定された経路上で実施するには、通信命令がコントローラ 29 から全ての P E へ発信される。P E は、この発信された P E 命令を受け取り、それを復号し、それぞれのそれらレジスタファイル 34 から選定済みデータを検索して取り出し、それをそれぞれのそれらクラススイッチ 86 に送る。図 3 A に示すように、スイッチイネーブル信号は、既にプログラムされている仮想 P E 番号およびコンフィギュレーション制御 44 出力と組合わされた、受信した命令から選定された通信情報に基づいて作られる。

#### 【 0 0 4 3 】

図 14 の構成図において、同じ  $4 \times 4$  マニフォールドアレイ 1000 に関する転置操作のための通信経路が示される。再び、アクティブなデータ経路は、経路に沿った方向性矢印によって示され、各マルチプレクサに対して、ある時間ではただ 1 つの入力がアクティブである。スイッチイネーブル信号は、図 10 ~ 図 13 に関して述べたのと同じ方法で形成される。図 15 は、 $4 \times 4$  マニフォールドアレイ 1500 における 4 個の独立した  $1 \times 4$  線形リング 1552、1554、1556、1558 のための通信経路を示す。 $1 \times 4$  線形リングは、行優先順を用いて  $2 \times 2$  から形成される。即ち、行優先順は、 $PE_{00}(A, B, C, \text{または}, D)$  から、 $PE_{01}(A, B, C, \text{または}, D)$  と  $PE_{10}(A, B, C, \text{または}, D)$  と  $PE_{11}(A, B, C, \text{または}, D)$  と  $PE_{00}(A, B, C, \text{または}, D)$  に至る線形経路である。 $PE_{00}$ 、01、10、11 の集合 A - D の各々は、P E の  $1 \times 4$  線形リングを構成する。

#### 【 0 0 4 4 】

図 16 の構成図は、本発明に係るマニフォールドアレイによって提供される融通性に関する更なる態様を示す。図 16 は、2 つのアレイとして構成される  $4 \times 4$  マニフォールドアレイ 1600 用の通信経路を示す。最上  $2 \times 2 \times 2$  アレイは「A」P E から成る。最下または「B」P E は、第 2 の  $2 \times 2 \times 2$  アレイを構成する。図 16 は、z 軸を介して通信する平面を用いた参照表記法(reference notation) (行)  $\times$  (列)  $\times$  (平面) を用いる。P E 間のこの種の通信は一般に、双方向性または単方向性どちらのポートが用いられるかに応じて、3 個または 6 個の通信ポート軸を必要とする。これとは対照的に、好ましいマニフォールドアレイの実施には、P E 当たり 1 つの入力ポートと 1 つの出力ポートを必要とするだけである。アレイ 1600 は、図 16 の相互接続スキームを作成するために、修正済みクラススイッチのスイッチ設定を用いて、図 11 ~ 14 に示すアレイ 1000 の最隣接通信のために用いられたのと同じ物理的配置の P E を備えた  $PE_{1652}$ 、1654、1656 のクラスタを使用できることが注記される。本発明に係るマニフォールドアレイ組織の多くの利点の 1 つは、コンフィギュレーション及び接続方法から強力な新規能力が生じることである。

#### 【 0 0 4 5 】

この強力な能力の更なる一例として、図 17 は、図に示すようにクラススイッチによって相互接続されたクラス 1752、1754、1756、1758 を備えたマニフォールドアレイ 1700 の 4 D 超立方体の具体化例を示す。図 17 において、最上 P E 番号、例えば 0, 0 又は 3, 1 は、図 1 及び 2 のように、トーラスにおける P E の位置を表し、下位の P E 番号、例えば 0000 又は 1001 は、超立方体におけるその位置またはアドレスを表す。標準超立方体 P E 間通信において、超立方体 P E 番号において、ただ 1 つのビットが変化するだけである。例えば図 1 において、 $PE_{0111}(PE_{12})$  は、 $PE_{0011}$ 、0110、1111、0101 と通信する。ここで各個別経路は、超立方体番号における 1 つのビットだけが変化している。図 17 において、同一クラスタ内に位置するのは P E とそれらの転置 P E だけでなく、P E s とそれらの超立方体補足 P E もその中に位置することが有利である。図 7 のアレイは、超立方体とその補集合の間の経路を含まない。

#### 【 0 0 4 6 】

図 18 は、本発明に係るマニフォールドアレイ 1800 を示す。このアレイにおいて、4 個の P E のクラス 1852、1854、1856、1858 の各々における P E は完全に

10

20

30

40

50

接続されている。N個のPEのN個のクラスタを有する $N \times N$ マニフォールドの場合、クラスタのグループ化は、1つのiと1つのjを選定し、次の論式を用いることによって形成され得る：任意のi,jおよび全ての“a” + {0,1,...,N-1}に関して： $PE_{(i+a) \bmod N, (j+N+a) \bmod N}$ 。例えばマニフォールドアレイ1800のような4D超立方体の場合には、クラスタノードは次のように適切にグレー符号化可能である： $PE_{G((i+a) \bmod N), G((j+N-a) \bmod N)}$ 。ここで、G(x)はxのグレーコード(交番2進コード)である。一般的な場合への拡張については、次に簡単に説明するマニフォールドアレイの数学的表現において検討することとする。

#### 【0047】

N = 4の場合における、超立方体ノードを持つ適当なマニフォールドアレイ1800を図18に示す。4×4マニフォールドアレイ1800は、超立方体補集合を接続する接続経路を含む。換言すれば、超立方体PEとその補足超立方体PEの間の通信経路がイネーブルされる。例えば、 $PE_{0111}$  ( $PE_{1,2}$ )は、 $PE_{1000}$  ( $PE_{3,0}$ )並びにそのクラスタの他のメンバと通信可能である。超立方体PEとその補足超立方体PEの間の通信の意味を考察すれば、この場合の4ステップに相当する最長経路用超立方体通信経路は1ステップに短縮される。この経路長短縮は、本発明に係る処理エレメントのマニフォールドアレイ組織に関する非常に効率的な超立方体型アルゴリズムの作成に関して大きい意味を持つ。その上、4個のPEで構成される4×4マニフォールドアレイクラスタは、従来技術における折畳みアレイの場合、同様の4D超立方体接続性のために8個のPEで構成されるクラスタを必要とする従来技術における実現と比較すると、PEとその超立方体補集合の間の通信リンクに関して低コストの解決方法を提供する。

#### 【0048】

PEにおける上部ラベルが4×4×2(行、列、平面)PE番号であり、底部PE番号が5D超立方体番号であるような5D超立方体1900を図19に示す。従来型5D超立方体の場合、各PEにおいて5個の双方向性、又は10個の単方向性接続ポートが必要とされる。5D超立方体が4×4×2マニフォールドアレイにマッピングされる図20に示すマニフォールドアレイ2000の場合、各PEにおいて、わずかに1個の単方向性または2個の双方向性ポートが必要とされる。更に、標準超立方体は、合計 $5N^2$  (N = 4)個の双方向性バス又は $10N^2$  (N = 4)個の単方向性バスを備えた $2^5$ 、即ち32個のPEを必要とする。図20の5D超立方体マニフォールドアレイ2000は、わずかに合計 $2N^2$  (N = 4)個の双方向性バス又は $4N^2$  (N = 4)個の、PEの全てのクラスタ間単方向性バスを必要とする。図20は、各クラスタの間に8個の送信および8個の受信経路を備えた単方向バスの場合を示す。図20に示すように、超立方体PEとそれらの補集合PEが、PEの同一クラスタ内に所在することが有利であることに注意されたい。その上、PEの各平面に関して、各PEとその最近隣接転置PEは、PEの同一クラスタ内に位置している。更に、各平面からの対応するクラスタエレメントと一緒にグループ化される。このグループ化は、図示されていない6D超立方体における、クラスタ当たり16個のPEに関しても真であることが保持される。

#### 【0049】

クラスタスイッチは、5Dの場合は4Dの場合と異なるが、同じレベルの相互接続性を提供する4Dの方法に類似した手法により組み立てられることに注意されたい。マニフォールドアレイ形成技法は、種々のサイズのクラスタを形成可能であることに注意されたい。クラスタサイズは用途および製品の必要条件に応じて選択される。

#### 【0050】

北、南、東、西、及び、Z軸入力/出力(I/O)ポートを備えた3Dトーラストポロジは、PEにつき6個の双方向性トライステート型リンクまたは12個の単方向リンクを必要とする。これは、 $N \times N \times N$ の3Dトーラスに関して、合計 $3(N^3)$ 個の双方向性トライステート型リンクおよび $6(N^3)$ 単方向リンクが必要とされることを意味する。4×4×4の3Dトーラスに関するマニフォールドアレイトポロジは、6D超立方体に関するマニフォールドアレイトポロジと区別できない。PEは、トーラスまたは超立方体の必要条

10

20

30

40

50

件に従ってラベル付けされる。更に、 $8 \times 8$  トーラスは、接続性必要条件の低下した  $4 \times 4 \times 4$  の 3 D トーラスのサブグラフと見なすことができる。3 D 立方体または 6 D 超立方体トポロジのマニフォールドアレイ実現に際して、P E は 1 つの送信ポート及び 1 つの受信ポートのみを必要とし、それぞれ、当該トポロジによって必要とされるポート数とは関係無い。3 D 立方体トポロジにおける必要位置へ P E を配置することは、当該 P E にとって外部メカニズムをスイッチングすることによって、適宜取り扱うことが可能である。マニフォールドアレイ 3 D トーラスに関しては、クラスタ間の配線複雑性は、現在一般的に必要とされる  $6(N^3)$  の代りに、 $3(N^3)$  リンクまたは  $(N^3)$  双方向性トラিসテート型リンクのわずかに 3 分の 1、および、単方向リンクに関してはわずかに  $2(N^3)$  だけを必要とするスイッチングメカニズムにおいて軽減される。これは、実現コストの実質的な減少を表す。

10

#### 【0051】

次に、本発明に係るマニフォールドアレイの様々な態様について数学的に説明する。超立方体は、次元につきサイズ 2 の正規トーラスである。例えば、1 つの 4 次元超立方体は、1 つの  $2 \times 2 \times 2 \times 2$  トーラスとみなされる。ただし、トーラスが比較的小さい次元である場合についての後続する討論においては埋込み（同相写像）を扱う。2 d 次元の超立方体は、次元 d および辺長 4 の正規トーラスに等価であるが、 $(2d + 1)$  次元の超立方体は  $(d + 1)$  次元および最終次元においてサイズ 2 である以外は、全次元における辺サイズ 4 のトーラスに等価である。最初に、d が自然数である場合、2 d 次元超立方体は、次元につきサイズ 4 の d 次元正規トーラスに等価である。

20

#### 【0052】

2 d 次元超立方体 H は、 $4^d = (2^2)^d$  に等しい  $2^{2d}$  個のノードから成る。前記ノード個数は、辺につきサイズ 4 の d 次元正規トーラス T のノード個数である。定義により、H の全てのノードは、各次元につき 1 ノードの割合で 2 d 個の他のノードに隣接する。T の全てのノードは、次元当たり 2 個の他のノードに隣接する。即ち、d 次元の正規トーラスにおいて、各ノードは他の 2 d 個のノードに隣接する。従って、H と T のノード数およびエッジ数は等しい。

#### 【0053】

それらの間における一対一の対応を定義するために、 $(i_1, i_2, \dots, i_d)$  を T のノードであるものとし、ここで  $i_j$  は次元 j に関するノード座標を表す。T は辺当たりサイズ 4 の d 次元正規トーラスであるので、 $i_j$  は、1 から d までの全ての j に関して 0 から 3 までの値をとる。次元 k において、このノードは、ノード  $(i_1, i_2, \dots, i_{k-1}, \dots, i_d)$  およびノード  $(i_1, i_2, \dots, i_{k+1}, \dots, i_d)$  に隣接する。ここで、演算  $i_{k-1}$  及び  $i_{k+1}$  はモジュロ(modulo) 4、即ち、トーラスの循環エッジをカバーするように、 $3 + 1 = 0$  及び  $0 - 1 = 3$  として実施されるものと仮定する。

30

#### 【0054】

T のノード  $(i_1, i_2, \dots, i_d)$  と H のノード  $(G(i_1), G(i_2), \dots, G(i_d))$  の間における一対一のマッピングについて考察することとする。ここに、 $G(0) = 00$ 、 $G(1) = 01$ 、 $G(2) = 11$ 、 $G(3) = 10$  は、2 数字グレーコード（交番 2 進符号）である。トーラスノードはタプル(tuple:集合)によりラベル付けされ、超立方体ノードは 2 進ストリングによりラベル付けされるが、 $(G(i_1), G(i_2), \dots, G(i_d))$  と表記する場合には、実際には、対応する 2 進ストリングの連結を意味する。この点および一対一マッピングの説明を明瞭にするために、3 次元正規トーラスからのノード  $(3, 1, 0)$  および  $(G(3), G(1), G(0))$  を連結することによって導出される 6 次元超立方体 100100 に関して対応するラベルについて考察することを提案する。

40

#### 【0055】

連続するグレーコードは 1 つの単一 2 進数字だけ異なるので、隣接するトーラスノードは同様に隣接超立方体ノードであり、その逆でもある。従って、H のノードとエッジの間、及び、T のノードとエッジの間の一対一マッピングが存在し、2 つのグラフが同じであることを意味する。従って、次元 2 d の超立方体は、次元につきサイズ 4 の d 次元正規トー

50

ラスに埋め込み可能である。

#### 【0056】

グレーコード（交番2進符号）およびグレーコードを用いた超立方体ノードのラベル表示スキームについての更なる定義に関しては、例えば、参考としてここに組み込み済みのF. Thomson Leighton「Introduction to Parallel Algorithms and Architectures: Arrays, Trees, Hypercubes」（並列アルゴリズム及びアーキテクチャ入門：アレイ、トリー、超立方体）Morgan Kaufmann、1992年、Ch. 3を参照されたい。

#### 【0057】

( $2d+1$ )次元の超立方体は、少なくともサイズ2を除く全ての次元に関してサイズ4の( $d+1$ )次元トラスに等価である。Leighton資料から、( $2d+1$ )次元の超立方体は、ページ393で検討されているように、接続されるそれらの対応するノードを持つ2つの $2d$ 次元超立方体とみなすことができることが分かる。しかし、 $2d$ 次元の超立方体は、辺につきサイズ4の $d$ 次元正規トラスに等価であるので、それらの対応するノードを接続することによる2つの $d$ 次元トラスの併合集合は、最後の次元に関してサイズ2の( $d+1$ )次元トラスである。

#### 【0058】

上記の検討から、 $2d$ 次元の超立方体は、次元 $d$ および辺長4の正規トラスに等価であることが理解されるはずである。同様に、( $2d+1$ )次元の超立方体は、サイズ2の最後の次元を除く全ての次元において辺サイズ4の( $d+1$ )次元トラスに等価である。

#### 【0059】

マニフォールドアレイグループ又はクラスタは、次に示すように、一般に、直径的に対面するノードを形成することが好ましい。 $d$ 次元超立方体の隣接するノードは、1つの単一2進数字だけが異なるので、それらのノードアドレスにおいて厳密に $d$ 数字だけ異なるノードは相互に最も遠く離れて所在し、それらは相互に直径的に対面する。最も遠く離れたノードのアドレスは、相互に2進補数である。従って、1つの所与ノードに直径的に対面するノードを、その補集合とも称する。

#### 【0060】

一例として、1つの2次元 $4 \times 4$ トラス、及び、図21Aに示すように、超立方体ノードラベルを付けた $4 \times 4$ の表として表記される対応する埋込み済み4次元超立方体について考察することとする。この表の行および列に沿った隣接エレメント間距離は1である。列2、3、4が1つの位置だけ上方回転された場合、第1と第2列の間の対応するエレメントの距離は2になる。列3及び4、次いで列4に対して同様に反復することにより求められる隣接列の対応エレメントを持つ列のエレメント間距離は2である。結果として得られる4Dマニフォールドアレイの表を図21Bに示す。

#### 【0061】

その表の各列は、4個のノードで構成される1つのグルーピング、又は、換言すれば、4次元超立方体上の最大距離は4であるので、直径的に対面する2対のノードを含むことが重要である点に注目されたい。直径的に対面する超立方体ノードが同じグループに属する場合には、その表の列がグループを定義する。

#### 【0062】

比較的高い次元のトラス、そして超立方体において、直径的に対面するノードのグルーピングは、最後の次元を除く各新規次元に沿った同じ回転によって達成される。

#### 【0063】

グループ形成のための置換を数学的に表記するために、テンソル多次元アレイを分解し、そのエレメントからベクトルを作るか、その逆を実施する2つの演算子が定義される。 $v_{ec}()$ 演算子は、単一のアーギュメント、即ち、テンソル $T$ をとり、テンソルの第1次元である列に沿って $T$ のエレメントを積み重ねることによりベクトルを返す。

#### 【数2】

$$T = \begin{bmatrix} 11 & 12 & 13 \\ 21 & 22 & 23 \\ 31 & 32 & 33 \end{bmatrix}$$

例えば、 $T$  が 2 次元テンソルである場合には、 $v = \text{vec}(T) = [11 \ 21 \ 31 \ 12 \ 22 \ 32 \ 13 \ 23 \ 33]^T$  である。一方、テンソル  $T$  は、アーギュメントとしてソース構造および次元のリストをとる演算子  $\text{reshape}()$  を用いることにより、ベクトル  $v$  から復元される。最後の例から、 $T$  は、 $\text{reshape}(v, 3, 3)$  を用いて再構築される。

【 0 0 6 4 】

2 つのマトリックス  $A$  と  $B$  のクロネッカー積は、 $A$  の対応エレメントによって基準化された  $B$  のコピーから成るブロックマトリックスである。即ち：

【 数 3 】

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix}$$

所要のグルーピングを表わすための多重アレイ  $T$  の操作は、マトリックスが順列マトリックスであり、行、列、及び、ベクトル当たり厳密に 1 つの単一ノンゼロエレメントを持つ直交マトリックスで、ベクトルが  $\text{vec}(T)$  である場合における、マトリックス・ベクトル積として定義される。先ず、マトリックス  $S$  を求めるために、サイズ 4 の上回転順列マトリックスが決定される。

【 数 4 】

$$S = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

【 0 0 6 5 】

同様に、マトリックス  $G$ 、即ち、実際にグルーピング置換を実施するブロック対角マトリックスが決定される。 $G$  の対角ブロックは  $S$  のべき乗（パワー）である。

【 数 5 】

$$G = \begin{bmatrix} S^0 & 0 & 0 & 0 \\ 0 & S^1 & 0 & 0 \\ 0 & 0 & S^2 & 0 \\ 0 & 0 & 0 & S^3 \end{bmatrix}$$

【 0 0 6 6 】

$T$  が  $4 \times 4$  トーラスである場合において、 $\text{reshape}(G \text{vec}(T), 4, 4)$  は

10

20

30

40

50

、上述の要求されるグルーピング特質を有する結果として得られるトーラスである。同様に、 $T$ が $4 \times 4 \times 4$ トーラスを表す場合において、要求されるグルーピングを定義する演算を次に示す：

【数 6】

$$\text{reshape}((G \otimes I_4) (I_4 \otimes G) \text{vec}(T), 4, 4, 4),$$

ここで、 $I_4$ は $4 \times 4$ 識別マトリックスである。

【0 0 6 7】

一般に、1次元につきサイズ4の $d$ 次元正規トーラス $T$ の場合における、グループ現示置換を次に示す：

【数 7】

$$\text{reshape}(\prod_{q=1}^{d-1} (I_{4^{q-1}} \otimes G \otimes I_{4^{d-q-1}}) \text{vec}(T), \underbrace{4, 4, \dots, 4}_d).$$

グルーピング置換えを $4 \times 4 \times 4$ トーラスに適用する例を次に示す：

( $I_4 \times G$ )による第1マトリックス乗算後の、4つの平面を次表に示す：

ここで $\times$ はクロネッカー積を示す。

【表 1】

000	110	220	330
100	210	320	030
200	310	020	130
300	010	120	230

001	111	221	331
101	211	321	031
201	311	021	131
301	011	121	231

002	112	222	332
102	212	322	032
202	312	022	132
302	012	122	232

003	113	223	333
103	213	323	033
203	313	023	133
303	013	123	233

( $G \times I_4$ )による第2マトリックス乗算後の結果を次表に示す：

ここで $\times$ はクロネッカー積を示す。

【表 2】

000	110	220	330
100	210	320	030
200	310	020	130
300	010	120	230

111	221	331	001
211	321	031	101
311	021	131	201
011	121	231	301

10

222	332	002	112
322	032	102	212
022	132	202	312
122	232	302	012

333	003	113	223
033	103	213	323
133	203	313	023
233	303	013	123

20

## 【 0 0 6 8 】

上回転の代りに下回転した場合には、直径的に対面するノードをまとめるという同じグルーピング特質が維持される。更に、ノード  $(i, j)$  がその転置対、ノード  $(j, i)$  と共にグルーピング化される場合には、そのグループはノードの対称対を含む。数学的には、下回転置換は上回転置換の転置である。同様に、順列マトリックスが直交する場合には、下回転置換は上回転置換の逆である。更に明確に、サイズ 4 の下回転順列マトリックスを次に示す：

30

## 【 数 8 】

$$S_4^T = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

40

## 【 0 0 6 9 】

ここに、サイズ 4 に関しては一切制約されないので、マトリックスのサイズについて言及することが必要であり、あらゆるサイズの 2 次元トーラス、又はより高い次元のトーラスの 2 次元サブグラフに回転を適用できる。同様に、対角線ブロックが  $S^T$  のべき乗（パワー）である場合には、全ての列に対応する回転を適用するために、マトリックス  $G$  の転置を決定する。

## 【 0 0 7 0 】

図 2 1 C に示す辺当たりサイズ 4 の 3 次元正規トーラスの前例に関して、第 1 マトリックスに  $(I_4 \times G^T)$  を乗算することによって得られる 4 つの平面を次に示す：

50

ここで×はクロネッカー積を示す。

【数 9】

$$G_4^T = \begin{bmatrix} S^0 & 0 & 0 & 0 \\ 0 & S^1 & 0 & 0 \\ 0 & 0 & S^2 & 0 \\ 0 & 0 & 0 & S^3 \end{bmatrix}^T$$

10

【表 3】

000	310	220	130
100	010	320	230
200	110	020	330
300	210	120	030

001	311	221	131
101	011	321	231
201	111	021	331
301	211	121	031

20

002	312	222	132
102	012	322	232
202	112	022	332
302	212	122	032

003	313	223	133
103	013	323	233
203	113	023	333
303	213	123	033

30

【 0 0 7 1 】

前述のグルーピングは、Z 軸からの斜視図として図 2 2 および図 2 3 に示される。第 2 マトリックスに  $(G^T \times I_4)$  を乗算することによって得られる結果を次に示す：

ここで×はクロネッカー積を示す。

40

【表 4】



000	310	220	130
100	010	320	230
200	110	020	330
300	210	120	030

131	001	311	221
231	101	011	321
331	201	111	021
031	301	211	121

10

222	132	002	312
322	232	102	012
022	332	202	112
122	032	302	212

313	223	133	003
013	323	233	103
113	023	333	203
213	123	033	303

20

## 【 0 0 7 2 】

この場合、直径的に対面するノードが第三次元、即ち、図 2 4 に示すように異なる平面内の同じ位置に沿ってグループ化されるばかりでなく、最初から 2 つの次元に関して対称的なノードも平面の第 2 の次元または行に沿ってグループ化される。図 2 4 においては、1 つのグルーピング 8 9 がハイライトされている。グルーピング 8 9 は、図のその次の集合に関する基準点として用いられる。図 2 5 は、クラスタの間で面が所有する接続性に基づいて A、B、C、D 平面を再整列する。図 2 6 は、サブクラスタの間の局部接続性に基づいて、各平面を  $2 \times 2$  サブクラスタに分割する。グループ 8 9 は、平面 A 内のサブクラスタとして示される。この段階においてマニフォールドアレイコアアーキテクチャは、図 2 7 A に示すように、識別された  $2 \times 2$  サブクラスタの各々を直接交換するために用いられる。  $4 \times 4 \times 4$  の P E ノード P E <sub>2,2,2</sub> 入力 (受信) 接続性を図 2 7 B に示す。マルチプレクサの追加集合は、ラベル  $x \times 1$ 、 $x \times 2$ 、 $x \times 3$  が追加されていることに注意されたい。  $4 \times 4 \times 4$  の各  $4 \times 4$  マニフォールドアレイ部分集合の場合には、追加された  $x \times \#$  型の 16 個のマルチプレクサの追加集合がある。これらのマルチプレクサは全て、ハイライトされた P E ノード 2、2、2 に関して、図 2 7 B に示されると同じ仕方において接続される。

30

40

## 【 0 0 7 3 】

ノードの同じグルーピングは、異なる置き換え順序を経て到達可能であることに注意することが重要である。以上に示したステップにおいて必要とされた置き換えを一緒に乗算する結果として得られる順列マトリックスを P とする。このマトリックスのあらゆる因数分解は、同じ結果を達成する異なる順序のステップに対応する。例えば、 $P = A_1 A_2 A_3$  を仮定する。ランダム順列マトリックス Q、R について考察することとする。順列 P のようなノードの同じグルーピングを達成する一連の異なる順列を得ることができる。例えば、 $P = A_1 Q^T Q A_2 R^T R A_3$  であるので、 $B_1 = A_1 Q^T$ 、 $B_2 = Q A_2 R^T$ 、及び、 $B_3 = R A_3$  と命名して、 $P = B_1 B_2 B_3$  を得ることができる。更に、グループのエレメント、又は、

50

グループの相対的順序付け、又は、両方の置換えを実施し、実質的に同じであるが異なって見えるノードのグルーピングに到達できる。

【 0 0 7 4 】

同様に、本発明に係るネットワークに基づくマニフォールドアレイの特質は、更に次に検討するように、多くの利点を有するネットワークを形成するためにノードを接続することにも有利に適用できる。

【 0 0 7 5 】

ネットワークの直径は、ノードの任意の対の間の最大距離である。ネットワークの直径は、2つのノード間通信に必要な最悪場合数のステップを含む。直径が小さければ小さい程、遠く離れたノード間通信に必要とされるステップ数は少なくなる。ネットワークの直径は小さいことが望ましい。d次元超立方体の場合、H'は相補ノードを接続するエッジをHに加えることによって生成される新規グラフである。s及びtはHの2つの相補ノードであり、vはHの他の任意のノードであるものとする。相補超立方体ノードの任意の対から任意の超立方体ノードvまでの距離の和は、超立方体の次元に等しいことが実証できる。即ち、相補ノードsとt、及び、ノードvを所与の1対であるとすれば、vからsとtを通る最短経路がある。

【 0 0 7 6 】

sからvまでの距離は、sとvの2進表現の差、例えばkである数字の個数に等しい。tはsの補集合数であるので、tとvの2進表現の差に相当する数字の個数は(d - k)に等しい。従って、sからvまでの距離はkであり、tからvまでの距離は(d - k)である。すなわち、2つの距離の和はdである。更に、sからvを経たtまでの経路は長さdであり、これが最短経路である。

【 0 0 7 7 】

更に、d次元超立方体の相補ノードを接続するエッジを追加すればグラフの直径を、dが偶数であれば半分に、dが奇数であれば(d + 1) / 2に減少させることができる。vをHのノードとして定義すれば、kおよび(d - k)は、Hの2つの相補ノードsとtからのそれぞれの距離である。一般性を失うことなしに、k < (d - k)と仮定する。そうすれば、Hの場合と同じ最短経路を使用できるので、H'におけるsからvまでの距離はkである。新規エッジを経てsを通る経路が相補ノードを接続するので、H'におけるtからvまでの距離は(k + 1)である。これは、H'のノードvとsのあらゆる対に関して、それらの距離は、dが偶数であればd / 2、又は、dが奇数であれば(d + 1) / 2を超過し得ないことを意味する。Hにおけるsからvまでの距離が(d + 1) / 2を越える場合には、H'におけるsの相補ノードtを通る最短経路は、d / 2未満の長さである。

【 0 0 7 8 】

d次元超立方体のネットワークの直径はdであり、相補ノード接続部の追加により、上記のように[d / 2]になる。前述の結果を下表に要約する。エッジ接続相補ノードだけが中央列において取り扱われることに注意されたい。第3列ラベル付けされたマニフォールドアレイは、本発明のこの態様に基づく構造に含まれるエッジの個数並びに2の一定ネットワーク直径を示す。

【 0 0 7 9 】

【表5】

10

20

30

40

	超立方体（マニフォルドアレイのサブグラフ）	相補エッジ付き超立方体（マニフォルドアレイのサブグラフ）	マニフォルドアレイ（本発明による）
ノード	$2^d$	$2^d$	$2^d$
エッジ	$d2^{d-1}$	$(d+1)2^{d-1}$	$2^{d-1}((4 * 3^{k-1})-1);$ for $d = 2k$
エッジ			$2^{d-1}((8 * 3^{k-1})-1);$ for $d = 2k + 1$
ネットワーク 直径	$d$	$\left\lceil \frac{d}{2} \right\rceil$	$2$

10

## 【 0 0 8 0 】

上記の表は、超立方体ネットワークの相補ノードを接続する超立方体より多くの  $2^{d-1}$  個のエッジを含むサブグラフが劇的な改良を起こすことを示す。ネットワークの直径は、超立方体に比較して、その元のサイズの半分に短縮される。本発明に従い上記の第 3 列に示すように、全個数のマニフォルドアレイエッジを備えている場合には、ネットワークの直径は、全ての  $d$  に関して一定な直径 2 に短縮される。超立方体および相補エッジを備えた超立方体は、マニフォルドアレイの適当なサブグラフである。

## 【 0 0 8 1 】

20

仮想ノードのエミュレーションは次のように実現可能である。より小さいネットワークによってエミュレートされることが必要な高次元ネットワークがあるものと仮定する。この必要条件是、多重仮想ノードが各物理的ノードによってエミュレートされなければならないことを意味する。エミュレーションが超立方体近傍並びに超立方体ネットワーク上の超立方体補集合および 2 次元トーラスネットワーク上のマトリックス転置を維持するように、仮想ネットワークを物理的ネットワークへマップするための幾つかの方法が以下に呈示される。超立方体エミュレーションは比較的簡単に達成可能である。q 次元の小さい方の超立方体にエミュレートすることが必要な d 次元超立方体を仮定する。次に、 $2^{d-q}$  個の仮想ノードが、各物理ノードによってエミュレートされなければならない。本発明に係る方法を非常に簡単に説明する方法は、ノードの 2 進アドレスについて考察することである。

30

d 次元超立方体ノードは、d 桁の数字による 2 進アドレスを必要とする。それら d 桁の数字のうちの q 桁の数字がエミュレーションを実施する物理ノードのアドレスを定義し、残りの  $(d - q)$  桁の数字が物理ノード内のローカル仮想ノード ID を定義する。

実際には、d 桁の数字アドレス

$$\overbrace{i_0 i_1 \dots i_{q-1} i_q \dots i_{d-1}}^{\text{仮想}}$$

$$\underbrace{i_0 i_1 \dots i_{q-1}}_{\text{物理的}} \underbrace{i_q \dots i_{d-1}}_{\text{ローカル}}$$

を持つ仮想ノード  $v$  に関して、このアドレスの最初の q ビットは、仮想アドレスのローカル ID セクションによって区別される  $2^{d-q}$  個の仮想ノードのグループをエミュレートする物理ノードの ID を示す。v のあらゆる隣接ノード  $w$  は、v のアドレスと単一の数字だけ異なる。この数字は、仮想 ID の最初の q 個の数字のなかのいずれか 1 つであり、従って物理ノードの近隣に所属するか、または、w のアドレスは  $(d - q)$  個のローカル数字のなかの 1 つだけ異なり、従って同じ物理ノードによってエミュレートされることを意味する。更に、仮想ノード  $v$  の補集合は、仮想アドレスの補数が物理アドレスおよびローカルアドレスの補数の連結に等しいので、v を賄う(host)物理ノードの補集合によってエミュレートされる。相補物理ノードは、マニフォルドアレイ内の同じクラスタに属するので、相補仮想ノードも同様に同一クラスタに属する。

40

## 【 0 0 8 2 】

50

一般に、仮想ノードIDは、必ずしも隣接していない物理およびローカルノードIDの2つの部分に有利に分割可能である。仮想ノードの隣接ノードは、IDの物理またはローカル部分のどちらかが異なる1つのノードIDを持つはずである。従って、仮想ノードは、それぞれ物理ノードの隣接ノード、又は同じ物理ノードによってエミュレートされる。更に、仮想IDの補集合はローカルおよび物理IDの補集合に等しいので、仮想ノードの補集合は、マニフォールドアレイ上の隣接でもある物理ノードの補集合によって常にエミュレートされる。

#### 【0083】

その代わりに、小さい方の超立方体が大きい方の超立方体によってエミュレートされる場合には、マニフォールドアレイネットワークは帰納的に定義されるので、全てがマニフォールドアレイのサブセットに予測通りに作用する。これは、前述の論理が保持される場合には、エミュレートされる超立方体に等しいサイズのマニフォールドアレイのサブグラフが存在することを意味する。

#### 【0084】

同じ概念がトーラスのエミュレーションに有効であるので、トーラスエミュレーションも同様に容易に扱うことができる。隣接して選定された仮想ノードID（次元当たり）は、物理およびローカルIDに相当する。物理ノードIDが仮想ノードIDの最上位ビットを含む場合には、仮想ノードがブロック分布している。そうではなくて、物理ノードIDが仮想ノードIDの最下位ビットを含む場合には、仮想ノードが循環分布している。ブロック分布の場合には、連続するIDを持つ仮想ノードのグループは、同じ物理ノードによってエミュレートされる。16個の仮想ノードが4個の物理ノード上にブロック分布している場合には、物理ノード0が仮想ノード0、1、2、及び、3をエミュレートし、物理ノード1が仮想ノード4、5、6、及び、7をエミュレートする等々である。循環分布の場合には、連続IDを持つノードのグループが異なる物理ノードによってエミュレートされる。16個の仮想ノードが4個の物理ノード上に循環分布している場合には、物理ノード0が仮想ノード0、4、8、及び、12をエミュレートし、物理ノード1が仮想ノード1、5、9、及び、13をエミュレートする等々である。

#### 【0085】

仮想アドレスに対して1を加算または減算することにより、指定された次元に沿ったこのノードの隣接ノードが見付けられる。この加算/減算は、仮想アドレスのローカルまたは物理部分のどちらか、又は、両方に1を加算/減算することに等価であり、同じ物理ノード又は隣接物理ノードのどちらかによって隣接仮想ノードがエミュレートされることを保証する。転置仮想ノードが、同じ物理ノード、又は隣接物理ノードによってエミュレートされることを保証するために、仮想アドレスの物理およびローカルセクションの割当は、全ての次元に関して同じでなければならない。すなわち、仮想ノードのブロック ブロックまたは循環 循環分布は、転置の近隣性を保存する。

#### 【0086】

並列マシンにおける超立方体マニフォールドアレイの例に戻って、データの配置は、アルゴリズムの高性能計算に関して最高の重要性を持つ。アレイプロセッサにおいては、処理エレメント間のデータ移動に起因する待ち時間を最小限化するために、最初データを適当なPE内に配置し、計算期間中に、直接接続されたPE間で移動させる。従って、全アルゴリズムに関して全体の通信待ち時間を最小限化するために、アルゴリズムがその計算段階を経て進行するにつれて、データの移動が最適化される必要がある。マニフォールドアレイの能力を実証するために、完全なシャッフルアルゴリズム及びPEとその超立方体補集合間の通信アルゴリズムが、図28A～図30に示す4×4マニフォールドアレイ2800上で調査されるはずである。テンソル積代数は、完全なシャッフルアルゴリズムをマニフォールドアレイプロセッサ上にマップするために用いられる。

#### 【0087】

テンソル積代数はクロネッカー積とも呼ばれ、数学方程式を目的のマシンアーキテクチャへのアルゴリズム的コーディングに適したマトリックス形式にマッピングするための方法

10

20

30

40

50

を表す。例えば、J. Granata、M. Conner、R. Tolimieri「The Tensor Product: A Mathematical Programming Language for FFTs and other Fast DSP operations」(テンソル積:FFTおよび高速DSP演算のための数学的プログラム言語)IEEE SPマガジン、pp. 40 - 48、1992年1月、及び、J. R. Johnson、R. W. Johnson、D. Roodriguez、および、R. Tolimieriによる「A Methodology for Designing, Modifying, and Implementing Fourier Transform Algorithms on Various Architectures」(各種アーキテクチャにおけるフーリエ変換アルゴリズムの設計、修正、および実行のための方法論)回路システム信号プロセスVol. 9、No. 4、pp. 449 - 500、1990年を参照のこと。これらの論文は両方とも参考としてここに組み込み済みである。

10

【0088】

マニフォールドアレイ完全シャッフル例のためのテンソル表記法による完全シャッフルは、J. R. Johnson等による参考資料のp. 472に用いられている

$$P^{2^n}_{2^n-1}$$

によって表される順列マトリックスとして定義されている。順列マトリックスは、特定マシン組織に対してロード又はストアされるべきデータにアクセスするためのアドレッシングメカニズムを定義すると一般的に解釈される。一般に、順列マトリックスは、その次に意図された計算的オペレーションに関して、データを適切な場所に置くために必要なデータの移動を表す。従って、データの移動、即ち順列マトリックスの目標アーキテクチャへのマッピングを最適化することが重要である。マニフォールドアレイ組織に関する、単一命令多重データ処理(SIMD)超立方体マシンとしての演算、即ち、完全シャッフルは、アレイ内にデータが適切に配置されれば容易に実行され得る。

20

$$P^{2^n}_{2^n-1}$$

$n = 5$  ( $P^{32}_{16}$ )の完全シャッフル例は、図28～図33を用いて記述される。図28Aは、完全シャッフルアルゴリズムを記述するために用いられるバス構造および演算ユニットを示す。図28Aには、多重コントローラ、メモリユニット0～3、および、特殊目的FIFOバッファが含まれる。この組織の好ましい実施形態では、メモリユニット、コントローラ、および、FIFOバッファがPEのアレイと同じチップ上に配置される。ただし、本発明は更に一般的であり、単一チップの実現を越えるものであることを理解されたい。マニフォールドアレイコンセプトは、マイクロプロセッサチップPEのアレイによって、ケーブル付きバス、外部メモリ、および、外部コントローラと共に容易に使用可能である。この検討のために、この種マシンのファミリ全体に互ってスケーラブルであると定義される単一アーキテクチャ及びマシン組織を可能にする単一チップ高性能アレイプロセッサが記述される。

30

【0089】

従って、低コストのために、コントローラ、メモリユニット、例えばFIFOのようなデータバッファ、および、PEは、全て単一チップに含まれる。コントローラは、それらのアレイクラスタ、メモリ、および、I/O機能へ、例えばメモリアドレス及びロード/ストア信号のような制御信号を介して、または、命令バス上で、例えば、PEに送られるディスパッチされた命令を介して制御を提供する。コントローラは、それぞれSIMDマシンにおける1つの一般的な機能ユニットであり、後続アルゴリズムのサポートという観点からのみ記述される。図28～図30において、特殊目的FIFOバッファとして示されるデータバッファは、一般に、メモリ/直接メモリアクセス(DMA)インタフェースユニットに同様に組み込まれ、この場合にも、ここではただ一般的にのみ記述される。

40

50

## 【0090】

図28Aは、クラスタ化されたPEの再構成可能なトポロジをサポートするために、多重コントローラがどのようにして適宜実行されるかを実証する。図28Aにおいて、コントローラは、4個のPEの各クラスタと連携する。他のスキームが使用可能であるが、コントローラ0はマスタコントローラとみなされる。命令バスI0は、それ自身のクラスタ、並びに、他のコントローラの各々と連携した命令スイッチ(ISW)に接続される。図28Aにおける各ISWは、命令バスI0のための出力C、または、コントローラの入力命令バスI1、I2、または、I3からそれぞれのISWの出力Cへの接続経路をイネーブルする。ISWは、マスタコントローラによって直接または間接的に供給される制御信号を用いて構成される。マスタコントローラも同様に、システム含まれ、かつこの種の情報を提供するように命じられた、例えばホストプロセッサのような他のプロセッサから、この情報を直接的または間接的に受け取る。この完全シャッフルの例のために、ISWはI0を全ての命令バス経路に接続するように設定される。これによって、コントローラ0は、図28Aに示すように、4個全てのPEクラスタ2852、2854、2856、2858のために、単一コントローラとして作動する。例えば、メモリサブシステムをデータソースへ接続する外部インタフェース経路を図28Aに示す。

10

## 【0091】

一例として、図28Bの底部において始まる、FIFOアドレス当たり8個のデータアイテムで構成されるグループとしてオンチップFIFOへ受信される一連のデータアイテムを示す。このアドレスは、FIFO内の各列の最上部に表示される。図28Bに示すように、データ項目{0-7}の第1グループはFIFO-0に記憶され、その次のグループ{8-15}はFIFO-1に記憶される、等々。この例において、FIFOは、次の図29に関連して次に記述される様式においてデータをロードするためのコントローラ、即ちこの例においてはコントローラ0、または、ローカルバッファ制御機能によってイネーブルされる4個のマルチプレクサに、各列データアイテムに対して1つずつ合計8個の出力経路を供給する。図28Bにおいて、図示される各バス、D0、D1、D2、D3はトライステート双方向性バス、又は、個別的にロードおよび個別的にストアするバスであり得る。1つ又は複数のバスは、意図した用途に適合する任意のデータ幅、一般に8ビット、16ビット、32ビット、または、64ビットであり。ただし、他の幅であっても差し支えない。

20

30

## 【0092】

説明を明瞭にするために、PEクラスタ間の相互接続は図示しない。図3A~3Cの基礎メモリブロックは、超立方体パターンにおけるデータ配置をサポートし、クラスタへのデータインターフェイス帯域を増大するために、図28の4x4マニフォールドアレイ2800におけるN=4のメモリブロックに拡大されている。メモリブロックは0-3にラベル付けされ、バス経路はクラスタ当たり1つに整理されている。例えばデータバスD0を備えたメモリ0は、クラスタA2852PE{(0,0)、(3,1)、(1,3)、(2,2)}へ接続され、データバスD1を備えたメモリ1は、クラスタB2854PE{(3,2)、(0,1)、(2,3)、(1,0)}に接続される等々。他のバス構造が可能であり、この典型的な記述によって排除されないことに注意されたい。

40

## 【0093】

FIFOバッファにデータがロードされた場合、図29に示すように、プロセスのその次のステップは、データを内部メモリ0-3(M0、M1、M2、M3)へロードする。この記述のために、データは32ビット、データバスは32ビットであるものと仮定する。次に示すシーケンスにおいて、一時に4個のデータアイテムを並列ロードするために、表記法メモリ データアイテム又はMx-aを用いた8回のFIFOからメモリへのロードサイクルが用いられる。図29に示すメモリユニットへロードされるデータパターンを生成するために用いられるシーケンスは、1<sup>st</sup>(M0-0、M1-1、M2-3、M3-2)、2<sup>nd</sup>(M0-6、M1-4、M2-5、M3-7)というように、これが8<sup>th</sup>(M0-31、M1-30、M2-28、M3-29)まで続く。例えば、メモリ2(M2)は

50

、図 29 に示す囲まれたアイテム 57 によって示される F I F O ラインからデータがロードされる。ここで、メモリデータは P E アレイ 2800 にロードされなければならない。メモリブロックはアドレスされ、コントローラ 0 によって制御される。P E ヘデータをロードするために、コントローラ 0 は、命令を P E ヘディスパッチし、それらのデータバスからデータがロードされるべきであること、および、P E 内部のどの場所へ当該データがロードされるべきかを通知する。次に、コントローラ 0 は、同期して、アドレスをメモリブロックに供給する。この場合には 4 個のメモリ 0 - 3 の各々に関して同じアドレスである。このアドレスと同期して、次に、コントローラ 0 は、アドレスされたロケーションからデータを読み出すため、またそのデータをメモリユニット自体のデータバスに配置するために、メモリユニットにとって必要な信号を生成する。

10

#### 【 0094 】

同期を保って、適当な P E s は、それらのデータバスからデータを取り、コントローラ 0 からディスパッチされた命令による指定に従って、それをロードする。コントローラ 0 は適当な P E の選択を識別する。この選択は、コントローラが P E ヘ送るディスパッチされた命令内の識別を介して、または、P E 内に位置するプログラム可能なイネーブル / ディスイネーブルビットを介して、等、幾つかの方法で実施できる。コントローラは、この結果を達成するために、一連の P E ロード命令を命令バスを介して各 P E ヘディスパッチする。32 データをロードするために、表記法：メモリユニット データアイテム P E # によって示される順序で並列に、ロードサイクル当たり 4 個のデータアイテムで、合計 8 回のメモリから P E へのロードサイクルが用いられる。図 30 に示す P E にロードされるデータパターンを生成するために用いられる順序を次に示す：1<sup>st</sup> (M0 - 0 - P E<sub>0,0</sub>, M1 - 1 - P E<sub>0,1</sub>, M2 - 3 - P E<sub>0,2</sub>, M3 - 2 - P E<sub>0,3</sub>) , 2<sup>nd</sup> (M0 - 6 - P E<sub>1,3</sub>, M1 - 4 - P E<sub>1,0</sub>, M2 - 5 - P E<sub>1,1</sub>, M3 - 7 - P E<sub>1,2</sub>) , 3<sup>rd</sup> (M0 - 9 - P E<sub>3,1</sub>, M1 - 11 - P E<sub>3,2</sub>, M2 - 10 - P E<sub>3,3</sub>, M3 - 8 - P E<sub>3,0</sub>) , 4<sup>th</sup> (M0 - 15 - P E<sub>2,2</sub>, M1 - 14 - P E<sub>2,3</sub>, M2 - 12 - P E<sub>2,0</sub>, M3 - 13 - P E<sub>2,1</sub>) , 図 30 に示すように 32 個のデータアイテムが P E アレイにロードされるまで継続する。データアイテムが正しい順序で読みとられた場合には、このロードパターンは完全シャッフルを実施する。完全シャッフル演算を順々に 32 のデータリストに実行することを図 31 に示す。

20

#### 【 0095 】

$X = (P^{32}_{16}) (P^{32}_{16}) (P^{32}_{16}) (P^{32}_{16}) (P^{32}_{16})$  X が知られている。  
この方程式は、マニフォールドアレイ上での通信事例を示すために用いられる。X の第 1 順列、即ち 32 エレメントベクトルは、図 30 に示すように、ロード操作によって達成される。次の順列は、4 つの隣接方向の各々に関して次のリストに定義済みであるように、P E の隣接対の間のス北ワップ演算によって実施される。

30

- ・東スワップ {0,0 & 0,1}, {0,2 & 0,3}, {1,0 & 1,1}, {1,2 & 1,3},  
                  {2,0 & 2,1}, {2,2 & 2,3}, {3,0 & 3,1}, {3,2 & 3,3},
- ・南スワップ {0,0 & 1,0}, {2,0 & 3,0}, {0,1 & 1,1}, {2,1 & 3,1},  
                  {0,2 & 1,2}, {2,2 & 3,2}, {0,3 & 1,3}, {2,3 & 3,3},
- ・西スワップ {0,0 & 0,3}, {0,1 & 0,2}, {1,0 & 1,3}, {1,1 & 1,2},  
                  {2,0 & 2,3}, {2,1 & 2,2}, {3,0 & 3,3}, {3,1 & 3,2},
- ・北スワップ {0,0 & 3,0}, {1,0 & 2,0}, {0,1 & 3,1}, {1,1 & 2,1},  
                  {0,2 & 3,2}, {1,2 & 2,2}, {0,3 & 3,3}, {1,3 & 2,3},

40

#### 【 0096 】

スワップ演算は、指定された P E 間におけるレジスタデータ値の交換を引き起こす。交換するレジスタ値の選択は、各 P E において受信されたディスパッチ命令に定義されている。この例の場合には、一方の P E におけるレジスタ R 1 は、もう一方の P E におけるレジスタ R 2 と交換またはスワップされる。図 31 は完全なシャッフルシーケンスを示し、P E の超立方体番号およびそれらに含まれるレジスタ R 1 及び R 2 を列形式においてリストする。スワップ ( 方向 ) 命令によって分離された各列は、スワップ演算の結果を示す。図

50

31に示すように、完全なシャッフルは、対構成されたPEの間におけるただ1回の最隣接データ移動の単一サイクルのみを必要とする各スワップ演算において実施される。

【0097】

この記述を更に拡張するために図32及び33が提供される。図32は、PEにディスパッチされた北スワップ命令の完了に際して得られるレジスタ結果を示す。図33は、PEにディスパッチされた南スワップ命令の完了に際して得られるレジスタ結果を示す。西スワップと東スワップは同様の仕方において処理される。この記述された事例の重要性は、データが指示通りの超立方体パターンでロード可能であれば、データの完全シャッフルを必要とする多くのアルゴリズムにとって、マニフォールドアレイ2000において非常に高速の処理が得られることである。

10

【0098】

最後に、マニフォールドアレイ超補足事例について記述する。この例においては、上記の完全シャッフル事例において指示されたようなスワップコマンドを用いて、超立方体PEと上述したそれらの補足超立方体PEとの間におけるレジスタ値の交換が実施される。超立方体ノード当たり単一のPEが存在するものと仮定すれば、超補足集合は超立方体マシンにおける最長経路に最適短縮を提供する。図18は、接続されたPE間における簡単な交換が1つの単一サイクルにおいて発生可能にするために超補足集合によって用いられる経路を示す。

【0099】

本発明は特定の好ましい実施形態および典型的なアプリケーションについて記述したが、本発明が多数のアプリケーションに適用可能であり、添付特許請求の範囲によってのみ限定されることが理解されるはずである。一例として、本実施形態は処理エレメントのクラスタを扱うが、ノードのクラスタも考慮対象とされる。この種のノードは、記憶されている多重ブロックのデータへの同時アクセスを可能にするタイル状のメモリシステムを形成するためのメモリエレメントであってもよい。更に、ノードは、接続ポート、入力/出力デバイス等であってもよい。一例として再度記述すれば、ノードは、通信ネットワークにおける複数の通信チャネルを接続するものでも良い。

20

【図面の簡単な説明】

【図1A】 従来の4×4最隣接接続型トーラス処理エレメント(PE)アレイのブロック図である。

30

【図1B】 図1Aの従来のトーラス接続経路が、どのようにT送信およびR受信結線を含むかを例示する図である。

【図1C】 図1Aの従来のトーラス接続経路が、どのようにBの2双方向性結線を含むかを例示する図である。

【図1D】 第2の従来の4×4最隣接接続型トーラスPEアレイの構成図である。

【図2A】 従来の折り畳まれたPEアレイの構成図である。

【図2B】 従来の折り畳まれたPEアレイの構成図である。

【図3A】 本発明に係るPEアレイ内で適宜使用される処理エレメントの構成図である。

【図3B】 本発明に係るPEアレイ内で適宜使用される処理エレメントの構成図である。

40

【図3C】 本発明に係るPEアレイ内で適宜使用される処理エレメントの構成図である。

【図3D】 本発明に係るクラスタスイッチ制御の更なる態様を示す図である。

【図4】 マニフォールドアレイ内PEのクラスタリングおよびPEのクラスタ間通信を示す構成図である。

【図5】 クラスタスイッチの更なる詳細を示す、図4の個別PEクラスタの構成図である。

【図6】 改良されたクラスタスイッチを備えた、本発明に係る改良されたPEクラスタの構成図である。

50

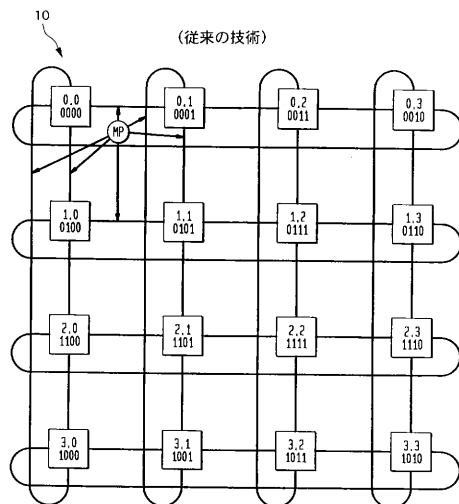


- 【図 7】 P E クラスタ間の相互接続経路をより詳細に示すブロック図である。
- 【図 8 A】 バッファを含まないクラスタスイッチの具体例を示すブロック図である。
- 【図 8 B】 バッファを含まないクラスタスイッチを用いる、P E クラスタ間の相互接続経路をより詳細に示すブロック図である。
- 【図 8 C】 バッファを含まないクラスタスイッチを用いる、P E クラスタ間の相互接続経路をより詳細に示すブロック図である。
- 【図 9 A】 大型アレイを形成するもう他の直交  $2 \times 2$  クラスタへの相互接続経路を示す  $2 \times 2$  クラスタのブロック図である。
- 【図 9 B】  $2 \times 2$  マニフォールドアレイのブロック図である。
- 【図 10】  $4 \times 4$  トーラス用の東側の通信経路を示す構成図である。 10
- 【図 11】  $4 \times 4$  トーラス用の西側の通信経路を示す構成図である。
- 【図 12】  $4 \times 4$  トーラス用の北側の通信経路を示す構成図である。
- 【図 13】  $4 \times 4$  トーラス用の南側の通信経路を示す構成図である。
- 【図 14】  $4 \times 4$  マニフォールドアレイの転置通信経路を示すブロック図である。
- 【図 15】  $4 \times 4$  マニフォールドアレイ上の 4 個の独立した  $1 \times 4$  線形リングを示す構成図である。
- 【図 16】 アレイコンフィギュレーションにおける  $z$  軸送信演算用の通信経路を示すブロック図である。
- 【図 17】  $4 \times 4$  マニフォールドアレイの超立方体ノードラベルを例示するブロック図である。 20
- 【図 18】  $4 \times 4$  マニフォールドアレイの超立方体補集合通信を示す構成図である。
- 【図 19】 5 D 超立方体を示す構成図である。
- 【図 20】 マニフォールドアレイにマッピングされた 5 D 超立方体を例示するブロック図である。
- 【図 21 A】 トーラスに埋め込まれた 4 D 超立方体のノードエレメントを示す表を示す図である。
- 【図 21 B】 埋め込まれた超立方体ノードの配置を示す改良されたマニフォールドアレイの表を示す図である。
- 【図 21 C】 頂部 P E ラベル ( P E - $x, y$  ) 付き  $8 \times 8$  の 2 D トーラス、中央 P E ラベル (  $x, y, z$  ) 付き 3 D 立方体、および、底部 P E ラベル  $G_x G_y G_z = d_5 d_4 d_3 d_2 d_1 d_0$  ラベル付き 6 D 超立方体の  $4 \times 4 \times 4$  の表現図である。 30
- 【図 22】 列の 1 D 下方回転後における  $4 \times 4 \times 4$  表現図である。
- 【図 23】 図 22 のノードの  $4 \times 4 \times 4$   $z$  平面表現図である。
- 【図 24】 列の 1 D 下方回転後における  $4 \times 4 \times 4$   $z$  平面表現図である。
- 【図 25】 レイアウト接続性に関する  $z$  平面表現の再順序付けを示す図である。
- 【図 26】  $z$  平面の  $2 \times 2$  サブクラスタへの分離を示す図である。
- 【図 27 A】  $2 \times 2$  サブクラスタの 4 個の  $4 \times 4$  マニフォールドアレイへの相互接続を示す図である。
- 【図 27 B】  $4 \times 4 \times 4$  P E ノード入力 ( 受信 ) 接続性の一例を示す図である。
- 【図 28 A】 クラスタ構成当たり 1 つの単一コントローラおよび典型的インタフェースを示す  $4 \times 4$  マニフォールドアレイの構成図である。 40
- 【図 28 B】 その外部インターフェイス内へ 32 個のデータアイテムを受け取る図 28 A の  $4 \times 4$  多重コントローラマニフォールドアレイを示す図である。
- 【図 29】 一例における 32 個のデータアイテムの 4 個のメモリコントローラへのローディングを示す図 28 A の  $4 \times 4$  多重コントローラマニフォールドアレイを示す図である。
- 【図 30】 各クラスタにおける個別 P E への 32 個のデータアイテムのロード配分を例示する、図 28 A の  $4 \times 4$  多重コントローラマニフォールドアレイを示す図である。
- 【図 31】 図 28 A ~ 図 30 に関する完全シャッフル例の各ステップの後における 32 個の典型的データをリストする表を示す図である。
- 【図 32】 スワップ北側通信演算を実行するために P E 間でデータが通る経路および完 50

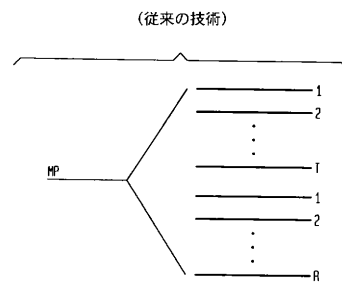
全シャッフル例の通信演算完了時における P E レジスタにおける結果を示す  $4 \times 4$  マニフォールドアレイを示す図である。

【図 3 3】 スワップ南側通信演算を実行するために P E の間でデータが通る経路および完全シャッフル例の通信演算完了時における P E レジスタ内の結果を示す  $4 \times 4$  マニフォールドアレイを示す図である。

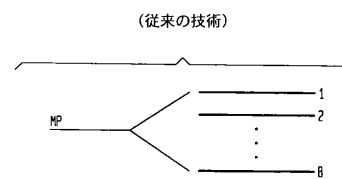
【図 1 A】



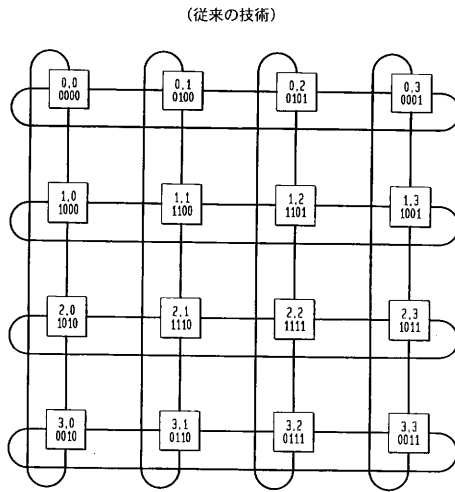
【図 1 B】



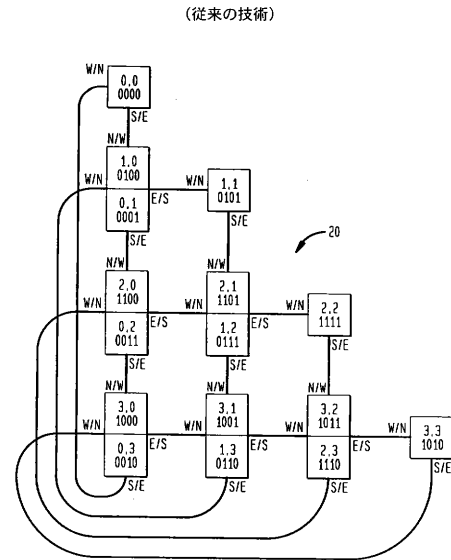
【図 1 C】



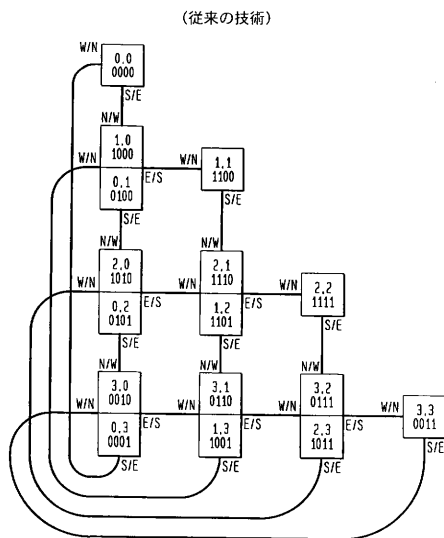
【図 1 D】



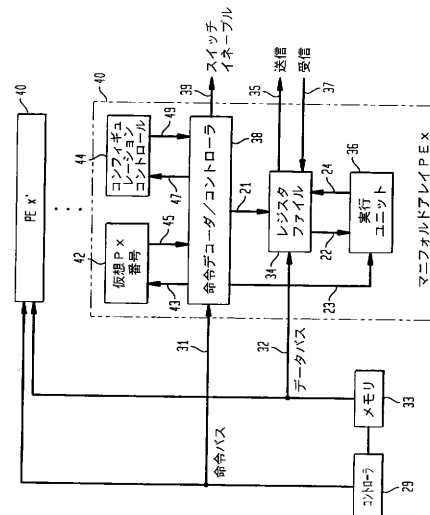
【図 2 A】



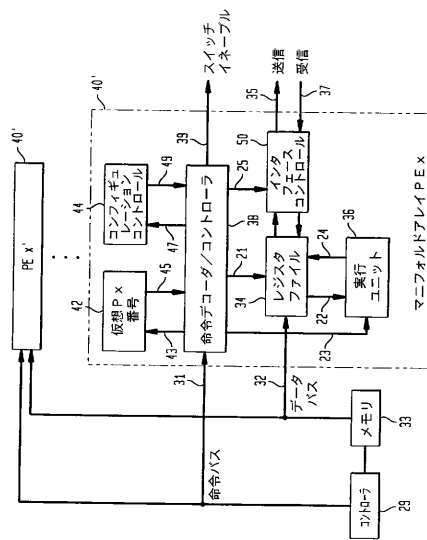
【図 2 B】



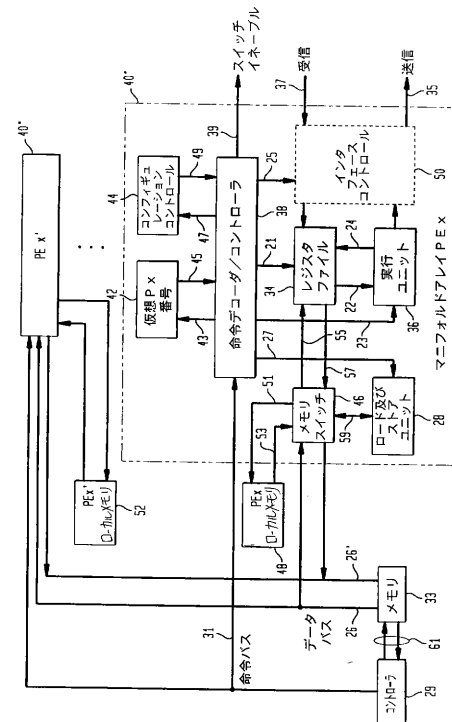
【図 3 A】



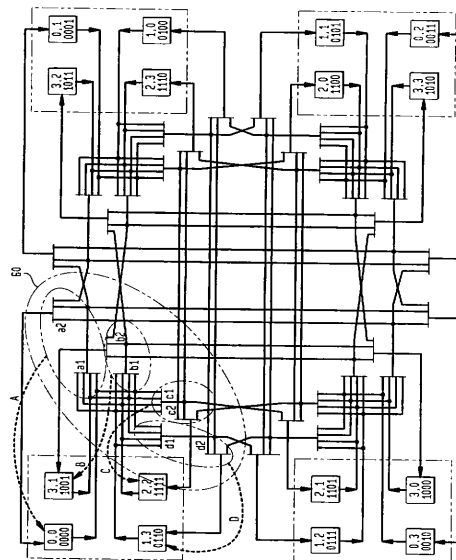
【 図 3 B 】



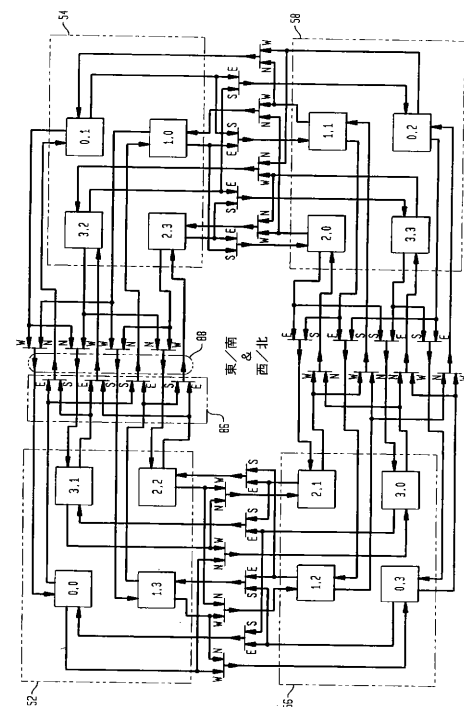
【 図 3 C 】



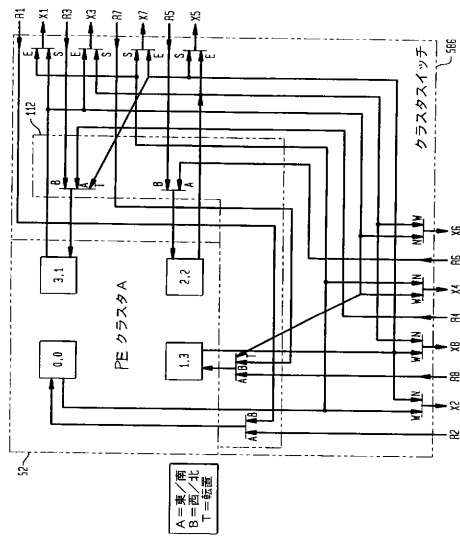
【 図 3 D 】



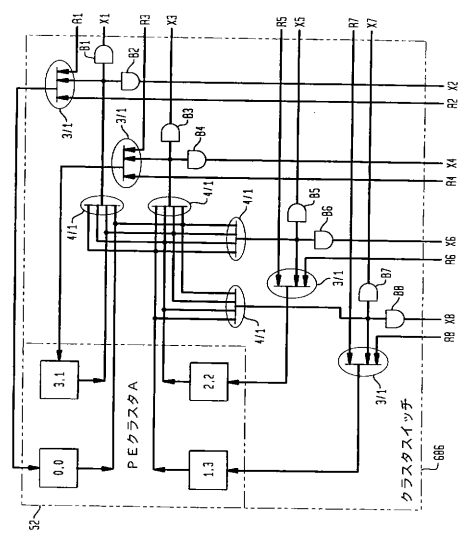
【圖 4】



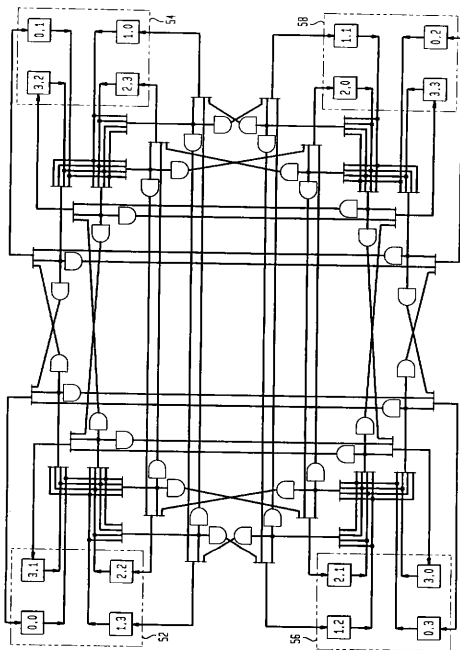
【図 5】



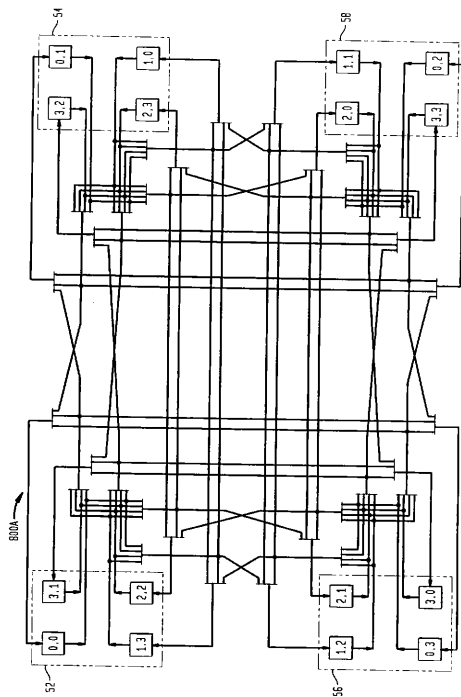
【図 6】



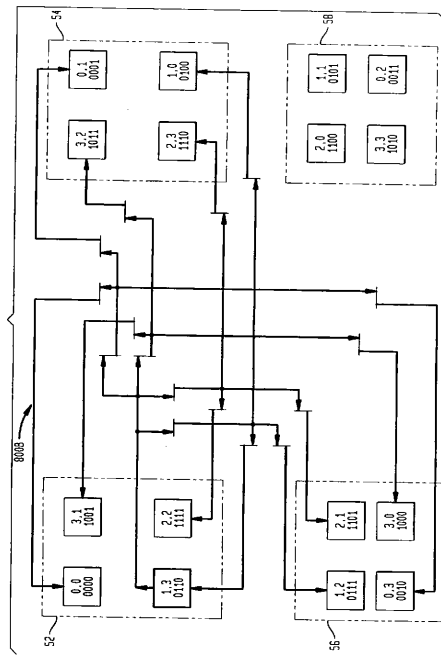
【図 7】



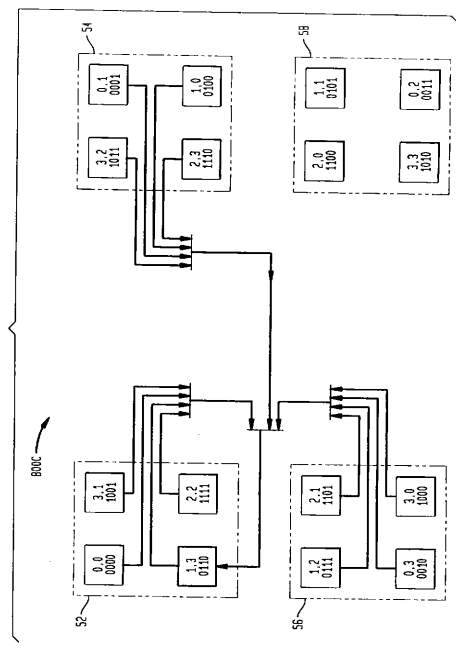
【図 8 A】



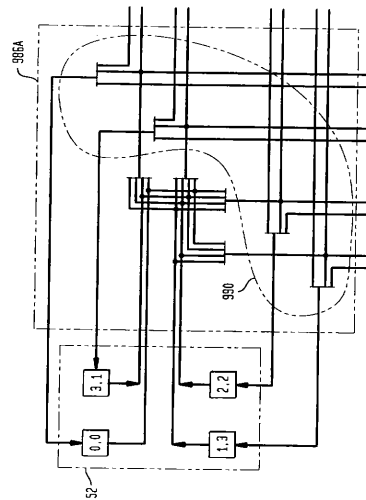
【図 8 B】



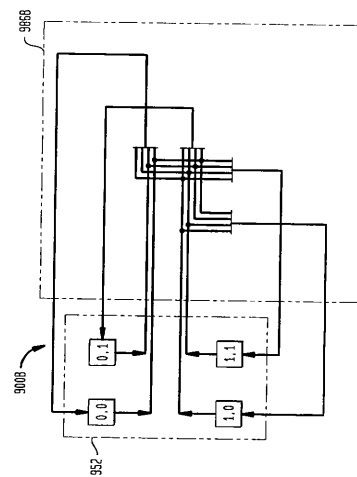
【図 8 C】



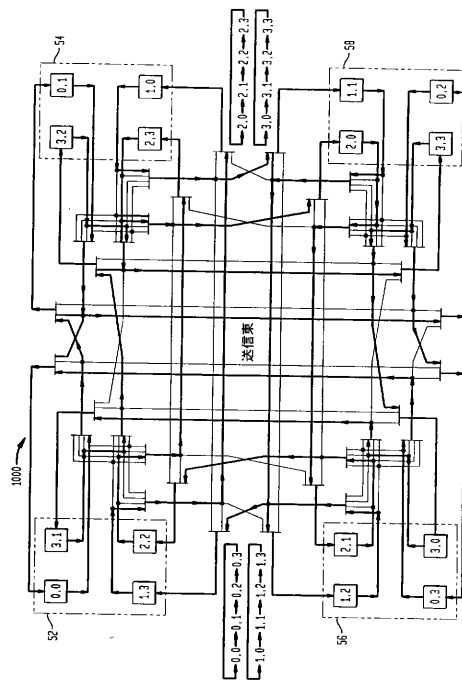
【図 9 A】



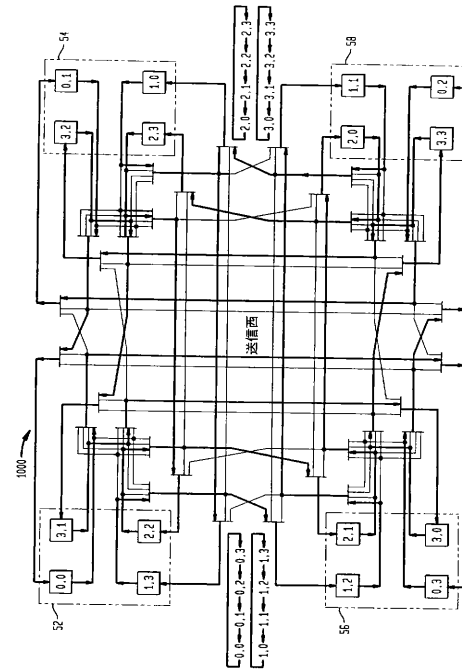
【図 9 B】



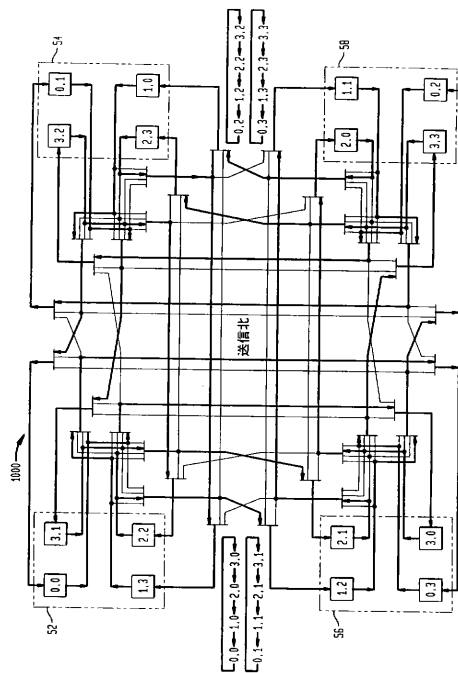
【図 10】



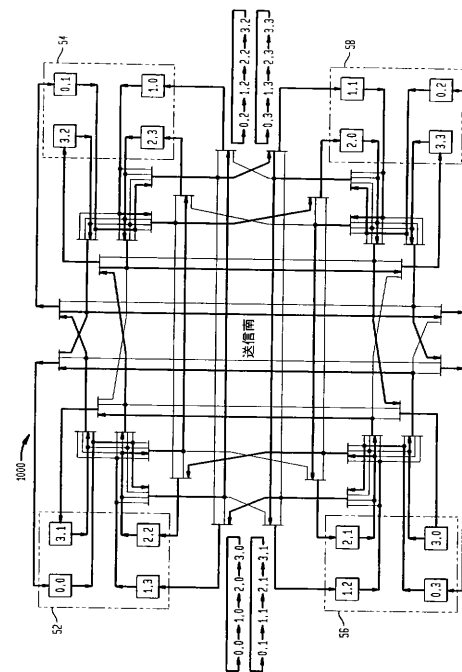
【図 11】



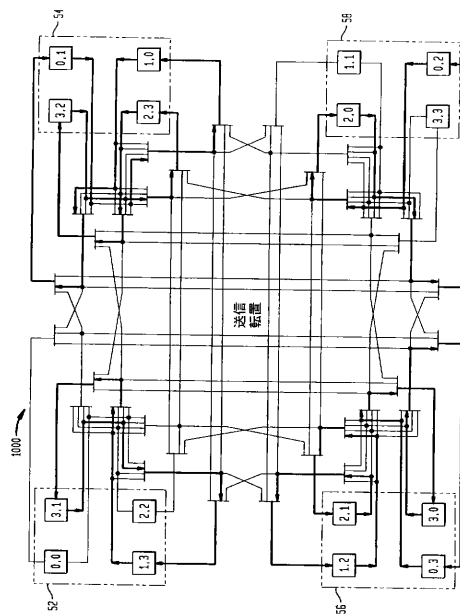
【図 12】



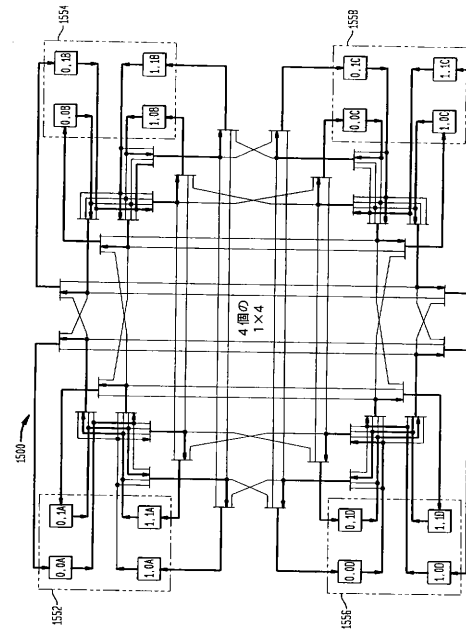
【図 13】



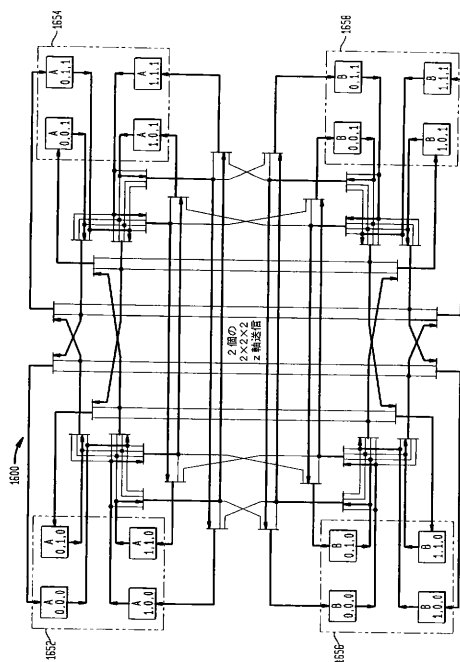
【図 14】



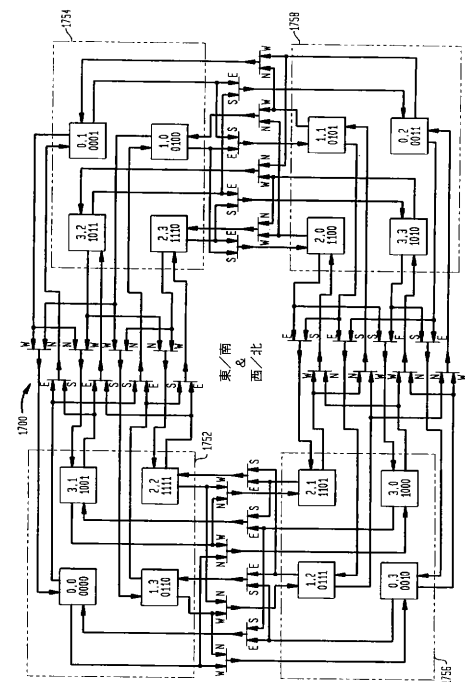
【図 15】



【図 16】

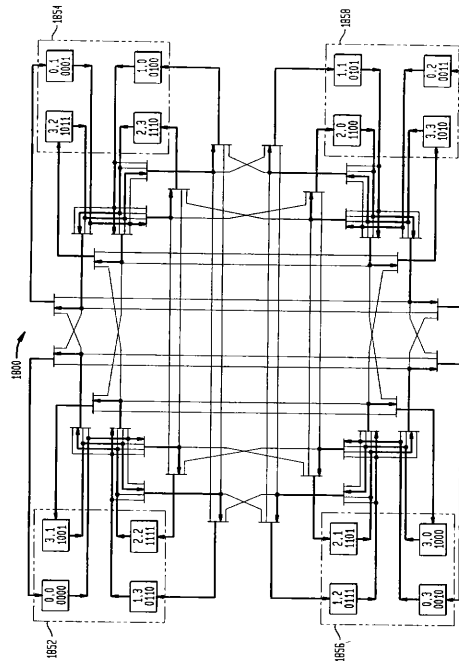


【図 17】

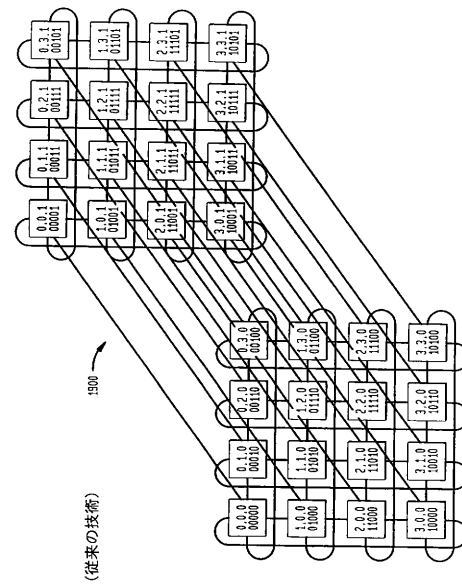




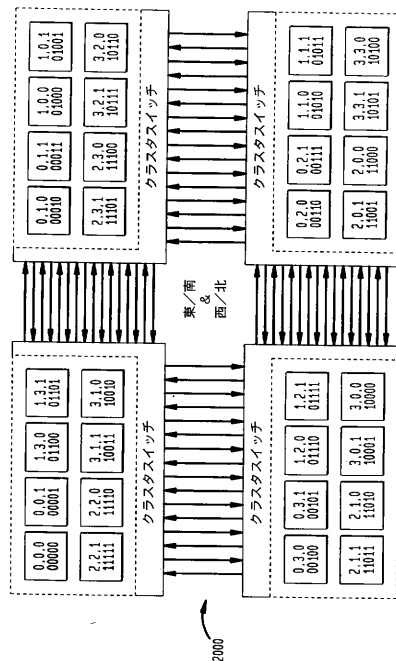
【図 18】



【図 19】



【図 20】



【図 21 A】

1	2	3	4
PE-0.0 0000	PE-0.1 0001	PE-0.2 0011	PE-0.3 0010
PE-1.0 0100	PE-1.1 0101	PE-1.2 0111	PE-1.3 0110
PE-2.0 1100	PE-2.1 1101	PE-2.2 1111	PE-2.3 1110
PE-3.0 1000	PE-3.1 1001	PE-3.2 1011	PE-3.3 1010

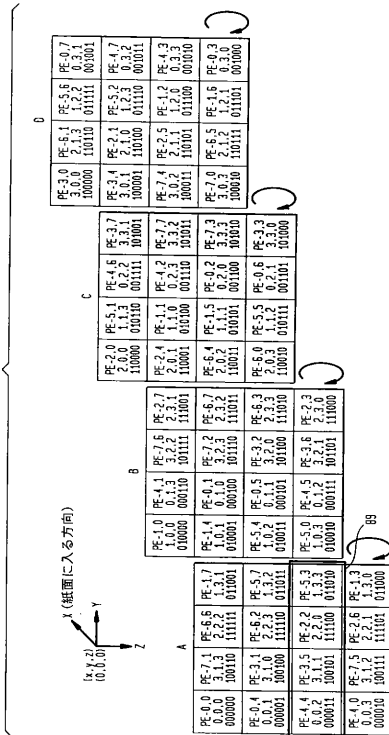
【 図 2 1 C 】

Figure 1 shows a 3D coordinate system with axes X, Y, and Z. The Z-axis is vertical, the X-axis is horizontal to the right, and the Y-axis is diagonal down and to the left. A point is labeled (x, y, z) with values (0, 0, 0).

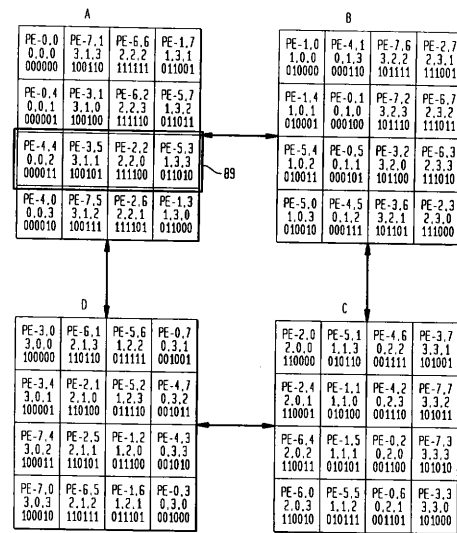
【 図 2 3 】

[illegible]

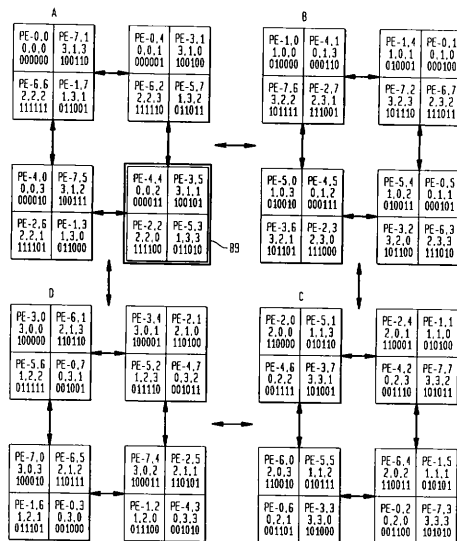
【図 24】



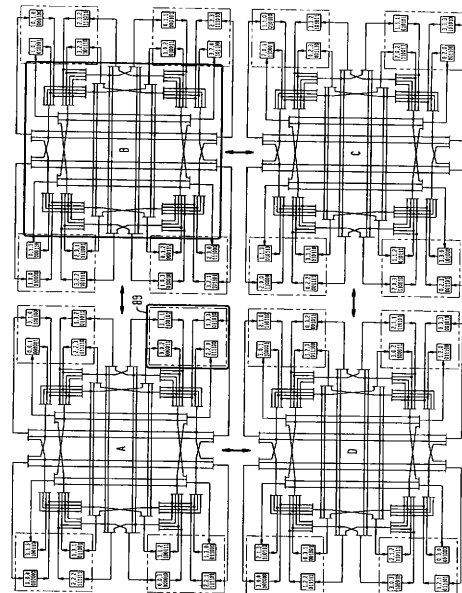
【図 25】



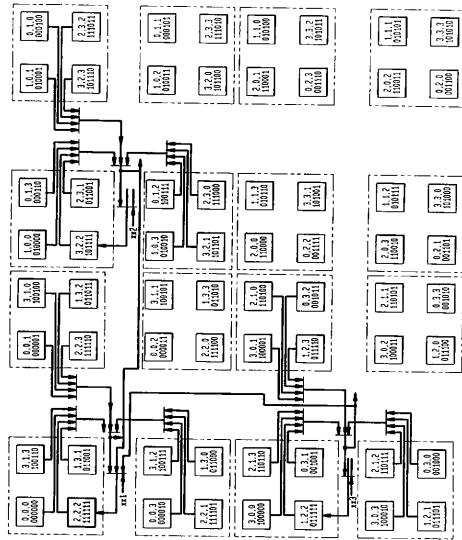
【図 26】



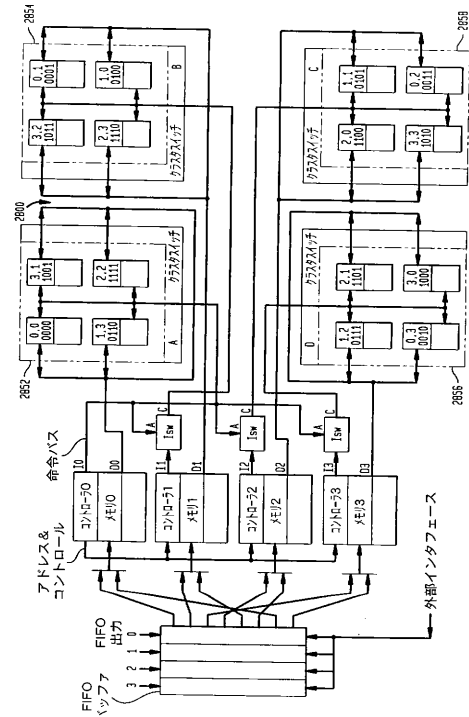
【図 27 A】



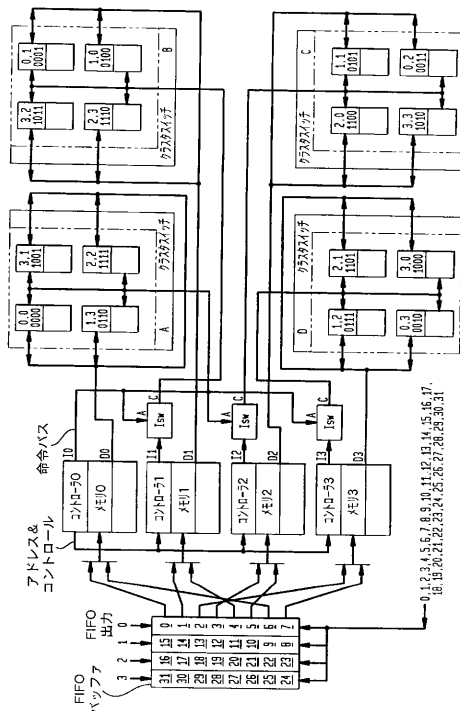
【図 27B】



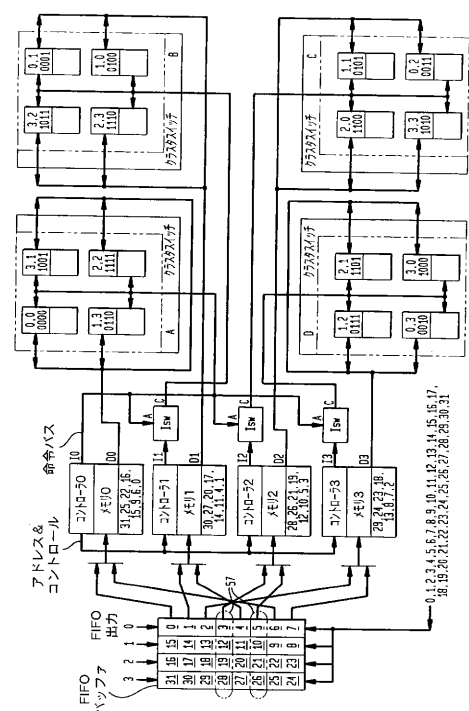
【図 28A】



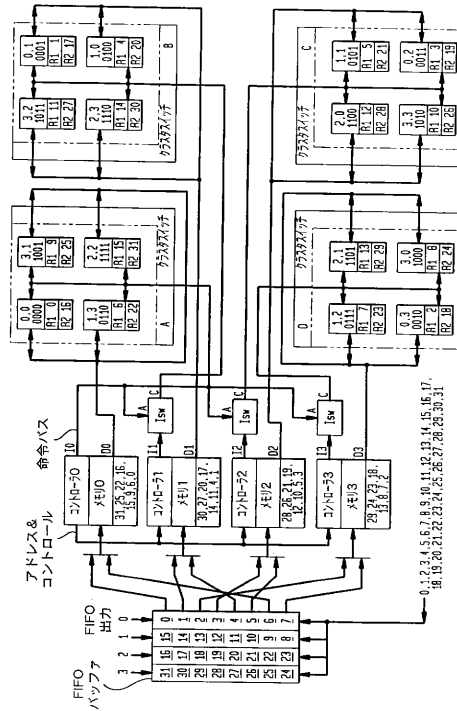
【図 28B】



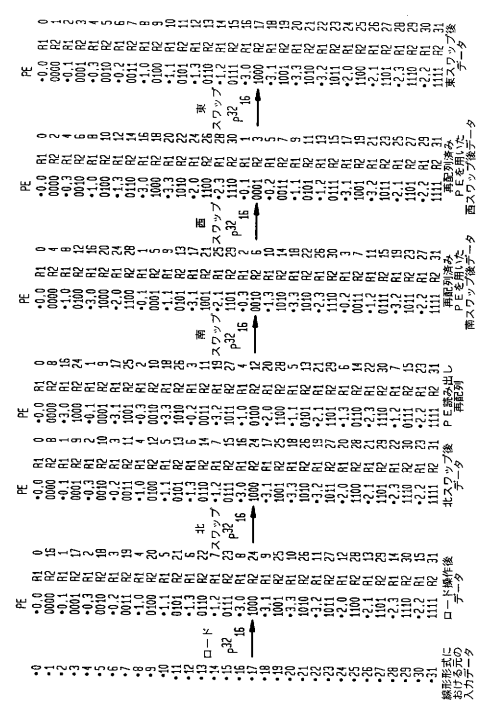
【図 29】



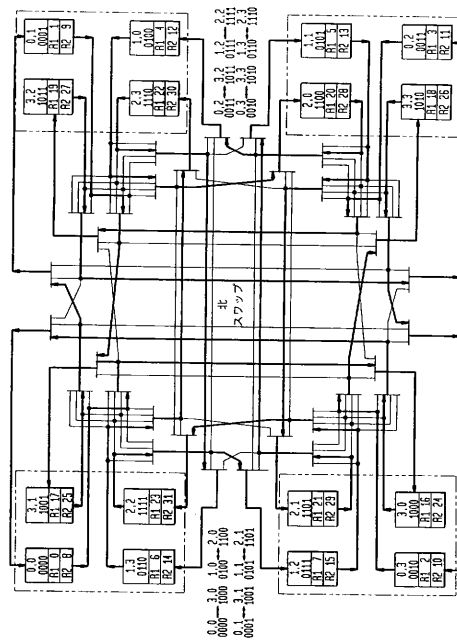
【図 30】



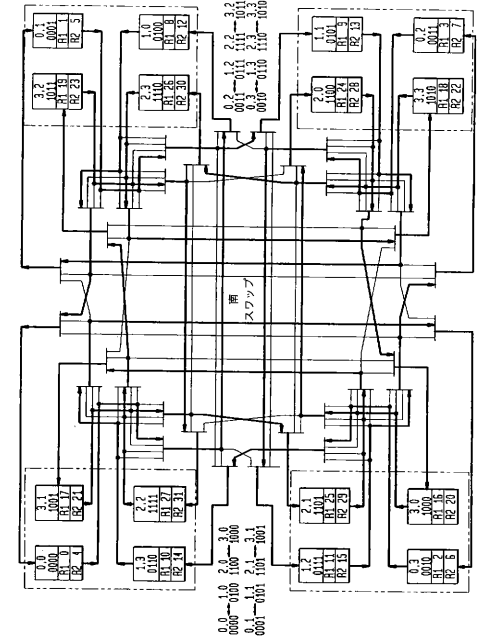
【図 31】



【図 32】



【図 33】



---

フロントページの続き

- (72)発明者 ペカネック ジェラルド ジー  
アメリカ合衆国 ノースカロライナ州 27511 キャリー, ストーンライヒ ドライブ 1  
07
- (72)発明者 ピットィアニス ニコス ピー  
アメリカ合衆国 ノースカロライナ州 27514 チャペル ヒル, ファリントン ロード  
6205 アpartment エイ11
- (72)発明者 バリー エドウィン エフ  
アメリカ合衆国 ノースカロライナ州 27511 キャリイ, ラークホール コート 120  
8
- (72)発明者 ドラベNSTOTT トーマス エル  
アメリカ合衆国 ノースカロライナ州 27514 チャペル ヒル, ファリントン ロード  
6123 アpartment エム9

審査官 久保 正典

- (56)参考文献 特開平07-152722(JP,A)  
特開昭63-059651(JP,A)  
特開平03-186963(JP,A)  
特開平04-138553(JP,A)  
特表2002-507300(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F15/80  
G06F15/16-15/177