

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3599197号

(P3599197)

(45) 発行日 平成16年12月8日(2004.12.8)

(24) 登録日 平成16年9月24日(2004.9.24)

(51) Int. Cl.<sup>7</sup>

G06F 15/173

F I

G06F 15/173

請求項の数 20 (全 27 頁)

(21) 出願番号	特願平6-524304	(73) 特許権者	クレイ・リサーチ・インコーポレーテッド アメリカ合衆国 55121 ミネソタ州 イーガン、ローン・オーク・ドライブ 655エイ番
(86) (22) 出願日	平成6年4月12日(1994.4.12)	(74) 代理人	弁理士 青山 稜
(65) 公表番号	特表平8-511886	(74) 代理人	弁理士 河宮 治
(43) 公表日	平成8年12月10日(1996.12.10)	(72) 発明者	スティーブンソン、ブリッキー・エイ アメリカ合衆国 54729 ウィスコン シン州 チッペワ・フォールズ、ウェスト ・レイク・ドライブ 6409番
(86) 国際出願番号	PCT/US1994/003994		
(87) 国際公開番号	W01994/025920		
(87) 国際公開日	平成6年11月10日(1994.11.10)		
審査請求日	平成13年4月12日(2001.4.12)		
(31) 優先権主張番号	08/055,814		
(32) 優先日	平成5年4月30日(1993.4.30)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 待ち時間が可変の、プロセッサをメモリに接続する相互接続ネットワーク

## (57) 【特許請求の範囲】

## 【請求項1】

マルチプロセッサ・システムにおける複数の中央処理装置(CPU)(100)と複数のメモリ・モジュール(600)との間でパケットの経路を設定するための相互接続ネットワーク(300)であって、

上記相互接続ネットワーク(300)は、

(a1) 前記複数のCPU(100)を前記複数のメモリ・モジュール(600)に接続するための出力経路ネットワーク(302)と、

(a2) 前記複数のメモリ・モジュール(600)を前記複数のCPU(100)に接続するための入力経路ネットワーク(304)とを備え、

前記出力経路ネットワークは、

(b1) 第1のクロスバー・スイッチ(320)と、

(b2) 前記複数のCPU(100)への接続のための入力バッファ(310)と、

(b3) 前記複数のメモリ・モジュール(600)への接続のための第2のクロスバー・スイッチ(330)とを備え、

前記入力バッファ(310)は前記第1のクロスバー・スイッチ(320)に接続され、パケットを前記CPU(100)から前記第1のクロスバー・スイッチ(320)へ転送し、

前記第2のクロスバー・スイッチ(330)は前記第1のクロスバー・スイッチ(320)に接続され、パケットを前記第1のクロスバー・スイッチ(320)から前記複数のメモリ・モジュール(600)へ転送し、

10

20

前記第 1 及び第 2 のクロスバー・スイッチ (320,330) の少なくとも一方は前記複数の CPU (100) から前記スイッチを介して転送されるパケットをキューイングするための先入れ先出し (FIFO) キュー構造体 (321,331) を含み、

前記入力経路ネットワークは、

(c1) 第 3 のクロスバー・スイッチ (360) と、

(c2) 前記複数のメモリ・モジュール (600) への接続のための入力バッファ (350) と、

(c3) 前記複数の CPU (100) への接続のための第 4 のクロスバー・スイッチ (370) とを

備え、

前記入力バッファ (350) は前記第 3 のクロスバー・スイッチ (360) に接続され、パケットを前記メモリ・モジュール (600) から前記第 3 のクロスバー・スイッチ (360) へ転送し、

10

前記第 4 のクロスバー・スイッチは前記第 3 のクロスバー・スイッチ (360) に接続され、パケットを前記第 3 のクロスバー・スイッチ (360) から前記複数の CPU (100) へ転送し、

前記第 3 及び第 4 のクロスバー・スイッチ (360,370) の少なくとも一方は前記複数のメモリ・モジュール (600) から前記スイッチを介して転送されるパケットをキューイングするための FIFO キュー構造体を含むことを特徴とする相互接続ネットワーク。

#### 【請求項 2】

前記出力経路ネットワーク (302) の前記入力バッファ (310) は、前記第 1 のクロスバー・スイッチ (320) へ転送されるパケットをキューイングするための FIFO キュー構造体 (331) を備えたことを特徴とする請求項 1 記載の相互接続ネットワーク。

20

#### 【請求項 3】

前記第 1 及び第 2 のクロスバー・スイッチ (320,330) の一方は、パケットの送信中にパケットを一時的に格納するための出力要求バッファをさらに含むことを特徴とする請求項 1 又は 2 記載の相互接続ネットワーク。

#### 【請求項 4】

前記第 1 及び第 2 のクロスバー・スイッチ (320,330) は、前記パケットの複数のメッセージ・ステアリング・ビットを、前記クロスバー・スイッチへの入力を表すビットで置換するための手段をさらに含むことを特徴とする請求項 1、2 又は 3 記載の相互接続ネットワーク。

30

#### 【請求項 5】

前記入力経路ネットワーク (304) の前記入力バッファ (350) は、前記第 3 のクロスバー・スイッチ (360) へ転送されるパケットをキューイングするための FIFO キュー構造体 (351) を備えたことを特徴とする請求項 1、2、3 又は 4 記載の相互接続ネットワーク。

#### 【請求項 6】

前記第 3 及び第 4 のクロスバー・スイッチ (360,370) の一方は、パケットの送信中にパケットを一時的に格納するための出力要求バッファをさらに含むことを特徴とする請求項 1 乃至 5 のうちのいずれか 1 つに記載の相互接続ネットワーク。

#### 【請求項 7】

請求項 1 記載の相互接続ネットワーク (300) によって、複数のメモリ・モジュールに接続された複数の CPU (100) を備えたマルチプロセッサ・システムにおいて使用するためのメモリ・モジュール (600) であって、

40

複数の入力と、

複数の出力と、

複数のメモリ・バンク (640a - 640p) と、

前記複数の入力を前記複数のメモリ・バンクに接続する第 1 のルーティング・ネットワークと、

前記複数のメモリ・バンクを前記複数の出力に接続する第 2 のルーティング・ネットワークとを備え、

前記第 1 のルーティング・ネットワークは、前記メモリ・バンクの 1 つへ転送されるパケ

50

ットをキューイングするためのFIFOキュー構造体(331)を有する第1のスイッチを含む複数のスイッチ(610)を備え、

前記第2のルーティング・ネットワークは、前記メモリ・バンクの1つから転送されるパケットをキューイングするためのFIFOキュー構造体(331)を有する第2のスイッチを含む複数のスイッチ(610)を備えたことを特徴とするメモリ・モジュール。

【請求項8】

前記第1のルーティング・ネットワークにおける前記FIFOキュー構造体(331)はN個の入力を備え、

前記数字Nは指定されたスイッチ間の通信の往復時間をカバーするように選択されたことを特徴とする請求項7記載のメモリ・モジュール(600)。

10

【請求項9】

前記各メモリ・バンクは複数のメモリ・デバイスを備え、

前記各メモリ・デバイスはアクセス時間とサイクル時間とを有し、

前記各メモリ・バンクは、複数のメモリ・デバイスの前記アクセス時間とサイクル時間とに依存して前記メモリ・バンクの速度を調整するための遅延手段をさらに備えたことを特徴とする請求項7又は8記載のメモリ・モジュール(600)。

【請求項10】

前記第2のルーティング・ネットワークにおける前記FIFOキュー構造体(331)はパイプライン方式のメモリ部分を扱うために十分な深さを有することを特徴とする請求項7、8又は9記載のメモリ・モジュール(600)。

20

【請求項11】

クロスバー・スイッチ(320)を備えたネットワーク・モジュールとの組み合わせ装置において、

前記クロスバー・スイッチ(320)は、前記第1のルーティング・ネットワークの前記第1のスイッチに接続され、前記第1のスイッチから受信されたパケットをバッファリングするためのFIFOキュー構造体(331)を含むことを特徴とする請求項7乃至10のうちのいずれか1つに記載のメモリ・モジュール(600)。

【請求項12】

前記クロスバー・スイッチ(320)は、パケットの送信中にパケットを一時的に格納するために使用される出力要求バッファをさらに含むことを特徴とする請求項11記載の組み合わせ装置。

30

【請求項13】

前記パケットはメッセージ・ステアリング・ビットを備え、

前記クロスバー・スイッチは、複数のメッセージ・ステアリング・ビットを、前記複数の入力の1つを表すビットで置換するための手段をさらに含むことを特徴とする請求項11又は12記載の組み合わせ装置。

【請求項14】

クロスバー・スイッチ(320)を備えたネットワーク・モジュールとの組み合わせ装置において、

前記クロスバー・スイッチ(320)は、前記第2のルーティング・ネットワークの前記第2のスイッチに接続され、前記第2のスイッチから受信されたパケットをバッファリングするためのFIFOキュー構造体(331)を含むことを特徴とする請求項7乃至10のうちのいずれか1つに記載のメモリ・モジュール(600)。

40

【請求項15】

前記クロスバー・スイッチ(320)は、パケットの送信中にパケットを一時的に格納するための出力要求バッファをさらに含むことを特徴とする請求項14記載の組み合わせ装置。

【請求項16】

前記クロスバー・スイッチ(320)は、前記パケットの複数のメッセージ・ステアリング・ビットを、前記複数のメモリ・バンクの1つを表すビットで置換するための手段をさらに含むことを特徴とする請求項14又は15記載の組み合わせ装置。

50

## 【請求項 17】

第1のクロスバー・スイッチと、第2のクロスバー・スイッチとを備えたネットワーク・モジュールとの組み合わせ装置において、

上記第1のクロスバー・スイッチは、前記第1のルーティング・ネットワークの前記第1のスイッチに接続され、前記第1のスイッチから受信されたパケットをバッファリングするためのFIFOキュー構造体を含み、

上記第2のクロスバー・スイッチは、前記第2のルーティング・ネットワークの前記第2のスイッチに接続され、前記第2のスイッチから受信されたパケットをバッファリングするためのFIFOキュー構造体を含むことを特徴とする請求項7乃至10のうちのいずれか1つに記載のメモリ・モジュール。

10

## 【請求項 18】

請求項1乃至6のうちのいずれか1つに記載の相互接続ネットワーク(300)と、

前記相互接続ネットワークに接続された複数のCPU(100)と、

前記相互接続ネットワークに接続された複数のメモリ・モジュール(600)とを備えたことを特徴とするマルチプロセッサ計算システム。

## 【請求項 19】

出力経路ネットワーク(302)と入力経路ネットワーク(304)とを有する相互接続ネットワークによって複数のCPU(100)に接続された複数のメモリ・モジュールのうちの1つのメモリ・モジュール(600)からの読み出し方法であって、

上記出力経路ネットワーク(302)は前記メモリ・モジュール(600)に接続され、入力と、出力と、前記入力を前記出力に接続する複数の第1のステージとを備え、前記各第1のステージはスイッチとバッファとを備え、

20

上記入力経路ネットワーク(304)は前記メモリ・モジュール(600)に接続され、入力と、出力と、前記入力を前記出力に接続する複数の第2のステージとを備え、前記各第2のステージはスイッチとバッファとを備え、

上記方法は、

メッセージ・ステアリング・ビットを有する読み出しコマンド・パケットを前記出力経路ネットワークの前記入力へ転送するステップと、

前記読み出しコマンド・パケットを前記出力経路ネットワークの第1のステージにおいてバッファに格納するステップと、

30

前記出力経路ネットワークの次のステージへの経路が空いているとき、前記読み出しコマンド・パケットを前記出力経路ネットワークの前記次のステージへ転送し、前記メッセージ・ステアリング・ビットを修正して前記読み出しコマンド・パケットがどこから来たかを示すステップと、

前記出力経路ネットワークのすべてのステージが通過されたとき、前記読み出しコマンド・パケットを前記メモリ・モジュールへ転送し、前記メッセージ・ステアリング・ビットを修正して前記読み出しコマンド・パケットがどこから来たかを示すステップと、

前記読み出しコマンド・パケットを前記メモリ・モジュールにおけるバッファに格納するステップと、

前記読み出しコマンド・パケットによって指示されたメモリ・ロケーションから読み出しを行ってデータを獲得するステップと、

40

メッセージ・ステアリング・ビットを含むデータをデータ・パケットにして転送するステップと、

前記データ・パケットを前記入力経路ネットワークの第1のステージにおいてバッファに格納するステップと、

前記入力経路ネットワークの次のステージへの経路が空いているとき、前記データ・パケットを前記入力経路ネットワークの前記次のステージへ転送し、前記メッセージ・ステアリング・ビットを修正して前記データ・パケットがどこから来たかを示すステップとを含むことを特徴とする読み出し方法。

## 【請求項 20】

50

出力経路ネットワーク(302)と入力経路ネットワーク(304)とを有する相互接続ネットワークによって複数のCPU(100)に接続された複数のメモリ・モジュールのうちの1つのメモリ・モジュール(600)への書き込み方法であって、

上記出力経路ネットワーク(302)は前記メモリ・モジュール(600)に接続され、入力と、出力と、前記入力を前記出力に接続する複数の第1のステージとを備え、前記各第1のステージはスイッチとバッファとを備え、

上記入力経路ネットワーク(304)は前記メモリ・モジュール(600)に接続され、入力と、出力と、前記入力を前記出力に接続する複数の第2のステージとを備え、前記各第2のステージはスイッチとバッファとを備え、

上記方法は、

書き込みコマンド・パケットを前記出力経路ネットワークの前記入力へ転送するステップを含み、前記書き込みコマンド・パケットはアドレスとメッセージ・ステアリング・ビットと書き込みデータとを含み、

上記方法は、

前記書き込みコマンド・パケットを前記出力経路ネットワークの第1のステージにおいてバッファに格納するステップと、

前記出力経路ネットワークの次のステージへの経路が空いているとき、前記書き込みコマンド・パケットを前記出力経路ネットワークの前記次のステージへ転送し、前記メッセージ・ステアリング・ビットを修正して前記書き込みコマンド・パケットがどこから来たかを示すステップと、

前記出力経路ネットワークのすべてのステージが通過されたとき、前記書き込みコマンド・パケットを前記メモリ・モジュールへ転送し、前記メッセージ・ステアリング・ビットを修正して前記書き込みコマンド・パケットがどこから来たかを示すステップと、

前記書き込みコマンド・パケットを前記メモリ・モジュールにおけるバッファに格納するステップと、

前記書き込みデータを前記アドレスによって指示されたロケーションに格納するステップとを含むことを特徴とする書き込み方法。

#### 【発明の詳細な説明】

##### <発明の背景>

多くのデータ処理のタスクは、データの整然とした配列に対する大規模な算術操作を必要とする。一般に、この種の操作すなわち「ベクトル」処理は、データ集合の連続する各要素に対して同一操作を繰り返し実行することが必要となる。整然とした配列データを扱うときに処理速度およびハードウェア効率を向上させるために、ベクトル計算機が開発されている。ベクトル計算機は、そのハードウェア機構により整然とした配列データを扱うものであり、したがって、スカラ計算機よりも高い演算速度(speed of operation)が得られる。

スカラ計算機およびベクトル計算機の双方における処理速度(computer processing speed)および処理効率は、マルチ処理技術(multiprocessing techniques)を利用して更に向上させることができる。マルチ処理は、主記憶のようなシステム資源を共有する数百ないし数千のプロセッサを使用する必要がある。異なるジョブの独立したタスクまたは単一ジョブの関連するタスクは、複数のプロセッサ上で実行される。各プロセッサはそれ自身の命令セット従い、それぞれの命令を並列に実行する。プロセッサの数を増やしてそれらを並列に動作させることにより、より多くの作業をより短い時間で行うことができる。

マルチ処理は作業速度(performance speed)を向上させることができるが、その向上は使用されるプロセッサの数と線形の関係にはならない。これは、主として二つの要因すなわちオーバヘッドとロックアウト(lock out)によるものである。プロセッサとプロセッサの機能との間で調整をとるのに要する同期と制御のレベルが増大するため、マルチプロセッサの環境の下では著しいオーバヘッドが生じる。全プロセッサ間での通信および全プロセッサの制御は、マルチ処理システムに性能低下を生じさせる。数個のプロセッサが協調して一つのタスクを実行しているときには、データの依存性とプロセッサ間でのデータ

10

20

30

40

50

の転送は不可避である。或るプロセッサが他のプロセッサから転送されるべきデータを待たなければならないときに、プロセッサのアイドル時間が生じる。このプロセッサのアイドル時間はシステム性能を低下させることになる。

マルチプロセッサ・システムの性能低下の他の重大な原因は、プロセッサのロックアウトすなわちブロッキング (blocking) であり、複数のプロセッサが共通の資源を共有することに関連している。これは、一つのプロセッサが、他のプロセッサが既に使用している共有資源、例えば共有メモリにアクセスをしようとするときに生じる。そのプロセッサは、このときその共有資源の使用を阻止され (ブロックされ)、他のプロセッサが停止するまで待機しなければならない。この場合もまた、プロセッサのアイドル時間が発生し、システムの性能が低下する。

プロセッサのメモリに対するインタフェイスは、オーバヘッドとロックアウトの概念と緊密につながっており、それは計算機の全体的な性能にも影響を与える。マルチプロセッサのインタフェイスの一例を、BBNシステム・アンド・テクノロジー社 (BBN Systems and Technologies Corporation) によって設計されたモナーク (Monarch) 並列マルチ処理計算機に見ることができる。モナークはスカラ方式で単一スレッドのマルチ処理アーキテクチャとなっており、プロセッサとメモリとの間の通信に回線交換技術を使用している。この回線交換技術によれば、全てのプロセッサはメモリへ至る同一の経路を共有する。モナークの設計における或るプロセッサがメモリを要求するとき、そのプロセッサ網からメモリへの経路全体が開放され、そのメモリとそのプロセッサとの間の通信が完了するまで開放されたままとなる。この方式は、その回線交換網を経由してメモリを参照しようとしている他のプロセッサをブロックすることができ、メモリ参照の転送速度を制限し、その結果、プロセッサのアイドル時間が長くなる。したがって、このような設計をマルチプロセッサによるマルチスレッドのベクトル処理において採用するのは実用的ではなく、このようなベクトル処理ではプロセッサとメモリとの間で本来的に大量のデータを転送しなければならない。

マルチプロセッサのメモリインタフェイスの他の例を、ホライゾン (HORIZON) ルーティング方式に見ることができる。このホライゾンのインタフェイス・ネットワークは、自暴自棄ルーティング (desperation routing) またはホットポット・ルーティング (hot potato routing) と呼ばれる方式を使用している。ホライゾンの自暴自棄ルーティングは、複数の入力と同数の出力を有する多段ネットワークである。この方式は、どのネットワークのサイクルにおいても各入力に対し一つの出力への経路を設定する必要がある。例えば、4個の入力情報 (input references) が存在し、その4個の入力情報が同一の出力に行くことを欲している場合、その4個の入力情報のうち一つが正しい出力へ行き、他の全ての入力は欲していない他のいずれかの出力へ行く。これは、4個の入力のうち3個は格段に長い経路をとってそのネットワークを通過することを意味する。このホライゾンの自暴自棄ネットワークは、これらの他の三つの情報 (references) は結局は望ましい入力へ戻って、望ましい出力に到達する別の機会を持つように、ネットワークに経路が定められる。したがって、それらの情報がネットワークの中で永久に失われるということではなく、ホライゾンのルーティング方式は、そのネットワーク内に最も長く存在している情報が最高の優先度を持つような機構を有しており、これにより、それらの情報はいつかは同一の出力に対する競合する情報に打ち勝つことになる。このようなルーティング方式によれば、一つの情報が要求する終点への可能な経路を複数持ち、多くの情報は、その目的地へ到達する前にネットワーク内を通行するのに非常に長い期間を費やすことになりかねない、ということが当業者には容易にわかるであろう。したがって、このホライゾンの自暴自棄ルーティング方式をマルチ処理の計算機で使用するの望ましくない。

マルチプロセッサのメモリインタフェイスの更に他の例が、ムラタ (Murata) らによって1992年4月1日に出願されたヨーロッパ特許出願第0 510 821号において説明されている。この出願は、共有メモリ・モジュールと複数のプロセッサとの間でデータを転送するためにアドレスバスに接続されたスイッチを使用することを教示している。このアドレスバスは、特定のプロセッサのキャッシュ内に存在するデータのアドレスに対するメモリ参照

10

20

30

40

50

について調べる (snoop) のために使用することができ、一方、このスイッチは、メモリ・モジュールとプロセッサとの間での高速のデータ転送路を提供する。

マルチプロセッサのメモリインタフェイスの他の例が、ハイデルバーガ (Heidelberger) らによって1990年3月2日に出願されたヨーロッパ特許出願第0 392 184号において説明されている。この出願は、メモリとプロセッサ双方のキャッシュの使用について述べている。このメモリ・キャッシュは、共有メモリに対するアクセス時間を短縮するために共有メモリデータをキャッシュするのに使用される。

マルチ処理システムにおける他の重要な概念は、異なる使用者の要求に合わせて種々の異なるサイズにすることである。例えば、最大システムは1024個のプロセッサを持つことができるが、512個のプロセッサを有する構成、256個のプロセッサを有する構成、または他の或る構成というように縮小化されたシステムを利用できるようにするのが望ましい。最大システムを構成する基本的なビルディングブロック (building blocks) を修正することなく使用して最小システムを作り上げることができることが重要であり、逆もまた同様である。したがってスケラブルなシステムは、格段に柔軟性に富み、このようなシステムの展開によって使用者の変化する要求に合わせるができる。

以上より、プロセッサをメモリに接続する相互接続ネットワークであって、例えばプロセッサが競合することなくメモリ参照を発行できるようにした相互接続ネットワークに対する要求が当業界にあり、それは、ネットワーク内でのメモリ参照の間での競合を低減し、どの1回の参照についてもネットワーク内で費やす時間を低減し、その結果、プロセッサのアイドル時間が減少し、システムの性能が向上する。また、モジュラーな相互接続ネットワークであって、その相互接続ネットワークを構成する個々のモジュールの再設計を必要とせずに、任意の個数のプロセッサおよび異なるサイズのメモリを有するマルチ処理システムに適合するように容易に規模を変更することができる相互接続ネットワークに対する要求もある。

#### < 発明の要約 >

上記技術における限界を克服するために、および本明細書を読んで理解すれば明らかになるであろう他の限界を克服するために、本発明は、プロセッサをメモリに接続するマルチ処理用の相互接続ネットワークを提供する。この相互接続ネットワークは、種々の個数のプロセッサおよび種々のメモリサイズを有するように構成されたシステムでの使用に適合したものである。すなわち、この相互接続ネットワークのモジュール性により、単純な複製で任意の規模のシステムを構築できるようになる。

この相互接続ネットワークは、ネットワーク・モジュールおよびメモリ・モジュールを含む。ネットワーク・モジュールおよびメモリ・モジュールは、一連の  $n \times m$  スイッチで構成されていて、各  $n \times m$  スイッチは  $n$  個の入力のうちの一つから  $m$  個の出力のうちの一つへの経路を設定する。このスイッチは、相互接続ネットワークにおけるメッセージの競合が低減されるように設計されている。このスイッチ、したがってメモリとネットワークのモジュールは高いモジュール性を有しており、これにより、同一の構成要素を利用して任意の規模のマルチ処理システムを実質的に構築できるようになる。

また、メッセージの競合を低減するための機構および相互接続ネットワークを実現するために必要な相互接続の数が与えられる。最後に、相互接続ネットワーク内の特定の経路上のデータが正しいか否かを判定するために健全性コード (sanity code) が使用される。この健全性コードは、メモリ・バンクもしくはネットワーク経路の故障を検出し、または、システムのグレードを下げるためにも使用される。

本発明は、請求項1 (ネットワーク)、請求項8 (読み出し方法)、請求項9 (書き込み方法)、請求項10 (計算機システム)、および請求項17 (メモリ・モジュール) において定義されている。

#### 【図面の簡単な説明】

図面において、同一符号はいくつかの図を通して同一構成要素を示している。

図1は、4個のCPUを有するマルチ処理システムの一例のブロック図を示す。

図2は、図1に示した4個のCPUを有するマルチ処理システムのより詳細なブロック図で

10

20

30

40

50

ある。

図3は、32個のCPUを有するマルチ処理システムの一列のブロック図を示す。

図4は、図3に示した32個のCPUのシステムにおける、8個のCPUと4個のネットワーク・モジュールとの間の接続のより詳細なブロック図を示す。

図5は、図3に示した32個のCPUのシステムにおける、ネットワーク・モジュールと8個のメモリ・モジュールとの間の接続のより詳細なブロック図を示す。

図6は、メモリ・モジュールのブロック図を示す。

図7は、図6に示したメモリ・モジュールにおける出力経路 (outgoing path) のより詳細なブロック図を示す。

図8は、図6に示したメモリ・モジュールにおける入力経路 (return path) のより詳細なブロック図を示す。

10

図9は、ネットワーク・モジュールのブロック図を示す。

図10は、図9に示したネットワーク・モジュールにおける出力経路のより詳細なブロック図を示す。

図11は、図9に示したネットワーク・モジュールにおける入力経路のより詳細なブロック図を示す。

図12は、レディ/レジューム (READY/RESUME) のハンドシェイクのための制御回路の詳細構成を示す。

図13は、健全性コードの制御回路の詳細構成を示す。

図14~15は、本発明の相互接続ネットワークを使用することができるマルチ処理システムのいくつかの例についての種々の構成を示す表である。

20

図16は、1×2スイッチの遅延チェーン (delay chain) の詳細を示す。

< 好ましい実施形態の詳細な説明 >

以下の詳細な説明では、本文の一部を構成すると共に、例として本発明の特定の実施形態を示した添付図面を参照することにする。しかし、他の実施形態を利用してよく、本発明の範囲を逸脱しない限り構造的または論理的な変更を施すこともできる。したがって、以下の詳細な説明は限定された意味で解されるべきではなく、本発明の範囲は添付した請求の範囲によって定めるべきである。

図1を参照しつつ、マルチ処理システムの一列の簡略ブロック図を説明する。図1は、4個のメモリ・モジュール600a-dから成る共通メモリに接続された4個のCPU100を示す。

30

図2は、図1のシステムのより詳細なブロック図を示す。4個のCPU100a-dのそれぞれは、ポート0とポート1の二つのポートを持っている。ポート0とポート1のそれぞれは、4個のメモリ・モジュール600a-dのうちの一つに経路が設定され、CPU1個当たり全部で8出力に対し、メモリ・モジュール1個当たり8入力に経路が設定される。

マルチ処理システムの他の例を図3に示す。図3は、32個のCPUを有する最大規模のシステムである。8個のCPUから成る4個のCPU群100a-h、100i-p、100q-x、100y-ffが、4個のネットワーク・モジュールから成る4個のネットワーク・モジュール群300a-d、300e-h、300i-l、300m-pを介して、8個のメモリ・モジュールから成る4個のメモリ・モジュール群600a-h、600i-p、600q-x、600y-ffに接続されている。プロセッサをメモリに接続するこの相互接続ネットワークを使用して構築することができる4個、8個、および32個のCPUのシステム例を以下において更に詳細に説明する。

40

図1~3は、ネットワーク・モジュール300およびメモリ・モジュール600のアーキテクチャのモジュール性により、如何にして、単純な複製で任意の規模のシステムの構築が可能となるのかを示している。1~4個のCPUを有するシステムにおいては、それらのCPUは図1および2に示すようにメモリ・モジュール600に直結される。4個よりも多くのCPUを有するシステムにおいては、それらのCPUは、まずネットワーク・モジュール300に接続され、そのネットワーク・モジュールがメモリ・モジュール600に接続される。最小のマルチ処理システムを構成するものと同じのネットワーク・モジュールおよびメモリ・モジュールを、修正することなく、最大のシステムを構築するために使用することができる。モジュール性と拡張性は本相互接続ネットワークの設計に本来的に備わっているものであるた

50

め、それらはシステム実現のために選択される技術に依存せず、極めて柔軟性に富み、その結果、容易に規模を変更できるマルチ処理システムが得られる、ということが当業者には容易にわかるであろう。

4個、8個、および32個のCPUを有するマルチ処理システムの例を説明する前に、メモリの相互接続ネットワークを構成するメモリ・モジュール600およびネットワーク・モジュール300について詳細に説明する。一般に、メモリ・モジュール600およびネットワーク・モジュール300は、 $n \times m$ スイッチの集合から構成される。これらの $n \times m$ スイッチは、種々の形態に接続されてメモリおよびネットワークのモジュールを形成し、それらは、プロセッサと共通メモリとの間にメッセージの経路を設定する。相互接続ネットワークを構築するために同一の $n \times m$ スイッチ群が繰り返し使用されるように、そのスイッチが作製される。したがって、その結果得られる構造はモジュール性が高く、小規模のマルチ処理システムを構築するために使用されるものと同一の構成要素が大規模なシステムにおける接続にも使用できるようになっている。

<メモリ・モジュール>

次にメモリ・モジュール600について詳細に説明する。図6は、メモリ・モジュール600のブロック図を示す。出力経路(outgoing path)604および入力経路(return path)606の双方が示されている。各メモリ・モジュール600は、8個の入力0~7と16個のメモリ・バンク640a~640pとを備えている。実現の容易化のため、このメモリ・モジュールは、論理的に二つの独立部分に分離されている。入力0~3はメモリ・バンク640a~640h(上半部)へ経路が設定されており、一方、入力4~7はメモリ・バンク640i~640p(下半部)へ経路が設定されている。これら二つの部分は、完全に分離されていて、互いに影響されることはない。物理的に、これら二つの部分は、一つのプリント回路基板の上部側と底部側に配置されるのが好ましい。しかし、本相互接続ネットワークは、この物理的な実現方法に限定されないと理解すべきである。すなわち、多くの他の物理的な所定の実現が当業者によって容易に考えられるであろう。

入力0~7のそれぞれは、まず $1 \times 1$ スイッチ610によってバッファリングされる。4個のスイッチ610から成る2個のスイッチ群は、2個の $4 \times 4$ スイッチ620のうちの一つの入力へ接続される。 $2 \times 2$ スイッチ420は、4個の入力のそれぞれから4個の $1 \times 2$ スイッチ630の一つへ経路を設定する。 $1 \times 2$ スイッチ630のそれぞれは、一つの入力から2個のメモリ・バンク640の一つへ経路を設定する。

各メモリ・モジュール600の入力経路606において、2個のメモリ・バンクは $2 \times 1$ スイッチ650へ経路が設定されている。4個の $2 \times 1$ スイッチ650は、 $4 \times 4$ スイッチ660の入力へ接続され、その後、それらは4個の $1 \times 1$ スイッチ670の一つへ経路が設定され、全部で8個の出力0~7へ経路が設定されている。

図7は、メモリ・モジュール600の出力経路側604のより詳細なブロック図を示す。上半部および下半部の双方が示されている。メモリ・モジュール600の入力0~7のそれぞれは、CPU100からそのメモリ・モジュールへのメモリ参照をバッファリングする $1 \times 1$ スイッチを有している。各 $1 \times 1$ スイッチ610におけるバッファの数は、クロスバー・スイッチ間での通信の往復時間によって変わる。各クロックの周期がその通信時間に付加されるため、各 $1 \times 1$ スイッチ610内のバッファを更に追加する必要がある。 $1 \times 1$ スイッチ610の目的は、メッセージ・パケットを多重に送ることができるように、クロスバー・スイッチ間の通信時間をカバーすることにある。したがって、図6に示した本メモリ・モジュールには、6個のバッファA~Fが設けられていて、 $n \times m$ スイッチがクロックの1周期の期間で通信を行う。

$1 \times 1$ スイッチ610におけるバッファA~F、および本相互接続ネットワークを構築するために使用される全ての $n \times m$ スイッチにおけるバッファは、先入れ先出し(FIFO)のキュー構造体(queue structure)を有している。入力0~7のうちの一つを経由してメモリ・モジュール600に入るメッセージは、その経路用の対応するバッファAに他のメッセージが既に存在するということがなければ、そのバッファAに保持される。バッファAが占有されている場合には、そのメッセージは、バッファBが使用可能であれば、バッファ

10

20

30

40

50

Aの代えてバッファBに保持される。バッファAがいったん空になると、バッファBに存在するメッセージはバッファAへ移される。別のメッセージがアクセスしようとしたときバッファAおよびバッファBの双方が満杯であれば、そのメッセージはバッファCに入る。このようにして入力バッファは、先入れ先出しキューとして動作する。別のメッセージがその入力バッファ610に入ろうとしたときバッファA～Fが全て満杯であれば、そのメッセージは、メモリ・モジュールに入る前に、FIFOバッファの一つが解放されるまで単に待機するだけである。入力保持バッファ(input holding buffer)621が他のメッセージを受け取る用意ができているときは、1×1スイッチ610のバッファA内のメッセージは、4×4スイッチの入力保持バッファ621の中に保持される。

いったんメッセージが入力保持バッファ621に保持されると、バッファ621のバッファAは2個のメッセージ・ステアリング・ビット(message steering bit)を解読してそのメッセージの経路を4個の出力要求バッファ(output request buffers)のいずれに設定すべきかを決定する。2個のメッセージ・ステアリング・ビットを解読した後、バッファAは、最初のケット内の2個のメッセージ・ステアリング・ビットを、8個の入力0～7のうちいずれの入力からそのメッセージがメモリ・モジュール600に入ったかを表す2個のビットに置き換える。このようにして、メッセージがメモリ・バンク640へ行く途中においてメモリ・モジュールの層を順次通過して行く間に、発信元のCPUへ戻すリターン・アドレス(return address)がそのメッセージ・ステアリング・ビットに組み込まれる。データが適切なメモリ・バンクから取り出された後、その取り出されたデータを発信元のプロセッサまで案内するために、プロセッサ・リターン・アドレス情報をそのメッセージ・ステアリング・ビットから得ることができるようになる。同様に、一旦そのメッセージがプロセッサに戻ると、その参照がどのメモリ・バンクから来たのかを示すリターン・アドレスをそのメッセージ・ステアリング・ビットから得ることができるようになる。この方法では、データ中に誤りがあれば、そのプロセッサは誤りのあるメッセージがどのメモリ・バンクから来たのかを知り、これにより、その誤りがどこで生じたのかを正確に指摘するための助けとなる。

このビット置換方式は、プロセッサ・リターン・アドレスをメッセージとともに送る必要がなくなるため、システムにおける相互接続およびメッセージ長を低減する。その代わりに、本発明の好ましい実施形態におけるビット置換方式は、出力経路上でメッセージを導き、自動的にプロセッサまたはメモリ・バンクのリターン・アドレスを生成するために、同じビットを使用して相互接続を行うだけである。

入力保持バッファ621のバッファA内で解読されたメッセージ・ステアリング・ビットは、メッセージが入力保持バッファAを出た後に4個の出力要求バッファ622a-dのうちのいずれに行くのかを決定する。出力要求バッファ622の目的は、相互接続ネットワーク内におけるメッセージ競合およびブロッキングを低減することにある。例えば、群バッファ(group buffer)623aが満杯であれば、出力要求バッファ内で待機しているメッセージはブロックされ(封鎖され)、その経路が空くのを待っているメッセージが存在することになる。しかし、例えば、異なる群バッファ623b、cまたはdに向かう別のメッセージが、入力0からメモリ・モジュール600に入ると、そのメッセージは群バッファ623aに対して待機しているメッセージによってブロックされることなく、適切な出力要求バッファ622の中へ入ることができる。このように、他のいずれかの群バッファ623に向かう異なるメッセージが入力0に入った場合、そのメッセージは、適切な出力要求バッファ622に入っていくだけであるため、群バッファ623へのアクセスをブロックされることがない。このようにして、より多くのメッセージはクロスバー・ネットワークにおけるブロックを「避けて通る」ことができる。相互接続ネットワークを通過する特定の経路が1×1スイッチ610を通る全ての通路でブロックされている場合にのみ、その経路に入ったメッセージは全ての出力経路からブロックされる。もし、4×4スイッチ620内に出力要求バッファ610が存在しなければ、4個の出力は全てブロックされるであろう。出力要求バッファ622を使用することにより、4個の出力の経路のうち一つのみがブロックされる。出力要求バッファは、各n×mスイッチおよびネットワークにおける、メッセージ競合およびブロッキ

10

20

30

40

50

ングの可能性とメモリ・モジュール自体とを大いに低減する、ということが当業者には容易に理解できるであろう。

メッセージが出力要求バッファ622に入った後は、そのメッセージは適切な群バッファ623 a - dに導かれる状態になる。各群バッファ623は、各入力経路に対応する出力要求バッファ622からメッセージを探す。例えば、群バッファ623aは、入力0、1、2、および3から4×4スイッチ620にそれぞれ入った適格なメッセージを、その対応する出力要求バッファ622a、622e、622iおよび622mから探す。適切な出力要求バッファ622内に適格なメッセージが一つだけ存在する場合は、群バッファ623はそのメッセージに対応する1×2スイッチ630へ送るだけである。好ましい実施形態では、適切な出力要求バッファ622内に適格なメッセージが複数存在すれば、群バッファ623は、最後に選ばれた出力要求バッファからラウンドロビンの順序での次のメッセージが送られる。例えば、群バッファ623aが転送のために出力要求バッファ622bを選択した場合において、もし、出力要求バッファ622cが適格なメッセージを有していて、適格なメッセージが複数存在すれば、その群バッファは、転送のために出力要求622cを次に選択するであろう。

一旦或るメッセージが4×4スイッチ620の群バッファ623を通過すると、そのメッセージは、対応する1×2スイッチ630へ進む。好ましい実施形態では、1×2スイッチ630に入るメッセージは、2個のメモリ・バンク640のうちの一つに経路が設定される。

或るメッセージが1×2スイッチ630のバッファAに入った後は、バッファAはメッセージ・ステアリング・ビットを解読して、そのメッセージを2個のバンク入力バッファ(bank input buffer)632aまたは632bのうちの一つに導く。各バンク入力バッファ632は、異なるメモリ・バンク640に対応づけられている。例えば、1×2スイッチ630aに入るメッセージは、メモリ・バンク04へ経路を設定されるべきものであって、バンク入力バッファ632bへ導かれる。この段階において、メッセージは、バンク制御論理チップにより選択されて適切なメモリチップにアクセスできる状態になる。

メモリ・モジュール600から要求元のCPUへ戻るために、参照は、上述のようにプロセッサからメモリへの出力経路側参照がメモリ・モジュールの出力経路側604を通過して進むのと同様にして、そのメモリ・バンクを出て、メモリ・モジュール600の適切な入力経路側606を通過して発信元のプロセッサに戻る。メモリ・モジュール600の入力経路側606の詳細を図8に示す。この場合もまた、上半部と下半部の双方が示されている。入力経路参照(return reference)は、メモリ・バンク616を出た後、2×1スイッチ650に入る。各2×1スイッチ650は8個のバッファFIFO652を備えている。このスイッチにおけるFIFOは、パイプライン方式のメモリ部分の可能性に対応できるように8段の深さがなければならない(これについては下記において詳細に論じられる)。各2×1スイッチ650は、2個の対応するメモリ・バンクからメッセージを受け取り、そのメッセージを8個のバッファFIFOによってバッファリングし、そのメッセージの経路を4×4スイッチ660の入力へと設定する。4×4スイッチ660は、図5を参照しつつ上記において説明された4×4スイッチ620と同じように動作する。リターン・メッセージ(return message)は、4×4スイッチ660から適切な1×1出力バッファ670へと経路が設定される。1×1スイッチ670は、図5を参照しつつ上記において説明された1×1スイッチ610と同じように動作する。このリターン・メッセージは、対応する出力0～7からメモリ・モジュール600を出て行く。

<ネットワーク・モジュール>

図9は、ネットワーク・モジュール300のブロック図を示す。各ネットワーク・モジュールは、16個の入力0～15と16個の出力0～15とを有している。実現の容易化のため、各ネットワーク・モジュールは二つの部分に論理的に分離されている。入力0～7はメモリ・モジュール600a-h(上半部)へと経路が設定され、入力8～15はメモリ・モジュール600i-p(下半部)へと経路が設定されている。メモリ・モジュールに関し上記で説明したように、これら二つの部分は完全に分離されていて、互いに影響されることはない。物理的に、これら二つの部分は、プリント回路基板の上部側と底部側に実装されるのが好ましい。しかし、本相互接続ネットワークは、ここに述べられた物理的な実現方法に限定されるものではなく、多くの他の物理的な所定の実現が当業者にとっては明白であろう。

10

20

30

40

50

或る特定のネットワーク・モジュール（図3参照）に接続された8個のプロセッサのそれぞれは、メモリに対して二つのポートを持っていて、それらは、16個の1×1スイッチ310a-pを介してネットワーク・モジュール300の入力へ接続されている。例えば、CPU0に対する二つのポートは、1×1スイッチ310aおよび1×1スイッチ310iを介して接続されている。そして各1×1スイッチ310は4個の4×4スイッチ320a-dのうちの一つに接続されている。そして各4×4スイッチは、その4個の入力のそれぞれから4個の出力の一つへと経路を設定し、出力側は、8個の2×2スイッチ330a-hで構成されている。そして8個の2×2スイッチ330a-hのそれぞれは、16個のメモリ・モジュール600a-pの一つへの入力を有している。

メモリ・モジュール300の入力経路304は、出力経路302上におけるものと同様のこれらの10  
スイッチから構成されている。16個の1×1スイッチ350a~350pは、16個のメモリ・モジュール600a-pの一つに接続されている。1×1スイッチ350は、1×1スイッチ310aと同じように動作し、それらと同様の構造となっており、上記において図6~8を参照しつつ1×1スイッチ610および670について論じられたのと同様である。4×4スイッチ360a-dは、それぞれ、4個の1×1スイッチ350から成るスイッチ群からの入力を受け取り、その4個の入力のそれぞれに対し、4個の2×2スイッチ370の一つへ経路を設定する。4×4スイッチ360は、ネットワーク・モジュール300の入力側302における4×4スイッチ320と同じように動作する。同様に、2×2スイッチ370は、ネットワーク・モジュール300の入力側302における2×2スイッチ330と同じように動作する。一つの2×2スイッチからの各出力は、8個のプロセッサ100a-hのうち2個のプロセッサに接続された20  
出力のうちの一つである。このようにして、図9に示すように、CPU100aのポート0は2×2スイッチ370aに接続され、CPU100aのポート1は2×2スイッチ370eに接続される。図10は、ネットワーク・モジュール300の出力経路302のより詳細なブロック図を示す。図10は出力経路の上半部のみを示している。すなわち、入力310a~310hは示されているが310i~310pは示されていない。しかし、図10に示されていないネットワーク・モジュールの下半部は、図10および11を参照しつつ説明される上半部と同じように動作する。図10は、1×1スイッチ310、4×4スイッチ320、および2×2スイッチ330のそれぞれの内部構造を示している。図10に示された詳細より、1×1スイッチ310は、図6~8を参照しつつ上記において説明された1×1スイッチ610aおよび1×1スイッチ670と同じように動作し、それらと同様の構造となっている。また、4×4スイッチ320は、図6~8を参照30  
しつつ説明された4×4スイッチ660および620と同じように動作し、それらと同様の構造となっている。

2個の4×4スイッチ320aおよび320bからの4個の出力のそれぞれは、4個の2×2スイッチ330a-dのうちの一つへ経路が設定されている。2×2スイッチ330のそれぞれは、4個のバッファを有する2個のFIFO331aおよび331bを有しており、これらのFIFOのそれぞれは、メッセージの経路を2個のモジュール・バッファ332aおよび332bへ設定している。FIFO331におけるバッファAは、メッセージ・ステアリング・ビットを解読して、そのメッセージに対し適切なモジュール・バッファ332へ経路を設定する。各モジュール・バッファは、メモリ・モジュールの一つに対する入力となっている。

図11は、ネットワーク・モジュールを通過する入力経路304を示す。図10と同様に、図1140  
は、ネットワーク・モジュール300における入力経路の上半部のみを示している。図11に示された詳細より、1×1スイッチ350は、図10に示された1×1スイッチ310、図7に示された1×1スイッチ610、および図8に示された1×1スイッチ670と同じように動作し、それらと同様の構造となっている。同様に、4×4スイッチ360は、図10に示された4×4スイッチ320、図8に示された4×4スイッチ660、および図7に示された4×4スイッチ620と同じように動作し、それらと同様の構造となっている。最後に、2×2スイッチ370は、図10を参照しつつ上記において説明された2×2スイッチ330と同じように動作し、それと同様の構造となっている。

<マルチ処理システムの一例についての詳細な説明>

上記の説明より当業者にとっては、メモリ・モジュール600およびネットワーク・モジュ 50

ール300の設計におけるモジュール性は明白であろう。最小のマルチ処理システムを構成するものと同じの  $n \times m$  スイッチを、修正することなく、最大のシステムを構築するために使用することができる。したがって、モジュール性と拡張性は本ネットワーク・モジュールおよびメモリ・モジュールの設計に本来的に備わっているものであるため、それらはシステム実現のために選択される技術に依存せず、また、極めて柔軟性に富み、容易に規模を変更できるマルチ処理システムを得ることができる、ということが当業者には容易に理解できるであろう。

本発明の相互接続ネットワークを用いて構築することができるマルチ処理システムの三つの例を詳細に説明する。

再び図1を参照しつつ、4個のCPUを有するマルチ処理システムを説明する。図2は、各CPU100からの8個の出力のそれぞれに対し共通メモリ602への経路を設定する方法の詳細を示している。図2に示された接続により、各CPUは共通メモリ602内の全てのメモリ・バンクにアクセスできるようなる。図2は、例えば、CPU100aのポート0および1がメモリ・モジュール600aの入力0および4にそれぞれ接続されていることを示している。図6は、入力0がメモリ・バンク640a-hへのアクセス経路を有し、入力4がメモリ・バンク640i-pへのアクセス経路を有することを示している。このように、各メモリ・バンクは、各CPU100のポート0またはポート1のいずれかに対応づけられる。CPU100aのポート0および1は、同様にしてメモリ・モジュール600b-dにも接続されている。CPU100b-dも同様にして各メモリ・モジュール600に接続され、これにより、各CPUが共通メモリ602内のメモリ・バンク640の全てにアクセス経路を有するようになっている。

次に図3を参照しつつ、CPUが32個のマルチ処理システムを説明する。32個のCPU100a-ffとネットワーク・モジュール300a-pとの間の接続の細部を図4により詳細に示す。図4は、8個のCPUと、ネットワーク・モジュール300a、300e、300iおよび300mの間の、それらのCPUに関連する相互接続とを示している。CPU100a-ffは、各CPUがメモリ・モジュール600の各メモリ・バンクへのアクセス経路を有するように、ネットワーク・モジュール300に接続されている。このようにして、CPU100aのポート0は、ネットワーク・モジュール300aの入力0、ネットワーク・モジュール300eの入力0、ネットワーク・モジュール300iの入力0、およびネットワーク・モジュール300mの入力0に接続される。図9に示すように、ネットワーク・モジュール300において、CPU100aのポート0は入力0経由でメモリ・モジュール600a-hへアクセスする経路を有し、CPU100aのポート1は入力8経由でメモリ・モジュール600i-pへアクセスする経路を有している。図3に示された32個のCPUのシステムにおける残りのプロセッサも同様にして接続され、各CPUが共通メモリ602内の全てのメモリ・バンクへの経路を有するようになっている。

図5は、図3に示された32個のCPUのシステムの一例における、ネットワーク・モジュール300a-dとメモリ・モジュール600a-hとの間の接続をより詳細に示している。各メモリ・モジュール600a-hは、4個のネットワーク・モジュール300a-dのそれぞれにおける2個の出力に接続されている。このネットワーク・モジュールとメモリ・モジュールは、32個のプロセッサのいずれからの参照も共通メモリ602内のいずれのメモリ・バンクにも到達できるように接続されている。

例えば、CPU100dからの参照に対しメモリ・モジュール600hのメモリ・バンク640eへ経路を設定するために(図3参照)、CPU100dのポート0はネットワーク・モジュール300aの入力3へ経路が設定される(図4参照)。ネットワーク・モジュール300aからは、メッセージがネットワーク・モジュール300aの出力7より出される(図5に示されている)。この出力はメモリ・モジュール600hの入力0に接続されている。次に図6を参照すると、各メモリ・モジュール600の入力0は、それぞれのバンク640eへの経路として、 $4 \times 4$  スイッチ620aおよび $1 \times 2$  スイッチ630cを経てバンク640eへ到達するという経路が設定されている。再び図3を見ると、CPU100qからの参照に対し、ネットワーク・モジュール300gを経由して、または、いずれかのCPUとメモリ・バンクとの間で上述の経路と同様の相互接続ネットワーク内の経路に沿って、メモリ・モジュール600iへの経路を設定することができる。

4個よりも多いが32個よりも少ないプロセッサを有するシステムは、図3の32個のCPUのシステムを変形したものをを用いて構築することができる。例えば、8個のCPUのシステムを、ネットワーク・モジュール300aを介してメモリ・モジュール600a-hに接続されたCPU100a-hで構築することができる。図4および5に示すように、8個のCPUのシステム内の各CPUの各ポートに対し8個のメモリ・モジュール600a-hへの経路を提供するためには、各CPU100a-hの出力0および1のみが必要である。したがって、各CPUの出力2~7は無効なものとなる。

8個のCPUのマルチ処理システム内のネットワーク・モジュール300aは、図5に示すように8個のメモリ・モジュール600a-hのそれぞれに接続されている。ただ一つのネットワーク・モジュールのみが必要であるため、各メモリ・モジュールの入力0および1のみが使用される。残りの入力接続されず、したがって残りの経路におけるデータの無効である。

8個のCPUと8個のメモリ・モジュールを有するマルチ処理システムのようなシステムでは、CPUの全ての出力が使用されているわけではなく、そのようなシステムに対しては、相互接続ネットワークにおける各経路上のデータが有効か否かを示す機構が設けられている。使用されていないそれらの出力に対しては、ノイズまたは他の無効な情報が所定のプロセッサからメモリ経路へ送られる。もし、メモリ・バンクがこれらの無効な信号を或るメモリ位置に対して読み出すべき又は書き込むべき命令またはデータと解釈すれば、エラーが発生する可能性がある。

有効なデータが各信号線に存在することを保証するために、プロセッサの有効な出力のそれぞれは、シリアルなビット・ストリームを相互接続ネットワークの最初のスイッチに送り込む。このシリアルなビット・ストリームは、その出力からのデータが有効か否かを示すために、その中に「健全性コード(sanity code)」が埋め込まれている。この健全性コードは、プロセッサの有効な出力のそれぞれに相互に接続された最初のネットワーク・スイッチへシリアルにかつ連続的に送り出される6ビットのパターンである。したがって、図2に示した各CPUに対しては、CPUの有効な出力のそれぞれが、それぞれのメモリ・モジュール600の最初のスイッチへ健全性コードを送る。図3に示した各CPUに対しては、CPUの有効な出力のそれぞれが、それぞれのネットワーク・モジュール300の最初のスイッチへ健全性コードを送る。この健全性コードの目的は、相互接続ネットワーク内のモジュール間のインタフェースを保護することである。正しい健全性コードが受け取られると、その入力がアクセス経路を有している8個の出力の全てに、その受け取ったものが同報送信される。このようにして、例えば、図7では、問題となっているCPUの出力がメモリ・モジュール600の入力2に接続されていれば、1×1スイッチ610cが健全性コードが正しいか否かを調べ、正しい健全性コードを受け取ってあれば、それを、入力2がアクセス経路を有している8個のメモリ・バンク640a-hへ同報送信する。

健全性コードを監視するための制御回路を図13に示す。各モジュールの最初のスイッチにおける健全性検査回路700は、受け取った健全性コードを監視して、モジュールのその入力に入って来たデータが有効か否かを判定する。メモリ・モジュール600に対しては、各入力0~7に対応する各1×1スイッチ610a-h内に健全性検査回路700が存在する。ネットワーク・モジュール300に対しては、各入力0~15に対応する各1×1スイッチ310a-p内に健全性検査回路700が存在する。

健全性検査回路700は、対応するスイッチのFIFOをイネーブルにするものとして動作する。健全性コードが正しく受信されなければ、そのスイッチ内のFIFOは使用可能とはされず、そのスイッチに入る如何なる情報もFIFOバッファの中へは格納されない。

無効なCPUによって又はネットワーク・モジュールもしくはメモリ・モジュールの未接続の入力によって生じるノイズが健全性コードと同じになるということが極めてまれであるように、健全性コードを選択しなければならない。したがって、6ビットの健全性コードの適切な一例は、

010011

である。

10

20

30

40

50

上述のように、この6ビットの健全性コードは、CPUの有効な出力によってのみ、シリアルにかつ連続的に送り出される。その6ビットのコードがシリアルに16回連続して受け取られたと健全性検査回路が判定すると、一つのメモリ・バンクが動作可能とされて、相互接続ネットワークの所定の経路上のデータに注意を払うようになるだけである。如何なるときであっても正しい健全性コードが受け取られない場合は、そのメモリ・バンクは、正しい6ビットのコードが再びシリアルに16回連続して受け取られるまで、その経路上のデータを無視する。

結果として得られる96ビットのコード(6ビットのコードが16回繰り返されたもの)は、非常に明瞭な周波数とその中に組み込まれた信号パターン(signalling patterns)を有する信号である。これらの周波数および信号パターンがノイズと同じとなるのは極めてまれである。

本相互接続ネットワークのモジュール性により、小規模システムを構築するために使用されるのと同じのモジュールを、修正を加えずに、最大規模のシステムを構築するために使用することができる。図13および14は、本発明の相互接続ネットワークを使用して構築することができるマルチ処理システムのいくつかの例を記載した表を示している。図13は、1~4個のCPUを有するマルチ処理システムについてのモジュール数、セクション構成(section configurations)、およびアドレス・ビットの割り付け(address bit layout)を記載した表を示している。図14は、8~32個のCPUを有するマルチ処理システムの例に対する表を示している。図13および14に示したシステムおよびこの中で説明されたシステムの例は、決して、本相互接続システムを使用して構築できる唯一のマルチ処理システムではない。それどころか、無限の種類構成が本発明の範囲を逸脱することなく構築可能であろう。

#### <パケット化メッセージ>

好ましい実施形態では、書込情報(write references)は、2個の連続したパケットでメモリ・モジュール600へCPUによって送られる。読出情報(read references)は、単一のパケットで転送される。書込情報に対しては、1番目のパケットが必要なアドレスおよび制御情報の全てと書込データの半分とを収容している。2番目のパケットは、書込データの後半部を収容している。共通メモリへの書込については時間は本質的なものではないため、この方法は、ネットワーク・モジュールおよびメモリ・モジュールの出力経路において必要とされる相互接続を最小化する。

好ましい実施形態では、2番目のパケットは1番目のパケットに追隨して、直後のクロックの周期でスイッチに入ってそれを通過し、1番目のパケットと同一の相互接続経路に沿って進む。このため、各データ移動にはクロックの2周期の期間を要する。このように、2番目のパケットは、1番目のパケットに対してちょうどクロックの1周期分だけ遅れて、メモリ・モジュールおよびネットワーク・モジュールの各 $n \times m$ スイッチ内に存在する。

書込情報をパケット化して、それらにネットワーク・モジュールおよびメモリ・モジュールを順次巡回させることにより、各 $n \times m$ スイッチの実現に要する相互接続の数が1/2に低減される、ということが当業者には容易にわかるであろう。現在の設計では、チップ内のゲートを利用し得る相互接続よりも、比率としては多くのゲートがチップに含まれている。したがって、相互接続は、多くのチップ設計において乏しい資源である。例えば、4個の独立した入力の経路および4個の独立した出力の経路を有する $4 \times 4$ スイッチと、120ビット幅のメモリ参照とを想定されたい。そのようなクロスバー・スイッチを作製するのに必要な相互接続の数は、この場合、120の8倍の数となる。これは極めて大きな数であって、単一チップに適合する数よりもはるかに大きく、実際、数個のチップに適合する数よりもはるかに大きい、ということが当業者には容易にわかるであろう。

メモリ参照をパケット化することにより、120ビット幅のメモリ参照が1/2に減少する。これにより、システムの全ての相互接続も1/2に減少する。これは極めて大きな意味を持ち得る低減である。本発明において使用されるこのパケット化方法は、各スイッチを単一チップ内に実現できるという点で、システム設計を大幅に簡潔にするものである、というこ

10

20

30

40

50

とが当業者には容易にわかるであろう。

しかし、読出情報は書込情報のようにパケット化されない。出力経路上において、読出情報は制御情報と所望の読出アドレスのみを有している。読み出しの待ち時間をできるだけ短くすることが重要であるため、メモリ・モジュールおよびネットワーク・モジュールの入力経路606および304は、それぞれ、その幅が完全に1ワード分となっている。このため、各読出情報は、書込情報の出力経路のように2クロックを必要とするのではなく、各クロック毎に転送される。また、相互接続を更に節約するために、目的地コード(すなわち、プロセッサ・リターン・アドレス)を出力経路の読出情報の書込データのフィールドの中に入れて送る。この方法は、相互接続ネットワークの出力経路上で必要となる相互接続の数を更に節約する、ということが当業者には容易にわかるであろう。

10

#### <ハンドシェイク・プロトコル>

レディ/レジューム(READY/RESUME)のハンドシェイク・プロトコルにより、連続する $n \times m$ スイッチ層の間での円滑な通信が達成される。図12は、 $4 \times 4$ スイッチ620に接続された $1 \times 1$ スイッチ610の制御回路の更なる詳細を示す。本ハンドシェイク・プロトコルは、図12に示した $1 \times 1$ スイッチおよび $4 \times 4$ スイッチを特に参照して説明されるが、スイッチ610および620は単に例に過ぎず、これらは、図6~11を参照しつつ上記において説明されたスイッチのいずれかの代表として考えている。

本相互接続ネットワークにおける各 $n \times m$ スイッチは、そのスイッチの各出力に対応づけられたカウンタを備えている。したがって、例えば、図12に示したスイッチ610は、その単一の出力の経路に対応づけられた1個のカウンタ612を備えている。 $4 \times 4$ スイッチ620は4個のカウンタ624a-dを備えていて、それらのカウンタのそれぞれは、4個の群バッファ623a-dのうちの一つに対応している。目的スイッチ内のFIFOバッファの数に等しい値を最大のカウンタ値としている。したがって、例えば、カウンタ612は、 $4 \times 4$ スイッチ620のFIFO621a内のバッファA~Dの数である4までカウントする。

20

これらのカウンタは、2個のスイッチの間で伝送されるレディおよびレジュームのハンドシェイク信号により、値が増加または減少する。レディ信号はカウンタの値を増加させる。レジューム信号はカウンタの値を減少させる。したがって、 $1 \times 1$ スイッチ610は、 $1 \times 1$ スイッチ610のカウンタ612の値が4よりも小さい間は更にメッセージを $4 \times 4$ スイッチ620に送信できることを知る。カウンタ612の値が4であれば、 $1 \times 1$ スイッチは、 $4 \times 4$ スイッチ620内のFIFOバッファ621aが満杯であり、 $1 \times 1$ スイッチ610から更なるメッセージの受理はできないことを知る。

30

カウンタの値は0から始まる。一つのメッセージが $1 \times 1$ スイッチ610内のFIFOのバッファAから送り出される毎に、レディ信号もそのメッセージの内容とともに $4 \times 4$ スイッチ620へ送られる。このレディ信号は二つの目的を持っている。第1にレディ信号は、一つのメッセージがバッファAから送り出される毎にカウンタ612に入力される。レディ信号はカウンタ612の値を増加させる。第2にレディ信号は、FIFO621aにより、 $1 \times 1$ スイッチ610から送られるメッセージ・データに対する有効信号として使用される。このカウンタの値は0から始まるため、 $1 \times 1$ スイッチ610は、4個のメッセージを送ることができ、それらの送過後は $4 \times 4$ スイッチ620から送られるレジューム信号を待たなければならない。

40

$4 \times 4$ スイッチ620内のFIFO621aのバッファAからメッセージが出て行くと、必ずレジューム信号が $1 \times 1$ スイッチ610へ送り返される。このレジューム信号はカウンタの値を減らし、FIFO621a内で一つのスペースが空いたことを示す。

レディ/レジュームのハンドシェイクは、この中で説明された各 $n \times m$ スイッチに対して同じように機能する。このようにして本相互接続ネットワーク内の各スイッチ層の間での円滑な通信が達成される。

#### <メモリのタイミング>

好ましい相互接続ネットワークは、広範囲のメモリ・アクセス時間およびサイクル時間での使用に適應できる。制御回路は、各メモリ・モジュール600の $1 \times 2$ スイッチ630内に配置されている。このスイッチの更なる詳細を図16に示す。図16は、FIFOバッファ631、ア

50

ドレス経路633、制御ビット経路634、および遅延チェイン (delay chain) 636を示しており、これらの全ては、 $1 \times 2$  スイッチ630内に存在する。遅延チェイン636は、マルチ処理システムに取り付けられたメモリの速度とサイクル時間に応じて、 $1 \times 2$  スイッチ630における遅延を制御する。マルチ処理システムに取り付けられたメモリのアクセス時間は、8位置セレクタ637aによって $1 \times 2$  スイッチ630へ知らされる。セレクタ637aは、制御タグ (control tag) が入力635a - hのいずれに入るかを制御することにより、遅延チェイン636aの長さを制御する。この制御タグは、遅延チェインの8入力635a - hのうちの一つに入る。この制御タグは、目的地コード、リターン・アドレス、および他の制御情報を有している。セレクタ637aは、遅延チェイン636a内のいずれの場所に制御タグが入るかを制御する。

10

このタグが一旦遅延チェイン636aに入ると、このタグは、遅延チェイン636aの最後に到達するまで、そのチェインを流れていく。その後、この制御タグは、メモリ・バンクから出て来るメッセージ・データを採取してそのメッセージをFIFO652における最初の空いているバッファにロードする時点になると、メモリ・バンクの出力側における $2 \times 1$  スイッチ650に信号を送るために使用される。

メモリのサイクル時間は、メモリのアクセス時間と同様に制御される。8位置のセレクタ637bが、まさに上述のようにして同様の遅延チェイン636bを制御する。しかし、セレクタ637bおよび遅延チェイン636bはメモリのサイクル時間を制御するため、この遅延チェインの出力は、 $1 \times 2$  スイッチ630のバッファ632aおよび632bまで送られて、メッセージがその $1 \times 2$  スイッチを出てメモリ・バンクへ行く頻度を制御する。

20

更なるセレクタ637cは、マルチ処理システムに取り付けられたメモリ部分がパイプライン方式のメモリ部分か、パイプライン方式ではないメモリ部分かを示す2位置のセレクタである。このように、本相互接続ネットワークは、異なるアクセス時間およびサイクル時間を有する多くの異なる種類のメモリを本相互接続ネットワークにおいて使用することを想定している。例えば、SRAMまたはDRAMを使用することができる。どのようなメモリが使用されるかに応じて、セレクタ637はそのメモリの種類に対応するように設定される。

本相互接続ネットワークの設計はマルチ処理システムのスケラビリティ (scalability) を大きく増大させる、ということが当業者には容易に認識できるであろう。ネットワーク・モジュールおよびメモリ・モジュール内の各 $n \times m$  スイッチはその隣接部と通信するだけでよいため、この相互接続ネットワークは、極めてスケラビリティに富んでいる。

30

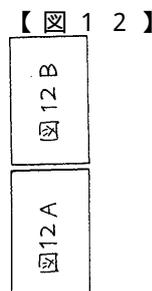
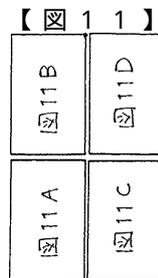
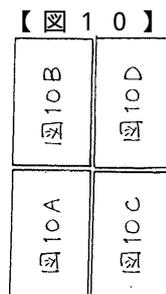
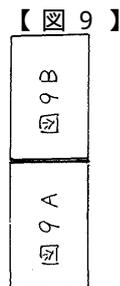
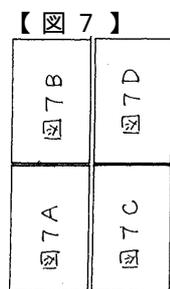
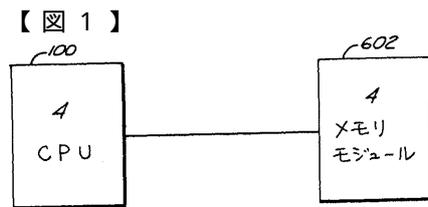
個々の $n \times m$  スイッチのそれぞれは、相互接続ネットワーク全体について、または、2つ若しくは3つ前方の層で起こっている事象についての大域的な情報を持つ必要はない。各 $n \times m$  スイッチは、CPUとメモリとの間で調停を行いデータを送るために、局所的な情報を持つだけでよい。同様に、CPUは、メモリ参照を送出する前に、相互接続ネットワークについて又はメモリの応答時間についての大域的な知識を持つ必要はない。この方法において、CPUは、そのCPUからメモリ・バンクまでのデータ経路全体が空いていることを確かめる必要はない。本相互接続ネットワークを使用すると、CPUは、簡単に「火をつけて忘れる (fire and forget)」ことができる。本相互接続ネットワークは大域的な調停回路 (arbitrator) を必要としないことから、CPUや、 $n \times m$  スイッチ、オペレーティングシステムについての如何なる再設計も必要としないので、単純な複製によって極めて容易にシステムの規模を変えることができる、ということを当業者は理解するであろう。異なる構成における構成要素を接続し直すだけで、アーキテクチャの規模が直ちに変更される。本相互接続ネットワークの他の利点は、この相互接続ネットワークは使用されるCPUの種類に依存せず、またメモリの速度にも依存しないということである。メモリの応答時間に依存しない如何なる種類のプロセッサの構成も、本発明のシステムに適應できる。このようなプロセッサ構成の一つは、同時係属中で本出願人に譲渡され、1994年12月28日に発行された、オーバーリン (Oberlin) らによる「SCALABLE PARALLEL VECTOR COMPUTER SYSTEM (スケラブルな並列ベクトル計算機システム)」という名称のEP公報第0553158号の中に見ることができる。参考までにこの出願をここに挙げておく。また、たとえメモリの全てのバンクが異なる速度で動作しても、本相互接続ネットワークは正しく機能し、複雑な

50

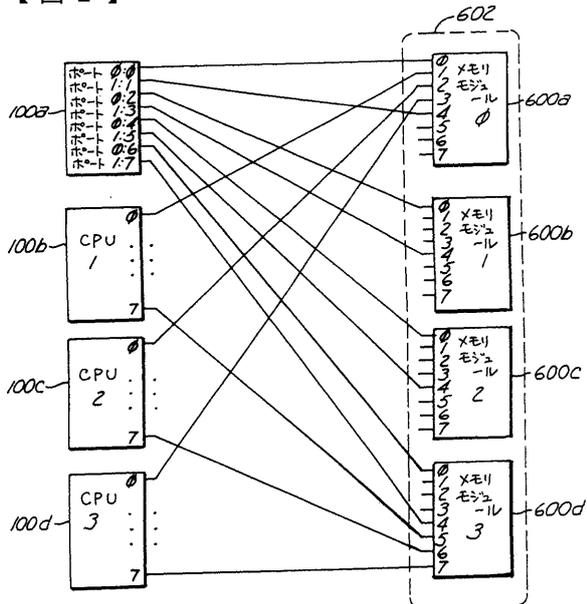
再設計を行う必要はない。したがって本相互接続ネットワークの設計はこの点において極めて柔軟性に富む、ということが当業者はわかるであろう。このように、本相互接続ネットワークで使用されるCPUは、時間に影響を受けやすいものではなく、したがって、メモリの待ち時間がより長い場合であっても効率よく動作する。本相互接続ネットワークの設計は、大量のデータがCPUとメモリとの間で移動できるように、バンド幅を広げるものである。このような設計方式により、極めて広いバンド幅で極めて高いスループットの計算機が得られる、ということを当業者は容易に理解するであろう。すなわち、これは、並列度の高いベクトル処理のタスクに特によく適合した設計である。

ここでは特定の実施形態が説明されているが、同一の目的を達成することを意図している如何なる構成も、その説明された特定の実施形態の代わりに用いることができる、ということを当業者は理解するであろう。例えば、本発明のアーキテクチャは、異なる数のプロセッサ、異なる量のメモリ、または、プロセッサ1個当たりの異なる数のプログラム・スレッド (program threads) で実現されてもよい。また、本発明のアーキテクチャは、使用されるプロセッサの特定の種類の種類、メモリの速度、または、この詳細な説明において開示された特定の論理設計のいずれにも依存しない、ということが当業者にはわかるであろう。本発明の範囲を逸脱することなく、異なる種類のプロセッサを使用できるであろう。この出願は、本発明の改造または変形をカバーするものである。したがって、この発明は請求の範囲によってのみ限定されることは明かである。

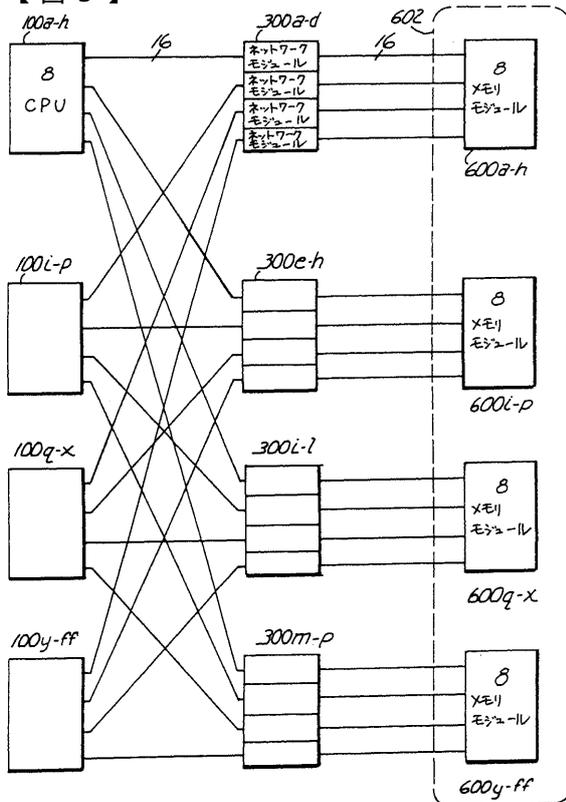
10



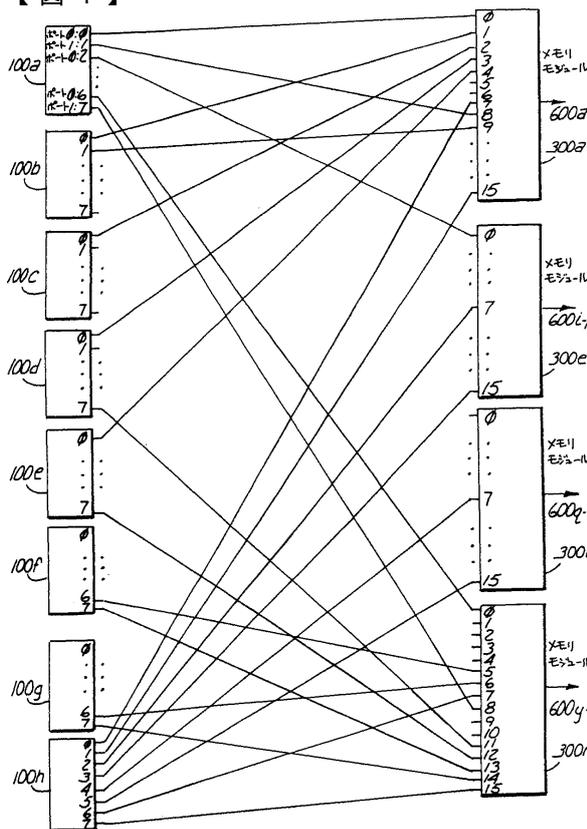
【図2】



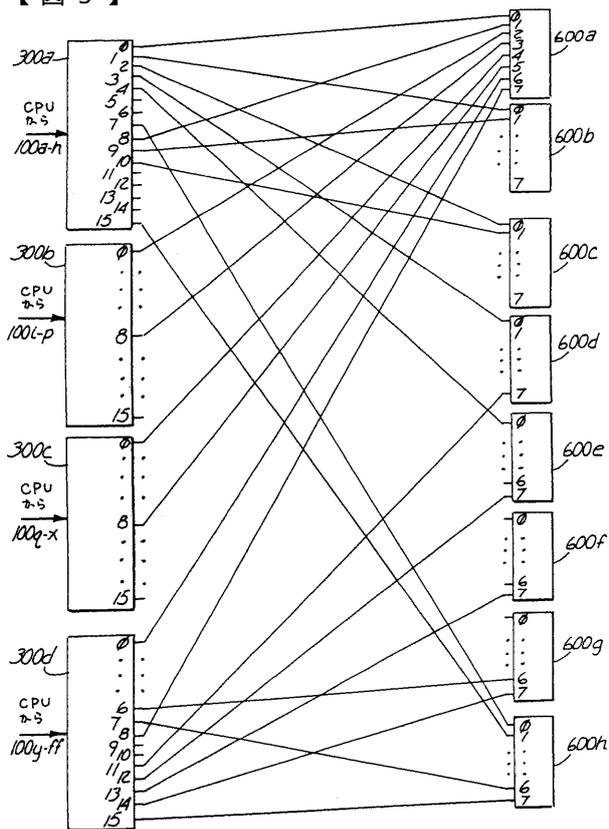
【図3】



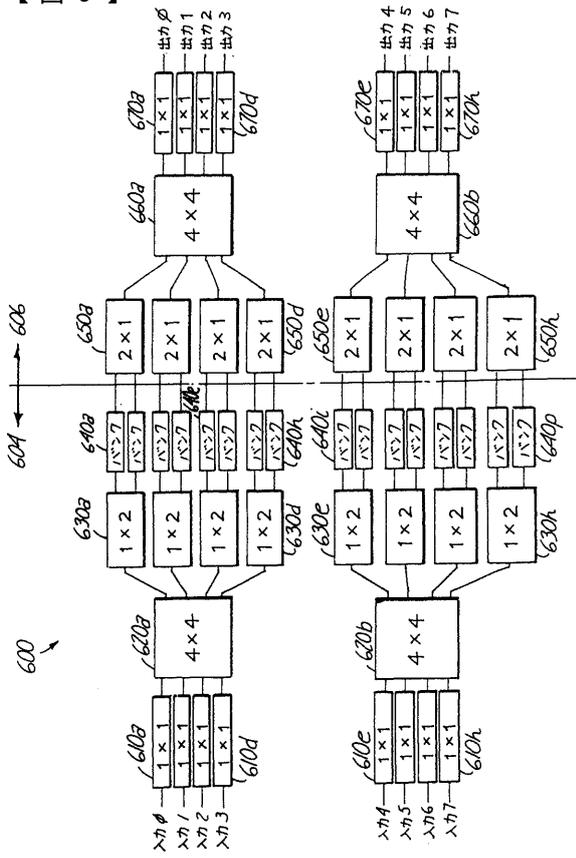
【図4】



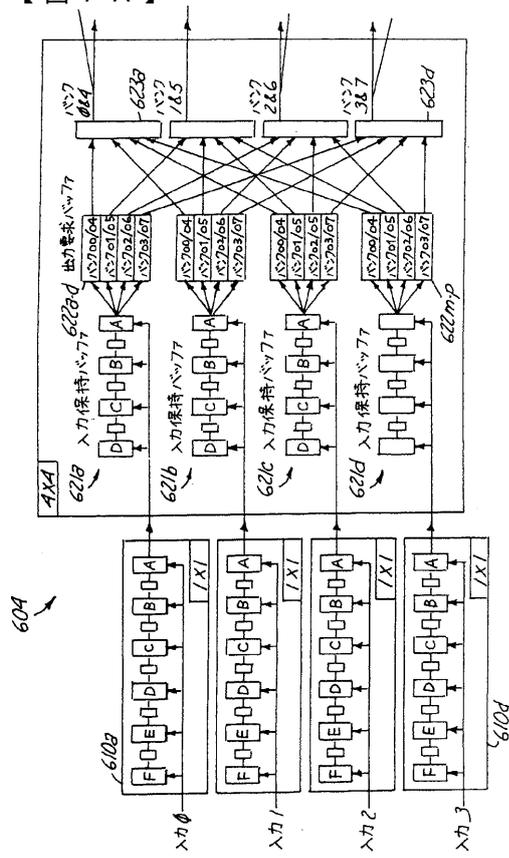
【図5】



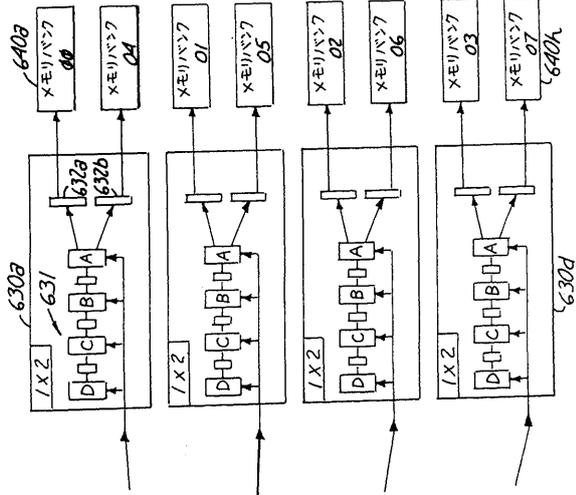
【図 6】



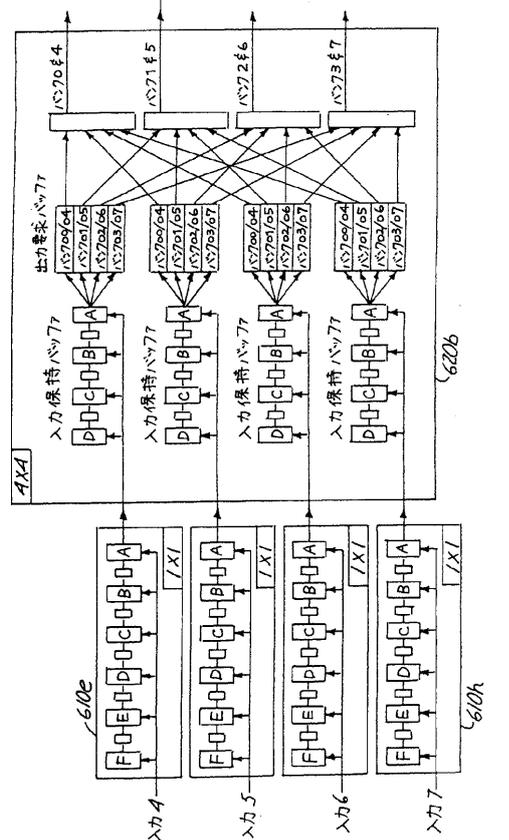
【図 7 A】



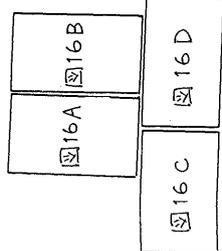
【図 7 B】

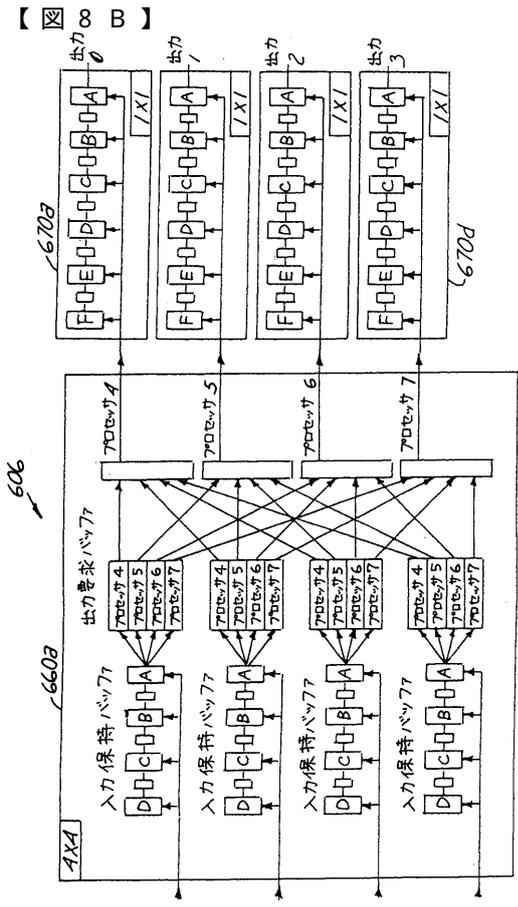
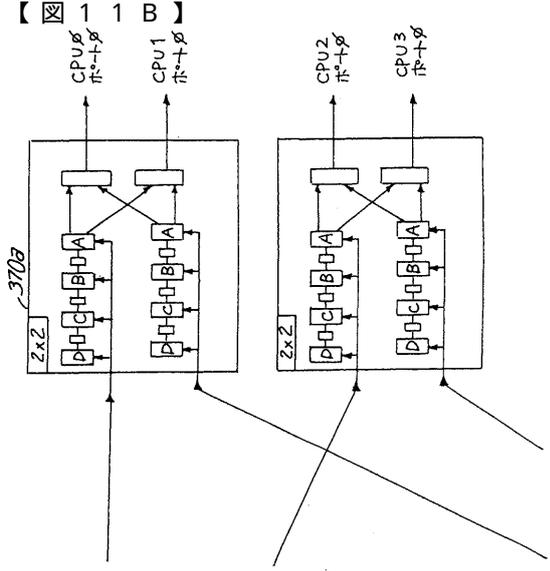
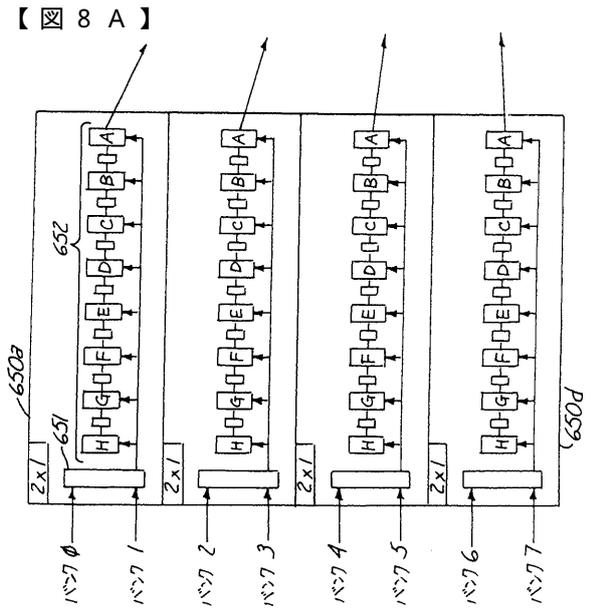
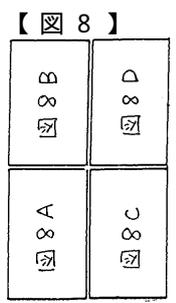
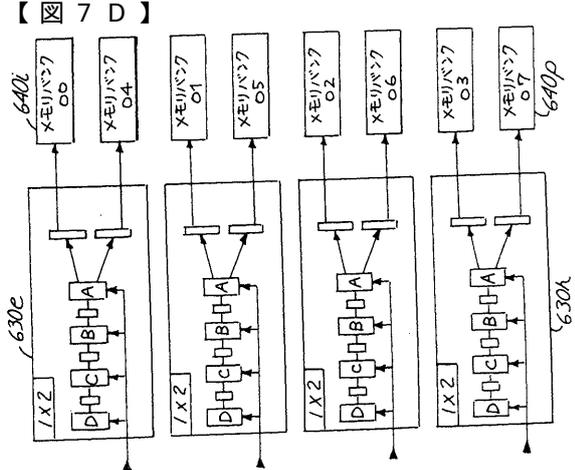


【図 7 C】

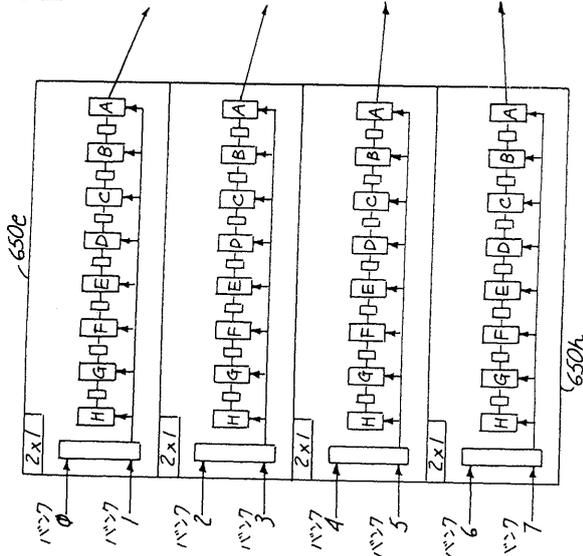


【図 16】

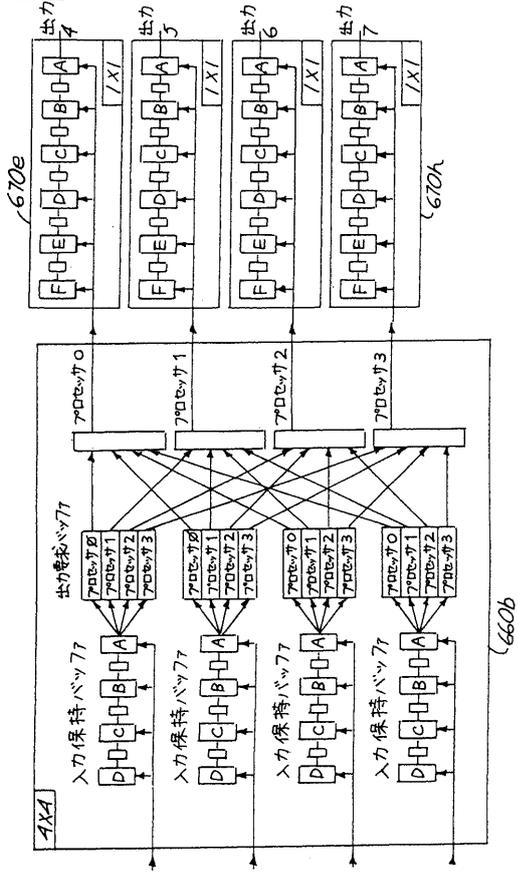




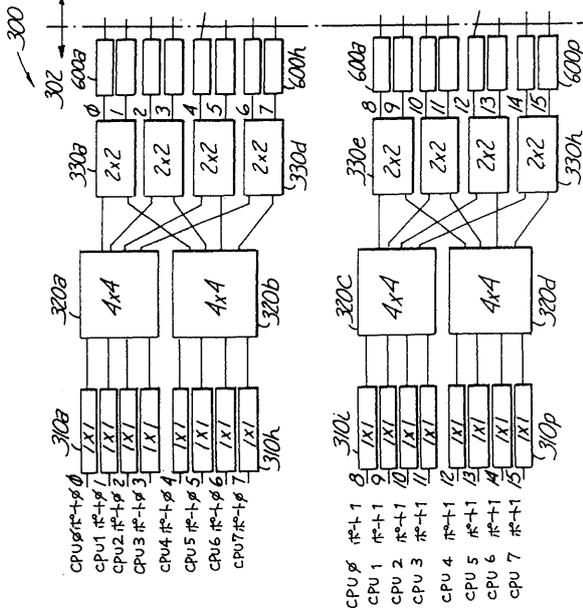
【図 8 C】



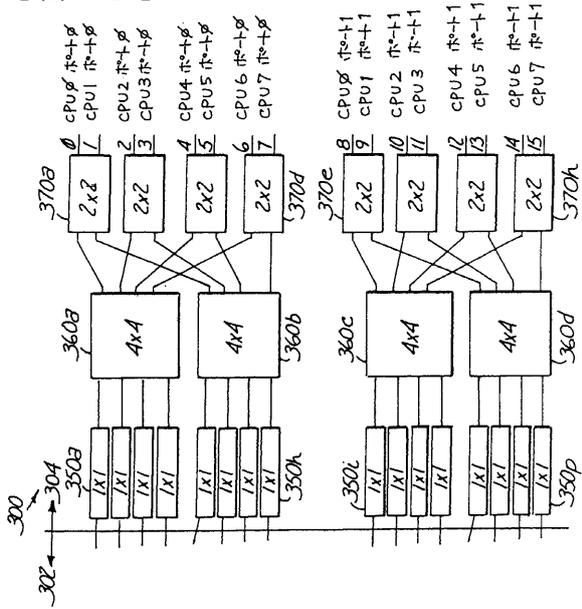
【図 8 D】



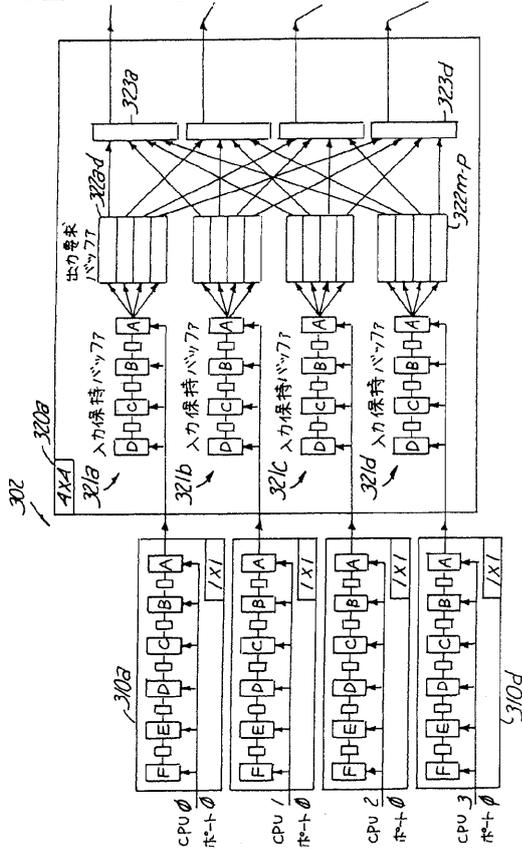
【図 9 A】



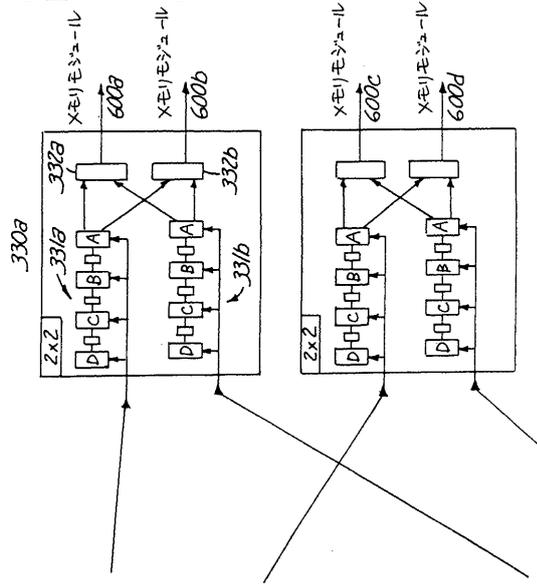
【図 9 B】



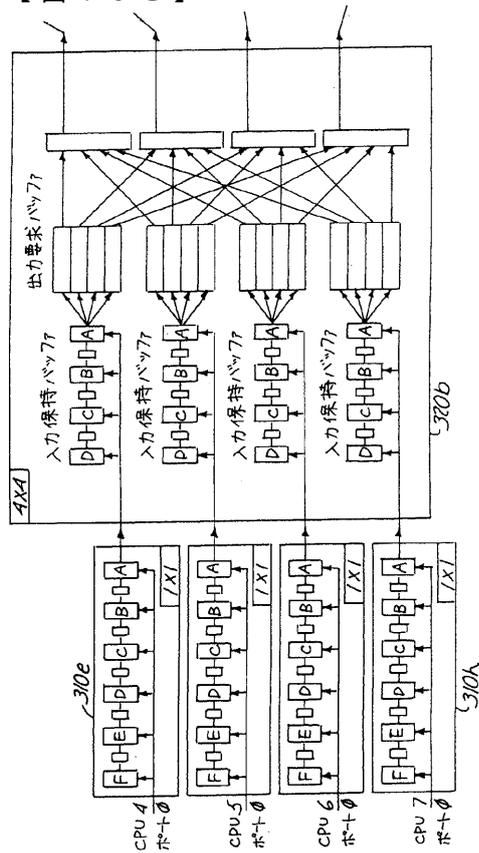
【図10A】



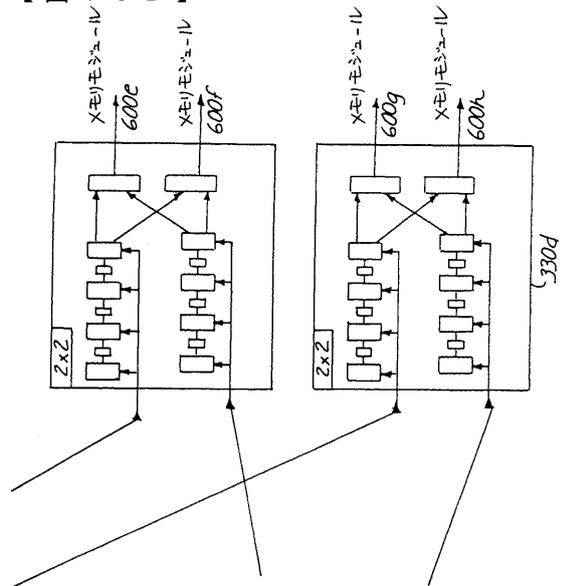
【図10B】



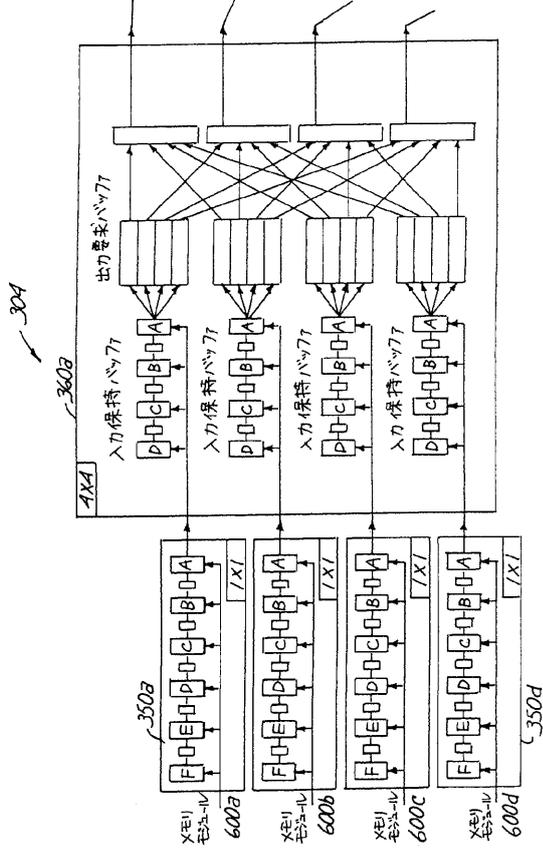
【図10C】



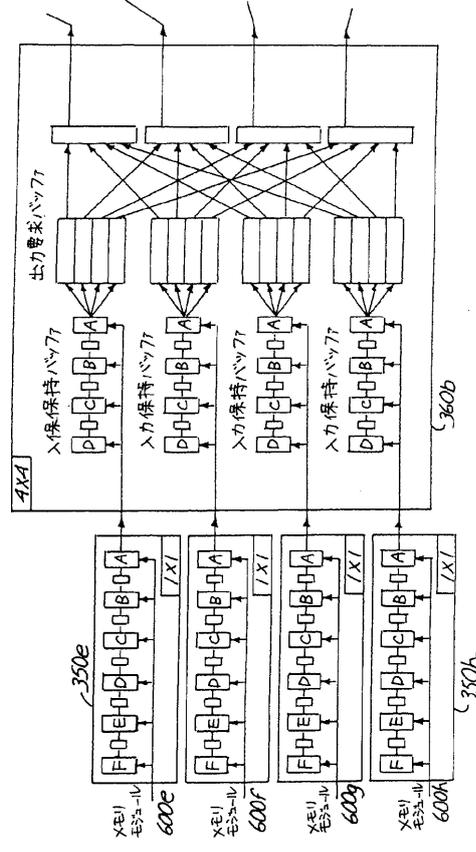
【図10D】



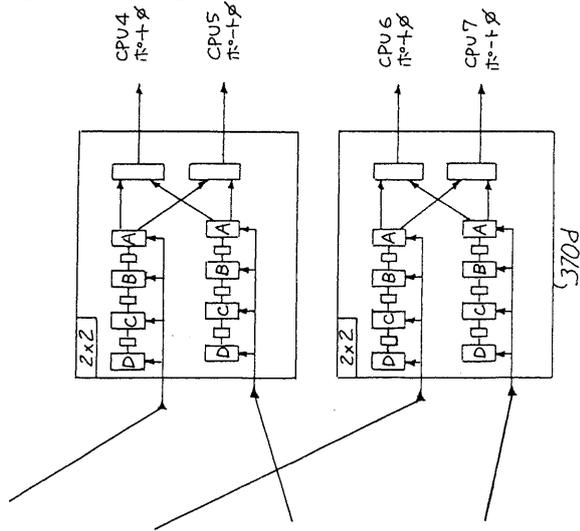
【 図 1 1 A 】



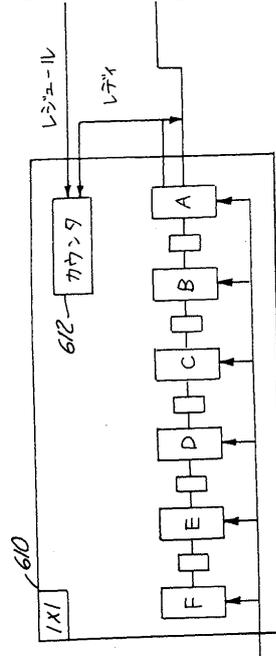
【 図 1 1 C 】

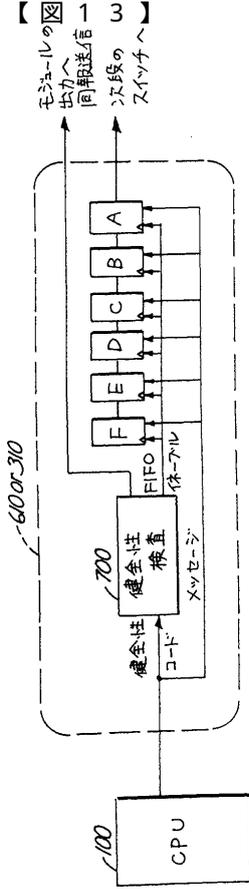


【 図 1 1 D 】



【 図 1 2 A 】

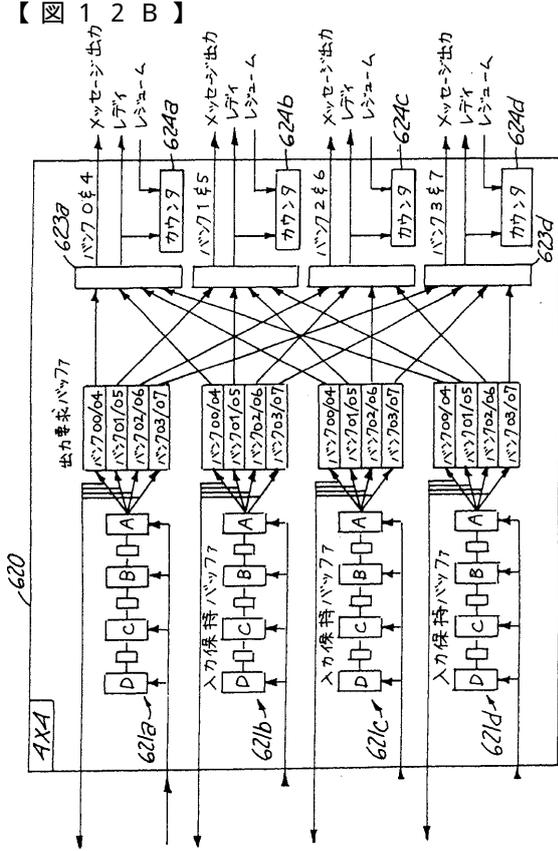




【 図 1 4 】

(--送判別目的に0で表されるビット)

モジュール数		セッション構成			アドレスビットの割り付け		
プロセッサ	ネットワーク	メモリ	セッション数	セッション数	セッション	セッション	セッション
1654	0	4	8	1	00.01.02	...	03.04.05...06-31
1654	0	4	8	1	00.01.02	...	03.04...05-30
1654	0	4	4	1	00.01...	...	02.03.04...05-30
1654	0	4	4	1	00.01...	...	02.03...04-29



【 図 1 5 】

(--送判別目的に0で表されるビット)

モジュール数		セッション構成			アドレスビットの割り付け		
プロセッサ	ネットワーク	メモリ	セッション数	セッション数	セッション	セッション	セッション
8	16	8	2	128	00.01.02	...	04.05.06...07-32
8	16	8	2	64	00.01.02	...	04.05...06-31
8	16	8	4	2	00.01...	...	03.04.05...06-31
8	16	8	4	2	00.01...	...	03.04...05-30
8	16	16	8	4	00.01.02	...	05.06.07...08-33
8	16	16	4	128	00.01.02	...	05.06...07-32
8	16	16	4	128	00.01...	...	04.05.06...07-32
8	16	16	4	64	00.01...	...	04.05...06-31
8	32	16	8	512	00.01.02	...	06.07.08...09-34
8	32	16	8	256	00.01.02	...	06.07...08-33
8	32	16	8	256	00.01...	...	05.06.07...08-33
8	32	16	4	128	00.01...	...	05.06...07-32



---

フロントページの続き

(72)発明者 ロッキー、ピーター・ジー  
アメリカ合衆国 54729 ウィスコンシン州 チッペワ・フォールズ、ピー・オー・ボックス  
212、ルート3番

審査官 鳥居 稔

(56)参考文献 特開平01-298459(JP,A)  
特開平04-328653(JP,A)  
特表平06-502033(JP,A)  
欧州特許出願公開第00339270(EP,A1)  
欧州特許出願公開第00553158(EP,A1)

(58)調査した分野(Int.Cl.<sup>7</sup>, DB名)  
G06F 15/173  
H04Q  
H04L