

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3704367号
(P3704367)

(45) 発行日 平成17年10月12日(2005.10.12)

(24) 登録日 平成17年7月29日(2005.7.29)

(51) Int. Cl.⁷

G06F 15/173

F I

G06F 15/173 650M

請求項の数 3 (全 9 頁)

| | | | |
|-----------|-----------------------|-----------|---|
| (21) 出願番号 | 特願平6-1913 | (73) 特許権者 | 000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号 |
| (22) 出願日 | 平成6年1月13日(1994.1.13) | (74) 代理人 | 100075096 弁理士 作田 康夫 |
| (65) 公開番号 | 特開平7-210528 | (74) 代理人 | 100068504 弁理士 小川 勝男 |
| (43) 公開日 | 平成7年8月11日(1995.8.11) | (72) 発明者 | 田中 輝雄 東京都国分寺市東恋ヶ窪1丁目280番地 株式会社 日立製作所 中央 研究所内 |
| 審査請求日 | 平成13年1月10日(2001.1.10) | (72) 発明者 | 保田 淑子 東京都国分寺市東恋ヶ窪1丁目280番地 株式会社 日立製作所 中央 研究所内 |

最終頁に続く

(54) 【発明の名称】 スイッチ回路

(57) 【特許請求の範囲】

【請求項1】

多段ネットワークを構成するm入力n出力のデータ転送用スイッチ回路であって、各入力ポート、出力ポート対応にそれぞれ入力バッファ、出力バッファを有し、各入力ポートからの転送データ間の競合調停を出力ポート対応に行い、前記出力ポートの競合で選択されなかった前記転送データは、通過し始めた複数ブロックからなる転送メッセージのすべてのブロックが通過して前記出力ポートが空くまで、前記入力バッファで待たせられる制御を行うデータ転送用スイッチにおいて、前記出力ポートの競合で選択されなかった転送メッセージを一時、保持するFIFOバッファを更に有し、転送メッセージが通過中の入力ポートを解放することを特徴とするスイッチ回路。

【請求項2】

請求項1において、前記FIFOバッファ上に前記転送メッセージが一時保持されているときの前記出力ポートの競合調停には、前記FIFOバッファ上のメッセージをm+1番目の入力として扱うスイッチ回路。

【請求項3】

請求項2において、前記FIFOバッファ上に前記転送メッセージが一時保持されているときの前記出力ポートの競合調停には、前記FIFOバッファ上のメッセージの優先度を高くするスイッチ回路。

【発明の詳細な説明】

【0001】

10

20

【産業上の利用分野】

本発明は、多数台のプロセサからなる並列プロセサなどを接続するネットワークを構成するスイッチ回路に関する。

【0002】**【従来の技術】**

科学技術計算処理分野において、多数台の要素プロセサを結合し、台数分の性能向上をねらう並列プロセサが商用化しつつある。並列プロセサでは、この多数台のプロセサを接続するプロセサ間結合方式が重要である。現在、このプロセサ間結合方式は、いくつかの m 入力 n 出力スイッチ(m, n は整数)を複数接続した多段ネットワークが一般的である。

【0003】

図4に代表的な多段ネットワークであるベネスネットワークを示す。もちろん、格子結合やハイパキューブ結合などのキューブ系結合や、トリー結合なども m 入力 n 出力スイッチを複数接続した多段ネットワークである。図4のベネスネットワークでは、4入力4出力のスイッチを3段接続することにより、16台のプロセサを接続している。このうち、最初の一段分は冗長なスイッチであり、このような冗長な段を持つことにより、転送先プロセサが重ならない時のすべての組み合わせのデータ転送をネットワーク内で衝突なしに行う経路を確保することができる。

【0004】

なお、以下では、プロセサ間のデータ転送をメッセージ転送と呼ぶ。メッセージとは、メッセージ情報として、転送すべきデータにあて先のプロセサ番号、転送データ量(メッセージ長)などを転送データに付加したもので、このメッセージ内のメッセージ情報を用いて、メッセージは能動的に各スイッチ内で経路を判断して、受信先のプロセサに送るべきデータを転送する。

【0005】

図5にメッセージの構成の例を示す。転送先プロセサ番号およびメッセージ長などのメッセージ構成情報は、一般にメッセージの最初の部分に保持されている。メッセージは、物理的な信号線数の制約により、複数回に分けて送られる。そのため、各部分を送るごとにデータ線の有効性を示す制御信号が必要となる。また、メッセージの転送先情報はメッセージの先頭のみしかないので、通過し始めたメッセージがすべて通りきるまで通過中の経路はかならず確保しておく必要がある。

【0006】

並列計算機の対象とする応用は、もともと大規模科学技術計算処理が考えられていたが、非数値処理として、推論処理などの知識処理の分野やデータベースやオンライントランザクション処理などと応用範囲が拡大されつつある。また、このような分野では、大規模科学技術計算処理と違って、一回に転送する転送データ量はあまり長くないと考えられる。さらに、最近実用化されだした、すべてのプロセサ内の記憶装置のアドレス空間を一元化する分散型共有記憶方式の並列計算機では、いままで以上に他のプロセサへのアクセスが必要となり、このデータアクセス単位は、一つのデータ(数バイト単位)からキャッシュのライン相(数10~数100バイト単位)と非常に小さい。

【0007】

一方、ハードウェアの状況を見ると、LSI高集積化は非常に進み、1チップあたりのゲート数は飛躍的に高まっている。しかし、信号ピン数は、ゲート数程には集積化が進んでおらず、ゲート/ピン比率は高まる方向にある。

【0008】**【発明が解決しようとする課題】**

ベネスネットワークでは、冗長なスイッチを設けることにより、転送先プロセサが重ならない時のすべての組み合わせのデータ転送をネットワーク内で衝突なしに行う経路を確保することができる。しかし、定形的な科学技術計算処理での静的なプロセサ間のメッセージ転送(つまり、事前にメッセージ転送手段が分かっている場合)を除いては、転送先プロセサが重ならないことを保証する事は困難である。また、仮りに転送先プロセサが重

10

20

30

40

50

ならないとしても各プロセサが独立に動作するような（いわゆるM I M D型：Multiple Instruction Mutiple Data）並列計算機システムでは、各プロセサの処理の進み方に差が生じる。このような場合、どうしてもネットワーク内のスイッチ上でメッセージ間で競合が発生し、メッセージの転送時間がどんどん延びてしまう。

【0009】

たとえば、同じ入力ポートから続けて異なる転送先プロセサへの複数のメッセージ転送があったとする。もし、最初のメッセージが他の入力ポートからのメッセージとの競合調停の結果、待つことになった場合、その次のメッセージが要求する転送先プロセサにつながる出力ポートが空いていても、この次のメッセージを送ることができない。このことは特に小さな単位のデータを転送する（つまり、メッセージ長が短い）場合、その転送時間が性能に与える影響は大きい。

10

【0010】

本発明の目的は、このような場合にも、競合待ちしているメッセージのある入力ポートを解放し、その次のメッセージを転送することを可能とする多段スイッチネットワークを構成するスイッチを提供することにある。

【0011】

【課題を解決するための手段】

上記課題を解決するためには、ネットワーク上のスイッチ内にメッセージの転送単位より、大きいF I F O（First In First Out）バッファを設けて、競合により待たせられたメッセージを一時的にこのF I F Oバッファに退避することにより、入力ポートを解放し、次のメッセージの競合調停への参加および転送を可能とする。

20

【0012】

【作用】

上記、F I F Oバッファを設けることにより、入力ポートの解放を可能とし、次のメッセージの転送を早めることができる。また、F I F Oに一時的に退避したメッセージは、優先度を高くして、次の選択に用いるようにすれば、特にF I F Oに退避したメッセージが必要以上に待たれることはない。

【0013】

【実施例】

図1に本発明のスイッチ構成の一例を示す。図中、10Aないし10Dは各入力ポートに対応する入力バッファ、11はメッセージ選択処理部、12はセレクタ、13は本発明により採用したF I F Oバッファ、14Wないし14Zはセレクタ、15Wないし15Zは各出力ポートに対応する出力バッファである。また、線L20Aないし線L20Dは、前段のスイッチないしプロセサからのメッセージの入力、線L21Aないし線L21Dは、入力バッファ10Aないし10Dがメッセージの待合せにより、これ以上メッセージの格納が不可能になったことを前段のスイッチないしプロセサに対して知らせるメッセージ送信抑止出力信号、線L23Wないし線L23Zは、後段のスイッチないしプロセサへのメッセージの出力、線L24Wないし線L24Zは、後段のスイッチないしプロセサからのメッセージ送信抑止信号である。

30

【0014】

前段のスイッチないしプロセサからの入力メッセージを線L20Aあるいは線L20Dを介して受け取った入力バッファ10Aあるいは10Dは、入力メッセージから転送先プロセサ番号およびメッセージ長などのメッセージ構成情報を取り出し、線L22Aあるいは線L22Dを介してメッセージ選択処理部11に送り、また、メッセージ本体をセレクタ14Wあるいは14Zに送る。

40

【0015】

図1では、入力ポート4，出力ポート4としたが、もちろん、これ以外での数でも構わない。さらに、入力ポート数と出力ポート数は異なっても良い。

【0016】

メッセージ競合調停部11では、複数の入力ポートからのメッセージ転送要求から、出力

50

ポートごとに、線 L 2 9 W ないし線 L 2 9 Z を介してセレクタ 1 4 W ないし 1 4 Z を切り替えることにより、メッセージ転送要求のある一つの入力ポートを選択する。この選択の時には、各出力ポートに対応する出力バッファ 1 5 W ないし 1 5 Z の使用状況（出力バッファがメッセージ取り込み可能か否か）の情報が線 L 3 0 W ないし線 L 3 0 Z を介してメッセージ競合調停部 1 1 に対して送られており、この情報もメッセージ選択の時に用いられる。

【 0 0 1 7 】

メッセージ競合調停部 1 1 でセレクタ 1 4 W ないし 1 4 Z の切り替えで選択されたメッセージは出力バッファ 1 5 W ないし 1 5 Z に蓄えられ、線 L 2 3 W ないし線 L 2 3 Z を介して、後段のプロセサあるいはスイッチの入力ポートに転送される。線 L 2 4 W ないし線 L 2 4 Z は後段からのメッセージ転送抑止信号であり、この信号が出ているときは、後段へのメッセージ転送は抑止される。

10

【 0 0 1 8 】

このスイッチには、本発明で用いる F I F O バッファ 1 3 が設けられており、複数の入力ポートから同一の出力ポートへの転送要求による競合のため、一時待機を余儀なくされたメッセージをセレクタ 1 2 を介して取り込むことができる。この時の制御として、メッセージ競合調停部 1 1 から

- (1) 線 L 2 5 を介してセレクタ 1 2 の切り替え、
 - (2) 線 L 2 6 から F I F O バッファへの書き込み場所の指示、
 - (3) 線 L 2 7 から F I F O バッファからの読みだし場所、
- の指示が送られてくる。

20

【 0 0 1 9 】

F I F O バッファ 1 3 からの出力としてのメッセージは、線 L 2 8 を介して、各出力ポートに対応するセレクタ 1 4 W ないし 1 4 Z に送られると同時に、メッセージ競合調停部 1 1 にもメッセージ内の転送先プロセサ番号情報を送り、メッセージ競合調停部 1 1 での出力ポートごとの入力ポート選択候補の一つとして扱う。この時、F I F O からの要求はできるだけ優先度を高くする方がよい。

【 0 0 2 0 】

図 2 にメッセージ競合調停部 1 1 の内部構成を示す。図において、4 0 は O R 回路、4 1 W ないし 4 1 Z は出力ポート対応競合調停部、4 2 は F I F O バッファ制御部である。

30

【 0 0 2 1 】

出力ポート対応競合調停部 4 1 W ないし 4 1 Z は、それぞれ出力ポートごとに対応するメッセージ競合調停部である。これらの競合調停部はすべて同じ構成をしている。入力は、各入力バッファ 1 0 A ないし 1 0 D (図 1) から線 L 2 2 A ないし線 L 2 2 D を介したメッセージ送信要求と、各出力バッファ 1 5 W ないし 1 5 Z (図 2) から線 L 3 0 W ないし線 L 3 0 Z を介した出力バッファの使用状況情報である。一方、出力は、線 L 2 9 W ないし線 L 2 9 Z を介したセレクタ

1 4 W ないし 1 4 Z (図 1) の切り替え信号と、線 L 1 9 A ないし線 1 9 D を介した各入力バッファ 1 0 A ないし 1 0 D (図 1) への競合調停結果の回答である。

【 0 0 2 2 】

出力ポート対応競合調停部 4 1 W ないし 1 4 Z は、後述するように F I F O バッファを第 5 の優先度の高い入力とし、(1) F I F O バッファ制御部 4 2 から F I F O バッファ上のメッセージ情報を受け取ること、(2) 競合調停に負けたメッセージ転送要求のある入力ポートに関するメッセージ情報を F I F O バッファ制御部 4 2 に転送すること以外は、通常のスイッチのメッセージ競合調停をそのまま用いることができるので、ここでは、詳細な記述は行わない。

40

【 0 0 2 3 】

次に、F I F O バッファ制御部 4 2 について説明する。F I F O バッファ制御部 4 2 は、F I F O バッファ 1 3 (図 1) の制御を行う。

【 0 0 2 4 】

50

入力は、各出力ポート対応競合調停部 4 1 Wないし 4 1 Zからのメッセージ情報として、線 L 5 0 Wないし線 5 0 Zを介して F I F Oバッファへのメッセージ格納要求が、線 L 5 1 Wないし線 5 1 Zを介してメッセージ格納要求を行う送信元の入力ポート番号が、線 L 5 2 Wないし線 5 2 Zを介して送り込んでいるメッセージの終了情報が、線 L 5 3 Wないし線 5 3 Zを介してメッセージの有効情報が、それぞれ送られてくる。

【 0 0 2 5 】

これらの情報を元に、F I F Oバッファ制御部 4 2は、取り込むべきメッセージを決定し、線 L 2 5を介して、セレクト 1 2 (図 1)を該当する入力バッファ 1 0 Aあるいは 1 0 Dからの入力に切り替え、線 L 2 6を介して F I F Oバッファ 1 3 (図 1)の書き込むべき場所を設定する。

10

【 0 0 2 6 】

さらに、F I F Oバッファ制御部 4 2は、F I F Oバッファ 1 3 (図 1)上に格納してあるメッセージを送りだすために、メッセージ内のメッセージ情報としての転送先プロセッサ番号を線 L 2 8を介して取り込み、線 L 5 5 Wないし線 L 5 5 Zを介して出力ポート対応競合調停部 4 1 Wないし 4 1 Zに対して、調停の依頼をし、その結果を線 L 5 4 Wないし線 L 5 4 Zにより受信し、それに応じて、線 L 2 7を介して F I F Oバッファ 1 3 (図 1)の読みだすべき場所を指示する。

【 0 0 2 7 】

F I F Oバッファ制御部の詳細な構成と動作を図 3を用いて説明する。図 3において、5 0は競合調停部、5 1, 5 2は状態を保持するレジスタ、5 3, 5 4はセレクト、5 5は AND回路、5 6は F I F Oバッファの書き込みポイント WP、5 7は F I F Oバッファの読みだしポイント RP、5 8は OR回路、5 9は F I F Oバッファの空き領域の計算部、6 0は要求信号生成部である。

20

【 0 0 2 8 】

まず、出力ポート対応競合調停部 4 1 Wないし 4 1 Z (図 2)からの要求により、F I F Oバッファ 1 3 (図 1)へのメッセージの取り込みについて説明する。出力ポート対応競合調停部 4 1 Wないし 4 1 Z (図 2)からのメッセージ情報として、線 L 5 0 Wないし線 5 0 Zを介して、F I F Oバッファへのメッセージ格納要求が、線 L 5 1 Wないし線 L 5 1 Zを介して、メッセージ格納要求を行う送信元の入力ポート番号が入力され、競合調停部 5 0に送られる。競合調停部

30

5 0では、要求のあった入力ポートの中から、一つを選択する。このとき、F I F Oバッファがメッセージ転送途中に、バッファの容量があふれてはいけないので、この競合調停時の入力ポート選択の基準のひとつとして、F I F Oバッファの空き領域情報を後述する F I F Oバッファの空き領域の計算部 5 9から読み込む。それ以外の競合調停部 5 0の動作は、通常の競合調停と同様であり、ここでは省略する。

【 0 0 2 9 】

競合調停の結果は、まず、新しく F I F Oバッファ 1 3 (図 1)に書き込むメッセージがあることを線 L 7 1を介して、レジスタ 5 1ないし 5 2に送り込む。レジスタ 5 1では、また線 L 7 0を介して、新しく選択された入力ポート番号を、線 L 7 1の指示により格納する。このレジスタ 5 1の情報は、線 L 2 5を介して、セレクト 1 2に送られる。もう一つのレジスタ 5 2は、現在、F I F Oバッファ 1 3 (図 1)への読み込みが行われている否かの状態を示す。

40

【 0 0 3 0 】

このレジスタ 5 2の S E T (F I F Oバッファへの読み込み開始)は、線 L 7 1を介して行われ、R E S E T (F I F Oバッファへの読み込み終了)は、出力ポート対応競合調停部 4 1 Wないし 4 1 Z (図 2)から線 L 5 2 Wないし L 5 2 Zを介して送り込まれたメッセージ情報の一つであるメッセージの終了指示からセレクト 5 3で選択された情報を用いて行う。

【 0 0 3 1 】

F I F Oバッファ 1 3 (図 1)の書き込む場所の設定は次のように行う。出力ポート対応

50

競合調停部 4 1 W ないし 4 1 Z (図 2) から線 L 5 3 W ないし L 5 3 Z を介して送り込まれたメッセージ情報の一つであるメッセージの有効情報からセクタ 5 4 で選択された情報を、レジスタ 5 2 の内容により現在 F I F O への読み込みを行っていることを確認しつつ、F I F O バッファの書き込みポイント

WP 5 6 に送る。WP 5 6 では、保持している WP の値を + 1 (メッセージの 1 サイクルでの転送単位を 1 とする) し、その結果を線 L 2 6 を介して、F I F O バッファ 1 3 (図 1) に送り込む。

【 0 0 3 2 】

次に、F I F O バッファ 1 3 (図 1) 上に格納してあるメッセージを出力バッファ 1 5 W ないし 1 5 Z (図 1) に送り出すための処理について説明する。まず、F I F O バッファ 1 3 (図 1) 上に格納してあるメッセージからメッセージ情報としての転送先プロセッサ番号を、線 L 2 8 を介して取り込み、要求信号生成部 6 0 に送り込む。要求信号生成部 6 0 は、F I F O バッファの空き領域の計算部 5 9 により F I F O バッファ 1 3 (図 1) 上にメッセージが格納されていることを確認し、線 L 5 5 W ないし線 L 5 5 Z を介して、出力ポート対応競合調停部

4 1 W ないし 4 1 Z (図 2) に対して、調停の依頼をする。出力ポート対応競合調停部 4 1 W ないし 4 1 Z (図 2) では、F I F O バッファ上のメッセージを各入力バッファ上のメッセージと同様に扱うが、競合時の優先度は、F I F O バッファの方を高くするべきである。

【 0 0 3 3 】

この調停結果は、線 L 5 4 W ないし線 L 5 4 Z により受信し、O R 回路 5 8 を介して、F I F O バッファの読みだしポイント R P 5 7 に送り込む。R P 5 7 は保持している R P の値を + 1 (メッセージの 1 サイクルでの転送単位を 1 とする) し、その結果を線 L 2 7 を介して、F I F O バッファ 1 3 (図 1) に送り込む。

【 0 0 3 4 】

F I F O バッファの空き領域の計算部 5 8 は、書き込みポイント WP 5 6 と読みだしポイント R P 5 7 を入力とし、F I F O バッファ 1 3 (図 1) 上に、メッセージが格納されているか、あるいは、次に格納すべきメッセージのためにどの程度の余裕があるかを計算し、その結果を、それぞれ要求信号生成部 6 0 と競合調停部 5 0 に送り込む。

【 0 0 3 5 】

本発明の F I F O バッファは、一つの F I F O 構成としたが、複数の F I F O を準備することも考えられる。また、さらに、この F I F O を出力ポート対応に設置することにより、一つの出力ポートが複数の入力ポートからの複数のメッセージを同時に受け付けることも可能になる。また、F I F O の容量は大きいほど効果的である。これらは、すべてスイッチを構成する L S I あるいは L S I 群のゲート規模により決定される。

【 0 0 3 6 】

【発明の効果】

本発明によれば、競合調停により、待つことを強いられたメッセージを一時的に格納する F I F O バッファに格納することにより、入力ポートを解放することができ、次のメッセージを競合調停の対象とすることができる。そのため、空いている出力ポートを少なくすることができ、スイッチの稼働率を高めることができる。

【図面の簡単な説明】

【図 1】スイッチの一実施例を示すブロック図。

【図 2】メッセージ調停制御部の一実施例を示すブロック図。

【図 3】F I F O バッファ制御部の一実施例を示すブロック図。

【図 4】並列計算機のプロセッサ間ネットワーク一実施例を示すブロック図。

【図 5】メッセージの構成の一実施例を示す説明図。

【符号の説明】

5 0 ... 競合調停部、5 1, 5 2 ... レジスタ、5 3, 5 4 ... セクタ、5 5 ... A N D 回路、5 6 ... 書き込みポイント WP、5 7 ... 読みだしポイント R P、5 8 ... O R 回路、5 9 ... F

10

20

30

40

50

I F Oバッファ空き領域計算部、60...要求信号生成部。

【 図 1 】

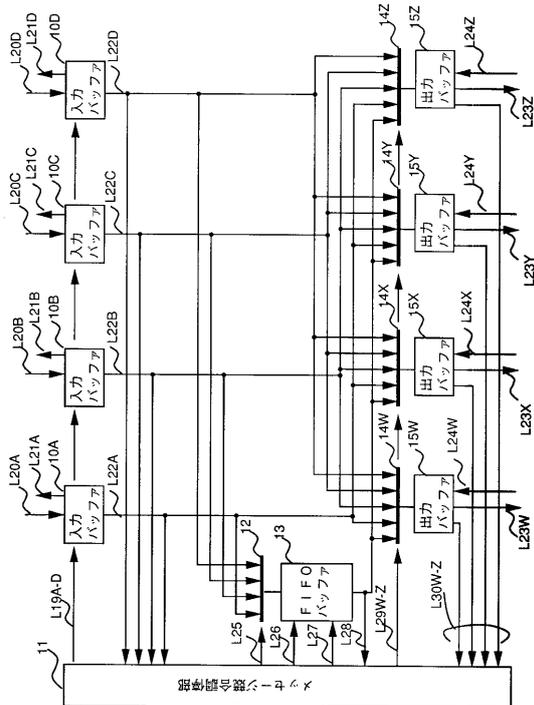


図 1

【 図 2 】

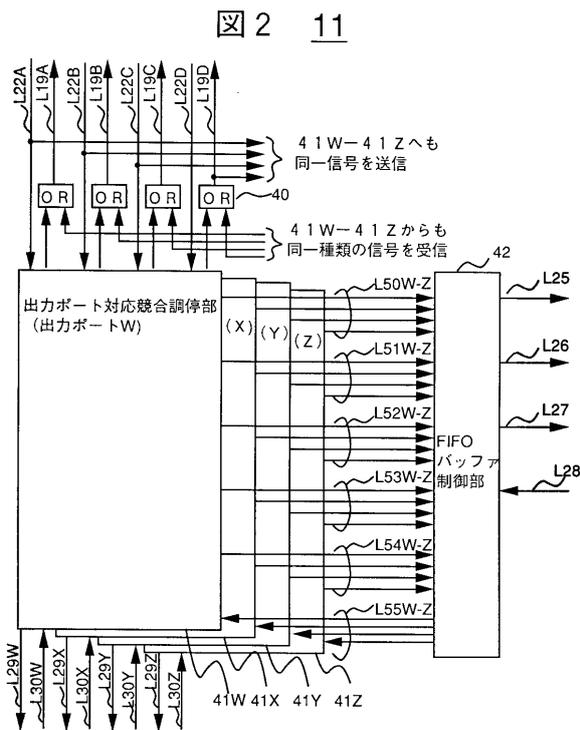
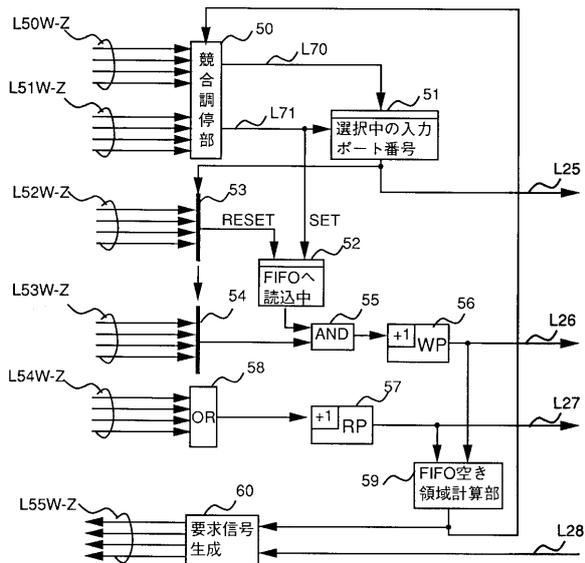


図 2 11

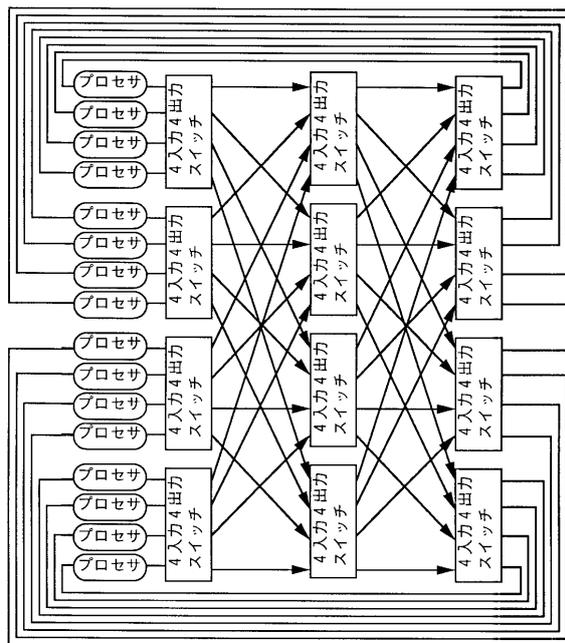
【 図 3 】

図 3 4 2



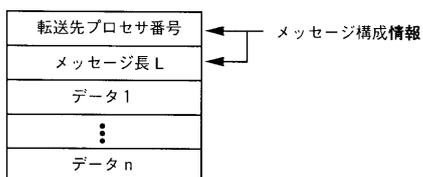
【 図 4 】

図 4



【 図 5 】

図 5



フロントページの続き

審査官 鳥居 稔

(56)参考文献 特公平04 - 053358 (JP, B2)

(58)調査した分野(Int.Cl.⁷, DB名)

G06F 15/16-173

G06F 13/38

H04L