

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5404947号
(P5404947)

(45) 発行日 平成26年2月5日(2014.2.5)

(24) 登録日 平成25年11月8日(2013.11.8)

(51) Int. Cl. F I
G06F 3/06 (2006.01) G O 6 F 3/06 3 O 1 J
G06F 13/10 (2006.01) G O 6 F 13/10 3 4 O B
 G O 6 F 3/06 3 O 1 X

請求項の数 16 (全 59 頁)

(21) 出願番号	特願2013-64768 (P2013-64768)	(73) 特許権者	000003078
(22) 出願日	平成25年3月26日(2013.3.26)		株式会社東芝
(62) 分割の表示	特願2010-252336 (P2010-252336) の分割		東京都港区芝浦一丁目1番1号
原出願日	平成22年11月10日(2010.11.10)	(74) 代理人	100108855 弁理士 蔵田 昌俊
(65) 公開番号	特開2013-145592 (P2013-145592A)	(74) 代理人	100109830 弁理士 福原 淑弘
(43) 公開日	平成25年7月25日(2013.7.25)	(74) 代理人	100088683 弁理士 中村 誠
審査請求日	平成25年3月26日(2013.3.26)	(74) 代理人	100103034 弁理士 野河 信久
早期審査対象出願		(74) 代理人	100095441 弁理士 白根 俊郎
		(74) 代理人	100075672 弁理士 峰 隆司

最終頁に続く

(54) 【発明の名称】 ストレージ装置

(57) 【特許請求の範囲】

【請求項1】

第1の入力ポートと、第1の不揮発性メモリと、前記第1の不揮発性メモリを制御する第1の制御部と、を含む第1のメモリノードと、

前記第1のメモリノードと接続される第2の入力ポートと、第2の不揮発性メモリと、前記第2の不揮発性メモリを制御する第2の制御部と、いずれか1つは前記第1の入力ポートと接続される複数の第1の出力ポートと、を含む第2のメモリノードと、

前記複数の第1の出力ポートのうちの1つに接続される第3の入力ポートを含み、前記第1乃至第2のメモリノードとは異なる複数の第3のメモリノードと、を具備するストレージ装置であって、

前記第2のメモリノードは、前記ストレージ装置内における前記第2のメモリノードの位置情報である第1の位置情報を記憶し、

前記第1の制御部は、送信先アドレスと送信元アドレスを少なくとも含むパケットを前記第2のメモリノードへ出力し、

前記第2の入力ポートは、前記第1のメモリノードから出力された前記パケットを受信し、

前記第2の制御部は、前記パケット内の前記送信先アドレスと前記送信元アドレス、前記第1の位置情報を少なくとも含む情報に基づいて、少なくとも前記複数の第3のメモリノード及び前記第1のメモリノードの中から前記パケットの出力先を決定するように構成されているストレージ装置。

【請求項 2】

第 1 の不揮発性メモリと、前記第 1 の不揮発性メモリを制御する第 1 の制御部と、を含む第 1 のメモリノードと、

前記第 1 のメモリノードと接続される第 1 の入力ポートと、第 2 の不揮発性メモリと、前記第 2 の不揮発性メモリを制御する第 2 の制御部と、複数の第 1 の出力ポートと、を含む第 2 のメモリノードと、

前記複数の第 1 の出力ポートのうちの 1 つに接続される第 2 の入力ポートを含み、前記第 1 乃至第 2 のメモリノードとは異なる複数の第 3 のメモリノードと、

前記第 1 乃至第 3 のメモリノードとは異なる 1 以上の第 4 のメモリノードと、を具備するストレージ装置であって、

前記第 1 の制御部は、送信先アドレスと送信元アドレスを少なくとも含むパケットを前記第 2 のメモリノードへ出力し、

前記第 1 の入力ポートは、前記第 1 のメモリノードから出力された前記パケットを受信し、

前記第 2 の制御部は、前記第 1 のメモリノードと前記複数の第 3 のメモリノードと前記 1 以上の第 4 のメモリノードのうち、前記ストレージ装置内における位置情報が前記送信元アドレスと一致する第 5 のメモリノードと、前記ストレージ装置内における位置情報が前記送信先アドレスと一致する第 6 のメモリノードとを結ぶ第 1 の直線で区切られる 2 つの領域のうち第 1 の領域に前記第 2 のメモリノードが属する場合に、前記複数の第 3 のメモリノードのうち第 1 の方向に存在する第 7 のメモリノードに前記パケットを出力し、前記第 1 の領域とは異なる第 2 の領域に前記第 2 のメモリノードが属する場合に、前記複数の第 3 のメモリノードのうち第 2 の方向に存在する第 8 のメモリノードに前記パケットを出力するストレージ装置。

【請求項 3】

前記第 2 の制御部は、前記ストレージ装置内における前記第 2 のメモリノードの位置情報が前記送信先アドレスと一致する場合に、前記パケットに含まれる情報を、前記第 2 の不揮発性メモリに記録する請求項 2 に記載のストレージ装置。

【請求項 4】

前記第 2 の制御部は、前記パケット内の前記送信先アドレス、前記送信元アドレス、前記第 1 の位置情報、及び前記複数の第 1 の出力ポートにおける前記パケットとは別のパケットの出力に使用されているかを示す占有情報を少なくとも含む情報に基づいて、少なくとも前記複数の第 3 のメモリノード及び前記第 1 のメモリノードの中から前記パケットの出力先を決定する請求項 1 に記載のストレージ装置。

【請求項 5】

前記第 2 の制御部は、前記複数の第 1 の出力ポートにおける前記パケットとは別のパケットの出力に使用されているかを示す占有情報を少なくとも含む情報に基づいて、前記パケットを、前記複数の第 3 のメモリノードのいずれかに出力する請求項 2 に記載のストレージ装置。

【請求項 6】

前記第 2 の制御部は、前記第 1 の領域に前記第 2 のメモリノードが属する場合に、前記複数の第 1 の出力ポートのうち前記第 7 のメモリノードと接続される出力ポートが前記パケットとは別のパケットの出力用に使用されている場合、前記第 8 のメモリノードに前記パケットを出力し、前記第 2 の領域に前記第 2 のメモリノードが属する場合に、前記複数の第 1 の出力ポートのうち前記第 8 のメモリノードと接続される出力ポートが前記パケットとは別のパケットの出力用に使用されている場合、前記第 7 のメモリノードに前記パケットを出力する請求項 2 に記載のストレージ装置。

【請求項 7】

前記第 2 の制御部は、前記第 1 の領域において、前記第 1 の直線と直行し、かつ前記第 6 のメモリノードを通過する第 2 の直線で区切られる領域であって、前記第 5 のメモリノードに近い領域を第 3 の領域、他方を第 4 の領域とし、前記第 2 の領域において、前記第

10

20

30

40

50

2の直線で区切られる領域のうち、前記第5のメモリノードに近い領域を第5の領域、他方を第6の領域とした場合に、前記第3の領域に前記第2のメモリノードが属する場合、前記第7のメモリノードに前記パケットを出力し、前記第4の領域に前記第2のメモリノードが属する場合、前記複数の第3のメモリノードのうち前記第2の方向と逆の方向に存在するメモリノードに前記パケットを出力し、

前記第5の領域に前記第2のメモリノードが属する場合、前記第8のメモリノードに前記パケットを出力し、前記第6の領域に前記第2のメモリノードが属する場合、前記複数の第3のメモリノードのうち前記第1の方向と逆の方向に存在するメモリノードに前記パケットを出力する請求項6に記載のストレージ装置。

【請求項8】

前記第1乃至第3のメモリノードとは異なる1以上の第4のメモリノードをさらに具備し、

前記第2の制御部は、前記第1のメモリノードと前記複数の第3のメモリノードと前記1以上の第4のメモリノードのうち、前記ストレージ装置内における位置情報が前記送信先アドレスと一致するメモリノードへの転送経路が最も短いメモリノードに前記パケットを出力する請求項1に記載のストレージ装置。

【請求項9】

前記第2のメモリノードは前記第1の入力ポートの他にコンピュータと接続可能な1以上の入力ポートを含む請求項1に記載のストレージ装置。

【請求項10】

前記第1乃至第3のメモリノードとは異なる1以上の第4のメモリノードをさらに具備し、

前記複数の第3のメモリノードのうち前記第2のメモリノードに隣接しない第1の非隣接メモリノードと、前記1以上の第4のメモリノードのうち前記第1の非隣接メモリノードと接続され前記第1の非隣接メモリノードに隣接しない1以上の第2の非隣接メモリノードと、前記第2のメモリノードは夫々メモリノード同士の相対的な物理位置によって決定する付加アドレスを持ち、

前記第2の制御部は、前記第2のメモリノードの付加アドレス、前記第1の非隣接メモリノードの付加アドレス、前記第2の非隣接メモリノードの付加アドレスの少なくともいずれか一つを含む情報に基づいて、少なくとも前記複数の第3のメモリノード及び前記第1のメモリノードの中から前記パケットの出力先を決定する請求項1に記載のストレージ装置。

【請求項11】

前記第1乃至第3のメモリノードとは異なる1以上の第4のメモリノードと、前記第1のメモリノードと前記第2のメモリノードと前記複数の第3のメモリノードと前記1以上の第4のメモリノードのうちいずれか1つ以上と接続されるリレーと、をさらに具備し、

前記リレーは、前記リレーと接続されるメモリノードの中で、前記第1のメモリノードと前記第2のメモリノードと前記複数の第3のメモリノードと前記1以上の第4のメモリノードのうち前記ストレージ装置内における位置情報が前記送信先アドレスと一致するメモリノードへの転送経路が最も短いメモリノードに前記パケット送信することを特徴とする請求項1に記載のストレージ装置。

【請求項12】

前記第1のメモリノードと前記第2のメモリノードと前記複数の第3のメモリノードと前記1以上の第4のメモリノードのうちいずれか1つ以上と接続されるリレーと、をさらに具備し、

前記リレーは、前記リレーと接続されるメモリノードの中で、前記第6のメモリノードへの転送経路が最も短いメモリノードに前記パケット送信することを特徴とする請求項2に記載のストレージ装置。

【請求項13】

10

20

30

40

50

前記第 1 乃至第 3 のメモリノードとは異なる 1 以上の第 4 のメモリノードをさらに具備し、

前記第 1 のメモリノードと前記第 2 のメモリノードと前記複数の第 3 のメモリノードと前記 1 以上の第 4 のメモリノードのうち、他のメモリノードと接続されない入力ポートを少なくとも 1 以上有するメモリノードに接続可能なコンピュータの総数を N_c とし、前記第 1 のメモリノードと前記第 2 のメモリノードと前記複数の第 3 のメモリノードと前記 1 以上の第 4 のメモリノードの総数を N_{node} とした場合、 $N_c < N_{node}$ を満たすことを特徴とする請求項 1 に記載のストレージ装置。

【請求項 14】

前記第 1 のメモリノードと前記第 2 のメモリノードと前記複数の第 3 のメモリノードと前記 1 以上の第 4 のメモリノードのうち、他のメモリノードと接続されない入力ポートを少なくとも 1 以上有するメモリノードに接続可能なコンピュータの総数を N_c とし、前記第 1 のメモリノードと前記第 2 のメモリノードと前記複数の第 3 のメモリノードと前記 1 以上の第 4 のメモリノードの総数を N_{node} とした場合、 $N_c < N_{node}$ を満たすことを特徴とする請求項 2 に記載のストレージ装置。

【請求項 15】

$N_c < 0.1 \times N_{node}$ の関係式を満たすことを特徴とする請求項 13 または 14 に記載のストレージ装置。

【請求項 16】

前記第 2 のメモリノードは前記パケットに含まれる key をアドレスに変換するアドレス変換器をさらに具備し、

前記第 2 の制御部は、key-value 型データの各レコードを保有し、前記送信先アドレスを前記アドレス変換器により変換された前記アドレスとし、前記パケットに value を含めて出力する請求項 1 または 2 に記載のストレージ装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、転送機能を有するメモリノードを相互に接続したストレージ装置及びデータ処理方法に関わり、例えばストレージ装置におけるデータパケットの転送制御方式に関わる。

【背景技術】

【0002】

容量を容易に拡張できるストレージ装置として、転送機能を有するメモリノードを相互に接続したストレージ装置が考えられる。各メモリノードは、自身のメモリノード宛のデータパケットを受信した場合は読み出しもしくは書き込みなどの所定の処理を行う。一方、各メモリノードは、自身のメモリノード宛でないパケットを受信した場合は、受信パケットを適切な他のメモリノードに転送する。各メモリノードによって適切な転送が繰り返されることにより、データパケットは目的のメモリノードに到達できる。

【0003】

各メモリノードは、メモリと、転送機能を有するコントローラ、複数のポートを有している。各メモリノードは、パケットの転送先を示した経路指定表（ルーティングテーブル）を維持・管理し、それに従ってパケットの転送を行う。経路指定表を管理すれば、物理的な位置に関わらず任意の論理的なパケット転送ネットワークを構築することが可能である。

【0004】

しかし、容量拡張のため新たなメモリノードを追加する場合、もしくは故障等の理由により既存のメモリノードを除去する場合には、各メモリノードの経路指定表を更新する必要があり、その手続きは煩雑である。経路指定表の維持・管理コストは、特にメモリノード数が大規模になった場合に膨大になり、それが容量の拡張性に制限を課すことになる。

【0005】

10

20

30

40

50

また、転送機能を有するメモリノードを相互に接続したストレージ装置において、複数のデータを複数のメモリノードに書き込み・読み出しを行う場合において、複数のデータを同一配線で同時に通信させることは一般に困難であるため、データの転送待機が生じやすい。このようなデータの転送待機は、データの書き込み・読み出しに必要な時間の増加を招く。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】米国特許出願公開第2009/0216924号明細書

【発明の概要】

10

【発明が解決しようとする課題】

【0007】

メモリノードが経路指定表を管理する必要がなく、効率的にパケットを転送することができるストレージ装置、及びデータ処理方法を提供する。

【課題を解決するための手段】

【0008】

一実施態様のストレージ装置は、第1の入力ポートと、第1の不揮発性メモリと、前記第1の不揮発性メモリを制御する第1の制御部と、を含む第1のメモリノードと、前記第1のメモリノードと接続される第2の入力ポートと、第2の不揮発性メモリと、前記第2の不揮発性メモリを制御する第2の制御部と、いずれか1つは前記第1の入力ポートと接続される複数の第1の出力ポートと、を含む第2のメモリノードと、前記複数の第1の出力ポートのうちの1つに接続される第3の入力ポートを含み、前記第1乃至第2のメモリノードとは異なる複数の第3のメモリノードとを具備するストレージ装置であって、前記第2のメモリノードは、前記ストレージ装置内における前記第2のメモリノードの位置情報である第1の位置情報を記憶し、前記第1の制御部は、送信先アドレスと送信元アドレスを少なくとも含むパケットを前記第2のメモリノードへ出力し、前記第2の入力ポートは、前記第1のメモリノードから出力された前記パケットを受信し、前記第2の制御部は、前記パケット内の前記送信先アドレスと前記送信元アドレス、前記第1の位置情報を少なくとも含む情報に基づいて、少なくとも前記複数の第3のメモリノード及び前記第1のメモリノードの中から前記パケットの出力先を決定するように構成されている。

20

30

【図面の簡単な説明】

【0009】

【図1】第1実施形態のストレージ装置の構成を示す図である。

【図2】第1実施形態におけるメモリノードの構成を示す図である。

【図3】第1実施形態におけるメモリノードの配置例を示す図である。

【図4】第1実施形態のストレージ装置における転送アルゴリズム1を示す図である。

【図5A】第1実施形態における転送アルゴリズム1によるパケットの転送過程を示す図である。

【図5B】第1実施形態における転送アルゴリズム1によるパケットの転送過程を示す図である。

40

【図5C】第1実施形態における転送アルゴリズム1によるパケットの転送過程のフロー図である。

【図5D】第1実施形態における転送アルゴリズム1によるパケットの転送過程の具体例を示す図である。

【図6】第1実施形態のストレージ装置を含むストレージシステムの構成を示す図である。

【図7】第1実施形態のストレージシステムにおける書き込み動作を示す図である。

【図8】第1実施形態のストレージシステムにおける読み出し動作を示す図である。

【図9】第1実施形態のストレージ装置におけるメモリノードの自動アドレス取得方式を示す図である。

50

- 【図10】第2実施形態のストレージ装置における転送アルゴリズム2を示す図である。
- 【図11A】第2実施形態における転送アルゴリズム2によるパケットの転送過程を示す図である。
- 【図11B】第2実施形態における転送アルゴリズム2によるパケットの転送過程を示す図である。
- 【図11C】第2実施形態における転送アルゴリズム2によるパケットの転送過程のフロー図である。
- 【図11D】第2実施形態における転送アルゴリズム2によるパケットの転送過程の具体例を示す図である。
- 【図12A】転送アルゴリズム1によるパケットの転送過程の具体例を示す図である。 10
- 【図12B】転送アルゴリズム2によるパケットの転送過程の具体例を示す図である。
- 【図13】第3実施形態のストレージ装置における転送アルゴリズム3を示す図である。
- 【図14A】第3実施形態における転送アルゴリズム3によるパケットの転送過程を示す図である。
- 【図14B】第3実施形態における転送アルゴリズム3によるパケットの転送過程を示す図である。
- 【図14C】第3実施形態における転送アルゴリズム3によるパケットの転送過程のフロー図である。
- 【図14D】第3実施形態における転送アルゴリズム3によるパケットの転送過程の具体例を示す図である。 20
- 【図15】第3実施形態のストレージ装置における転送アルゴリズム4を示す図である。
- 【図16】第3実施形態のストレージ装置における転送アルゴリズム5を示す図である。
- 【図17】第4実施形態のストレージシステムの構成を示す図である。
- 【図18】第4実施形態のストレージシステムにおけるバイパス転送発生率と渋滞発生率を示す図である。
- 【図19】第5実施形態のストレージ装置の構成を示す図である。
- 【図20】第6実施形態のストレージシステムの構成を示す図である。
- 【図21】第6実施形態のストレージシステムの他の構成例を示す図である。
- 【図22】第7実施形態のストレージシステムの構成を示す図である。
- 【図23】第7実施形態のストレージシステムに対する比較例を示す図である。 30
- 【図24】第7実施形態におけるパケットのヘッダ部に記録されたアドレス情報を示す図である。
- 【図25】第7実施形態のストレージシステムにおける書き込み動作を示す図である。
- 【図26】第7実施形態のストレージシステムの他の構成例を示す図である。
- 【図27】第8実施形態のストレージシステムの構成を示す図である。
- 【図28】第8実施形態のストレージシステムに対する比較例を示す図である。
- 【図29】第8実施形態におけるパケットのヘッダ部に記録されたアドレス情報を示す図である。
- 【図30】第8実施形態のストレージシステムにおける書き込み動作を示す図である。
- 【図31】第8実施形態のストレージシステムの他の構成例を示す図である。 40
- 【図32】第8実施形態のストレージシステムの他の構成例を示す図である。
- 【図33A】第9実施形態のストレージシステムの構成を示す図である。
- 【図33B】ストレージシステムにおいて転送待機が発生する読み出し動作を示す図である。
- 【図33C】ストレージシステムにおいて転送待機が発生する読み出し動作を示す図である。
- 【図33D】ストレージシステムにおいて転送待機が発生する読み出し動作を示す図である。
- 【図34A】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作を示す図である。 50

【図34B】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作を示す図である。

【図34C】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作を示す図である。

【図34D】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作を示す図である。

【図34E】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作を示す図である。

【図35A】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

10

【図35B】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

【図36A】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

【図36B】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

【図36C】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

【図36D】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

20

【図36E】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

【図36F】第9実施形態のストレージシステムにおいて転送待機の発生を回避する読み出し動作の他の例を示す図である。

【図37A】第10実施形態のストレージシステムの構成を示す図である。

【図37B】ストレージシステムにおいて転送待機が発生する書き込み動作を示す図である。

【図37C】ストレージシステムにおいて転送待機が発生する書き込み動作を示す図である。

【図37D】ストレージシステムにおいて転送待機が発生する書き込み動作を示す図である。

30

【図37E】ストレージシステムにおいて転送待機が発生する書き込み動作を示す図である。

【図37F】ストレージシステムにおいて転送待機が発生する書き込み動作を示す図である。

【図37G】ストレージシステムにおいて転送待機が発生する書き込み動作を示す図である。

【図38A】第10実施形態のストレージシステムにおいて転送待機の発生を回避する書き込み動作を示す図である。

【図38B】第10実施形態のストレージシステムにおいて転送待機の発生を回避する書き込み動作を示す図である。

40

【図38C】第10実施形態のストレージシステムにおいて転送待機の発生を回避する書き込み動作を示す図である。

【図38D】第10実施形態のストレージシステムにおいて転送待機の発生を回避する書き込み動作を示す図である。

【図38E】第10実施形態のストレージシステムにおいて転送待機の発生を回避する書き込み動作を示す図である。

【発明を実施するための形態】

【0010】

以下、図面を参照して実施形態について説明する。なお、以下の説明において、同一の

50

機能及び構成を有する構成要素については、同一符号を付し、重複説明は必要な場合にのみ行う。

【 0 0 1 1 】

[第 1 実施形態]

第 1 実施形態は、転送機能を有する複数のメモリノードを相互に接続したストレージ装置を備え、各メモリノードが効率的にデータパケットを転送する転送方式を有する。

【 0 0 1 2 】

[1] ストレージ装置の構成

図 1 は、第 1 実施形態のストレージ装置の構成を示す図であり、複数のメモリノードの物理的な配置法と、それに対するアドレスの割り当て法の一例を示す。

10

【 0 0 1 3 】

図示するように、ストレージ装置 1 0 はデータ転送機能を有するメモリノード 1 1 を複数備える。各々のメモリノード 1 1 は、正方格子の格子点に置かれる。格子点の座標である格子座標 (x , y) に位置するメモリノードの論理アドレスは、位置座標と一致して (x , y) とする。すなわち、メモリノード 1 1 の論理アドレスは、物理アドレス (格子座標 (x , y)) に一致する。

【 0 0 1 4 】

各メモリノード 1 1 は、4 つの入力ポート 1 2 と 4 つの出力ポート 1 3 を持つ。各メモリノードは、隣接する 4 つのメモリノードと入力ポート 1 2 及び出力ポート 1 3 を介して相互に接続される。具体的には、対向する 2 つの隣接したメモリノードは、互いの入力ポート 1 2 と出力ポート 1 3 とが接続される。

20

【 0 0 1 5 】

図 2 に、各メモリノード 1 1 の構成を示す。メモリノード 1 1 は、入力ポート 1 2、入力ポートバッファ 1 2 A、出力ポート 1 3、出力ポートバッファ 1 3 A、セクタ 1 4、パケットコントローラ 1 5、メモリ 1 6、メモリコントローラ 1 7、MPU 1 8、及びローカルバス 1 9 を有する。

【 0 0 1 6 】

入力ポート 1 2 に入力されたパケットは、入力ポートバッファ 1 2 A に一時的に記憶される。セクタ 1 4 には入力ポートバッファ 1 2 A からパケットが入力され、またパケットコントローラ 1 5 から制御信号が入力される。セクタ 1 4 は、制御信号に従い、入力されたパケットからいずれかのパケットを選択して出力ポートバッファ 1 3 A に出力する。出力ポートバッファ 1 3 A は、セクタ 1 4 から出力されたパケットを一時的に記憶すると共に、出力ポート 1 3 に出力する。パケットコントローラ 1 5 は、セクタ 1 4 の出力を制御する。パケットとは、送信先アドレス及び送信元アドレスを少なくとも含むヘッダ部と、データ部とからなる転送データの単位である。

30

【 0 0 1 7 】

メモリ 1 6 は、データを記憶する複数のメモリセルを有する。メモリ 1 6 は、例えば NAND 型フラッシュメモリ等からなる。メモリコントローラ 1 7 は、メモリ 1 6 への書き込み、読み出し、及び消去の動作を制御する。MPU 1 8 は、メモリノード内で必要な演算処理を行う。ローカルバス 1 9 は、入力ポートバッファ 1 2 A、パケットコントローラ 1 5、メモリコントローラ 1 7、及び MPU 1 8 間を相互に接続し、これらの間の信号伝送を行う。

40

【 0 0 1 8 】

メモリノード 1 1 が受信したパケットは、入力ポート 1 2 を介して入力ポートバッファ 1 2 A に格納される。パケットコントローラ 1 5 は、パケットが含む送信先 (宛先) アドレスとそのメモリノード自身 (以下、自ノード) のアドレスの 2 つの情報に基づいて、受信したパケットが自ノード宛であるか否かを判断する。

【 0 0 1 9 】

もし、自ノード宛であれば、パケットコントローラ 1 5 は、自ノードのメモリ 1 6 に対する書き込み、読み出し、もしくはその他所定の処理を行う。もし、自ノード宛でなけれ

50

ば、パケットの送信先アドレスと自ノードのアドレスの2つの情報に基づいて、パケットを転送する隣接するメモリノードを決定し、セクタ14により対応する出力ポートバッファ13Aにパケットを出力する。

【0020】

図1に示したストレージ装置では、正方格子の格子点にメモリノードが配置されたが、本実施形態はこの例に限定されるものではない。これを図3を参照して、例を挙げて説明する。

【0021】

図3(a)に示す構成が図1に示した例に相当する。より一般的には、本実施形態において、各メモリノードは格子点に配置される。格子とは、幾何学においては平面上の点のうち、 x 座標、 y 座標がともに整数である点を言う。ここでは、 x 方向と y 方向の単位ベクトル、 e_x と e_y の長さが異なる場合、即ち x 方向と y 方向の繰り返し周期が異なる場合も含む。この例を図3(b)に示す。

10

【0022】

また、 x 方向と y 方向の単位ベクトルが直交しない場合、即ち x 軸と y 軸が直交しない場合も含む。この例を図3(c)に示す。また、メモリノードの相互接続数は4とは限らない。相互接続数が6の例を図3(d)に示す。

【0023】

本実施形態においては、格子の定義に関わらず、いずれの場合においても、格子座標(x, y)に位置するメモリノードの論理アドレスは、位置座標と一致して(x, y)とする。すなわち、メモリノードの論理アドレスは、物理アドレス(格子座標(x, y))に一致する。

20

【0024】

また、より一般的には、本実施形態においては、各格子点に2以上のメモリノードの組が配置される場合を含む。各格子点に2つのメモリノードが配置される例を、図3(e)に示す。この例では、格子点(1, 0)などに2つのメモリノードが配置される。これを、(1, 0, 0)と(1, 0, 1)とする。即ち、1つのメモリノードの座標は(x, y, z)の3つの整数の組によって表される。この場合においても、ノード座標(x, y, z)に位置するメモリノードの論理アドレスは、位置座標と一致して(x, y, z)とする。すなわち、メモリノードの論理アドレスは、物理アドレス(格子座標(x, y, z))に一致する。さらに、格子点が配置される平面が三次元的に曲げられたり、もしくは折りたたまれた場合も含む。

30

【0025】

隣接したメモリノードとは、図3(a)~図3(e)において以下のような位置関係にあるメモリノードを指す。図3(a)~図3(c)において、例えば自ノードが座標(1, 1)に存在する場合、自ノードに隣接したメモリノードは座標(0, 1)、(1, 2)、(2, 1)、(1, 0)に存在する4つのメモリノードを指す。また図3(d)において、例えば自ノードが座標(1, 1)に存在する場合、自ノードに隣接したメモリノードは座標(0, 1)、(0, 2)、(1, 2)、(2, 1)、(2, 0)、(1, 0)に存在する6つのメモリノードを指す。さらに、図3(e)において、例えば自ノードが座標(1, 0, 1)に存在する場合、自ノードに隣接したメモリノードとは座標(0, 1, 0)、(1, 1, 0)、(1, 0, 0)に存在する3つのメモリノードを指す。

40

【0026】

[2] ストレージ装置の転送アルゴリズム1

転送アルゴリズム1では、パケットが含む送信先アドレス、及び自ノードのアドレスの2つの情報に基づいて、転送先のメモリノードを決定する。

【0027】

図4は、第1実施形態のストレージ装置における転送アルゴリズム1を示す図である。図4を参照して、パケットの送信先アドレスと自ノードのアドレスの2つの情報に基づいて、転送先の隣接ノードを決定する方法の一例を示す。これを転送アルゴリズム1とする

50

。

【 0 0 2 8 】

パケットを受け取ったメモリノードは、パケットの送信先ノード(to)と自ノード(PP: Present position)との距離が最も小さくなる隣接するメモリノードにパケットを転送する。

。

【 0 0 2 9 】

図 5 A ~ 図 5 D を参照して、転送アルゴリズム 1 に従うパケットの転送過程の具体例を示す。

【 0 0 3 0 】

図 5 A に示すように、送信先ノード(to)のアドレスを (x_{to}, y_{to}) 、自ノード(PP)のアドレスを (x_{pp}, y_{pp}) とし、さらに $dx = x_{to} - x_{pp}$ 、 $dy = y_{to} - y_{pp}$ とする。転送先を示す方向として、図 5 B に示すように、 y が増加する方向をN(North)、 x が増加する方向をE(East)、 y が減少する方向をS(South)、 x が減少する方向をW(West)とする。

10

【 0 0 3 1 】

a に対してその絶対値をあらわす記号を $|a|$ として、 $|dx| > |dy|$ ならば x 方向へ、 $|dx| < |dy|$ ならば y 方向へ進む。 x 方向へ進む場合は、 $dx > 0$ ならEへ、 $dx < 0$ ならWに転送する。同様に、 y 方向へ進む場合は、 $dy > 0$ ならNへ、 $dy < 0$ ならSに転送する。

【 0 0 3 2 】

図 5 C に、ストレージ装置における転送アルゴリズム 1 のフローを示す。転送アルゴリズム 1 はパケットコントローラ 15 に記憶されており、パケットコントローラ 15 により実行される。

20

【 0 0 3 3 】

まず、パケットコントローラ 15 は、 $dx = x_{to} - x_{pp}$ 、 $dy = y_{to} - y_{pp}$ を算出する(ステップ S 1)。続いて、パケットコントローラ 15 は、 dx が0であるか否かを判定する(ステップ S 2)。 dx が0であるとき、 $y_{to} > y_{pp}$ が成り立つか否かを判定する(ステップ S 3)。 $y_{to} > y_{pp}$ が成り立つとき、パケットをNへ転送する(ステップ S 4)。一方、 $y_{to} > y_{pp}$ が成り立たないとき、パケットをSへ転送する(ステップ S 5)。

【 0 0 3 4 】

次に、ステップ S 2 において dx が0でないとき、パケットコントローラ 15 は、 dy が0であるか否かを判定する(ステップ S 6)。 dy が0であるとき、 $x_{to} > x_{pp}$ が成り立つか否かを判定する(ステップ S 7)。 $x_{to} > x_{pp}$ が成り立つとき、パケットをEへ転送する(ステップ S 8)。一方、 $x_{to} > x_{pp}$ が成り立たないとき、パケットをWへ転送する(ステップ S 9)。

30

【 0 0 3 5 】

次に、ステップ S 6 において dy が0でないとき、すなわち dx と dy が0でないとき、パケットコントローラ 15 は、 $dx > 0$ かつ $dy > 0$ が成り立つか否かを判定する(ステップ S 10)。 $dx > 0$ かつ $dy > 0$ が成り立つとき、 $dx > dy$ が成り立つか否かを判定する(ステップ S 11)。 $dx > dy$ が成り立つとき、パケットをEへ転送する(ステップ S 12)。一方、 $dx > dy$ が成り立たないとき、パケットをNへ転送する(ステップ S 13)。

40

【 0 0 3 6 】

次に、ステップ S 10 において $dx > 0$ かつ $dy > 0$ が成り立たないとき、パケットコントローラ 15 は、 $dx < 0$ かつ $dy > 0$ が成り立つか否かを判定する(ステップ S 14)。 $dx < 0$ かつ $dy > 0$ が成り立つとき、 $(-1) \cdot dx > dy$ が成り立つか否かを判定する(ステップ S 15)。 $(-1) \cdot dx > dy$ が成り立つとき、パケットをWへ転送する(ステップ S 16)。一方、 $(-1) \cdot dx > dy$ が成り立たないとき、パケットをNへ転送する(ステップ S 17)。

【 0 0 3 7 】

次に、ステップ S 14 において $dx < 0$ かつ $dy > 0$ が成り立たないとき、パケットコ

50

ントローラ15は、 $dx < 0$ かつ $dy < 0$ が成り立つか否かを判定する(ステップS18)。 $dx < 0$ かつ $dy < 0$ が成り立つとき、 $dx > dy$ が成り立つか否かを判定する(ステップS19)。 $dx > dy$ が成り立つとき、パケットをSへ転送する(ステップS20)。一方、 $dx > dy$ が成り立たないとき、パケットをWへ転送する(ステップS21)。

【0038】

次に、ステップS18において $dx < 0$ かつ $dy < 0$ が成り立たないとき、パケットコントローラ15は、 $dx > (-1) \cdot dy$ が成り立つか否かを判定する(ステップS22)。 $dx > (-1) \cdot dy$ が成り立つとき、パケットをEへ転送する(ステップS23)。一方、 $dx > (-1) \cdot dy$ が成り立たないとき、パケットをSへ転送する(ステップS24)。

10

【0039】

以上の処理により、メモリノードに入力されたパケットは、送信先ノード(to)と自ノード(PP)との距離が最も小さくなるような隣接するメモリノードに転送される。

【0040】

図5Dに、パケットの転送過程の具体例を示す。転送アルゴリズム1においては $|dx| = |dy|$ となるまでは $|dx| > |dy|$ ならばx方向へ、 $|dx| < |dy|$ ならばy方向へ連続して進み、 $|dx| = |dy|$ となった後はx方向とy方向へ交互に進む。

【0041】

例えば、ケース1では、 $dx = 0$ かつ $dy > 0$ であるため、送信先ノード(to)に到達するまでNへ進む。ケース4では、 $dx > 0$ かつ $dy > 0$ であり、 $dx < dy$ であるため、 $dx = dy$ となるまでNへ連続して進み、 $dx = dy$ となった後はEとNへ交互に進む。

20

【0042】

[3] ストレージ装置を含むストレージシステム

図6は、第1実施形態のストレージ装置を含むストレージシステムの構成を示す図である。

【0043】

ストレージシステム20は、ストレージ装置をクライアントが使用するためのシステムであり、以下のような構成を備える。

【0044】

ストレージ装置10はゲートウェイサーバを介してクライアントと接続される。ストレージ装置10内部の通信規格とゲートウェイサーバ21A、21Bの通信規格が異なる場合は、両者の間にアダプタ22A、22Bをそれぞれ設置しても良い。

30

【0045】

詳述すると、ストレージ装置10の外周部に配置されたメモリノード(1,4)は、アダプタ22A、ゲートウェイサーバ21Aを介してクライアント31Aと接続される。同様に、メモリノード(1,1)は、アダプタ22B、ゲートウェイサーバ21Bを介してクライアント31B1、31B2と接続される。なお、“メモリノード(x,y)”は、アドレス(x,y)のメモリノードであることを表す。以降も同様である。

【0046】

ゲートウェイサーバ21A、21Bは、例えばコンピュータから構成され、ストレージ装置10内部のメモリノード11と同じ規約に基づいてアドレスを持つ。図6では、ゲートウェイサーバ21Aはアドレス(0,4)を持ち、ゲートウェイサーバ21Bはアドレス(0,1)を持つ。

40

【0047】

次に、ストレージシステムにおける書き込み動作を説明する。

【0048】

図7は、ストレージ装置10にクライアントがファイルを書き込む手続きを示す。ここでは、クライアント31Aがストレージ装置10にファイルを書き込む場合を述べる。

【0049】

50

クライアント31Aは、ファイルとファイルIDをゲートウェイサーバ21Aに送信する((1)参照)。ファイルIDは、ファイルを一意に特定できる識別子である。ファイルIDとしては、¥strage_system ¥home ¥cliantA ¥file1.txtといった所定のファイルシステムにおけるフルパスのファイル名を使うことができる。

【0050】

次に、ゲートウェイサーバ21Aは、ファイルを規定サイズのデータパケットに分割し、各パケットにパケットIDを割り振る。続いて、ファイルIDと、分割されたパケットのパケットIDをファイルテーブルに書き込む((2)参照)。パケットIDは、パケットを一意に特定できる識別子である。パケットIDとしては、・・・¥file1.txt~1、・・・¥file1.txt~2というように、「ファイルID+連続する番号」を割り当てることができる。

10

【0051】

次に、ゲートウェイサーバ21Aは、パケットIDの情報に基づき、そのパケットを書き込むメモリノードのアドレス(以下、書き込みノードアドレス)を決定する((3)参照)。この際、大規模分散ファイルシステムで使用されるコンシステントハッシング(文献1参照)と呼ばれるノード決定手法を使っても良い。コンシステントハッシングは、ノードアドレスのハッシュ値、およびパケットIDのハッシュ値の両方を使って、書き込みアドレスを決定する点に特徴がある。

【0052】

[文献1]:丸山不二夫/首藤一幸編、「雲の世界の向こうをつかむクラウドの技術」、株式会社アスキー・メディアワークス、2009年11月6日、p.88、ISBN978-4-04-868064-6

20

次に、ゲートウェイサーバ21Aは、書き込みアドレスを送信先アドレス、ゲートウェイサーバ21Aのアドレスを送信元アドレスとして、書き込みパケットをストレージ装置10に送信する((4)、(5)参照)。

【0053】

メモリノード(1,4)に送信されたパケットは、ストレージ装置10内で転送アルゴリズム1に従って適切な転送が繰り返されることにより、送信先アドレスのメモリノードに到達する((6)参照)。目的のメモリノードでは、受け取ったパケットを自ノードのメモリ16に書き込む((7)参照)。その後、書き込み完了パケットをゲートウェイサーバ21Aへ返信する((8)参照)。

30

【0054】

次に、ストレージシステムにおける読み出し動作を説明する。

【0055】

図8は、ストレージ装置10からクライアントがファイルを読み出す手続きを示す。ここでは、クライアント31Aがストレージ装置10からファイルを読み出す場合を述べる。

【0056】

クライアント31Aは、読み出し要求(ファイルID)をゲートウェイサーバ21Aに送信する((1)参照)。

【0057】

ゲートウェイサーバ21Aは、ファイルテーブルからファイルIDに対応するパケットIDを取得する((2)参照)。続いて、ゲートウェイサーバ21Aは、パケットIDの情報に基づき、読み出しを行うメモリノードのアドレス(以下、読み出しノードアドレス)を決定する((3)参照)。この際、大規模分散ファイルシステムで使用されるコンシステントハッシングと呼ばれるノード決定手法を使っても良い。

40

【0058】

次に、ゲートウェイサーバ21Aは、読み出しノードアドレスを送信先アドレス、ゲートウェイサーバ21Aのアドレスを送信元アドレスとして、読み出しパケットをストレージ装置10に送信する((4)、(5)参照)。

【0059】

50

メモリノード(1, 4)に送信されたパケットは、ストレージ装置10内で転送アルゴリズム1に従って適切な転送が繰り返されることにより、送信先アドレスのメモリノードに到達する(6)参照)。目的のメモリノードでは、読み出しパケットに従い、自ノードのメモリ16から目的のデータを読み出す(7)参照)。その後、読み出したデータをデータパケットとしてゲートウェイサーバ21Aに返信する(8)参照)。

【0060】

図6に示したシステムにおいては、書き込みデータのランダム化処理(Rand処理)、書き込みデータへの誤り訂正符号(ECC: error-correcting code)の付与、及びECCを用いた読み出しデータの誤り検出および訂正の機能(ECC機能)をゲートウェイサーバが実施してよい。この場合、各メモリノードのメモリコントローラはランダム化処理とECC機能を持つ必要がない。さらに、メモリノード間の通信の際にECC処理を実施する必要がない。このため、メモリノード一つ当たりのコストを低減することが可能となる。

【0061】

[4]ストレージ装置の拡張性

ストレージ装置に新たなメモリノードを追加する際の方式について説明する。

【0062】

図9は、ストレージ装置に新たなメモリノードを追加した際の自動アドレス取得方式を示す。

【0063】

本実施形態のストレージ装置10においては、格子座標(x, y)に位置するメモリノードの論理アドレスは位置座標と一致して(x, y)とするため、追加するメモリノードは隣接ノードにアドレスを問い合わせることにより、自ノードのアドレスを簡易に割り出すことができる。

【0064】

例えば、追加するメモリノードがW方向にある隣接ノードにアドレスを問い合わせる場合、隣接ノードのアドレスを(x, y)とすると、追加ノードのアドレスは(x+1, y)となる。同様に、N, E, Sにある隣接ノード(x, y)に問い合わせた場合、自ノードのアドレスはそれぞれ(x, y-1), (x-1, y), (x, y+1)となる。

【0065】

図9に示した場合においては、メモリノード(4, 4)のE方向に新しいメモリノードを追加した状況を例示している。追加するメモリノードは、toとfromがnoneであるaddress_requestパケットをW方向にある隣接ノード(4, 4)に送信する。address_requestパケットを受信したメモリノード(4, 4)は、fromを自ノードアドレスとし、toを(5, 4)としたaddress_answerパケットをE方向にあるメモリノードの出力ポートバッファにセットする。追加ノードは、address_answerパケットを受信することにより、自ノードのアドレスが、(5, 4)であることを割り出すことができる。

【0066】

また、別の方式として、手動アドレス設定方式も考えられる。手動アドレス設定方式では、ストレージ装置10の運用者が予め追加するメモリノードのアドレスを割り出し、それをプリセットしてから、新規のメモリノードを追加する。

【0067】

前述した自動アドレス取得方式は、各メモリノードのパケットコントローラがAddress_requestパケットに対して回答する機能を有することが必要である。このため、パケットコントローラ1つ当たりのコストが高価になる傾向がある。一方、追加するメモリノードのアドレスをプリセットする必要がないため、メモリノードを追加する時の手続きが簡易化され、運用コストを低下させることができる。

【0068】

以上説明したように第1実施形態によれば、メモリノードの論理的なアドレスと物理的な位置(物理アドレス)が一致するため、各メモリノードが経路指定表を管理する必要の

10

20

30

40

50

ない効率的なデータ転送方式を実現することができる。このため、ストレージ装置の拡張性に優れる、すなわちストレージ装置にメモリノードを容易に追加することが可能である。

【 0 0 6 9 】

[第 2 実施形態]

第 2 実施形態のストレージ装置について説明する。第 2 実施形態は、第 1 実施形態が有する転送アルゴリズム 1 と異なる転送アルゴリズム 2 を有する。転送アルゴリズム 2 では、パケットが含む送信先アドレス及び送信元アドレス、自ノードのアドレスの 3 つの情報に基づいて、転送先のメモリノードを決定する。なお、このストレージ装置は、図 1 及び図 2 と同様な構成を有するため、記載を省略する。

10

【 0 0 7 0 】

[1] ストレージ装置の転送アルゴリズム 2

図 1 0 は、第 2 実施形態のストレージ装置における転送アルゴリズム 2 を示す図である。図 1 0 を参照して、パケットの送信先アドレス及び送信元アドレス、自ノードのアドレスの 3 つの情報に基づいて、転送先の隣接ノードを決定する方法の一例を示す。これを転送アルゴリズム 2 とする。

【 0 0 7 1 】

図 1 0 に示すように、送信先ノード(to)と送信元ノード(from)を結ぶ直線で区切られる 2 つの領域のどちらに自ノード(PP)が存在するかを判断し、その領域に割り当てられた方向の隣接するメモリノードにパケットを転送する。

20

【 0 0 7 2 】

図 1 1 A ~ 図 1 1 D を参照して、転送アルゴリズム 2 に従うパケットの転送過程の具体例を示す。

【 0 0 7 3 】

図 1 1 A に示すように、送信先ノード(to)のアドレスを(x_{to} , y_{to})、送信元ノード(from)のアドレスを(x_{from} , y_{from})、自ノード(PP)のアドレスを(x_{pp} , y_{pp})とする。さらに、 $dx = x_{to} - x_{from}$ 、 $dy = y_{to} - y_{from}$ 、 $Dx = x_{pp} - x_{from}$ 、 $Dy = y_{pp} - y_{from}$ とする。転送先を示す方向として、図 1 1 B に示すように、 y が増加する方向を N、 x が増加する方向を E、 y が減少する方向を S、 x が減少する方向を W とする。

【 0 0 7 4 】

$y = (dx / dy) \cdot x$ という送信先ノード(to)と送信元ノード(from)を結ぶ直線の式に基づき、 Dy と $(dy / dx) \cdot Dx$ のどちらが大きいかを比較することで、その直線で区切られる 2 つの領域のどちらに自ノード(PP)が存在しているかを判断する。

30

【 0 0 7 5 】

$dx > 0$ かつ $dy > 0$ の場合を例に転送方向の割り当て例を説明する。 Dy と $(dy / dx) \cdot Dx$ を比較した場合に、 Dy が $(dy / dx) \cdot Dx$ よりも大きくなる領域を A 領域、もう一方の領域を B 領域とする。転送において、A 領域の第一優先方向には E を割り当て、B 領域の第一優先方向には N を割り当てる。転送アルゴリズム 2 においては、送信先ノード(to)と送信元ノード(from)を結ぶ直線に沿ってパケットが進むように転送が行われる。

40

【 0 0 7 6 】

図 1 1 C に、ストレージ装置における転送アルゴリズム 2 のフローを示す。転送アルゴリズム 2 はパケットコントローラ 1 5 に記憶されており、パケットコントローラ 1 5 により実行される。

【 0 0 7 7 】

まず、パケットコントローラ 1 5 は、 $dx = x_{to} - x_{from}$ 、 $dy = y_{to} - y_{from}$ 、 $Dx = x_{pp} - x_{from}$ 、 $Dy = y_{pp} - y_{from}$ を算出する (ステップ S 3 1)。続いて、パケットコントローラ 1 5 は、 dx が 0 であるか否かを判定する (ステップ S 3 2)。 dx が 0 であるとき、 $y_{to} > y_{pp}$ が成り立つか否かを判定する (ステップ S 3 3)。 $y_{to} > y_{pp}$ が成り立つとき、パケットを N へ転送する (ステップ S 3 4)。一方、 $y_{to} > y_{pp}$ が成り立たないとき、パケッ

50

トをSへ転送する(ステップS35)。

【0078】

次に、ステップS32において dx が0でないとき、パケットコントローラ15は、 dy が0であるか否かを判定する(ステップS36)。 dy が0であるとき、 $x_{to} > x_{pp}$ が成り立つか否かを判定する(ステップS37)。 $x_{to} > x_{pp}$ が成り立つとき、パケットをEへ転送する(ステップS38)。一方、 $x_{to} > x_{pp}$ が成り立たないとき、パケットをWへ転送する(ステップS39)。

【0079】

次に、ステップS36において dy が0でないとき、すなわち dx と dy が0でないとき、パケットコントローラ15は、 $Dy \cdot dx > dy \cdot Dx$ が成り立つか否かを判定する(ステップS40)。 $Dy \cdot dx > dy \cdot Dx$ が成り立つとき、 $dx > 0$ かつ $dy > 0$ が成り立つか否かを判定する(ステップS41)。 $dx > 0$ かつ $dy > 0$ が成り立つとき、パケットをEへ転送する(ステップS42)。一方、 $dx > 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy > 0$ が成り立つか否かを判定する(ステップS43)。 $dx < 0$ かつ $dy > 0$ が成り立つとき、パケットをNへ転送する(ステップS44)。

10

【0080】

次に、ステップS43において $dx < 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy < 0$ が成り立つか否かを判定する(ステップS45)。 $dx < 0$ かつ $dy < 0$ が成り立つとき、パケットをWへ転送する(ステップS46)。 $dx < 0$ かつ $dy < 0$ が成り立たないとき、パケットをSへ転送する(ステップS47)。

20

【0081】

次に、ステップS40において $Dy \cdot dx > dy \cdot Dx$ が成り立たないとき、 $dx > 0$ かつ $dy > 0$ が成り立つか否かを判定する(ステップS48)。 $dx > 0$ かつ $dy > 0$ が成り立つとき、パケットをNへ転送する(ステップS49)。 $dx > 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy > 0$ が成り立つか否かを判定する(ステップS50)。 $dx < 0$ かつ $dy > 0$ が成り立つとき、パケットをWへ転送する(ステップS51)。 $dx < 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy < 0$ が成り立つか否かを判定する(ステップS52)。 $dx < 0$ かつ $dy < 0$ が成り立つとき、パケットをSへ転送する(ステップS53)。 $dx < 0$ かつ $dy < 0$ が成り立たないとき、パケットをEへ転送する(ステップS54)。

30

【0082】

以上の処理により、メモリノードに入力されたパケットは、送信先ノード(to)と送信元ノード(from)を結ぶ直線に沿ってパケットが進むように、隣接するメモリノードに転送される。

【0083】

図11Dに、パケットの転送過程の具体例を示す。転送アルゴリズム2においては、送信先ノード(to)と送信元ノード(from)を結ぶ直線に沿ってパケットが進むように転送が行われる。例えば、ケース3では、 $dx < 0$ かつ $dy > 0$ であり、送信先ノード(to)と送信元ノード(from)を結ぶ直線からできるだけ離れないように、WまたはNへ進む。ケース4では、 $dx > 0$ かつ $dy > 0$ であり、送信先ノード(to)と送信元ノード(from)を結ぶ直線からできるだけ離れないように、NまたはEへ進む。

40

【0084】

図12A及び図12Bを参照して、転送アルゴリズム1に対する転送アルゴリズム2の優位性を示す。

【0085】

図12A及び図12Bは、2つの送信元ノード(from1)と(from2)からそれぞれ送信先ノード(to)宛てのパケットを送信した場合のパケットの転送過程を例示している。図12Aは転送アルゴリズム1による転送過程を示し、図12Bは転送アルゴリズム2による転送過程を示す。

【0086】

50

図12Aに示す転送アルゴリズム1では、2つのパケットが合流する地点で、合流待ちの渋滞(jamで示す)が発生している。一方、図12Bに示す転送アルゴリズム2では渋滞は発生していない。したがって、転送アルゴリズム2では、転送アルゴリズム1よりも渋滞発生確率を減少させることができる。

【0087】

以上説明したように第2実施形態によれば、第1実施形態に比べて渋滞発生確率を減少させることができ、多くのクライアントが同時に接続した場合でも、応答速度を維持することができるストレージ装置を提供可能である。

【0088】

また、第1実施形態と同様に、メモリノードの論理アドレスと物理アドレスが一致するため、各メモリノードが経路指定表を管理する必要のない効率的なデータ転送方式を実現することができる。このため、ストレージ装置の拡張性に優れる、すなわちストレージ装置にメモリノードを容易に追加することが可能である。その他の構成及び効果は、前述した第1実施形態と同様である。

【0089】

[第3実施形態]

第3実施形態のストレージ装置について説明する。第3実施形態は、第1,第2実施形態が有する転送アルゴリズム1,2と異なる転送アルゴリズム3を有する。転送アルゴリズム3では、パケットが含む送信先アドレス及び送信元アドレス、自ノードのアドレス、及び自ノードの出力ポート占有情報の4つの情報に基づいて、転送先のメモリノードを決定する。なお、このストレージ装置は、図1及び図2と同様な構成を有するため、記載を省略する。

【0090】

[1]ストレージ装置の転送アルゴリズム3

図13は、第3実施形態のストレージ装置における転送アルゴリズム3を示す図である。図13を参照して、パケットの送信先アドレス及び送信元アドレス、自ノードのアドレス、自ノードの出力ポート占有情報の4つの情報に基づいて、転送先の隣接ノードを決定する方法の一例を示す。これを転送アルゴリズム3とする。

【0091】

図13に示すように、送信先ノード(to)と送信元ノード(from)を結ぶ直線で区切られる2つの領域のどちらに自ノード(PP)が存在するかを判断し、その領域に割り当てられた2つの方向のうち、自ノードの出力ポートの占有情報によって決まる方向の隣接ノードにパケットを転送する。各領域には、第一優先方向と、第二優先方向を割り当てる。第一優先方向の出力ポートバッファが別のパケットにより占有されていた場合に、第二優先方向が選択される。

【0092】

図13、図14A~図14Dを参照して、転送アルゴリズム3に従うパケットの転送過程の具体例を示す。

【0093】

図14Aに示すように、送信先ノード(to)のアドレスを (x_{to}, y_{to}) 、送信元ノード(from)のアドレスを (x_{from}, y_{from}) 、自ノードのアドレス(PP)を (x_{pp}, y_{pp}) とする。さらに、 $dx = x_{to} - x_{from}$ 、 $dy = y_{to} - y_{from}$ 、 $Dx = x_{pp} - x_{from}$ 、 $Dy = y_{pp} - y_{from}$ とする。

【0094】

図13では、 $dx > 0$ かつ $dy > 0$ の場合を例示している。 Dy と $(dy/dx) \cdot Dx$ を比較した場合に、 Dy が $(dy/dx) \cdot Dx$ よりも大きくなる領域をA領域、もう一方の領域をB領域とする。転送において、A領域の第一優先方向にはEを割り当て、B領域の第一優先方向にはNを割り当てる。また、A領域の第二優先方向(bypass方向)にはNを割り当て、B領域の第二優先方向にはEを割り当てる。例えば、図示されているようにパケットを転送するメモリノードがB領域に属す場合において、第一優先方向であるN方向の出力ポートバッファが別のパケットに占有されていれば、そのメモリノードは第二優先

10

20

30

40

50

方向であるE方向にパケットを転送する。

【0095】

図14Cに、ストレージ装置における転送アルゴリズム3のフローを示す。転送アルゴリズム3はパケットコントローラ15に記憶されており、パケットコントローラ15により実行される。なお、図14Bに示すように、Nへパケットを出力するための出力ポートバッファをOPBNとし、Eへパケットを出力するための出力ポートバッファをOPBE、Wへパケットを出力するための出力ポートバッファをOPBW、Sへパケットを出力するための出力ポートバッファをOPBSとする。

【0096】

これら出力ポートバッファOPBN、OPBE、OPBW、OPBSが空いているか（パケットが記憶できるか）、あるいはパケットに占有されているかの情報（出力ポート占有情報）は、以下のようにパケットコントローラ15に記憶される。パケットコントローラ15は、出力ポートバッファと入力ポートバッファの総数分だけバッファ占有フラグビットを有する。パケットコントローラ15は、出力ポートバッファにパケットが書き込まれた場合、そのバッファに対応するバッファ占有フラグビットを“1”にし、出力ポートバッファからパケットが出力された場合、そのバッファに対応するバッファ占有フラグビットを“0”にする。パケットコントローラ15は、バッファ占有フラグビットを評価することで、対応する出力ポートバッファが空いているか否か（パケットに占有されているか否か）を判断できる。

【0097】

まず、パケットコントローラ15は、 $dx = x_{to} - x_{from}$ 、 $dy = y_{to} - y_{from}$ 、 $Dx = x_p - x_{from}$ 、 $Dy = y_{pp} - y_{from}$ を算出する（ステップS61）。続いて、パケットコントローラ15は、 dx が0であるか否かを判定する（ステップS62）。 dx が0であるとき、 $y_{to} > y_{pp}$ が成り立つか否かを判定する（ステップS63）。 $y_{to} > y_{pp}$ が成り立つとき、パケットをNへ転送する（ステップS64）。一方、 $y_{to} > y_{pp}$ が成り立たないとき、パケットをSへ転送する（ステップS65）。

【0098】

次に、ステップS62において dx が0でないとき、パケットコントローラ15は、 dy が0であるか否かを判定する（ステップS66）。 dy が0であるとき、 $x_{to} > x_{pp}$ が成り立つか否かを判定する（ステップS67）。 $x_{to} > x_{pp}$ が成り立つとき、パケットをEへ転送する（ステップS68）。一方、 $x_{to} > x_{pp}$ が成り立たないとき、パケットをWへ転送する（ステップS69）。

【0099】

次に、ステップS66において dy が0でないとき、すなわち dx と dy が0でないとき、パケットコントローラ15は、 $Dy \cdot dx > dy \cdot Dx$ が成り立つか否かを判定する（ステップS70）。 $Dy \cdot dx > dy \cdot Dx$ が成り立つとき、 $dx > 0$ かつ $dy > 0$ が成り立つか否かを判定する（ステップS71）。 $dx > 0$ かつ $dy > 0$ が成り立つとき、Eへ出力する出力ポートバッファOPBEが空いているか否か、すなわち別のパケットにより占有されていないかを判定する（ステップS72）。出力ポートバッファOPBEが空いているとき、パケットをEへ転送する（ステップS73）。一方、出力ポートバッファOPBEが空いていないとき、パケットをNへ転送する（ステップS74）。

【0100】

次に、ステップS71において $dx > 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy > 0$ が成り立つか否かを判定する（ステップS75）。 $dx < 0$ かつ $dy > 0$ が成り立つとき、Nへ出力する出力ポートバッファOPBNが空いているか否かを判定する（ステップS76）。出力ポートバッファOPBNが空いているとき、パケットをNへ転送する（ステップS77）。一方、出力ポートバッファOPBNが空いていないとき、パケットをWへ転送する（ステップS78）。

【0101】

次に、ステップS75において $dx < 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ か

10

20

30

40

50

つ $dy < 0$ が成り立つか否かを判定する (ステップ S 7 9)。 $dx < 0$ かつ $dy < 0$ が成り立つとき、Wへ出力する出力ポートバッファOPBWが空いているか否かを判定する (ステップ S 8 0)。出力ポートバッファOPBWが空いているとき、パケットをWへ転送する (ステップ S 8 1)。一方、出力ポートバッファOPBWが空いていないとき、パケットをSへ転送する (ステップ S 8 2)。

【0102】

次に、ステップ S 7 9において $dx < 0$ かつ $dy < 0$ が成り立たないとき、Sへ出力する出力ポートバッファOPBSが空いているか否かを判定する (ステップ S 8 3)。出力ポートバッファOPBSが空いているとき、パケットをSへ転送する (ステップ S 8 4)。一方、出力ポートバッファOPBSが空いていないとき、パケットをEへ転送する (ステップ S 8 5)。

10

【0103】

次に、ステップ S 7 0において $Dy \cdot dx > dy \cdot Dx$ が成り立たないとき、 $dx > 0$ かつ $dy > 0$ が成り立つか否かを判定する (ステップ S 8 6)。 $dx > 0$ かつ $dy > 0$ が成り立つとき、Nへ出力する出力ポートバッファOPBNが空いているか否かを判定する (ステップ S 8 7)。出力ポートバッファOPBNが空いているとき、パケットをNへ転送する (ステップ S 8 8)。一方、出力ポートバッファOPBNが空いていないとき、パケットをEへ転送する (ステップ S 8 9)。

【0104】

次に、ステップ S 8 6において $dx > 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy > 0$ が成り立つか否かを判定する (ステップ S 9 0)。 $dx < 0$ かつ $dy > 0$ が成り立つとき、Wへ出力する出力ポートバッファOPBWが空いているか否かを判定する (ステップ S 9 1)。出力ポートバッファOPBWが空いているとき、パケットをWへ転送する (ステップ S 9 2)。一方、出力ポートバッファOPBWが空いていないとき、パケットをNへ転送する (ステップ S 9 3)。

20

【0105】

次に、ステップ S 9 0において $dx < 0$ かつ $dy > 0$ が成り立たないとき、 $dx < 0$ かつ $dy < 0$ が成り立つか否かを判定する (ステップ S 9 4)。 $dx < 0$ かつ $dy < 0$ が成り立つとき、Sへ出力する出力ポートバッファOPBSが空いているか否かを判定する (ステップ S 9 5)。出力ポートバッファOPBSが空いているとき、パケットをSへ転送する (ステップ S 9 6)。一方、出力ポートバッファOPBSが空いていないとき、パケットをWへ転送する (ステップ S 9 7)。

30

【0106】

次に、ステップ S 9 4において $dx < 0$ かつ $dy < 0$ が成り立たないとき、Eへ出力する出力ポートバッファOPBEが空いているか否かを判定する (ステップ S 9 8)。出力ポートバッファOPBEが空いているとき、パケットをEへ転送する (ステップ S 9 9)。一方、出力ポートバッファOPBEが空いていないとき、パケットをSへ転送する (ステップ S 1 0 0)。

【0107】

以上の処理により、第一優先方向において渋滞が発生している場合には、第二優先方向へパケットを転送することにより、渋滞を回避しながら送信先ノード(to)と送信元ノード(from)を結ぶ直線に沿って、隣接するメモリノードにパケットが転送される。

40

【0108】

図14Dに、パケットの転送過程の具体例を示す。ケース1は渋滞が発生しない例であり、送信先ノード(to)と送信元ノード(from)を結ぶ直線に沿ってパケットが転送される。ケース2, 3は、渋滞を回避する例を示す。ケース2では、パケットをNへ転送しようとしているとき、N方向に渋滞が発生している。このため、ケース3に示すように、パケットをEへ転送して渋滞を回避している。ケース4, 5は、渋滞を回避する他の例を示す。ケース4では、パケットをEへ転送しようとしているとき、E方向に渋滞が発生している。このため、ケース5に示すように、パケットをNへ転送して渋滞を回避している。

50

【 0 1 0 9 】

転送アルゴリズム 3 においては、第一優先方向で渋滞が発生している場合には、第二優先方向へパケットを転送することにより、渋滞を回避することができる。この際、渋滞がそもそも発生しなかった場合と比べて、送信先ノードにパケットが到達するまでの転送回数は変わらない。

【 0 1 1 0 】

[2] ストレージ装置の転送アルゴリズム 4

図 1 5 は、ストレージ装置における転送アルゴリズム 4 を示す図である。図 1 5 を参照して、パケットの送信先アドレス及び送信元アドレス、自ノードのアドレス、自ノードの出力ポート占有情報の 4 つの情報に基づいて、転送先の隣接ノードを決定する方法の別の第 1 例を示す。これを転送アルゴリズム 4 とする。

10

【 0 1 1 1 】

図 1 5 に示すように、送信先ノード(to)と送信元ノード(from)を結ぶ直線、及びその直線と送信先ノード(to)の位置で直交する直線で区切られる 4 つの領域を定義する。それら 4 つの領域のうち、どの領域に自ノード(PP)が存在するかを判断し、自ノード(PP)が存在する領域に割り当てられた第一優先と第二優先の 2 つの方向のうち、自ノード(PP)の出力ポートの占有情報によって決まる方向の隣接ノードにパケットを転送する。

【 0 1 1 2 】

アルゴリズム 3 と比較したアルゴリズム 4 の優位性は次の通りである。アルゴリズム 3 に対しては、送信先ノード(to)と送信元ノード(from)を結ぶ直線を対角線とする長方形により規定される領域の外にパケットを転送することを禁止する制限を設ける必要がある。一方、アルゴリズム 4 に対しては、その制限を設ける必要がない。

20

【 0 1 1 3 】

[3] ストレージ装置の転送アルゴリズム 5

図 1 6 は、ストレージ装置における転送アルゴリズム 5 を示す図である。図 1 6 を参照して、パケットの送信先アドレス及び送信元アドレス、自ノードのアドレス、自ノードの出力ポート占有情報の 4 つの情報に基づいて、転送先の隣接ノードを決定する方法の別の第 2 例を示す。これを転送アルゴリズム 5 とする。

【 0 1 1 4 】

図 1 6 に示すように、送信先ノード(to)と送信元ノード(from)を結ぶ直線、その直線と送信先ノード(to)の位置で直交する直線、送信先ノード(to)の位置で x 方向と y 方向に伸びる 2 つの直線の 4 つの直線で区切られる 8 つの領域を定義する。それら 8 つの領域のうち、どの領域に自ノード(PP)が存在するかを判断し、自ノード(PP)が存在する領域に割り当てられた第一優先と第二優先の 2 つの方向のうち、自ノード(PP)の出力ポートの占有情報によって決まる方向の隣接ノードにパケットを転送する。前記送信先ノード(to)の位置で x 方向と y 方向に伸びる 2 つの直線は、例えば送信先ノード(to)を通りメモリノードの配列方向に沿った 2 つの直線を含む。

30

【 0 1 1 5 】

アルゴリズム 4 と比較したアルゴリズム 5 の優位性は次の通りである。アルゴリズム 4 では、送信先ノード(to)と送信元ノード(from)を結ぶ直線を対角線とする長方形により規定される領域の外でパケットを第二優先方向(bypass方向)に転送させた場合、送信先ノード(to)にパケットが到達するまでの転送回数が少なくとも 1 つ増加する。一方、アルゴリズム 5 の場合では転送回数の増加がない。

40

【 0 1 1 6 】

以上説明したように第 3 実施形態によれば、第 1 , 第 2 実施形態に比べて渋滞発生確率を減少させることができ、多くのクライアントが同時に接続した場合でも、応答速度を維持することができるストレージ装置を提供可能である。

【 0 1 1 7 】

また、第 1 実施形態と同様に、メモリノードの論理アドレスと物理アドレスが一致するため、各メモリノードが経路指定表を管理する必要のない効率的なデータ転送方式を実現

50

することができる。このため、ストレージ装置の拡張性に優れる、すなわちストレージ装置にメモリノードを容易に追加することが可能である。その他の構成及び効果は、前述した第1実施形態と同様である。

【0118】

[第4実施形態]

第4実施形態では、転送機能を有するメモリノードを複数備えたストレージ装置と、ストレージ装置に接続された制御コンピュータからなるストレージシステムにおいて、パケット転送時の渋滞が発生しない、つまり渋滞フリーとなるシステムの運用条件について説明する。

【0119】

[1]ストレージシステムの構成

図17は、第4実施形態のストレージシステムの構成を示す図である。

【0120】

図示するように、ストレージシステムは、複数のメモリノード11を含むストレージ装置10、及びストレージ装置10に接続された複数の制御コンピュータ41を備える。ストレージ装置10は、図1に示したように、データ転送機能を有するメモリノードが相互に複数接続された構成を有する。ストレージ装置10の外周部に配置されたメモリノード11には、制御コンピュータ41が接続されている。

【0121】

[2]ストレージシステムの渋滞発生フリーの運用条件

図17を参照して、ストレージシステムにおけるパケット転送シミュレーションの枠組みを説明する。ストレージ装置10に接続された制御コンピュータ41から送信されたリクエストパケットが送信先アドレスのメモリノードに到達するまでの過程、および目的のメモリノードから返送されたデータパケットが元の制御コンピュータ41に到着するまでの過程をシミュレートする。パケットの転送アルゴリズムとしては、図13を用いて説明したバイパス機能付ルーティングアルゴリズム(転送アルゴリズム3)を使用する。

【0122】

ストレージ装置10のメモリノード数を $N_{node} = N_x \times N_y$ とし、ストレージ装置10に接続されている制御コンピュータ数を N_c とする。メモリノード間で1回のパケット転送が行われる時間を単位ステップ時間とする。制御コンピュータ41が、単位ステップ時間

あたりに、リクエストパケットを送信する確率を R_r とする。 R_r は最大で1である。バイパス転送発生率を R_{bypass} 、渋滞発生率を R_{jam} とする。

【0123】

バイパス転送発生率 R_{bypass} は、メモリノードが1つのパケットの転送を試みた際に、バイパス転送により渋滞が回避された確率を示す。一方、渋滞発生率 R_{jam} は、メモリノードが1つのパケットの転送を試みた際に、第一優先方向の出力ポートバッファのみならず、第二優先方向の出力ポートバッファも別のパケットにより占有されていたために、パケット転送を実施できず渋滞が発生し、パケットが滞留した確率を示す。また、負荷係数を R_{load} とする。負荷係数 R_{load} は、1つのメモリノードに存在しているパケット数の平均である。

【0124】

図18は、バイパス転送発生率 R_{bypass} と負荷係数 R_{load} の関係、ならびに渋滞発生率 R_{jam} と負荷係数 R_{load} の関係を示す。多数のプロットは、メモリノード数 N_{node} 、制御コンピュータ数を N_c 、確率 R_r を様々に変化させた場合の結果を示している。

【0125】

渋滞発生率 R_{jam} は負荷係数 R_{load} が0.2より小さい場合ほぼゼロであり、渋滞フリーのシステム運用条件は負荷係数 $R_{load} < 0.2$ であることが分かる。渋滞発生率 R_{jam} は負荷係数 R_{load} が0.2以上になると増加し始める。しかし、負荷係数 $R_{load} = 2$ であっても渋滞発生率 R_{jam} は0.05程度であり、実用上ほとんど問題にならない。一方、渋滞発生率 R_{jam} は負荷係数 R_{load} が2より大きくなると急増する。これは、負荷係数 R_l

10

20

30

40

50

ρ_{oad} が 2 より大きくなると連鎖的に渋滞が発生するハングアップ現象が起こるためである。

【 0 1 2 6 】

バイパス転送発生率 R_{bypass} は負荷係数 R_{load} の増加に対して、渋滞発生率 R_{jam} よりも一桁早く立ち上がっている。つまり、バイパス転送アルゴリズムは渋滞発生要件を一桁程度改善する効果がある。

【 0 1 2 7 】

次に、このシミュレーション結果に基づき、渋滞フリーもしくは渋滞発生が実用上問題とならない、制御コンピュータ数とメモリノード数の関係を求める。平均パケット滞在ステップ時間を S_{avg} とする。平均パケット滞在ステップ時間 S_{avg} は、制御コンピュータ 4 1 がリクエストパケットを送信してからそのリクエストパケットが送信先アドレスのメモリノードに到達するまでのステップ時間 $S_{request}$ と、目的のメモリノードから返送されたデータパケットが元の制御コンピュータに到着するまでのステップ時間 S_{dat} の和である。ステップ時間 $S_{request}$ と S_{dat} は、渋滞が発生しない場合 ($N_x / 2 + N_y / 2$) である。従って、平均パケット滞在ステップ時間 S_{avg} は ($N_x + N_y$) であり、これはおよそ $2 \cdot N_{node}$ である。

【 0 1 2 8 】

一方、制御コンピュータ 4 1 が、1 ステップ時間あたりに送信するリクエストパケットの総数は、 $N_c \times R_r$ である。従って、ストレージシステム内に存在するパケット総数は平均で $N_c \times R_r \times 2 \cdot N_{node}$ である。このため、 $R_{load} \sim N_c \times R_r \times 2 \cdot N_{node} / N_{node} = N_c \times R_r \times 2 / N_{node}$ である。ここで、渋滞発生要件を R_{load}^{limit} とすると、好ましいシステム運用条件は、 $R_{load} < R_{load}^{limit}$ となる。従って、渋滞フリーもしくは渋滞発生が実用上問題とならない最大の制御コンピュータ数 N_c^{max} は、 $N_c^{max} < R_{load}^{limit} \times N_{node} / (R_r \times 2)$ となる。 R_r は最大 1 であるから、より厳しい条件は、 $N_c^{max} < R_{load}^{limit} \times N_{node} / 2$ である。

【 0 1 2 9 】

前述したように、パケット転送時に渋滞発生フリーの運用条件は $R_{load}^{limit} = 0.2$ であり、渋滞が実用上問題とならない運用条件は $R_{load}^{limit} = 2$ である。従って、好ましい最大の制御コンピュータ数 N_c^{max} は $N_c^{max} < N_{node}$ であり、より好ましくは $N_c^{max} < 0.1 \times N_{node}$ である。よって、 $N_c < N_{node}$ という関係が成立している。その他の構成及び効果は、前述した第 1 実施形態と同様である。

【 0 1 3 0 】

[第 5 実施形態]

第 5 実施形態は、図 1 に示した転送機能を有する複数のメモリノードを相互に接続したストレージ装置を備え、各メモリノードが分散処理機能を有する。

【 0 1 3 1 】

[1] ストレージ装置の構成

図 1 9 は、第 5 実施形態のストレージ装置の構成を示す図である。

【 0 1 3 2 】

図示するように、ストレージ装置は、図 1 に示した複数のメモリノード 1 1 を含むストレージ装置 1 0 を備えると共に、各メモリノード 1 1 はアドレス変換器 4 2 を備え、分散処理機能を有している。

【 0 1 3 3 】

[2] ストレージ装置の分散処理機能

ここでは、図 1 9 を参照して、apple と orange というキーワードを共に含む URL を見つける AND 検索処理を例に挙げて、分散処理機能を説明する。

【 0 1 3 4 】

ここで、メモリノード (1 , 3) と (4 , 3) がそれぞれ apple と orange の転置ファイルを保持しているものとする。転置ファイルとは、キーワードの一つ一つについて作成される索引ファイルで、例えば apple というキーワードに対応する転置ファイルには apple を

10

20

30

40

50

含む全 URL のリストが格納される。

【 0 1 3 5 】

ゲートウェイサーバは、AND 検索命令を apple と orange の転置ファイルを管理するメモリノードに送信する。転置ファイルを管理するメモリノードは、転置ファイルの内容を key-value 型データにマッピングする。ここで、key と value はそれぞれ URL とその転置ファイルのキーワードとする。

【 0 1 3 6 】

各メモリノードは、自らが保有する key-value 型データの各レコードについて、key をアドレス変換器 4 2 によりアドレスに変換し、そのアドレスに value を含むパケットを送信する。アドレス変換器 4 2 は、あるルールに従い、自ら key からアドレスを算出する場合もあるし、key をアドレスに変換する機能を有するサーバに問い合わせ、key に対応するアドレスを得ても良い。

10

【 0 1 3 7 】

例えば、アドレス変換器 4 2 は、次のアドレッシングルールに従い、key からアドレスを算出する。

【 0 1 3 8 】

アドレス = hash(key) mod N

ここで、hash() は暗号学的なハッシュ関数、N はストレージ装置 1 0 内のメモリノード数、mod N は N を法とする剰余演算を表す。

【 0 1 3 9 】

20

また例えば、アドレス変換器 4 2 は、key に対応するアドレスを別のサーバに問い合わせ、そして別のサーバがコンシステントハッシングにより key をアドレスに変換して元のアドレス変換器 4 2 に回答する。

【 0 1 4 0 】

図 1 9 に示す例では、apple と orange の転置ファイルは共に URL 1 を含むとし、URL 1 はアドレス (2 , 1) に変換されたものとする。この場合、メモリノード (2 , 1) が、URL 1 が AND 検索式を満たすかどうかを判定する役割を担うことになる。メモリノード (2 , 1) には、value が apple と orange である 2 つのパケットが届く。このため、メモリノード (2 , 1) は URL 1 が AND 検索条件を満たしていることが分かる。

【 0 1 4 1 】

30

メモリノード (2 , 1) はその旨をクエリーを出したゲートウェイサーバに通知する。同様の判定を転置ファイルに記載された URL すべてについて多数のメモリノードが行うことにより、クエリーを出したゲートウェイサーバは AND 検索条件を満たした URL のリストを得ることができる。

【 0 1 4 2 】

以上の AND 検索演算を、単一メモリノードで行う場合は、次のコードで示す演算を行う必要がある。

【 数 1 】

```
for(i=0;i<Napple;i++)
    for(j=i;j<Norange;j++)
        if(URLlistapple[i]==URLlistorange[j])
            ... (1)
```

40

【 0 1 4 3 】

ここでは、apple と orange の転置ファイルがそれぞれ格納する URL の数を、それぞれ N_{apple} と N_{orange} とし、== は一致演算を示す。メモリノードは一致演算を、N_{apple} × N_{orange} / 2 回だけ繰り返す必要がある。本実施形態のストレージ装置では、この一致演算を数多くのメモリノードで分散して実施することができる。

【 0 1 4 4 】

以上説明したように第 5 実施形態によれば、第 1 実施形態に記載したストレージ機

50

能だけでなく、分散処理（分散コンピューティング）機能も併せ持つことができる。その他の構成及び効果は、前述した第1実施形態と同様である。

【0145】

〔第6実施形態〕

第6実施形態は、転送機能を有するメモリノードを複数備えたストレージ装置と、ストレージ装置の一部のメモリノードが、隣接するメモリノード（以下、隣接ノード）と接続された入出力ポート（以下、隣接ポート）以外の入出力ポート（以下、非隣接ポート）を有し、前記一部のメモリノードが、前記非隣接ポートにより、制御コンピュータもしくは隣接しないメモリノード（以下、非隣接ノード）と接続された構成を有する。

【0146】

〔1〕ストレージシステムの構成

図20は、第6実施形態のストレージシステムの構成を示す図である。

【0147】

図示するように、ストレージシステムは、複数のメモリノード11を含むストレージ装置10、及びストレージ装置10に接続された複数のゲートウェイサーバ21A、21Bを備える。ストレージ装置10は、図1に示したように、データ転送機能を有するメモリノードが相互に複数接続された構成を有する。ストレージ装置10の外周部に配置されたメモリノード（1,4）には、アダプタ22Aを介してゲートウェイサーバ21Aが接続されている。ストレージ装置10の中央部に配置されたメモリノード（3,5）は、隣接ノードと接続された入出力ポート以外の入出力ポート（非隣接ポート）を有し、この非隣接ポートにはアダプタ22Bを介してゲートウェイサーバ21Bが接続されている。

【0148】

〔2〕ストレージシステムの packets 転送

図20を参照して、隣接ノードと接続された入出力ポート以外の入出力ポート（非隣接ポート）を有するメモリノードを導入した効果について述べる。ここで、ストレージ装置10のx方向のメモリノード数を N_x 、y方向のメモリノード数を N_y とする。

【0149】

ゲートウェイサーバの接続先がストレージ装置10の外周部に配置されたメモリノードに限定される場合、ゲートウェイサーバから目的のメモリノードに到達するまでのパケットの平均転送回数は、 $N_x / 2 + N_y / 2$ である。

【0150】

一方、図20に示すように、ストレージ装置10の中央部に位置するメモリノード（3,5）の入出力ポート数を増やし、追加した入出力ポート（前記非隣接ポート）によって、メモリノード（3,5）をゲートウェイサーバ21Bと接続する。すると、ゲートウェイサーバ21Bから目的のメモリノードまでのパケットの平均転送回数は、 $N_x / 4 + N_y / 4$ となる。

【0151】

このように、ストレージ装置10内の一部のメモリノードの入出力ポート数を増やし、追加した入出力ポートでメモリノードとゲートウェイサーバを接続すると、パケットが目的のメモリノードに到達するまでの平均転送回数を低減させることができる。

【0152】

また、図21に示すように、追加した入出力ポートによってゲートウェイサーバに接続されたメモリノードをストレージ装置10に複数配置しても良い。すなわち、メモリノード（4,7）は隣接ノードと接続された入出力ポート以外の入出力ポート（非隣接ポート）を有し、この非隣接ポートはアダプタ22Cを介してゲートウェイサーバ21Cに接続されている。

【0153】

さらに、ストレージ装置10内のメモリノードに接続された複数のゲートウェイサーバ21A、21B、21C、21D同士を、サーバとメモリノードとを接続する配線とは別の配線23で接続しても良い。また、ゲートウェイサーバ間、例えばゲートウェイサーバ

10

20

30

40

50

2 1 A , 2 1 B 間にアダプタ 2 2 E を接続しても良い。これにより、配線 2 3 を使用して、ゲートウェイサーバ間、及びゲートウェイサーバとメモリノード間でパケット転送を行うことができる。例えば、ゲートウェイサーバ 2 1 A ゲートウェイサーバ 2 1 B ゲートウェイサーバ 2 1 C メモリノード (4 , 7) の順でパケットを転送してもよい。また、ゲートウェイサーバ 2 1 A ゲートウェイサーバ 2 1 D メモリノード (5 , 4) の順でパケットを転送してもよい。

【 0 1 5 4 】

以上説明したように第 6 実施形態によれば、ゲートウェイサーバとストレージ装置内のメモリノードとの間でパケット転送を行う際に、転送時間を短縮させることができる。その他の構成及び効果は、前述した第 1 実施形態と同様である。

10

【 0 1 5 5 】

[第 7 実施形態]

第 7 実施形態は、転送機能を有するメモリノードを複数備えたストレージ装置と、ストレージ装置の一部のメモリノードが、隣接ノードと接続された入出力ポート以外の入出力ポート (非隣接ポート) を有し、前記一部のメモリノードが、前記非隣接ポートにより、制御コンピュータもしくは隣接しないメモリノード (非隣接ノード) と接続された構成を有する。

【 0 1 5 6 】

[1] ストレージシステムの構成

図 2 2 は、第 7 実施形態のストレージシステムの構成を示す図である。

20

【 0 1 5 7 】

図示するように、ストレージシステムは、複数のメモリノード 1 1 を含むストレージ装置 1 0、及びストレージ装置 1 0 に接続されたゲートウェイサーバ 2 1 A を備える。ストレージ装置 1 0 は、図 1 に示したように、データ転送機能を有するメモリノードが相互に複数接続された構成を有する。ストレージ装置 1 0 の外周部に配置されたメモリノード (7 , 1) には、アダプタ 2 2 A を介してゲートウェイサーバ 2 1 A が接続されている。

【 0 1 5 8 】

ストレージ装置 1 0 の内部に配置されたメモリノード (3 , 3)、(3 , 8)、(8 , 3)、(8 , 8) は、隣接ノードと接続された入出力ポート以外の入出力ポート (非隣接ポート) を有する。非隣接ポート間は配線 2 4 により接続されている。例えば、メモリノード (3 , 3) と (3 , 8) 間、メモリノード (3 , 8) と (8 , 8) 間、メモリノード (8 , 8) と (8 , 3) 間、及びメモリノード (8 , 3) と (3 , 3) 間が、配線 2 4 により相互に接続されている。

30

【 0 1 5 9 】

非隣接ノードに接続されたメモリノードは、物理的な位置によって決まるアドレスに加えて、非隣接ノードと接続されたメモリノード同士の相対的な物理位置によって決まる付加アドレスを持つ。非隣接ノードと接続されたメモリノードは、自ノード以外のメモリノードを宛先とするパケットを受信した場合、自ノードのパケットコントローラが、自ノードの付加アドレス、もしくは自ノード以外の、非隣接ノードに接続されたメモリノードの付加アドレスのいずれか一つを少なくとも含む情報に基づいてパケットの出力ポートを決定する。

40

【 0 1 6 0 】

[2] ストレージシステムのパケット転送

まず、隣接しないメモリノード (非隣接ノード) 同士の接続が、ストレージ装置 1 0 内のパケット転送回数に与える影響について説明する。前述したように、一部の非隣接ノードが非隣接ポートによって相互接続されたストレージ装置を図 2 2 に、全メモリノードの入出力ポート数と同じで、非隣接ノード間の接続が無いストレージ装置を図 2 3 に示す。

【 0 1 6 1 】

これらストレージ装置において、アドレス (7 , 0) のゲートウェイサーバ 2 1 A から、メモリノード (7 , 8) にデータパケットを送付するケースを考える。

50

【 0 1 6 2 】

非隣接ノード間の接続が無い場合(図23)、最短のデータ転送経路は、ゲートウェイサーバ21A(7,0)メモリノード(7,1)(7,2)(7,3)(7,4)(7,5)(7,6)(7,7)(7,8)である。よって、ゲートウェイサーバ21Aからメモリノード(7,8)にパケット転送を行う際の(ゲートウェイサーバとメモリノード間、及びメモリノードとメモリノード間の)総パケット転送回数は8回である。

【 0 1 6 3 】

他方、非隣接ノード間が相互接続されている場合(図22)、ゲートウェイサーバ21Aメモリノード(7,1)(7,2)(7,3)(8,3)(8,8)(7,8)が最短経路になり、6回の転送でパケットをゲートウェイサーバ21Aから送信先メモリノード(7,8)に送ることができる。

10

【 0 1 6 4 】

このように、一部の隣接しないメモリノード同士を増設した入出力ポートによって接続し、それら入出力ポート間の接続配線を用いてパケット転送を行うことにより、送信先ノードに到達するまでのパケット転送回数を低減でき、転送時間を短縮することができる。

【 0 1 6 5 】

但し、図22に示すように、非隣接ノード間が接続されたストレージ装置では、メモリノードとメモリノード間、或いはメモリノードとゲートウェイサーバ21A間の物理的な距離と、それらのデータ転送上の距離が必ずしも一致しない点に注意が必要である。

20

【 0 1 6 6 】

前述したように、アドレス(7,0)のゲートウェイサーバ21Aから、メモリノード(7,8)のデータパケットを送付する場合において、非隣接ノード間の接続が存在しないストレージ装置(図23)では、パケットの最低総転送回数(8回)は、パケット送信元と送信先の物理的な位置で決まるアドレス(ゲートウェイサーバ21A(7,0)、メモリノード(7,8))の差に等しい。

【 0 1 6 7 】

他方、前述したように、非隣接ノード間が接続された図22のストレージ装置では、同じ物理アドレスを有するゲートウェイサーバ21A(7,0)~送信先メモリノード(7,8)間のパケット通信を、6回の転送で行うことができる。このように、非隣接ノード間が接続されたストレージ装置では、物理的な最短距離と、データ転送上の最短距離が異なるため、物理的な位置によって決まるアドレスのみでパケットの最短転送経路を決定することはできない。

30

【 0 1 6 8 】

そこで、非隣接ノード間が接続されたストレージ装置では、これら非隣接ノードに接続されたメモリノード同士の相対的な物理位置を反映したアドレス(以下、サブアドレスと記す)を非隣接ノードに接続されたメモリノードに追加付与し、全メモリノードに与えられた、物理位置を反映したアドレス(以下、メインアドレスと記す)とサブアドレスの両者からデータ転送経路・パケット転送先を決定すれば、前述したアルゴリズムに従って効率良くパケット転送を行うことができる。なお、図22では、メインアドレスを丸括弧で、サブアドレスを四角括弧で表示している。また、前記非隣接ノードに接続されたメモリノードは、例えば図22に示したメモリノード(3,3)、(3,8)、(8,3)、(8,8)を指す。

40

【 0 1 6 9 】

但し、前述したように効率良くパケット転送を行うためには、パケットのヘッダー部に必要なアドレス情報が全て書かれていること、転送先決定に用いられる、パケットの送信元アドレス及び送信先アドレスを適宜書き換えること等が必要である。これらについて、次に説明する。

【 0 1 7 0 】

非隣接ノードに接続されたメモリノードを有するストレージシステムでパケット転送を

50

行う場合、望ましいパケットのヘッダー部のアドレス情報の一例を図 2 4 に示す。

【 0 1 7 1 】

ここで、最終送信先ノードと送信元ノードは、それぞれパケットを最終的に届ける宛先のノードと、パケットを最初に作成・発信するノードを意味する。他方、一時送信先ノードアドレスと一時送信元ノードアドレスは、各メモリノードでパケット転送先決定に用いられるアドレスであり、パケットの通信過程で書き換えられる。このときの書き換えルールについては後述する。

【 0 1 7 2 】

なお、一時送信先ノード/送信元アドレスのタイプは、そのアドレスが、メインアドレスとサブアドレスのどちらであるか区別するためのものである。また、第 1 中継ノード、第 2 中継ノードは、パケット送信時に経由すべき非隣接ノードに接続されたメモリノードであり、送信元ノードに最も近い非隣接ノードに接続されたメモリノードが第 1 中継ノード、最終送信先ノードに最も近い非隣接ノードに接続されたメモリノードが第 2 中継ノードである。

【 0 1 7 3 】

図 2 2 に示したように、隣接しないメモリノードと一部のメモリノードが接続されたストレージ装置を含むシステムで、非隣接ノード間の接続を経由したパケット転送を行う場合、パケットを受信したメモリノードにおけるパケット転送先決定やパケットのヘッダー部の修正等は、例えば次のようなルールで行うことが好ましい。

【 0 1 7 4 】

(a) 追加の入出力ポート（非隣接ポート）を有しないメモリノードがパケットを受信した場合、パケットコントローラがパケットのヘッダー部に記録されたアドレス情報を調べ、

(i) 最終送信先ノードのアドレスが自ノードのアドレスと一致する場合には、パケット転送を行わない。

【 0 1 7 5 】

(ii) 最終送信先ノードのアドレスが自ノードのアドレスと異なる場合には、一時送信元ノード、一時送信先ノード、自ノードのメインアドレスを参照して転送先を決定し、隣接したメモリノードにパケットを送信する。

【 0 1 7 6 】

(b) 非隣接ノードに接続されたメモリノードがパケットを受信した場合、パケットコントローラはパケットのヘッダー部に記録されたアドレス情報を調べ、

(i) 自ノードのアドレスが最終送信先ノードのアドレスと一致する場合には、パケット転送を行わない。

【 0 1 7 7 】

(ii) 自ノードのアドレスが、最終送信先ノードのアドレスと異なり、

(1) 第 1 中継ノードのアドレスと一致する場合は、パケットのヘッダー部の一時送信先、一時送信元のアドレスを、それぞれ第 2 中継ノードのサブアドレス、自ノードのサブアドレスに書き換える。さらに、一時送信元ノード、一時送信先ノード、自ノードのサブアドレスを参照して転送先を決定し、他の非隣接ノードに接続されたメモリノードにパケットを送信する。

【 0 1 7 8 】

(2) 第 2 中継ノードのアドレスと一致する場合は、一時送信先アドレスを最終送信先ノードのメインアドレス、一時送信元アドレスを自ノードのメインアドレスに変更する。さらに、一時送信元ノード、一時送信先ノード、自ノードのメインアドレスを参照して転送先を決定し、隣接した他のメモリノードにパケットを転送する。

【 0 1 7 9 】

(3) 第 1 中継ノード、第 2 中継ノードのいずれのアドレスとも異なり、さらにパケットのヘッダー部に書かれた一時送信先アドレスと一時送信元アドレスがサブアドレスである場合には、一時送信元ノード、一時送信先ノード、自ノードのサブアドレスを参照し

10

20

30

40

50

て転送先を決定し、他の非隣接ノードに接続されたメモリノードにパケットを転送する。

【0180】

(4) 第1中継ノード、第2中継ノードのいずれのアドレスとも異なり、さらに一時送信先アドレスと一時送信元アドレスがメインアドレスである場合には、一時送信元ノード、一時送信先ノード、自ノードのメインアドレスを参照して転送先を決定し、他の隣接したメモリノードにパケットを転送する。

【0181】

次に、非隣接ノード間の接続を併用してパケット送信を行う例として、図22に示したシステムで、クライアントがゲートウェイサーバ21Aを介してストレージ装置にファイルを書き込む手続きを説明する。

【0182】

図25は、図22に示したストレージシステムにおける書き込み動作を示す図である。

【0183】

まず、クライアントは、ファイルとファイルIDとをゲートウェイサーバ21Aに送信する((1)参照)。ファイルIDは、ファイルを一意に特定できる識別子である。

【0184】

次に、ゲートウェイサーバ21Aは、ファイルを規定サイズのデータに分割し、分割した各データに分割データIDを割り振る。さらに、ファイルIDと分割データIDをファイルテーブルに書き込む。分割データIDは、分割されたデータを一意に特定できる識別子である((2)参照)。

【0185】

次に、ゲートウェイサーバ21Aは、分割データIDの情報に基づき、分割データを書き込む書き込み先のメモリノード(以下、書き込みノード)のアドレスを決定する(図22ではメインアドレス(7,8)のメモリノード)。

【0186】

次に、下記手順で、ゲートウェイサーバ21Aから書き込みノードにパケットを転送する際に、通信時間(パケット転送回数)が最小となる経路を求める：

1. ゲートウェイサーバ21Aに最も近い、非隣接ノードに接続されたメモリノード(第1中継ノード)と、書き込みノードに最も近接した非隣接ノードに接続されたメモリノード(第2中継ノード)のそれぞれのアドレスを調べる。図22で(7,8)に書き込む場合、ゲートウェイサーバ21Aと書き込みノードのそれぞれの最近接の非隣接ノードに接続されたメモリノードは、メインアドレス(8,3)、(8,8)のメモリノードである。

【0187】

2. 非隣接ノード間の接続を含む最短経路と、含まない最短経路を求める。図22でメモリノード(7,8)に書き込む場合、前者は、ゲートウェイサーバ21Aメモリノード(7,1)(7,2)(7,3)(8,3)([2,1])(8,8)([2,2])(7,8)であり、後者は、ゲートウェイサーバ21Aメモリノード(7,1)(7,2)(7,3)(7,4)(7,5)(7,6)(7,7)(7,8)である。

【0188】

3. 隣接ノードとのみ接続されたメモリノード間の経路についてはメインアドレスで、非隣接ノードに接続されたメモリノード間の経路は、サブアドレスを元にして、それぞれの経路でパケット送信を行う際に生じる転送回数を算出する。図22でメモリノード(7,8)に書き込む場合、非隣接ノードに接続されたメモリノードを介する最短経路の転送回数は6回であり、非隣接ノードに接続されたメモリノードを介さない最短経路では、転送回数は8回である。

【0189】

4. パケット転送回数が小さい方の経路をデフォルト経路として決定する((3)参照)。以下では、非隣接ノード間の接続を含む経路の方が、含まない経路よりもパケット転送回数が少ないケースについて説明する。

10

20

30

40

50

【 0 1 9 0 】

・アドレスデータ、書き込み命令等からなるヘッダ部を、書き込みデータに付加した書き込みパケットをゲートウェイサーバ21Aで作成する。なお、パケットのヘッダ部に記録されたアドレス情報の一時送信先は第1中継ノードのメインアドレス(図22では(8,3))とし、一時送信元はパケット送信元のメインアドレス(図22では(7,0))とする。

【 0 1 9 1 】

・作成した書き込みパケットを、ゲートウェイサーバ21Aからゲートウェイサーバ21Aに接続されたメモリノードに転送する((4)参照)。

【 0 1 9 2 】

・前述したアルゴリズムに従って、一時送信先メモリノード(図22ではメインアドレス(8,3)の第1中継ノード)に到達するまで、メモリノード間で書き込みパケットの転送を行う(各メモリノードは一時送信先、一時送信元、自ノードのメインアドレスを参照して転送先を決定する)((5)参照)。

【 0 1 9 3 】

・書き込みパケットを受け取った第1中継ノードは、パケットのヘッダ部を読み、パケットの最終送信先ノードが他のメモリノードであり、第1中継ノードが自ノードであることから、一時送信先を第2中継ノードのサブアドレス(図22では[2,2])、一時送信元を自ノードのサブアドレス(図22では[2,1])に変更する((6)参照)。さらに、隣接した非隣接ノードに接続されたメモリノードにパケットを転送する((7)参照)。

【 0 1 9 4 】

・前述したアルゴリズムに従い、一時送信先メモリノードである第2中継ノードに到達するまで、非隣接ノードに接続されたメモリノード間で書き込みパケットの転送を行う(各メモリノードは一時送信先、一時送信元、自ノードのサブアドレスを参照して転送先を決定する)。

【 0 1 9 5 】

・書き込みパケットを受信した第2中継ノードは、パケットのヘッダ部を読み、最終送信先が他のメモリノードで、第2中継ノードが自ノードであることから、一時送信先を最終送信先ノードのメインアドレス(図22では書き込みノード(7,8))、一時送信元を自ノードのメインアドレス(図22では(8,8))に変更する((8)参照)。さらに、隣接したメモリノードに書き込みパケットを転送する。

【 0 1 9 6 】

・一時送信元ノード、一時送信先ノード、自ノードのメインアドレスを参照して、前述したアルゴリズムによりメモリノード間で転送が繰り返され、最終送信先である書き込みノード(図8ではメインアドレス(7,8))に書き込みパケットが到達する((9)参照)。

【 0 1 9 7 】

・パケットを受信した書き込みノードは、最終送信先が自ノードであることから、パケットを転送せず、ヘッダ部のアドレス情報や書き込みデータ等を自ノードのメモリ16に書き込む((10)参照)。その後、宛先や中継先を逆にした書き込み完了報告のパケットを作成する。このパケットでは、発信元ノードを書き込みノードの(7,8)とし、最終送信先ノードをゲートウェイサーバ21Aの(7,0)、第1中継ノードをメインアドレス(8,8)(サブアドレス[2,2])のノード、第2中継ノードをメインアドレス(8,3)(サブアドレス[2,1])のノードとする。その後、書き込み完了報告のパケットを前述した手順と同様の手続きで、ゲートウェイサーバ21Aへ返信する((11)参照)。全分割データの書き込みが終わった後、ゲートウェイサーバ21Aはクライアントにファイル書き込みの終了報告を行う。

【 0 1 9 8 】

前述した例では、隣接ノードと接続された入出力ポート以外の入出力ポート(非隣接ボ

10

20

30

40

50

ート)を有し、隣接しないメモリノード同士が非隣接ポートによって相互接続されたメモリノードを有するストレージ装置にクライアントがファイルを書き込む手順を説明したが、メモリノードに書き込まれたデータの読み出し・消去やメモリの空き容量確認などのため、読み出し命令・消去命令や空き容量の返信命令等をゲートウェイサーバからストレージ装置のメモリノードに送る場合も、非隣接ノード間の接続を介して通信を行うことでパケット転送時間を短縮することができる。さらに、読み出したデータや問合せ結果などをメモリノードからゲートウェイサーバに送る場合などにおいても、非隣接ノード間の接続を介して通信を行うことでパケット転送時間を短縮することができる。

【0199】

前述したように、隣接しないノード間の接続を介してデータ送信を行うことでパケットの平均転送回数を減らすことが可能である。なお、追加入出力ポート(非隣接ポート)によって非隣接ノードと接続されるメモリノードは、ストレージ装置内で均一に分散して配置されていることが好ましい。

【0200】

例えば、 10×10 のメモリノードから構成されている図22に示したストレージ装置では、 5×5 のメモリノードを含む4つの領域に分割(境界を点線で示す)し、各領域の中心に位置するメモリノード(3, 3)、(8, 3)、(3, 8)、(8, 8)同士を増設ポート(非隣接ポート)で接続している。この場合には、どの 5×5 のメモリノードブロックにも、1つの追加入出力ポート(非隣接ポート)を有するメモリノードが存在する。

【0201】

このように、含まれるメモリノード数が同じであるような、複数の領域にストレージ装置を分割し、その中心のメモリノードを追加入出力ポート付きのメモリノードとして相互接続することで、追加入出力ポート付きのメモリノードをストレージ装置内で均一に分散配置させることができる。

【0202】

なお、例えば、 $a1 \times b1$ のノードから構成されたストレージ装置を $1 \times m$ ずつの領域 a 個に分割し、各領域の中心に位置するメモリノード同士を追加入出力ポートにより接続する場合、 $(c1 + d, e1 + f)$ ($a = c1, b = e1, l = d, m = f$)で表わされるメモリノードに最も近く、各領域の中心に位置する追加入出力ポート付きのメモリノードアドレスは、 $(c1 + \text{round}(l/2, 0), e1 + \text{round}(m/2, 0))$ で表わされる。 $(\text{round}(A, 0))$ は、 A の小数点以下を四捨五入する関数。

【0203】

さて、ストレージ装置において、隣接しない一部のメモリノード間を接続することで生じる課題は、その接続を介して通信を行うパケットの数が増えると、パケット通信の渋滞が生じやすくなることである。

【0204】

これに対して、隣接したメモリノード間よりも、隣接しないメモリノード間のパケット通信速度を高くすることで、このパケット転送の渋滞を緩和させることができる。但し、この場合には、ゲートウェイサーバにおいて、パケット転送のデフォルト経路決定時に行う、非隣接ノード間の接続を介する最短経路と、非隣接ノード間の接続を経由しない最短経路の総パケット通信時間の算出に注意が必要である。総パケット通信時間は、メモリノード間の通信時間とパケット転送回数の積に等しい。従って、隣接したメモリノード間と隣接しないメモリノード間のデータ通信速度が等しい場合には、前述したように、各経路の転送回数でデータ通信時間を比較することができる。しかし、パケット送信速度が異なる場合には、転送回数だけで総パケット通信時間を比べることができない。

【0205】

通信速度の逆数(速度 = 距離/時間であるから、これは通信時間に比例)と転送回数の積を参照することで、非隣接ノード同士の接続を経由する経路と、非隣接ノード同士の接続を経由しない経路の総パケット通信時間を正確に比較することができる。しかし、前述

10

20

30

40

50

した例のように、隣接したノード間の接続と非隣接ノード間の接続がパケット通信経路で混在している場合、この計算は煩雑なものとなる。

【0206】

隣接ノード間と非隣接ノード間でパケット通信速度が異なるストレージ装置では、付与するアドレスのステップをパケット通信速度に逆比例させ、そのアドレス差を転送回数と見なすのも一つの方法である。

【0207】

隣接ノード間よりも非隣接ノード間の方が、パケット通信速度が10倍高い場合のアドレス付与例を図26に示す。図26では、隣接しないノードと接続されたメモリノードには、それら物理位置を反映させながら、あるメモリノードと、そのメモリノードと接続された非隣接ノードとの間で1ずつ異なるようにサブアドレスを付与している(四角括弧で表示)。さらに、図26に示す全メモリノードには、互いの物理的な位置関係を反映させながら、隣り合うメモリノード間で、10ずつ異なるメインアドレスを与えている(丸括弧で表示)。

10

【0208】

このようなストレージ装置で、非隣接ノードに接続されたメモリノード間でパケット送信を行う場合は各メモリノードのサブアドレスの差を計算し、隣接メモリノード間でパケット送信を行う場合はメインアドレスの差を計算し、それらを転送回数と見なす。転送回数の算出ルールをこのように定めれば、隣接ノード間と非隣接ノード間で接続速度が異なる場合においても、転送回数だけでパケット通信時間を見積もり、比較することができる。

20

【0209】

以上説明したように第7実施形態によれば、ゲートウェイサーバとストレージ装置内のメモリノードとの間でパケット転送を行う際に、転送時間を短縮させることができる。その他の構成及び効果は、前述した第1実施形態と同様である。

【0210】

[第8実施形態]

第8実施形態では、パケットの平均転送回数を減らすため、ゲートウェイサーバとストレージ装置との間にスイッチングリレーを追加したストレージシステムについて説明する。

30

【0211】

[1]ストレージシステムの構成

図27は、第8実施形態のストレージシステムの構成を示す図である。

【0212】

図示するように、ストレージシステムは、複数のメモリノード11を含むストレージ装置10、ストレージ装置10に接続されたスイッチングリレー81、及びスイッチングリレー81に接続されたゲートウェイサーバ21Aを備える。

【0213】

ストレージ装置10は、図1に示したように、データ転送機能を有するメモリノードが相互に複数接続された構成を有する。スイッチングリレー81は、ストレージ装置10の一端側(図27において左端)に配置されたメモリノード(1,1)、(1,2)、(1,3)、(1,4)、(1,5)、(1,6)、(1,7)、(1,8)、(1,9)の全てに接続されている。スイッチングリレー81には、アダプタ22Aを介してゲートウェイサーバ21Aが接続されている。

40

【0214】

スイッチングリレー81は、パケットのヘッダー部に記録されたアドレス情報に従い、受信したパケットを指定された送信先ノードに転送する。スイッチングリレー81には、ストレージ装置内のメモリノードとは異なるアドレスが付与されている(図27では四角括弧で表示)。

【0215】

50

〔 2 〕ストレージシステムの packets 転送

第 8 実施形態のストレージシステムにおいて packets を転送する手順について説明する。

【 0 2 1 6 〕

ゲートウェイサーバ 2 1 A から送り出された packets は、アダプタ 2 2 A を通ってスイッチングリレー 8 1 に入る。スイッチングリレー 8 1 に入った packets は、スイッチングリレー 8 1 と接続された、メモリノード (1 , 1)、(1 , 2)、(1 , 3)、(1 , 4)、(1 , 5)、(1 , 6)、(1 , 7)、(1 , 8)、(1 , 9) のいずれか一つに送られ、その後、送信先アドレスのメモリノードまで転送される。

【 0 2 1 7 〕

逆に、ストレージ装置 1 0 内のメモリノードから発信された packets は、メモリノード (1 , 1) ~ (1 , 9) のいずれか一つのメモリノードに送られ、スイッチングリレー 8 1 とアダプタ 2 2 A を介してゲートウェイサーバ 2 1 A に転送される。

【 0 2 1 8 〕

図 2 7 に示したシステムにおいて、ゲートウェイサーバ 2 1 A からスイッチングリレー 8 1 とメモリノード (1 , 9) を介して、メモリノード (5 , 9) に packets を送る場合は以下となる。最短経路はスイッチングリレー 8 1 → メモリノード (1 , 9) → (2 , 9) → (3 , 9) → (4 , 9) → (5 , 9) で、ストレージ装置 1 0 内の転送回数は 4 回である。

【 0 2 1 9 〕

他方、図 2 8 に示すように、スイッチングリレー 8 1 を介さずに、ゲートウェイサーバ 2 1 A がストレージ装置 1 0 のメモリノード (1 , 4) に接続されているシステムでは、ゲートウェイサーバ 2 1 A からメモリノード (5 , 9) に packets を送る場合は以下となる。ゲートウェイサーバ 2 1 A → メモリノード (1 , 4) → (1 , 5) → (1 , 6) → (1 , 7) → (1 , 8) → (1 , 9) → (2 , 9) → (3 , 9) → (4 , 9) → (5 , 9) が最短経路の一つであり、ストレージ装置 1 0 内で packets の転送は最低 9 回必要である。

【 0 2 2 0 〕

このように、ゲートウェイサーバ 2 1 A とストレージ装置 1 0 との間にスイッチングリレー 8 1 を導入すると、packets の転送回数を減らし、転送時間を短縮することができる。

【 0 2 2 1 〕

なお、ゲートウェイサーバ 2 1 A からストレージ装置 1 0 に packets を送信するケースについて説明したが、ストレージ装置 1 0 内のメモリノードに保存されたデータをゲートウェイサーバ 2 1 A に送信する場合においても、ゲートウェイサーバ 2 1 A とストレージ装置 1 0 との間にスイッチングリレー 8 1 が存在するシステム(図 2 7)の方が、アダプタを介してゲートウェイサーバ 2 1 A とストレージ装置 1 0 が直接接続されたシステム(図 2 8)よりも、一般に、packets の平均転送回数が少なく、転送時間が短い。

【 0 2 2 2 〕

但し、スイッチングリレーを導入したシステムにおける packets 転送回数は、スイッチングリレーとストレージ装置間で packets の送受信を行う際に中継するメモリノードに依存する。例えば、前述したように、図 2 7 に示したシステムでゲートウェイサーバ 2 1 A がメモリノード (5 , 9) との間で packets 転送を行う場合、中継するメモリノードがメモリノード (1 , 9) であれば、ストレージ装置 1 0 内の転送回数は最低 4 回である。しかし、他のメモリノードを中継した場合には、5 回以上の転送が必要である。

【 0 2 2 3 〕

このように、packets の転送回数を最小限にするためには、スイッチングリレーに直接接続されたメモリノードのうち、送信先ノードに最も近いメモリノードを中継メモリノードとして選択することが重要である。packets 転送時間が最短となる、中継メモリノードのアドレス算出手順については後述する。

10

20

30

40

50

【 0 2 2 4 】

なお、スイッチングリレーを介さずにゲートウェイサーバとストレージ装置が接続されたシステムにおいて、ゲートウェイサーバとストレージ装置内のメモリノードとの間でパケット転送を行う場合、送信元と送信先は、ゲートウェイサーバ、メモリノードのいずれかであり、これはパケット転送が行われている間変更されない。

【 0 2 2 5 】

他方、スイッチングリレーを介してゲートウェイサーバとストレージ装置が接続されたシステムにおいて、パケット転送を行う場合には、まず中継メモリノードを宛先としてパケット転送が行われる。パケットが中継メモリノードに到達した後は、メモリノード(ゲートウェイサーバからメモリノードにパケットを送信する場合)或いはスイッチングリレー(メモリノードからゲートウェイサーバにパケット送信する場合)を送信先としてパケット転送が行われる。すなわち、パケットが中継メモリノードに到達する前後で、パケットの送信先の変更が必要である。

10

【 0 2 2 6 】

スイッチングリレーを使用したシステムにおいて、パケットのヘッダ部に記録されたアドレス情報の例を図 2 9 に示す。ここで、最終送信先ノードと送信元ノードは、それぞれパケットを最終的に届ける宛先のメモリノード、パケットを最初に作成・発信するメモリノードを意味する。他方、一時送信先ノードアドレスと一時送信元ノードアドレスは、各メモリノードでパケット転送先の決定に用いられるアドレスであり、スイッチングリレーに接続された中継メモリノードで書き換えられる(書き換えルールについては後述する)。

20

【 0 2 2 7 】

なお、アドレスタイプとは、そのアドレスが、ストレージ装置 1 0 内のメモリノードのアドレスであるか、スイッチングリレー 8 1 のアドレスであるかを、区別するためのものである。また、中継ノードアドレスは、パケット転送時に経由すべき、スイッチングリレー 8 1 と接続されたメモリノードのアドレスである。

【 0 2 2 8 】

前述したように、スイッチングリレーに接続されたメモリノードでは、パケットのヘッダ部に記録されたアドレス情報の書き換えが行われる。パケットのアドレス情報が図 2 9 に示したものである場合、書き換えルールは例えば次のようなものである。

30

【 0 2 2 9 】

1. 最終送信先ノードが自ノード以外のメモリノードで、送信元がスイッチングリレーである場合、一時送信先を最終送信先ノードに、一時送信元を自ノードに変更する。

【 0 2 3 0 】

2. 最終送信先ノードがスイッチングリレーである場合、一時送信先をスイッチングリレーに、一時送信元を自ノードに変更する。

【 0 2 3 1 】

3. 最終送信先が自ノード以外のメモリノードで、送信元もメモリノードである場合、一時送信先、一時送信元共に変更しない。

【 0 2 3 2 】

例えば、図 2 7 に示したシステムで、クライアントがストレージ装置 1 0 にファイルを書き込む手続きは次のようになる。

40

【 0 2 3 3 】

図 3 0 は、図 2 7 に示したストレージシステムにおける書き込み動作を示す図である。

【 0 2 3 4 】

まず、クライアントは、ファイルとファイル ID とをゲートウェイサーバ 2 1 A に送信する((1) 参照)。ファイル ID は、ファイルを一意に特定できる識別子である。

【 0 2 3 5 】

次に、ゲートウェイサーバ 2 1 A は、ファイルを規定サイズ of データに分割し、分割した各データに分割データ ID を割り振る。さらに、ファイル ID と分割データ ID をファ

50

イルテーブルに書き込む。分割データIDは、分割されたデータを一意に特定できる識別子である（（2）参照）。

【0236】

次に、ゲートウェイサーバ21Aは、分割データIDの情報に基づき、分割データを書き込む書き込み先のメモリノード（書き込みノード）のアドレスを決定する（図27ではアドレス（5，9））。さらに、スイッチングリレー81に接続されたメモリノードのうち、書き込みノードに最も近いメモリノードを中継メモリノードとして、そのアドレスを算出する（図27ではアドレス（1，9））（（3）参照）。

【0237】

次に、ゲートウェイサーバ21Aは、書き込みデータに、前述したアドレス情報を含むヘッダー部等を付け加えた書き込みパケットを作成する。ここで、一時送信先アドレスは中継ノードアドレス（図27ではアドレス（1，9））、一時送信元ノードアドレスは、スイッチングリレー81のアドレス（図27では[1]）とする。その後、パケットをスイッチングリレー81に送信する（（4）参照）。

【0238】

書き込みパケットを受け取ったスイッチングリレー81は、指定された一時送信先メモリノード（中継メモリノード（1，9））にこの書き込みパケットを送信する（（5）参照）。

【0239】

スイッチングリレー81から書き込みパケットを受信した中継メモリノードは、パケットのヘッダー部を読み、最終送信先ノードがストレージ装置10内の他のメモリノードであることから、以下のようにヘッダー部を書き換えた書き込みパケットを作成する。ヘッダー部において、一時送信元アドレスを自ノードアドレスに、一時送信先ノードアドレスを、最終送信先である書き込みノードのアドレス（図27では（5，9））に変更する（（6）参照）。

【0240】

その後、書き込みパケットを隣接したメモリノードに転送する。ストレージ装置10内で適切な転送を繰り返すことにより、書き込みパケットを書き込みノード（図27ではメモリノード（5，9））に到達させる（（7）参照）。

【0241】

書き込み先のメモリノードでは、受け取ったパケットの書き込みデータ、パケット送信元や中継ノードのアドレス等を自ノードのメモリ16に書き込む（（8）参照）。その後、書き込み完了報告のパケットを作成し、上記と逆の経路でゲートウェイサーバ21Aに返信する（（9）参照）。書き込み完了報告のパケットでは、ヘッダー部のアドレス情報の最終送信先ノードはスイッチングリレーとし、送信元は書き込みノード、中継ノードはゲートウェイサーバ21Aから書き込みノードにパケット転送を行った際と同じメモリノード、一時送信先ノードは中継ノード、一時送信元ノードは書き込みノードとする。前記パケットは、中継メモリノードにおいて、一時送信先をスイッチングリレーに、一時送信元を中継ノードに、ヘッダー部のアドレス情報が書き換えられる（（10）参照）。

【0242】

全分割データの書き込みが終わった後、ゲートウェイサーバ21Aは、クライアントにファイル書き込みの終了報告を行う。

【0243】

なお、クライアントがストレージ装置10にファイルを書き込む場合だけでなく、書き込まれたデータの読み出し・消去やメモリの空き容量の確認などのため、読み出し・消去命令や空き容量の返信命令等をストレージ装置10のメモリノードに送る場合も、前述した手続きに従ってスイッチングリレー81に接続されたメモリノードを介してパケットを転送することができる。さらに、書き込みや消去等の命令実施完了報告、メモリノードから読み出したデータなどをメモリノードからゲートウェイサーバ21Aに送信する際などにも、同様に、前述した手続きに従ってスイッチングリレー81に接続されたメモリノード

10

20

30

40

50

ドを介してパケットを転送することができる。これにより、転送時間を短縮することができる。

【0244】

なお、ここまでは、図27に示したように、左端のメモリノードのみがスイッチングリレー81と接続されたストレージシステムについて説明したが、図31に示すように、全ての外周部のメモリノードがスイッチングリレー81, 82, 83, 84と接続されたストレージシステムにおいても、スイッチングリレーを導入しないシステムより、パケットの平均転送回数が少なく、転送時間を短くすることができる。図31に示すストレージシステムでは、ストレージ装置の左端に配置されたメモリノードにスイッチングリレー81が接続され、上端に配置されたメモリノードにスイッチングリレー82、右端に配置されたメモリノードにスイッチングリレー83、下端に配置されたメモリノードにスイッチングリレー84がそれぞれ接続されている。また、全てのスイッチングリレーは、アダプタ22Aを介してゲートウェイサーバ21Aと接続されている。

10

【0245】

さらに、図32に示すように、隣接ノードと接続された入出力ポート以外の入出力ポート（非隣接ポート）を有する複数のメモリノードが、非隣接ポートによってスイッチングリレー81と接続されたストレージ装置を含むストレージシステムにおいても、スイッチングリレーを導入しないシステムより、パケットの平均転送回数が少なく、転送時間を短くすることができる。図32に示すストレージシステムでは、ストレージ装置10内のメモリノード(3, 2)、(3, 7)、(8, 2)、(8, 7)にスイッチングリレー81

20

【0246】

さて、前述したように、スイッチングリレーに接続されているメモリノードを介してパケット転送を行う場合、パケット転送時間（転送回数）が最小となるメモリノードを中継メモリノードとして選択することが重要である。以下では、その中継メモリノードのアドレス算出について説明する。

【0247】

まず、図27に示したように、左端のメモリノード全てがスイッチングリレー81に接続されたストレージ装置10とゲートウェイサーバ21Aとの間でデータ転送を行う場合、ストレージ装置10内のアドレス(x, y)のパケット送信先/送信元メモリノードにとって、パケット転送距離（転送回数）が最小となる、スイッチングリレー81に接続されたメモリノードは、アドレス(1, y)のメモリノードである。

30

【0248】

角に配置されたメモリノードのアドレスが(a0, b0)で、横c0個、縦d0個のメモリノードから成るストレージ装置において、外周部の全メモリノードがスイッチングリレーに接続されている場合、ストレージ装置10内のアドレス(x, y)のメモリノードがゲートウェイサーバ21Aとの間でパケット転送を行うとき、転送時間（転送回数）が最小となる中継メモリノードのアドレスは、以下の通りである。ここで、図31に示すように、スイッチングリレーと接続されたメモリノードのアドレスは、(a0, y)、(a0 + c0, y)、(x, b0)、(x, b0 + d0)で表わされる(xはa0 ~ a0 + c0の任意の整数値、yはb0 ~ b0 + d0の任意の整数値)。

40

【0249】

$\min(x - a_0, a_0 + c_0 - x) \min(y - b_0, b_0 + d_0 - y)$ かつ $\min(x - a_0, a_0 + c_0 - x) = x - a_0$ の場合は $(x - a_0, y)$;

$\min(x - a_0, a_0 + c_0 - x) \min(y - b_0, b_0 + d_0 - y)$ かつ $\min(x - a_0, a_0 + c_0 - x) = a_0 + c_0 - x$ の場合は $(a_0 + c_0 - x, y)$;

$\min(x - a_0, a_0 + c_0 - x) \min(y - b_0, b_0 + d_0 - y)$ かつ $\min(y - b_0, b_0 + d_0 - y) = y - b_0$ の場合は $(x, y - b_0)$;

$\min(x - a_0, a_0 + c_0 - x) \min(y - b_0, b_0 + d_0 - y)$ かつ $\min(y - b_0, b_0 + d_0 - y) = b_0 + d_0 - y$ の場合は $(x, b_0 + d_0 - y)$ 。

50

【0250】

なお、 $\min(x, y)$ は、二つの引数 x, y のうち、小さい方を与える関数である。

【0251】

他方、図32に示すように、スイッチングリレー81と接続されたメモリノードのアドレスが $(am + b, cn + d)$ (m, n は整数) で与えられる場合 (図32では、 $a = 5$ 、 $b = 3$ 、 $c = 5$ 、 $d = 2$ 、 $m = 0$ または 1 、 $n = 0$ または 1)、任意のメモリノード (アドレス (x, y)) に対して、パケット転送時間 (転送回数) が最小となる、スイッチングリレー81に接続されるメモリノード (中継メモリノード) は、 $(a(\text{round}(x/a, 0) + b, c(\text{round}(y/c, 0) + d))$ で与えられる。ここで、 $\text{round}(u, 0)$ は、数値 u を小数点で四捨五入する関数である。

10

【0252】

以上説明したように第8実施形態によれば、ストレージ装置内でメモリノード間のパケット転送回数を低減させることができる。その他の構成及び効果は、前述した第1実施形態と同様である。

【0253】

[第9実施形態]

第9実施形態は、複数のデータを複数のメモリノードに保存するデータ処理において、ゲートウェイサーバとのデータ転送時間が互いに異なるメモリノードを、複数データの保存先として選択するデータ処理手順を備える。

【0254】

[1] ストレージシステムの構成

図33Aは、第9実施形態のストレージシステムの構成を示す図である。

20

【0255】

図示するように、ストレージシステムは、複数のメモリノード11を含むストレージ装置10、及びストレージ装置10に接続された複数のゲートウェイサーバ21Aを備える。ストレージ装置10は、図1に示したように、データ転送機能を有するメモリノードが相互に複数接続された構成を有する。ストレージ装置10の外周部に配置されたメモリノード(1, 4)には、アダプタ22Aを介してゲートウェイサーバ21Aが接続されている。

【0256】

[2] ストレージシステムのデータ処理方法

第9実施形態のストレージシステムにおけるデータ処理手順について説明する。

30

【0257】

図33A～図33D及び図34A～図34Eは、ストレージ装置10がアダプタ22Aを介してゲートウェイサーバ21Aと接続されたストレージシステムである。ストレージ装置は、隣接したメモリノード11が相互接続された複数のメモリノードから構成される。

【0258】

各メモリノード11は、送信されたパケットが自ノード宛であれば受け取り、他のメモリノードが宛先の場合は、隣接するメモリノードにパケットを転送する。このデータ転送機能により、ゲートウェイサーバ21Aと指定メモリノードとの間でパケット通信を行うことができる。

40

【0259】

但し、パケット転送に必要な転送回数は、メモリノードにより異なる。例えば、ゲートウェイサーバ21Aとのパケット通信に必要な最低転送回数は、アドレス(1, 4)のメモリノードは0回、アドレス(1, 5)、(2, 4)、(1, 3)のメモリノードは1回、アドレス(1, 6)、(2, 5)、(3, 4)、(2, 3)、(1, 2)のメモリノードは2回である。図33A～図33D及び図34A～図34Eでは、同一の最低転送回数を持つメモリノードを同一のハッチングで表している。

【0260】

50

またここでは、ストレージ装置 10 のメモリノード内のパケット転送時間とメモリノード間のパケット送信時間が、メモリノードに依らず一定とする。この場合、総パケット転送時間はどのメモリノード間でも同じである。なお、パケット転送時間は、入力ポートにパケットを受信した後、パケットのヘッダ部に記録されたアドレスから自ノード宛かどうかを判断し、出力ポートからパケットを出力するまでの時間である。総パケット転送時間は、メモリノードがパケットを受け取ってから隣接ノードにパケットを送信し、そのパケットが隣接ノードに到達するまでの時間である。

【0261】

さて、一つのファイルを 3 つに分割し、分割後のデータの順番に合わせて $ID = 1$ 、 $ID = 2$ 、 $ID = 3$ と ID を付与した 3 つのデータを、3 つのメモリノードにそれぞれ保存する。その後、メモリノードから 3 つのデータを読み出して、ゲートウェイサーバに送信するデータ処理について考える。

10

【0262】

データを保存する 3 つのメモリノードが、ゲートウェイサーバとのパケット通信に必要な最低転送回数が同じになるメモリノードである場合と、最低転送回数が全て異なるメモリノードの場合のそれぞれについて、読み出されたデータのストレージ装置内の転送過程を比較する。なお、読み出し命令は一斉に行われ、読み出されたデータの送信は同時に開始されること、従って 3 つのデータの転送も全く同じタイミングで進行するものとする。

【0263】

ゲートウェイサーバ 21A とのパケット通信に必要な最低転送回数が 3 回であるメモリノード (1, 1)、(2, 2)、(3, 5) に保存された 3 つのデータが、各メモリノードから読み出された後、ゲートウェイサーバ 21A に向けて転送されていく過程の一例を図 33A ~ 図 33D に示す。

20

【0264】

図 33A は、各メモリノード内のメモリからデータが読み出された直後 (転送前) のストレージ装置 10 の状態を示し、図 33B は、図 33A から一定時間経過し、3 つのデータ全てが、保存されていたメモリノードの隣接ノードに転送された状態を示す。図 33B に示すように、図 33A においてメモリノード (1, 1)、(2, 2)、(3, 5) に保存されていた 3 つのデータは、メモリノード (1, 2)、(2, 3)、(2, 5) にそれぞれ転送されている。

30

【0265】

図 33C、図 33D は、それぞれデータの転送が 2 回、3 回行われた状態を示す。図 33C に示すように、図 33B においてメモリノード (1, 2)、(2, 3)、(2, 5) に保存されていた 3 つのデータは、メモリノード (1, 3)、(2, 4)、(1, 5) にそれぞれ転送されている。さらに、図 33D に示すように、図 33C においてメモリノード (1, 3)、(2, 4)、(1, 5) に保存されていた 3 つのデータは、メモリノード (1, 4) にそれぞれ転送されている。なお図 33A ~ 図 33D は、保存されていたメモリノードからゲートウェイサーバ 21A へ向けてデータ転送を行う過程であるため、時間の経過と共に、ゲートウェイサーバ 21A とのパケット通信に必要な、最低転送回数が少ないメモリノードにデータが転送されている。

40

【0266】

さて、図 33A ~ 図 33C に示したストレージ装置 10 では、3 つのデータはそれぞれ別のメモリノードに存在しているが、図 33D に示したストレージ装置 10 では、全てのデータがアドレス (1, 4) のメモリノード上に位置している。これは、3 つのデータが同時にアドレス (1, 4) のメモリノードに到達した状態を表している。

【0267】

但し、メモリノードの一時保存用メモリ (入力ポートバッファ) の記憶容量サイズが、複数のデータを保存できるほど大きくない場合には、メモリノード (1, 4) では一度に一つのデータしか受け入れることができない。一時保存用メモリは、自ノード宛でないデータを受信した場合に、そのデータを転送するまでの間、保管するためのメモリである。

50

【 0 2 6 8 】

この場合には、受信した一つのデータのアダプタ 2 2 A への転送が終了するまで、他のデータはアドレス (1 , 4) のメモリノードに転送できず、隣接ノードで待機する必要がある。

【 0 2 6 9 】

メモリノードの一時保存用メモリのサイズが十分大きく、3つのデータを同時に一時保存できる場合においても、アドレス (1 , 4) のメモリノードからゲートウェイサーバ 2 1 A に3つのデータを一度に送信することは通常できないため、一つのデータがゲートウェイサーバ 2 1 A に送信されるまでの間、他のデータはアドレス (1 , 4) のメモリノードで待機する必要がある。このようなデータ転送の待機が生じると、全データの転送に要する時間が増加する。

10

【 0 2 7 0 】

また、図 3 3 A ~ 図 3 3 D からわかるように、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が同じであるメモリノードの数は、転送回数が少なくなるにつれて減少する。従って、図 3 3 A に示すように、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が同じで異なるメモリノードに複数データを保存し、同時に読み出して転送を開始すると、時間が経過してデータがゲートウェイサーバ 2 1 A に近づくほど、転送先のメモリノードの数が減少し、複数データが同一のメモリノードに転送される確率が高まる。このように、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が同じであるメモリノードに複数データを保存すると、読み出し後の転送時に無駄な待機時間が発生しやすくなる。

20

【 0 2 7 1 】

図 3 4 A ~ 図 3 4 E、図 3 5 A、及び図 3 5 B を参照して、このような転送時の待機時間を低減するためのデータ処理手順を説明する。

【 0 2 7 2 】

図 3 4 A ~ 図 3 4 E は、メモリノードから読み出されたデータがゲートウェイサーバ 2 1 A に向けて転送されていく過程の一例を示す図である。図 3 4 A ~ 図 3 4 E は、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が 1 回であるメモリノード (1 , 5) に ID = 2 のデータを、最低転送回数が 2 回であるメモリノード (1 , 2) に ID = 1 のデータを、3 回であるメモリノード (4 , 4) に ID = 3 のデータを保存させた後、これらデータを同時に読み出してゲートウェイサーバ 2 1 A に転送する過程を示している。

30

【 0 2 7 3 】

この場合も、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が少ないメモリノードへとデータは転送されている。しかし、図 3 4 A ~ 図 3 4 E の場合、全ての過程で、図 3 3 D に示したような同一メモリノードへの複数データの転送が生じていない。これは、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が異なるメモリノードにデータを保存したためである。この場合にも、最低転送回数が少ないメモリノードへとデータ転送が行われるが、最低転送回数が異なるメモリノードから全データの転送が開始されるため、どの時間においても各データは最低転送回数が異なるメモリノードに位置する。

40

【 0 2 7 4 】

このように、ゲートウェイサーバ 2 1 A との packet 通信に必要な最低転送回数が異なるメモリノードにデータを保存することにより、読み出しデータの転送時における複数データの同一メモリノードへの同時転送が回避され、データ転送時間を短縮することができる。

【 0 2 7 5 】

なお、当然ではあるが、ゲートウェイサーバとの packet 通信に必要な最低転送回数が異なるメモリノードに保存したデータを読み出してゲートウェイサーバに転送する場合、データを保存したメモリノードがゲートウェイサーバに近い順番でゲートウェイサーバに

50

データが到達する。図34Aでは、ゲートウェイサーバ21Aに最も近いメモリノードにID=2のデータを、次にサーバ21Aに近いメモリノードにID=1のデータを、最もサーバ21Aから遠いメモリノードにID=3のデータを保存しているため、これらメモリノードから読み出し・転送されたデータは、図34C~図34Eに示したように、ID=2のデータ、ID=1のデータ、ID=3のデータの順番でゲートウェイサーバ21Aに到達する。

【0276】

前述したように、これらデータは一つのファイルを3分割したものであり、分割後のデータの順番に合わせてID=1、ID=2、ID=3とIDを付与されている。このため、分割データからファイルを再構成するためには、ゲートウェイサーバ21Aに到達したデータをID順に入れ替える必要がある。

10

【0277】

他方、図35Aでは、ID=1のデータをゲートウェイサーバ21Aからの最低転送回数1回のメモリノードに、ID=2のデータを最低転送回数2回のメモリノードに、ID=3のデータを最低転送回数3回のメモリノードに保存している。このため、これらのデータをメモリノードから読み出し、ゲートウェイサーバ21Aに転送すると、図35Bに示すように、ID=1のデータ、ID=2のデータ、ID=3のデータの順番でゲートウェイサーバ21Aに到達する。従って、ファイルの再構成を行う際に、データの並べ替えを行う必要が無い。

【0278】

20

このように、データの読み出し順番が意味を持つ複数のデータを、ゲートウェイサーバとのパケット通信に必要な最低転送回数が異なるメモリノードに保存する場合、その順番に合わせて、ゲートウェイサーバに近いメモリノードから遠いメモリノードへ順に保存する。これにより、ゲートウェイサーバ21Aに転送されたデータの並び替え作業を省略することができる。

【0279】

さて、全て等価なメモリノードから構成されたストレージ装置の単一メモリノードとゲートウェイサーバが、アダプタを介して接続されているストレージシステムで、ストレージ装置に複数データを保存する手順について説明したが、図36Aに示すように、隣接ノードと接続された入出力ポート以外の入出力ポート（非隣接ポート）を有し、非隣接ポートによって相互接続された複数のメモリノードを有するストレージ装置においても、ゲートウェイサーバ21Bとのパケット通信に必要な最低転送回数が異なるメモリノードにデータを保存することで、保存データの転送時に無駄な待機時間が発生するのを回避することができる。

30

【0280】

さらに、図36B~図36Dに示すように、スイッチングリレーを介してストレージ装置10とゲートウェイサーバ21Aが接続されている場合においても、ゲートウェイサーバ21Aとのパケット通信に必要な最低転送回数が異なるメモリノードにデータを保存することで、保存データの転送時に無駄な待機時間が発生するのを回避することができる。なお、図36A~図36Dでは、ID=1、ID=2、ID=3のデータを、パケット通信に必要な最低転送回数が1回、2回、3回のメモリノードにそれぞれ保存した例を示している。

40

【0281】

なお前述したように、メモリノードのアドレスのハッシュ値とパケットIDのハッシュ値から、パケットデータを保存するメモリノードのアドレスを決める方法（コンシステントハッシング）がある。これは例えば、 $0 \sim 2^{16} - 1$ の整数値を取るID空間を考え、メモリノードとパケットのそれぞれに対して、前者はアドレス、後者はパケットIDで暗号的ハッシュ関数SHA-1を計算（計算結果は $0 \sim 2^{16} - 1$ のいずれかの整数値になる）することで、それらをID空間のいずれかのIDに割り当てる。次に、各パケットに対して、ハッシュ値のIDからID空間を時計回りにたどり、最初に行き当たった

50

メモリノードを、そのパケットを保存するメモリノードとして決定するものである。

【0282】

本方法は、高い均一性でパケットをメモリノードに分散保存できるだけでなく、メモリノードを増減させても、変更が必要なパケットが少ない(すなわち、メモリノードのスケラビリティが高い)特長を有している。しかし、この方法でパケットを保存するメモリノードを決めると、ゲートウェイサーバとの通信時に必要な転送回数が同じメモリノードにパケットが保存されうる(前述したように、これは、データ読み出し時にパケット衝突を引き起こす)。

【0283】

ゲートウェイサーバとの通信でパケット転送回数が異なるメモリノードに、N個のパケットを均一に分散させて保存するためには、例えばパケットIDのハッシュ値mod Nを計算して(modは割り算の余りを返す関数)、そのパケットを保存するメモリノードの転送回数(0 ~ (N - 1)のいずれか)を決める(計算の結果、転送回数が一致するパケットがあれば、異なる値になるように調整する)。さらに、各パケットに対して、決められた転送回数を有するメモリノードの中から、そのパケットを保存するメモリノードを前述したコンシステントハッシングにより決めれば良い。

10

【0284】

なお、このようなパケット保存先のメモリノードを決定する手順では、決められた転送回数を有するメモリノードのアドレス把握が必要である。ある転送回数を有するメモリノードのアドレスは次のように表わされる。

20

【0285】

図33A ~ 図35Bのように、ゲートウェイサーバ21Aが、アダプタ22Aを介してストレージ装置(ここでは、メモリノードのアドレスは全て正の整数値とする)の端に位置するメモリノード(1, a0)と接続されているシステムにおいて、メモリノード(1, a0)との間でパケット通信を行う場合に、最低転送回数がn回であるメモリノードのアドレスは、(1 + b, a0 + (n - b))(n - b = 0)、及び(1 + c, a0 - (n - c))(a0 - 1 - n - c = 0)で表わされる。

【0286】

図36A、図36D示したように、増加ポート(非隣接ポート)によって、隣接しないメモリノード同士やメモリノードとスイッチングリレーが接続されたストレージ装置(メモリノードのアドレスは正の整数値とする)において、増加ポートを有するノード(a0, b0)との間でパケット通信を行う場合に、最低転送回数がn回であるメモリノードのアドレスは、(a0 + c, b0 + (n - c))(n - c = 0)、(a0 + d, b0 - (n - d))(n - d = 0, b0 - 1 - n - d)、(a0 - e, b0 + (n - e))(min(a0 - 1, n) - e = 0)、(a0 - f, b0 - (n - f))(a0 - 1 - f = 0, b0 - 1 - n - f = 0)で表わされる。

30

【0287】

図36Bに示すように、ストレージ装置(全メモリノードのアドレスは正の整数値とする)の左端(アドレス(1, y))のメモリノードがスイッチングリレーと接続されているシステムにおいて、アドレス(1, a0)の中継メモリノードとの間でパケット通信を行う場合に、最低転送回数がn回であるメモリノードのアドレスは、(1, a0 + n)で表わされる。

40

【0288】

なお、データ転送時間が同じメモリノードが相互接続されたストレージ装置を含むシステムについて説明してきたが、図36Eに示すようなデータ転送時間が同じメモリノードがツリー状に接続されたストレージシステムや、図36Fに示すようなデータ転送時間の異なるメモリノードから成るストレージシステムなど、サーバとのデータ通信時間が異なるメモリノードが存在するストレージ装置を含むシステムであれば、本実施形態を適用することができる。

【0289】

50

図36Eに示すシステムは、データ転送時間が同じメモリノードがツリー状に接続されたストレージシステムであり、ゲートウェイサーバ21Aとの通信時間が同じメモリノードを同一のハッチングで示している。図36Fに示すシステムでは、メモリノード間とメモリノード～スイッチングリレー間のデータ通信速度は同じであり、さらにアドレス(a, l)、(a, m)、(a, n)、(b, l)、(b, m)、(b, n)のメモリノード間、及びこれらメモリノード～スイッチングリレー81間のデータ通信速度に対して、アドレス(,)、(,)、(,)のメモリノード間、及びこれらメモリノード～スイッチングリレー81間のデータ通信速度は2倍、アドレス(A, B)とスイッチングリレー81間の通信速度は4倍の場合を示している。なお、スイッチングリレー81とのデータ転送時間が同じメモリノードを同一のハッチングで示している。但し、メモリノード内でのデータ転送時間は、メモリノード間及びスイッチングリレー～メモリノード間のデータ転送時間より十分に小さいと仮定する。

10

【0290】

以上説明したように第9実施形態によれば、複数のメモリノードに保存したデータを各メモリノードから読み出した後、ゲートウェイサーバに送信する際の転送時間を短縮することができる。これにより、高速なデータ読み出しが可能となる。その他の構成及び効果は、前述した第1実施形態と同様である。

【0291】

[第10実施形態]

第10実施形態は、ゲートウェイサーバから複数のデータを複数のメモリノードに送信するデータ処理において、データ通信時間が長いメモリノードを宛先とするデータから、データ通信時間が短いメモリノードを宛先とするデータの順に、データ転送を行うデータ処理手順を備える。

20

【0292】

[1]ストレージシステムの構成

図37Aは、第10実施形態のストレージシステムの構成を示す図である。

【0293】

図示するように、ストレージシステムは、複数のメモリノード11を含むストレージ装置10、及びストレージ装置10に接続された複数のゲートウェイサーバ21Aを備える。ストレージ装置10は、図1に示したように、データ転送機能を有するメモリノードが相互に複数接続された構成を有する。ストレージ装置10の外周部に配置されたメモリノード(1, 4)には、アダプタ22Aを介してゲートウェイサーバ21Aが接続されている。

30

【0294】

[2]ストレージシステムのデータ処理方法

第10実施形態のストレージシステムにおけるデータ処理手順について説明する。

【0295】

図37A～図37G及び図38A～図38Eは、第9実施形態と同様に、ストレージ装置がアダプタ22Aを介してゲートウェイサーバ21Aと接続されたストレージシステムである。ストレージ装置は、隣接したメモリノード11が相互接続された複数のメモリノードから構成される。

40

【0296】

ここでは、ゲートウェイサーバ21Aから、ID=1、ID=2、ID=3の3つのデータを、それぞれアドレス(1, 5)、(1, 2)、(4, 4)のメモリノードに送信する場合の、ゲートウェイサーバ21Aからのデータ送信順番と、全データ送信に必要な時間の関係について考える。なお、メモリノードの一時保存用メモリ(入力ポートバッファ)にパケットは1つしか保存できず、また一度に1つのパケットしか、ゲートウェイサーバ21Aからアドレス(1, 4)のメモリノードに送信できないものとする。

【0297】

図37A～図37Gでは、ゲートウェイサーバ21Aに近いメモリノードを宛先とする

50

データから遠いメモリノードを宛先とするデータの順番（すなわち、ID = 1のデータ、ID = 2のデータ、ID = 3のデータの順）で、ゲートウェイサーバ21Aからストレージ装置にパケットを送信した場合のパケット転送過程を示している。

【0298】

この場合当然のことながら、図37B～図37Cに示すように、ゲートウェイサーバ21Aに最も近い、アドレス(1,5)のメモリノード宛のデータ(ID = 1)送信は最も早く終了する。

【0299】

しかし、ゲートウェイサーバ21Aから最も遠い、アドレス(4,4)のメモリノードを宛先としたデータ(ID = 3)は、図37Dに示すように、2つのデータ(ID = 1、ID = 2)が、アドレス(1,4)のメモリノードから隣接メモリノードに転送されてしまうまで、アドレス(1,4)に送信することができず、それまでの間、ゲートウェイサーバ21Aで待機する必要がある。その待機が終わり、アドレス(1,4)のメモリノードに送信されても、図37D～図37Gに示すように、アドレス(4,4)までデータが到達するためには最低3回の転送が必要である。

10

【0300】

他方、図38A～図38Eでは、ゲートウェイサーバ21Aから最も遠いメモリノードを宛先とするデータから、近いメモリノードを宛先とするデータの順番（すなわちID = 3のデータ、ID = 2のデータ、ID = 1のデータの順）で、ゲートウェイサーバ21Aからストレージ装置にパケットを送信した場合のパケット転送過程を示している。

20

【0301】

この場合においても、図38B～図38Dに示すように、ID = 3のデータとID = 2のデータがアドレス(1,4)のメモリノードから隣接ノードに転送されてしまうまで、ID = 1のデータはアドレス(1,4)のメモリノードに送信できず、ゲートウェイサーバ21Aで待機する必要がある。

【0302】

しかし、後から発信されるデータ(ID = 1)の宛先であるアドレス(1,5)のメモリノードは、ゲートウェイサーバ21Aに近く、少ない転送回数で宛先のメモリノードに到達する。他方、先にゲートウェイサーバ21Aから送信されたデータ(ID = 3)は、宛先がゲートウェイサーバ21Aから離れているが、他のデータがゲートウェイサーバ21Aで待機している間に転送が開始されているため、その分早く宛先に到達する。これらの結果、図38Eに示すように、3つのデータは、同時に宛先のメモリノードに到達する。

30

【0303】

図37A～図37Gと図38A～図38Eの比較からわかるように、宛先のメモリノードがゲートウェイサーバ21Aから離れたデータから、ゲートウェイサーバ21Aから近いデータの順番で、ゲートウェイサーバ21Aからデータ送信を行うことで、全てのデータ送信に必要な時間を最小にすることができる。

【0304】

なお、データ転送時間が同じメモリノードが相互接続されたストレージ装置を含むシステムについて説明してきたが、データ転送時間が同じメモリノードがツリー状に接続されたストレージシステム(図36E)や、データ転送時間の異なるメモリノードから成るストレージシステム(図36F)など、サーバとのデータ通信時間が異なるメモリノードが存在するストレージ装置に複数のデータを送信する場合には、本実施形態を同様に適用することができる。

40

以上説明したように第10実施形態によれば、ゲートウェイサーバとの通信時間が異なる複数のメモリノードに複数のデータを送る際、必要となる通信時間を最小にすることができる。その他の構成及び効果は、前述した第1実施形態と同様である。

【0305】

以上述べたように第1～第10実施形態によれば、メモリノードが経路指定表を管理す

50

する必要がなく、効率的にパケットを転送することができるストレージ装置、及びデータ処理方法を提供することができる。

【0306】

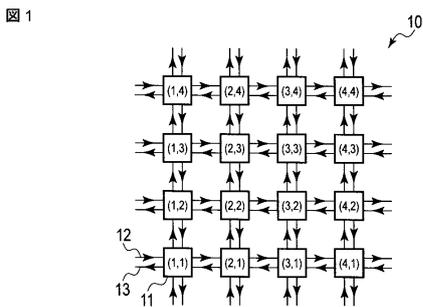
本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

【符号の説明】

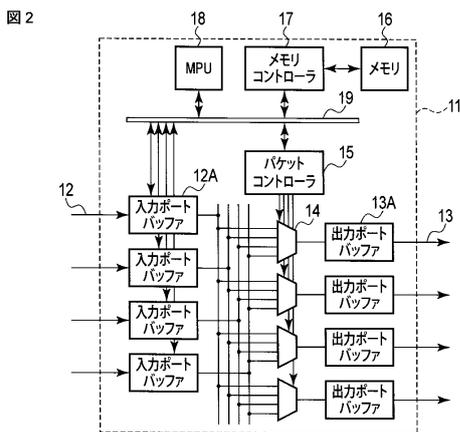
【0307】

10...ストレージ装置、11...メモリノード、12...入力ポート、12A...入力ポートバッファ、13...出力ポート、13A...出力ポートバッファ、14...セレクタ、15...パケットコントローラ、16...メモリ、17...メモリコントローラ、18...MPU、19...ローカルバス、20...ストレージシステム、21A, 21B, 21C, 21D...ゲートウェイサーバ、22A, 22B, 22C, 22D, 22E...アダプタ、31A, 31B1, 31B2...クライアント、41...制御コンピュータ、42...アドレス変換器、23, 24...配線、81, 82, 83, 84...スイッチングリレー。

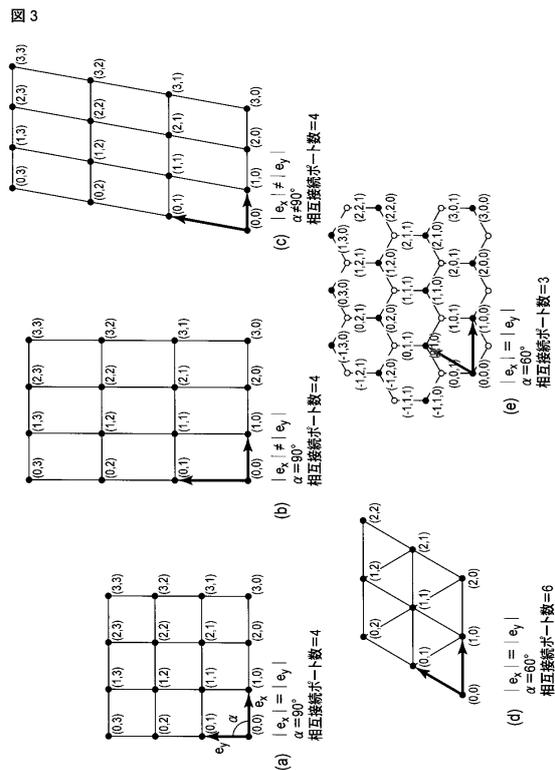
【図1】



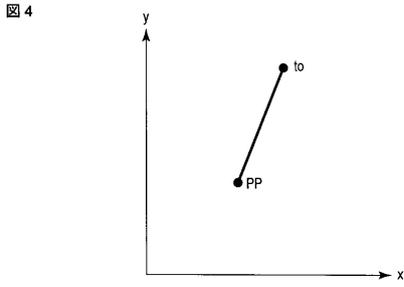
【図2】



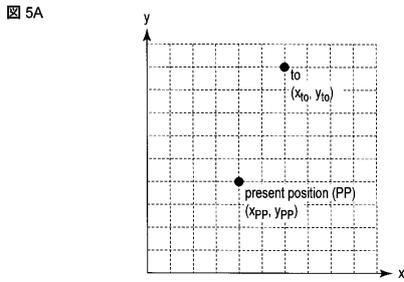
【図3】



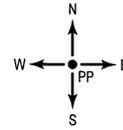
【 図 4 】



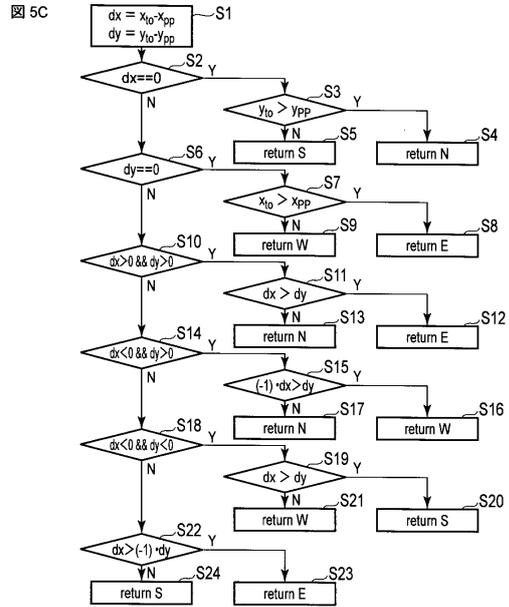
【 図 5 A 】



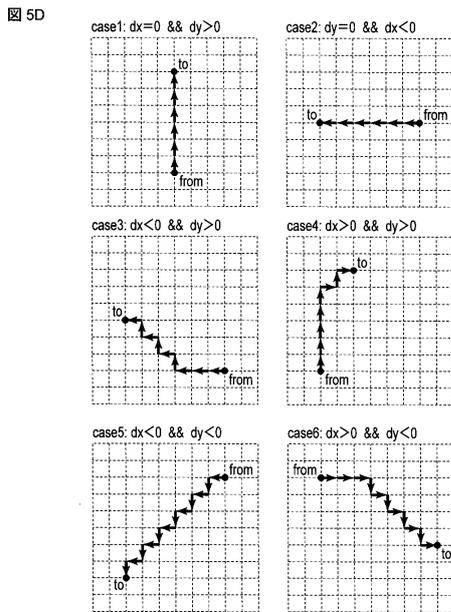
【 図 5 B 】



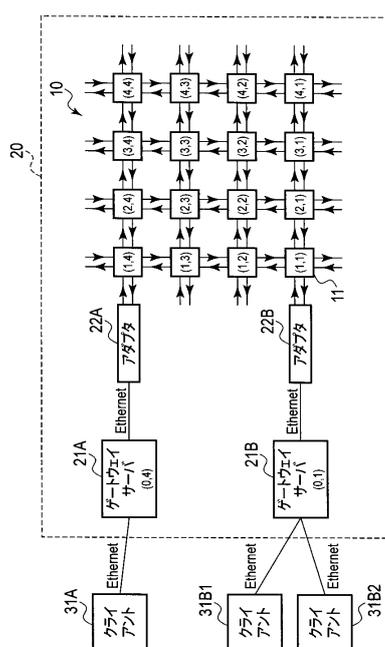
【 図 5 C 】



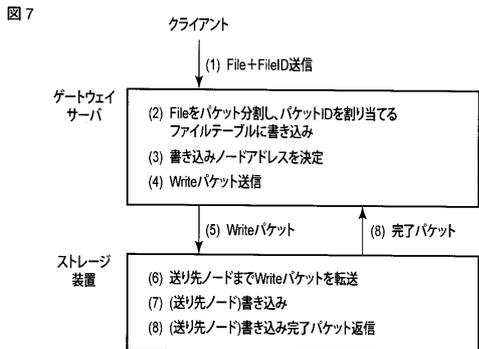
【 図 5 D 】



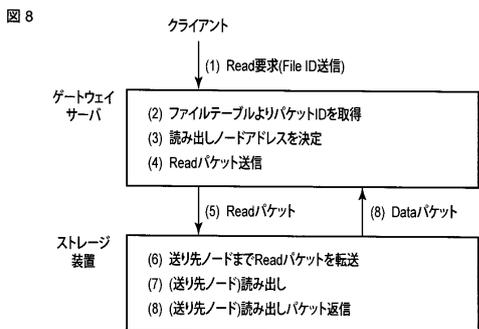
【 図 6 】



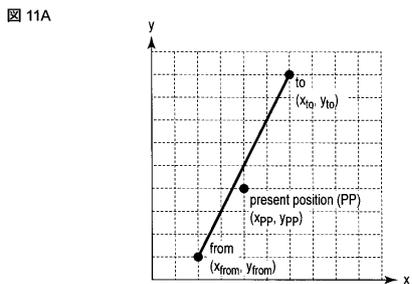
【 図 7 】



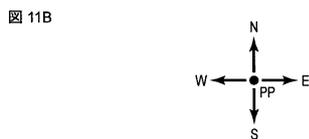
【 図 8 】



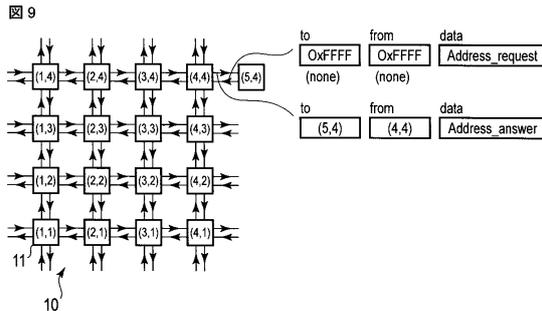
【 図 1 1 A 】



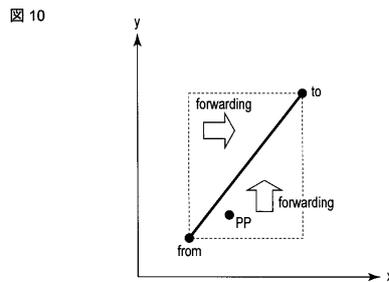
【 図 1 1 B 】



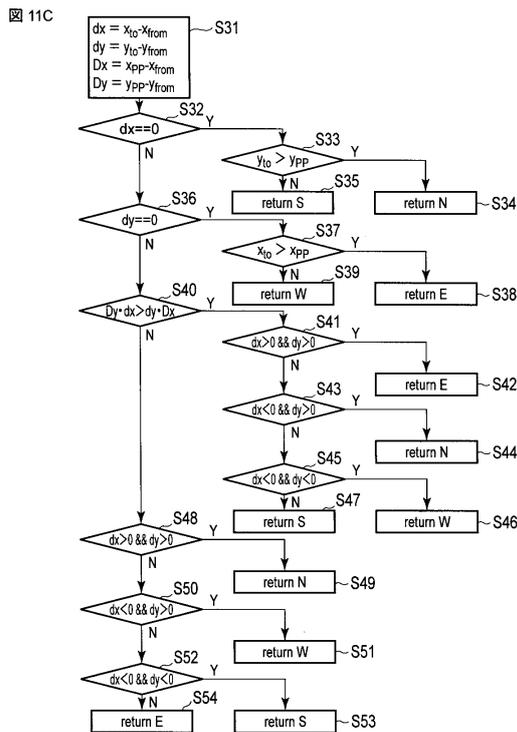
【 図 9 】



【 図 1 0 】

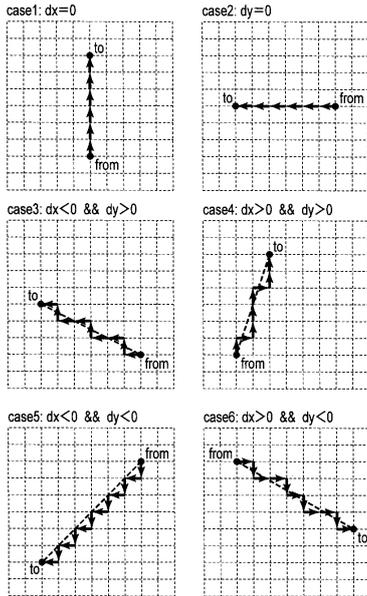


【 図 1 1 C 】



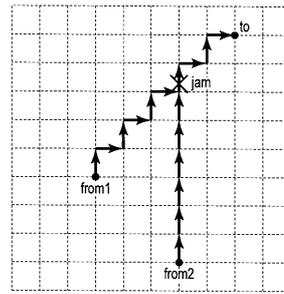
【 1 1 D 】

11D



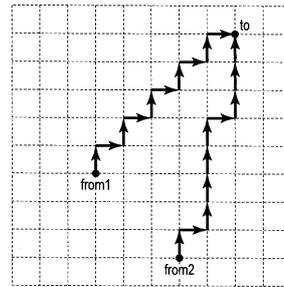
【 1 2 A 】

12A



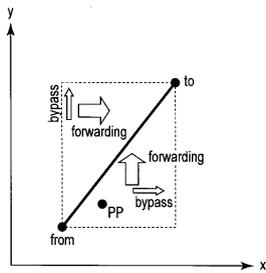
【 1 2 B 】

12B



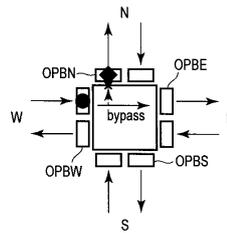
【 1 3 】

13



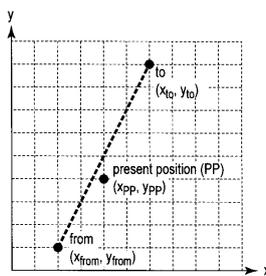
【 1 4 B 】

14B

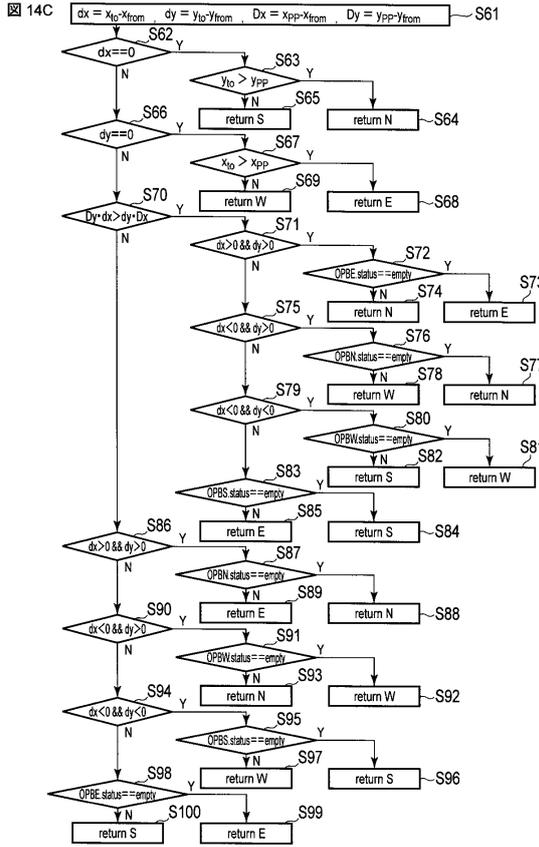


【 1 4 A 】

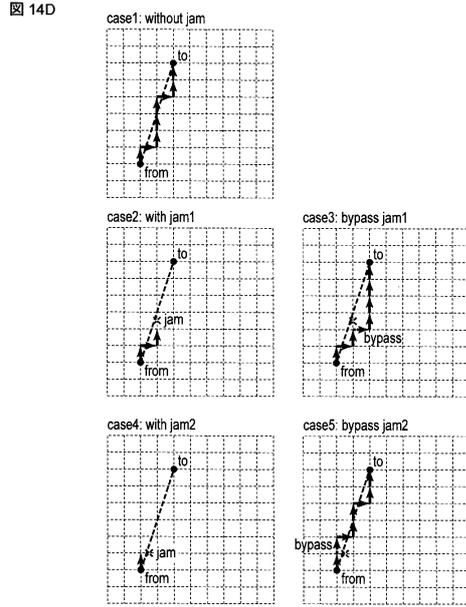
14A



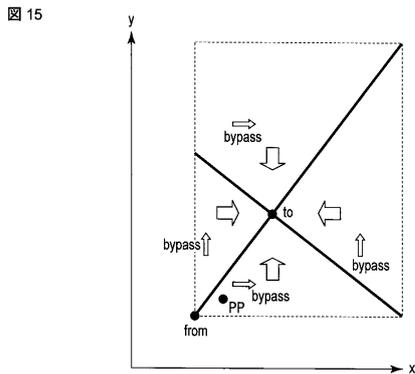
【 図 1 4 C 】



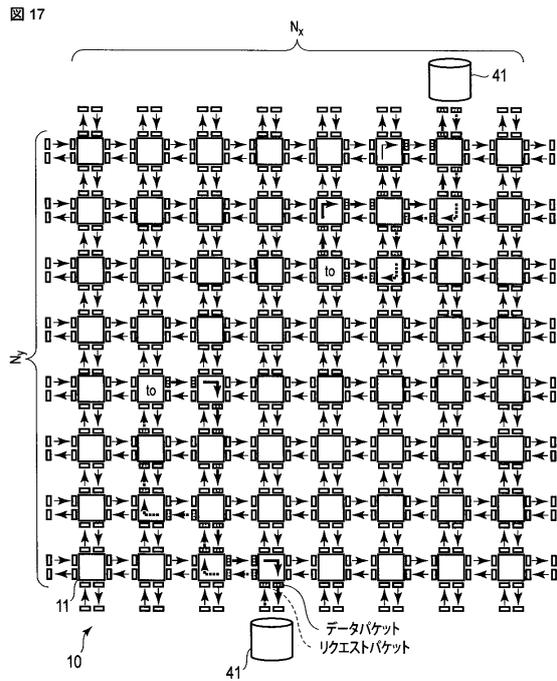
【 図 1 4 D 】



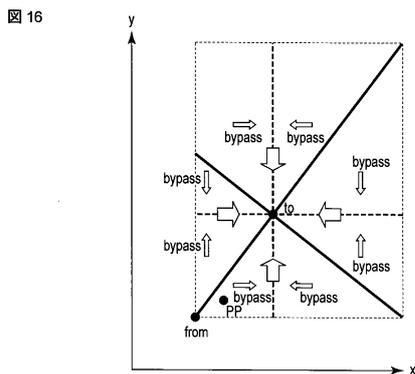
【 図 1 5 】



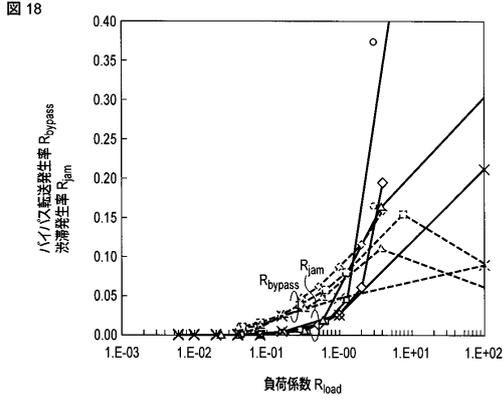
【 図 1 7 】



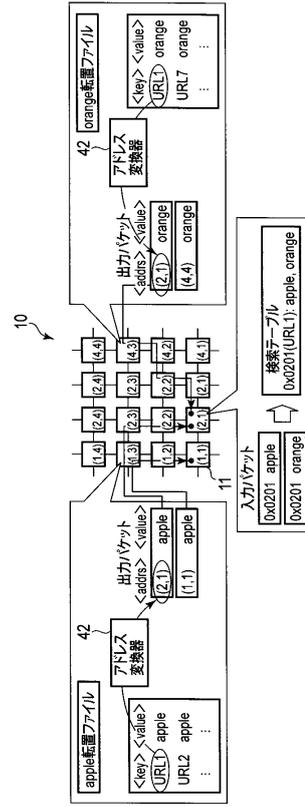
【 図 1 6 】



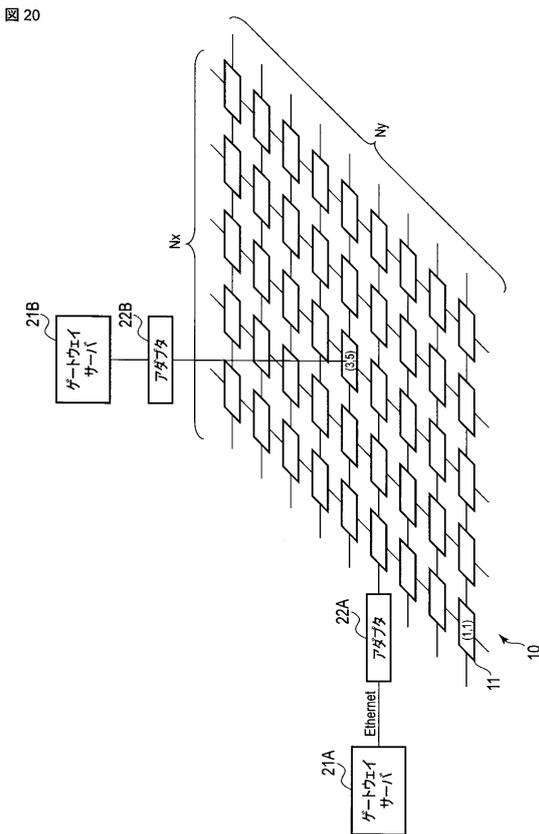
【 図 18 】



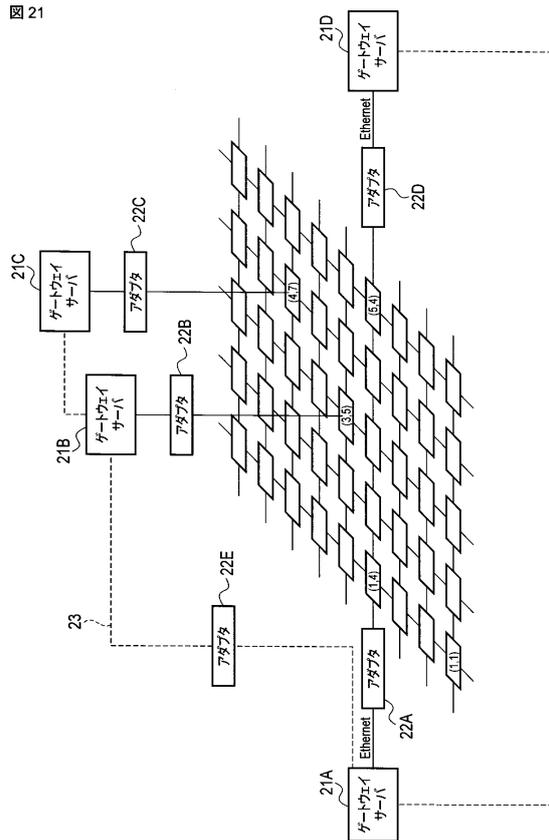
【 図 19 】



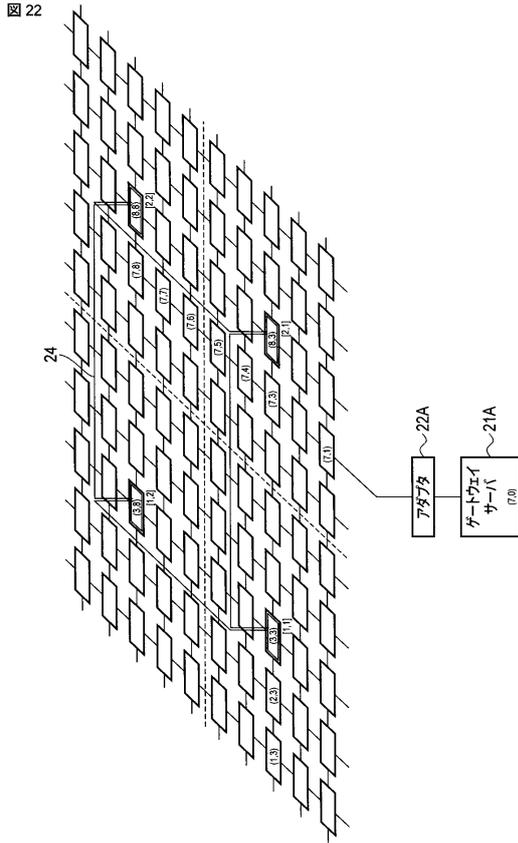
【 図 20 】



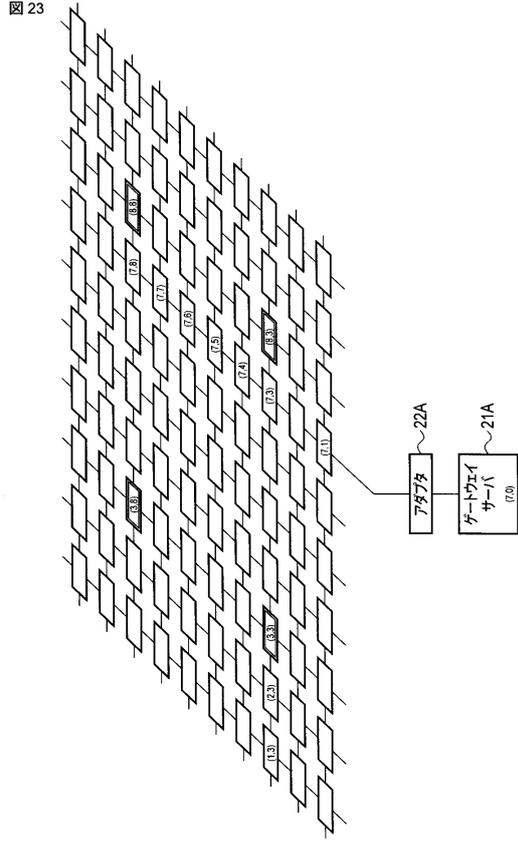
【 図 21 】



【図 2 2】



【図 2 3】

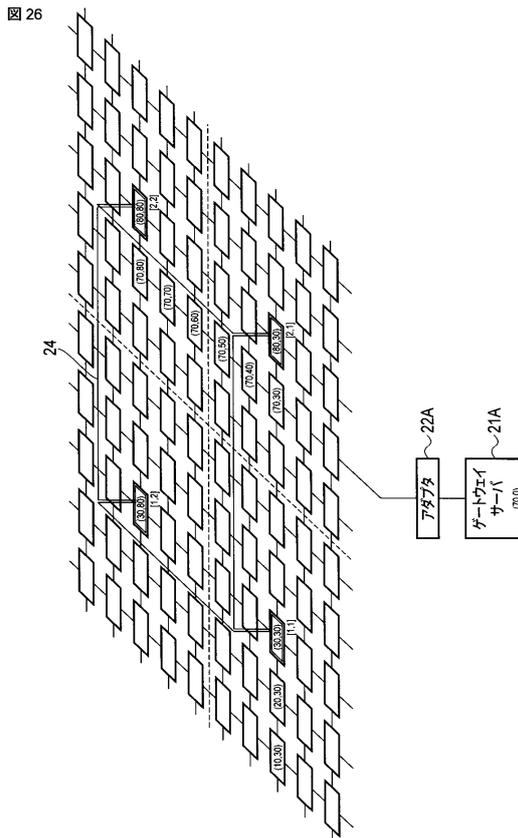


【図 2 4】

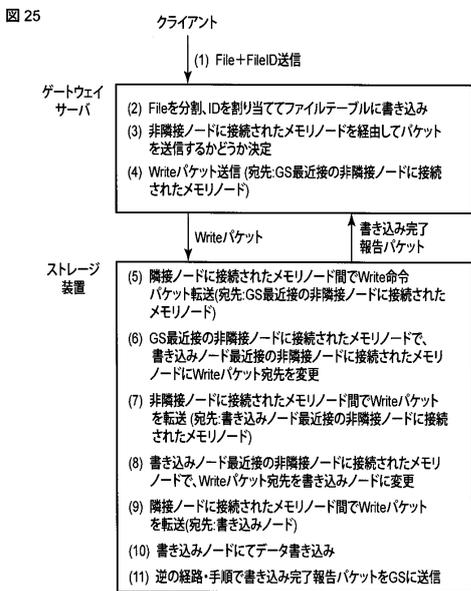
図 24

一時送信先 ノード アドレス 及び タイプ	一時送信元 ノード アドレス 及び タイプ	第1中継 ノード サブ アドレス	第2中継 ノード サブ アドレス	最終送信先 ノード メイン アドレス	送信元 ノード メイン アドレス
-----------------------------------	-----------------------------------	---------------------------	---------------------------	-----------------------------	---------------------------

【図 2 6】

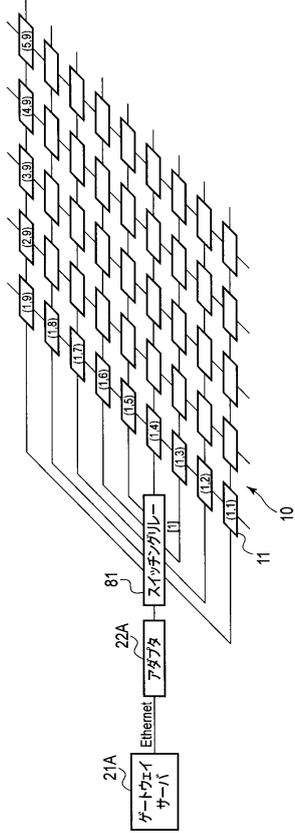


【図 2 5】



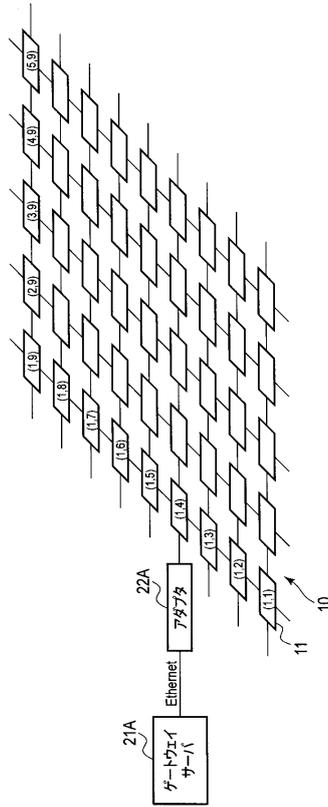
【図 27】

図 27



【図 28】

図 28



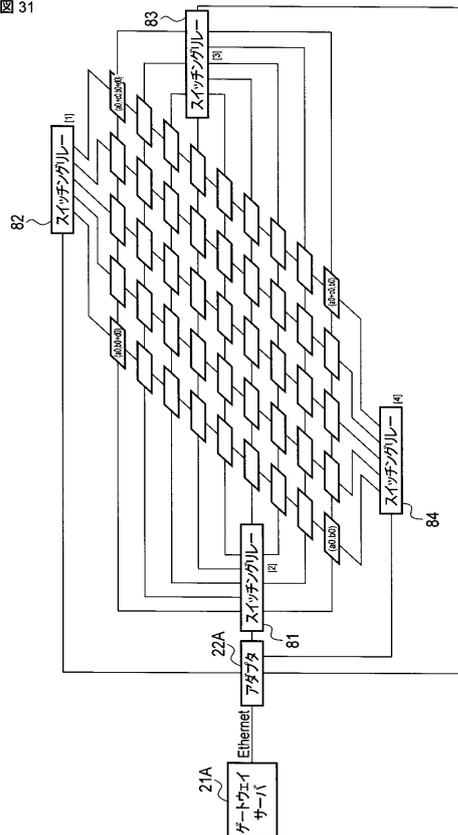
【図 29】

図 29

一時送信先 ノード アドレス 及び タイプ	一時送信元 ノード アドレス 及び タイプ	中継 ノード アドレス	最終送信先 ノード アドレス 及び タイプ	送信元 ノード アドレス 及び タイプ
-----------------------------------	-----------------------------------	-------------------	-----------------------------------	---------------------------------

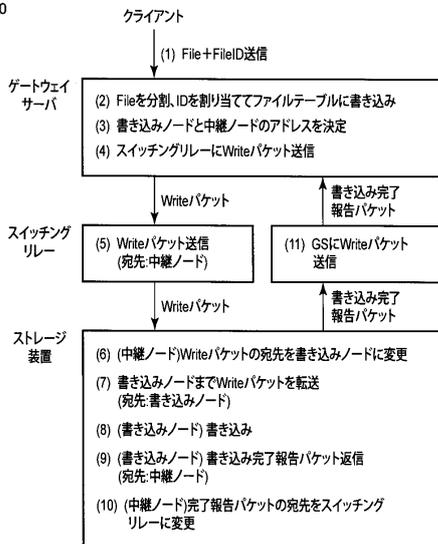
【図 31】

図 31

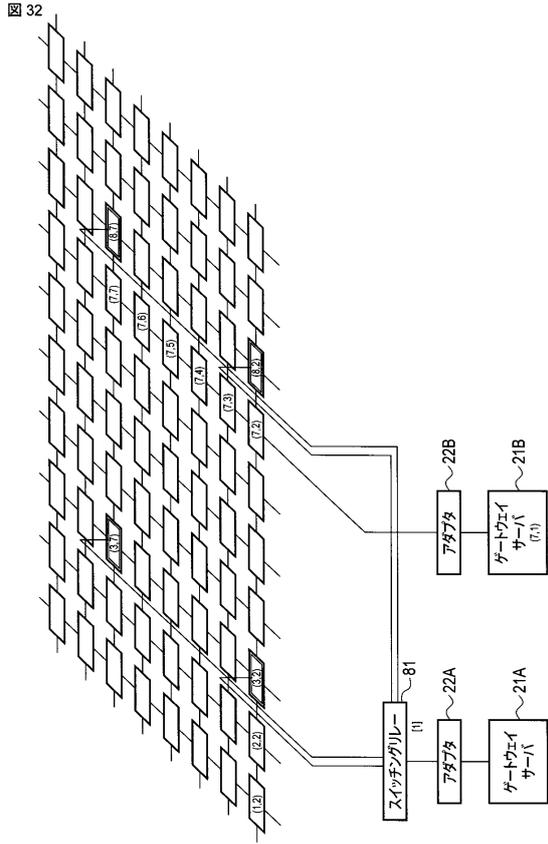


【図 30】

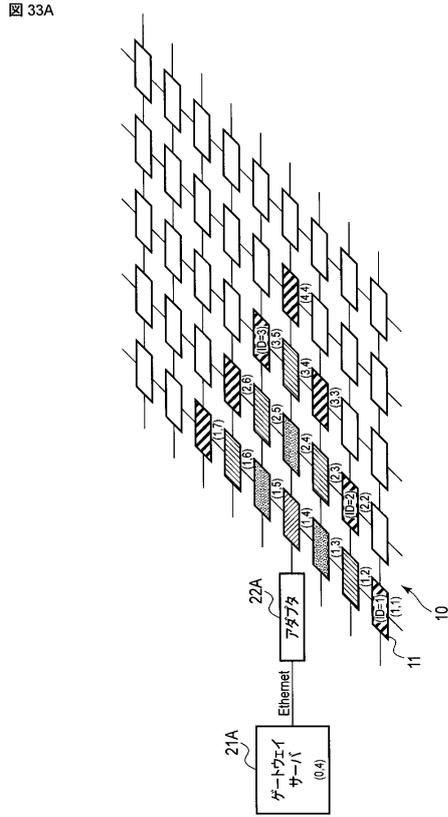
図 30



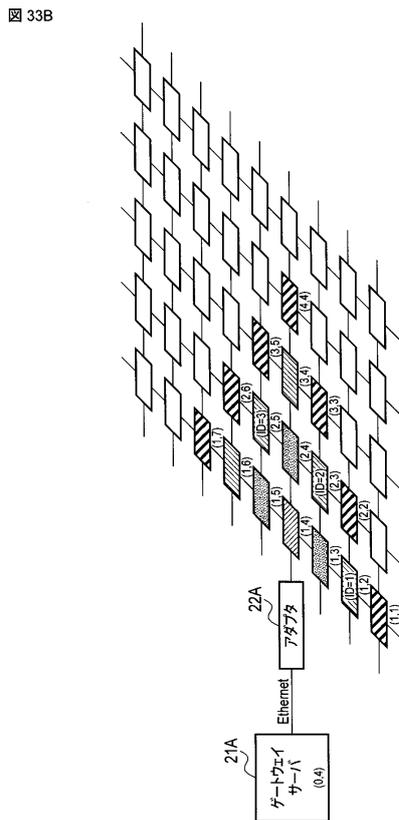
【 図 3 2 】



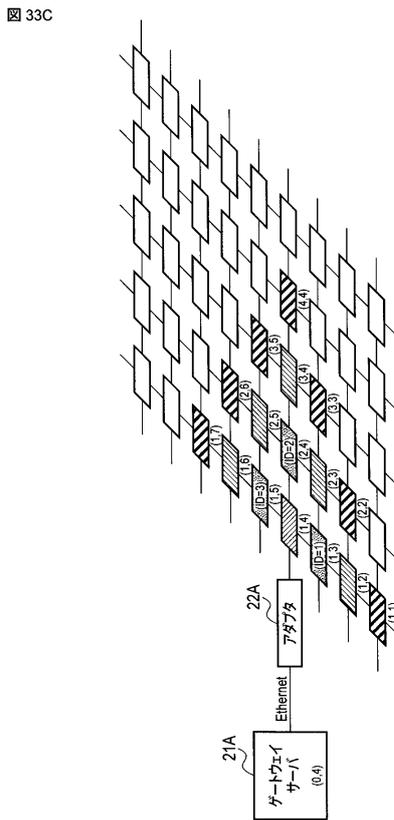
【 図 3 3 A 】



【 図 3 3 B 】

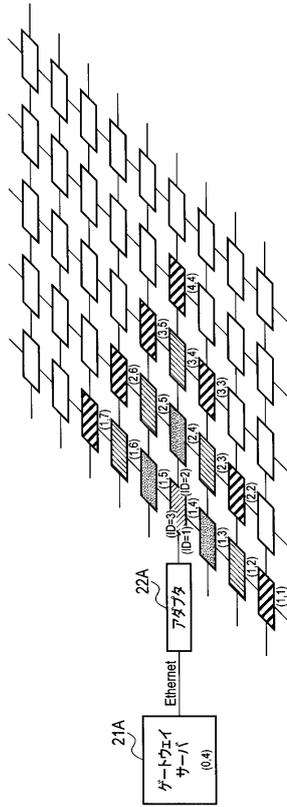


【 図 3 3 C 】



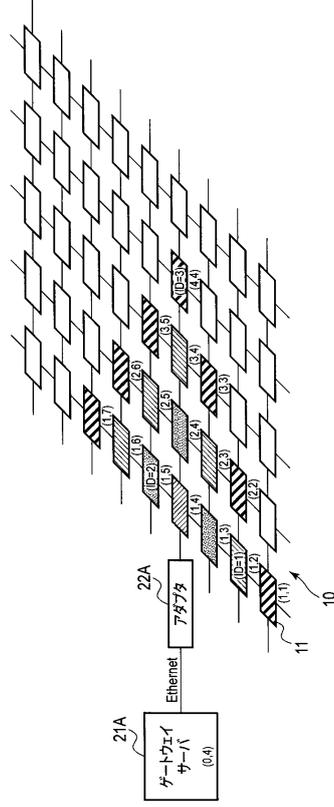
【図 33 D】

図 33D



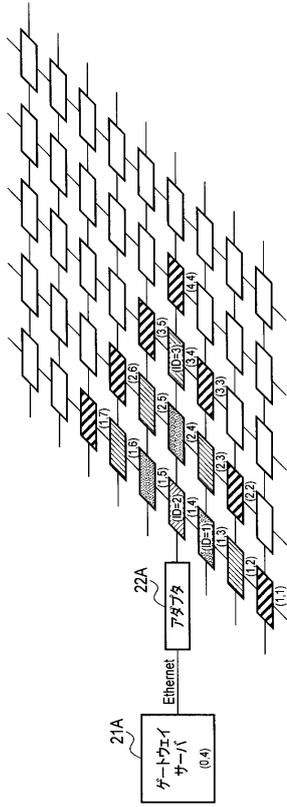
【図 34 A】

図 34A



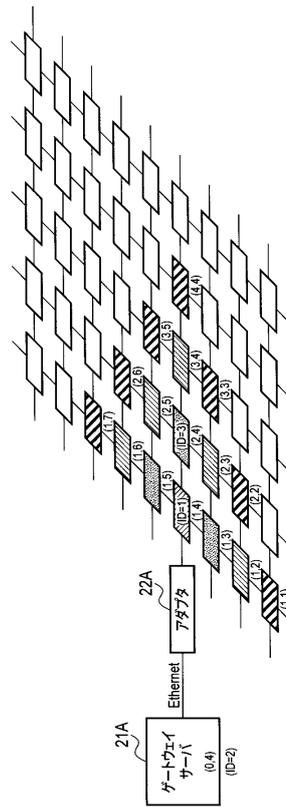
【図 34 B】

図 34B



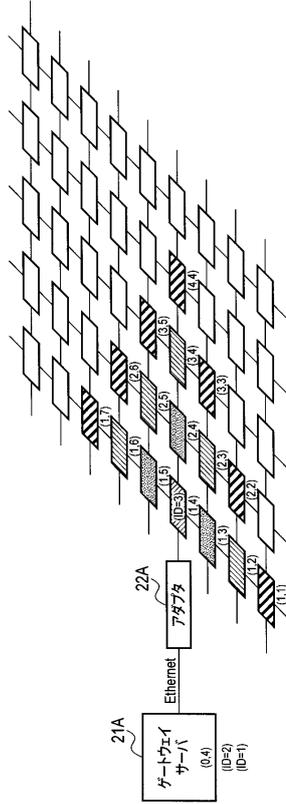
【図 34 C】

図 34C



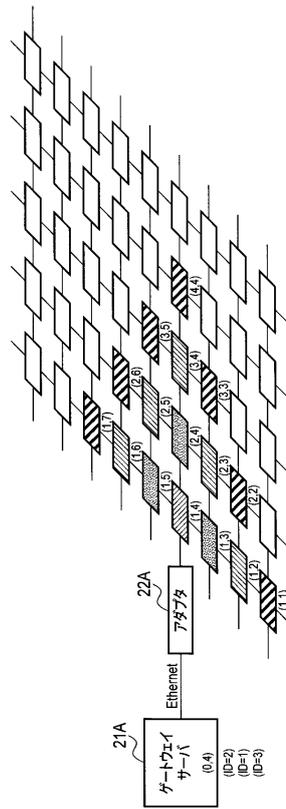
【 図 3 4 D 】

図 34D



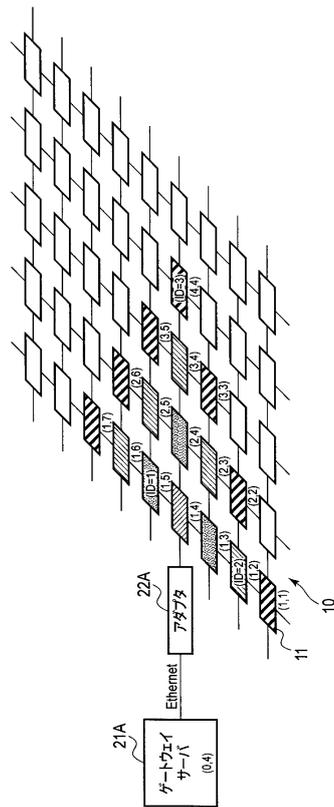
【 図 3 4 E 】

図 34E



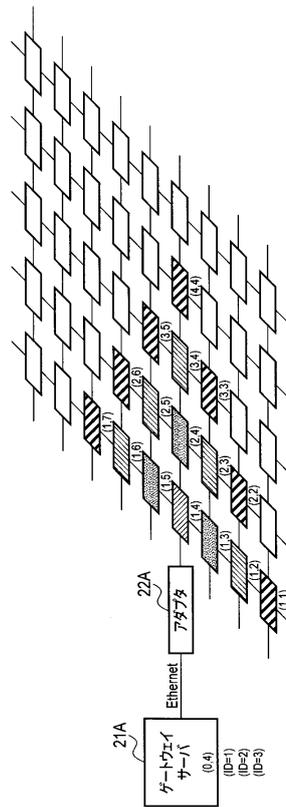
【 図 3 5 A 】

図 35A

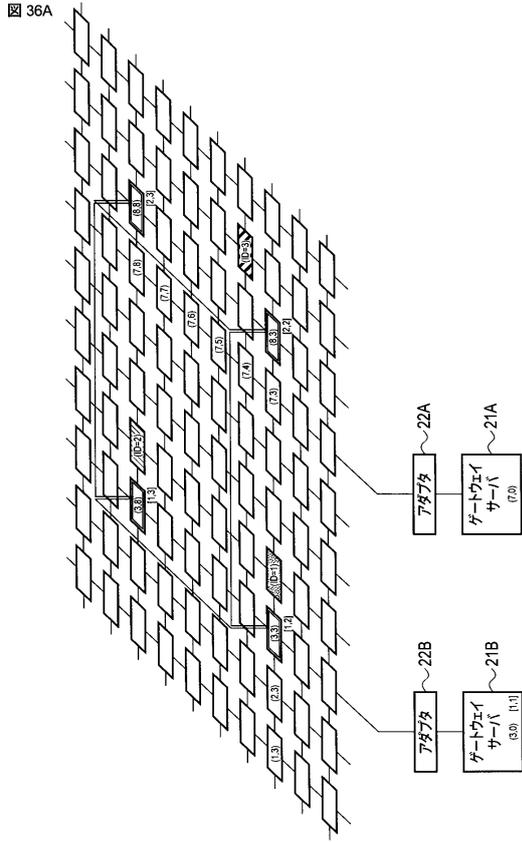


【 図 3 5 B 】

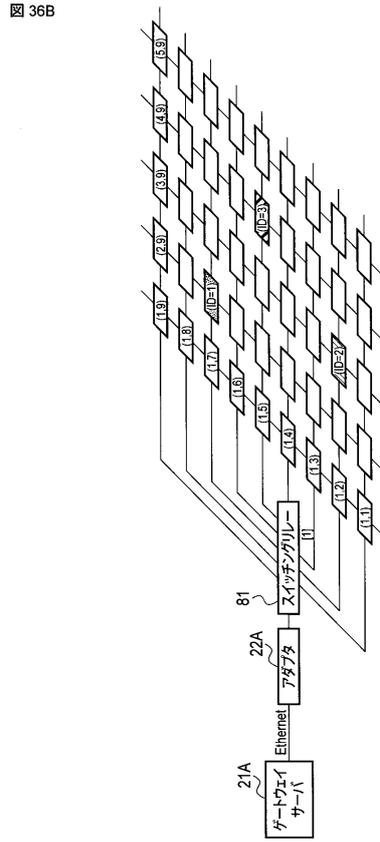
図 35B



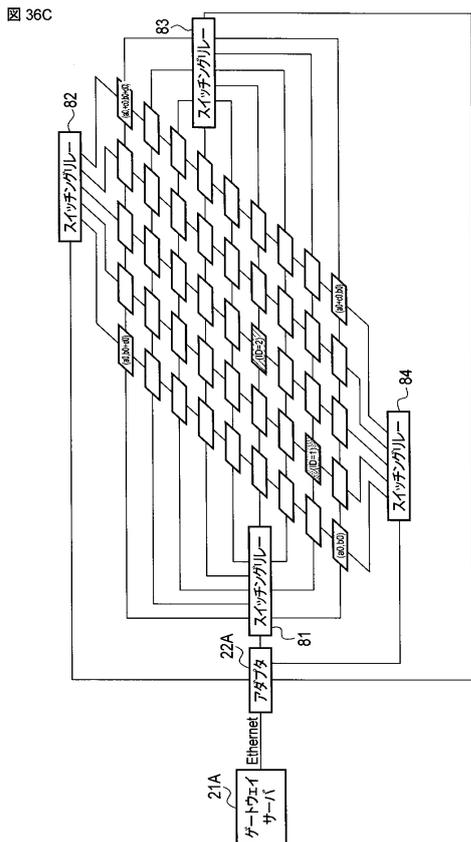
【図 36 A】



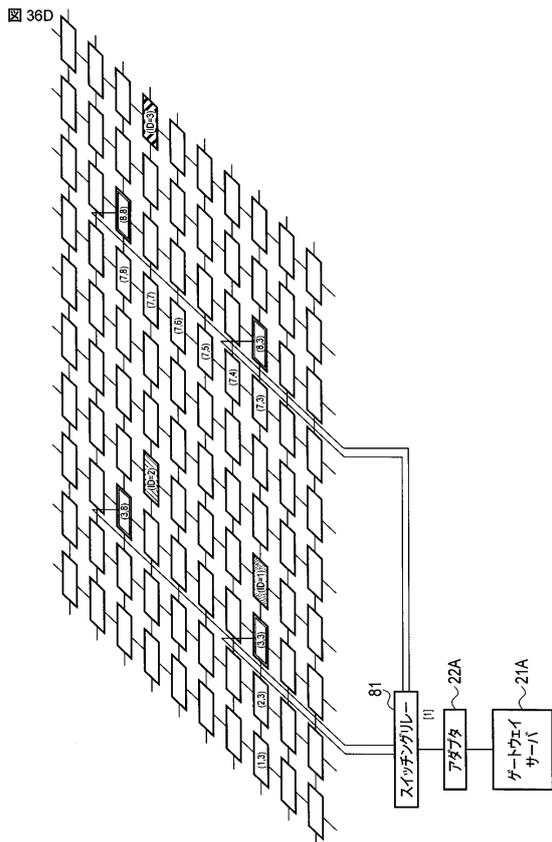
【図 36 B】



【図 36 C】

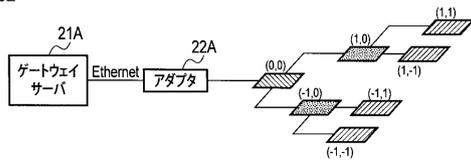


【図 36 D】



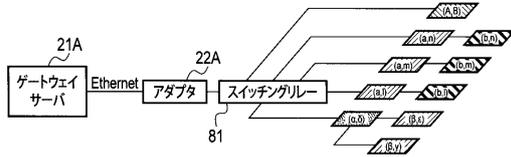
【 図 3 6 E 】

図 36E



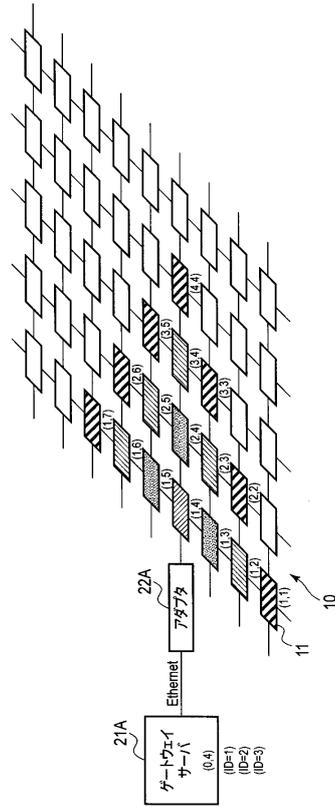
【 図 3 6 F 】

図 36F



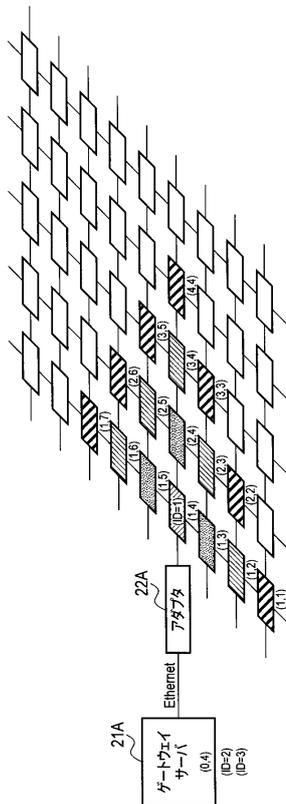
【 図 3 7 A 】

図 37A



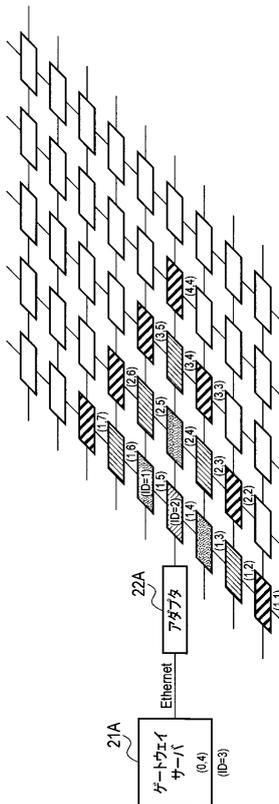
【 図 3 7 B 】

図 37B



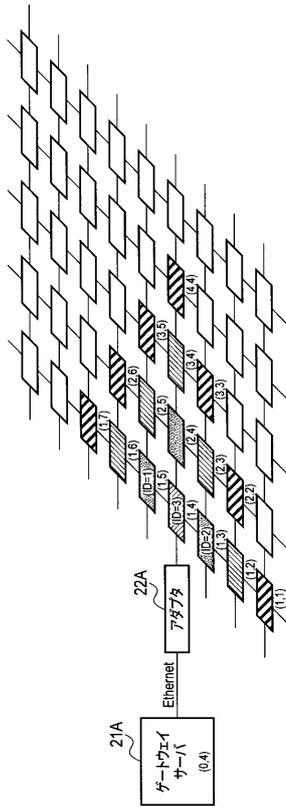
【 図 3 7 C 】

図 37C



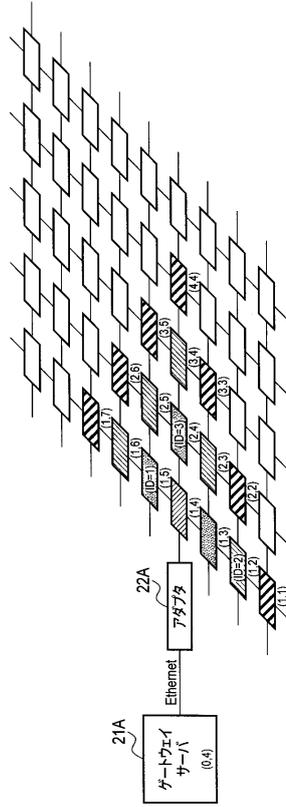
【図 37D】

図 37D



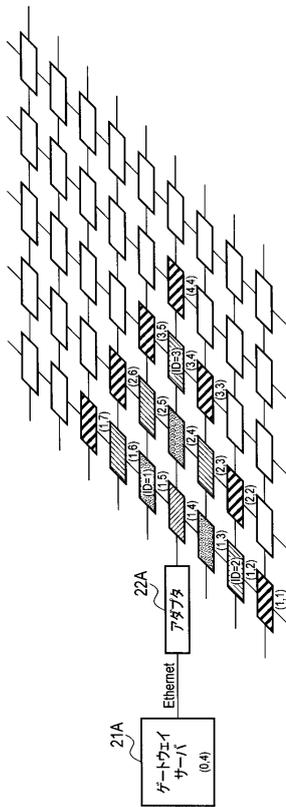
【図 37E】

図 37E



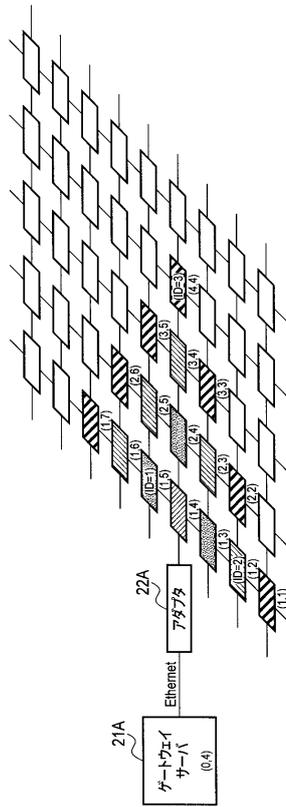
【図 37F】

図 37F



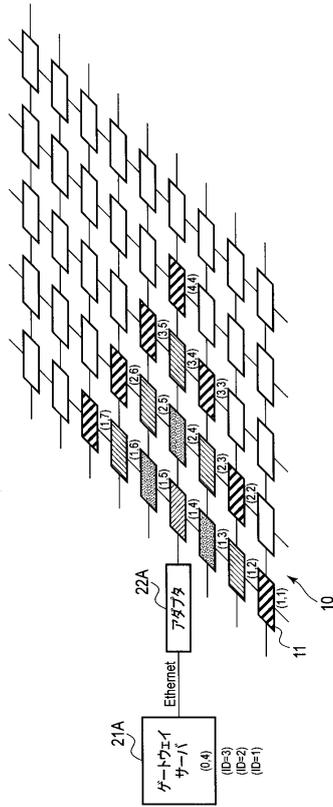
【図 37G】

図 37G



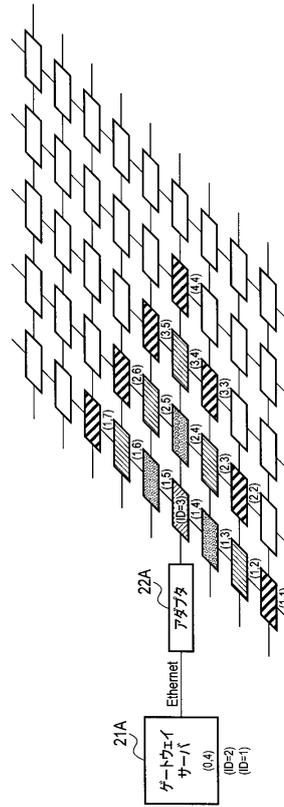
【 図 38 A 】

図 38A



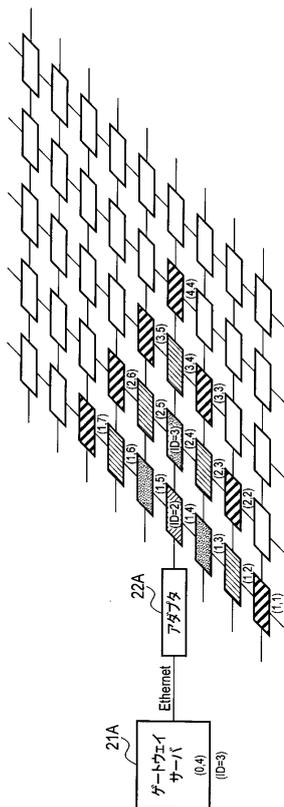
【 図 38 B 】

図 38B



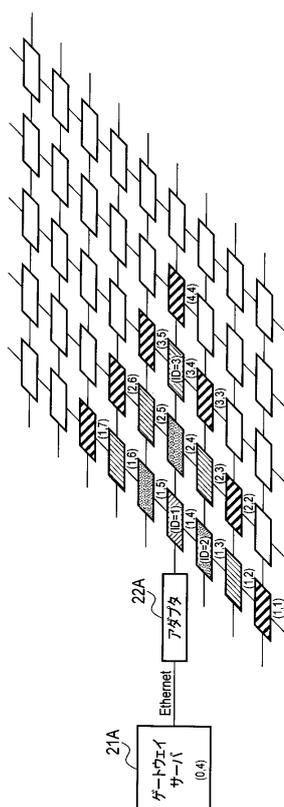
【 図 38 C 】

図 38C



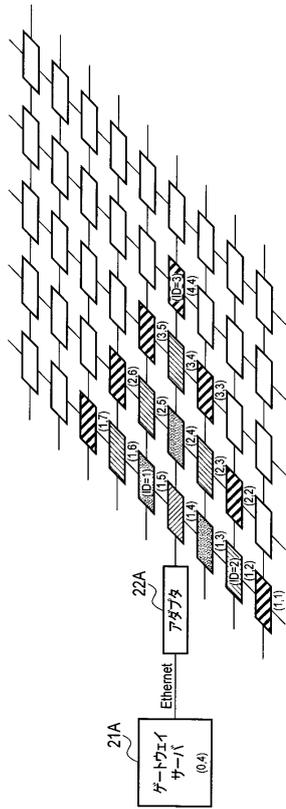
【 図 38 D 】

図 38D



【 図 3 8 E 】

図 38E



フロントページの続き

- (74)代理人 100119976
弁理士 幸長 保次郎
- (74)代理人 100153051
弁理士 河野 直樹
- (74)代理人 100140176
弁理士 砂川 克
- (74)代理人 100158805
弁理士 井関 守三
- (74)代理人 100172580
弁理士 赤穂 隆雄
- (74)代理人 100179062
弁理士 井上 正
- (74)代理人 100124394
弁理士 佐藤 立志
- (74)代理人 100112807
弁理士 岡田 貴志
- (74)代理人 100111073
弁理士 堀内 美保子
- (74)代理人 100134290
弁理士 竹内 将訓
- (72)発明者 辰村 光介
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 木下 敦寛
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 西野 弘剛
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 鈴木 正道
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 西 義史
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 丸亀 孝生
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 栗田 貴宏
東京都港区芝浦一丁目1番1号 株式会社東芝内

審査官 稲葉 崇

- (56)参考文献 再公表特許第2008/149784(JP, A1)
特開2003-174460(JP, A)
特表2008-537265(JP, A)
特開2009-037273(JP, A)
特許第4385387(JP, B2)
特開2010-218364(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08
G06F 13/10 - 13/14
G06F 12/00

G 0 6 F 1 5 / 1 6 - 1 5 / 1 7 7
G 0 6 F 1 5 / 8 0