

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5055942号
(P5055942)

(45) 発行日 平成24年10月24日(2012.10.24)

(24) 登録日 平成24年8月10日(2012.8.10)

(51) Int.Cl.

F I

G 0 6 F 1 5 / 1 7 3 (2 0 0 6 . 0 1)

G 0 6 F 1 5 / 1 7 3 6 4 0 C

請求項の数 5 (全 24 頁)

<p>(21) 出願番号 特願2006-281703 (P2006-281703) (22) 出願日 平成18年10月16日(2006.10.16) (65) 公開番号 特開2008-97528 (P2008-97528A) (43) 公開日 平成20年4月24日(2008.4.24) 審査請求日 平成21年7月10日(2009.7.10)</p>	<p>(73) 特許権者 000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号 (74) 代理人 100092152 弁理士 服部 毅巖 (72) 発明者 安島 雄一郎 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 審査官 清木 泰</p>
---	--

最終頁に続く

(54) 【発明の名称】 計算機クラスタ

(57) 【特許請求の範囲】

【請求項1】

格子状に構成された相互結合網の各格子点に配置される計算機クラスタにおいて、
 格子点の位置を示す外部アドレスと前記計算機クラスタ内における相対位置を示す内部
 アドレスとから成る宛先アドレスを含むパケットの処理を行う演算部と、

前記計算機クラスタの外部の通信相手と前記パケットを送受信する第1の通信部と、
 前記計算機クラスタの内部の通信相手と前記パケットを送受信する第2の通信部と、

前記演算部、前記第1の通信部および前記第2の通信部と接続されており、取得した前
 記パケットを、前記宛先アドレスが自己のアドレスと一致する場合は前記演算部へ出力し
 、前記外部アドレスが前記自己のアドレスと異なり前記パケットを転送すべき方向が前記
 第1の通信部が送信可能な方向である場合は前記第1の通信部に転送させ、前記内部アド
 レスのみが前記自己のアドレスと異なる場合および前記外部アドレスが前記自己のアドレ
 スと異なり前記パケットを転送すべき方向が前記第1の通信部が送信可能な方向でない場
 合は前記第2の通信部に転送させるスイッチ部と、

を有する第1のノード、第2のノード、第3のノードおよび第4のノードと、

前記第1のノード、前記第2のノード、前記第3のノードおよび前記第4のノードの前
 記第2の通信部を相互に接続する内部通信網と、

前記第1のノードの前記第1の通信部と格子のX軸方向に隣接する2つの前記計算機ク
 ラスタとを接続する第1の外部通信網と、

前記第2のノードの前記第1の通信部と格子のY軸方向に隣接する2つの前記計算機ク

10

20

ラストとを接続する第 2 の外部通信網と、

前記第 3 のノードの前記第 1 の通信部と格子の X 軸方向に隣接する 2 つの前記計算機クラスタとを接続する第 3 の外部通信網と、

前記第 4 のノードの前記第 1 の通信部と格子の Y 軸方向に隣接する 2 つの前記計算機クラスタとを接続する第 4 の外部通信網と、

を有することを特徴とする計算機クラスタ。

【請求項 2】

前記内部通信網は、全ての前記第 2 の通信部を双方向リンクによって直接接続した完全結合網であることを特徴とする請求項 1 記載の計算機クラスタ。

【請求項 3】

前記第 1 の外部通信網は、X 軸の負側に隣接する前記計算機クラスタから前記第 1 のノードの前記第 1 の通信部への片方向リンクと、前記第 1 のノードの前記第 1 の通信部から X 軸の正側に隣接する前記計算機クラスタへの片方向リンクとで構成され、

前記第 2 の外部通信網は、Y 軸の負側に隣接する前記計算機クラスタから前記第 2 のノードの前記第 1 の通信部への片方向リンクと、前記第 2 のノードの前記第 1 の通信部から Y 軸の正側に隣接する前記計算機クラスタへの片方向リンクとで構成され、

前記第 3 の外部通信網は、X 軸の正側に隣接する前記計算機クラスタから前記第 3 のノードの前記第 1 の通信部への片方向リンクと、前記第 3 のノードの前記第 1 の通信部から X 軸の負側に隣接する前記計算機クラスタへの片方向リンクとで構成され、

前記第 4 の外部通信網は、Y 軸の正側に隣接する前記計算機クラスタから前記第 4 のノードの前記第 1 の通信部への片方向リンクと、前記第 4 のノードの前記第 1 の通信部から Y 軸の負側に隣接する前記計算機クラスタへの片方向リンクとで構成される、

ことを特徴とする請求項 1 記載の計算機クラスタ。

【請求項 4】

前記第 1 の通信部は、前記パケットを一時的に保持する複数の仮想チャネルバッファを有しており、受信した前記パケットを前記宛先アドレスに応じて前記複数の仮想チャネルバッファに振り分けて格納し、

前記スイッチ部は、前記複数の仮想チャネルバッファに格納された前記パケットを順次取得して、スイッチング処理を行う、

ことを特徴とする請求項 1 記載の計算機クラスタ。

【請求項 5】

前記相互結合網は 3 次元の格子状に構成されており、

前記演算部、前記第 1 の通信部、前記第 2 の通信部および前記スイッチ部を有する第 5 のノードおよび第 6 のノードと、

前記第 5 のノードの前記第 1 の通信部と格子の Z 軸方向に隣接する 2 つの前記計算機クラスタとを接続する第 5 の外部通信網と、

前記第 6 のノードの前記第 1 の通信部と格子の Z 軸方向に隣接する 2 つの前記計算機クラスタとを接続する第 6 の外部通信網と、

を更に有し、

前記内部通信網は、前記第 1 のノード、前記第 2 のノード、前記第 3 のノード、前記第 4 のノード、前記第 5 のノードおよび前記第 6 のノードの前記第 2 の通信部を相互に接続する、

ことを特徴とする請求項 1 記載の計算機クラスタ。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は複数のノードを備える計算機クラスタに関し、特に格子状に構成された相互結合網の各格子点に配置される計算機クラスタに関する。

【背景技術】

【0002】

10

20

30

40

50

近年、ハードウェア技術の進歩に伴い、演算装置の計算能力が飛躍的に向上している。一方で、大規模な科学技術計算や大量のマルチメディアデータの処理など、1つの演算装置の計算能力を大きく超える計算能力を必要とする処理も依然として存在する。これに対し、複数の演算装置を並列に動作させる並列処理の技術が知られている。並列処理を行う場合、演算装置と通信装置とから成るノードを相互にリンク接続した相互結合網を構築する。相互結合網では、ノード間でパケットを送受信することでデータ処理が進行する。相互結合網の結合形態としては、完全結合型、ツリー型、スター型、リング型、メッシュ型、トラス型、ハイパーキューブ型などが存在する。

【0003】

ここで、多数のノードを含む相互結合網の場合、リンク数の増大を抑止するため、格子状の結合形態であるメッシュ型やトラス型が採用されることが多い。格子状の相互結合網では、パケットは幾つかのノードによってリレー方式で中継されて宛先のノードに届けられる。一方、格子状の相互結合網では、膨大な数のノードを相互結合すると、格子の1辺に並ぶノードの数が大きくなり、パケットが宛先のノードへ到達するまでの中継回数が増大するという問題がある。

【0004】

この問題に対し、相互結合網を階層化する技術が知られている。すなわち、相互結合網の各格子点に、1つのノードに代えて、複数のノードを備える計算機クラスタを配置する。計算機クラスタの内部では、複数のノードを、リング型、トラス型、ハイパーキューブ型などで相互接続する(例えば、特許文献1, 2参照)。このように相互結合網を階層化することで、パケットの平均中継回数を抑止することができ、システム全体の計算能力を向上させることができる。

【特許文献1】特開平6-35873号公報

【特許文献2】特開平7-191947号公報

【発明の開示】

【発明が解決しようとする課題】

【0005】

しかし、上記特許文献1, 2記載の技術では、演算装置の動作速度が向上するに伴い、計算機クラスタの内部でのパケット転送がシステム全体の計算能力のボトルネックになるという問題がある。すなわち、計算機クラスタの内部でもパケットはリレー方式で中継されるため、中継によるオーバーヘッドやノード間を結合するリンクの伝送能力が、計算能力の向上を妨げる原因となる。

【0006】

一方、計算機クラスタの内部での中継回数を減らすために、補助的なリンクやスイッチ装置を設けると、計算機クラスタの内部の構成が非対称になるという問題がある。単純な格子状の相互結合網では全てのノードの中継規則は同一であるが、補助的なリンクやスイッチ装置を設けると、中継規則がノードによって異なってしまうためである。この場合、相互結合網を実現する回路構成が複雑になってしまう。

【0007】

本発明はこのような点に鑑みてなされたものであり、格子状に構成された相互結合網の各格子点に配置される計算機クラスタであって、内部におけるパケットの中継回数を削減しつつ、内部構成を対称化できる計算機クラスタを提供することを目的とする。

【課題を解決するための手段】

【0008】

本発明では、上記課題を解決するために、図1に示すような計算機クラスタが提供される。本発明に係る計算機クラスタは、格子状に構成された相互結合網の各格子点に配置される。

【0009】

計算機クラスタ10は、第1のノード11、第2のノード12、第3のノード13、第4のノード14、内部通信網15、第1の外部通信網16、第2の外部通信網17、第3

10

20

30

40

50

の外部通信網 18 および第 4 の外部通信網 19 を有する。第 1 のノード 11 は、格子点の位置を示す外部アドレスと計算機クラスタ内における相対位置を示す内部アドレスとから成る宛先アドレスを含むパケットの処理を行う演算部 11a と、計算機クラスタ 10 の外部の通信相手とパケットを送受信する第 1 の通信部 11b と、計算機クラスタ 10 の内部の通信相手とパケットを送受信する第 2 の通信部 11c と、演算部 11a、第 1 の通信部 11b および第 2 の通信部 11c と接続されており、取得したパケットを、宛先アドレスが自己のアドレスと一致する場合は演算部 11a へ出力し、外部アドレスが自己のアドレスと異なりパケットを転送すべき方向が第 1 の通信部 11b が送信可能な方向である場合は第 1 の通信部 11b に転送させ、内部アドレスのみが自己のアドレスと異なる場合および外部アドレスが自己のアドレスと異なりパケットを転送すべき方向が第 1 の通信部 11b が送信可能な方向でない場合は第 2 の通信部 11c に転送させるスイッチ部 11d と、を有する。第 2 のノード 12、第 3 のノード 13 および第 4 のノード 14 も、第 1 のノード 11 と同様に、演算部、第 1 の通信部、第 2 の通信部およびスイッチ部を有する。内部通信網 15 は、第 1 のノード 11、第 2 のノード 12、第 3 のノード 13 および第 4 のノード 14 の第 2 の通信部を相互に接続する。第 1 の外部通信網 16 は、第 1 のノード 11 の第 1 の通信部 11b と格子の X 軸方向に隣接する 2 つの計算機クラスタとを接続する。第 2 の外部通信網 17 は、第 2 のノード 12 の第 1 の通信部と格子の Y 軸方向に隣接する 2 つの計算機クラスタとを接続する。第 3 の外部通信網 18 は、第 3 のノード 13 の第 1 の通信部と格子の X 軸方向に隣接する 2 つの計算機クラスタとを接続する。第 4 の外部通信網 19 は、第 4 のノード 14 の第 1 の通信部と格子の Y 軸方向に隣接する 2 つの計算機クラスタとを接続する。

【0010】

このような計算機クラスタ 10 によれば、第 1 のノード 11 により、第 1 の外部通信網 16 を介して、相互結合網の X 軸方向へ向かうパケットが中継される。第 2 のノード 12 により、第 2 の外部通信網 17 を介して、相互結合網の Y 軸方向へ向かうパケットが中継される。第 3 のノード 13 により、第 3 の外部通信網 18 を介して、相互結合網の X 軸方向へ向かうパケットが中継される。第 4 のノード 14 により、第 4 の外部通信網 19 を介して、相互結合網の Y 軸方向へ向かうパケットが中継される。また、第 1 のノード 11、第 2 のノード 12、第 3 のノード 13 および第 4 のノード 14 により、内部通信網 15 を介して、計算機クラスタ 10 内でパケットが転送される。そして、パケットが宛先アドレスで示されるノードに到達すると、ノードの演算部によって、パケット処理が行われる。

【発明の効果】

【0011】

本発明によれば、複数のノードをクラスタ化して、第 1 のノードおよび第 3 のノードが X 軸方向へ向かうパケットを中継し、第 2 のノードおよび第 4 のノードが Y 軸方向へ向かうパケットを中継することとした。これにより、格子状に構成された相互結合網におけるパケットの平均中継回数を削減することができる。また、パケットの中継処理を複数のノードに均等に分散させることができ、特定のノードやリンクにパケットが集中することによる計算能力の低下を防止できる。また、各ノードの構成を均一にすることができ、計算機クラスタを実現する回路を単純化できる。

【発明を実施するための最良の形態】

【0012】

以下、本発明の実施の形態を図面を参照して詳細に説明する。まず、本実施の形態の概要を説明し、その後、本実施の形態の具体的な内容を説明する。

図 1 は、本実施の形態の概要を示す図である。本実施の形態に係る並列計算機システムは、2次元のトーラス型相互結合網の各格子点に、複数のノードを備える計算機クラスタを配置したものである。すなわち、階層的な相互結合網で構成した並列計算機システムである。なお、トーラス型相互結合網とは、メッシュ型相互結合網の各座標軸の端点同士を結合してループ状にしたものである。

【0013】

図1は、本実施の形態に係る並列計算機システムの一部を示している。図1に示す並列計算機システムは、計算機クラスタ10、20、30、40を有している。計算機クラスタ10は、格子点(0,0)に配置されている。計算機クラスタ20は、格子点(1,0)に配置されている。計算機クラスタ30は、格子点(0,1)に配置されている。計算機クラスタ40は、格子点(1,1)に配置されている。隣接する格子点に配置された計算機クラスタ同士がリンクで結合されている。

【0014】

計算機クラスタ10は、第1のノード11、第2のノード12、第3のノード13、第4のノード14、内部通信網15、第1の外部通信網16、第2の外部通信網17、第3の外部通信網18および第4の外部通信網19を有している。計算機クラスタ20、30、40も、計算機クラスタ10と同様の構成で実現される。

10

【0015】

第1のノード11、第2のノード12、第3のノード13および第4のノード14は、それぞれ、データ処理の機能とパケットを送受信する機能とを備えている。内部通信網15は、第1のノード11、第2のノード12、第3のノード13および第4のノード14が送受信するパケットを、計算機クラスタ10の内部で相互に伝送する。

【0016】

第1の外部通信網16は、X軸の負方向に隣接する計算機クラスタから到来するパケットを第1のノード11へ伝送すると共に、第1のノード11が送信するパケットをX軸の正方向に隣接する計算機クラスタ20に向かって伝送する。第2の外部通信網17は、Y軸の負方向に隣接する計算機クラスタから到来するパケットを第2のノード12へ伝送すると共に、第2のノード12が送信するパケットをY軸の正方向に隣接する計算機クラスタ30に向かって伝送する。

20

【0017】

第3の外部通信網18は、X軸の正方向に隣接する計算機クラスタ20から到来するパケットを第3のノード13へ伝送すると共に、第3のノード13が送信するパケットをX軸の負方向に隣接する計算機クラスタに向かって伝送する。第4の外部通信網19は、Y軸の正方向に隣接する計算機クラスタ30から到来するパケットを第4のノード14へ伝送すると共に、第4のノード14が送信するパケットをY軸の負方向に隣接する計算機クラスタに向かって伝送する。

30

【0018】

第1のノード11は、内部に、演算部11a、第1の通信部11b、第2の通信部11cおよびスイッチ部11dを有している。第2のノード12、第3のノード13および第4のノード14も、第1のノード11と同様の構成で実現される。

【0019】

演算部11aは、指定されたプログラムを実行する。また、演算部11aは、必要に応じてパケットの入出力を行う。演算部11aが出力するパケットには、格子点の座標(外部アドレス)と計算機クラスタ内におけるノードの位置を示す番号(内部アドレス)とから成る宛先アドレスが含まれる。

【0020】

第1の通信部11bは、第1の外部通信網16を介して、計算機クラスタ10の外部の通信相手とパケットを送受信する。具体的には、第1の通信部11bは、計算機クラスタ10の外部からパケットを受信すると、受信したパケットをスイッチ部11dへ出力する。また、第1の通信部11bは、スイッチ部11dからパケットを取得すると、取得したパケットを計算機クラスタ10の外部へ送信する。

40

【0021】

第2の通信部11cは、内部通信網15を介して、計算機クラスタ10の内部の通信相手とパケットを送受信する。具体的には、第2の通信部11cは、第2のノード12、第3のノード13および第4のノード14のいずれかのノードからパケットを受信すると、受信したパケットをスイッチ部11dへ出力する。また、第2の通信部11cは、スイッ

50

子部 1 1 d からパケットを取得すると、パケットの宛先アドレスに応じて、取得したパケットを第 2 のノード 1 2、第 3 のノード 1 3 および第 4 のノード 1 4 のいずれかのノードへ送信する。

【 0 0 2 2 】

スイッチ部 1 1 d は、演算部 1 1 a、第 1 の通信部 1 1 b および第 2 の通信部 1 1 c が出力するパケットを取得する。スイッチ部 1 1 d は、取得したパケットから宛先アドレスを抽出し、宛先アドレスに応じてパケットの出力先を決定する。そして、スイッチ部 1 1 d は、決定した出力先へパケットを出力する。

【 0 0 2 3 】

具体的には、スイッチ部 1 1 d は、宛先アドレスが第 1 のノード 1 1 のアドレスと一致する場合は、出力先を演算部 1 1 a に決定する。宛先アドレスが第 2 のノード 1 2、第 3 のノード 1 3 もしくは第 4 のノード 1 4 の場合は、出力先を第 2 の通信部 1 1 c に決定する。外部アドレスが第 1 のノード 1 1 のアドレスと異なり、かつ、パケットを転送すべき方向が第 1 の通信部 1 1 b が送信可能な方向である場合は、出力先を第 1 の通信部 1 1 b に決定する。外部アドレスが第 1 のノード 1 1 のアドレスと異なり、かつ、パケットを転送すべき方向が第 1 の通信部 1 1 b が送信可能な方向でない場合は、出力先を第 2 の通信部 1 1 c に決定する。

【 0 0 2 4 】

このような計算機クラスタ 1 0 によれば、第 1 のノード 1 1 により、第 1 の外部通信網 1 6 を介して、X 軸の正方向へ向かうパケットが中継される。第 2 のノード 1 2 により、第 2 の外部通信網 1 7 を介して、Y 軸の正方向へ向かうパケットが中継される。第 3 のノード 1 3 により、第 3 の外部通信網 1 8 を介して、X 軸の負方向へ向かうパケットが中継される。第 4 のノード 1 4 により、第 4 の外部通信網 1 9 を介して、Y 軸の負方向へ向かうパケットが中継される。また、第 1 のノード 1 1、第 2 のノード 1 2、第 3 のノード 1 3 および第 4 のノード 1 4 により、内部通信網 1 5 を介して、計算機クラスタ 1 0 内でパケットが転送される。そして、パケットが宛先アドレスで示されるノードに到達すると、ノードの演算部によって、パケット処理が行われる。

【 0 0 2 5 】

これにより、格子状に構成された相互結合網におけるパケットの平均中継回数を削減することができる。また、パケットの中継処理を複数のノードに均等に分散させることができ、特定のノードやリンクにパケットが集中することによる計算能力の低下を防止できる。また、各ノードの構成を均一にすることができ、計算機クラスタを実現する回路を単純化できる。

【 0 0 2 6 】

例えば、図 1 に示した 1 6 個のノード、すなわち、第 1 のノード 1 1、2 1、3 1、4 1、第 2 のノード 1 2、2 2、3 2、4 2、第 3 のノード 1 3、2 3、3 3、4 3 および第 4 のノード 1 4、2 4、3 4、4 4 を、単層の格子状に配置すると、1 辺の長さが 4 の 2 次元の格子が構成される。この場合、第 1 のノード 1 1 から第 3 のノード 4 3 へパケットを送信するには、6 回の中継処理を行う必要がある。一方、図 1 に示した構成の相互結合網では、第 1 のノード 1 1、第 1 のノード 2 1、第 2 のノード 2 2、第 2 のノード 4 2、第 3 のノード 4 3 の順に、4 回の中継処理を行えばよい。

【 0 0 2 7 】

[第 1 の実施の形態]

最初に、第 1 の実施の形態を、図面を参照して詳細に説明する。

図 2 は、第 1 の実施の形態のシステム構成を示す図である。第 1 の実施の形態に係る並列計算機システムは、2 次元のトラス型相互結合網の各格子点に、複数の計算ユニットを備えるクラスタサーバを配置したものである。図 2 に示す並列計算機システムは、9 個のクラスタサーバ 1 0 0、1 0 0 a、1 0 0 b、1 0 0 c、1 0 0 d、1 0 0 e、1 0 0 f、1 0 0 g、1 0 0 h を有している。

【 0 0 2 8 】

10

20

30

40

50

クラスタサーバ100は、格子点(0, 0)に配置されている。クラスタサーバ100aは、格子点(1, 0)に配置されている。クラスタサーバ100bは、格子点(2, 0)に配置されている。クラスタサーバ100cは、格子点(0, 1)に配置されている。クラスタサーバ100dは、格子点(1, 1)に配置されている。クラスタサーバ100eは、格子点(2, 1)に配置されている。クラスタサーバ100fは、格子点(0, 2)に配置されている。クラスタサーバ100gは、格子点(1, 2)に配置されている。クラスタサーバ100hは、格子点(2, 2)に配置されている。

【0029】

隣接する格子点に配置されたクラスタサーバ同士が、双方向通信が可能な通信ケーブルで接続されている。また、各座標軸の端点にあるクラスタサーバ同士も、双方向通信が可能な通信ケーブルで接続されている。これにより、任意のクラスタサーバ間でパケットを送受信することができる。

10

【0030】

例えば、クラスタサーバ100dからクラスタサーバ100hへ向かうパケットは、クラスタサーバ100d、クラスタサーバ100e、クラスタサーバ100hの順に中継される。また、クラスタサーバ100hからクラスタサーバ100へ向かうパケットは、クラスタサーバ100h、クラスタサーバ100f、クラスタサーバ100の順に中継される。

【0031】

図3は、第1の実施の形態のクラスタサーバの構成を示す図である。なお、図3にはクラスタサーバ100の構成を示したが、クラスタサーバ100a, 100b, 100c, 100d, 100e, 100f, 100g, 100hも同様の構成で実現できる。クラスタサーバ100は、計算ユニット110, 120, 130, 140および通信パネル150を有している。

20

【0032】

計算ユニット110, 120, 130, 140は、演算機能とパケットのスイッチ機能とを備える装置である。個々の計算ユニットには、クラスタサーバ100の内部における相対位置を示すクラスタ内番号が割り当てられている。具体的には、計算ユニット110に番号0、計算ユニット120に番号1、計算ユニット130に番号2、計算ユニット140に番号3が割り当てられている。

30

【0033】

個々の計算ユニットを一意に識別するアドレスは、クラスタサーバ100が配置された格子点の座標とクラスタ内番号とで構成される。すなわち、計算ユニット110のアドレスは(0, 0, 0)、計算ユニット120のアドレスは(0, 0, 1)、計算ユニット130のアドレスは(0, 0, 2)、計算ユニット140のアドレスは(0, 0, 3)である。

【0034】

計算ユニット110, 120, 130, 140は、双方向通信が可能な6本の通信ケーブルによって、相互に接続されている。すなわち、個々の計算ユニットは、他の3つの計算ユニットと3本の通信ケーブルによって直接接続されている。これにより、クラスタサーバ100内の任意の2つの計算ユニット間で、パケットを送受信することができる。

40

【0035】

通信パネル150は、通信ポート151, 152, 153, 154を有している。通信ポート151, 152, 153, 154は、他のクラスタサーバと通信を行うための通信ケーブルを接続するインターフェースである。

【0036】

具体的には、通信ポート151には、X軸の負方向に隣接するクラスタサーバ100bと通信を行うための通信ケーブルが接続される。通信ポート152には、X軸の正方向に隣接するクラスタサーバ100aと通信を行うための通信ケーブルが接続される。通信ポート153には、Y軸の負方向に隣接するクラスタサーバ100fと通信を行うための通

50

信ケーブルが接続される。通信ポート 154 には、Y 軸の正方向に隣接するクラスタサーバ 100c と通信を行うための通信ケーブルが接続される。

【0037】

通信ポート 151, 152, 153, 154 は、通信ケーブルを介して受信したパケットを、各計算ユニットへ出力する。また、計算ユニットから取得したパケットを、通信ケーブルを介して送信する。ここで、パケットの出力先となる計算ユニットと、パケットの取得元となる計算ユニットとは、別の計算ユニットである。

【0038】

図 4 は、第 1 の実施の形態の通信ケーブルの接続関係を示す図である。通信ポート 151, 152, 153, 154 に接続される通信ケーブルは、より詳細には、上り方向（クラスタサーバ 100 へ入力される方向）のパケットを伝送するケーブルと、下り方向（クラスタサーバ 100 から出力する方向）のパケットを伝送するケーブルとで構成される。

10

【0039】

ここで、通信ポート 151 は、上り方向のケーブルで受信したパケットを計算ユニット 110 へ出力すると共に、計算ユニット 130 から取得したパケットを下り方向のケーブルで送信する。通信ポート 153 は、上り方向のケーブルで受信したパケットを計算ユニット 130 へ出力すると共に、計算ユニット 110 から取得したパケットを下り方向のケーブルで送信する。

【0040】

同様に、通信ポート 152 は、上り方向のケーブルで受信したパケットを計算ユニット 120 へ出力すると共に、計算ユニット 140 から取得したパケットを下り方向のケーブルで送信する。通信ポート 154 は、上り方向のケーブルで受信したパケットを計算ユニット 140 へ出力すると共に、計算ユニット 120 から取得したパケットを下り方向のケーブルで送信する。

20

【0041】

このように、上り方向の通信経路と下り方向の通信経路とを分離し交差させることで、個々の計算ユニットは、特定の一方方向のパケットの中継のみを担当することができる。すなわち、計算ユニット 110 は、X 軸の正方向の中継のみを担当する。計算ユニット 120 は、Y 軸の正方向の中継のみを担当する。計算ユニット 130 は、X 軸の負方向の中継のみを担当する。計算ユニット 140 は、Y 軸の負方向の中継のみを担当する。

30

【0042】

図 5 は、第 1 の実施の形態の計算ユニットの構成を示す図である。なお、図 5 には計算ユニット 110 の構成を示したが、計算ユニット 120, 130, 140 も同様の構成で実現できる。計算ユニット 110 は、演算回路 111、入力回路 112、出力回路 113、受信回路 114, 116a, 116b, 116c、送信回路 115, 117a, 117b, 117c およびスイッチ回路 118 を有している。

【0043】

演算回路 111 は、入力回路 112 が保持するパケットを順次取得し、データ処理を行う。また、演算回路 111 は、他の計算ユニットへデータを送信する必要がある場合、宛先アドレスを含むパケットを生成して、出力回路 113 へ出力する。

40

【0044】

入力回路 112 は、内部に F I F O (First In First Out) 型のバッファメモリを有している。入力回路 112 は、スイッチ回路 118 からパケットを取得すると、取得したパケットをバッファメモリの最後尾に格納する。また、入力回路 112 は、演算回路 111 からの要求に回答して、バッファメモリの先頭からパケットを取り出し、演算回路 111 へ出力する。

【0045】

出力回路 113 は、内部に F I F O 型のバッファメモリを有している。出力回路 113 は、演算回路 111 からパケットを取得すると、取得したパケットをバッファメモリの最後尾に格納する。また、出力回路 113 は、スイッチ回路 118 からの要求に回答して、

50

バッファメモリの先頭からパケットを取り出し、スイッチ回路 118 へ出力する。

【0046】

受信回路 114 は、内部に複数の F I F O 型のバッファメモリを仮想チャネルバッファとして有している。受信回路 114 は、通信ポート 151 からのパケット、すなわち、クラスタサーバ 100 b が送信したパケットを取得する。受信回路 114 は、パケットを取得すると、パケットの宛先アドレスに基づいて仮想チャネルバッファを 1 つ選択し、選択した仮想チャネルバッファの最後尾に取得したパケットを格納する。また、受信回路 114 は、スイッチ回路 118 からの要求に応答して、仮想チャネルバッファの先頭からパケットを取り出し、スイッチ回路 118 へ出力する。これにより、1 つの通信経路の上に複数の仮想チャネルを構成できる。

10

【0047】

このように、受信回路 114 が複数の仮想チャネルを構成できることで、デッドロックの発生を回避することができる。例えば、クラスタサーバ 100 (座標 (0, 0)) がクラスタサーバ 100 a (座標 (1, 0)) へ、クラスタサーバ 100 a がクラスタサーバ 100 d (座標 (1, 1)) へ、クラスタサーバ 100 d がクラスタサーバ 100 c (座標 (1, 0)) へ、クラスタサーバ 100 c がクラスタサーバ 100 へ、それぞれパケットを送信する場合を考える。ここで、受信回路にバッファメモリが 1 つしかなく、経路上の全てのバッファメモリが満杯の場合、デッドロックが発生する。これに対し、複数の仮想チャネルを構成できるようにすることで、デッドロックの発生を回避できる。

【0048】

送信回路 115 は、スイッチ回路 118 からパケットを取得すると、取得したパケットを通信ポート 153 へ出力する。これにより、パケットがクラスタサーバ 100 a へ送信される。

20

【0049】

受信回路 116 a, 116 b, 116 c は、それぞれ内部に F I F O 型のバッファメモリを有している。受信回路 116 a は、計算ユニット 120 からのパケットを取得する。受信回路 116 b は、計算ユニット 130 からのパケットを取得する。受信回路 116 c は、計算ユニット 140 からパケットを取得する。受信回路 116 a, 116 b, 116 c は、取得したパケットをバッファメモリの先頭に格納する。また、受信回路 116 a, 116 b, 116 c は、スイッチ回路 118 からの要求に応答して、バッファメモリの先頭からパケットを取り出し、スイッチ回路 118 へ出力する。

30

【0050】

送信回路 117 a, 117 b, 117 c は、スイッチ回路 118 からパケットを取得する。送信回路 117 a は、取得したパケットを計算ユニット 120 へ出力する。送信回路 117 b は、取得したパケットを計算ユニット 130 へ出力する。送信回路 117 c は、取得したパケットを計算ユニット 140 へ出力する。

【0051】

スイッチ回路 118 は、出力回路 113 および受信回路 114, 116 a, 116 b, 116 c のバッファメモリを監視する。バッファメモリにパケットが格納されている場合、スイッチ回路 118 は、出力回路 113 および受信回路 114, 116 a, 116 b, 116 c からパケットを順次取得する。そして、スイッチ回路 118 は、パケットの宛先アドレスに基づいて、パケットの出力先を決定する。パケットの出力先は、入力回路 112 および送信回路 115, 117 a, 117 b, 117 c のいずれかである。

40

【0052】

図 6 は、第 1 の実施の形態のパケット出力先の決定方法を示す図である。なお、図 6 にはスイッチ回路 118 によるパケット出力先の決定方法を示したが、他の計算ユニットのスイッチ回路も同様の方法でパケット出力先を決定することができる。その場合、図 6 に示した宛先アドレスを適宜読み替えばよい。

【0053】

スイッチ回路 118 は、パケットの宛先アドレスが計算ユニット 110 のアドレスと一

50

致する場合、更に、パケットの取得元が出力回路 1 1 3 か否かを判断する。パケットの取得元が出力回路 1 1 3 の場合、パケットを破棄する。これは、自分自身へパケットを送信することは、エラーと判断できるからである。一方、パケットの取得元が出力回路 1 1 3 でない場合、パケットを入力回路 1 1 2 へ出力する。

【 0 0 5 4 】

宛先アドレスが計算ユニット 1 2 0 の場合、スイッチ回路 1 1 8 は、パケットを送信回路 1 1 7 a へ出力する。宛先アドレスが計算ユニット 1 3 0 の場合、パケットを送信回路 1 1 7 b へ出力する。宛先アドレスが計算ユニット 1 4 0 の場合、パケットを送信回路 1 1 7 c へ出力する。

【 0 0 5 5 】

宛先アドレスと計算ユニット 1 1 0 のアドレスとで X 座標が異なり、X 軸の正方向へ転送することで距離が縮まる場合、スイッチ回路 1 1 8 は、パケットを送信回路 1 1 5 へ出力する。宛先アドレスと計算ユニット 1 1 0 のアドレスとで X 座標が異なり、X 軸の負方向へ転送することで距離が縮まる場合、パケットを送信回路 1 1 7 b へ出力する。

【 0 0 5 6 】

宛先アドレスと計算ユニット 1 1 0 のアドレスとで X 座標が同一で Y 座標が異なり、Y 座標の正方向へ転送することで距離が縮まる場合、スイッチ回路 1 1 8 は、パケットを送信回路 1 1 7 a へ出力する。宛先アドレスと計算ユニット 1 1 0 のアドレスとで X 座標が同一で Y 座標が異なり、Y 座標の負方向へ転送することで距離が縮まる場合、パケットを送信回路 1 1 7 c へ出力する。

【 0 0 5 7 】

なお、図 6 に示した方法で、X 軸方向への転送を Y 軸方向への転送よりも優先するように統一しているのは、デッドロックの発生を抑制するためである。なお、Y 軸方向への転送を X 軸方向への転送よりも優先するように統一しても、同様の効果を得られる。

【 0 0 5 8 】

このように、第 1 の実施の形態に係る並列計算機システムによれば、宛先アドレスで示される計算ユニットに到達するまでのパケットの中継回数を削減することができる。また、2 次元座標空間における 4 方向へのパケットの送信を 4 つの計算ユニットで均等に分担するため、特定の回路や通信ケーブルにパケットが集中することを防止できる。

【 0 0 5 9 】

また、クラスタサーバの内部では完全結合網によりパケットを直接転送するため、内部の転送に関してデッドロック防止のための仮想チャネルバッファを設ける必要がなく、回路規模の増大を抑制できる。また、各計算ユニットの構成を均一にすることができ、クラスタサーバを実現する構成を単純化できる。

【 0 0 6 0 】

[第 2 の実施の形態]

次に、第 2 の実施の形態を、図面を参照して詳細に説明する。前述の第 1 の実施の形態との相違点を中心に説明し、同様の事項については説明を省略する。

【 0 0 6 1 】

図 7 は、第 2 の実施の形態のシステム構成を示す図である。第 2 の実施の形態に係る並列計算機システムは、2 次元のトラス型相互結合網の各格子点に、複数の PC サーバを備える PC クラスタを配置したものである。PC サーバとは、通常は個人で使用するパーソナルコンピュータをサーバコンピュータとして使用したものである。図 7 に示す並列計算機システムは、9 個の PC クラスタ 2 0 0 , 2 0 0 a , 2 0 0 b , 2 0 0 c , 2 0 0 d , 2 0 0 e , 2 0 0 f , 2 0 0 g , 2 0 0 h を有している。

【 0 0 6 2 】

PC クラスタ 2 0 0 は、格子点 (0 , 0) に配置されている。PC クラスタ 2 0 0 a は、格子点 (1 , 0) に配置されている。PC クラスタ 2 0 0 b は、格子点 (2 , 0) に配置されている。PC クラスタ 2 0 0 c は、格子点 (0 , 1) に配置されている。PC クラスタ 2 0 0 d は、格子点 (1 , 1) に配置されている。PC クラスタ 2 0 0 e は、格子点

10

20

30

40

50

(2, 1)に配置されている。PCクラスタ200fは、格子点(0, 2)に配置されている。PCクラスタ200gは、格子点(1, 2)に配置されている。PCクラスタ200hは、格子点(2, 2)に配置されている。

【0063】

隣接する格子点に配置されたPCクラスタ同士が、双方向通信が可能な3本の通信ケーブルで接続されている。また、各座標軸の端点にあるPCクラスタ同士も、双方向通信が可能な3本の通信ケーブルで接続されている。

【0064】

図8は、第2の実施の形態のPCクラスタの構成を示す図である。なお、図8にはPCクラスタ200の構成を示したが、PCクラスタ200a, 200b, 200c, 200d, 200e, 200f, 200g, 200hも同様の構成で実現できる。PCクラスタ200は、PCサーバ210, 220, 230, 240, 250, 260およびスイッチ270を有している。

【0065】

PCサーバ210, 220, 230, 240, 250, 260は、演算機能とパケットのスイッチ機能とを備えるコンピュータである。個々のPCサーバには、PCクラスタ200の内部における相対位置を示すクラスタ内番号が割り当てられている。具体的には、PCサーバ210に番号0、PCサーバ220に番号1、PCサーバ230に番号2、PCサーバ240に番号3、PCサーバ250に番号4、PCサーバ260に番号5が割り当てられている。

【0066】

個々のPCサーバを一意に識別するアドレスは、PCクラスタ200が配置された格子点の座標とクラスタ内番号とで構成される。すなわち、PCサーバ210のアドレスは(0, 0, 0)、PCサーバ220のアドレスは(0, 0, 1)、PCサーバ230のアドレスは(0, 0, 2)、PCサーバ240のアドレスは(0, 0, 3)、PCサーバ250のアドレスは(0, 0, 4)、PCサーバ260のアドレスは(0, 0, 5)である。

【0067】

PCサーバ210, 220, 230, 240, 250, 260は、それぞれ3つの通信ポートを有している。すなわち、PCサーバ210は、通信ポート219a, 219b, 219cを有している。PCサーバ220は、通信ポート229a, 229b, 229cを有している。PCサーバ230は、通信ポート239a, 239b, 239cを有している。PCサーバ240は、通信ポート249a, 249b, 249cを有している。PCサーバ250は、通信ポート259a, 259b, 259cを有している。PCサーバ260は、通信ポート269a, 269b, 269cを有している。

【0068】

通信ポート219c, 229c, 239c, 249c, 259c, 269cには、スイッチ270と通信を行うための通信ケーブルが接続されている。PCサーバ210, 220, 230, 240, 250, 260は、転送先のPCサーバを指定してスイッチ270へパケットを出力する。スイッチ270は、取得したパケットを指定されたPCサーバへ転送する。

【0069】

通信ポート219a, 219b, 229a, 229b, 239a, 239b, 249a, 249b, 259a, 259b, 269a, 269bには、それぞれ、PCクラスタ200の外のPCクラスタと通信を行うための通信ケーブルが接続される。

【0070】

具体的には、通信ポート219a, 239a, 259aには、X軸の負方向に隣接するPCクラスタ200b内のPCサーバと通信を行うための通信ケーブルが接続される。通信ポート219b, 239b, 259bには、X軸の正方向に隣接するPCクラスタ200a内のPCサーバと通信を行うための通信ケーブルが接続される。通信ポート229a, 249a, 269aには、Y軸の負方向に隣接するPCクラスタ200f内のPCサー

10

20

30

40

50

バと通信を行うための通信ケーブルが接続される。通信ポート 2 2 9 b , 2 4 9 b , 2 6 9 b には、Y 軸の正方向に隣接する P C クラスタ 2 0 0 c 内の P C サーバと通信を行うための通信ケーブルが接続される。

【 0 0 7 1 】

なお、第 2 の実施の形態では、構成要素間を接続する場合は、全て双方向通信が可能な通信ケーブルを用いる。したがって、第 1 の実施の形態の構成と異なり、通信ポートにおいて上り方向の通信経路と下り方向の通信経路とを分離し交差させることは行わない。

【 0 0 7 2 】

図 9 は、第 2 の実施の形態の P C サーバの構成を示す図である。なお、図 9 には P C サーバ 2 1 0 の構成を示したが、P C サーバ 2 2 0 , 2 3 0 , 2 4 0 , 2 5 0 , 2 6 0 も同様の構成で実現できる。P C サーバ 2 1 0 は、演算回路 2 1 1、入力回路 2 1 2、出力回路 2 1 3、受信回路 2 1 4 a , 2 1 4 b , 2 1 6、送信回路 2 1 5 a , 2 1 5 b , 2 1 7 およびスイッチ回路 2 1 8 を有している。

10

【 0 0 7 3 】

演算回路 2 1 1、入力回路 2 1 2 および出力回路 2 1 3 は、図 5 に示した第 1 の実施の形態の演算回路 1 1 1、入力回路 1 1 2 および出力回路 1 1 3 とそれぞれ同様の機能を備えている。

【 0 0 7 4 】

受信回路 2 1 4 a , 2 1 4 b は、それぞれ内部に複数の F I F O 型のバッファメモリを仮想チャネルバッファとして有している。受信回路 2 1 4 a は、通信ポート 2 1 9 a からのパケット、すなわち、P C クラスタ 2 0 0 b からのパケットを取得する。受信回路 2 1 4 b は、通信ポート 2 1 9 b からのパケット、すなわち、P C クラスタ 2 0 0 a からのパケットを取得する。複数の仮想チャネルバッファを用いて仮想チャネルを実現する方法は、第 1 の実施の形態の方法と同様である。

20

【 0 0 7 5 】

送信回路 2 1 5 a は、スイッチ回路 2 1 8 からパケットを取得すると、取得したパケットを通信ポート 2 1 9 b へ出力する。これにより、パケットが P C クラスタ 2 0 0 a へ送信される。送信回路 2 1 5 b は、スイッチ回路 2 1 8 からパケットを取得すると、取得したパケットを通信ポート 2 1 9 a へ出力する。これにより、パケットが P C クラスタ 2 0 0 b へ送信される。

30

【 0 0 7 6 】

受信回路 2 1 6 は、内部に F I F O 型のバッファメモリを有している。受信回路 2 1 6 は、通信ポート 2 1 9 c からのパケット、すなわち、スイッチ 2 7 0 が転送したパケットを取得する。そして、受信回路 2 1 6 は、取得したパケットをバッファメモリの先頭に格納する。

【 0 0 7 7 】

送信回路 2 1 7 は、スイッチ回路 2 1 8 から、転送先の P C サーバが指定されたパケットを取得する。そして、送信回路 2 1 7 は、取得したパケットを通信ポート 2 1 9 c へ出力する。これにより、パケットがスイッチ 2 7 0 を介して指定された P C サーバへ転送される。

40

【 0 0 7 8 】

スイッチ回路 2 1 8 は、出力回路 2 1 3 および受信回路 2 1 4 a , 2 1 4 b , 2 1 6 のバッファメモリを監視する。バッファメモリにパケットが格納されている場合、スイッチ回路 2 1 8 は、出力回路 2 1 3 および受信回路 2 1 4 a , 2 1 4 b , 2 1 6 からパケットを順次取得する。そして、スイッチ回路 2 1 8 は、パケットの宛先アドレスに基づいて、パケットの出力先を決定する。パケットの出力先は、入力回路 2 1 2 および送信回路 2 1 5 a , 2 1 5 b , 2 1 7 のいずれかである。

【 0 0 7 9 】

図 1 0 は、第 2 の実施の形態のパケット出力先の決定方法を示す図である。なお、図 1 0 にはスイッチ回路 2 1 8 によるパケット出力先の決定方法を示したが、他の P C サーバ

50

のスイッチ回路も同様の方法でパケット出力先を決定することができる。その場合、図 10 に示した宛先アドレスを適宜読み替えばよい。

【 0 0 8 0 】

スイッチ回路 2 1 8 は、パケットの宛先アドレスが P C サーバ 2 1 0 のアドレスと一致する場合、更に、パケットの取得元が出力回路 2 1 3 か否かを判断する。パケットの取得元が出力回路 2 1 3 の場合、パケットを破棄する。一方、パケットの取得元が出力回路 2 1 3 でない場合、パケットを入力回路 2 1 2 へ出力する。宛先アドレスが P C サーバ 2 2 0 , 2 3 0 , 2 4 0 , 2 5 0 , 2 6 0 の場合、宛先アドレスの P C サーバを転送先に指定してパケットを送信回路 2 1 7 へ出力する。

【 0 0 8 1 】

宛先アドレスと P C サーバ 2 1 0 のアドレスとで X 座標が異なり、X 軸の正方向へ転送することで距離が縮まる場合、スイッチ回路 2 1 8 は、パケットを送信回路 2 1 5 a へ出力する。宛先アドレスと P C サーバ 2 1 0 のアドレスとで X 座標が異なり、X 軸の負方向へ転送することで距離が縮まる場合、パケットを送信回路 2 1 5 b へ出力する。宛先アドレスと P C サーバ 2 1 0 のアドレスとで X 座標が同一で Y 座標が異なる場合、スイッチ回路 2 1 8 は、P C サーバ 2 2 0 を転送先に指定してパケットを送信回路 2 1 7 へ出力する。

【 0 0 8 2 】

なお、図 1 0 に示した方法では、X 軸方向への転送を Y 軸方向への転送よりも優先するように統一しているが、Y 軸方向への転送を X 軸方向への転送よりも優先するように統一してもよい。

【 0 0 8 3 】

このように、第 2 の実施の形態に係る並列計算機システムによれば、宛先アドレスで示される P C サーバに到達するまでのパケットの中継回数を削減することができる。また、2 次元座標空間における 4 方向へのパケットの送信を 6 つの P C サーバで分担するため、特定の回路や通信ケーブルにパケットが集中することを防止できる。特に、1 つの座標軸に対して複数の P C サーバを割り当てたため、P C クラスタ間のパケット転送能力を向上させることができる。

【 0 0 8 4 】

また、各 P C サーバの構成を均一にすることができ、P C クラスタを実現する構成を単純化できる。また、P C クラスタ内でのパケットの転送にスイッチを用いたため、各 P C サーバが有する受信回路および送信回路の個数を少なくでき、回路規模の増大を抑制できる。

【 0 0 8 5 】

[第 3 の実施の形態]

次に、第 3 の実施の形態を、図面を参照して詳細に説明する。前述の第 1 の実施の形態および第 2 の実施の形態との相違点を中心に説明し、同様の事項については説明を省略する。

【 0 0 8 6 】

図 1 1 は、第 3 の実施の形態のシステム構成を示す図である。第 3 の実施の形態に係る並列計算機システムは、3 次元のトーラス型相互結合網の各格子点に、複数の計算ユニットを備えるクラスタサーバを配置したものである。図 1 1 は、1 辺が 1 6 個のクラスタサーバで構成される 3 次元の並列計算機システムの一部を示している。第 3 の実施の形態に係る並列計算機システムは、クラスタサーバ 3 0 0 , 3 0 0 a , 3 0 0 b , 3 0 0 c , 3 0 0 d , 3 0 0 e , 3 0 0 f を有している。

【 0 0 8 7 】

クラスタサーバ 3 0 0 は、格子点 (0 , 0 , 0) に配置されている。クラスタサーバ 3 0 0 a は、格子点 (1 , 0 , 0) に配置されている。クラスタサーバ 3 0 0 b は、格子点 (1 5 , 0 , 0) に配置されている。クラスタサーバ 3 0 0 c は、格子点 (0 , 1 , 0) に配置されている。クラスタサーバ 3 0 0 d は、格子点 (0 , 1 5 , 0) に配置されてい

10

20

30

40

50

る。クラスタサーバ300eは、格子点(0, 0, 1)に配置されている。クラスタサーバ300fは、格子点(0, 0, 15)に配置されている。

【0088】

隣接する格子点に配置されたクラスタサーバ同士が、双方向通信が可能な3本の通信ケーブルで接続されている。また、各座標軸の端点にあるクラスタサーバ同士も、3本の双方向通信が可能な通信ケーブルで接続されている。

【0089】

図12は、第3の実施の形態のクラスタサーバの構成を示す図である。なお、図12にはクラスタサーバ300の構成を示したが、クラスタサーバ300a, 300b, 300c, 300d, 300e, 300fも同様の構成で実現できる。クラスタサーバ300は、計算ユニット310, 310a, 310b, 320, 320a, 320b, 330, 330a, 330bおよび通信パネル340を有している。

10

【0090】

計算ユニット310, 310a, 310b, 320, 320a, 320b, 330, 330a, 330bは、演算機能とパケットのスイッチ機能とを備える装置である。個々の計算ユニットには、クラスタサーバ300の内部における相対位置を示すクラスタ内番号が割り当てられている。具体的には、計算ユニット310に番号0、計算ユニット310aに番号1、計算ユニット310bに番号2、計算ユニット320に番号3、計算ユニット320aに番号4、計算ユニット320bに番号5、計算ユニット330に番号6、計算ユニット330aに番号7、計算ユニット330bに番号8が割り当てられている。

20

【0091】

個々の計算ユニットを一意に識別するアドレスは、クラスタサーバ300が配置された格子点の座標とクラスタ内番号とで構成される。すなわち、計算ユニット310のアドレスは(0, 0, 0, 0)、計算ユニット310aのアドレスは(0, 0, 0, 1)、計算ユニット310bのアドレスは(0, 0, 0, 2)、計算ユニット320のアドレスは(0, 0, 0, 3)、計算ユニット320aのアドレスは(0, 0, 0, 4)、計算ユニット320bのアドレスは(0, 0, 0, 5)、計算ユニット330のアドレスは(0, 0, 0, 6)、計算ユニット330aのアドレスは(0, 0, 0, 7)、計算ユニット330bのアドレスは(0, 0, 0, 8)である。

【0092】

計算ユニット310, 310a, 310b, 320, 320a, 320b, 330, 330a, 330bは、双方向通信が可能な通信ケーブルによって、相互に直接接続されている。すなわち、個々の計算ユニットは、他の8個の計算ユニットと8本の通信ケーブルによって接続されている(完全結合網)。

30

【0093】

通信パネル340は、通信ポート341a, 341b, 341c, 342a, 342b, 342c, 343a, 343b, 343c, 344a, 344b, 344c, 345a, 345b, 345c, 346a, 346b, 346cを有している。

【0094】

具体的には、通信ポート341a, 341b, 341cには、X軸の負方向に隣接するクラスタサーバ300bと通信を行うための通信ケーブルがそれぞれ接続されている。通信ポート341a, 341b, 341cは、それぞれ、計算ユニット310, 310a, 310bとの間でパケットの入出力を行う。

40

【0095】

通信ポート342a, 342b, 342cには、X軸の正方向に隣接するクラスタサーバ300aと通信を行うための通信ケーブルがそれぞれ接続される。通信ポート342a, 342b, 342cは、それぞれ、計算ユニット310, 310a, 310bとの間でパケットの入出力を行う。

【0096】

通信ポート343a, 343b, 343cには、Y軸の負方向に隣接するクラスタサーバ

50

バ 3 0 0 d と通信を行うための通信ケーブルが接続される。通信ポート 3 4 3 a , 3 4 3 b , 3 4 3 c は、それぞれ、計算ユニット 3 2 0 , 3 2 0 a , 3 2 0 b との間でパケットの入出力を行う。

【 0 0 9 7 】

通信ポート 3 4 4 a , 3 4 4 b , 3 4 4 c には、Y 軸の正方向に隣接するクラスタサーバ 3 0 0 c と通信を行うための通信ケーブルが接続される。通信ポート 3 4 4 a , 3 4 4 b , 3 4 4 c は、それぞれ、計算ユニット 3 2 0 , 3 2 0 a , 3 2 0 b との間でパケットの入出力を行う。

【 0 0 9 8 】

通信ポート 3 4 5 a , 3 4 5 b , 3 4 5 c には、Z 軸の負方向に隣接するクラスタサーバ 3 0 0 f と通信を行うための通信ケーブルが接続される。通信ポート 3 4 5 a , 3 4 5 b , 3 4 5 c は、それぞれ、計算ユニット 3 3 0 , 3 3 0 a , 3 3 0 b との間でパケットの入出力を行う。

【 0 0 9 9 】

通信ポート 3 4 6 a , 3 4 6 b , 3 4 6 c には、Z 軸の正方向に隣接するクラスタサーバ 3 0 0 e と通信を行うための通信ケーブルが接続される。通信ポート 3 4 6 a , 3 4 6 b , 3 4 6 c は、それぞれ、計算ユニット 3 3 0 , 3 3 0 a , 3 3 0 b との間でパケットの入出力を行う。

【 0 1 0 0 】

図 1 3 は、第 3 の実施の形態の計算ユニットの構成を示す図である。なお、図 1 3 には計算ユニット 3 1 0 の構成を示したが、計算ユニット 3 1 0 a , 3 1 0 b , 3 2 0 , 3 2 0 a , 3 2 0 b , 3 3 0 , 3 3 0 a , 3 3 0 b も同様の構成で実現できる。計算ユニット 3 1 0 は、演算回路 3 1 1、入力回路 3 1 2、出力回路 3 1 3、受信回路 3 1 4 a , 3 1 4 b , 3 1 6 a ~ 3 1 6 h、送信回路 3 1 5 a , 3 1 5 b , 3 1 7 a ~ 3 1 7 h およびスイッチ回路 3 1 8 を有している。

【 0 1 0 1 】

演算回路 3 1 1、入力回路 3 1 2 および出力回路 3 1 3 は、図 5 に示した第 1 の実施の形態の演算回路 1 1 1、入力回路 1 1 2 および出力回路 1 1 3 とそれぞれ同様の機能を備えている。

【 0 1 0 2 】

受信回路 3 1 4 a , 3 1 4 b は、それぞれ内部に複数の F I F O 型のバッファメモリを仮想チャネルバッファとして有している。受信回路 3 1 4 a は、通信ポート 3 4 1 a からのパケット、すなわち、クラスタサーバ 3 0 0 b からのパケットを取得する。受信回路 3 1 4 b は、通信ポート 3 4 2 a からのパケット、すなわち、クラスタサーバ 3 0 0 a からのパケットを取得する。複数の仮想チャネルバッファを用いて仮想チャネルを実現する方法は、第 1 の実施の形態の方法と同様である。

【 0 1 0 3 】

送信回路 3 1 5 a は、スイッチ回路 3 1 8 からパケットを取得すると、取得したパケットを通信ポート 3 4 2 a へ出力する。これにより、パケットがクラスタサーバ 3 0 0 a へ送信される。送信回路 3 1 5 b は、スイッチ回路 3 1 8 からパケットを取得すると、取得したパケットを通信ポート 3 4 1 a へ出力する。これにより、パケットがクラスタサーバ 3 0 0 b へ送信される。

【 0 1 0 4 】

受信回路 3 1 6 a ~ 3 1 6 h は、それぞれ内部に F I F O 型のバッファメモリを有している。受信回路 3 1 6 a は、計算ユニット 3 1 0 a からのパケットを取得する。受信回路 3 1 6 b は、計算ユニット 3 1 0 b からのパケットを取得する。受信回路 3 1 6 c は、計算ユニット 3 2 0 からのパケットを取得する。受信回路 3 1 6 d は、計算ユニット 3 2 0 a からのパケットを取得する。受信回路 3 1 6 e は、計算ユニット 3 2 0 b からのパケットを取得する。受信回路 3 1 6 f は、計算ユニット 3 3 0 からのパケットを取得する。受信回路 3 1 6 g は、計算ユニット 3 3 0 a からのパケットを取得する。受信回路 3 1 6 h

10

20

30

40

50

は、計算ユニット330bからのパケットを取得する。

【0105】

送信回路317a~317hは、スイッチ回路318からパケットを取得する。送信回路317aは、取得したパケットを計算ユニット310aへ出力する。送信回路317bは、取得したパケットを計算ユニット310bへ出力する。送信回路317cは、取得したパケットを計算ユニット320へ出力する。送信回路317dは、取得したパケットを計算ユニット320aへ出力する。送信回路317eは、取得したパケットを計算ユニット320bへ出力する。送信回路317fは、取得したパケットを計算ユニット330へ出力する。送信回路317gは、取得したパケットを計算ユニット330aへ出力する。送信回路317hは、取得したパケットを計算ユニット330bへ出力する。

10

【0106】

スイッチ回路318は、出力回路313および受信回路314a, 314b, 316a~316hのバッファメモリを監視する。バッファメモリにパケットが格納されている場合、スイッチ回路318は、出力回路313および受信回路314a, 314b, 316a~316hからパケットを順次取得する。そして、スイッチ回路318は、パケットの宛先アドレスに基づいて、パケットの出力先を決定する。パケットの出力先は、入力回路312および送信回路315a, 315b, 317a~317hのいずれかである。

【0107】

図14は、第3の実施の形態のパケット出力先の決定方法を示す図である。なお、図14にはスイッチ回路318によるパケット出力先の決定方法を示したが、他の計算ユニットのスイッチ回路も同様の方法でパケット出力先を決定することができる。その場合、図14に示した宛先アドレスを適宜読み替えばよい。

20

【0108】

スイッチ回路318は、パケットの宛先アドレスが計算ユニット310のアドレスと一致する場合、更に、パケットの取得元が出力回路313か否かを判断する。パケットの取得元が出力回路313の場合、パケットを破棄する。一方、パケットの取得元が出力回路313でない場合、パケットを入力回路312へ出力する。宛先アドレスが計算ユニット310a, 310b, 320, 320a, 320b, 330, 330a, 330bの場合、それぞれ、パケットを送信回路317a~317hへ出力する。

【0109】

宛先アドレスと計算ユニット310のアドレスとでX座標が異なり、X軸の正方向へ転送することで距離が縮まる場合、スイッチ回路318は、パケットを送信回路315aへ出力する。宛先アドレスと計算ユニット310のアドレスとでX座標が異なり、X軸の負方向へ転送することで距離が縮まる場合、パケットを送信回路315bへ出力する。

30

【0110】

宛先アドレスと計算ユニット310のアドレスとでX座標が同一でY座標が異なる場合、スイッチ回路318は、パケットを送信回路317cへ出力する。これにより、パケットが計算ユニット320へ転送される。宛先アドレスと計算ユニット310のアドレスとでX座標およびY座標が同一でZ座標が異なる場合、スイッチ回路318は、パケットを送信回路317fへ出力する。これにより、パケットが計算ユニット330へ転送される。

40

【0111】

なお、図14に示した方法では、X軸方向への転送をY軸方向への転送よりも優先し、Y軸方向への転送をZ軸方向の転送よりも優先するように統一しているが、優先する座標軸を変更してもよい。

【0112】

このように、第3の実施の形態に係る並列計算機システムによれば、宛先アドレスで示される計算ユニットに到達するまでのパケットの中継回数を削減することができる。また、3次元座標空間における6方向へのパケットの送信を9つの計算ユニットで分担するため、特定の回路や通信ケーブルにパケットが集中することを防止できる。特に、1つの座

50

標軸に対して複数の計算ユニットを割り当てたため、クラスタサーバ間のパケット転送能力を向上させることができる。

【0113】

また、クラスタサーバの内部では完全結合網によりパケットを直接転送するため、内部の転送に関してデッドロック防止のための仮想チャネルバッファを設ける必要がなく、回路規模の増大を抑止できる。また、各計算ユニットの構成を均一にすることができ、クラスタサーバを実現する構成を単純化できる。

【0114】

なお、第1の実施の形態および第3の実施の形態では、クラスタサーバ内でパケットを転送する通信網として完全結合網を採用し、第2の実施の形態では、PCクラスタ内でパケットを転送する通信網としてスター型結合網を採用した。これに対し、他の形状の通信網を考えられる。

10

【0115】

図15は、クラスタ内でパケットを転送する他の通信網の例を示す模式図である。図15に示すクラスタは、3次元のトーラス型相互結合網の各格子点に配置されるものである。このクラスタは、6個の計算ユニット410、420、430、440、450、460で構成される。

【0116】

計算ユニット410は、X軸の正方向へのパケットの中継を担当する。計算ユニット420は、Y軸の正方向へのパケットの中継を担当する。計算ユニット430は、Z軸の正方向へのパケットの中継を担当する。計算ユニット440は、X軸の負方向へのパケットの中継を担当する。計算ユニット450は、Y軸の負方向へのパケットの中継を担当する。計算ユニット460は、Z軸の負方向へのパケットの中継を担当する。

20

【0117】

ここで、計算ユニット410、420、430が3本の通信ケーブルで相互に接続されている。また、計算ユニット440、450、460が3本の通信ケーブルで相互に接続されている。また、計算ユニット410と計算ユニット440、計算ユニット420と計算ユニット450、計算ユニット430と計算ユニット460がそれぞれ通信ケーブルで接続されている。

【0118】

このように、正方向同士および負方向同士を相互に接続したのは、パケットの中継において、正方向もしくは負方向への中継が連続して行われる可能性が高いと考えられるためである。これにより、クラスタ内でのパケットの平均中継回数を抑えつつ、通信ケーブルの本数および計算ユニット毎の通信ポートの個数を小さくすることができる。

30

【0119】

例えば、6個の計算ユニットを全て相互に接続した場合、15本の通信ケーブルと、計算ユニット毎に5個の通信ポートとを設ける必要がある。これに対し、図15に示した構成では、9本の通信ケーブルと、計算ユニット毎に3個の通信ポートとを設ければよい。

【0120】

以上、本発明の計算機クラスタを図示の実施の形態に基づいて説明したが、本発明はこれに限定されるものではなく、各部の構成は同様の機能を有する任意の構成のものに置換することができる。また、本発明に他の任意の構成物や工程が付加されていてもよい。また、本発明は、前述した実施の形態のうちの任意の2以上の構成(特徴)を組み合わせたものであってもよい。

40

【0121】

以上説明した実施の形態の主な技術的特徴は、以下の付記の通りである。

(付記1) 格子状に構成された相互結合網の各格子点に配置される計算機クラスタにおいて、

格子点の位置を示す外部アドレスと前記計算機クラスタ内における相対位置を示す内部アドレスとから成る宛先アドレスを含むパケットの処理を行う演算部と、

50

前記計算機クラスタの外部の通信相手と前記パケットを送受信する第1の通信部と、
 前記計算機クラスタの内部の通信相手と前記パケットを送受信する第2の通信部と、
 前記演算部、前記第1の通信部および前記第2の通信部と接続されており、取得した前記パケットを、前記宛先アドレスが自己のアドレスと一致する場合は前記演算部へ出力し、前記外部アドレスが前記自己のアドレスと異なり前記パケットを転送すべき方向が前記第1の通信部が送信可能な方向である場合は前記第1の通信部に転送させ、前記内部アドレスのみが前記自己のアドレスと異なる場合および前記外部アドレスが前記自己のアドレスと異なり前記パケットを転送すべき方向が前記第1の通信部が送信可能な方向でない場合は前記第2の通信部に転送させるスイッチ部と、

を有する第1のノード、第2のノード、第3のノードおよび第4のノードと、

前記第1のノード、前記第2のノード、前記第3のノードおよび前記第4のノードの前記第2の通信部を相互に接続する内部通信網と、

前記第1のノードの前記第1の通信部と格子のX軸方向に隣接する2つの前記計算機クラスタとを接続する第1の外部通信網と、

前記第2のノードの前記第1の通信部と格子のY軸方向に隣接する2つの前記計算機クラスタとを接続する第2の外部通信網と、

前記第3のノードの前記第1の通信部と格子のX軸方向に隣接する2つの前記計算機クラスタとを接続する第3の外部通信網と、

前記第4のノードの前記第1の通信部と格子のY軸方向に隣接する2つの前記計算機クラスタとを接続する第4の外部通信網と、

を有することを特徴とする計算機クラスタ。

【0122】

(付記2) 前記内部通信網は、全ての前記第2の通信部を双方向リンクによって直接接続した完全結合網であることを特徴とする付記1記載の計算機クラスタ。

(付記3) 前記内部接続網は、受信した前記パケットを振り分けるスイッチ装置と、前記スイッチ装置と前記第1のノード、前記第2のノード、前記第3のノードおよび前記第4のノードの前記第2の通信部とをそれぞれ接続する4つの双方向リンクとで構成されることを特徴とする付記1記載の計算機クラスタ。

【0123】

(付記4) 前記第1の外部通信網は、X軸の負側に隣接する前記計算機クラスタから前記第1のノードの前記第1の通信部への片方向リンクと、前記第1のノードの前記第1の通信部からX軸の正側に隣接する前記計算機クラスタへの片方向リンクとで構成され、

前記第2の外部通信網は、Y軸の負側に隣接する前記計算機クラスタから前記第2のノードの前記第1の通信部への片方向リンクと、前記第2のノードの前記第1の通信部からY軸の正側に隣接する前記計算機クラスタへの片方向リンクとで構成され、

前記第3の外部通信網は、X軸の正側に隣接する前記計算機クラスタから前記第3のノードの前記第1の通信部への片方向リンクと、前記第3のノードの前記第1の通信部からX軸の負側に隣接する前記計算機クラスタへの片方向リンクとで構成され、

前記第4の外部通信網は、Y軸の正側に隣接する前記計算機クラスタから前記第4のノードの前記第1の通信部への片方向リンクと、前記第4のノードの前記第1の通信部からY軸の負側に隣接する前記計算機クラスタへの片方向リンクとで構成される、

ことを特徴とする付記1記載の計算機クラスタ。

【0124】

(付記5) 前記第1の外部通信網は、前記第1のノードの前記第1の通信部とX軸方向に隣接する2つの前記計算機クラスタとをそれぞれ接続する2つの双方向リンクで構成され、

前記第2の外部通信網は、前記第2のノードの前記第1の通信部とY軸方向に隣接する2つの前記計算機クラスタとをそれぞれ接続する2つの双方向リンクで構成され、

前記第3の外部通信網は、前記第3のノードの前記第1の通信部とX軸方向に隣接する2つの前記計算機クラスタとをそれぞれ接続する2つの双方向リンクで構成され、

10

20

30

40

50

前記第 4 の外部通信網は、前記第 4 のノードの前記第 1 の通信部と Y 軸方向に隣接する 2 つの前記計算機クラスタとをそれぞれ接続する 2 つの双方向リンクで構成される、
ことを特徴とする付記 1 記載の計算機クラスタ。

【 0 1 2 5 】

(付記 6) 前記第 1 の通信部は、前記パケットを一時的に保持する複数の仮想チャネルバッファを有しており、受信した前記パケットを前記宛先アドレスに応じて前記複数の仮想チャネルバッファに振り分けて格納し、

前記スイッチ部は、前記複数の仮想チャネルバッファに格納された前記パケットを順次取得して、スイッチング処理を行う、

ことを特徴とする付記 1 記載の計算機クラスタ。

10

【 0 1 2 6 】

(付記 7) 前記相互結合網は 3 次元の格子状に構成されており、

前記演算部、前記第 1 の通信部、前記第 2 の通信部および前記スイッチ部を有する第 5 のノードおよび第 6 のノードと、

前記第 5 のノードの前記第 1 の通信部と格子の Z 軸方向に隣接する 2 つの前記計算機クラスタとを接続する第 5 の外部通信網と、

前記第 6 のノードの前記第 1 の通信部と格子の Z 軸方向に隣接する 2 つの前記計算機クラスタとを接続する第 6 の外部通信網と、

を更に有し、

前記内部通信網は、前記第 1 のノード、前記第 2 のノード、前記第 3 のノード、前記第 4 のノード、前記第 5 のノードおよび前記第 6 のノードの前記第 2 の通信部を相互に接続する、

20

ことを特徴とする付記 1 記載の計算機クラスタ。

【図面の簡単な説明】

【 0 1 2 7 】

【図 1】本実施の形態の概要を示す図である。

【図 2】第 1 の実施の形態のシステム構成を示す図である。

【図 3】第 1 の実施の形態のクラスタサーバの構成を示す図である。

【図 4】第 1 の実施の形態の通信ケーブルの接続関係を示す図である。

【図 5】第 1 の実施の形態の計算ユニットの構成を示す図である。

30

【図 6】第 1 の実施の形態のパケット出力先の決定方法を示す図である。

【図 7】第 2 の実施の形態のシステム構成を示す図である。

【図 8】第 2 の実施の形態の PC クラスタの構成を示す図である。

【図 9】第 2 の実施の形態の PC サーバの構成を示す図である。

【図 10】第 2 の実施の形態のパケット出力先の決定方法を示す図である。

【図 11】第 3 の実施の形態のシステム構成を示す図である。

【図 12】第 3 の実施の形態のクラスタサーバの構成を示す図である。

【図 13】第 3 の実施の形態の計算ユニットの構成を示す図である。

【図 14】第 3 の実施の形態のパケット出力先の決定方法を示す図である。

【図 15】クラスタ内でパケットを転送する他の通信網の例を示す模式図である。

40

【符号の説明】

【 0 1 2 8 】

1 0 , 2 0 , 3 0 , 4 0 計算機クラスタ

1 1 , 2 1 , 3 1 , 4 1 第 1 のノード

1 2 , 2 2 , 3 2 , 4 2 第 2 のノード

1 3 , 2 3 , 3 3 , 4 3 第 3 のノード

1 4 , 2 4 , 3 4 , 4 4 第 4 のノード

1 5 内部通信網

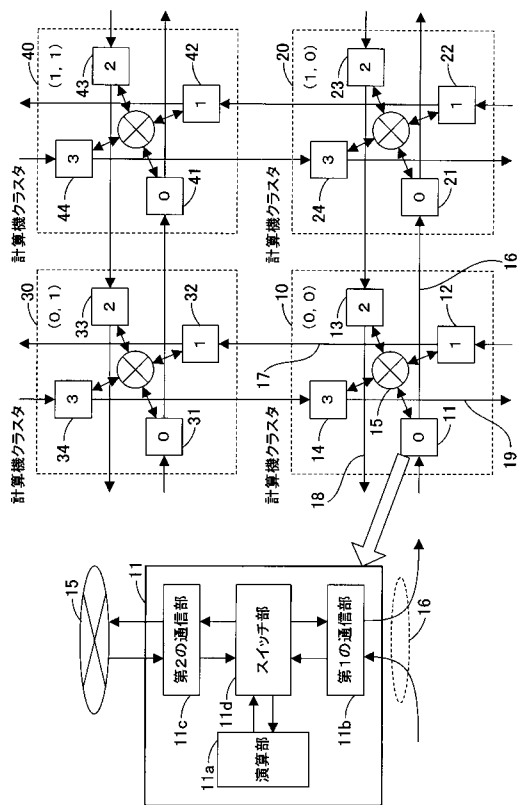
1 6 第 1 の外部通信網

1 7 第 2 の外部通信網

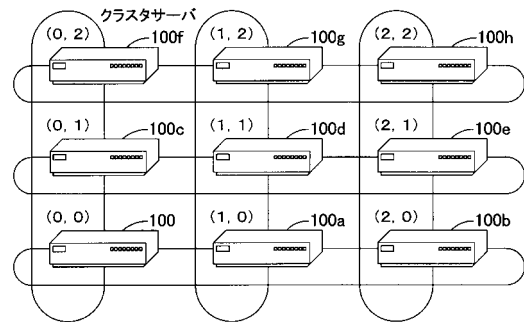
50

- 1 8 第 3 の外部通信網
- 1 9 第 4 の外部通信網
- 1 1 a 演算部
- 1 1 b 第 1 の通信部
- 1 1 c 第 2 の通信部
- 1 1 d スイッチ部

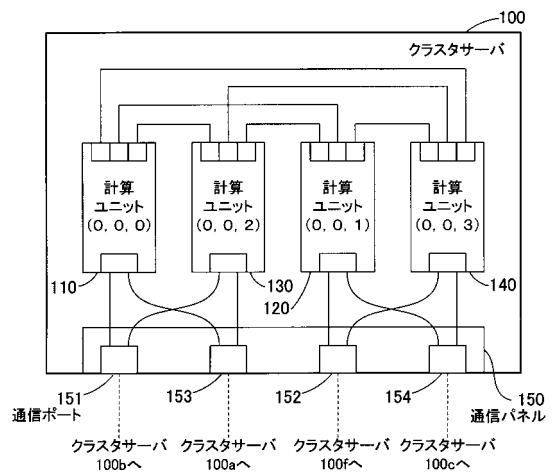
【 図 1 】



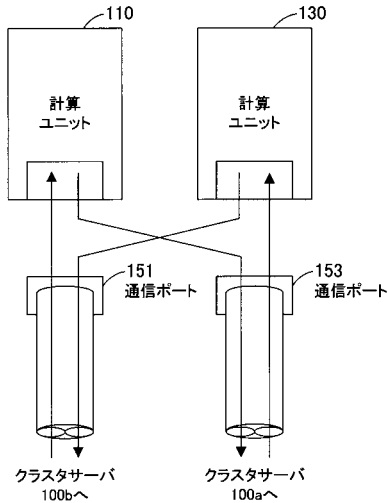
【 図 2 】



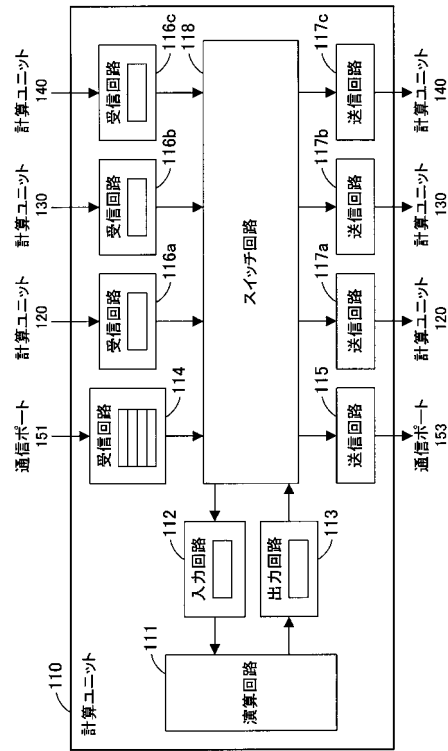
【 図 3 】



【 図 4 】



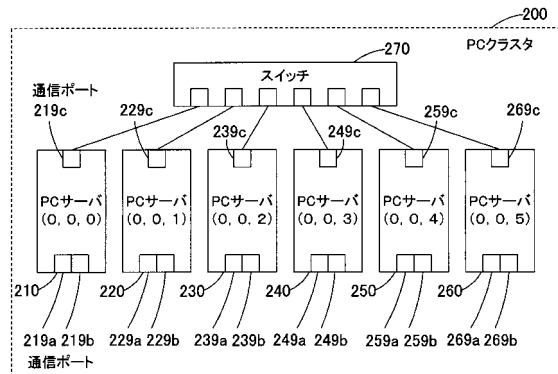
【 図 5 】



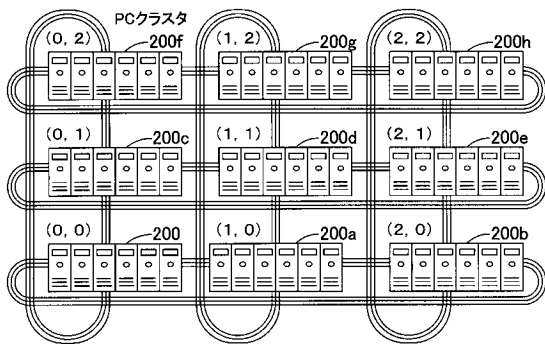
【 図 6 】

宛先アドレス		出力先	
自己	(0, 0, 0)	出力回路113から	破棄
		出力回路113以外から	入力回路112
クラスタ内	(0, 0, 1)	送信回路117a	
	(0, 0, 2)	送信回路117b	
	(0, 0, 3)	送信回路117c	
X軸正方向	(1, *, *)	送信回路115	
X軸負方向	(2, *, *)	送信回路117b	
Y軸正方向	(0, 1, *)	送信回路117a	
Y軸負方向	(0, 2, *)	送信回路117c	

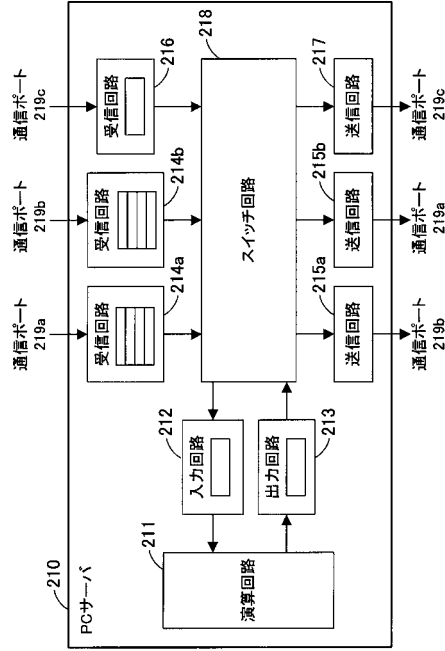
【 図 8 】



【 図 7 】



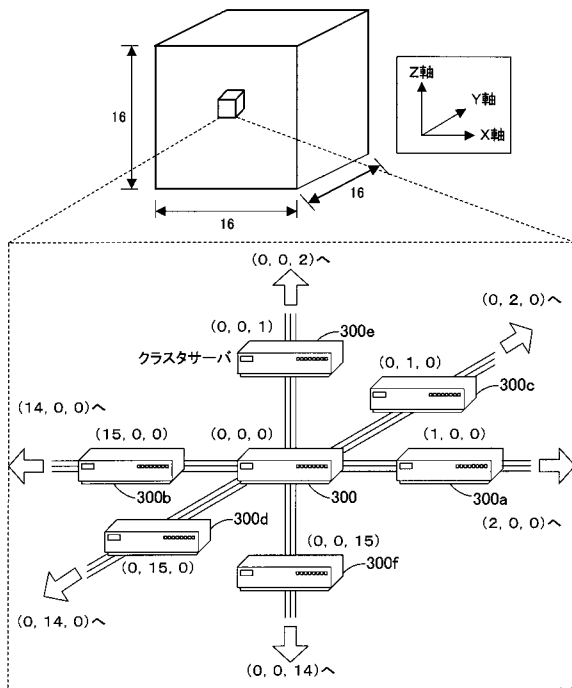
【図9】



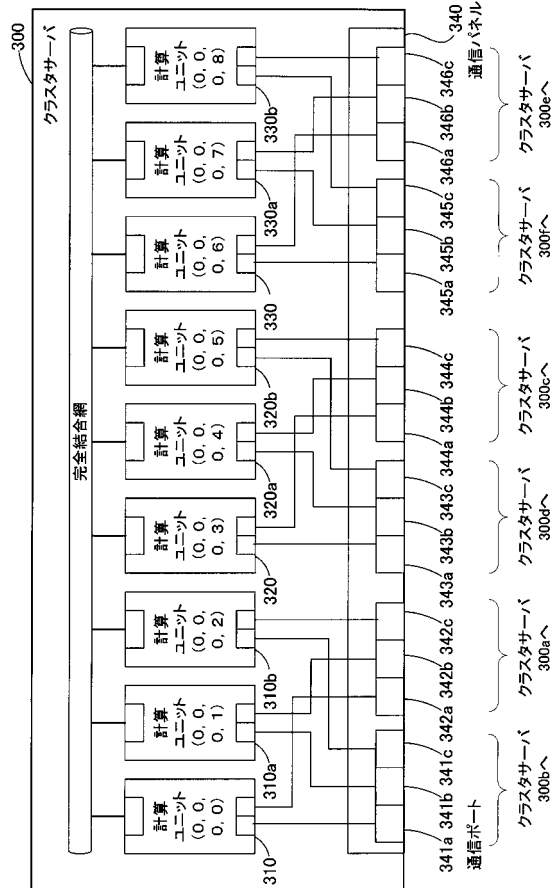
【図10】

宛先アドレス		出力先	
自己	(0, 0, 0)	出力回路213から	破棄
		出力回路213以外から	入力回路212
クラスタ内	(0, 0, 1)	送信回路217 → PCサーバ220	
	(0, 0, 2)	送信回路217 → PCサーバ230	
	(0, 0, 3)	送信回路217 → PCサーバ240	
	(0, 0, 4)	送信回路217 → PCサーバ250	
	(0, 0, 5)	送信回路217 → PCサーバ260	
X軸正方向	(1, *, *)	送信回路215a	
X軸負方向	(2, *, *)	送信回路215b	
Y軸方向	(0, 1, *)	送信回路217 → PCサーバ220	
	(0, 2, *)		

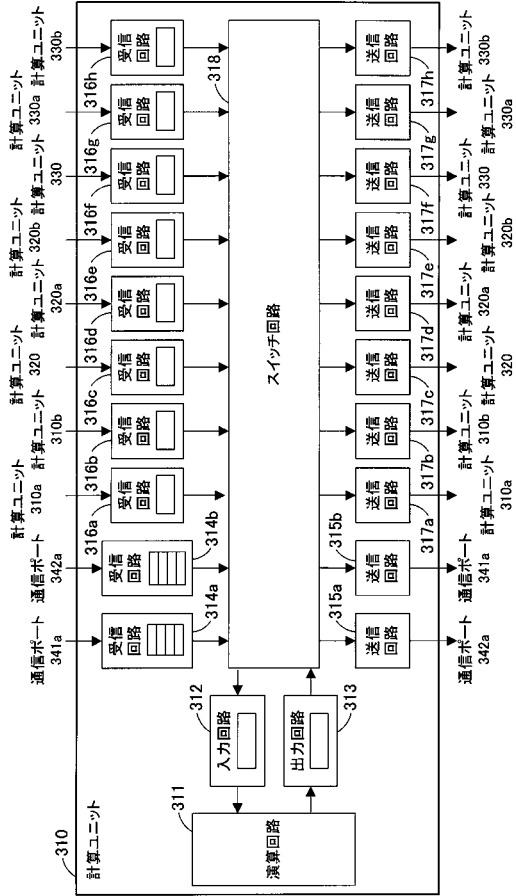
【図11】



【図12】



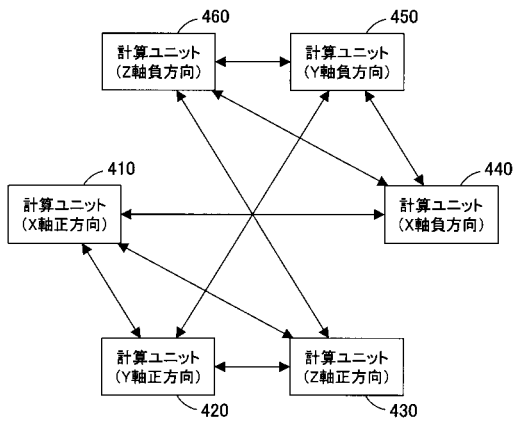
【図13】



【図14】

宛先アドレス		出力先	
自己	(0, 0, 0, 0)	出力回路313から	破棄
		出力回路313以外から	入力回路312
クラスタ内	(0, 0, 0, 1)	送信回路317a	
	(0, 0, 0, 2)	送信回路317b	
	(0, 0, 0, 3)	送信回路317c	
	(0, 0, 0, 4)	送信回路317d	
	(0, 0, 0, 5)	送信回路317e	
	(0, 0, 0, 6)	送信回路317f	
	(0, 0, 0, 7)	送信回路317g	
	(0, 0, 0, 8)	送信回路317h	
X軸正方向	(1, *, *, *)	送信回路315a	
	(8, *, *, *)	送信回路315a	
X軸負方向	(9, *, *, *)	送信回路315b	
	(15, *, *, *)	送信回路315b	
Y軸方向	(0, 1, *, *)	送信回路317c	
	(0, 15, *, *)	送信回路317c	
Z軸方向	(0, 0, 1, *)	送信回路317f	
	(0, 0, 15, *)	送信回路317f	

【図15】



フロントページの続き

- (56)参考文献 特開平9 - 160893 (JP, A)
特開平7 - 191947 (JP, A)
特開平7 - 99491 (JP, A)
特開平6 - 266684 (JP, A)
特開平6 - 35873 (JP, A)
特開平4 - 253256 (JP, A)
特開昭64 - 4856 (JP, A)
米国特許出願公開第2006/0176888 (US, A1)
田邊昇, 菅野伸一, 小柳滋, Wavefront Array 動作が可能な汎用超並列マシン向け結合網アーキテクチャ, 電子情報通信学会論文誌, 日本, 社団法人電子情報通信学会, 1995年 2月25日, Vol:J-78-D-I, No:2, Pages:99-107
水戸部理, 吉永努, 曾和将容, 適応ルーティングを用いたPCクラスタ用ネットワークスイッチの提案, 並列処理シンポジウムJSPP2002, 日本, 社団法人情報処理学会, 2002年 5月29日, Pages:187-188

(58)調査した分野(Int.Cl., DB名)

G06F15/16 - 15/177
H04L12/28
H04L12/40 - 12/417
H04L12/44 - 12/46
H04W8/26
H04W24/00
H04W72/04 - 72/04
H04W74/08
H04W84/12
H04W88/08