

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7193733号
(P7193733)

(45)発行日 令和4年12月21日(2022.12.21)

(24)登録日 令和4年12月13日(2022.12.13)

(51)国際特許分類 F I
H 0 4 L 49/201 (2022.01) H 0 4 L 49/201
H 0 4 L 47/60 (2022.01) H 0 4 L 47/60

請求項の数 6 (全25頁)

(21)出願番号	特願2019-77726(P2019-77726)	(73)特許権者	000005223 富士通株式会社
(22)出願日	平成31年4月16日(2019.4.16)		神奈川県川崎市中原区上小田中4丁目1番1号
(65)公開番号	特開2020-178180(P2020-178180 A)	(74)代理人	110002918 弁理士法人扶桑国際特許事務所
(43)公開日	令和2年10月29日(2020.10.29)	(72)発明者	関澤 龍一 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
審査請求日	令和4年1月11日(2022.1.11)	審査官	佐々木 洋

最終頁に続く

(54)【発明の名称】 通信制御プログラム、通信制御方法および情報処理装置

(57)【特許請求の範囲】

【請求項1】

コンピュータに、

複数のノードと複数の第1の中継装置と複数の第2の中継装置とを含み、前記複数のノードそれぞれが前記複数の第1の中継装置の1つと接続され、前記複数の第1の中継装置それぞれが前記複数の第2の中継装置の一部である2以上の第2の中継装置と接続されたシステムにおいて、接続された前記2以上の第2の中継装置が同一でない異なる第1の中継装置に接続された異なるノードが、異なるグループに分類されるように、前記システムに含まれる前記複数のノードを複数のグループに分類し、

前記複数のグループそれぞれから代表ノードを選択し、

前記複数のグループに対応する複数の代表ノードの間で実行される第1のブロードキャスト通信の通信順序を、1つの代表ノードが最初の送信元ノードとして動作し、データを受信した代表ノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定し、

前記複数のグループそれぞれについて、前記第1のブロードキャスト通信の後に当該グループに含まれる2以上のノードの間で実行される第2のブロードキャスト通信の通信順序を、当該グループの代表ノードが最初の送信元ノードとして動作し、データを受信したノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定する、

処理を実行させる通信制御プログラム。

【請求項 2】

前記複数のノードの分類では、接続された前記 2 以上の第 2 の中継装置が同一である異なる第 1 の中継装置に接続された異なるノードを、同じグループに分類する、

請求項 1 記載の通信制御プログラム。

【請求項 3】

前記複数のグループそれぞれの代表ノードは、当該グループに含まれる前記 2 以上のノードのうち、配置されたプロセスの識別番号が最小のノードである、

請求項 1 記載の通信制御プログラム。

【請求項 4】

前記第 1 のブロードキャスト通信および前記第 2 のブロードキャスト通信では、2 以上の送信元ノードが異なるノードに並列にデータを送信するフェーズを繰り返すことで、データを受信済みであるノードの数をフェーズ毎に 2 倍にする、

請求項 1 記載の通信制御プログラム。

【請求項 5】

コンピュータが、

複数のノードと複数の第 1 の中継装置と複数の第 2 の中継装置とを含み、前記複数のノードそれぞれが前記複数の第 1 の中継装置の 1 つと接続され、前記複数の第 1 の中継装置それぞれが前記複数の第 2 の中継装置の一部である 2 以上の第 2 の中継装置と接続されたシステムにおいて、接続された前記 2 以上の第 2 の中継装置が同一でない異なる第 1 の中継装置に接続された異なるノードが、異なるグループに分類されるように、前記システム

に含まれる前記複数のノードを複数のグループに分類し、

前記複数のグループそれぞれから代表ノードを選択し、

前記複数のグループに対応する複数の代表ノードの間で実行される第 1 のブロードキャスト通信の通信順序を、1 つの代表ノードが最初の送信元ノードとして動作し、データを受信した代表ノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定し、

前記複数のグループそれぞれについて、前記第 1 のブロードキャスト通信の後に当該グループに含まれる 2 以上のノードの間で実行される第 2 のブロードキャスト通信の通信順序を、当該グループの代表ノードが最初の送信元ノードとして動作し、データを受信したノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するよ

うに決定する、
通信制御方法。

【請求項 6】

複数のノードと複数の第 1 の中継装置と複数の第 2 の中継装置とを含み、前記複数のノードそれぞれが前記複数の第 1 の中継装置の 1 つと接続され、前記複数の第 1 の中継装置それぞれが前記複数の第 2 の中継装置の一部である 2 以上の第 2 の中継装置と接続されたシステムにおいて、前記複数のノードの間のブロードキャスト通信の通信順序を示す通信制御データを記憶する記憶部と、

前記ブロードキャスト通信の通信順序を決定する処理部と、

を有し、前記処理部は、

接続された前記 2 以上の第 2 の中継装置が同一でない異なる第 1 の中継装置に接続された異なるノードが、異なるグループに分類されるように、前記システムに含まれる前記複数のノードを複数のグループに分類し、

前記複数のグループそれぞれから代表ノードを選択し、

前記複数のグループに対応する複数の代表ノードの間で実行される第 1 のブロードキャスト通信の通信順序を、1 つの代表ノードが最初の送信元ノードとして動作し、データを受信した代表ノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定し、

前記複数のグループそれぞれについて、前記第 1 のブロードキャスト通信の後に当該グループに含まれる 2 以上のノードの間で実行される第 2 のブロードキャスト通信の通信順

10

20

30

40

50

序を、当該グループの代表ノードが最初の送信元ノードとして動作し、データを受信したノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定する、

情報処理装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は通信制御プログラム、通信制御方法および情報処理装置に関する。

【背景技術】

【0002】

複数の情報処理装置をノードとして含む並列処理システムがある。並列処理システムは、同一のジョブに属する複数のプロセスを複数のノードに割り振り、それら複数のプロセスを複数のノードにより並列に実行する。ジョブの中ではノード間で通信を行うことがある。ノード間で通信を行うユーザプログラムは、MPI (Message Passing Interface) ライブラリなどの通信ライブラリを利用して実装されることがある。ノード間の通信としては、ジョブに使用される複数のノードが一斉にデータ送信に参加するコレクティブ通信 (集合通信または集団通信と言うこともある) がある。コレクティブ通信には、1つのノードから複数の他のノードに同一データをコピーするブロードキャスト通信が含まれる。

【0003】

ところで、並列処理システムが多数のノードを含む場合、単一のスイッチなど単一の中継装置に全てのノードを直接接続することは難しい。そのため、複数のノードおよび複数の中継装置の接続形態を示すネットワークトポロジ (単にトポロジと言うことがある) が問題となる。1つのノードと別の1つのノードとの間の通信が、2以上の中継装置を経由することがある。並列処理システムのトポロジの選択では、ノード間の通信経路の冗長性や、中継装置の個数などのコストが考慮されることがある。

【0004】

並列処理システムの1つとして、多層フルメッシュトポロジをもつ多層フルメッシュシステムが提案されている。提案の多層フルメッシュシステムは、複数のノードと複数のLeaf (葉) スイッチと複数のSpine (背骨) スイッチを含み、複数の層 (レイヤ) を形成する。各ノードは何れか1つのLeafスイッチに接続され、各Leafスイッチは何れか1つの層に属し、各Spineスイッチは複数の層を貫通する。

【0005】

層内では、2以上のLeafスイッチがフルメッシュトポロジにより接続されている。Leafスイッチのペア毎に、他のLeafスイッチを経由しない通信経路が存在する。ただし、Leafスイッチのペア毎に、その間に1つのSpineスイッチが配置されている。よって、層内では、1つのLeafスイッチは別の1つのLeafスイッチと、1つのSpineスイッチを経由して通信することになる。このSpineスイッチは複数の層を接続している。よって、当該1つのLeafスイッチは別の層に属する1つのLeafスイッチとも、1つのSpineスイッチを経由して通信できる。

【0006】

なお、ツリー状に接続された複数の中継装置を用いて1つの送信端末から複数の受信端末にデータを配信するデータ配信システムが提案されている。提案のデータ配信システムは、各受信端末の属性情報を、ツリーの末端からルートに向かって転送することで送信端末に集約する。送信端末が属性条件を指定したパケットを出力すると、各中継装置は、属性条件に合致する受信端末が存在する方向にパケットを転送する。

【先行技術文献】

【特許文献】

【0007】

【文献】特開2018-26657号公報

特開2018-185650号公報

10

20

30

40

50

【発明の概要】**【発明が解決しようとする課題】****【0008】**

ブロードキャスト通信のアルゴリズムには、始点ノードからコピーされたデータを受信したノードが、その後は送信元ノードとして動作することで、並列にデータを送信するノードが増加していくものがある。例えば、ブロードキャスト通信のアルゴリズムとして *Binomial Tree* (二項木) アルゴリズムがある。 *Binomial Tree* アルゴリズムでは、第1フェーズにおいて、プロセス0がプロセス1にデータを送信し、第2フェーズにおいて、プロセス0がプロセス2にデータを送信すると共に、プロセス1がプロセス3にデータを送信する。これにより、送信元ノード数が2の累乗で増加する。

10

【0009】

しかし、並列処理システムのトポロジによっては、ブロードキャスト通信において並列にデータを送信するノードが増加すると、通信が競合するリスクが高くなるという問題がある。プロセス0がプロセス2にデータを送信し、これと並列にプロセス1がプロセス3にデータを送信するにあたり、2つの通信経路が同じリンクを使用することになる場合、通信が競合し得る。通信が競合すると、パケットの送信待ちが発生する、1つのリンクの通信帯域が分割されるなどにより、通信遅延が発生して通信時間が長くなる。

【0010】

例えば、前述の多層フルメッシュシステムでは、2つの *Leaf* スイッチの間に存在する最短経路の個数は、当該2つの *Leaf* スイッチに共通に接続されている *Spine* スイッチの個数に相当する。同じ層に属する2つの *Leaf* スイッチの間にある最短経路は1つのみである。よって、1つの *Leaf* スイッチの配下にプロセス0, 1が配置され、別の1つの *Leaf* スイッチの配下にプロセス2, 3が配置されている場合、プロセス0からプロセス2への通信とプロセス1からプロセス3への通信が競合することがある。

20

【0011】

1つの側面では、本発明は、ノード間の通信の競合を抑制できる通信制御プログラム、通信制御方法および情報処理装置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】**【0012】**

1つの態様では、コンピュータに以下の処理を実行させる通信制御プログラムが提供される。複数のノードと複数の第1の中継装置と複数の第2の中継装置とを含み、複数のノードそれぞれが複数の第1の中継装置の1つと接続され、複数の第1の中継装置それぞれが複数の第2の中継装置の一部である2以上の第2の中継装置と接続されたシステムにおいて、接続された2以上の第2の中継装置が同一でない異なる第1の中継装置に接続された異なるノードが、異なるグループに分類されるように、システムに含まれる複数のノードを複数のグループに分類する。複数のグループそれぞれから代表ノードを選択する。複数のグループに対応する複数の代表ノードの間で実行される第1のブロードキャスト通信の通信順序を、1つの代表ノードが最初の送信元ノードとして動作し、データを受信した代表ノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定する。複数のグループそれぞれについて、第1のブロードキャスト通信の後

に当該グループに含まれる2以上のノードの間で実行される第2のブロードキャスト通信の通信順序を、当該グループの代表ノードが最初の送信元ノードとして動作し、データを受信したノードが送信元ノードに加わることで並列にデータを送信する送信元ノードが増加するように決定する。

30

40

【0013】

また、1つの態様では、コンピュータが実行する通信制御方法が提供される。また、1つの態様では、記憶部と処理部とを有する情報処理装置が提供される。

【発明の効果】**【0014】**

1つの側面では、ノード間の通信の競合を抑制できる。

50

【図面の簡単な説明】

【0015】

【図1】第1の実施の形態の情報処理システムの例を説明する図である。

【図2】第2の実施の形態の多層フルメッシュシステムの例を示す図である。

【図3】多層フルメッシュシステムの配線例を示す図である。

【図4】サーバのハードウェア例を示すブロック図である。

【図5】スイッチのハードウェア例を示すブロック図である。

【図6】Binomial Treeアルゴリズムの例を示すシーケンス図である。

【図7】ブロードキャスト通信の競合例を示す図である。

【図8】プロセスの配置例を示す図である。

10

【図9】二段階ブロードキャスト通信の例を示すシーケンス図である。

【図10】グループ間通信の競合回避例を示す図である。

【図11】グループ内通信の競合回避例を示す図である。

【図12】サーバとジョブスケジューラの機能例を示す図である。

【図13】プロセス配置テーブルの例を示す図である。

【図14】通信手順テーブルの例を示す図である。

【図15】通信手順決定の処理例を示すフローチャートである。

【図16】ブロードキャスト通信の処理例を示すフローチャートである。

【発明を実施するための形態】

【0016】

20

以下、本実施の形態を図面を参照して説明する。

〔第1の実施の形態〕

第1の実施の形態を説明する。

【0017】

図1は、第1の実施の形態の情報処理システムの例を説明する図である。

第1の実施の形態の情報処理システムは、並列に情報処理を行う複数のノードの間のブロードキャスト通信を制御する。ブロードキャスト通信を行う対象のシステムは、例えば、多層フルメッシュトポロジをもつ多層フルメッシュシステムである。ただし、対象のシステムは、後述する構成を備えていれば多層フルメッシュシステムでなくてもよい。

【0018】

30

第1の実施の形態の情報処理システムは、情報処理装置10を含む。情報処理装置10は、ブロードキャスト通信を行う対象のシステムを制御するジョブスケジューラなどの制御装置でもよいし、対象のシステムに含まれるノードの1つであってもよい。

【0019】

情報処理装置10は、記憶部および処理部を有する。記憶部は、RAM (Random Access Memory) などの揮発性メモリでもよいし、HDD (Hard Disk Drive) やフラッシュメモリなどの不揮発性ストレージでもよい。処理部は、例えば、CPU (Central Processing Unit)、GPU (Graphics Processing Unit)、DSP (Digital Signal Processor) などのプロセッサである。ただし、処理部は、ASIC (Application Specific Integrated Circuit) やFPGA (Field Programmable Gate Array) などの特定用途の電子回路を含んでもよい。プロセッサは、メモリに記憶されたプログラムを実行する。複数のプロセッサの集合を「マルチプロセッサ」または単に「プロセッサ」と言うことがある。

40

【0020】

ブロードキャスト通信を行う対象のシステムは、ノード11, 12, 13, 14, 15, 16, 17, 18を含む複数のノードを有する。また、対象のシステムは、中継装置21, 22, 23, 24, 25, 26, 27, 28を含む複数の中継装置を有する。中継装置21, 22, 23, 24は、下位の中継装置(第1の中継装置)である。中継装置25, 26, 27, 28は、上位の中継装置(第2の中継装置)である。中継装置21, 22, 23, 24, 25, 26, 27, 28は、接続関係に応じてデータを転送する。

50

【 0 0 2 1 】

複数のノードはそれぞれ、複数の第 1 の中継装置のうちの 1 つと接続される。複数の第 1 の中継装置はそれぞれ、複数の第 2 の中継装置の一部である 2 以上の第 2 の中継装置と接続される。図 1 の例では、ノード 1 1 , 1 2 が中継装置 2 1 に接続され、ノード 1 3 , 1 4 が中継装置 2 2 に接続され、ノード 1 5 , 1 6 が中継装置 2 3 に接続され、ノード 1 7 , 1 8 が中継装置 2 4 に接続される。中継装置 2 1 が中継装置 2 5 , 2 6 に接続され、中継装置 2 2 が中継装置 2 5 , 2 6 に接続され、中継装置 2 3 が中継装置 2 6 , 2 7 に接続され、中継装置 2 4 が中継装置 2 6 , 2 7 に接続される。

【 0 0 2 2 】

情報処理装置 1 0 は、対象のシステム上で実行されるブロードキャスト通信について、複数のノードの間の通信順序を決定する。ブロードキャスト通信は、例えば、複数のフェーズに分割して実行される。1 つのフェーズの中で、2 以上のノードが並列にデータを送信することがある。ブロードキャスト通信の通信順序の決定では、例えば、各フェーズにおいて何れのノードが送信元ノードとして動作するかが決定される。

10

【 0 0 2 3 】

まず、情報処理装置 1 0 は、対象のシステムに含まれる複数ノードを複数のグループに分類する。このとき、情報処理装置 1 0 は、複数の第 1 の中継装置を、接続された 2 以上の第 2 の中継装置の同一性に基いて複数のグループに分類する。接続された 2 以上の第 2 の中継装置が一致しない第 1 の中継装置を、異なるグループに振り分けるようにする。情報処理装置 1 0 は、あるグループに分類した第 1 の中継装置に接続されたノードを、当該グループに属するノードと判定する。1 つのノードは何れか 1 つのグループに属する。

20

【 0 0 2 4 】

図 1 の例では、中継装置 2 1 と中継装置 2 2 は共に、上位の中継装置として中継装置 2 5 , 2 6 と接続されている。中継装置 2 3 と中継装置 2 4 は共に、上位の中継装置として中継装置 2 6 , 2 7 と接続されている。そこで、中継装置 2 1 と中継装置 2 2 は同一のグループに分類してもよい。また、中継装置 2 3 と中継装置 2 4 は同一のグループに分類してもよい。一方、中継装置 2 1 と中継装置 2 3 , 2 4 は異なるグループに分類する。また、中継装置 2 2 と中継装置 2 3 , 2 4 は異なるグループに分類する。

【 0 0 2 5 】

ここでは、情報処理装置 1 0 は、中継装置 2 1 に接続されたノード 1 1 , 1 2 および中継装置 2 2 に接続されたノード 1 3 , 1 4 を、グループ 3 1 に分類する。また、情報処理装置 1 0 は、中継装置 2 3 に接続されたノード 1 5 , 1 6 および中継装置 2 4 に接続されたノード 1 7 , 1 8 を、グループ 3 2 に分類する。ただし、ノード 1 1 , 1 2 とノード 1 3 , 1 4 を異なるグループに分類することも許容される。また、ノード 1 5 , 1 6 とノード 1 7 , 1 8 を異なるグループに分類することも許容される。

30

【 0 0 2 6 】

次に、情報処理装置 1 0 は、複数のグループそれぞれから代表ノードを選択する。ここでは、各グループから 1 つのノードが代表ノードとして選択されればよく、何らかの選択基準を予め定めておけばよい。例えば、情報処理装置 1 0 は、グループ内の 2 以上のノードのうち、実行するプロセスの識別番号が最小のノードを選択する。図 1 の例では、情報処理装置 1 0 は、グループ 3 1 からノード 1 1 を代表ノードとして選択する。また、情報処理装置 1 0 は、グループ 3 2 からノード 1 5 を代表ノードとして選択する。

40

【 0 0 2 7 】

次に、情報処理装置 1 0 は、複数のグループに対応する複数の代表ノードの間で実行されるブロードキャスト通信 3 3 (第 1 のブロードキャスト通信) の通信順序を決定する。また、情報処理装置 1 0 は、複数のグループそれぞれについて、ブロードキャスト通信 3 3 の後に当該グループに含まれる 2 以上のノードの間で実行されるブロードキャスト通信 3 4 (第 2 のブロードキャスト通信) の通信順序を決定する。

【 0 0 2 8 】

ブロードキャスト通信 3 3 とブロードキャスト通信 3 4 は、それぞれ全体のブロードキ

50

キャスト通信の一部として実行され、異なる段階として区別されて実行される。ブロードキャスト通信 33 はグループ間通信であり、ブロードキャスト通信 34 はグループ内通信である。ブロードキャスト通信 33 では、グループ毎に 1 つの代表ノードのみが通信に参加する。複数のグループのブロードキャスト通信 34 は並列に実行してもよい。

【0029】

情報処理装置 10 は、ブロードキャスト通信 33 が以下の条件を満たすように、複数のグループに対応する複数の代表ノードの間の通信順序を決定する。複数の代表ノードのうち、1 つの代表ノードが最初の送信元ノードとして動作する。データ受信した代表ノードが以降は送信元ノードとして動作することで、並列にデータを送信する送信元ノードが増加する。例えば、フェーズが 1 つ進む毎に、並列にデータを送信する送信元ノードの数が 2 倍になる。ここで決定されるブロードキャスト通信 33 の通信順序は、Binomial Tree アルゴリズム (二項木アルゴリズム) に基づいていてもよい。

10

【0030】

また、情報処理装置 10 は、ブロードキャスト通信 34 が以下の条件を満たすように、グループ内の 2 以上のノードの間の通信順序を決定する。グループ内の 2 以上のノードのうち、代表ノードが最初の送信元ノードとして動作する。代表ノードには、上記のブロードキャスト通信 33 によってデータがコピーされている。データ受信したノードが以降は送信元ノードとして動作することで、並列にデータを送信する送信元ノードが増加する。例えば、フェーズが 1 つ進む毎に、並列にデータを送信する送信元ノードの数が 2 倍になる。ここで決定されるブロードキャスト通信 34 の通信順序は、ブロードキャスト通信 33 と同様に、Binomial Tree アルゴリズムに基づくものであってもよい。

20

【0031】

情報処理装置 10 は、ブロードキャスト通信 33 の通信順序とブロードキャスト通信 34 の通信順序とを結合して、全体のブロードキャスト通信の通信手順を決定する。情報処理装置 10 は、決定した通信手順を示す通信制御情報を生成して記憶する。情報処理装置 10 がノードの 1 つである場合、情報処理装置 10 は、生成した通信制御情報を参照してブロードキャスト通信を実行してもよい。情報処理装置 10 が制御装置である場合、情報処理装置 10 は、生成した通信制御情報を複数ノードに配布してもよい。

【0032】

第 1 の実施の形態の情報処理システムによれば、接続された上位の中継装置が同一でないような異なる下位の中継装置に接続された異なるノードが、異なるグループに分類されるように、ノードがグループ分けされ、グループ毎に代表ノードが選択される。そして、全体のブロードキャスト通信が、代表ノード間のブロードキャスト通信とグループ内のブロードキャスト通信に分けて実行される。代表ノード間のブロードキャスト通信およびグループ内のブロードキャスト通信それぞれでは、並列にデータを送信する送信元ノードが段階的に増加するアルゴリズムに従って通信順序が決定される。

30

【0033】

代表ノード間のブロードキャスト通信およびグループ内のブロードキャスト通信それぞれでは、異なるノードペアの間のデータ送信が並列に実行されるため、ブロードキャスト通信が高速化される。また、上記のグループ内通信では、下位の中継装置と上位の中継装置との間の通信経路が冗長化されているため、通信競合が抑制される。また、グループ間通信に参加するノードはグループ毎に 1 つであるため、グループ間の通信経路の冗長度が低い場合であっても通信競合が抑制される。従って、ブロードキャスト通信の全体を通じて通信競合を抑制でき、通信遅延を抑制して通信時間を短縮できる。

40

【0034】

[第 2 の実施の形態]

次に、第 2 の実施の形態を説明する。

図 2 は、第 2 の実施の形態の多層フルメッシュシステムの例を示す図である。

【0035】

第 2 の実施の形態の多層フルメッシュシステムは、複数のサーバおよび複数のスイッチ

50

を含み、それら複数のサーバおよび複数のスイッチが多層フルメッシュトポロジで接続された並列処理システムである。サーバは、ユーザプログラムを実行可能なノードであり、コンピュータや情報処理装置とすることもできる。

【0036】

スイッチは、サーバ間で送信されるデータを中継する通信装置である。後述するように、スイッチはLeafスイッチとSpineスイッチとに分類される。LeafスイッチとSpineスイッチは、同様のハードウェアをもつスイッチであってもよい。第2の実施の形態では、説明を簡単にするため、スイッチのポート数が6であるとする。ただし、スイッチのポート数は、8や10や36など6より大きい偶数であってもよい。

【0037】

多層フルメッシュシステムは、複数の層を形成する。各サーバは、何れか1つのLeafスイッチに接続される。各Leafスイッチは、何れか1つの層に属する。各Spineスイッチは、複数の層を貫通しており、複数の層のLeafスイッチに接続される。

【0038】

層内では、複数のLeafスイッチがフルメッシュトポロジを形成する。よって、Leafスイッチのペア毎に、他のLeafスイッチを経由しない最短経路が存在する。2つのLeafスイッチの間には、複数の層を貫通するSpineスイッチが配置される。よって、同じ層に属する2つのLeafスイッチは、1つのSpineスイッチを経由する通信経路によって通信することができる。異なる層に属する2つのLeafスイッチも、1つのSpineスイッチを経由する通信経路によって通信することができる。Leafスイッチは、データをその宛先に応じて最短経路で転送するよう設定される。

【0039】

ポート数が6である第2の実施の形態では、多層フルメッシュシステムは3つの層を形成する。各層は4つのLeafスイッチを含む。各Leafスイッチには、3つのサーバと3つのSpineスイッチが接続される。各Spineスイッチには、層毎に2つのLeafスイッチが接続され、3つの層の合計で6つのLeafスイッチが接続される。多層フルメッシュシステムは6つのSpineスイッチを含む。

【0040】

一般に、ポート数が p (p は6以上の偶数)であるスイッチを使用すると、多層フルメッシュシステムは $p/2$ 個の層を形成する。各層は $p/2 + 1$ 個のLeafスイッチによって $p/2 + 1$ 角形を形成する。多層フルメッシュシステムは、 $p^2(p+2)/8$ 個のサーバと $3p(p+2)/8$ 個のスイッチを含む。 $p=8$ の場合、多層フルメッシュシステムは、5角形の4層を形成し、80個のサーバと30個のスイッチを含む。 $p=10$ の場合、多層フルメッシュシステムは、6角形の5層を形成し、150個のサーバと45個のスイッチを含む。 $p=36$ の場合、多層フルメッシュシステムは、19角形の18層を形成し、6156個のサーバと513個のスイッチを含む。

【0041】

第2の実施の形態の多層フルメッシュシステムは、層41, 42, 43を形成する。層41は、Leafスイッチ200, 210, 220, 230を含む。Leafスイッチ200, 210, 220, 230にはそれぞれ3つのサーバが接続される。

【0042】

Leafスイッチ200とLeafスイッチ210の間にSpineスイッチ240が配置される。Leafスイッチ200とLeafスイッチ220の間にSpineスイッチ241が配置される。Leafスイッチ200とLeafスイッチ230の間にSpineスイッチ242が配置される。Leafスイッチ210とLeafスイッチ220の間にSpineスイッチ243が配置される。Leafスイッチ210とLeafスイッチ230の間にSpineスイッチ244が配置される。Leafスイッチ220とLeafスイッチ230の間にSpineスイッチ245が配置される。

【0043】

層42, 43も、Leafスイッチ200, 210, 220, 230に対応するLeaf

10

20

30

40

50

fスイッチを含む。Spineスイッチ240, 241, 242, 243, 244, 245は、層41, 42, 43を貫通しており層41, 42, 43の間で共通である。

【0044】

例えば、層42は、Leafスイッチ200, 220, 230に対応するLeafスイッチ201, 221, 231を含む。Leafスイッチ201とLeafスイッチ231の間にSpineスイッチ242が配置される。Leafスイッチ221とLeafスイッチ231の間にSpineスイッチ245が配置される。層43は、Leafスイッチ202, 222, 232を含む。Leafスイッチ202とLeafスイッチ232の間にSpineスイッチ242が配置される。Leafスイッチ222とLeafスイッチ232の間にSpineスイッチ245が配置される。

10

【0045】

また、第2の実施の形態の多層フルメッシュシステムは、ジョブスケジューラ300を含む。ジョブスケジューラ300は、ユーザからジョブ要求を受け付け、ジョブに使用するサーバ(ノード)を選択するサーバ装置である。ジョブスケジューラ300は、コンピュータや情報処理装置ということもできる。ジョブは、ユーザプログラムから起動される複数のプロセスを含む。ユーザプログラムは、MPIライブラリなどの通信ライブラリを用いることがある。複数のプロセスには、ランクと呼ばれる非負整数の識別番号が付与される。1つのサーバには1つのプロセスが配置される。ジョブスケジューラ300は、プロセスの配置を決定し、サーバに対してプロセス配置に関する情報を通知する。

【0046】

ジョブスケジューラ300とサーバとの間の通信には、上記のLeafスイッチやSpineスイッチを含むデータ用ネットワークを使用してもよいし、データ用ネットワークとは異なる管理用ネットワークを使用してもよい。

20

【0047】

図3は、多層フルメッシュシステムの配線例を示す図である。

図3は、図2の多層フルメッシュシステムに含まれるサーバとLeafスイッチとSpineスイッチの間の配線を、図2とは異なる形式で表現したものである。

【0048】

多層フルメッシュシステムは、Spineスイッチ240, 241, 242, 243, 244, 245(SpineスイッチA, B, C, D, E, F)を含む。

30

また、多層フルメッシュシステムは、Leafスイッチ200, 201, 202(Leafスイッチa1, a2, a3)を含む。Leafスイッチ200, 201, 202はそれぞれ、Spineスイッチ240, 241, 242の3つのSpineスイッチに接続されている。Leafスイッチ200には、サーバ100, 101, 102が接続されている。Leafスイッチ201には、サーバ103, 104, 105が接続されている。Leafスイッチ202には、サーバ106, 107, 108が接続されている。

【0049】

また、多層フルメッシュシステムは、Leafスイッチ210, 211, 212(Leafスイッチb1, b2, b3)を含む。Leafスイッチ210, 211, 212はそれぞれ、Spineスイッチ240, 243, 244の3つのSpineスイッチに接続されている。Leafスイッチ210には、サーバ110, 111, 112が接続されている。Leafスイッチ211には、サーバ113, 114, 115が接続されている。Leafスイッチ212には、サーバ116, 117, 118が接続されている。

40

【0050】

また、多層フルメッシュシステムは、Leafスイッチ220, 221, 222(Leafスイッチc1, c2, c3)を含む。Leafスイッチ220, 221, 222はそれぞれ、Spineスイッチ241, 243, 245の3つのSpineスイッチに接続されている。Leafスイッチ220には、サーバ120, 121, 122が接続されている。Leafスイッチ221には、サーバ123, 124, 125が接続されている。Leafスイッチ222には、サーバ126, 127, 128が接続されている。

50

【0051】

また、多層フルメッシュシステムは、Leafスイッチ230, 231, 232 (Leafスイッチd1, d2, d3)を含む。Leafスイッチ230, 231, 232はそれぞれ、Spineスイッチ242, 244, 245の3つのSpineスイッチに接続されている。Leafスイッチ230には、サーバ130, 131, 132が接続されている。Leafスイッチ231には、サーバ133, 134, 135が接続されている。Leafスイッチ232には、サーバ136, 137, 138が接続されている。

【0052】

このように、各Leafスイッチには、上位スイッチとして3つのSpineスイッチが接続されている。層41, 42, 43の間の対応する位置にあるLeafスイッチは、同一のSpineスイッチに接続されている。第2の実施の形態では、接続されている3つのSpineスイッチが全て同一であるLeafスイッチおよびその配下のサーバを、「層間グループ」または単に「グループ」と言うことがある。

10

【0053】

Leafスイッチ200, 201, 202およびその配下のサーバ100, 101, 102, 103, 104, 105, 106, 107, 108は、1つのグループ(グループa)を形成する。Leafスイッチ210, 211, 212およびその配下のサーバ110, 111, 112, 113, 114, 115, 116, 117, 118は、1つのグループ(グループb)を形成する。Leafスイッチ220, 221, 222およびその配下のサーバ120, 121, 122, 123, 124, 125, 126, 127, 128は、1つのグループ(グループc)を形成する。Leafスイッチ230, 231, 232およびその配下のサーバ130, 131, 132, 133, 134, 135, 136, 137, 138は、1つのグループ(グループd)を形成する。

20

【0054】

図4は、サーバのハードウェア例を示すブロック図である。

サーバ100は、CPU151、RAM152、HDD153、画像インタフェース154、入力インタフェース155、媒体リーダー156およびHCA (Host Channel Adapter) 157を有する。上記ユニットはバスに接続されている。他のサーバやジョブスケジューラ300も、サーバ100と同様のハードウェアを有する。

【0055】

CPU151は、プログラムの命令を実行するプロセッサである。CPU151は、HDD153に記憶されたプログラムやデータの少なくとも一部をRAM152にロードし、プログラムを実行する。なお、CPU151は複数のプロセッサコアを備えてもよく、サーバ100は複数のプロセッサを備えてもよい。複数のプロセッサの集合を「マルチプロセッサ」または単に「プロセッサ」と言うことがある。

30

【0056】

RAM152は、CPU151が実行するプログラムやCPU151が演算に使用するデータを一時的に記憶する揮発性の半導体メモリである。なお、サーバ100は、RAM以外の種類のメモリを備えてもよく、複数のメモリを備えてもよい。

【0057】

HDD153は、OS (Operating System) やミドルウェアやアプリケーションソフトウェアなどのソフトウェアのプログラム、および、データを記憶する不揮発性ストレージである。なお、サーバ100は、フラッシュメモリやSSD (Solid State Drive) など他の種類のストレージを備えてもよく、複数のストレージを備えてもよい。

40

【0058】

画像インタフェース154は、CPU151からの命令に従って、サーバ100に接続された表示装置161に画像を出力する。表示装置161として、CRT (Cathode Ray Tube) ディスプレイ、液晶ディスプレイ (LCD: Liquid Crystal Display)、有機EL (OEL: Organic Electro-Luminescence) ディスプレイ、プロジェクタなど、任意の種類の表示装置を使用することができる。また、サーバ100に、プリンタなど表示

50

装置 161 以外の出力デバイスが接続されてもよい。

【0059】

入力インタフェース 155 は、サーバ 100 に接続された入力デバイス 162 から入力信号を受け付ける。入力デバイス 162 として、マウス、タッチパネル、タッチパッド、キーボードなど、任意の種類の入力デバイスを使用することができる。また、サーバ 100 に複数種類の入力デバイスが接続されてもよい。

【0060】

媒体リーダ 156 は、記録媒体 163 に記録されたプログラムやデータを読み取る読み取り装置である。記録媒体 163 として、フレキシブルディスク (FD: Flexible Disk) や HDD などの磁気ディスク、CD (Compact Disc) や DVD (Digital Versatile Disc) などの光ディスク、半導体メモリなど、任意の種類記録媒体を使用することができる。媒体リーダ 156 は、例えば、記録媒体 163 から読み取ったプログラムやデータを、RAM 152 や HDD 153 などの他の記録媒体にコピーする。読み取られたプログラムは、例えば、CPU 151 によって実行される。なお、記録媒体 163 は可搬型記録媒体であってもよく、プログラムやデータの配布に用いられることがある。また、記録媒体 163 や HDD 153 を、コンピュータ読み取り可能な記録媒体と言うことがある。

10

【0061】

HCA 157 は、InfiniBand の通信インタフェースである。HCA 157 は、全二重通信が可能であり、データの送信と受信を並列に行える。HCA 157 は、Leaf スイッチ 200 に接続される。ただし、サーバ 100 は、HCA 157 に代えてまたは HCA 157 に加えて、他の通信規格の通信インタフェースを有してもよい。

20

【0062】

図 5 は、スイッチのハードウェア例を示すブロック図である。

Leaf スイッチ 200 は、CPU 251、RAM 252、ROM 253 および通信ポート 254、255、256、257、258、259 を有する。他の Leaf スイッチや Spine スイッチも、Leaf スイッチ 200 と同様のハードウェアを有する。

【0063】

CPU 251 は、通信制御プログラムを実行するプロセッサである。CPU 251 は、通信制御プログラムに従い、受信されたパケットをその宛先に応じた通信ポートに出力する。CPU 251 は、ROM 253 に記憶された通信制御プログラムの少なくとも一部を RAM 252 にロードし、通信制御プログラムを実行する。ただし、通信制御の少なくとも一部を、専用のハードウェア回路を用いて実装することもできる。

30

【0064】

RAM 252 は、CPU 251 が実行する通信制御プログラムや通信制御に使用するデータを一時的に記憶する揮発性の半導体メモリである。データには、パケットの宛先と出力先の通信ポートとを対応付けたルーティング情報が含まれる。ROM 253 は、通信制御プログラムを記憶する不揮発性ストレージである。ただし、Leaf スイッチ 200 は、フラッシュメモリなど書き換え可能な不揮発性ストレージを備えてもよい。

【0065】

通信ポート 254、255、256、257、258、259 は、InfiniBand の通信インタフェースである。通信ポート 254、255、256、257、258、259 は、全二重通信が可能であり、データの送信と受信を並列に行える。通信ポート 254 は、サーバ 100 に接続される。通信ポート 255 は、サーバ 101 に接続される。通信ポート 256 は、サーバ 102 に接続される。通信ポート 257 は、Spine スイッチ 241 に接続される。通信ポート 258 は、Spine スイッチ 242 に接続される。通信ポート 259 は、Spine スイッチ 243 に接続される。ただし、Leaf スイッチ 200 は、通信ポート 254、255、256、257、258、259 に代えてまたは通信ポート 254、255、256、257、258、259 に加えて、他の通信規格の通信インタフェースを有してもよい。

40

【0066】

50

次に、多層フルメッシュシステム上のブロードキャスト通信について説明する。

同一のジョブに属する複数のプロセスは、それら複数のプロセスが一斉にデータ送信に参加するコレクティブ通信を行うことがある。ユーザプログラムがMPIライブラリのコレクティブ通信の命令を呼び出すことで、一斉のデータ送信を開始できる。コレクティブ通信の1つの種類として、ブロードキャスト通信がある。ブロードキャスト通信では、ランク0のプロセスなど特定のプロセスがもつデータを、他の全てのプロセスにコピーする。1ノード1プロセスを仮定すると、ブロードキャスト通信は、あるサーバ(ノード)から他の全てのサーバ(ノード)にデータをコピーするものであると言える。ブロードキャスト通信のアルゴリズムの1つに、BinomialTreeアルゴリズムがある。

【0067】

図6は、BinomialTreeアルゴリズムの例を示すシーケンス図である。

ここでは、サーバ110がランク0のプロセスを実行し、サーバ111がランク1のプロセスを実行し、サーバ120がランク2のプロセスを実行し、サーバ121がランク3のプロセスを実行するものとする。また、サーバ130がランク4のプロセスを実行し、サーバ131がランク5のプロセスを実行し、サーバ100がランク6のプロセスを実行し、サーバ101がランク7のプロセスを実行するものとする。また、ランク0のプロセスがランク1, 2, 3, 4, 5, 6, 7のプロセスに同一データを渡すことを考える。

【0068】

BinomialTreeアルゴリズムでは、あるフェーズでデータを受信したサーバが、次以降のフェーズではデータの送信元として動作する。これにより、同一データを保持しているサーバが2の累乗の速度で増加する。

【0069】

フェーズt1では、サーバ110がサーバ111にデータを送信する(S10)。これにより、サーバ110, 111が同一データを保持する。フェーズt2では、サーバ110がサーバ120にデータを送信し(S11)、これと並列にサーバ111がサーバ121にデータを送信する(S12)。これにより、サーバ110, 111, 120, 121が同一データを保持する。フェーズt3では、サーバ110がサーバ130にデータを送信し(S13)、これと並列にサーバ111がサーバ131にデータを送信する(S14)。更に、これと並列にサーバ120がサーバ100にデータを送信し(S15)、これと並列にサーバ121がサーバ101にデータを送信する(S16)。

【0070】

BinomialTreeアルゴリズムのフェーズ数は、ブロードキャスト通信に参加するプロセスの数をNとすると、 $O(\log_2 N)$ である。具体的には、フェーズ数nは、 $2^{n-1} < N \leq 2^n$ を満たす自然数である。例えば、プロセス数N=8の場合は通信フェーズ数n=3である。プロセス数N=36の場合はフェーズ数n=6である。プロセス数N=80の場合はフェーズ数n=7である。

【0071】

ただし、単純なBinomialTreeアルゴリズムでは、フェーズの進行に伴ってデータ通信の並列度が増加する。よって、プロセスの配置状況によっては、データ通信が競合するリスクが高くなる。同じフェーズにおいて2つのデータ通信が同じリンクを同じ方向に使用する場合、データ通信が競合している(衝突している)と言える。データ通信が競合すると、パケットの送信待ちが発生する、1つのリンクの通信帯域が分割されるなどにより、通信遅延が発生して通信時間が長くなるおそれがある。

【0072】

図7は、ブロードキャスト通信の競合例を示す図である。

前述のフェーズt2では、サーバ110がサーバ120にデータを送信し、サーバ111がサーバ121にデータを送信する。サーバ110からサーバ120への最短経路と、サーバ111からサーバ121への最短経路は共に、Leafスイッチ210、Spineスイッチ243、Leafスイッチ220を順に経由するものであり、競合している。

【0073】

10

20

30

40

50

また、フェーズ t 3 では、サーバ 1 1 0 がサーバ 1 3 0 にデータを送信し、サーバ 1 1 1 がサーバ 1 3 1 にデータを送信する。サーバ 1 1 0 からサーバ 1 3 0 への最短経路と、サーバ 1 1 1 からサーバ 1 3 1 への最短経路は共に、Leaf スイッチ 2 1 0、Spine スイッチ 2 4 4、Leaf スイッチ 2 3 0 を順に経由するものであり、競合している。
【0074】

また、フェーズ t 3 では、サーバ 1 2 0 がサーバ 1 0 0 にデータを送信し、サーバ 1 2 1 がサーバ 1 0 1 にデータを送信する。サーバ 1 2 0 からサーバ 1 0 0 への最短経路と、サーバ 1 2 1 からサーバ 1 0 1 への最短経路は共に、Leaf スイッチ 2 2 0、Spine スイッチ 2 4 1、Leaf スイッチ 2 0 0 を順に経由するものであり、競合している。
【0075】

そこで、第 2 の実施の形態の多層フルメッシュシステムは、通信競合が生じないようにブロードキャスト通信の手順を規定する。具体的には、各グループから 1 つの代表プロセスを選択し、選択した代表プロセスの間で Binomial Tree アルゴリズムを実行する。そして、各グループの代表プロセスにデータがコピーされた後、グループの内部で代表プロセスを始点とする Binomial Tree アルゴリズムを実行する。以下、第 2 の実施の形態のブロードキャスト通信の手順を説明する。

【0076】

図 8 は、プロセスの配置例を示す図である。

第 2 の実施の形態のブロードキャスト通信の手順を説明するにあたり、36 個のサーバのうち 32 個のサーバに、32 個のプロセスを配置することを考える。4 つのグループそれぞれに、32 個のプロセスのうち 8 個のプロセスが配置される。

【0077】

グループ a のサーバ 1 0 0, 1 0 1, 1 0 2, 1 0 3, 1 0 4, 1 0 5, 1 0 6, 1 0 7 に、ランク 0, 4, 8, 12, 16, 20, 24, 28 のプロセスが配置される。グループ b のサーバ 1 1 0, 1 1 1, 1 1 2, 1 1 3, 1 1 4, 1 1 5, 1 1 6, 1 1 7 に、ランク 1, 5, 9, 13, 17, 21, 25, 29 のプロセスが配置される。グループ c のサーバ 1 2 0, 1 2 1, 1 2 2, 1 2 3, 1 2 4, 1 2 5, 1 2 6, 1 2 7 に、ランク 2, 6, 10, 14, 18, 22, 26, 30 のプロセスが配置される。グループ d のサーバ 1 3 0, 1 3 1, 1 3 2, 1 3 3, 1 3 4, 1 3 5, 1 3 6, 1 3 7 に、ランク 3, 7, 11, 15, 19, 23, 27, 31 のプロセスが配置される。

【0078】

図 9 は、二段階ブロードキャスト通信の例を示すシーケンス図である。

まず、グループ a, b, c, d それぞれから代表プロセスが選択される。代表プロセスは、例えば、グループ内で最もランクが小さいプロセスである。ここでは、グループ a からランク 0 のプロセスが選択され、グループ b からランク 1 のプロセスが選択され、グループ c からランク 2 のプロセスが選択され、グループ d からランク 3 のプロセスが選択される。ただし、代表プロセスを他の基準で選択してもよい。また、ここでは 4 つのプロセスが同一の層に配置されているが、異なる層に配置されたプロセスが混在してもよい。

【0079】

代表プロセスが選択されると、まず代表プロセスの間で Binomial Tree アルゴリズムが実行される。ランク 0 のプロセスのデータをブロードキャストする場合、フェーズ t 1 では、サーバ 1 0 0 がサーバ 1 1 0 にデータを送信する (S 2 0)。フェーズ t 2 では、サーバ 1 0 0 がサーバ 1 2 0 にデータを送信し (S 2 1)、これと並列にサーバ 1 1 0 がサーバ 1 3 0 にデータを送信する (S 2 2)。これにより、代表プロセスが配置されたサーバ 1 0 0, 1 1 0, 1 2 0, 1 3 0 が同一データを保持する。

【0080】

代表プロセスの間でのデータコピーが完了すると、グループ a, b, c, d それぞれの内部で Binomial Tree アルゴリズムが実行される。ここでは、グループ a の通信手順を説明する。グループ a と並列にグループ b, c, d でも同様の通信が行われる。

【0081】

10

20

30

40

50

フェーズ t 3 では、サーバ 1 0 0 がサーバ 1 0 1 にデータを送信する (S 2 3)。フェーズ t 4 では、サーバ 1 0 0 がサーバ 1 0 2 にデータを送信し (S 2 4)、これと並列にサーバ 1 0 1 がサーバ 1 0 3 にデータを送信する (S 2 5)。フェーズ t 5 では、サーバ 1 0 0 がサーバ 1 0 4 にデータを送信し (S 2 6)、これと並列にサーバ 1 0 1 がサーバ 1 0 5 にデータを送信する (S 2 7)。更に、これと並列にサーバ 1 0 2 がサーバ 1 0 6 にデータを送信し (S 2 8)、これと並列にサーバ 1 0 3 がサーバ 1 0 7 にデータを送信する (S 2 9)。これにより、グループ a に属するサーバ 1 0 0 , 1 0 1 , 1 0 2 , 1 0 3 , 1 0 4 , 1 0 5 , 1 0 6 , 1 0 7 が同一データを保持する。

【 0 0 8 2 】

次に、グループ間通信とグループ内通信の通信経路を説明する。

10

図 1 0 は、グループ間通信の競合回避例を示す図である。

上記のフェーズ t 1 のステップ S 2 0 では、サーバ 1 0 0 から、Leaf スイッチ 2 0 0 と Spine スイッチ 2 4 0 と Leaf スイッチ 2 1 0 を経由して、サーバ 1 1 0 にデータが送信される。上記のフェーズ t 2 のステップ S 2 1 では、サーバ 1 0 0 から、Leaf スイッチ 2 0 0 と Spine スイッチ 2 4 1 と Leaf スイッチ 2 2 0 を経由して、サーバ 1 2 0 にデータが送信される。また、上記のフェーズ t 2 のステップ S 2 2 では、サーバ 1 1 0 から、Leaf スイッチ 2 1 0 と Spine スイッチ 2 4 4 と Leaf スイッチ 2 3 0 を経由して、サーバ 1 3 0 にデータが送信される。

【 0 0 8 3 】

このように、各グループから代表プロセスを選択し、代表プロセスが配置されたサーバのみが通信を行うようにすると、グループ間通信において競合が生じない。これは、異なるグループに属する Leaf スイッチの間にはフルメッシュの通信経路が存在するためである。フルメッシュの通信経路が存在することから、グループ a とグループ c の間の通信経路は、グループ b とグループ d の間の通信経路とリンクを共有しない。代表プロセスの配置されたサーバが異なる層に跨がっていても同様である。

20

【 0 0 8 4 】

図 1 1 は、グループ内通信の競合回避例を示す図である。

上記のフェーズ t 3 のステップ S 2 3 では、サーバ 1 0 0 から、Leaf スイッチ 2 0 0 を経由してサーバ 1 0 1 にデータが送信される。

【 0 0 8 5 】

30

上記のフェーズ t 4 のステップ S 2 4 では、サーバ 1 0 0 から、Leaf スイッチ 2 0 0 を経由してサーバ 1 0 2 にデータが送信される。上記のフェーズ t 4 のステップ S 2 5 では、サーバ 1 0 1 から、Leaf スイッチ 2 0 0 と Spine スイッチ 2 4 0 と Leaf スイッチ 2 0 1 を経由して、サーバ 1 0 3 にデータが送信される。

【 0 0 8 6 】

上記のフェーズ t 5 のステップ S 2 6 では、サーバ 1 0 0 から、Leaf スイッチ 2 0 0 と Spine スイッチ 2 4 1 と Leaf スイッチ 2 0 1 を経由して、サーバ 1 0 4 にデータが送信される。上記のフェーズ t 5 のステップ S 2 7 では、サーバ 1 0 1 から、Leaf スイッチ 2 0 0 と Spine スイッチ 2 4 2 と Leaf スイッチ 2 0 1 を経由して、サーバ 1 0 5 にデータが送信される。上記のフェーズ t 5 のステップ S 2 8 では、サーバ 1 0 2 から、Leaf スイッチ 2 0 0 と Spine スイッチ 2 4 0 と Leaf スイッチ 2 0 2 を経由して、サーバ 1 0 6 にデータが送信される。上記のフェーズ t 5 のステップ S 2 9 では、サーバ 1 0 3 から、Leaf スイッチ 2 0 1 と Spine スイッチ 2 4 1 と Leaf スイッチ 2 0 2 を経由して、サーバ 1 0 7 にデータが送信される。

40

【 0 0 8 7 】

なお、図 1 1 はグループ a の内部のデータ通信を表しているが、グループ b , c , d の内部のデータ通信も同様である。ただし、グループ b は、上位スイッチとして Spine スイッチ 2 4 0 , 2 4 3 , 2 4 4 を使用する。グループ c は、上位スイッチとして Spine スイッチ 2 4 1 , 2 4 3 , 2 4 5 を使用する。グループ d は、上位スイッチとして Spine スイッチ 2 4 2 , 2 4 4 , 2 4 5 を使用する。

50

【 0 0 8 8 】

このように、グループ内通信では競合が生じない。これは、グループ内ネットワークが Fat Tree トポロジに相当するためである。Fat Tree トポロジは、Tree トポロジに含まれる上位の通信装置を多重化することで、異なる下位の通信装置の間の通信経路を多重化し、トラフィックの混雑を軽減するネットワークトポロジである。

【 0 0 8 9 】

第2の実施の形態の多層フルメッシュシステムでは、3つの Leaf スイッチそれぞれがもつ Spine スイッチ側のリンクは3本であり、サーバ側のリンクと同数である。また、1つの Leaf スイッチから別の1つの Leaf スイッチに到達する通信経路は3つある。3つの Leaf スイッチと3つの Spine スイッチの間には合計で9つの通信経路が存在することになり、3つの Leaf スイッチに接続されたサーバと同数である。よって、1つのサーバに対して Leaf スイッチと Spine スイッチの間の1つの通信経路を割り当てれば、9つのサーバは競合なしにデータ通信を行うことができる。

10

【 0 0 9 0 】

次に、サーバとジョブスケジューラの機能について説明する。

図12は、サーバとジョブスケジューラの機能例を示す図である。

サーバ100は、通信手順決定部171、通信手順記憶部172およびブロードキャスト実行部173を有する。通信手順記憶部172は、例えば、RAM152またはHDD153の記憶領域を用いて実現される。通信手順決定部171およびブロードキャスト実行部173は、例えば、CPU151が実行するプログラムを用いて実現される。他のサーバもサーバ100と同様のモジュールを有する。

20

【 0 0 9 1 】

通信手順決定部171は、ジョブスケジューラ300から、ジョブに属する複数のプロセスの配置を示すプロセス配置情報を受信する。プロセス配置情報は、例えば、プロセスのランクとプロセスが配置されたサーバを識別するノードIDとを対応付ける。通信手順決定部171は、受信したプロセス配置情報に基づいて、ブロードキャスト通信における複数のプロセスの間の通信手順を決定する。通信手順決定部171は、決定した通信手順を示す通信手順情報を生成し、通信手順記憶部172に格納する。

【 0 0 9 2 】

ブロードキャスト通信の通信手順は、MPIライブラリなどの通信ライブラリの初期化時に決定される。通信ライブラリを使用するユーザプログラムが複数のサーバに配置され、それら複数のサーバでユーザプログラムが起動されると、通信ライブラリが初期化される。通信ライブラリの初期化時には、サーバ間で通信が行われることがある。通信手順の決定に使用するプロセス配置情報を、ジョブスケジューラ300から受信する代わりに、サーバ間の通信によって収集するようにしてもよい。また、ブロードキャスト通信の通信手順を、通信ライブラリの初期化時に行う代わりに、ブロードキャスト通信を初めてユーザプログラムから要求されたときに行うようにしてもよい。

30

【 0 0 9 3 】

ここで決定される通信手順は、ランク0のプロセスを始点とするブロードキャスト通信の通信手順である。ランク0のプロセス以外のプロセスがもつデータをブロードキャストする場合、ランク0のプロセスにデータを渡せばよい。ただし、ランク0のプロセス以外のプロセスを始点とするブロードキャスト通信の通信手順を決定することも可能である。また、通信手順決定部171が生成する通信手順情報は、ジョブに含まれる全てのプロセスの間の通信手順を示す。複数のサーバは、同一のプロセス配置情報と同一のブロードキャスト通信アルゴリズムを使用すれば、同一の通信手順情報を生成することになる。ただし、生成する通信手順情報を、サーバ100の通信手順のみを示すようにしてもよい。

40

【 0 0 9 4 】

通信手順記憶部172は、通信手順決定部171が生成した通信手順情報を記憶する。通信手順情報は、ブロードキャスト通信の各フェーズにおいて、データを送信する際の送信先プロセスのランクと、データを受信する際の送信元プロセスのランクとを示す。

50

【 0 0 9 5 】

ブロードキャスト実行部 1 7 3 は、ブロードキャスト通信を開始する命令をユーザプログラムが呼び出すと、通信手順記憶部 1 7 2 に記憶された通信手順情報に基づいてブロードキャスト通信を実行する。ブロードキャスト実行部 1 7 3 は、通信手順情報が示す複数のフェーズを 1 つずつ実行する。あるフェーズで送信元プロセスが指定されている場合、ブロードキャスト実行部 1 7 3 は、送信元プロセスからデータを受信する。あるフェーズで送信先プロセスが指定されている場合、ブロードキャスト実行部 1 7 3 は、保持しているデータをコピーして送信先プロセスに送信する。

【 0 0 9 6 】

データ送信では、ブロードキャスト実行部 1 7 3 は、送信先プロセスが配置されたサーバのアドレスとデータ本体とを含むパケットを生成し、H C A 1 5 7 を介して L e a f スイッチ 2 0 0 にパケットを出力する。各プロセスが配置されたサーバのアドレスは、通信ライブラリの初期化時に把握される。

10

【 0 0 9 7 】

ジョブスケジューラ 3 0 0 は、プロセス配置決定部 3 7 1 を有する。プロセス配置決定部 3 7 1 は、例えば、C P U が実行するプログラムを用いて実現される。

プロセス配置決定部 3 7 1 は、ユーザからジョブ要求を受け付け、受け付けたジョブ要求に応じてジョブに含まれる複数のプロセスの配置を決定する。起動するプロセスの数は、ユーザからのジョブ要求で指定される。プロセス配置決定部 3 7 1 は、例えば、同一のジョブに属する複数のプロセスが、できる限りグループ a , b , c , d に均等に配置されるようにプロセス配置を決定する。プロセス配置決定部 3 7 1 は、決定したプロセス配置を示すプロセス配置情報を、ジョブで使用する複数のサーバに送信する。

20

【 0 0 9 8 】

図 1 3 は、プロセス配置テーブルの例を示す図である。

プロセス配置テーブル 1 7 4 は、通信手順決定部 1 7 1 がジョブスケジューラ 3 0 0 から受信するプロセス配置情報を示す。プロセス配置テーブル 1 7 4 は、ランクとノード I D とを対応付ける。ランクは、ジョブに含まれる複数のプロセスを識別する非負整数の識別番号である。ノード I D は、プロセスが配置されたサーバを識別する識別子である。ノード I D が、パケットの宛先を示す通信アドレスを兼ねてもよい。

【 0 0 9 9 】

図 1 4 は、通信手順テーブルの例を示す図である。

送信手順テーブル 1 7 5 および受信手順テーブル 1 7 6 は、通信手順決定部 1 7 1 により生成されて通信手順記憶部 1 7 2 に格納される。

30

【 0 1 0 0 】

送信手順テーブル 1 7 5 は、フェーズとランクの組に対して、当該フェーズで当該ランクのプロセスがデータを送信する際の送信先プロセスのランクが登録される。送信先プロセスが存在しない場合、すなわち、当該フェーズで当該ランクのプロセスがデータを送信しない場合、「 - 1 」などランクに使用されない所定の数値が登録される。例えば、フェーズ t 1 でランク 0 のプロセスがランク 1 のプロセスにデータを送信する場合、フェーズ t 1 とランク 0 の組に対して「 1 」が登録される。

40

【 0 1 0 1 】

受信手順テーブル 1 7 6 は、フェーズとランクの組に対して、当該フェーズで当該ランクのプロセスがデータを受信する際の送信元プロセスのランクが登録される。送信元プロセスが存在しない場合、すなわち、当該フェーズで当該ランクのプロセスがデータを受信しない場合、「 - 1 」などランクに使用されない所定の数値が登録される。例えば、フェーズ t 1 でランク 1 のプロセスがランク 0 のプロセスからデータを受信する場合、フェーズ t 1 とランク 1 の組に対して「 0 」が登録される。

【 0 1 0 2 】

ブロードキャスト通信を行う場合、各サーバは、送信手順テーブル 1 7 5 および受信手順テーブル 1 7 6 から、当該サーバに配置されたプロセスのランクに対応する行を読み出

50

し、読み出した行に含まれる数値を左側から右側に向かって順に参照すればよい。

【0103】

次に、サーバ100の処理手順について説明する。

図15は、通信手順決定の処理例を示すフローチャートである。

(S30)通信手順決定部171は、ジョブに含まれる複数のプロセスそれぞれが配置されたサーバ(ノード)を特定し、特定したサーバが属するグループを判定する。各グループは、接続されているSpineスイッチが同一である複数のLeafスイッチおよびそれらLeafスイッチに接続された複数のサーバによって形成される。

【0104】

第2の実施の形態の多層フルメッシュシステムは、グループa, b, c, dを含む。図2の例では、グループaは、層41, 42, 43それぞれの四角形の左上に位置するLeafスイッチおよびサーバを含む。グループbは、層41, 42, 43それぞれの四角形の左下に位置するLeafスイッチおよびサーバを含む。グループcは、層41, 42, 43それぞれの四角形の右下に位置するLeafスイッチおよびサーバを含む。グループdは、層41, 42, 43それぞれの四角形の右上に位置するLeafスイッチおよびサーバを含む。なお、通信手順決定部171は、多層フルメッシュシステムのトポロジ、すなわち、多層フルメッシュシステムに含まれるグループの定義を予め知っている。

10

【0105】

(S31)通信手順決定部171は、ジョブに含まれる複数のプロセスが単一グループに閉じているか、すなわち、ブロードキャスト通信に参加する複数のサーバが全て同一のグループに属するか判断する。単一グループに閉じている場合はステップS35に進み、単一グループに閉じていない場合はステップS32に進む。

20

【0106】

(S32)通信手順決定部171は、グループ毎に当該グループに配置されたプロセスの中から、ランクが最小のプロセスを代表プロセスとして選択する。なお、代表プロセスはグループ毎に1つに決まればよく、他の基準によって代表プロセスを選択してもよい。ただし、ランク0のプロセスのように、始点となるプロセスが選択されるようにする。

【0107】

(S33)通信手順決定部171は、ステップS32で選択された複数のグループに対応する複数の代表プロセスを、ランクの昇順(小さい順)にソートする。なお、ここでは複数の代表プロセスが一定の順に並べばよく、ランクの昇順でなくてもよい。ただし、ランク0のプロセスのように、始点となるプロセスが先頭になるようにする。

30

【0108】

(S34)通信手順決定部171は、ソートした複数の代表プロセスの間でBinomial Treeを生成する。ランクの昇順に代表プロセスがソートされている場合、Binomial Treeアルゴリズムでは、ランクの小さい代表プロセスから優先的にデータを受信することになる。データを受信済みの代表プロセスは2の累乗で増加する。通信手順決定部171は、Binomial Treeに従い、全ての代表プロセスがデータのコピーを受信するまでの代表プロセス間の通信手順を決定する。

【0109】

(S35)通信手順決定部171は、ジョブに使用される1以上のグループそれぞれについて、当該グループに配置されたプロセスをランクの昇順(小さい順)にソートする。なお、ここではグループ内のプロセスが一定の順に並べばよく、ランクの昇順でなくてもよい。ただし、ステップS32で選択された代表プロセスが先頭になるようにする。

40

【0110】

(S36)通信手順決定部171は、1以上のグループそれぞれについて、当該グループ内のソート済みのプロセスの間でBinomial Treeを生成する。ランクの昇順にプロセスがソートされている場合、代表プロセスから開始して、ランクの小さいプロセスから優先的にデータを受信することになる。データを受信済みのプロセスは2の累乗で増加する。通信手順決定部171は、Binomial Treeに従い、グループ内の全

50

でのプロセスがデータのコピーを受信するまでのプロセス間の通信手順を決定する。

【0111】

(S37) 通信手順決定部171は、ステップS34で決定したグループ間の通信手順の後ろに、ステップS36で決定したグループ内の通信手順を結合し、始点のプロセスがもつデータを全てのプロセスにコピーするまでの全体の通信手順を決定する。通信手順決定部171は、決定した全体の通信手順に基づいて、送信手順テーブル175および受信手順テーブル176を生成し、通信手順記憶部172に格納する。

【0112】

図16は、ブロードキャスト通信の処理例を示すフローチャートである。

(S40) ブロードキャスト実行部173は、ユーザプログラムからブロードキャスト通信が指示されると、送信手順テーブル175と受信手順テーブル176を取得する。

10

【0113】

(S41) ブロードキャスト実行部173は、未実行のフェーズのうちフェーズ番号の小さい方から優先的に次のフェーズを選択する。最初はフェーズt1が選択される。

(S42) ブロードキャスト実行部173は、受信手順テーブル176から、ステップS41で選択したフェーズとサーバ100に配置されたプロセスのランクとの組に対応する数値を読み出す。ブロードキャスト実行部173は、読み出した数値が送信元ランクを表しているか、すなわち、該当する送信元ランクが受信手順テーブル176に登録されているか判断する。読み出した数値が「-1」である場合、送信元ランクが登録されていないことになる。送信元ランクが登録されている場合はステップS43に進み、送信元ランクが登録されていない場合はステップS44に進む。

20

【0114】

(S43) ブロードキャスト実行部173は、送信元ランクが示す相手プロセスからデータを受信できるように待機し、データを受信する。例えば、ブロードキャスト実行部173は、相手プロセスに対応する受信バッファを定期的に確認し、受信バッファにデータが到着している場合には到着したデータを取り出す。データ受信は、以下のステップS44, S45と並列に実行でき、ステップS46までに実行されればよい。ブロードキャスト実行部173は、受信したデータを保持しておく。

【0115】

(S44) ブロードキャスト実行部173は、送信手順テーブル175から、ステップS41で選択したフェーズとサーバ100に配置されたプロセスのランクとの組に対応する数値を読み出す。ブロードキャスト実行部173は、読み出した数値が送信先ランクを表しているか、すなわち、該当する送信先ランクが送信手順テーブル175に登録されているか判断する。読み出した数値が「-1」である場合、送信先ランクが登録されていないことになる。送信先ランクが登録されている場合はステップS45に進み、送信先ランクが登録されていない場合はステップS46に進む。

30

【0116】

(S45) ブロードキャスト実行部173は、保持しているデータのコピーを、送信先ランクが示す相手プロセスに送信する。送信データはパケットに分割され、各パケットには相手プロセスが配置されたサーバのアドレスが付加される。保持しているデータは、サーバ100に配置されたプロセスが始点プロセスである場合はオリジナルデータであり、始点プロセスでない場合は前フェーズまでに他のサーバから受信したデータである。

40

【0117】

(S46) ブロードキャスト実行部173は、送信手順テーブル175や受信手順テーブル176に規定された全フェーズを実行したか判断する。全フェーズを実行した場合はブロードキャスト通信を終了し、未実行のフェーズがある場合はステップS41に戻る。

【0118】

ここで、第2の実施の形態のブロードキャスト通信のフェーズ数について説明する。単純なBinomial Treeアルゴリズムのフェーズ数は、36プロセスの場合は6フェーズであり、80プロセスの場合は7フェーズである。一方、第2の実施の形態では、

50

36プロセスの場合、グループ間通信の最小フェーズ数が2フェーズであり、グループ内通信の最小フェーズ数が4フェーズであるため、合計で6フェーズである。80プロセスの場合、グループ間通信の最小フェーズ数が3フェーズであり、グループ内通信の最小フェーズ数が4フェーズであるため、合計で7フェーズである。

【0119】

グループ間のプロセス数の偏りが小さければ、第2の実施の形態のブロードキャスト通信は、単純なBinomial Treeアルゴリズムのフェーズ数と同一かまたはそれに近いフェーズ数で実行できる。ブロードキャスト通信の効率の観点から、ジョブスケジューラ300は、複数のグループにできる限り均等にプロセスを配置することが好ましい。

【0120】

なお、第2の実施の形態では、接続されているSpineスイッチが同一である複数のLeafスイッチおよびそれら複数のLeafスイッチの配下のサーバから、1つのグループを形成した。これに対して、1つのLeafスイッチおよび当該Leafスイッチの配下のサーバから、1つのグループを形成することもできる。

【0121】

この場合、Leafスイッチ毎に1つの代表プロセスが選択される。グループ間通信については、複数のLeafスイッチに対応する複数の代表プロセスの間で1つのBinomial Treeが形成される。グループ内通信については、各Leafスイッチの配下にある複数のプロセスの間で1つのBinomial Treeが形成される。例えば、図2, 3の多層フルメッシュシステムでは12個のグループが形成される。このようにサーバをグループ化しても、複数のLeafスイッチの間のグループ間通信では通信競合が生じない。また、Leafスイッチ配下のグループ内通信でも通信競合が生じない。

【0122】

また、第2の実施の形態では二段階でブロードキャスト通信を行ったが、グループを階層化して三段階でブロードキャスト通信を行ってもよい。接続されているSpineスイッチが同一である複数のLeafスイッチおよびそれら複数のLeafスイッチの配下のサーバから、1つの大グループを形成する。また、1つのLeafスイッチおよび当該Leafスイッチの配下のサーバから、1つの小グループを形成する。

【0123】

この場合、大グループ毎に上位代表プロセスが選択され、更にLeafスイッチ毎に下位代表プロセスが選択される。第1段階として、複数の大グループに対応する複数の上位代表プロセスの間で1つのBinomial Treeが形成される。第2段階として、複数のLeafスイッチに対応する複数の下位代表プロセスの間で1つのBinomial Treeが形成される。第3段階として、各Leafスイッチの配下にある複数のプロセスの間で1つのBinomial Treeが形成される。例えば、図2, 3の多層フルメッシュシステムでは、4個の大グループと12個の小グループが形成される。サーバが多い場合、このようにLeafスイッチ単位でグループを形成することも有用である。

【0124】

第2の実施の形態の多層フルメッシュシステムによれば、多層フルメッシュトポロジが採用される。多層フルメッシュトポロジでは、単純なTreeトポロジと比べて上位の通信装置が冗長化され、下位の通信装置の間の通信経路が冗長化される。よって、トラフィックの混雑を抑制することができる。また、単純なFat Treeトポロジと比べて通信装置の個数を削減でき、システム構築コストを削減できる。また、第2の実施の形態の多層フルメッシュシステムによれば、Binomial Treeに従い、データをコピー済みのノードがフェーズ数に対して2の累乗で増加するようにブロードキャスト通信が実行される。よって、ブロードキャスト通信を高速に実行できる。

【0125】

また、接続されているSpineスイッチの集合が同一であるLeafスイッチおよびその配下のノードがグループ化され、各グループから代表ノードが選択される。そして、代表ノード間のデータ送信が優先的に実行され、代表ノードを始点とするグループ内のデ

10

20

30

40

50

ータ送信がその後実行される。ここで、複数のグループの間にはフルメッシュの通信経路が存在するため、通信に参加するノードがグループ毎に1つであれば、複数のノードが並列通信を行っても通信競合は生じない。また、グループ内のネットワークポロジはF a t T r e eに相当するため、グループ内のノード同士の閉じた通信であれば、複数のノードが並列通信を行っても通信競合は生じない。よって、通信競合を抑制でき、通信遅延を抑制してブロードキャスト通信の所要時間を短縮できる。

【符号の説明】

【 0 1 2 6 】

1 0 情報処理装置

1 1 , 1 2 , 1 3 , 1 4 , 1 5 , 1 6 , 1 7 , 1 8 ノード

2 1 , 2 2 , 2 3 , 2 4 , 2 5 , 2 6 , 2 7 , 2 8 中継装置

3 1 , 3 2 グループ

3 3 , 3 4 ブロードキャスト通信

10

20

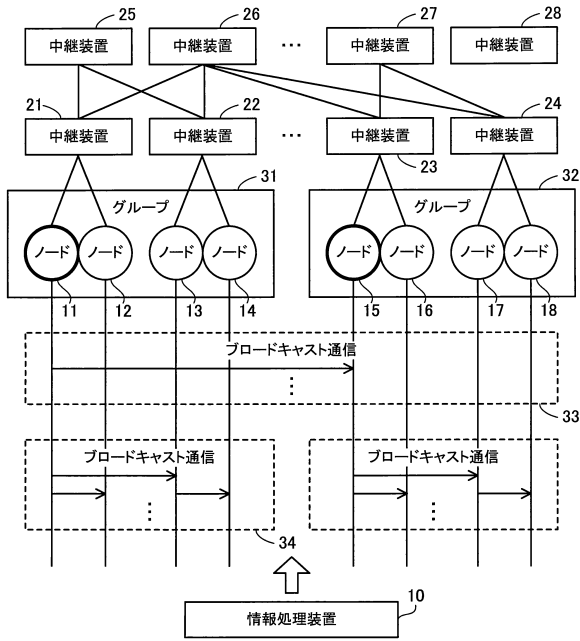
30

40

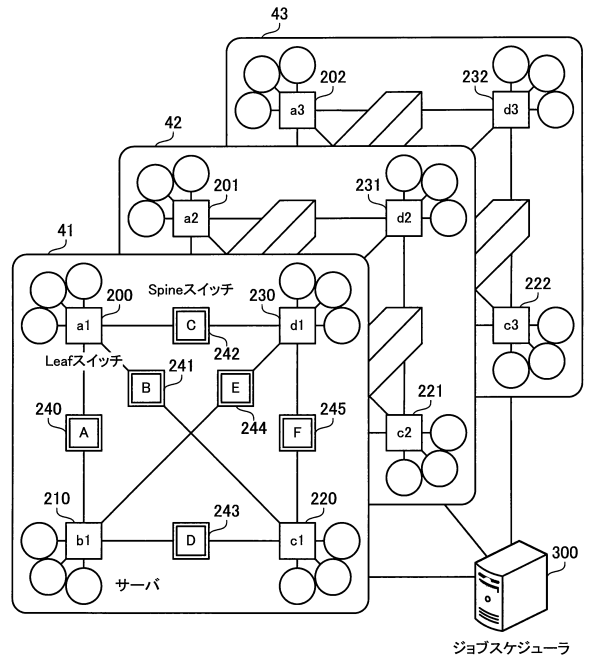
50

【図面】

【図 1】



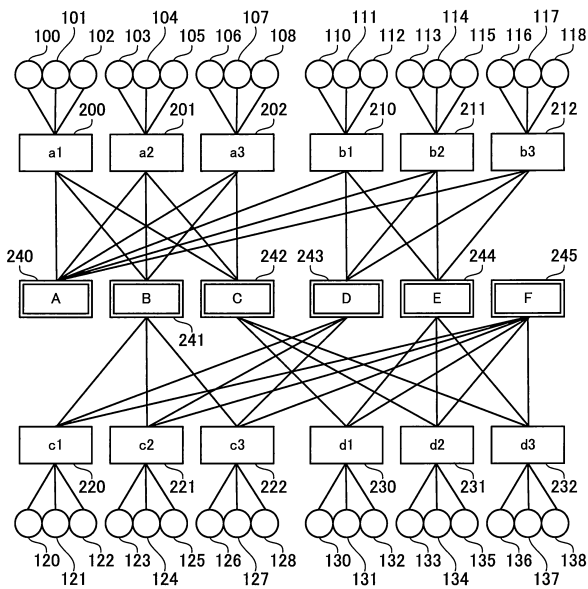
【図 2】



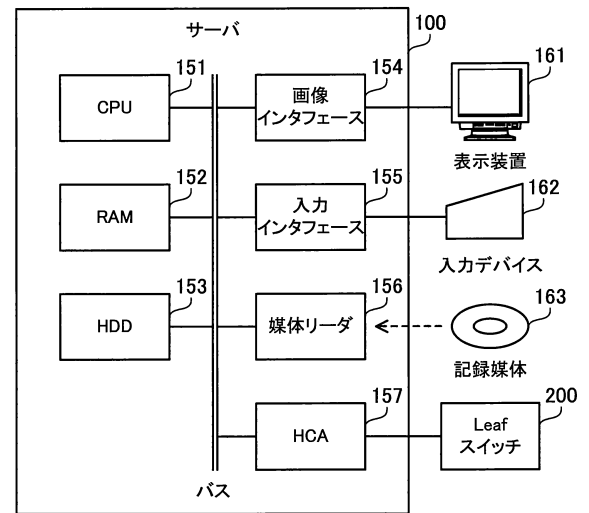
10

20

【図 3】



【図 4】

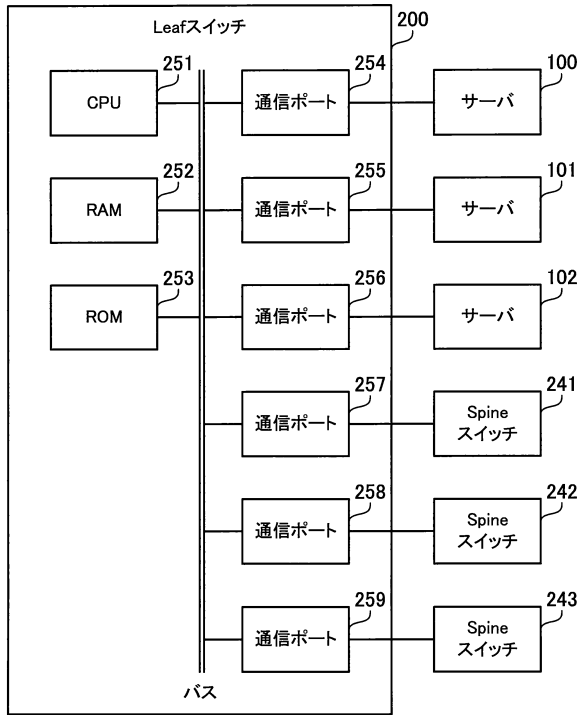


30

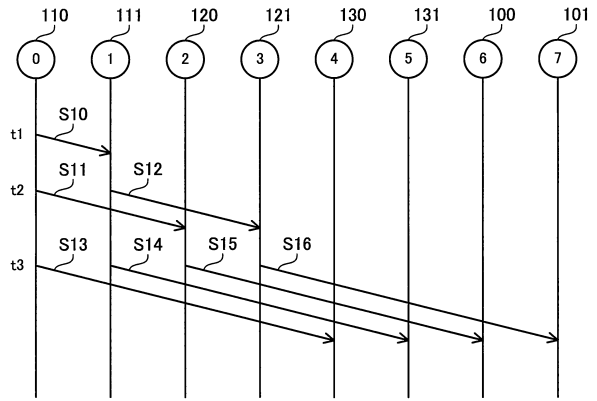
40

50

【 図 5 】



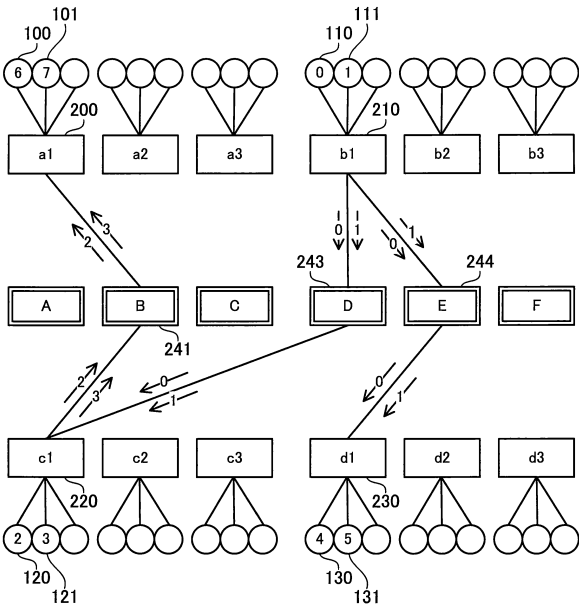
【 図 6 】



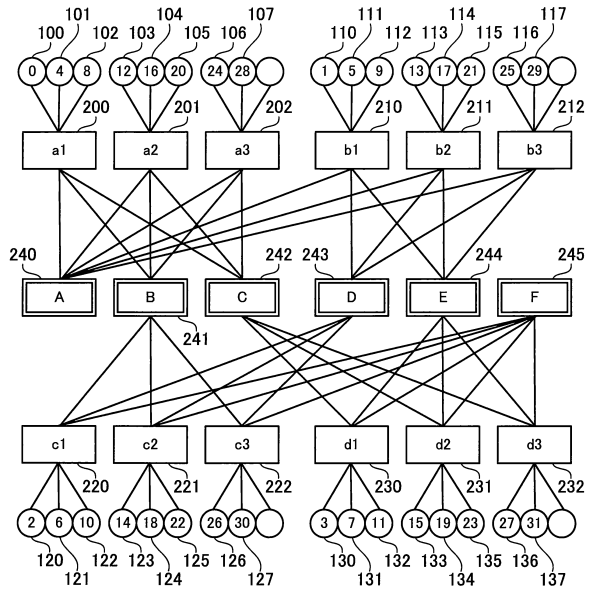
10

20

【 図 7 】



【 図 8 】

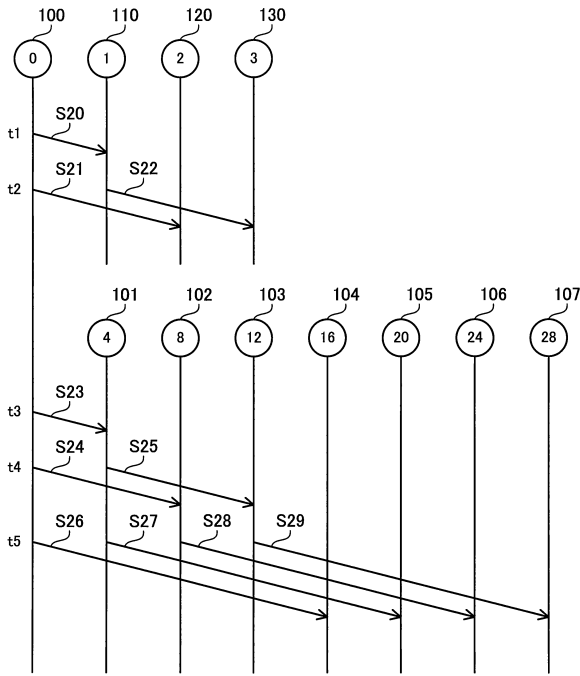


30

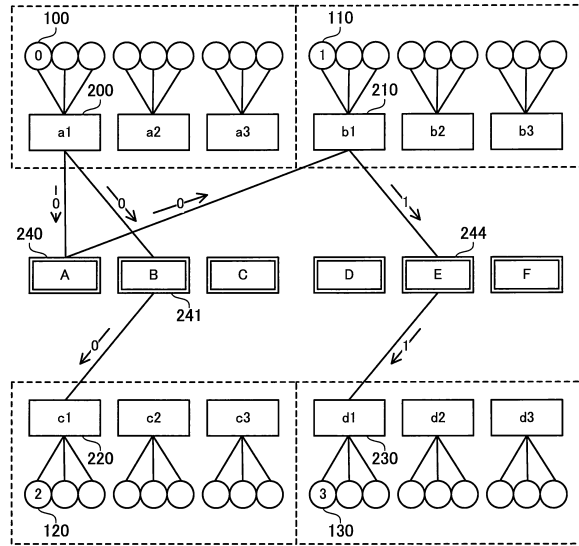
40

50

【図 9】



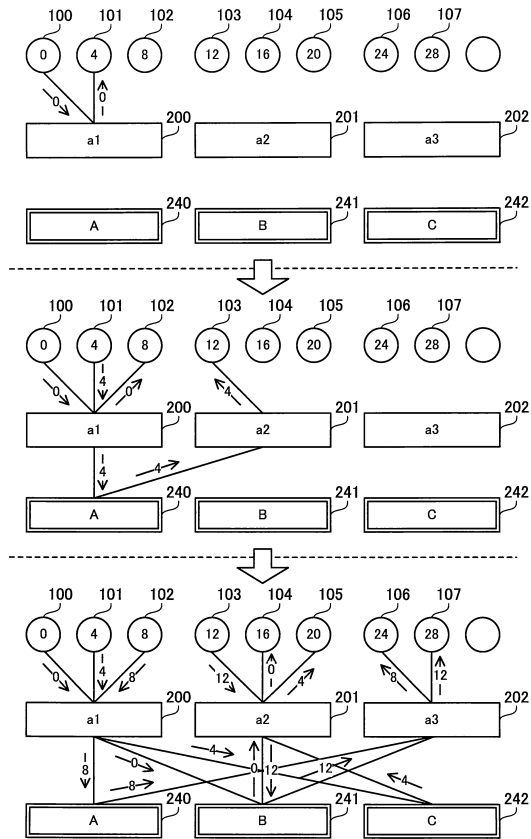
【図 10】



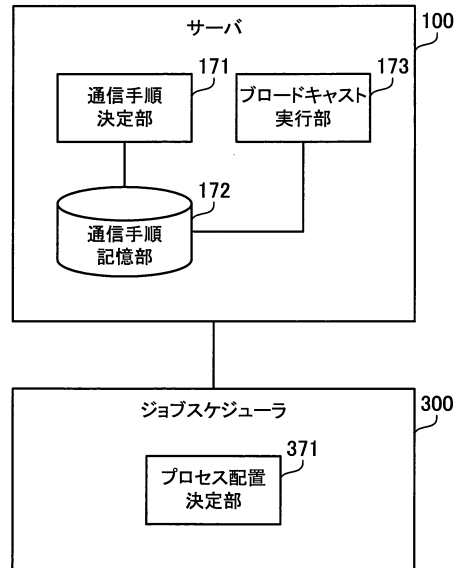
10

20

【図 11】



【図 12】



30

40

50

【 図 1 3 】

プロセス配置テーブル 174

ランク	ノードID
0	n0
1	n9
2	n18
3	n27
4	n1
...	...

【 図 1 4 】

通信手順記憶部 172

送信手順テーブル 175

ランク	t1	t2	t3	t4	t5
0	1	2	4	8	16
1	-1	3	5	9	17
2	-1	-1	6	10	18
3	-1	-1	7	11	19
4	-1	-1	-1	12	20
...

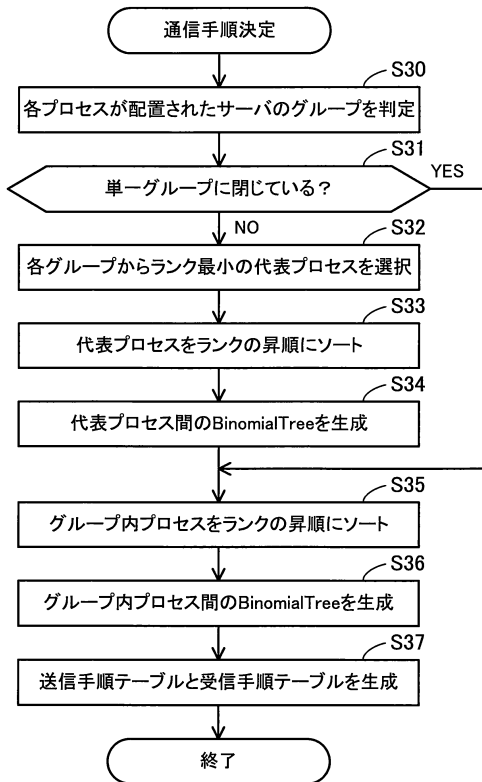
受信手順テーブル 176

ランク	t1	t2	t3	t4	t5
0	-1	-1	-1	-1	-1
1	0	-1	-1	-1	-1
2	-1	0	-1	-1	-1
3	-1	1	-1	-1	-1
4	-1	-1	0	-1	-1
...

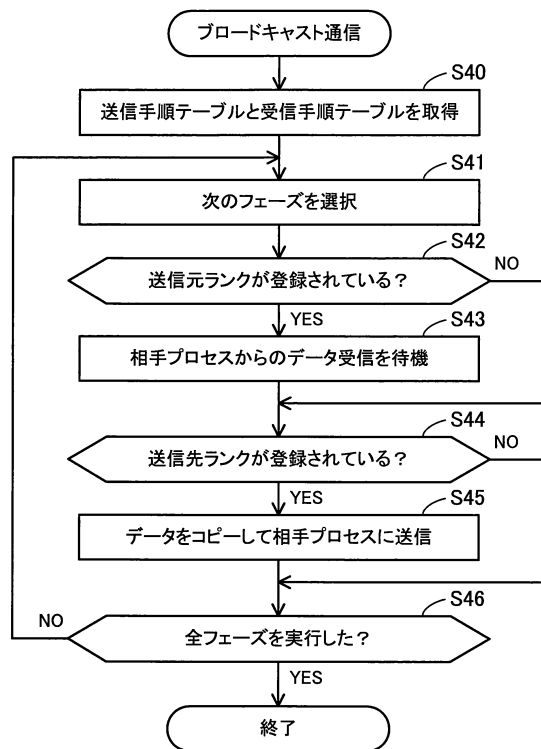
10

20

【 図 1 5 】



【 図 1 6 】



30

40

50

フロントページの続き

- (56)参考文献 特表 2020 - 512638 (JP, A)
特開 2018 - 185650 (JP, A)
特表 2009 - 519504 (JP, A)
特開 2009 - 193255 (JP, A)
米国特許出願公開第 2014 / 0156890 (US, A1)
米国特許出願公開第 2015 / 0160965 (US, A1)
千葉 立寛 ほか, グリッド環境におけるマルチレーンを用いた MPI コレクティブ通信アルゴリズム MPI Collective Operations Algorithm by Using Multi-lane for Grid Environment, 情報処理学会論文誌 第 48 巻 No. SIG8 (ACS18) IPSJ, 日本, 社団法人情報処理学会 Information Processing Society of Japan, 2007年05月15日, 第48巻, pp. 104-113
- (58)調査した分野 (Int.Cl., DB名)
H04L 49/201
H04L 47/60