

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6492080号
(P6492080)

(45) 発行日 平成31年3月27日 (2019. 3. 27)

(24) 登録日 平成31年3月8日 (2019. 3. 8)

(51) Int. Cl.

H04L 12/44 (2006.01)

F I

H04L 12/44

Z

請求項の数 13 (全 18 頁)

(21) 出願番号 特願2016-531743 (P2016-531743)
 (86) (22) 出願日 平成26年7月18日 (2014. 7. 18)
 (65) 公表番号 特表2016-525853 (P2016-525853A)
 (43) 公表日 平成28年8月25日 (2016. 8. 25)
 (86) 国際出願番号 PCT/US2014/047280
 (87) 国際公開番号 W02015/017145
 (87) 国際公開日 平成27年2月5日 (2015. 2. 5)
 審査請求日 平成29年6月21日 (2017. 6. 21)
 (31) 優先権主張番号 61/859, 650
 (32) 優先日 平成25年7月29日 (2013. 7. 29)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 14/226, 288
 (32) 優先日 平成26年3月26日 (2014. 3. 26)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 502303739
 オラクル・インターナショナル・コーポレ
 イション
 アメリカ合衆国カリフォルニア州9406
 5レッドウッド・シティ、オラクル・パ
 ークウェイ500
 (74) 代理人 110001195
 特許業務法人深見特許事務所
 (72) 発明者 ボグダンスキー、バルトシュ
 ノルウェー、エヌー0275 オスロ、エ
 イチ・0203、ホフ・テラス、15
 (72) 発明者 ヨンセン、ビョルン・ダグ
 ノルウェー、エヌー0687 オスロ、ビ
 ルベルクグレンダ、9

最終頁に続く

(54) 【発明の名称】 ミドルウェアマシン環境におけるマルチホーム型ファットツリー・ルーティングをサポートする
 ためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

ネットワーク環境のサブネットにおいてマルチホーム型ルーティングをサポートするた
 めの方法であって、

前記サブネットにおいてリーフスイッチ上のスイッチポートに関連付けられたエンドノ
 ードを提供するステップを含み、前記エンドノードは複数のポートを有し、前記方法はさ
 らに、

前記エンドノードが有する各々の前記ポートのためにルーティングを実行するステップ
 と、

前記ルーティングのために、前記エンドノードが有する前記複数のポートが相互に独立
 した経路を取ることを保証するステップとを含む、方法。

【請求項 2】

前記エンドノードが有する前記複数のポートは、第1ポートおよび第2ポートを含み、
 前記サブネットは、複数の前記エンドノードと、前記複数のエンドノードに接続される
 複数の前記リーフスイッチと、前記複数のリーフスイッチに接続される複数のルートスイ
 ッチとを備え、前記第1ポートおよび前記第2ポートのそれぞれは、前記複数のリーフス
 イッチのうちの異なるリーフスイッチに接続され、

前記複数のルートスイッチ、前記複数のリーフスイッチおよび前記複数のエンドノード
 の間の接続はファットツリー・トポロジーに基づいている、請求項1に記載の方法。

【請求項 3】

10

20

前記エンドノードは、複数のポートを介して前記ファットツリー・トポロジーの2つ以上の部分に接続されるマルチホーム型ノードである、請求項2に記載の方法。

【請求項4】

前記エンドノードの前記複数のポートにおける1つのポートに関連付けられた経路上の各々のスイッチにマーク付けするステップと、

前記エンドノードの前記複数のポートにおける別のポートに関連付けられた別の経路が、前記マーク付けされたスイッチを用いることを防止するステップとをさらに含む、請求項1～3のいずれか1項に記載の方法。

【請求項5】

冗長なスイッチがない場合に、前記エンドノードの前記複数のポートにおける別のポートに関連付けられた別の経路が、1つ以上のマーク付けされたスイッチを用いることを可能にするステップと、

相互に独立したリンクが前記1つ以上のマーク付けされたスイッチ上に存在する場合に、同じエンドノードの複数のポートのうちの異なる目標ポートのための独立したリンクを選択するステップとをさらに含む、請求項4に記載の方法。

【請求項6】

前記エンドノードの前記複数のポートの前記ルーティングを完了した後に各々のマーク付けされたスイッチのマーク付けを解除するステップをさらに含む、請求項4または5に記載の方法。

【請求項7】

前記エンドノードからの相互に独立した各々の経路を異なるスパインスイッチに関連付けるステップをさらに含む、請求項1～6のいずれか1項に記載の方法。

【請求項8】

前記エンドノードの前記複数のポートの前記ルーティングを完了した後に、前記エンドノードをルーティングされたエンドノードとしてマーク付けするステップをさらに含む、請求項1～7のいずれか1項に記載の方法。

【請求項9】

前記エンドノードが別のリーフスイッチに遭遇した場合に、前記エンドノードが再びルーティングされることを防止するステップをさらに含む、請求項8に記載の方法。

【請求項10】

前記エンドノードをパラメータとして採用するルーティングアルゴリズムを用いるステップをさらに含む、請求項1～9のいずれか1項に記載の方法。

【請求項11】

システムによって実行されたときに請求項1～10のいずれか1項に記載の方法を前記システムに実行させる命令を含むコンピュータプログラム。

【請求項12】

命令を有するコンピュータプログラムであって、前記命令が実行されると、ネットワーク環境のサブネットにおいてリーフスイッチ上のスイッチポートに関連付けられ、複数のポートを有するエンドノードを提供するステップと、

前記エンドノードの各々の前記ポートのためにルーティングを実行するステップと、
前記エンドノードの前記複数のポートが相互に独立した経路を取ることを保証するステップとをシステムに実行させる、コンピュータプログラム。

【請求項13】

ネットワーク環境のサブネットにおいてマルチホーム型ルーティングをサポートするためのシステムであって、

1つ以上のマイクロプロセッサと、

前記1つ以上のマイクロプロセッサ上で実行されるサブネットマネージャとを含み、前記サブネットマネージャは、

前記サブネットにおいてリーフスイッチ上のスイッチポートに、複数のポートを有する
エンドノードを関連付け、

10

20

30

40

50

前記エンドノードの各々の前記ポートのためにルーティングを実行し、
前記エンドノードの前記複数のポートが前記ルーティングのために相互に独立した経路
を取ることを保証するように動作する、システム。

【発明の詳細な説明】

【技術分野】

【0001】

著作権表示：

この特許文献の開示の一部は、著作権保護の対象となる題材を含んでいる。著作権の所有者は、特許商標庁の包袋または記録に掲載されるように特許文献または特許情報開示を誰でも複製できることに對して異議はないが、その他の点ではすべての如何なる著作権をも保有する。

10

【0002】

発明の分野：

本発明は概してコンピュータシステムに関し、特にネットワーク環境に関する。

【背景技術】

【0003】

背景：

ファットツリー・トポロジは、高性能コンピューティング（HPC：high performance computing）クラスタと、インフィニバンド（IB：InfiniBand）技術に基づいたクラスタとのために用いられる。たとえば、ファットツリー・トポロジは、MilkyWay-2などの最速のスーパーコンピュータにおいて用いられる。また、ファットツリーIBシステムは、Stampede、TGCC CurieおよびSuperMUCなどの大型の設備を含む。

20

【0004】

これらは、概して、本発明の実施形態が対処するように意図された分野である。

【発明の概要】

【課題を解決するための手段】

【0005】

概要：

この明細書中に記載されるシステムおよび方法は、ファットツリーまたは同様のトポロジを用いるインフィニバンドアーキテクチャに基づき得るネットワーク環境において、マルチホーム型ルーティングをサポートすることができる。システムは、ネットワークアプリケーションにおいてリーフスイッチ上のスイッチポートに関連付けられたエンドノードを提供することができる。さらに、システムは、エンドノード上の複数のポートの各々のためにルーティングを実行し、エンドノード上の複数のポートが相互に独立した経路を取ることを保証することができる。

30

【図面の簡単な説明】

【0006】

【図1】ネットワーク環境におけるファットツリー・ルーティングをサポートする例を示す図である。

40

【図2】本発明の実施形態に従った、ネットワーク環境におけるマルチホーム型ルーティングをサポートする例を示す図である。

【図3】本発明の実施形態に従った、ネットワーク環境におけるファットツリー・ルーティングをサポートするために冗長性を提供する例を示す図である。

【図4】本発明の実施形態に従った、ネットワーク環境におけるマルチホーム型ルーティングをサポートするための例示的なフローチャートである。

【図5】本発明の実施形態が実現され得るコンピュータシステムを示すブロック図である。

【図6】本発明の実施形態を実現するためのシステムを示すブロック図である。

【発明を実施するための形態】

50

【 0 0 0 7 】

詳細な説明：

本発明は、添付の図面において、限定ではなく例示のために示されており、図中同様の参照符号は同様の要素を指している。この開示において「或る」、「1つの」または「いくつかの」実施形態と言及する場合、これは必ずしも同じ実施形態である必要はなく、このように言及する場合、少なくとも1つを意味している。

【 0 0 0 8 】

以下の本発明の説明では、高性能ネットワークについての一例としてインフィニバンド（IB）ネットワークを用いている。他のタイプの高性能ネットワークを制限なく使用できることが当業者にとって明らかになるだろう。また、以下の本発明の説明では、ファブリック・トポロジーについての一例としてファットツリー・トポロジーを用いている。他のタイプのファブリック・トポロジーを制限なく使用できることが当業者にとって明らかになるだろう。

【 0 0 0 9 】

この明細書中に記載されるシステムおよび方法により、ネットワーク環境におけるマルチホーム型ルーティングをサポートすることができる。

【 0 0 1 0 】

インフィニバンドアーキテクチャ

インフィニバンドアーキテクチャ（IBA：InfiniBand Architecture）は2層型トポロジー分割をサポートする。下層においては、IBネットワークはサブネットと称される。この場合、サブネットは、スイッチおよび二地点間リンクを用いて相互接続される1組のホストを含み得る。上層においては、IBファブリックは、ルータを用いて相互接続することができる1つ以上のサブネットを構成する。

【 0 0 1 1 】

さらに、サブネット内のホストおよびスイッチは、ローカル識別子（LID：local identifier）を用いてアドレス指定することができ、単一のサブネットは49151LIDに制限することができる。サブネット内においてのみ有効なローカルアドレスであるLIDに加えて、各々のIBデバイスは、その不揮発性メモリに焼きつけられる64ビットのグローバル固有識別子（GUID：global unique identifier）を有し得る。GUIDは、IB層3（L3）アドレスであるグローバル識別子（GID：global identifier）を形成するために用いることができる。GIDは、IPv6のような128ビットアドレスを形成するように64ビットサブネット識別子（ID：identifier）を64ビットGUIDと連結することによって作成することができる。たとえば、さまざまなポートGUIDを、IBファブリックに接続されたポートに割当てることができる。

【 0 0 1 2 】

加えて、サブネットマネージャ（SM：subnet manager）は、IBファブリックにおけるルーティングテーブル計算を実行する役割を果たすことができる。ここで、IBネットワークのルーティングは、ローカルサブネットにおけるすべてのソースと宛先対との間において完全な接続性、デッドロック自由性および適切なロードバランシングを獲得することを目的としている。

【 0 0 1 3 】

サブネットマネージャは、ネットワーク初期化時間にルーティングテーブルを計算することができる。さらに、ルーティングテーブルは、最適な性能を保証するために、トポロジーが変化するたびに更新することができる。通常の動作中、サブネットマネージャは、トポロジーの変化をチェックするために、ネットワークに対して周期的な光掃引を実行することができる。光掃引中に変化が発見された場合、または、ネットワークの変化を信号で伝えるメッセージ（トラップ）がサブネットマネージャによって受信された場合、サブネットマネージャは、発見された変化に応じてネットワークを再構成することができる。

【 0 0 1 4 】

たとえば、サブネットマネージャは、たとえば、リンクがダウンした場合、デバイスが

10

20

30

40

50

追加された場合、または、リンクが排除された場合など、ネットワークトポロジーが変化した場合に、ネットワークを再構成することができる。再構成するステップは、ネットワーク初期化中に実行されるステップを含み得る。さらに、再構成には、ネットワークの変化が起こった局所的な範囲が含まれてもよく、この局所的な範囲は、サブネットに制限されている。また、大型のファブリックをルータでセグメント化することにより、再構成範囲が制限される可能性がある。

【 0 0 1 5 】

加えて、IBネットワークは、無損失のネットワーキング技術に基づいたものであって、いくつかの条件下ではデッドロックし易くなる可能性がある。たとえば、デッドロックが起こる可能性があるIBネットワークにおいては、バッファまたはチャネルなどのネットワークリソースが共有されており、パケットドロップが許容されていない。ここで、デッドロックを発生させる必須条件として、周期的な信頼依存性の生成が挙げられる。これは、周期的な信頼依存性がデッドロックを発生させる可能性があることを意味している。他方で、これは、周期的な信頼依存性が存在しているときは常にデッドロックが起こるであろうことを意味しているわけではない。

【 0 0 1 6 】

ファットツリー・ルーティング

ファットツリー・トポロジーは、高性能相互接続をサポートするためのさまざまな利点を提供することができる。これらの利点は、デッドロック自由性、固有のフォールトトレランスおよび完全な二分帯域幅を含み得る。デッドロック自由性があれば、デッドロックの回避を特別に考慮することなく、ツリー構造を用いることによってファットツリーをルーティングすることが可能となる。固有のフォールトトレランスがあれば、個々のソース宛先対の間に複数の経路が存在することによりネットワークの障害に対処することが容易になる。完全な二分帯域幅があれば、ネットワークは、二等分にした当該ネットワーク間における通信を最高速度で維持することができる。

【 0 0 1 7 】

さらに、ファットツリー・ルーティングアルゴリズムを用いることにより、基礎をなすファットツリー・トポロジーの効率的な使用をサポートすることができる。以下のアルゴリズム 1 は、例示的なファットツリー・ルーティングアルゴリズムである。

【 0 0 1 8 】

【表 1】

アルゴリズム 1 *route_to_cns()* 機能

要求: アドレス指定が完了

保証: すべての *hca_ports* がルーティングされる

```

1: for swleaf = 0 to max_leaf_sw do
2:     for swleaf.port = 0 to max_ports do
3:         hca_lid = swleaf.port -> remote_lid
4:         swleaf.routing_table[hca_lid] = swleaf.port
5:         route_downgoing_by_going_up()
6:     end for
7: end for

```

【 0 0 1 9 】

上に示されるように、ルーティング機能である *route_to_cns()* は、一連のリーフスイッチにわたって繰り返すことができる (1 ~ 7 行目)。選択されたリーフスイッチごとに、ルーティング機能は、たとえば、ポート番号付けシーケンスにおいて、選択されたリーフスイッチに接続される各々のエンドノードポートに向けてルーティングす

ることができる（２～６行目）。

【００２０】

さらに、特定のＬＩＤに関連付けられるエンドノードポートをルーティングする場合、ルーティング機能は、ネットワークトポロジーにおいて１レベル上がって下降経路をルーティングすることができ、各々のスイッチポートをルーティングする場合、ルーティング機能が下がり、上昇経路をルーティングすることができる。このプロセスは、ルートスイッチレベルに達するまで繰り返すことができる。その後、すべてのノードに向かう経路がルーティングされ、ファブリックにおけるすべてのスイッチのリニアフォワーディングテーブル（ＬＦＴ：linear forwarding table）に挿入される。

【００２１】

たとえば、`route__downgoing__by__going__up()`機能（５行目）は反復機能であってもよく、この反復機能は、経路のバランスをとることができ、`route__upgoing__by__going__down()`機能呼び出し、この機能が、ファットツリーにおける上り経路をスイッチを介してルーティングするが、このルーティングの宛先からは、`route__downgoing__by__going__up()`機能が呼び出されている。

【００２２】

さらに、`route__to__cns()`機能に関していくつかの潜在的な欠点が生じる可能性がある。第一に、`route__to__cns()`機能は記憶されておらず、エンドポートがどのエンドノードに属しているかを考慮することなくエンドポートをルーティン

【００２３】

図１は、ネットワーク環境においてファットツリー・ルーティングをサポートする例を示す。図１に示されるように、１つ以上のエンドノード１０１～１０４はネットワークファブリック１００に接続することができる。ネットワークファブリック１００は、ファットツリー・トポロジーに基づき得るものであって、複数のリーフスイッチ１１１～１１４および複数のスパインスイッチまたはルートスイッチ１３１～１３４を含む。加えて、ネットワークファブリック１００は、スイッチ１２１～１２４などの１つ以上の中間スイッチを含み得る。

【００２４】

また図１に示されるように、エンドノード１０１～１０４の各々は、マルチホーム型ノードであってもよく、すなわち、複数のポートを介してネットワークファブリック１００の２つ以上の部分に接続される単一ノードであってもよい。たとえば、ノード１０１は、ポートＨ１およびＨ２を含み得る。ノード１０２はポートＨ３およびＨ４を含み得る。ノード１０３はポートＨ５およびＨ６を含み得る。ノード１０４はポートＨ７およびＨ８を含み得る。

【００２５】

加えて、各々のスイッチは複数のスイッチポートを有し得る。たとえば、ルートスイッチＳ１ １３１はスイッチポート１、２を有し得る。ルートスイッチＳ２ １３２はスイッチポート３、４を有し得る。ルートスイッチＳ３ １３３はスイッチポート５、６を有し得る。ルートスイッチＳ４ １３４はスイッチポート７、８を有し得る。

【００２６】

リーフスイッチベースでルーティングするアルゴリズム１などのファットツリー・ルーティングアルゴリズムを用いると、独立したルートが異なる２ポートノード１０１～１０４に割当てられることが保証されなくなる。たとえば、ポートＨ１、Ｈ２、Ｈ５およびＨ６は、各々のスイッチ（図示せず）上のポート１に接続することができ、ポートＨ３、Ｈ４、Ｈ７およびＨ８は、各々のスイッチ（図示せず）上のポート２に接続される。ここで、４つのエンドポートを介してルーティングし、リーフスイッチ１１３およびスイッチ１２３を通過した後、ファットツリー・ルーティングアルゴリズムは、ノード１０１上で対

10

20

30

40

50

をなすエンドポートH 1およびH 2からの2つの経路を、(それぞれ、スイッチポート1、2を介して)同じ最左端のルートスイッチにS 1 1 3 1に割り当て得る。同様に、他のエンドポートの対、たとえば、ノード1 0 2上のH 3およびH 4、ノード1 0 3上のH 5およびH 6、ならびにノード1 0 4上のH 7およびH 8は、同じルートスイッチ(すなわち、S 2 1 3 2~S 4 1 3 4のそれぞれ)を介してルーティングされてもよい。

【0027】

これにより、結果として、ユーザにとって不所望な動作が生じる可能性がある。図1に示されるように、エンドノード1 0 1は、当該エンドノード1 0 1が、固有の物理的なフォールトトレランス(すなわち、さまざまなリーフスイッチ1 1 1および1 1 3に接続された2つのエンドポート)を有する可能性があったとしても、ルートスイッチS 1 1 3 1において単一障害点を被る可能性がある。加えて、物理的な配線に応じて、同様の問題が起こる可能性があり、単一障害点がファットツリー・トポロジにおける他のスイッチ上で発生する可能性がある。

【0028】

さらに、ネットワークファブリック1 0 0内においては、同じノード上における異なるポートへのトラフィックが単一のリンクを介してルーティングされてもよい。このため、このような単一のリンクは、エンドポートの組についての付加的な単一障害点と、性能ボトルネックとの両方を呈する可能性がある(なぜなら、異なるエンドポートを対象としたトラフィックが有効に利用できるのが単一の共有リンクの帯域幅だけであるかもしれないからである)。

【0029】

マルチホーム型ファットツリー・ルーティング

本発明の実施形態に従うと、システムは、ファットツリーにおけるマルチホーム型ノードのために独立したルートを提供することができ、そのため、単一障害点によって完全な機能停止がもたらされる可能性がなくなる。

【0030】

図2は、本発明の実施形態に従った、ネットワーク環境におけるマルチホーム型ルーティングをサポートする例を示す。図2に示されるように、ネットワーク環境は、複数のエンドノード(たとえば、ノード2 0 1~2 0 4)を含み得る。これら複数のエンドノードは各々、1つ以上のポート(たとえば、ポートH 1~H 8)を含み得る。加えて、複数のエンドノード2 0 1~2 0 4は、ファットツリー・トポロジにあり得るネットワークファブリック2 0 0に接続することができる。また、ネットワークファブリック2 0 0は、複数のスイッチ、たとえばリーフスイッチ2 1 1~2 1 4およびルートスイッチS 1 2 3 1およびS 2 2 3 2を含み得る。

【0031】

本発明の実施形態に従うと、m F t r e eアルゴリズムなどのマルチホーム型ファットツリー・ルーティングアルゴリズムは、たとえば、サブネットマネージャ(S M)2 1 0によって、ファットツリー・ルーティングを実行するために用いることができる。図2に示される例においては、m F t r e eアルゴリズムは、相互に冗長なやり方でルーティングされることがあり得るノード2 0 1上のポートH 1およびポートH 2からの経路を識別することができる。なぜなら、両方のポートが単一のエンドノード1上に位置しているからである。

【0032】

さらに、m F t r e eアルゴリズムは、経路を実際に冗長にすることを保証することができる。たとえば、ノード2 0 1上のポートH 1からの経路は、リーフスイッチ2 1 1を通過し、最終的にルートスイッチ2 3 1に到達し得る。図2に示されるように、システムは、経路における(濃い陰影で示される)スイッチにマーク付けすることができる。そして、システムは、ノード2 0 1上のポートH 2からの経路を決定するために、マーク付けされたスイッチを使用しないようにすることができる。これにより、ノード2 0 1上のポートH 2からの経路は、(たとえば、リーフスイッチ2 1 3を介する)冗長経路を通過す

10

20

30

40

50

ることができ、最終的に、（軽く陰影を付けて示される）別のルートスイッチ 2 3 2 に到達する。

【 0 0 3 3 】

ノード 2 0 1 のためのルーティングステップが完了すると、アルゴリズムは、（太線で示されるように）ルーティングされたノードにマーク付けすることができる。これにより、アルゴリズムがノード 2 0 1 の別のポートに達したときに、このノード 2 0 1 についてルーティングステップが繰り返されないようにする。こうして、システムは、単一障害点が、マルチポートノードの完全な機能停止をもたらさないようにすることを保証することができる。

【 0 0 3 4 】

加えて、ファットツリー・ルーティングアルゴリズムは、インフィニバンド（IB）ファットツリー・トポロジーの性能、スケーラビリティ、利用可能性および予測可能性を改善させることができる。

【 0 0 3 5 】

以下のアルゴリズム 2 は、例示的なマルチホーム型ファットツリー・ルーティングアルゴリズムである。

【 0 0 3 6 】

【表 2】

アルゴリズム 2 *route_multihomed_cns()* 機能

要求: アドレス指定が完了

保証: すべての *hca_ports* が独立したスパインを介してルーティングされる

1: for *swleaf* = 0 to *leaf_sw_num* do

2: for *swleaf.port* = 0 to *max_ports* do

3: *hca_node* = *swleaf.port* -> *remote_node*

4: if *hca_node.routed* == *true* then

5: *continue*

6: end if

7: *route hcas(hca_node)*

8: end for

9: end for

【 0 0 3 7 】

上に示されるように、アルゴリズム 2 は、マルチホーム型ルーティングアルゴリズムであって、すべてのリーフスイッチに対して繰り返すことができ、さらに、各々のリーフスイッチのためのすべてのリーフスイッチポートに対して繰り返すことができる（1～9行目）。これにより、アルゴリズム 2 は、アルゴリズム 1 と同様に決定論的になり得る。

【 0 0 3 8 】

さらに、アルゴリズム 2 は、スイッチポートに関連付けられたエンドノードを発見するために、リーフスイッチ上のスイッチポートを採用することができる（3行目）。単にリーフスイッチに接続された遠隔ポートの L I D を採用するアルゴリズム 1 とは異なり、アルゴリズム 2 は、ルーティング演算を実行するためのパラメータとしてエンドノードを採用することができる（7行目）。

【 0 0 3 9 】

以下のアルゴリズム 3 は、ファットツリーにおける単一のエンドノードをルーティングするための例示的なアルゴリズムである。

【 0 0 4 0 】

【表 3】

 アルゴリズム 3 *route_hcas(hca)* 機能

要求: ルーティングされるべきノード

保証: *hca_lid* を有するノードに属するすべての *hca_ports* がルーティングされる

```

1: for hca_node.port = 0 to port_num do
2:   hca_lid = hca_node.port -> lid
3:   swleaf = hca_node.port -> remote_node
4:   swleaf.port = hca_node.port -> remote_port_number
5:   swleaf.routing_table[hca_lid] = swleaf.port
6:   route_downgoing_by_going_up()
7: end for
8: hca_node.routed = true
9: clear_redundant_flag()

```

10

【0041】

上に示されるように、アルゴリズム 3 は、選択されたエンドノード上のすべてのポートに対して繰り返すことができる（1～7行目）。たとえば、アルゴリズム 3 は、*route_downgoing_by_going_up()* 機能の修正バージョンを用いて、選択されたエンドノード上の各々のポートをルーティングすることができる（6行目）。選択されたエンドノード上のすべてのポートがルーティングされる場合、ルーティングアルゴリズムは、選択されたエンドノードをルーティングされるものとしてマーク付けすることができ（8行目）、これにより、別のリーフスイッチ上に到達したときにエンドノードがルーティングされないようにする。また、アルゴリズム 3 は、さまざまな状況でシステムの性能を改善させることができる（たとえば、アルゴリズム 3 は、2ポートノードのためにループ反復のうち半分を保存しておくことができる）。

20

【0042】

加えて、アルゴリズム 3 は、単一のホストチャンネルアダプタ（HCA: host channel adapter）上に複数のポートがあるシナリオと、2つ以上の HCA 上に複数のポートがあるとシナリオとの両方に適用され得る。アルゴリズムはさまざまな方法を用いて、同じ論理ノード上の単一の HCA または複数の HCA 上のポート（またはいずれかのエンドポート）を識別することができる。ここで、ノードは、物理サーバもしくは仮想サーバ、IO デバイス、または、1つ以上の HCA ポートを介して IB ファブリックに接続された何なるの種類のエンドノードであってもよい。

30

【0043】

さらに、アルゴリズム 3 は、選択されたノード上の各々のポートをルーティングことができ、フラグを用いて経路上の各々のスイッチにマーク付けすることができる。こうして、アルゴリズム 3 は、同じエンドノード上のさまざまなポートのためにさまざまなスイッチを選択することができる。その後、アルゴリズムは、次のノードにまで進み得るように、すべてのスイッチ上のフラグをひっくり返すことができる。

40

【0044】

加えて、システムは、*clear_redundant_flag()* 機能においてスイッチ冗長性フラグをクリアすることによって最適化することができる（9行目）。スイッチ冗長性フラグをクリアするための最適な方法としては、特定のスイッチが経路上にあったか否かに関わらずすべてのスイッチに対して繰り返されるループをこの機能で用いるのではなく、経路上にあるスイッチのリストを作成して、*clear_redundant_flag()* 機能がリスト中のそれらスイッチに対してのみ繰り返されることを保証する方法が挙げられる。

50

【 0 0 4 5 】

以下のアルゴリズム 4 は、ファットツリーにおける単一のエンドノードポートをルーティングするための例示的なアルゴリズムである。

【 0 0 4 6 】

【表 4】

 アルゴリズム 4 *route_downgoing_by_going_up()* 機能

要求: 現在のホップスイッチ

保証: 上向きの冗長スイッチを発見するために最大限の努力が払われる

保証: 経路上のスイッチを冗長とマーク付けする

10

1: *groupmin* = 0

2: *redundant_group* = 0

3: for *port_group* = 0 to *port_group_num* do

4: if *groupmin* == 0 then

5: if *groupmin* -> *remote_node.redundant* then

6: *groupmin* = *port_group*

7: end if

20

8: else if *port_group.cntdown* < *groupmin.cntdown* then

9: *groupmin* = *port_group*

10: if *groupmin* -> *remote_node.redundant* then

11: *min_redundant_group* = *groupmin*

12: end if

13: end if

14: end for

15: if *groupmin* == 0 then

30

16: *fallback_normal_routing(hca_lid)*

17: else if *groupmin* -> *remote_node.redundant* then

18: *groupmin* = *min_redundant_group*

19: *groupmin* -> *remote_node.redundant* = false

20: end if

【 0 0 4 7 】

上に示されるように、アルゴリズム 4 における *route_downgoing_by_going_up()* 機能の修正バージョンでは、冗長性が主要な問題として処理される。最低位の下方向カウンタを有するポートグループが選択されているアルゴリズム 1 とは異なり、アルゴリズム 4 は、当該アルゴリズム 4 がエンドノードに属する他のいずれのポートもルーティングしない場合（すなわち、冗長フラグが真である場合）、エンドノードの上向きノードを次のホップとして選択し得るだけである。

40

【 0 0 4 8 】

ここで、次のエンドノードがルーティングされる前に、冗長フラグがクリアされる。極めて過度にサブスクライブされたファブリックにおいて起こり得るような冗長なノードが存在しない場合、または、リンク障害がある場合、*m F t r e e* は通常のファットツリー・ルーティングにまで後退し、この場合、ユーザは、アルゴリズム 1 において提示されるルーティング機能と同様にこのルーティング機能が動作することを認識し得る。

【 0 0 4 9 】

50

加えて、2つのスイッチ間に代替的なスイッチおよび平行なリンクが存在しない場合、上述のアルゴリズム4は、性能/負荷の分散およびリンクレベル冗長性の両方をサポートするために、同じエンドノード上のさまざまな目標ポートのためにさまざまなリンクを選択することができる。他方では、性能が優先される場合、両方のリンクが同じレベルの負荷を有する場合にのみ、同じエンドノード上の異なる目標ポートに対して別個のリンクが選択され得る。

【0050】

図3は、ネットワーク環境においてファットツリー・ルーティングをサポートするために冗長性を提供する例を示す。図3に示されるように、1つ以上のマルチホーム型エンドノード301~304はネットワークファブリック300に接続することができる。ネットワークファブリック300は、ファットツリー・トポロジーに基づくものであって、複数のリーフスイッチ311~312と、複数のスパインスイッチまたはルートスイッチS1 331~S4 334とを含み得る。加えて、ネットワークファブリック300は、スイッチ321~324などの1つ以上の中間スイッチを含み得る。

【0051】

図3にも示されるように、ノード301はポートH1およびH2を含み得る。ノード302はポートH3およびH4を含み得る。ノード303はポートH5およびH6を含み得る。また、ノード304はポートH7およびH8を含み得る。加えて、ルートスイッチS1 331はスイッチポート1、2を有し得る。ルートスイッチS2 332はスイッチポート3、4を有し得る。ルートスイッチS3 333はスイッチポート5、6を有し得る。ルートスイッチS4 334はスイッチポート7、8を有し得る。

【0052】

さらに、mFtreeアルゴリズムなどのマルチホーム型ルーティングは、ノード上の各々のポートへの経路が排他的となるように、各々のマルチホーム型ノード301~304をルーティングすることができる。すなわち、mFtreeアルゴリズムは、独立した経路を介してマルチホーム型ノード上の各ポートに到達可能となるように保証する。加えて、mFtreeアルゴリズムはネットワーク性能を改善させることができる。

【0053】

さらに、単一のマルチホーム型エンドノードの場合、mFtreeアルゴリズムは、同じエンドノードに属するいずれのポート対に至る経路によっても単一のリンクが共有されないように保証することができる。また、ネットワークファブリック300において同じ宛先ノード上に異なるソースポートから異なるポートに至るコンカレントなトラフィックがある場合、mFtreeアルゴリズムは、代替的な経路が存在する場合に、コンカレントなトラフィックが如何なる中間リンクも共有しないことを保証することができる。

【0054】

このように、mFtreeアルゴリズムを用いると、ファブリック300におけるスパインスイッチS1 331などの単一デバイスが故障しても、ノード301の切断が引き起こされる可能性はない。なぜなら、異なるポートへの経路が単一のスパインスイッチS1 331においては収束しないからである。

【0055】

加えて、mFtreeアルゴリズムは、別個の独立したエンティティと同じノード上の各々のポートを処理する。これにより、mFtreeアルゴリズムは、ポートベースではなくノードベースでルーティングすることができ、mFtreeアルゴリズムは、異なるエンドノードが有し得るさまざまな特徴をアドレス指定することができる。

【0056】

図4は、本発明の実施形態に従った、ネットワーク環境においてマルチホーム型ルーティングをサポートするための例示的なフローチャートを示す。図4に示されるように、ステップ401において、システムは、ネットワークファブリックにおいてリーフスイッチ上のスイッチポートに関連付けられたエンドノードを提供することができる。この場合、エンドノードは複数のポートに関連付けられている。次いで、ステップ402において、

システムは、エンドノード上の各々の上記ポートのためにルーティングを実行することができる。さらに、ステップ403において、システムは、エンドノード上の複数のポートが相互に独立した経路を取ることを保証することができる。

【0057】

図5は、本発明の実施形態が実現され得るコンピュータシステムのブロック図を示す。コンピュータシステム500は、情報を通信するためのバスまたは他の通信メカニズムと、バスと結合されて情報を処理するためのハードウェアプロセッサとを含む。ハードウェアプロセッサは、たとえば、汎用のマイクロプロセッサであってもよい。

【0058】

コンピュータシステム500はまた、プロセッサによって実行されるべき命令および情報を格納するための、バスに結合されたランダムアクセスメモリ(RAM: random access memory)または他の動的記憶装置などのメインメモリを含む。メインメモリはまた、プロセッサによって実行される命令の実行中に一時変数または他の中間情報を格納するために用いられてもよい。このような命令がプロセッサにアクセス可能な非一時的な記憶媒体に格納されると、コンピュータシステム500が、命令で指定された動作を実行するようにカスタマイズされた専用マシンとなる。

【0059】

コンピュータシステム500はさらに、プロセッサのために静的情報および命令を格納するための、バスに結合された読取専用メモリ(ROM: read only memory)または他の静的記憶装置を含む。情報および命令を格納するための、バスに結合された磁気ディスクまたは光ディスクなどの記憶装置が提供される。

【0060】

コンピュータシステム500は、情報をコンピュータユーザに対して表示するための、陰極線管(CRT: cathode ray tube)などのディスプレイを含み得るか、または、当該ディスプレイに対してバスを介して結合され得る。英数字および他のキーもしくはカーソル制御を含む入力装置は、プロセッサに対して情報およびコマンド選択を伝えるためにバスに結合される。

【0061】

コンピュータシステム500はこの明細書中に記載される技術を実現し得るが、この場合、コンピュータシステムと組み合わせることで、コンピュータシステム500を専用マシンにするかまたは専用マシンにするようプログラムする、カスタマイズされたハードワイヤード論理、1つ以上のASICもしくはFPGA、ファームウェアおよび/またはプログラム論理を用い得る。一実施形態に従うと、この明細書中に記載される技術は、プロセッサがメインメモリに含まれる1つ以上の命令の1つ以上のシーケンスを実行することに応じて、コンピュータシステム500によって実行される。このような命令は、記憶装置などの別の記憶媒体からメインメモリに読み込まれてもよい。メインメモリに含まれる命令のシーケンスを実行することにより、この明細書中に記載されるプロセスステップをプロセッサに実行させる。代替的な実施形態においては、ハードワイヤード回路が、ソフトウェア命令の代わりに、またはソフトウェア命令と組み合わせで用いられてもよい。

【0062】

コンピュータシステム500はまた、バスに連結された通信インターフェイスを含む。通信インターフェイスは、統合サービスデジタル網(ISDN: integrated services digital network)カード、ケーブルモデム、衛星モデム、または対応するタイプの電話線にデータ通信接続するためのモデムであってもよい。別の例として、通信インターフェイスは、互換性のあるローカルエリアネットワーク(LAN: local area network)にデータ通信接続するためのLANカードであってもよい。ワイヤレスリンクが実現されてもよい。このような如何なる実現例においても、通信インターフェイスは、さまざまなタイプの情報を表わすデジタルデータストリームを搬送する電気信号、電磁信号または光信号を送受信する。

【0063】

図5に示されるように、記憶装置は、ネットワークファブリックにおいてリーフスイッチ上のスイッチポートにエンドノードを関連付けるように構成され得るサブネットマネージャを含んでもよい。この場合、エンドノードは、複数のポートに関連付けられている。サブネットマネージャはさらに、エンドノード上の各々の上記ポートのためにルーティングを実行し、エンドノード上の複数のポートが相互に独立した経路を取ることを保証するように構成され得る。

【0064】

いくつかの実施形態においては、サブネットマネージャはさらに、エンドノード上の複数のポートにおける或るポートに関連付けられた経路上の各々のスイッチにマーク付けし、エンドノード上の複数のポートにおける別のポートに関連付けられた別の経路が、マーク付けされたスイッチを用いることを防止するように構成されてもよい。

10

【0065】

いくつかの実施形態においては、サブネットマネージャはさらに、冗長なスイッチがない場合に、エンドノード上の複数のポートにおける別のポートに関連付けられた別の経路が、1つ以上のマーク付けされたスイッチを用いることを可能にし、平行なリンクが1つ以上のマーク付けされたスイッチ上に存在する場合に、同じエンドノード上の異なる目標ポートのために独立したリンクを選択するように構成されてもよい。

【0066】

いくつかの実施形態においては、サブネットマネージャはさらに、エンドノード上の複数のポートのルーティングを完了した後、各々のマーク付けされたスイッチのマーク付けを解除するように構成されてもよい。

20

【0067】

いくつかの実施形態においては、サブネットマネージャはさらに、エンドノード上の複数のポートのルーティングを完了した後、エンドノードをルーティングされたエンドノードとしてマーク付けするように構成されてもよい。

【0068】

いくつかの実施形態においては、サブネットマネージャはさらに、エンドノードが別のリーフスイッチに遭遇したときに、エンドノードが再びルーティングされることを防止するように構成されてもよい。

【0069】

30

図6は、本発明の実施形態を実現するためのシステム600のブロック図を示す。システム600のブロックは、本発明の原理を実行するために、ハードウェア、ソフトウェア、またはハードウェアとソフトウェアとの組み合わせによって実現されてもよい。上述のとおり本発明の原理を実現するために、図6において記載されたブロックが組み合わされ得るかまたはサブブロックに分解され得ることが当業者によって理解される。したがって、この明細書中における説明は、この明細書中に記載される機能ブロックの実現可能な任意の組み合わせまたは分解またはさらなる定義をサポートし得る。

【0070】

図6に示されるように、システム600は、ネットワークファブリックにおいてリーフスイッチ上のスイッチポートに関連付けられたエンドノードを提供するように構成された提供ユニット601を含み得る。エンドノードは、複数のポートに関連付けられている。システム600はさらに、エンドノード上の各々の上記ポートのためにルーティングを実行するように構成された実行ユニット602を含み得る。さらに、システム600は、エンドノード上の複数のポートが相互に独立した経路を取ることを保証するように構成された保証ユニット603をさらに含み得る。

40

【0071】

図6には示されないが、いくつかの実施形態においては、システム600はさらに、ネットワークファブリックをファットツリー・トポロジーに基づいたものにすることを可能にするように構成された第1の可能ユニットを含み得る。

【0072】

50

いくつかの実施形態においては、システム 600 はさらに、エンドノードが、複数のポートを介してファットツリー・トポロジーの 2 つ以上の部分に接続されるマルチホーム型ノードとなることを可能にするように構成された第 2 の可能ユニットを含み得る。

【0073】

いくつかの実施形態においては、システム 600 はさらに、エンドノード上の複数のポートにおける或るポートに関連付けられた経路上の各々のスイッチにマーク付けするように構成された第 1 のマーク付けユニットと、エンドノード上の複数のポートにおける別のポートに関連付けられた別の経路が当該マーク付けされたスイッチを用いることを防止するように構成された第 1 の防止ユニットとを含み得る。

【0074】

いくつかの実施形態においては、システム 600 はさらに、冗長なスイッチがない場合に、エンドノード上の複数のポートにおける別のポートに関連付けられた別の経路が 1 つ以上のマーク付けされたスイッチを用いることを可能にするように構成された第 3 の可能ユニットと、平行なリンクが 1 つ以上のマーク付けされたスイッチ上に存在する場合に、同じエンドノード上の異なる目標ポートのために独立したリンクを選択するように構成された選択ユニットとを含み得る。

【0075】

いくつかの実施形態においては、システム 600 はさらに、エンドノード上の複数のポートのルーティングを完了した後、各々のマーク付けされたスイッチのマーク付けを解除するように構成されたマーク付け解除ユニットを含み得る。

【0076】

いくつかの実施形態においては、システム 600 はさらに、エンドノードからの相互に独立した各々の経路を異なるスパインスイッチに関連付けるように構成された関連付けユニットを含み得る。

【0077】

いくつかの実施形態においては、システム 600 はさらに、エンドノード上の複数のポートのルーティングを完了した後、エンドノードをルーティングされたエンドノードとしてマーク付けするように構成された第 2 のマーク付けユニットを含み得る。

【0078】

いくつかの実施形態においては、システム 600 はさらに、エンドノードが別のリーフスイッチに遭遇したときに、エンドノードが再びルーティングされることを防止するように構成された第 2 の防止ユニットを含み得る。

【0079】

いくつかの実施形態においては、システム 600 はさらに、ルーティングアルゴリズムがエンドノードをパラメータとして採用することを可能にするように構成された第 4 の可能ユニットを含み得る。

【0080】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段は、ネットワークファブリックにおいてリーフスイッチ上のスイッチポートに関連付けられたエンドノードを提供することを含む。エンドノードは複数のポートに関連付けられている。当該手段はさらに、エンドノード上の各々の上記ポートのためにルーティングを実行することと、エンドノード上の複数のポートが上記ルーティングのために相互に独立した経路を取ることを保証することを含む。

【0081】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段においては、ネットワークファブリックはファットツリー・トポロジーに基づいている。

【0082】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段においては、エンドノードは、複数のポートを介してファットツリー・トポロジーの 2 つ以上の部分に接続されたマルチホーム型ノードである。

10

20

30

40

50

【 0 0 8 3 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、エンドノード上の複数のポートにおける或るポートに関連付けられた経路上の各々のスイッチにマーク付けすることと、エンドノード上の複数のポートにおける別のポートに関連付けられた別の経路が、マーク付けされたスイッチを用いることを防止することを含む。

【 0 0 8 4 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、冗長なスイッチがない場合に、エンドノード上の複数のポートにおける別のポートに関連付けられた別の経路が、1つ以上のマーク付けされたスイッチを用いることを可能にすることと、平行なリンクが1つ以上のマーク付けされたスイッチ上に存在する場合に、同じエンドノード上の異なる目標ポートのために独立したリンクを選択することを含む。

10

【 0 0 8 5 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、エンドノード上の複数のポートのルーティングを完了した後に各々のマーク付けされたスイッチのマーク付けを解除することをさらに含む。

【 0 0 8 6 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、エンドノードからの相互に独立した経路を異なるスパインスイッチに関連付けることを含む。

20

【 0 0 8 7 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、エンドノード上の複数のポートのルーティングを完了した後、エンドノードをルーティングされたエンドノードとしてマーク付けすることを含む。

【 0 0 8 8 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、エンドノードが別のリーフスイッチに遭遇した場合に、エンドノードが再びルーティングされることを防止することを含む。

【 0 0 8 9 】

ネットワーク環境においてマルチホーム型ルーティングをサポートするための手段はさらに、エンドノードをパラメータとして採用するルーティングアルゴリズムを用いることを含む。

30

【 0 0 9 0 】

本発明は、本開示の教示に従いプログラムされた、1つ以上のプロセッサ、メモリ、および/またはコンピュータ読取可能な記録媒体を含む、従来の汎用もしくは専用デジタルコンピュータ、コンピュータ化デバイス、マシン、またはマイクロプロセッサを1つ以上用いて、適宜実装し得る。適切なソフトウェアコーディングは、熟練したプログラマーが本開示の教示に基づいて容易に準備できるものである。これはソフトウェア技術における当業者には明らかであろう。

40

【 0 0 9 1 】

実施形態によっては、本発明は、本発明のプロセスのうちいずれかを実行するためにコンピュータをプログラムするのに使用できる命令が格納された記録媒体または(1つまたは複数の)コンピュータ読取可能な媒体であるコンピュータプログラムプロダクトを含む。この記録媒体の例は、フロッピー(登録商標)ディスク、光ディスク、DVD、CD-ROM、マイクロドライブ、および光磁気ディスクを含む、任意の種類のディスク、ROM、RAM、EPROM、EEPROM、DRAM、VRAM、フラッシュメモリデバイス、磁気もしくは光カード、ナノシステム(分子メモリICを含む)、または、命令および/またはデータを格納するのに適した任意の種類の媒体もしくはデバイスを含み得るものの、これらに限定されない。

50

【 0 0 9 2 】

本発明のこれまでの記載は例示および説明を目的として提供されている。すべてを網羅するまたは本発明を開示された形態そのものに限定することは意図されていない。当業者には数多くの変更および変形が明らかであろう。実施の形態は、本発明の原理およびその実際の応用を最もうまく説明することによって他の当業者がさまざまな実施の形態および意図している特定の用途に適したさまざまな変形を理解できるようにするために、選択され説明されている。

【 図 1 】

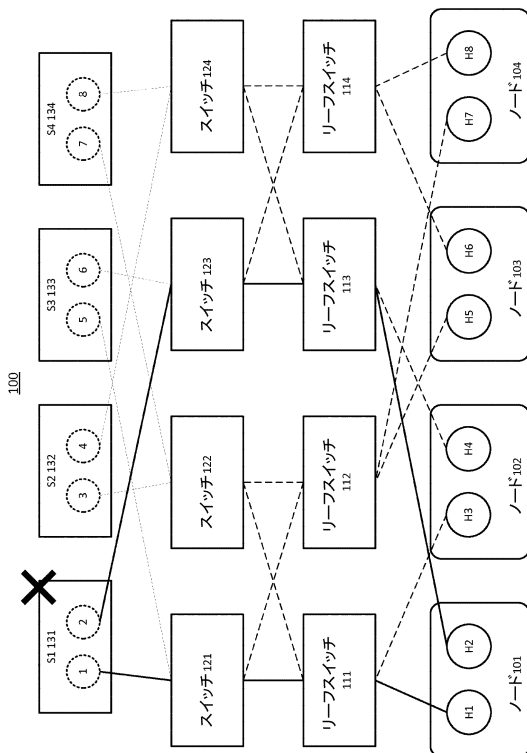


FIGURE 1

【 図 2 】

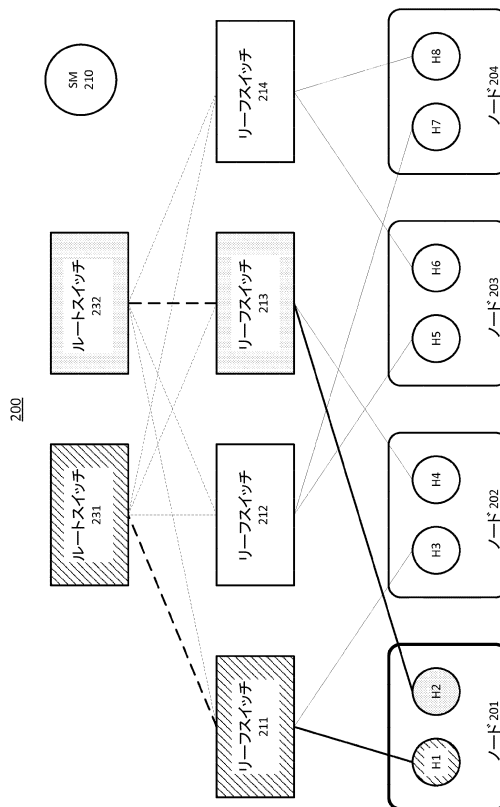


FIGURE 2

【 図 3 】

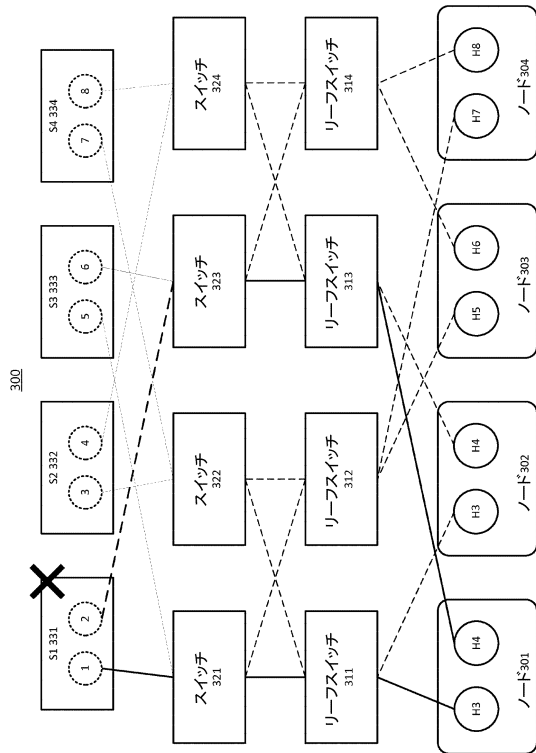


FIGURE 3

【 図 4 】

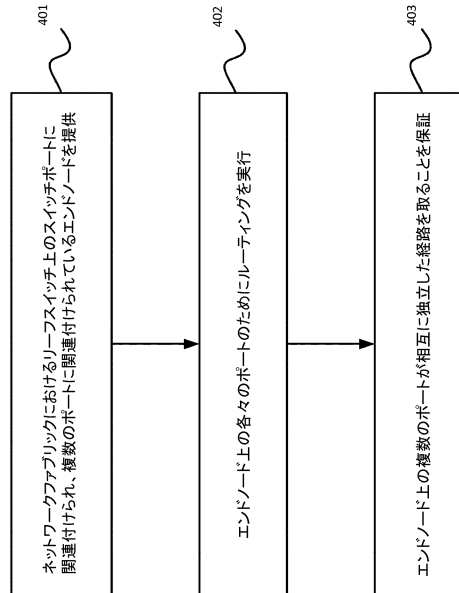


FIGURE 4

【 図 5 】

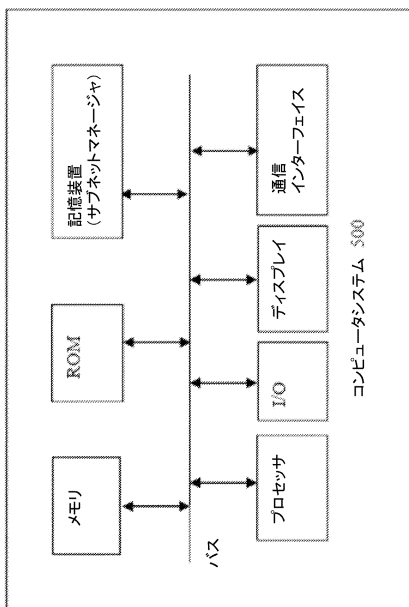


FIGURE 5

【 図 6 】

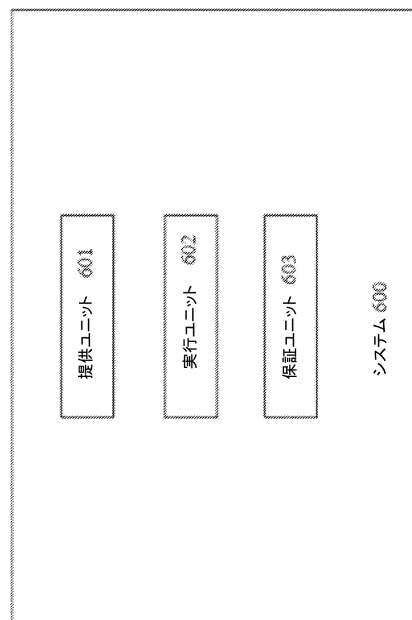


FIGURE 6

フロントページの続き

審査官 野元 久道

(56)参考文献 米国特許出願公開第2012/0300669(US, A1)

特開2007-318449(JP, A)

特開2013-051647(JP, A)

(58)調査した分野(Int.Cl., DB名)

H04L 12/44