

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6018217号
(P6018217)

(45) 発行日 平成28年11月2日(2016.11.2)

(24) 登録日 平成28年10月7日(2016.10.7)

(51) Int.Cl.

F I

H04L 12/753 (2013.01)

H04L 12/753

請求項の数 11 (全 18 頁)

(21) 出願番号	特願2014-541332 (P2014-541332)	(73) 特許権者	502303739
(86) (22) 出願日	平成24年11月9日 (2012.11.9)		オラクル・インターナショナル・コーポレ イション
(65) 公表番号	特表2014-533071 (P2014-533071A)		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
(43) 公表日	平成26年12月8日 (2014.12.8)		
(86) 国際出願番号	PCT/US2012/064479	(74) 代理人	110001195
(87) 国際公開番号	W02013/071130		特許業務法人深見特許事務所
(87) 国際公開日	平成25年5月16日 (2013.5.16)	(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
審査請求日	平成27年9月10日 (2015.9.10)		
(31) 優先権主張番号	61/557,722	審査官	安藤 一暁
(32) 優先日	平成23年11月9日 (2011.11.9)		
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	13/653,303		
(32) 優先日	平成24年10月16日 (2012.10.16)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 ファットツリートポロジーにおけるスイッチ間のデッドロックフリールーティングを提供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】

1つ以上のマイクロプロセッサに対して機能するミドルウェアマシン環境における複数のスイッチ間のルーティングをサポートする方法であって、

ルーティングアルゴリズムを用いて、前記ミドルウェアマシン環境におけるスイッチ間トラフィックのためのルーティングを実行すること、を含み、前記ルーティングアルゴリズムは、前記ミドルウェアマシン環境におけるファットツリートポロジーにおいてソーススイッチから宛先スイッチへのスイッチ間トラフィックをルーティングするアップ/ダウンターンモデルを含み、前記ファットツリーは複数のリーフスイッチとその上方のレベルの複数のスイッチとを含み、各リーフスイッチは終端ノードに接続されて、

前記複数のリーフスイッチから、特定のリーフスイッチを、ハブスイッチとして選択すること、

前記ハブスイッチの兄弟スイッチに関連付けられているルーティングテーブルが、前記宛先スイッチへの経路を含むことを判断することと、

前記ソーススイッチと前記ハブスイッチのそれぞれに、前記兄弟スイッチを通して、前記宛先スイッチへの経路を設定することと、

前記ソーススイッチと宛先スイッチと前記ハブスイッチを含む経路で、前記ハブスイッチと関連付けられたルーティングテーブルを更新することを含む、方法。

【請求項 2】

前記ミドルウェアマシン環境内の複数のスイッチをファットツリートポロジー内に存在

させることをさらに含む、請求項 1 に記載の方法。

【請求項 3】

前記特定のリーフスイッチは前記ファットツリートポロジー内の複数のスイッチ宛先に到達することができる、請求項 1 または 2 に記載の方法。

【請求項 4】

前記アップ/ダウンターンモデルが機能しないとき、スイッチ間トラフィックのダウン/アップターンを行ない、前記ファットツリートポロジーにおけるデッドロックを防止するために、すべてのダウン/アップターンを、前記ハブスイッチをサブツリールートノードとするデッドロックフリーの単一のサブツリーにローカライズすることをさらに含む、請求項 1 から 3 のいずれか 1 項に記載の方法。

10

【請求項 5】

前記ファットツリートポロジーにおける複数のスイッチ宛先に到達可能な 1 つ以上のリーフスイッチを検知するために、前記ファットツリートポロジーにおける複数のリーフスイッチを通して繰返すことをさらに含む、請求項 1 ~ 4 のいずれか 1 項に記載の方法。

【請求項 6】

前記ハブスイッチに前記兄弟スイッチがあるか否かを確認するとともに、宛先スイッチへの経路が存在しないときに前記兄弟スイッチに関連付けられたルーティングテーブルが前記経路を含むか否かを確認することをさらに含む、請求項 1 ~ 5 のいずれか 1 項に記載の方法。

20

【請求項 7】

前記設定することは、

前記兄弟スイッチが存在し前記兄弟スイッチに関連付けられたルーティングテーブルが前記宛先スイッチへの経路を含むとき、前記兄弟スイッチを通して前記宛先スイッチへの経路を前記ソーススイッチおよび前記サブツリールートスイッチに設定することを含む、請求項 1 ~ 6 のいずれか 1 項に記載の方法。

【請求項 8】

前記宛先スイッチに対する出力ポートを前記ハブスイッチに対する出力ポートと同一にすることをさらに含む、請求項 1 ~ 7 のいずれか 1 項に記載の方法。

【請求項 9】

スイッチ間ルーティングにおけるホップ計算のために再帰関数を用いることをさらに含む、前記再帰関数は、前記経路を構成する 1 つ以上のスイッチに対して繰返される、請求項 1 ~ 8 のいずれか 1 項に記載の方法。

30

【請求項 10】

1 つ以上のコンピュータシステムによって実行されると前記 1 つ以上のコンピュータシステムに請求項 1 ~ 9 のいずれか 1 項に記載の方法を実行させるプログラム命令を含むコンピュータプログラム。

【請求項 11】

ミドルウェアマシン環境におけるスイッチ間トラフィック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上で実行する複数のスイッチを備え、前記複数のスイッチは、

40

ルーティングアルゴリズムを用いて、前記ミドルウェアマシン環境におけるスイッチ間トラフィックのためのルーティングを実行するステップを実行するように動作し、前記ルーティングアルゴリズムは、前記ミドルウェアマシン環境におけるファットツリートポロジーにおいてソーススイッチから宛先スイッチへのスイッチ間トラフィックをルーティングするアップ/ダウンターンモデルを含み、前記ファットツリーは複数のリーフスイッチとその上方のレベルの複数のスイッチとを含み、各リーフスイッチは終端ノードに接続されて、

さらに、

前記複数のリーフスイッチから、特定のリーフスイッチを、ハブスイッチとして選択す

50

るステップと、

前記ハブスイッチの兄弟スイッチに関連付けられているルーティングテーブルが、前記宛先スイッチへの経路を含むことを判断するステップ、

前記ソーススイッチと前記ハブスイッチのそれぞれに、前記兄弟スイッチを通して、前記宛先スイッチへの経路を設定するステップと、

前記ソーススイッチと宛先スイッチと前記ハブスイッチを含む経路で、前記ハブスイッチと関連付けられたルーティングテーブルを更新するステップとを実行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】

10

【0001】

著作権に関する注意

本特許文献の開示の一部には、著作権保護の対象となるものが含まれている。著作権者は、この特許文献または特許開示の何者かによる複製が、特許商標庁の特許ファイルまたは記録にある限り、それに対して異議を唱えないが、そうでなければ、いかなる場合もすべての著作権を留保する。

【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、具体的にはミドルウェアマシン環境におけるスイッチ間通信のサポートに関する。

20

【背景技術】

【0003】

背景

ファットツリートポロジは、現代の高性能計算 (high performance computing) (HPC) クラスタ、および、インフィニバンド (Infiniband) (IB) 技術に基づいたクラスタに使用される。ファットツリーにとっては、その他多くのトポロジと同様、ネットワークリソースを効率的に利用するためにはルーティングアルゴリズムが有益である。しかしながら、既存のルーティングアルゴリズムは、スイッチ間通信においては限界がある。既存のルーティングアルゴリズムの中で、効率的なシステム管理にとって有益な、デッドロックフリーで完全接続のスイッチ間通信をサポートするものはない。これらが、本発明の実施の形態が取組もうとしている一般的な分野である。

30

【発明の概要】

【課題を解決するための手段】

【0004】

概要

本明細書に記載のシステムおよび方法は、ミドルウェアマシン環境における複数のスイッチ間のパケットのルーティングをサポートすることができる。ミドルウェアマシン環境は、インターネットプロトコル (IP) に基づく管理トラフィックを、ミドルウェアマシン環境におけるインフィニバンド上での IP (IP over Infiniband) (IPOIB) 通信を可能にすることで、サポートすることができる。複数のスイッチは、ルーティングアルゴリズムを用いて、ミドルウェアマシン環境におけるスイッチ間トラフィックのためのルーティングを実行することができる。次に、このルーティングアルゴリズムを用いて、ミドルウェアマシン環境内のスイッチを、宛先に到達できないスイッチ間トラフィックのためのハブスイッチとして選択することができる。さらに、ソーススイッチと宛先スイッチとの間にハブスイッチを介した経路が存在するとき、ハブスイッチと関連付けられたルーティングテーブルを更新することができる。

40

【0005】

本発明のある代表的な実施の形態において、ミドルウェアマシン環境におけるスイッチ間トラフィックルーティングをサポートするためのシステムが提供される。このシステムは、複数のスイッチを含み得る。これら複数のスイッチのうちの少なくとも1つは、ルー

50

ティング部と、選択部と、更新部とを含み得る。ルーティング部は、第1のルーティングアルゴリズムを用いて、ミドルウェアマシン環境におけるスイッチ間トラフィックのためのルーティングを実行し得る。選択部は、第1のルーティングアルゴリズムを用いて、ミドルウェアマシン環境内のスイッチを、宛先に到達できないスイッチ間トラフィックのためのハブスイッチとして選択し得る。更新部は、ソーススイッチと宛先スイッチとの間にハブスイッチを介した経路が存在するとき、ハブスイッチと関連付けられたルーティングテーブルを更新する機能を果たし得る。

【0006】

ある実施の形態において、複数のスイッチはファットツリートポロジを形成し得る。選択部によって、ファットツリートポロジ内のリーフスイッチを、ハブスイッチとして選択してもよく、リーフスイッチは、ファットツリートポロジ内の複数のスイッチ宛先に到達することができる。ファットツリートポロジにおいて、第1のルーティングアルゴリズムは、アップ/ダウナーンモデルを用いてスイッチ間トラフィックをルーティングしてもよい。アップ/ダウナーンモデルが機能しないとき、スイッチ間トラフィックはダウン/アップターンを行ない、ファットツリートポロジにおけるデッドロックを防止するために、すべてのダウン/アップターンは、ハブスイッチをサブツリールートノードとするデッドロックフリーの単一のサブツリーにローカライズされる。

【0007】

別の実施の形態において、上記複数のスイッチのうちの少なくとも1つのスイッチは繰返し部をさらに含み得る。繰返し部は、ファットツリートポロジにおける複数のスイッチ宛先に到達可能な1つ以上のリーフスイッチを検知するために、ファットツリートポロジにおける複数のリーフスイッチを通して繰返す。一例では、上記複数のスイッチのうちの少なくとも1つのスイッチは確認部をさらに含み得る。確認部は、ハブスイッチに兄弟スイッチがあるか否かを確認するとともに、宛先スイッチへの経路が存在しないときに上記兄弟スイッチに関連付けられたルーティングテーブルがこの経路を含むか否かを確認してもよい。一例では、上記複数のスイッチのうちの少なくとも1つのスイッチは経路設定部をさらに含み得る。経路設定部は、上記兄弟スイッチが存在しこの兄弟スイッチに関連付けられたルーティングテーブルが宛先スイッチへの経路を含むとき、上記兄弟スイッチを通して宛先スイッチへの経路をソーススイッチおよびサブツリールートスイッチに設定する。別の例では、上記複数のスイッチのうちの少なくとも1つのスイッチは、宛先スイッチに対する出力ポートがハブスイッチに対する出力ポートと同一になるように設定するためのポート設定部をさらに含み得る。

【図面の簡単な説明】

【0008】

【図1】ミドルウェアマシン環境においてデッドロックが生じるときの代表的なシナリオを説明する図である。

【図2】本発明のある実施の形態に従い、シングルコアのファットツリートポロジにおいてスイッチ間のデッドロックフリールーティングを提供することを説明する図である。

【図3】本発明のある実施の形態に従い、マルチコアのファットツリートポロジにおいてスイッチ間のデッドロックフリールーティングを提供することを説明する図である。

【図4】本発明のある実施の形態に従い、ファットツリートポロジにおいてスイッチ間のデッドロックフリールーティングを提供するための代表的なフローチャートを示す。

【図5】本発明に従うスイッチの代表的な実施の形態を示す。

【発明を実施するための形態】

【0009】

詳細な説明

本明細書に記載のシステムおよび方法は、ミドルウェアマシン環境における複数のスイッチ間のパケットのルーティングをサポートすることができる。ミドルウェアマシン環境は、インターネットプロトコル(IP)に基づく管理トラフィックを、ミドルウェアマシン環境におけるインフィニバンド上でのIP(IPoIB)通信を可能にすることで、サ

10

20

30

40

50

ポートすることができる。

【 0 0 1 0 】

インフィニバンド (I B) ネットワークおよびサブネット管理

本発明のある実施の形態に従うと、 I B ネットワークはサブネットと呼ぶことができる。 1 つのサブネットは、スイッチおよびポイント間リンクを用いて相互接続された一組のホストからなる。 I B ファブリックは 1 つ以上のサブネットを構成することができ、各サブネットはルータを用いて相互接続することができる。 1 つの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 (local identifier) (L I D) を用いてアドレス指定され、単一のサブネットは 4 9 1 5 1 個の L I D に限定される。

【 0 0 1 1 】

1 つの I B サブネットは少なくとも 1 つのサブネットマネージャ (subnet manager) (S M) を有することができる。 S M は、サブネット内の、スイッチ、ルータ、およびホストチャネルアダプタ (host channel adapter) (H C A) 上にあるすべての I B ポートの構成を含むネットワークを、初期化し起動する役割を果たす。初期化時、 S M は検出状態で始動する。検出状態において、 S M は、すべてのスイッチおよびホストを検出するためにネットワークのスweepを行なう。検出状態中に、 S M は、他の S M を検出することができ、どの S M がマスタ S M になるべきかについて交渉することができる。検出状態が終了すると、 S M はマスタ状態になる。マスタ状態において、 S M は、 L I D の割当、スイッチの構成、ルーティングテーブルの計算および導入、ならびにポートの構成を進める。マスタ状態が終了した時点で、サブネットは立ち上げられて使用できる状態となり、 S M は、サブネットが構成された後の変化についてネットワークをモニタリングする役割を果たす。

【 0 0 1 2 】

加えて、 S M は、完全接続性、デッドロックフリー状態、およびソースと宛先のすべての対の間での適切な負荷分散を維持するルーティングテーブルを計算する役割を果たすことができる。ルーティングテーブルはネットワークの初期化時に計算することができ、このプロセスを、トポロジが変化する度に、ルーティングテーブルを更新し最適性能を保証するために繰返すことができる。

【 0 0 1 3 】

通常動作中、 S M は、リンクが機能停止したとき、デバイスが追加されたとき、またはリンクが削除されたとき等の、トポロジ変化事象について確認するために、定期的にネットワークを軽くスweepする。 S M が、軽くスweepを行なっている間に変化事象を検出した場合またはネットワークの変化を知らせるメッセージ (トラップ) を受けた場合、 S M は、検出された変化に従ってネットワークを再構成することができる。この再構成も、初期化中に用いられるステップを含むことができる。

【 0 0 1 4 】

I B は、フロー制御をバーチャルレーン (virtual lane) (V L) 毎に行なうことができる、ロスのないネットワーキング技術である。 V L は、別々のバッファリング、フロー制御、および輻輳管理リソースを備えた、同一の物理リンク上の論理チャネルである。 V L の概念によって、物理トポロジの上にバーチャルネットワークを構築することが可能になる。これらバーチャルネットワークまたは層を、効率的なルーティング、デッドロックの回避、フォールトトレランス、およびサービス差別化等のさまざまな目的に使用することができる。

【 0 0 1 5 】

ファットツリールーティング

本発明のある実施の形態に従うと、ファットツリートポロジは階層ネットワークトポロジであり、たとえば、平衡のとれたファットツリーでは、すべての層のリンク容量を等しくすることができる。さらに、ファットツリートポロジは、複数のルートを持つ 1 つのツリーを構築することによって実現でき、たとえば、 X G F T 表記を用いて記載できる、 m ポート n ツリー定義または k 分木 n ツリー定義によって実現できる。

【 0 0 1 6 】

ファットツリールーティングアルゴリズムは、利用可能なネットワークリソースを活用することができる。ファットツリールーティングアルゴリズムは、2つの段階、すなわちパケットをソースから送る上昇段階と、パケットを宛先に向けて送るときの下降段階とを含むことができる。これら2つの段階の間の移行期は、最下部にある共通の祖先で発生する。この祖先は、ソースおよび宛先双方に、その下向きのポートを通して到達できるスイッチである。このようなルーティングの実装はデッドロックフリーを保証し、この実装はまた、同じ宛先に向かうすべての経路が同一のルート（トップ）ノードで収束することで、この宛先に向かうすべてのパケットが下向きの1つの専用経路を辿ることを保証する。すべての宛先に対して専用下降経路を設けることによって、下降段階における競合を効果的になくす（上昇段に移す）ことができ、そのため、宛先が異なるパケットは、その経路上のスイッチのうちの2分の1における出力ポートについてしか競合しない。加えて、オーバーサブスクリプション状態のファットツリーにおける下降経路は、専用経路ではなく、数個の宛先が共有できるものである。

10

【 0 0 1 7 】

スイッチ間通信

ファットツリートポロジーでは、終端ノード間の通信に加えて、スイッチ間通信が可能である。スイッチ間通信から生じるトラフィックがほとんどないとき、スイッチ間通信は無視してもよい。言い換えると、IBスイッチは、ルーティングアルゴリズムの観点から安全に無視できる透明なデバイスとみなし得る。

20

【 0 0 1 8 】

一方、スイッチ間通信トラフィックは常に無視できる訳ではない。こういった場合、先のセクションで説明したファットツリールーティングを用いて、制限のあるスイッチ間通信のルーティングを行なうことができる。なぜなら、いずれか2つのスイッチの間の経路を検知して利用できる最初のポートからトラフィックを送るためには、最下部にある共通の先祖を検知する必要があるからである。さらに、ファットツリールーティング方式は、最下部にある共通の先祖を持たないスイッチの間は接続しない、すなわち、先のセクションで説明したファットツリールーティング方式は、ファットツリートポロジーにおいて完全接続を構築しないことがある。

【 0 0 1 9 】

終端ノードおよびスイッチ双方を含むファットツリートポロジーにおける完全接続は、さまざまな理由から有益である可能性がある。たとえば、IB診断ツールは、LIDルーティングに依拠しており、基本的なファブリック管理およびモニタリングのために完全接続を必要とする。さらに、ベンチマーキングに使用されるperf test等の進化したツールは、基礎をなすLIDルーティングに依拠し完全接続を必要とするIPOIBを採用することがある。IPOIBは、IBネットワーク上でTCP/IPトラフィックを実行できるようにするカプセル化方法である。加えて、IPOIBは、非インフィニバンドを意識した管理およびモニタリングプロトコルならびにSNMPまたはSSHのようなアプリケーションが必要とする場合がある。このように、ファブリック内のすべてのスイッチ間の完全接続は有益であり得る。なぜなら、スイッチが、任意のトラフィックを生成することができ、他の終端ノードのようにアクセスされることができ、埋込まれたSMを含むことが多いからである。

30

40

【 0 0 2 0 】

ファットツリートポロジーにおいて、デッドロックは、バッファまたはチャネル等のネットワークリソースが共有されていることが原因で生じ得る。

【 0 0 2 1 】

図1は、ミドルウェアマシン環境においてデッドロックが生じるときの代表的なシナリオを説明する図である。図1に示されるように、ミドルウェアマシン環境におけるファットツリートポロジー100では、スイッチ間通信の対が2つあり（スイッチ0～スイッチ3およびスイッチ3～スイッチ6）ノード間通信の対が2つある（ノードB～ノードDお

50

よびノードD～ノードA)、混合通信パターンが可能である。ここで、スイッチ0からスイッチ3へのトラフィックはスイッチ4、8、および5を介して導かれ、スイッチ3からスイッチ6へのトラフィックはスイッチ7および9を介して導かれる。また、ノードBからノードDへのトラフィックはスイッチ8、5、3、7、および9を介して導かれ、ノードDからノードAへのトラフィックはスイッチ9、6、0、4、および8を介して導かれる。

【0022】

図1の点線で示されるように、上記混合通信パターンがファットツリートポロジ－100内に存在するとき、循環クレジット従属性またはクレジットループ101が作成される。さらに、デッドロックは、循環クレジット従属性が作成されたときに発生し得る。これは、循環クレジット従属性が存在するときには常にデッドロックがあるということを意味しているのではない。言い換えると、循環クレジット従属性を作成することで、デッドロックが発生する可能性が生じる。たとえば、minhop等のルーティングアルゴリズムを使用するとき、デッドロックは、ノード間およびスイッチ間トラフィックがある3段のファットツリーにおいて生じ得る。その理由として、minhopが5ノード以上のリングをデッドロックフリーで解決することができないことが挙げられ、3段のファットツリーは6ノードのリングを含み得る。このように、すべてのスイッチの間でスイッチ間通信が行なわれるときにファットツリートポロジ－においてデッドロックフリールーティングを提供できるアルゴリズムが必要である。

【0023】

デッドロックは、利用可能な物理リソースをセグメントに分けるVLを用いることによって回避できる。しかしながら、ファットツリーにおいてデッドロックを回避するためにVLを使用することは、効率的でない可能性がある。なぜなら、サービス品質(Quality of Service)(QoS)のようなその他の目的のために追加のVLを使用することがあるからである。

【0024】

たとえば、IBのためのLASHアルゴリズムは、完全接続とデッドロックフリーをサポートできる。LASHアルゴリズムは、VLを用いることによってクレジットループ101を切断し、ファブリック内のすべてのノード間の完全接続を提供する。LASHは、経路を割当てるときにファットツリーの特性を活用しないので、通常のファットツリールーティングよりもネットワーク性能が悪く、加えて、LASHはルーティング時間が長い。DFSSSPは、ファットツリーにおける全対全のスイッチ間通信をサポートするのに使用し得るもう1つのルーティングプロトコルである。しかしながら、DFSSSPは非HCAノード間のクレジットループを切断しない。このことは、スイッチ間通信がデッドロックフリーではない場合があることを意味する。

【0025】

このように、相互接続ネットワーク内でデッドロックが生じると、このデッドロックは、ネットワークの一部またはすべてにおける通信を妨げる可能性がある。したがって、相互接続ネットワークのシステム全体の観点から、ファットツリートポロジ－においてサポートされる完全接続は、デッドロックフリーであることが好ましい。

【0026】

シングルコアのファットツリートポロジ－におけるデッドロックフリールーティング
本発明のある実施の形態に従い、ミドルウェアマシン環境におけるノード間の完全接続をサポートすることができ、1つのVLを用いるだけでファブリック全体をデッドロックフリーにすることができる。

【0027】

図2は、本発明のある実施の形態に従い、シングルコアのファットツリートポロジ－においてスイッチ間のデッドロックフリールーティングを提供することを説明する図である。図2に示されるように、ミドルウェアマシン環境におけるシングルコアのファットツリートポロジ－200は、複数のスイッチを含むことができる。さらに、ミドルウェアマシ

ン環境は、ファットツリートポロジ－２００内のリーフスイッチ２０２に接続するサーバノード等の複数の終端ノード２２０を含むことができる。

【００２８】

アップ/ダウンターンモデルを用いることにより、ミドルウェアマシン環境におけるファットツリールーティングを提供することができる。たとえば、スイッチＡ１およびＡ２はＸ１またはＸ２を経由してアップすることによって接続性を得ることができる。一方、アップ/ダウンターンモデルは、最下部の共通の先祖がないスイッチたとえばＡ１とＢ１との間の通信には有用でないことがある。

【００２９】

本発明のある実施の形態に従い、代替の手法を採用することによって、アップ/ダウンターンモデルを用いて通信を確立することができないスイッチ間の接続性を得ることができる。図２に示されるように、上下逆のツリー、すなわちサブツリー２１０（影を付けた領域に示される）を、ファットツリートポロジ－２００内に形成することができる。サブツリー２１０は、ファットツリートポロジ－２００内のリーフスイッチであるサブツリールートノードＳ_nすなわちハブスイッチ２０１を有することができる。完全接続は、ファットツリートポロジ－２００内のハブスイッチ２０１を通して確立することができる。

【００３０】

サブツリー２１０内において、スイッチＸ１とＹ１はサブツリールートＳ_n２０１を経由して通信することができ、Ｘ１とＸ２はスイッチＡ１を経由して通信できる。さらに、サブツリー２１０の外にあるスイッチは、ファットツリートポロジ－２００内の他のスイッチと、サブツリー２１０を通して通信することができる。これは、最短でない（または最適に及ばない）経路の使用につながる。ファブリック内の、アップ/ダウンターンモデルを用いて別のスイッチに到達することが不可能なスイッチは、先ずパケットをサブツリー２１０に送ればよい。次に、これらのパケットを、宛先スイッチへの経路を有するスイッチに到達するまで、サブツリー２１０内で送ればよい。言い換えると、これらのパケットは、ハブスイッチ２０１でＵターンする、すなわち、下降し折返して上昇し続いて上昇してから宛先に向かって下降する。ミドルウェアマシン環境におけるサブツリー２１０は、ファブリック内のすべてのダウン/アップターンを、１つのデッドロックフリーツリーにローカライズすることができる。

【００３１】

たとえば、ミドルウェアマシン環境における第１のスイッチたとえばＡ１から第２のスイッチたとえばＢ２まで１つ以上のパケットをルーティングすることを要求されたとき、システムは先ず、サブツリールートスイッチすなわちハブスイッチ２０１に関連付けられたルーティングテーブルが第２のスイッチＢ２への経路を含むか否か判断することができる。図２において点線で示されるように、ハブスイッチ２０１は、Ｂ１、Ｙ１またはＹ２を通る経路を介して第２のスイッチＢ２に接続できる。次に、第２のスイッチＢ２への経路を、第１のスイッチＡ１に関連付けられたルーティングテーブルに挿入することができる。加えて、Ａ１に関連付けられたルーティングテーブルにおいて、スイッチＡ２に対する出力ポートをハブスイッチ２０１に対する出力ポートと同一になるように設定することができる。

【００３２】

上下逆のサブツリー２１０のルートは、リーフスイッチのうちのいずれでもよい。さらに、ミドルウェアマシン環境におけるデッドロックフリーの完全接続を得るために、ファットツリートポロジ－２００内のすべてのスイッチがともに１つのリーフスイッチを通信用を選択することができる。

【００３３】

本発明のある実施の形態に従い、すべてのＵターンがサブツリー２１０内でのみ行なわれ上昇（アップ）から下降（ダウン）へのいずれかのターンによってトラフィックがサブツリー２１０から出るようにすることによって、ファットツリートポロジ－２００内をデッドロックフリーにすることができる。また、サブツリー２１０の外部ではＵターンがで

10

20

30

40

50

きないので、サブツリー 2 1 0 から出たトラフィックは循環してサブツリー 2 1 0 の中に
戻ることができないかもしれない。

【 0 0 3 4 】

このように、シングルコアのファットツリー 2 0 0 内では、ソーススイッチをハブス
witch 2 0 1 に接続しハブスイッチ 2 0 1 をすべての宛先スイッチに接続する経路が常に存
在し得る。さらに、このトポロジー内のリーフスイッチを通して接続性を得るための通信
方式は、ファットツリールーティングを用いて互いに到達することができるスイッチ間の
上方向から下方向への経路を変更しない。

【 0 0 3 5 】

添付の付録 A は、ミドルウェアマシン環境におけるスイッチ間のルーティングのサポー
トおよび本開示を通して記載されているプラットフォームのさまざまな側面に関するさら
に他の情報を提供する。付録 A における情報は、例示を目的として提供されており、本発
明の実施の形態すべてを限定するものであると解釈されてはならない。

【 0 0 3 6 】

マルチコアのファットツリートポロジーにおけるデッドロックフリールーティング

本発明のある実施の形態に従うと、複雑なマルチコアのファットツリーまたは不規則な
ファットツリーでは、2つのスイッチ間に常に接続性があるとは限らないであろう。こう
いった場合、ベストエフォートの手法を用いることによってスイッチ間通信をルーティン
グすることができる。

【 0 0 3 7 】

図 3 は、本発明のある実施の形態に従い、マルチコアのファットツリートポロジーにお
いてスイッチ間のデッドロックフリールーティングを提供することを説明する図である。
図 3 に示されるように、ミドルウェアマシン環境におけるマルチコアのファットツリー
ポロジ 3 0 0 は、複数のスイッチを含むことができる。さらに、このミドルウェアマシ
ン環境は、ファットツリートポロジ 3 0 0 内のリーフスイッチに接続するサーバノード
等の複数の終端ノード 3 2 0 を含むことができる。

【 0 0 3 8 】

同様に、上下逆のツリーすなわちサブツリー 3 1 0 (影を付けた領域に示される)をフ
ァットツリートポロジ 3 0 0 内に形成することができる。サブツリー 3 1 0 は、ファッ
トツリートポロジ 3 0 0 内のリーフスイッチであるサブツリールートノード S n すなわ
ちハブスイッチ 3 0 1 を有することができる。図 2 とは異なり、このサブツリールートス
イッチすなわちハブスイッチ 3 0 1 はすべての宛先への経路を有さないかもしれない。

【 0 0 3 9 】

加えて、このミドルウェアマシン環境には兄弟スイッチ 3 0 2 が存在し得る。この兄弟
スイッチは、水平リンク 3 3 0 を通してサブツリールートスイッチ S n すなわちハブス
イッチ 3 0 1 に接続できる。

【 0 0 4 0 】

一例において、第 1 のスイッチたとえば B 2 は、ミドルウェアマシン環境において 1 つ
以上のパケットを第 2 のスイッチたとえば B 3 にルーティングすることを要求できる。サ
ブツリールートスイッチ 3 0 1 には、B 3 に接続する経路がない。なぜなら、B 3 は、他
のノードと水平リンクのみを介して接続する別のファットツリーの中にあるからである。

【 0 0 4 1 】

システムはまず、サブツリールートスイッチ 3 0 1 が、そのルーティングテーブルが B
3 への経路を含んでいる兄弟スイッチ 3 0 2 を有するか否か確認することができる。シス
テムは次に、上記兄弟スイッチ 3 0 2 を通して、第 2 のスイッチ B 3 への経路を、第 1 の
スイッチ B 2 およびサブツリールートスイッチ 3 0 1 に設定することができる。結果とし
て、このシステムを、スイッチ B 2 からのパケットをスイッチ B 3 までたとえば Y 1 また
は Y 2、B 1、サブツリールートノード 3 0 1、およびその兄弟ノード 3 0 2 を介してル
ーティングする(点線で示される)ように構成できる。

【 0 0 4 2 】

10

20

30

40

50

S F T R E E アルゴリズム

本発明のある実施の形態に従い、ルーティングアルゴリズム、すなわち S F T R E E アルゴリズムを用いて、スイッチ間のデッドロックフリーの完全接続を得ることができる。S F T R E E アルゴリズムは、2つのステップ、すなわち、リーフスイッチのうちの1つであってもよいサブツリールート S_n を検知する第1のステップと、スイッチ間ルーティングを実行する第2のステップとを含むことができる。

【0043】

以下のアルゴリズム1は、サブツリールートを検知するための疑似コードを含み、サブツリーを確立する。

【0044】

10

【表1】

アルゴリズム1 サブツリールート関数を選択

Require: ルーティングテーブルが生成されている。

Ensure: サブツリールート(sw_{root})が選択される。

```

1: found = false
2: for  $sw_{leaf} = 0$  to  $max\_leaf\_sw$  do
3:   if found == true then
4:     break
5:   end if
6:    $sw_{root} = sw_{leaf}$ 
7:   found = true
8:   for  $dst = 1$  to  $max\_dst\_addr$  do
9:     if  $sw_{root}.routing\_table[dst] == no\_path$  then
10:      found = false
11:      break
12:    end if
13:  end for
14: end for
15: if found = false then
16:    $sw_{root} = get\_leaf(0)$ 
17: end if

```

20

30

【0045】

アルゴリズム1は、トポロジ内のすべてのリーフスイッチをトラバースし、各リーフスイッチについて、そのルーティングテーブルにおいてスイッチ宛先アドレスのうちのいずれかが到達不能と記されているか否か確認する。ルーティングテーブル内のすべてのスイッチ宛先が到達可能と記されている場合は、最初に遭遇したリーフスイッチをサブツリールートスイッチとして選択する。さらに、ファットツリーにおける潜在的なサブツリー数は、完全接続のリーフスイッチの数と同一である可能性がある。これらリーフスイッチのうちの1つのみをサブツリールートとして選択すればよく、この特定のリーフスイッチにおいてルートがあるサブツリーは1つのみであってもよい。

40

【0046】

以下のアルゴリズム2は、ファットツリートポロジにおけるスイッチ間ルーティングを実行するための疑似コードを含む。

【0047】

【表 2】

アルゴリズム2 スイッチ間ルーティング関数

Require: サブツリールート(sw_{root})**Ensure:** 最悪の場合でも sw_{root} によって各 sw_{src} が各 sw_{dst} に到達する。

```

1: if found or  $sw_{root}.routing\_table[dst] \neq no\_path$  then
2:    $port = sw_{src}.routing\_table[sw_{root}.addr]$ 
3:    $sw_{src}.routing\_table[dst] = port$ 
4:    $get\_path\_length(sw_{src}, null, dst, sw_{root}, hops)$ 
5:    $set\_hops(sw_{src}, hops)$ 
6:   return true
7: else if ( $sw_{sib} = get\_sibling\_sw(sw_{root}) \neq null$ ) then
8:   if  $sw_{src}.routing\_table[sw_{sib}.addr] \neq no\_path$  then
9:     if  $sw_{sib}.routing\_table[dst] \neq no\_path$  then
10:       $port = get\_port\_to\_sibling(sw_{sib})$ 
11:       $sw_{root}.routing\_table[dst] = port$ 
12:       $hops = get\_hops(sw_{sib}, dst)$ 
13:       $set\_hops(sw_{root}, hops + 1)$ 
14:       $hops = 0$ 
15:       $port = sw_{src}.routing\_table[sw_{sib}.addr]$ 
16:       $sw_{src}.routing\_table[dst] = port$ 
17:       $get\_path\_length(sw_{src}, null, dst, sw_{root}, hops)$ 
18:       $set\_hops(sw_{src}, hops)$ 
19:      return true
20:   end if
21: end if
22: end if
23: print 'SW2SW failed for  $sw_{src}$  and  $sw_{dst}$ .'
24: return false

```

10

20

【0048】

スイッチ間ルーティング関数がコールされたとき、第1のステップは、アルゴリズム2のライン1に示されているように、サブツリールート(sw_{root})のルーティングテーブルが宛先スイッチ(sw_{dst})への経路を含むか否か判断することである。含む場合、宛先スイッチへの経路をソーススイッチ(sw_{src})のルーティングテーブルに挿入し、宛先スイッチに対する出力ポートは、サブツリールートへの経路に対する出力ポートと同一である(ライン2~6)。すべてのスイッチに対してスイッチ間ルーティング関数がコールされるので、最終的に、到達不能な各宛先スイッチに対するすべての経路はサブツリーに収束する。言い換えると、選択されたサブツリールートは、到達不能な宛先スイッチすべてにとって新たなターゲットである。ダウン/アップターンは、ソーススイッチおよび宛先スイッチ双方への上向きの経路を有するスイッチであってもよい、サブツリー内に位置する可能な第1のスイッチで、生じるであろう。

30

40

【0049】

アルゴリズム2の第2の部分(ライン7~24)は、ベストエフォートの手法が用いられサブツリールートがすべての宛先への経路を有していない複雑なマルチコアまたは不規則なファットツリーの場合に対応する。一方、このアルゴリズム2の第2の部分は、シングルコアのファットツリーではスキップしてもよい。なぜなら、ソーススイッチからサブツリールートへの経路は常に存在し得るものでありサブツリールートはすべての宛先スイッチへの経路を有し得るからである。

【0050】

アルゴリズム2は、サブツリールートに、このサブツリールートが図3に示されるように水平リンクを通して接続される真隣の「隣人」があるか否か確認する(ライン7)。次

50

に、アルゴリズム 2 は、兄弟とも呼ばれるこの隣人が、宛先スイッチへの経路を有するかどうかを確認する（ライン 8）。兄弟が存在しこの兄弟が宛先スイッチへの経路を有していれば、兄弟スイッチを通して、宛先スイッチへの経路を、サブツリールートおよび起点となるソーススイッチ双方に設定する（ライン 9 ~ 19）。言い換えると、到達不能な宛先スイッチにとってのターゲットは依然としてサブツリールートであってもよく、サブツリールートはパケットをその兄弟スイッチを介して宛先に送ることができる。

【 0 0 5 1 】

上記 2 つのステップがいずれも真値とともに関数からリターンしない場合、S F T R E E アルゴリズムは 2 つのスイッチ間の経路を検知していないかもしれない。たとえば、この事態が生じ得るのは、リンク障害が複数あるためまたはノード間の接続が非常に不規則であるためにファットツリールーティングを実行することが推奨されないトポロジの場合である。このようなトポロジの一例は、大きさが異なる 2 つのファットツリーが非対称に相互接続されているトポロジである可能性があり、たとえば大きい方のツリー上の 2、3 の中間段スイッチが小さい方のツリー上のルートに接続されている。O p e n S M におけるさまざまなコマンドラインパラメータを用いて、ファットツリールーティングをこういったトポロジ上で強制的に実行することができる。S F T R E E アルゴリズムは、性能だけでなく接続に関しても最適に及ばないルーティングを与えるかもしれない。

【 0 0 5 2 】

本発明のある実施の形態に従い、O p e n S M は、宛先へのすべての経路にホップカウントを記してもよい。ファットツリールーティングアルゴリズムは、ツリー内のスイッチランキングを用いてホップの数を計算することができる。また、カウントは、主ルーティング関数の単純なカウンタを用いて行なうことができる。これらのカウント方式は、スイッチ間ルーティングが、ファットツリーアルゴリズムによって実施される主ルーティングから完全に独立しているときは、ジグザグ経路（すなわち経路が最短ホップ経路に従っていない）が存在し得るので、信頼できないことがある。

【 0 0 5 3 】

スイッチ間ルーティングの間、ホップ計算、たとえば `get_path_length`（ライン 17）に対して再帰関数を用いることができる。この関数は、経路を構成する一連のスイッチに対して繰返して用いられることができ、この関数は、宛先への適切なホップカウントを有するノードに達したとき、ホップカウントを、経路上のスイッチのルーティングテーブルに書込むことができる。

【 0 0 5 4 】

加えて、サブツリールートを、サブネットマネージャノード（エンドノードであってその上でサブネットマネージャが実行している）がサブツリールートと記されたスイッチに接続されないように、選択することができる。このようにスイッチ間トラフィックをサブネットマネージャから切離すことによって、重要な場所でボトルネックが生じないようにすることができる。

【 0 0 5 5 】

図 4 は、本発明のある実施の形態に従い、ファットツリートポロジにおいてスイッチ間のデッドロックフリールーティングを提供するための代表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、ミドルウェアマシン環境は、ルーティングアルゴリズムを用いて、ミドルウェアマシン環境におけるスイッチ間トラフィックのためのルーティングを実行することができる。次に、ステップ 402 において、ミドルウェアマシン環境は、ルーティングアルゴリズムを用いて、ミドルウェアマシン環境内のあるスイッチを、宛先に到達できないスイッチ間トラフィックのためのハブスイッチとして選択することができる。さらに、ステップ 403 で、ミドルウェアマシン環境は、ソーススイッチと宛先スイッチとの間にハブスイッチを介した経路が存在するとき、ハブスイッチと関連付けられたルーティングテーブルを更新することができる。

【 0 0 5 6 】

図 5 は本発明に従うスイッチの代表的な実施の形態を示す。このスイッチは、ミドルウ

10

20

30

40

50

エアマシン環境におけるスイッチ間トラフィックルーティングをサポートし複数のスイッチを含む上記システムのうちいずれかに含まれるものであってもよい。図5を参照して、代表的なスイッチ10は、上記のように本発明の原理に従って構成された、ルーティング部11と、選択部13と、更新部14とを含み得る。部11、13、および14を含むがこれらに限定されないスイッチ10の機能部は、本発明の原理を実施するハードウェア、ソフトウェア、またはハードウェアとソフトウェアの組合せによって実現し得る。図5に示される各機能ブロックを組合わせてまたはサブユニットに分割して上記本発明の原理を実現し得ることが、当業者に理解される。したがって、本明細書における説明は、本明細書に記載の機能ブロックの可能な組合せまたは分割またはさらに他の定義をサポートし得る。

10

【0057】

ルーティング部11は、先に開示されているように、第1のルーティングアルゴリズムを用いて、ミドルウェアマシン環境におけるスイッチ間トラフィックのルーティングを実行してもよい。選択部13は、第1のルーティングアルゴリズムを用いて、ミドルウェアマシン環境内のスイッチを、宛先に到達できないスイッチ間トラフィックのためのハブスイッチとして選択してもよい。ソーススイッチと宛先スイッチとの間にハブスイッチを介した経路が存在するとき、更新部14はこのハブスイッチと関連付けられたルーティングテーブルを更新してもよい。

【0058】

ある代表的な実施の形態において、システムの複数のスイッチはファットツリートポロジを形成する。選択部13がファットツリートポロジ内のリーフスイッチをハブスイッチとして選択してもよく、このリーフスイッチはファットツリートポロジ内の複数のスイッチ宛先に到達できる。第1のルーティングアルゴリズムは、アップ/ダウンターンモデルを用いることによってファットツリートポロジ内のスイッチ間トラフィックをルーティングしてもよい。スイッチ間トラフィックは、アップ/ダウンターンモデルが機能しないときは、ダウン/アップターンを行ない、すべてのダウン/アップターンは、ファットツリートポロジにおけるデッドロックを防止するために、ハブスイッチをサブツリールートノードとする単一のデッドロックフリーサブツリーに、ローカライズされる。

20

【0059】

ある代表的な実施の形態において、スイッチ10は任意で、ファットツリートポロジにおける複数のスイッチ宛先に到達できる1つ以上のリーフスイッチを検知するために、ファットツリートポロジ内の複数のリーフスイッチを通して繰返すための繰返し部12を含んでもよい。したがって、選択部13は、この1つ以上のリーフスイッチからハブスイッチを選択してもよい。

30

【0060】

ある代表的な実施の形態において、スイッチ10は確認部15をさらに含んでもよい。確認部15は、ハブスイッチが兄弟スイッチを有しているか否かを確認するとともに、宛先スイッチへの経路が存在しないとき、兄弟スイッチと関連付けられたルーティングテーブルがこの宛先スイッチへの経路を含むか否かを確認してもよい。

【0061】

ある代表的な実施の形態において、スイッチ10は経路設定部16をさらに含んでもよい。兄弟スイッチが存在しこの兄弟スイッチと関連付けられたルーティングテーブルが宛先スイッチへの経路を含むとき、経路設定部16は、兄弟スイッチを通して、宛先スイッチへの経路を、ソーススイッチおよびサブツリールートスイッチに設定してもよい。

40

【0062】

ある代表的な実施の形態において、スイッチ10は、宛先スイッチに対する出力ポートがハブスイッチに対する出力ポートと同一になるように設定するためのポート設定部17をさらに含んでもよい。

【0063】

本発明は、1つ以上のプロセッサ、メモリ、および/または本開示の教示に従いプログ

50

ラムされたコンピュータ読取可能な記憶媒体を含む、従来の汎用または専用デジタルコンピュータ、コンピューティングデバイス、マシン、またはマイクロプロセッサを1つ以上用いて、適宜実現し得る。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が本開示の教示に基づいて容易に準備できる。これはソフトウェア技術の当業者には明らかであろう。

【0064】

実施の形態によっては、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憶媒体またはコンピュータ読取可能な媒体（複数の媒体）であるコンピュータプログラムプロダクトを含む。この記憶媒体は、フロッピーディスク（登録商標）、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。

【0065】

本発明に関する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。変形は、開示された特徴および/または要素の任意の組合せを含み得る。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当業者が本発明の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。本発明の範囲は、以下の特許請求の範囲およびその均等物によって定められることが意図されている。

【0066】

付録A

以下は、SFTREEアルゴリズムがデッドロックフリーであることを示す証拠である。この証拠には関連用語の代表的な定義も含まれる。

【0067】

定義1．スイッチ S_n は、ランク n のスイッチを意味する。ルートスイッチのランクは $n = 0$ である。

【0068】

定義2．下方向のチャンネルは、トラフィックが s_{n-1} から s_n に流れる s_{n-1} と s_n の間のリンクを意味する。上方向のチャンネルは、トラフィックが s_n から s_{n-1} に流れる s_n と s_{n-1} の間のリンクを意味する。

【0069】

定義3．経路は、チャンネルによって接続された一連のスイッチであると定義される。経路はソーススイッチで始まり宛先スイッチで終わる。

【0070】

定義4．Uターンは、下方向のチャンネルの次に上方向のチャンネルが続くターンである。

定義5．チャンネル c_i とチャンネル c_j との間のチャンネル従属性は、パケットを保持しているチャンネル c_i がチャンネル c_j の使用を要求するときに生じる。

【0071】

定義6．チャンネル従属性グラフ $G = (V, E)$ は、頂点 V がネットワーク N のチャンネルでありエッジ E がチャンネル対 (c_i, c_j) であり c_i から c_j への（チャンネル）従属性が存在する有向グラフである。

【0072】

定義7．サブツリーは、サブツリーの単一のルートである1つのリーフスイッチ s に収束するファットツリー内の上下逆の論理ツリー構造を意味する。サブツリーは、そのルートから拡張しそのリーフはすべてファットツリー内のトップレベルスイッチである。上方向から下方向へのターンによってサブツリー構造の外に出る。

【 0 0 7 3 】

以下の定理はルーティング関数のデッドロックフリーにとって十分な条件を示す。

定理 1 . ネットワーク N に対するルーティング確定関数 R は、チャンネル従属性グラフ G に循環がない場合に限り、デッドロックフリーである。

【 0 0 7 4 】

次に補助定理を示すことができる。この補助定理の後に S F T R E E アルゴリズムのデッドロックフリー定理が続く。

【 0 0 7 5 】

補助定理 8 ファットツリー内に少なくとも 1 つのサブツリーが存在する。これにより、このサブツリー内の U ターンのみを用いてスイッチ間の完全接続を構築できる。

10

【 0 0 7 6 】

証拠 . ファットツリーにおいて、アップ / ダウン経路または水平兄弟経路を用いて、ファブリック内の任意のリーフから任意の他のノード (すべてのスイッチを含む) に到達することができる。その理由は、すべての終端ノードがその他すべての終端ノードと通信できること、および終端ノードが直接リーフスイッチに接続されていることにある。すなわち、トポロジー内においてリーフスイッチはその他任意のノードに到達することができる。結果として、リンク障害が含まれていなければ、いずれのリーフスイッチもサブツリーのルートとして機能できる。

【 0 0 7 7 】

補助定理 9 . デッドロックを生じさせずに 1 つのサブツリー内で無制限の数の U ターンが可能である。

20

【 0 0 7 8 】

証拠 . サブツリーは論理ツリー構造でありその中で 2 種類のターンが可能である。これらのターンは、U ターンすなわち下方向から上方向へのターンと、通常の上方向から下方向へのターンである。U ターンは、サブツリー (およびファットツリー) 内の最上部のスイッチでは不可能である。なぜなら、これらのスイッチには上向きの出力ポートがないからである。定義 7 に従い、上方向から下方向への任意のターンによってサブツリーの外に出る。

【 0 0 7 9 】

サブツリーは循環がない連結グラフであるため、それ自体がデッドロックフリーである。いかなるデッドロックにもサブツリー外部の U ターンが関わっているはずである。その理由は、U ターンの次には上方向から下方向へのターンが続いてサブツリーから出ることにある。サブツリーから始まりアップ / ダウンターンを通して進むすべてのトラフィックは、ターンが行なわれるとサブツリーから外に出て、宛先への下方向のチャンネルのみを通る。このような従属性は、再びサブツリーに入ることはできず、その他の U ターンに決して到達できず、したがって循環を形成できない。

30

【 0 0 8 0 】

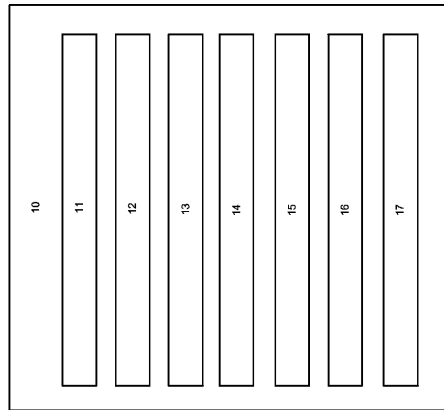
定理 2 . サブツリー法を用いた S F T R E E アルゴリズムはデッドロックフリーである。

【 0 0 8 1 】

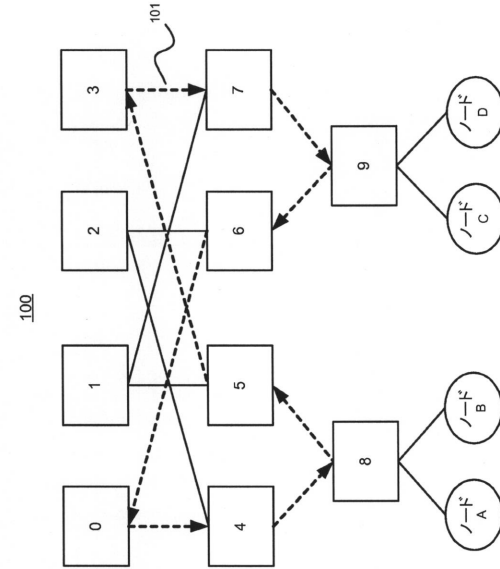
40

証拠 . 補助定理 8 および補助定理 9 に従うと、ファットツリー内のデッドロックは、サブツリーの外部で U ターンがあった場合にしか生じない。このような状況が生じ得ない理由は、すべての U ターンがサブツリー内で生じる S F T R E E アルゴリズムの設計にある。サブツリーから出たトラフィックが循環してサブツリーに戻ることはない。なぜなら、トラフィックは、一旦サブツリーから出ると、下向きのチャンネルのみを通して宛先に送られ、よって、トポロジー内では循環クレジット従属性が発生しないと考えられるからである。

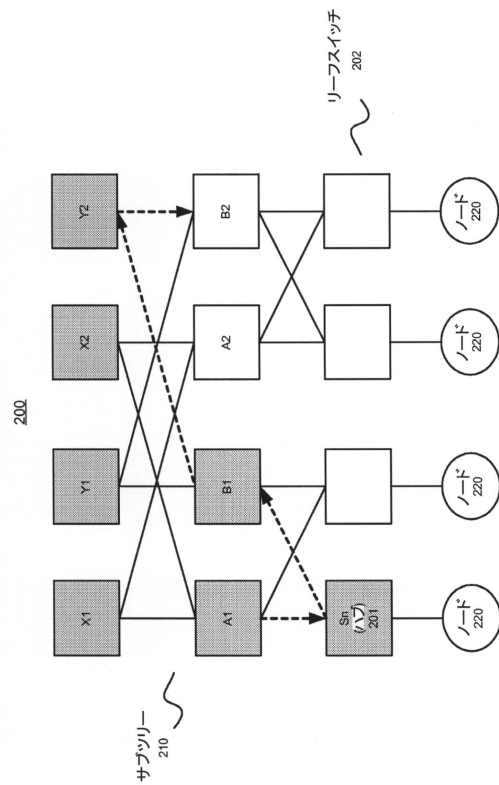
【図 5】



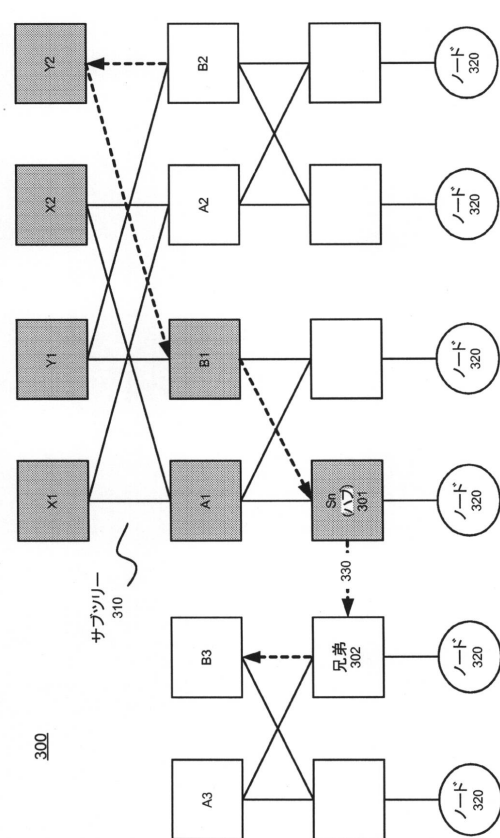
【図 1】



【図 2】



【図 3】



【図 4】

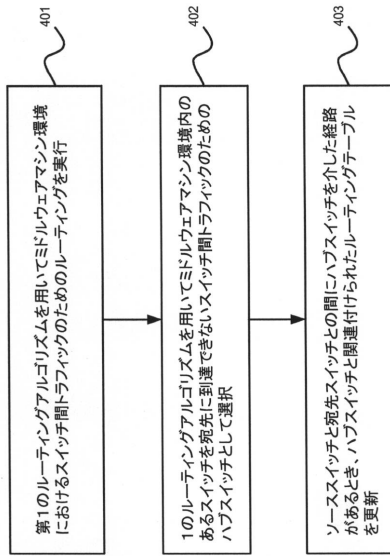


FIGURE 4

フロントページの続き

(56)参考文献 特開2011-223419(JP, A)

Frank Olaf Sem-Jacobsen 他, Dynamic Fault Tolerance in Fat Trees, IEEE TRANSACTIONS ON COMPUTERS, IEEE Computer Society, 2011年 4月, Vol.60, No.4, p.508-525

(58)調査した分野(Int.Cl., DB名)

H04L 12/753