

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5860670号
(P5860670)

(45) 発行日 平成28年2月16日(2016.2.16)

(24) 登録日 平成27年12月25日(2015.12.25)

(51) Int.Cl. F I
H04L 12/715 (2013.01) H04L 12/715

請求項の数 10 外国語出願 (全 27 頁)

(21) 出願番号	特願2011-241239 (P2011-241239)	(73) 特許権者	593096712 インテル コーポレーション
(22) 出願日	平成23年11月2日(2011.11.2)		アメリカ合衆国 95054 カリフォル ニア州 サンタ クララ ミッション カ レッジ ブールバード 2200
(65) 公開番号	特開2012-105265 (P2012-105265A)	(74) 代理人	100068755 弁理士 恩田 博宣
(43) 公開日	平成24年5月31日(2012.5.31)	(74) 代理人	100105957 弁理士 恩田 誠
審査請求日	平成26年10月22日(2014.10.22)	(74) 代理人	100142907 弁理士 本田 淳
(31) 優先権主張番号	61/410, 641	(72) 発明者	マイク パーカー アメリカ合衆国 98164 ワシントン 州 シアトル フィフス アベニュー 9 01
(32) 優先日	平成22年11月5日(2010.11.5)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 Dragonflyプロセッサ相互接続ネットワークにおけるテーブル駆動型ルーティング

(57) 【特許請求の範囲】

【請求項1】

Dragonflyネットワーク内のターゲットノードから宛先ノードまでの複数のネットワークパスから、1セットのルーティングテーブルに基づきネットワークパスを選択することによって、データを適応的にルーティングするように動作する少なくとも1つのルータを備え、前記Dragonflyネットワークは、複数のルータを備え、前記ルータの各々は、複数のルータグループのうちに対応する1つに含まれており、各ルータグループは、対応するリンクを通じて前記複数のルータグループにおける他のルータグループの各々に接続されており、前記ルーティングテーブルは、前記複数のルータグループにおけるルータグループ内でデータをルーティングするように使用される1以上のローカルテーブルを含み、前記ルーティングテーブルは、前記複数のルータグループにおけるルータグループ間においてデータをルーティングするように使用される1以上のグローバルテーブルをさらに備える、Dragonflyプロセッサ相互接続ネットワークを含むマルチプロセッサコンピュータシステム。

【請求項2】

前記ルーティングテーブルは、前記ターゲットノードと前記宛先ノードとの間において適応的ルーティングを提供するように用いられる、

請求項1記載のマルチプロセッサコンピュータシステム。

【請求項3】

前記ルーティングテーブルは、最小テーブルと非最小テーブルを備える、

10

20

請求項 1 記載のマルチプロセッサコンピュータシステム。

【請求項 4】

前記非最小テーブルは、規制されたポートリストを含む、
請求項 3 記載のマルチプロセッサコンピュータシステム。

【請求項 5】

適応的にルーティングすることは、経路の選択において、近隣ルータからのネットワーク輻輳情報と、前記近隣ルータからのネットワークリンク失敗情報とのうちの 1 または複数を使用することからなる、

請求項 1 記載のマルチプロセッサコンピュータシステム。

【請求項 6】

マルチプロセッサコンピュータシステムの動作方法であって、前記動作方法は、

Dragonfly ネットワーク内のターゲットノードから宛先ノードまでの複数のネットワークパスから、1 または複数のルーティングテーブルに基づきネットワークパスを選択することによって、データを適応的にルーティングするステップ

を備え、前記 Dragonfly ネットワークは、複数のルータを備え、前記ルータの各々は、複数のルータグループのうちの対応する 1 つに含まれており、各ルータグループは、対応するリンクを通じて前記複数のルータグループにおける他のルータグループの各々に接続されており、前記ルーティングテーブルは、前記複数のルータグループにおけるルータグループ内でデータをルーティングするように使用される 1 以上のローカルテーブルを含み、前記ルーティングテーブルは、前記複数のルータグループにおけるルータグループ間においてデータをルーティングするように使用される 1 以上のグローバルテーブルをさらに備える、マルチプロセッサコンピュータシステムの動作方法。

【請求項 7】

前記ルーティングテーブルは、前記ターゲットノードと前記宛先ノードとの間において適応的にルーティングを提供するように用いられる

請求項 6 記載のマルチプロセッサコンピュータシステムの動作方法。

【請求項 8】

前記ルーティングテーブルは、最小テーブルと非最小テーブルとを備える、

請求項 6 記載のマルチプロセッサコンピュータシステムの動作方法。

【請求項 9】

前記非最小テーブルは、規制されたポートリストを含む、

請求項 8 記載のマルチプロセッサコンピュータシステムの動作方法。

【請求項 10】

適応的にルーティングすることは、経路の選択において、近隣ルータからのネットワーク輻輳情報と、前記近隣ルータからのネットワークリンク失敗情報とのうちの 1 または複数を使用することからなる、

請求項 6 記載のマルチプロセッサコンピュータシステムの動作方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、概してコンピュータ相互接続ネットワークに関し、より詳細には、一実施形態において、Dragonfly トポロジープロセッサ相互接続ネットワークでのテーブル駆動型ルーティングに関する。

【0002】

(制限された著作権放棄)

本特許文書の開示の一部には、著作権保護が請求される要素が含まれている。著作権所有者は、本特許文書または本特許開示を任意の人物がファクシミリ複製することについては、これらが米国特許商標庁のファイルあるいは記録に記載されているため異論を唱えないが、その他の権利は全て所有するものである。

【背景技術】

10

20

30

40

50

【0003】

長期間にわたり、コンピュータシステムは、データ伝送を行う場合にネットワーク接続に依存してきた。これは、データ伝送が1つのコンピュータシステムから別のコンピュータシステムへの伝送、1つのコンピュータコンポーネントから別のコンピュータコンポーネントへの伝送、または同じコンピュータ内の1つのプロセッサから別のプロセッサへの伝送のいずれであろうと同じである。多くのコンピュータネットワークは、複数のコンピュータ化した素子どうしを相互にリンクし、またネットワーク上で送信したメッセージが目的の受信者に到着したことの検証、メッセージの整合性の確認、メッセージをネットワーク上の目的の受信者にルーティングする方法といった様々な機能を含む。

【0004】

プロセッサ相互接続ネットワークは、データを1つのプロセッサから別のプロセッサに、または1つのプロセッサグループから別のプロセッサグループに転送するマルチプロセッサコンピュータシステムにおいて使用される。相互接続リンクの数を、数百台または数千台のプロセッサを備えたコンピュータシステムに合わせて変更でき、またシステムパフォーマンスを、プロセッサ相互接続ネットワークの効率に基づき大幅に変更できる。接続の数、送信側の処理ノードと受信側の処理ノードの間の中間ノードの数、そして接続の速度またはタイプは全て、相互接続ネットワークパフォーマンスにおける1つの要素である。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2004-294568号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

同様に、ネットワークトポロジーや、処理ノードどうしの結合に使用する接続パターンによってもパフォーマンスが影響を受けるので、積極的な研究の余地がある。数十台のプロセッサを備えたシステム内で各ノードを相互に直接つなぐことは実用的でなく、プロセッサの数が数千台に達した場合にはほとんど不可能となる。

【0007】

さらに、特に長距離接続または高速の光ファイバリンクを必要とする場合には、通信インターフェース、ケーブル、その他の要素のコストによって、不十分な設計または非効率なプロセッサ相互接続ネットワークのコストが大幅にかさんでしまう。そのため、プロセッサ相互接続ネットワークの設計者にとっては、総リンク数、プロセッサ相互接続ネットワークのコストおよび複雑性を制御しながら、様々な処理ノード間に高速かつ効率的な通信を提供することが課題になっている。

【0008】

したがって、ネットワークのトポロジーまたはマルチプロセッサコンピュータシステムにおいて、どのように1つの処理ノードを別のノードにどのようにしてつなぐか決定するために使用される方法が関心範囲となる。

【課題を解決するための手段】

【0009】

本発明は1例において、複数のプロセッサノードと複数のルータとを設けたDragonflyプロセッサ相互接続ネットワークを実装したマルチプロセッサコンピュータシステムを備える。ルータは、例えばローカルルーティングテーブルおよびグローバルルーティングテーブル、最小および非最小ルーティングテーブルのような、1または複数のルーティングテーブルに基づき、Dragonflyネットワーク内のターゲットノードから宛先ノードまでの複数のネットワークパスの中からネットワークパスを選択することによって、データをルーティングするように動作する。

【図面の簡単な説明】

【 0 0 1 0 】

【図 1】本発明の例示的な一実施形態と一致する Dragonfly ネットワークポロジのブロック図である。

【図 2】本発明の例示的な一実施形態と一致する、ノードにおける Dragonfly ネットワークのスケラビリティを、様々なルータ基数について示したグラフである。

【図 3】本発明の例示的な一実施形態と一致する Dragonfly ネットワークポロジを示すブロック図である。

【図 4】本発明の或る例示的な実施形態と一致する Dragonfly ネットワークポロジグループのブロック図である。

【図 5】本発明の例示的な一実施形態と一致する、仮想チャネルを使った最小ルーティングおよび非最小ルーティングを示す Dragonfly ネットワークのブロック図である。

10

【図 6】本発明の例示的な一実施形態と一致する、様々なルーティングアルゴリズムについてのレイテンシ対供給負荷を、様々なトラフィックパターンを使用して示したグラフである。

【図 7】本発明の例示的な一実施形態と一致する、中間ノードからのバックプレッシャを使用する、グローバルチャネル経由の適応型ルーティングを示す、Dragonfly ネットワークのノードグループ線図である。

【図 8 A】従来のクレジットフロー制御を示すノード線図である。

【図 8 B】本発明の例示的な一実施形態と一致する、クレジットラウンドトリップレイテンシ追跡を示すノード線図である。

20

【図 9】本発明の例示的な一実施形態と一致するルータ構成を示す。

【図 10】本発明の例示的な一実施形態と一致する、Dragonfly プロセッサ相互接続ネットワークにおけるノードグループを示す。

【図 11】本発明の例示的な一実施形態と一致する、Dragonfly プロセッサ相互接続ネットワーク内のいくつかのノードグループ間での接続を示す。

【図 12】本発明の例示的な一実施形態と一致する、Dragonfly プロセッサ相互接続ネットワークルータのためのルータテーブル構成を示す。

【発明を実施するための形態】

【 0 0 1 1 】

30

以下の本発明の例示的な実施形態の詳細な説明では、特定の例を図面と例証の方法により参照する。これらの例は、当業者が本発明を実施するために十分詳細に説明され、どのようにして本発明を様々な目的または実施形態に応用できるか例証する役割を果たす。本発明の他の実施形態が存在し、それらは本発明の範囲内に入り、また、本発明の課題または範囲から逸脱しない限り、論理的、機械的、電気的な変更、さらにその他の変更を加えることができる。ここで記述している本発明の様々な実施形態の特徴または制限は、これらが援用されている例示的な実施形態にとって必須であるが、本発明全体を制限することではなく、また、本発明、およびその要素、操作、用途のいかなる参照も本発明全体を限定することではなく、これら例示的な実施形態を定義するためだけに貢献する。したがって、以下の詳細な説明は、添付の請求の範囲によってのみ定義される本発明の範囲を限定するものではない。

40

【 0 0 1 2 】

相互接続ネットワークは、ハイエンドルータおよびスイッチ用のスイッチングファブリックとしての、マルチプロセッサ内のプロセッサおよびメモリを接続するために、そして、I/O デバイスを接続するために幅広く使用されている。マルチプロセッサコンピュータシステム内のプロセッサおよびメモリのパフォーマンスが向上を続けるに従って、相互接続ネットワークのパフォーマンスがシステム全体のパフォーマンスを決定する中心的役割となっている。ネットワークのレイテンシと帯域幅は、遠隔メモリアクセスレイテンシおよび帯域幅の大部分を確立する。

【 0 0 1 3 】

50

一般に、優れた相互接続ネットワークは、利用可能な技術の能力と制約とから大きく離れない範囲で設計される。例えば少数のポートを維持してポート毎の帯域幅を増加させるのではなく、増加した帯域幅を使ってルータ毎のポートの数を増やす高基数ルータの使用は、ルータピン帯域幅を増加させることによって、動機付けられた。高基数ネットワークを採用した最初のシステムの1つである、クレイ社のBlack Widowシステムは、折り返しクロス(Folded Cross)トポロジの変異型と基数64(radix-64)ルータを使用する。これは従来の低基数3次元トラスネットワークからの重要な新発展である。近年、経済的な光信号通信の登場によって、長距離チャンネルを伴ったトポロジが可能である。しかしこれらの長距離型の光チャンネルは、短距離の電気チャンネルよりも依然として遥かに高額である。そこで、ルータをグループ化してネットワークの有効基数をさらに増加させる先進の光信号通信技術を利用したDragonflyトポロジが導入された。

10

【0014】

ネットワークのパフォーマンスとコストの両方は、相互接続ネットワークのトポロジによって大きく左右される。ネットワークコストは、チャンネルコスト、特に長距離型でグローバルなキャビネット間チャンネルのコストによってそのほとんどが占められる。そのためグローバルチャンネルの数を減らすことによって、ネットワークコストも大幅に低減できる。パフォーマンスを低下させずにグローバルチャンネルの数を減らすには、平均パケットがトラバースするグローバルチャンネルの数を減らす必要がある。Dragonflyトポロジは、最小ルーティングを使用することによって、各パケットがトラバースするグローバルチャンネルの数を1にまで低減する。

20

【0015】

(Dragonflyトポロジの例)

この1つのグローバル直径(global diameter)を達成するには、約2N(ここで、Nはネットワークのサイズである)という非常に高基数のルータを使用する。基数64のルータは既に導入されており、また基数128(radix-128)も実行可能であるが、従来の非常に高基数のルータ技術を使用して各パケットが1グローバルホップのみに限定される場合には、8K~1Mノードの規模のマシンを作るために、これよりも遥かに高い数百または数千という基数が必要となる。各ノードにつき数百台または数千台のポートを設けずに、ルータによってこの非常に高基数の恩恵を達成するために、Dragonflyネットワークトポロジは、サブネットワーク内に接続したルータグループを、非常に高基数の仮想ルータとして使用することを提案する。この非常に高い有効基数によって、全ての最小ルータが最大1本のグローバルチャンネルをトラバースするネットワークを構築することが可能になる。さらに先進の光信号通信技術の能力を利用することによって、グローバルチャンネルの物理長も増す。

30

【0016】

Dragonflyトポロジ上の幅広いトラフィックパターンにおいて優れたパフォーマンスを達成するには、グローバルチャンネルにかけて効率的な負荷分散が可能なルーティングアルゴリズムを選択することが必要となる。グローバル適応型ルーティング(UGAL)は、ルーティング決定を行うソースルータにおいてグローバルチャンネルの負荷を利用できる場合に、このような負荷分散を実行することが可能である。しかしDragonflyトポロジを用いた場合、ソースルータは、当該のグローバルチャンネルに接続されていないことがほとんどである。したがって適応型ルーティング決定は、遠隔または間接的な情報に基づき行われる。

40

【0017】

従来のUGAL(ローカルキュー占有率を使用してルーティング決定を行う)を使用している場合には、この決定の間接的な性質は、レイテンシとスループット両方の劣化を招く。我々は、DragonflyネットワークトポロジのUGALルーティングアルゴリズムに対しこの制限を克服する2つの修正を加えることを提案したところ、グローバル情報を使用した理想的な実現に近いパフォーマンス結果が得られた。UGAL(UGAL

50

- V C H) に選択的な仮想チャネル分別を追加することで、最小パスおよび非最小パス間でのローカルチャネル共有による帯域幅の劣化が排除される。クレジットラウンドトリップレイテンシを、グローバルチャネル輻輳の感知と、この輻輳情報の上り伝播 (U G A L - C R) の両方に使用すると、輻輳の感知にキュー占有率のみを使用した場合よりも遥かに力強いバックプレッシャが提供され、レイテンシの劣化が排除される。

【 0 0 1 8 】

高基数ネットワークでは、低基数ネットワークと比べてネットワークの直径は縮小するが、ケーブルは長距離化する。信号通信技術の進化と近年におけるアクティブ光ケーブルの開発によって、長距離ケーブルを使用した高基数トポロジーの実現が容易化される。

【 0 0 1 9 】

相互接続ネットワークは、パッケージング階層に組み込まれる。最下位レベルでは、ルータは、回路基板を介して接続され、次にこれらがバックプレーンまたはミッドプレーンを介して接続されている。1または複数のバックプレーンがキャビネット内に梱包され、且つ複数のキャビネットが電力ケーブルまたは光ケーブルによって接続されることによって、完全なシステムが形成される。多くの場合、ネットワークコストの大部分は、グローバル (キャビネット間) ケーブルとこれに関連するトランシーバによって占められる。ネットワークコストを最小化するためには、トポロジーは、利用可能な相互接続技術の特性、例えばコストおよびパフォーマンスと一致している必要がある。

【 0 0 2 0 】

距離が増すに従って表皮効果と誘電体吸収による信号減衰が直線的に増加するため、電力ケーブルの最大帯域幅は、ケーブルの長さが長くなると低下する。一般的な高性能信号通信速度 (1 0 ~ 2 0 G b / s) と技術パラメータの場合、電気信号通信パスは、回路基板内で約 1 m、ケーブル内で 1 0 m に制限される。これよりも長い距離では、信号通信速度を低減するかまたはリピータを挿入して、減衰を克服しなければならない。

【 0 0 2 1 】

光信号通信の歴史上、高コストが原因で、光信号通信の使用は非常に長距離や、コスト度外視でパフォーマンスを求める用途に限定されてきた。光ケーブルの固定費用は高いが、データを長距離にわたって銅ケーブルのデータ速度の数倍の速度で伝送する能力のために、光ケーブルの単位距離当たりのコストは、電力ケーブルよりも低くなる。現在の技術を使用して利用できるデータに基づくと、損益分岐点は 1 0 m である。1 0 m 未満の距離では、電気信号通信の方が安価である。1 0 m を超えると、光信号通信がより経済的である。D r a g o n f l y トポロジーは、このコストと距離の関係を利用する。グローバルケーブルの数を減らすことで、光信号通信の高い固定間接費が最小化され、グローバルケーブルを長くすることで、光ファイバのより低い単位当たりのコストの利点が最大化される。

【 0 0 2 2 】

D r a g o n f l y は、本数の少ない長距離のグローバルケーブルを使用するため、1 k ノード以上のネットワークでは D r a g o n f l y のドルコストも平坦化したバタフライ (F l a t t e n e d B u t t e r f l y) と好意的に比較され、4 k ノードまでのものについては約 1 0 % の節約、さらに 4 k ノード以上のものについては約 2 0 % のコスト節約を、平坦化したバタフライトポロジーに関連して示している。これに比較すると、折り返しクロスおよび 3 次元トラスネットワークは、大きなネットワークの直径をサポートするために多数のケーブルが必要であることで問題を抱える。1 k ノードだけのネットワークでは、D r a g o n f l y のコストは 3 次元トラスネットワークのコストの 6 2 %、折り返しクロスネットワークの 5 0 % である。このネットワークコストの低減はネットワーク消費電力の低減と直接関連するので、大型ネットワークにとっては、また、環境に優しいことが望ましい設置にとっては大きな利点となる。

【 0 0 2 3 】

ここで提示している D r a g o n f l y ネットワークの例示的な実施形態は、1 つのルータグループを仮想ルータとして使用することによって、どのようにネットワークの有効

10

20

30

40

50

基数が増加するか、さらに、これによってどのようにネットワークの直径、コスト、レイテンシが低減するか示している。Dragonflyトポロジは、ネットワークにおけるグローバルケーブルの本数を減らしながらグローバルケーブルの長さを増加させるので、先進のアクティブ光ケーブル（固定費用は高いが、単位長さ当たりのコストは電力ケーブルよりも低い）を使用した実現に特に適している。Dragonflyネットワークは、グローバルチャンネルにアクティブ光ケーブルを使用することによって、同じ帯域幅を用いる平坦化したバタフライと比べて20%、やはり同じ帯域幅を用いる折り返しクロスネットワークと比べて52%のコスト低減が可能である。

【0024】

Dragonflyネットワークトポロジの1例を示すために、Dragonflyトポロジと、後述するルーティングアルゴリズムの例との説明において次の記号を使用する：

- N ネットワーク端末の数、
- p 各ルータに接続された端末の数、
- a 各グループにおけるルータの数、
- k ルータの基数、
- $k_{_}$ グループの有効基数（または仮想ルータ）、
- h 他のグループに接続するために使用された各ルータ内のチャンネルの数、
- g システム内に存在するグループの数、
- q 出力ポートのキュー深度、
- q_{vc} それぞれの出力VCのキュー深度、
- H ホップカウント、
- Out_i ルータ出力ポートi。

【0025】

Dragonflyトポロジは、図1に示すように、ルータ104、105、106、グループ101、102、103、システム、の3レベルで構成された階層ネットワークである。ルータレベルにて、各ルータはp個のノードに接続し、a-1個のローカルチャンネルは同じグループ内の他のルータに接続し、h個のグローバルチャンネルは他のグループ内のルータに接続している。したがって各ルータの基数（または度数）は、 $k = p + a + h - 1$ と定義される。グループは、図1の符号101で示すように、ローカルチャンネルで形成されたグループ間相互接続ネットワークを介して接続したa個のルータで構成されている。各グループは端末へのap個の接続を有し、グローバルチャンネルへのah個の接続を有しており、また、1グループ内の全てのルータは、基数 $k' = a(p + h)$ を伴う仮想ルータとして集合的に機能する。この非常に高い基数 $k' \gg k$ により、非常に小さいグローバル直径でのシステムレベルネットワークの実現が可能になる（任意の2つのノード間の最小パス上に、最大数の高価なグローバルチャンネルを設けている）。最大で $g = ah + 1$ までの数のグループ（ $N = ap(a h + 1)$ 個の端末）を1のグローバル直径に接続することができる。これに対し、基数k個のルータで直接構築したシステムレベルネットワークでは、より大きなグローバル直径が必要となる。

【0026】

最大サイズ（ $N = ap(a h + 1)$ ）のDragonflyでは、各1対のグループの間には厳密に1つの接続しか存在しない。これよりも小規模のDragonflyでは、各グループ以外との接続以外のグローバル接続の方が、他のグループとの接続よりも多い。これらの追加的なグローバル接続は、複数のグループにかけて分布しており、この複数のグループは、少なくとも $ah + 1g$ 本のチャンネルで接続した各1対のグループ毎にまとめられている。

【0027】

Dragonflyパラメータa、p、hは、任意の値であってよい。しかしチャンネル負荷を分散させるために、この例のネットワークは、 $a = 2p = 2h$ を有する。各パケットは、その経路に沿って、1本のグローバルチャンネルと1本の端末チャンネルの合計2本の

10

20

30

40

50

ローカルチャネル（グローバルチャネルの各終端部分につき1本ずつ）をトラバースするので、この比率によってバランスが保たれる。グローバルチャネルは高価であるため、或る実施形態では、ローカルチャネルと端末チャネルを過剰供給して高価なグローバルチャネルが十分に利用される状態を保つ方法によって、この2：1の比率からの逸脱を行っている。つまり、このような例では、 $a = 2h$ 、 $2p = 2h$ となるようにネットワークのバランスをとっている。

【0028】

バランスのとれたDragonflyのスケラビリティを図2に示す。有効基数を増加することでDragonflyトポロジは高度にスケラブルとなり、基数64ルータを用いた場合、トポロジは、わずか3のホップ数でネットワークの直径256kノード以上にまでスケラリングされる。図1のグループ内およびグループ間ネットワークには、任意のネットワークを使用できる。本明細書で提示する例は、両方のネットワークについて、1次元の平坦化したバタフライまたは完全接続型トポロジを使用している。図3に、Dragonflyの単純な例を示す。ここでは $p = h = 2$ （ルータ1つにつき処理ノード2つ、そして、各ルータ内には他グループに接続したチャネル2本が設けられている）、 $a = 4$ （各グループ内にルータ4つ）であり、これが、 $k = 7$ （基数7）のルータを用いて、 $N = 7^2$ （ネットワーク内の72ノード）にスケラリングされる。図3のグループ G_0 が8つのグローバル接続と8つのノード接続を有するので、仮想ルータを使用することによって、有効基数が $k = 7$ から $k' = 16$ にまで増加する。

【0029】

グループ内ネットワークにより高次元のトポロジを使用することによって、グローバル基数 k' をさらに増加することができる。さらに、このようなネットワークは、グループ内でのパッケージ化の局所参照性も利用する。例えば図4の符号401が示す2次元の平坦化したバタフライは、図5に示すグループと同じ k' を有するが、より多くの帯域幅をローカルルータに提供することによって、パッケージ化局所参照性を利用する。図4の符号402では、3次元の平坦化したバタフライを使用して、有効基数を $k' = 16$ から $K' = 32$ に増加することによって、図1のものと同じ $k = 7$ のルータを使用してトポロジを最大 $N = 1056$ にまでスケラリングできるようにしている。

【0030】

Dragonflyのような高基数ネットワークの端末帯域幅を増加するために、チャネルスライシングを採用することができる。チャネルを幅広くするのではなく（この場合、ルータ基数が減少してしまう）、複数のネットワークを並列接続して容量を追加することができる。同様に或る実施形態では、Dragonflyトポロジは、ネットワーク容量を追加するために、並列ネットワークを利用することもできる。これに加えて、ここまで説明したDragonflyネットワークでは、ネットワーク内の全てのノードに均等な帯域幅を仮定した。しかしこのような均等な帯域幅が不要な場合には、いくつかのグループからグループ間チャネルを除去することによって、帯域幅のテーパリングを実現することが可能である。

【0031】

（Dragonflyルーティングの例）

様々な最小および非最小ルーティングアルゴリズムは、Dragonflyトポロジを使用して実現できる。ローカル情報を使用するグローバル適応型ルーティングの或る実施形態は、中間負荷にて、スループットの制限と非常に高いレイテンシを招く。これらの問題を克服するために、我々はグローバル適応型ルーティングに、理想的なグローバル適応型ルーティングの実現にアプローチするパフォーマンスを提供できる新たなメカニズムを導入する。

【0032】

Dragonflyにおける、グループ G_s 内のルータ R_s に取り付けたソースノード s からグループ G_d 内のルータ R_d に取り付けた宛先ノード d までの最小ルーティングが、1本のグローバルチャネルをトラバースし、これは次の3つのステップで達成される：

10

20

30

40

50

ステップ1: $G_s _ = G_d$ であり且つ R_s が G_d と接続していない場合、 G_s 内で R_s から R_a (G_d へのグローバルチャネルを有するルータ) までをルーティングする。

【0033】

ステップ2: $G_s _ = G_d$ である場合、グローバルチャネルを、 R_a から G_d 内のルータ R_b に到達するまでトラバースする。

ステップ3: $R_b _ = R_d$ である場合、 G_d 内で R_b から R_d までルーティングする。

【0034】

この最小ルーティングは、負荷分散トラフィックには上手く適応するが、これと対抗するトラフィックパターン上では満足なパフォーマンスが得られない。対抗するトラフィックパターンを負荷分散させるためには、ヴァリエント (Variant) のアルゴリズムをシステムレベルに適用して、各パケットをまずランダムに選択した中間グループ G_i に、次にその最終宛先 d にルーティングすることができる。ヴァリエントのアルゴリズムをグループに適用することによって、グローバルチャネルとローカルチャネルの両方の上の負荷は、十分に分散される。このランダム型の非最小ルーティングは、最大で2本のグローバルチャネルをトラバースし、また、次の5つのステップを必要とする：

ステップ1: $G_s _ = G_i$ であり且つ R_s が G_i と接続していない場合、 G_s 内で R_s から R_a (G_i へのグローバルチャネルを有するルータ) までをルーティングする。

【0035】

ステップ2: $G_s _ = G_i$ である場合、グローバルチャネルを、 R_a から G_i 内のルータ R_x に到達するまでトラバースする。

ステップ3: $G_i _ = G_d$ であり且つ R_x が G_d と接続していない場合、 G_i 内で R_x から R_y (G_d へのグローバルチャネルを有するルータ) までをルーティングする。

【0036】

ステップ4: $G_i _ = G_d$ である場合、グローバルチャネルを、 R_y から G_d 内のルータ R_b に到達するまでトラバースする。

ステップ5: $R_b _ = R_d$ である場合、 G_d 内で R_b から R_d までをルーティングする。

【0037】

ルーティングデッドロックを防ぐためには、図5に示すように、最小ルーティングに2本の仮想チャネル (VC) を採用し、また、非最小ルーティングの場合は3本の仮想チャネルが必要である。これらの仮想ルータを指定することで、ルーティングによって生じるチャネル依存が全て排除される。いくつかの用途では、プロトコルデッドロックを回避するために追加の仮想チャネルが必要になることがある。例えば共有メモリシステムでは、メッセージの要求と応答のために仮想チャネルの別個のセットが必要になることがある。

【0038】

次のような Dragonfly トポロジーのための様々なルーティングアルゴリズムが評価されてきた：

最小 (MIN) : 先述したように最小パスを経る。

【0039】

バリエント (VAL) [32] : 先述したランダム型の非最小ルーティング。

ユニバーサルグローバル適応型負荷分散 [29] : ネットワークを負荷分散するために、(UGAL-G, UGAL-L) UGAL が、MIN と VAL の中からパケットバイパケットに基づき選択する。この選択は、ネットワーク遅延を推定するためにキュー長とホップカウントを使用し、遅延が最も小さいパスを選択することによって行われる。我々は、次の2つのバージョンの UGAL を実現する。

【0040】

UGAL-L : 現在のルータノードでのローカルキュー情報を使用する。

UGAL-G : G_s 内の全てのグローバルチャネルのためのキュー情報を使用する (他のルータ上のキュー長がわかっていると仮定する)。ローカルチャネルではなくグロー

10

20

30

40

50

バルチャネルの負荷分散が必要なので、これは実現が困難である一方で、UGALの理想的な実現を表すものである。

【0041】

図6に示すように、BENIGNパターンと、これに対抗する合成トラフィックパターンとの両方を使用して、異なるルーティングアルゴリズムを評価する。符号601の均等なランダムトラフィックと、これに対抗する符号602のトラフィックとの両方を使用して、4つのルーティングアルゴリズムについてのレイテンシ対供給負荷を示す。合成トラフィックパターンを使用することで、ネットワークを十分に評価するためにトポロジーとルーティングアルゴリズムを強調できるようになる。図6の符号601で示すように、均等ランダム型(UR)のようなBENIGNトラフィックでは、MINで十分に低レイテンシと高スループットを提供できる。VALは、その負荷分散によってグローバルチャネル上の負荷が2倍になるので、ネットワーク容量の約半分を達成する。UGAL-GとUGAL-Lの両方はMINのスループットと似ているが、これらの方が飽和付近でのレイテンシが若干高い。この若干高いレイテンシは、並列またはグリーディ割当の使用によって生じる。並列またはグリーディ割当の使用では、各ポートにてルーティング決定は、並列に行われる。逐次割当を使用することによって、より複雑なアロケータの犠牲の上にレイテンシは、短縮する。

10

【0042】

Dragonflyでの適応型ルーティングは、ルータ出力ではなく、グローバルチャネル、グループ出力のバランスを取る必要があるため、課題が伴う。これによって、間接的なルーティング問題が生じる。各ルータは、グローバルチャネルの状態に間接的にのみ依存するローカル情報だけを使用して、使用するグローバルチャネルを選ぶ。先行技術によるグローバル適応型ルーティング方法は、ネットワーク輻輳を正確に推定するために、ローカルキュー情報、ソースキュー、出力キューを使用する。これらのケースでは、開始させた経路上の輻輳を直接示すローカルキューは、グローバル輻輳の正確なプロキシである。しかしDragonflyトポロジーを使用した場合には、ローカルキューは、ローカルチャネルにかかるバックプレッシャを介してグローバルチャネル上の輻輳を感知するだけである。ローカルチャネルを過剰提供した場合、ソースルータが輻輳を感知する以前に、過負荷状態の最小経路上で著しい数のパケットをエンキューしなければならない。これによって、先に図6の符号602で示したように、スループットとレイテンシが低下してしまう。

20

30

【0043】

UGAL-Lに伴うスループットの問題は、1本のローカルチャネルで最小および非最小トラフィックの両方を扱うことが原因で生じる。例えば図7では、R1のパケットは、gc7を使用する最小パスと、gc6を使用する非最小パスとを設けている。両パスは、同じR1からR2までのローカルチャネルを共用する。両パスが同じローカルキューを共用し(したがって、同じキュー占有率を有する)、最小パスの方が非最小パスよりも短いため(グローバルホップ:1対2)、常に最小チャネルが選択される。これは、たとえパスが飽和状態にある時でも同じである。これによって最小グローバルチャネルが過負荷状態となるため、この最小チャネルと同じルータを共用している非最小グローバルチャネルは、利用されなくなってしまう。UGAL-Gを使用することで、最小チャネルが優先され、負荷は、全ての他のグローバルチャネルにかけて均等に分散される。これに対しUGAL-Lを使用することで、最小グローバルチャネルを含んでいるルータの非最小チャネルが利用されなくなり、その結果、ネットワークスループットは、低下する。

40

【0044】

この制限を克服するために、それぞれの仮想チャネル(UGAL-LVC)を使用することによって、キュー占有率を最小成分と非最小成分に分けるようにUGALアルゴリズムを修正する:

```
if ( q m v c H m q n m v c H n m )
    route minimally ;
```

50

```
else
```

```
    route nonminimally;
```

ここで、添字 m は最小パスを、 nm は非最小パスを示す。図 5 の仮想チャネル割当を使用すれば、 $q_m \ v_c = q(V \ C1)$ および $q_{nm} \ v_c = q(V \ C0)$ となる。

【0045】

比較すると、UGAL-LVC は WC トラフィックパターン内の UGAL-G のスループットと一致するが、UR トラフィックでは、スループットの制限により、スループットが約 30% 低下する。ほとんどのトラフィックを非最小で送信する必要がある WC トラフィックの場合、最小キューが負荷の大きな状態となっているため、UGAL-LVC が上手く機能する。しかし、ほとんどのトラフィックを最小で送信する必要がある時に負荷分散されたトラフィックを用いると、それぞれの仮想チャネルがチャネル輻輳の正確な表示を提供せず、その結果、スループットが低下してしまう。

10

【0046】

この制限を克服するために、我々は、最小パスと非最小パスが同一の出力ポートから開始する場合にのみキュー占有率を最小成分と非最小成分に分けるように、UGAL アルゴリズムをさらに修正する。我々のハイブリッド修正版の UGAL ルーティングアルゴリズム (UGAL-LVCH) は、次のとおりである：

```
    if ( q_m H m   q_n m H n m   &&   Out m _ = Out n m   ) || ( q_m
v_c H m   q_n m   v_c H n m   &&   Out m = Out n m )
        route minimally;
```

20

```
else
```

```
    route nonminimally;
```

UGAL-LVC と比較すると、UGAL-LVCH が提供するスループットは、WC トラフィックパターン上のものと同じであるが、UR トラフィック上の UGAL-G のスループットに一致し、それ故に、飽和状態に近い 0.8 の供給負荷での高い方のレイテンシのほぼ 2 倍となる。WC トラフィックでは、UGAL-LVCH は、中間レイテンシも UGAL-G のもの比べて高くなる。

【0047】

この UGAL-L の高い中間レイテンシは、輻輳を感知する前に、ソースと輻輳ポイントとの間のチャネルバッファが最小経路で送られたパケットで充填されることで生じる。我々の研究では、非最小経路で送られたパケットは UGAL-G と比較可能なレイテンシ曲線を有する一方で、最小経路で送られたパケットは著しく高いレイテンシに遭遇することが示されている。入力バッファが増加すると、最小経路で送られたパケットのレイテンシが増加し、パケットのレイテンシは、バッファの深さに比例する。レイテンシ分布のヒストグラムは、次の 2 つの明白な分布を示している。1 つは、非最小パケットについてのレイテンシが低い大きな分布であり、他の 1 つは、パケット数が制限されているが、最小パケットについてのレイテンシが遥かに高い分布である。

30

【0048】

UGAL-L に伴うこの問題を理解するために、図 7 に示す Dragonfly グループの例において、R1 のパケットは、gc0 を介して最小のルーティング、または gc7 を介して非最小のルーティングのいずれを行うかについてのグローバル適応型ルーティング決定を行っているとは仮定する。ルーティング決定は、グローバルチャネル利用の負荷分散を行う必要があり、またチャネル利用をグローバルチャネル q_0 、 q_3 に関連したキューから得られることが理想的である。しかし、 q_0 、 q_3 キュー情報は R0、R2 でしか利用できず、R1 ではまだ利用できる状態にないので、ルーティング決定は、R1 で利用できるローカルキュー情報を介して間接的にしか行えない。

40

【0049】

この例では、 q_1 は q_0 の状態を反映し、 q_2 は q_3 の状態を反映する。 q_0 または q_3 のどちらかがフル状態である場合には、図 7 の矢印で示すように、フロー制御によって q_1 と q_2 にバックプレッシャが提供される。その結果、安定した状態の測定が得られ、

50

スループットの正確な測定にこれらのローカルキュー情報を使用できるようになる。スループットは、レイテンシが無限大になると（またはキュー占有率が無限大になると）供給負荷として定義されるため、このローカルキュー情報で十分である。しかし q_1 が $g_c 0$ の輻輳を反映でき、且つ R_1 がパケットを非最小経路で送ることができるようにするには、 q_0 は、完全にフル状態になる必要がある。したがってローカル情報を使用するには、数個のパケットを犠牲にして正確な輻輳を決定することを要し、その結果、最小で送られているパケットのレイテンシが遥かに高くなってしまふ。負荷が増加するに従って、最小経路で送られたパケットのレイテンシは増加し続けるが、より多くのパケットが非最小で送られるようになり、その結果、飽和までの平均レイテンシは、低下する。

【 0 0 5 0 】

ローカルキューがグローバル輻輳の優れた推定を提供できるようにするには、グローバルキューが完全にフル状態となり、ローカルキューに向けて力強いバックプレッシャを提供する必要がある。バックプレッシャの力強さはバッファの深さに反比例し、バッファが深いほどバックプレッシャの伝播に時間がかかり、バッファが浅いほど遥かに力強いバックプレッシャが提供される。バッファサイズが減少するに従って、バックプレッシャが力強くなるため、中間負荷におけるレイテンシは、低下する。しかし使用するバッファの数を減らすと、ネットワークスループットの低下という犠牲を払うことになる。

【 0 0 5 1 】

高い中間レイテンシを克服するために、我々は、クレジットラウンドトリップレイテンシを使用して、高速な輻輳の感知とレイテンシの低減を図ることを提案する。図 8 に示すクレジットベースのフロー制御では、下流バッファについてクレジットカウントが維持される。パケットが下流に送られるに従い、適切なクレジットカウントが減少してゆくが、パケットが下流ルータを離れるとクレジットは上流に送り戻され、クレジットカウントが増加する。クレジットが戻るためのレイテンシはクレジットラウンドトリップレイテンシ（ t_{crt} ）と呼ばれ、下流に輻輳がある場合にはパケットが直ぐに処理されないため、結果として t_{crt} が増加する。

【 0 0 5 2 】

図 8 を参照すると、図 8 A では従来のクレジットフロー制御を示している。パケットが下流に送信されると（1）、出力クレジットカウントが減少し（2）、クレジットが上流に送り戻される（3）。図 8 B では、このスキームが、ネットワーク内の輻輳を推定するためにクレジットラウンドトリップレイテンシを使用するように修正される。減少している出力クレジットカウント（2）に加えて、CTQ で示すクレジットタイムキュー内にタイムスタンプが押し込まれる。クレジットを上流に送り戻す前に（4）、クレジットが遅延し（3）、また、下流のクレジットが受信されると（5）、クレジットカウント並びにクレジットラウンドトリップレイテンシ t_{crt} がアップデートされる。

【 0 0 5 3 】

t_{crt} の値を、グローバルチャネルの輻輳を推定するために使用できる。我々は、上流クレジットを遅延させるためにこの情報を使用することで、バックプレッシャを力強くして、輻輳情報をより高速に上流に伝播できるようにした。各出力 O について、 $t_{crt}(O)$ が測定され、量 $t_d(O) = t_{crt}(O) - t_{crt0}$ がレジスタに記憶される。次に、クレジットを即座に上流に送り戻す代わりに、フリットを出力 O に送信する場合には、クレジットは $t_d(O) - \min[t_d(o)]$ だけ遅延される。グローバルチャネル上で送信されるクレジットは遅延しない。これによってこのメカニズムに周期ループが存在しないことが保証され、グローバルチャネルを十分に利用できるようになる。

【 0 0 5 4 】

戻されるクレジットの遅延によって、力強いバックプレッシャを作るための、より浅いバッファが提供される。しかしバッファ全体が利用され、高負荷においてスループットの低下が生じないことを保証するために、全ての出力にかけてクレジットを t_d の変動分だけ遅延させる必要がある。我々は、この変動を、 $\min[t_d(o)]$ 値を求め、差分を使用することによって推定する。クレジットを遅延させることで、上流ルータが輻輳を（

10

20

30

40

50

キューが充填されるのを待つ時と比べて)より高速で観察し、より優れたグローバル適応型ルーティング決定が得られるようになる。

【0055】

WCおよびURトラフィックの両方について、クレジットレイテンシ(UGAL-LCR)を使用したUGAL-Lルーティングアルゴリズム評価が、深さ16と256のバッファを使用して調査される。UGAL-LCRはUGAL-Lと比べてレイテンシを著しく低下させ、UGAL-Gのレイテンシ近くにまでする。WCトラフィックについて、UGAL-LCRは、深さ16のバッファでレイテンシを最大35%低下させ、また、深さ256のバッファで中間レイテンシを最大で20分の1に低下させた(低下率はUGAL-Lとの比較)。UGAL-Lとは違い、UGAL-LCRの中間レイテンシはバッファサイズと無関係である。URトラフィックについて、UGAL-LCRは、飽和付近にてレイテンシを、UGAL-LVCHと比較し最大50%低下させる。しかし、UGAL-LCRとUGALLVCHの両方共、そのローカル情報が不正確なためにいくつかの packets が非最小でルーティングされてしまうので、URトラフィックを伴うUGAL-Gのスループットには達しない。

10

【0056】

このスキームを実現した結果、各ルータ側に必要なのは次の3つの特徴であるため、複雑性にかかる間接費は最小であった。

- ・tcrtを測定するために個々のクレジットの追跡、
- ・td値を記憶するためのレジスタ、
- ・クレジットを戻す際の遅延メカニズム。

20

【0057】

必要なtd記憶量は、O(k)レジスタのみが必要であるため最小量である。クレジットはデータフリット上にピギーバックされて戻されることが多く、また、次の上流でのデータフリットの送信を待つために、クレジットを遅延させる必要がある。提案したこのメカニズムに必要なことは、さらなる遅延を追加することだけである。

【0058】

個々のクレジットの追跡については、従来から、クレジットは、クレジットフロー制御においてクレジットのプールとして追跡されており、つまり、各出力仮想チャネルにつき1つのクレジットカウンタが維持され、このクレジットカウンタは、クレジットが受信されると増加する。UGAL-LCRの実現には、各クレジットを個別に追跡する必要がある。これは、図8Bに示すように、フリットが送信される度に、クレジットタイムスタンプキュー(CTQ)を使ってキューの末尾にタイムスタンプを押し、該当するクレジットが到着したらキューの先頭からタイムスタンプを取り出すことによって実行される。フリットとクレジットの割合は1:1であり、オーダリングを維持するので、ラウンドトリップクレジットレイテンシの測定には単純なキューで十分である。キューの深さはデータバッファの深さに比例していなければならないが、キューサイズは、輻輳の測定に不正確な情報を利用するために縮小される(例えば、データバッファサイズの1/4のサイズのキューを設ければ、輻輳の測定を行うために、4個のクレジットのうちの1個のみを追跡すればよくなる)。

30

40

【0059】

Dragonflyトポロジーにかかるコストも、平坦化したバタフライや他のトポロジーのコストと比較して優れている。平坦化したバタフライトポロジーは、中間ルータとチャネルを除去することによってButterflyのネットワークコストを低減する。その結果、平坦化したバタフライは、バランスの取れたトラフィック上での折り返しクロスと比べてコストが約50%低減する。Dragonflyトポロジーは、ルータの有効基数を増加して、さらなるコスト低減およびネットワークのスケラビリティを向上することによって、平坦化したバタフライを拡張する。

【0060】

各々64k個のノードを接続したDragonflyネットワークと平坦化したバタフ

50

ライネットワークを比較したところ、平坦化したバタフライがグローバルチャネル用ルータポートの50%を使用する一方で、Dragonflyはグローバル接続用ポートの25%を使用することを示した。平坦化したバタフライは2つの次元を追加する必要があるが、Dragonflyは1次元である。さらにDragonflyでは、グループサイズを増加してネットワークのスケーリングが可能なので、より優れたスケーラビリティが得られるのに対し、平坦化したバタフライでは次元を追加する必要がある。Dragonflyは、ホップカウントがほぼ同一な状態であれば、より長い方のグローバルケーブルを相殺してグローバルケーブルの本数を減少させることによって、先進の信号通信技術に適合する、よりコスト効果的なトポロジーを提供する。

【0061】

ここで説明しているDragonflyネットワークの様々な実施形態も、Dragonflyが呈する間接的な適応型ルーティングの課題を克服する、グローバル適応型ルーティングの2つの新たなバリエーションを備える。一般に、Dragonflyルータは、同一グループ内の別のルータに取り付けられているグローバルチャネルの状態に基づきルーティング決定を行う。この遠隔チャネルの状態を推論するためにローカルキュー占有率を使用する従来のグローバル適応型ルーティングアルゴリズムは、スループットとレイテンシを低下させてしまう。そこで我々は、仮想チャネル分別の選択的な使用を導入することによって、帯域幅の減少を克服する。さらに我々は、チャネル輻輳を感知してこれを信号通信するために、クレジットラウンドトリップレイテンシも使用する。この2つの技術を組み合わせることによって、遠隔チャネル状態を完全に知得した理想的なアルゴリズムのパフォーマンスに近づくよう試みるグローバル適応型ルーティングアルゴリズムが得られる。

【0062】

(Dragonflyネットワークにおける革新的な適応型ルーティング)

ここでは、輻輳リンクまたはダウンしたリンクに基づき複数の正当な経路の中から選択を行うよう動作可能な、デッドロックを回避する適応型ルーティングを提供することによって、Dragonflyプロセッサ相互接続ネットワークのための向上したルーティング方法を提案する。この適応型ルーティング方法は、向上したルーティングパフォーマンスを提供することと、ダウンしたリンクまたはトラフィックの多いリンクを許容することを従来の方法よりも上手く行い、さらに帯域幅に悪影響を与えるクレジットを保留するのではなく、チャネル上の輻輳を明快に通信させる。

【0063】

或る実施形態では、ネットワーク経路は、例えば複数の異なる次元でのルーティングのような複数の最小経路からまず選択され、次に任意で、例えば輻輳リンクやダウンしたリンクを回避するためにランダム選択したホップを使用することによって、1または複数の非最小経路からさらに選択される。

【0064】

1例では、ルーティングの選択は、テーブルによって提示され、ネットワークの構成および状態に応じて、特定の経路、あるいは、最小経路または非最小経路に偏る可能性がある。例えば経路の選択が最高の効率のデフォルトによって最小ルーティングに偏るが、この偏りは、追加のトラフィックを任意または不必要に受信することから特定のネットワークリンクを保護するために、非最小ルーティングへの偏りに切り替わることもある。

【0065】

或る実施形態では、輻輳情報は、例えば出力キュー内のメッセージ数のカウントのような要素から予測される次のリンク輻輳を導出し、伝送中のクレジットやメッセージのような要因に基づき受信バッファ輻輳の推定を確立することによって利用される。ノードは、潜在的な受信側ノードに、平均的な「次のリンク」出力の輻輳について問い合わせ、ノードが輻輳リンクまたはダウンしたリンクの回避に基づきルーティング決定を行えるようにすることができる。

【0066】

10

20

30

40

50

図9は、本発明の例示的な実施形態と一致するDragonflyネットワークルータを示す。ここに示すルータブロックは、それぞれが入力/出力の対に対応した48個のタイルを備えている。タイルは8×6の行列に編成されているため、特定のタイルにおける入来パケットデータ(incoming packet data)を、8列のうちの一つにつながる行にかけてルーティングし、次に、8列を上って、または下って、6行のうちの一つにルーティングし、適切なタイルに到達させて出力させる。さらなる実施形態では、チャンネルは、複数の仮想チャンネルと、仮想チャンネル伝送中切り換えと、SECCEDのようなエラー修復と、さらに、ネットワークパフォーマンスを向上するために必要に応じて仮想チャンネルへの動的割当を含む入力バッファリングとを特徴とする。

【0067】

再び図9の例を参照すると、タイルのうち40個は外部ネットワークリンクに接続しており、8個はプロセッサノード域内のプロセッサコアに接続している。各タイルは、入力キュー、サブスイッチ、列バッファを備えている。入力キューは、シリアライザ/デシリアライザインターフェースからネットワークに送られたパケットを受信し、このパケットをどのようにルーティングするか決定する。パケットは、行バス上で、適当な列のサブスイッチへと送信される。サブスイッチは、このパケットを受信すると、これを適切な仮想チャンネルへと切り換えてから、6列のバスのうちの一つを介して適切な行内の列バッファに送る。列バッファは、列内の6個のタイルからのパケットデータを収集し、これらをネットワークチャンネル上で送信する。

【0068】

この例におけるDragonflyネットワークトポロジーは、2層の平坦化したバタフライトポロジーで構成された階層ネットワークである。第1層は、コンピュータキャビネットやシャーシのようなローカルグループ内の全てのルータチップを接続する2次元の平坦化したバタフライである。各グループは、非常に高基数のルータとして扱われ、また、単一次元の平坦化したバタフライ(all-to-all)は、全てのグループを接続して、ここで示すDragonflyトポロジーの例の第2層を形成する。

【0069】

グループ内の第1の次元は(便宜上「緑色」次元と呼ぶ)シャーシ内の16個のルータを接続する。グループ内の第2の次元は(同様に「黒色」次元と呼ぶ)、2キャビネットから成るグループ内の6個のシャーシを接続する。これは、図10のネットワーク「グループ」に示すネットワーク構成に反映されている。この図10のネットワーク「グループ」には、各シャーシにつき16個のルータ(16個の列で示す)で構成された6つのシャーシ(6つの行で示す)が図示されている。

【0070】

図10に示したようなグループは、図11に示すように、「青色」次元のリンクを使用してさらに相互結合する。これらのグループ間の「青色」リンクは、各グループを他の各グループと接続するものであり、接続可能な数は、この例では各グループにつき最大で240の青色リンク、または各システムにつき241グループである。各リンクは、例えば1つのリンクまたは1本の光ケーブルにつき4ポートというように、複数のポートを備えることができる。したがって4つのポートは、1本のケーブルでグループの各対に接続することになる。グループ数がより少ないシステムでは、各グループにつき240個の青色ポートのうち未使用ポートを使用することによって、構成グループ間に追加の帯域幅を提供できる。これによって例えば、120個のグループで構成され、且つグループの各対を接続する8個のポートを提供しているネットワーク内部の各グループの対につき2つのリンクを設けることができる。

【0071】

このネットワークでは、ソースノードからターゲットノードまでルーティングされるパケットは、図9、図10、図11に示した次元のうち少なくとも一つの次元、しかしおそらくは3つ全ての次元をトラバースする。3つ全ての次元をトラバースするルーティングパスは、まず緑色次元にルーティングされ、次に黒色次元に送られて、ターゲットグルー

10

20

30

40

50

ブとつながっているグループ内の適切なノードに達し、その後青色次元にルーティングされて目的のターゲットグループに到達する。次にパケットは、グループ内の緑色次元と黒色次元とにルーティングされ、ターゲットグループ内の目的のターゲットノードに到達することによって、ターゲット到達までに、3つの次元における5つのルーティングを辿ったことになる。

【 0 0 7 2 】

一実施形態では、このネットワークは、適応型ルーティングと決定論的なルーティングの両方をサポートする。決定論的なルーティングは、ネットワーク輻輳に関係なく、所与のパケットをネットワーク上の定義済みの経路上に送信する。複数の決定論的なパスが利用可能である場合には、複数のパス間のトラフィックを分布させるために、宛先ノード、アドレス、または他の同様の特征に基づき、決定論的なトラフィックをハッシュすることができる。或る実施形態では、ソース ターゲット間の全てのパケットは同一の決定論的なパスを使用するので、同一のソース ターゲット間を移動するパケットはターゲットに順番どおりに到達する。

10

【 0 0 7 3 】

適応型ルーティングによって、パケットは、ネットワーク内の輻輳レベルに基づき、複数の異なる経路を使用できる。或る実施形態では、適応型ルーティングを使用すると、パケットは元の順番と違うバラバラの順番で到着してもよく、また輻輳のために最小パスの回避が指示された場合には、非最小パスを使用してもよい。

20

【 0 0 7 4 】

Dragonflyにおける最小ルーティングは、パケットが所与の次元の最大1つのリンクをトラバースする際に生じる。したがって、例えば図10に示すようなグループ内の最小ルーティングは、「緑色」次元の最大1つのホップと、「黒色」次元の1つのホップとを使用する。異なるグループにあるノードどうし間の最小パスは、各グループにおける、緑色次元の最大1つのホップと、黒色次元の1つのホップとを使用し、さらに追加の1つのホップを使用して、ソースグループからターゲットグループに移動する。

【 0 0 7 5 】

最初に黒色または緑色次元のどちらかをトラバースすることができるので、ソースグループと宛先グループとの両方に複数の最小パスが存在する。グループ間に複数のリンクが存在する場合には、1つのパスにおける、ソースグループと宛先グループのどちらかの黒色または緑色次元でのホップを0にすることによって、最小パスの完了に必要な総ホップ数を5未満に減らすことができる。

30

【 0 0 7 6 】

非最小ルーティングでは、ソースグループまたはターゲットグループ内の黒色または緑色いずれかの次元でのホップを複数にすることによって、総ホップ数を5以上にすることができる。ルータまでの最小パスまたは利用可能なパスに輻輳が発生している状況では、さらにホップを追加することによって、ターゲットへのメッセージ伝送速度を向上させる一方で、既に輻輳状態にあるネットワークリンクをさらに輻輳させないようにすることが望ましい。さらなる実施形態では、既に輻輳しているリンク周辺で同じパスを繰り返しルーティングした結果、輻輳したネットワーク領域をさらに作ってしまうことを回避するために、例えばパス選択をランダム化またはハッシュすることによって、トラフィックを利用可能なリンクにかけて拡散するべく試みる。

40

【 0 0 7 7 】

このような実施形態の1つでは、図10に示したようなグループから中間ノードを1つ選択することで、メッセージをまず最小経路でこの中間ノードまでルーティングし、次に中間ノードからグループ内の最終ノードまでルーティングできるようになる。これによって、緑色次元と黒色次元の各々で最大2ホップ、またはグループ内の最小ルーティングにおいてはその2倍のホップ数となる。ルーティングは、ソースグループ内で非最小、ターゲットグループ内で非最小、またはソースグループおよびターゲットグループの両方において非最小であってよい。

50

【 0 0 7 8 】

メッセージを、ソースグループとターゲットグループにおいて最小経路で、しかし両グループ間のリンクの輻輳を回避するために、ソースグループとターゲットグループの間の中間グループを経由させてルーティングする場合には、非最小ルーティングは、グループ間においても生じ得る。ソース、中間、ターゲットグループにおけるルーティングは、さらに、各グループでの輻輳に応じて最小または非最小であってよい。

【 0 0 7 9 】

一実施形態では、所与のパケットまたはメッセージに使用するルーティングのタイプは、パケットヘッダ内のルーティング制御フィールドによって決定される。例えばルーティング制御記号は、パケット順序を保つことが望ましい時には、決定論的な非最小ハッシュ化ルーティングを使用すべきであることを示してもよい。パケットは、ターゲットノードをハッシュとして使用することによって、利用可能な複数のパスにかけて分布される。トラフィックは非最小ルーティングされるが、パケットをグループ内の様々な中間ノード間に分布させることによって、ホットスポットまたは輻輳は、減少する。

【 0 0 8 0 】

決定論的な最小ハッシュ化ルーティングは、最小パス上でパケットのハッシュを提供するが、これは、緑色次元より前に黒色次元を、あるいは黒色次元より前に緑色次元というように、別の最小パス上でのルーティングを許可することによって所与グループ内のホップ数を減少させる。その結果、特定の状況においてはネットワークの重大な輻輳が発生する可能性があるため、これはグローバルトラフィックが特に均等に分布している場合を除き望ましくないかもしれない。

【 0 0 8 1 】

決定論的な最小非ハッシュ化ルーティングは、1つの決定論的な最小パスを全トラフィックに用いるが、これはパケット順序付けを提供する一方で、利用可能なパス間への優れた帯域幅や負荷分布を提供しない。このようなルーティングは、制御メッセージやレイテンシが問題になるメッセージといった、頻繁でないまたは小サイズのメッセージに使用できる。

【 0 0 8 2 】

順序付けが不要な場合には、適応型ルーティングをデフォルトルーティングタイプとして使用できる。パケットは、最小でのルーティングを試みるが、ネットワーク輻輳を回避するために、グループ内またはグループ間の非最小パスが使用されてもよい。或る実施形態では、適応型ルーティングは、ルーティングの選択を考慮するために、2つ以上の最小ポートおよび2つ以上の非最小ポートを提供するルーティングテーブルを使用することによって行われる。各ノードについて輻輳値が計算されることによって、または図9に示すルータタイルのようなルータ内のタイルの数が計算されることによって、同ルータ内の別のタイルに分布される。この例では適応型ルーティングアルゴリズムは、利用可能な2つの最小パスと2つの非最小パスを考慮し、この中から、輻輳値に基づき、また任意で様々な構成したバイアスに基づき選択を行う。

【 0 0 8 3 】

さらなる実施形態では、下流ポート輻輳、推定される遠端リンク輻輳、近端リンク輻輳のような要素から、ポート輻輳値を導出する。特定の例では、2ビットの下流ポート輻輳情報は、ルータチップ内において各タイルに対応した外部チャネル上で伝播され、定期的に更新される。これらのビットは、送信ルータチップにおいて、チップ上の下流ポートの輻輳のビューを組合せることによって生成される。この2ビットの輻輳値に組み合わせられた下流ポートは、各タイルにおけるMMR構成可能マスクを介して選択される。これらの下流ポートの輻輳値を合計し、3つのプログラム可能な閾値と比較する。合計が最高閾値よりも大きい場合は、輻輳は2'b11である。合計が最高閾値未満であり、中間閾値よりも大きい場合は、輻輳は2'b10である。合計が中間閾値未満であり、最低閾値よりも大きい場合は、輻輳は2'b01である。あるいは、合計が最低閾値未満の場合は、輻輳は2'b00である。

10

20

30

40

50

【 0 0 8 4 】

チャンネルの受信側にて、この2ビット値は、4ビット幅の下流輻輳リマッピングテーブルによって4つのエントリ内にインデックスされることによって、4ビット値にマッピングされる。推定された遠端リンク輻輳の計算は、過去にチャンネルラウンドリップレイテンシよりも長い距離で送信され、未だ確認応答されていないフリットの数を追跡し、フリット送信とその確認応答受信との相対値を調整することによって行われる。この計算を行うために使用するメカニズムは、5ビット幅、32エントリの深さの遅延チェーンである。MMR構成可能なサイクル数(1~31)については、ルータは、この遅延チェーンの末尾位置に送信されたフリットの数をカウントする。この遅延の後に、全ての値がシフトされる。予想されるチャンネル上の未だ確認応答されていないフリット(送信済みで、確認応答が予想されるもの)の総数は、このチェーン内の値の合計である。この値を、未だ確認応答されていないクレジットカウントと比較する。未だ確認応答されていないクレジットの総数から、予想されるチャンネル上のフリットを引くと、遠隔入力キューに記憶されるフリットの推定数が得られる。

10

【 0 0 8 5 】

推定される遠端輻輳は、10ビット数として計算される。この数値をマッピングテーブルに従って4ビットインデックスに変換し、次にこの4ビット数を、16エントリ遠端輻輳リマッピングテーブル内にインデックスすることによって、別のプログラム可能な4ビット値にリマッピングする。

20

【 0 0 8 6 】

近端リンク輻輳は、列バッファ内にキューイングされ、リンク上で送信されるのを待っている状態のフリットを合計することによって計算される。この合計も10ビット値であり、マッピングテーブルに従って4ビット値に変換される。次にこの4ビット数を、16エントリ近端輻輳リマッピングテーブル内にインデックスすることによって、別のプログラム可能な4ビット値にリマッピングする。

20

【 0 0 8 7 】

このリマッピングされた4ビットの下流ポート輻輳値と、リマッピングされた4ビットの遠端リンク輻輳値と、さらにリマッピングされた4ビットの近端リンク輻輳値とを互いに組み合わせることによって、1つの4ビット輻輳値をタイル毎に生成する。この組み合わせは、4ビット3入力符号なし飽和加算として行う。この4ビット輻輳値がチップ上の全ての他のタイルへと伝播されることによって、タイルが適応型のインフォームドチョイスを行うことが補助される。

30

【 0 0 8 8 】

チップ上のnタイルの各々からチップ上の他の全てのタイルに「リンクアライブ(link alive)」信号は、配信される。nタイルの各々に配信されるこのリンクアライブ信号は、該当するタイルと接続しているルータとの間にシリアルリンクが確立しているか否かを示す。リンクがアライブ状態でないポートは、ポート選択の観点から無効であると考慮される。これによって、ルータは、最近失敗し且つまだソフトウェアによってルーティングテーブルから除去されていないリンクを適応的に回避できるようになる。

40

【 0 0 8 9 】

リンクアライブ信号は、全てのネットワークタイルを接続する2線式シリアルチェーンを介してルータ周囲に伝播される。各タイルは、そのリンク状態情報を、適切なビットタイミングでシリアルチェーン上に配置する。輻輳論理に示されている全てのポートが無効である場合には、そのパケットは破棄される。この場合では、紛失パケットにタイムアウトを設定するか否かはエンドポイントハードウェア次第であり、また、再送信するか、エラーを適切と扱うかはより高レベルのソフトウェア次第である。

【 0 0 9 0 】

各入力キューでは、2つの最小ポート候補と2つの非最小ポート候補の間で適応的な選択を行うために、配信輻輳値が使用される。これらの輻輳値を使用する前に、選択した2つの最小ポートおよび非最小ポート輻輳値にバイアス値を適用する。最初に、値の最も大

50

きな部分にゼロを2つプリペンドすることで、値を論理的に6ビット値にまで拡張する。適応型ルーティングの制御タイプ(適応型0、適応型1、適応型2、適応型3)を使用して、4つのエントリバイアステーブルの中から1組のバイアスを選択する。各エントリは、最小ポートおよび非最小ポート輻輳値の各々に達するまで左にどれくらいシフトすればよいのか決定する、2ビットシフト値の1対を有する。6ビットに拡張された輻輳値は、0ビット、1ビット、または2ビットシフトされる。このフィールドの符号化は、次のとおりである: 2' b 0 0 = 左に0ビットシフト($\times 1$)、2' b 0 1 = 左に1ビットシフト($\times 2$)、2' b 1 0 = 左に2ビットシフト($\times 4$)、2' b 1 1 = リザーブ。

【0091】

各バイアスMMRは、拡張された最小および非最小輻輳値6ビットに加えて、追加の6ビット値の1対を含んでいる。この追加は飽和加算として実行され、6ビット数が得られる。最も低い輻輳に対応したポートが選ばれる。最小ポートと非最小ポートが同値である場合、ルータは、最小ポートを優先する。非最小として提示された2つのポート間、または、最小として提示された2つのポート間が同値の場合には、選択は、任意であり、あらゆる適切な方法で行われる。

【0092】

(Dragonflyネットワークにおけるテーブル駆動型ルーティングメカニズム)
ここで挙げるルーティング例では、パケットまたはメッセージをルーティングするために利用できるパスを決定するため、Dragonflyネットワーク構成にルーティングフレキシビリティを提供するために、様々なテーブルを使用する。グループ内またはグループ間にルーティングを提供するために、また、最小および非最小ルーティングパスのために、各種テーブルが存在する。

【0093】

ここに挙げた例示的なルータアーキテクチャにおけるルーティング構造は、次の4つの別個のテーブルセットに分割される。グローバル非最小(GN)テーブルセットと、グローバル最小(GM)テーブルと、ローカル非最小(LN)テーブルセットと、ローカル最小(LM)テーブルとである。この特定の例の論理フローを図12に示す。

【0094】

グローバルテーブルは、現在のグループがターゲットグループでない場合に、どのように遠隔グループにルーティングするか決定するために使用される。これらのテーブルは、ローカルグループから出るための出口である特定の光学ポートにルーティングするために使用される。ローカルテーブルは、現在のグループ内の特定のルータチップにルーティングするために使用される。これらのテーブルは、ローカルルーティングのためグループ内での「アップ」または「ダウン」ルーティングに使用されたり、あるいは中間グループ内での「アップ」ルーティングに使用されたりする。最小テーブルは、最小のローカルまたはグローバル経路を指定する。これらの最小テーブルは、「ダウン」ルーティングする時、または、適応型ルーティングのケースでは、「アップ」過程で最小パスの使用を試みる時に使用される。非最小テーブルは、非最小パスを指定し、「アップ」ルーティングする時のみ使用される。非最小テーブルはまた、「アップ」ルーティングの停止時を決定するための「ルート(root)検出」メカニズムを提供する。

【0095】

グローバル非最小テーブルセットは、非最小トラフィックを中間グループにルーティングするために使用される。このテーブルセットは、「安全な」中間グループへと続くポートのリストを含んでおり、この「安全な」中間グループとは、他の全てのグループに接続している中間グループのことである。(健全なネットワークでは、全てのグループが安全である。特に健全なネットワークにおいては、トラフィックをターゲットグループに接続していない可能性のある中間グループに送信することを回避するように、テーブルをプログラムしなければならない。)このテーブルセットは、3つのテーブルで構成されている。第1テーブルは、現在の(ソース)グループから出るためには緑色次元のどのランクをトラバースするか選択する。第2テーブルは、黒色次元をトラバースするよう選択する。

10

20

30

40

50

第3テーブルは、現在のルーチップをオンにしておくために光学ポートを選択する。

【0096】

これらのテーブルは、固定優先順序で階層配列される。緑色次元テーブルは優先順位が1番高く、青色次元テーブルは一番低い。各テーブルには、Aries（ルーチ兼ファイアウォール）をオンにしておくための複数のポート番号が、または現在のテーブルがその優先順位を下げて、優先順位階層における次のテーブルを考慮すべきであることを示す特別値がリストされている。最下位優先順位（青色）テーブル上の特別値を参照した場合、エラー状態が生じる。各テーブルは128エントリで構成されており、エントリの各々は6ビットポート番号または特別値6'b11xxxxである。各テーブルは、16x8エントリで編成されており、8エントリ毎の各ブロックにつき7ビットのECCが付加されている。

10

【0097】

このテーブルは、システム内の他の全てのグループに安全にルーティングすることが可能な中間グループに確実に続いている他のルーチップまたは光学ポート番号への経路しが含まべきでない。このテーブルはさらに、非最小トラフィックをシステム内の複数のグループにかけてほぼ均等に分布させるメカニズムを提供する。各テーブルにはエントリが128個あるので、有効基数18次元を用いた場合でも、各ポートが7回または8回リストされ、その次元で2つのポート間に最大14.3%の不安定さが生じる。この不安定さは、グループ全体を通してテーブルを複数回コピーし、不安定なポートを変化させることによって最小化する。

20

【0098】

グローバル決定論的なルーティングの場合、このテーブルセットは、ターゲットtgtid（おそらくはローカルポート番号）を含むハッシュ値によって、またオプションでパケットヘッダ（パケットアドレスからのもの）からのハッシュフィールドによって、インデックスされる。各テーブルは、互いに異なるインデックスを得る。グローバル適応型ルーティングの場合は、テーブルから、それぞれが8エントリで構成されている複数ブロックのうちの1つは、ランダムに選択される。次にこの同じ8エントリ構成のブロックから第2のエントリは、ランダムに選択される。2つのポートを互いに、またグローバル最小テーブルからの2つのエントリと比較することによって、パケットをどのパスでルーティングするか決定する。

30

【0099】

pタイルにおける緑色テーブルは、一般に、8回リストされた15個の緑色ポートと、8個の特別値を有する。さらにpタイルでは、黒色テーブルは、約7回リストされた15個の黒色ポートを、特別値を含んだ約21のエントリと共に有する。青色テーブルは、それぞれ約13回リストされた光学ポートの各々を有する。

【0100】

緑色のnタイルポートは、一般に、緑色テーブル内のエントリ全てを特別値として有する。黒色および青色テーブルは、pタイルの場合と同じ比率で構成される。黒色nタイルポートは、一般に、緑色および黒色テーブル内の全てのエントリを特別値として有する。青色テーブルは、pタイルの場合と同じ比率で構成される。

40

【0101】

グローバル最小テーブルは、現在のグループからターゲットグループまでの直接パスを決定するために使用される。このテーブルは256エントリで構成され、各エントリのビット幅は81ビットである。各エントリは、全ポートセットと規制されたポートセットとの2つの部分に分割されている。全ポートセットは、8個の6ビットポートエントリと、3ビットのモジュロ指定子とで構成される。モジュロフィールドは、関連するエントリにおける有効なポートの総数を示す。モジュロ指定子は、モジュロ-1として符号化される。つまり、モジュロフィールド内に7の値があれば、8のモジュロ演算ということである。規制されたポートセットは、4個の6ビットポートと、2ビットのモジュロ指定子で構成されている。81ビットエントリの各々が8ビットのECCを有する。

50

【0102】

このテーブルは、ターゲットグループ番号で編成されている。各ターゲットグループは、システムのサイズに応じて、テーブル内の1、2、4、8、16、32、64、または128エントリの「ブロック」に対応している。241グループを持ったシステムは、テーブル内のブロック1個につき1エントリを有する（エントリのうち15個は使用されない）。65～128グループを有するシステムは、各ブロックにつき2つのエントリを使用する。33～64グループを有するシステムは4つのエントリを使用する。同様に続く。グループ番号、並びに0～7の追加のランダム（適応型ルーティング）またはハッシュ（決定論的なルーティング）ビットを使用して、テーブル内にインデックスを定義する。各エントリは、関連するターゲットグループに最小で接続するポートであって、ルーティングにおける現在のポイントから到達可能なAries（ルータ兼ファイアウォール）に続いているポートのリスト、または青色リンクを介してターゲットグループに直接続いているポートのリストを含む。

10

全ポートセットは、グループ内で（pタイルまたは光学nタイルのどちらかにて）他のグループへの最小ルーティングを丁度開始した時に使用されたり、あるいは中間グループ内において非最小でルーティングし且つローカル非最小テーブルにルート（root）が検出された際に任意のタイルにおいて使用されたりする（以下を参照）。テーブルのこのサイドには、インデックスが指定するグループに最小接続している利用可能な光学ポートまでの利用可能なパスが全て挙げられている。規制されたポートセットは、全ポートセットテーブルについて言及するルート（root）検出の場合とルート（root）注入の場合を除いて、グループ内でのルーティングに使用される。テーブルのこの半分は、最小ルーティングを行っているとは仮定した場合に、グループネットワーク内の現在地点から正当であるネットワーク内のパスのみを示す。

20

【0103】

規制されたポートリストの重要な目的は、パケットがその出発点の方向に戻らないようにすることである。緑色ポートでは、規制されたテーブルエントリは、通常、黒色ポートと青色ポートしかリストしないはずである。黒色ポートでは、規制されたテーブルエントリは、通常、青色ポートしかリストしないはずである。

【0104】

規制されたセットにリストされた全てのポートが無効である場合には、このことは、パケットが正当な最小パスから分岐したことを、適応型ルーティング論理に示している。この場合、適応型ルーティング論理は、非最小選択肢から1つを選ぶ。（これは、決定論的または最小にルーティングしたトラフィックでは絶対に発生しないはずである。なぜならテーブルは一貫した方法で書き込まれているはずであり、パケットが宛先にルーティングできないポイントに到着することはあり得ないからである。これが発生した場合には、ルータがエラーフラグを立て、そのパケットを破棄する。）

30

タイル内でルーティングされた正当な規制されたポートがない場合には、mod値は、任意の値に設定されうる。経路テーブルは、グループ番号に関連した全てのエントリに特別値6'b11xxxxを含んでいなければならない。正当な経路が1つしかない場合には、ポートリストは、少なくとも2回リストされている正当な経路と、これに合わせて2またはそれ以上に設定したmod値とを含んでいなければならない。

40

【0105】

決定論的なルーティングでは、関連するインデックス内の有効なエントリの数によってハッシュのモジュロを計算することで、全ポートセットまたは規制されたポートセットにおける有効なエントリのうちの1つが選択される。上述のケースと同様に、適応型ルーティングは、乱数と、N-1の第2モジュロとのモジュロを計算することによって第1の数字+1に可算し、テーブル中に第2のランダムであるが固有のエントリのオフセットを取得することを除いて、テーブルから2つのエントリを選択する。

【0106】

50

ローカル非最小テーブルセットは、ローカルグループ内のルーチップを、このグループ内の非最小ルーティングのルート (r o o t) として使用する目的で選ぶために使用される。このテーブルは、ソースグループおよびターゲットグループが同一である場合に、非最小ルーティングのために使用される。さらにテーブルは、中間グループ内での非最小ルーティングのためにも使用される。このテーブルセットは、青色テーブルが存在しないことを除けば、グローバル非最小テーブルと同様の構造である。

【 0 1 0 7 】

ローカル非最小テーブルは、適応型ルーティングの場合にはランダムにインデックスされ、また非最小決定論的なルーティングの場合にはハッシュによってインデックスされる。グローバル非最小テーブルと類似し、適応型ルーティングの場合には、このテーブルから2つのエントリが生成され、比較される。設計におけるRAMマクロ総数を低減するために、これらのテーブルをRAM内のグローバル非最小テーブルと物理的に組み合わせる。

10

【 0 1 0 8 】

このテーブルは、このタイルから到達可能な A r i e s (ルータ兼ファイアウォール) をリストしている。これら A r i e s は、ローカル非最小ルーティングに安全に使用できる。健全なネットワークでは、pタイルと青色(光学)タイルは、グループ内の全ての A r i e s (ルータ兼ファイアウォール) をほぼ均等にリストしなければならない。緑色テーブルのエントリのほぼ 1 5 / 1 6 には緑色ポートがリストされるべきであり、また約 1 / 6 は、緑色次元が既に条件を満たしており、黒色テーブルを使用すべきであるということを示す特別値を含んでいる必要がある。同様に黒色テーブルのエントリの約 5 / 6 は、黒色ポートをリストすべきであり、また約 1 / 6 は、黒色次元が条件を既に満たしたことを示す特別値を含んでいる必要がある。緑色テーブルと黒色テーブルの両方における特別値は、ルート (r o o t) に到達したことと(「ルート (r o o t) 検出」)、パケットをこの地点からダウンルーティングすべきであることを示す。

20

【 0 1 0 9 】

緑色タイルは、特別値(緑色次元の条件が満たされたことを示す)で緑色テーブルを充填すべきであり、到達可能な6個の A r i e s (ルータ兼ファイアウォール) を(自身を含み、特別値を使用する。)黒色テーブルに均等にリストすべきである。黒色タイルは、緑色テーブルと黒色テーブルの両方を特別ルート (r o o t) 検出値によって充填すべきである。pタイルと光学タイルとは、全テーブルセットを必要とすべきである。nタイルは緑色テーブルなしでも技術上問題はないが、しかし、ここで示すルーチップの例は、フレキシビリティを持たせるためにnタイルを実現している。

30

【 0 1 1 0 】

ローカル最小テーブルは、ターゲットグループ内での最小ルーティング(「ダウンルーティング」)に使用され、さらにターゲットグループ内で適応的に「アップルーティング」を行う際にも使用される。このテーブルは、128エントリを有する。各エントリは52ビット幅であり、8個の6ビットポート番号と、「分岐」ビットと、テーブルのこのラインで有効なエントリの数を示す m o d 値とで構成されている。ターゲットグループ内のパスが最小パスから分岐しているため、そのパスは適応型アップルーティングを行う時には最小パスとして使用できず、したがってダウンルーティングにしか使用できないことを、分岐したビットは示す。これは、規制されたセットにおける全てのポートが無効であるグローバル最小テーブルのケースと類似している。

40

【 0 1 1 1 】

このテーブルは、グループ内の「ターゲット」 A r i e s 番号で編成されている。各ローカル A r i e s 番号は、グループのサイズに応じて、テーブル内の1、2、4、8、または16エントリで構成されているブロックに対応している。65~128個の A r i e s (ルータ兼ファイアウォール) を持ったグループは、1つのローカル A r i e s 番号につき1エントリのブロックサイズを使用する。33~64個の A r i e s (ルータ兼ファイアウォール) を有するサイズのグループは、2のブロックサイズを使用する、などであ

50

る。ローカル A r i e s 番号、並びに 0 ~ 4 の追加の乱数 (適応型ルーティングの場合) またはハッシュ (決定論的なルーティング) ビットは、テーブルにインデックスを定義するために使用される。各エントリは、関連するローカル A r i e s へと続くポートのリストを含む。

【 0 1 1 2 】

決定論的なルーティングでは、関連インデックス内の有効なエントリの数によってハッシュのモジュロを計算することによって、テーブル内の有効なエントリの中から 1 つは、選択される。上述のケースと同様に、適応型ルーティングは、乱数と、N - 1 の第 2 モジュロとのモジュロを計算することによって第 1 の数字 + 1 に可算し、テーブル中に第 2 のランダムであるが固有のエントリのオフセットを取得することを除いて、テーブルから 2

10

【 0 1 1 3 】

グローバル非最小テーブルは、ソースグループのみにおいて、別のグループに向かうトラフィックのために使用される。グローバル非最小テーブルとローカル非最小テーブルとは、決して同時に使用されることはない。そのため必要な R A M 総数を減らすために、グローバル非最小緑色テーブルは、ローカル非最小緑色テーブルと同じ R A M に記憶される。グローバル非最小黒色テーブルは、ローカル非最小黒色テーブルと同じ R A M に記憶される。グローバルテーブルは、これら 2 つの R A M それぞれの下方インデックス値部分に記憶される。

【 0 1 1 4 】

(結論)

上述した例は、D r a g o n f l y ネットワークにおけるルーティングが、ネットワーク輻輳やトラフィックタイプのような要素に基づきネットワークパスを選択することが可能な適応型ルーティングと、最小および非最小ルーティングや、ローカルルーティングおよびグローバルルーティングを含む様々なルーティング用のルーティングテーブルとを使用して、どのように向上させられるか説明する。

20

【 0 1 1 5 】

適応型ルーティングは、輻輳リンクまたはダウンしたリンクに基づき複数の正当な経路を選択し、また、チャンネル上の輻輳を明快に通信させることによって向上したルーティングパフォーマンスと許容性を提供する、デッドロックを回避するルーティングを提供する

30

【 0 1 1 6 】

輻輳情報は、出力キュー内のメッセージ数をカウントすることや、伝送中のクレジットまたはメッセージのような要因から受信側バッファ輻輳推定を確立することといった要素から予想される、次のリンク輻輳に基づく。ノードは潜在的な受信側ノードに、平均的な「次のリンク」出力輻輳について問い合わせることができるため、輻輳リンクやダウンしたリンクの回避に基づきルーティング決定を行えるようになる。これ以外にもさらに、例えば、ルーティングパスを選ぶ際に、決定論的なハッシュまたは乱数を使用してトラフィックを拡散させるといった特徴が複数提供され、これらは輻輳を回避するためにトラフィックを拡散させる上で役立つ。

40

【 0 1 1 7 】

1 例では、ルーティングの選択はテーブルによって提示され、ネットワークの構成および状態に応じて、特定の経路に、あるいは、最小経路または非最小経路に偏る可能性がある。例えば、経路の選択が最高の効率のデフォルトによって最小ルーティングに偏るが、この偏りが、追加のトラフィックを任意または不必要に受信することから特定のネットワークリンクを保護するために、非最小ルーティングへの偏りに切り替わることもある。さらなる例では、ルーティングテーブルは、ローカルルーティングテーブルおよびグローバ

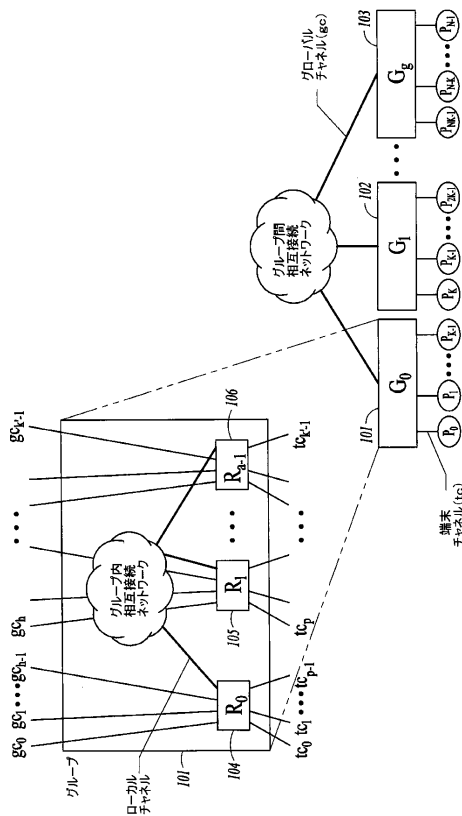
50

ルルーティングテーブルを備えたテーブルと、最小パスおよび非最小パスとを含む。

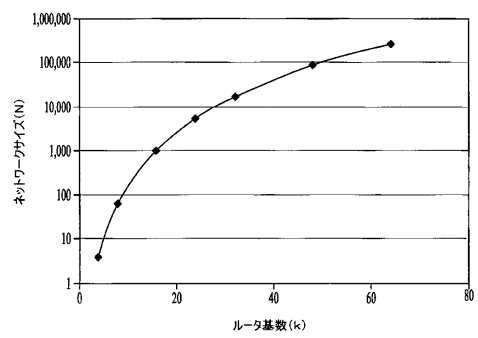
【0118】

特定の実施形態をここで例証および説明したが、当業者は、同じ目的を達成すると推定されるあらゆる配列はここに示した特定の実施形態の代用となることを理解するであろう。本出願は、ここで説明した本発明の例示的な実施形態のあらゆる改造または応用を包括することを意図する。本発明は、特許請求の範囲、およびその均等物の全範囲によってのみ限定されることが意図されている。

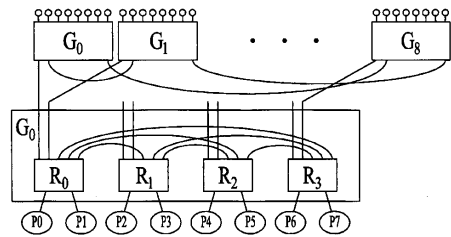
【図1】



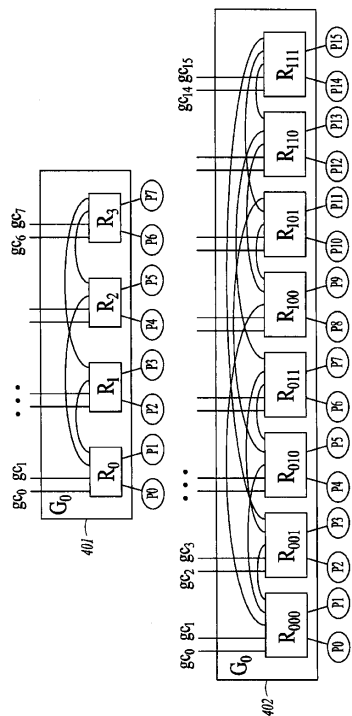
【図2】



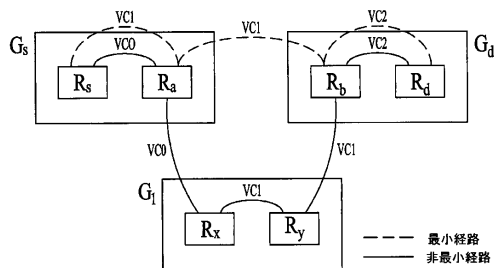
【図3】



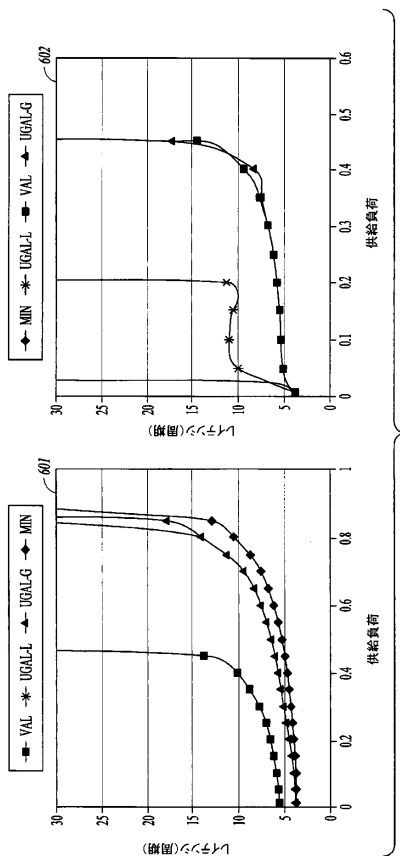
【 図 4 】



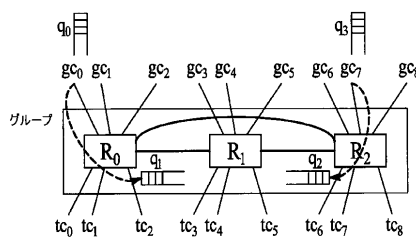
【 図 5 】



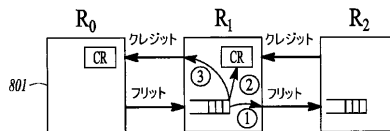
【 図 6 】



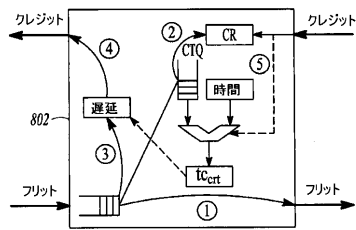
【 図 7 】



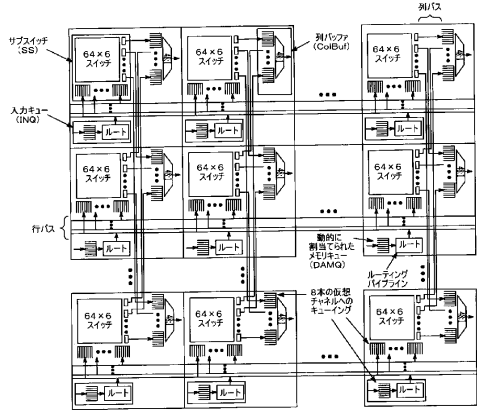
【 図 8 A 】



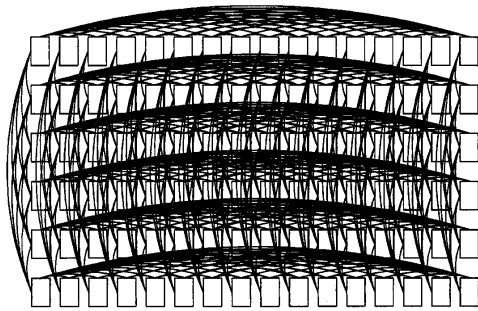
【 図 8 B 】



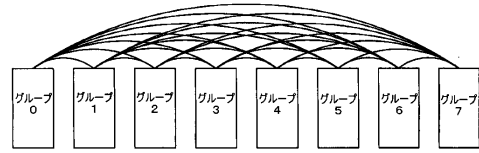
【図9】



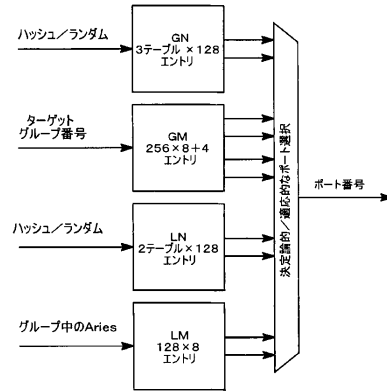
【図10】



【図11】



【図12】



フロントページの続き

- (72)発明者 スティーブ スコット
アメリカ合衆国 98164 ワシントン州 シアトル フィフス アベニュー 901
- (72)発明者 アルバート チェン
アメリカ合衆国 98164 ワシントン州 シアトル フィフス アベニュー 901
- (72)発明者 ロバート アルバーソン
アメリカ合衆国 98164 ワシントン州 シアトル フィフス アベニュー 901

審査官 宮島 郁美

- (56)参考文献 米国特許出願公開第2010/0049942(US, A1)
特開2006-012112(JP, A)
米国特許出願公開第2008/0285562(US, A1)
米国特許出願公開第2008/0123679(US, A1)

- (58)調査した分野(Int.Cl., DB名)
H04L12/00-12/26, 12/50-12/955
G06F15/16-15/177