

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6254617号
(P6254617)

(45) 発行日 平成29年12月27日(2017.12.27)

(24) 登録日 平成29年12月8日(2017.12.8)

(51) Int.Cl.	F I
GO6F 15/173 (2006.01)	GO6F 15/173 683D
HO4L 12/721 (2013.01)	HO4L 12/721 Z
	GO6F 15/173 660B

請求項の数 19 (全 15 頁)

(21) 出願番号	特願2015-558103 (P2015-558103)	(73) 特許権者	591016172
(86) (22) 出願日	平成26年2月12日(2014.2.12)		アドバンスト・マイクロ・ディバイス・
(65) 公表番号	特表2016-513321 (P2016-513321A)		インコーポレイテッド
(43) 公表日	平成28年5月12日(2016.5.12)		ADVANCED MICRO DEVI
(86) 国際出願番号	PCT/US2014/016045		CES INCORPORATED
(87) 国際公開番号	W02014/127012		アメリカ合衆国、94088-3453
(87) 国際公開日	平成26年8月21日(2014.8.21)		カリフォルニア州、サニibel、ピー・
審査請求日	平成29年2月6日(2017.2.6)		オウ・ボックス・3453、ワン・エイ・
(31) 優先権主張番号	13/766, 115		エム・ディ・プレイス、メイル・ストップ
(32) 優先日	平成25年2月13日(2013.2.13)		・68 (番地なし)
(33) 優先権主張国	米国 (US)	(74) 代理人	100108833
早期審査対象出願			弁理士 早川 裕司
		(74) 代理人	100111615
			弁理士 佐野 良太

最終頁に続く

(54) 【発明の名称】 改良3Dトラス

(57) 【特許請求の範囲】

【請求項1】

それぞれ複数のノードを含む複数のトラスと、
前記複数のトラスの各々の少なくとも1つのノードに接続されたホストと、
を備え、
各ノードは、データトラフィックを送受信するように構成された計算装置であって、リンクによって同一トラス内の1つ以上の他のノードに接続されており、
前記ホストは、
第1の基準を有するデータトラフィックを伝搬するために第1のトラスを選択し、第2の基準を有するデータトラフィックを伝搬するために第2のトラスを選択することと

10

、
前記第1のデータトラフィックを、選択した第1のトラス内のノードへ配送することであって、1つのトラス内のノード間のデータトラフィックは、他のトラス内のデータトラフィックと混ざらない、ことと、

前記ホストを前記第1のトラスに接続する第1のノードから受信したデータトラフィックを、前記ホストを前記第2のトラスに接続する第2のノードへ再配送することと、
を行うように構成されている、
並列トラスネットワークインターコネクトシステム。

【請求項2】

前記並列トラスネットワークインターコネクトシステム内の各トラスは、三次元ト

20

ーラスである、請求項 1 に記載の並列トーラスネットワークインターコネクトシステム。

【請求項 3】

前記ホストは、前記第 1 及び前記第 2 のデータトラフィックの各々のサービス品質 (QoS) に基づいて、前記第 1 及び前記第 2 のトーラスを選択するように構成されている、請求項 1 に記載の並列トーラスネットワークインターコネクトシステム。

【請求項 4】

前記ホストは、前記トーラスの何れかにおける輻輳に基づいて、前記第 1 及び前記第 2 のトーラスを選択するように構成されている、請求項 1 に記載の並列トーラスネットワークインターコネクトシステム。

【請求項 5】

前記ホストは、前記第 1 及び前記第 2 のデータトラフィックのタイプに基づいて、前記第 1 及び前記第 2 のトーラスを選択するように構成されている、請求項 1 に記載の並列トーラスネットワークインターコネクトシステム。

【請求項 6】

前記ホストは、前記第 1 及び前記第 2 のデータトラフィックに関連付けられたサービスクラス (COS) に基づいて、前記第 1 及び前記第 2 のトーラスを選択するように構成されている、請求項 1 に記載の並列トーラスネットワークインターコネクトシステム。

【請求項 7】

前記ホストは、前記第 1 及び前記第 2 のデータトラフィック及び前記第 1 及び前記第 2 のトーラスの各々のセキュリティレベルに基づいて、前記第 1 及び前記第 2 のトーラスを選択するように構成されている、請求項 1 に記載の並列トーラスネットワークインターコネクトシステム。

【請求項 8】

それぞれ複数のノードを含む複数のトーラスを有する並列トーラスネットワークインターコネクトシステムにおいて、データトラフィックを伝搬する方法であって、

前記複数のトーラスの各々の少なくとも 1 つのノードに接続されたホストにて、第 1 の基準を有する第 1 のデータトラフィックと、第 2 の基準を有する第 2 のデータトラフィックとを受信することと、

前記第 1 のデータトラフィックを受信するために、前記ホストによって第 1 のトーラス内のノードを選択し、前記第 2 のデータトラフィックを受信するために、前記ホストによって第 2 のトーラス内のノードを選択することと、

前記ホストによって、前記第 1 のデータトラフィックを、前記第 1 のトーラス内で選択したノードへ配送し、前記第 2 のデータトラフィックを、前記第 2 のトーラス内で選択したノードへ配送することであって、1 つのトーラス内のノード間のデータトラフィックは、他のトーラス内のデータトラフィックと混ざらない、ことと、

前記ホストを前記第 1 のトーラスに接続する前記第 1 のトーラス内の前記ノードから受信したデータトラフィックを、前記ホストを前記第 2 のトーラスに接続する前記第 2 のトーラス内の前記ノードへ再配送することと、を備える、

方法。

【請求項 9】

前記並列トーラスネットワークインターコネクト内の各トーラスは、三次元トーラスである、請求項 8 に記載の方法。

【請求項 10】

前記選択することは、前記第 1 及び前記第 2 のデータトラフィックの各々のサービス品質 (QoS) に基づいて、前記第 1 及び前記第 2 のノードを選択することを含む、請求項 8 に記載の方法。

【請求項 11】

前記選択することは、前記第 1 及び前記第 2 のデータトラフィックの各々のサービスクラス (COS) に基づいて、前記第 1 及び前記第 2 のノードを選択することを含む、請求項 8 に記載の方法。

10

20

30

40

50

【請求項 1 2】

前記選択することは、前記トラスの何れかにおける輻輳に基づいて、前記第 1 及び前記第 2 のノードを選択することを含み、請求項 8 に記載の方法。

【請求項 1 3】

前記選択することは、前記第 1 及び前記第 2 のデータトラフィック及び前記第 1 及び前記第 2 のトラスのセキュリティレベルに基づいて、前記第 1 及び前記第 2 のノードを選択するように構成されている、請求項 8 に記載の方法。

【請求項 1 4】

それぞれ複数のノードを含む複数のトラスを有する並列トラスネットワークインターコネクトシステムにおいてデータトラフィックをコンピュータに伝搬させるためのプログラムを記憶するコンピュータ可読記憶媒体であって、

前記プログラムは、前記コンピュータに、

前記複数のトラスの各々の少なくとも 1 つのノードに接続されたホストにて、第 1 の基準を有する第 1 のデータトラフィックと、第 2 の基準を有する第 2 のデータトラフィックとを受信するための受信機能と、

前記第 1 のデータトラフィックを受信するために、前記ホストによって第 1 のトラス内のノードを選択し、前記第 2 のデータトラフィックを受信するために、前記ホストによって第 2 のトラス内のノードを選択するための選択機能と、

前記ホストによって、前記第 1 のデータトラフィックを、前記第 1 のトラス内で選択したノードへ配送し、前記第 2 のデータトラフィックを、前記第 2 のトラス内で選択したノードへ配送するための配送機能であって、1 つのトラス内のノード間のデータトラフィックは、他のトラス内のデータトラフィックと混ざらない、配送機能と、

前記ホストを前記第 1 のトラスに接続する前記第 1 のトラス内の前記ノードから受信したデータトラフィックを、前記ホストを前記第 2 のトラスに接続する前記第 2 のトラス内の前記ノードへ再配送するための再配送機能と、を実現させる、

コンピュータ可読記憶媒体。

【請求項 1 5】

前記選択機能は、前記第 1 及び前記第 2 のデータトラフィックの各々のサービス品質に基づいて、前記各ノードを選択することを含み、請求項 1 4 に記載のコンピュータ可読記憶媒体。

【請求項 1 6】

前記選択機能は、前記トラスの何れかにおける輻輳に基づいて、前記各ノードを選択することを含み、請求項 1 4 に記載のコンピュータ可読記憶媒体。

【請求項 1 7】

前記選択機能は、前記第 1 及び前記第 2 のデータトラフィックの各々のタイプに基づいて、前記各ノードを選択することを含み、請求項 1 4 に記載のコンピュータ可読記憶媒体。

【請求項 1 8】

前記選択機能は、前記第 1 及び前記第 2 のデータトラフィックの各々に関連付けられたサービスクラス (C o S) に基づいて、前記各ノードを選択することを含み、請求項 1 4 に記載のコンピュータ可読記憶媒体。

【請求項 1 9】

前記選択機能は、前記第 1 及び前記第 2 のデータトラフィック及び前記第 1 及び前記第 2 のトラスの各々のセキュリティレベルに基づいて、前記各ノードを選択することを含み、請求項 1 4 に記載のコンピュータ可読記憶媒体。

【発明の詳細な説明】**【技術分野】****【0 0 0 1】**

実施形態は、概してネットワークトラフィックの最適化に関し、より具体的には並列トラスインターコネクトを用いたネットワークトラフィックの最適化に関する。

10

20

30

40

50

【背景技術】

【0002】

トーラスは、並列コンピュータネットワークにおいて処理ノードを接続するためのネットワークポロジである。トーラスは、N次元のフィールドアレイに配置されてもよく、そこでは、処理ノード（ノードとも呼ばれる）は、リンクを用いて最寄の処理ノードに接続されている。

【0003】

従来のトーラスネットワークポロジでは、トーラスインターコネクトの帯域幅に限りがある。帯域幅に限りがあるのは、トーラス内のノードのサブセットに接続する各ホストが、帯域幅のごく一部を受け取るからである。これにより、ホストを、より多くのノードを介してトーラスファブリックに接続すると、トーラス内の他のノードと、これらのノードに接続された他のホストと、から帯域幅が奪われる。

10

【0004】

トーラスインターコネクトのノード間では、リンクによってデータトラフィックが伝搬される。トーラスインターコネクト内のリンクが輻輳したり故障したりすると、当該リンクを用いるノード間のデータトラフィックの経路が変更される。この経路変更は、トーラスネットワークにおけるトラフィックの待ち時間に影響を及ぼす。例えば、経路を変更されたトラフィックが、変更された経路に沿って終点ノードに到達する場合には、より長い時間がかかる可能性がある。また、データトラフィックの経路が別の経路に変更されると、当該別の経路で輻輳が増大して、当該別の経路を流れるように予定されていたトラフィックにも影響が及ぶ。

20

【発明の概要】

【課題を解決するための手段】

【0005】

データトラフィックの流れを最適化するシステム及び方法が提供される。並列トーラスインターコネクトにおいて複数のトーラスが接続されている。各トーラスは、複数のノードを含む。トーラス内のノードは、リンクを用いて相互に接続されている。ネットワーク内のホストはノードのサブセットに接続されており、サブセット内のノードは、異なるトーラスに関連付けられている。ホストは、ノードのサブセットからノードを選択することによって、並列トーラスインターコネクトへパケットを送信する。パケットは、トーラス内で前記ノードと複数のノードとを結ぶリンクを用いて送信されるが、複数のトーラス間では送信されない。

30

【0006】

添付の図面を参照しながら、実施形態のさらなる特徴及び利点、並びに、種々の実施形態の構造及び動作を詳述する。実施形態は、本明細書で説明する特定の実施形態に限定されないことに留意されたい。かかる実施形態は、専ら例示目的で本明細書に示されている。本明細書に記載された教示に基づくさらなる実施形態が、当業者に明らかになるであろう。

【0007】

本書に組み込まれて本明細書の一部を構成する添付の図面は、実施形態を示すものであって、本明細書の記載とともに、実施形態の原理を説明して、当業者による実施形態の製造及び実施を可能にするものである。図面を参照しながら様々な実施形態を以下に説明するが、全体を通じて、同様の要素には同様の参照番号が使われている。

40

【図面の簡単な説明】

【0008】

【図1A】一実施形態による三次元トーラスのブロック図である。

【図1B】一実施形態による並列トーラスインターコネクトのブロック図である。

【図2】一実施形態による並列トーラスインターコネクトを通じてデータトラフィックを伝搬する方法のフローチャートである。

【図3】一実施形態による並列トーラスインターコネクトにおけるノードの例示的な物理

50

的配置を示すブロック図である。

【図4】実施形態を実装し得るコンピュータシステムのブロック図である。

【発明を実施するための形態】

【0009】

添付の図面を参照しながら実施形態を説明する。全体的に、参照番号の左端の数字は、要素が最初に現れる図面を表している。

【0010】

以降の詳細な説明では、「one embodiment（一実施形態）」、「an embodiment（一実施形態）」、「an example embodiment（例示的な実施形態）」等は、説明される実施形態が、特定の特徵、構造または特性を含み得ることを意味するが、全ての実施形態が、必ずしも当該特定の特徵、構造または特性を含むとは限らない。また、かかる語句が、必ずしも同じ実施形態を指すとは限らない。さらに、一実施形態に関して特定の特徵、構造または特性が説明される場合には、明示されているか否かに関わらず、他の実施形態においてかかる特徵、構造、または特性が当業者の知識の範囲内で実現され则认为られる。

【0011】

「embodiments（複数の実施形態）」という語句は、全ての実施形態が、記載された特徵、利点または動作モードを含むことを要していない。本開示の範囲から逸脱することなく、代替の実施形態を案出可能である。また、本発明に関連する詳細が不明瞭にならないように、本開示における周知の要素が詳しく説明されない場合があり、省かれる場合もある。さらに、本書で用いられている用語は、専ら特定の実施形態を説明することを目的としており、本開示を制限するものではない。例えば、本書で用いられている単数形「a（或る）」、「an（或る）」及び「the（その）」は、複数形を含まないことが文脈から明らかな場合を除いて、複数形も含む。さらにまた、本書で用いられている用語「comprise（備える）」、「comprising（備える）」、「includes（含む）」及び/または「including（含む）」は、記載された特徵、整数、ステップ、操作、要素及び/又はコンポーネントの存在を明示するものであるが、1つ以上の他の特徵、整数、ステップ、操作、要素、コンポーネント及び/又はこれらの集合の存在や追加を否定するものではないことが理解されよう。

【0012】

従来のトーラスインターコネクトは、単一トーラスネットワークトポロジとして実装されている。単一トーラスネットワークトポロジには、いくつかの限界がある。第一に、単一トーラスネットワークトポロジでは、帯域幅に限りがある。例えば、各ホストは、トーラス内のいくつかのノードに接続する。各トーラスが有限の合計帯域幅を有しているので、各接続は、トーラス内の合計利用帯域幅のごく一部を受け取る。ノード間のデータトラフィックが、トーラス内で割り当てられた帯域幅よりも多くの帯域幅を要する場合には、2つのノードが、複数のリンクを用いて互いに接続される場合がある。しかし、この場合には、トーラス内の他のノードで使用可能な帯域幅が減ってしまう。また、トーラス内の複数のノードに接続するために追加のリンクがデータトラフィックの発信元に割り当てられると、他のホストがトーラスへ接続するために使用可能な帯域幅も減る。

【0013】

第二に、従来のトーラスインターコネクトでは、ノード間のリンクの故障や輻輳がネットワークに悪影響を及ぼす。例えば、従来のトーラスインターコネクトのノード間には、様々な種類のデータトラフィックが流れている。リンク又はノードの何れかで故障や輻輳が生じると、そのリンクを流れる全てのデータトラフィックに対して故障や輻輳の影響が及ぶ場合がある。輻輳や故障を是正するために、従来のトーラスインターコネクトでは、データトラフィックを他のノードに差し向けることがある。これにより、ネットワーク全体で遅延が生じる場合がある。この問題は、データトラフィックが、リンクやノードの故障や輻輳のために満たすことのできないサービス品質（QoS）やサービスクラス（CoS）に関連付けられる場合に、とりわけ顕著となる。

【 0 0 1 4 】

従来のトーラスインターコネクトにおいてリンクの故障や輻輳を是正する初歩的な方法として、冗長リンクの形成がある。輻輳、サービス品質またはサービスクラスの劣化の発生確率が冗長リンクによって低減することもあるが、上述したように、従来のトーラスインターコネクトでは、冗長リンクによって帯域幅が減少することもある。

【 0 0 1 5 】

後述する並列トーラスインターコネクトは、上述した制限に対する解決策である。

【 0 0 1 6 】

図 1 A は、一実施形態による三次元 (3 D) トーラスのブロック図 1 0 0 A である。ブロック図 1 0 0 A は、並列トーラスインターコネクトに含まれ得るトーラス 1 0 2 を含んでいる。ブロック図 1 0 0 A の例示的なトーラス 1 0 2 は、それぞれ 3 つのノード 1 0 4 のリングで構成された 2 7 個のノード 1 0 4 を含んでいるが、実装は、この実施形態に限定されない。リングは、3 つの直交次元 (X、Y、Z) で形成され得る。一実施形態において、各ノード 1 0 4 は、3 つの異なるリングのメンバーであり、リングは各次元に 1 つずつある。図 1 では、各ノード 1 0 4 の相対的位置がタプル (x、y、z) で示されており、x、y、z は、X、Y、Z 座標軸におけるノード 1 0 4 の論理的位置を表す。また、各ノード 1 0 4 は、接続すなわちリンク 1 0 6 によって 6 つの隣接ノード 1 0 4 に接続されている。一実施形態において、リンク 1 0 6 は双方向接続であってよい。

【 0 0 1 7 】

一実施形態において、トーラス 1 0 2 はネットワークを表す。ネットワークは、データトラフィックを搬送し、サービスやアプリケーションへのアクセスを提供するネットワークであってよい。ネットワークは、ローカルエリアネットワーク (L A N)、メトロポリタンエリアネットワーク、及び / 又は、インターネット等のワイドエリアネットワーク (W A N) を含み得るが、これらに限定されない。

【 0 0 1 8 】

ノード 1 0 4 は、トーラス 1 0 2 における接続点である。一実施形態において、ノード 1 0 4 は、リンク 1 0 6 を介してデータトラフィックを送信、受信及び転送できる計算装置であってよい。例示的な計算装置については、図 4 にて詳述する。ノード 1 0 4 は、クライアント、サーバー及びピアノードを含むネットワークの一部であってよい。一実施形態において、ピアノードは、クライアントノードまたはサーバーノードであってよい。非限定的な実施例において、クライアントは、上述した計算装置であって、ネットワーク上でデータを要求し、受信したデータを処理し、表示する。サーバーは、上述した電子機器であって、データを蓄積し、クライアントへ配送する。

【 0 0 1 9 】

一実施形態において、トーラス 1 0 2 は、メッシュとして組み立てられてもよい。メッシュにおいて、ノード 1 0 4 は、自身のデータを捕捉し、頒布するほか、他のノード 1 0 4 からのデータトラフィックを中継する。

【 0 0 2 0 】

図 1 A に示されたトーラス 1 0 2 は、X、Y、Z 座標空間内の 3 D 配列であるが、ノード 1 0 4 は、ネットワーク内の各ノード 1 0 4 の他ノード 1 0 4 に対する位置を示す論理的次元を表しており、必ずしも各ノード 1 0 4 の物理的配置を示す物理的次元を表すものではない。例えば、サーバーとして機能するトーラス 1 0 2 のネットワークトポロジは、ラックやバックプレーンの 1 つ以上の段に物理的に配置されたネットワーク内のノード 1 0 4 に対してファブリックインターコネクトを配線することによって実装できる。つまり、トーラス 1 0 2 内の所定のノード 1 0 4 の相対的位置は、ノード 1 0 4 を含む電子ノードの物理的位置ではなく、当該ノード 1 0 4 が接続されたノード 1 0 4 によって規定され得る。いくつかの実施形態において、トーラス 1 0 2 は、トーラスネットワークトポロジを実装するように、ファブリックインターコネクトによってともに配線された複数のソケットを備える。各ノード 1 0 4 は、ファブリックインターコネクトで使用するソケットに接続するように構成された現地交換可能装置 (F R U) (後述する) を備えており

、トーラス１０２内のノード１０４の位置は、ＦＲＵが挿入されたソケットによって決まる。

【００２１】

いくつかの実施形態において、ノード１０４間のリンク１０６は、例えば、接続された処理ノード間で差動対シグナリングを利用する１つ以上の高速ポイント・ツー・ポイント・シリアル通信リンクを含む。例えば、ノード１０４間の双方向接続は、１つ以上のペリフェラルコンポーネントインターコネクトエクスプレス（ＰＣＩｅ）リンク若しくは例えば×１　ＰＣＩｅリンク、×４　ＰＣＩｅリンク、×８　ＰＣＩｅリンク、×１６　ＰＣＩｅリンク等の外部ＰＣＩｅリンク、又は、１０ギガビットイーサネット（登録商標）（ＧｂＥ）アタッチメントユニットインターフェース（ＸＡＵＩ）インターフェースを含み得る。別の実施形態において、ノード１０４間のリンク１０６は、例えばイーサネット（登録商標）、ポイント・ツー・ポイント（ＰＰＰ）、高水準データリンク制御（ＨＤＬＣ）プロトコル、及び、アドバンスドデータ通信制御手順（ＡＤＣＣＰ）プロトコルインターフェースを含み得る。

【００２２】

図１Ｂは、一実施形態による並列トーラスインターコネクト１０１のブロック図１００Ｂである。並列トーラスインターコネクト１０１は、図１Ａで説明した複数のトーラス１０２Ａ～１０２Ｃを含んでいる。並列トーラスインターコネクト１０１では、異なるトーラス１０２からの複数のノード１０４がホスト１０８に接続されている。ホスト１０８は、データトラフィックを、トーラス１０２（例えば並列トーラスインターコネクト１０１内のトーラス１０２Ａ～１０２Ｃ）へ配送する計算装置である。ホストが単一のトーラス内のノードに接続する従来のトーラスインターコネクトとは違って、ホスト１０８は、複数のトーラス１０２内のノード１０４に接続する。図１の例において、ホスト１０８は、トーラス１０２Ａ内のノード（０、０、０）と、トーラス１０２Ｂ内のノード（２、２、０）と、トーラス１０２Ｃ内のノード（２、０、０）と、に接続している。ホスト１０８から異なるトーラス１０２へ至る複数の接続が、並列トーラスインターコネクト１０１を形成する。

【００２３】

一実施形態において、ホスト１０８は、ソケットを用いて、異なるトーラス１０２Ａ～１０２Ｃのノード１０４に接続する。一実施形態において、ソケットは、ソケットアドレス（例えば、インターネットプロトコルすなわちＩＰ）と、ポート番号と、を含み得る。並列トーラスインターコネクト１０１において、それぞれ異なるトーラス１０２に属するノード１０４は、同じソケットアドレスと、異なるポート番号と、を用いてホスト１０８に接続し得る。

【００２４】

一実施形態において、各トーラス１０２Ａ，１０２Ｂ，１０２Ｃは、独立し並列した他のトーラスのレプリカである。図１Ｂのトーラス１０２Ａ～１０２Ｃが並列接続される場合には、各トーラス１０２は独自のエコシステムを形成するネットワークを表す。つまり、１つのトーラス１０２内のデータトラフィックは、当該トーラス１０２内のノード１０４を通過するため、他のトーラス１０２内のデータトラフィックと混ざらない。一実施形態において、データトラフィックがホスト１０８を通過し、ホスト１０８が他のトーラスへデータトラフィックを再配送する場合には、データトラフィックは、トーラス１０２Ａ～１０２Ｃ間を移動する。

【００２５】

一実施形態において、ノード１０４間でやり取りされるデータトラフィックは、パケットに分割される。パケットは、並列トーラスインターコネクト１０１内の何れかのトーラス１０２において、始点ノードと終点ノードとの間の経路に沿って配信される。一実施形態において、始点ノードは、トーラス１０２へパケットを送信するホスト１０８に接続するノード１０４である。一実施形態において、終点ノードは、パケット内のデータを受信し、蓄積し、表示するノード１０４であるが、パケットをさらに伝搬することができない

10

20

30

40

50

。経路は、ゼロ個、１個または２個以上の中間ノードを含み得る。一実施形態において、各ノード１０４は、ファブリックインターコネクタへのインターフェースを含む。インターフェースは、ファブリックインターコネクタの対応するリンクに接続されたノードのポート間でパケットを配信するリンク層スイッチを実装している。

【００２６】

一実施形態において、異なるトーラス１０２内のノード１０４に接続されたホスト１０８は、特定のトーラス１０２を選択して、データトラフィックを伝搬する。一実施例において、ホスト１０８は、データトラフィックの種類に基づき、又は、種々のデータタイプに対する所定のＱoS要件に基づき、トーラス１０２を選択する。例えば、トーラス１０２Ａは、「ゴールド」ＱoSタイプを有するデータトラフィックを伝搬してもよいし、トーラス１０２Ｂは、「シルバー」ＱoSタイプを有するデータトラフィックを伝搬してもよいし、トーラス１０２Ｃは、「ブロンズ」ＱoSタイプを有するデータトラフィックを伝搬してもよい。「ゴールド」、「シルバー」及び「ブロンズ」ＱoSタイプは、データトラフィックが始点ノードから終点ノードへ到達するのにかかる保証時間の上限を指定する。別の実施例において、ホスト１０８は、トーラス１０２における輻輳に基づいてトーラス１０２を選択する。例えば、トーラス１０２Ａでデータトラフィックの輻輳が生じた場合には、ホスト１０８は、トーラス１０２Ｂまたは１０２Ｃを用いてデータトラフィックを送信してもよい。ホスト１０８は、並列トーラスインターコネクタ１０１にわたって特定のＱoSを有するデータトラフィックの配送を制御するのに対し、各トーラス１０２内のノード１０４は、トーラス１０２内で特定のＱoSを有するデータトラフィックの伝搬を制御する。

【００２７】

別の実施形態において、ホスト１０８はＱoSのタイプに基づいてトーラス１０２を選択する。例示的なＱoSは、データトラフィックにおいて表現される特定のコンフィデンシャルティグループ、カスタマアソシエーション等を含んでもよい。ＱoSのタイプは、各ＱoS内で予め設定されてもよい。一実施形態において、ＱoSのタイプは、データトラフィックの種類を区別するのに用いられるデータまたは音声プロトコルに含まれてもよい。

【００２８】

別の実施例において、ホスト１０８は、ホスト１０８で予め設定されたアルゴリズムに従い、一部または全ての並列トーラス１０２にわたってデータトラフィックを配送する。このアルゴリズムによって、例えば、ホスト１０８は、各トーラス１０２内でネットワーク輻輳を監視することができる。この実施形態では、ホスト１０８が、並列トーラスインターコネクタ１０１におけるトラフィックの輻輳やリンクの故障に応じて、データトラフィックを他の並列トーラス１０２へ送り直すので、ノードやリンクの故障の影響が低減する。

【００２９】

トーラス１０２が並列トーラスインターコネクタ１０１に接続される場合には、並列トーラスインターコネクタ１０１の帯域幅が、トーラス１０２の数に比例して増加する。例えば、ネットワークの帯域幅は、並列トーラスインターコネクタ１０１に追加されるトーラス１０２の各々の帯域幅によって、比例して増加する。

【００３０】

並列トーラスインターコネクタ１０１を管理するソフトウェアの拡張性も並列トーラスインターコネクタ１０１の利点である。例えば、並列トーラスインターコネクタ１０１用の管理ソフトウェアは、並列トーラスインターコネクタ１０１に追加されたトーラス１０２の各々と、追加されたトーラス１０２へのデータトラフィック配送と、を管理するように拡張され得る。一実施形態において、並列トーラスインターコネクタにトーラス１０２が追加される場合には、追加の前にトーラス１０２内のノード１０４に接続されていたホスト１０８間の帯域幅が増加する。

【００３１】

一実施形態において、管理ソフトウェアは、各トーラス102上でQoSを管理する。
一実施形態において、管理ソフトウェアは、ホスト108上で作動し、データトラフィックを並列トーラスインターコネクト101へ配送する。一実施例において、管理ソフトウェアは、上述したように、QoSに基づいてデータトラフィックを各トーラス102へ配送してもよい。別の実施例において、管理ソフトウェアは、セキュリティレベルに基づいてデータトラフィックを配送してもよい。例えば、1つのセキュリティレベルに関連付けられたデータトラフィックは、1つのトーラス102へ配送されてもよく、他のセキュリティレベルを有するデータトラフィックは、他のトーラス102へ配送されてもよい。このように、異なるセキュリティレベルを有するデータトラフィックは、単一のトーラスにわたって移送されることがない。また、特定のセキュリティレベルを有するデータトラフィックを伝搬するトーラス102は、追加的なセキュリティ対策を含んでもよい。当業者であれば、セキュリティレベルが、アプリケーションによって、又は、データを送受信するアプリケーションを使用するユーザによって設定され得ることを理解するであろう。

【0032】

さらなる実施形態において、ホスト108は、データトラフィックを、並列トーラスインターコネクト101内の特定のトーラス102へ送信することが制限されてもよい。例えば、ホスト108は、データを、並列トーラスインターコネクト101内のトーラス102のサブセットへ配送することが制限されてもよい。データ配送を制限するための方法は、トーラスのサブセット内のノード104にホスト108を接続することである（図示せず）。別の実施形態において、ホスト108は、並列トーラスインターコネクト101内のノード104に物理的に接続されてもよいが、接続されたノード104へのデータ送信を開始及び停止するタイミングは、管理ソフトウェアによって判断される。

【0033】

図2は、一実施形態による並列トーラスインターコネクトを通じてデータトラフィックを伝搬する方法のフローチャート200である。

【0034】

工程202では、ホストは、データトラフィックを受信する。例えば、ホスト108は、並列トーラスインターコネクト101へ配送されるデータトラフィックを受信する。

【0035】

工程204では、ホストは、データトラフィックを受信するノードを選択する。例えば、ホスト108は、並列トーラスインターコネクト101内のノード104のサブセットに接続されており、サブセットに含まれるノード104は、異なるトーラス102に関連付けられている。例えば、ホスト108は、トーラス102Aのノード(0、0、0)と、トーラス102Bのノード(2、2、0)と、トーラス102Cのノード(2、0、0)とに接続されてもよい。ホスト108は、データトラフィックを例えばパケットとして受信する場合に、データトラフィックを受信するノード104を、ノードのサブセットから選択する。上述したように、この選択は、例えば並列トーラスインターコネクト101のトーラス102における輻輳、データトラフィックで指定されるQoSのタイプ、又は、データトラフィックに関連付けられたセキュリティレベルに基づいて行われてもよい。例えば、ホスト108は、上記に基づいて、トーラス102Aのノード(0、0、0)、トーラス102Bのノード(2、2、0)、又は、トーラス102Cのノード(2、0、0)を選択し得る。

【0036】

工程206では、データトラフィックが、選択されたノードへ伝搬される。例えば、ホスト108は、パケットを、工程204で選択されたノード104へ伝搬する。

【0037】

図3は、一実施形態による並列トーラスインターコネクトにおけるノードの例示的な物理的配置を示すブロック図300である。図示された実施例において、ファブリックインターコネクトは、1つ以上のインターコネクト302を含み、インターコネクト302は、1つ以上のプラグインソケット304の段または集合体を有する。インターコネクト3

02は、例えば、バックプレーン、プリント配線基板、マザーボード、ケーブル若しくは他の自在配線又はこれらの組み合わせ等のような固定又は自在インターコネクトを含み得る。さらに、インターコネクト302は、電気信号伝達、光信号伝達又はこれらの組み合わせを実装し得る。各プラグインソケット304は、1つ以上のFRU（例えばFRU306～311）をインターコネクト302に接続するように動作するカードエッジソケットを備える。各FRUは、個々のトーラス102に関連付けられたノード104を表す。

【0038】

各FRUは、PCB上に配置されたコンポーネントを含む。これにより、コンポーネントは、PCBの金属層を介して相互に接続され、FRUによって表されるノードの機能を提供する。例えば、FRU306はPCB312を含み、PCB312には、1つ以上のプロセッサコア322を備えるプロセッサ320と、例えばDRAMデュアルインラインメモリーモジュール（DIMM）等の1つ以上のメモリーモジュール324と、ファブリックインターフェースデバイス326と、が実装されている。各FRUはソケットインターフェース330を含み、ソケットインターフェース330は、プラグインソケット304を介してFRUをインターコネクト302に接続するように動作する。

【0039】

インターコネクト302は、プラグインソケット304間にデータ通信経路を提供する。これにより、インターコネクト302は、FRUをリング状に接続し、リングを2D又は3Dトーラスネットワークトポロジータ（例えば、図1Bのトーラスネットワーク100B状）に接続するように動作する。FRUは、対応するファブリックインターフェース（例えば、FRU306のファブリックインターフェースデバイス326）を通るデータ通信経路を利用する。ソケットインターフェース330には、プラグインソケット304の対応する電気接点に電氣的に接続する電気接点（例えばカードエッジピン）が設けられており、X次元リング（例えば、ピン0, 1のリングX__INポート332、及び、ピン2, 3のリングX__OUTポート334）、Y次元リング（例えば、ピン4, 5のリングY__INポート336、及び、ピン6, 7のリングY__OUTポート338）、Z次元リング（例えば、ピン8, 9のリングZ__INポート340及びピン10, 11のリングZ__OUTポート342）のポートインターフェースとして機能する。図示された実施例において、各ポートは、例えばPCIEレーンの入力ポートまたは出力ポートを備える差動送信器である。当業者であれば、さらなるレーンやさらなるポートに対応するために、ポートがさらなるTX/RX信号ピンを含み得ることを理解するであろう。

【0040】

図4は、実施形態を実装し得るコンピュータシステムのブロック図400である。

【0041】

種々の実施形態は、ソフトウェア、ファームウェア、ハードウェア又はこれらの組み合わせで実装され得る。図4は、開示される実施形態又はその一部を、コンピュータ可読コードとして実装可能な例示的なコンピュータシステム400を示している。例えば、ここで説明したフローチャートの方法をシステム400で実装できる。この例示的なコンピュータシステム400に関して種々の実施形態を説明する。この説明を読んだ当業者であれば、他のコンピュータシステム及び/又はコンピュータアーキテクチャを用いて実施形態を実装する方法が明白となるであろう。

【0042】

コンピュータシステム400は、1つ以上のプロセッサ（例えばプロセッサ410）を含む。プロセッサ410は、専用のプロセッサであってもよいし、汎用のプロセッサであってもよい。例示的なプロセッサは、中央処理装置（CPU）を用いてデータを処理する。CPUは、コンピュータプログラムやアプリケーションの命令を遂行するプロセッサである。例えば、CPUは、算術演算と論理演算と入出力操作とを実行することによって、命令を遂行する。一実施形態において、CPUは、コンピュータプログラムやアプリケーションの意思決定コードを含む制御命令を遂行し、電子装置内の他のプロセッサ（例えばグラフィックス処理装置（GPU））に対して処理を任せる。GPUは、他の例示的なプ

ロセッサであって、電子装置上で計算が集中するアプリケーションを迅速に処理するための専用の電子回路である。GPUは、大きいデータブロック（例えば、コンピュータグラフィックスアプリケーション、画像、ビデオにおいて一般的な計算が集中するデータ）を効率良く並列処理するように、高度に並列化された構造を有している。GPUは、処理対象のデータをCPUから受信してもよいし、処理対象のデータを、処理済みのデータや操作から生成してもよい。一実施形態において、GPUは、ハードウェアを用いてデータを並列処理するハードウェアベースのプロセッサである。

【0043】

プロセッサ410は、通信基盤420（例えば、バスまたはネットワーク）に接続されている。

10

【0044】

また、コンピュータシステム400は、メインメモリー430（好ましくはランダムアクセスメモリー（RAM））を含み、さらには二次メモリー440を含み得る。二次メモリー440は、例えば、ハードディスクドライブ450、取外可能記憶ドライブ460及び/又はメモリースティックを含み得る。取外可能記憶ドライブ460は、フロッピー（登録商標）ディスクドライブ、磁気テープドライブ、光ディスクドライブ、フラッシュメモリー等を有してもよい。取外可能記憶ドライブ460は、周知の方法で、取外可能記憶装置470からの読み取り、及び/又は、取外可能記憶装置470への書き込みを行う。取外可能記憶装置470は、取外可能記憶ドライブ460によって読み取られ、取外可能記憶ドライブ460によって書き込まれるフロッピー（登録商標）ディスク、磁気テープ、光ディスク等を有してもよい。当業者によって理解されるように、取外可能記憶装置470は、コンピュータソフトウェア及び/又はデータを蓄積するコンピュータ可用記憶媒体を含む。

20

【0045】

別の実装において、二次メモリー440は、コンピュータシステム400がコンピュータプログラムや他の命令を読み込むための他の類似手段を含み得る。かかる手段は、例えば、取外可能記憶装置470と、インターフェース（図示せず）と、を含み得る。かかる手段の例は、プログラムカートリッジ及びカートリッジインターフェース（例えばビデオゲーム装置に見られるもの等）、取外可能メモリーチップ（EPROM又はPROM等）及び関連ソケット、並びに、取外可能記憶装置470からコンピュータシステム400へソフトウェアやデータを転送することの可能な他の取外可能記憶装置470及びインターフェースを含み得る。

30

【0046】

また、コンピュータシステム400は、通信・ネットワークインターフェース480を含み得る。通信・ネットワークインターフェース480は、コンピュータシステム400と外部装置との間でソフトウェアやデータを転送できるようにする。通信・ネットワークインターフェース480は、モデムと、通信ポートと、PCMCIAスロット及びカードと、などを含み得る。通信・ネットワークインターフェース480を介して転送されるソフトウェア及びデータは信号形式であって、信号は、通信・ネットワークインターフェース480によって受信できる電子信号、電磁信号、光信号、又は、その他の信号であり得る。これらの信号は、通信経路485を通して通信・ネットワークインターフェース480へ提供される。通信経路485は、信号を搬送し、ワイヤー若しくはケーブル、光ファイバ、電話回線、携帯電話リンク、RFリンク、又は、その他の通信チャンネルを用いて実装され得る。

40

【0047】

通信・ネットワークインターフェース480は、コンピュータシステム400が、LAN、WAN、インターネット等の通信ネットワーク又は媒体で通信できるようにする。通信・ネットワークインターフェース480は、有線又は無線接続で遠隔地のサイトやネットワークに繋がり得る。

【0048】

50

本書で使われている用語「コンピュータプログラム媒体」と、「コンピュータ可用媒体」と、「コンピュータ可読媒体」とは、取外可能記憶装置 470、取外可能記憶ドライブ 460、ハードディスクドライブ 450 に組み込まれたハードディスク等の媒体を指す総称である。通信経路 485 で搬送される信号もここで説明するロジックを具現化できる。コンピュータプログラム媒体と、コンピュータ可用媒体と、コンピュータ可読媒体とは、メインメモリ 430 や二次メモリ 440 等のメモリを指すこともあり、メモリはメモリ半導体（例えば、DRAM 等）であってよい。コンピュータプログラム製品は、コンピュータシステム 400 へソフトウェアを提供する手段である。

【0049】

コンピュータプログラム（コンピュータ制御ロジックとも呼ばれる）は、メインメモリ 430 及び／又は二次メモリ 440 に蓄積される。また、コンピュータプログラムは、通信・ネットワークインターフェース 480 を介して受信され得る。かかるコンピュータプログラムが実行されると、コンピュータシステム 400 は、ここで説明した実施形態を実施する。具体的に述べると、コンピュータプログラムが実行されると、プロセッサ 410 は、実施形態のプロセス（例えば、上述したフローチャートの方法のステップ）を実施する。つまり、かかるコンピュータプログラムは、コンピュータシステム 400 のコントローラに対応する。実施形態がソフトウェアを用いて実装される場合には、このソフトウェアは、コンピュータプログラム製品に蓄積されてもよく、例えば取外可能記憶ドライブ 460、インターフェース、ディスクドライブ 450、又は、通信・ネットワークインターフェース 480 を用いてコンピュータシステム 400 に読み込まれてもよい。

【0050】

また、コンピュータシステム 400 は、キーボード、モニター、ポインティング装置等の入力／出力／表示装置 490 を含み得る。

【0051】

実施形態は、例えば一般的なプログラミング言語（C、C++ 等）、Verilog HDL、VHDL、Altera HDL（AHDL）等のハードウェア記述言語（HDL）、又は、他のプログラミング及び／若しくはスキマティックキャプチャツール（サーキットキャプチャツール等）を用いて実現可能である。プログラムコードは、半導体、磁気ディスク、光ディスク（CD-ROM、DVD-ROM 等）を含む周知のコンピュータ可読媒体に格納することができる。コードは、インターネットやイントラネットを含む通信ネットワークで送信可能である。上述したシステム及び手法によって実現される機能、並びに／又は、上述したシステム及び手法によって提供される構造は、コア（CPU コア及び／又は GPU コア）において表現され得る。機能及び／又は構造は、プログラムコードで具現化され、集積回路の製造の一部としてハードウェアに変換され得ることが理解されよう。

【0052】

また、実施形態は、ソフトウェアを含むコンピュータプログラム製品を対象とし、ソフトウェアは、何らかのコンピュータ可用媒体に蓄積される。かかるソフトウェアが 1 つ以上のデータ処理装置で実行されると、データ処理装置が、ここで説明したように動作し、又は、上述したように、かかるソフトウェアが、ここで説明した実施形態を遂行する電子装置（例えば、ASIC 又はプロセッサ）の合成及び／又は製造を可能にする。実施形態は、何らかのコンピュータ可用媒体又はコンピュータ可読媒体と、現在または将来の何らかのコンピュータ可用媒体又はコンピュータ可読媒体とを使用する。コンピュータ可用媒体又はコンピュータ可読媒体の例は、一次記憶装置（例えば、ある種のランダムアクセスメモリ）、二次記憶装置（例えば、ハードドライブ、フロッピー（登録商標）ディスク、CD-ROM、ZIP ディスク、テープ、磁気式記憶装置、光学式記憶装置、MEMS、ナノテクノロジーによる記憶装置等）、並びに、通信媒体（例えば、有線及び無線通信ネットワーク、ローカルエリアネットワーク、ワイドエリアネットワーク、イントラネット等）を含むが、これらに限定されない。

【0053】

「概要」及び「要約」ではなく「詳細な説明」が請求項を解釈するためのものであることを理解されたい。「概要」及び「要約」は、発明者によって検討された１つ以上の例示的实施形態を記述し得るが、全ての例示的な実施形態を記述するものではなく、実施形態や添付の請求項を制限するものではない。

【 0 0 5 4 】

特定の機能及びこれらの関係の実装を例示する機能的な構成ブロックを参照しながら実施形態を説明した。ここでは、説明の都合上、これらの機能的構成ブロックの境界が恣意的に定められている。特定の機能及びこれらの関係が適切に遂行されるのであるならば、別の境界を定めることもできる。

【 0 0 5 5 】

前述した特定の実施形態の説明によって実施形態の本質が十分に明らかにされるため、他の人々は、当該技術の知識を応用し、過度の実験を行わずとも、本開示の一般概念から逸脱することなく、かかる特定の実施形態を様々な用途に向けて容易に修正及び／又は適用できる。かかる適用及び修正は、本書の教示と指導に基づいて開示された実施形態の均等物の趣旨及び範囲内にある。本書の用語や表現は、制限ではなく説明を目的としており、本明細書の用語や表現が教示と指導を踏まえて当業者によって解釈されるべきものであることを理解されたい。

【 0 0 5 6 】

実施形態の幅と範囲は上記の例示的な実施形態によって制限されず、専ら以降の請求項とその均等物に基づいて規定されるべきものである。

10

20

【 図 1 A 】

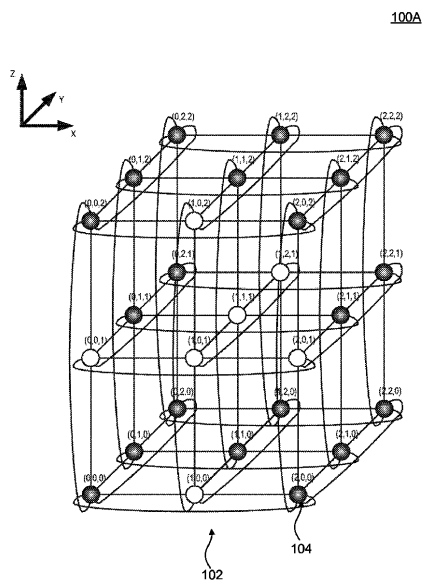


FIG. 1A

【 図 1 B 】

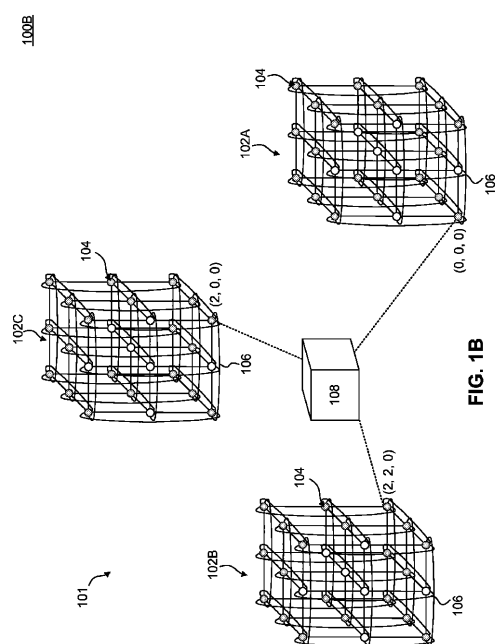
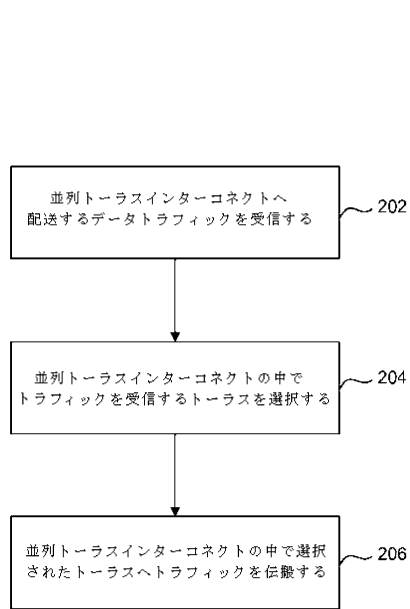
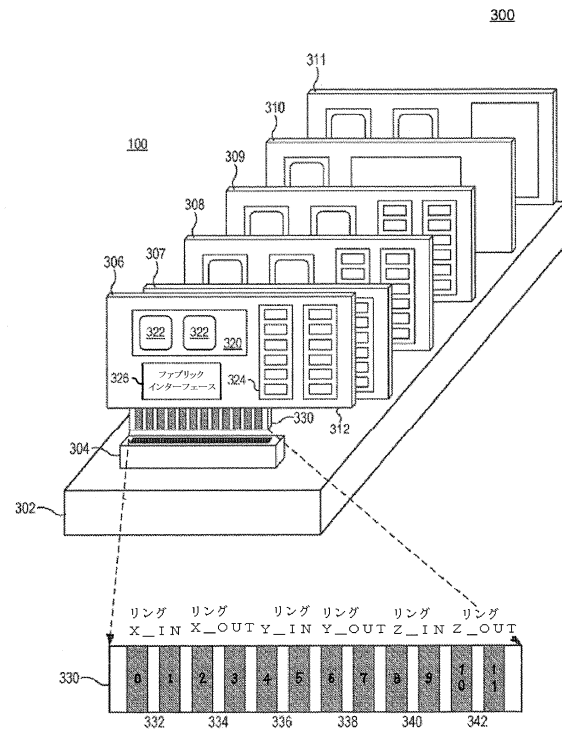


FIG. 1B

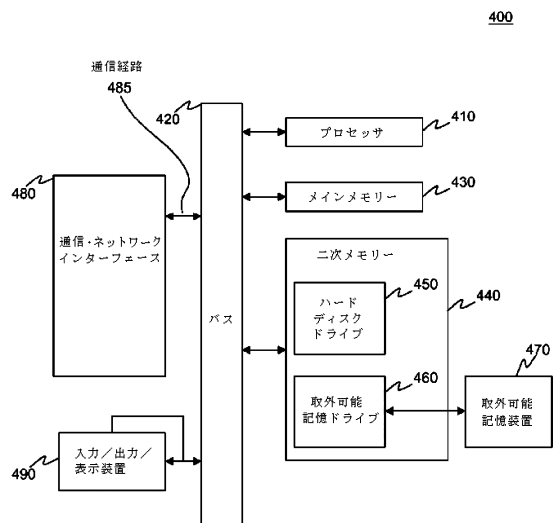
【図 2】



【図 3】



【図 4】



フロントページの続き

(74)代理人 100162156

弁理士 村雨 圭介

(72)発明者 ジャン - フィリップ フリッカー

アメリカ合衆国 9 4 0 4 0 カリフォルニア州、マウンテン ビュー、エイックラー コート
1 2 2 5

審査官 漆原 孝治

(56)参考文献 国際公開第2 0 0 8 / 1 1 4 4 4 0 (W O , A 1)

特表2 0 0 8 - 5 3 6 3 7 2 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 5 / 1 7 3

H 0 4 L 1 2 / 7 2 1