

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6436394号  
(P6436394)

(45) 発行日 平成30年12月12日 (2018.12.12)

(24) 登録日 平成30年11月22日 (2018.11.22)

(51) Int. Cl.	F I
<b>H04L 12/42 (2006.01)</b>	H04L 12/42 Z
<b>G06F 13/14 (2006.01)</b>	G06F 13/14 310H
<b>H04L 12/771 (2013.01)</b>	H04L 12/771
<b>H04L 12/931 (2013.01)</b>	H04L 12/931
<b>G06F 13/10 (2006.01)</b>	G06F 13/14 310F
請求項の数 18 (全 23 頁) 最終頁に続く	

(21) 出願番号	特願2015-118773 (P2015-118773)	(73) 特許権者	504407000
(22) 出願日	平成27年6月12日 (2015.6.12)		パロ アルト リサーチ センター イン
(65) 公開番号	特開2016-10157 (P2016-10157A)		コーポレイテッド
(43) 公開日	平成28年1月18日 (2016.1.18)		アメリカ合衆国 カリフォルニア州 94
審査請求日	平成30年6月12日 (2018.6.12)		304 パロ アルト カイオーテ ヒル
(31) 優先権主張番号	14/313, 922		ロード 3333
(32) 優先日	平成26年6月24日 (2014.6.24)	(74) 代理人	100092093
(33) 優先権主張国	米国 (US)		弁理士 辻居 幸一
(31) 優先権主張番号	14/512, 341	(74) 代理人	100082005
(32) 優先日	平成26年10月10日 (2014.10.10)		弁理士 熊倉 禎男
(33) 優先権主張国	米国 (US)	(74) 代理人	100067013
早期審査対象出願			弁理士 大塚 文昭
		(74) 代理人	100086771
			弁理士 西島 孝喜
		最終頁に続く	

(54) 【発明の名称】 統合ストレージ、処理、およびネットワークスイッチを組み込むネットワークスイッチング構造を有するコンピューティングシステムフレームワーク、および、それらを作成および使用する方

(57) 【特許請求の範囲】

【請求項 1】

統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムであって、

複数の処理ノードと、

複数のモジュール内ポートであって、各モジュール内ポートは前記処理ノードの1つと関連付けられる、モジュール内ポートと、

複数のモジュール間ポートであって、各モジュール間ポートは前記処理ノードの1つと関連付けられる、モジュール間ポートと、

複数の処理モジュールであって、各処理モジュールは、前記処理ノードのうちの少なくとも1つを備え、モジュール内ネットワークを形成するために、各処理ノードは当該処理ノードが備えられる前記処理モジュール内の前記処理ノードのすべての残りの処理ノードに接続される、処理モジュールと、

前記処理モジュールのうちの1つの処理モジュールの1つのモジュール間ポートと前記処理モジュールのうちの他の1つの処理モジュールの他の1つのモジュール間ポートとの間に、少なくとも1つの接続を備える、モジュール間ネットワークであって、前記モジュール間ネットワークは、前記処理モジュール  $(2^P * N + 1)$  から  $(2^P * N + 2^{P-1})$  のモジュール間ポート  $P$  を、前記処理モジュール  $(2^P * N + 2^{P-1} + 1)$  から  $(2^P * N + 2^P)$  のモジュール間ポート  $P$  と、それぞれ接続する方法に従って提供される接続を備え、 $N$  は  $[0, 1, \dots, (M / 2^P) - 1]$  であり、 $P$  は整数であり、前記処理モジュール

10

20

の数はMであり、前記処理モジュールの各々は少なくともP個のモジュール間ポートを有する、前記モジュール間ネットワークと、

を備える、コンピューティングシステム。

【請求項2】

前記モジュール内ネットワークは、

バス、リング、スター、メッシュ、およびクロスバースイッチ、

のうちの少なくとも1つを備える、請求項1に記載のシステム。

【請求項3】

前記処理ノードは各々、

物理ノード及び仮想ノードの少なくとも1つを備える、

請求項1に記載のシステム。

10

【請求項4】

前記処理ノードは各々、

処理要素、メモリコントローラ、メモリ、ストレージコントローラ、ストレージデバイス、及び、モジュール内ポート及びモジュール間ポートへのインタフェースの少なくとも1つを備える、

請求項1に記載のシステム。

【請求項5】

線形的に接続された処理モジュールの列であって、前記列の最初と最後が相互に結合されてリングを形成し、前記接続された処理モジュールの各々から2つのモジュール間ポートを使用する、前記列を備える、更なるモジュール間ネットワーク、

をさらに備える、請求項1に記載のシステム。

20

【請求項6】

前記リングの2つの非隣接処理モジュールの間の接続であって、前記2つの処理モジュールのうちの1つに配置され、前記線形接続に使用されない、1つのモジュール間ポートを介する、前記2つの非隣接処理モジュールのうちの他方に配置され、前記線形接続に使用されない、別のモジュール間ポートへの接続をさらに備える、請求項5に記載のシステム。

【請求項7】

前記リングの処理モジュールNと前記リングの処理モジュールN+Sとの間の、前記処理モジュールNの1つの未使用モジュール間ポートを、前記処理モジュールN+Sの別の未使用モジュール間ポートと接続することによる接続であって、Mは前記リングの前記処理モジュールの数であり、SはMの整数の約数であり、Nは $[x, x+S, \dots, x+(M-S)(S\text{ずつ増加})]$ であり、xは $[1, 2, \dots, S-1]$ である、接続と、

30

前記リングの処理モジュールNと前記リングの処理モジュールN+S-1との間の、前記処理モジュールNの1つの未使用モジュール間ポートを、前記処理モジュールN+S-1の別の未使用モジュール間ポートと接続することによる接続であって、Mは前記リングの前記処理モジュールの数であり、SはMの整数の約数であり、Nは $[x, x+S, \dots, x+(M-S)(S\text{ずつ増加})]$ であり、xは $[1, 2, \dots, S-1]$ である、接続と、

前記リングの処理モジュールNと前記リングの処理モジュールN+S-rとの間の、前記処理モジュールNの1つの未使用モジュール間ポートを、前記処理モジュールN+S-rの別の未使用モジュール間ポートと接続することによる接続であって、Mは前記リングの前記処理モジュールの数であり、SはMの整数の約数であり、Nは $[x, x+S, \dots, x+(M-S)(S\text{ずつ増加})]$ であり、xは $[1, 2, \dots, S-1]$ であり、rは $[2, 3, \dots, S-1]$ である、接続と、

40

のうちの少なくとも1つをさらに備える、請求項5に記載のシステム。

【請求項8】

データバケットを、前記処理ノードのうちの1つから前記処理ノードのうちの他の1つへ確立されたモジュール間接続を介して転送する、送信モジュールと、を更に備える請求項1に記載のシステム。

50

## 【請求項 9】

データパケットトラフィックパターンを測定する、監視モジュールと、  
前記確立されたモジュール間接続を前記データパケットトラフィックパターンに基づいて変更する、修正モジュールと、  
をさらに備える、請求項 8 に記載のシステム。

## 【請求項 10】

統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムを生成する又は使用方法であって、

複数の処理モジュールを提供するステップであって、各処理モジュールは複数の処理ノードを含み、各処理ノードは当該処理ノードが備えられる前記処理モジュール内の前記処理ノードのすべての残りの処理ノードに接続される、複数の処理モジュールを提供するステップと、

前記処理ノードの少なくとも 1 つのためにモジュール間ポートを提供するステップと、  
前記モジュール間ポートの各々のために、モジュール間接続を作成するために、前記処理ノードが備えられる処理モジュールとは異なる前記処理モジュールにおける処理ノードの他の 1 つを決定するステップと、

前記処理モジュール間のモジュール間接続を確立することにより、前記決定に基づいて前記モジュール間接続を確立するステップであって、前記処理モジュールの数は  $M$  であり、各処理モジュールは少なくともモジュール間ポート  $P$  を有し、処理モジュール  $(2^P * N + 1)$  から  $(2^P * N + 2^{P-1})$  のモジュール間ポート  $P$  を処理モジュール  $(2^P * N + 2^{P-1} + 1)$  から  $(2^P * N + 2^P)$  の更なるモジュール間ポート  $P$  に各々接続し、 $N$  は  $[0, 1, \dots, (M / 2^P) - 1]$  であり、 $P$  は整数である、確立するステップと、

を備える方法。

## 【請求項 11】

バス、リング、スター、メッシュ、およびクロスバースイッチのうちの少なくとも 1 つを介するモジュール内接続を形成するステップを更に備える、

請求項 10 に記載の方法。

## 【請求項 12】

前記処理ノードの各々を物理ノード及び仮想ノードの 1 つとして提供するステップをさらに備える、

請求項 10 に記載の方法。

## 【請求項 13】

前記処理ノードは各々は、

処理要素、メモリコントローラ、メモリ、ストレージコントローラ、ストレージデバイス、及び、モジュール内ポート及びモジュール間ポートへのインタフェースの少なくとも 1 つを備える、

請求項 10 に記載の方法。

## 【請求項 14】

残りの処理モジュールの列を線形的に接続することと、

リングを形成するために前記列の最初と最後を接続することと、を含む、さらなるモジュール間接続を確立するステップをさらに備え、

前記接続は、前記接続された処理モジュールの各々から 2 つのモジュール間ポートを用いて提供される、

請求項 10 に記載の方法。

## 【請求項 15】

前記リングにおける前記処理モジュールの 1 つにおける 1 つの使用されないモジュール間ポートを前記リングにおける他の 1 つの隣接しない処理モジュールにおける他の 1 つの使用されないモジュール間ポートに接続するステップをさらに備える、請求項 14 に記載の方法。

## 【請求項 16】

前記リングの処理モジュールNと前記リングの処理モジュールN + S との間の、前記処理モジュールNの1つの未使用モジュール間ポートを、前記処理モジュールN + S の別の未使用モジュール間ポートと接続することによる接続であって、Mは前記リングの前記処理モジュールの数であり、SはMの整数の約数であり、Nは $[x, x + S, \dots, x + (M - S)(S \text{ ずつ増加})]$ であり、xは $[1, 2, \dots, S - 1]$ である、接続を提供するステップと、

前記リングの処理モジュールNと前記リングの処理モジュールN + S - 1 との間の、前記処理モジュールNの1つの未使用モジュール間ポートを、前記処理モジュールN + S - 1の別の未使用モジュール間ポートと接続することによる接続であって、Mは前記リングの前記処理モジュールの数であり、SはMの整数の約数であり、Nは $[x, x + S, \dots, x + (M - S)(S \text{ ずつ増加})]$ であり、xは $[1, 2, \dots, S - 1]$ である、接続を提供するステップと、

10

前記リングの処理モジュールNと前記リングの処理モジュールN + S - r との間の、前記処理モジュールNの1つの未使用モジュール間ポートを、前記処理モジュールN + S - rの別の未使用モジュール間ポートと接続することによる接続であって、Mは前記リングの前記処理モジュールの数であり、SはMの整数の約数であり、Nは $[x, x + S, \dots, x + (M - S)(S \text{ ずつ増加})]$ であり、xは $[1, 2, \dots, S - 1]$ であり、rは $[2, 3, \dots, S - 1]$ である、接続を提供するステップと、

のうちの少なくとも1つをさらに備える、請求項14に記載の方法。

【請求項17】

20

データパケットを、前記処理ノードのうちの1つから前記処理ノードのうちの他の1つへ前記確立されたモジュール間接続を介して転送するステップ、をさらに備える請求項10に記載の方法。

【請求項18】

データパケットトラフィックパターンを測定するステップと、

前記確立されたモジュール間接続を前記データパケットトラフィックパターンに基づいて変更するステップと、

をさらに備える、請求項17に記載の方法。

【発明の詳細な説明】

【技術分野】

30

【0001】

本出願は、一般的には、データセンタで使用されるコンピューティングシステムフレームワークに関し、より詳細には、統合ストレージ、処理、およびネットワークスイッチング構造を作成および使用するシステムおよび方法に関する。

【背景技術】

【0002】

データセンタは急速に発展しており、その発展速度は加速すると思われる。急速な発展は、高まる要求により促され、データセンタのコンポーネントのコスト削減により可能となる。データセンタは、主に、処理ノード、ストレージノード、および、処理ノードおよびストレージノードを接続するネットワークから構成される。処理ノードおよびストレージノードは両方とも、以前より小さく安価になっており、エネルギー効率も良くなっているため、データセンタは、より狭い空間に多くの処理およびストレージノードを詰め込んで、データ処理およびストレージの高まる要求に応えることができるようになってきている。処理ノードは、さらに多くのデータを以前より高頻度で消費し、ストレージノードとの間でデータの検索および保存を行うので、ネットワークは、さらに多くのデータを以前より高速で、増加する接続間で送信しなければならない。結果として、処理およびストレージノードのコスト低下との関連において、ネットワークのコストが重要になる。1つの予測では、大抵はネットワークスイッチおよびケーブルから成るネットワークのコストは、新しいデータセンタの約50%としている。

40

【発明の概要】

50

## 【発明が解決しようとする課題】

## 【0003】

従来のデータセンタにおいて、処理ノードは、典型的には、単一の主ネットワークを介して接続される。二次ネットワークがある場合、主に管理目的で使用されるが、本明細書においては述べない。各処理ノードは、ハードディスクまたは固体ディスクなど、局所的に取り付けられた1つ以上の長期ストレージデバイスを有してもよい。処理ノードは、その長期ストレージデバイスにアクセスして内部的な要求を満たし、システム全体の分散ストレージシステムを代表する場合が多い。各々が1つ以上の長期ストレージデバイスを有する多くの処理ノードが、処理モジュールに包括される。データセンタの計算能力は、主として、処理モジュールを追加することにより向上する。この構築フレームワークでは、追加された全処理ノードが主高速ネットワークに依存して既存の処理ノードと通信するため、主高速ネットワークへの要求が高くなる。高速ネットワークの能力は、追加される処理ノードの数に比例して向上しなければならない。処理ノードが高速かつ安価になっている一方で、高速接続は高価になっているため、高速ネットワークのコストは、データセンタの計算能力の低価格での向上を妨げる障害となっている。

10

## 【0004】

例えば、データセンタで使用される従来のネットワークアーキテクチャは、3層のネットワークスイッチ（すなわち、低層から高層へ、アクセス、アグリゲート、およびコア層）から成る多重根のツリートポロジに従う3階層システムである。アクセス層は、直接的にルートサーバに到達し、アグリゲート層内に相互接続する。アクセス層スイッチは、データセンタをインターネットへ接続する役割も果たす、コア層スイッチにより、最終的に相互に接続する。3階層レガシーシステムは増大化が困難であり、3階層の上の層ほど、非常に超過傾向にある。加えて、耐障害性、エネルギー効率、および断面帯域幅が問題になる。

20

## 【0005】

さらなる例において、Fat Treeデータセンタアーキテクチャが、従来の3階層データセンタネットワークアーキテクチャが直面する、超過および断面帯域幅の問題に対処しようと試みている。Fat Treeトポロジは、1:1の超過割合および完全な二分帯域幅を提案する。しかしながら、Fat Treeトポロジは、3階層レガシーシステムよりも非常に多い数のネットワークスイッチを利用し、同様に増大化が困難である。

30

## 【0006】

さらなる例において、DCellアーキテクチャは、1つのサーバが直接的に多くの他のサーバと接続される、サーバ中心のハイブリッドアーキテクチャを採用する。DCellトポロジは、複数のレベルに配置したセルの再帰的に構築された階層構造に依存し、ここで高いレベルのセルは低い層の複数のセルを包含し、セル内部でサーバは、サーバ自体のスイッチに割り当てられる。容易に拡張可能である一方、断面帯域幅およびネットワーク待ち時間は、DCellアーキテクチャにおいて重要な課題である。加えて、DCellは、拡張性を達成するため、各サーバに複数のネットワークインタフェースを要する。

## 【0007】

Facebookは、Open Compute Projectを企画しており、エネルギー効率およびコスト効率の両方を満たすデータセンタサーバの開発を志している。Open Compute Projectが進める取り組みは、虚飾のないハードウェア設計、ブロックを構築するオープンボルトのストレージ、機械的な搭載システム、および高いディスク密度を含む。これらの取り組みの結果が、構築および稼働するのに、従来のサーバハードウェアと比較して、38%のエネルギー効率向上および24%の安価を達成した、虚飾のないサーバからなるデータセンタである。しかしながら、Open Compute Projectにおいて実践される解決策は、処理ノードの包括の最適化ということになる。処理機能とストレージ機能との間の基本的な二分法は、計算専用の処理ノードと保管専用のストレージデバイスとの間にもたらされる、ネットワークトラフィックと連動して、変化しないままである。

40

50

## 【 0 0 0 8 】

Nutanixは、高速ストレージ (Server Flash) および低速ストレージ (Hard Disk Storage) を局所的に処理ノードに組み込む、Nutanix Virtual Computing Platformを開発して、データセンタの計算の速度および効率を向上させてきた。しかしながら、基本的なネットワークの改良は明らかにされていない。

## 【 0 0 0 9 】

したがって、データセンタにおける高い帯域幅および低い遅延接続の要求を満たし、一方で高い計算能力への高まり続ける必要性に適応するために、より多くの処理ノードを追加する必要性が残されたままである。好適には、新しい処理ノードが既存のコンピュータシステムに追加される際、新しい処理ノードはネットワーク機能性を包含しており、ネットワークスイッチなどの専用ネットワーク機器をインストールする必要があるが、ほとんどないか、または全くなく、したがって、計算能力およびネットワーク性能は処理ノードの追加に伴って向上する。このデータセンタネットワークアーキテクチャのパラダイムは、それ自体のネットワーク機能性を包含するデータ処理ノード構造を介して、データ処理、保管、および送信を統合することを想定しており、したがって、ネットワークスイッチの必要性を除去または最小化し、一方で計算能力、保管性能、低コスト、拡張容易性、および頑健性、および低エネルギー消費への要求を満たす。

## 【 0 0 1 0 】

さらに、このパラダイムは、ネットワークスイッチを完全に取り除く代わりに、統合ストレージ、処理、およびネットワークスイッチング構造に組み込んでもよい。ネットワークスイッチは、ネットワークにおいてハブと同じ場所を占有する。ハブと異なり、ネットワークスイッチは、単にパケットを全ポートまで繰り返すのではなく、適切に各データパケットを検査および処理する。ネットワークスイッチは、各ネットワークセグメントに存在するノードのネットワークアドレスをマップし、必要なトラフィックだけ、スイッチを通過させる。スイッチによりパケットが受信されると、スイッチはハードウェアアドレスの送信先および送信元を調査して、ネットワークセグメントおよびアドレスの表と比較する。セグメントが同じ場合、パケットは廃棄、すなわち「フィルタリング」される。セグメントが異なる場合、パケットは正しいセグメントへ「転送」される。加えて、スイッチは、不良または正しく整理していないパケットの拡散を、転送しないことにより防ぐ。したがって、ネットワークスイッチを統合ストレージ、処理、およびネットワークスイッチング構造に組み込むことで、データトラフィックフローの制御が容易になり、データセンタコンピュータシステムの全体的な性能が向上する。さらに重要なことに、既存のデータセンタは広範囲にスイッチを使用するため、それ自体のネットワーク機能を包含する処理ノードで既存のデータセンタを増大させること、または、既存のデータセンタを処理ノードへ移行することは、ネットワークスイッチが処理ノードを介して影響を受けるネットワーク機能に大きく依存する接続スキームへ効率的に組み込まれ得る場合、容易になり、コスト効率が良くなるであろう。

## 【 0 0 1 1 】

したがって、ネットワークスイッチを、少なくとも処理ノードの一部が、それ自体のネットワーク機能を包含するデータセンタネットワークシステムへ効率的に組み込む、システムおよび方法の必要性が残っている。

## 【 課題を解決するための手段 】

## 【 0 0 1 2 】

統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムフレームワークを作成および使用するシステムおよび方法が、提供される。処理ノードは、物理的または仮想的にかかわらず、モジュール内ポート、モジュール間ポート、およびローカルストレージデバイスと関連付けられる。複数の処理ノードがモジュール内ポートを介して連結され、処理モジュールを形成する。複数の処理モジュールがモジュール間ポートを介して、さらに接続され、コンピューティングシステムを形成する。

いくつかのモジュール間接続スキームが説明されており、各々が既存のネットワークパケットルーティングアルゴリズムで使用するよう適合され得る。各処理ノードは、直接的に接続される隣接ノードの状態を追跡すればよく、システム内の残りの処理ノードと高速接続される必要性は排除される。結果として、専用ネットワークスイッチング機器は必要なく、ネットワーク性能は、処理ノードが追加されるにつれて自然に向上する。さらに、ネットワークスイッチはネットワーク接続に組み込まれてもよく、それによりネットワークトラフィック制御が容易になる。

【 0 0 1 3 】

1つの実施形態は、統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムフレームワークを提供する。システムは4個以上の処理ノードを含む。システムは4個以上のモジュール内ポートをさらに含み、各モジュール内ポートは処理ノードの1つと一意的に関連付けられる。システムは複数のモジュール間ポートをさらに含み、各モジュール間ポートは処理ノードの1つと関連付けられる。システムは複数の処理モジュールをさらに含み、各処理モジュールは複数の処理ノードから一意的に選択される2つ以上の処理ノードを備え、ここで処理ノードの各々は、処理モジュールのうちの1つにのみ備えられる。システムは複数のモジュール内ネットワークをさらに含み、ここで各処理モジュール内の処理ノードは完全に相互接続されている。最後に、システムは、1つの処理モジュールの1つのモジュール間ポートと別の処理モジュールの別のモジュール間ポートとの間に、少なくとも1つの接続を備える、モジュール間ネットワークを含む。接続は、ケーブルまたはネットワークスイッチを介して、または両方を介して、なされてよい。処理モジュールは、ネットワークスイッチを模倣してもよい。

【 0 0 1 4 】

さらなる実施形態は、統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムフレームワークを作成および使用方法を、提供する。方法は以下を含む：4個以上の処理ノードを提供すること；4個以上のモジュール内ポートであって、各モジュール内ポートは処理ノードの1つと一意的に関連付けられる、モジュール内ポートを提供すること；複数のモジュール間ポートであって、モジュール間ポートの各々は処理ノードの1つと関連付けられる、モジュール間ポートを提供すること；複数の処理モジュールであって、各処理モジュールは複数の処理ノードから一意的に選択される2つ以上の処理ノードを備え、処理ノードの各々は処理モジュールのうちの1つのみに備えられる、処理モジュールを提供すること；複数のモジュール内ネットワークであって、各処理モジュール内の処理ノードは完全に相互接続される、モジュール内ネットワークを提供すること；および、1つの処理モジュールの1つのモジュール間ポートと別の処理モジュールの別のモジュール間ポートとの間に、少なくとも1つの接続を備える、モジュール間ネットワークを提供すること。さらに、モジュール内ネットワーク接続およびモジュール間接続は、ケーブルまたはネットワークスイッチを介して、または両方を介して、なされ得る。処理モジュールは、ネットワークスイッチを模倣してもよい。

【 0 0 1 5 】

さらなる実施形態は、統合ストレージ、処理、およびネットワークスイッチを組み込むネットワークスイッチング構造を有するコンピューティングシステムフレームワークを、提供する。システムは4個以上の処理ノードを含む。システムは4個以上のモジュール内ポートをさらに含み、各モジュール内ポートは処理ノードの1つと一意的に関連付けられる。システムは複数のモジュール間ポートをさらに含み、各モジュール間ポートは処理ノードの1つと関連付けられる。システムは複数の処理モジュールをさらに含み、各処理モジュールは複数の処理ノードから一意的に選択される2つ以上の処理ノードを備え、処理ノードの各々は処理モジュールのうちの1つのみに備えられる。システムは複数のモジュール内ネットワークをさらに含み、各処理モジュール内の処理ノードは完全に相互接続されている。システムは、1つの処理モジュールの1つのモジュール間ポートと別の処理モジュールの別のモジュール間ポートとの間に、少なくとも1つの接続を備える、モジュール間ネットワークをさらに含む。最後に、システムは、モジュール内ポートおよびモジュ

ール間ポートのうちの少なくとも1つと動作可能に接続される、少なくとも1つのネットワークスイッチを含む。接続は、ケーブルまたはネットワークスイッチを介して、または両方を介して、なされてよい。処理モジュールは、ネットワークスイッチを模倣してもよい。

【0016】

本発明の他の実施形態が、以降の詳細な説明から、当業者に容易に明らかになるであろう。ここで、本発明の実施を考慮した最良の手法を図示することにより、本発明の実施形態が説明される。理解されるように、本発明は他の異なる実施形態でも可能であり、その個別の詳細は、全てが本発明の精神および範囲から逸脱することなく、様々な明確な点において修正可能である。したがって、図および詳細な説明は、本質的に例示と見なされ、

10

【図面の簡単な説明】

【0017】

【図1】図1は、1つの実施形態による、統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムフレームワークを示すブロック図である。

【図2】図2は、1つの実施形態による、統合ストレージ、処理、およびネットワークスイッチを組み込むネットワークスイッチング構造を有する、コンピューティングシステムフレームワークを示すブロック図である。

【図3】図3は、例として、処理モジュールの各々に4個のモジュール間ポート( $P = 4$ )を有する16個の処理モジュール( $M = 16$ )を含む、コンピューティングシステムフレームワークとして図示される、二分スパニング接続スキームを示す図である。

20

【図4】図4は、例として、8個の処理モジュール( $M = 8$ )を含むコンピューティングシステムフレームワークとして図示される、リングシステム接続スキームを示す図である。

【図5】図5は、例として、ショートカットを有するリングシステム接続スキームを示す図である。

【図6】図6は、例として、 $S$ 個ごとのホップを有するリングシステム接続スキームを示す図である。

【図7】図7は、例として、 $S$ 個ごとの分割ホップを有するリングシステム接続スキームを示す図である。

30

【図8】図8は、例として、 $S$ 個ごとの調整可能分割ホップを有するリングシステム接続スキームを示す図である。

【図9】図9は、例として、ネットワークスイッチにより実現されるショートカットを有するリングシステム接続スキームを示す図である。

【図10】図10は、例として、ネットワークスイッチにより実現される $S$ 個ごとのホップを有するリングシステム接続スキームを示す図である。

【図11】図11は、例として、ネットワークスイッチにより実現される $S$ 個ごとの分割ホップを有するリングシステム接続スキームを示す図である。

【図12】図12は、例として、ネットワークスイッチにより実現される $S$ 個ごとの調整可能分割ホップを有するリングシステム接続スキームを示す図である。

40

【図13】図13は、例として、ブロックまたは外部デバイスを、ショートカットを有するリングシステム接続スキームと接続するケーブルリンクを示す図である。

【図14】図14は、例として、ブロックまたは外部デバイスを、 $S$ 個ごとのホップを有するリングシステム接続スキームと接続するケーブルリンクを示す図である。

【図15】図15は、例として、ブロックまたは外部デバイスを、 $S$ 個ごとの分割ホップを有するリングシステム接続スキームと接続するケーブルリンクを示す図である。

【図16】図16は、例として、ブロックまたは外部デバイスを、 $S$ 個ごとの調整可能分割ホップを有するリングシステム接続スキームと接続するケーブルリンクを示す図である。

50



【図 17】図 17 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現されるショートカットを有するリングシステム接続スキームへの接続を示す図である。

【図 18】図 18 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現される S 個ごとのホップを有するリングシステム接続スキームへの接続を示す図である。

【図 19】図 19 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現される S 個ごとの分割ホップを有するリングシステム接続スキームへの接続を示す図である。

【図 20】図 20 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現される S 個ごとの調整可能分割ホップを有するリングシステム接続スキームへの接続を示す図である。

10

【発明を実施するための形態】

【0018】

1 つの実施形態において、コンピューティングシステムフレームワークは、2 レベル構成に体系化された複数の処理ノードを組み込む。第 1 のレベルで、複数の処理ノードは処理モジュールを形成し；第 2 のレベルで、複数の処理モジュールはコンピューティングシステムを形成する。図 1 は、1 つの実施形態による、統合ストレージ、処理、およびネットワークスイッチング構造を有するコンピューティングシステムフレームワーク (10) を示すブロック図である。

20

【0019】

処理ノード (1) は、物理ノードまたは仮想ノードであってよい。1 つの実施形態において、物理ノードは、いくつかの仮想ノードを実現してもよい。以降の説明において、処理ノードは、物理ノードおよび仮想ノードの両方を指す。

【0020】

処理ノード (1) は、処理要素、メモリコントローラ、メモリ、ストレージコントローラ、1 つ以上のストレージデバイス、および、モジュール内およびモジュール間ポートへのインタフェース、であり得る。物理的な処理ノードにおいて、それらは物理デバイスである。仮想的な処理要素において、これらのデバイスの機能はエミュレートされる。全ノードが一意的なネットワークアドレスを有する。処理ノード (1) は、典型的には、処理ノードの内部の必要性を満たすか、または、システム全体の分散ストレージシステムの代理の役割を果たしてもよい、ノード自体のストレージデバイスと接続される。

30

【0021】

複数の処理ノードは、処理モジュール (2) を構成する。各処理ノード (1) は、典型的には、モジュール内ポート (3) と称される高速データ転送ポートを備える。処理モジュールにおける各処理ノードは、典型的には、同じ処理モジュール内の各々または全ての他の処理ノードと、モジュール内ポート (3) を介して接続される。言い換えれば、各処理モジュール内の処理ノードは、完全に相互接続される。1 つの処理モジュール内のモジュール内ポートを介するこれらの接続は、モジュール内接続 (4) と称される。モジュール内ネットワークは、処理モジュール内の各処理ノードから、同じ処理モジュール内の残りの処理モードへと形成される接続 (単数または複数) からなる。モジュール内接続 (4) により、同じ処理モジュール内の 2 つの処理ノード間でのデータパケット交換が、単一ステップで完結される。1 つの実施形態において、モジュール内ポートは、PCI Express Serial Rapid IO を使用して、または、低いエラー率および短い相互接続距離を利用する他の技術を使用して実現され、非常に小型化された、高い帯域幅の信頼性がある安価な接続を提供する。さらなる実施形態において、モジュール内接続は、リング、バス、スター、メッシュ、および 1 つ以上のクロスバスイッチの集合を含む、少なくとも 1 つの任意の技術およびトポロジを使用して、達成される。

40

【0022】

リングトポロジにおいて、全ての処理ノードは、通信する目的で正確に 2 つの隣接ノード

50

ドを有する。全メッセージが、リングを介して同じ方向に進んでもよく、進まなくてもよい。バストポロジにおいて、バスネットワーク（コンピュータのシステムバスと混同されるべきでない）は、共通の基幹を使用して、全処理ノードまたはデバイスを接続する。基幹は、処理ノードまたはデバイスがインタフェースコネクタに付着または接触する、共有の通信媒体として機能する。スターネットワークは、残りの処理ノードと接続する「ハブ」と呼ばれる中心接続点が特徴的である。メッシュトポロジにおいて、処理ノードは、多くの冗長な相互接続と接続される。ハイブリッドトポロジは、任意の2つ以上のネットワークトポロジの組み合わせを備える。

#### 【0023】

モジュール内ポート（3）に加えて、処理ノードは、別の処理モジュールに配置される処理ノードと別のモジュール間ポートを介して接続する、モジュール間ポート（5）と称されるネットワークポートを、さらに備えてもよい。2つの処理モジュール間の、2つの別々の処理モジュールからの2つの処理ノードに配置される2つのモジュール間ポートを介するこれらの接続は、それぞれ、モジュール間接続（6）と称される。データパケットは、2つの処理モジュールがモジュール間接続を介して直接的に接続される場合、1つの処理モジュールから別の処理モジュールへ、1つのステップで送信されてよい。モジュール間接続は、典型的には、2つのモジュール間ポートを連結するケーブルを介して、すなわち、略してケーブルリンクを介して、実現される。1つの実施形態において、モジュール間ポートは、イーサネット（Ethernet）（登録商標）などの従来のネットワークポートである。さらなる実施形態において、モジュール間ネットワークポートは、モジュールの外部から物理的にアクセス可能である。

#### 【0024】

コンピューティングシステム（10）は、モジュール間ポートを介して達成されるモジュール間接続のネットワークを介して相互接続される、複数の処理モジュールを含む。処理モジュールは、複数の内部接続された処理ノードを同様に含む。1つの実施形態において、1つの処理ノードは1つのモジュール間ポートを所有し、1つのモジュール間ポートは1つのモジュール間接続をなす。したがって、処理モジュールは、典型的には、処理モジュール内に存在するモジュール間ポートの数を上回らない、限定数のモジュール間接続をなす。したがって、構築されるコンピューティングシステムは、処理ノードが主ネットワーク内の他の全ての処理ノードと直接的に接続されず、それらを追跡しないため、処理ノードの数が増える際、コンピューティングシステム内の高速トラフィックへの圧力を軽減する。代わりに、各処理ノードは、その直接的に接続された隣接ノードの状態を追跡すればよい。したがって、構築されるコンピューティングシステムは、さらに、異なる処理モジュールに配置される処理ノードの中のデータトラフィックの制御を、データパケットを限定数のモジュール内とモジュール間との対のステップを介して送信することにより、簡略化する。

#### 【0025】

コンピューティングシステムフレームワークは、最少で2つの層の接続を含む：すなわち、モジュール内接続およびモジュール間接続である。以下に説明されるように、より多い層のモジュール間接続が、モジュール間接続の接続スキームに応じて、提供されてもよい。以下に説明されるように、接続スキームを介して相互接続される処理モジュールの群は、ブロックと称される。同じスキームまたは異なるスキームのいずれから形成されるかにかかわらず、複数のブロックは、モジュール間接続を介してさらに接続され、上層のブロックを形成してもよい。ブロックまたは上層のブロックは、最終的にシステムと一体化する。考察の目的で、以下を含む3階層システムが参照されることに留意されたい：1）処理ノードから形成される処理モジュール、2）処理モジュールから形成されるブロック、および、3）ブロックから形成されるシステム。ただし、複数のブロックのレベルを設けることにより、3階層より多いシステムも可能である。したがって、別に明示されない限り、システムへの各参照は同等に適用され、任意の数のブロックの層に置換可能である。

## 【 0 0 2 6 】

中間にデータパケットトラフィック制御の追加の層が、少なくとも1つのネットワークスイッチを含んで達成されてよい。ネットワークスイッチはマルチポートブリッジであり、すなわち、Open Systems Interconnectionモデルの階層2で作動する実要素である。図2は、1つの実施形態による、統合ストレージ、処理、およびネットワークスイッチを組み込むネットワークスイッチング構造を有するコンピューティングシステムフレームワークを示すブロック図である。少なくとも1つのネットワークスイッチが、コンピューティングシステムで使用され得る。ネットワークスイッチは、処理モジュール内の処理ノード、ブロック内の処理モジュール、システム内のブロック、または、それらの組み合わせの間の接続を達成するために使用され得る。

10

## 【 0 0 2 7 】

コンピューティングシステムフレームワーク(20)は3つのレベルを包含するが、当業者により認識されるように、さらに多くのレベルも可能である。第1のレベルにおいて、複数の処理ノード(1)は処理モジュール(2)を形成する；第2のレベルにおいて、複数の処理モジュールはブロックを形成する；(3)第3のレベルにおいて、複数のブロックはシステムを形成する。処理モジュール(2)内で、処理ノードの各々に配置されるモジュール内ポート(3)は、モジュール内接続(4)を形成する。2つの処理モジュール(2)の間に、モジュール間ポート(5)は、モジュール間接続(6)を形成する。モジュール間接続は、さらに、ブロック間を接続してもよい(図示せず)。

20

## 【 0 0 2 8 】

加えて、モジュール内接続の少なくとも1つは、ネットワークスイッチまたは回路、およびネットワークスイッチで使用されるプロトコルを使用して、実現され得る。ネットワークスイッチまたは回路、およびネットワークスイッチで使用されるプロトコルを介して実現されるモジュール内接続は、ノードレベルネットワークスイッチ接続と称され、したがって、使用されるネットワークスイッチは、ノードレベルネットワークスイッチ(17)と称される。モジュール間接続の少なくとも1つは、ネットワークスイッチまたは回路、およびネットワークスイッチで使用されるプロトコルを使用して実現され得る。ネットワークスイッチまたは回路、およびネットワークスイッチで使用されるプロトコルを介して実現されるモジュール間接続は、ネットワークスイッチが提供する接続の層に応じて、モジュールレベルネットワークスイッチ接続またはブロックレベルネットワークスイッチ接続と称される。ネットワークスイッチがブロック内の処理モジュールの中で接続を提供する場合、ネットワークスイッチは、モジュールレベルネットワークスイッチ(18)と称される；ネットワークスイッチが2つのブロックに存在する処理モジュールの中で接続を提供する場合、ネットワークスイッチは、ブロックレベルネットワークスイッチ(19)と称される。

30

## 【 0 0 2 9 】

さらなる実施形態において、モジュール間接続は、ケーブルリンクとネットワークスイッチとのハイブリッドで実現される。ネットワークスイッチを含めることで、必要なケーブルの長さが減り、拡張性の問題に対処でき、および、データパケットトラフィック制御の追加的な層が提供される。

40

## 【 0 0 3 0 】

さらなる実施形態において、ノードレベルネットワークスイッチ、モジュールレベルネットワークスイッチ、およびブロックレベルネットワークスイッチは、個々にまたは結合して、またはケーブルリンクと結合して、コンピューティングシステムに存在して接続を形成し、データパケット転送を促進、修正、および変更してもよい。接続は、リング、バス、スター、メッシュ、ツリー、またはハイブリッドなどのトポロジを使用して、実現されてよい。1つの実施形態において、ツリートポロジは、ノードレベルネットワークスイッチをルートとして有して使用される。ツリートポロジは、複数のスタートポロジを、まとめてバス上へ統合する。このバス/スターハイブリッド手法は、バスまたはスターより、非常に優れた将来のネットワーク拡張性を支援する。ハイブリッドトポロジにおいて、

50

2つ以上のネットワークトポロジの組み合わせが使用される。

【0031】

処理モジュールは複数の処理ノードを包含し、各々がモジュール間ポートを包含する。したがって、処理モジュールはマルチポートデバイスであり、ネットワークスイッチとして使用され得る。1つの実施形態において、処理モジュールはネットワークスイッチの役割を果たす。そのようなスイッチは、処理モジュール模擬スイッチと称される。

【0032】

処理モジュール模擬スイッチは、ノードレベルネットワークスイッチ、モジュールレベルネットワークスイッチ、ブロックレベルネットワークスイッチ、または、それらの組み合わせのレベルで、存在し得る。ノードレベル、モジュールレベル、およびブロックレベルに存在する処理モジュール模擬スイッチは接続を形成し、データパケット転送を促進、修正、および変更してもよい。コンピューティングシステムは、さらに、ノードレベル、モジュールレベル、ブロックレベル、または、それらの組み合わせで、ネットワークスイッチと処理モジュール模擬スイッチとの両方を実現してもよい。

【0033】

モジュール間接続スキームの例が、以下に説明される。これらの例は例示的であり、限定を意味するものではない。当業者により認識されるように、他の構成、トポロジ、配列、および、接続、ポート、および処理モジュールの置換が可能である。

【0034】

図を簡略化するために、以下の例は、各処理モジュールがP個のモジュール間ポートを有する、M個の処理モジュールを含むコンピューティングシステムを想定しており、ここで、MとPは両方とも整数である。Pがゼロに等しい場合、処理モジュール間の接続はない。このようなシステムは、厳格に制限された有用性を有する。Pが1に等しい場合、システムは連結されたモジュールの組を含むであろう。このようなシステムは、処理を行う際、1つのモジュールより2つのモジュールの方が効率的であろうという点において、全く連結されない処理モジュールを有する場合より、わずかに優れた有用性を有する。Pが1より大きい場合、システムは、連結された処理モジュールの大きなネットワークを形成し、より好都合である。モジュール間接続の数が増えると、転送データパケットにおけるコンピューティングシステムの利点が大きくなる。

【0035】

一般的に、各モジュールと関連付けられるモジュール間ポートの数は、同等である必要はない。各モジュール間ポートは1つの処理ノードと関連付けられ、各モジュール間ポートは、別の処理モジュールの別のモジュール間ポートと、切断または接続されるかのいずれかである。例えば、コンピューティングシステムは、各処理モジュールが最少でP個のモジュール間ポートを有する、M個の処理モジュールを有してもよい。このようなシステムにおいて、Pを超えるモジュール間ポートは接続されないままであってよく、または、他のモジュール間ポートと接続されていてもよい。

【0036】

1つの処理モジュールのモジュール間ポートを他の処理モジュールのモジュール間ポートと無作為に接続

1つの実施形態において、1つの処理モジュールに配置されるモジュール間ポートは、他の処理モジュールに配置されるモジュール間ポートと無作為に接続される。このスキームの利点は、製造および保守が簡易であり、したがって、安価なことである。このスキームの不利な点は、接続の形成が無作為であることが原因で、一部の処理モジュールが接続されない一方で、一部の他の処理モジュールが過度に接続されてしまうリスクである。しかしながら、同じ無作為性は、現代のデータセンタ設定においてその傾向にあるように、処理モジュールの数が大きい場合、処理モジュールの過度な接続または少ない接続から生じる重大な問題は起きないことを保証している。大数の法則によると、膨大な数の同じ実験の試行から得られる結果の平均は、期待値に近くなり、試行が多く行われるほど近づく傾向にある。したがって、Mが大きくなるにつれて、特定の処理モジュールと接続される

10

20

30

40

50

異なる処理モジュールの数は、ますます  $P$  に近づく と期待される。

#### 【 0 0 3 7 】

さらなる実施形態において、処理モジュールの各々におけるモジュール間ポートの数は異なってもよい。さらなる実施形態において、モジュール間ポートの一部は、切断されている。

#### 【 0 0 3 8 】

二分スパニング接続スキームを使用して 1 つの処理モジュールのモジュール間ポートを他の処理モジュールのモジュール間ポートと接続

1 つの実施形態において、二分スパニングシステムが、1 つの処理モジュールに配置されるモジュール間ポートを、別の処理モジュールに配置されるモジュール間ポートと接続するために使用される。このスキームでは、処理モジュールの数  $M$  は 2 の累乗であり、処理モジュールは最初に  $M / 2$  個の群に分割される。結果として、群ごとに 2 つの処理モジュールが割り当てられる。 $M / 2$  個の群の各々の内部で、2 つの処理モジュールは、第 1 のモジュール間ポートを介して接続される。次に、処理モジュールは、 $M / 4$  個の群に分割され、群ごとに 4 つの処理モジュールとなる。 $M / 4$  個の群の各々の内部で、4 つの処理モジュールは、第 2 のモジュール間ポートを介して、1 と 3 および 2 と 4 が、それぞれ接続される。次に、処理モジュールは  $M / 8$  個の群に分割され、群ごとに 8 つの処理モジュールとなる。 $M / 8$  個の群の各々の内部で、8 つの処理モジュールは、第 3 のモジュール間ポートを介して、1 と 5、2 と 6、3 と 7、および 4 と 8 が、それぞれ接続される。接続は、このようなパターンで、使用されるモジュール間ポートの数が  $P$  または  $\log_2 M$  のどちらか小さい方に到達するまで、継続的に増大する。

#### 【 0 0 3 9 】

二分スパニング接続スキームは、例を用いて、より良好に図示される。図 3 は、例として、処理モジュールの各々に 4 つのモジュール間ポート ( $P = 4$ ) を有する 16 個の処理モジュール ( $M = 16$ ) を含む、コンピューティングシステムフレームワークとして図示される、二分スパニング接続スキームを示す図である。処理モジュールは、モジュール 1、2、... 16 まですべて省略されている。モジュール間ポートは、1、2、3、および 4 と番号が付けられている。当業者は、処理モジュールおよびモジュール間接続ポートの数が例示目的に過ぎず、限定の意味はないことを理解し得る。別の数の  $M$  および  $P$  も、可能である。 $M$  個のモジュールの各々において、モジュール内接続を介して接続される多様な数の処理ノードが存在する。

#### 【 0 0 4 0 】

したがって、二分スパニングスキームでは、処理モジュールの以下の対が、モジュール間ポート 1 を介して接続される：1 と 2、3 と 4、5 と 6、7 と 8、9 と 10、11 と 12、13 と 14、および 15 と 16 (図 3)。処理モジュールの以下の対が、モジュール間ポート 2 を介して連結される：1 と 3、2 と 4、5 と 7、6 と 8、9 と 11、10 と 12、13 と 15、および 14 と 16 (図 3)。処理モジュールの以下の対が、モジュール間ポート 3 を介して連結される：1 と 5、2 と 6、3 と 7、4 と 8、9 と 13、10 と 14、11 と 15、および 12 と 16 (図 3)。最後に、処理モジュールの以下の対が、モジュール間ポート 4 を介して連結される：1 と 9、2 と 10、3 と 11、4 と 12、5 と 13、6 と 14、7 と 15、および 8 と 16 (図 3)。

#### 【 0 0 4 1 】

したがって、二分スパニングスキームは、 $M$  個の処理モジュールおよび各処理モジュールにおける  $P$  個以上のモジュール間ポートを有するシステムとして、一般化され得る。各モジュール間ポートは、1 つの処理ノードと正確に関連付けられる。 $M$  は、2 の累乗である。モジュール間接続は、以下の規則にしたがって実現される：

i . モジュール  $2N + 1$  のポート 1 は、モジュール  $2N + 2$  のポート 1 と接続され、ここで、 $N$  は  $[0, 1, \dots, (M / 2) - 1]$  である。

ii . モジュール  $4N + 1$  および  $4N + 2$  のポート 2 は、モジュール  $4N + 3$  および  $4N + 4$  のポート 2 と、それぞれ接続され、ここで、 $N$  は  $[0, 1, \dots, (M / 4) - 1]$

」である。

$i i i$  . 一般的に、モジュール  $(2^P * N + 1)$  から  $(2^P * N + 2^{P-1})$  のポート  $P$  は、モジュール  $(2^P * N + 2^{P-1} + 1)$  から  $(2^P * N + 2^P)$  のポート  $P$  と、それぞれ接続され、ここで、 $N$  は  $[0, 1, \dots, (M / 2^{P+1}) - 1]$  であり、 $P$  は  $[1, 2, \dots, P]$  である。

#### 【0042】

したがって、二分スパニング接続スキームは、最初に、全ての処理モジュールが第1のモジュール間ポートを介して対となって接続され、続いて、対の各々が第2のモジュール間ポートを介して、さらに別の対と対になり、対の対を形成して、未使用のモジュール間ポートがなくなるまで、または、全ての処理モジュールが対となり1つの対になるまで、対になり続ける、接続スキームを表す。したがって、1つの実施形態において、処理モジュール内の接続は、次のようになされる：1、処理モジュールの1つを処理モジュールの別の1つと対にして、対になった処理モジュールを接続し、対になった処理モジュール内の処理ノードのモジュール間ポートを介して、接続対を形成する；2、接続対の1つを接続対の別の1つと対にして、対を接続し、新規に接続された対の未使用のモジュール間ポートを介して、新規の接続対を形成する；および、3、新規の接続対を、全てのモジュール間ポートが利用されるまで、または、全ての処理モジュールが接続されるまで、対にして接続する。

#### 【0043】

モジュール間ポートの数  $P$  が  $\log_2 M$  と等しい場合、コンピューティングシステムは完全に接続される可能性があり、すなわち、各処理モジュールは、直接的または間接的に、別の処理モジュールと接続されている。接続の総数は、 $(M * P) / 2$  に等しい。

#### 【0044】

$P$  が  $\log_2 M$  より小さい場合、コンピューティングシステムは、完全には接続されない。例えば、 $P$  が  $\log_2 M - 1$  と等しい場合、コンピューティングシステムは、各々の半分は内部的に接続されるが、他の半分は接続しない、2つの半分から構成される。 $P$  が  $\log_2 M - 2$  と等しい場合、コンピューティングシステムは、各々の4分割部分は内部的に接続されるが、他の4分割部分は接続しない、4つの4分割部分から構成される。 $P$  が  $\log_2 M - 3$  と等しい場合、コンピューティングシステムフレームワークは、内部的に接続される処理モジュールの集合の8つの等しいセクションに分割される。コンピュータシステムの性能は、一般的に、処理モジュールが良好に接続されると向上するため、 $\log_2 M$  に近い  $P$  を有するのが好都合である。

#### 【0045】

$P$  が  $\log_2 M$  より大きい場合、コンピューティングシステムは、 $\log_2 M$  個のモジュール間ポートを使用して完全に接続されてよく、残りのモジュール間ポートは、処理モジュールごとに  $P - \log_2 M$  個だけ、未接続のままであるか、または、モジュール間ポートの中で追加的な接続を形成するかのいずれかであり得る。1つの実施形態において、追加的な接続は、コンピューティングシステムの中で無作為に形成される。さらなる実施形態において、追加的な接続は、二分スパニングスキームにより接続される2つのコンピューティングシステムにより形成される。

#### 【0046】

さらなる実施形態において、処理モジュールの各々のモジュール間ポートの数は、相互に異なる可能性がある。二分スパニングスキームでは、モジュール間ポートの数が処理モジュールの中で等しい必要はない。1つの実施形態において、二分スパニングスキームにおける処理モジュールの中での異なるモジュール間ポートの数は、コンピューティングシステム内の処理モジュールのモジュール間ポートの最小数を特定して、その最小数をコンピューティングシステム内の全ての処理モジュールのモジュール間ポートの数として割り当てることにより、実現され得る。

#### 【0047】

リングシステム接続スキームを使用して1つの処理モジュールのモジュール間ポートを

他の処理モジュールのモジュール間ポートと接続

リングシステムは、ブロックまたはシステム内の全ての処理モジュールを、最初と最後が相互に連結される、線形に接続された処理モジュールの列状に接続し、接続された処理モジュールの各々からの2つのモジュール間ポートを使用して、閉鎖した円、すなわち、リングを形成する、共通の特徴を有する、接続スキームである。全てのモジュールが正確に2つのポートを利用し、一方が前のモジュールと接続して一方が次のモジュールと接続する、接続システムは、基本リングシステムと称される。基本リングシステムは、リングに存在する処理モジュール間の追加的な接続を導入することにより、修正されてもよい。リングシステムは、明示されない限り、以下に説明するように、基本および修正されたリングシステムの両方を指す。

10

#### 【0048】

基本リングシステム。1つの実施形態において、1つの処理モジュールに配置されるモジュール間ポートは、他の処理モジュールに配置されるモジュール間ポートと、リングシステムを使用して接続される。図4は、例として、8個の処理モジュール( $M = 8$ )を含むコンピューティングシステムフレームワークとして図示される、リングシステム(40)接続スキームを示す図である。M個の処理モジュールを含むコンピューティングシステムフレームワークにおいて、リングシステムは、Nが $[1, 2, \dots, M - 1]$ で、処理モジュールNのモジュール間ポート1を処理モジュールN + 1のモジュール間ポート2と接続し、処理モジュールMのモジュール間ポート1を処理モジュール1のモジュール間ポート2と接続することにより、実現される。ポート番号またはモジュール番号は、リングシステムから逸脱せずに、自明に交換され得る。したがって、処理モジュールごとに2つのモジュール間ポートを使用して、リングシステムは、コンピューティングシステムの全ての処理モジュールを頭尾方式で線形的に接続し、最後の処理モジュールを最初の処理モジュールと連結することにより、線形配列を円形に閉じる。このスキームは、処理モジュールごとに2つのモジュール間ポートのみを使用して、コンピューティングシステムフレームワーク内の全ての処理モジュールをまとめて接続する利点を有する。しかしながら、システムの処理モジュールの数が大きくなる傾向にあるため、基本リングシステムにおいて1つの処理モジュールから別のモジュールへ移動するためのホップの数が大きくなる可能性がある(パケットがリングに沿っていずれかの方向に移動する場合、ホップの最大数は $M / 2$ ; パケットがリングに沿って1つの方向にのみ移動する場合、ホップの最大数はM)。

20

30

#### 【0049】

ショートカットを有するリングシステム。さらなる実施形態において、コンピューティングシステムフレームワークは、基本リングシステムで実現される。さらに、リングにおける互いに隣接しない2つの処理モジュール間で、2つの処理モジュールの1つに配置され線形接続に使用されない1つのモジュール間ポートを介して、2つの処理モジュールの他方に配置され線形接続に使用されない別のモジュール間ポートへ、接続がなされる。図5は、例として、ショートカットを有するリングシステム接続スキームを示す図である。線形的に接続される処理モジュール(2)により形成されるリング(41)において、2つの隣接しない処理モジュールが、モジュール間接続を介して接続され、ショートカット(42)を形成している。ショートカット(42)は、一部の待ち時間を減らす可能性がある。しかしながら、ショートカット(42)は、データパケットトラフィックパターンと良好に適合しても、しなくてもよい。

40

#### 【0050】

1つの実施形態において、ショートカットは無作為に形成される。別の実施形態において、複数のショートカットが形成される。さらに別の実施形態において、ショートカットは、データトラフィックパターンにしたがって、2つの処理モジュールを選択することにより形成される。さらに別の実施形態において、ショートカットは、データトラフィックパターンにしたがって、2つの処理モジュールを選択することにより修正される。

50

## 【 0 0 5 1 】

S 個ごとのホップを有するリングシステム。さらなる実施形態において、コンピューティングシステムフレームワークは、M 個の処理モジュールを有する基本リングシステムで実現される。さらに、処理モジュール N のポート P は、処理モジュール N + S のポート P + 1 と接続され、ここで、S は M の整数の約数であり、N は  $[x, x + S, \dots, x + (M - S)(S \text{ ずつ増加})]$  であり、x は  $[0, 1, 2, \dots, S - 1]$  である。これらの接続により、データパケットは、サイズ S ごとにリングの周囲を移動できる。図 6 は、例として、S 個ごとのホップを有するリングシステム接続スキームを示す図である。例示目的で、5 個ごとのホップが示される。ステップのサイズ S は、 $M / 5$  に等しい。実線 (5 1) の五角形は、 $x = 0$  の 5 個ごとのホップを表す。破線 (5 2) の五角形は、 $x = 1$  の 5 個ごとのホップを表す。x の値は 0 から S - 1 の範囲である一方で、x の値の少なくとも 1 つが使用されるため、S 個ごとのホップの少なくとも 1 つの列がリングに存在する。より大きい値の x が使用されるほど、より多くのホップの列が作成され、潜在的な障害がさらに削減されるが、他のショートカットに利用できるポートが少なくなる。

10

## 【 0 0 5 2 】

S 個ごとの分割ホップを有するリングシステム。さらなる実施形態において、コンピューティングシステムフレームワークは、M 個の処理モジュールを有する基本リングシステムで実現される。さらに、モジュール N のポート P は、モジュール N + S - 1 のポート P と接続され、ここで、S は M の整数の約数であり、モジュール N の数は  $[x, x + S, \dots, x + (M - S)(S \text{ ずつ増加})]$  であり、x は  $[0, 1, \dots, S - 1]$  である。これらの接続により、データパケットは、それぞれ長さ S - 1 および 1 ごとのペアにおいて、リングの周囲を移動できる。図 7 は、例として、S 個ごとの分割ホップを有するリングシステムを示す図である。例示目的で、5 個ごとの分割ホップが示される。ステップのサイズ S は、 $M / 5$  に等しい。分割ホップの最初の部分はサイズ S - 1 を有し、長い実直線 (7 1) により表される。分割ホップの第 2 の部分はリングに沿っており、1 のサイズを有する (7 1)。破線の五角形は、分割ステップを介する結果的なホップを表す。接続スキームは、ポート 1 および 2 で構築されるリングを利用することにより、各モジュールの 1 つの追加的なポートを使用するのみであるという利点を有する。この接続スキームにおいて、3 個より多いポートが必要な処理モジュールはないが、待ち時間が発生する。

20

## 【 0 0 5 3 】

S 個ごとの調整可能分割ホップを有するリングシステム。さらなる実施形態において、コンピューティングシステムフレームワークは、M 個の処理モジュールを有する基本リングシステムで実現される。さらに、モジュール N のポート P は、モジュール N + S - r のポート P と接続されており、ここで、S は M の整数の約数であり、N は  $[x, x + S, \dots, x + (M - S)(S \text{ ずつ増加})]$  であり、x は  $[0, 1, \dots, M / S - 1]$  である。これらの接続により、データパケットは、それぞれ長さ S - r および r ごとのペアにおいて、リングの周囲を移動できる。図 8 は、例として、S 個ごとの調整可能分割ホップを有するリングシステムを示す図である。例示目的で、5 個ごとのホップが示される。ステップのサイズ S は、 $M / 5$  に等しい。破線の五角形は、S 個ごとの集合ホップを表す。実直線は、S - r 個のホップと r 個のホップを表し、S - r 個のホップおよび r 個のホップが一緒に、S 個ごとの調整可能分割ホップ (8 1) を形成する。r は  $1 < r < S$  の値を有する整数である。この接続スキームは、ポート 1 および 2 で構築されるリングを利用することにより、各モジュールの 1 つのポートを使用するのみであるという利点を有する。

30

40

## 【 0 0 5 4 】

ネットワークスイッチを接続スキームへ組み込み。上述した接続スキームは、少なくとも部分的に、ネットワークスイッチまたは処理モジュール模擬スイッチのうちの少なくとも 1 つ、または、それらの組み合わせの使用を介して、実現されてよい。ネットワークスイッチまたは処理モジュール模擬スイッチは、ケーブルリンクまたは直接リンクのうちの少なくとも 1 つに置換されてもよい。さらに、ネットワークスイッチまたは処理モジュール模擬スイッチが、上述したように、ケーブルリンクまたは直接リンクを介して実現され

50



る接続スキームに加えて、導入されてもよい。ネットワークスイッチは、パケットトラフィックパターンに応じてショートカットを変更または増加することにより、他の有用性の中で、追加的な柔軟性を提供する。

【0055】

1つ以上のネットワークスイッチをシステムに組み込むか、あるいは、ケーブルリンクを介して接続する、いくつかの例が、以下に説明される。当業者に既知のように、他のスキーム、トポロジ、構成、実装、セットアップ、および組み合わせが可能である。

【0056】

1つの実施形態において、リングシステム接続スキーム内のショートカットは、1つ以上のネットワークスイッチで形成される。図9は、例として、ネットワークスイッチにより実現されるショートカットを有するリングシステム接続スキームを示す図である。ネットワークスイッチ(91)は、この実施形態において示されるように、ケーブルリンクに置換されてもよい。さらに、ネットワークスイッチ(91)は、補完ケーブルリンクであってもよい(図示せず)。

【0057】

さらなる実施形態において、S個ごとのホップを有するリングシステム接続スキーム内のショートカットが、1つ以上のネットワークスイッチにより実現される。図10は、例として、ネットワークスイッチにより実現されるS個ごとのホップを有するリングシステム接続スキームを示す図である。例示目的で、5個ごとのホップが示される。ステップのサイズSは、 $M/5$ と等しい。リングのサイズに応じて、2つ以上のネットワークスイッチが使用されてもよいが、1つのネットワークスイッチ(91)が例として表されている。処理モジュール $N+0*S$ 、 $N+1*S$ 、 $N+2*S$ 、...、 $N-S$ は、ネットワークスイッチを介して順に接続される。データパケットがリングを横断して両方向に切り替わる場合、最大ホップ数は $S+1$ (起点の処理モジュールからネットワークスイッチへの最短接続まで $S/2$ ホップ、スイッチを通過するのに1ホップ、および、ネットワークスイッチと接続される処理モジュールから最終目的地の処理モジュールまで $S/2$ ホップ)である。リングおよびネットワークスイッチは、ホイールの外周に配置される処理モジュール、中央に配置されるネットワークスイッチ、および、ネットワークスイッチからスポークを形成するS番目ごとの処理モジュールへの接続を伴って、ホイールとスポークのパラダイムを形成する。Sが大きい場合、ショートカットはケーブルリンクで実現されてもよい。これにより、ネットワークスイッチのポートの数を削減しながら、ネットワーク遅延時間に対処できる。

【0058】

さらなる実施形態において、S個ごとの分割ホップを有するリングシステム接続スキーム内のショートカットは、1つ以上のネットワークスイッチを介して実現される。図11は、例として、ネットワークスイッチにより実現されるS個ごとの分割ホップを有するリングシステム接続スキームを示す図である。例示目的で、5個ごとのホップが示される。ステップのサイズSは、 $M/5$ と等しい。処理モジュールNおよび処理モジュール $N+S-1$ はネットワークスイッチと接続されており、ここで、Sはモジュールの数Mの整数の約数、Nは $[1, 1+S, \dots, 1+(M-S)]$ (Sずつ増加)である。リングおよびネットワークスイッチは、ホイールの外周に配置される処理モジュール、中央に配置されるネットワークスイッチ、および、ネットワークスイッチからスポークを形成するS番目およびS番目-1ごとの処理モジュールへの接続を伴って、ホイールとスポークのパラダイムを形成する。加えて、接続スキームは、どちらもネットワークスイッチと接続される、2つの隣接した処理モジュール間のリングリンクを排除してもよい。この接続スキームの利点の1つは、処理モジュールのポートの必要最大数が2個に制限されることである。最後に、1つ以上のネットワークスイッチ(図示せず)が、以下で説明されるように、別のリングまたは外部デバイスを接続するために使用される可能性がある。

【0059】

さらなる実施形態において、S個ごとの調整可能分割ホップを有するリングシステム接

10

20

30

40

50

続スキーム内のショートカットは、1つ以上のネットワークスイッチにより実現される。図12は、例として、ネットワークスイッチにより実現されるS個ごとの調整可能分割ホップを有するリングシステム接続スキームを示す図である。例示目的で、5個ごとのホップが示される。ステップのサイズSは、 $M/5$ と等しい。処理モジュールNおよび $N+S-r$ はネットワークスイッチにより接続されており、ここで、SはMの整数の約数であり、Nは $[1, 1+S, \dots, 1+(M-S)(S\text{ずつ増加})]$ であり、rは $[2, 3, \dots, M/S-1]$ である。データパケットがリングを横断して両方向に切り替わってよい場合、最大ホップ数は、 $MAX((S-r), r)+1$  (起点の処理モジュールから、ネットワークスイッチへの最短接続まで $MAX((S-r), r)/2$ ホップ、スイッチを通過するのに1ホップ、および、ネットワークスイッチと接続される処理モジュールから最終目的地の処理モジュールまで $MAX((S-r), r)/2$ ホップ)である。リングおよびネットワークスイッチは、ホイールの外周に配置される処理モジュール、中央に配置されるネットワークスイッチ、および、ネットワークスイッチからスポークを形成するS番目およびS番目-rごとの処理モジュールへの接続を伴って、ホイールとスポークのパラダイムを形成する。加えて、1つ以上のネットワークスイッチ(図示せず)が、以下で説明されるように、別のリングまたは外部デバイスを接続するために使用される可能性がある。

#### 【0060】

上記に説明した例は、全ての可能性を網羅していない。ネットワークスイッチは、二分スパニング、ランダム接続、または、それらの接続スキームの組み合わせなど、他の接続スキームに組み込まれ得る。上述した例において、ネットワークスイッチは、ブロック内の2つの処理モジュールを接続するため、モジュールレベルで配置される。しかしながら、ネットワークスイッチは、ノードレベルで実現され得る。さらに、ネットワークスイッチは、以下で説明するように、ブロックレベルでも実現され得る。

#### 【0061】

最後に、ネットワークスイッチは、処理モジュール模擬スイッチに置換され得る。

#### 【0062】

コンピューティングシステムを処理モジュールおよびブロックで構築

データセンタで使用されるようなコンピューティングシステムは、処理モジュールのクラスタから構築され得る。典型的には、処理モジュールは、上記で説明した接続スキームの1つを介して接続され、ブロックを形成する。コンピューティングシステムフレームワークは、複数のブロックから、または、各層が複数のブロックを包含する、複数のブロックの層からさえ、形成され得る。したがって、構築されるコンピュータシステムは、拡張可能であり、機動的で、耐障害性を有する。

#### 【0063】

1つの実施形態において、2つのブロック間の接続は、1つのブロックの1つの処理モジュールに配置されるモジュール間ポートを、他のブロックの別の処理モジュールに配置されるモジュール間ポートと連結するケーブルを介して、実現され得る。図13~図16は、例として、ブロックまたは外部デバイスを、ショートカットを有するリングシステム接続スキーム(図13)、S個ごとのホップを有するリングシステム接続スキーム(図14)、S個ごとの分割ホップを有するリングシステム接続スキーム(図15)、および、S個ごとの調整可能分割ホップを有するリングシステム接続スキーム(図16)と接続する、ケーブルリンクを示す図である。ブロックの処理モジュールは、リングまたは他の接続スキームを使用して、相互接続し得る。外部デバイスは、ネットワークスイッチであり得る。他の接続スキームおよび組み合わせが、可能である。

#### 【0064】

コンピューティングシステムを処理モジュールおよびブロックで構築、および1つ以上のネットワークスイッチを組み込み

別の実施形態において、2つのブロック間の接続は、ブロック内で実装されるネットワークスイッチを介して実現され得る。図17は、例として、ブロックまたは外部デバイス

10

20

30

40

50

から、ネットワークスイッチにより実現されるショートカットを有するリングシステム接続スキームへの接続を示す図である。図 18 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現される S 個ごとのホップを有するリングシステム接続スキームへの接続を示す図である。図 19 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現される S 個ごとの分割ホップを有するリングシステム接続スキームへの接続を示す図である。図 20 は、例として、ブロックまたは外部デバイスから、ネットワークスイッチにより実現される S 個ごとの調整可能分割ホップを有するリングシステム接続スキームへの接続を示す図である。ブロックの処理モジュールは、リングまたは他の接続スキームを使用して、相互接続され得る。外部デバイスは、ネットワークスイッチであり得る。他の接続スキームおよび組み合わせが、可能である。

10

**【0065】**

したがって、2つのブロックは、ケーブルリンク、ネットワークスイッチ、または、それらの組み合わせで接続され得る。コンピューティングシステムフレームワークは、複数のブロックを含んでよい。さらに、コンピューティングシステムフレームワークは、ブロックの各層が複数の相互接続されたブロックを包含する、複数のブロックの層を包含してもよい。

**【0066】**

ハイブリッドシステムを使用して1つの処理モジュールのモジュール間ポートを他の処理モジュールのモジュール間ポートと接続

1つの実施形態において、1つの処理モジュールに配置されるモジュール間ポートは、ハイブリッドシステムを使用して、他の処理モジュールに配置されるモジュール間ポートと接続される。ハイブリッドシステムは、2つ以上の接続スキームの組み合わせを指す。接続スキームを介して相互接続される処理モジュールの群は、上述したように、ブロックと称される。システムは、同じまたは異なるスキームから形成される、複数のブロックを含んでもよい。さらなる実施形態において、処理モジュールは、各々が二分スパニングスキームを使用して内部的に接続される、多数のブロックにグループ化される。ブロックは、リングスキームを使用して相互接続される。さらなる実施形態において、処理モジュールのブロックは、ラックまたはラック式の内部に設置されてもよい。最後に、確立されたモジュール間接続は、ランダム接続、二分スパニングスキーム、リングスキーム、ハイブリッドシステム、および、それらの組み合わせのうちの少なくとも1つを含む。

20

30

**【0067】**

コンピューティングシステムを調整して1つのモジュールのモジュール間ポートを他のモジュールのモジュール間ポートと接続

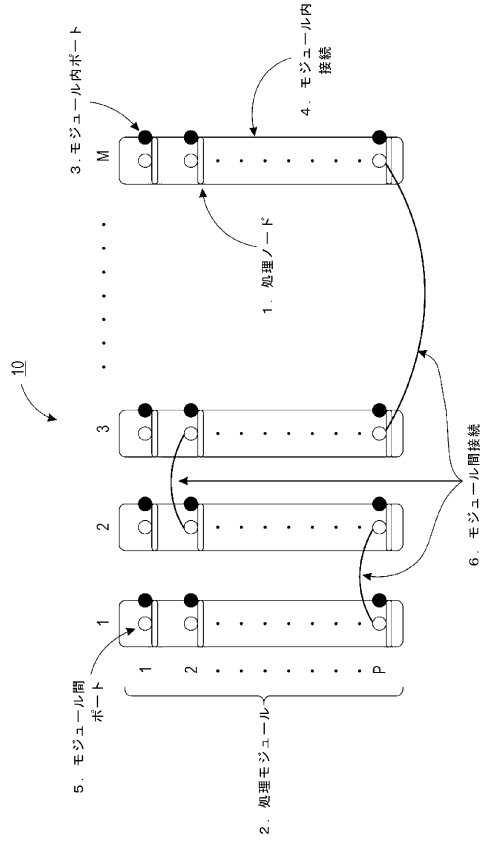
上記に説明した接続スキームは、調整システムを使用して、さらに改良または修正されてもよい。1つの実施形態において、コンピューティングシステムフレームワークは、データパケットを1つの処理ノードから別の処理ノードへ、確立されたモジュール間接続を介して転送する必要がある、タスクまたは割り当てを行う。データパケットトラフィックパターンが測定され、トラフィック障害が特定され、モジュール間接続がパケットトラフィックパターンに基づいて変更されて、トラフィックの流れを最適化する。

**【0068】**

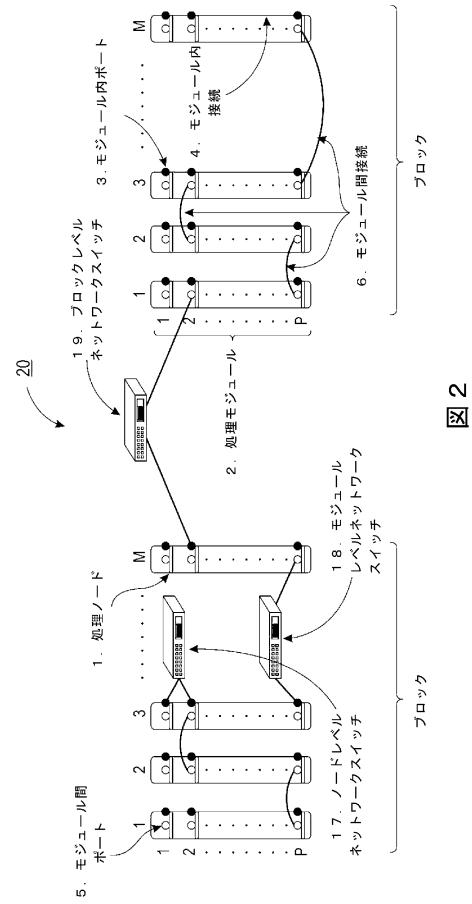
40

説明された接続スキームおよびネットワークプロトコルは、適合ルーティングアルゴリズム、非適合ルーティングアルゴリズム、データルーティング、マルチパスルーティング、および階層的ルーティングを含む、ネットワークパケットルーティングアルゴリズムに適応され得る。

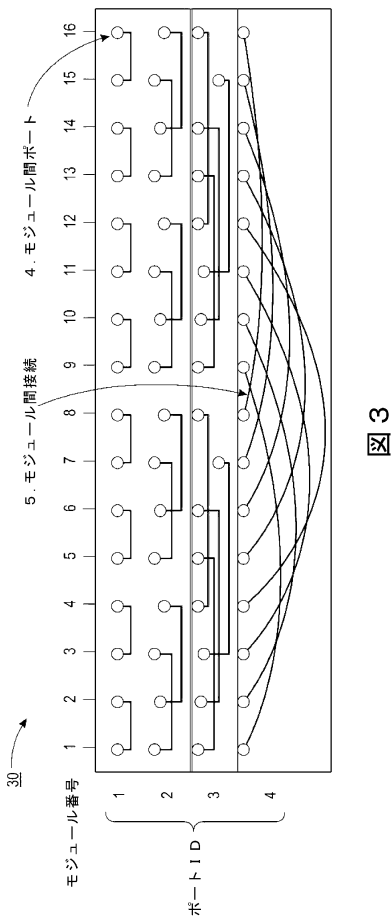
【図 1】



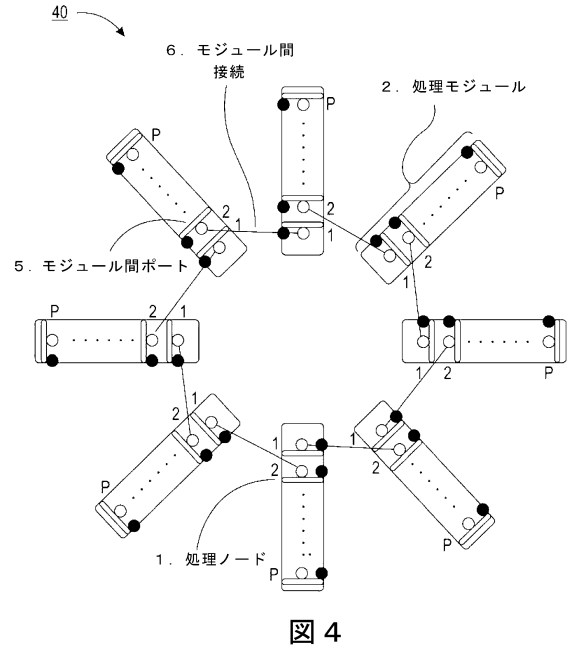
【図 2】



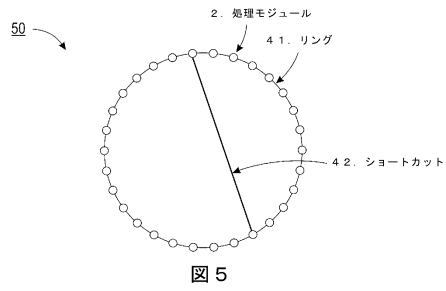
【図 3】



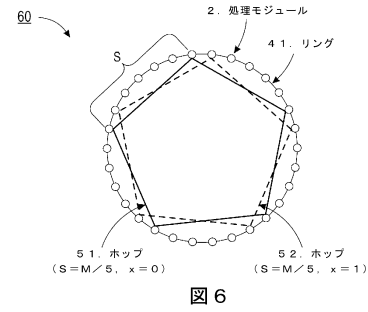
【図 4】



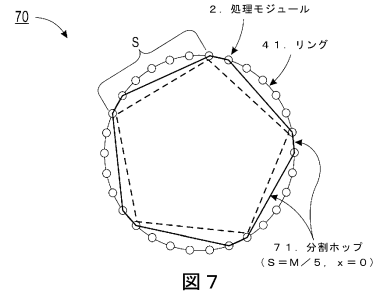
【図 5】



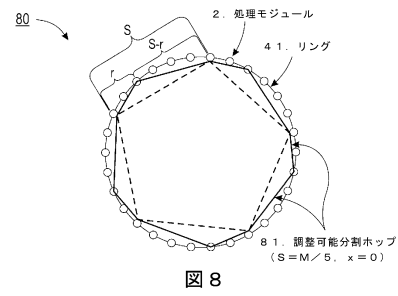
【図 6】



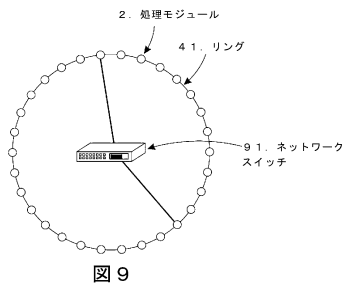
【図 7】



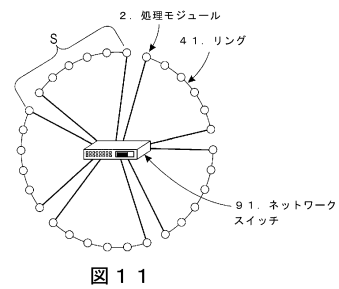
【図 8】



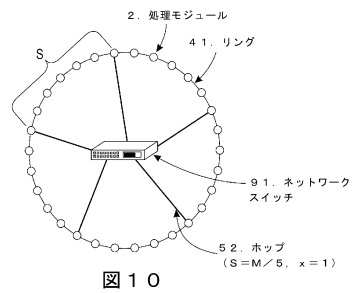
【図 9】



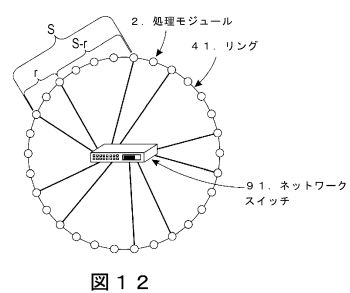
【図 11】



【図 10】



【図 12】



【図 13】

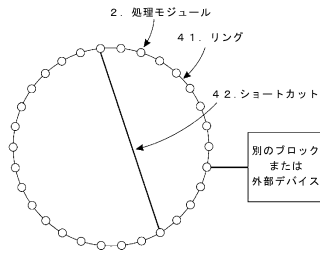


図 13

【図 15】

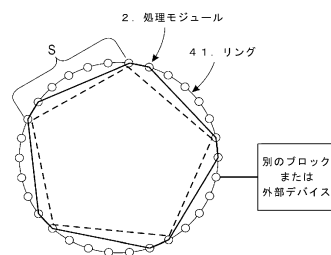


図 15

【図 14】

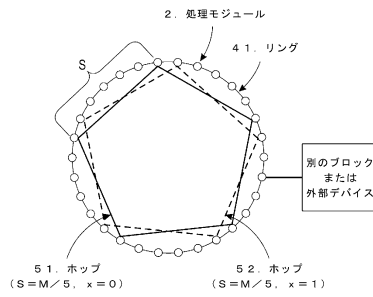


図 14

【図 16】

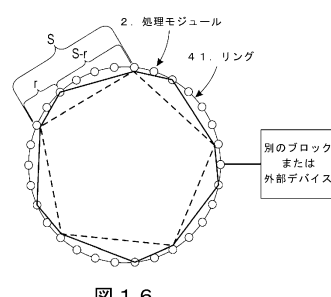


図 16

【図 17】

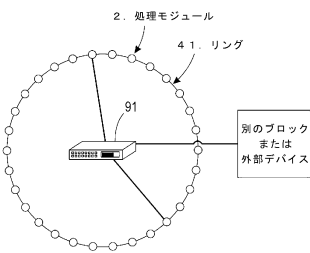


図 17

【図 19】

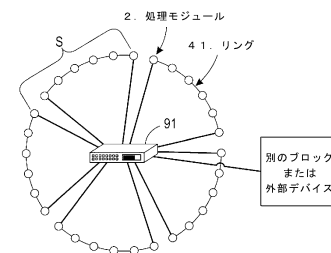


図 19

【図 18】

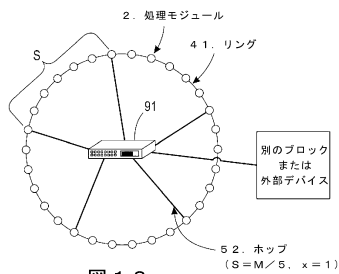


図 18

【図 20】

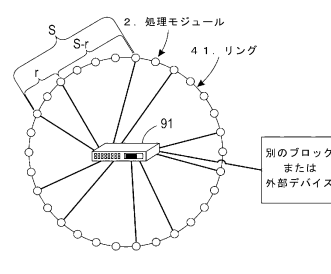


図 20

## フロントページの続き

(51)Int.Cl. F I  
G 0 6 F 13/10 3 4 0 A

(74)代理人 100109070  
弁理士 須田 洋之

(74)代理人 100109335  
弁理士 上杉 浩

(74)代理人 100120525  
弁理士 近藤 直樹

(74)代理人 100139712  
弁理士 那須 威夫

(72)発明者 ダニエル・デイヴィス  
アメリカ合衆国 カリフォルニア州 9 4 3 0 6 パロアルト スタンフォード・アベニュー 3  
7 4

審査官 鈴木 肇

(56)参考文献 特開2001-356847(JP,A)  
米国特許出願公開第2013/0073814(US,A1)  
米国特許出願公開第2014/0032731(US,A1)  
米国特許出願公開第2013/0250802(US,A1)  
米国特許出願公開第2011/0258340(US,A1)

(58)調査した分野(Int.Cl.,DB名)  
H 0 4 L 1 2 / 0 0 - 1 2 / 9 5 5  
H 0 4 L 1 3 / 0 0 - 1 3 / 1 8  
H 0 4 L 2 9 / 0 0 - 2 9 / 1 4  
G 0 6 F 1 3 / 1 0 - 1 3 / 1 4

(54)【発明の名称】統合ストレージ、処理、およびネットワークスイッチを組み込むネットワークスイッチング構造  
を有するコンピューティングシステムフレームワーク、および、それらを作成および使用する方  
法