

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5973570号
(P5973570)

(45) 発行日 平成28年8月23日 (2016. 8. 23)

(24) 登録日 平成28年7月22日 (2016. 7. 22)

(51) Int. Cl.	F I
HO 4 L 12/715 (2013. 01)	HO 4 L 12/715
HO 4 L 12/717 (2013. 01)	HO 4 L 12/717
HO 4 L 12/725 (2013. 01)	HO 4 L 12/725

請求項の数 16 (全 22 頁)

(21) 出願番号	特願2014-524465 (P2014-524465)	(73) 特許権者	598036300
(86) (22) 出願日	平成24年7月26日 (2012. 7. 26)		テレフオンアクチーボラゲット エルエム
(65) 公表番号	特表2014-525692 (P2014-525692A)		エリクソン (パブル)
(43) 公表日	平成26年9月29日 (2014. 9. 29)		スウェーデン国 スtockホルム エスー
(86) 国際出願番号	PCT/IB2012/053833		1 6 4 8 3
(87) 国際公開番号	W02013/021304	(74) 代理人	100079108
(87) 国際公開日	平成25年2月14日 (2013. 2. 14)		弁理士 稲葉 良幸
審査請求日	平成27年6月26日 (2015. 6. 26)	(74) 代理人	100109346
(31) 優先権主張番号	13/208, 251		弁理士 大貫 敏史
(32) 優先日	平成23年8月11日 (2011. 8. 11)	(72) 発明者	ヤダバリ, キラン
(33) 優先権主張国	米国 (US)		アメリカ合衆国, カリフォルニア州 9 5
			1 3 5, サン ノゼ, デラクロワ コート
			4 2 8 0

最終頁に続く

(54) 【発明の名称】 分割アーキテクチャ・ネットワークにおけるOSPFの実装

(57) 【特許請求の範囲】

【請求項 1】

分割アーキテクチャ・ネットワークの複数の領域のうちの1つのための複数のコントローラのうちの1つとして機能するネットワーク要素内に実装される方法であって、前記コントローラは、前記分割アーキテクチャ・ネットワークの前記領域に制御プレーンを提供し、前記コントローラは、分割アーキテクチャ・ネットワークの前記領域にデータ・プレーンを提供する複数のスイッチからはリモートにあり、前記コントローラは、前記分割アーキテクチャ・ネットワークの他の領域の他のコントローラ、及び前記分割アーキテクチャ・ネットワークを含むネットワークの従来型ルータに、限定領域内リンク・コスト・データを提供することによって、前記分割アーキテクチャ・ネットワークの前記複数領域にまたがる最適化ルーティングを容易にするものであって、前記限定領域内リンク・コスト・データは、すべての内部リンク・コスト・データを提供することなく、前記コントローラの領域のそれぞれの可能な最短パスのトラバースに関するコストを提供し、

前記コントローラの前記領域内の各境界スイッチを含む前記分割アーキテクチャ・ネットワーク内の前記コントローラの前記領域のトポロジを習得するステップであって、前記コントローラの前記領域内の各境界スイッチは、前記コントローラの前記領域を、前記分割アーキテクチャ・ネットワークの別の領域、又は前記ネットワーク内の前記従来型ルータのうちの1つにリンクさせる、少なくとも1つの外部ポートを有する、習得するステップと、

前記コントローラの前記領域内の各境界スイッチ・ペア間の最短パスを計算するステッ

10

20

プと、

前記コントローラのルーティング・テーブル内の各境界スイッチ・ペア間の各最短パスのコストを記憶するステップと、

前記分割アーキテクチャ・ネットワーク内の各近隣コントローラ又は前記ネットワーク内の近隣従来型ルータを、ハロー・プロトコルを使用して識別するステップであって、各近隣コントローラは、前記コントローラの前記領域の少なくとも1つの外部ポートを通じてアクセス可能な、前記分割アーキテクチャ・ネットワークの別の領域内のスイッチを制御する、識別するステップと、

リンク状態データベースを各近隣コントローラ及び近隣従来型ルータと交換するステップであって前記リンク状態データベースは、各境界スイッチ・ペア間の各最短パスの前記コストを含む、交換するステップと、

前記コントローラをツリーのルートとして備える前記ネットワークについて最短パス・ツリーを計算するステップと、

前記最短パス・ツリーに従って転送を実装するために、前記コントローラの前記領域のスイッチ内の転送テーブルを更新するステップと、
を含む、方法。

【請求項2】

各近隣コントローラを識別するステップが、

前記コントローラの前記領域の各外部ポートでハロー・パケットを送信するステップと、

、
前記コントローラの前記領域の少なくとも1つの外部ポートを通じて、各近隣コントローラからハロー・パケットを受信するステップと、
を更に含む、請求項1に記載の方法。

【請求項3】

リンク状態データを交換するステップが、各近隣コントローラにリンク状態アドバタイズを送信するステップを更に含み、前記リンク状態アドバタイズは前記コントローラの各外部リンクへのコストを含む、請求項1に記載の方法。

【請求項4】

リンク状態データを交換するステップが、各近隣コントローラにリンク状態アドバタイズを送信するステップを更に含み、前記リンク状態アドバタイズは、前記領域内の各境界スイッチ・ペアに関する前記最短パスの前記コストを含む、請求項1に記載の方法。

【請求項5】

更新されたリンク状態データを各近隣コントローラにアドバタイズするステップを更に含む、請求項1に記載の方法。

【請求項6】

前記更新リンク状態データをアドバタイズするステップが、更新されたリンク状態データを伴うリンク状態アドバタイズ・メッセージを各近隣コントローラに送信するステップを更に含む、請求項5に記載の方法。

【請求項7】

転送テーブルを更新するステップが、OpenFlowプロトコルを使用して前記コントローラの前記領域内の各スイッチの転送テーブルを更新するステップを更に含む、請求項1に記載の方法。

【請求項8】

前記コントローラでオープン・ショーテスト・パス・ファースト(OSPF)プロトコル・データ構造を開始するステップを更に含む、請求項1に記載の方法。

【請求項9】

分割アーキテクチャ・ネットワークの複数の領域のうちの1つについての複数のコントローラのうちの1つとして機能する、ネットワーク要素であって、前記コントローラは、前記分割アーキテクチャ・ネットワークの前記領域に制御プレーンを提供し、前記コントローラは、前記分割アーキテクチャ・ネットワークの前記領域にデータ・プレーンを提供

10

20

30

40

50

する複数のスイッチからはリモートにあり、前記コントローラは、前記分割アーキテクチャ・ネットワークの他の領域の他のコントローラ、及び前記分割アーキテクチャ・ネットワークを含むネットワークの従来型ルータに、限定領域内リンク・コスト・データを提供することによって、前記分割アーキテクチャ・ネットワークの前記複数領域にまたがる最適化ルーティングを容易にし、前記限定領域内リンク・コスト・データは、すべての内部リンク・コスト・データを提供することなく、前記コントローラの前記領域のそれぞれの可能な最短パスのトラバースに関するコストを提供し、

ネットワークを介してデータを受信するように構成された入口モジュールと、

前記ネットワークを介してデータを送信するように構成された出口モジュールと、

前記入口モジュール及び前記出口モジュールに結合されたネットワーク・プロセッサであって、コントローラ・モジュール、トポロジ習得モジュール、最短パス計算モジュール、近隣発見モジュール、及びリンク状態管理モジュールを備える、モジュール・セットを実行するように構成される、ネットワーク・プロセッサと、

前記トポロジ習得モジュールは、前記コントローラの前記領域内の各境界スイッチを含む前記分割アーキテクチャ・ネットワーク内の前記コントローラの前記領域のトポロジを決定するように構成され、各境界スイッチは、前記コントローラの前記領域を前記ネットワークの別の領域又は前記ネットワーク内の従来型ルータにリンクさせる、少なくとも1つの外部ポートを有し、

前記コントローラ・モジュールは、前記分割アーキテクチャ・ネットワーク内のコントローラの前記領域に制御プレーン機能を提供するように構成され、

前記最短パス計算モジュールは、近隣コントローラと共有することになる前記コントローラの前記領域内の各境界スイッチと前記従来型ルータとの間の最短パスを識別すること、及び、前記コントローラを前記ツリーのルートとして備える前記ネットワークについて最短パス・ツリーを計算することを、実行するように構成され、

前記近隣発見モジュールは、ハロー・プロトコルを使用して前記ネットワーク内の各近隣コントローラ及び従来型ルータを識別するように構成され、

前記リンク状態管理モジュールは、リンク状態データベースを、前記ネットワーク内の各近隣コントローラ及び従来型ルータと交換するように構成され、

前記リンク状態データベースは、前記コントローラの前記領域内の各境界スイッチ・ペアの間の各最短経路のコストを含み、

前記ネットワーク・プロセッサに通信可能に結合されたルーティング・テーブル記憶デバイスであって、前記コントローラの前記領域について、及び前記コントローラの前記領域の境界スイッチ間の、最短パス情報を含む、コントローラ・モジュールに関するルーティング・テーブルを記憶するように構成される、ルーティング・テーブル記憶デバイスと、
を備える、ネットワーク要素。

【請求項 10】

前記近隣発見モジュールが、前記領域の各外部ポートでハロー・パケットを送信すること、及び、前記領域の少なくとも1つの外部ポートを通じて、各近隣コントローラからハロー・パケットを受信することを、実行するように構成される、請求項9に記載のネットワーク要素。

【請求項 11】

前記リンク状態管理モジュールが、各近隣コントローラにリンク状態アドバタイズを送信するように更に構成され、前記リンク状態アドバタイズは前記コントローラの各外部リンクへのコストを含む、請求項9に記載のネットワーク要素。

【請求項 12】

前記リンク状態管理モジュールが、各近隣コントローラにリンク状態アドバタイズを送信するように更に構成され、前記リンク状態アドバタイズは、前記コントローラの前記領域内の各境界スイッチ・ペアに関する前記最短パスの前記コストを含む、請求項9に記載のネットワーク要素。

10

20

30

40

50

【請求項 13】

前記リンク状態管理モジュールが、更新されたリンク状態データを各近隣コントローラにアドバタイズするように更に構成される、請求項9に記載のネットワーク要素。

【請求項 14】

前記リンク状態管理モジュールが、前記更新されたリンク状態データを伴うリンク状態アドバタイズ・メッセージを、各近隣コントローラに送信するように更に構成される、請求項9に記載のネットワーク要素。

【請求項 15】

前記コントローラ・モジュールが、前記OpenFlowプロトコルを使用して前記コントローラの前記領域内の各スイッチの転送テーブルを更新するように更に構成される、請求項9に記載のネットワーク要素。

10

【請求項 16】

前記最短パス計算モジュールが、前記コントローラでオープン・ショーテスト・パス・ファースト (OSPF) プロトコル・データ構造を開始するように更に構成される、請求項9に記載のネットワーク要素。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、分割アーキテクチャ・ネットワークにおけるパスの探知又はルーティングに関する。具体的に言えば、本発明の実施形態は、分割アーキテクチャ・ネットワークを含むネットワーク及び従来のルータにおいてデータをルーティングするために、オープン・ショーテスト・パス・ファースト・プロセスを実装するための方法及びシステムに関する。

20

【背景技術】

【0002】

分割アーキテクチャ・ネットワーク設計は、制御プレーン及び転送又はデータ・プレーンとも呼ばれるネットワークの制御構成要素と転送構成要素との間に、分離を導入する。分割アーキテクチャ・ネットワークは、キャリアグレード・ネットワークのアクセス/集約ドメイン、モバイル・バックホール、クラウド・コンピューティング、及びマルチレイヤ (L3 & L2 & L1、光伝送ネットワーク (OTN)、波長分割多重 (WDM)) サポートで使用可能であり、これらはすべてネットワーク・アーキテクチャの構築ブロックの1つである。

30

【0003】

転送 (データ) プレーンと制御プレーンの両方を同じボックス内に統合する従来のネットワーク・アーキテクチャとは異なり、分割アーキテクチャはこれら2つの機能を切り離し、転送要素 (スイッチ) とは異なる物理位置にあり得るサーバ (コントローラ) 上で制御プレーンを実行する。分割アーキテクチャは、転送プラットフォームの機能及びハードウェアを簡略化し、ネットワークのインテリジェンス及び管理を、スイッチを監視するコントローラ・セットに集約する。従来のネットワーク・アーキテクチャにおける転送及び制御プレーンの密結合の結果として、通常、かなり複雑な制御プレーン及び手間のかかるネットワーク管理が生じている。これにより、新しいネットワーキング・デバイスの製作には費用が掛かり、これらデバイスの潜在的な展開のための新しいプロトコル及び技術に参入するためには高い障壁が設けられている。回線速度、ポート密度、及びパフォーマンスは急速に改善されているにもかかわらず、これらの機能を管理するためのネットワーク制御プレーン機構の進歩はかなり遅れている。

40

【0004】

分割アーキテクチャ・ネットワークにおいて、コントローラはスイッチから情報を収集し、適切な転送決定を計算して、これをスイッチに配信する。コントローラ及びスイッチは、制御プレーン・プロトコルを使用して情報の通信及び交換を行う。こうしたプロトコルの例がOpenFlowであり、コントローラと通信するためのオープンで標準的な方法をスイ

50

ッチに提供する。図1は、スイッチとコントローラとの間のOpenFlowインターフェースの概要を示す図である。OpenFlowコントローラは、転送テーブル（フロー・テーブル）を構成するためのセキュアなチャネルを用いてOpenFlowスイッチと通信する。

【0005】

OpenFlowスイッチ内の転送テーブルには、パケット・ヘッダ内のフィールドに関する合致を定義する規則、規則によって定義された合致の検出時に実行されるアクション、及び、データ・プレーン内のデータ・パケットの処理に関する統計の収集からなる、エントリが存在する。着信データ・パケットが特定の規則に合致する場合、関連付けられたアクションがデータ・パケット上で実行される。規則は、例えば、Ethernet MACアドレス、IPアドレス、IPプロトコル、TCP/IPポート番号、並びに着信ポート番号等の、プロトコル・スタック内のいくつかのヘッダからのキー・フィールドを含む。同様の特徴を有するデータ・パケット・セットはフローとして管理することができる。フローは、データ・パケット内の任意数の使用可能フィールド又はその組み合わせを使用して定義することができる。望ましくないフィールドについてのワイルドカードを使用して、使用可能フィールドのサブセットに関する合致に規則を制約することも可能である。

【0006】

分割アーキテクチャの制御プレーンとデータ・プレーンの分離により、ネットワーク制御論理を修正するタスクが容易になり、開発者が多様な新しいプロトコル及び管理アプリケーションを構築可能なプログラマチック・インターフェースが提供される。このモデルでは、データ・プレーン及び制御プレーンは独立に展開及びスケーリング可能であり、データ・プレーン要素のコストは削減される。

【発明の概要】

【発明が解決しようとする課題】

【0007】

分割アーキテクチャ・ネットワークの複数の領域のうちの1つのための複数のコントローラのうちの1つとして機能するネットワーク要素内に実装される方法を説明する。

【課題を解決するための手段】

【0008】

コントローラは、分割アーキテクチャ・ネットワークの領域に制御プレーンを提供し、ここでコントローラは、分割アーキテクチャ・ネットワークの領域にデータ・プレーンを提供する複数のスイッチからはリモートにある。コントローラは、分割アーキテクチャ・ネットワークの他の領域の他のコントローラ、及び分割アーキテクチャ・ネットワークを含むネットワークの従来型ルータに、限定領域内リンク・コスト・データを提供することによって、分割アーキテクチャ・ネットワークの複数領域にまたがる最適化ルーティングを容易にする。限定領域内リンク・コスト・データは、すべての内部リンク・コスト・データを提供することなく、コントローラの領域のそれぞれの可能な最短パスのトラバースに関するコストを提供する。方法は、コントローラの領域内の各境界スイッチ（border switch）を含む分割アーキテクチャ・ネットワーク内のコントローラの領域のトポロジを習得することを含むステップを含み、コントローラの領域内の各境界スイッチは、コントローラの領域を、分割アーキテクチャ・ネットワークの別の領域、又はネットワーク内の従来型ルータのうちの1つにリンクさせる、少なくとも1つの外部ポートを有する。ステップは、コントローラの領域内の各境界スイッチ・ペア間の最短パスを計算することを含む。ステップは、コントローラのルーティング・テーブル内の各境界スイッチ・ペア間の各最短パスのコストを記憶することを含む。ステップは、分割アーキテクチャ・ネットワーク内の各近隣コントローラ又はネットワーク内の近隣従来型ルータを、ハロー・プロトコル（hello protocol）を使用して識別することを含み、ここで各近隣コントローラは、コントローラの領域の少なくとも1つの外部ポートを通じてアクセス可能な、分割アーキテクチャ・ネットワークの別の領域内のスイッチを制御する。ステップは、リンク状態データベースを各近隣コントローラと交換することを含み、リンク状態データベースは、各境界スイッチ・ペア間の各最短パスのコストを含む。ステップは、コントローラをツリー

のルートとして備えるネットワークについて最短パス・ツリーを計算すること、及び、最短パス・ツリーに従って転送を実装するために、コントローラの領域のスイッチ内の転送テーブルを更新することを、含む。

【0009】

ネットワーク要素は、分割アーキテクチャ・ネットワークの複数の領域のうちの1つについての複数のコントローラのうちの1つとして機能する。コントローラは、分割アーキテクチャ・ネットワークの領域に制御プレーンを提供し、ここでコントローラは、分割アーキテクチャ・ネットワークの領域にデータ・プレーンを提供する複数のスイッチからはリモートにある。コントローラは、分割アーキテクチャ・ネットワークの他の領域の他のコントローラ、及び分割アーキテクチャ・ネットワークを含むネットワークの従来型ルータに、限定領域内リンク・コスト・データを提供することによって、分割アーキテクチャ・ネットワークの複数領域にまたがる最適化ルーティングを容易にする。限定領域内リンク・コスト・データは、すべての内部リンク・コスト・データを提供することなく、コントローラの領域のそれぞれの可能な最短パスのトラバースに関するコストを提供する。ネットワーク要素は、ネットワークを介してデータを受信するように構成された入口（ingress）モジュールと、ネットワークを介してデータを送信するように構成された出口（egress）モジュールとを備える。ネットワーク要素は、入口モジュール及び出口モジュールに結合されたネットワーク・プロセッサも含み、ネットワーク・プロセッサは、コントローラ・モジュール、トポロジ習得モジュール、最短パス計算モジュール、近隣発見モジュール、及びリンク状態管理モジュールを備える、モジュール・セットを実行するように構成される。トポロジ習得モジュールは、コントローラの領域内の各境界スイッチを含む分割アーキテクチャ・ネットワーク内のコントローラの領域のトポロジを決定するように構成され、各境界スイッチは、コントローラの領域をネットワークの別の領域又はネットワーク内の従来型ルータにリンクさせる、少なくとも1つの外部ポートを有する。コントローラ・モジュールは、分割アーキテクチャ・ネットワーク内のコントローラの領域に制御プレーン機能を提供するように構成される。最短パス計算モジュールは、近隣コントローラと共有することになるコントローラの領域内の各境界スイッチと従来型ルータとの間の最短パスを識別すること、及び、コントローラをツリーのルートとして備えるネットワークについて最短パス・ツリーを計算することを、実行するように構成される。近隣発見モジュールは、ハロー・プロトコルを使用して分割アーキテクチャ・ネットワーク内の各近隣コントローラを識別するように構成され、リンク状態管理モジュールは、リンク状態データベースを各近隣コントローラと交換するように構成される。リンク状態データベースは、コントローラの領域内の各境界スイッチ・ペア間の、各最短経路のコストを含み、ネットワーク要素は、ネットワーク・プロセッサに通信可能に結合されたルーティング・テーブル記憶デバイスも含む。ルーティング・テーブル記憶デバイスは、コントローラの領域について、及びコントローラの領域の境界スイッチ間の、最短パス情報を含む、コントローラ・モジュールに関するルーティング・テーブルを記憶するように構成される。

【0010】

本発明は、添付の図面の図中で、限定的ではなく単なる例示として示されており、図面内の同じ参照番号は同様の要素を示す。本開示における「実施形態」又は「一実施形態」という異なる言い方は、必ずしも同じ実施形態を指しておらず、こうした言い方は少なくとも1つを意味することに留意されたい。更に、特定の機能、構造、又は特徴が実施形態に関連して説明される場合、明示的に説明されるか否かにかかわらず、こうした機能、構造、又は特徴を他の実施形態に関連して実行することは、当業者の知識の範囲内であるものと考えられる。

【図面の簡単な説明】

【0011】

【図1】単純なOpenFlowネットワークに関する例示アーキテクチャの一実施形態を示す図である。

【図2】オープン・ショーテスト・パス・ファースト（OSPF）パケット・ヘッダを示

10

20

30

40

50

す図である。

【図 3】OSPF ハロー・パケットを示す図である。

【図 4】リンク状態アドバタイズ (advertisement) (LSA) ヘッダ・フォーマットを示す図である。

【図 5】ルータ LSA メッセージを示す図である。

【図 6】OSPF 外部リンク・コスト情報が交換される、例示ネットワークを示す図である。

【図 7】OSPF 内部リンク・コスト情報が交換される、例示ネットワークを示す図である。

【図 8】複数コントローラを備えた例示の複数領域 OpenFlow ネットワークの一実施形態を示す図である。

10

【図 9】OSPF プロセスを実装するネットワーク要素の一実施形態を示す図である。

【図 10】分割アーキテクチャ領域を備えたネットワーク内の OSPF ルーティングに関するプロセスの一実施形態を示すフローチャートである。

【図 11】OSPF リンク状態アドバタイズ・ヘッダ・フォーマットを示す図である。

【図 12】ルータ LSA メッセージを示す図である。

【図 13】分割アーキテクチャ OSPF ハロー・メッセージを示す図である。

【図 14】ネットワーク内で分割アーキテクチャ・コントローラから従来型ルータへハロー・メッセージを送信するプロセスの、一実施形態を示す図である。

【発明を実施するための形態】

20

【0012】

以下の説明では、多数の特定の細部が示される。しかしながら、本発明の実施形態はこれら特定の細部なしで実施可能であることを理解されよう。その他の場合、良く知られた回路、構造、及び技法は、この説明に関する理解を不明瞭にしないために詳細には示されていない。しかしながら、当業者であれば、こうした特定の細部なしに本発明が実施できることを理解されよう。当業者であれば、過度の実験なしに、含められた説明を用いて適切な機能を実装することができよう。

【0013】

フローチャートの動作は、図 7 ~ 図 9 及び図 14 の例示の実施形態を参照しながら説明される。しかしながら、図 10 のフローチャートの動作は、図 7 ~ 図 9 及び図 14 を参照しながら考察される実施形態以外の本発明の実施形態によって実行可能であり、図 7 ~ 図 9 及び図 14 を参照しながら考察される実施形態は、図 10 のフローチャートを参照しながら考察される動作とは異なる動作を実行可能であることを理解されたい。

30

【0014】

図内に示される技法は、1 つ又は複数の電子デバイス (例えばエンド・ステーション、ネットワーク要素、サーバ、又は同様の電子デバイス) 上に記憶され、実行される、コード及びデータを使用して実装可能である。こうした電子デバイスは、非一過性の機械読み取り可能又はコンピュータ読み取り可能記憶媒体 (例えば、磁気ディスク、光ディスク、ランダム・アクセス・メモリ、読み取り専用メモリ、フラッシュ・メモリ・デバイス、及び相変化メモリ) 等の、非一過性の機械読み取り可能又はコンピュータ読み取り可能媒体を使用して、コード及びデータを記憶及び (内部的に、及び/又は、ネットワークを介して他の電子デバイスと) 通信する。加えてこうした電子デバイスは、通常、1 つ又は複数の記憶デバイス、ユーザ入力/出力デバイス (例えば、キーボード、タッチ・スクリーン、及び/又はディスプレイ)、及びネットワーク接続等の、1 つ又は複数の他の構成要素に結合された、1 つ又は複数のプロセッサのセットを含む。プロセッサのセットと他の構成要素との結合は、通常、1 本又は複数のバス及びブリッジ (バス・コントローラとも呼ばれる) を介する。記憶デバイスは、1 つ又は複数の非一過性の機械読み取り可能又はコンピュータ読み取り可能記憶媒体、及び、非一過性の機械読み取り可能又はコンピュータ読み取り可能通信媒体を表す。したがって、所与の電子デバイスの記憶デバイスは、通常、その電子デバイスの 1 つ又は複数のプロセッサのセット上で実行するための、コード及

40

50

び／又はデータを記憶する。もちろん、本発明の実施形態の１つ又は複数の部分は、ソフトウェア、ファームウェア、及び／又はハードウェアの異なる組み合わせを使用して実装可能である。

【 0 0 1 5 】

本明細書で使用される場合、ネットワーク要素（例えば、ルータ、スイッチ、ブリッジ、又は同様のネットワーキング・デバイス）は、ネットワーク上の他の機器（例えば、他のネットワーク要素、エンド・ステーション、又は同様のネットワーキング・デバイス）を通信可能に相互接続するハードウェア及びソフトウェアを含む、１つのネットワーキング機器である。いくつかのネットワーク要素は、複数のネットワーキング機能（例えば、ルーティング、ブリッジング、スイッチング、レイヤ２集約、セッション境界制御、マルチキャスト、及び／又は加入者管理）に対するサポートを提供する、及び／又は、複数のアプリケーション・サービス（例えばデータ収集）に対するサポートを提供する、「複数サービス・ネットワーク要素」である。

【 0 0 1 6 】

分割アーキテクチャ領域

単一のアクセス／集約ネットワークは、複数の従来型ルータと協働する複数の個別の分割アーキテクチャ領域で構成され得る。本明細書で使用される場合、分割アーキテクチャ領域とは、ドメインに類似した、別々のルーティングを伴う分割アーキテクチャ・ネットワークのセクションである。これは、ネットワークの堅固さ又は制御プレーンのスケラビリティに関する、広範な地理的領域にわたる管理を簡略化するために実行することができる。各分割アーキテクチャ領域は、別々のコントローラによって管理可能である。特定のアプリケーションに応じて、これら別個の分割アーキテクチャ領域のコントローラは、分割アーキテクチャ・ネットワークの適切な管理のために、何らかの情報を共有及び交換する必要がある。

【 0 0 1 7 】

オープン・ショーテスト・パス・ファースト・ルーティング

従来型ネットワーク及び分割アーキテクチャ・ネットワークは、どちらも、スイッチとネットワークによってサービスが提供される他のデバイスとの間の経路を計算しなければならない。オープン・ショーテスト・パス・ファースト（ＯＳＰＦ）は、内部ゲートウェイ・ルーティング・プロトコルである。ＯＳＰＦ（ＲＦＣ 2328に定義）は、ルータがその近隣のリンク状態情報をルーティング・ドメイン内のすべてのノードにブロードキャストする、リンク状態プロトコルである。この情報を使用して、あらゆるルータがドメイン内にネットワーク全体のトポロジ・マップを構築する。各ルータは、全ネットワーク・トポロジを反映するリンク状態データベースを維持する。このトポロジ・マップ及びリンク・コスト・メトリクスに基づき、ルータは、ダイクストラ法アルゴリズムを使用してすべての他のルータへの最短パスを決定する。次にこの情報を使用し、インターネット・プロトコル（ＩＰ）パケットの転送に使用されるルーティング・テーブルが作成される。

【 0 0 1 8 】

ＯＳＰＦは、ルーティング・ドメインを、管理しやすいように領域境界ルータ（ＡＢＲ）によって分離されている複数の「領域」に分割することができる。各領域は、ＯＳＰＦネットワークのコア又はバックボーン用に予約された識別子「０」を備える、３２ビット識別子、通常は領域内のメイン・ルータのＩＰアドレスによって、識別される。リンク状態情報のブロードキャストは領域に限定され、領域を越えて送信されることはない。ＯＳＰＦプロセスにおいて、各ルータ（スイッチ）は、それが属する各領域に対するＯＳＰＦプロトコル・プロセスの別々のコピーを実行する。ルータは、異なる領域に属する複数のインターフェースを有する場合、プロセスの複数のコピー、各インターフェースについて１つを実行する。起動時、ルータのＯＳＰＦプロセスは、すべてのルーティング・プロトコル・データ構造を開始し、そのルータのアクティブ・インターフェースに関する情報を下位レイヤから取得する。その後ルータは、ＯＳＰＦのハロー・プロトコル・パケットを使用して、その近隣を検出する。ルータはその近隣にハロー・パケットを送信し、その近

隣からハロー・パケットを受信する。ブロードキャスト及びポイント・ツー・ポイント・ネットワークでは、ハロー・パケットはマルチキャスト・アドレス 2 2 4 . 0 . 0 . 5 に送信される。非ブロードキャスト・ネットワークでは、ユーザ構成は近隣を開発するために必要である。ブロードキャスト及び非ブロードキャスト多重アクセス・ネットワーク (NBMA) では、ネットワークのセグメントに「指定ルータ」及び「バックアップ指定ルータ」を選択するために、ハロー・プロトコルが使用される。

【 0 0 1 9 】

近隣を検出すると、ルータはその新しく検出された近隣を用いて「隣接 (adjacencies)」を確立するように試行することになる。隣接は、領域内のルーティング情報の分布を決定する。ルーティング更新は、隣接上でのみ送受信される。隣接を確立すると、ルータはその「リンク状態データベース」を、隣接の他端上の対応するルータと同期させることになる。ブロードキャスト及びNBMAネットワークの場合、指定ルータはどのルータが隣接になるかを決定する。

10

【 0 0 2 0 】

ルータは、リンク状態アドバタイズ (LSA) を使用して、「リンク状態」とも呼ばれるその状態を定期的にアドバタイズする。リンク状態は、ルータの状態が変化した時にもアドバタイズされる。LSAは、アドバタイズ・ルータの隣接も含む。ルータは、領域全体にそのLSAをフラッディングする (flood)。フラッディング・アルゴリズムは、領域内のすべてのルータが同じ正確なリンク状態データベース確実に有するように、情報の信頼性を保証する。リンク状態データベースは、領域に属する各ルータによって発せられるLSAの集合からなる。

20

【 0 0 2 1 】

各ルータは、このデータベースを使用して、それ自体をルートとして備える最短パス・ツリーを計算する。次に最短パス・ツリーを使用して、ルータのルーティング・テーブルが作成される。OSPFメッセージは、TCP又はUDP等の移送レイヤ・プロトコルを使用せずに、直接伝送され、プロトコル番号 8 9 と共にIPパケットに封入される。加えてOSPFは、それ自体の誤り検出及び訂正を使用する。

【 0 0 2 2 】

OSPFは、以下の5つの異なるパケット・タイプを定義する。ハロー・パケット：OSPFのハロー・プロトコル・パケットは、近隣関係を発見及び維持するために使用される。データベース記述パケット：これらのパケットは、隣接を形成するために使用される。リンク状態要求パケット：これらのパケットは、隣接ルータ間でリンク状態データベースをダウンロードするために使用される。リンク状態更新パケット：これらのパケットは、OSPFの信頼できる更新機構に使用される。単一のリンク状態更新パケットは、いくつかのルータのLSAを含む。リンク状態Ackパケット：これらのパケットは、OSPFの信用できる更新機構に関するリンク状態更新パケットと共に使用される。OSPFプロトコル・パケット (ハロー・パケットを除く) は、隣接を介してのみ送信される。したがってすべてのOSPFプロトコル・パケットは、ソース・アドレスとして1つのルータのIPアドレス、及び、宛先アドレスとして他のルータのIPアドレスを用いて、単一のインターネット・プロトコル (IP) ホップを移動する。

30

40

【 0 0 2 3 】

RFC 2328に定義されたOSPFは、以下の5つの異なるタイプのLSAを指定する。ルータLSA：これらのLSAは、領域内のすべてのルータによって送信される。各LSAは、領域に対するルータのインターフェースの状態を含む。ルータLSAは単一の領域全体のみでフラッディングされる。ネットワークLSA：これらのLSAは、ブロードキャスト又はNBMAネットワークの指定ルータによって発信される。このLSAは、ネットワークに接続されたルータのリストを含む。ルータLSAと同様に、このLSAも単一の領域全対のみでフラッディングされる。ネットワーク要約LSA：このLSAは、領域境界ルータ (ABR) によって発信され、各LSAは、領域外部であるが依然として自律システム (AS) 内部にある宛先ネットワークへの経路を記述する。境界要約LS

50

A：ネットワーク要約LSAと同様に、このLSAは境界領域ルータ（ABR）によって発信され、AS境界ルータへの経路を記述する。AS外部LSA：このLSAは、AS境界ルータによって発信され、AS全体でフラディングされる。これらのLSAのそれぞれが、他のAS内の宛先への経路を記述する。OSPFパケット・ヘッダ・フォーマットは、図2に示される通りである。ハロー・パケット・フォーマットは図3に示されている。LSAヘッダ・フォーマットは図4に示される通りである。ルータLSAパケット・フォーマットは図5に示される通りである。

【0024】

分割アーキテクチャ・ネットワークにおけるOSPF

分割アーキテクチャ・ネットワークにおいて、ルーティング・メッセージは、すべての他の制御メッセージと同様に、コントローラ間で交換される。分割アーキテクチャにおいてOSPFを実装するための単純な方法の1つは、内部コストと、外部コストに関する交換情報のみと、すなわち、異なるネットワーク領域を接続するリンクのコストを、無視することである。このように実行することによって、従来のOSPFメッセージはコントローラ間で交換可能であり、各分割アーキテクチャ領域は単一のノードとみなすことができる。

【0025】

図6に示されたネットワークでは、例えばコントローラBが、10単位のコストを伴う領域Cへのリンクを有することを告知する。しかしながら、この情報に基づいて探知された経路は、準最適である。この実装では、内部コストが他のコントローラに提供されないため、コントローラは実際のエンド・ツー・エンド・コストに基づいて決定することはできない。その結果、1つの領域内の内部コストが内部パスによってかなり異なる場合、準最適な経路が生じる。

【0026】

代替実施形態において、上記の準最適経路の問題は、コントローラ間で交換される情報に内部コストを含めることによって解決できる。各コントローラは、内部リンク・コストを計算し、その境界スイッチのうちの任意の2つの間の最短パス（すなわち最も安価なパス）を探知することができる。例えば、再度図6を参照すると、コントローラCは、それ自体の領域とコントローラFの領域との間のリンク・コストを告知した場合、境界スイッチS1とS2の間の最短パス・コストを追加し、それをコントローラBに送信されるコスト情報に追加する。同様に、コントローラCはスイッチS2とS3の間の最短パス・コストを、領域CとFの間のコストに追加して、この情報をコントローラDに送信する。

【0027】

この代替実施形態は問題も発生させる。この実施形態は、同じリンクについて複数の（及び恐らくは異なる）リンク・コスト・メッセージを送信する必要がある。例えば、領域Fに到達するためにコントローラCによって告知されるコストは、この情報がコントローラBに送信される場合とコントローラDに送信される場合とでは、それぞれが異なる内部パス・コストに基づいているため、異なる（例えば、S1とS2の間の内部最短パス・コストがS2とS3の間のコストと異なる場合、コントローラCはそれ自体の領域とコントローラFの領域との間の2つの異なるリンク・コストを告知することになる）。最終的に、これら2つの矛盾するメッセージが他のコントローラによって受信されることになり、これらのメッセージは、2つの別々のメッセージとして解釈されるのではなく、第2のメッセージは第1の受信されたメッセージの更新として解釈されることになる。

【0028】

例えば、複数のコストを同じリンクに起因させることができるようにOSPFプロトコルを修正することによって、この複数コスト問題を解決しようとする他の実施形態は、分割アーキテクチャ内のコントローラによって探知される最短パスが、最適でない可能性がある。例えば、図7に示されたシナリオでは、各リンクの番号がいずれかの方向のリンクのコストを示し、準最適経路が決定されることになる。領域Aから領域Fへの最短パスを計算する際には、以下のステップが実行される。（1）CはB及びDに、Fに到達するた

めのコスト $2(1 +)$ を告知し、(2) B は A に、C に到達するためのコスト $11(10 + 1)$ を告知し、(3) D は B に、C に到達するためのコスト $21(1 + 20)$ を告知し (領域 D 内の境界スイッチ S4 と S5 の間には 1 つの内部パスのみが存在することに留意されたい)、(4) B は A に、C に到達するためのコスト $2(1 + 1)$ を告知する。

【0029】

コントローラ A が領域 F へのその最短パスを計算する場合、上記情報に基づいて、コストが最も小さい 14 であることから、他のパスよりもパス A - B - C - F を選択する。しかしながら、図 7 に示されるすべてのコスト情報に基づく、A と F の間の最適なルーティングは、A - B - D - E - D - C - F であり、コストは 10 のみである。コントローラ A は、コントローラ間で交換される情報がこうした選択にとって不十分であることから、このパスは選択され得ない。したがって、最適なルーティング・ソリューションの場合、コントローラ間で交換される OSPF メッセージは、以下でより詳細に考察されるように、内部リンク・コストに関する更なる情報を伝搬する必要がある。

【0030】

抽出された (abstracted) 内部領域パス・コストを備える OSPF

本発明の実施形態は、従来技術の欠点を回避するための方法及びシステムを提供する。従来技術及び前述の分割アーキテクチャ・ネットワークにおける OSPF の単純な実施は、最短パスが必ずしも正確に決定されない、及び / 又は、非効率的又は拡張性がない経路を識別するために過剰な情報が提供される、準最適ルーティング・ソリューションを提供する。

【0031】

本発明の実施形態は、従来技術のこれらの欠点を克服する。本発明の実施形態は、境界スイッチの各ペア間の内部領域パス・コストを、コスト値に関連付けられた直接リンクとして抽出する。このソリューションは、最適なパスを提供し、効率的に実行可能であり、拡張性があり、非分割アーキテクチャにおける従来のルータとの後方互換性がある。

【0032】

一実施形態において、OSPF ルーティング・プロトコルは、ネットワーク内の任意の 2 つの転送要素間に最適な (最短の) パスを確立するために、分割アーキテクチャにおけるコントローラ間で実装される。実施形態は、内部領域パス・コスト及び内部領域リンク・コストの両方の必要な情報を、複数領域分割アーキテクチャ・ネットワーク全体に公開する。前述のように、実施形態は、最適性、効率性、拡張性、及び後方互換性を提供する。

【0033】

ルーティング・プロトコルは、各分割アーキテクチャ・コントローラが、任意の他の宛先に到達するために最適なパスを個別に計算できるようにする、十分かつ正確な情報を提供する。これは、従来のドメイン内ルーティング・プロトコル (OSPF、IS - IS) が従来型ネットワーク内で提供する、最適性特性である。しかしながら、最適パスの定義は、分割アーキテクチャ・ネットワークで使用される場合はわずかに異なる。従来型ネットワークでは、これは最小のルータ間コストを伴うパスである。分割アーキテクチャ・ネットワークにおいて、最適パスは、最小ルータ間コストとルータ内コスト (又は SA 領域内コスト) とを伴うパスである。

【0034】

分割アーキテクチャ OSPF (SA - OSPF) ルーティング・プロトコルの実施形態は、各ルータが最短パス決定を独立に計算できるようにするものである。OSPF は高速で収束する。言い換えれば、分割アーキテクチャ OSPF プロトコルは、従来型ネットワーク OSPF と比べて、追加の収束オーバーヘッドを招かない。分割アーキテクチャ OSPF プロトコルは、何百ものスイッチを備える大規模ネットワークに拡張する。拡張性は、交換されるメッセージの数のオーバーヘッドとコントローラに関する記憶要件の両方によって量子化することができる。増分の展開は、ネットワーク・プロトコルに関する任意の新しい提案の適合にとって重要である。分割アーキテクチャ OSPF プロトコルは、従来型

ルータとの後方互換性がある。この特性により、実際のネットワーク環境におけるその潜在的な使用率を増加させる。

【 0 0 3 5 】

分割アーキテクチャ O S P F の実施形態は、コスト値に関連付けられた直接リンクとして、境界スイッチの任意のペア間の内部領域パス・コストを抽出する。この機能を容易にする分割アーキテクチャ O S P F のいくつかの態様が存在する。内部領域コストは、後方互換性を保証するために従来の O S P F メッセージの形で組み込まれる。領域内コストは、領域間コストとは別に、すべての領域におけるすべてのコントローラに伝搬される。これにより、最短パス計算の最適性を保証する。これは、各コントローラが、領域間及び領域内の両方の、全ネットワークの完全なピクチャを有することを保証する。こうした十分な領域がある場合、従来の O S P F と同様に、最短パスの計算はダイクストラ法アルゴリズムを使用して容易に実行可能である。

10

【 0 0 3 6 】

分割アーキテクチャ O S P F は、どの内部情報が外部コントローラに伝搬されることになるかが慎重に管理されるという点で、拡張可能である。1つの生来の手法は、分割アーキテクチャ・スイッチの任意のペア間の内部コストを送信することである。しかしながら、これは拡張性の問題を生じさせる。更にこの情報のほとんどが、外部コントローラでの決定に有用ではない。分割アーキテクチャ O S P F は、境界分割アーキテクチャ・スイッチの任意のペアの合計集約コストのみを伝搬する。

【 0 0 3 7 】

20

図 8 は、分割アーキテクチャ・ネットワークの例示の一実施形態を示す図である。例示の分割アーキテクチャ・ネットワークは、別々の分割アーキテクチャ領域 (S A) 8 0 1 A ~ C に分割される。各領域 8 0 1 A ~ C はスイッチ・セットを含む。同じ領域内のすべてのスイッチは、単一の論理コントローラ 8 0 3 A ~ C によって制御される。一実施形態において、S A は、冗長性の目的で、1 次コントローラ及びバックアップ・コントローラとして実装可能である。

【 0 0 3 8 】

各 S A におけるスイッチは、分割アーキテクチャ・ネットワークのデータ・プレーンを実装可能な、任意のタイプのルータ、スイッチ、又は同様のネットワーキング・デバイスとすることができる。スイッチは、境界分割アーキテクチャ・スイッチ及び内部分割アーキテクチャ・スイッチを含むことができる。境界分割アーキテクチャ・スイッチは、異なる S A 領域内の別のスイッチに接続するインターフェースで、分割アーキテクチャ機能をサポートする。境界分割アーキテクチャ・スイッチは、通常、単一の S A 領域のコントローラによって制御される。他の実施形態において、境界分割アーキテクチャ・スイッチは、複数の S A 内に存在することが可能であり、各それぞれの S A コントローラによって制御されるインターフェースを有する。内部分割アーキテクチャ・スイッチは、分割アーキテクチャ・プロトコルをサポートする。これは、その領域内のコントローラによって制御される。そのすべての近隣は同じ S A 領域内にある。

30

【 0 0 3 9 】

スイッチは、リンク・セットを介して互いに通信している。これらのリンクは、有線又は無線の通信媒体、及びそれらの任意の組み合わせを含む、任意のタイプの通信媒体とすることができる。リンクは、内部リンク又は外部リンクのいずれかとして分類可能である。内部リンクは、S A 内の 2 つのスイッチ間のリンクであり、これらのスイッチは、同じ S A 領域に属する境界スイッチ又は内部 S A 領域スイッチのいずれかとしてすることができる。外部リンクは、異なる S A 領域に属する 2 つの S A スイッチ間のリンクである。この場合、どちらの S A スイッチも境界 S A スイッチである。

40

【 0 0 4 0 】

リンク状態アドバタイズ (L S A) 8 0 5 は、分割アーキテクチャ O S P F を実装する L S A の例示セットである。各 L S A は、境界スイッチのペア及びこれら 2 つの境界スイッチの間の S A 領域内をトラバースする関連付けられたコストの抽出である。これらの L

50

S Aは各コントローラによって生成され、隣接するコントローラに伝送される。例示のL S Aセット8 0 5は、S A領域8 0 1 Bに関するコントローラ8 0 3 BのL S Aのセットである。

【0041】

図9は、コントローラを実装するネットワーク要素の一実施形態を示す図である。一実施形態において、コントローラ9 0 1はルータ、スイッチ、又は同様のネットワーキング・デバイスである。コントローラ9 0 1は、入口モジュール9 0 3、出口モジュール9 0 5、ネットワーク・プロセッサ9 0 7、及び記憶デバイス9 1 1を含むことができる。入口モジュール9 0 3は、物理及びリンク・レベルで着信データ・トラフィックを処理し、このデータを更なる処理のためにネットワーク・プロセッサに提供する。同様に、出口モジュール9 0 5は、物理及びリンク・レベルで発信データ・トラフィックを処理し、接続されたネットワークを介してこれを他のデバイスに伝送する。これら2つのモジュール機能がまとまって、ネットワークを介した他のデバイスとの通信を実行可能にする。

10

【0042】

ネットワーク・プロセッサ9 0 7は、ネットワークのデータ・プレーンを支配するネットワークの制御プレーンに関する機能のそれぞれを含むネットワーク要素の機能を実行する、処理デバイス又は処理デバイスのセットである。ネットワーク・プロセッサ9 0 7は、近隣発見モジュール9 1 3、O S P Fモジュール9 1 5、トポロジ習得モジュール9 1 7、リンク状態管理モジュール9 1 9、及びOpenFlowコントローラ9 2 1等のコントローラ・モジュールを含む、モジュール・セットを実行することができる。

20

【0043】

加えて、ネットワーク・プロセッサ9 0 7は、記憶デバイス9 1 1内に記憶されたデータにアクセスすることができる。記憶デバイス9 1 1に記憶されたデータは、ルーティング・テーブル9 2 3及びリンク状態データベース9 2 5を含むことができる。他の実施形態において、記憶デバイス9 1 1は、任意数の別々のローカル又は分散型の記憶デバイス、及び、これらのデバイス全体にわたって記憶されたデータの任意の配置構成を含むことができる。ネットワーク・プロセスによって実行される他のモジュールは、記憶デバイス9 1 1からロードすること、又は記憶デバイス9 1 1上に記憶することも可能である。

【0044】

近隣発見モジュール9 1 3は、近隣コントローラのそれぞれに関する情報を取得して、コントローラによって管理されるS A領域のスイッチ間での適切な通信及びそれらスイッチの構成を実行可能にするために、ハロー・プロトコル又は同様のプロトコルを使用して、ネットワーク内の他のデバイスと通信するためのプロトコルを管理することができる。任意のハロー・プロトコル又はプロセスを使用して、S A領域に関して隣接するコントローラ及びスイッチを識別することができる。

30

【0045】

トポロジ習得モジュール9 1 7は、内部でコントローラが動作するネットワークのトポロジを決定するために、近隣発見モジュール9 1 3によって収集された情報を使用する。このトポロジ情報は、ネットワークを介して最適なルートを計算するために、O S P Fモジュール9 1 5によって使用される。トポロジ情報は、ネットワーク内のリンク・コストを追跡及び決定するために、リンク状態管理モジュール9 1 9によっても使用される。

40

【0046】

O S P Fモジュール9 1 7は、ネットワーク内のソース又は発信デバイスから宛先デバイスの間の最適な経路を計算する。O S P Fモジュール9 1 7は、ルーティング・テーブル・セット9 2 3内にルーティング情報を記憶することができる。O S P Fモジュール9 1 7は、トポロジ習得モジュール9 1 7によって生成されたトポロジ情報を使用して、ネットワーク要素間の経路を計算することができる。いくつかの実施形態において、O S P Fモジュール9 1 7は、他の各デバイスへの経路を確立するために、ルートにコントローラを備えたネットワークに関するスパニング・ツリーを計算することもできる。

【0047】

50

リンク状態管理モジュール 9 1 9 は、リンク状態データベース 9 2 5 内のネットワーク又はネットワークの領域に関するリンク状態情報を管理する。リンク状態管理モジュール 9 1 9 は、コントローラと他のコントローラの間並びにコントローラの S A 領域内の接続に関する情報を提供する隣接コントローラに広められることになる、リンク状態アドバタイズを生成することもできる。この情報は、ネットワーク内の他のコントローラに送信される各境界スイッチ・ペアに関するリンク状態アドバタイズ・セットとして、パッケージングすることができる。

【 0 0 4 8 】

コントローラ 9 2 1 は、S A 領域の制御プレーンを管理するための、任意のタイプの分割アーキテクチャ・コントローラとすることができる。コントローラは、分割アーキテクチャ・ネットワークを管理するための、OpenFlow プロトコル又は同様のプロトコルを実装することができる。コントローラ・モジュール 9 2 1 は、データ・プレーン上でのパケットの転送を構成するために、S A 領域内のスイッチと通信することができる。コントローラ 9 2 1 は、近隣情報、リンク状態アドバタイズを交換するため、及び同様の情報をピアに提供するために、他のコントローラとも通信する。

【 0 0 4 9 】

図 1 0 は、分割アーキテクチャ・ネットワーク内の O S P F を実装及びサポートするためのコントローラの動作の一実施形態を示すフローチャートである。プロセスは、わかりやすくするために、コントローラによって実行されていると考えて説明されるが、プロセスは、特定のコントローラ・モジュール（例えば OpenFlow コントローラ）と共にコントローラの構成要素（例えば O S P F モジュール）によって実行されることが可能である。プロセスは、コントローラが活動化又はリセットされた時点で開始することができる。他の実施形態において、プロセスは連続的又は定期的である。コントローラは、O S P F ルーティング・プロトコル関連データ構造等のデータを記憶及び操作するために使用されることになる、データ構造セットを開始することができる（ブロック 1 0 0 1 ）。

【 0 0 5 0 】

コントローラは、境界スイッチである S A 領域内のスイッチを識別することを含む、その割り当てられた S A 領域のトポロジを習得する（ブロック 1 0 0 3 ）。トポロジは、リンク状態アドバタイズ又はリンク状態データベースの他のコントローラ及び従来型ルータとの交換を通じて、並びに同様の機構を通じて、習得することができる。異なる S A 領域内のコントローラが S A - O S P F プロトコルと通信可能である。交換されるメッセージは、ハロー・メッセージ、データベース記述、リンク状態要求、リンク状態更新、及びリンク状態肯定応答を含む、従来の O S P F メッセージと同様である。特に図 8 では、リンク状態更新メッセージ 8 0 5 が示されている。単一のリンク状態更新パケットは、いくつかのリンク状態アドバタイズ（L S A ）を含むことができる。S A - O S P F の基本ルーティング・アルゴリズムの別々のコピーは、各領域内で実行する。複数の領域へのインターフェースを有するルータは、アルゴリズムの複数のコピーを実行する。このデータは、S A 領域及び S A 領域が位置する境界ネットワークのトポロジカル・マップにコンパイルされる。

【 0 0 5 1 】

コントローラ（例えばコントローラの O S P F モジュール）は、コントローラの S A 領域内の境界スイッチの各ペア間の最短パスを計算する（ブロック 1 0 0 5 ）。これは、スイッチ間の他の領域内経路を計算することと共に実行可能である。境界スイッチ・ペア・コストは、他のコントローラ及び従来型ルータと共有され、それらのコントローラ及び従来型ルータがコントローラの S A 領域を横断する最適なパスを決定できるようにするものである。最短パスは、S A 領域の習得トポロジを使用して計算される。一実施形態において、ダイクストラ法アルゴリズムは、任意の境界スイッチ・ペア間の最短パスを計算するために使用される。各ペアについて計算された最短パスのコストは、その後、トポロジカル・データと共に、又は O S P F ルーティング・データと共に、記憶される（ブロック 1 0 0 7 ）。

【 0 0 5 2 】

次にコントローラは、近隣コントローラを識別する（ブロック 1 0 0 9）。一実施形態において、コントローラは、近隣を獲得するために S A - O S P F のハロー・プロトコルを使用することができる。コントローラは、その近隣の領域又は近隣の従来型ルータ内の他のコントローラにハロー・パケットを送信し、次に他のコントローラ及び従来型ルータからハロー・パケットを受信する（ブロック 1 1 1 1）。ブロードキャスト及びポイント・ツー・ポイント・ネットワークでは、コントローラは、そのハロー・パケットをマルチキャスト・アドレスに送信することによって、その近隣領域内のコントローラ又は従来型ルータを動的に検出する。非ブロードキャスト・ネットワークでは、近隣コントローラを発見するために、いくつかの構成情報を提供することができる。

10

【 0 0 5 3 】

次にコントローラは、その新しく獲得された近隣コントローラのいくつかと共に、隣接の形成を試行することになる。リンク状態データベースは、隣接コントローラと L S A メッセージを交換すること（ブロック 1 1 1 5）によって、隣接コントローラのペア間で同期される（ブロック 1 1 1 3）。各 L S A メッセージは、その外部リンク全体のコストを含む。コントローラは、近隣に接続しているリンクのコストをこの特定の近隣に送信するのみならず、リンクのコストをすべての近隣にも送信する。外部リンクのコストの送信に加えて、コントローラは内部リンクのコストも送信することになる。L S A メッセージは、2つの境界スイッチの S A スイッチ I D を、2つの関連インターフェース I D と共に使用して構築することができる。コントローラは、リンク状態とも呼ばれるその S A 領域の状態を、定期的にアドバタイズすることができる。リンク状態は、領域の（コントローラの）状態が変化した場合にもアドバタイズされる。領域の隣接は、その L S A の内容に反映される。この隣接とリンク状態との間の関係により、プロトコルが適時に障害を検出できるようになる。

20

【 0 0 5 4 】

L S A は、複数の S A 領域にまたがってフラッディングされる。フラッディング・プロセスは信頼できるものであり、ネットワーク全体のすべてのコントローラが確実に同じリンク状態データベースを有することを保証する。このデータベースは、領域間トポロジ及び領域内トポロジの両方を含む各領域によって発信された L S A の集合からなる。

【 0 0 5 5 】

各コントローラは、このリンク状態データベースから、それ自体をルートとして備えるネットワーク全体について最短パス・ツリーを計算する（ブロック 1 1 1 7）。次にこの最短パス・ツリーは、O S P F プロトコルに関するルーティング・テーブルを生み出す。コントローラはこのルーティング情報に基づき、計算された最短パスを使用して、S A 領域を介してデータ・パケットをルーティングするように、その S A 領域内のスイッチをプログラミングする（ブロック 1 1 1 9）。スイッチの更新は、セキュア・チャネルを介する制御プレーン・プロトコル（例えば OpenFlow）を使用するか、又は同様のプロセスを使用して、達成可能である。

30

【 0 0 5 6 】

ネットワークに変化が生じた場合、コントローラによって更新リンク状態アドバタイズを送信することができる（ブロック 1 1 2 1）。L S A メッセージは、隣接するコントローラがそれらのルーティングを必要に応じて更新できるように、コントローラの S A 領域内の変化を通知するために隣接のコントローラに送信される。

40

【 0 0 5 7 】

分割アーキテクチャ O S P F は、従来型ルータがネットワーク内の他のルータには単一のルータのように見えるという点で、従来の O S P F とは異なる。しかしながら S A 領域は、ネットワーク内の他の領域／ルータには N 個のルータのように見えることになり、ここで N は、領域の境界分割アーキテクチャ・スイッチの数である。相違は、従来の O S P F は領域あたり単一のルータ上で単一のプロセスを実行し、単一のルータのように見える一方で、S A - O S P F は S A 領域の単一のコントローラ上で単一のプロセスを実行し、

50

他のコントローラにはN個のルータのように見えるという点である。

【0058】

一実施形態において、単一のSA-OSPFプロセスは、従来型ルータにおいてOSPFプロセスによって送信されるあらゆるハロー・パケットについて、N個の境界SAスイッチのそれぞれを表すために、N個の異なるハロー・パケットを送信する。これらのハロー・パケットは、他のルータ/SA領域に接続されたインターフェース上で送信される。同様に、単一のSA-OSPFプロセスは、従来型ルータにおいてOSPFプロセスによって送信されるあらゆるルータLSAパケットについて、N個の境界SAスイッチのそれぞれを表すために、N個の異なるルータLSAパケットを送信する。これらのLSAパケットは、他のルータ/SA領域に接続されたインターフェース上で送信される。

10

【0059】

端的に言えば、SA領域に関するコントローラは、N個のOSPFコントローラ・プロセスを有するかのようにOSPFを管理し、ここでNは、SA領域内の境界分割アーキテクチャ・スイッチの数である。これにより、接続された従来型ルータ上で実行中のOSPFとの後方互換性を保証する。したがってコントローラは、N個のハロー・パケット、N個のLSAメッセージ等を送信する。従来型ネットワークにおいて同じ数のルータの場合、SA-OSPFは、分割アーキテクチャ・ネットワーク内で従来のOSPFに比べて少ない数のプロトコル・パケットを使用する。同じSA領域に属する境界SAスイッチ間でOSPFプロトコル・パケットを送信する必要がないことから、節約が生じる。

20

【0060】

LSAメッセージ

図11は、OSPFリンク状態アドバタイズ・ヘッダ・フォーマットを示す図である。SA-OSPFに後方互換性があることを保証するために、標準OSPFメッセージ・フォーマットと同じフォーマット（すなわち、同じフィールド・シーケンス及びフィールド長さ）が使用される。しかしながら、いくつかのフィールドに対する値の割り当ては異なる可能性がある。

【0061】

図12は、SA-OSPF LSAメッセージ、具体的に言えばルータLSAタイプ・メッセージを示す図である。ルータLSAは領域内の各スイッチによって生成される。これは、領域内の領域インターフェースの状態を記述する。このメッセージはすべてのコントローラによって発信される。このLSAは、複数ドメインSAネットワークへのルータのインターフェースの集合状態を記述する。これは、すべてのコントローラ全体にフラッディングされる。

30

【0062】

一実施形態において、ルータLSAメッセージは、自己スイッチID、自己インターフェースID、近隣インターフェースID、及び近隣スイッチIDの、メッセージの4つのフィールドにおいて、従来型ルータLSAメッセージと異なる（図12に図示）。スイッチIDは、分割アーキテクチャ内の各スイッチに割り当てられる固有の32ビット識別子である。このIDは、複数領域SAドメイン内のすべてのスイッチの中で固有でなければならない。一実施形態において、スイッチの最高位IPアドレスがスイッチIDに使用される。自己スイッチIDはアドバタイズ・スイッチの識別子であり、近隣スイッチIDは近隣スイッチの識別子である。インターフェースIDは、スイッチの各インターフェースに割り当てられる32ビット識別子である。この識別子は単一のスイッチにおいてのみ固有であればよい。自己インターフェースIDはアドバタイズ・スイッチのインターフェースIDであり、近隣スイッチIDは（近隣スイッチの）近隣インターフェースのインターフェースIDである。

40

【0063】

ハロー・メッセージの処理

OSPFハロー・プロトコルは、ルータが近隣ルータを伴う隣接を確立及び維持できるようにする機構である。隣接ルータはハロー・メッセージを交換する。ブロードキャスト

50

及びポイント・ツー・ポイント・ネットワークにおいて、ルータは、ハロー・メッセージをマルチキャスト・アドレスに送信することによって、その近隣ルータを動的に検出することができる。ルータはハロー・パケットを受信すると、近隣内に隣接を形成する。リンク状態データベースは、隣接ルータのペア並びにルーティング更新間で同期される。隣接が確立された後、依然としてルータは、活動中であることを示すためにハロー・メッセージを定期的に交換する必要がある。要約すると、ハロー・メッセージの目的は、近隣関係を発見及び維持することである。

【 0 0 6 4 】

図 1 3 は、S A - O S P F におけるハロー・メッセージを示す図である。S A - O S P F プロトコルにおいて、ハロー・メッセージは未変更とすることができる。しかしながら、分割アーキテクチャ・コントローラは、近隣ルータとそれに接続する分割アーキテクチャ・スイッチとの間で隣接を維持するために、追加のハロー・メッセージを構築する必要がある。

10

【 0 0 6 5 】

図 1 4 は、分割アーキテクチャ・コントローラと従来型ルータとの間で交換されるハロー・メッセージを示す図である。S A - O S P F において、コントローラは、境界スイッチの代わりに第 1 に O S P F ハロー・メッセージを構築する。これは、スイッチ I D のフィールドで境界スイッチの最高位 I P アドレスを使用する。パケットは第 1 にスイッチ S 1 に投入され、その後、出口リンク上の隣接ルータに送信される。

【 0 0 6 6 】

20

このように、分割アーキテクチャ・ネットワークにおいて O S P F を実装するための方法、システム、及び装置を説明した。上記の説明は例示的であり、制約的ではないことが意図されることを理解されよう。当業者であれば、上記の説明を読み、これを理解することで、多くの他の実施形態が明らかとなろう。したがって本発明の範囲は、添付の特許請求の範囲、並びにこうした特許請求の範囲が権利を認められる等価物の全範囲を参照しながら、決定されるものとする。

【図 6】

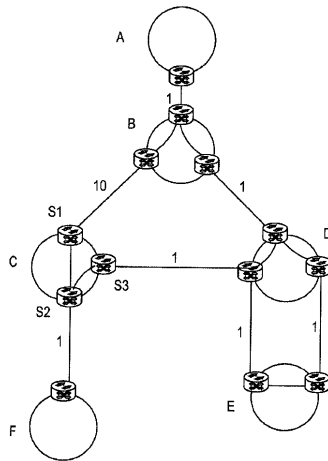


FIG. 6

【図 7】

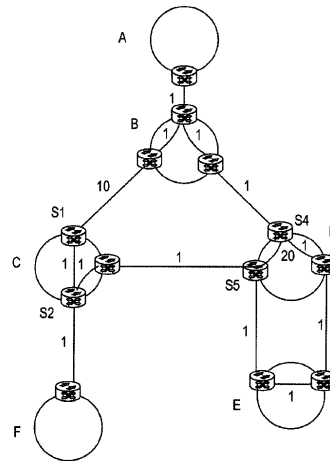
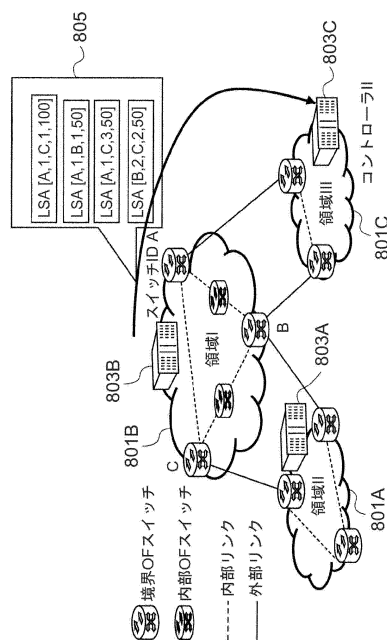
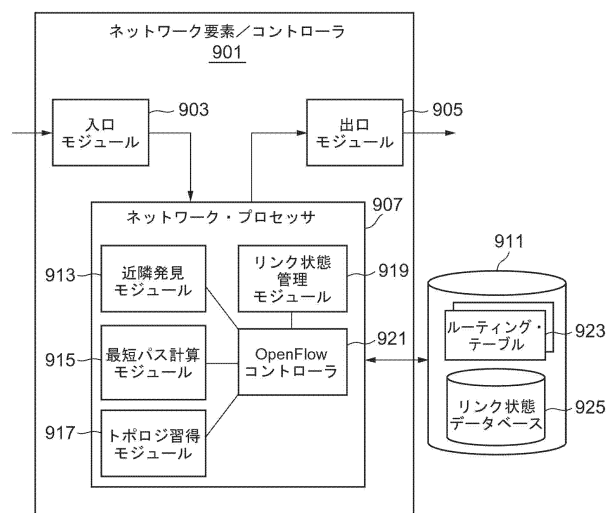


FIG. 7

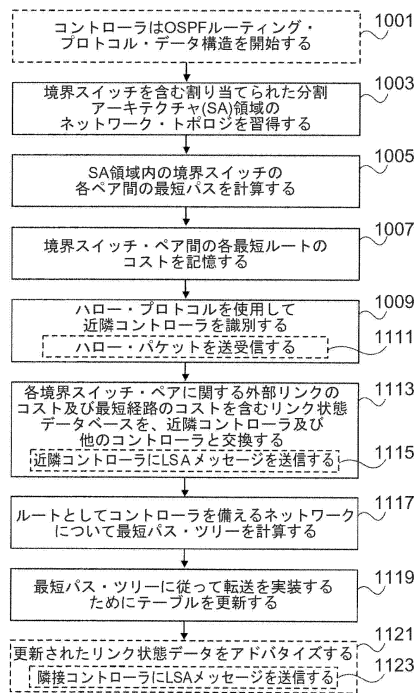
【図 8】



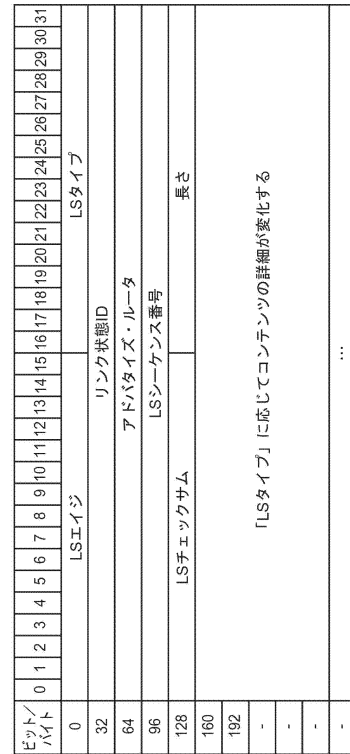
【図 9】



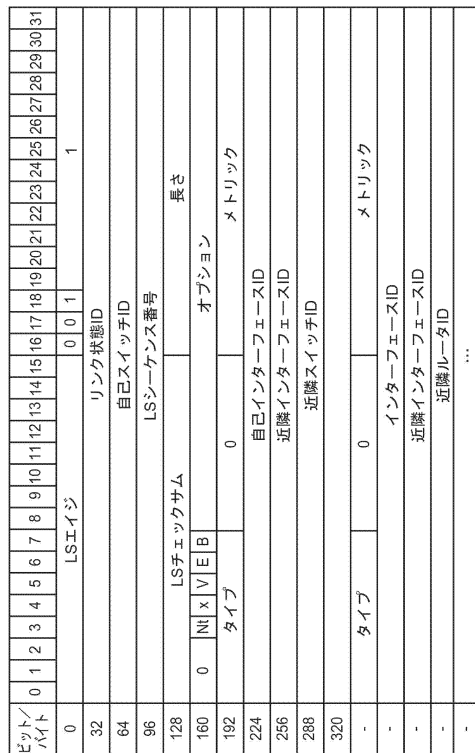
【 図 1 0 】



【 図 1 1 】



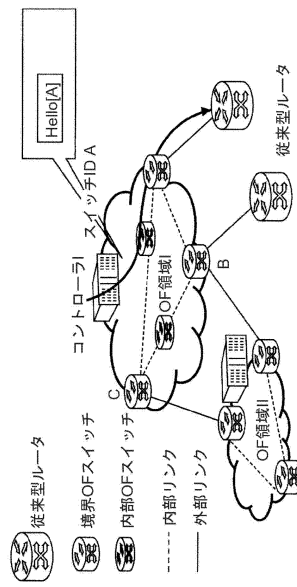
【 図 1 2 】



【 図 1 3 】



【図 14】



フロントページの続き

(72)発明者 ベヘシュチ - サバレ, ネダ
アメリカ合衆国, カリフォルニア州 95134, サン ノゼ, パルミラ ドライブ 3500,
ユニット 1026

(72)発明者 チャン, イン
アメリカ合衆国, カリフォルニア州 94555, フリーモント, ミッデイ コモン 5474

審査官 安藤 一道

(56)参考文献 特表2010-541426(JP, A)
米国特許出願公開第2009/0138577(US, A1)
欧州特許出願公開第02355423(EP, A1)
特開2010-199882(JP, A)
特表2011-524728(JP, A)
国際公開第2010/096552(WO, A1)

(58)調査した分野(Int.Cl., DB名)

H04L 12/715

H04L 12/717

H04L 12/725