

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2015-138987

(P2015-138987A)

(43) 公開日 平成27年7月30日(2015.7.30)

(51) Int.Cl.	F I	テーマコード (参考)
HO4L 12/717 (2013.01)	HO4L 12/717	5K030
HO4L 12/775 (2013.01)	HO4L 12/775	
HO4L 12/703 (2013.01)	HO4L 12/703	

審査請求 未請求 請求項の数 8 O L (全 18 頁)

(21) 出願番号	特願2014-7773 (P2014-7773)	(71) 出願人	000004237
(22) 出願日	平成26年1月20日 (2014.1.20)		日本電気株式会社
			東京都港区芝五丁目7番1号
		(74) 代理人	100103090
			弁理士 岩壁 冬樹
		(74) 代理人	100124501
			弁理士 塩川 誠人
		(72) 発明者	小川 英輝
			東京都港区芝五丁目7番1号 日本電気株式会社内
		Fターム(参考)	5K030 GA12 HD03 KA01 LB07 LE03 LE10 MA12 MB01 MD02

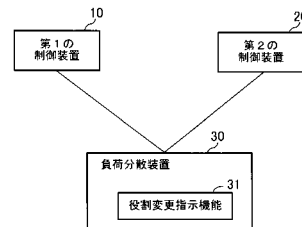
(54) 【発明の名称】 通信システムおよび通信システムにおけるサービス復旧方法

(57) 【要約】

【課題】 ビザンチン型故障の発生を検出した場合に、早期に故障状態を解消しパケット転送を再開可能とする通信システムおよび通信システムにおけるサービス復旧方法を提供する。

【解決手段】 フロー制御機能と仮想ネットワークの設定情報とをもとに複数のパケット転送装置を制御する第1の制御装置10と、動作実績があるフロー制御機能と、運用実績がある仮想ネットワークの設定情報とを有し、第1の制御装置10に故障が発生した場合に切り替え先となる第2の制御装置20と、複数のパケット転送装置と各制御装置との間に配置された負荷分散装置30とを備え、負荷分散装置30は、第1の制御装置10においてビザンチン型故障が発生した場合に、複数のパケット転送装置の制御を第2の制御装置20に実行させる役割変更指示機能31を有する。

【選択図】 図8



【特許請求の範囲】**【請求項 1】**

フロー制御機能と仮想ネットワークの設定情報とをもとに複数のパケット転送装置を制御する第 1 の制御装置と、

動作実績があるフロー制御機能と、運用実績がある仮想ネットワークの設定情報とを有し、前記第 1 の制御装置に故障が発生した場合に切り替え先となる第 2 の制御装置と、

前記複数のパケット転送装置と各制御装置との間に配置された負荷分散装置とを備え、

前記負荷分散装置は、前記第 1 の制御装置においてビザンチン型故障が発生した場合に、前記複数のパケット転送装置の制御を前記第 2 の制御装置に実行させる役割変更指示機能を有する

ことを特徴とする通信システム。

【請求項 2】

負荷分散装置の役割変更指示機能は、第 1 の制御装置の CPU 使用率が所定の閾値を超えたときに、当該第 1 の制御装置においてビザンチン型故障が発生したと判断する

請求項 1 に記載の通信システム。

【請求項 3】

通常運用時、第 1 の制御装置は、複数のパケット転送装置から送信される全てのメッセージを受信するマスタとして動作し、第 2 の制御装置は、前記複数のパケット転送装置からトポロジ検出に必要なメッセージのみを受信するスレーブとして動作し、

負荷分散装置の役割変更指示機能は、第 1 の制御装置においてビザンチン型故障が発生した場合に、前記第 1 の制御装置をスレーブとして動作させ、第 2 の制御装置をマスタとして動作させる

請求項 1 または請求項 2 に記載の通信システム。

【請求項 4】

各制御装置は、負荷分散装置からの役割変更指示に基づいて、自装置をマスタとして動作させるか、スレーブとして動作させるかを管理する役割管理機能を有する

請求項 3 に記載の通信システム。

【請求項 5】

第 2 の制御装置の役割管理機能は、スレーブ動作中に、トポロジ検出に必要なメッセージを受信すると、当該メッセージをもとにマスタとして動作を開始する際に必要なフローを作成する

請求項 4 に記載の通信システム。

【請求項 6】

各制御装置は、パケット転送装置から受信したメッセージがトポロジ検出に必要なメッセージであると判断した場合に、他の制御装置に当該メッセージを通知する

請求項 1 に記載の通信システム。

【請求項 7】

第 1 の制御装置が、複数の制御装置によりクラスタ構成された

請求項 1 から請求項 6 のうちのいずれか 1 項に記載の通信システム。

【請求項 8】

フロー制御機能と仮想ネットワークの設定情報とをもとに複数のパケット転送装置を制御する第 1 の制御装置と、動作実績があるフロー制御機能と、運用実績がある仮想ネットワークの設定情報とを有し、前記第 1 の制御装置に故障が発生した場合に切り替え先となる第 2 の制御装置と、前記複数のパケット転送装置と各制御装置との間に配置された負荷分散装置とを備えた通信システムにおいて、

前記負荷分散装置が、前記第 1 の制御装置においてビザンチン型故障が発生した場合に、前記複数のパケット転送装置の制御を前記第 2 の制御装置に実行させる

ことを特徴とするサービス復旧方法。

【発明の詳細な説明】**【技術分野】**

10

20

30

40

50

【0001】

本発明は、複数のパケット転送装置を制御する制御装置を含む通信システムおよび通信システムにおけるサービス復旧方法に関する。

【背景技術】

【0002】

オープンフローなどのSDN (Software-Defined Networking) 分野では、ネットワーク (NW) におけるパケット転送を単一のコントローラにより制御する。このような分野においては、これまでもコントローラの可用性を高めるための技術が提案されてきた。例えば、特許文献1には、オープンフローネットワークにおいて、複数配置されたコントローラの負荷を均等化する技術が記載されている。

10

【0003】

オープンフロー (OpenFlow) では、データプレーンとコントロールプレーンが分離され、OFC (OpenFlow Controller) がネットワークを集中的に管理する。このような構成において可用性をより高めるためには、OFCに故障が発生した場合に、速やかに業務を復旧するための手段を備えている必要がある。しかし、コンピュータシステムで発生する故障、特にビザンチン型故障からの復旧に関する技術についての提案は現時点では少ない。

【0004】

一般的に、コンピュータシステムで発生する故障について、「沈黙型故障」と「ビザンチン型故障」の二つに分類することができる。

20

【0005】

沈黙型故障は、他ノードからの要求に対し応答を返却できない状態に陥った場合を指す。一方、ビザンチン型故障は、入力に対する出力が正しかったり誤っていたりして、入力に対する出力が信用できない状態に陥った場合を指す。

【0006】

OFCにおける沈黙型故障とは、OFS (OpenFlow Switch) や管理クライアントなどからの要求に対し、応答を返却できない状態であると考えられる。OFCにおける沈黙型故障は、致命的な障害が発生した場合や、過剰な負荷が掛かった場合などに発生する。この沈黙型故障からの復旧手段については、既に考案・実用化されている。例えば、クラスタソフトウェア (クラスタSW) によりOFCを冗長化しておき、沈黙型故障発生時にはスタンバイノードにフェイルオーバーして処理を継続する方法である。

30

【0007】

一方で、OFCにおけるビザンチン型故障とは、なんらかの要因によりOFSに意図しないフローを設定し、意図しないパケット転送処理が行われるようになった状態であると考えられる。

【0008】

ビザンチン型故障対策の詳細に移る前に、先ず背景となる技術として、(I) シングルOFC構成における通常運用処理、(II) シングルOFC構成における再起動処理、(III) クラスタOFC構成におけるフェイルオーバー処理、について説明する。

【0009】

(I) シングルOFC構成における通常運用処理

40

【0010】

シングルOFC構成、すなわちOFCが冗長化されていない場合のシステム構成について説明する。図9は、オープンフローシステムにおけるシングルOFC構成の一例を示すブロック図である。OFC100は、複数のOFS101~103を管理する。OFS101~103は、それぞれ固有のフローテーブル113~115を保持する。各フローテーブルには、有限個のフローエントリが格納される。フローエントリには、パケット転送ルールが記述される。OFC100は、コンフィグ110、トポロジ111、および、フローテーブル112を保持する。

【0011】

50

OFC100は、配下のOFSの数だけフローテーブルを保持する。また、OFC100が保持するフローテーブル112とOFS101～103が保持するフローテーブル113～115の内容は、常時同期される。また、OFC100のフローテーブル112は、OFCの再起動等の処理において揮発する。

【0012】

コンフィグ110は、コンフィグを示す情報(コンフィグ情報)であって、仮想ネットワークの設定が記述されている。オープンフローネットワーク管理者は、適宜、OFCを操作することで、仮想ネットワークの設定を変更することができる。つまり、コンフィグは、OFCの運用中に変更される情報である。また、通常、コンフィグはセーブすることができる。従って、コンフィグは、OFCの再起動等の処理において揮発しない。

10

【0013】

トポロジ111は、トポロジを示す情報(トポロジ情報)であって、OFCがトポロジ検出機能を用いて作成する情報である。トポロジ111は、スイッチ増設や冗長化など、ネットワーク構成を変更した場合に更新される。また、トポロジ111は、OFCの再起動等の処理において揮発する。

【0014】

ここで、OFC100の基本的な動作を説明する。図10は、OFC100の構成の一例を示すブロック図である。

【0015】

OFC100は、通信手段201と、OF(OpenFlow)プロトコル処理手段202と、トポロジ検出手段203と、コンフィグ管理手段205と、経路管理手段207とを含む。

20

【0016】

OFプロトコル処理手段202は、通信手段201を用いて、複数のOFSとの間でオープンフロープロトコルに則った通信を行う。

【0017】

トポロジ検出手段203は、OFプロトコル処理手段202を用いて、ネットワークのトポロジを示す情報(図10に示すトポロジ204)を作成する。トポロジ検出手段203は、LLDP(Link Layer Discovery Protocol)パケットを配下の全てのOFS宛に送信する。OFSは、受け取ったLLDPパケットを全ポートから送出する。OFSのポートにケーブルが接続され、当該ケーブルが他のOFSに接続されている場合、接続先OFSはLLDPパケットのPACKET_INメッセージをOFC100に対して通知する。OFC100は、このPACKET_INメッセージを受け取り、OFSポート間の接続情報を得る。OFC100は、全ての接続情報をまとめることで、ネットワークのトポロジを取得できる。トポロジ検出手段203は、定期的にLLDPを送出することで、OFSの追加やケーブルの抜き差しによる、トポロジの変化を検出できる。そのたびに、トポロジ検出手段203は、トポロジを更新する。

30

【0018】

コンフィグ管理手段205は、ユーザによる仮想ネットワークの設定変更操作に応じて、コンフィグを示す情報(図10に示すコンフィグ206)を更新する。

40

【0019】

経路管理手段207は、経路計算手段209と、経路設定手段210と、オーディット手段211とを含む。

【0020】

オープンフローによりフロー制御を行う場合、始めに少なくともブロードキャストやマルチキャストされたパケットを適切に転送するためのフローをOFSに設定しておく必要がある。このフロー(以下、BC/MC(Broadcast/Multicast)フローと呼称する)は、トポロジやコンフィグの変更に合わせて、更新される。経路計算手段209は、システム開始時や、トポロジまたはコンフィグの更新時にBC/MCフローを計算し、フローテーブル(図10に示すフローテーブル208)に設定する。

50

【 0 0 2 1 】

フローテーブルにフローエントリが追加された場合、経路設定手段 2 1 0 は、O F S に対し、フローエントリの追加を通知する。

【 0 0 2 2 】

O F S にパケットが到着したとき、当該パケットのヘッダに記された宛先アドレスや送信元アドレスなどの情報（ヘッダ情報）が適合条件にマッチするフローエントリが存在する場合、O F S は、そのフローエントリに記述された通りにパケットを処理する。当該処理をハード転送と呼ぶ。フローエントリには、パケット転送ルールとして、上記適合条件や、該当パケットをあるポートから送出する、などの処理内容が記述されている。

【 0 0 2 3 】

一方で、O F S にパケットが到着したとき、当該パケットのヘッダ情報にマッチするフローエントリが存在しない場合、O F S は O F C 1 0 0 に対して、当該パケットの処理方法を問い合わせる。O F S からの問い合わせメッセージを受け取った経路管理手段 2 0 7 は、トポロジ、コンフィグ、および、当該パケットのヘッダ情報を、経路計算手段 2 0 9 に入力し、当該パケットの処理方法を定義するフローエントリを導出する。経路管理手段 2 0 7 がこのフローエントリをフローテーブル 2 0 8 に追加したのち、経路設定手段 2 1 0 は、O F S のフローテーブルに当該フローエントリを追加する。O F S は、追加されたフローエントリに従って、転送や廃棄等の処理を行う。当該処理をソフト転送と呼ぶ。ソフト転送は、フローエントリ導出処理を含む分、ハード転送に比べて性能が悪い。

【 0 0 2 4 】

また、O F S のフローテーブルについて、一定時間以上ヒットしなかったフローエントリを削除する機能がある。この機能では、O F S が、既定時間ヒットしないフローエントリを識別し削除する。その後、O F S は、フローエントリを削除した旨を O F C 1 0 0 に通知する。当該通知を受けた経路管理手段 2 0 7 は、フローテーブルから当該フローエントリを削除する。

【 0 0 2 5 】

このようにして、O F C のフローテーブルと O F S のフローテーブルは同期を保ったまま更新される。しかし、O F C と O F S の間の接続が切れた場合などは、テーブルの同期が失われる。従って、O F C と O F S とが再度接続した際には、オーディット手段 2 1 1 によって、再度テーブルを同期させる処理が実行される。この処理をオーディット（A u d i t）処理と呼ぶ。オーディット手段 2 1 1 は、O F C 側で持つフローテーブルを主とし、従である O F S 側フローテーブルに生じた差分を修正する。

【 0 0 2 6 】

以上のように、オープンフローシステムがフロー制御を行うためには、O F C がコンフィグ、トポロジ、フローテーブルを保持し、且つ、O F S が O F C と同期されたフローテーブルを保持している必要がある。

【 0 0 2 7 】

（ I I ） シングル O F C 構成における再起動処理

【 0 0 2 8 】

（ I ） のシングル O F C 構成において、O F C を再起動する場合の処理について説明する。図 1 1 は、シングル O F C 構成において、沈黙型故障が発生した際の O F C 再起動によるサービス停止時間、つまり O F C が再起動してからパケット転送を再開するまでの時間を示す説明図である。

【 0 0 2 9 】

O F C の O S を再起動すると、O F C と O F S の間の通信が切断される。切断を検出した O F S は、O F C への接続をリトライし続ける。O F C が起動すると、O F C が保持するコンフィグ、トポロジ、フローテーブルは、何れも空の状態になる。従って、O F C では、まず、セーブされているコンフィグのリロード（コンフィグリロード処理 3 0 0）、トポロジの作成（トポロジ情報作成処理 3 0 1）を行う必要がある。これらの処理が完了すると、O F C は、B C / M C フローエントリの導出（B C / M C フローエントリ導出処

10

20

30

40

50

理 302) が可能になる。OFC は、この処理が完了したのち、OFS からの接続要求を受け付ける。OFC と OFS の通信が切断されている間に、OFC と OFS のそれぞれのフローテーブルに差異が発生するため、OFC は、オーディット (オーディット処理 303) を行う。このオーディット処理により、再起動前に設定されていたフローエントリは削除され、BC/MC フローのみが登録された状態となる。

【0030】

以上のように、シングル OFC 構成において、OFC を再起動する場合、上記の処理 300 ~ 303 が完了した後にオープンフローによるフロー制御が可能となる。なお、サービス再開後は、ソフト転送となる。

【0031】

(III) クラスタ OFC 構成におけるフェイルオーバー処理

【0032】

OFC では、沈黙型故障に対応するため、既にクラスタ方式が実用化されている。ここで、クラスタ OFC 構成、すなわち OFC がクラスタ構成 (冗長化) されたシステム構成について説明する。図 12 は、オープンフローシステムにおけるクラスタ OFC 構成の一例を示すブロック図である。図 12 を参照し、クラスタ OFC 構成でのフェイルオーバー処理について説明する。

【0033】

OFC (ACT) 400、OFC (SBY) 401 の 2 台の OFC ノードは、クラスタミドルウェアによりクラスタシステムを構成する。なお、“ACT (ACTIVE)” は、現用系であることを表す。“SBY (STANDBY)” は、待機系であることを表す。OFC (ACT) 400 および OFC (SBY) 401 には、それぞれ固有のアドレスが割り当てられる。また、クラスタシステムには、固有の仮想アドレスが割り当てられる。

【0034】

始めに、クラスタミドルウェアは、OFC (ACT) 400 のアドレスに仮想アドレスを対応付ける。OFS 403 ~ 405 は仮想アドレスに接続することで、結果的に OFC (ACT) 400 と接続される。こうして、OFC (ACT) 400 と OFS 403 ~ 405 によりオープンフローによるフロー制御が行われる。OFC (ACT) 400 は、コンフィグ 410、トポロジ 411、フローテーブル 412 を保持する。これらの情報は、上記「(I) シングル OFC 構成における通常運用処理」と同様に更新される。また、OFC (SBY) 401 も、コンフィグ 413、トポロジ 414、フローテーブル 415 を保持する。OFC (SBY) 401 が保持するこれらの情報は、クラスタミドルウェアが備えるミラーディスク機能やメモリ同期機能により、OFC (ACT) 400 が保持する、対応する情報と常に同期される。

【0035】

OFC (ACT) 400 と OFC (SBY) 401 は、LAN ケーブルにより接続されていて、定期的に通信 (ハートビート通信) することで、対向ノードが正常に稼働していることを確認する。

【0036】

沈黙型故障が発生した場合、OFC (ACT) 400 は、ハートビート通信を継続できず、タイムアウトなどにより通信が切断する。ハートビート通信の切断を検出したクラスタミドルウェアは、仮想アドレスと OFS との間で確立していた通信を切断する。OFS は、通信切断を検出すると、再度 OFC と接続するべく、仮想アドレスへの通信確立を繰り返しリトライする。

【0037】

一方で、クラスタミドルウェアは、OFC (ACT) 400 側でサービスを停止、OFC (SBY) 401 側でサービスを開始させ、仮想アドレスを OFC (SBY) 401 のアドレスに対応付ける。これにより、OFS は、OFC (SBY) 401 と通信を確立できるようになる。OFC (SBY) 401 のコンフィグ 413、トポロジ 414、および

10

20

30

40

50

、フローテーブル415は、OFC(ACT)400側のものと同期されているため最新の情報になっている。従って、OFC(SBY)401において、(I)のケースで必要だった、コンフィグロードやトポロジ検出、BC/MCフローエントリの導出、といった処理は不要である。つまり、OFC(SBY)401では、OFSとOFCの間の通信が切断されていた間に生じた差分を解消するためのオーディット処理のみが必要となり、オーディット処理完了後にはパケット転送可能となる。なお、サービス再開時からハードウェア転送になる。

【0038】

図13は、クラスタOFC構成において、沈黙型故障が発生した際のフェイルオーバー処理におけるサービス停止時間、つまりノード切り替えを開始してからパケット転送を再開するまでの時間を示す説明図である。図13に示すように、クラスタOFC構成では、クラスタミドルウェアによるノード切り替え処理500と、オーディット処理501だけで、サービス再開できる。

10

【先行技術文献】

【特許文献】

【0039】

【特許文献1】再表WO2013/114490号

【発明の概要】

【発明が解決しようとする課題】

【0040】

上記(I)~(III)に記載した、基本的なOFCの動作を踏まえた上で、OFCにおけるビザンチン型故障について考える。

20

【0041】

OFCにおけるビザンチン型故障は、例えば、システム管理者が意図した通りにパケット転送されることもあれば、意図しない形でパケット転送されることもあるという、フロー制御が信用できない状態に陥った状態である。つまり、意図しないフローエントリが導出され、フローテーブルに登録されてしまった状態である。このような状態は、以下のような場合に発生し得る。

【0042】

(A)不正なコンフィグが設定されてしまった場合。例えば、コンフィグ設定作業にミスがあった場合や、内部犯行者により意図的にセキュリティホールを含むような論理ネットワークが設定された場合。

30

【0043】

(B)OFC機能(例えば、経路導出機能やトポロジ検出機能)が不正である場合。例えば、OFCに潜在的なバグがある場合や、不正使用者によりOFC機能が改ざんされた場合。

【0044】

(I)~(III)で説明したような構成では、ビザンチン型故障から復旧することが難しい。(A)の場合、一度OFCを再起動し、修正したコンフィグを再度読み込ませてサービスを再開することにより、OFSに設定された不正なフローが一掃でき、故障状態を解消できる。しかし、少なくとも図11に示す程度のサービス停止時間が発生する。(B)の場合、バグや改ざんを取り除いたOFCを新たに用意し、コンフィグをロードするなど、必要となる処理を行うことにより、サービスを再開できる。しかし、少なくともバグや改ざんを除去する時間と図11に示す時間とを合計したサービス停止時間が発生する。何れの場合も、図11に示す程度のサービス停止時間が発生する。

40

【0045】

そこで、本発明は、ビザンチン型故障の発生を検出した場合に、早期に故障状態を解消しパケット転送を再開可能とする通信システムおよび通信システムにおけるサービス復旧方法を提供することを目的とする。

【課題を解決するための手段】

50

【 0 0 4 6 】

本発明による通信システムは、フロー制御機能と仮想ネットワークの設定情報とをもとに複数のパケット転送装置を制御する第1の制御装置と、動作実績があるフロー制御機能と、運用実績がある仮想ネットワークの設定情報とを有し、第1の制御装置に故障が発生した場合に切り替え先となる第2の制御装置と、複数のパケット転送装置と各制御装置との間に配置された負荷分散装置とを備え、負荷分散装置は、第1の制御装置においてビザンチン型故障が発生した場合に、複数のパケット転送装置の制御を第2の制御装置に実行させる役割変更指示機能を有することを特徴とする。

【 0 0 4 7 】

本発明によるサービス復旧方法は、フロー制御機能と仮想ネットワークの設定情報とをもとに複数のパケット転送装置を制御する第1の制御装置と、動作実績があるフロー制御機能と、運用実績がある仮想ネットワークの設定情報とを有し、第1の制御装置に故障が発生した場合に切り替え先となる第2の制御装置と、複数のパケット転送装置と各制御装置との間に配置された負荷分散装置とを備えた通信システムにおいて、負荷分散装置が、第1の制御装置においてビザンチン型故障が発生した場合に、複数のパケット転送装置の制御を第2の制御装置に実行させることを特徴とする。

【 発明の効果 】

【 0 0 4 8 】

本発明によれば、通信システムにおいて、ビザンチン型故障の発生を検出した場合に、早期に故障状態を解消しパケット転送を再開することができる。

【 図面の簡単な説明 】

【 0 0 4 9 】

【 図 1 】 本発明による通信システムの第1の実施形態の構成を示すブロック図である。

【 図 2 】 アドレス変換表の一例を示す説明図である。

【 図 3 】 第1の実施形態におけるOFCの構成を示すブロック図である。

【 図 4 】 ロードバランサにおけるOFCの切り替え動作を示すフローチャートである。

【 図 5 】 ステップS5において更新されたアドレス変換表の一例を示す説明図である。

【 図 6 】 第1の実施形態の通信システムにおいてビザンチン型故障が発生した際の復旧に伴うサービス停止時間を示す説明図である。

【 図 7 】 本発明による通信システムの第2の実施形態の構成を示す説明図である。

【 図 8 】 本発明による通信システムの最小構成を示すブロック図である。

【 図 9 】 オープンフローシステムにおけるシングルOFC構成の一例を示すブロック図である。

【 図 1 0 】 OFCの構成の一例を示すブロック図である。

【 図 1 1 】 シングルOFC構成において、沈黙型故障が発生した際のOFC再起動によるサービス停止時間を示す説明図である。

【 図 1 2 】 オープンフローシステムにおけるクラスタOFC構成の一例を示すブロック図である。

【 図 1 3 】 クラスタOFC構成において、沈黙型故障が発生した際のフェイルオーバー処理におけるサービス停止時間を示す説明図である。

【 発明を実施するための形態 】

【 0 0 5 0 】

実施形態 1 .

以下、本発明の第1の実施形態を図面を参照して説明する。

【 0 0 5 1 】

本実施形態では、オープンフロープロトコル Ver 1 . 2 で実装されたマルチコントローラ機能を用いて、現用系でオープンフロー制御を行いつつ、待機系でも常時トポロジ情報とBC/MCフローエントリーの導出を行う通信システムを例として説明する。

【 0 0 5 2 】

本実施形態では、クラスタ構成されたOFC（以下、第1のOFCと呼称する）とは別

10

20

30

40

50

に、ビザンチン型故障が発生した場合に切り替え先となるOFC（以下、第2のOFCと呼称する）を配置する。また、OFSと、第1のOFC/第2のOFCとの間にロードバランサ（LB）を配置し、ロードバランサの操作によりOFSの接続先OFCを切り替えられるようにする。

【0053】

第2のOFCとして、経路導出やトポロジ検出機能にバグや改ざんがないOFC、つまり信用がおけるOFC機能を有するOFCを用意する。例えば、以下のような考え方で第2のOFCを用意する。

【0054】

- ・第1のOFCを構成するOFCノードにバグが見つかった場合、そのバグを改修したOFCを用意し、第2のOFCとして使用する。或いは、既に十分な動作実績があり、品質が安定している、古いバージョンのOFCを第2のOFCとして使用する、など。
- ・第1のOFCより厳しいセキュリティポリシーを適用し、機能改ざんなどの不正操作に対する対策が施されたOFCを第2のOFCとして使用する。例えば、第2のOFCにはより限定されたユーザのみアクセスできるように設定する、第2のOFCをより保護されたネットワークに配置する、など。

10

【0055】

また、第2のOFCには、信用がおけるコンフィグ情報をロードしておく。例えば、以下のような考え方で信用がおけるコンフィグを定義する。

【0056】

- ・十分運用実績があるコンフィグ
- ・更新できないように設定したコンフィグ
- ・シンプルな構成とし、容易に改ざんを見抜けるようなコンフィグ

20

【0057】

以上のように、本実施形態では、第2のOFCが、常に最も信用がおけるOFC機能と、最も信用がおけるコンフィグ情報とを有するように管理する。第1のOFCと第2のOFCは非対称になる。

【0058】

なお、トポロジ情報は、OFCで管理できるものではなく、LLDPパケットをOFSに送信しその応答を得て生成する情報である。本実施形態では、このトポロジ情報を、第1のOFC、第2のOFCで常に最新の情報を取得できるように、通信システムを構成する。それにより、第2のOFCでは、信用がおけるOFC機能およびコンフィグ情報と、最新のトポロジ情報とから、常にBC/MCフローエントリを導出できる状態とすることができる。その結果、通信システムは、切り替え時のサービス停止時間を短縮できる。

30

【0059】

図1は、本発明による通信システムの第1の実施形態の構成を示すブロック図である。なお、図1において点線で示す構成要素以外の要素については、図12に示す要素と同様であるため、詳細な説明を省略する。

【0060】

図1に示すように、通信システムは、第1のOFC600と、第2のOFC603と、ロードバランサ604とを備える。なお、図1には、ロードバランサ604に3台のOFS（OFS608～610）が接続されているが、OFSは、ロードバランサ604にいくつ接続されていてもよい。

40

【0061】

第1のOFC600は、クラスタミドルウェアにより冗長化構成されたOFCノードである。第1のOFC600は、OFC（ACT）601とOFC（SBY）602とを含む。なお、第1のOFC600は、2台以上のOFCで冗長化構成されていてもよい。

【0062】

OFC（ACT）601が保持する、コンフィグ611、トポロジ612、フローテーブル613は、それぞれ、OFC（SBY）602が保持する、コンフィグ614、トポ

50

ロジ 6 1 5、フローテーブル 6 1 6 と同期されている。

【 0 0 6 3 】

本実施形態では、第 1 の O F C 6 0 0 が沈黙型故障に対処することに加えて、第 2 の O F C 6 0 3 が、ビザンチン型故障に対処する。第 2 の O F C 6 0 3 は、前述したように、信用できる O F C 機能とコンフィグ情報とを有する。なお、第 2 の O F C は 1 台に限定されず、複数台配置されていてもよい。

【 0 0 6 4 】

O F S 6 0 8 ~ 6 1 0 が O F C に接続する際、各 O F S は、複数の O F C と通信可能な方法で接続する。各 O F S は、例えば、オープンフロープロトコル V e r 1 . 2 規格で定義されているマルチコントローラ機能などを用いて、複数の O F C と通信する。

10

【 0 0 6 5 】

マルチコントローラ機能では、それぞれの O F C には役割が与えられる。「マスタ (M a s t e r) 」の役割が与えられた O F C は、O F S からのメッセージを受け取り、そのメッセージに応答し、O F S が持つフローエントリを更新する権限を有する。一方、「スレーブ (S l a v e) 」の役割が与えられた O F C は、O F S からのメッセージを受け取ることができるが、メッセージに対する応答を返さない。つまり、スレーブの O F C は、O F S のフローテーブルを更新することはできない。

【 0 0 6 6 】

ロードバランサ 6 0 4 は、後述する役割変更指示機能 6 0 5 を有する。なお、役割変更指示機能 6 0 5 は、例えば、プログラムに従って動作するコンピュータの C P U によって実現される。当該プログラムは、例えば、ロードバランサ 6 0 4 の記憶装置 (図示せず) に記憶される。ロードバランサ 6 0 4 の C P U は、そのプログラムを読み込み、そのプログラムに従って、役割変更指示機能 6 0 5 として動作する。

20

【 0 0 6 7 】

ロードバランサ 6 0 4 は、マスタの O F C とスレーブの O F C に、それぞれ仮想アドレスを与える。また、ロードバランサ 6 0 4 は、O F C の実際のアドレス (以下、実アドレスという) との対応付けを、アドレス変換表 6 0 6 で管理する。ロードバランサ 6 0 4 の N W アドレス変換機能 6 0 7 は、アドレス変換表 6 0 6 を参照して、パケットヘッダの送信先・宛先情報を書き換える。図 2 は、アドレス変換表の一例を示す説明図である。

【 0 0 6 8 】

初期状態では、アドレス変換表 6 0 6 は、図 2 に示すようになっている。O F S は、ロードバランサ 6 0 4 の仮想アドレス 1 にアクセスすることで、結果的に第 1 の O F C 6 0 0 にアクセスできる。また、ロードバランサ 6 0 4 の仮想アドレス 2 にアクセスすることで、結果的に第 2 の O F C 6 0 3 にアクセスできる。

30

【 0 0 6 9 】

各 O F S は、接続する O F C のアドレス情報として、マスタとスレーブの二つのアドレス情報をもつ。各 O F S は、マスタの O F C のアドレスとしてロードバランサ 6 0 4 の仮想アドレス 1 を設定し、スレーブの O F C として仮想アドレス 2 を設定する。こうすることで、第 1 の O F C 6 0 0 と O F S の間でオープンフロープロトコルによる通信が成立し、フロー制御が可能になる。一方、第 2 の O F C 6 0 3 では、O F S から発行されたオープンフローメッセージを受信可能になる。

40

【 0 0 7 0 】

本実施形態における O F C (第 1 の O F C 6 0 0 の O F C (A C T) 6 0 1、O F C (S B Y) 6 0 2、第 2 の O F C 6 0 3) は、「役割管理機能」と、「拡張 OF プロトコル処理手段」とを有する。

【 0 0 7 1 】

図 3 は、第 1 の実施形態における O F C の構成を示すブロック図である。図 2 を参照して、O F C の機能について説明する。なお、図 3 において点線で示す構成要素以外の要素については、図 1 0 に示す要素と同様であるため、詳細な説明を省略する。

【 0 0 7 2 】

50

図3に示すOFC700は、図10に示すOFC100に含まれる構成要素に加え、拡張OFプロトコル処理手段701と、役割管理機能702とを有する。

【0073】

役割管理機能702は、当該OFC(OFC700)に、マスタとしての動作を行わせるか、スレーブとしての動作を行わせるかを管理する機能である。

【0074】

OFC700がマスタとして動作するように設定された場合、役割管理機能702は、拡張OFプロトコル処理手段701に対し、全てのオープンフローメッセージを受け付け、必要となる処理を行い、応答を返却するよう指示する。

【0075】

一方、OFC700がスレーブとして動作するように設定された場合、役割管理機能702は、拡張OFプロトコル処理手段701に対し、受信したオープンフローメッセージの内、LLDPのPACKET_INメッセージ以外を破棄するよう指示する。こうすることで、スレーブ側、つまり「スレーブ」の役割が与えられたOFC700では、トポロジ検出に必要なメッセージのみを受信できるようになる。LLDPのPACKET_INメッセージを受け取ったスレーブ側では、トポロジが更新されると、それに伴ってBC/MCフローエントリが導出され、フローテーブルに設定される。

【0076】

なお、拡張OFプロトコル処理手段701および役割管理機能702は、例えば、プログラムに従って動作するコンピュータのCPUによって実現される。当該プログラムは、例えば、OFC700の記憶装置(図示せず)に記憶される。OFC700のCPUは、そのプログラムを読み込み、そのプログラムに従って、拡張OFプロトコル処理手段701および役割管理機能702として動作する。また、拡張OFプロトコル処理手段701および役割管理機能702が別々のハードウェアで実現されていてもよい。

【0077】

次に、本実施形態の動作を説明する。

【0078】

まず、ビザンチン型故障が発生した場合に、ロードバランサ604が、コントローラを第1のOFCから第2のOFCに切り替える動作を説明する。

【0079】

なお、ビザンチン型故障の発生を検出する手段は様々あり得るが、例として以下のような場合に、ビザンチン型故障の発生を検出することができる。

【0080】

- ・OFC(ACT)のCPU使用率が異常に上がった場合。この場合、なんらかの要因で、多数の意図しないフロー設定が行われた可能性がある。
- ・OFSのフローテーブルに設定されたフローエントリの数が異常に増加した場合。この場合、多数の意図しないフローが設定された可能性がある。

【0081】

ビザンチン型故障が発生した場合、ロードバランサ604は、以下のようにしてOFCを切り替える。図4は、ロードバランサ604におけるOFCの切り替え動作を示すフローチャートである。

【0082】

まず、システム管理者は、ロードバランサ604の役割変更指示機能605を使用して、ロードバランサ604に対して、各コントローラの役割を変更するように指示する。つまり、役割変更指示機能605は、システム管理者から操作部(図示せず)等を介して、各コントローラの役割の変更指示を入力する(ステップS1)。

【0083】

なお、ロードバランサ604が、ビザンチン型故障の発生を検出して、ステップS2以降の処理を開始するようにしてもよい。例えば、ロードバランサ604が、OFCのCPU使用率を監視可能な場合には、当該CPU使用率が所定の閾値を超えたときに、ビザン

10

20

30

40

50

チン型故障が発生したと判断することができる。また例えば、ロードバランサ604が、OFSのフローエントリの数を取得可能な場合には、当該フローエントリの数が所定の閾値を超えたときに、ビザンチン型故障が発生したと判断することができる。

【0084】

ロードバランサ604は、ロードバランサ604とOFSの間で確立していた通信を切断する(ステップS2)。このとき、ロードバランサ604は、OFSからの接続要求の受け入れを停止し、ステップS6が完了するまで、OFSからの接続要求に応じない。

【0085】

役割変更指示機能605は、第1のOFC600の役割管理機能に対し、スレーブとして動作するよう通知する(ステップS3)。これ以降、第1のOFC600は、LLDPのPACKET_INメッセージのみを受信ようになる。

10

【0086】

役割変更指示機能605は、第2のOFC603の役割管理機能に対し、マスタとして動作するよう通知する(ステップS4)。これ以降、第2のOFC603は、すべてのオープンフローメッセージを受信ようになる。

【0087】

役割変更指示機能605は、アドレス変換表606を更新し、マスタの実アドレスとして、第2のOFC603の物理アドレスを設定する。また、スレーブの実アドレスとして、第1のOFC600の仮想アドレスを設定する(ステップS5)。図5は、ステップS5において更新されたアドレス変換表606の一例を示す説明図である。

20

【0088】

ロードバランサ604は、OFSからの接続要求の受け入れを再開する(ステップS6)。

【0089】

OFSと第2のOFC603との間で接続が確立され、オープンフロープロトコルによる通信が成立する(ステップS7)。OFSと第1のOFC600の間でも接続が確立され、第1のOFC600はLLDPのPACKET_INメッセージのみを受信ようになる。

【0090】

第2のOFC603は、オーディット処理を実行する(ステップS8)。これにより、第2のOFC603のフローテーブル619が、OFSのフローテーブル620~622に同期される。つまり、OFSにはBC/MCフローのみが設定された状態となり、OFCの切り替え前に設定されていた意図しないフローは一掃される。

30

【0091】

上記ステップS1~S8の処理により、図1に示す通信システムは、信用できるOFCとコンフィグ情報により、フロー制御を再開できるようになる。なお、サービス再開直後は、ソフト転送となる。

【0092】

以上に説明したように、本実施形態では、マスタのOFCにおいてビザンチン型故障が発生した場合に、最も信用がおけるOFC機能と最も信用がおけるコンフィグ情報とを有するスレーブのOFCをマスタとして動作させる。それにより、上記(A)に示すコンフィグ情報の不正、上記(B)に示すOFC機能の不正、という二つの原因に由来するビザンチン型故障からの復旧が可能となる。

40

【0093】

また、本実施形態では、スレーブのOFCが、トポロジ検出に必要なメッセージを受信し、BC/MCフローエントリを導出する。それにより、OFCの役割を「スレーブ」から「マスタ」に切り替える際、当該OFCにおいてBC/MCフローエントリ導出処理等が不要となり、サービス停止時間を短縮できる。ビザンチン型故障からの復旧に伴うサービス停止時間につき、特に対策しない場合は少なくとも図11に示す程度のサービス停止時間が発生するが、本実施形態によれば、図6に示す程度のサービス停止時間となる。

50

図6は、第1の実施形態の通信システムにおいてビザンチン型故障が発生した際の復旧に伴うサービス停止時間を示す説明図である。

【0094】

また、第2のOFCを、第1のOFCとは離れた場所、例えば、第1のOFCと異なるラック、異なるフロア、異なるDCに配置することで、本発明をディザスタリカバリ対策としても応用できる。

【0095】

実施形態2.

以下、本発明の第2の実施形態を図面を参照して説明する。

【0096】

本実施形態では、マルチコントローラ機能を使用しないOFSを制御するOFCを含む通信システムを例にする。マルチコントローラ機能を使用しないOFSは、単一のコントローラにしかメッセージを送信できない。つまり、第1のOFCのみでトポロジ情報を検出できる。図7は、本発明による通信システムの第2の実施形態の構成を示す説明図である。

【0097】

図7に示す通信システムの構成は、第1の実施形態と同様である。

【0098】

ただし、ロードバランサ1107は、役割変更指示機能を有さない。また、第1のOFC1100のOFC(ACT)1101、OFC(SBY)1102、第2のOFC1103はそれぞれ、役割管理機能および拡張OFプロトコル処理手段の代わりに、トポロジ変更イベント送受信機能1104、1105、1106を有する。

【0099】

ロードバランサ1107は、NWアドレス変換機能1108により、単一の仮想アドレスを管理する。つまり、ロードバランサ1107は、アドレス変換表1109で、当該単一の仮想アドレスとOFCの実アドレスとの対応付けを管理する。図7に示す例では、当該単一の仮想アドレスとして、「LBの仮想IPアドレス」が設定されている。また、OFCの実アドレスとして、「第1のOFCの仮想アドレス」が設定されている。OFS1108~1110は、ロードバランサ1107の当該単一の仮想アドレスに接続する。

【0100】

ロードバランサ1107は、アドレス変換表1109に従い、第1のOFC1100または第2のOFC1103にパケットを中継する。

【0101】

OFS1110~1112と接続されたOFCは、オープンフロー制御を行う。また、各OFCは、受信したオープンフローメッセージがLLDPのPACKET_INメッセージだった場合は、トポロジ変更イベント送受信機能を用いて、対向のOFCノードにメッセージを通知する。これにより、オープンフローメッセージを直接受信していないOFCノードでもトポロジ情報を生成することが可能になり、結果、常時BC/MCフローエントリを導出することが可能となる。

【0102】

このように、ロードバランサが保持するアドレス変換表の実アドレスに、切り替え先OFCのアドレスを設定することで、ノード切り替えを実行できる。つまり、本実施形態によれば、OFCのOFプロトコル処理手段等にバグや改ざんなどがなく、対向ノードに伝達するメッセージに誤りが含まれる可能性がなければ、図7に示すような簡易な構成により、第1の実施形態と同様の効果を得ることができる。

【0103】

なお、各OFCのトポロジ変更イベント送受信機能(トポロジ変更イベント送受信機能1104、1105、1106)は、例えば、プログラムに従って動作するコンピュータのCPUによって実現される。当該プログラムは、例えば、各OFCの記憶装置(図示せず)に記憶される。各OFCのCPUは、そのプログラムを読み込み、そのプログラムに

10

20

30

40

50

従って、トポロジ変更イベント送受信機能（トポロジ変更イベント送受信機能 1104、1105、1106）として動作する。

【0104】

以上、各実施形態においてオープンフローに適用される通信システムを例にしたが、本発明はオープンフロー以外にも適用可能である。例えば、各実施形態における OFC は、オープンフローにおけるコントローラ以外の制御装置であってもよい。また例えば、各実施形態における OFS は、オープンフローにおけるスイッチ以外のパケット転送装置であってもよい。つまり、通信ネットワーク上の各パケット転送装置を制御装置が集中管理する構成の通信システムであれば、本発明を適用することができる。

【0105】

次に、本発明の概要を説明する。図 8 は、本発明による通信システムの最小構成を示すブロック図である。本発明による通信システムは、フロー制御機能（例えば、オープンフローにおける OFC 機能）と仮想ネットワークの設定情報（例えば、オープンフローにおけるコンフィグ情報）とをもとに複数のパケット転送装置を制御する第 1 の制御装置 10（図 1 に示す第 1 の OFC 600 に相当）と、動作実績があるフロー制御機能と、運用実績がある仮想ネットワークの設定情報とを有し、第 1 の制御装置 10 に故障が発生した場合に切り替え先となる第 2 の制御装置 20（図 1 に示す第 2 の OFC 603 に相当）と、複数のパケット転送装置と各制御装置との間に配置された負荷分散装置 30（図 1 に示すロードバランサ 604 に相当）とを備え、負荷分散装置 30 は、第 1 の制御装置 10 においてビザンチン型故障が発生した場合に、複数のパケット転送装置の制御を第 2 の制御装置 20 に実行させる役割変更指示機能 31（図 1 に示すロードバランサ 604 における役割変更指示機能 605 に相当）を有する。

【0106】

そのような構成によれば、第 1 の制御装置においてビザンチン型故障が発生した場合に、最も信用が与えられるフロー制御機能と最も信用が与えられるネットワークの設定情報とを有する第 2 の制御装置を動作させることができる。それにより、例えば、オープンフローにおいて、コンフィグ情報の不正、OFC 機能の不正、という二つの原因に由来するビザンチン型故障からの復旧が可能となる。従って、通信システムにおいて、ビザンチン型故障の発生を検出した場合に、早期に故障状態を解消しパケット転送を再開することができる。

【0107】

また、負荷分散装置 30 の役割変更指示機能 31 は、第 1 の制御装置 10 の CPU 使用率が所定の閾値を超えたときに、当該第 1 の制御装置においてビザンチン型故障が発生したと判断してもよい。そのような構成によれば、ビザンチン型故障をより確実に検出することができ、より早期に故障状態を解消しパケット転送を再開することができる。

【0108】

また、通常運用時、第 1 の制御装置 10 は、複数のパケット転送装置から送信される全てのメッセージを受信するマスターとして動作し、第 2 の制御装置 20 は、複数のパケット転送装置からトポロジ検出に必要なメッセージのみを受信するスレーブとして動作し、負荷分散装置 30 の役割変更指示機能 31 は、第 1 の制御装置 10 においてビザンチン型故障が発生した場合に、第 1 の制御装置 10 をスレーブとして動作させ、第 2 の制御装置 20 をマスターとして動作させてもよい。そのような構成によれば、第 2 の制御装置は、マスターとしてフロー制御を開始する際に必要なブロードキャストやマルチキャストされたパケットを適切に転送するための情報（例えば、オープンフローにおける BC/MC フロー）を、スレーブ動作中に計算し保持しておくことができる。

【0109】

また、各制御装置は、負荷分散装置 30 からの役割変更指示に基づいて、自装置をマスターとして動作させるか、スレーブとして動作させるかを管理する役割管理機能（図 3 に示す役割管理機能 702 に相当）を有していてもよい。そのような構成によれば、負荷分散装置は、役割変更指示を出力するだけで各制御装置の役割を変更することが可能となる。

【0110】

10

20

30

40

50

また、第2の制御装置20の役割管理機能は、スレーブ動作中に、トポロジ検出に必要なメッセージ（例えば、オープンフローにおけるLLDPのPACKET_INメッセージ）を受信すると、当該メッセージをもとにマスタとして動作を開始する際に必要なフロー（例えば、オープンフローにおけるBC/MCフロー）を作成してもよい。そのような構成によれば、第2の制御装置20は、マスタとしてフロー制御を開始する際に、マスタとして動作を開始する際に必要なフローを作成する必要がないので、サービス停止時間をより短縮することができる。

【0111】

また、各制御装置は、パケット転送装置から受信したメッセージがトポロジ検出に必要なメッセージであると判断した場合に、他の制御装置に当該メッセージを通知するトポロジ変更イベント送受信機能（図7に示すトポロジ変更イベント送受信機能1104～1106に相当）を有してもよい。そのような構成によれば、より簡易な構成で、ビザンチン型故障の発生を検出した場合の、故障状態の早期解消およびパケット転送の再開を行うことができる。

10

【0112】

また、第1の制御装置10が、複数の制御装置（例えば、図1に示すOFC（ACT）601、OFC（SBY）602）によりクラスタ構成されていてもよい。そのような構成によれば、通信システムを、ビザンチン型故障だけでなく、沈黙型故障にも対応させることができる。

20

【符号の説明】

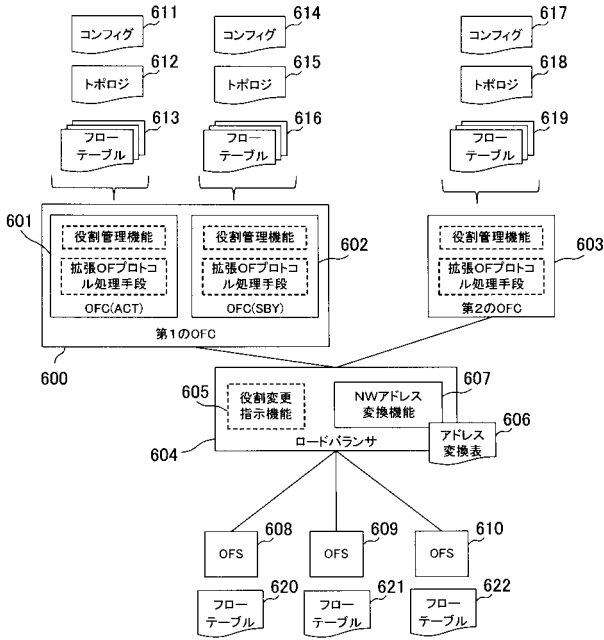
【0113】

- 10 第1の制御装置
- 20 第2の制御装置
- 30 負荷分散装置
- 31、605 役割変更指示機能
- 100、700 OFC
- 101～103、403～405、608～610、1110～1112 OFS
- 110、206、410、413、611、614、617 コンフィグ
- 111、204、411、414、612、615、618 トポロジ
- 112、113～115、208、412、415、417～419、613、616
、619、620～622 フローテーブル
- 201 通信手段
- 202 OFプロトコル処理手段
- 203 トポロジ検出手段
- 205 コンフィグ管理手段
- 207 経路管理手段
- 209 経路計算手段
- 210 経路設定手段
- 211 オーディット手段
- 400、601、1101 OFC（ACT）
- 401、602、1102 OFC（SBY）
- 600、1100 第1のOFC
- 603、1103 第2のOFC
- 604、1107 ロードバランサ
- 607、1108 NWアドレス変換機能
- 606、1109 アドレス変換表
- 701 拡張OFプロトコル処理手段
- 702 役割管理機能
- 1104～1106 トポロジ変更イベント送受信機能

30

40

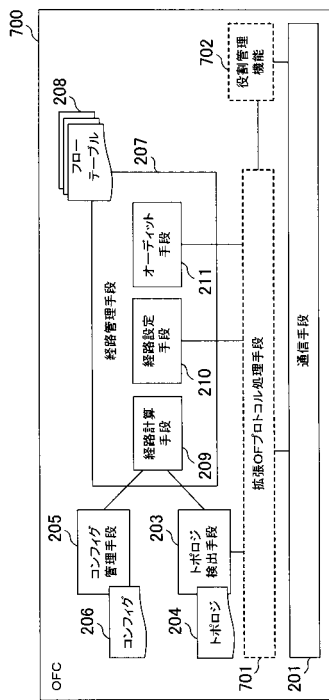
【 図 1 】



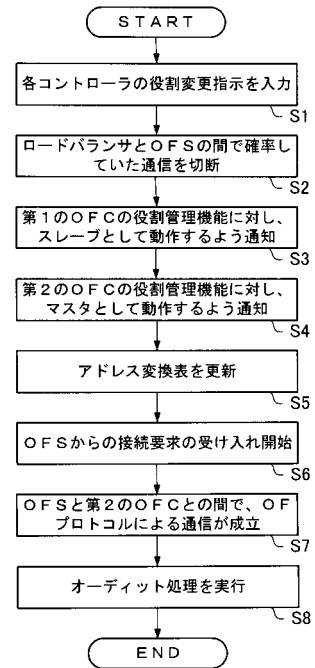
【 図 2 】

OFCの役割	仮想アドレス	実アドレス
Master	LBの仮想アドレス1	第1のOFCの仮想アドレス
Slave	LBの仮想アドレス2	第2のOFCの物理アドレス

【 図 3 】



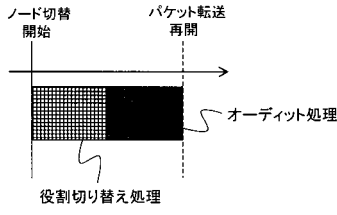
【 図 4 】



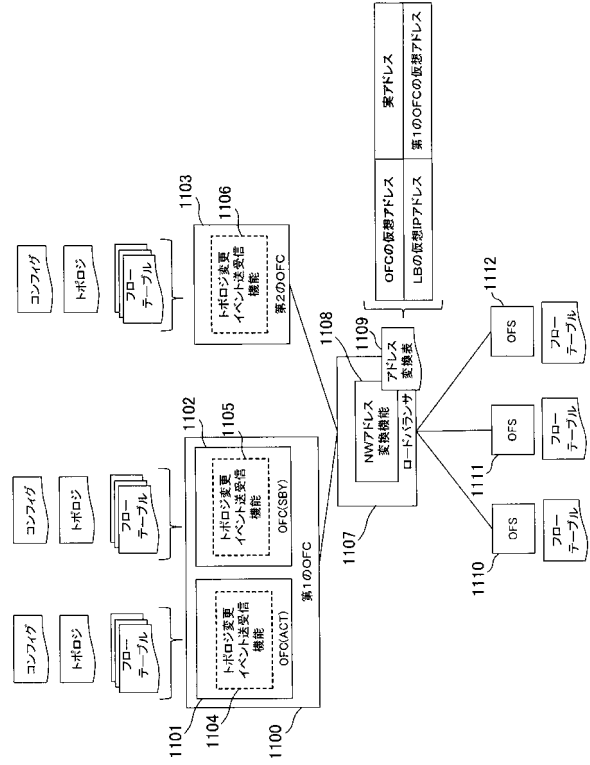
【 図 5 】

OFCの役割	仮想アドレス	実アドレス
Master	LBの仮想アドレス1	第2のOFCの物理アドレス
Slave	LBの仮想アドレス2	第1のOFCの仮想アドレス

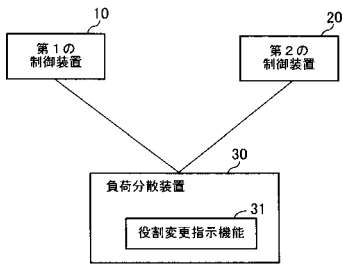
【 図 6 】



【 図 7 】

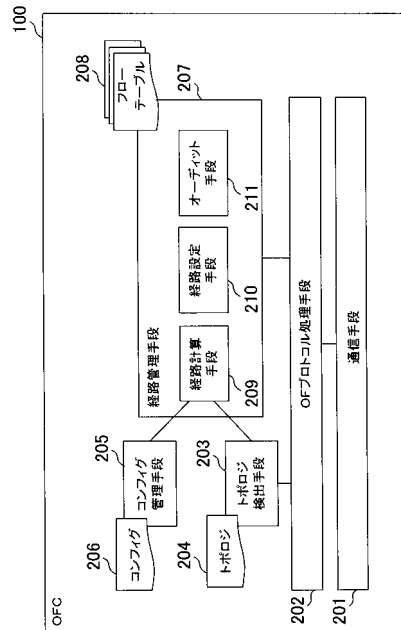
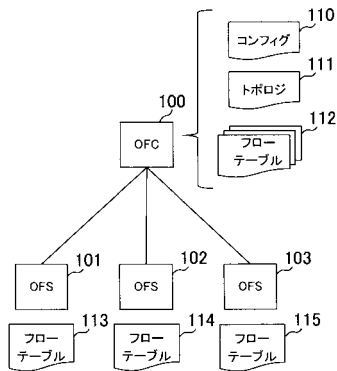


【 図 8 】

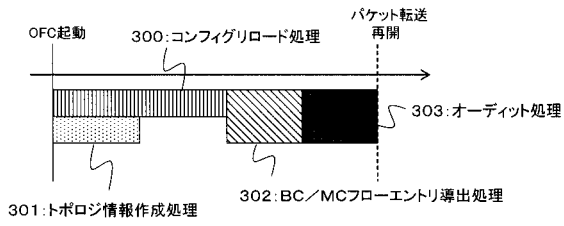


【 図 10 】

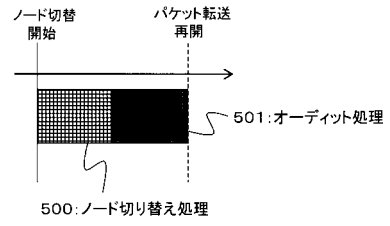
【 図 9 】



【 図 1 1 】



【 図 1 3 】



【 図 1 2 】

