

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5872731号  
(P5872731)

(45) 発行日 平成28年3月1日(2016.3.1)

(24) 登録日 平成28年1月22日(2016.1.22)

(51) Int.Cl.		F I			
HO 4 L	12/28	(2006.01)	HO 4 L	12/28	2 0 0 M
HO 4 L	12/713	(2013.01)	HO 4 L	12/713	
HO 4 L	12/70	(2013.01)	HO 4 L	12/70	1 0 0 Z

請求項の数 15 (全 14 頁)

(21) 出願番号	特願2015-501687 (P2015-501687)	(73) 特許権者	501113353
(86) (22) 出願日	平成25年2月28日 (2013.2.28)		シマンテック コーポレーション
(65) 公表番号	特表2015-511099 (P2015-511099A)		Symantec Corporation
(43) 公表日	平成27年4月13日 (2015.4.13)		アメリカ合衆国, カリフォルニア州 94
(86) 国際出願番号	PCT/US2013/028346		043, マウンテン ビュー, エリス ス
(87) 国際公開番号	W02013/142023		トリート 350
(87) 国際公開日	平成25年9月26日 (2013.9.26)	(74) 代理人	100147485
審査請求日	平成27年7月1日 (2015.7.1)		弁理士 杉村 憲司
(31) 優先権主張番号	13/425, 127	(74) 代理人	100134119
(32) 優先日	平成24年3月20日 (2012.3.20)		弁理士 奥町 哲行
(33) 優先権主張国	米国 (US)		
早期審査対象出願			

最終頁に続く

(54) 【発明の名称】 クラスタの複数のノードのそれぞれに対してリンクの障害の検出を伝えるためのコンピュータ実装方法、非一時的なコンピュータ可読媒体およびコンピュータシステム

(57) 【特許請求の範囲】

【請求項1】

クラスタの複数のノードのそれぞれに対してリンクの障害の検出を伝えるためのコンピュータ実装方法であって、該方法が、

該クラスタの特定のノードの特定のリンクの障害を含む、リンクダウン事象を検出する工程と、

該クラスタの各ノードに通信可能に連結される別のリンクであってリンクダウン事象通知メッセージをブロードキャストするための専用のリンクを使用して、該クラスタの該複数のノードに該特定のリンクの該障害の通知を伝播する工程であって、これにより該ノードが該リンクダウン事象を平行して処理する、工程と、

を含み、該特定のリンクの該障害の該通知の該伝播は、該クラスタの該ノードが、対応するハートビートの期限切れにより該リンクダウン事象を知るよりも前に、該通知を受け取るように実行される、コンピュータ実装方法。

【請求項2】

前記クラスタの前記複数のノードに前記特定のリンクの前記障害の通知を伝播する工程が、前記クラスタの前記ノードにメッセージをブロードキャストする工程であって該メッセージは前記特定のリンクの前記障害を前記クラスタの前記ノードに通知する、工程を更に含む、請求項1に記載の方法。

【請求項3】

前記別のリンクを使用して前記クラスタの該ノードに該メッセージをブロードキャスト

する工程を更に含む、請求項 2 に記載の方法。

【請求項 4】

前記クラスタの各ノードに通信可能に連結された別のリンクを確保する工程であって、該リンクがリンクダウン事象通知メッセージのブロードキャスト専用である、工程が、該別のリンクを、前記クラスタのユーザーに公開されていないプライベートリンクとして確保する工程を更に含む、請求項 3 に記載の方法。

【請求項 5】

前記リンクダウン事象を処理する前に前記ブロードキャストされたメッセージが前記クラスタの前記ノードに受信されたことを確実にする工程を更に含む、請求項 2 に記載の方法。

10

【請求項 6】

前記リンクダウン事象を処理する前に前記ブロードされたキャストメッセージが、前記クラスタの前記ノードにより受信されたことを確実にする工程が、前記リンクダウン事象を処理する前に前記クラスタの前記ノードのそれぞれから前記ブロードキャストされたメッセージを受信したことを確認する確認応答を受け取るのを待つ工程を更に含む、請求項 5 に記載の方法。

【請求項 7】

所与の時間内に前記ブロードキャストされたメッセージの受信を確認する確認応答を受信していないことに応答して、前記ブロードキャストされたメッセージを再送信する工程を更に含む、請求項 5 に記載の方法。

20

【請求項 8】

前記リンクの障害検出に対応して、前記リンクの該障害の前記通知を伝播する前に特定の猶予期間待つ工程を更に含む、請求項 1 に記載の方法。

【請求項 9】

前記クラスタの前記複数のノードに前記特定のリンクの前記障害の通知を伝播する工程が、前記クラスタの前記ノードのすべてにアクセス可能な中央集中型コンピューティングデバイスに、前記特定のリンクの前記障害の通知を送信する工程を更に含む、請求項 1 に記載の方法。

【請求項 10】

前記クラスタの各ノードが、リンクの障害の通知のために、特定の周波数で、前記中央集中型コンピューティングデバイスのポーリングを行い、これにより前記クラスタの前記ノードが、前記中央集中型コンピューティングデバイスのポーリングにより前記特定のリンクの前記障害を知る、請求項 9 に記載の方法。

30

【請求項 11】

前記中央集中型コンピューティングデバイスが、前記クラスタの前記ノードに、前記特定のリンクの前記障害の前記通知を送信する、請求項 9 に記載の方法。

【請求項 12】

前記クラスタの前記特定のノードの前記特定のリンクの障害を検出する工程が、前記特定のノードのオペレーティングシステムから、前記特定のリンクの該障害の通知を受け取る工程を更に含む、請求項 1 に記載の方法。

40

【請求項 13】

クラスタの複数のノードそれぞれにリンクの障害の検出を伝えるための、コンピュータプログラムを格納する少なくとも 1 つの非一時的なコンピュータ可読媒体であって、該コンピュータプログラムが、

該クラスタの特定のノードの特定のリンクの障害を含む、リンクダウン事象を検出するためのプログラムコードと、

該クラスタの各ノードに通信可能に連結される別のリンクであってリンクダウン事象通知メッセージをブロードキャストするための専用のリンクを使用して、該クラスタの該複数のノードに該特定のリンクの該障害の通知を伝播し、これにより該ノードが該リンクダウン事象を平行して処理するためのプログラムコードと、

50

を含み、該特定のリンクの該障害の該通知の該伝播は、該クラスタの該ノードが、対応するハートビートの期限切れにより該リンクダウン事象を知るよりも前に、該通知を受け取るように実行される、コンピュータ可読媒体。

【請求項 14】

前記クラスタの前記複数のノードに前記特定のリンクの前記障害の通知を伝播するための前記プログラムコードが、前記クラスタの前記ノードにメッセージをブロードキャストするためのプログラムコードを更に含み、該メッセージは、前記特定のリンクの前記障害を前記クラスタの前記ノードに通知する、請求項 13 に記載のコンピュータプログラム。

【請求項 15】

クラスタの複数のノードのそれぞれに対してリンクの障害の検出を伝播するためのコンピュータシステムであって、該コンピュータシステムが、

プロセッサと、

コンピュータメモリと、

該クラスタの特定のノードの特定のリンクの障害を含むリンクダウン事象を検出するための手段と、

該クラスタの各ノードに通信可能に連結される別のリンクであってリンクダウン事象通知メッセージをブロードキャストするための専用のリンクを使用して、該クラスタの該複数のノードに該特定のリンクの該障害の通知を伝播し、これにより該ノードが該リンクダウン事象を平行して処理する、手段と、

を含み、該特定のリンクの該障害の該通知の該伝播は、該クラスタの該ノードが、対応するハートビートの期限切れにより該リンクダウン事象を知るよりも前に、該通知を受け取るように実行される、コンピュータシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、概ね、コンピュータクラスタの管理に関し、より具体的には、リンク障害のクラスタ全体の一貫した検出に関する。

【背景技術】

【0002】

高可用性クラスタ（別名、HAクラスタ、又はフェイルオーバークラスタ）は、ダウンタイムを最小限に抑え、サーバアプリケーションの実行をサポートする、コンピュータ（ノード）群である。高可用性クラスタは、個々のシステムコンポーネントに障害が生じた場合にも連続したサービスを提供するために、冗長なコンピュータリソース群を使用する。より具体的には、高可用性クラスタは、複数のサーバ、複数のネットワーク接続、冗長なデータ記憶装置などを提供することにより、単一の障害ポイントを排除する。クラスタがない場合、特定のアプリケーションを実行しているサーバに障害が生じた場合、このアプリケーションは、このサーバが復帰するまで利用できないことになる。高可用性クラスタリングでは、サーバ（又は、例えばネットワークアダプタ、記憶装置などの、これにより使用されている特定のコンピューティングリソース）の障害が検出される。障害が生じているサーバで実行中であつたアプリケーションを使用できる状態に維持できるようにするため、工程が自動的に行われる。これは、異なるネットワークリソース（例えばネットワークアダプタ）を使用するか、又は別のコンピューティングシステム（すなわちクラスタの別のノード）上でアプリケーションを自動的に再起動する、サーバとアプリケーションの再起動の形態である場合がある。このプロセスは「フェイルオーバー」と呼ばれる。高可用性クラスタはしばしば、企業データベース、重要な業務用アプリケーション、電子商取引ウェブサイトなどの重要なサーバアプリケーション用に使用される。そのようなアプリケーションでは、特に、銀行業務及び通信などの特定の業種において、たとえ短時間であってもダウンタイムが起こると非常に高価なものとなる場合がある。

【0003】

コンピュータ記憶装置においては、論理ボリューム管理が、大容量記憶装置のスペースを割り当てる柔軟な方法である。特に、ボリュームマネージャは、下位の物理的パーティションを連結させ、整合させ、又は組み合わせて、より大きな仮想パーティションにすることができる。よって管理者は、可能性としてはシステム利用を中断することなく、論理ボリュームのリサイズ又は移動を行うことができる。クラスタボリュームマネージャは、クラスタの複数のノードにわたってボリューム論理管理を拡張し、これにより各ノードは同じ論理ボリュームレイアウトを認識し、全ノードの全ボリュームリソースの同じ状態を認識する。クラスタボリューム管理の下では、ディスク又はボリューム構成に対して、クラスタ内のどのノードから変更が行われても、クラスタの全ノードにより認識される。クラスタレベルのボリューム管理をサポートするため、新しいノードがクラスタに加わったとき、及び既存のノードがクラスタから離れるときに、クラスタが再構成される。

10

【発明の概要】

【発明が解決しようとする課題】

【0004】

クラスタは、データ用ノードと管理通信との間の相互接続（リンク）を使用している。リンクに障害が生じた場合は、修正処置を行う必要がある。フェイルオーバー及びその他のクラスタ処置からの復帰の効率と適時性は、クラスタがリンク障害を検出及び処理可能なスピードに依存する。個々のノードは、リンクのいずれかに障害が生じているかどうかを検出するのに、オペレーティングシステムのサービスを使用することができる。2つのノードがスイッチなしで（クロスオーバーケーブルを使って）直接接続されているとき、両方のノードがリンクダウン通知を受け取り、これにより両ノードが平行してこの事象に対処することができる。しかしながら、2つ以上のノードがスイッチを使用して連結されている場合は、ローカルノードだけがリンク障害をリアルタイムで知る。よって、クラスタのネットワークポロジーにより、典型的には単一のノードのみ、又はローカライズされたノードのサブセットのみが、リンク障害を即座に知る。従来、他のノードはその後に、クラスタ内の各ノードの状況を監視するのに使用されるクラスタ全体のハートビート機構を介して、接続障害を知る。ハートビート機構は、個々のノードがローカルリンク障害を知るのにかかる時間に比べ、比較的遅い。クラスタの適切な管理は、修正処置を行う前に、クラスタの各ノードがリンク障害を認識することに依存し得る。ハートビート機構を介して各ノードが障害を知るまでの時間によって、ダウンした接続に必要なクラスタ再構成又はフェイルオーバー処置を遅らせることになり、これにより、クラスタ管理（例えばフェンシングアービトレーション（arbitration）判断の実行）に更なる問題を引き起こす。更に、ノード又はクラスタが異なるタイミングでリンク障害を知ると、不正又は望ましくない修正処置が生じる場合がある。

20

30

【0005】

これらの問題に対処することが望ましいと考えられる。

【課題を解決するための手段】

【0006】

リンク障害の通知は、クラスタ内の複数のノードのそれぞれに伝えられる。例えば特定のノードのオペレーティングシステムを介して、クラスタの特定のノードの特定のリンクの障害が検出される。特定のリンクの障害の通知は、クラスタの複数のノードに伝播され、これによりノードがリンクダウン事象を平行して処理する。リンク障害の通知の伝播は、クラスタのノードが、対応するハートビートの期限切れによりリンクダウン事象を知るよりも前に、通知を受け取るように実行される。いくつかの実施形態において、リンクの障害の通知は、リンクが即座に復帰した場合、特定の猶予期間待った後に限り、伝播される。

40

【0007】

一実施形態において、リンク障害の通知は、特定のリンクの障害をクラスタのノードに通知するメッセージをブロードキャストすることにより、伝播される。別のリンクが確保され、これはクラスタの各ノードに接続され、リンクダウン事象通知メッセージをブロー

50

ドキャストするための専用となる。この別のリンクは、クラスタのノードにメッセージをブロードキャストするのに使用することができる。この別のリンクは、クラスタのユーザーに公開されないプライベートリンクの形態であり得るが、必ずしもプライベートリンクである必要はない。一実施形態において、ブロードキャストメッセージがクラスタのノードに受信されたことを確実にしてから、リンクダウン事象が処理される。これは、例えば、クラスタのノードそれぞれからブロードキャストメッセージを受信したことを確認する確認応答を受け取るのを待つことによりなされてもよい。一実施形態において、ブロードキャストメッセージの受信を確認する確認応答が、所与の時間以内に各ノードから受信されない場合、そのブロードキャストメッセージは再送信される。

【0008】

10

別の一実施形態において、リンク障害の通知は、ノードのすべてにアクセス可能な中央集中型コンピューティングデバイスに通知を送信することにより、クラスタのノードに伝播される。クラスタのノードはリンク障害の通知のため特定の周波数で中央集中型コンピューティングデバイスのポーリングを行い、そのポーリングから特定のリンクの障害を知ることができる。一実施形態において、ポーリングの代わりに（又はこれに加えて）、中央集中型コンピューティングデバイスは、特定のリンクの障害の通知を、クラスタのノードに送信する。

【0009】

この「発明の概要」及び後述の「発明を実施するための形態」に記述される特徴及び利点は、すべてを包含したものではなく、特に、本特許の図面、明細書、及び請求項の見地から、関連分野の当業者には、数多くの付加的な特徴及び利点が明らかとなる。更に、本明細書に使用されている表現は、読みやすさと説明の目的のために主に選択されており、発明の主題を限定又は制限するために選択されたものではなく、そのような発明の主題を決定するには、請求項に依ることが必要である。

20

【図面の簡単な説明】

【0010】

【図1】いくつかの実施形態による、リンク障害管理システムが実装され得る代表的なネットワークアーキテクチャのブロック図である。

【図2】いくつかの実施形態による、リンク障害管理システムを実装するのに好適なコンピュータシステムのブロック図である。

30

【図3】いくつかの実施形態による、例示的なクラスタトポグラフィにおいてリンクの障害を検出する、リンク障害管理システムのブロック図である。

【図4】いくつかの実施形態による、ブロードキャストメッセージを介してリンク障害の通知をクラスタのノードに伝播する、リンク障害管理システムのブロック図である。

【図5】いくつかの実施形態による、中央集中型コンピューティングデバイスを介してリンク障害の通知をクラスタのノードに伝播する、リンク障害管理システムのブロック図である。

【図6】一実施形態によるリンク障害管理システムのオペレーションのフローチャートである。

【図7】別の一実施形態によるリンク障害管理システムのオペレーションのフローチャートである。

40

【0011】

図には、単に例示目的のため、様々な実施形態が示されている。当業者には、本明細書に記述されている原理から逸脱することなく、本明細書に示されている構造及び方法の別の実施形態を採用し得ることが、下記の議論から容易に理解されよう。

【発明を実施するための形態】

【0012】

図1は、リンク障害管理システム101が実装され得る代表的なネットワークアーキテクチャ100を示すブロック図である。図示のネットワークアーキテクチャ100は、複数のクライアント103A、103B及び103N、並びに複数のサーバ105A及び1

50

05Nを含む。図1において、リンク障害管理システム101は、サーバ105A上にあるものとして図示されている。これは単に一例であり、様々な実施形態においてこのシステム101の様々な機能性が、サーバ105、クライアント103で実例を示すことができ、あるいは、複数のクライアント103及び/又はサーバ105間に分散できることが理解されよう。

#### 【0013】

クライアント103及びサーバ105は、例えば図2及び下記に示すもののような、コンピュータシステム210を使用して実装することができる。クライアント103及びサーバ105は、図2に関連して下記で述べられるように、例えばネットワークインタフェース248又はモデム247を介して、ネットワーク107に通信可能に連結されている。クライアント103は、例えばウェブブラウザ又はその他のクライアントソフトウェア(図示なし)を用いて、サーバ105上のアプリケーション及び/又はデータにアクセス可能である。

10

#### 【0014】

図1には3つのクライアント及び2つのサーバが一例として示されているが、実際にはもっと多くの(又はもっと少ない)数のクライアント103及び/又はサーバ105を配備することができる。一実施形態において、ネットワーク107はインターネットの形態である。他の実施形態において、他のネットワーク107又はネットワークベースの環境を使用することができる。

#### 【0015】

図2は、リンク障害管理システム101を実装するのに好適なコンピュータシステム210のブロック図である。クライアント103とサーバ105の両方が、そのようなコンピュータシステム210の形態で実装され得る。図示のように、コンピュータシステム210の1つの構成要素がバス212である。バス212は、コンピュータシステム210の他の構成要素、例えば、少なくとも1つのプロセッサ214、システムメモリ217(例えばランダムアクセスメモリ(RAM)、リードオンリーメモリ(ROM)、フラッシュメモリ)、入出力(I/O)コントローラ218、スピーカシステム220などの外部音声装置に通信可能に連結された音声出力インタフェース222、ディスプレイ画面224などの外部ビデオ出力装置に通信可能に連結されたディスプレイアダプタ226、1つ以上のインタフェース(例えばシリアルポート230、ユニバーサルシリアルバス(USB)のレセプタクル230、パラレルポート(図示なし)など)、キーボード232に通信可能に連結されたキーボードコントローラ233、少なくとも1つのハードディスク244(又はその他の形態の磁気媒体)に通信可能に連結された記憶装置インタフェース234、フロッピーディスク238(「フロッピー」は登録商標、以下同じ)を受容するよう構成されたフロッピーディスクドライブ237、ファイバーチャネル(FC)ネットワーク290と接続するよう構成されたホストバスアダプタ(HBA)インタフェースカード235A、SCSIバス239に接続するよう構成されたHBAインタフェースカード235B、光ディスク242を受容するよう構成された光ディスクドライブ240、バス212に(例えばUSBレセプタクル228を介して)連結されたマウス246(又はその他のポインティングデバイス)、バス212に(例えばシリアルポート230を介して)接続されたモデム247、及び、例えばバス212に直接連結されたネットワークインタフェース248と、通信可能に連結される。

20

30

40

#### 【0016】

他の構成要素(図示なし)は、同様の方法で接続され得る(例えば文書スキャナ、デジタルカメラ、プリンタなど)。逆に、図2に示されているすべての構成要素は必ずしも存在する必要はない。構成要素は、図2に示されるものとは異なる手法で相互接続することができる。

#### 【0017】

バス212は、プロセッサ214と、上述のように、ROM及び/又はフラッシュメモリ並びにRAMを含み得るシステムメモリ217との間のデータ通信を可能にする。RA

50

Mは典型的に、オペレーティングシステム及びアプリケーションプログラムがロードされるメインメモリである。ROM及び/又はフラッシュメモリは、他のコードと共に、特定の基本的なハードウェアオペレーションを制御するベーシックインプット・アウトプットシステム(BIOS)を含み得る。アプリケーションプログラムはローカルコンピュータ可読媒体(例えばハードディスク244、光ディスク242)に格納され、システムメモリ217にロードされ、そしてプロセッサ214により実行され得る。アプリケーションプログラムは、例えばネットワークインタフェース248又はモデム247を介して、離れた場所(すなわち、離れたところに配置されたコンピュータシステム210)からシステムメモリ217にロードすることもできる。図2において、リンク障害管理システム101は、システムメモリ217内にあるものとして図示されている。リンク障害管理システム101のはたらきは、図3に関連して下記で更に詳しく説明される。

10

**【0018】**

記憶装置インタフェース234は、1つ以上のハードディスク244(及び/又はその他の標準記憶媒体)に連結されている。(複数の)ハードディスク244は、コンピュータシステム210の一部であってもよく、又は、物理的に別であってもよく、他のインタフェースシステムを介してアクセスされてもよい。

**【0019】**

ネットワークインタフェース248及び/又はモデム247は、例えばインターネットなどのネットワーク107に、直接的又は間接的に、通信可能に連結され得る。そのような連結は、有線又は無線であり得る。

20

**【0020】**

図3は、いくつかの実施形態による、一例のクラスタ301トポグラフィにおいてリンク305の障害を検出する、リンク障害管理システム101を示す。上述のように、リンク障害管理システム101の機能性はクライアント103又はサーバ105上にあるか、又は、複数のコンピュータシステム210間で分散され得る。これには、リンク障害管理システム101の機能性がネットワーク107にわたるサービスとして提供されるクラウドベースのコンピューティング環境内が含まれる。リンク障害管理システム101は図3において単独のエンティティとして示されているが、図示されているリンク障害管理システム101は、機能性の集まりを示し、これは所望に応じて単一又は複数のモジュールとして実例を示すことができる(リンク障害管理システム101の具体的な複数モジュールの実例は図3及び4に示されている)。例示目的のため、リンク障害管理システム101は、図示のクラスタ301の各ノード303上にあり、かつノード303上で起こる機能性を管理しているものとして図示されている。実際には、リンク障害管理システム101は所望に応じて、中央集中化されてよく、又はクラスタ301の複数ノード303にわたって分散されていてもよい。

30

**【0021】**

リンク障害管理システム101のモジュールは、任意のコンピュータシステム210のシステムメモリ217(例えばRAM、ROM、フラッシュメモリ)内で実例を示すことができ(例えばオブジェクトコード又は実行可能画像として)、これによって、コンピュータシステム210のプロセッサ214がモジュールを処理する際、コンピュータシステム210がその関連する機能性を実行することが理解されよう。本明細書において使用される用語「コンピュータシステム」、「コンピュータ」、「クライアント」、「クライアントコンピュータ」、「サーバ」、「サーバコンピュータ」及び「コンピューティングデバイス」は、記述される機能性を実行するよう設定及び/又はプログラムされた、1つ以上のコンピュータを意味する。加えて、リンク障害管理システム101の機能性を実装するためのプログラムコードは、コンピュータ可読記憶媒体に格納することができる。この文脈において、例えば磁気又は光学記憶媒体など、任意の形態の有形のコンピュータ可読記憶媒体を使用することができる。本明細書において使用される用語「コンピュータ可読記憶媒体」は、下位の物理的媒体とは別体の電気信号は意味しない。

40

**【0022】**

50

図3に示すように、リンク障害管理システム101は、障害のあるリンク305の検出を、クラスタ301全体に迅速に伝えることを可能にする。明確にするため、図3は3つのノード303A、303B、及び303Cを含むクラスタ301を示す。この分野において、クラスタ301は典型的に、この数桁以上の数のノード303を含み得ることが理解されよう。ノード303は、ネットワークリンク305及びハブ307を使用して接続されている。図3に示されている例において、各ノード303は、別のリンク305で、3つの異なるハブ307A、307B、及び307Cに接続されている。別の実施形態において、クラスタ301のノード303を接続するネットワークトポロジを形成するのに、配備するリンク305及びハブ307の数は、これより多くても少なくてもよい。

#### 【0023】

図3に示すように、ノード303Aにあるリンク障害管理システム101のリンク障害検出モジュール309は、ノード303Aの任意のリンク305（すなわちリンク305A～C）が障害を生じた時点を検出する。（他のノード303それぞれの類似のモジュールは、それぞれのローカルリンク305の障害を検出する）。リンク障害検出モジュール309は、ローカルリンク305がダウンしたときに即座に通知を受け取るため、ローカルノード303のオペレーティングシステムサービスを利用することができる。例えば、リンク305Aが壊れた場合、ノード303Aのリンク障害検出モジュール309は、ノード303A上のオペレーティングシステム（図示なし）から即座に通知を受け取ることができる。しかしながら、リンク障害検出モジュール309はそのローカルノード303のリンク305の障害を検出するだけであるため、ノード303B、303Cはこのリンクダウン事象に気づいていないことになる。所与のノード（例えば303A）のオペレーティングシステムは、そのノード303にとってローカルであるリンク305の障害を検出するだけであるため、このようになる。

#### 【0024】

従来、クラスタ301の他のノード303は、ハートビートを監視することにより、ノード303Aのリンクの障害を知り得る。しかしながら、上述のように、ハートビートに依存すると、離れたノード303が、離れたリンク305に障害が生じたと結論づけられるまでに、過度に長い遅れが必要となる。図示の例において、ノード303Aはリンク305Aの障害を即座に知り得るが、ノード303B及び303Cが、リンク305Aに障害が生じたと結論づけられるようになるまでには、対応するハートビートが期限切れになるのを待たなければならない。上記で説明したように、このような状況は問題である。それは、クラスタ301の管理にとって、各ノード303がリンクダウン事象を平行して処理する（すなわち、ダウンしているリンク305に対応する）ことが重要であり、また同時に、ハートビート機構に依存して促進され得るよりも迅速にこれを成し遂げるのが望ましいためである。

#### 【0025】

図4～5に示すように、ローカルノード303のリンク障害管理システム101は、クラスタ301の他のノード303に、リンク305の障害の通知を伝播し、これによりクラスタ301のノード303がリンクダウン事象を平行して処理できるようにする。下記で詳しく記述されるように、異なる実施形態において、リンク障害管理システム101は、別の方法を使用して、クラスタ301のノード303に、リンク305の障害の通知伝播を実行する。これらの異なる実施形態において、リンク305の障害の通知伝播は、対応するハートビートが期限切れになる前に、クラスタ301のノード303がその通知を受け取るように実行されることが理解される。よって、リンク障害管理システム101によって、クラスタ301のノード303が、リンクダウン事象を平行に、ハートビート機構に依存して達成され得るよりも迅速に処理することができるようになる。

#### 【0026】

図4は、リンク障害管理システム101が、リンク305の障害に関する情報を、ブロードキャストメッセージ403を介して、クラスタ301のノード303に伝播する一実施形態を示す。一実施形態において、ローカルリンク（例えば図3に示される例示的なト

10

20

30

40

50



ポロジにおけるリンク 305 A) の障害を知ったことに対応して、そのローカルノード (例えばノード 303 A) 上にあるリンク障害管理システム 101 のブロードキャストモジュール 401 が、検出されたリンク 305 の障害のことをクラスタ 301 のノード 303 に通知する、ブロードキャストメッセージ 403 を生成かつ送信する。よって、クラスタ 301 の他のノード 303 はすべて、リンク 305 の障害の通知を同時に受け取り、よってこの事象の処理を平行して行うことができる。別の実施形態において、この目的にどのリンク 305 が利用できるかに応じて、別のリンク 305 を使用して、ブロードキャストメッセージ 403 を送信することができる。

【0027】

例えば、一実施形態において、クラスタ 301 内のいくつかのリンク 305 が、標準クラスタリンク 305 として構成され (例えば図 3 に示す例示的なトポロジにおいて、ハブ 307 A を介して相互接続しているリンク 305 A、305 D、及び 305 G、並びに、ハブ 307 B を介して相互接続しているリンク 305 B、305 E、及び 305 H)、ここにおいて一組のリンク 305 は、リンクダウン事象ブロードキャストメッセージ 403 のために確保される (例えば、ハブ 307 C を介して相互接続しているリンク 305 C、305 F 及び 305 I)。このシナリオにおいて、ノード 303 A のリンク障害検出モジュール 309 が、リンク 305 A 及び/又はリンク 305 B がダウンしたことを検出すると、専用のリンク 305 C を使用して、ブロードキャストメッセージ 403 をノード 303 B 及び 303 C へと送信することができる。一実施形態において、障害検出ブロードキャストメッセージ 403 専用の別のリンク 305 は、リンク障害管理システム 101 のプライベートリンクであり、クラスタ 301 のユーザーに対して公開されない。このシナリオにおいて、専用リンク 305 は、クラスタ 301 内のすべてのノード 303 に接続され、リンクダウン事象情報を送信する際にのみ使用される。これにより、リンク 305 の障害検出を示すブロードキャストメッセージ 403 の通信のための専用チャンネルが確保される。

【0028】

別の一実施形態において、これらのリンク 305 は、リンク障害管理システム 101 のプライベートリンクであるが、リンクダウン事象情報以外のトラフィックに使用することができる。更に別の一実施形態において、リンクダウン事象情報の送信に使用されるリンク 305 は、リンク障害管理システム 101 のためのプライベートリンクではなく、例えば利用可能な帯域幅に基づいて、他の当事者により他のトラフィックのために使用することができる。いくつかの実施形態において、リンクダウン事象情報を通信するための専用リンク 305 は利用できず、この場合、他のリンク 305 をこの目的のために使用することができる (例えば、高優先度リンク 305 又は低優先度リンクを含む既存の公共リンク 305 で、例えば他のものがダウンした場合など)。別の実施形態において、利用できるリンク 305 が何であれ、所望に応じてリンクダウン事象情報を送信するのに使用できることが理解されよう。

【0029】

リンクダウン事象メッセージ 403 がブロードキャストされると、ローカルノード (例えば図 3 の 303 A) のリンク障害管理システム 101 が工程を実施し、ブロードキャストメッセージ 403 をクラスタ 301 の他のノード 303 が確実に受け取ってから、リンクダウン事象の処理が行われる。これにより、クラスタ 301 の複数のノード 303 が、リンクダウン事象を同時に確実に処理することができる。一実施形態において、リンク障害管理システム 101 の確認応答受信モジュール 405 が、クラスタ 301 の他のノード 303 のそれぞれから、ブロードキャストメッセージ 403 の受信を確認するための、確認応答 (ACK) 407 の受信を待つ。他のノード 303 のそれぞれから ACK 407 を受信するということは、クラスタ内のすべてのノード 303 がリンク 305 の障害を知っていることを示す。ブロードキャストメッセージ 403 自体、又は 1 つ以上のノード 303 からの ACK 407 が、ネットワークにより中断されている場合、確認応答受信モジュール 405 は、クラスタ 301 の他のノード 303 それぞれからの ACK 407 を

10

20

30

40

50

受信しない。一実施形態において、所与の一定時間内に、クラスタ301の各ノード303からACK 407を受信していないことに対応して、ブロードキャストモジュール401はブロードキャストメッセージ403を再送信する。いくつかの例において、ブロードキャストメッセージ403は、1つ以上の予期されるACK 407を受信していないことに対応して、複数回再送信することができる。再送信までの待ち時間の長さ、並びに再送信の回数は、実施形態間で変えることができる設計パラメータであることが理解されよう。いかなる場合でも、これらのパラメータは典型的に、確認応答受信モジュール405がACK 407を待つ合計時間が、ノード303のステータス情報のクラスタ301の全体への伝播のために、ハートビート機構によって使用される時間スケールよりも更に短くなるように、設定される。やがてクラスタの他のノード303がハートビート機構を介してリンク305の障害を知ることが理解されよう。よって、ハートビート機構を介してリンク305の障害をノード303が知るのにかかる時間内に、すべてのACK 407が受信されなかった場合のバックアップとして、クラスタ301のノード303は、リンクダウン事象を従来の方法で知ることができる。リンクダウン事象に関する情報をブロードキャストするため、及びその確認応答受信に使用されるハンドシェーキングに使用されるプロトコル及びフォーマットの具体的な実装は、所望に応じて実施形態間で変えることができることが理解されよう（例えばアトミックブロードキャスト、2相コミットなど）。

#### 【0030】

いくつかの実施形態において、リンク障害検出モジュール309がリンク305の障害を検出したとき、ブロードキャストモジュール401は猶予期間待ってから、ブロードキャストメッセージ403を送信する。この猶予期間の目的は、リンクがダウンしたけれどもほぼ即座に元に戻った場合に（例えば不安定なリンク）、リンクダウン事象のブロードキャストを避けることである。猶予期間の長さは可変の設計パラメータであるが、典型的には、ハートビート機構の時間スケールに比べ、かなり短い。例えば、0.5秒、1秒、又は2秒の猶予期間が使用され得る。そのような実施形態において、障害が生じたリンク305が猶予期間内に元に戻った場合、ブロードキャストメッセージ403は送信されない。一方、猶予期間が経過し、リンク305が依然としてダウンしている場合は、ブロードキャストモジュール401がメッセージ403を他のノード303に送信する手続きを行う。

#### 【0031】

図5は別の実施形態を示し、ここにおいてリンク障害管理システム101は、中央集中型コンピューティングデバイス503を介してクラスタ301の他のノード303に、リンク305の障害の通知501を伝播する。この実施形態において、リンク障害検出モジュール309がリンク305の障害を検出した場合、他のノード303にメッセージ403をブロードキャストする代わりに（又はそれに加えて）、リンク障害管理システム101の送信モジュール505が、中央集中型コンピューティングデバイス503（例えば、サーバ105、ディスク244など）に通知501を送信する。この実施形態において、クラスタ301の各ノード303上のリンク障害管理システム101のポーリングモジュール507は、リンク305の障害の通知501のために、特定の周波数で中央集中型コンピューティングデバイス503のポーリングを行い、その過程でリンク305の障害を知る。使用する具体的なポーリング周波数は、可変の設計パラメータである。ポーリングの代わりに又はこれに加えて、中央集中型コンピューティングデバイス503は、リンクダウン事象に関する通知501を受信すると、その通知501をクラスタ301のノード303すべてに送信することができる。これらの実施形態は、例えば、リンクダウン事象のブロードキャストメッセージ403を送信するための、所与のノード303上のブロードキャストモジュール401が利用できるクラスタリンク305がない場合に、使用することができる。

#### 【0032】

図6は、一実施形態によるリンク障害管理システム101のオペレーションの工程を示

10

20

30

40

50

す。リンク障害検出モジュール309が、クラスタ301の特定のノード303の特定のリンク305の障害を検出する601。ブロードキャストモジュール401が、検出されたリンク305の障害を、クラスタ301のノード303を通知するメッセージ403をブロードキャストし603、これによりノード303は、対応するハートビートの期限切れによりリンクダウン事象を知るより前に、メッセージ403を受信する。確認応答受信モジュール405が、クラスタ301のノード303からのブロードキャストメッセージ403の受信を確認する確認応答407を受信する605。クラスタ301の各ノード303が、リンクダウン事象を平行に処理する607。

【0033】

図7は、別の一実施形態によるリンク障害管理システム101のオペレーションの工程を示す。リンク障害検出モジュール309が、クラスタ301の特定のノード303の特定のリンク305の障害を検出する701。送信モジュール505が、リンクダウン事象に関する通知501を、中央集中型コンピューティングデバイス503に送信する703。クラスタ301のノード303のポーリングモジュール507が、リンク305の障害の通知501のため、特定の周波数で中央集中型コンピューティングデバイス503のポーリングを行い705、これによって、クラスタ301のノード303が、対応するハートビートの期限切れによって知るよりも前に、ポーリングからリンクダウン事象のを知る。クラスタ301の各ノード303が、リンクダウン事象を平行に処理する707。

【0034】

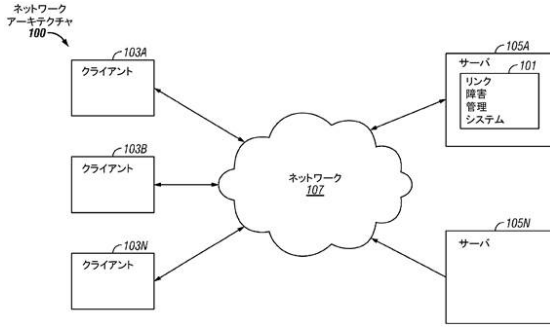
当該技術分野を知る業者には理解されるように、本発明は、その趣旨又は本質的特徴から逸脱することなく、他の具体的な形態で実施することができる。同様に、部分、モジュール、エージェント、マネージャー、構成要素、機能、手順、処置、階層、特徴、属性、方法論、データ構造、及びその他の態様の具体的な名称及び区分は、必須のものでも重要なものでもなく、本発明又はその特徴を実施するメカニズムは、別の名称、区分、及び/又は形式を有し得る。説明目的のための上記の記述は、特定の実施形態を参照して記述されている。しかしながら、上記の例示的議論は網羅的であることを意図したものではなく、開示されている正確な形態に限定することを意図したものでもない。上記の教示の見地から、数多くの改変及び変化が可能である。実施形態は、関連する原理及びその実際的な適用を最適に説明するよう選択かつ記述されており、これにより、他の当業者が、想到される特定の使用に好適であり得るような様々な改変を行うか行わないかにかかわらず、様々な実施形態を最適に利用できる。

10

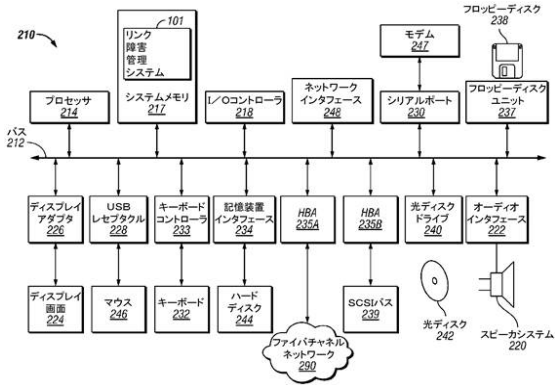
20

30

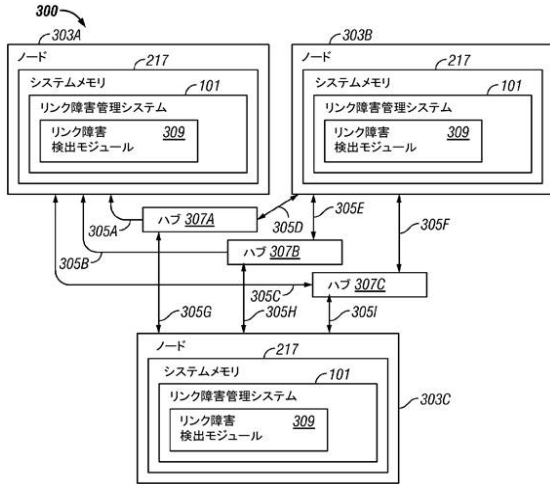
【図1】



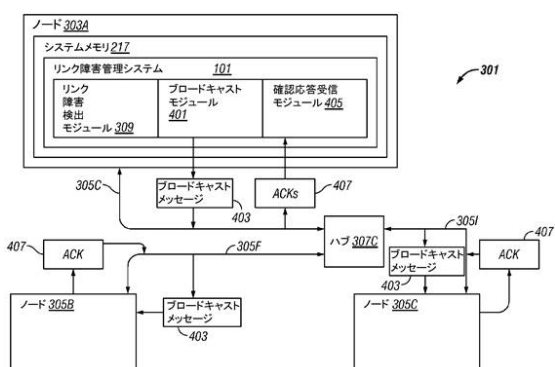
【図2】



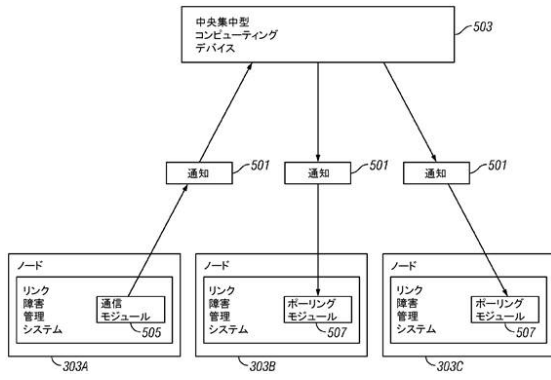
【図3】



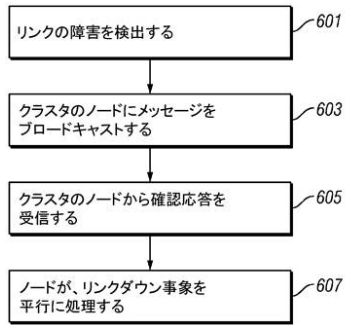
【図4】



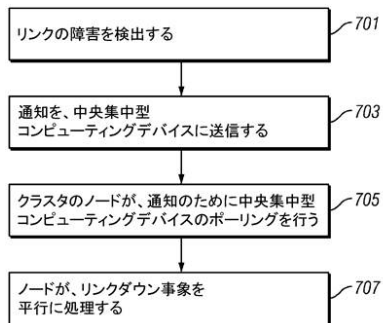
【 図 5 】



【 図 6 】



【 図 7 】



## フロントページの続き

- (72)発明者 カトカー・アモール  
インド マハラシュトラ州 411020 プネー ボポディー バウパティルロード ラビラヘ  
リテージ アウイング 303
- (72)発明者 アガーワル・オーム・プラカシュ  
インド マハラシュトラ州 411032 プネー ティンジャーナガー ロードナンバー10シ  
ー アサーヴァレジデンシー フラットナンバー8
- (72)発明者 セイカー・バービン  
アメリカ合衆国 カリフォルニア州 94087 サニーヴェール インヴァネスウェイ 756

審査官 衣嶋 文彦

- (56)参考文献 特開2005-124171(JP,A)  
特開2003-179629(JP,A)  
特表2004-519024(JP,A)  
堀田 勇樹 他, 耐故障並列計算を支援する自律的な故障検知機構, 情報処理学会論文誌, 20  
05年 8月15日, 第46巻, 第SIG12(ACS 11)号, p.236~244

- (58)調査した分野(Int.Cl., DB名)  
H04L 12/00~12/955