

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6081293号
(P6081293)

(45) 発行日 平成29年2月15日 (2017.2.15)

(24) 登録日 平成29年1月27日 (2017.1.27)

(51) Int. Cl.		F I			
G06F 13/14	(2006.01)	G06F 13/14	330A		
G06F 3/06	(2006.01)	G06F 3/06	301A		
G06F 13/10	(2006.01)	G06F 13/10	340A		
		G06F 13/10	330C		

請求項の数 20 (全 23 頁)

(21) 出願番号	特願2013-107281 (P2013-107281)	(73) 特許権者	510149482
(22) 出願日	平成25年5月21日 (2013.5.21)		グイェムウェア インコーポレイテッド
(65) 公開番号	特開2013-246821 (P2013-246821A)		VMware, Inc.
(43) 公開日	平成25年12月9日 (2013.12.9)		アメリカ合衆国 94304 カリフォル
審査請求日	平成28年5月9日 (2016.5.9)		ニア州 パロ アルト ヒルビュー アベ
(31) 優先権主張番号	13/479, 118		ニュー 3401
(32) 優先日	平成24年5月23日 (2012.5.23)	(74) 代理人	100105957
(33) 優先権主張国	米国 (US)		弁理士 恩田 誠
早期審査対象出願		(74) 代理人	100068755
			弁理士 恩田 博宣
		(74) 代理人	100142907
			弁理士 本田 淳

最終頁に続く

(54) 【発明の名称】 ファブリック分散リソースのスケジューリング

(57) 【特許請求の範囲】

【請求項1】

ストレージ・エリア・ネットワーク (SAN) における複数の仮想マシン (VM) 向けの、集中型の入出力 (I/O) 負荷ベースパス選択用システムであって、

コンピュータ装置に関連した記憶エリアであって、前記記憶エリアは前記 SAN のトポロジ記述を記憶し、前記トポロジ記述は少なくとも前記 SAN 用の閾値容量を含む、記憶エリアと、

プロセッサと、

を備え、前記プロセッサは、

前記仮想マシンを実行する複数のホストのうちの1つから、前記ホストが認識可能な少なくとも1つの論理的記憶装置にそれぞれ至る複数の記憶パスを識別することによって、初期の記憶パス選択を遂行する工程であって、前記識別には最大フロー・アルゴリズムを繰り返し用い、同最大フロー・アルゴリズムに用いる、前記ホストと前記論理的記憶装置との間の各エッジのウェイトを、前記識別を繰り返す毎に増大させる、工程と

前記遂行する工程において識別された前記記憶パスにおける複数の記憶パスが形成するリンクに基づき、識別された前記記憶パスのそれぞれについて I/O 負荷を求める工程と、

前記記憶パスのそれぞれについての求められた前記 I/O 負荷を前記閾値容量から減じて、識別された前記記憶パスのそれぞれに対して、それぞれの識別された前記記憶パスに適合する最小限の帯域幅を見いだす工程と、

前記最小限の帯域幅に基づき、前記複数のホストのそれぞれに対して、識別された前記記憶パスのうちの1つを選択する工程と、
を行うようにプログラムされる、システム。

【請求項2】

前記プロセッサはさらに、
前記SANのトポロジの変更を検出し、
検出された前記変更に応じて、前記初期の記憶パス選択を再度遂行するようにプログラムされる、
請求項1に記載のシステム。

【請求項3】

前記プロセッサはさらに、求められた前記I/O負荷に基づき前記SANのトポロジの変更を推奨するようにプログラムされる、
請求項1に記載のシステム。

【請求項4】

前記プロセッサはさらに、前記ホストのうち少なくとも1つに対して前記複数の記憶パスから別の記憶パスを選択することと、別のコントローラに対して前記論理的記憶装置を再マッピングするように、記憶プロセッサ・アレイに命じることと、前記論理的記憶装置から別の論理的記憶装置へとデータを移動することと、前記複数のホストのうちの第1のホストから前記複数のホストのうちの第2のホストへとデータストアを移動することと、のうち1または複数を前記SANの管理者からの入力なしで遂行することによって、前記SANの前記トポロジに対する推奨された前記変更を開始するようにプログラムされる、
請求項3に記載のシステム。

【請求項5】

前記記憶エリアはさらに、前記複数の記憶パスが形成するリンクに関連したリンク容量を記憶する、
請求項1に記載のシステム。

【請求項6】

前記プロセッサは、前記最小限の帯域幅および前記記憶エリアに記憶された前記リンク容量に基づき、識別された前記記憶パスのうち前記1つを選択するようにプログラムされる、
請求項5に記載のシステム。

【請求項7】

前記システムはさらに、
前記記憶パスの輻輳を改善するためにトポロジ変更を識別する手段と；
前記記憶パスの輻輳を改善するためにトポロジ変更を自動的に開始する手段とを備える、
請求項1に記載のシステム。

【請求項8】

分散型共有リソース記憶システムにおける集中型の入出力(I/O)負荷ベースパス選択の方法であって、

前記分散型共有リソース記憶システムにおいて動作するプロセッサが、複数のホストのうちの1つから、前記ホストが認識可能な少なくとも1つの論理的記憶装置にそれぞれ至る複数の記憶パスを検索する工程であって、前記検索には最大フロー・アルゴリズムを繰り返し用い、同最大フロー・アルゴリズムに用いる前記ホストと前記論理的記憶装置との間の各エッジのウェイトを、前記検索を繰り返す毎に増大させる、工程と

前記検索する工程において検索された前記記憶パスが形成する複数のリンクに基づき、検索された前記記憶パスのそれぞれについてI/O負荷を求める工程と；

前記記憶パスのそれぞれについての求められた前記I/O負荷を閾値容量から減じて、前記記憶パスのそれぞれに対して、それぞれの前記記憶パスについての最小限の帯域幅を見いだす工程と；

10

20

30

40

50

前記複数のホストのそれぞれに対して、前記記憶パスのうちの1つを前記最小限の帯域幅に基づき選択する記憶パス選択工程とを有する、方法。

【請求項9】

求められた前記I/O負荷に基づき、前記分散型共有リソース記憶システムにおけるトポロジ変更を選択する工程を含む、

請求項8に記載の方法。

【請求項10】

選択された前記トポロジ変更を開始する工程を含む、

請求項9に記載の方法。

10

【請求項11】

変更を検出するために時間を通じて前記I/O負荷を監視する工程を含む、

請求項8に記載の方法。

【請求項12】

前記記憶パスのうちの1つに関して、監視された前記I/O負荷における変更を検出すると、

前記記憶パスのうちの前記1つに対する更新されたI/O負荷を求める工程と、

前記記憶パスのそれぞれに対して、それぞれの前記記憶パスに適合する最小限の帯域幅を算定する工程と、

前記複数のホストのそれぞれに対して、前記記憶パスのうちの1つを算定された前記最小限の帯域幅に基づき選択する工程と、を含む、

20

請求項11に記載の方法。

【請求項13】

前記記憶パス選択工程はさらに、前記最小限の帯域幅に基づきバイナリ・サーチを遂行する工程を含む、

請求項8に記載の方法。

【請求項14】

対応する前記ホストに対して選択された記憶パスのそれぞれを識別する工程を含む、

請求項8に記載の方法。

30

【請求項15】

前記複数のホストのうちの1つ、ホスト・バス・アダプタ、アレイ・ポート、および論理装置番号(LUN)を識別することによって、前記複数の記憶パスのそれぞれを定義する工程を含む、

請求項8に記載の方法。

【請求項16】

コンピュータ実行可能命令を含む1または複数のコンピュータ可読記憶媒体であって、前記コンピュータ実行可能命令が実行されると、少なくとも1つのプロセッサは、

分散型共有リソース記憶システムにおいて動作するプロセッサが、複数のホストのうちの1つから、前記ホストが認識可能な少なくとも1つの論理的記憶装置にそれぞれ至る複数の記憶パスを検索する工程であって、前記検索には最大フロー・アルゴリズムを繰り返し用い、同最大フロー・アルゴリズムに用いる前記ホストと前記論理的記憶装置との間の各エッジのウェイトを、前記検索を繰り返す毎に増大させる、工程と、

40

検索された前記記憶パスの少なくとも1つにおける少なくとも1つの記憶パスに対して上限リンク容量を定義する工程と、

前記記憶パスのそれぞれについてI/O負荷を求める工程と、

前記記憶パスのそれぞれについての求められた前記I/O負荷を閾値容量から減じて、前記記憶パスのそれぞれに対して、それぞれの前記記憶パスおよび定義された前記上限リンク容量に適合する最小限の帯域幅を見いだす工程と、

前記複数のホストのそれぞれに対して、前記記憶パスのうちの1つを最小限の帯域幅に

50

基づき選択する工程と、
を行うことによって集中型の入出力（I/O）負荷ベースパス選択を遂行する、コンピュータ可読記憶媒体。

【請求項 17】

前記複数のホストは 1 または複数の仮想マシン（VM）を実行する、請求項 8 に記載の方法。

【請求項 18】

前記コンピュータ実行可能命令によって、前記プロセッサは、前記記憶パスにおけるリンクの過剰予約を防止するように前記上限リンク容量を定義する、

請求項 16 に記載のコンピュータ可読記憶媒体。

10

【請求項 19】

前記コンピュータ実行可能命令によって、前記プロセッサは、前記記憶パスの内部のスイッチに問い合わせる前記スイッチの最大帯域幅容量を得ることによって、前記上限リンク容量を定義する、

請求項 16 に記載のコンピュータ可読記憶媒体。

【請求項 20】

前記コンピュータ実行可能命令によって、前記プロセッサはさらに、前記最小限の帯域幅および前記スイッチの前記最大帯域幅容量に基づき、前記複数のホストのそれぞれに対して前記記憶パスのうち 1 つを選択する、

請求項 19 に記載のコンピュータ可読記憶媒体。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明概念は、一般にデータ記憶装置の分野に関する。より詳細には、本開示は、共有記憶システムおよび関連方法に関する。

【背景技術】

【0002】

共有記憶システムでは、複数のホストは、記憶装置の同一の組および/または記憶入出力（I/O）経路の同一の組を共有してよい。共有記憶システムの実例の 1 つには、多くの仮想マシン（VM: virtual machine）を実行する仮想データセンタがある。或る現存システムでは、ファブリック輻輳（congestion）を防止するのが困難である。この困難は、部分的には、ファブリック・トポロジ、I/O 負荷、および I/O 経路選択の複雑さによる場合がある。また、ファブリック・トポロジ、I/O 負荷、および I/O 経路選択の変更による困難が存在することもある。例えば、仮想マシン、仮想ディスク、ホスト、記憶装置（または他の記憶領域）、およびファイバ・チャンネル・リンクが、動的に付加されたり除去されたりする可能性がある。

30

【0003】

固体ディスク（SSD: solid-state disk）などの先進のディスク技術は、他のタイプのディスクよりも優れたランダム I/O 性能をもたらす。SSD は、キャッシュ記憶装置として、フロントエンド段として、および/または完全なスピンドル交換として使用されている。SSD を用いると、3 ギガバイト/秒ものランダム・ディスクの I/O トラフィックを達成することが可能になり得て、記憶装置または論理装置番号（LUN: logical unit number）当たりの I/O 帯域幅の向上をもたらす。さらに、バックアップ、クローニング、およびテンプレート配置など、高スループットの連続した入出力操作によって、ファブリック・リンクの飽和状態および/または故障が生じる恐れがある。

40

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】米国特許出願公開第 2012/0054329 号明細書

50

【特許文献2】米国特許出願公開第2010/0083262号明細書

【特許文献3】米国特許出願公開第2011/0119413号明細書

【特許文献4】米国特許出願公開第2011/0072208号明細書

【非特許文献】

【0005】

【非特許文献1】GULATI et al., "PARDA: Proportional Allocation of Resources for Distributed Storage Access", 7th USENIX Conference on File and Storage Technologies, 2009, 14 pages.

10

【非特許文献2】UNKNOWN, "VMware VROOM! Blog: 350,000 I/O operations per Second, One vSphere Host", Retrieved from <<http://blogs.vmware.com/performance/2009/05/350000-io-operations-per-second-one-vsphere-host-with-30-efds.html>>, May 18, 2009, 2 pages.

【非特許文献3】UNKNOWN, "Veritas Storage Foundation 5.0 Dynamic Multi-Pathing", Retrieved from <<http://eval.symantec.com/mktginfo/enterprise/white_papers/ent-whitepaper_vs_f_5.0_dynamic_multi_pathing_05-2007.en-us.pdf>>, May 2007, 41 pages.

20

【発明の概要】

【発明が解決しようとする課題】

【0006】

個々のホスト向け最適パスを手動で求めることによって動的事象にตอบสนองすることは、困難であり、信頼性が低く、エラーが発生しやすく、効果的な負荷バランシングをもたらす可能性が低い。さらに、現存システムには、マルチパス化、入出力制限、LUNパス選択技法の遂行、1つのホストから別のホストへの作業負荷の移動、または1つのLUNもしくはデータストアから別のLUNへとデータを移動することによって、負荷バランスを試みるものもある。しかし、このような現存システムは、LUN輻輳とリンク輻輳を区別しない。そのため、これらの現存システムは、現行パスが輻輳しているとき、LUNにアクセスするための代替パスを提案することができない。さらに、現存システムの多くはホスト・レベルで動作するものであり、したがって、全体最適をもたらすこと、あるいはボトルネックを改善するためのトポロジ変更または代替パスを推奨することができない。

30

【0007】

当技術のこれらの困難または他の困難に対処するシステムおよび方法の提供が望まれており、これらのシステムおよび方法は、当業者に、本明細書の教示および議論を考慮した後理解されるであろう。

【課題を解決するための手段】

40

【0008】

本発明によれば、添付の特許請求の範囲で説明される装置および方法が提供される。本発明の他の特徴は、従属請求項およびそれに続く記述から明らかになるであろう。

この概要は、以下で、より詳細に説明される概念の選択を紹介する。この概要は、基本的な特徴を列挙したり、要求された対象の範囲を限定したりするようには意図されていない。

【0009】

1例では、共有ストレージ・エリア・ネットワーク(SAN: storage area network)における、集中型の全体最適化された入出力(I/O)負荷ベースパス管理を有するシステムが提供される。複数のホストのうちの1つから、同ホストが認

50

識可能な少なくとも1つの論理的記憶装置にそれぞれ至る複数の記憶パスが識別される。少なくとも記憶パス上のI/O負荷および反復検索技法に基づき、それぞれのリンクまたは記憶パスに対する最小限の帯域幅が求められる。ホストと論理的記憶装置の各対に対して、記憶パスのうちの1つが最小限の帯域幅に基づき選択される。

【0010】

本明細書で説明される例示的システムは、記憶パスに沿って、かつ/または論理的記憶装置において、輻輳の確率を低減するために、SANトポロジおよびI/O負荷を監視して、トポロジの変更を識別し、変更を開始する。或る实例では、トポロジの変更には、ホスト上で動作している複数の仮想マシン(VM)と関連した移動する負荷および/またはデータストアが含まれる。

10

【0011】

1例では、エリア・ネットワーク(SAN)における複数の仮想マシン(VM)向けの、集中型の入出力(I/O)負荷ベースパス選択用システムが提供され、前記システムは、少なくともSAN用の閾値容量を含む同SANのトポロジ記述を記憶する、コンピュータ装置に関連した記憶エリアと、VMを実行する複数のホストのうちの1つから、同ホストが認識可能な少なくとも1つの論理的記憶装置にそれぞれ至る複数の記憶パスを識別することによって、初期のパス選択を遂行し、識別された記憶パスのそれぞれについてI/O負荷を求め、記憶パスのそれぞれについて求められたI/O負荷を閾値容量から減じて、識別された記憶パスのそれぞれに対して、それぞれの識別された記憶パスに適合する最小限の帯域幅を見だし、複数のホストのそれぞれに対して、識別された記憶パスのうち1つを最小限の帯域幅に基づき選択するようにプログラムされたプロセッサと、からなる。

20

【0012】

1例では、このプロセッサは、SANのトポロジの変更を検出し、検出された変更に応じて、初期のパス選択を再度遂行するようにさらにプログラムされる。

1例では、このプロセッサは、求められたI/O負荷に基づきSANのトポロジ変更を推奨するようにさらにプログラムされる。

【0013】

1例では、このプロセッサは、ホストのうち少なくとも1つに対して複数の記憶パスから別のパスを選択する工程と、様々なコントローラに対して論理的記憶装置を再マッピングするように、記憶プロセッサ・アレイに命じる工程と、この論理的記憶装置から別の論理的記憶装置へとデータを移動する工程と、複数のホストの中の第1のホストから複数のホストの中の第2のホストへとデータストアを移動する工程と、のうち1つまたは複数を経由してSANの管理者からの入力なしで遂行することによって、SANのトポロジに対する推奨された変更を開始するようにさらにプログラムされる。

30

【0014】

1例では、複数の記憶パスのそれぞれが複数のリンクからなり、記憶エリアは、複数のリンクに関連したリンク容量をさらに記憶する。

1例では、このプロセッサは、最小限の帯域幅および記憶エリアに記憶されたリンク容量に基づき、識別された記憶パスのうち1つを選択するようにプログラムされる。

40

【0015】

1例では、このシステムは、記憶パスの輻輳を修正するためにトポロジ変更を識別する手段と、記憶パスの輻輳を修正するためにトポロジ変更を自動的に開始する手段とを含む。

【0016】

1例では、分散型共有リソース記憶システムにおける集中型入出力(I/O)負荷ベースパス選択の方法が提供され、前記方法は、複数のホストのうちの1つから、同ホストが認識可能な少なくとも1つの論理的記憶装置にそれぞれ至る複数の記憶パスに、該分散型共有記憶システムにおいて動作するプロセッサでアクセスする工程と、アクセスされた記憶パスのそれぞれについてI/O負荷を求める工程と、記憶パスのそれぞれについての求

50

められた I / O 負荷を閾値容量から減じて、アクセスされた記憶パスのそれぞれに対して、それぞれのアクセスされた記憶パスについての最小限の帯域幅を見いだす工程と、複数のホストのそれぞれに対して、アクセスされた記憶パスのうちの 1 つを最小限の帯域幅に基づき選択する工程と、からなる。

【 0 0 1 7 】

1 例では、この方法は、求められた I / O 負荷に基づき、分散型共有リソース記憶システムにおけるトポロジ変更を選択する工程を含む。

1 例では、この方法は、選択されたトポロジ変更を開始する工程を含む。

【 0 0 1 8 】

1 例では、この方法は、変更を検出するために長時間にわたって I / O 負荷を監視する工程を含む。

1 例では、この方法は、記憶パスのうちの 1 つに関して、監視された I / O 負荷における変更を検出すると、記憶パスのうちの 1 つに対する更新された I / O 負荷を求める工程と、記憶パスのそれぞれに対して、それぞれの記憶パスに適合する最小限の帯域幅を算定する工程と、複数のホストのそれぞれに対して、記憶パスのうちの 1 つを、算定された最小限の帯域幅に基づき選択する工程と、を含む。

【 0 0 1 9 】

1 例では、複数のホストのそれぞれに対して記憶パスのうちの 1 つを選択する工程は、最小限の帯域幅に基づきバイナリ・サーチを遂行する工程をさらに含む。

1 例では、この方法は、各繰返しに対してエッジ・ウェイトを増す最大フロー・アルゴリズムを繰り返し遂行することによって、複数の記憶パスを検索する工程を含む。

【 0 0 2 0 】

1 例では、この方法は、対応するホストに対して選択されたパスのそれぞれを識別する工程を含む。

1 例では、この方法は、複数のホスト、ホスト・バス・アダプタ、アレイ・ポート、および論理装置番号 (L U N) のうちの 1 つを識別することによって、複数の記憶パスのそれぞれを定義する工程を含む。

【 0 0 2 1 】

1 例では、コンピュータ実行可能命令を含むコンピュータ可読記憶媒体が提供され、これらのコンピュータ実行可能命令が実行されると、少なくとも 1 つのプロセッサが、複数のホストのうちの 1 つから、同ホストが認識可能な少なくとも 1 つの論理的記憶装置にそれぞれ至る複数の記憶パスに、分散型共有記憶システムにおいて動作するプロセッサでアクセスする工程と、アクセスされた記憶パスの少なくとも 1 つにおける少なくとも 1 つのリンクに対して上限リンク容量を定義する工程と、アクセスされた記憶パスのそれぞれに対する I / O 負荷を求める工程と、記憶パスのそれぞれについての求められた I / O 負荷を閾値容量から減じて、記憶パスのそれぞれに対して、それぞれの記憶パスおよび定義された上限リンク容量に適合する最小限の帯域幅を見いだす工程と、複数のホストのそれぞれに対して、記憶パスのうちの 1 つを、算定された最小限の帯域幅に基づき選択する工程と、によって集中型の入出力 (I / O) 負荷ベースパス選択を遂行する。

【 0 0 2 2 】

1 例では、これらの工程は、複数の記憶パスの中のリンクの過剰な予約を防止するために上限リンク容量を定義する工程を含む。

1 例では、これらの工程は、記憶パスの内部のスイッチに問い合わせることでスイッチの最大帯域幅容量を得ることによって、上限リンク容量を定義する工程を含む。

【 0 0 2 3 】

1 例では、これらの工程は、最小限の帯域幅およびスイッチの最大帯域幅容量に基づき、複数のホストのそれぞれに対して記憶パスのうちの 1 つを選択する工程を含む。

1 例では、ストレージ・エリア・ネットワーク (S A N) における複数の仮想マシン (V M) 向けの、集中型の入出力 (I / O) 負荷ベースパス管理の方法が提供され、前記方法は、V M を実行する複数のホストのうちの 1 つから、同ホストが認識可能な少なくとも

10

20

30

40

50

1つの論理的記憶装置にそれぞれ至る複数の記憶パスに、プロセッサでアクセスする工程と、アクセスされた記憶パスのそれぞれについてI/O負荷を監視する工程と、監視されたI/O負荷が記憶パスの少なくとも1つにおける輻輳を表しているとして判断すると、1または複数のトポロジ変更を識別する工程と、識別された1または複数のトポロジ変更を、プロセッサによって開始する工程と、からなる。

【0024】

1例では、識別された1または複数のトポロジ変更を開始する工程は、SANの管理者からの入力なしでトポロジ変更を開始する工程を含む。

1例では、識別された1または複数のトポロジ変更を開始する工程は、複数の記憶パスから、ホストのうち少なくとも1つ向けに別のパスを選択する工程と、ホストのうち少なくとも1つに対して選択された別のパスを用いるように指示する工程とをさらに含む。

10

【0025】

1例では、識別された1または複数のトポロジ変更を開始する工程は、論理的記憶装置を管理している記憶プロセッサ・アレイに対して、論理的記憶装置を別のコントローラに再マッピングするように命令を送る工程を含む。

【0026】

1例では、識別された1または複数のトポロジ変更を開始する工程は、この論理的記憶装置から別の論理的記憶装置へとデータを移動する工程と、複数のホストの中の第1のホストから複数のホストの中の第2のホストへとデータストアを移動する工程と、のうち1または複数を含む。

20

【図面の簡単な説明】

【0027】

【図1】例示のホスト・コンピュータ装置のブロック図。

【図2】図1に示されたホスト・コンピュータ装置などのコンピュータ装置上にインスタンスが生成される仮想マシンのブロック図。

【図3】ネットワーク・スイッチによって記憶パスを介してディスク・アレイにアクセスするホストを含む例示のストレージ・エリア・ネットワーク(SAN)構成のブロック図。

【図4】ネットワーク・フロー解析に対する最適の記憶パス選択の例示のマッピングのブロック図。

30

【図5】集中型の入出力(I/O)記憶パスの選択および管理を実施するための例示のコンピュータ装置のブロック図。

【図6】SANの中の最適の記憶パスを選択するように、複数のホストを管理する集中型システムによって遂行される例示の方法の流れ図。

【図7】リンクまたは論理装置番号(LUN)の輻輳にตอบสนองしてトポロジ変更を開始するように、複数のホストを管理する集中型システムによって遂行される例示の方法の流れ図。

【発明を実施するための形態】

【0028】

対応する参照文字は、各図面を通じて該当部分を示す。

40

本明細書で説明される実施形態は、共有ストレージアレイを有し、SAN 512などのストレージ・エリア・ネットワーク(SAN)における複数のホスト(例えばクラスター)にわたって集中型の入出力(I/O)パス選択および管理を実施する。パス推奨は、ファブリック(例えばリンク)のI/O負荷状態に基づくものであり、例えばネットワーク・フロー・アルゴリズム(例えば最大フロー・アルゴリズム)の最上位で条件付き反復検索を用いて選択される。結果として、本開示の態様は、ホスト514からアレイ302までのパス(経路)上のどの単一ポイントにも過負荷をかけず、したがって全体最適をもたらすように、各ホストとLUNの対に対してI/Oパスを選択する。

【0029】

或る実施形態では、初期のパス選択推奨は、ディスク・アレイ302などのストレージ

50

アレイにおけるホスト・バス・アダプタ (HBA: host bus adapter) およびポート304にわたってパスの数のバランスをとるように、ホスト514ごとに記憶装置当たりの均一なI/O負荷を想定して、本明細書に説明される動作を遂行することによって求められる。ホスト514ごとに記憶装置当たりのI/O帯域幅消費(例えば負荷または要求)が長時間にわたって監視されながら、パス選択推奨は、再構成される。SAN 512トポロジが変更されるときも、パス選択推奨が再構成されてよい。或る実施形態では、リンク当たりの帯域幅制限が(既知のとき)強制される。

【0030】

本開示の態様は、ファブリックにおけるI/O輻輳の可能性を低減し、トポロジ変更および動的I/O状態に適合する、または対応する。集中的管理は、所与のトポロジに対して、各ホスト514上で認識可能な個々のLUN用のI/Oパスを、負荷バランスをとるやり方で計算し、トポロジ変更または他の改善を開始して、トポロジの変更およびリンク輻輳の他の徴候につながる事象に対処する。例えば、集中型のパス選択は、記憶プロセッサ306または他の記憶コントローラにわたってLUNまたはLUN負荷の不均衡などの輻輳シナリオを識別するばかりでなく、特定のホスト514用パスの変更、LUN侵入(例えばLUNとコントローラのマッピングの変更)の開始、および/またはI/O負荷の移動などの改善も提案する。

10

【0031】

或る実施形態は、共有ストレージアレイを有する多くの仮想マシン(VM)を実行する仮想データセンタにおいて実施される。共有ストレージアレイでは、複数のLUNは、複数の別々のホスト514によって同一の記憶パスを介してアクセスされる可能性がある。そのため、合計のスループットが単一パス上で利用可能なリンクの帯域幅を超えることがあり、ボトルネックが生じる。1例として、LUNが、8Gb/秒のファイバ・チャネルのリンク速度によって250MB/秒の能力がある場合、1000MB/秒の理論上の最大帯域幅は、4つのLUN上の同時の連続したI/Oによって飽和状態になる可能性がある。このような同時の連続したI/O動作は、例えば、バックアップ、クローニング、およびプレート配置の間に起こる可能性がある。本開示の態様により、パス選択は、リンク上で、またLUNにおいて、このような輻輳を低減することができる。

20

【0032】

さらに、集中型パス選択によって、ホスト514からアレイへのマルチホップパス上のすべてのリンクが、その最大のリンク容量で動作することが可能になり、ストレージアレイ・ポート304および/またはファブリックのリンクのいくつかは、ホスト514上のHBAよりも低いリンク速度で動作するときでさえ、最大のエンドツーエンド帯域幅が使用可能になる。例えば、HBAが8ギガビット/秒で動作する一方で、ストレージアレイ・ポート304は、4ギガビット/秒で動作しうる。

30

【0033】

さらに、本開示の態様は、複数のパスが単一の記憶ポート304に集中する可能性がある、かつ/または不均衡なLUN-コントローラのマッピング(例えば「インキャスト」問題)がある環境で動作可能である。例えば、複数のホスト514が、ディスク・アレイ302上のより少量の記憶ポート304に接続していることがある。本開示の態様は、このようなシステムの性能を改善する。

40

【0034】

次に、1または複数の仮想マシンを実行するホスト・コンピュータ装置100を含む例示の動作環境が説明される。しかし、本開示の態様は、このような環境に限定されることなく、非VM環境にも適用可能である。

【0035】

図1は、例示のホスト・コンピュータ装置100のブロック図である。ホスト・コンピュータ装置100は、命令を実行するためのプロセッサ102を含み、ホスト514のうちの一つと称されてよい。或る実施形態では、実行可能命令は、記憶装置104に記憶される。記憶装置104は、実行可能命令および/または他のデータなどの情報が記憶され

50

、かつ取り出され得る任意の装置である。例えば、記憶装置 104 は、1 または複数のランダム・アクセス・メモリ (RAM) モジュール、フラッシュ・メモリ・モジュール、ハード・ディスク、固体ディスク、および/または光ディスクを含んでよい。

【0036】

ホスト・コンピュータ装置 100 は、ユーザ 108 からデータを受け取るため、および/またはユーザ 108 にデータを表示するためのユーザ・インターフェース装置 110 を含んでよい。ユーザ 108 は、VMware の vCenter Server または他の管理装置など別のコンピュータ装置を介してホスト・コンピュータ装置 100 と間接的に対話してよい (interact)。ユーザ・インターフェース装置 110 は、例えばキーボード、ポインティング・デバイス、マウス、スタイラス、タッチ・パネル (例えばタ
10
ッチ・パッドまたはタッチ・スクリーン)、ジャイロスコープ、加速度計、位置検出器、および/またはオーディオ入力装置を含んでよい。或る実施形態では、ユーザ・インターフェース装置 110 は、ユーザ 108 からデータを受け取るように動作し、別の装置 (例えばプレゼンテーション装置) は、ユーザ 108 にデータを見せるように動作する。他の実施形態では、ユーザ・インターフェース装置 110 は、ユーザ 108 へのデータ出力お
20
よびユーザ 108 からのデータ受取りのどちらにも機能するタッチ・スクリーンなどの単一の構成要素を有する。このような実施形態では、ユーザ・インターフェース装置 110 は、ユーザ 108 に情報を表示するためのプレゼンテーション装置として動作する。このような実施形態では、ユーザ・インターフェース装置 110 は、ユーザ 108 に情報を伝
20
えることができる任意の構成要素を表す。例えば、ユーザ・インターフェース装置 110 は、限定されることなく、表示装置 (例えば液晶ディスプレイ (LCD)、有機発光ダイ
20
オード (OLED) ディスプレイ、または「電子インク」ディスプレイ) および/または音声出力装置 (例えばスピーカまたはヘッドホン) を含んでよい。或る実施形態では、ユーザ・インターフェース装置 110 は、ビデオ・アダプタおよび/またはオーディオ・ア
20
ダプタなどの出力アダプタを含む。出力アダプタは、プロセッサ 102 に対して動作可能に結合され、表示装置または音声出力装置などの出力装置に対して動作可能に結合されるように構成される。

【0037】

ホスト・コンピュータ装置 100 は、ネットワーク通信インターフェース 112 も含み、これによって、有線または無線の packets ネットワークなどの通信媒体を介して遠隔
30
装置 (例えば別のコンピュータ装置) と通信することができる。例えば、ホスト・コンピュータ装置 100 は、ネットワーク通信インターフェース 112 を介してデータを送受してよい。ユーザ・インターフェース装置 110 および/またはネットワーク通信インター
30
フェース 112 は、まとめて入力インターフェースと称されてよく、ユーザ 108 から情報を受け取るように構成されてよい。

【0038】

ホスト・コンピュータ装置 100 はさらに、記憶インターフェース 116 を含み、これによって、仮想ディスク・イメージ、ソフトウェア・アプリケーション、および/または本明細書に説明された方法とともに用いるのに適切なその他のデータを記憶する 1 または
40
複数のデータストア 316 と通信することができる。例示の実施形態では、記憶インターフェース 116 は、ホスト・コンピュータ装置 100 を、(例えば packets ネットワークを介して) SAN 512 などのストレージ・エリア・ネットワーク (例えばファイバ・
40
チャンネル・ネットワーク) におよび/またはネットワークに接続された記憶 (NAS: network-attached storage) システムに結合する。記憶インターフェース 116 は、ネットワーク通信インターフェース 112 と一体化されてよい。

【0039】

図 2 は、ホスト・コンピュータ装置 100 上にインスタンスが生成される仮想マシン 235₁、235₂、...、235_N のブロック図を示す。ホスト・コンピュータ装置 100 は、x86 アーキテクチャのプラットフォームなどのハードウェア・プラットフォーム 205 を含む。ハードウェア・プラットフォーム 205 は、プロセッサ 102、記憶装
50

置 1 0 4、ネットワーク通信インターフェース 1 1 2、ユーザ・インターフェース装置 1 1 0、およびプレゼンテーション装置 1 0 6 (図 1 に示されている) など他の入出力 (I / O) 装置を含んでよい。以下でハイパーバイザ 2 1 0 と称される仮想化ソフトウェア層が、ハードウェア・プラットフォーム 2 0 5 の最上位にインストールされる。

【 0 0 4 0 】

仮想化ソフトウェア層は、複数の仮想マシン (V M 2 3 5 ₁ ~ 2 3 5 _N) が同時にインスタンスを生成され、かつ実行され得る仮想マシン実行空間 2 3 0 に対応する。ハイパーバイザ 2 1 0 は、装置ドライバ層 2 1 5 を含む。ハイパーバイザ 2 1 0 は、ハードウェア・プラットフォーム 2 0 5 の物理的リソース (例えばプロセッサ 1 0 2、記憶装置 1 0 4、ネットワーク通信インターフェース 1 1 2、および / またはユーザ・インターフェース装置 1 1 0) を、「仮想」リソース V M 2 3 5 ₁ ~ 2 3 5 _N のそれぞれに対して、マッピングする。それによって、V M 2 3 5 ₁ ~ 2 3 5 _N のそれぞれが、それ自体の仮想ハードウェア・プラットフォーム (例えば仮想ハードウェア・プラットフォーム 2 4 0 ₁ ~ 2 4 0 _N の対応するもの) を有し、各仮想ハードウェア・プラットフォームが、それ自体のエミュレートされたハードウェア (V M 2 3 5 ₁ のプロセッサ 2 4 5、記憶装置 2 5 0、ネットワーク通信インターフェース 2 5 5、ユーザ・インターフェース装置 2 6 0 および他のエミュレートされた入出力デバイスなど) を有する。ハイパーバイザ 2 1 0 は、「V M 2 3 5 ₁ ~ 2 3 5 _N は、ハイパーバイザ 2 1 0 の予想外の終了および / または初期設定に際して自動的に再起動されること」などのハイパーバイザ 2 1 0 に関する方針に従って、V M 2 3 5 ₁ ~ 2 3 5 _N の実行を管理してよい (例えば、監視、開始、および / または終了してよい) 。それに加えて、またはその代わりに、ハイパーバイザ 2 1 0 は、ホスト・コンピュータ装置 1 0 0 とは別の装置から受け取った要求に基づき V M 2 3 5 ₁ ~ 2 3 5 _N の実行を管理してよい。例えば、ハイパーバイザ 2 1 0 は、管理装置から、ネットワーク通信インターフェース 1 1 2 を介して第 1 の V M 2 3 5 ₁ の実行開始を指示する実行指示を受け取り、これを実行して第 1 の V M 2 3 5 ₁ の実行を開始してよい。

【 0 0 4 1 】

或る実施形態では、第 1 の仮想ハードウェア・プラットフォーム 2 4 0 ₁ の記憶装置 2 5 0 は、ホスト・コンピュータ装置 1 0 0 の記憶装置 1 0 4 (例えばハード・ディスクまたは固体ディスク) に記憶された 1 または複数の仮想ディスク・イメージに関連づけられた、すなわち「マッピングされた」仮想ディスクを含む。仮想ディスク・イメージは、第 1 の V M 2 3 5 ₁ によって、単一のファイルまたはそれぞれがファイル・システムの一部を含んでいる複数のファイルの中で用いられるファイル・システム (例えばディレクトリおよびファイルの階層) を表す。それに加えて、またはその代わりに、仮想ディスク・イメージは、S A N 構成の中などの 1 または複数の遠隔コンピュータ装置 1 0 0 の記憶装置 1 0 4 に記憶されてよい。このような実施形態では、いかなる量の仮想ディスク・イメージも遠隔コンピュータ装置 1 0 0 によって記憶され得る。

【 0 0 4 2 】

例えば、装置ドライバ層 2 1 5 は、例えばホスト・コンピュータ装置 1 0 0 に接続されたローカル・エリア・ネットワーク (L A N) との間でデータを送受するネットワーク通信インターフェース 1 1 2 と相互作用する通信インターフェース・ドライバ 2 2 0 を含む。通信インターフェース・ドライバ 2 2 0 には、1 つの通信インターフェース (例えばネットワーク通信インターフェース 1 1 2) から他の通信インターフェース (例えば V M 2 3 5 ₁ ~ 2 3 5 _N の仮想通信インターフェース) が受け取る、物理的ネットワークにおけるデータ・パケットの同送をシミュレートする仮想ブリッジ 2 2 5 も含まれる。第 1 の V M 2 3 5 ₁ 用ネットワーク通信インターフェース 2 5 5 など、各 V M 2 3 5 ₁ ~ 2 3 5 _N 用の各仮想通信インターフェースには、仮想ブリッジ 2 2 5 がネットワーク通信インターフェース 1 1 2 からの着信データ・パケットの転送をシミュレートすることを可能にする、独自の仮想媒体アクセス制御 (M A C : M e d i a A c c e s s C o n t r o l) アドレスが割り当てられてよい。1 実施形態では、ネットワーク通信インターフェ

10

20

30

40

50

ース112は、(それ自体の物理的MACアドレス宛てのイーサネット(登録商標)・パケットばかりでなく)受け取ったイーサネット(登録商標)・パケットのすべてを仮想ブリッジ225に渡す「プロミキヤス・モード」に構成されたイーサネット(登録商標)・アダプタであり、仮想ブリッジ225は、同様に、これらのイーサネット(登録商標)・パケットをVM 235₁~235_Nへさらに転送することができる。この構成によって、宛先アドレスとして仮想MACアドレスを有するイーサネット(登録商標)・パケットは、このような仮想MACアドレスに対応する仮想通信インターフェースを有するホスト・コンピュータ装置100の中のVMに適切に到達することができる。

【0043】

仮想ハードウェア・プラットフォーム240₁は、任意のx86互換デスクトップ・オペレーティング・システム(例えばマイクロソフトWINDOWS(登録商標)のオペレーティング・システム、LINUX(登録商標)のオペレーティング・システム、SOLARIS(登録商標)のオペレーティング・システム、NETWARE、またはFREEBSD)が、第1のVM 235₁などのインスタンスを生成されたVM向けのアプリケーション270を実行するためのゲスト・オペレーティング・システム(OS)265としてインストールされるように、標準的なx86ハードウェア・アーキテクチャ相当品として機能してよい。仮想ハードウェア・プラットフォーム240₁~240_Nは、ハイパーバイザ210と対応するVM 235₁~235_Nとの間の動作を調整するために仮想システム・サポートを実施する仮想マシン・モニタ(VMM: virtual machine monitor)275₁~275_Nの一部であると見なされてよい。当業者は、図2の中の仮想化構成要素を説明するのに用いられた様々な用語、層、および分類は、それらの機能または本開示の趣旨もしくは範囲から逸脱することなく、別々に参照されてよいことを理解するであろう。例えば、仮想ハードウェア・プラットフォーム240₁~240_NはVMM 275₁~275_Nと分離していると考えられてもよく、VMM 275₁~275_Nはハイパーバイザ210と分離していると考えられてもよい。本開示の実施形態で用いられ得るハイパーバイザ210の実例の1つは、VMwareのESX商標のソフトウェアに構成要素として含まれており、VMware社から市販されている。

【0044】

図3は、ネットワーク・スイッチ516によって記憶パスを介してディスク・アレイ302にアクセスするホスト514を含む例示のストレージ・エリア・ネットワーク(SAN)アーキテクチャのブロック図である。一般に、I/Oパスは、ホスト・バス・アダプタ、ケーブル、SANスイッチ(例えばスイッチ516)、記憶ポート304、およびディスク・アレイ302におけるディスク・コントローラまたは記憶プロセッサ306を含んでよい。他の実施形態では、I/Oパスまたは他の記憶パスは、ホスト514、1つのホスト・バス・アダプタ、1つのアレイ・ポート304、およびLUNのうちの1つを識別することによって定義されてよい。記憶プロセッサ306当たり複数のポート304があってもよい。記憶プロセッサ306は、1または複数のLUNを構成して、これらを単一の記憶装置としてホスト514に表示することによって、アレイの中のディスクを仮想化する。ホスト514、スイッチ516、ポート304、コントローラ、およびLUNのそれぞれの独特な組合せは、様々なI/Oパスを表す。

【0045】

図3の実例では、ホストA、ホストB、およびホストCを含む3つのホスト514は、ネットワーク・スイッチ1およびネットワーク・スイッチ2を含むネットワーク・スイッチ516を介してディスク・アレイ302に接続する。或る実施形態では、ネットワーク・スイッチ516は、ファイバ・チャネル(FC)のプロトコルに適合する。例示のネットワーク・スイッチ516は、それぞれが、4Gb/秒のリンクを提供し得る。

【0046】

ホストA、B、およびCのそれぞれは、1または複数の仮想マシンを実行する。これらのホストのそれぞれは、少なくとも1つのホスト・バス・アダプタを含む。ホストAは、

10

20

30

40

50

ホスト・バス・アダプタ A 1 およびホスト・バス・アダプタ A 2 を含む。ホスト B は、ホスト・バス・アダプタ B 1 およびホスト・バス・アダプタ B 2 を含む。ホスト C は、ホスト・バス・アダプタ C 1 およびホスト・バス・アダプタ C 2 を含む。例示の H B A は、それぞれが 8 G b / 秒のリンクを提供し得る。

【 0 0 4 7 】

ディスク・アレイ 3 0 2 は、1 または複数のポート 3 0 4、1 または複数の記憶プロセッサ 3 0 6、および 1 または複数の L U N を含む。図 3 の実例では、ディスク・アレイ 3 0 2 は、3 つの L U N を管理する 1 つの記憶プロセッサ 3 0 6 をそれぞれ有する 2 つのポート 3 0 4 を含む。一方の記憶プロセッサ 3 0 6 が L U N 1、3、および 5 を管理し、他方の記憶プロセッサ 3 0 6 が L U N 0、2、および 4 を管理する。例示の L U N のそれぞれが、2 5 0 M B / 秒のリンクを提供し得る。

10

【 0 0 4 8 】

本開示の態様は、任意の I / O パス記述に対して動作可能である。1 例として、以下で式 (1) に示されるように、ホスト H_1 と L U N L_1 の間に I / O パスが定義され得る。

【 0 0 4 9 】

【 数 1 】

$$\{H_1, HBA_i, Sport_{k_1}, \dots, Sport_{k_n}, Aport_j, LUN_1\} \quad (1)$$

20

【 0 0 5 0 】

この実例では、 HBA_i はホスト 5 1 4 における各 H B A のうちの 1 つを示し、 $Sport_{k_n}$ はネットワーク・スイッチ 5 1 6 のうちの 1 つにおけるポートを示し、 $Aport_j$ はディスク・アレイ 3 0 2 におけるポート 3 0 4 と記憶プロセッサ 3 0 6 の組合せを示す。しかし、ネットワーク・スイッチ 5 1 6 に関する構造の細部が未知の場合 (例えば、マルチパス層がスイッチ・ポートを認識できない場合)、I / O パスは、次の式 (2) に示されるようにエンドポイントの観点から記述されてよい。

【 0 0 5 1 】

【 数 2 】

$$\{H_1, HBA_i, Aport_j, LUN_1\} \quad (2)$$

30

【 0 0 5 2 】

各ホスト 5 1 4 から L U N への複数のパスがあると、ディスク・アレイ 3 0 2 の全体の有用性が増す。複数の I / O パスは、ホスト 5 1 4 から、H B A、S A N スイッチ (例えばスイッチ 5 1 6)、または様々な要素を接続するケーブルを含むディスク・アレイ 3 0 2 までのあらゆる単一故障に対して防護する。例えば、連続した I / O 要求または作業負荷に関するキャッシュ・ヒットが一時的飽和をもたらした後に、L U N 性能が、I / O パスの高帯域幅によって制限されることがある。ホスト 5 1 4 は、応答して、または I / O パスのいずれかの要素が故障したとき、ホスト 5 1 4 にアクセスしているアプリケーションに対して劣化した性能または要素の故障をさらすのではなく、別の利用可能なパスに切り換える、すなわち転換する (*divert*) ことを選択してよい。

40

【 0 0 5 3 】

マルチパス化は、様々な方法で実施されてよい。例えば、アクティブ - アクティブのディスク・アレイは、任意のポート 3 0 4 上の単一の L U N に対する I / O 要求を同時に実行する。別の実例として、アクティブ - パッシブのディスク・アレイは、1 つの記憶プロセッサ 3 0 6 (アクティブ・コントローラまたは 1 次コントローラ) の 1 または複数のポート 3 0 4 の単一の L U N に対する I / O 要求を実行するが、これらの I / O 要求を、別

50

の記憶プロセッサ306（パッシブ・コントローラまたは2次コントローラ）のポート304のLUNへフェイルオーバーすることもできる。このようなフェイルオーバーの遂行は、LUN侵入（trespass）と称されることがある。このようなディスク・アレイは、I/O要求が任意の所与の時間に記憶プロセッサ306のうちの1つを通して流れていること確認するために、I/O要求を受け取ったポート304に基づきLUNに対するフェイルオーバーを起動する。

【0054】

ディスク・アレイを実施するさらに別の实例には、疑似アクティブ・アクティブがある。このようなディスク・アレイでは、LUNは、性能が不均一な（例えば様々な程度またはレベルの性能の）複数の記憶プロセッサ306を通じてアクセスされてよい。複数の記憶プロセッサ306を通るパスは、初期の装置発見段階の間に（例えばオペレーティング・システムによって）発見されてよい。発見されたパスは、別々の装置として出現し得るが、パス管理層は、パス間の共通の装置エンドポイントを識別してよく、動的なパス選択を可能にする。

10

【0055】

パス管理層は、パス選択の管理およびエラー処理の方針を実施する。例示の方針には、ラウンドロビン、最小限のキュー長さ、バランスのとれたパス、適応型パス選択、および単一アクティブパスが含まれる。ラウンドロビンの方針は、特にI/O要求によって要求されたデータがほぼ同一サイズであって、パスがHBA、ポート304、およびコントローラにわたって均一に配置されているとき、I/O要求をそれぞれのアクティブパスへ均一に配分する。

20

【0056】

最小限のキュー長さの方針は、ホスト514を要求しているホスト・バス・アダプタにおいて顕著な要求の量が最少のアクティブパスへ、各I/O要求をバス設定することによって、局所的負荷バランスを実施する。バランスのとれたパスの方針は、各I/O要求を、開始ブロック・アドレスに対応するパスに割り当てることができるように、装置のブロック・アドレス空間を、それぞれが各アクティブパスの量の領域に分割する。適応型パス選択の方針は、より高いスループットまたはより優れた性能を最近配送したパスに対してより高い優先順位を割り当てるように、サービス時間および最近のスループットに基づき周期的なパス優先順位を計算する。このような各方針は、様々なパスに対して（厳密な優先順位でなく）重みを割り当て、割り当てられた重みに比例したI/O要求を送る。このような各方針は、高度に動的な作業負荷または非常に様々な性能もしくはホップ・カウントを有するパスに対して動作可能である。単一アクティブパスの方針は、1つのパスを故障するまで用いて、次いでフェイルオーバーパスに切り換える。

30

【0057】

本開示の態様は、本明細書に説明されたものなど、既存のマルチバス技法に対して動作可能であるが、パス選択の集中的管理はさらに、I/Oパスにおける輻輳を検出して対応する。或る実施形態では、集中型のパス選択は、パス選択を、ネットワークの最大フロー解析へマッピングすることによって実施される。次に例示のネットワークの最大フロー解析を参照しながら説明するが、本開示の態様は、本明細書に説明されるように改良された任意の最大フロー解析に対して動作可能である。

40

【0058】

ネットワークの最大フロー解析は、グラフ $G = (V, E)$ 上で動作する。ここで、 V は頂点の組であり、 E はエッジと各エッジ (u, v) に対する容量 $c(u, v)$ の組である。例示のネットワークの最大フロー問題は、以下のようにエッジ (u, v) 当たりのフロー $f(u, v)$ を計算することによって、起点 s と宛先 t から最大フローを見つけるものと説明され得る。

1. すべてのエッジに関して $f(u, v) \leq c(u, v)$
2. ノード s および t を除いて、あるノードに入るフローの和は、そのノードを出るフローの和に等しい。

50

3. 起点 s から宛先 t までの全フローは、可能な最大である。

【0059】

本開示の態様は、フォード - フルカーソン (Ford - Fulker son) のアルゴリズムおよびエドモンド - カープ (Edmond - Karp) のアルゴリズムなどの、最大フロー問題を解決する、または最大フロー解析を遂行する様々な現存システムに対して動作可能である。或る実施形態では、集中型のパス選択のマッピングを可能にする複数起点、複数宛先の最大フロー問題の改良は、それぞれのホスト H_i および $LUN L_j$ に対して起点ノード $H_i L_j$ を生成することと、各 $LUN L_j$ に対して宛先ノード L_j を生成することと、すべてのノード $H_i L_j$ の前に単一の起点を付加し、エッジ上に無限の容量を有するすべての L_j のノードに単一のシンクを付加することと、各ホスト 514 および各 LUN に接続されたアレイ・ポート 304 においてホスト・パス・アダプタに対応する中間ノードを生成することを含む。このようなマッピングの 1 例が、次に図 4 で示される。

10

【0060】

図 4 は、ネットワーク・フロー解析に対する最適の記憶パス選択の例示のマッピングのブロック図である。この実例では、2つのホスト 514 ($H1$ および $H2$) がそれぞれ 2つの HBA を有し、 HBA のそれぞれ ($HBA1$ 、 $HBA2$ 、 $HBA3$ 、および $HBA4$) が 3つの LUN ($LUN1$ 、 $LUN2$ 、および $LUN3$) に接続されている。

【0061】

初期のパス選択は、すべての LUN にわたって I/O 負荷が均一であると想定し、各 HBA および記憶プロセッサ・ポート 304 を通るパスのバランスをとることによって遂行される。動作においては、起点へのエッジ (例えば起点ノード) およびシンクへのエッジ (宛先ノード) を除いたすべてのエッジ (例えばリンク) に、1の容量が割り当てられる。すべての $H_i L_j$ ノードが L_j ノードへの非ゼロパスを有するかどうか判断するために、最大フロー・アルゴリズムが遂行される。見つかったそれぞれのパスに関して、パスのすべてのエッジによって 1 に等しい容量が消費される。すべてのノードにパスがあるわけではない場合、エッジ・ウェイトが 1 だけ増加され、すべてのノードがこのようなパスを有するようになるまで最大フロー・アルゴリズムが (例えば反復して) 再実行される。すべてのノードがパスを有するようになった後、このような最大フロー・アルゴリズムの反復動作からもたらされる、すべてのホスト 514 からのすべての LUN 向けの初期パスは、各 HBA およびコントローラ・ポート 304 に対する最小ストレスを用いる。

20

30

【0062】

次に、利用可能なパスの数または量ばかりでなく、各ホスト 514 によって各 LUN のために消費される実際の I/O 帯域幅を用いることにも基づくパス割当てを可能にするために、負荷ベース I/O パス選択が遂行される。 $H_i L_j$ ノードから HBA ノードへのエッジおよびポート・ノードから LUN へのエッジの I/O 帯域幅に関して、特定の量の容量 C が定義されている。例えば、エッジ当たり C_{min} および C_{max} の 2つの容量値が、それぞれ、 HBA 当たりのホスト 514 ごとの平均的 I/O 帯域幅および最高の可能なネットワーク・スイッチ・リンク容量として初期化される。しかし、本開示の態様は、他の初期容量値を企図する。

40

【0063】

すべての $H_i L_j$ ノードが L_j ノードへの非ゼロパスを有するかどうか判断するために、最大フロー・アルゴリズムが遂行される。パスが見つかるたびに、ホスト H_i から L_j への要求が、パスの容量から減じられる。すべてのパスを満たす最小限の I/O 帯域幅または容量を見いだすために、バイナリ・サーチが遂行される。この最小限の I/O 帯域幅に対応するパスは、選択されたパスすなわちホスト 514 に推薦されるパスを表す。

【0064】

或る実施形態では、特定のパスが過剰に予約されることがある。リンクがその容量以上の要求に対応しないことを保証するために、ホスト・パス・アダプタ、ポート 304、または中間リンクの帯域幅もしくは他の容量が既知であれば、このリンク制限を、パス選択

50

アルゴリズムにおけるリンク容量の上限として用いてもよい。最大のリンク容量値は、（例えばキャッシュに入った記憶位置に対してアイドル期間中に大きなI/O要求を用いて）リンク帯域幅を発見すること、および/またはネットワーク・スイッチ516に対して情報を問い合わせることによって、見いだされ得る。特定のホスト514からLUNへの、リンク容量を超えないパスが見つからない場合には、トポロジ変更の提案および/またはリンク飽和の警報が、管理者502に送られてよい。図7を参照しながら以下でさらに説明されるように、トポロジ変更が開始されてよい。例えば、LUN侵入、特定のLUNからのデータ移動、およびLUNからポートへの再構成マッピングが、本明細書の説明のように改良された最大フロー・アルゴリズムを繰り返し実行することによって評価されてよい。

10

【0065】

図5は、SANの複数の仮想マシンなどに対して集中型の入出力（I/O）記憶パスの選択および管理を実施するための例示のコンピュータ装置504のブロック図である。コンピュータ装置504は、コンピュータ装置504と関連した動作および機能を実施するための命令を（例えばアプリケーション・プログラム、オペレーティング・システムの機能または両方として）実行する任意の装置を表す。例えば、コンピュータ装置504は、本明細書に説明されたように、分散型リソース・スケジューリングを管理するための命令を実行する。コンピュータ装置504は、任意のコンピュータ装置または処理ユニットを含んでよい。例えば、コンピュータ装置504は、1群の処理ユニットまたはクラウド・コンピューティング構成など他のコンピュータ装置を表してよい。或る実施形態では、管理者502、またはユーザ108などの他のユーザは、コンピュータ装置504の動作の態様を管理するためにコンピュータ装置504と対話する。

20

【0066】

コンピュータ装置504は、少なくとも1つのプロセッサ506および記憶エリア508を有する。プロセッサ506は、何らかの量の処理ユニットを含み、本開示の態様を実施するためのコンピュータ実行可能命令を実行するようにプログラムされる。これらの命令は、プロセッサ506もしくはコンピュータ装置504の内部で動作する複数のプロセッサによって遂行されてよく、またはコンピュータ装置504の外部のプロセッサによって遂行されてもよい。或る実施形態では、プロセッサ506は、分散型共有記憶システムの図に示されたものなどの命令を実行するようにプログラムされる。

30

【0067】

記憶エリア508は、コンピュータ装置504に関連した、またはコンピュータ装置504がアクセス可能な、何らかの量のコンピュータ可読媒体を含む。記憶エリア508またはその一部分は、コンピュータ装置504の内部、コンピュータ装置504の外部、またはその両方にある。

【0068】

図5の実例では、記憶エリア508は、SAN 512のトポロジ記述510を記憶する。トポロジ記述510は、ホスト514、スイッチ516などSAN 512の要素、およびディスク・アレイ302、ならびにその構成（例えばHBA、ポート304、記憶プロセッサ306など）およびその間のリンクを識別する。或る実施形態では、トポロジ記述510はさらに、SAN 512に関する閾値容量を含む。閾値容量は、例えばHBA当たりのホスト514ごとのI/O帯域幅を表すが、他の閾値容量も本開示の範囲内である。記憶エリア508は、リンクと関連したリンク容量も記憶してよい。

40

【0069】

或る実施形態では、コンピュータ装置504は、ネットワーク（図示せず）を介してSAN 512の要素にアクセスする。ネットワークは、コンピュータ装置504とSAN 512の要素の間を連絡するための任意の手段を表す。本開示の態様は、任意のネットワークのタイプまたは構成に対して動作可能である。

【0070】

図6は、SAN 512などのSANの中の最適の記憶パスを選択するように、複数の

50

ホスト 5 1 4 を管理する集中型システムによって遂行される例示の方法の流れ図である。方法 6 0 0 は、コンピュータ装置 5 0 4 (図 5 に示されている) による実行を参照しながら説明されているが、任意のコンピュータ装置によって遂行され得るように企図されている。さらに、図 6 に示された動作の実行は、V M 環境に限定されることなく、任意の複数起点および複数宛先の環境に適用可能である。また、コンピュータ実行可能命令を記憶する 1 または複数のコンピュータ可読記憶媒体が動作することにより、プロセッサ 5 0 6 は、図 6 に示された動作を遂行することによって集中型の I / O 負荷ベースパス選択を実施してよい。

【 0 0 7 1 】

ステップ 6 0 4 で、複数のホスト 5 1 4 のうちの 1 つから、同ホストが認識可能な少なくとも 1 つの論理的記憶装置にそれぞれ至る複数の記憶パスを識別するように、複数のホスト 5 1 4 に対して初期のパス選択が遂行される (図 4 を参照されたい) 。初期のパス選択は、新規のトポロジに対して、または I / O 負荷の変更に際して遂行されるが、或る実施形態は、パスを選択するのではなく、1 組の選択されたパスにアクセスする。

10

【 0 0 7 2 】

一般に、初期のパス選択は、各繰返しに対してエッジ・ウェイトを増す最大フロー・アルゴリズムを繰り返し遂行することによって、記憶パスを繰り返し検索することを含む。各繰返しで、ホスト - L U N 対の何らかの量が発見されてよく、または発見されないこともあり、ホスト - L U N 対がすべて見つかるまでエッジ・ウェイトが増加される。例えば、ステップ 6 0 2 で、各リンクに対して 1 のエッジ・ウェイトを定義することによって初期のパス選択が開始する。ステップ 6 0 6 で各ホスト 5 1 4 用に記憶パスが見つからなければ、ステップ 6 0 7 でエッジ・ウェイトが増加され、ステップ 6 0 4 でパス選択が遂行される。このように、各ホスト 5 1 4 からそれぞれの接続された L U N (例えば記憶装置) へのパスが見つかるまでパス選択を繰り返して、初期のパス選択を完了する。

20

【 0 0 7 3 】

或る実施形態では、例えばリンクの過剰予約を防止するために、任意選択で、1 または複数のリンクに対する上限リンク容量が定義されてよい。或る実施形態では、ネットワーク要素にリンク容量を問い合わせよく、またはテスト I / O 要求によってリンク容量を発見することができる。例えば、記憶パスを選択するのに使用されるスイッチ 5 1 6 の最大容量を得るために、記憶パス内のスイッチ 5 1 6 に問い合わせよく。

30

【 0 0 7 4 】

ステップ 6 0 8 では、見つかった各記憶パスについて、記憶パスに対する I / O 負荷が求められる。ステップ 6 1 0 では、各記憶パスに対して求めた I / O 負荷を閾値容量から減じる。ステップ 6 1 4 で、例えばバイナリ・サーチを遂行して、バイナリ・サーチからの各値に対して最大フロー・アルゴリズムを実行することによって、記憶パスのそれぞれに適合する、そうでなければ各ホスト 5 1 4 からそれぞれの接続された L U N への少なくとも 1 つのパスを見つけるのに必要とされる、最小限の帯域幅が見いだされる。

【 0 0 7 5 】

ステップ 6 1 6 で、ホスト 5 1 4 のそれぞれに対して、記憶パスのうちの 1 つが最小限の帯域幅に基づき選択される。各ホスト 5 1 4 向けに選択された記憶パスが、少なくとも対応するホスト 5 1 4 に対して、表示される、識別される、推奨される、または示される。

40

【 0 0 7 6 】

図 6 に示されたパス選択動作は、S A N 5 1 2 のトポロジの変更またはパスに沿った I / O 負荷の変更に応答して繰り返される。例えば、ステップ 6 1 8 で、トポロジおよび I / O の負荷が監視されてよい。ステップ 6 2 0 でトポロジ変更が検出された場合、または選択された記憶パスに沿った I / O の変更がステップ 6 2 2 で検出された場合、これに応答して、初期のパス選択が再度遂行されてよい。例えば、記憶パスのうちの 1 つに関して監視された I / O 負荷の変更が検出されるとすぐに、更新された I / O 負荷が求められ、記憶パスのそれぞれに適合する最小限の帯域幅が再度計算され、各ホスト 5 1 4 に対し

50

て、記憶パスのうちの1つが、再計算された最小の帯域幅に基づき選択される。

【0077】

図7を参照しながら次に説明されるように、トポロジ変更が推奨されてよく、或る実施形態では自動的に開始される。

図7は、リンクまたは論理装置番号(LUN)の輻輳にตอบสนองしてトポロジ変更を開始するように、複数のホスト514を管理する集中型システムによって遂行される例示の方法の流れ図である。方法700は、コンピュータ装置504(図5に示されている)による実行を参照しながら説明されているが、任意のコンピュータ装置によって遂行され得るように企図されている。さらに図7に示された動作の実行は、VM環境に限定されることなく、トポロジ変更が開始され得る任意の複数起点および複数宛先の環境に適用可能である。またコンピュータ実行可能命令を記憶する1または複数のコンピュータ可読記憶媒体が動作することにより、プロセッサ506は、図7に示された動作を遂行することによって集中型のI/O負荷ベースパス選択を実施してよい。

10

【0078】

ホスト514(例えば仮想マシンを実行する)のうちの1つから、LUNまたはホスト514から認識可能な他の論理的記憶装置のうちの1つにそれぞれ至る複数の記憶パスが、識別される、またはアクセスされる。ステップ704で、記憶パスのそれぞれに対するI/O負荷が監視される。ステップ706で、記憶パスの少なくとも1つにおいて監視されたI/O負荷が輻輳を示していると判断すると、ステップ708で、1または複数のトポロジ変更が識別される。例示のトポロジ変更は、それだけではないが、1つの論理的記憶装置から別の論理的記憶装置へデータを移動することと、複数のホスト514の第1のホストから複数のホスト514の第2のホストへデータストアを移動することと、論理的記憶装置をコントローラへ再マッピングすることを含む。

20

【0079】

コンピュータ装置504がSAN 512の中のエンドポイントを再構成することができる実施形態では、ステップ710で、識別された1または複数のトポロジ変更は、自動的に、かつ/または管理者502の入力なしで開始されてよい。例えば、記憶装置を別の記憶プロセッサ306またはコントローラに再マッピングするために、記憶装置を管理する記憶プロセッサ306のうちの1つに命令が送られてよい。或る実施形態では、トポロジ変更を開始することはさらに、記憶パスを再構成することと、ホスト514のうち少なくとも1つに対して別の記憶パスを用いるように指示することとを含んでよい。

30

追加の実例

本開示の態様が、1つのLUN当たり1つのパスを用いるアクティブ・パッシブ・モードでアクセスされているLUNを参照しながら本明細書で説明されてきたが、他の実施形態は、1つのLUN当たり複数のパスを用いるアクティブ・アクティブ・モードでLUNがアクセスされる環境を企図する。このような実施形態では、最大フロー・アルゴリズムが、k個のパスが許容されるLUNに対してk個のパスを見つけるように変更される。パスが見つかるたびに、要求は、1ではなくkだけインクリメントされる。

例示の動作環境

本明細書で説明された動作は、コンピュータ装置504などのコンピュータまたはコンピュータ装置によって遂行されてよい。コンピュータ装置は、メッセージおよび/または記憶したデータの交換を通じて互いに通信する。通信は、任意の有線または無線の接続にわたって、任意のプロトコルまたは機構を用いて起こり得る。コンピュータ装置は、(例えばネットワークおよび/またはデータバスの全体への)ブロードキャスト・メッセージ、(例えば複数の他のコンピュータ装置にアドレス指定された)マルチキャスト・メッセージ、および/または複数のユニキャスト・メッセージとしてメッセージを伝送してよく、これらのそれぞれが、個々のコンピュータ装置にアドレス指定される。さらに、或る実施形態では、メッセージは、ユーザ・データグラム・プロトコル(UDP: User Datagram Protocol)などの配信を保証しないネットワーク・プロトコルを用いて伝送される。したがってコンピュータ装置は、メッセージを伝送するとき、配信

40

50

不能のリスクを低減することができるように、メッセージの複数のコピーを伝送してよい。

【0080】

例示のコンピュータ読取り可能媒体は、フラッシュ・メモリ・ドライブ、デジタル多用途ディスク（DVD）、コンパクト・ディスク（CD）、フロッピー（登録商標）・ディスク、およびテープ・カセットを含む。限定的でない実例として、コンピュータ読取り可能媒体は、コンピュータ記憶媒体および通信媒体を含む。コンピュータ記憶媒体には、コンピュータ読取り可能命令、データ構造、プログラム・モジュールまたは他のデータなどの情報を記憶するための任意の方法または技術で実施された、揮発性媒体および不揮発性媒体、取外し可能媒体および固定型媒体が、含まれる。コンピュータ記憶媒体は、有形であり、伝搬されたデータ信号を除外し、通信媒体とは相互排他的である。或る実施形態では、コンピュータ記憶媒体はハードウェアで実施される。例示のコンピュータ記憶媒体は、ハード・ディスク、フラッシュ・ドライブ、および他の固体記憶装置を含む。それと対照的に、通信媒体は、一般に、コンピュータ読取り可能命令、データ構造、プログラム・モジュール、または搬送波もしくは他の移送機構などの変調データ信号における他のデータを具体化し、任意の情報配送媒体を含む。

10

【0081】

本開示の実施形態は、例示のコンピュータ・システム環境と関連して説明されてきたが、多数の他の汎用または専用のコンピュータ・システムの環境または構成で動作可能である。本開示の態様とともに用いるのに適切であり得る周知のコンピュータ・システムの実例、環境、および/または構成には、それだけではないが、モバイル・コンピューティング装置、パーソナル・コンピュータ、サーバ・コンピュータ、ハンドヘルドまたはラップトップの装置、マルチプロセッサ・システム、ゲーム機、マイクロプロセッサ・ベースのシステム、セット・トップ・ボックス、プログラム可能な家電、携帯電話、ネットワークPC、ミニコンピュータ、メインフレーム・コンピュータ、上記システムのうち任意ものを含む分散型コンピュータ環境または装置などが含まれる。

20

【0082】

本開示の実施形態は、1または複数のコンピュータまたは他の装置によって実行されるプログラム・モジュールなどのコンピュータ実行可能命令の一般的な状況において説明されてよい。コンピュータ実行可能命令は、1または複数のコンピュータ実行可能な構成要素またはモジュールに体系化され得る。一般に、プログラム・モジュールは、それだけではないが、特定のタスクを遂行する、または特定の抽象データ型を実施するルーチン、プログラム、オブジェクト、構成要素、およびデータ構造を含む。本開示の態様は、このような構成要素またはモジュールの任意の数および構成を用いて実施され得る。例えば、本開示の態様は、特定のコンピュータ実行可能命令または各図面に示されて本明細書で説明された特定の構成要素もしくはモジュールに限定されない。本開示の他の実施形態は、本明細書に示され説明されたものよりも、程度の差はあるが機能性を有する様々なコンピュータ実行可能命令または構成要素を含んでよい。

30

【0083】

本開示の態様は、本明細書に説明された命令を実行するようにプログラムされたとき、汎用コンピュータを専用コンピュータ装置に変換する。

40

本明細書で図示して説明した実施形態ならびに本明細書で具体的には説明されていないが本発明の態様の範囲内の実施形態は、記憶パスの輻輳を修正するためにトポロジ変更を識別する例示の手段と、記憶パスの輻輳を修正するためにトポロジ変更を自動的に開始する例示の手段とを構成する。例えば、本明細書で説明されたような動作を実行するようにプログラムされたプロセッサ506は、これらの例示の手段を構成する。

【0084】

図に示された様々な要素の機能の少なくとも一部分は、図中の他の要素、また図には示されていないエンティティ（例えばプロセッサ、ウェブ・サービス、サーバ、アプリケーション・プログラム、コンピュータ装置など）によって遂行されてよい。

50

【 0 0 8 5 】

或る実施形態では、図に示された動作は、コンピュータ可読媒体に符号化されたソフトウェア命令として、この動作を遂行するようにプログラムされた、もしくは設計されたハードウェアで、またはその両方で実施されてよい。例えば、本開示の態様は、チップまたは複数の相互接続された導電素子を含んでいる他の回路上のシステムとして実施されてよい。

【 0 0 8 6 】

本明細書で図示して説明した本開示の実施形態の動作の実行または遂行の順序は、別段の定めがない限り必須ではない。すなわち、これらの動作は、別段の定めがない限り任意の順序で遂行されてよく、本開示の実施形態は、本明細書で開示されたものに対して追加の動作を含んでよく、より少ない動作しか含まなくてもよい。例えば、特定の動作を、別の動作の以前に、同動作と同時に、または同動作の後に、実行する、または遂行するということは、本開示の態様の範囲に入るように企図されている。

10

【 0 0 8 7 】

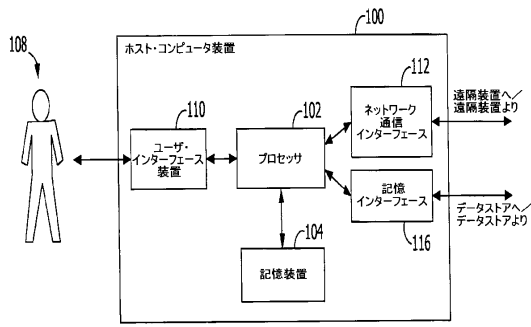
本開示の態様の要素またはその実施形態を紹介するとき、冠詞「1つの(a)」、「1つの(an)」、「この(the)」、および「前記(said)」は、同要素の1または複数があることを意味するように意図されている。用語「備える」、「含む」、および「有する」は、包含的であるように意図されており、列挙された要素とは別に追加の要素が存在し得ることを意味する。

【 0 0 8 8 】

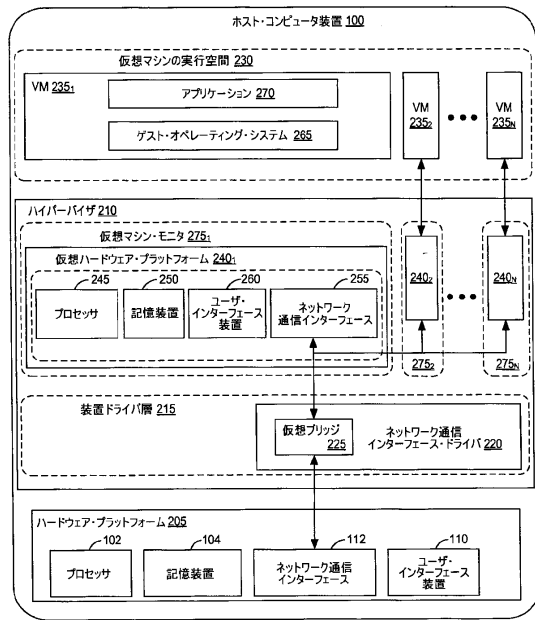
本開示の態様を詳細に説明してきたが、添付の特許請求の範囲で定義される本開示の態様の範囲から逸脱することなく、修正形態および変形形態が可能であることが明らかであろう。本開示の態様の範囲から逸脱することなく、上記の構成、結果、および方法に様々な変更を加えることができるので、上記の記述の中に含まれ、添付図面に示されていることは、すべて例示であって限定する意味ではないと解釈されるべきであることが意図されている。

20

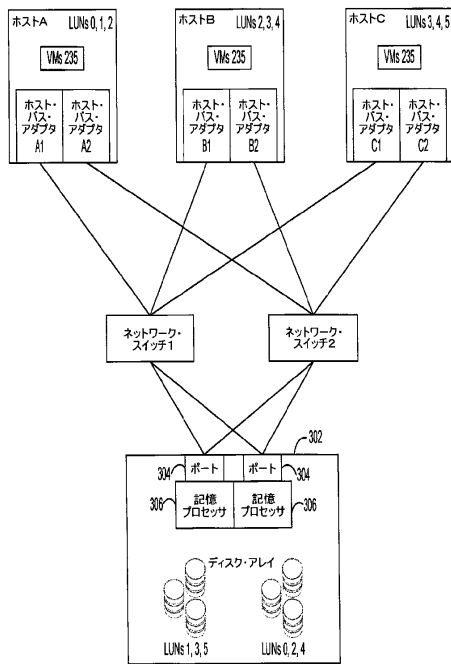
【図1】



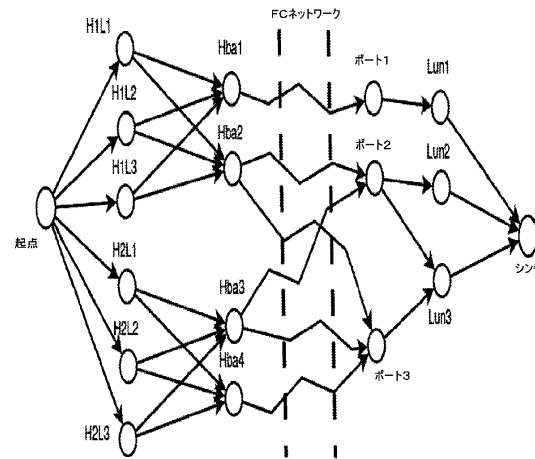
【図2】



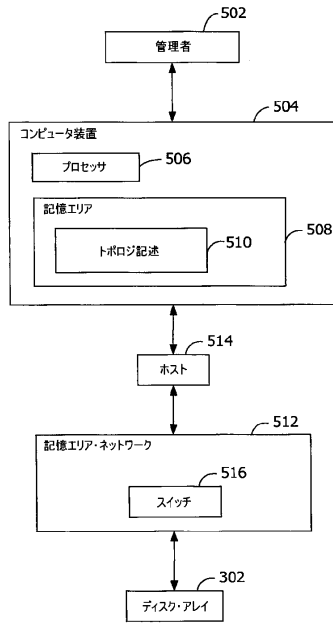
【図3】



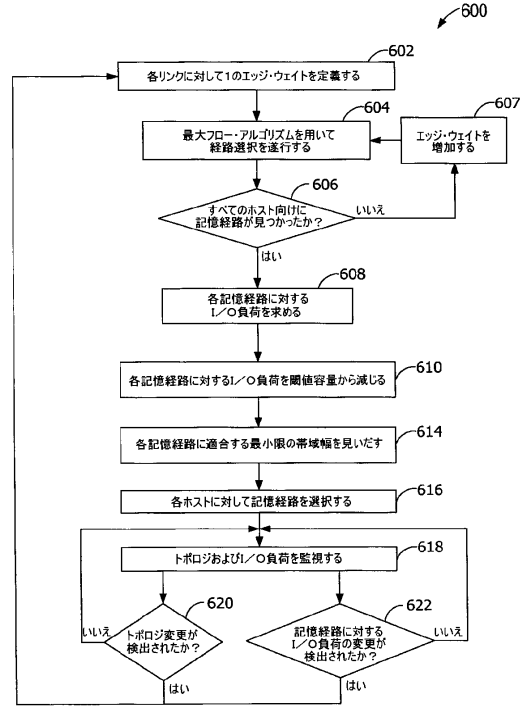
【図4】



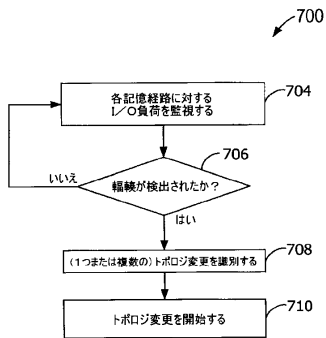
【図5】



【図6】



【図7】



フロントページの続き

(72)発明者 クリシュナ ラージ ラージャ

アメリカ合衆国 95051 カリフォルニア州 サンタ クララ ギャンプリン ドライブ 2
658

(72)発明者 アージェイ グラティ

アメリカ合衆国 94304 カリフォルニア州 パロ アルト ヒルビュー アベニュー 34
01

審査官 田上 隆一

(56)参考文献 米国特許出願公開第2009/0172666(US,A1)

米国特許出願公開第2008/0250178(US,A1)

特開2007-094681(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/14

G06F 3/06

G06F 13/10