

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2014-529804  
(P2014-529804A)

(43) 公表日 平成26年11月13日(2014.11.13)

(51) Int.Cl.		F I		テーマコード (参考)
<b>G06F 3/06 (2006.01)</b>		G06F 3/06	301Z	
<b>G06F 13/10 (2006.01)</b>		G06F 3/06	301F	
		G06F 13/10	340A	

審査請求 未請求 予備審査請求 未請求 (全 70 頁)

(21) 出願番号 特願2014-527263 (P2014-527263)  
 (86) (22) 出願日 平成24年8月22日 (2012. 8. 22)  
 (85) 翻訳文提出日 平成26年1月30日 (2014. 1. 30)  
 (86) 国際出願番号 PCT/US2012/051872  
 (87) 国際公開番号 WO2013/032810  
 (87) 国際公開日 平成25年3月7日 (2013. 3. 7)  
 (31) 優先権主張番号 13/219, 368  
 (32) 優先日 平成23年8月26日 (2011. 8. 26)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 510149482  
 ヴィエムウェア インコーポレイテッド  
 VMware, Inc.  
 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 3401  
 (74) 代理人 100105957  
 弁理士 恩田 誠  
 (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】 入力/出力オペレーションのためにオブジェクト・ストレージ・システムを構成すること

(57) 【要約】

ストレージ・システムが、ストレージ・オブジェクトとしてプロビジョンされる論理ストレージ・ボリュームをエクスポートする。これらのストレージ・オブジェクトは、ストレージ・システムにおいて構成されているプロトコル・トラフィックに関する論理エンドポイントを通じて、SCSIおよびNFSなどの標準的なプロトコルを使用して、接続されているコンピュータ・システムによってオン・デマンドでアクセスされる。ストレージ・システムにおいて入力/出力コマンド(I/O)が受信された場合には、そのI/Oから識別子が取り出され、論理ストレージ・ボリューム識別子へと変換され、その論理ストレージ・ボリューム識別子に対応する論理ストレージ・ボリュームによって参照されるストレージ・ロケーション上でI/Oが実行される。

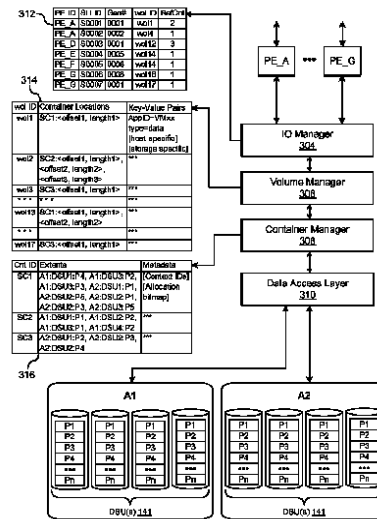


FIGURE 3

**【特許請求の範囲】****【請求項 1】**

ストレージ・システムにおいて入力/出力コマンド（I/O）を処理するための方法であって、

第 1 および第 2 の識別子を含む I/O をストレージ・システムにおいて受信することであって、該第 1 の識別子が、該 I/O を該ストレージ・システムへ導くために使用される、前記 I/O を受信すること、

該ストレージ・システムにおいて、該 I/O から該第 2 の識別子を取り出し、該第 2 の識別子を論理ストレージ・ボリューム識別子に変換すること、

該論理ストレージ・ボリューム識別子に対応する論理ストレージ・ボリュームによって参照されるストレージ・ロケーション上で I/O を実行すること、を含む方法。

10

**【請求項 2】**

前記 I/O が、前記ストレージ・システムにおいて構成されている LUN において受信され、前記第 1 の識別子が、該 LUN に関するワールド・ワイド・ネームである、請求項 1 に記載の方法。

**【請求項 3】**

前記論理ストレージ・ボリュームまたは前記 LUN におけるエラー状況を検知し、前記 I/O を発行したコンピュータ・システムにエラー・メッセージを送信すること

をさらに含む、請求項 2 に記載の方法。

20

**【請求項 4】**

前記 I/O が、前記ストレージ・システムにおいて構成されているマウント・ポイントにおいて受信され、前記第 1 の識別子が、前記ストレージ・システムおよび該マウント・ポイントの IP アドレスを含む、請求項 1 に記載の方法。

**【請求項 5】**

前記論理ストレージ・ボリュームまたは前記マウント・ポイントにおけるエラー状況を検知し、前記 I/O を発行したコンピュータ・システムにエラー・メッセージを送信すること

をさらに含む、請求項 4 に記載の方法。

**【請求項 6】**

論理ストレージ・ボリューム識別子への第 2 の識別子のマッピングを提供するデータ構造を保持することであって、該第 2 の識別子が、該データ構造を使用して該論理ストレージ・ボリューム識別子に変換される、前記保持すること

をさらに含む、請求項 1 に記載の方法。

30

**【請求項 7】**

ストレージ・システムを用意して入力/出力コマンド（I/O）を実行するための方法であって、

論理ストレージ・ボリュームを用意するための要求を受信すること、

該論理ストレージ・ボリュームに関する I/O が受信される際に経由することになる該ストレージ・システムにおいて構成されているプロトコル・エンドポイントを選択すること

、  
該論理ストレージ・ボリュームに関する I/O 内に含まれることになる識別子を生成し、該要求を発行したコンピュータ・システムに該識別子を返すこと、を含む方法。

40

**【請求項 8】**

プロトコル・パスを介して前記コンピュータ・システムに接続されている前記ストレージ・システムにおいて構成されている 1 つまたは複数のプロトコル・エンドポイントを特定することをさらに含み、

前記プロトコル・エンドポイントが、該 1 つまたは複数のプロトコル・エンドポイントから選択される、

請求項 7 に記載の方法。

**【請求項 9】**

前記 1 つまたは複数のプロトコル・エンドポイントが LUN である、請求項 8 に記載の方

50

法。

【請求項 10】

前記 1 つまたは複数のプロトコル・エンドポイントが、ネットワーク・アタッチト・ストレージ・システムのマウント・ポイントである、請求項 8 に記載の方法。

【請求項 11】

前記要求が、非プロトコル・パスを介して前記ストレージ・システムにおいて受信され、前記識別子が、該非プロトコル・パスを介して返される、請求項 8 に記載の方法。

【請求項 12】

IO 内に含まれている識別子の、論理ストレージ・ボリューム識別子へのマッピングを提供するデータ構造を保持することをさらに含み、

前記要求が、論理ストレージ・ボリューム識別子を含み、前記選択されたプロトコル・エンドポイントを識別して生成された識別子を該論理ストレージ・ボリューム識別子に関連付けるエントリが、該データ構造に加えらる、

請求項 7 に記載の方法。

【請求項 13】

前記データ構造が、それぞれのエントリに関するリファレンス・カウントを保持し、該リファレンス・カウントが、前記論理ストレージ・ボリュームを用意するための要求が行われるたびにインクリメントされ、前記論理ストレージ・ボリュームに関する IO を受信するために選択されるプロトコル・エンドポイントが、該エントリ内で識別される同じプロトコル・エンドポイントである、請求項 12 に記載の方法。

【請求項 14】

前記リファレンス・カウントが、前記論理ボリュームに関する IO を受信するために選択された前記プロトコル・エンドポイントから前記論理ストレージ・ボリュームを切り離すための要求が行われるたびにデクリメントされる、請求項 13 に記載の方法。

【請求項 15】

前記論理ストレージ・ボリュームに関する IO が受信される際に経由することになる前記ストレージ・システムにおいて構成されている新たなプロトコル・エンドポイントを選択し、前記論理ストレージ・ボリュームを該新たなプロトコル・エンドポイントに結合すること、

前記論理ストレージ・ボリュームに関する IO 内に含まれることになる新たな識別子を生成し、前記要求を発行したコンピュータ・システムに該新たな識別子を返すこと、

をさらに含む、請求項 7 に記載の方法。

【請求項 16】

前記論理ストレージ・ボリュームを古いプロトコル・エンドポイントから切り離して、前記論理ストレージ・ボリュームに関する新たな IO が、前記新たなプロトコル・エンドポイントを通じて受信されるようにすること

をさらに含む、請求項 15 に記載の方法。

【請求項 17】

ストレージ・システムを用意して、ホスト・コンピュータ上で稼働する仮想マシンによって発行される入力/出力コマンド (IO) を実行するための方法であって、

該仮想マシンに関するメタデータ論理ストレージ・ボリュームを用意するための要求を受信すること、

該メタデータ論理ストレージ・ボリュームに関する IO が受信される際に経由することになる該ストレージ・システムにおいて構成されているプロトコル・エンドポイントを選択すること、

該プロトコル・エンドポイントを通じて読み取り IO を受信し、データ論理ストレージ・ボリュームに関する識別子を含む読み取りデータを返すこと、

該仮想マシンに関する該データ論理ストレージ・ボリュームを用意するための要求を受信すること、

それぞれのデータ論理ストレージ・ボリュームに関して、論理ストレージ・ボリューム

10

20

30

40

50

に関する I O が受信される際に経由することになる該ストレージ・システムにおいて構成されているプロトコル・エンドポイントを選択すること、を含む方法。

【請求項 18】

前記メタデータ論理ストレージ・ボリュームが、前記仮想マシンに関する構成ファイルおよびログ・ファイルと、前記仮想マシンに関するスワップ・ファイルと、ディスク記述子ファイルとを含み、ファイルは、前記仮想マシンの仮想ディスクのそれぞれに関するものである、請求項 17 に記載の方法。

【発明の詳細な説明】

【背景技術】

【0001】

コンピュータ・システムが、とりわけ大規模なデータ・センターをサポートするというコンテキストにおいて、エンタープライズ・レベルへ拡張するにつれて、基礎をなすデータ・ストレージ・システムは、ストレージ・エリア・ネットワーク (SAN) またはネットワーク・アタッチト・ストレージ (NAS) を採用する場合が多い。従来からよく理解されているように、SAN または NAS は、複数の技術的能力および機能上の利点を提供し、それらは基本的に、データ・ストレージ・デバイスの仮想化と、トランスペアレントな、フォルト・トレラントな、フェイルオーバな、およびフェイルセーフなコントロールを伴う物理的なデバイスの冗長性と、地理的に分散され複製されたストレージと、クライアント中心のコンピュータ・システム管理から切り離された集中化された監督およびストレージ構成管理とを含む。

【0002】

アーキテクチャ上は、SAN ストレージ・システム (たとえば、ディスク・アレイなど) 内のストレージ・デバイスは、典型的にはネットワーク・スイッチ (たとえば、ファイバ・チャネル・スイッチなど) に接続され、次いで、それらのネットワーク・スイッチは、ストレージ・デバイス内のデータへのアクセスを必要とするサーバまたは「ホスト」に接続される。SAN 内のサーバ、スイッチ、およびストレージ・デバイスは、典型的には、ディスク・データ・ブロックのレベルでネットワークを介してデータを転送するスモール・コンピュータ・システム・インターフェース (SCSI) プロトコルを使用して通信する。対照的に、NAS デバイスは、典型的には、1 つまたは複数のストレージ・ドライブを内部に含み、イーサネットなどのネットワーク・プロトコルを通じてホスト (または中間スイッチ) に接続されるデバイスである。NAS デバイスはまた、ストレージ・デバイスを含むことに加えて、ネットワーク・ファイル・システム (NFS) またはコモン・インターネット・ファイル・システム (CIFS) など、ネットワークベースのファイル・システムに従って自分のストレージ・デバイスを事前にフォーマットする。したがって、SAN は、ディスク (LUN と呼ばれ、これについては、以降でさらに詳述する) をホストに公開し、次いでそれらのディスクは、フォーマットされてから、ホストによって利用されるファイル・システムに従ってマウントされる必要があるが、そうした SAN とは対照的に、NAS デバイスのネットワークベースのファイル・システム (これは、ホストのオペレーティング・システムによってサポートされる必要がある) は、NAS デバイスが、ホストのオペレーティング・システムにとってファイル・サーバとして見えるようにし、次いでホストは、その NAS デバイスを、たとえば、オペレーティング・システムによってアクセス可能なネットワーク・ドライブとして、マウントまたはマップすることができる。ストレージ・システム・ベンダーによる継続的なイノベーションおよび新製品のリリースに伴って、SAN ストレージ・システムと NAS ストレージ・システムとの間における明確な区別が薄れ続けており、実際のストレージ・システムの実施態様は、しばしば両方の特徴を呈しており、同じシステムにおいてファイルレベル・プロトコル (NAS) およびブロックレベル・プロトコル (SAN) の両方を提供しているということを認識されたい。たとえば、代替 NAS アーキテクチャにおいては、従来の NAS デバイスではなく、NAS 「ヘッド」または NAS 「ゲートウェイ」デバイスが、ホストにネットワーク接続される。そのような NAS ゲートウェイ・デバイスは、自分自身ではストレージ・

10

20

30

40

50

ドライブを含まず、外部のストレージ・デバイスが（たとえば、ファイバ・チャネル・インターフェースなどを介して）そのNASゲートウェイ・デバイスに接続されることを可能にする。そのようなNASゲートウェイ・デバイス（これは、従来のNASデバイスと同様の様式でホストによって知覚される）は、ファイルレベルのストレージ・アクセスのシンプルさを保持しながら、NASベースのストレージ・アーキテクチャのキャパシティーを（たとえば、SANによって、より伝統的にサポートされているストレージ・キャパシティー・レベルに）著しく増大させる能力を提供する。

**【0003】**

図1Aにおいて示されているストレージ・システム30など、SCSIおよびその他のブロック・プロトコルベースのストレージ・デバイスは、1つまたは複数のプログラムされたストレージ・プロセッサに相当するストレージ・システム・マネージャ31を利用して、そのストレージ・デバイス内のストレージ・ユニットまたはドライブを集約し、それらを、一意に識別可能な番号をそれぞれ伴う1つまたは複数のLUN (Logical Unit Number) 34として提示する。LUN34は、ネットワーク20（たとえば、ファイバ・チャネルなど）を介して物理ホスト・バス・アダプタ(HBA)11を通じて1つまたは複数のコンピュータ・システム10によってアクセスされる。コンピュータ・システム10内で、HBA11の上に、ストレージ・アクセス・アブストラクションが、ローレベル・デバイス・ドライバ・レイヤ12から始まってオペレーティング・システム固有のファイル・システム・レイヤ15で終わる一連のソフトウェア・レイヤを通じて特徴的に実装されている。デバイス・ドライバ・レイヤ12は、LUN34への基本的なアクセスを可能にし、典型的には、ストレージ・システムによって使用される通信プロトコル（たとえば、SCSIなど）に固有のものである。HBA11を通じて見えるLUN34のマルチパスの統合、およびその他のデータ・アクセス・コントロールおよび管理機能をサポートするために、データ・アクセス・レイヤ13が、デバイス・ドライバ・レイヤ12の上に実装されることが可能である。論理ボリューム・マネージャ14が、典型的にはデータ・アクセス・レイヤ13と従来のオペレーティング・システム・ファイル・システム・レイヤ15との間に実装され、HBA11を通じてアクセス可能なLUN34のボリューム指向の仮想化および管理をサポートする。複数のLUN34が集められて、1つの論理デバイスとしてファイル・システム・レイヤ15に提示されてファイル・システム・レイヤ15によって使用されるために論理ボリューム・マネージャ14のコントロールのもとで1つのボリュームとしてまとめて管理されることが可能である。

**【0004】**

ストレージ・システム・マネージャ31は、ストレージ・システム30内に存在する、図1Aにおいてスピンドル32と呼ばれている、物理的な、典型的にはディスク・ドライブベースのストレージ・ユニットの仮想化を実施する。論理的な観点からは、これらのスピンドル32のそれぞれは、固定されたサイズのエクステント33のシーケンシャル・アレイと考えられることが可能である。ストレージ・システム・マネージャ31は、LUN34として知られている一組の仮想SCSIデバイスへと分けられている連続した論理ストレージ・スペースを、コンピュータ・システム10など、接続されているコンピュータ・システムに公開することによって、ディスク・ドライブの実際のスピンドルおよびエクステントのアドレスへのターゲットとなる読み取りおよび書き込みオペレーションの複雑さを取り去る。それぞれのLUNは、そのようなLUNの存在と、コンピュータ・システム10へのそのようなLUNの提示とのおかげでコンピュータ・システム10によって使用されるために割り振られるいくらかのキャパシティーを表す。ストレージ・システム・マネージャ31は、それぞれのそのようなLUNに関する、エクステントの順序付けられたリストへのマッピングを含むメタデータを保持し、そのリストにおいては、それぞれのそのようなエクステントは、スピンドル/エクステントのペア<スピンドル#, エクステント#>として識別されることが可能であり、したがって、さまざまなスピンドル32のうちのいずれかに位置特定されることが可能である。

**【0005】**

10

20

30

40

50

図1Bは、ネットワーク21（たとえば、イーサネット）を介してネットワーク・インターフェース・カード（NIC）11'を経由して1つまたは複数のコンピュータ・システム10に接続される従来のNASまたはファイルレベル・ベースのストレージ・システム40のブロック図である。ストレージ・システム40は、1つまたは複数のプログラムされているストレージ・プロセッサに相当するストレージ・システム・マネージャ41を含む。ストレージ・システム・マネージャ41は、ストレージ・システム40内に存在する、図1Bにおいてスピンドル42と呼ばれている、物理的な、典型的にはディスク・ドライブベースのストレージ・ユニットの上にファイル・システム45を実装する。論理的な観点からは、これらのスピンドルのそれぞれは、固定されたサイズのエクステント43のシーケンシャル・アレイと考えられることが可能である。ファイル・システム45は、各自のマウント・ポイントを通じてアクセスされるファイル・システム・レベル・ボリューム44（以降では「FSボリューム」と呼ばれる）へと編成されることが可能であるディレクトリおよびファイルを含むネームスペースを、コンピュータ・システム10など、接続されているコンピュータ・システムに公開することによって、ディスク・ドライブの実際のスピンドルおよびエクステントのアドレスへのターゲットとなる読み取りおよび書き込みオペレーションの複雑さを取り去る。

10

## 【0006】

上述のストレージ・システムにおける進歩をもってさえ、それらのストレージ・システムは、仮想化されたコンピュータ・システムの特定のニーズを満たす上で十分にスケラブルではないということが広く認識されている。たとえば、サーバ・マシンのクラスタは、10,000個もの仮想マシン（VM: virtual machine）にサービス提供することができ、それぞれのVMは、複数の「仮想ディスク」および複数の「スナップショット」を使用し、それらはそれぞれ、たとえば、1つのファイルとして特定のLUNまたはFSボリューム上に格納されることが可能である。VMごとに2つの仮想ディスクおよび2つのスナップショットというスケール・ダウンされた推定においてさえ、VMが物理ディスクに直接接続される（すなわち、物理ディスクごとに1つの仮想ディスクまたはスナップショット）場合、ストレージ・システムがサポートするのは、合計60,000個の別々のディスクに達する。加えて、このスケールでのストレージ・デバイスおよびトポロジーの管理は困難であることがわかっている。結果として、本願明細書に援用する「Providing Multiple Concurrent Access to a File System」と題されている（特許文献1）に記載されているような、VMがさらに小さな組の物理ストレージ・エンティティ（たとえば、LUNベースのVMFSクラスタ化されたファイル・システムまたはFSボリューム）へと多重化されるデータストアというコンセプトが開発された。

20

30

## 【0007】

LUNまたはFSボリュームを採用している従来のストレージ・システムにおいては、複数のVMからのワークロードは、典型的には、単一のLUNまたは単一のFSボリュームによってサービス提供される。結果として、1つのVMワークロードからのリソース需要が、同じLUNまたはFSボリューム上の別のVMワークロードに提供されるサービス・レベルに影響を与えることになる。したがって、待ち時間、および1秒あたりの入力/出力処理（IO）、すなわちIOPS（input/output operations per second）など、ストレージに関する効率尺度は、所与のLUNまたはFSボリューム内のワークロードの数に応じて変わり、保証されることは不可能である。その結果として、LUNまたはFSボリュームを採用しているストレージ・システムに関するストレージ・ポリシーがVMごとに実行されることは不可能であり、サービス・レベル・アグリーメント（SLA）保証がVMごとに与えられることは不可能である。加えて、スナップショット、複製、暗号化、および重複排除など、ストレージ・システム・ベンダーによって提供されるデータ・サービスは、VMの仮想ディスクの粒度ではなく、LUNまたはFSボリュームの粒度で提供される。結果として、スナップショットは、ストレージ・システム・ベンダーによって提供されるデータ・サービスを使用してLUN全体ま

40

50

たはFSボリューム全体に関して作成されることが可能であるが、LUN、または仮想ディスクが格納されているファイル・システムとは別に、VMの単一の仮想ディスクに関するスナップショットが作成されることは不可能である。

【先行技術文献】

【特許文献】

【0008】

【特許文献1】米国特許第7,849,098号明細書

【発明の概要】

【課題を解決するための手段】

【0009】

1つまたは複数の実施形態が対象にしているストレージ・システムは、その中で実行されるワークロード同士を分離するように構成されており、それによって、SLA保証がワークロードごとに提供されることが可能であり、ストレージ・システムのデータ・サービスがワークロードごとに提供されることが可能であり、ストレージ・システムの抜本的な再設計は必要とされない。複数の仮想マシンに関する複数の仮想ディスクを格納するストレージ・システムにおいては、SLA保証が仮想ディスクごとに提供されることが可能であり、ストレージ・システムのデータ・サービスが仮想ディスクごとに提供されることが可能である。

【0010】

本発明の複数の実施形態によれば、ストレージ・システムは、論理的なストレージ・キャパシティの割り当て（本明細書においては、「ストレージ・コンテナ」と呼ばれる）から、ワークロードごとにストレージ・オブジェクトとしてプロビジョンされる論理ストレージ・ボリューム（本明細書においては、「仮想ボリューム」と呼ばれる）をエクスポートする。VMに関しては、そのVMの仮想ディスクおよびスナップショットのそれぞれに関して仮想ボリュームが作成されることが可能である。一実施形態においては、仮想ボリュームは、ストレージ・システムにおいて構成されているプロトコル・トラフィックに関する論理エンドポイント（「プロトコル・エンドポイント」として知られている）を通じて、SCSIおよびNFSなどの標準的なプロトコルを使用して、接続されているコンピュータ・システムによってオン・デマンドでアクセスされる。

【0011】

本発明の一実施形態による、ストレージ・システムにおいて入力/出力コマンド（IO）を処理するための方法は、第1および第2の識別子を含むIOをストレージ・システムにおいて受信することであって、第1の識別子が、IOをストレージ・システムへ導くために使用される、IOを受信すること、IOから第2の識別子を取り出し、その第2の識別子を論理ストレージ・ボリューム識別子に変換すること、その論理ストレージ・ボリューム識別子に対応する論理ストレージ・ボリュームによって参照されるストレージ・ロケーション上でIOを実行することを含む。

【0012】

本発明の一実施形態による、ストレージ・システムを用意してIOを実行するための方法は、論理ストレージ・ボリュームを用意するための要求を受信すること、論理ストレージ・ボリュームに関するIOが受信される際に経由することになるストレージ・システムにおいて構成されているプロトコル・エンドポイントを選択すること、論理ストレージ・ボリュームに関するIO内に含まれることになる識別子を生成し、要求を発行したコンピュータ・システムにその識別子を返すことを含む。

【0013】

本発明の一実施形態による、ストレージ・システムを用意してホスト・コンピュータ上で稼働する仮想マシンによって発行されるIOを実行するための方法は、仮想マシンに関するメタデータ論理ストレージ・ボリュームを用意するための要求を受信すること、メタデータ論理ストレージ・ボリュームに関するIOが受信される際に経由することになるストレージ・システムにおいて構成されているプロトコル・エンドポイントを選択すること

10

20

30

40

50

、そのプロトコル・エンドポイントを通じて読み取りI/Oを受信し、データ論理ストレージ・ボリュームに関する識別子を含む読み取りデータを返すこと、仮想マシンに関するデータ論理ストレージ・ボリュームを用意するための要求を受信すること、それぞれのデータ論理ストレージ・ボリュームに関して、論理ストレージ・ボリュームに関するI/Oが受信される際に経由することになるストレージ・システムにおいて構成されているプロトコル・エンドポイントを選択することを含む。

【0014】

本発明の複数の実施形態はさらに、コンピュータ・システムによって実行されたときに上述の方法のうちの一つをそのコンピュータ・システムに実行させる命令を格納している非一時的なコンピュータ可読ストレージ・メディアを含む。

10

【図面の簡単な説明】

【0015】

【図1A】ネットワークを介して一つまたは複数のコンピュータ・システムに接続されている従来のブロック・プロトコルベースのストレージ・デバイスのブロック図。

【図1B】ネットワークを介して一つまたは複数のコンピュータ・システムに接続されている従来のNASデバイスのブロック図。

【図2A】本発明の一実施形態による、仮想ボリュームを実装しているブロック・プロトコルベースのストレージ・システム・クラスタのブロック図。

【図2B】本発明の一実施形態による、仮想ボリュームを実装しているNASベースのストレージ・システム・クラスタのブロック図。

20

【図3】本発明の一実施形態による、仮想ボリュームを管理するための図2Aまたは図2Bのストレージ・システム・クラスタのコンポーネントのブロック図。

【図4】ストレージ・コンテナを作成するための方法工程の流れ図。

【図5A】SANベースのストレージ・システム上にホストされる仮想ボリュームを実装するように構成されているコンピュータ・システムの一実施形態のブロック図。

【図5B】NASベースのストレージ・システム上にホストされる仮想ボリュームのために構成されている図5Aのコンピュータ・システムのブロック図。

【図5C】SANベースのストレージ・システム上にホストされる仮想ボリュームを実装するように構成されているコンピュータ・システムの別の実施形態のブロック図。

【図5D】NASベースのストレージ・システム上にホストされる仮想ボリュームのために構成されている図5Cのコンピュータ・システムのブロック図。

30

【図6】本発明の一実施形態による、仮想ボリュームを管理するために使用されるコンポーネントおよび通信パスを示すコンピュータ環境の簡略化されたブロック図。

【図7】図2Aまたは図2Bのストレージ・システム・クラスタに対してコンピュータ・システムを認証するための方法工程の流れ図。

【図8】一実施形態による、仮想ボリュームを作成するための方法工程の流れ図。

【図9】Aは、コンピュータ・システムにとって利用可能であるプロトコル・エンドポイントを発見するための方法工程の流れ図、Bは、コンピュータ・システムが帯域内パスを介して接続されるプロトコル・エンドポイントをストレージ・システムが発見するための方法工程の流れ図。

40

【図10】一実施形態による、仮想ボリューム・バインド要求を発行および実行するための方法工程の流れ図。

【図11】Aは、一実施形態による、仮想ボリュームにI/Oを発行するための方法工程の流れ図、Bは、一実施形態による、仮想ボリュームにI/Oを発行するための方法工程の流れ図。

【図12】一実施形態による、ストレージ・システムにおいてI/Oを実行するための方法工程の流れ図。

【図13】一実施形態による、仮想ボリューム再バインド要求を発行および実行するための方法工程の流れ図。

【図14】仮想ボリュームのライフ・サイクルの概念図。

50



【図15】図2Aのストレージ・システムを使用する一実施形態による、VMをプロビジョンするための方法工程の流れ図。

【図16】Aは、VMをパワー・オンするための方法工程の流れ図、Bは、VMをパワー・オフするための方法工程の流れ図。

【図17】VMのvvolのサイズを拡張するための方法工程の流れ図。

【図18】ストレージ・コンテナ同士の間においてVMのvvolを移動させるための方法工程の流れ図。

【図19】テンプレートVMからVMをクローンするための方法工程の流れ図。

【図20】別の実施形態による、VMをプロビジョンするための方法工程の流れ図。

【図21】サンプル・ストレージ能力プロファイルと、プロファイル選択工程を含む、ストレージ・コンテナを作成するための方法とを示す図。

【図22】vvolを作成して、そのvvolに関するストレージ能力プロファイルを定義するための方法工程を示す流れ図。

【図23】スナップショットを作成するための方法工程を示す流れ図。

【発明を実施するための形態】

【0016】

図2Aおよび図2Bは、本発明の実施形態による、「仮想ボリューム」を実装しているストレージ・システム・クラスタのブロック図である。このストレージ・システム・クラスタは、1つまたは複数のストレージ・システム、たとえば、ストレージ・システム130<sub>1</sub>および130<sub>2</sub>（これらは、ディスク・アレイであることが可能であり、それぞれが、複数のデータ・ストレージ・ユニット(DSU)を有しており、それらのDSUのうちの1つは、図において141とラベル付けされている)と、本明細書に記載されている本発明の実施形態を可能にするためのストレージ・システム130のさまざまなオペレーションをコントロールするストレージ・システム・マネージャ131および132を含む。一実施形態においては、複数のストレージ・システム130が、分散型ストレージ・システム・マネージャ135を実装することができ、分散型ストレージ・システム・マネージャ135は、ストレージ・システム・クラスタのオペレーションを、あたかもそれらが単一の論理ストレージ・システムであるかのようにコントロールする。分散型ストレージ・システム・マネージャ135の運用ドメインは、同じデータ・センター内に、または複数のデータ・センターにわたってインストールされているストレージ・システムに及ぶことができる。たとえば、そのような一実施形態においては、分散型ストレージ・システム・マネージャ135は、ストレージ・システム・マネージャ131を含むことができ、ストレージ・システム・マネージャ131は、ストレージ・システム・マネージャ132と通信する場合に「マスター」マネージャとして機能し、ストレージ・システム・マネージャ132は、「スレーブ」マネージャとして機能するが、分散型ストレージ・システム・マネージャを実装するためのさまざまな代替方法が実施されることが可能であるということ認識されたい。DSUは、物理ストレージ・ユニット、たとえば、回転ディスクまたはソリッド・ステート・ディスクなどのディスクまたはフラッシュ・ベースのストレージ・ユニットに相当する。複数の実施形態によれば、ストレージ・システム・クラスタは、本明細書においてさらに詳述するように、「仮想ボリューム」(vvol: virtual volume)を作成して、コンピュータ・システム100<sub>1</sub>および100<sub>2</sub>などの接続されているコンピュータ・システムに公開する。コンピュータ・システム100内で稼働するアプリケーション(たとえば、自分の仮想ディスクにアクセスするVMなど)は、図2Aの実施形態におけるSCSI、および図2Bの実施形態におけるNFSなど、標準的なプロトコルを使用して、SCSIまたはNFSプロトコル・トラフィックに関する論理エンドポイント(「プロトコル・エンドポイント」(PE)として知られており、ストレージ・システム130内で構成されている)を通じて、オン・デマンドでvvolにアクセスする。アプリケーション関連のデータ・オペレーションに関する、コンピュータ・システム100からストレージ・システム130への通信パスは、本明細書においては「帯域内」パスと呼ばれる。コンピュータ・システム100のホスト・バス・アダプタ(

10

20

30

40

50

H B A ) と、ストレージ・システム 1 3 0 内で構成されている P E との間における通信パス、およびコンピュータ・システム 1 0 0 のネットワーク・インターフェース・カード ( N I C ) と、ストレージ・システム 1 3 0 内で構成されている P E との間における通信パスは、帯域内パスの例である。帯域内ではなく、かつ典型的には、管理オペレーションを実行するために使用される、コンピュータ・システム 1 0 0 からストレージ・システム 1 3 0 への通信パスは、本明細書においては「帯域外」パスと呼ばれる。コンピュータ・システム 1 0 0 と、ストレージ・システム 1 3 0 との間におけるイーサネット・ネットワーク接続など、帯域外パスの例は、図 6 において帯域内パスとは別に示されている。簡単にするために、コンピュータ・システム 1 0 0 は、ストレージ・システム 1 3 0 に直接接続されているように示されている。しかしながら、それらのコンピュータ・システム 1 0 0 は、複数のパス、およびスイッチのうちの 1 つまたは複数を通じてストレージ・システム 1 3 0 に接続されることが可能であるということを理解されたい。

10

## 【 0 0 1 7 】

分散型ストレージ・システム・マネージャ 1 3 5、または単一のストレージ・システム・マネージャ 1 3 1 もしくは 1 3 2 は、(たとえば、コンピュータ・システム 1 0 0 などの要求に応じて、) 物理的な D S U の論理的な集約に相当する論理的な「ストレージ・コンテナ」から v v o l を作成することができる。一般には、1 つのストレージ・コンテナは、複数のストレージ・システムにわたることができ、単一のストレージ・システム・マネージャ、または分散型ストレージ・システム・マネージャによって、多くのストレージ・コンテナが作成されることが可能である。同様に、単一のストレージ・システムは、多くのストレージ・コンテナを含むことができる。図 2 A および図 2 B においては、分散型ストレージ・システム・マネージャ 1 3 5 によって作成されたストレージ・コンテナ 1 4 2<sub>A</sub> は、ストレージ・システム 1 3 0<sub>1</sub> およびストレージ・システム 1 3 0<sub>2</sub> にわたるものとして示されており、その一方で、ストレージ・コンテナ 1 4 2<sub>B</sub> およびストレージ・コンテナ 1 4 2<sub>C</sub> は、単一のストレージ・システム (すなわち、ストレージ・システム 1 3 0<sub>1</sub> およびストレージ・システム 1 3 0<sub>2</sub> それぞれ) の中に含まれているものとして示されている。1 つのストレージ・コンテナは、複数のストレージ・システムにわたることができるため、ストレージ・システム管理者は、ストレージ・システムのいずれか 1 つのストレージ・キャパシティーを超えるストレージ・キャパシティーを自分の顧客にプロビジョンすることができるということ認識されたい。単一のストレージ・システム内に複数のストレージ・コンテナが作成されることが可能であるため、ストレージ・システム管理者は、単一のストレージ・システムを使用して複数の顧客にストレージをプロビジョンすることができるということさら認識されたい。

20

30

## 【 0 0 1 8 】

図 2 A の実施形態においては、それぞれの v v o l は、ブロック・ベースのストレージ・システムからプロビジョンされている。図 2 B の実施形態においては、N A S ベースのストレージ・システムが、D S U 1 4 1 の上にファイル・システム 1 4 5 を実装しており、それぞれの v v o l は、このファイル・システム内の 1 つのファイル・オブジェクトとしてコンピュータ・システム 1 0 0 に公開されている。加えて、以降でさらに詳細に説明するように、コンピュータ・システム 1 0 0 上で稼働するアプリケーションは、P E を通じた I O のために v v o l にアクセスする。たとえば、図 2 A および図 2 B において破線で示されているように、v v o l 1 5 1 および v v o l 1 5 2 は、P E 1 6 1 を介してアクセス可能であり、v v o l 1 5 3 および v v o l 1 5 5 は、P E 1 6 2 を介してアクセス可能であり、v v o l 1 5 4 は、P E 1 6 3 および P E 1 6 4 を介してアクセス可能であり、v v o l 1 5 6 は、P E 1 6 5 を介してアクセス可能である。ストレージ・コンテナ 1 4 2<sub>A</sub> 内の v v o l 1 5 3、およびストレージ・コンテナ 1 4 2<sub>C</sub> 内の v v o l 1 5 5 など、複数のストレージ・コンテナからの v v o l は、任意の所与の時点において P E 1 6 2 などの単一の P E を介してアクセス可能とすることができるということ認識されたい。P E 1 6 6 などの P E は、それらの P E を介してアクセス可能である v v o l がまったくなくても、存在することができるということさら認識されたい。

40

50

## 【 0 0 1 9 】

図 2 A の実施形態においては、ストレージ・システム 1 3 0 は、L U N をセットアップするための知られている方法を使用して、特別なタイプの L U N として P E を実装している。L U N と同様に、ストレージ・システム 1 3 0 は、W W N ( W o r l d W i d e N a m e ) として知られている一意の識別子をそれぞれの P E に提供する。一実施形態においては、P E を作成する際に、ストレージ・システム 1 3 0 は、その特別な L U N に関するサイズを指定しない。なぜなら、本明細書に記載されている P E は、実際のデータ・コンテナではないためである。そのような一実施形態においては、ストレージ・システム 1 3 0 は、P E 関連の L U N のサイズとしてゼロの値または非常に小さな値を割り振ることができ、それによって管理者は、以降でさらに論じるように、ストレージ・システムが L U N (たとえば、従来のデータ L U N および P E 関連の L U N ) のリストを提供するように要求する場合に、P E を迅速に識別することができる。同様に、ストレージ・システム 1 3 0 は、P E に対する L U N に関する識別番号として、2 5 5 よりも大きな数を L U N に割り振って、それらの L U N がデータ L U N ではないということを、人間にとってわかりやすい方法で示すことができる。P E と L U N との間において区別を行うための別の方法として、P E ビットが E x t e n d e d I n q u i r y D a t a V P D ページ ( ページ 8 6 h ) に加えられることが可能である。P E ビットは、L U N が P E である場合には 1 に設定され、L U N が通常のデータ L U N である場合には 0 に設定される。コンピュータ・システム 1 0 0 は、S C S I コマンド R E P O R T \_ L U N S を発行することによって帯域内パスを介して P E を発見することと、示されている P E ビットを調べることによって、それらの P E が、本明細書に記載されている実施形態による P E であるか、または従来のデータ L U N であるかを判定することが可能である。コンピュータ・システム 1 0 0 は、L U N が P E であるか、または従来の L U N であるかをさらに確認するために、L U N のサイズおよび L U N の番号のプロパティを任意選択で検査することができる。P E 関連の L U N を通常のデータ L U N から区別するために上述の技術のうちの任意の技術が使用されることが可能であるということを確認されたい。一実施形態においては、P E ビット技術が、P E 関連の L U N を通常のデータ L U N から区別するために使用される唯一の技術である。

10

20

## 【 0 0 2 0 】

図 2 B の実施形態においては、P E は、F S ボリュームに対してマウント・ポイントをセットアップするための知られている方法を使用してストレージ・システム 1 3 0 内に作成される。図 2 B の実施形態において作成されるそれぞれの P E は、I P アドレスおよびファイル・システム・パスによって一意に識別され、それらの I P アドレスおよびファイル・システム・パスは、従来、合わせて「マウント・ポイント」とも呼ばれている。しかしながら、従来のマウント・ポイントとは異なり、それらの P E は、F S ボリュームに関連付けられない。加えて、図 2 A の P E とは異なり、図 2 B の P E は、仮想ボリュームが所与の P E にバインドされていない限り、帯域内パスを介してコンピュータ・システム 1 0 0 によって発見可能ではない。したがって、図 2 B の P E は、帯域外パスを介してストレージ・システムによって報告される。

30

## 【 0 0 2 1 】

図 3 は、一実施形態による、仮想ボリュームを管理するための図 2 A または図 2 B のストレージ・システム・クラスタのコンポーネントのブロック図である。それらのコンポーネントは、一実施形態におけるストレージ・システム 1 3 0 において実行されるストレージ・システム・マネージャ 1 3 1 および 1 3 2 のソフトウェア・モジュール、または別の実施形態における分散型ストレージ・システム・マネージャ 1 3 5 のソフトウェア・モジュール、すなわち、入力 / 出力 ( I / O ) マネージャ 3 0 4、ボリューム・マネージャ 3 0 6、コンテナ・マネージャ 3 0 8、およびデータ・アクセス・レイヤ 3 1 0 を含む。本明細書における実施形態の説明においては、分散型ストレージ・システム・マネージャ 1 3 5 によって取られるあらゆるアクションは、実施形態に応じてストレージ・システム・マネージャ 1 3 1 またはストレージ・システム・マネージャ 1 3 2 によって取られること

40

50

が可能であるということを理解されたい。

【 0 0 2 2 】

図 3 の例においては、分散型ストレージ・システム・マネージャ 1 3 5 は、3 つのストレージ・コンテナ SC 1、SC 2、および SC 3 を DSU 1 4 1 から作成しており、DSU 1 4 1 のそれぞれは、P 1 から P n とラベル付けされているスピンドル・エクステントを有するように示されている。一般に、それぞれのストレージ・コンテナは、固定された物理サイズを有しており、DSU の特定のエクステントに関連付けられている。図 3 に示されている例においては、分散型ストレージ・システム・マネージャ 1 3 5 は、コンテナ・データベース 3 1 6 へのアクセスを有しており、コンテナ・データベース 3 1 6 は、それぞれのストレージ・コンテナに関して、そのコンテナ ID と、物理的なレイアウトの情報と、何らかのメタデータとを格納している。コンテナ・データベース 3 1 6 は、コンテナ・マネージャ 3 0 8 によって管理および更新され、コンテナ・マネージャ 3 0 8 は、一実施形態においては、分散型ストレージ・システム・マネージャ 1 3 5 のコンポーネントである。コンテナ ID とは、ストレージ・コンテナが作成されたときにそのストレージ・コンテナに与えられる汎用一意識別子である。物理的なレイアウトの情報は、所与のストレージ・コンテナに関連付けられている DSU 1 4 1 のスピンドル・エクステントから構成されており、<システム ID, DSU ID, エクステント番号> の順序付けられたリストとして格納されている。メタデータ・セクションは、何らかの一般的なメタデータ、および何らかのストレージ・システム・ベンダー固有のメタデータを含むことができる。たとえば、メタデータ・セクションは、ストレージ・コンテナにアクセスすることを許可されているコンピュータ・システムまたはアプリケーションまたはユーザの ID を含むことができる。別の例として、メタデータ・セクションは、ストレージ・コンテナのどの<システム ID, DSU ID, エクステント番号> のエクステントが既存の v v o l に既に割り当てられているか、およびどの<システム ID, DSU ID, エクステント番号> のエクステントが空いているかを示すためのアロケーション・ビットマップを含む。一実施形態においては、ストレージ・システム管理者は、別々のビジネス・ユニットに関して別々のストレージ・コンテナを作成することができ、それによって、別々のビジネス・ユニットの v v o l 同士は、同じストレージ・コンテナからはプロビジョンされない。v v o l 同士を分けるためのその他のポリシーが適用されることも可能である。たとえば、ストレージ・システム管理者は、クラウド・サービスの別々の顧客の v v o l 同士が別々のストレージ・コンテナからプロビジョンされるというポリシーを採用することができる。また、v v o l 同士が、自分たちの必要とされるサービス・レベルに従ってストレージ・コンテナからグループ化されてプロビジョンされることも可能である。加えて、ストレージ・システム管理者は、ストレージ・コンテナを作成すること、削除すること、およびその他の形で管理すること、たとえば、作成されることが可能であるストレージ・コンテナの数を定義すること、および、ストレージ・コンテナごとに設定されることが可能である最大の物理サイズを設定することなどが可能である。

【 0 0 2 3 】

また、図 3 の例においては、分散型ストレージ・システム・マネージャ 1 3 5 は、(要求を行っているコンピュータ・システム 1 0 0 のために) 複数の v v o l を、それぞれ別々のストレージ・コンテナからプロビジョンしている。一般に、v v o l は、固定された物理サイズを有することができ、またはシン・プロビジョニングされることが可能であり、それぞれの v v o l は、v v o l ID を有しており、v v o l ID とは、v v o l が作成されたときにその v v o l に与えられる汎用一意識別子である。それぞれの v v o l に関して、v v o l データベース 3 1 4 が、それぞれの v v o l ごとに、その v v o l ID と、その v v o l が作成されているストレージ・コンテナのコンテナ ID と、その v v o l のアドレス空間を含むそのストレージ・コンテナ内の<オフセット, 長さ> の値の順序付けられたリストとを格納している。v v o l データベース 3 1 4 は、ボリューム・マネージャ 3 0 6 によって管理および更新され、ボリューム・マネージャ 3 0 6 は、一実施形態においては、分散型ストレージ・システム・マネージャ 1 3 5 のコンポーネント

である。一実施形態においては、v v o l データベース 3 1 4 は、v v o l に関する少量のメタデータも格納する。このメタデータは、一組のキー/値のペアとして v v o l データベース 3 1 4 内に格納され、その v v o l の存在中はいつでも帯域外バスを介してコンピュータ・システム 1 0 0 によって更新されることおよびクエリーされることが可能である。格納されるキー/値のペアは、3つのカテゴリーへと分かれる。第1のカテゴリーは、よく知られているキーであり、特定のキーの定義（ひいては、それらの値の解釈）は、公に利用可能である。1例は、仮想ボリューム・タイプ（たとえば、仮想マシンの実施形態においては、v v o l が V M のメタデータを含むか、または V M のデータを含むか）に対応するキーである。別の例は、A p p I D であり、これは、v v o l 内にデータを格納したアプリケーションの I D である。第2のカテゴリーは、コンピュータ・システム固有のキーであり、コンピュータ・システムまたはその管理モジュールは、特定のキーおよび値を仮想ボリュームのメタデータとして格納する。第3のカテゴリーは、ストレージ・システム・ベンダー固有のキーであり、これらによって、ストレージ・システム・ベンダーは、仮想ボリュームのメタデータに関連付けられている特定のキーを格納することができる。ストレージ・システム・ベンダーがそのメタデータに関してこのキー/値の格納を使用する1つの理由は、これらのキーのすべてが、v v o l に関する帯域外チャンネルを介してストレージ・システム・ベンダーのプラグインおよびその他の拡張にとって容易に利用可能であることである。キー/値のペアに関する格納オペレーションは、仮想ボリュームの作成およびその他の工程の一部であり、したがって格納オペレーションは、適度に高速になるはずである。ストレージ・システムはまた、特定のキー上で提供される値に対する厳密な一致に基づいて仮想ボリュームの検索を可能にするように構成される。

10

20

#### 【 0 0 2 4 】

I O マネージャ 3 0 4 は、P E と v v o l との間における現在有効な I O 接続パスを記憶する接続データベース 3 1 2 を保持するソフトウェア・モジュール（また、特定の実施形態においては、分散型ストレージ・システム・マネージャ 1 3 5 のコンポーネント）である。図 3 に示されている例においては、7つの現在有効な I O セッションが示されている。それぞれの有効なセッションは、関連付けられている P E I D と、セカンダリー・レベル識別子（S L L I D ）と、v v o l I D と、この I O セッションを通じて I O を実行している別々のアプリケーションの数を示すリファレンス・カウント（R e f C n t ）とを有する。（たとえば、コンピュータ・システム 1 0 0 による要求に応じて）分散型ストレージ・システム・マネージャ 1 3 5 によって P E と v v o l との間における有効な I O セッションを確立する工程は、本明細書においては「バインド」工程と呼ばれる。それぞれのバインドごとに、分散型ストレージ・システム・マネージャ 1 3 5 は、（たとえば、I O マネージャ 3 0 4 を介して）接続データベース 3 1 2 にエントリーを加える。その後分散型ストレージ・システム・マネージャ 1 3 5 によって I O セッションを取り壊す工程は、本明細書においては「アンバインド」工程と呼ばれる。それぞれのアンバインドごとに、分散型ストレージ・システム・マネージャ 1 3 5 は、（たとえば、I O マネージャ 3 0 4 を介して）I O セッションのリファレンス・カウントを1ずつデクリメントする。I O セッションのリファレンス・カウントがゼロである場合には、分散型ストレージ・システム・マネージャ 1 3 5 は、（たとえば、I O マネージャ 3 0 4 を介して）その I O 接続パスに関するエントリーを接続データベース 3 1 2 から削除することができる。前述したように、一実施形態においては、コンピュータ・システム 1 0 0 は、バインド要求およびアンバインド要求を生成して、帯域外バスを介して分散型ストレージ・システム・マネージャ 1 3 5 へ送信する。あるいは、コンピュータ・システム 1 0 0 は、アンバインド要求を生成して、オーバーローディングが存在しているエラー・パスによって帯域内バスを介して送信することができる。一実施形態においては、リファレンス・カウントが 0 から 1 に変わった場合、またはその逆の場合には、世代番号は、単調に増加する番号、またはランダムに生成された番号に変更される。別の実施形態においては、世代番号は、ランダムに生成された番号であり、接続データベース 3 1 2 から R e f C n t 列が取り除かれており、それぞれのバインドごとに、バインド要求が、既にバインドされている v v o

30

40

50

1 に対してなされた場合でさえ、分散型ストレージ・システム・マネージャ 1 3 5 は、（たとえば、I O マネージャ 3 0 4 を介して）接続データベース 3 1 2 にエントリーを加える。

#### 【0025】

図 2 A のストレージ・システム・クラスタにおいては、I O マネージャ 3 0 4 は、接続データベース 3 1 2 を使用して P E を通じて受信されたコンピュータ・システム 1 0 0 からの I O 要求（I O）を処理する。I O が P E のうちの 1 つにおいて受信された場合には、I O マネージャ 3 0 4 は、その I O の対象であった v v o l を特定する目的で、その I O 内に含まれている P E I D および S L L I D を識別するために、その I O を解析する。次いで、接続データベース 3 1 4 にアクセスすることによって、I O マネージャ 3 0 4 は、解析された P E I D および S L L I D に関連付けられている v v o l I D を取り出すことができる。図 3 および後続の図においては、簡単にするために、P E I D は、P E \_ A、P E \_ B などと示されている。一実施形態においては、実際の P E I D は、P E の W W N である。加えて、S L L I D は、S 0 0 0 1、S 0 0 0 2 などと示されている。実際の S L L I D は、接続データベース 3 1 2 内の所与の P E I D に関連付けられている複数の S L L I D のうちの任意の一意の番号として、分散型ストレージ・システム・マネージャ 1 3 5 によって生成される。v v o l I D を有する仮想ボリュームの論理アドレス空間と、D S U 1 4 1 の物理ロケーションとの間におけるマッピングは、v v o l データベース 3 1 4 を使用してボリューム・マネージャ 3 0 6 によって、およびコンテナ・データベース 3 1 6 を使用してコンテナ・マネージャ 3 0 8 によって実行される。D S U 1 4 1 の物理ロケーションが得られると、データ・アクセス・レイヤ 3 1 0（一実施形態においては、やはり分散型ストレージ・システム・マネージャ 1 3 5 のコンポーネント）が、これらの物理ロケーション上で I O を実行する。

10

20

#### 【0026】

図 2 B のストレージ・システム・クラスタにおいては、I O は、P E を通じて受信され、そのようなそれぞれの I O は、N F S ハンドル（または類似のファイル・システム・ハンドル）を含み、その N F S ハンドルに対して、その I O は発行されている。一実施形態においては、そのようなシステムのための接続データベース 3 1 2 は、P E I D としてストレージ・システムの N F S インターフェースの I P アドレスを、および S L L I D としてファイル・システム・パスを含む。S L L I D は、ファイル・システム 1 4 5 内の v v o l のロケーションに基づいて生成される。v v o l の論理アドレス空間と、D S U 1 4 1 の物理ロケーションとの間におけるマッピングは、v v o l データベース 3 1 4 を使用してボリューム・マネージャ 3 0 6 によって、およびコンテナ・データベース 3 1 6 を使用してコンテナ・マネージャ 3 0 8 によって実行される。D S U 1 4 1 の物理ロケーションが得られると、データ・アクセス・レイヤが、これらの物理ロケーション上で I O を実行する。図 2 B のストレージ・システムに関しては、コンテナ・データベース 3 1 2 は、所与の v v o l に関してコンテナ・ロケーション・エントリー内にファイル：<オフセット，長さ>エントリーの順序付けられたリストを含むことができる（すなわち、v v o l は、ファイル・システム 1 4 5 内に格納されている複数のファイル・セグメントから構成されることが可能である）ということ認識されたい。

30

40

#### 【0027】

一実施形態においては、接続データベース 3 1 2 は、揮発性メモリ内に保持され、その一方で、v v o l データベース 3 1 4 およびコンテナ・データベース 3 1 6 は、D S U 1 4 1 などの永続ストレージ内に保持される。その他の実施形態においては、データベース 3 1 2、3 1 4、3 1 6 のすべてが永続ストレージ内に保持されることが可能である。

#### 【0028】

図 4 は、ストレージ・コンテナを作成するための方法工程 4 1 0 の流れ図である。一実施形態においては、これらの工程は、ストレージ管理者のコントロールのもとで、ストレージ・システム・マネージャ 1 3 1、ストレージ・システム・マネージャ 1 3 2、または分散型ストレージ・システム・マネージャ 1 3 5 によって実行される。上述したように、

50

ストレージ・コンテナは、物理的なD S Uの論理的な集約に相当し、複数のストレージ・システムからの物理的なD S Uにわたることができる。工程4 1 1において、ストレージ管理者は、(分散型ストレージ・システム・マネージャ1 3 5などを介して、)ストレージ・コンテナの物理的なキャパシティーを設定する。クラウドまたはデータ・センター内では、この物理的なキャパシティーは、たとえば、顧客によってリースされる物理ストレージの量に相当することができる。本明細書において開示されているストレージ・コンテナによって提供される柔軟性として、別々の顧客の複数のストレージ・コンテナが、1人のストレージ管理者によって同じストレージ・システムからプロビジョンされることが可能であり、また、たとえば、いずれか1つのストレージ・デバイスの物理的なキャパシティーが、顧客によって要求されているサイズを満たすのに十分ではないケースにおいて、または1つのv v o lの物理ストレージ・フットプリントが、必然的に複数のストレージ・システムにわたることになる複製などのケースにおいて、単一の顧客のための1つのストレージ・コンテナが、複数のストレージ・システムからプロビジョンされることが可能である。工程4 1 2において、ストレージ管理者は、ストレージ・コンテナにアクセスするための許可レベルを設定する。マルチテナント・データ・センターにおいては、たとえば、顧客は、自分にリースされているストレージ・コンテナにアクセスすることしかできない。工程4 1 3において、分散型ストレージ・システム・マネージャ1 3 5は、ストレージ・コンテナに関する一意の識別子を生成する。次いで工程4 1 4において、分散型ストレージ・システム・マネージャ1 3 5は、(たとえば、一実施形態においては、コンテナ・マネージャ3 0 8を介して、)D S U 1 4 1の空いているスピンドル・エクステントを、工程4 1 1において設定された物理的なキャパシティーを満たすのに十分な量でストレージ・コンテナに割り当てる。上述したように、いずれか1つのストレージ・システムの空きスペースが、物理的なキャパシティーを満たすのに十分ではないケースにおいては、分散型ストレージ・システム・マネージャ1 3 5は、複数のストレージ・システムからD S U 1 4 1のスピンドル・エクステントを割り当てることができる。パーティションが割り当てられた後に、分散型ストレージ・システム・マネージャ1 3 5は、(たとえば、コンテナ・マネージャ3 0 8を介して、)一意のコンテナIDと、<システム番号, D S U ID, エクステント番号>の順序付けられたリストと、ストレージ・コンテナにアクセスすることを許可されているコンピュータ・システムのコンテキストIDとでコンテナ・データベース3 1 6を更新する。

10

20

30

#### 【0 0 2 9】

本明細書に記載されている実施形態によれば、ストレージ能力プロファイル、たとえば、S L Aまたはサービス品質(Q o S : q u a l i t y o f s e r v i c e)は、v v o lごとに(たとえば、要求を行っているコンピュータ・システム1 0 0のために)分散型ストレージ・システム・マネージャ1 3 5によって構成されることが可能である。したがって、別々のストレージ能力プロファイルを有する複数のv v o lが、同じストレージ・コンテナの一部であることが可能である。一実施形態においては、システム管理者は、ストレージ・コンテナの作成時に、新たに作成されたv v o lに関するデフォルトのストレージ能力プロファイル(または、複数の可能なストレージ能力プロファイル)を定義して、コンテナ・データベース3 1 6のメタデータ・セクション内に格納した。ストレージ・コンテナ内で作成されている新たなv v o lに関してストレージ能力プロファイルが明示的に指定されない場合には、その新たなv v o lは、そのストレージ・コンテナに関連付けられているデフォルトのストレージ能力プロファイルを引き継ぐことになる。

40

#### 【0 0 3 0】

図5 Aは、図2 Aのストレージ・システム・クラスタ上にホストされる仮想ボリュームを実装するように構成されているコンピュータ・システムの一実施形態のブロック図である。コンピュータ・システム1 0 1は、従来の、典型的にはサーバクラスである、ハードウェア・プラットフォーム5 0 0上に構築されることが可能であり、ハードウェア・プラットフォーム5 0 0は、1つまたは複数の中央処理装置(C P U)5 0 1と、メモリ5 0 2と、1つまたは複数のネットワーク・インターフェース・カード(N I C)5 0 3と、

50

1つまたは複数のホスト・バス・アダプタ(HBA)504とを含む。HBA504は、コンピュータ・システム101が、ストレージ・デバイス130内に構成されているPEを通じて仮想ボリュームにIOを発行することを可能にする。図5Aにおいてさらに示されているように、オペレーティング・システム508は、ハードウェア・プラットフォーム500の上にインストールされており、複数のアプリケーション512<sub>1</sub>~512<sub>N</sub>が、オペレーティング・システム508の上で実行される。オペレーティング・システム508の例としては、よく知られているコモディティ・オペレーティング・システム、たとえばMicrosoft Windows、Linuxなどのうちの任意のものが含まれる。

#### 【0031】

本明細書に記載されている実施形態によれば、それぞれのアプリケーション512は、自分に関連付けられている1つまたは複数のvvolを有しており、アプリケーション512によるオペレーティング・システム508への「CREATE DEVICE」コールに従ってオペレーティング・システム508によって作成されたvvolのブロック・デバイス・インスタンスにIOを発行する。ブロック・デバイス名と、vvol IDとの間における関連付けは、ブロック・デバイス・データベース533内に保持される。アプリケーション512<sub>2</sub>~512<sub>N</sub>からのIOは、ファイル・システム・ドライバ510によって受信され、ファイル・システム・ドライバ510は、それらのIOをブロックIOに変換し、それらのブロックIOを仮想ボリューム・デバイス・ドライバ532に提供する。その一方で、アプリケーション512<sub>1</sub>からのIOは、ファイル・システム・ドライバ510を迂回するように示されており、仮想ボリューム・デバイス・ドライバ532に直接提供され、これが意味するのは、アプリケーション512<sub>1</sub>が、自分のブロック・デバイスに直接、ロー・ストレージ・デバイス(raw storage device)として、たとえば、データベース・ディスク、ログ・ディスク、バックアップ・アーカイブ、およびコンテンツ・リポジトリとして、「Providing Access to a Raw Data Storage Unit in a Computer System」と題されている米国特許第7,155,558号(その全内容を本願明細書に援用する)に記載されている様式でアクセスするということである。仮想ボリューム・デバイス・ドライバ532は、ブロックIOを受信した場合には、ブロック・デバイス・データベース533にアクセスして、そのIO内で指定されているブロック・デバイス名と、そのブロック・デバイス名に関連付けられているvvolへのIO接続バスを定義するPE ID(PE LUNのWWN)およびSLIDとの間におけるマッピングを参照する。ここで示されている例においては、「archive」というブロック・デバイス名は、アプリケーション512<sub>1</sub>に関して作成されたvvol12のブロック・デバイス・インスタンスに対応しており、「foo」、「database」、および「log」というブロック・デバイス名は、アプリケーション512<sub>2</sub>~512<sub>N</sub>のうちの1つまたは複数に関してそれぞれ作成されたvvol1、vvol16、およびvvol17のブロック・デバイス・インスタンスに対応する。ブロック・デバイス・データベース533内に格納されているその他の情報としては、ブロック・デバイスがアクティブであるか否かを示すそれぞれのブロック・デバイスに関するアクティブ・ビット値と、CIF(commands-in-flight:処理中コマンド)値とが含まれる。「1」というアクティブ・ビットは、IOがブロック・デバイスに発行されることが可能であるということの意味する。「0」というアクティブ・ビットは、ブロック・デバイスが非アクティブであり、IOがブロック・デバイスに発行されることは不可能であるということの意味する。CIF値は、いくつかのIOが処理中であるか、すなわち、発行されたが完了されていないかの表示を提供する。ここで示されている例においては、「foo」というブロック・デバイスは、アクティブであり、いくつかの処理中コマンドを有している。「archive」というブロック・デバイスは、非アクティブであり、さらに新しいコマンドを受け入れないであろう。しかしながら、このブロック・デバイスは、2つの処理中コマンドが完了するのを待っている。「database」というブロック・デバイスは、非アクティブであり、未処理のコマ

10

20

30

40

50



ンドはない。最後に、「log」というブロック・デバイスは、アクティブであるが、アプリケーションは現在、このデバイスに対する未処理のI/Oを有していない。仮想ボリューム・デバイス・ドライバ532は、いつでも自分のデータベース533からこれらのようなデバイスを除去することを選択することができる。

#### 【0032】

上述のマッピングを実行することに加えて、仮想ボリューム・デバイス・ドライバ532は、データ・アクセス・レイヤ540にロー・ブロックレベルI/O (raw block level I/O)を発行する。データ・アクセス・レイヤ540は、コマンド・キューイングおよびスケジューリング・ポリシーをロー・ブロックレベルI/Oに適用するデバイス・アクセス・レイヤ534と、プロトコルに準拠したフォーマットでロー・ブロックレベルI/Oをフォーマットして、それらのロー・ブロックレベルI/Oを、帯域内バスを介してPEへ転送するためにHBA504に送信する、HBA504のためのデバイス・ドライバ536とを含む。SCSIプロトコルが使用される実施形態においては、vvol情報は、SAM-5 (SCSI Architecture Model-5)において指定されているように、SCSI LUNデータ・フィールド (これは、8バイト構造である) 内にエンコードされる。PE IDは、最初の2バイト (これは、従来はLUN ID用に使用されている) 内にエンコードされ、vvol情報、とりわけSLCIDは、残っている6バイト (の一部) を利用して、SCSIセカンド・レベルLUN ID内にエンコードされる。

10

#### 【0033】

図5Aにおいてさらに示されているように、データ・アクセス・レイヤ540はまた、ストレージ・システムから帯域内バスを通じて受信されるI/Oエラーを取り扱うためのエラー・ハンドリング・ユニット542を含む。一実施形態においては、エラー・ハンドリング・ユニット542によって受信されたI/Oエラーは、I/Oマネージャ304によってPEを通じて伝搬される。I/Oエラー・クラスの例としては、コンピュータ・システム101とPEとの間におけるバス・エラーと、PEエラーと、vvolエラーとが含まれる。エラー・ハンドリング・ユニット542は、検知されたすべてのエラーを上述のクラスへと分類する。PEへのバス・エラーに出くわし、PEへの別のバスが存在する場合には、データ・アクセス・レイヤ540は、PEへの別のバスに沿ってI/Oを送信する。I/OエラーがPEエラーである場合には、エラー・ハンドリング・ユニット542は、PEを通じてI/Oを発行しているそれぞれのブロック・デバイスに関するエラー状況を示すために、ブロック・デバイス・データベース533を更新する。I/Oエラーがvvolエラーである場合には、エラー・ハンドリング・ユニット542は、vvolに関連付けられているそれぞれのブロック・デバイスに関するエラー状況を示すために、ブロック・デバイス・データベース533を更新する。エラー・ハンドリング・ユニット542は、アラームまたはシステム・イベントを発行することもでき、それによって、エラー状況を有するブロック・デバイスへのさらなるI/Oは、拒否されることになる。

20

30

#### 【0034】

図5Bは、図2Aのストレージ・システム・クラスタの代わりに図2Bのストレージ・システム・クラスタとインターフェースを取るように構成されている図5Aのコンピュータ・システムのブロック図である。この実施形態においては、データ・アクセス・レイヤ540は、NFSクライアント545と、NIC503のためのデバイス・ドライバ546とを含む。NFSクライアント545は、ブロック・デバイス名を、PE ID (NASストレージ・システムのIPアドレス) と、ブロック・デバイスに対応するNFSファイル・ハンドルであるSLCIDとにマップする。このマッピングは、図5Bにおいて示されているように、ブロック・デバイス・データベース533内に格納される。「アクティブ」および「CIF」の列は、依然として存在するが、図5Bに示されているブロック・デバイス・データベース533においては示されていないということに留意されたい。以降で説明するように、NFSファイル・ハンドルは、NASストレージ・システム内のファイル・オブジェクトを一意に識別し、バインド工程中に生成されることが可能である

40

50

。あるいは、`vvol`をバインドしたいという要求に回答して、`NAS`ストレージ・システムは、`PE ID`および`SLID`を返し、通常の帯域内メカニズム（たとえば、`ルックアップ`または`readdirplus`）を使用して`vvol`が開かれると、`NFS`ファイル・ハンドルが与えられることになる。`NFS`クライアント`545`はまた、仮想ボリューム・デバイス・ドライバ`532`から受信されたロー・ブロックレベル`IO`を`NFS`ファイルベースの`IO`に変換する。次いで、`NIC503`のためのデバイス・ドライバ`546`は、プロトコルに準拠したフォーマットで`NFS`ファイルベースの`IO`をフォーマットして、それらの`NFS`ファイルベースの`IO`を、帯域内パスを介して`PE`のうちの1つへ転送するために、`NFS`ハンドルとともに、`NIC503`へ送信する。

#### 【0035】

図5Cは、仮想ボリュームを実装するように構成されているコンピュータ・システムの別の実施形態のブロック図である。この実施形態においては、コンピュータ・システム102は、仮想化ソフトウェア（ここでは、ハイパーバイザ560として示されている）を伴って構成されている。ハイパーバイザ560は、ハードウェア・プラットフォーム550の上にインストールされており、ハードウェア・プラットフォーム550は、`CPU551`、`メモリ552`、`NIC553`、および`HBA554`を含み、仮想マシン実行スペース570をサポートし、仮想マシン実行スペース570内では、複数の仮想マシン（`VM`）`5711 ~ 571N`が、同時にインスタンス化されて実行されることが可能である。1つまたは複数の実施形態においては、ハイパーバイザ560および仮想マシン`571`は、`VMware, Inc.` [米国カリフォルニア州パロアルト（`Palo Alto`）所在]によって販売されている`VMware vSphere`（登録商標）製品を使用して実装される。それぞれの仮想マシン`571`は、仮想ハードウェア・プラットフォーム573を実装しており、仮想ハードウェア・プラットフォーム573は、アプリケーション579を実行することができるゲスト・オペレーティング・システム（`OS`）`572`のインストレーションをサポートする。ゲスト`OS572`の例としては、よく知られている`コモディティ・オペレーティング・システム`、たとえば`Microsoft Windows`、`Linux`などのうちの任意のものが含まれる。それぞれのインスタンスにおいて、ゲスト`OS572`は、ネイティブ・ファイル・システム・レイヤ（図5Cにおいては示されていない）、たとえば、`NTFS`または`ext3FS`タイプのファイル・システム・レイヤのいずれかを含む。これらのファイル・システム・レイヤは、仮想ハードウェア・プラットフォーム573とインターフェースを取って、ゲスト`OS572`の視点からは、データ・ストレージ`HBA`にアクセスするが、このデータ・ストレージ`HBA`は、実際には、ゲスト`OS572`の実行を可能にするためにディスク・ストレージ・サポートの外見（実際には、仮想ディスク、すなわち仮想ディスク`575A ~ 575X`）を提供する仮想ハードウェア・プラットフォーム573によって実装されている仮想`HBA574`である。特定の実施形態においては、仮想ディスク`575A ~ 575X`は、ゲスト`OS572`の視点からは、仮想マシンに接続するための`SCSI`標準、または、`IDE`、`ATA`、および`ATAPI`を含む、当技術分野における標準的な技術者に知られているその他の任意の適切なハードウェア接続インターフェース標準をサポートするように見えることが可能である。ゲスト`OS572`の視点からは、ファイル・システム関連のデータ転送およびコントロール・オペレーションを実施するためにそのようなゲスト`OS572`によって開始されたファイル・システム・コールは、最終的な実行のために仮想ディスク`575A ~ 575X`へ回送されるように見えるが、実際には、そのようなコールは、処理され、仮想`HBA574`を通じて補助的な仮想マシン・モニタ（`VMM`）`5611 ~ 561N`に渡され、`VMM5611 ~ 561N`は、ハイパーバイザ560とのオペレーションを調整するために必要とされる仮想システム・サポートを実施する。とりわけ、`HBAエミュレータ562`は、データ転送およびコントロール・オペレーションがハイパーバイザ560によって正しく取り扱われることを機能的に可能にし、ハイパーバイザ560は最終的に、そのようなオペレーションを、自分のさまざまなレイヤを通じて、ストレージ・システム130へ接続している`HBA554`に渡す。

10

20

30

40

50

## 【0036】

本明細書に記載されている実施形態によれば、それぞれのVM571は、自分に関連付けられている1つまたは複数のvvolを有しており、VM571によるハイパーバイザ560への「CREATE DEVICE」コールに従ってハイパーバイザ560によって作成されたvvolのブロック・デバイス・インスタンスにIOを発行する。ブロック・デバイス名と、vvol IDとの間における関連付けは、ブロック・デバイス・データベース580内に保持される。VM571<sub>2</sub>～571<sub>N</sub>からのIOは、SCSI仮想化レイヤ563によって受信され、SCSI仮想化レイヤ563は、それらのIOを、仮想マシン・ファイル・システム(VMFS)ドライバ564によって理解されるファイルIOへと変換する。次いで、VMFSドライバ564は、それらのファイルIOをブロックIOに変換し、それらのブロックIOを仮想ボリューム・デバイス・ドライバ565に提供する。その一方で、VM571<sub>1</sub>からのIOは、VMFSドライバ564を迂回するように示されており、仮想ボリューム・デバイス・ドライバ565に直接提供され、これが意味するのは、VM571<sub>1</sub>が、自分のブロック・デバイスに直接、ロー・ストレージ・デバイスとして、たとえば、データベース・ディスク、ログ・ディスク、バックアップ・アーカイブ、およびコンテンツ・リポジトリとして、米国特許第7,155,558号に記載されている様式でアクセスするということである。

10

## 【0037】

仮想ボリューム・デバイス・ドライバ565は、ブロックIOを受信した場合には、ブロック・デバイス・データベース580にアクセスして、そのIO内で指定されているブロック・デバイス名と、そのブロック・デバイス名に関連付けられているvvolへのIOセッションを定義するPE IDおよびSLIDとの間におけるマッピングを参照する。ここで示されている例においては、「database」および「log」というブロック・デバイス名は、VM571<sub>1</sub>に関してそれぞれ作成されたvvol1およびvvol4のブロック・デバイス・インスタンスに対応しており、「vmdk2」、「vmdkn」および「snapn」というブロック・デバイス名は、VM571<sub>2</sub>～571<sub>N</sub>のうちの1つまたは複数に関してそれぞれ作成されたvvol12、vvol16、およびvvol17のブロック・デバイス・インスタンスに対応する。ブロック・デバイス・データベース580内に格納されているその他の情報としては、ブロック・デバイスがアクティブであるか否かを示すそれぞれのブロック・デバイスに関するアクティブ・ビット値と、CIF(commands-in-flight:処理中コマンド)値とが含まれる。「1」というアクティブ・ビットは、IOがブロック・デバイスに発行されることが可能であるということの意味する。「0」というアクティブ・ビットは、ブロック・デバイスが非アクティブであり、IOがブロック・デバイスに発行されることは不可能であるということの意味する。CIF値は、いくつのIOが処理中であるか、すなわち、発行されたが完了されていないかの表示を提供する。

20

30

## 【0038】

上述のマッピングを実行することに加えて、仮想ボリューム・デバイス・ドライバ565は、データ・アクセス・レイヤ566にロー・ブロックレベルIOを発行する。データ・アクセス・レイヤ566は、コマンド・キューイングおよびスケジューリング・ポリシーをロー・ブロックレベルIOに適用するデバイス・アクセス・レイヤ567と、プロトコルに準拠したフォーマットでロー・ブロックレベルIOをフォーマットして、それらのロー・ブロックレベルIOを、帯域内バスを介してPEへ転送するためにHBA554に送信する、HBA554のためのデバイス・ドライバ568とを含む。SCSIプロトコルが使用される実施形態においては、vvol情報は、SAM-5(SCSI Architecture Model-5)において指定されているように、SCSI LUNデータ・フィールド(これは、8バイト構造である)内にエンコードされる。PE IDは、最初の2バイト(これは、従来はLUN ID用に使用されている)内にエンコードされ、vvol情報、とりわけSLIDは、残っている6バイト(の一部)を利用して、SCSIセカンド・レベルLUN ID内にエンコードされる。図5Cにおいてさらに

40

50

示されているように、データ・アクセス・レイヤ 5 6 6 はまた、エラー・ハンドリング・ユニット 5 6 9 を含み、エラー・ハンドリング・ユニット 5 6 9 は、エラー・ハンドリング・ユニット 5 4 2 と同じ様式で機能する。

【 0 0 3 9 】

図 5 D は、図 2 A のストレージ・システム・クラスタの代わりに図 2 B のストレージ・システム・クラスタとインターフェースを取るように構成されている図 5 C のコンピュータ・システムのブロック図である。この実施形態においては、データ・アクセス・レイヤ 5 6 6 は、NFS クライアント 5 8 5 と、NIC 5 5 3 のためのデバイス・ドライバ 5 8 6 とを含む。NFS クライアント 5 8 5 は、ブロック・デバイス名を、PE ID (IP アドレス) と、ブロック・デバイスに対応する S L L I D (NFS ファイル・ハンドル) とにマップする。このマッピングは、図 5 D において示されているように、ブロック・デバイス・データベース 5 8 0 内に格納される。「アクティブ」および「CIF」の列は、依然として存在するが、図 5 D に示されているブロック・デバイス・データベース 5 8 0 においては示されていないということに留意されたい。以降で説明するように、NFS ファイル・ハンドルは、NFS 内のファイル・オブジェクトを一意に識別し、一実施形態においてはバインド工程中に生成される。NFS クライアント 5 8 5 はまた、仮想ボリューム・デバイス・ドライバ 5 6 5 から受信されたロー・ブロックレベル I/O を NFS ファイルベースの I/O に変換する。次いで、NIC 5 5 3 のためのデバイス・ドライバ 5 8 6 は、プロトコルに準拠したフォーマットで NFS ファイルベースの I/O をフォーマットして、それらの NFS ファイルベースの I/O を、帯域内パスを介して PE のうちの 1 つへ転送するために、NFS ハンドルとともに、NIC 5 5 3 へ送信する。

10

20

【 0 0 4 0 】

図 5 A ~ 図 5 D におけるコンポーネントを説明するために使用されているさまざまな用語、レイヤ、および分類は、それらの機能、または本発明の趣旨もしくは範囲から逸脱することなく、別の形で参照されることが可能であるということに認識されたい。たとえば、VMM 5 6 1 は、VM 5 7 1 とハイパーバイザ 5 6 0 との間における別個の仮想化コンポーネントとみなされることが可能である (ハイパーバイザ 5 6 0 は、そのような概念においては、それ自体が仮想化「カーネル」コンポーネントとみなされることが可能である)。なぜなら、それぞれのインスタンス化された VM ごとに別々の VMM が存在するためである。あるいは、それぞれの VMM 5 6 1 は、その VMM 5 6 1 の対応する仮想マシンのコンポーネントであるとみなされることが可能である。なぜなら、そのような VMM は、その仮想マシンに関するハードウェア・エミュレーション・コンポーネントを含むためである。そのような代替概念においては、たとえば、仮想ハードウェア・プラットフォーム 5 7 3 として記載されている概念レイヤは、VMM 5 6 1 と、および VMM 5 6 1 内へ合併されることが可能であり、それによって、仮想ホスト・バス・アダプタ 5 7 4 は、図 5 C および図 5 D から除去される (すなわち、その仮想ホスト・バス・アダプタ 5 7 4 の機能は、ホスト・バス・アダプタ・エミュレータ 5 6 2 によって実施されるためである)。

30

【 0 0 4 1 】

図 6 は、本発明の一実施形態による、vvol を管理するために使用されるコンポーネントおよび通信パスを示すコンピュータ環境の簡略化されたブロック図である。前述したように、I/O プロトコル・トラフィックのための通信パスは、帯域内パスと呼ばれ、図 6 においては破線 6 0 1 として示されており、破線 6 0 1 は、コンピュータ・システムのデータ・アクセス・レイヤ 5 4 0 を (コンピュータ・システムにおいて提供されている HBA または NIC を通じて)、ストレージ・システム 1 3 0 において構成されている 1 つまたは複数の PE と接続している。vvol を管理するために使用される通信パスは、帯域外パス (前に定義したように、「帯域内」ではないパス) であり、図 6 においては実線 6 0 2 として示されている。本明細書に記載されている実施形態によれば、vvol は、管理サーバ 6 1 0 において提供されているプラグイン 6 1 2、および / またはコンピュータ・システム 1 0 3 のそれぞれにおいて提供されているプラグイン 6 2 2 を通じて管理され

40

50

ることが可能であり、図 6 においては、コンピュータ・システム 103 のうちの 1 つのみが示されている。ストレージ・デバイス側では、管理インターフェース 625 が、ストレージ・システム・マネージャ 131 によって構成されており、管理インターフェース 626 が、ストレージ・システム・マネージャ 132 によって構成されている。加えて、管理インターフェース 624 が、分散型ストレージ・システム・マネージャ 135 によって構成されている。それぞれの管理インターフェースは、プラグイン 612、622 と通信する。管理コマンドの発行および取り扱いを容易にするために、特別なアプリケーション・プログラミング・インターフェース (API) が開発されている。一実施形態においては、プラグイン 612、622 の両方が、特定のストレージ・システム・ベンダーからのストレージ・ハードウェアと通信するようにカスタマイズされているということを認識されたい。したがって、管理サーバ 610 およびコンピュータ・システム 103 は、別々のストレージ・システム・ベンダーに関するストレージ・ハードウェアと通信する場合には、別々のプラグインを採用することになる。別の実施形態においては、任意のベンダーの管理インターフェースと対話する単一のプラグインが存在することができる。これは、ストレージ・システム・マネージャが、(たとえば、コンピュータ・システムおよび/または管理サーバによって発行されているという理由で、) よく知られているインターフェースに合わせてプログラムされることを必要とすることになる。

10

#### 【0042】

管理サーバ 610 はさらに、コンピュータ・システムを管理するためのシステム・マネージャ 611 を伴って構成されている。一実施形態においては、コンピュータ・システムは、仮想マシンを実行しており、システム・マネージャ 611 は、コンピュータ・システムにおいて稼働している仮想マシンを管理する。仮想マシンを管理するシステム・マネージャ 611 の 1 例は、VMware, Inc. によって販売されている vSphere (登録商標) 製品である。示されているように、システム・マネージャ 611 は、コンピュータ・システム 103 からリソース使用レポートを受信するために、およびコンピュータ・システム 103 において稼働しているアプリケーション上でさまざまな管理オペレーションを開始するために、(管理サーバ 610 およびコンピュータ・システム 103 の両方において適切なハードウェア・インターフェースを通じて、) コンピュータ・システム 103 において稼働しているホスト・デーモン (hostd) 621 と通信する。

20

30

#### 【0043】

図 7 は、認証関連の API を使用して図 2 A または図 2 B のストレージ・システム・クラスタに対してコンピュータ・システムを認証するための方法工程の流れ図である。これらの方法工程は、コンピュータ・システムが自分のセキュア・ソケット・レイヤ (SSL) 証明書をストレージ・システムに送信することによって認証を要求したときに、開始される。工程 710 において、ストレージ・システムは、認証クレデンシャル (たとえば、ユーザ名およびパスワード) を求めるプロンプトを、認証を要求しているコンピュータ・システムに発行する。工程 712 において認証クレデンシャルを受信すると、ストレージ・システムは、工程 714 において、それらの認証クレデンシャルを、格納されているクレデンシャルと比較する。正しいクレデンシャルが提供された場合には、ストレージ・システムは、認証されたコンピュータ・システムの SSL 証明書をキー・ストア内に格納する (工程 716)。正しくないクレデンシャルが提供された場合には、ストレージ・システムは、SSL 証明書を無視して、適切なエラー・メッセージを返す (工程 718)。認証された後に、コンピュータ・システムは、SSL リンクを介してストレージ・システムに管理コマンドを発行するために API を呼び出すことができ、また、SSL 証明書内に含まれている一意のコンテキスト ID が、どのコンピュータ・システムがどのストレージ・コンテナにアクセスすることができるかを定義することなどの特定のポリシーを実施するために、ストレージ・システムによって使用される。いくつかの実施形態においては、コンピュータ・システムのコンテキスト ID が、それらのコンピュータ・システムに与えられた許可を管理する際に使用されることが可能である。たとえば、ホスト・コンピュー

40

50

タは、v v o l を作成することを許可されることが可能であるが、そのv v o l を削除すること、もしくはそのv v o l をスナップショットすることを許可されないことが可能であり、またはホスト・コンピュータは、v v o l のスナップショットを作成することを許可されることが可能であるが、そのv v o l をクローンすることを許可されないことが可能である。加えて、許可は、認証されたコンピュータ・システムにログインされるユーザのユーザレベル特権に従って変わることが可能である。

#### 【0044】

図8は、仮想ボリューム作成APIコマンドを使用して仮想ボリュームを作成するための方法工程の流れ図である。一実施形態においては、コンピュータ・システム103は、最小IOPSおよび平均待ち時間など、特定のサイズおよびストレージ能力プロファイルを有するv v o l を作成したいという要求を自分のアプリケーションのうちの1つから工程802において受信した場合には、帯域外バス602を介してストレージ・システムに仮想ボリューム作成APIコマンドを発行する。それに応じて、コンピュータ・システム103は、工程804において、1つのストレージ・コンテナを(コンピュータ・システム103および要求を行っているアプリケーションがアクセスすることを許可されていて、かつ要求に対応するのに十分な空きキャパシティーを有するストレージ・コンテナの中から)選択し、プラグイン622を介してストレージ・システムに仮想ボリューム作成APIコマンドを発行する。このAPIコマンドは、ストレージ・コンテナIDと、v v o l サイズと、v v o l のストレージ能力プロファイルとを含む。別の実施形態においては、このAPIコマンドは、一組のキー/値のペアを含み、アプリケーションは、ストレージ・システムが、その一組のキー/値のペアを、新たに作成されたv v o l とともに格納することを必要とする。別の実施形態においては、管理サーバ610は、帯域外バス602を介してストレージ・システムに仮想ボリューム作成APIコマンドを(プラグイン612を介して)発行する。

#### 【0045】

工程806において、ストレージ・システム・マネージャは、管理インターフェース(たとえば、管理インターフェース624、625、または626)を介して、v v o l を生成したいという要求を受信し、コンテナ・データベース316内の選択されたストレージ・コンテナのメタデータ・セクションにアクセスして、コンピュータ・システム103とアプリケーションとを含む要求コンテキストが、その選択されたストレージ・コンテナ内にv v o l を作成するのに十分な許可を有していることを確かめる。一実施形態においては、許可レベルが十分でない場合には、エラー・メッセージがコンピュータ・システム103に返される。許可レベルが十分である場合には、工程810において、一意のv v o l IDが生成される。次いで工程812において、ストレージ・システム・マネージャは、コンテナ・データベース316のメタデータ・セクション内のアロケーション・ビットマップをスキャンして、選択されたストレージ・コンテナの空いているパーティションを特定する。ストレージ・システム・マネージャは、選択されたストレージ・コンテナの空いているパーティションを、要求されているv v o l サイズに対応するのに十分なだけ割り当て、コンテナ・データベース316のストレージ・コンテナのメタデータ・セクション内のアロケーション・ビットマップを更新する。ストレージ・システム・マネージャはまた、新たなv v o l エントリでv v o l データベース314を更新する。新たなv v o l エントリは、工程810において生成されたv v o l IDと、新たに割り当てられたストレージ・コンテナ・エクステンツの順序付けられたリストと、キー/値のペアとして表されている新たなv v o l のメタデータとを含む。次いで工程814において、ストレージ・システム・マネージャは、v v o l IDをコンピュータ・システム103へ送信する。工程816において、コンピュータ・システム103は、v v o l IDを、v v o l の作成を要求したアプリケーションに関連付ける。一実施形態においては、それぞれのアプリケーションに関して、1つまたは複数のv v o l 記述子ファイルが保持され、v v o l の作成を要求したアプリケーションに関して保持されているv v o l 記述子ファイル内にv v o l IDが書き込まれる。

10

20

30

40

50

## 【 0 0 4 6 】

図 2 A および図 2 B において示されているように、すべての v v o l が P E に接続されているわけではない。P E に接続されていない v v o l は、対応するアプリケーションによって発行された I O に気づかない。なぜなら、その v v o l には I O セッションが確立されていないためである。I O が v v o l に発行されることが可能になる前に、その v v o l はバインド工程を経て、その結果として、その v v o l は特定の P E にバインドされることになる。v v o l が P E にバインドされると、その v v o l がその P E からアンバインドされるまで、I O がその v v o l に発行されることが可能である。

## 【 0 0 4 7 】

一実施形態においては、バインド要求は、コンピュータ・システム 1 0 3 によって、バインド仮想ボリューム A P I を使用して、帯域外パス 6 0 2 を介してストレージ・システムに発行される。バインド要求は、( v v o l I D を使用して、 ) バインドされることになる v v o l を識別し、それに応じてストレージ・システムは、その v v o l を、コンピュータ・システム 1 0 3 が帯域内パスを介して接続される P E にバインドする。図 9 A は、コンピュータ・システムが帯域内パスを介して接続される P E をそのコンピュータ・システムが発見するための方法工程の流れ図である。S C S I プロトコルベースのストレージ・デバイスにおいて構成されている P E は、標準的な S C S I コマンド、R E P O R T \_ L U N S を使用して、帯域内パスを介して発見される。N F S プロトコルベースのストレージ・デバイスにおいて構成されている P E は、A P I を使用して、帯域外パスを介して発見される。図 9 A の方法工程は、コンピュータ・システムによって、それぞれの接続されているストレージ・システムごとに実行される。

## 【 0 0 4 8 】

工程 9 1 0 において、コンピュータ・システムは、接続されているストレージ・システムが S C S I プロトコルベースであるか、または N F S プロトコルベースであるかを判定する。ストレージ・システムが S C S I プロトコルベースである場合には、S C S I コマンド、R E P O R T \_ L U N S が、帯域内のコンピュータ・システムによってストレージ・システムに発行される(工程 9 1 2)。次いで工程 9 1 3 において、コンピュータ・システムは、ストレージ・システムからの応答、とりわけ、返される P E I D のそれぞれに関連付けられている P E ビットを調べて、P E 関連の L U N と、従来のデータ L U N との間における区別を行う。ストレージ・システムが N F S プロトコルベースである場合には、利用可能な P E の I D を得るために、A P I コールが、帯域外のコンピュータ・システムによって、プラグイン 6 2 2 から管理インターフェース(たとえば、管理インターフェース 6 2 4、6 2 5、または 6 2 6)に発行される(工程 9 1 4)。工程 9 1 3 および 9 1 4 の後に続く工程 9 1 6 において、コンピュータ・システムは、ストレージ・システムによって返された P E 関連の L U N の P E I D、または管理インターフェースによって返された P E I D を、バインド工程中に使用するために格納する。S C S I プロトコルベースのストレージ・デバイスによって返された P E I D はそれぞれ、W W N を含み、N F S プロトコルベースのストレージ・デバイスによって返された P E I D はそれぞれ、I P アドレスおよびマウント・ポイントを含むということを認識されたい。

## 【 0 0 4 9 】

図 9 B は、所与のコンピュータ・システム 1 0 3 が帯域内パスを介して接続される P E を、ストレージ・システム・マネージャ 1 3 1、またはストレージ・システム・マネージャ 1 3 2、または分散型ストレージ・システム・マネージャ 1 3 5 (以降では、「ストレージ・システム・マネージャ」と呼ばれる)が発見するための方法工程の流れ図である。そのような P E をストレージ・システム・マネージャによって発見することにより、ストレージ・システムは、コンピュータ・システムからのバインド要求に応答して、要求を行っているコンピュータ・システムに、有効な P E I D を返すことができ、コンピュータ・システムは、その P E I D 上に実際に接続されることが可能である。工程 9 5 0 において、ストレージ・システム・マネージャは、管理インターフェースおよびプラグイン 6 2 2 を介してコンピュータ・システム 1 0 3 に帯域外の「D i s c o v e r \_ T o p o l

10

20

30

40

50

ogy」APIコールを発行する。コンピュータ・システム103は、自分のシステムIDと、図9Aの流れ図を介して自分が発見したすべてのPE IDのリストとを返す。一実施形態においては、ストレージ・システム・マネージャは、管理インターフェースおよびプラグイン612を介して管理サーバ610に「Discover\_Topology」APIコールを発行することによって、工程950を実行する。そのような一実施形態においては、ストレージ・システムは、複数のコンピュータ・システムIDと、関連付けられているPE IDとを含む応答を、管理サーバ610が管理するそれぞれのコンピュータ・システム103ごとに1つ受信することになる。次いで工程952において、ストレージ・システム・マネージャは、工程950からの結果を処理する。たとえば、ストレージ・システム・マネージャは、自分の現在のコントロール下にはないすべてのPE IDのリストを消去する。たとえば、Discover\_Topologyコールを発行しているときにストレージ・システム・マネージャ135によって受信される特定のPE IDは、同じコンピュータ・システムに接続されている別のストレージ・システムに対応している可能性がある。同様に、受信される特定のPE IDは、その後ストレージ・システム管理者によって削除されたさらに古いPEに対応している可能性がある、といった具合である。工程954において、ストレージ・システム・マネージャは、処理された結果を、その後のバインド要求中に使用するためにキャッシュする。一実施形態においては、ストレージ・システム・マネージャは、図9Bの工程を定期的に行うして、コンピュータ・システムおよびネットワーク・トポロジーの進行中の変化で自分のキャッシュされた結果を更新する。別の実施形態においては、ストレージ・システム・マネージャは、新たなvvol作成要求を受信するたびに、図9Bの工程を実行する。さらに別の実施形態においては、ストレージ・システム・マネージャは、図7の認証工程を実行した後に、図9Bの工程を実行する。

#### 【0050】

図10は、バインド仮想ボリュームAPIを使用して仮想ボリューム・バインド要求を発行および実行するための方法工程の流れ図である。一実施形態においては、コンピュータ・システム103は、自分のアプリケーションのうちの1つが、まだPEにバインドされていないvvolに関連付けられているブロック・デバイスへのIOアクセスを要求した場合には、帯域外パス602を介してストレージ・システムにバインド要求を発行する。別の実施形態においては、管理サーバ610は、VMのパワー・オン、および1つのストレージ・コンテナから別のストレージ・コンテナへのvvolの移行を含む特定のVM管理オペレーションに関連したバインド要求を発行する。

#### 【0051】

まだPEにバインドされていないvvolに関連付けられているブロック・デバイスへのIOアクセスをアプリケーションが要求している上述の例について続けると、コンピュータ・システム103は、工程1002において、ブロック・デバイス・データベース533（または580）から、そのvvolのvvol IDを特定する。次いで工程1004において、コンピュータ・システム103は、そのvvolをバインドしたいという要求を、帯域外パス602を通じてストレージ・システムに発行する。

#### 【0052】

ストレージ・システム・マネージャは、工程1006において、管理インターフェース（たとえば、管理インターフェース624、625、または626）を介して、vvolをバインドしたいという要求を受信し、次いで、vvolがバインドされることになるPEを選択すること、選択されたPEに関するSLIDおよび世代番号を生成すること、ならびに（たとえば、IOマネージャ304を介して）接続データベース312を更新することを含む工程1008を実行する。vvolがバインドされることになるPEの選択は、接続（すなわち、コンピュータ・システム103への既存の帯域内接続を有するPEのみが、選択に利用可能である）と、利用可能なPEを通る現在のIOTrafficなどのその他の要因とに従って行われる。一実施形態においては、ストレージ・システムは、図9Bの方法に従ってコンピュータ・システム103がそのストレージ・システムに送信



したPEの処理がなされてキャッシュされたリストから選択を行う。SLLIDの生成は、図2Aのストレージ・システム・クラスタを採用している実施形態と、図2Bのストレージ・システム・クラスタを採用している実施形態との間において異なる。前者のケースにおいては、選択されたPEに関して一意であるSLLIDが生成される。後者のケースにおいては、vvolに対応するファイル・オブジェクトへのファイル・パスが、SLLIDとして生成される。選択されたPEに関してSLLIDおよび世代番号が生成された後に、接続データベース312は、vvolに対する新たに生成されたIOセッションを含むように更新される。次いで工程1010において、選択されたPEのID、生成されたSLLID、および世代番号が、コンピュータ・システム103に返される。任意選択で、図2Bのストレージ・システム・クラスタを採用している実施形態においては、一意のNFSファイル・ハンドルが、vvolに対応するファイル・オブジェクトに関して生成され、選択されたPEのID、生成されたSLLID、および世代番号とともにコンピュータ・システム103に返されることが可能である。工程1012において、コンピュータ・システム103は、ストレージ・システムから返されたPE ID、SLLID（および任意選択で、NFSハンドル）、および世代番号を含めるようにブロック・デバイス・データベース533（または580）を更新する。とりわけ、ストレージ・システムから返されたPE ID、SLLID（および任意選択で、NFSハンドル）、および各組の世代番号は、新たなエントリとしてブロック・デバイス・データベース533（または580）に加えられることになる。世代番号は、リプレイ攻撃を防ぐために使用されるということを認識されたい。したがって、リプレイ攻撃が懸念されない実施形態においては、世代番号は使用されない。

10

20

**【0053】**

同じvvolにIOを発行したいと望む別のアプリケーションによって開始された同じvvolへのその後のバインド要求に対して、ストレージ・システム・マネージャは、そのvvolを同じまたは別のPEにバインドすることができる。そのvvolが同じPEにバインドされる場合には、ストレージ・システム・マネージャは、その同じPEのIDと、以前に生成されたSLLIDとを返し、接続データベース312内に格納されているこのIO接続パスのリファレンス・カウントをインクリメントする。その一方で、そのvvolが別のPEにバインドされる場合には、ストレージ・システム・マネージャは、新たなSLLIDを生成し、その別のPEのIDと、新たに生成されたSLLIDとを返し、そのvvolへのこの新たなIO接続パスを新たなエントリとして接続データベース312に加える。

30

**【0054】**

アンバインド仮想ボリュームAPIを使用して、仮想ボリューム・アンバインド要求が発行されることが可能である。アンバインド要求は、vvolがそれまでにバインドされた際に経由したIO接続パスのPE IDおよびSLLIDを含む。しかしながら、アンバインド要求の処理は、助言的なものである。ストレージ・システム・マネージャが、すぐに、または遅れてvvolをPEからアンバインドするのは自由である。アンバインド要求は、PE IDおよびSLLIDを含むエントリのリファレンス・カウントをデクリメントするように接続データベース312を更新することによって、処理される。リファレンス・カウントがゼロまでデクリメントされた場合には、そのエントリは、削除されることが可能である。このケースにおいては、そのvvolは、引き続き存在するが、その所与のPE IDおよびSLLIDを使用したIOにそれ以上利用することはできないということに留意されたい。

40

**【0055】**

VMの仮想ディスクを実装しているvvolのケースにおいては、このvvolに関するリファレンス・カウントは、少なくとも1となる。VMがパワー・オフされ、それに関連してアンバインド要求が発行された場合には、リファレンス・カウントは、1だけデクリメントされることになる。リファレンス・カウントがゼロである場合には、vvolエントリは、接続データベース312から除去されることが可能である。一般に、エント

50

リーを接続データベース 312 から除去することは有益である。なぜなら、I/O マネージャ 304 は、より少ないデータを管理し、SLLID を再利用することもできるためである。そのような利点は、ストレージ・システムによって格納されている vvol の総数が多い（たとえば、数百万個程度の vvol）が、アプリケーションによってアクティブにアクセスされている vvol の総数が少ない（たとえば、数万個の VM）場合に、顕著になる。加えて、vvol がいずれの PE にもバインドされていない場合には、ストレージ・システムは、その vvol を DSU141 内のどこに格納するかを選択する際に、より大きな柔軟性を有する。たとえば、ストレージ・システムは、いくつかの DSU141 が、より高速なデータ・アクセスを提供し、その他の DSU141 が、（たとえば、ストレージ・コストを節約するために）より低速なデータ・アクセスを提供する、非対称的な、階層的な DSU141 を伴って実装されることが可能である。1 実施態様においては、vvol がいずれの PE にもバインドされていない場合には（これは、接続データベース 312 内のその vvol のエントリのリファレンス・カウントをチェックすることによって判定されることが可能である）、ストレージ・システムは、その vvol を、より低速なおよび/またはより安価なタイプの物理ストレージに移行させることができる。次いで、その vvol が PE にバインドされると、ストレージ・システムは、その vvol を、より高速なタイプの物理ストレージに移行させることができる。そのような移行は、vvol データベース 314 内の所与の vvol を構成するコンテナ・ロケーションの順序付けられたリストの 1 つまたは複数の要素を変更すること、およびコンテナ・データベース 316 のメタデータ・セクション内の対応するエクステンション・アロケーション・ビットマップを更新することによって達成されることが可能であるということ認識されたい。

#### 【0056】

vvol を PE にバインドすることおよびアンバインドすることによって、ストレージ・システム・マネージャは、vvol の活性を決定することができる。ストレージ・システム・マネージャは、この情報を利用して、I/O サービスを提供しない（パッシブな）vvol および I/O サービスを提供する（アクティブな）vvol に関してストレージ・システム・ベンダー固有の最適化を実行することができる。たとえば、ストレージ・システム・マネージャは、vvol が、特定のしきい値の時間を超えてパッシブな状態にとどまっている場合には、その vvol を、待ち時間が少ない（高いコストの）SSD から、待ち時間が中程度の（低いコストの）ハード・ドライブへ移動させるように構成されることが可能である。

#### 【0057】

図 11A および図 11B は、一実施形態による、仮想ボリュームに I/O を発行するための方法工程の流れ図である。図 11A は、アプリケーションからの I/O を直接ロー・ブロック・デバイスに発行するための方法工程 1100 の流れ図であり、図 11B は、アプリケーションからの I/O を、ファイル・システム・ドライバを通じて発行するための方法工程 1120 の流れ図である。

#### 【0058】

方法 1100 は、工程 1102 において開始し、工程 1102 では、図 5A ~ 図 5B において示されているアプリケーション 512、または図 5C ~ 図 5D において示されている VM 571 などのアプリケーションが、ロー・ブロック・デバイスに I/O を発行する。工程 1104 において、仮想ボリューム・デバイス・ドライバ 532 または 565 が、アプリケーションによって発行された I/O からロー・ブロックレベル I/O を生成する。工程 1106 において、ロー・ブロック・デバイスの名前が、仮想ボリューム・デバイス・ドライバ 532 または 565 によって PE ID および SLLID に（そしてまた、図 2B のストレージ・デバイスを採用している実施形態においては、NFS クライアント 545 または 585 によって NFS ハンドルに）変換される。工程 1108 において、データ・アクセス・レイヤ 540 または 566 が、PE ID および SLLID を（そしてまた、図 2B のストレージ・デバイスを採用している実施形態においては、NFS ハンドルを）ロー・ブロックレベル I/O へとエンコードすることを実行する。次いで工程 1110 にお

いて、HBA/NICが、ロー・ブロックレベルIOを発行する。

【0059】

図5A～図5Bにおいて示されているアプリケーション512など、VM以外のアプリケーションに関しては、方法1120は、工程1121において開始する。工程1121においては、アプリケーションが、vvolベースのブロック・デバイス上に格納されているファイルにIOを発行する。次いで工程1122において、ファイル・システム・ドライバ、たとえばファイル・システム・ドライバ510が、ファイルIOからブロックレベルIOを生成する。工程1122の後には、工程1126、1128、および1130（これらは、工程1106、1108、および1110と同じである）が実行される。

【0060】

図5C～図5Dにおいて示されているVM571など、VMアプリケーションに関しては、方法1120は、工程1123において開始する。工程1123においては、VMが、自分の仮想ディスクにIOを発行する。次いで工程1124において、このIOは、たとえばSCSI仮想化レイヤ563によって、ファイルIOに変換される。次いで、ファイル・システム・ドライバ、たとえばVMFSドライバ564が、工程1125において、ファイルIOからブロックレベルIOを生成する。工程1125の後には、工程1126、1128、および1130（これらは、工程1106、1108、および1110と同じである）が実行される。

【0061】

図12は、一実施形態による、ストレージ・システムにおいてIOを実行するための方法工程の流れ図である。工程1210において、コンピュータ・システムによって発行されたIOが、ストレージ・システムにおいて構成されているPEのうちの一つを通じて受信される。そのIOは、工程1212においてIOマネージャ304によって解析される。工程1212の後には、ストレージ・システム・クラスタが、図2Aにおいて示されているタイプのものである場合には、IOマネージャ304によって工程1214aが実行され、ストレージ・システム・クラスタが、図2Bにおいて示されているタイプのものである場合には、IOマネージャ304によって工程1214bが実行される。工程1214aにおいては、IOマネージャ304は、解析されたIOからSLLIDを抽出し、接続データベース312にアクセスして、PE IDと、抽出されたSLLIDとに対応するvvol IDを特定する。工程1214bにおいては、IOマネージャ304は、解析されたIOからNFSハンドルを抽出し、PE IDと、SLLIDとしてのNFSハンドルとを使用して、vvolを識別する。工程1214aおよび1214bの後には、工程1216が実行される。工程1216においては、IOが実行されることになる物理ストレージ・ロケーションを得るために、vvolデータベース314およびコンテナ・データベース316が、それぞれボリューム・マネージャ306およびコンテナ・マネージャ308によってアクセスされる。次いで工程1218において、データ・アクセス・レイヤ310が、工程1216において得られた物理ストレージ・ロケーション上でIOを実行する。

【0062】

いくつかの状況においては、アプリケーション（アプリケーション512またはVM571）、管理サーバ610、および/またはストレージ・システム・マネージャは、特定のPEに対するあるvvolのバインディングが、問題（そのPEが、あまりにも多くのバインディングでオーバーロード状態になっている場合など）を経験していると判定する可能性がある。そのような問題を解決するための方法として、バインドされているvvolは、IOコマンドがそのvvolへ導かれている間でさえ、ストレージ・システム・マネージャによって別のPEに再バインドされることが可能である。図13は、再バインドAPIを使用した、一実施形態による、vvol再バインド要求を発行および実行するための方法工程1300の流れ図である。

【0063】

示されているように、方法1300は、工程1302において開始し、工程1302で

10

20

30

40

50

は、ストレージ・システム・マネージャは、vvolが、そのvvolが現在バインドされている第1のPEとは異なる第2のPEにバインドされるべきであると判定する。工程1304において、ストレージ・システム・マネージャは、vvolを再バインドしたいという要求を、vvolにIOを発行するアプリケーションを実行しているコンピュータ・システム（たとえば、コンピュータ・システム103）に、帯域外バスを介して発行する。工程1306において、コンピュータ・システム103は、ストレージ・システム・マネージャから再バインド要求を受信し、それに応じて、vvolを新たなPEにバインドしたいという要求を発行する。工程1308において、ストレージ・システム・マネージャは、再バインド要求を受信し、それに応じて、vvolを新たなPEにバインドする。工程1310において、ストレージ・システム・マネージャは、図10に関連して上述したように、vvolが現在やはりバインドされている新たなPEのIDと、vvolにアクセスするためのSLIDとをコンピュータ・システムに送信する。

10

#### 【0064】

工程1312において、コンピュータ・システムは、ストレージ・システム・マネージャから新たなPE IDおよびSLIDを受信する。ブロック・デバイス・データベース533または580において、新たなPE接続のアクティブ・ビットが、はじめは1に設定され、これが意味するのは、新たなPEを介したvvolのための新たなIOセッションが確立されたということである。コンピュータ・システムはまた、第1のPE接続のアクティブ・ビットを0に設定し、これが意味するのは、このPE接続を通じてそのvvolにそれ以上IOが発行されることは不可能であるということである。認識されたいこととして、このPE接続は、非アクティブ化された際にすぐにアンバインドされるべきではない。なぜなら、処理中である、すなわち、発行されたが完了されていない可能性がある、このPE接続を通じたそのvvolへのIOが存在する可能性があるためである。したがって、工程1314において、コンピュータ・システムは、ブロック・デバイス・データベース533または580にアクセスして、第1のPE接続を通じてそのvvolに発行されたすべての「処理中コマンド」（CIF）が完了されているか、すなわち、CIF=0であるかを確認する。コンピュータ・システムは、工程1318を実行する前に、CIFがゼロになるのを待つ。その間に、そのvvolへのさらなるIOが、新たなPEを通じて発行される。なぜなら、新たなPE接続のアクティブ・ビットが既に1に設定されているためである。CIFがゼロに達しない場合には、工程1318が実行され、工程1318では、第1のPE接続をアンバインドしたいという要求が、ストレージ・システム・マネージャに発行される。次いで工程1320において、ストレージ・システム・マネージャは、そのvvolを第1のPEからアンバインドする。また、コンピュータ・システムは、工程1324において、新たなPEを通じてそのvvolにさらなるIOをすべて発行する。

20

30

#### 【0065】

図14は、一実施形態による、仮想ボリュームのライフ・サイクルの概念図である。図14において示されているすべてのコマンド、すなわち、作成、スナップショット、クローン、バインド、アンバインド、拡張、および削除は、vvol管理コマンド・セットを形成しており、図6に関連して上述したプラグイン612、622を通じてアクセス可能である。示されているように、「vvolを作成する」、「vvolをスナップショットする」、または「vvolをクローンする」というコマンドのうちのいずれかの結果としてvvolが生成された場合には、その生成されたvvolは、「パッシブな」状態にとどまり、パッシブな状態では、そのvvolは、特定のPEにバインドされておらず、したがってIOを受信することはできない。加えて、vvolがパッシブな状態にあるときに、「vvolをスナップショットする」、「vvolをクローンする」、または「vvolを拡張する」というコマンドのうちのいずれかが実行された場合には、オリジナルのvvolおよび（もしあれば）新たに作成されたvvolは、パッシブな状態にとどまる。やはり示されているように、パッシブな状態にあるvvolがPEにバインドされた場合には、そのvvolは、「アクティブな」状態に入る。逆に、アクティブなvvolが

40

50

PEからアンバインドされた場合には、そのvvolがいずれのさらなるPEにもバインドされていないと仮定すると、そのvvolは、パッシブな状態に入る。vvolがアクティブな状態にあるときに、「vvolをスナップショットする」、「vvolをクローンする」、「vvolを拡張する」、または「vvolを再バインドする」というコマンドのうちのいずれかが実行された場合には、オリジナルのvvolは、アクティブな状態にとどまり、(もしあれば)新たに作成されたvvolは、パッシブな状態にとどまる。

#### 【0066】

上述したように、1つのVMは、複数の仮想ディスクを有することができ、それぞれの仮想ディスクごとに別々のvvolが作成される。VMはまた、そのVMの構成について記述するメタデータ・ファイルを有する。それらのメタデータ・ファイルとしては、VM構成ファイル、VMログ・ファイル、ディスク記述子ファイル、すなわち、VMのための仮想ディスクのそれぞれに関するファイル、VMスワップ・ファイルなどが含まれる。仮想ディスクに関するディスク記述子ファイルは、仮想ディスクに関連する情報、たとえば、その仮想ディスクのvvol ID、その仮想ディスクのサイズ、その仮想ディスクがシン・プロビジョニングされているかどうか、および、その仮想ディスクに関して作成された1つまたは複数のスナップショットのIDなどを含む。VMスワップ・ファイルは、ストレージ・システム上におけるVMのスワップ・スペースを提供する。一実施形態においては、これらのVM構成ファイルは、vvol内に格納され、このvvolは、本明細書においてはメタデータvvolと呼ばれる。

#### 【0067】

図15は、一実施形態による、VMをプロビジョンするための方法工程の流れ図である。この実施形態においては、管理サーバ610と、VMをホストしているコンピュータ・システム、たとえば、図5Cにおいて示されているコンピュータ・システム102(以降では、「ホスト・コンピュータ」と呼ばれる)と、図2Aのストレージ・システム・クラスタ、とりわけストレージ・システム・マネージャ131、132、または135とが使用される。示されているように、ストレージ・システム・マネージャは、工程1502において、VMをプロビジョンしたいという要求を受信する。これは、VM管理者が、管理サーバ610への適切なユーザ・インターフェースを使用して、特定のサイズおよびストレージ能力プロファイルを有するVMをプロビジョンするためのコマンドを管理サーバ610に発行するとき生成される要求であることが可能である。それに応じて、工程1504において、管理サーバ610は、図8に関連して上述した様式で、VMのメタデータを含めるためのvvol(以降では、「メタデータvvol」と呼ばれる)を作成するための方法を開始し、それに従ってストレージ・システム・マネージャは、工程1508において、メタデータvvolを作成し、そのメタデータvvolのvvol IDを管理サーバ610に返す。工程1514において、管理サーバ610は、VMをホストしているコンピュータ・システムへ戻るメタデータvvolのvvol IDを登録する。工程1516において、ホスト・コンピュータは、図10に関連して上述した様式で、メタデータvvolをPEにバインドするための方法を開始し、それに従ってストレージ・システム・マネージャは、工程1518において、メタデータvvolをPEにバインドし、PE IDおよびSLLIDをホスト・コンピュータに返す。

#### 【0068】

工程1522において、ホスト・コンピュータは、ホスト・コンピュータのオペレーティング・システムへの「CREATE DEVICE」コールを使用して、メタデータvvolのブロック・デバイス・インスタンスを作成する。次いで工程1524において、ホスト・コンピュータは、ブロック・デバイスの上にファイル・システム(たとえば、VMFS)を作成し、それに応答して、ファイル・システムID(FSID)が返される。ホスト・コンピュータは、工程1526において、返されたFSIDを有するファイル・システムをマウントし、VMのメタデータを、そのファイル・システムに関連付けられている名前空間内に格納する。メタデータの例としては、VMログ・ファイル、ディスク記述子ファイル、すなわち、VMのための仮想ディスクのそれぞれに関するファイル

、およびVMスワップ・ファイルが含まれる。

【0069】

工程1528において、ホスト・コンピュータは、図8に関連して上述した様式で、VMの仮想ディスクのそれぞれに関するvvol（そのようなそれぞれのvvolは、本明細書においては「データvvol」と呼ばれる）を作成するための方法を開始し、それに従ってストレージ・システム・マネージャは、工程1530において、データvvolを作成し、そのデータvvolのvvol IDをホスト・コンピュータに返す。工程1532において、ホスト・コンピュータは、データvvolのIDを、仮想ディスクに関するディスク記述子ファイル内に格納する。この方法は、VMの仮想ディスクのうちのすべてに関してデータvvolが作成された後にメタデータvvolがアンバインドされること（図示せず）に伴って、終了する。

10

【0070】

図16Aは、図15に関連して説明した様式でVMがプロビジョンされた後にVMをパワー・オンするための方法工程の流れ図である。図16Bは、VMがパワー・オンされた後にVMをパワー・オフするための方法工程の流れ図である。これらの2つの方法は、VMのためのホスト・コンピュータによって実行される。

【0071】

工程1608においてVMパワー・オン・コマンドを受信すると、そのVMに対応するメタデータvvolのIDが、工程1610において取り出される。次いで工程1612において、メタデータvvolは、図10に関連して上述したようなバインド工程を経る。工程1614において、ファイル・システムがメタデータvvol上にマウントされ、それによって、工程1616において、データvvolに関するメタデータ・ファイル、とりわけディスク記述子ファイルを読み取ることができ、データvvol IDを得ることができる。次いで工程1618において、データvvolは、図10に関連して上述したように、1つずつバインド工程を経る。

20

【0072】

工程1620においてVMパワー・オフ・コマンドを受信すると、そのVMのデータvvolが、ブロック・デバイス・データベース（たとえば、図5Cのブロック・デバイス・データベース580）において非アクティブとしてマークされ、ホスト・コンピュータは、それらのデータvvolのそれぞれに関連付けられているCIFがゼロに達するのを待つ（工程1622）。それぞれのデータvvolに関連付けられているCIFがゼロに達した際に、ホスト・コンピュータは、工程1624において、そのデータvvolをアンバインドするようストレージ・システムに要求する。すべてのデータvvolに関連付けられているCIFがゼロに達した後に、工程1626において、メタデータvvolが、ブロック・デバイス・データベースにおいて非アクティブとしてマークされる。次いで工程1628において、メタデータvvolに関連付けられているCIFがゼロに達したときに、ホスト・コンピュータは、工程1630において、メタデータvvolがアンバインドされるよう要求する。

30

【0073】

図17および図18は、VMを再プロビジョンするための方法工程の流れ図である。ここで示されている例においては、図17は、VMのvvol、とりわけVMの仮想ディスクに関するデータvvolのサイズを拡張するための、ホスト・コンピュータ上で実行される方法工程の流れ図であり、図18は、ストレージ・コンテナ同士の間においてVMのvvolを移動させるための、ストレージ・システムにおいて実行される方法工程の流れ図である。

40

【0074】

VMの仮想ディスクに関するデータvvolのサイズを拡張するための方法が、工程1708において開始し、工程1708では、ホスト・コンピュータが、VMがパワー・オンされているかどうかを判定する。ホスト・コンピュータは、工程1708において、VMがパワー・オンされていないと判定した場合には、工程1710において、そのVMに

50

対応するメタデータ `vvol` の ID を取り出す。次いで工程 1712 において、メタデータ `vvol` に関するバインド工程が、ホスト・コンピュータによって開始される。バインドの後に、工程 1714 において、ホスト・コンピュータが、ファイル・システムをメタデータ `vvol` 上にマウントし、仮想ディスクに対応するデータ `vvol` の ID を、仮想ディスクに関するディスク記述子ファイル（これは、メタデータ `vvol` 上にマウントされたファイル・システム内のファイルである）から取り出す。次いで工程 1716 において、ホスト・コンピュータは、拡張 `vvol` API コールを工程 1716 においてストレージ・システムへ送信し、その拡張 `vvol` API コールは、データ `vvol` の ID と、データ `vvol` の新たなサイズとを含む。

**【0075】**

VM がパワー・オンされている場合には、ホスト・コンピュータは、工程 1715 において、拡張されることになる VM の仮想ディスクのデータ `vvol` の ID を取り出す。図 16A の方法から認識されたいこととして、この ID は、VM の仮想ディスクに関連付けられているディスク記述子ファイルから入手されることが可能である。次いで工程 1716 において、ホスト・コンピュータは、拡張 `vvol` API コールを工程 1716 においてストレージ・システムへ送信し、その拡張 `vvol` API コールは、データ `vvol` の ID と、データ `vvol` の新たなサイズとを含む。

**【0076】**

拡張 `vvol` API コールの結果、`vvol` データベースおよびコンテナ・データベース（たとえば、図 3 の `vvol` データベース 314 およびコンテナ・データベース 316）は、`vvol` の増大されたアドレス空間を反映するようにストレージ・システムにおいて更新される。拡張 `vvol` API コールが完了したという肯定応答を受信すると、ホスト・コンピュータは、工程 1718 において、VM の仮想ディスクに関するディスク記述子ファイルを新たなサイズで更新する。次いで工程 1720 において、ホスト・コンピュータは、VM がパワー・オンされているかどうかを判定する。VM がパワー・オンされていない場合には、ホスト・コンピュータは、工程 1722 において、ファイル・システムをアンマウントし、メタデータ `vvol` をアンバインドしたいという要求をストレージ・システムへ送信する。その一方で、VM がパワー・オンされている場合には、この方法は終了する。

**【0077】**

現在 PE にバインドされている VM の `vvol` をソース・ストレージ・コンテナから宛先ストレージ・コンテナへ移動させるための方法（この場合、ソース・ストレージ・コンテナおよび宛先ストレージ・コンテナの両方が、同じストレージ・システム・マネージャの範囲内にある）が、工程 1810 において開始し、工程 1810 では、ソース・ストレージ・コンテナおよび宛先ストレージ・コンテナ（それぞれ、SC1 および SC2）のコンテナ ID と、移動されることになる `vvol` の `vvol` ID とが受信される。次いで工程 1812 において、`vvol` データベース（たとえば、図 3 の `vvol` データベース 314）、およびコンテナ・データベース（たとえば、図 3 のコンテナ・データベース 316）のエクステント・アロケーション・ビットマップが、下記のように更新される。はじめに、ストレージ・システム・マネージャが、SC1 内の `vvol` エクステントをコンテナ・データベース 316 内の SC1 のエントリーから除去し、次いで、コンテナ・データベース 316 内の SC2 のエントリーを修正することによって、これらのエクステントを SC2 に割り振る。一実施形態においては、ストレージ・システムは、SC1 における（`vvol` ストレージ・エクステントの除去に起因する）ストレージ・キャパシティーの喪失を、新たなスピンドル・エクステントを SC1 に割り振ることによって補うこと、および SC2 における（`vvol` ストレージ・エクステントの追加に起因する）ストレージ・キャパシティーの増大を、いくつかの使用されていないスピンドル・エクステントを SC2 から除去することによって調整することが可能である。工程 1814 において、ストレージ・システム・マネージャは、現在バインドされている PE が `vvol` の新たなロケーションに IO を最適にサービス提供することができるかどうかを判定する。現在の PE

10

20

30

40

50

がvvolの新たなロケーションにIOをサービス提供することができない場合の1例は、ストレージ管理者が、ストレージ・システム・マネージャを、別々のPEを別々の顧客ひいては別々のストレージ・コンテナからのvvolに割り振るように静的に構成している場合である。現在のPEがvvolにIOをサービス提供することができない場合には、vvolは、工程1815において、図13に関連して上述した再バインド工程（および接続データベース、たとえば、図3の接続データベース312に対する関連する変更）を経る。工程1815の後には、工程1816が実行され、工程1816では、移動が首尾よく完了した旨の肯定応答が、ホスト・コンピュータに返される。工程1814において、現在のPEがvvolの新たなロケーションにIOをサービス提供することができるストレージ・システム・マネージャが判定した場合には、工程1815は迂回され、次いで工程1816が実行される。

10

## 【0078】

互換性がないストレージ・コンテナ同士の間において、たとえば、別々の製造業者のストレージ・デバイス内に作成されたストレージ・コンテナ同士の間においてvvolが移動される場合には、コンテナ・データベース316、vvolデータベース314、および接続データベース312に対する変更に加えて、ストレージ・コンテナ同士の間においてデータの移動が実行される。一実施形態においては、2008年5月29日に出願された「Offloading Storage Operations to Storage Hardware」と題されている米国特許出願第12/129,323号（その全内容を本願明細書に援用する）に記載されているデータ移動技術が採用される。

20

## 【0079】

図19は、テンプレートVMからVMをクローンするための、ホスト・コンピュータおよびストレージ・システムにおいて実行される方法工程の流れ図である。この方法は、工程1908において開始し、工程1908では、ホスト・コンピュータが、新たなVMに関するメタデータvvolを作成したいという要求をストレージ・システムへ送信する。1910において、ストレージ・システムは、図8に関連して上述した方法に従って新たなVMに関するメタデータvvolを作成し、新たなメタデータvvol IDをホスト・コンピュータに返す。次いで工程1914において、クローンvvol APIコールが、テンプレートVMに属するすべてのデータvvol IDに関して、ホスト・コンピュータから帯域外パス601を介してストレージ・システムに発行される。工程1918において、ストレージ・システム・マネージャが、テンプレートVMのデータvvolと、新たなVMのデータvvolとに互換性があるか否かをチェックする。別々の製造業者のストレージ・システム内に作成されたストレージ・コンテナ同士の間においてクローニングが行われる場合には、データvvol同士に互換性がない可能性があるということを認識されたい。互換性がある場合には、工程1919が実行される。工程1919において、ストレージ・システム・マネージャは、新たなデータvvol IDを生成すること、コンテナ・データベース316内のアロケーション・ビットマップを更新すること、および新たなvvolエントリをvvolデータベース314に加えることによって、新たなデータvvolを作成し、テンプレートVMのデータvvol内に格納されているコンテンツを新たなVMのデータvvolにコピーする。工程1920において、ストレージ・システム・マネージャは、新たなデータvvol IDをホスト・コンピュータに返す。新たなデータvvol IDを受信することは、データvvolのクローニングがエラーなく完了した旨の確認をホスト・コンピュータに提供する。次いで工程1925において、ホスト・コンピュータは、メタデータ・ファイル、とりわけディスク記述子ファイルを、新たに生成されたデータvvol IDで更新するために、新たなVMのメタデータvvolにIOを発行する。ホスト・コンピュータによってストレージ・システムに発行されたIOは、工程1926においてストレージ・システムによって実行され、その結果として、新たなVMのディスク記述子ファイルが、新たに生成されたデータvvol IDで更新される。

30

40

## 【0080】

50



工程 1918 において、テンプレート VM のデータ `vvol` と、新たな VM のデータ `vvol` とに互換性がないとストレージ・システム・マネージャが判定した場合には、エラー・メッセージが、ホスト・コンピュータに返される。このエラー・メッセージを受信すると、ホスト・コンピュータは、工程 1921 において、新たなデータ `vvol` を作成するために、作成 `vvol` API コールをストレージ・システムに発行する。工程 1922 において、ストレージ・システム・マネージャは、新たなデータ `vvol` ID を生成すること、コンテナ・データベース 316 内のアロケーション・ビットマップを更新すること、および新たな `vvol` エントリーを `vvol` データベース 314 に加えることによって、新たなデータ `vvol` を作成し、新たなデータ `vvol` ID をホスト・コンピュータに返す。工程 1923 において、ホスト・コンピュータは、2009 年 1 月 21 日に  
10  
出願された「Data Mover for Computer System」と題されている米国特許出願第 12 / 356 , 694 号（その全内容を本願明細書に援用する）に記載されている技術に従って、データの移動を実行する（工程 1923）。工程 1923 の後には、工程 1925 および 1926 が、上述のように実行される。

#### 【0081】

図 20 は、別の実施形態による、VM をプロビジョンするための方法工程の流れ図である。この実施形態においては、管理サーバ 610 と、VM をホストしているコンピュータ・システム、たとえば、図 5D において示されているコンピュータ・システム 102（以降では、「ホスト・コンピュータ」と呼ばれる）と、図 2B のストレージ・システム・クラスタ、とりわけストレージ・システム・マネージャ 131、またはストレージ・システム・マネージャ 132、またはストレージ・システム・マネージャ 135 とが使用される  
20  
。示されているように、VM をプロビジョンしたいという要求が、工程 2002 において受信される。これは、VM 管理者が、管理サーバ 610 への適切なユーザ・インターフェースを使用して、特定のサイズおよびストレージ能力プロファイルを有する VM をプロビジョンするためのコマンドを管理サーバ 610 に発行するときに生成される要求であることが可能である。それに応じて、工程 2004 において、管理サーバ 610 は、図 8 に関連して上述した様式で、VM のメタデータを含めるための `vvol`、とりわけメタデータ `vvol` を作成するための方法を開始し、それに従ってストレージ・システム・マネージャは、工程 2008 において、メタデータ `vvol`（これは、NAS デバイス内のファイルである）を作成し、メタデータ `vvol` ID を管理サーバ 610 に返す。工程 202  
30  
0 において、管理サーバ 610 は、ホスト・コンピュータへ戻るメタデータ `vvol` の `vvol` ID を登録する。工程 2022 において、ホスト・コンピュータは、メタデータ `vvol` ID に関するバインド要求をストレージ・システムに発行し、それに応答して、ストレージ・システムは、工程 2023 において、IP アドレスおよびディレクトリ・パスをそれぞれ PE ID および SLL ID として返す。工程 2024 において、ホスト・コンピュータは、指定された IP アドレスおよびディレクトリ・パスにおいてディレクトリをマウントし、そのマウントされたディレクトリ内にメタデータ・ファイルを格納する。NFS を使用する実施形態においては、NFS クライアント 545 または 585 が、そのようなディレクトリに NFS 要求を発行するために、所与の IP アドレスおよびディレクトリ・パスを NFS ハンドルへと変換することができる。  
40

#### 【0082】

工程 2026 において、ホスト・コンピュータは、図 8 に関連して上述した様式で、VM の仮想ディスクのそれぞれに関するデータ `vvol` を作成するための方法を開始し、それに従ってストレージ・システム・マネージャは、工程 2030 において、データ `vvol` を作成し、そのデータ `vvol` の `vvol` ID をホスト・コンピュータに返す。工程 2032 において、ホスト・コンピュータは、データ `vvol` の ID を、仮想ディスクに関するディスク記述子ファイル内に格納する。この方法は、VM の仮想ディスクのうちのすべてに関してデータ `vvol` が作成された後にメタデータ `vvol` がアンバインドされること（図示せず）に伴って、終了する。

#### 【0083】

10

20

30

40

50

図 8 に関連して上述したように、新たな `vvol` がストレージ・コンテナから作成され、その新たな `vvol` に関してストレージ能力プロファイルが明示的に指定されない場合には、その新たな `vvol` は、ストレージ・コンテナに関連付けられているストレージ能力プロファイルを引き継ぐことになる。ストレージ・コンテナに関連付けられているストレージ能力プロファイルは、いくつかの別々のプロファイルのうちの一つから選択されることが可能である。たとえば、図 2 1 において示されているように、それらの別々のプロファイルとしては、プロダクション (`prod`) プロファイル 2 1 0 1、開発 ( `dev` ) プロファイル 2 1 0 2、およびテスト・プロファイル 2 1 0 3 (ここでは「プロファイル 2 1 0 0」と総称される) が含まれる。その他の多くのプロファイルが定義されることも可能であるということ認識されたい。示されているように、特定のプロファイルのそれぞれのプロファイル・エントリは、固定タイプまたは可変タイプのものであり、1 つの名前と、それに関連付けられている 1 つまたは複数の値とを有している。固定タイプのプロファイル・エントリは、固定された数の選択可能なアイテムを有している。たとえば、プロファイル・エントリ「複製」は、真または偽であるように設定されることが可能である。対照的に、可変タイプのプロファイル・エントリは、事前に定義された選択肢を有していない。その代わりに、可変タイプのプロファイル・エントリに関しては、デフォルトの値およびある範囲の値が設定され、ユーザは、その範囲内にある任意の値を選択することができる。値がまったく指定されない場合には、デフォルトの値が使用される。図 2 1 に示されている例示的なプロファイル 2 1 0 0 においては、可変タイプのプロファイル・エントリは、コンマによって区切られている 3 つの数字を有している。第 1 の数字は、指定された範囲の下限であり、第 2 の数字は、指定された範囲の上限である。第 3 の数字は、デフォルトの値である。したがって、プロダクション・プロファイル 2 1 0 1 において定義されているストレージ能力プロファイルを引き継ぐ `vvol` は、複製されることになり (複製。値 = 真)、複製に関する目標復旧時間 (`RTO: recovery time objective`) は、0 . 1 時間から 2 4 時間の範囲内で定義されることが可能であり、デフォルトは 1 時間である。加えて、この `vvol` に関しては、スナップショットが可能である (スナップショット。値 = 真)。保持されるスナップショットの数は、1 から 1 0 0 の範囲内であり、デフォルトは 1 であり、スナップショットの頻度は、1 時間に 1 回から 2 4 時間に 1 回の範囲内であり、デフォルトは 1 時間に 1 回である。「スナップ引き継ぎ」の列は、派生 `vvol` である新たな `vvol` を作成するために所与の `vvol` がスナップショットされる場合に所与のプロファイル属性 (およびその値) が派生 `vvol` に伝搬されるべきかどうかを示す。プロダクション・プロファイル 2 1 0 1 の例においては、最初の 2 つのプロファイル・エントリ (複製および `RTO`) のみが、プロダクション・プロファイル 2 1 0 1 を有する所与の `vvol` のスナップショット `vvol` に伝搬されることが可能である。スナップショット `vvol` のその他のすべての属性の値は、そのプロファイルにおいて指定されているデフォルトの値に設定されることになる。言い換えれば、所与の `vvol` に関するこれらのその他の属性のあらゆるカスタマイゼーション (たとえば、スナップショット頻度のデフォルトではない値) は、それらの対応する「スナップ引き継ぎ」の列が偽であるため、スナップショット `vvol` に伝搬されない。このプロファイルは、どの属性値が所与の `vvol` のそれぞれクローンおよびレプリカに伝搬されるかをコントロールする「クローン引き継ぎ」 (図示せず) および「レプリカ引き継ぎ」 (図示せず) などのその他の列も含む。

【 0 0 8 4 】

図 4 の方法に従ってストレージ・コンテナが作成される場合には、そのストレージ・コンテナから作成される `vvol` に関して定義することができるストレージ能力プロファイルのタイプが設定されることが可能である。図 2 1 における流れ図は、図 4 において示されているストレージ・コンテナを作成するための方法を、工程 4 1 2 と工程 4 1 3 との間に工程 2 1 1 0 が挿入された状態で示している。工程 2 1 1 0 において、ストレージ管理者は、作成されているストレージ・コンテナに関するプロファイル 2 1 0 0 のうちの一つまたは複数を選択する。たとえば、1 人の顧客のために作成された 1 つのストレージ・コ

10

20

30

40

50

ンテナが、プロダクション・プロファイル 2 1 0 1 およびデベロップメント・プロファイル 2 1 0 2 に関連付けられることが可能であり、それによって、プロダクション・タイプのものである `vv01` は、場合によってデフォルトの値または顧客によって指定された値を伴うプロダクション・プロファイル 2 1 0 1 において定義されているストレージ能力プロファイルを引き継ぐことになり、デベロップメント・タイプのものである `vv01` は、場合によってデフォルトの値または顧客によって指定された値を伴うデベロップメント・プロファイル 2 1 0 2 において定義されているストレージ能力プロファイルを引き継ぐことになる。

#### 【 0 0 8 5 】

図 2 2 は、`vv01` を作成して、その `vv01` に関するストレージ能力プロファイルを定義するための、ストレージ・システム・マネージャ 1 3 1、1 3 2、または 1 3 5 によって実行される方法工程を示す流れ図である。図 2 2 の方法工程、とりわけ工程 2 2 1 0、2 2 1 2、2 2 1 8、および 2 2 2 0 はそれぞれ、図 8 において示されている工程 8 0 6、8 1 0、8 1 2、および 8 1 4 に対応する。加えて、図 2 2 の方法工程は、工程 2 2 1 4、2 2 1 5、および 2 2 1 6 を含み、これらの工程は、作成されている `vv01` に関するストレージ能力プロファイルを定義する。

10

#### 【 0 0 8 6 】

工程 2 2 1 4 において、ストレージ・システム・マネージャは、ストレージ能力プロファイルにおいて使用されることになる値が、`vv01` を作成したいという要求内で指定されているかどうかを判定する。ストレージ能力プロファイルにおいて使用されることになる値が指定されていない場合には、ストレージ・システム・マネージャは、工程 2 2 1 5 において、`vv01` のストレージ・コンテナに関連付けられているストレージ能力プロファイルを、デフォルトの値を伴う `vv01` のストレージ能力プロファイルとして採用する。ストレージ能力プロファイルにおいて使用されることになる値が指定されている場合には、ストレージ・システム・マネージャは、工程 2 2 1 6 において、`vv01` のストレージ・コンテナに関連付けられているストレージ能力プロファイルを、デフォルトの値の代わりに指定されている値を伴う `vv01` のストレージ能力プロファイルとして採用する。

20

#### 【 0 0 8 7 】

一実施形態においては、`vv01` のストレージ能力プロファイルは、キー/値のペアとして `vv01` データベース 3 1 4 内に格納される。いったん `vv01` のストレージ能力プロファイルが定義されて、キー/値のペアとして `vv01` データベース 3 1 4 内に格納されると、図 2 1 の例示的なプロファイルにおいて示されているように複製およびスナップショット関連の属性および値がこのプロファイルの一部である限り、ストレージ・システムは、ホスト・コンピュータによって発行されるさらなる命令を伴わずに、その `vv01` に関する複製およびスナップショットを実行することができる。

30

#### 【 0 0 8 8 】

図 2 3 は、親 `vv01` からスナップショットを作成するための、ストレージ・システム・マネージャ 1 3 1、1 3 2、または 1 3 5 によって実行される方法工程を示す流れ図である。一実施形態においては、所与の `vv01` のストレージ能力プロファイル内のスナップショット定義に従ってスナップショットをスケジュールするためのスナップショット・トラッキング・データ構造が採用される。スナップショットのためのスケジュールされた時刻に達すると、ストレージ・システム・マネージャは、工程 2 3 1 0 において、スナップショット・トラッキング・データ構造から `vv01 ID` を取り出す。次いで工程 2 3 1 2 において、ストレージ・システム・マネージャは、スナップショットに関する一意の `vv01 ID` を生成する。ストレージ・システム・マネージャは、工程 2 3 1 5 において、親 `vv01` (すなわち、スナップショット・トラッキング・データ構造から取り出された `vv01 ID` を有する `vv01`) のストレージ能力プロファイルを、スナップショット `vv01` のストレージ能力プロファイルとして採用する。これは、ストレージ・システムによって駆動される自動化されたプロファイル駆動型のスナップショット工程であるため、ユーザには、スナップショット `vv01` のストレージ能力プロファイルにおいて使

40

50

用されることになるカスタムの値を指定するための機会はないということに留意されたい。工程 2318 において、ストレージ・システム・マネージャは、コンテナ・データベース 316 内のアロケーション・ビットマップを更新すること、およびスナップショット vvol に関する新たな vvol エントリを vvol データベース 314 に加えることによって、親 vvol のストレージ・コンテナ内にスナップショット vvol を作成する。次いで工程 2320 において、ストレージ・システム・マネージャは、親 vvol に関する次なるスナップショットを生成するための時刻をスケジュールすることによって、スナップショット・トラッキング・データ構造を更新する。ストレージ・システム・マネージャは、スナップショット・トラッキング・データ構造を保持することと、スケジュールされたスナップショットを命じるストレージ能力プロファイルを有するすべての vvol に関して図 23 の方法工程を実行することとを同時に行わなければならないということを確認されたい。

10

20

30

40

50

**【0089】**

上述の様式でスナップショットが作成された後に、vvol データベース 314 内に格納されているキー/値のペアは、スナップショット vvol がタイプ = スナップショットのものであるということを示すように更新される。また、スナップショットに関して世代番号が保持され、その世代番号が、スナップショットがとられるたびにインクリメントされるか、または日にち + 時刻に等しくなるように設定される実施形態においては、世代番号は、キー/値のペアとして格納される。スナップショット vvol の親 vvol ID も、キー/値のペアとしてスナップショット vvol エントリ内に格納される。結果として、ホスト・コンピュータは、特定の vvol ID に対応するスナップショットを求めて vvol データベース 314 にクエリーを行うことができる。ホスト・コンピュータは、特定の vvol ID および特定の世代番号に対応するスナップショットを求めて vvol データベースにクエリーを発行することも可能である。

**【0090】**

本明細書に記載されているさまざまな実施形態は、コンピュータ・システム内に格納されているデータを含むさまざまなコンピュータ実施オペレーションを採用することができる。たとえば、これらのオペレーションは、通常、必須ではないが、物理量の物理的な操作を必要とする場合があり、これらの量は、電気信号または磁気信号の形態を取ることができ、そうした形態では、これらの量、またはこれらの量の表示は、格納されること、転送されること、結合されること、比較されること、またはその他の形で操作されることが可能である。さらに、そのような操作はしばしば、作り出すこと、識別すること、判定すること、または比較することなどの用語で呼ばれる。1つまたは複数の実施形態の一部を形成する、本明細書に記載されているあらゆるオペレーションは、有用なマシン・オペレーションであることが可能である。加えて、1つまたは複数の実施形態はまた、これらのオペレーションを実行するためのデバイスまたは装置に関する。その装置は、特定の必要とされる目的のために特別に構築されることが可能であり、または、コンピュータ内に格納されているコンピュータ・プログラムによって選択的にアクティブ化または構成される汎用コンピュータであることが可能である。とりわけ、さまざまな汎用マシンが、本明細書における教示に従って書かれたコンピュータ・プログラムとともに使用されることが可能であり、または、必要とされるオペレーションを実行するためのさらに専門化された装置を構築することが、より好都合である可能性がある。

**【0091】**

本明細書に記載されているさまざまな実施形態は、ハンドヘルド・デバイス、マイクロプロセッサ・システム、マイクロプロセッサベースのまたはプログラム可能な家庭用電化製品、ミニコンピュータ、メインフレーム・コンピュータなどを含むその他のコンピュータ・システム構成とともに実施されることが可能である。

**【0092】**

1つまたは複数の実施形態は、1つもしくは複数のコンピュータ・プログラムとして、または1つもしくは複数のコンピュータ可読メディアにおいて具体化される1つもしくはは

複数のコンピュータ・プログラム・モジュールとして実装されることが可能である。コンピュータ可読メディアという用語は、その後コンピュータ・システムに入力されることが可能であるデータを格納することができる任意のデータ・ストレージ・デバイスを指し、コンピュータ可読メディアは、コンピュータ・プログラムがコンピュータによって読み取られることを可能にする様式でそれらのコンピュータ・プログラムを具体化するための任意の既存のまたはその後開発されるテクノロジーに基づくことができる。コンピュータ可読メディアの例としては、ハード・ドライブ、ネットワーク・アタッチト・ストレージ(NAS)、読み取り専用メモリ、ランダムアクセス・メモリ(たとえば、フラッシュ・メモリ・デバイス)、CD(Compact Disc)、CD-ROM、CD-R、またはCD-RW、DVD(Digital Versatile Disc)、磁気テープ、ならびにその他の光学式および非光学式のデータ・ストレージ・デバイスが含まれる。コンピュータ可読メディアは、ネットワークに結合されたコンピュータ・システムを介して分散されることも可能であり、それによってコンピュータ可読コードは、分散された様式で格納および実行される。

10

#### 【0093】

1つまたは複数の実施形態について、理解を明確にするためにいくらか詳細に説明したが、特許請求の範囲の範囲内で特定の変更および修正が行われることが可能であるということは明らかであろう。たとえば、SCSIが、SANデバイスのためのプロトコルとして採用され、NFSが、NASデバイスのためのプロトコルとして使用される。ファイバ・チャンネルなど、SCSIプロトコルに対する任意の代替手段が使用されることが可能であり、CIFS(Common Internet File System)プロトコルなど、NFSプロトコルに対する任意の代替手段が使用されることが可能である。したがって、記載されている実施形態は、限定的なものではなく例示的なものとみなされるべきであり、特許請求の範囲の範囲は、本明細書において与えられている詳細に限定されるものではなく、特許請求の範囲の範囲および均等物の中で修正されることが可能である。特許請求の範囲において、要素および/または工程は、特許請求の範囲において明示的に記載されていない限り、オペレーションのいかなる特定の順序も意味するものではない。

20

#### 【0094】

加えて、記載されている仮想化方法は一般に、仮想マシンが、特定のハードウェア・システムと整合するインターフェースを提示すると想定しているが、記載されているそれらの方法は、いかなる特定のハードウェア・システムにも直接対応しない仮想化に関連して使用されることが可能である。ホストされる実施形態、ホストされない実施形態として、またはそれらの両者の間における区別をあいまいにする傾向がある実施形態として実施される、さまざまな実施形態による仮想化システムが、すべて想定されている。さらに、さまざまな仮想化オペレーションは、全体的にまたは部分的にハードウェアで実施されることが可能である。たとえば、ハードウェアの実施態様は、非ディスク・データをセキュアにするためのストレージ・アクセス要求の修正のためのルックアップ・テーブルを採用することができる。

30

#### 【0095】

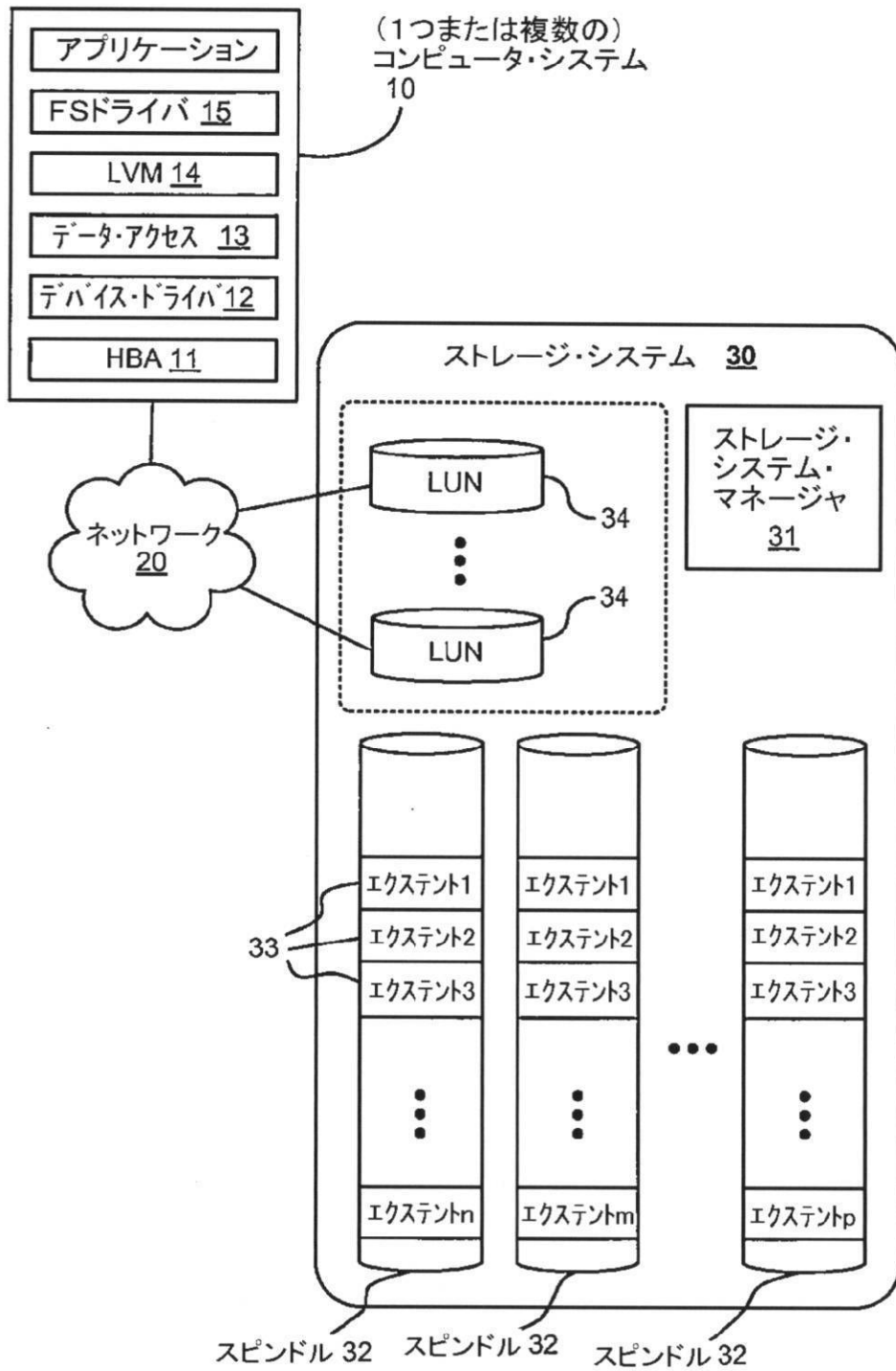
仮想化の度合いにかかわらず、多くの変形、修正、追加、および改良が可能である。したがって仮想化ソフトウェアは、仮想化機能を実行するホスト、コンソール、またはゲスト・オペレーティング・システムのコンポーネントを含むことができる。単一のインスタンスとして本明細書に記載されているコンポーネント、オペレーション、または構造のために、複数のインスタンスが提供されることが可能である。最後に、さまざまなコンポーネント、オペレーション、およびデータストア同士の間における境界は、いくらか任意のものであり、特定のオペレーションは、特定の例示的な構成のコンテキストにおいて示されている。機能のその他の割り当ても想定されており、本明細書に記載されている実施形態の範囲内に収まるることができる。一般に、例示的な構成において別々のコンポーネントとして提示されている構造および機能は、結合された構造またはコンポーネントとして実装されることが可能である。同様に、単一のコンポーネントとして提示されている構造お

40

50

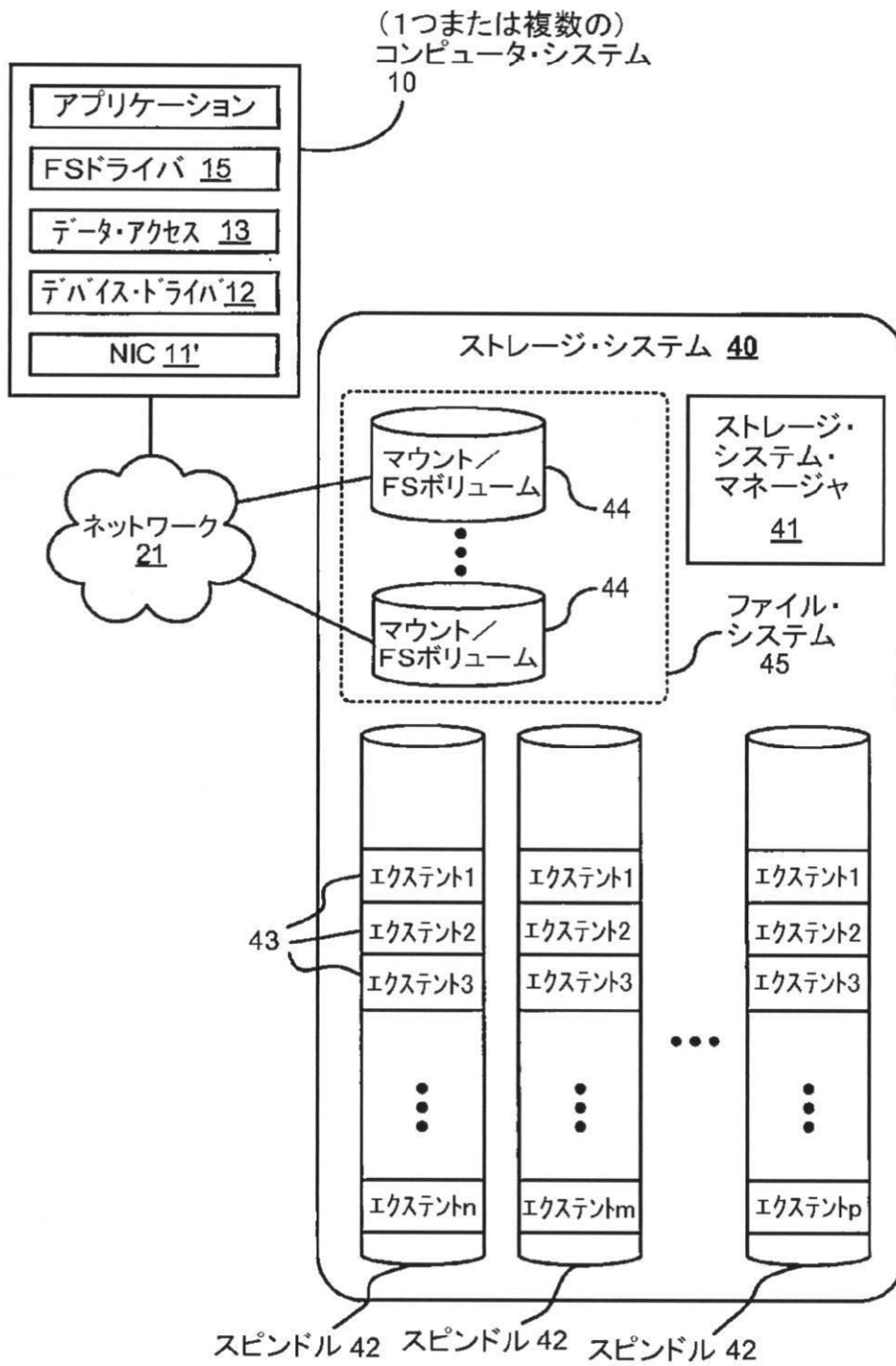
よび機能は、別々のコンポーネントとして実装されることが可能である。これらおよびその他の変形、修正、追加、および改良は、添付の（１つまたは複数の）特許請求の範囲の範疇内に収まることができる。

【図 1 A】



(従来技術)

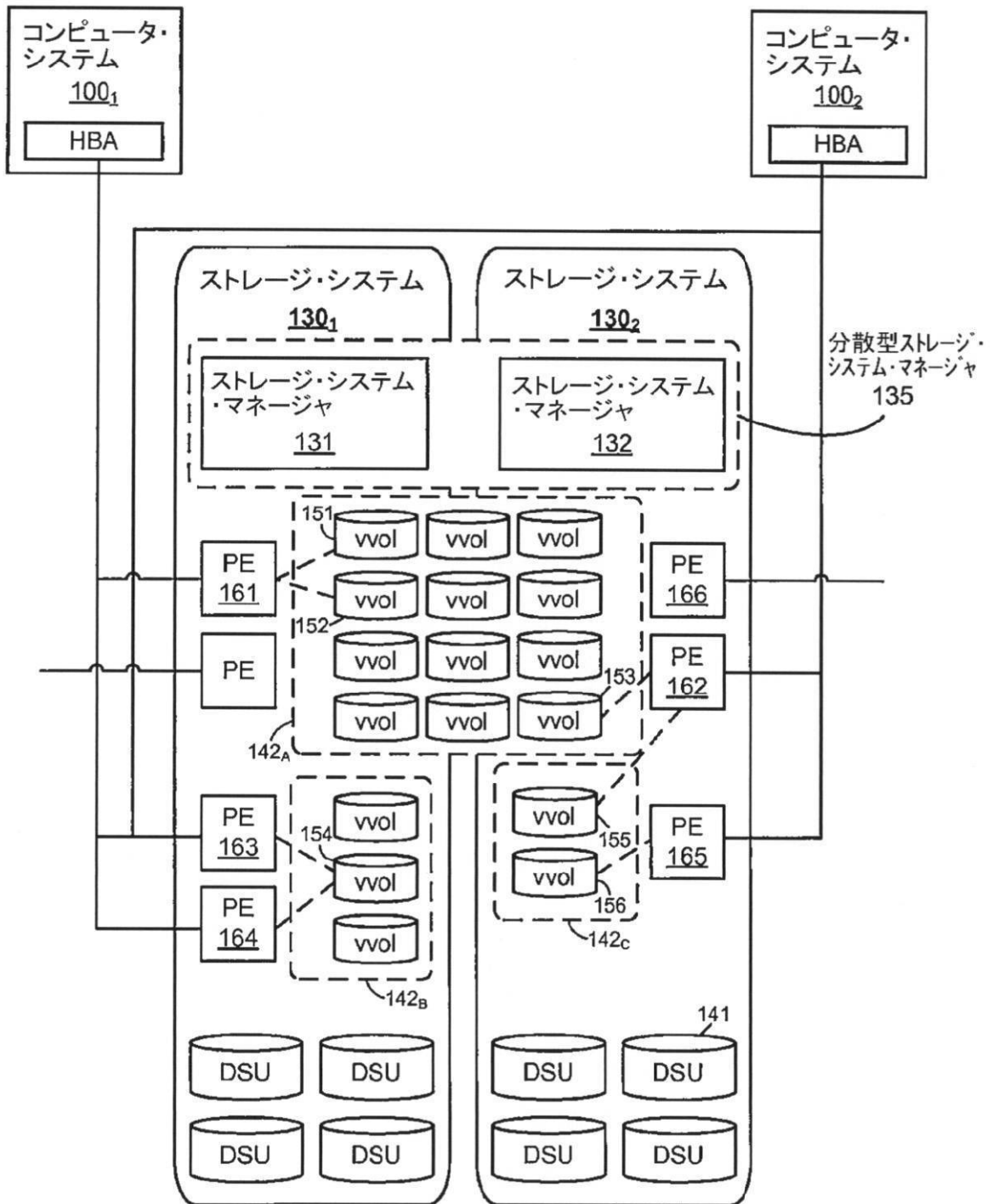
【図1B】



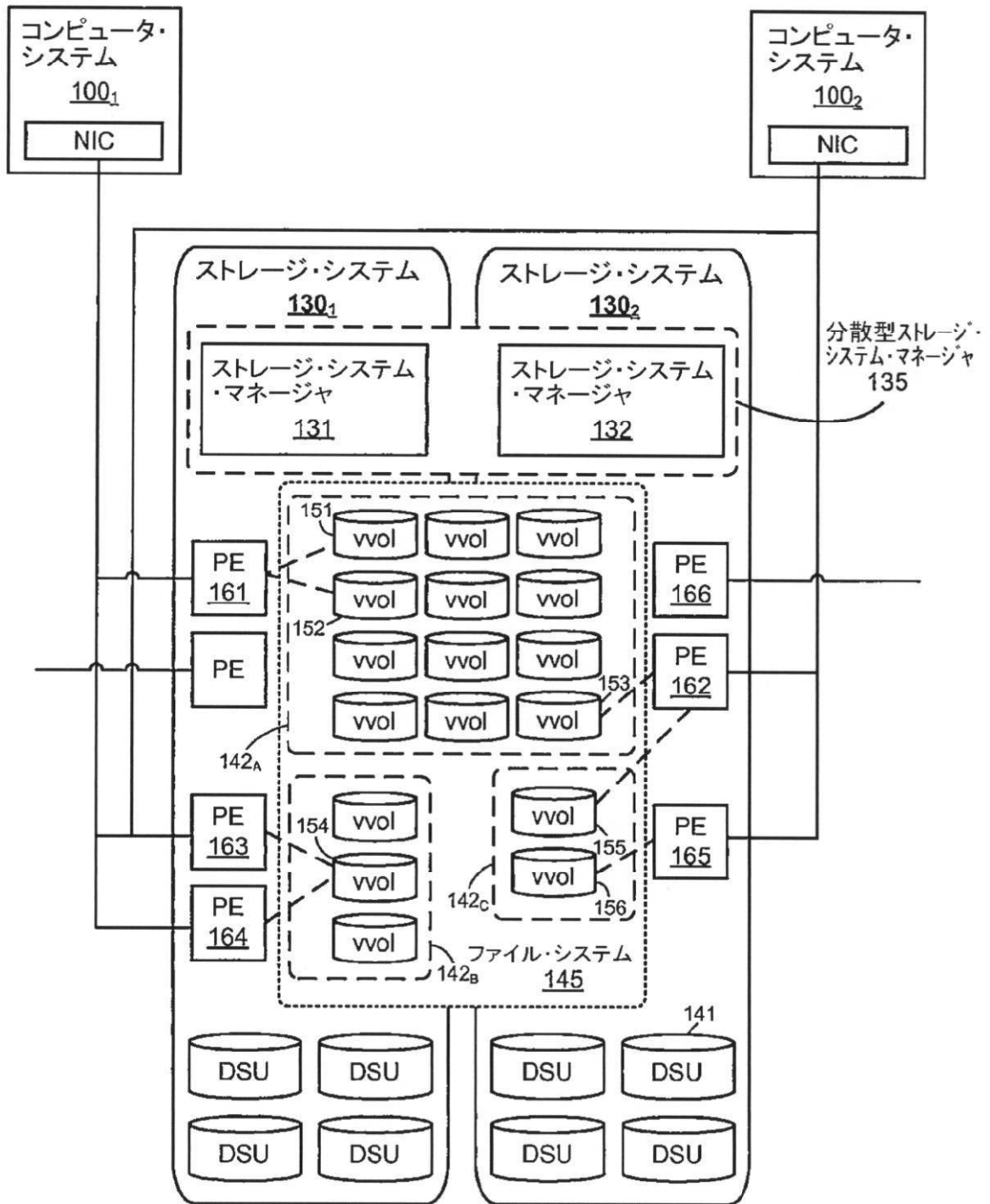
(従来技術)



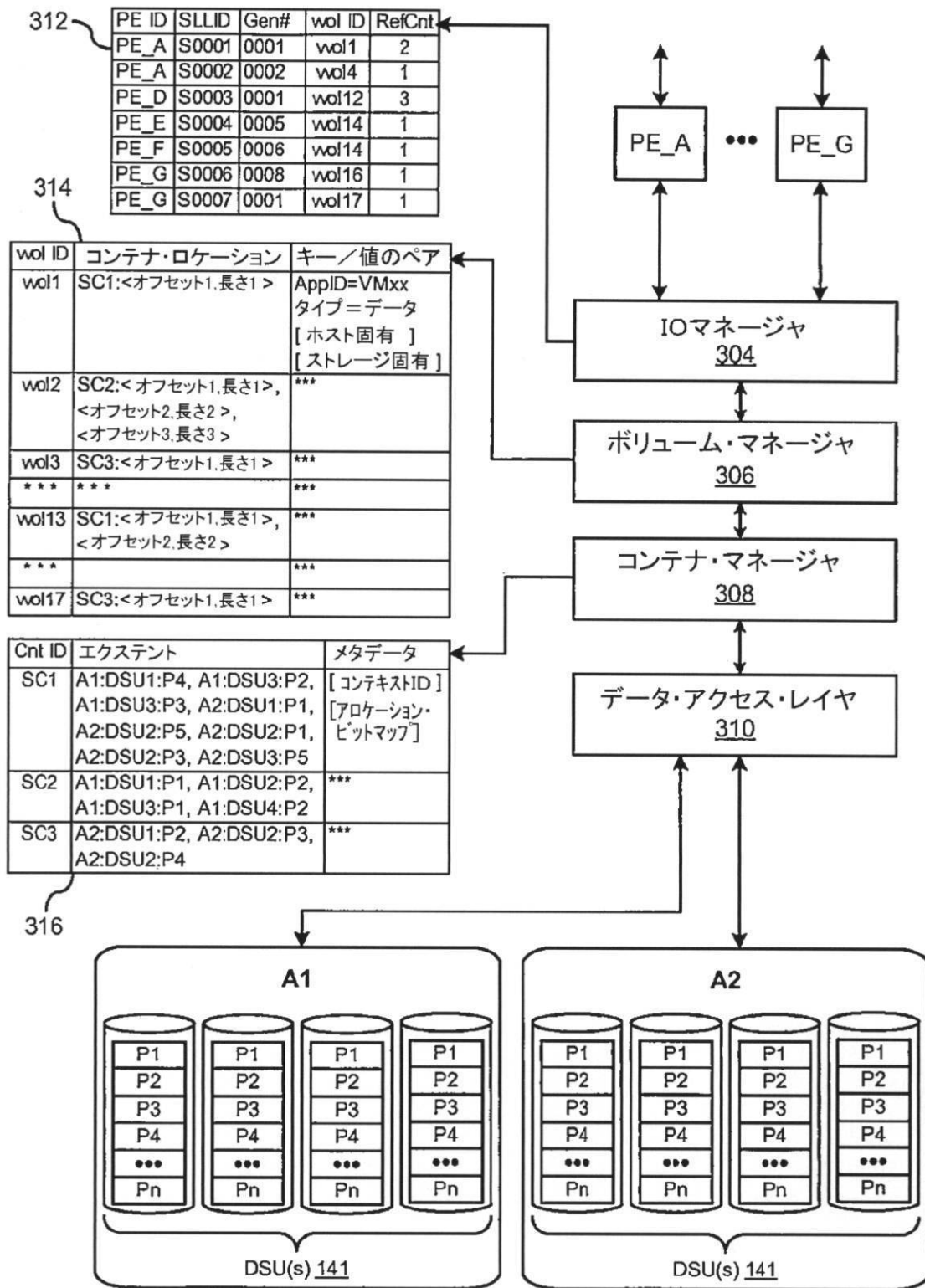
【図 2 A】



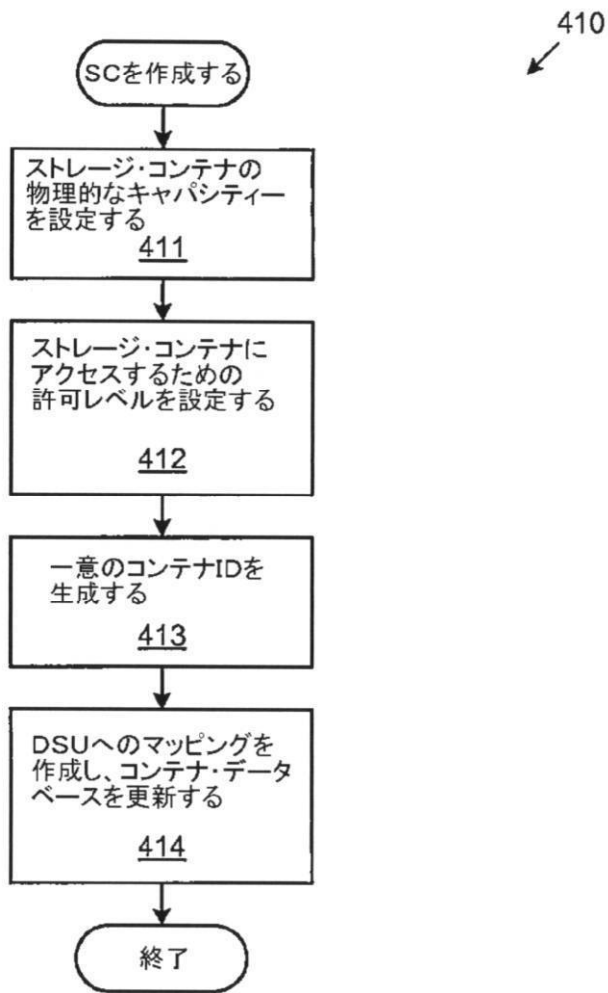
【図 2 B】



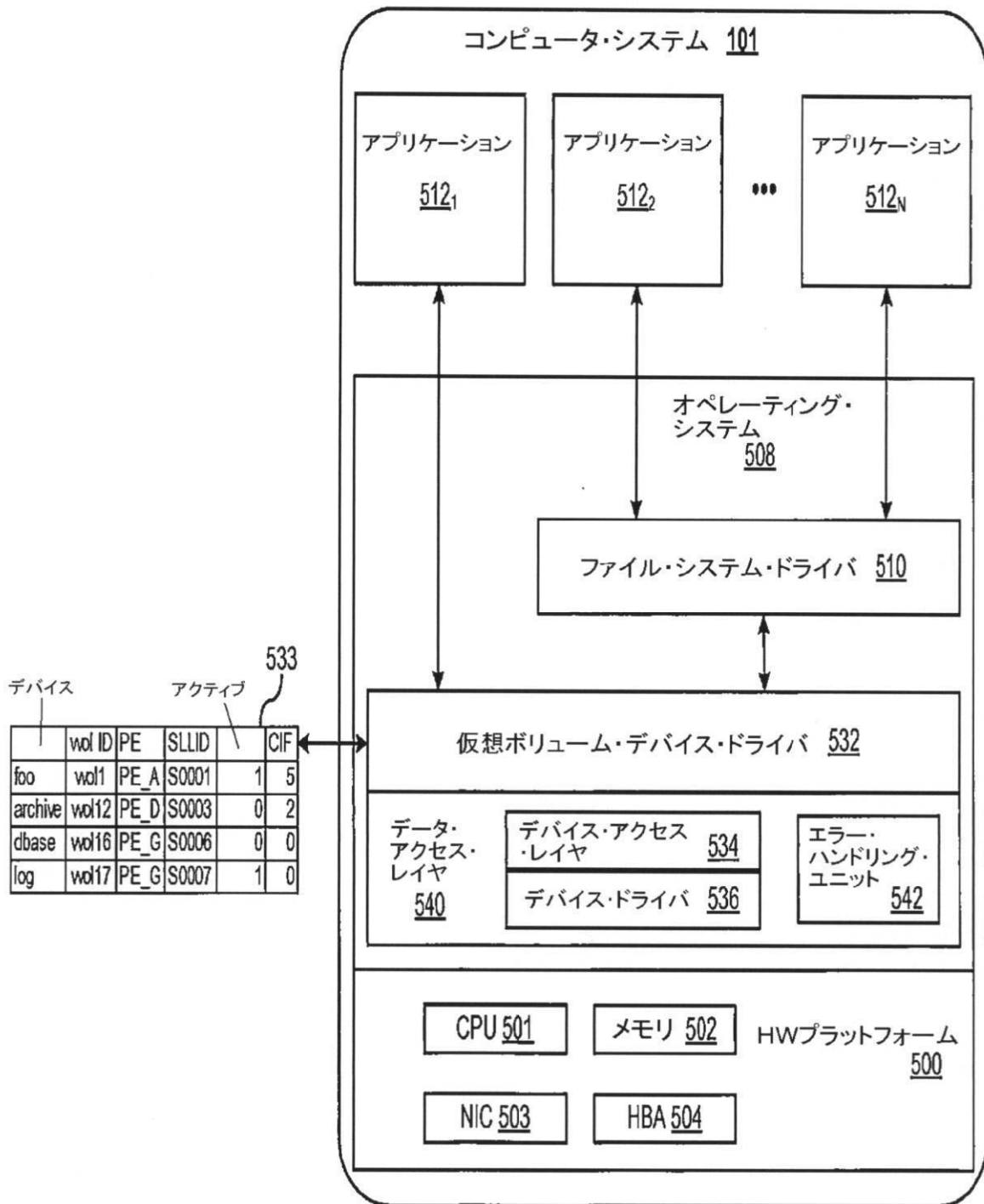
【 図 3 】



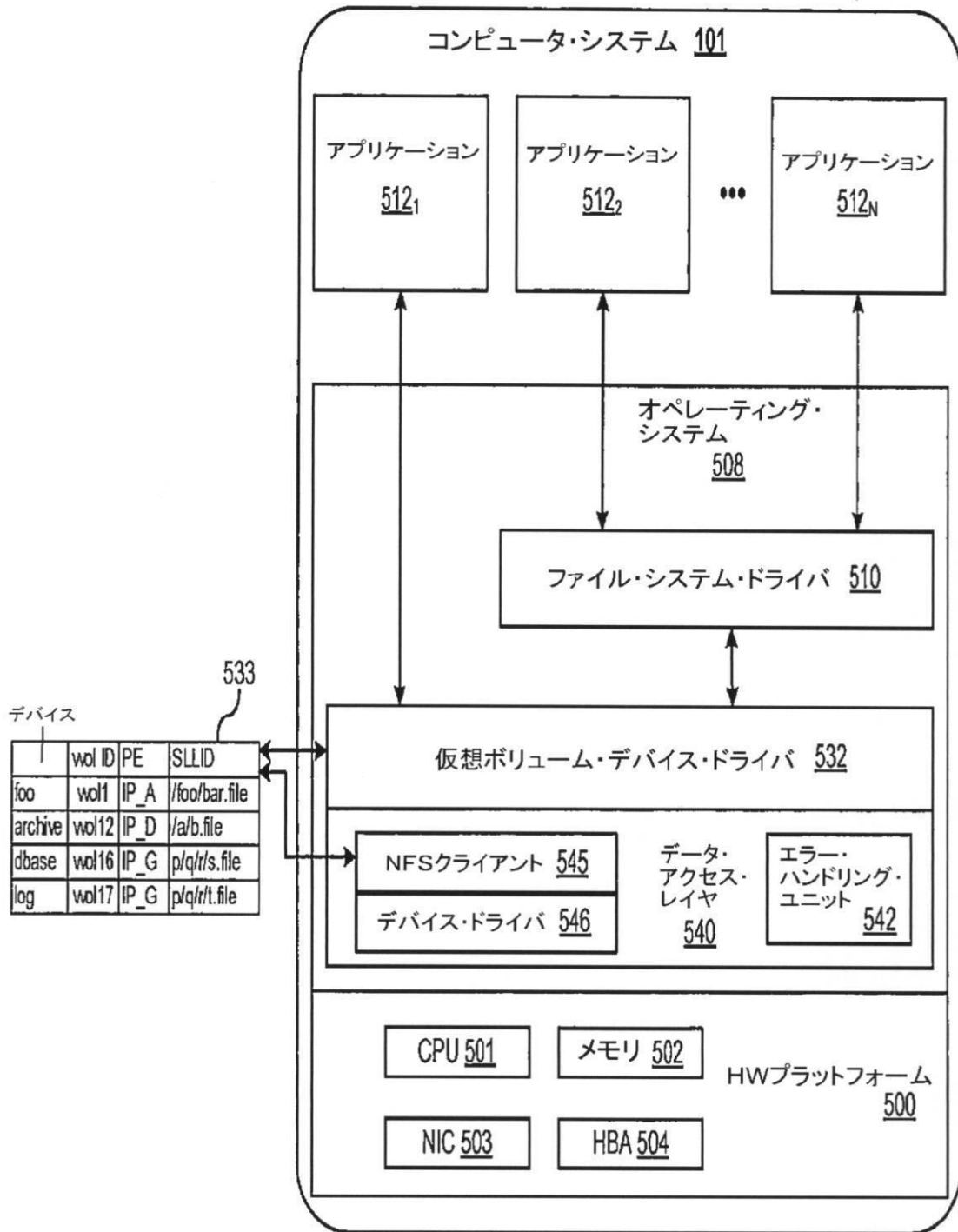
【 図 4 】



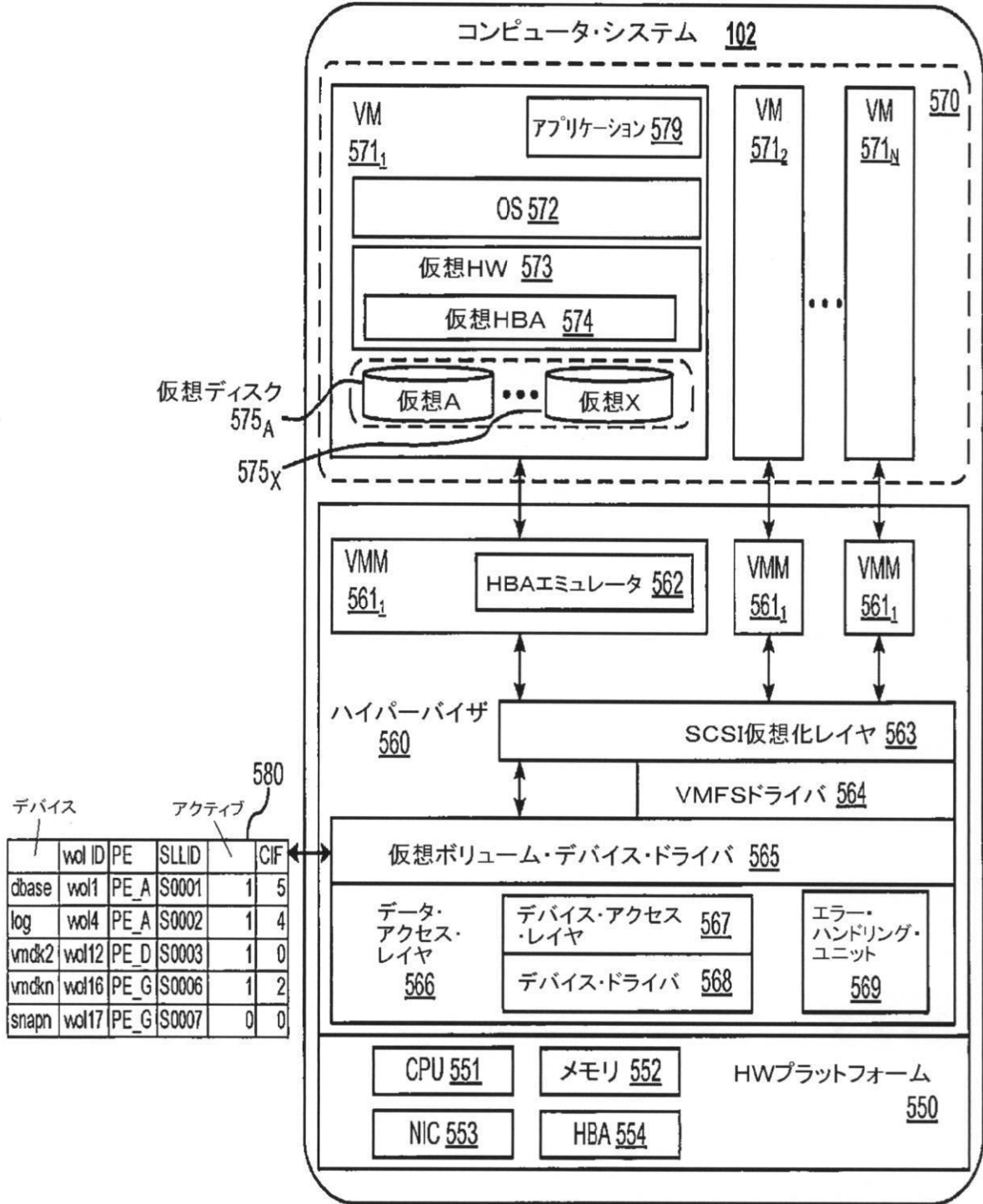
【図5A】



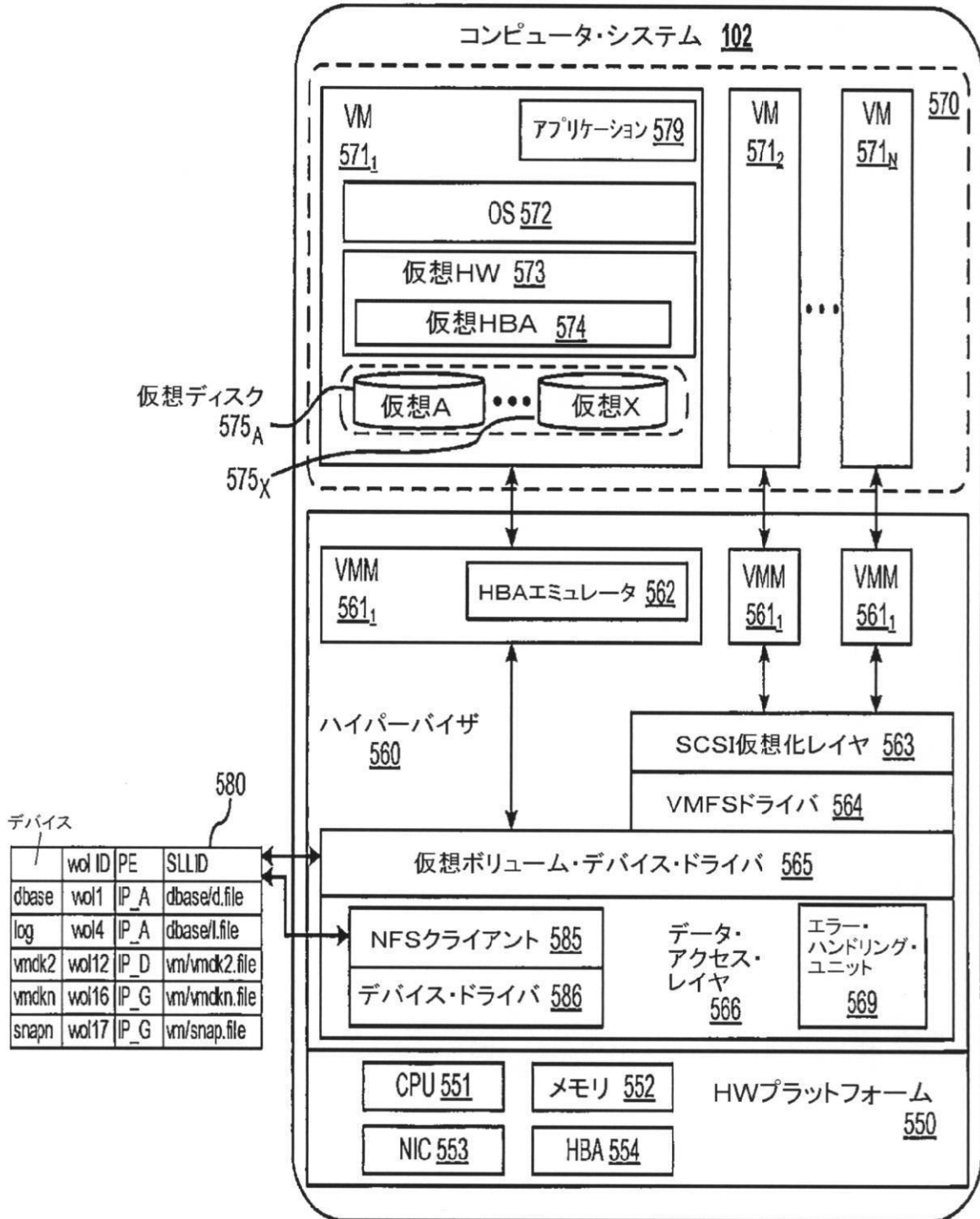
【図5B】



【図5C】

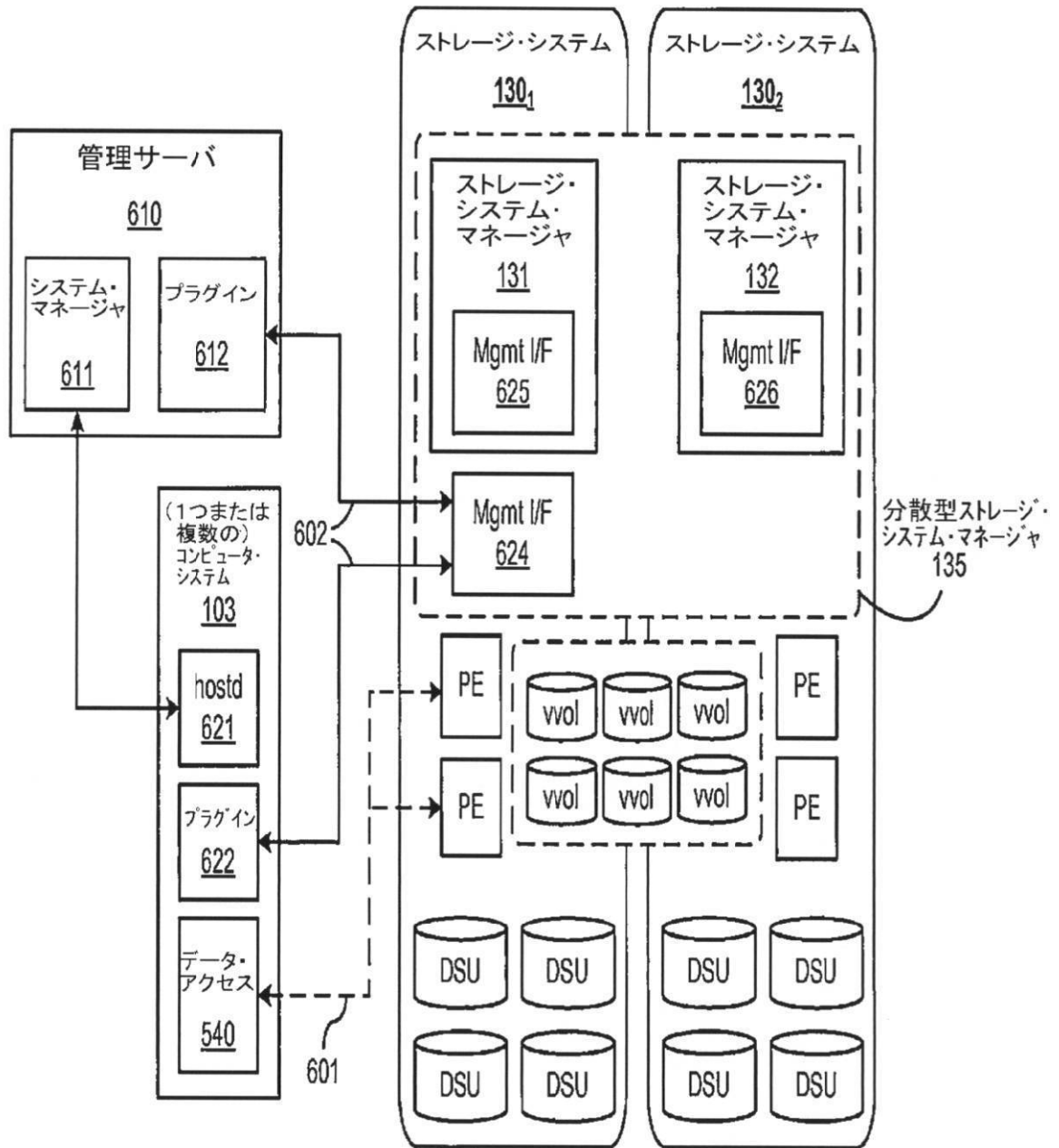


【図5D】

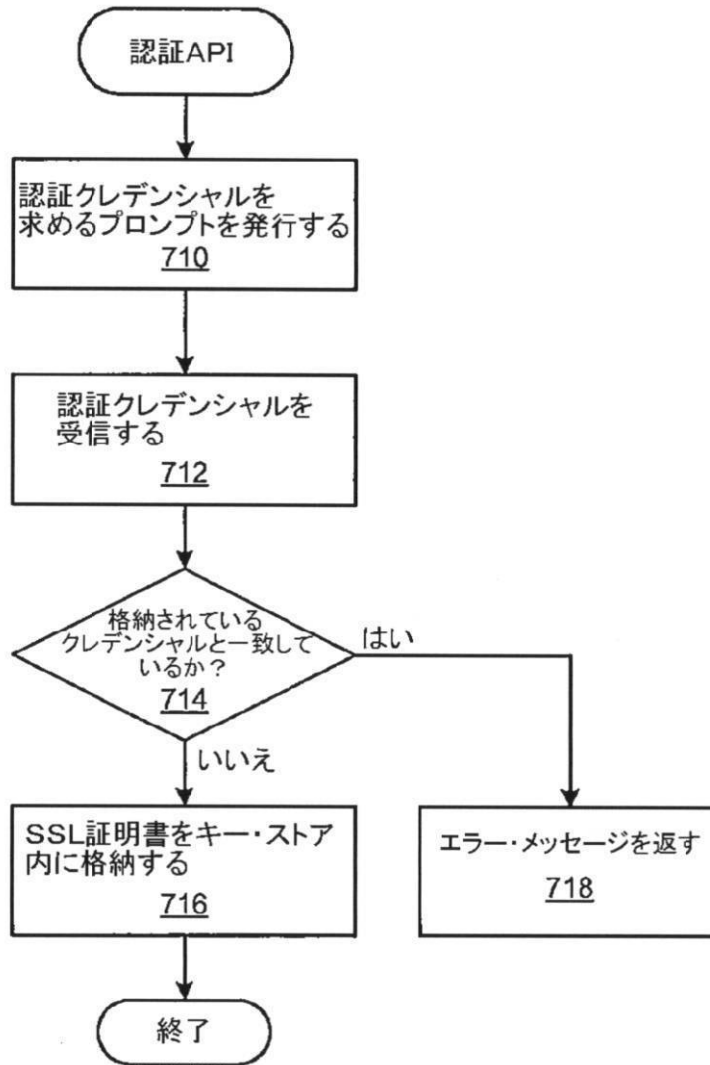




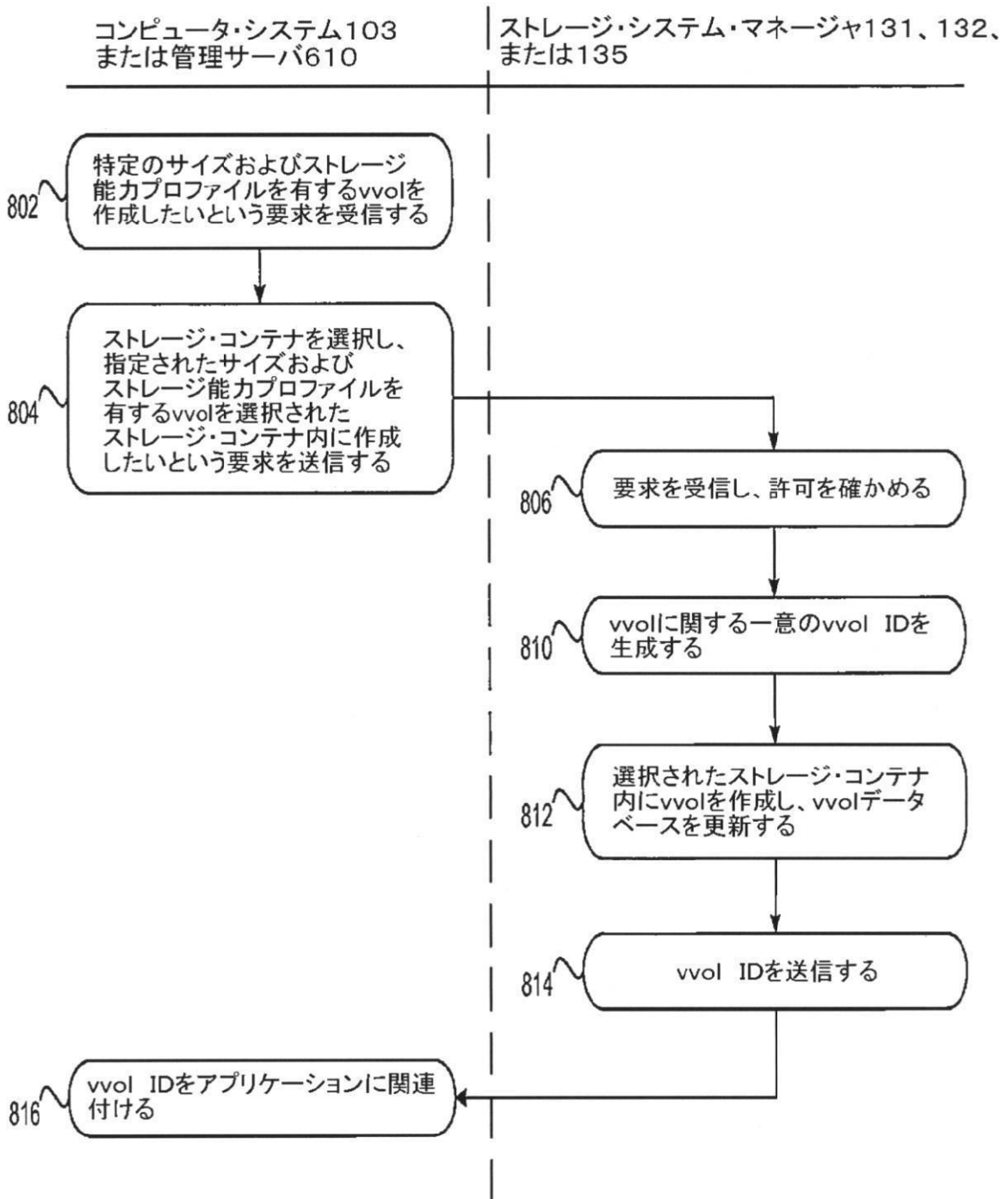
【 図 6 】



【 図 7 】



【図8】



【 図 9 】

図9A

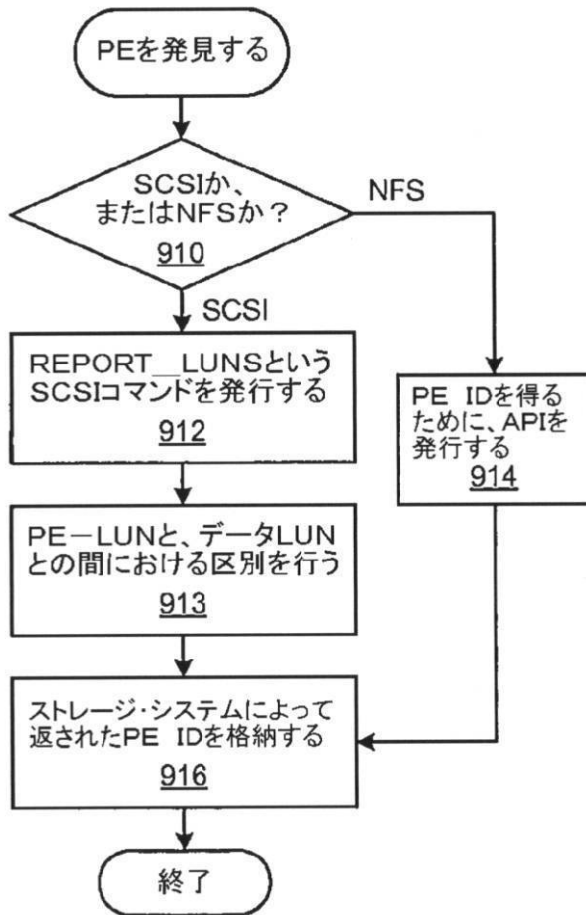
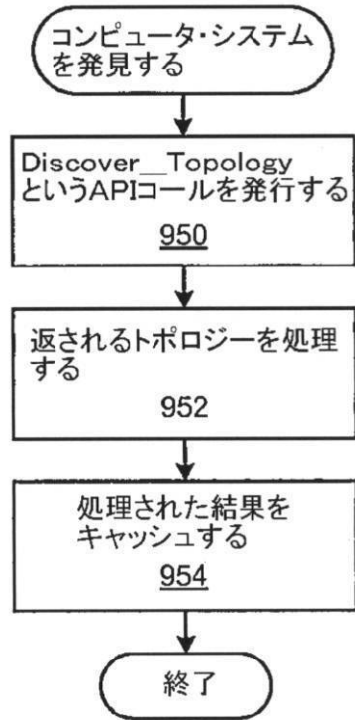
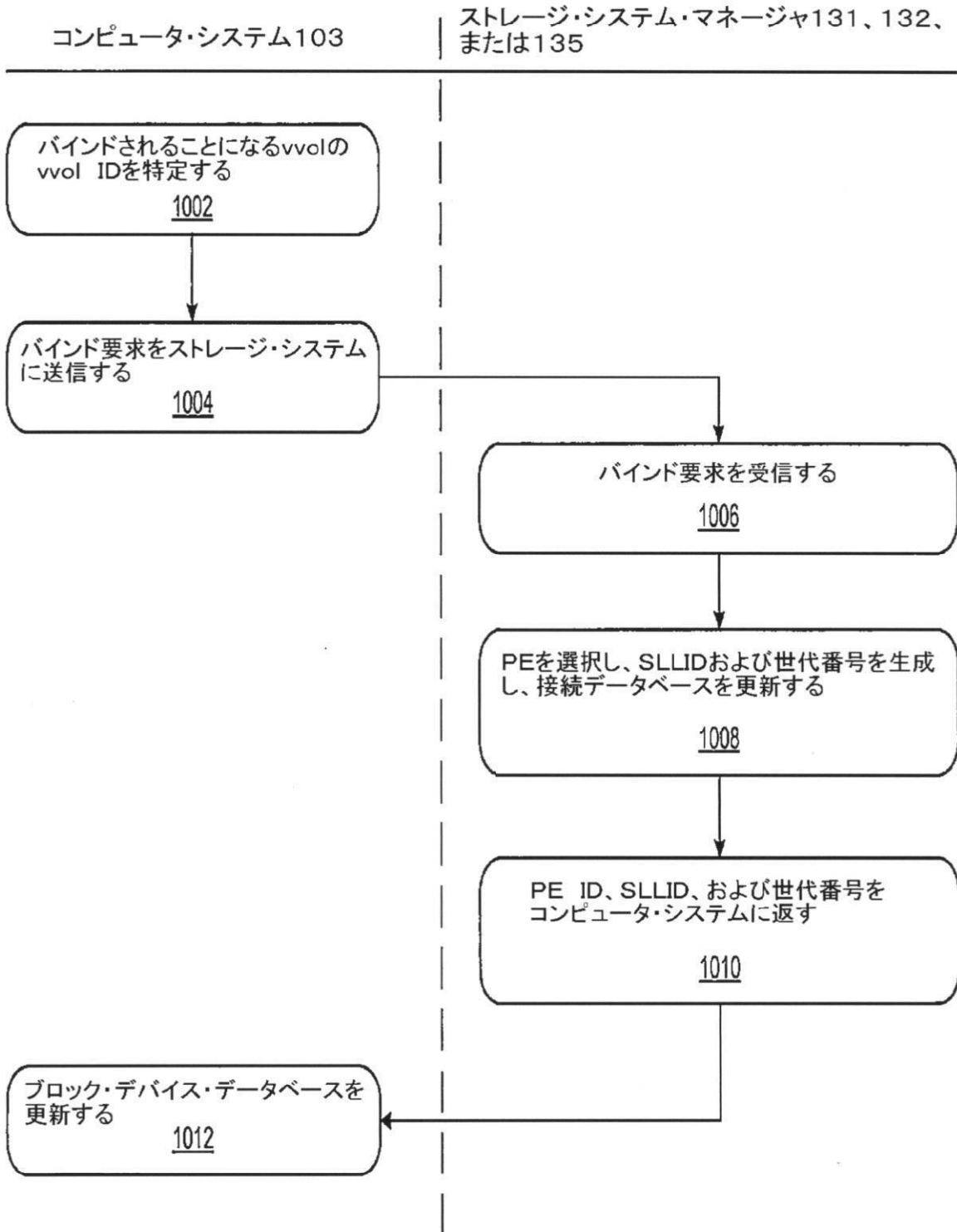


図9B



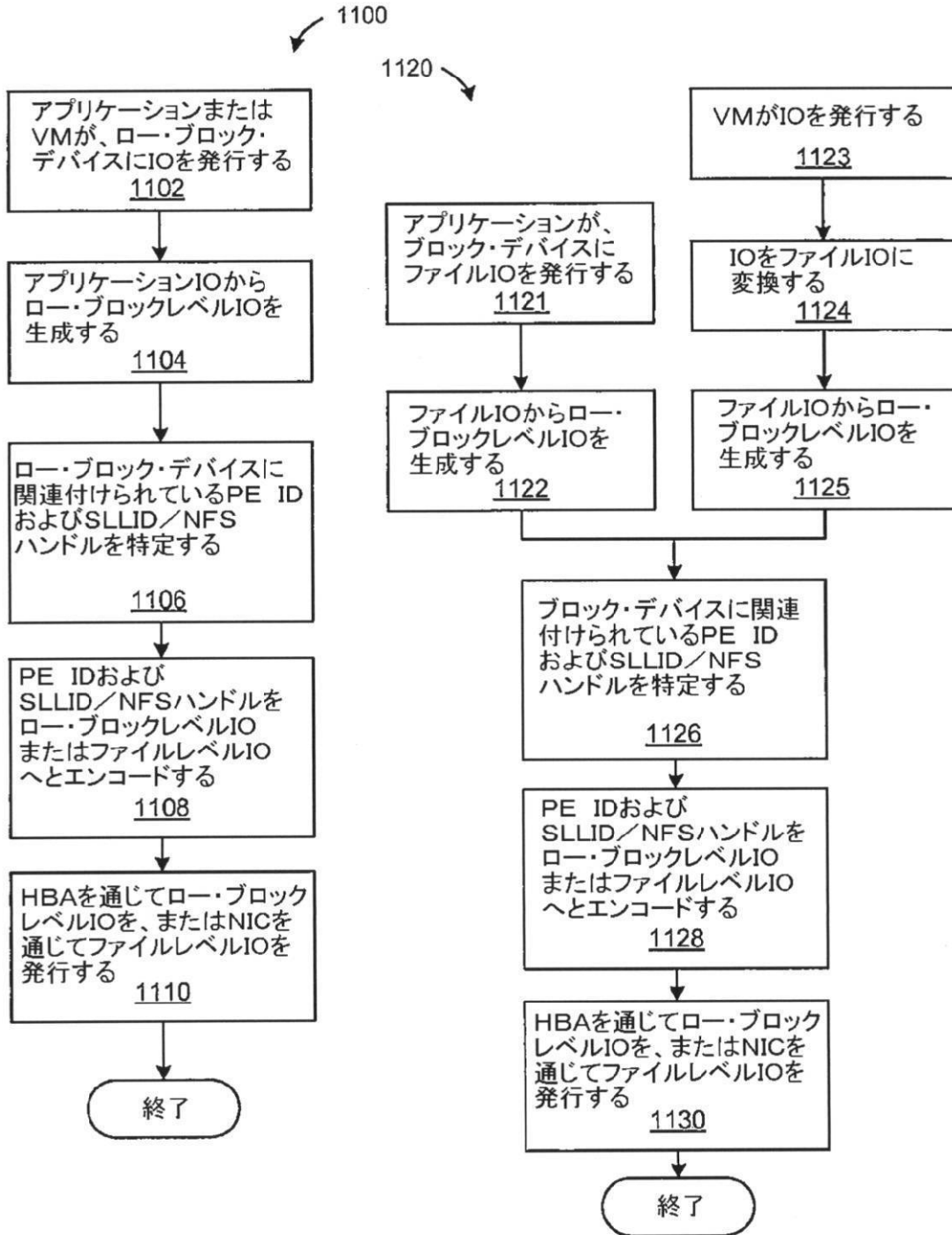
【図10】



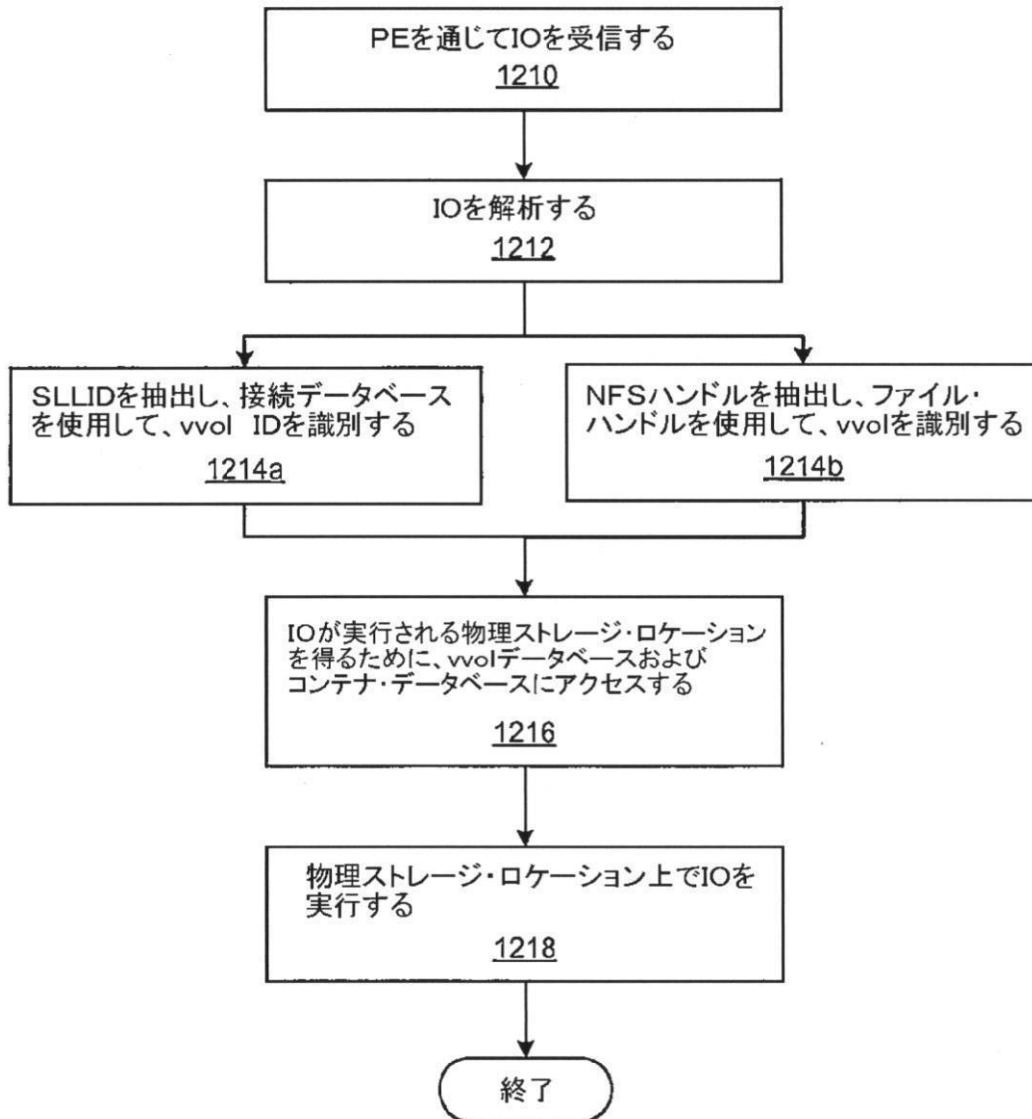
【図 11】

図11A

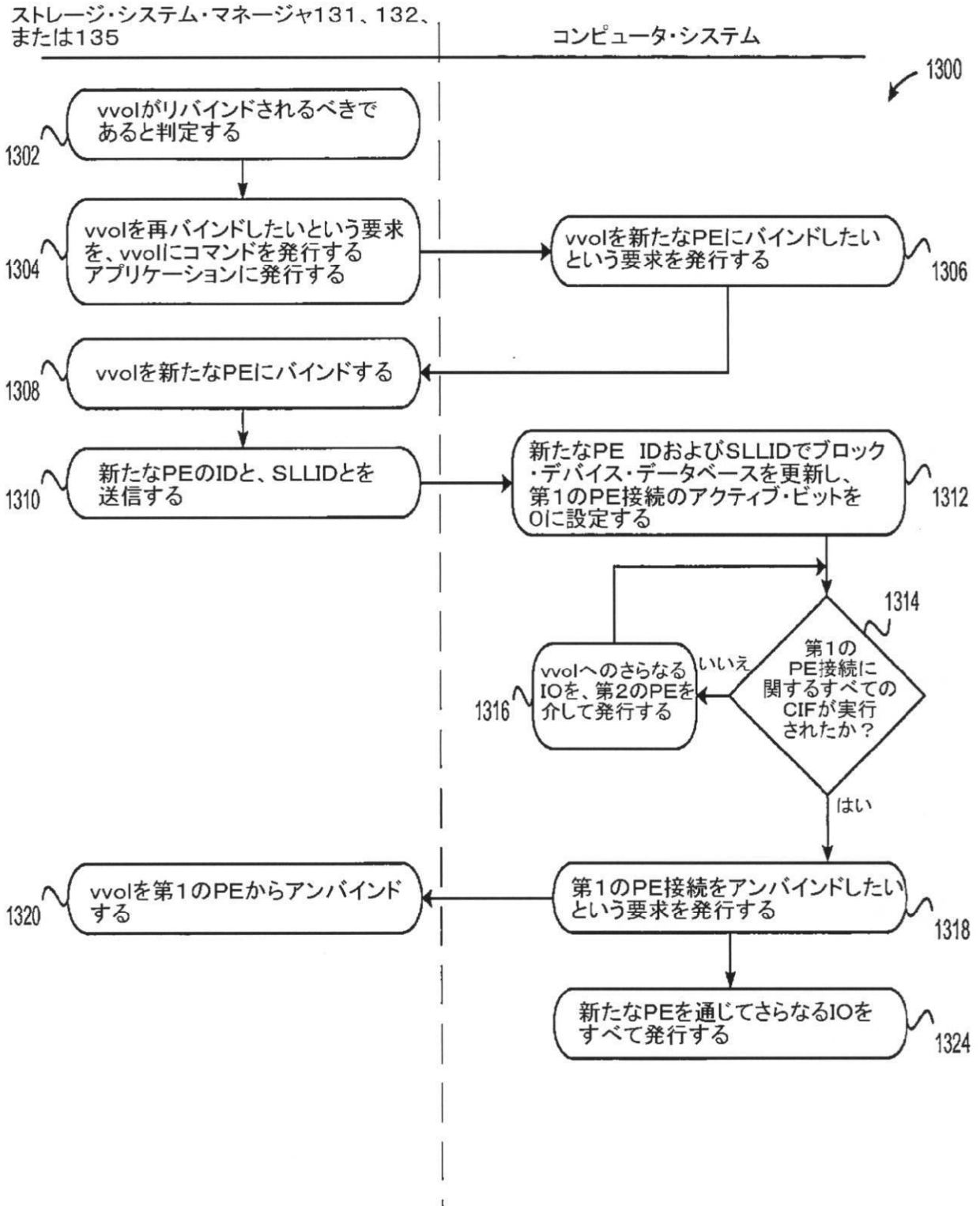
図11B



【 図 1 2 】



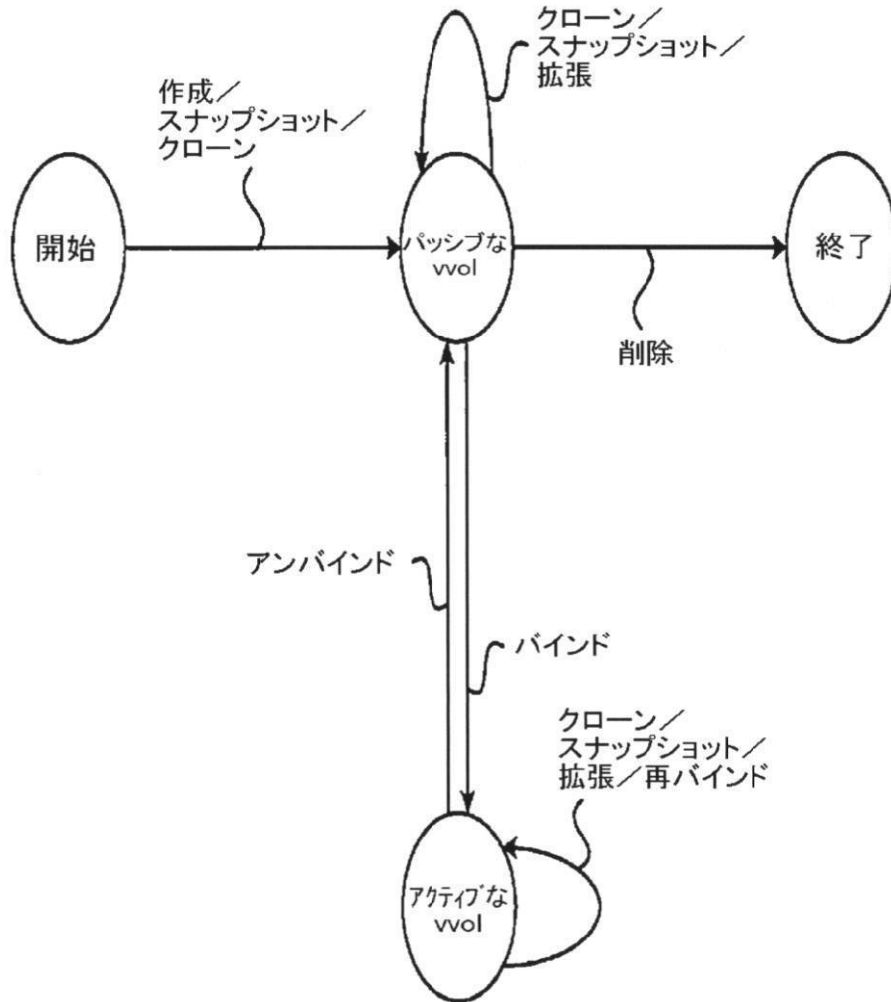
【図13】



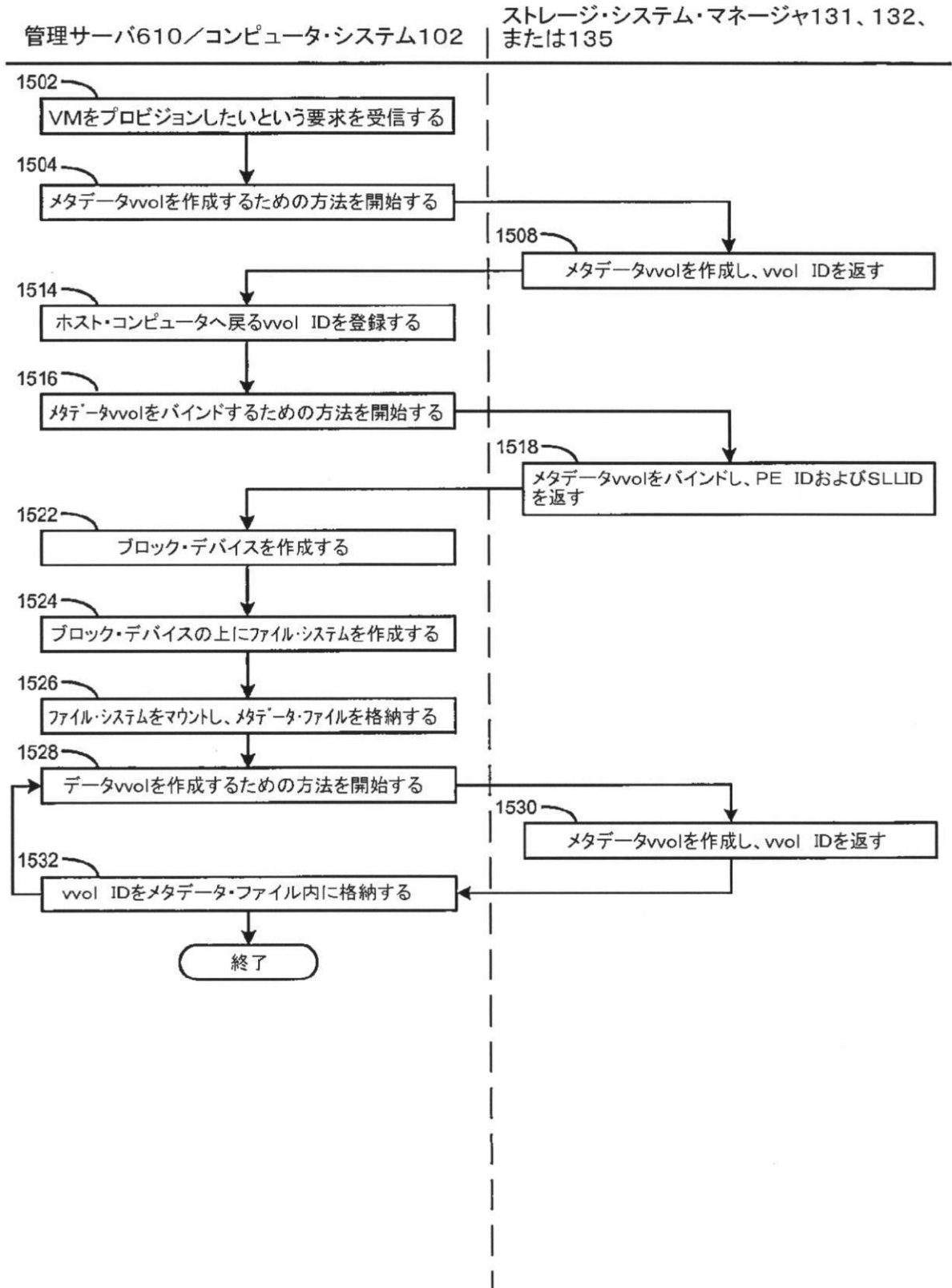


【 図 1 4 】

vvolのライフ・  
サイクルおよび  
帯域外コマンド・セット  
1400



【図15】



【 図 1 6 】

図16A

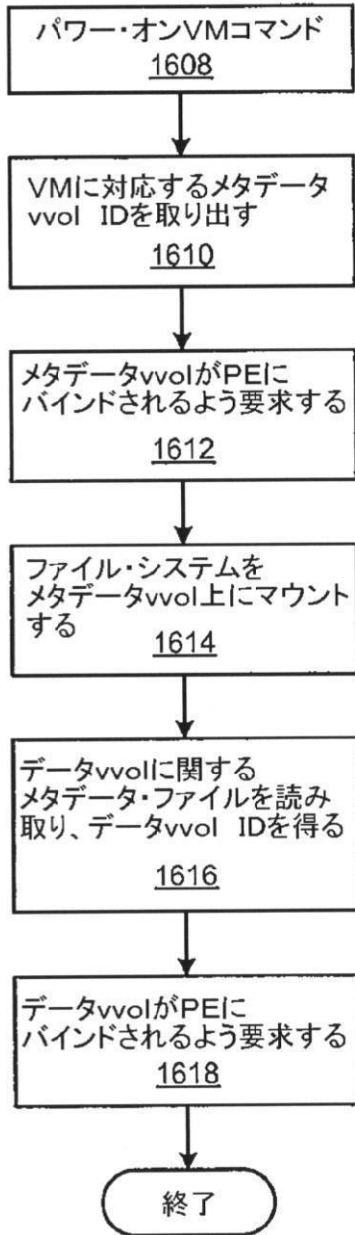
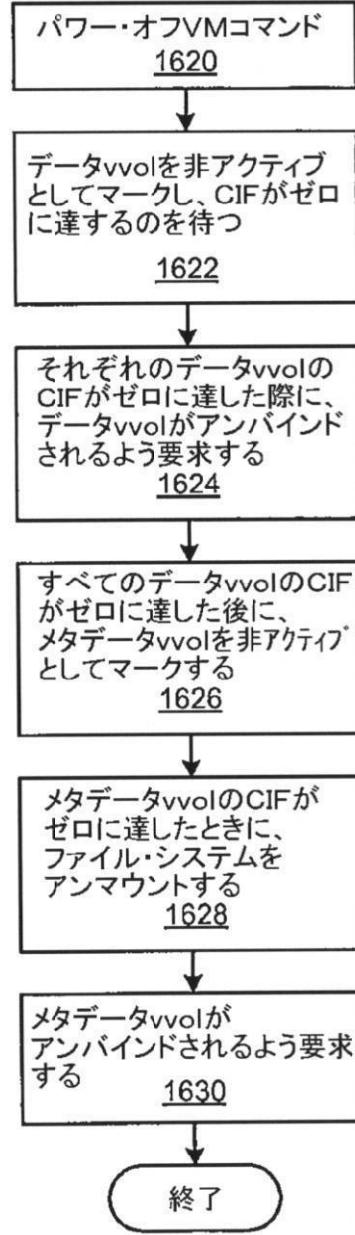
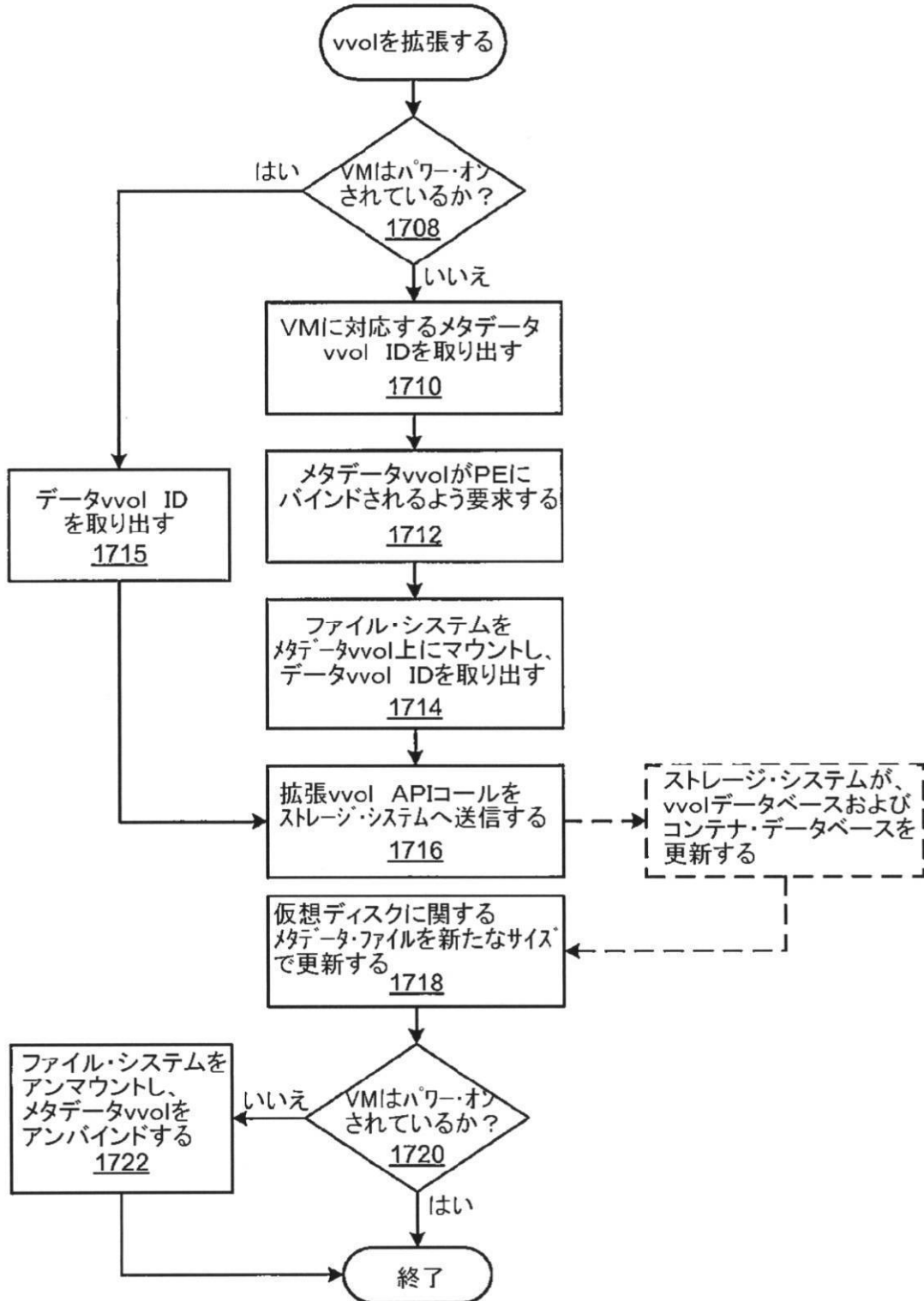


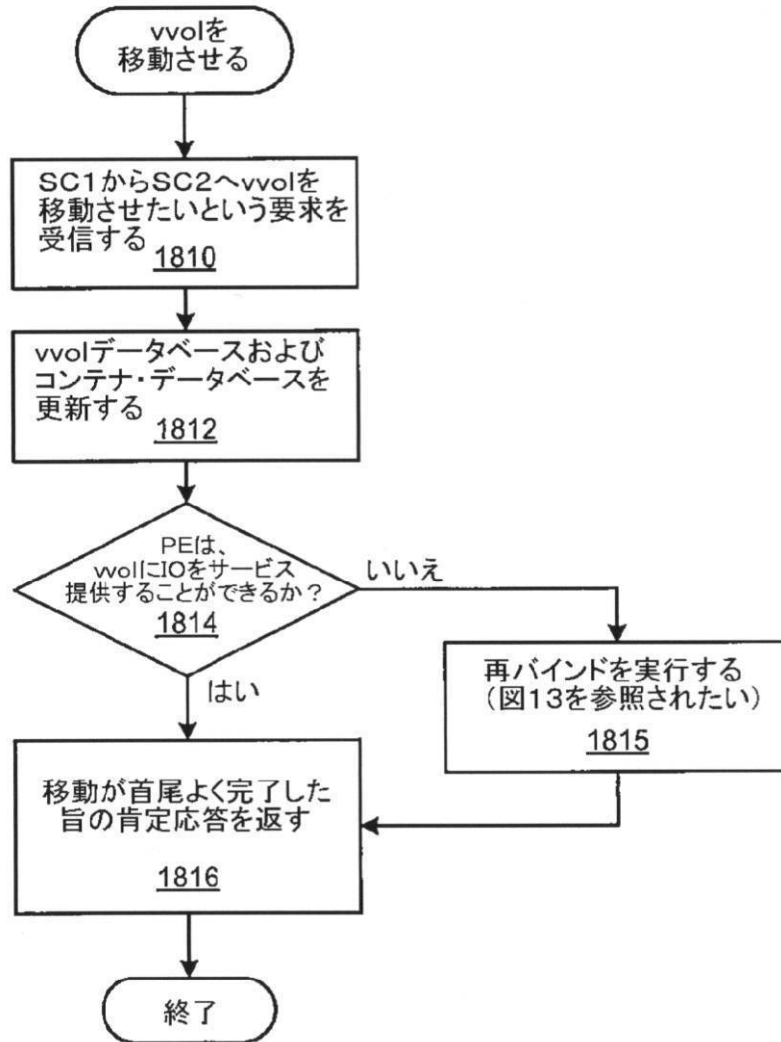
図16B



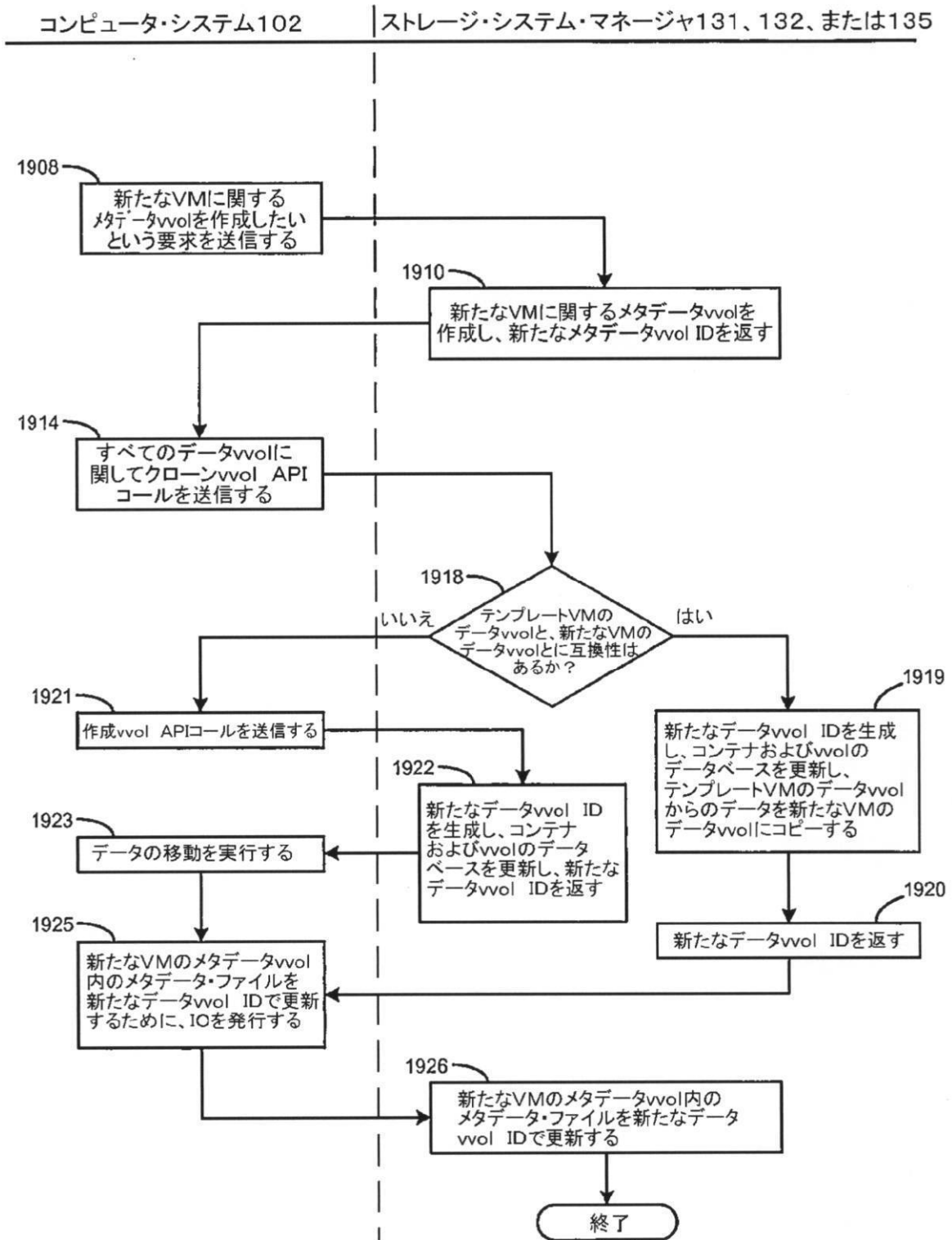
【図17】



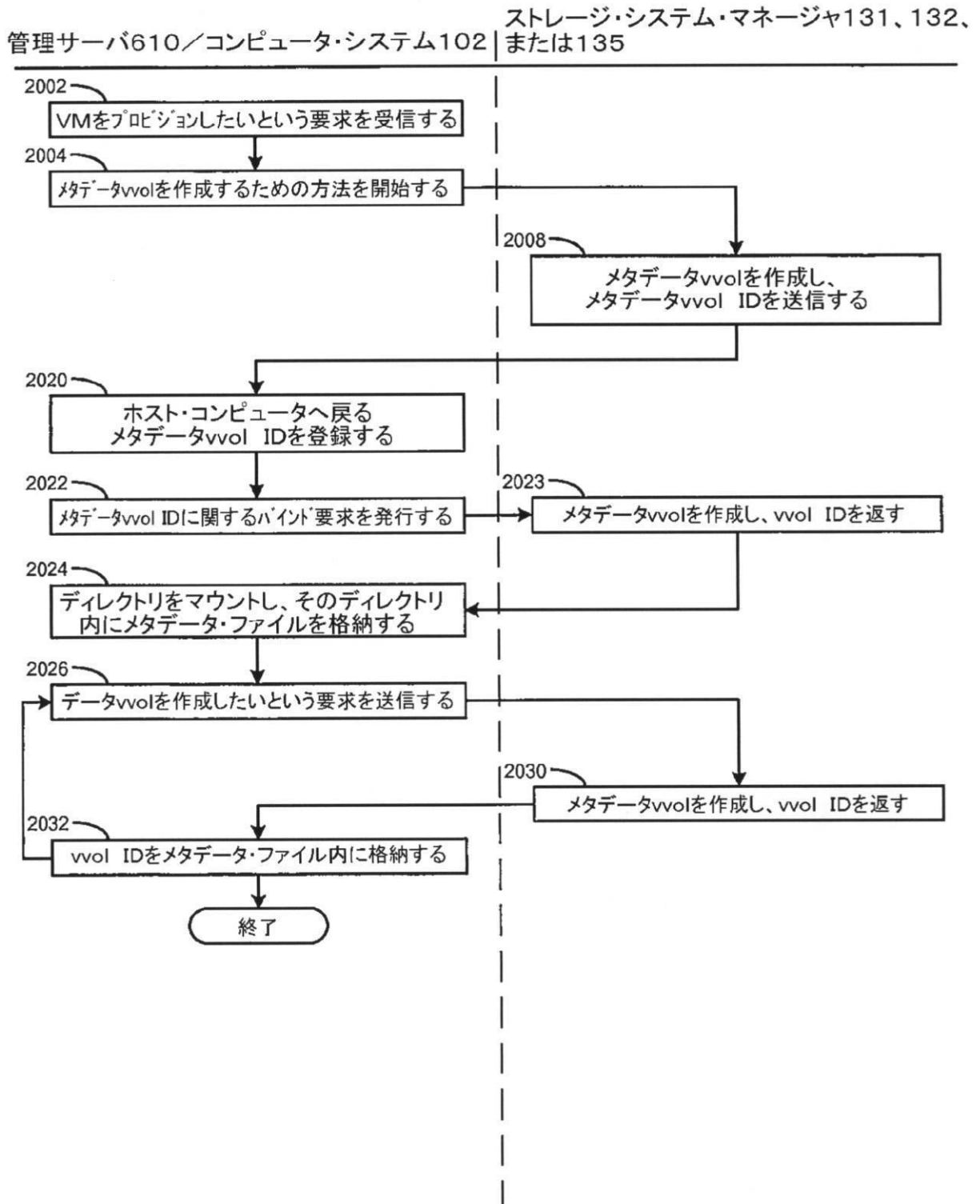
【図18】



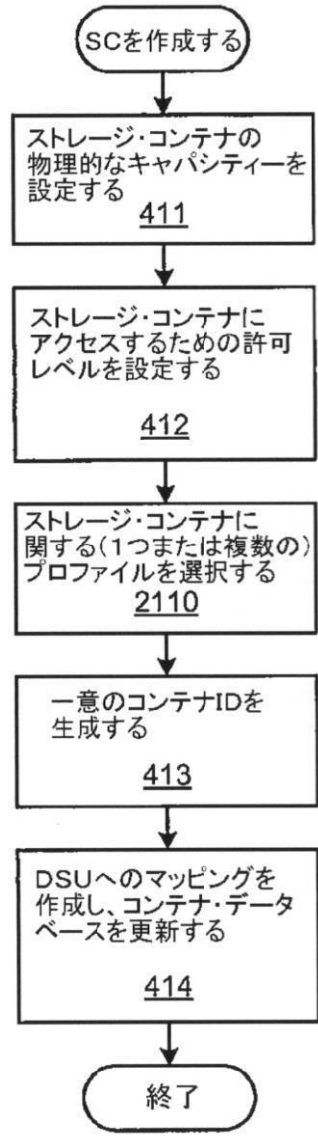
【図19】



【図 20】

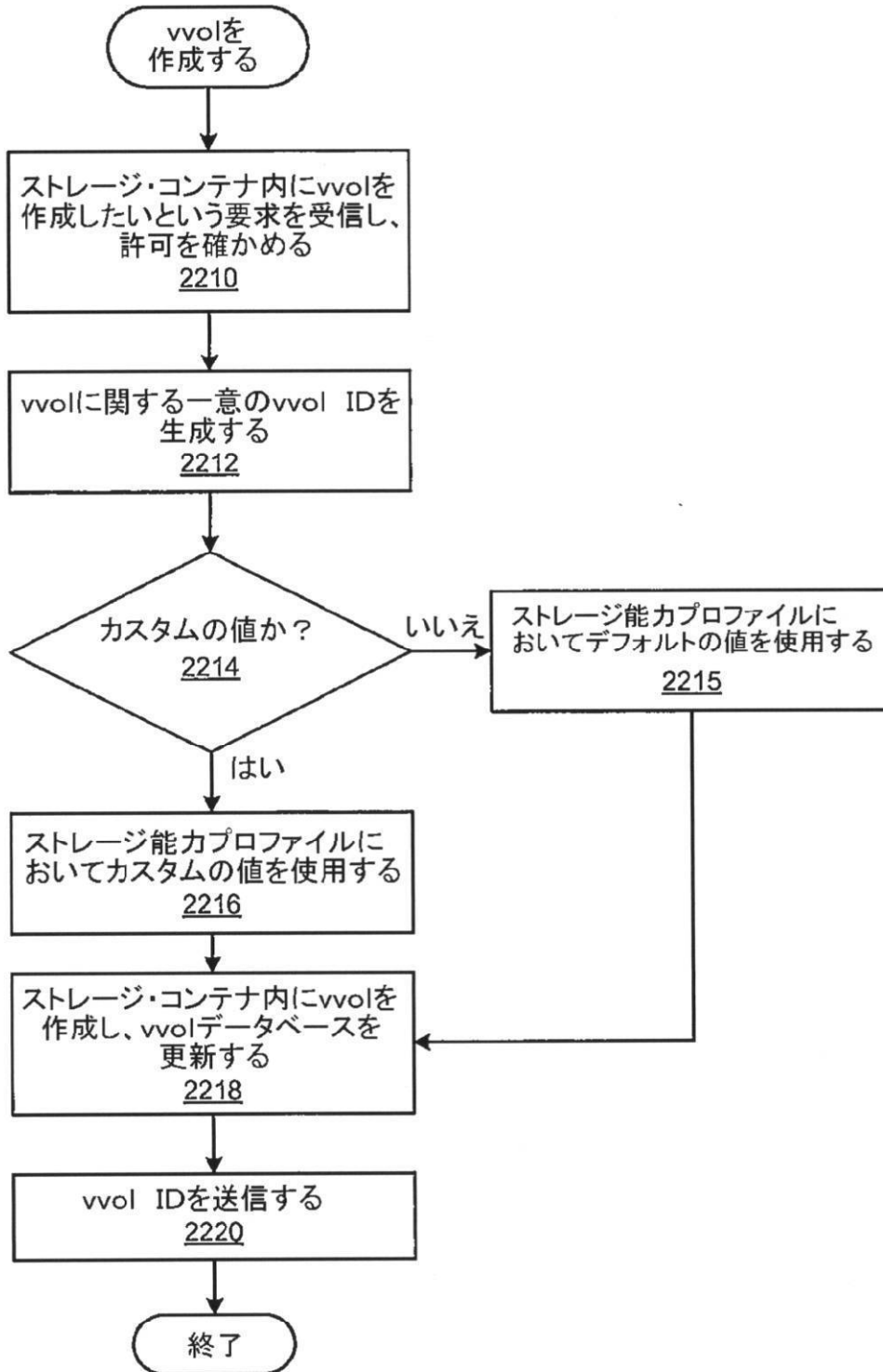


【図 2 1】

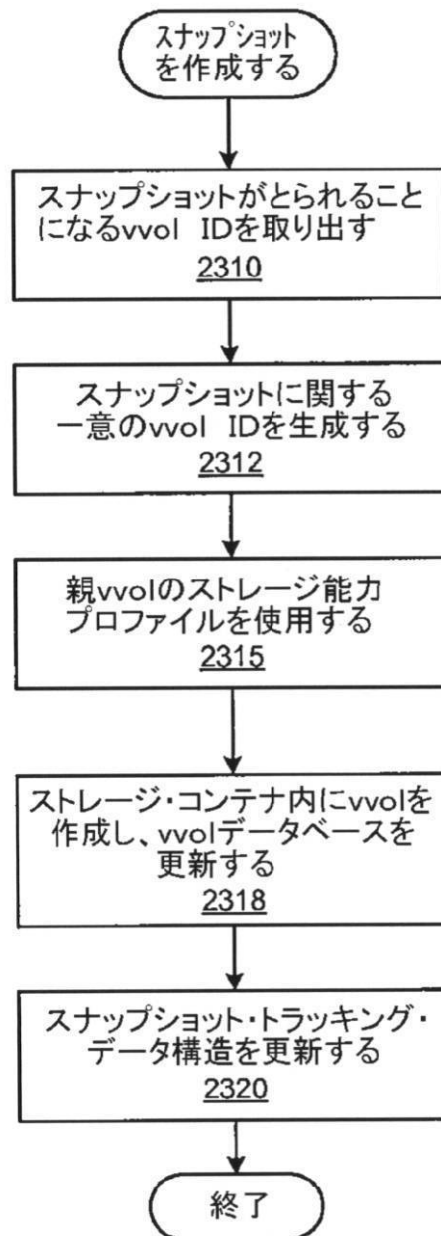




【図 2 2】



【図 23】



## 【 国際調査報告 】

## INTERNATIONAL SEARCH REPORT

International application No

PCT/US2012/051872

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> INV. G06F3/06 G06F9/00 H04L29/06 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) G06F H04L		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal, WPI Data		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	WO 2011/070605 A1 (HITACHI LTD [JP]; SUGINO SHOJI [JP]; URABE KIICHIRO [JP]; TANINAKA HIR) 16 June 2011 (2011-06-16) paragraph [0011] paragraphs [0014] - [0024] figure 1 paragraphs [0025] - [0038] figure 2 paragraph [0039] figure 3 paragraphs [0042] - [0044] figure 4 paragraphs [0045] - [0055] figures 5,6 paragraph [0062] figure 8 ----- -/--	1-18
<input checked="" type="checkbox"/>	Further documents are listed in the continuation of Box C.	<input checked="" type="checkbox"/> See patent family annex.
* Special categories of cited documents : "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed		"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family
Date of the actual completion of the international search 27 November 2012		Date of mailing of the international search report 04/12/2012
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer De Ceulaer, Bart

## INTERNATIONAL SEARCH REPORT

International application No PCT/US2012/051872
---

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	WO 01/80013 A1 (STOREAGE NETWORKING TECHNOLOGI [IL]; NAHUM NELSON [IL]) 25 October 2001 (2001-10-25) abstract page 10, line 10 - page 12, line 5 figure 1 -----	1-18
A	US 2010/153947 A1 (HARUMA TOSHIYUKI [JP]) 17 June 2010 (2010-06-17) paragraphs [0076] - [0084] figure 2 paragraph [0090] figure 3G -----	17,18

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No

PCT/US2012/051872

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 2011070605 A1	16-06-2011	US 2011258377 A1 WO 2011070605 A1	20-10-2011 16-06-2011
WO 0180013 A1	25-10-2001	AU 4679901 A CA 2405405 A1 EP 1282861 A1 US 2003236945 A1 WO 0180013 A1	30-10-2001 25-10-2001 12-02-2003 25-12-2003 25-10-2001
US 2010153947 A1	17-06-2010	JP 4705982 B2 JP 2010140273 A US 2010153947 A1	22-06-2011 24-06-2010 17-06-2010

## フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN

(特許庁注：以下のものは登録商標)

1. イーサネット
2. Linux

(72)発明者 バト、ラジェシュ  
 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 3401

(72)発明者 アチャリヤ、サンジェイ  
 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 3401

(72)発明者 バガーニ、サティヤム ピー .  
 アメリカ合衆国 95117 カリフォルニア州 サンノゼ ダフォディル ウェイ 893

(72)発明者 シー、チャオ - チャン  
 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 3401