

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4224077号  
(P4224077)

(45) 発行日 平成21年2月12日(2009.2.12)

(24) 登録日 平成20年11月28日(2008.11.28)

(51) Int.Cl. F 1  
G 0 6 F 3 / 0 6 (2006.01) G 0 6 F 3 / 0 6 3 0 1 Z

請求項の数 7 (全 30 頁)

(21) 出願番号	特願2006-103256 (P2006-103256)	(73) 特許権者	000003078 株式会社東芝
(22) 出願日	平成18年4月4日(2006.4.4)		東京都港区芝浦一丁目1番1号
(65) 公開番号	特開2007-279898 (P2007-279898A)	(73) 特許権者	301063496 東芝ソリューション株式会社
(43) 公開日	平成19年10月25日(2007.10.25)		東京都港区芝浦一丁目1番1号
審査請求日	平成18年4月4日(2006.4.4)	(74) 代理人	100058479 弁理士 鈴江 武彦
		(74) 代理人	100091351 弁理士 河野 哲
		(74) 代理人	100088683 弁理士 中村 誠
		(74) 代理人	100108855 弁理士 蔵田 昌俊

最終頁に続く

(54) 【発明の名称】 ストレージシステム

(57) 【特許請求の範囲】

【請求項1】

物理ボリュームを有する少なくとも2つのノードを含む複数のノードから構成され、前記少なくとも2つのノードが有する物理ボリュームを利用して、ホストに対して複数の論理エクステントから構成される論理ボリュームを提供するストレージシステムにおいて、

前記複数のノードの少なくとも2つは、前記ホストに対して前記論理ボリュームを提供し、当該ホストからのアクセスを受け付けるコントローラを有し、

前記少なくとも2つのノードは、前記論理ボリュームを構成する論理エクステントを物理エクステントにマッピングして前記物理ボリューム内に格納するストレージを有し、

前記複数のノードの少なくとも2つは、前記論理ボリュームを構成する論理エクステントのうち、自身が担当すべき論理エクステントを格納しているストレージの所在に関する前記コントローラからの問い合わせに対して回答する所在管理サービスを実行する所在管理サーバを有し、

前記コントローラは、前記論理ボリュームを構成する論理エクステントに対応する所在管理サーバを示す第1のマッピングテーブルを有し、

前記所在管理サーバは、前記論理ボリュームを構成する論理エクステントのうち、当該サーバが担当すべき論理エクステントについて、その論理エクステントがどのノードのストレージに格納されているかを示す第2のマッピングテーブルを有し、

前記ストレージは、論理エクステントを、当該ストレージが有する物理ボリュームのどの物理エクステントに対応させて格納したかを示す第3のマッピングテーブルを有し、

10

20

前記コントローラは、前記ホストから前記論理ボリューム内の任意の論理エクステントをアクセス先として指定するアクセス要求が与えられた場合、当該コントローラが有する前記第1のマッピングテーブルを参照することにより当該アクセス先の論理エクステントに対応する所在管理サーバを特定し、特定された所在管理サーバに対して、当該論理エクステントがどのストレージに格納されているかを問い合わせ、

前記所在管理サーバは前記コントローラから問い合わせを受けた場合、当該サーバが有する前記第2のマッピングテーブルを参照することにより問い合わせられたストレージを特定して、当該ストレージを示す情報を前記問い合わせに対する回答として問い合わせ元の前記コントローラに通知し、

前記コントローラは、前記所在管理サーバから前記問い合わせに対する回答として通知された情報の示すストレージに対して前記ホストから要求された論理エクステントへのアクセスを要求し、

前記ストレージは、前記コントローラから論理エクステントへのアクセスが要求された場合、当該ストレージが有する前記第3のマッピングテーブルを参照することにより、当該ストレージが有する前記物理ボリューム内で前記要求された論理エクステントに対応させられている物理エクステントを特定して、当該特定された物理エクステントに対してアクセスしてアクセス要求元の前記コントローラに応答し、

論理エクステントを構成する物理エクステント内のデータを前記ストレージをまたがって別の物理エクステントに移動またはコピーするエクステント移動またはコピーが発生した場合、移動元及び移動先またはコピー先のストレージが有する前記第3のマッピングテーブルは当該エクステント移動またはコピーが反映されるように更新され、当該論理エクステントを格納しているストレージの所在の問い合わせに対して回答する所在管理サーバが有する前記第2のマッピングテーブルは、当該論理エクステントを格納しているストレージとして移動先のストレージ、またはコピー元及びコピー先のストレージを示すように更新される

ことを特徴とするストレージシステム。

#### 【請求項2】

前記所在管理サーバが有する前記第2のマッピングテーブルによって示されている、ある論理エクステントがどのノードのストレージに格納されているかを示す情報を別の所在管理サーバが有する前記第2のマッピングテーブルに移動する所在管理サービス移動が発生した場合、移動元及び移動先の前記第2のマッピングテーブル、及び全ての前記コントローラが有する前記第1のマッピングテーブルは、当該所在管理サービス移動が反映されるように更新される

ことを特徴とする請求項1記載のストレージシステム。

#### 【請求項3】

前記所在管理サーバが有する前記第2のマッピングテーブルによって示されている、ある論理エクステントがどのノードのストレージに格納されているかを示す情報を別の所在管理サーバが有する前記第2のマッピングテーブルにコピーする所在管理サービスコピーが発生した場合、コピー先の前記第2のマッピングテーブル、及び全ての前記コントローラが有する前記第1のマッピングテーブルは、当該所在管理サービスコピーが反映されるように更新される

ことを特徴とする請求項1記載のストレージシステム。

#### 【請求項4】

前記第3のテーブルは、前記論理ボリュームの作成時の初期状態において当該論理ボリュームを構成する論理エクステントが格納される物理ボリュームのみを示し、

前記ストレージは、前記コントローラから論理エクステントへのアクセスが要求されて、当該ストレージが有する前記第3のマッピングテーブルを参照した結果、当該ストレージが有する前記物理ボリューム内で前記要求された論理ボリュームに物理エクステントが対応付けられていないことが判別された場合、前記物理ボリューム内の空き物理エクステントを当該論理エクステントに割り付けて前記第3のマッピングテーブルを更新する

ことを特徴とする請求項 1 記載のストレージシステム。

【請求項 5】

前記第 2 のテーブルは、前記論理ボリュームの作成時の初期状態において、当該サーバに対応する論理エクステントに対してストレージの未割り付けの状態にあることを示し、

前記第 3 のテーブルは、前記論理ボリュームの作成時の初期状態において、当該論理ボリュームを構成する論理エクステントに対して物理ボリューム及び物理エクステントの未割り付けの状態にあることを示し、

前記所在管理サーバは、前記コントローラからの問い合わせに応じて当該サーバが有する前記第 2 のマッピングテーブルを参照した結果、前記アクセス先の論理エクステントに対してストレージの未割り付けの状態にあることが判別された場合、空き物理エクステントを有するストレージを当該論理エクステントに割り付けて、前記第 2 のマッピングテーブルを更新し、

前記ストレージは、前記コントローラから論理エクステントへのアクセスが要求されて、当該ストレージが有する前記第 3 のマッピングテーブルを参照した結果、当該ストレージが有する前記物理ボリューム内で前記要求された論理ボリュームに物理エクステントが対応付けられていないことが判別された場合、当該ストレージの前記物理ボリューム内の空き物理エクステントを当該論理エクステントに割り付けて、前記第 3 のマッピングテーブルを更新する

ことを特徴とする請求項 1 記載のストレージシステム。

【請求項 6】

前記コントローラは、前記問い合わせ及び当該問い合わせに対する前記所在管理サーバからの回答に関する情報を格納するキャッシュを有し、

前記コントローラは、前記所在管理サーバに論理エクステントがどのストレージに格納されているかを問い合わせる前に当該問い合わせに対する回答に関する情報が前記キャッシュに格納されているかを調べ、格納されている場合には当該キャッシュに格納されている回答に関する情報の示すストレージに対して前記ホストから要求された論理エクステントへのアクセスを要求する

ことを特徴とする請求項 1 記載のストレージシステム。

【請求項 7】

前記所在管理サーバが有する前記第 2 のマッピングテーブルの示す論理エクステントが格納されているストレージと、当該所在管理サーバが有するキャッシュに格納されている回答に関する情報の示す当該論理エクステントが格納されているストレージとが一致しているかを調べて、不一致の場合に当該キャッシュに格納されている当該情報を無効化する無効化手段を更に有することを特徴とする請求項 6 記載のストレージシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、物理ボリュームを有する少なくとも 1 つのノードを含む複数のノードから構成されるストレージシステムに関する。

【背景技術】

【0002】

近年、相互に接続された複数のストレージ装置（例えばディスクストレージ装置）を有するストレージシステムが開発されている。この種のストレージシステムの 1 つとして、ストレージクラスタ（ストレージクラスタシステム）が知られている。

【0003】

ストレージクラスタの特徴は、複数のストレージ装置（以下、ノードと称する）の各々に存在する記憶領域（以下、物理ボリュームと称する）を、当該ノードを超えて統合する（束ねる）ことで、ノード間にまたがった横断的・仮想的な少なくとも 1 つの記憶領域（以下、論理ボリュームと称する）をホストに提供することにある。ここでのボリュームは、ブロック単位でアクセスされるブロックボリュームであるとする。

10

20

30

40

50

## 【 0 0 0 4 】

ストレージクラスタを構成するノードは、常時稼働のストレージ装置である。つまり、ノードの参加・離脱は、運用中の計画的参加・離脱や、故障による突然の離脱による程度である。このためストレージクラスタでは、例えば毎日の電源 ON / OFF によるような、頻繁なノードの参加・離脱はない。

## 【 0 0 0 5 】

ストレージクラスタにおいては、データ（の断片）を物理ボリューム内や、別のノードの物理ボリュームをまたがって、（自動/手動を問わず）適切に移動/コピーすることは、頻繁に起きる。これは、ホストからの論理ボリュームへのアクセス透過性を保障しながら、ストレージクラスタを構成するノードの突発的停止やディスク故障などに備えたデータの冗長化や、アクセス性能・信頼性の向上などの目的による。

10

## 【 0 0 0 6 】

一方、ホストは、ストレージクラスタ内の任意の1つのノードに対してパスを張ってアクセスするだけで、複数のノードをまたがって構成された論理ボリューム全体に対してアクセスすることができる。ここで、ホストがパスを張っているノードで障害が発生して、当該障害が発生したノード（障害ノード）が停止したもとする。この場合でも、障害ノードに格納されているデータがストレージクラスタ内の他のノード（生存ノード）に冗長化されていれば、ホストは任意の生存ノードに対してパスを張り直すことで、論理ボリュームに対するアクセスを継続することができる。

20

## 【 0 0 0 7 】

このようなストレージクラスタを構築するには、論理ボリューム上のデータが、どのノードのどの物理ボリューム上の、どの位置に配置されているかを管理するための仕組みが必要になる。

## 【 0 0 0 8 】

一般に、ストレージクラスタのノードは、機能的にコントローラとストレージとから構成される。コントローラはホストとのパスを管理しアクセスを受け付ける。ホストはどのコントローラに対してパスを張っても、同じ論理ボリュームにアクセスすることができる。ストレージはデータを蓄える物理ボリュームを有する。

## 【 0 0 0 9 】

ボリューム（ブロックボリューム）は、エクステントと呼ばれる、固定長のブロックの集合により構成される。各ノードのストレージは1つの物理ボリュームを管理する。一方、各ノードのコントローラは、物理ボリュームを構成する物理エクステントと、論理ボリュームの論理エクステントとをマッピングテーブルにより対応付けて管理する。これにより、論理ボリューム（論理ボリューム空間）が構成される。このようなマッピングテーブルによる物理エクステントと論理エクステントとの対応付けは、例えば特許文献1にも記載されている。

30

## 【 0 0 1 0 】

上述のマッピングテーブルを、論理エクステント - 物理エクステントマッピングテーブル（L P M T : Logical Extent to Physical Extent Mapping Table）と呼ぶ。L P M T は、論理ボリュームの論理エクステントが、どのノード（ストレージ）の物理ボリュームの、どの物理エクステントから構成されているか、その対応付けを管理するのに用いられる。

40

## 【 0 0 1 1 】

ストレージクラスタ内の全てのノードのコントローラは、それぞれ同一内容のL P M T を保持する。また各ノードは、L P M T により全論理ボリューム分を管理する。これらの理由は、ホストからコントローラへのパスが、別のコントローラに切り替わった際に、直ちに同一の論理ボリュームへのアクセスを継続可能とするためである。

【特許文献1】特開平9 - 6 2 4 5 2号公報

【発明の開示】

【発明が解決しようとする課題】

50

## 【 0 0 1 2 】

上記したように従来のストレージクラスタ（ストレージシステム）では、当該クラスタ内の各ノードは、L P M T（論理エクステント - 物理エクステントマッピングテーブル）により全論理ボリューム分を管理する必要がある。ところが、L P M Tのサイズは物理ボリュームの数と、その物理ボリュームを構成するエクステントの量によってはとても大きくなり、コントローラ上の大量のメモリを消費することになる。

## 【 0 0 1 3 】

ストレージクラスタでは、当該クラスタ内のあるノード（以下、第1のノードと称する）の例えば第1の物理エクステントに格納されていた第1の論理エクステントを、別のノード（以下、第2のノードと称する）の例えば第2の物理エクステントにマイグレーション（移動）することが発生する。このようなマイグレーションは、第1のノードのストレージの負荷を軽減するなどの目的で行われる。また、上記特許文献1に記載されているように、あるノードのストレージに対するアクセス性能を向上させるために、当該ノードのある物理エクステントに格納されていた第1の論理エクステントを、当該ノードの別の物理エクステントにマイグレーションすることもある。このマイグレーションにより、ストレージクラスタのアクセス性能及び信頼性を向上することが可能となる。

## 【 0 0 1 4 】

このようにストレージクラスタでは、ストレージクラスタのアクセス性能・信頼性向上のために、論理エクステントを同じストレージ内の別の物理エクステントへ、また別のストレージの物理エクステントへの移動/コピーが発生する。この場合、ストレージクラスタの各コントローラは互いに同期をとって、自身が管理するL P M Tを一斉（同時）に更新する必要がある。しかし、各コントローラのL P M Tが同一となるように同時に更新することは、L P M T管理の処理の複雑さを招き、アクセス性能に悪影響を及ぼす。

## 【 0 0 1 5 】

本発明は上記事情を考慮してなされたものでその目的は、各ノードが管理する情報の更新作業量を抑えてアクセス性能を向上させることができるストレージシステムを提供することにある。

## 【課題を解決するための手段】

## 【 0 0 1 6 】

本発明の1つの観点によれば、物理ボリュームを有する少なくとも1つのノードを含む複数のノードから構成され、前記少なくとも1つのノードが有する物理ボリュームを利用して、ホストに対して複数の論理エクステントから構成される論理ボリュームを提供するストレージシステムが提供される。このシステムにおいて、前記複数のノードの少なくとも2つは、前記ホストに対して前記論理ボリュームを提供し、当該ホストからのアクセスを受け付けるコントローラを有する。また、前記少なくとも1つのノードは、前記論理ボリュームを構成する論理エクステントを物理エクステントにマッピングして前記物理ボリューム内に格納するストレージを有し、前記複数のノードの少なくとも1つは、前記論理ボリュームを構成する論理エクステントを格納しているストレージの所在に関する前記コントローラからの問い合わせに対して回答する所在管理サービスを実行する所在管理サーバを有する。また、前記コントローラは、前記論理ボリュームを構成する論理エクステントに対応する所在管理サーバを示す第1のマッピングテーブルを有する。前記所在管理サーバは、当該サーバに対応する論理エクステントについて、その論理エクステントがどのノードのストレージに格納されているかを示す第2のマッピングテーブルを有し、前記ストレージは、論理エクステントを、当該ストレージが有する物理ボリュームのどの物理エクステントに対応させて格納したかを示す第3のマッピングテーブルを有する。前記コントローラは、前記ホストから前記論理ボリューム内の任意の論理エクステントをアクセス先として指定するアクセス要求が与えられた場合、当該コントローラが有する前記第1のマッピングテーブルを参照することにより当該アクセス先の論理エクステントに対応する所在管理サーバを特定し、特定された所在管理サーバに対して、当該論理エクステントがどのストレージに格納されているかを問い合わせる。前記所在管理サーバは前記コントロ

10

20

30

40

50

ーラから問い合わせを受けた場合、当該サーバが有する前記第2のマッピングテーブルを参照することにより問い合わせられたストレージを特定して、当該ストレージを示す情報を前記問い合わせに対する回答として問い合わせ元の前記コントローラに通知する。前記コントローラは、前記所在管理サーバから前記問い合わせに対する回答として通知された情報の示すストレージに対して前記ホストから要求された論理エクステントへのアクセスを要求する。前記ストレージは、前記コントローラから論理エクステントへのアクセスが要求された場合、当該ストレージが有する前記第3のマッピングテーブルを参照することにより、当該ストレージが有する前記物理ボリューム内で前記要求された論理エクステントに対応させられている物理エクステントを特定して、当該特定された物理エクステントに対してアクセスしてアクセス要求元の前記コントローラに応答する。

10

【発明の効果】

【0017】

本発明によれば、第1乃至第3のマッピングテーブルを用いることにより、ノード（コントローラ）が管理する情報の更新作業量を抑えてアクセス性能を向上させることができる。

【発明を実施するための最良の形態】

【0018】

以下、本発明の実施の形態につき図面を参照して説明する。

[第1の実施形態]

図1は本発明の第1の実施形態に係るストレージクラスタ構成のストレージシステムの概略構成を示すブロック図である。図1のストレージシステム（以下、ストレージクラスタと称する）は、複数、例えば3台のストレージ装置（以下、ノードと称する）1-1, 1-2及び1-3（#1, #2及び#3）から構成される。ノード1-1, 1-2及び1-3は、ファイバチャネル、或いはiSCSI（Small Computer System Interface）のようなインタフェース2によって、ストレージクラスタを利用するホスト（ホストコンピュータ）と接続されている。

20

【0019】

ノード1-1, 1-2及び1-3は、それぞれ、コントローラ11-1, 11-2及び11-3と、所在管理サーバ12-1, 12-2及び12-3と、ストレージ13-1, 13-2及び13-3とから構成される。

30

【0020】

コントローラ11-i（ $i = 1, 2, 3$ ）はホストに対して論理ボリュームを提供する。ここでは、ストレージ13-1, 13-2及び13-3の後述する物理ボリュームを統合することで、論理ボリュームLVol=1がホストに提供されるものとする。コントローラ11-iはホストからのアクセス要求を受け付ける。アクセス要求は、アクセスされるべき論理ボリューム内の論理ブロックを示すアドレス（論理ブロックアドレス）LBAを含む。コントローラ11-iは、要求された論理ブロックを含む論理エクステントがどのノードのストレージに格納されているか、後述するLMMT112aを参照して所在管理サーバ12-j（ $j = 1, 2, 3$ ）に問い合わせる。コントローラ11-iは、問い合わせに対して所在管理サーバ12-jから通知されたノードのストレージに対してアクセスを行う。

40

【0021】

図2は、コントローラ11-iの構成を示すブロック図である。コントローラ11-iは、ターゲットインタフェース111、サーバ特定部112及び入出力マネージャ（IOMネージャ）113を含む。

【0022】

ターゲットインタフェース111は、ホストからのアクセス要求を受け付け、要求された論理ボリュームの論理ブロックアドレスLBA（の示す論理ブロック）が位置する論理エクステントを特定する。ここでは、論理エクステントを識別する論理エクステントID（LEID）が、（ $LEID = LBA / \text{エクステントサイズ}$ ）の演算によって特定される。ターゲットインタフェース111は、特定された論理ボリューム内の論理エクステント

50

がどのストレージに格納されているかをサーバ特定部 1 1 2 を介して所在管理サーバ 1 2 -j に問い合わせる。この問い合わせは、論理ボリューム / 論理エクステントを識別する論理ボリューム ID (LVol) / 論理エクステント ID (LEID) を含む。ターゲットインタフェース 1 1 1 は、所在管理サーバ 1 2 -i に対する問い合わせに応じてサーバ特定部 1 1 2 を介して通知されたストレージに対して I/O マネージャ 1 1 3 を介してアクセスし、ホストに対して応答する。

#### 【 0 0 2 3 】

サーバ特定部 1 1 2 は、ターゲットインタフェース 1 1 1 からの問い合わせを受け渡す所在管理サーバ 1 2 -j を特定する。以下の説明では、ターゲットインタフェース 1 1 1 からの問い合わせに含まれている論理ボリューム ID 及び論理エクステント ID で示される論理エクステントをターゲット論理エクステントと呼ぶこともある。所在管理サーバ 1 2 -j は、このターゲット論理エクステントを格納するストレージの所在を管理するサーバとして特定される。

10

#### 【 0 0 2 4 】

サーバ特定部 1 1 2 は、ターゲット論理エクステントの所在を管理する所在管理サーバを特定するために、第 1 のマッピングテーブルとしての論理エクステント - 所在管理サーバマッピングテーブル (LMMT: LogicalExtent to ManagementServer Mapping Table) 1 1 2 a を有する。LMMT 1 1 2 a は、論理エクステントをどのストレージが格納しているか、その所在に関して責任を負う (つまり所在を管理する) 所在管理サーバを提供するノードの ID (MS) を保持する。サーバ特定部 1 1 2 は、コントローラ 1 1 -1 からの問い合わせを特定された所在管理サーバに受け渡す。サーバ特定部 1 1 2 は、ターゲットインタフェース 1 1 1 からの上記問い合わせに対する所在管理サーバ 1 2 -j からの応答 (ターゲット論理エクステントが格納されているストレージのノード ID (PVol)) を当該ターゲットインタフェース 1 1 1 に受け渡す。サーバ特定部 1 1 2 はまた、キャッシュ 1 1 2 b を有する。このキャッシュ 1 1 2 b については、後述する第 5 の実施形態で説明する。

20

#### 【 0 0 2 5 】

I/O マネージャ 1 1 3 は、コントローラ 1 1 -i とストレージ 1 3 -j との間のインタフェースを提供する。

#### 【 0 0 2 6 】

再び図 1 を参照すると、所在管理サーバ 1 2 -i (i = 1, 2, 3) は、コントローラ 1 1 -j (j = 1, 2, 3) によって問い合わせられた論理エクステントを格納しているストレージのノード ID を調べて、そのノード ID を当該コントローラ 1 1 -j に通知する。所在管理サーバ 1 2 -i は、自身が担当する論理エクステントについて、その論理エクステントがどのノード ID のストレージ (物理ボリューム) に格納されているかを管理するための、第 2 のマッピングテーブルとしての論理エクステント - ストレージマッピングテーブル (LSMT: Logical Extent to Storage Mapping Table) LSMT 1 2 0 -i を有する。つまり所在管理サーバ 1 2 -1, 1 2 -2 及び 1 2 -3 は、それぞれ LSMT 1 2 0 -1, 1 2 0 -2 及び 1 2 0 -3 を有する。

30

#### 【 0 0 2 7 】

ストレージ 1 3 -1, 1 3 -2 及び 1 3 -3 は、それぞれ物理ボリューム 1 3 0 -1, 1 3 0 -2 及び 1 3 0 -3 (PVol = 1, PVol = 2 及び PVol = 3) を有する。ストレージ 1 3 -i (i = 1, 2, 3) は論理エクステントを物理エクステントにマッピングして物理ボリューム 1 3 0 -i 内に格納する。ストレージ 1 3 -i は、コントローラ 1 1 -i からの入出力要求 (I/O 要求) を受け付ける。この I/O 要求は、アクセスされるべき論理ボリューム内の論理エクステントの ID (LEID) を含む。ストレージ 1 3 -i は、この論理エクステントと対応付けられている物理エクステントを調べる。ストレージ 1 3 -i は、この物理エクステントに対して I/O を行い、コントローラ 1 1 -i に対して応答する。

40

#### 【 0 0 2 8 】

図 3 は、ストレージ 1 3 -i の構成を示す。ストレージ 1 3 -i は、I/O ドライバ 1 3 1、

50

論理エクステント - 物理エクステント管理部 (LP 管理部) 132 及びディスク 133 を含む。

【0029】

I/O ドライバ 131 は、ストレージ 13-i とコントローラ 11-i とのインタフェースを提供する。

【0030】

LP 管理部 132 は、論理エクステントと物理エクステントとの対応を管理する。LP 管理部 132 は、論理エクステントが、自身の物理ボリューム 130-i のどの物理エクステントに対応させて格納されているかを示す、第 3 のマッピングテーブルとしての論理エクステント - 物理エクステントマッピングテーブル (LPMT: LogicalExtent to PhysicalExtent Mapping Table) 132a を保持する。このように、本実施形態で適用される LPMT 132a は、[背景技術] で述べた LPMT と異なり、全論理ボリューム分のマッピング情報を持たないことに注意されたい。また、LPMT 132a の内容は、各ストレージ 13-i 毎に異なることにも注意されたい。

10

【0031】

ディスク 133 は、物理ボリューム 130-i (PVol = i) を有する。物理ボリューム 130-i は物理エクステント単位で管理される。また、物理ボリューム 130-i へのアクセスの単位はブロックであり、物理エクステントは、複数のブロックから構成される。

【0032】

また、図 1 のストレージクラスタは、当該クラスタの構成の計画・管理を司る構成管理部 (図示せぬ) を有する。この構成管理部は、ボリューム作成時に、論理ボリューム / 論理エクステントの格納先となる物理ボリューム / 物理エクステントを決定する。構成管理部はまた、エクステントや後述する所在管理サービスの移動 / コピー (冗長化) のための移動元 / 移動先を決定する。構成管理部は更に、オンデマンドでエクステント割り付けの割り付け先を決定する。この構成管理部が、ノード 1-1 ~ 1-3 から独立に設けられても、当該ノード 1-1 ~ 1-3 のいずれかに設けられても構わない。

20

【0033】

本実施形態において、構成管理部は、以下の手順で論理ボリュームを作成する。

(a) 構成管理部は、論理ボリュームを構成する論理エクステントの各々を、どのノード ID のストレージに格納するか決定する。構成管理部は、論理エクステントの各々に、決定されたストレージの物理エクステントを割り付ける。構成管理部は、ストレージ 13-1 ~ 13-3 毎に、そのストレージの物理エクステントと当該物理エクステントが割り付けられた論理エクステントとの対応を示す LPMT 132a を生成して、そのストレージに保持する。

30

【0034】

(b) 構成管理部は、論理エクステントが割り付けられた物理エクステントを格納するストレージ (物理ボリューム) から、その論理エクステントの所在を管理するサービス (所在管理サービス) をノード 1-1 ~ 1-3 のうちのいずれのノードの所在管理サーバに担当させるかを決定する。構成管理部は、所在管理サーバ 12-1 ~ 12-3 毎に、その所在管理サーバが担当する論理エクステントと物理ボリュームとの対応を示す LSMT 120-1 ~ 120-3 を生成して、その所在管理サーバ 12-1 ~ 12-3 内に設定する。なお、後述する第 2 の実施形態のように、ストレージクラスタ内の全ての所在管理サーバ 12-1 ~ 120-3 に LSMT が設定される必要はない。

40

【0035】

(c) 構成管理部は、LSMT 112a を生成して全コントローラ 11-1 ~ 11-3 に配布する。LSMT 112a は先に述べたように、論理ボリュームを構成する論理エクステント毎に、その論理エクステントが割り付けられている物理エクステントを格納するストレージ (物理ボリューム) をどの所在管理サーバに問い合わせればよいかを示す。

【0036】

このようにして、論理ボリュームが作成される。すると、作成された論理ボリュームに

50



対してホストからアクセスすることが可能となる。ホストからの論理ボリュームに対するアクセスは、ホストからノード1-1~1-3のうちの任意のノード1-iのコントローラ11-iに対してパスを張ることで行われる。

【0037】

以下、例えばホストからノード1-1( $i = 1$ の場合)に対してパスが張られた後に、当該ホストからアクセス要求が発行された場合の処理の流れについて、図4及び図5を参照して説明する。図4は処理の手順を示すフローチャート、図5は情報の流れを示す図である。図5における矢印A1~A8は、図4のフローチャートにおけるステップS1~S8にそれぞれ対応する。

【0038】

まず、ホストからノード1-1のコントローラ11-1に含まれているターゲットインタフェース111に対してアクセス要求(I/O要求)が発行されたものとする。また、このアクセス要求が、ある論理ボリューム $LVol = LVolx$ (例えば $LVol = 1$ )のある論理ブロックアドレス $LBA = LBAx$ に対するアクセスを要求しているものとする。

【0039】

ターゲットインタフェース111は、ホストからのアクセス要求を受け付ける(ステップS0)。するとターゲットインタフェース111は、ホストからのアクセス要求で指定されるアクセスされるべき論理ボリューム $LVol = LVolx$ 上の $LBA = LBAx$ から、当該 $LBAx$ が位置する論理エクステントの $ID = LEIDx$ を次式

$$LEIDx = LBAx / \text{エクステントサイズ}$$

により求める(ステップS1)。ここで、論理ボリューム $LVol = LVolx = 1$ 上の論理エクステントIDが $LEIDx$ の論理エクステントを、論理エクステント( $LVolx, LEIDx$ )と表現する。

【0040】

ターゲットインタフェース111は、論理ボリューム名( $ID : LVol = LVolx = 1$ )及び求められた論理エクステント $ID = LEIDx$ をサーバ特定部112に渡すことにより、当該論理エクステント $ID = LEIDx$ の論理エクステント( $LVolx, LEIDx$ )が、どのノードのストレージに格納されているかを問い合わせる(ステップS2)。

【0041】

サーバ特定部112は、ターゲットインタフェース111によって渡された論理エクステント $ID = LEIDx$ の論理エクステント( $LVolx, LEIDx$ )がどのノードのストレージに格納されているかを管理している所在管理サーバを提供するノードのID( $MS = MIDx$ )を、 $LMMT112a$ に基づいて特定する(ステップS3)。ここでは、サーバ特定部112と同じノード1-1、つまり所在管理サーバ12-1を提供するノード1-1のID( $MS = MIDx = 1$ )が特定されたものとする。本実施形態において、ノード1-1~1-3のストレージ13-1~13-3は物理ボリューム130-1~130-3を有する。つまり、ノード1-1~1-3のストレージ13-1~13-3は、それぞれ1つの物理ボリュームを有する。このため本実施形態では、所在管理サーバ12-1~12-3を提供するノード1-1~1-3のID(MS)と当該ノード1-1~1-3のストレージ13-1~13-3が有する物理ボリューム130-1~130-3のID(PVol)とは一致し、1~3である。

【0042】

サーバ特定部112は、特定されたノード1-1の所在管理サーバ12-1に対して、ターゲットインタフェース111からの問い合わせを受け渡す(ステップS4)。即ち、サーバ特定部112は所在管理サーバ12-1に対して、論理エクステント $ID = LEIDx$ の論理エクステント(ターゲット論理エクステント)( $LVolx, LEIDx$ )に対応する物理エクステントを格納しているストレージのノードID( $PVol) = SIDx$ を問い合わせる。

【0043】

所在管理サーバ12-1は、どのノード $ID = SIDx$ のストレージが、論理エクステン

10

20

30

40

50

ト (LVolx, LEIDx) に対応する物理エクステントを格納しているかを、LSMT120-1に基づいて調べて、そのノードID = SIDxを、問い合わせ元のサーバ特定部112に対して通知する(ステップS5)。

【0044】

サーバ特定部112は、所在管理サーバ12-1から通知されたノードID = SIDxを、ターゲットインタフェース111からの問い合わせに対する回答として当該ターゲットインタフェース111に渡す(ステップS6)。ここでは、ストレージ13-1のノードID = SIDxがターゲットインタフェース111に渡されたものとする。

【0045】

ターゲットインタフェース111は、サーバ特定部112によって渡されたノードID = SIDxで示されるノード1-1のストレージ13-1に対し、論理エクステント(LVolx, LEIDx)に対応する物理エクステントへのアクセスをIOMネージャ113を介して要求する(ステップS7)。

10

【0046】

ストレージ13-1のIODライバ131は、要求された論理エクステント(LVolx, LEIDx)に、当該ストレージ13-1(物理ボリューム130-1)内のどのID (= PEIDx)の物理エクステント(以下、物理エクステントPEIDxと表現する)が対応付けられているかを、LPMT132aに基づいて調べる(ステップS8)。このステップS8において、IODライバ131は、論理エクステント(LVolx, LEIDx)に対応する物理エクステントPEIDxに対してアクセスし、コントローラ11-1のIOMネージャ113に対して応答する。ターゲットインタフェース111は、ストレージ13-1のIODライバ131からの応答をIOMネージャ113から受け取って、ホストに返す(ステップS9)。

20

【0047】

以上の説明では、ホストからのアクセス要求を受け付けたターゲットインタフェースと、ターゲット論理エクステントがどのストレージに格納されているかを管理している所在管理サーバと、当該ターゲット論理エクステントを物理エクステント内に格納しているストレージとが、いずれも同一のノード1-1に属する場合を想定している。

【0048】

図6は、上述のターゲットインタフェース、所在管理サーバ及びストレージが、それぞれ別々のノードに属する場合の情報の流れを示す。図6における矢印B1~B8は、図4のフローチャートにおけるステップS1~S8にそれぞれ対応する。

30

【0049】

図7は、図1のシステムにおいて、論理ボリュームLVolがLVol = 1であって、当該論理ボリュームLVol = 1を、LEIDがLEID = 1~LEID = 4の4つの論理エクステントから構成した場合のシステム構成を示す。また、図8は図7の状態におけるLMMT112aのデータ構造例を、図9は図7の状態におけるLSMT120-i、例えばノード1-1のLSMT120-1のデータ構造例を、図10は図7の状態における例えばノード1-1のLPMT132aのデータ構造例を示す。

【0050】

40

LMMT112aは、論理ボリュームLVol = 1を構成する論理エクステントの数に一致するエントリを有する。LMMT112aの各エントリは、論理ボリュームのID(論理ボリューム番号LVol)と、当該LVolで示される論理ボリュームに含まれる論理エクステントのID(LEID)と、当該LVol及びLEIDで示される論理エクステント(LVol, LEID)をどのストレージ(物理ボリューム)が格納しているか、その所在を管理する所在管理サーバ12-jを提供するノードのID(MS)とを保持する。図8の例では、論理ボリュームLVol = 1を構成する論理エクステントが、LEID = 1~LEID = 4の4つであることが示されている。

【0051】

LSMT120-i(120-1)は、所在管理サーバ12-i(12-1)が管理する論理工

50

クステントの数に一致するエントリを有する。L S M T 1 2 0 - i ( 1 2 0 - 1 ) の各エントリは、論理ボリュームの I D ( L V o l ) と、当該 L V o l で示される論理ボリュームに含まれる論理エクステントの I D ( L E I D ) と、当該 L V o l 及び L E I D で示される論理エクステントが格納されているストレージ ( 物理ボリューム ) のノード I D ( 物理ボリューム番号 P V o l ) とを保持する。

【 0 0 5 2 】

L P M T 1 3 2 a は、物理ボリューム 1 3 0 - i ( 1 3 0 - 1 ) を構成する物理エクステントの数に一致するエントリを有する。L P M T 1 3 2 a の各エントリは、物理ボリューム 1 3 0 - i を構成する物理エクステントの I D ( P E I D ) と、当該 P E I D で示される物理エクステントが割り付けられる論理エクステントを含む論理ボリュームの I D ( L V o l ) と、当該論理エクステントの I D ( L E I D ) と、当該論理エクステントをどのストレージ ( 物理ボリューム ) が格納しているか、その所在を管理する所在管理サービス ( 所在管理サーバ 1 2 - j ) を提供するノードの I D ( M S ) とを保持する。この L P M T 1 3 2 a 内のエントリの M S は、所在管理サービスの逆引きに用いられる。図 1 0 の例では、物理ボリューム 1 3 0 - i ( 1 3 0 - 1 ) を構成する物理エクステントが、P E I D = 1 ~ P E I D = 3 の 3 つであることが示されている。

10

【 0 0 5 3 】

図 7 において、コントローラ 1 1 - 1 の L M M T 1 1 2 a のエントリから所在管理サーバ 1 2 - 1 ~ 1 2 - 3 の L S M T 1 2 0 - 1 ~ 1 2 0 - 3 へ向かう一点鎖線の矢印は、L M M T 1 1 2 a のエントリの情報によって参照される L S M T 1 2 0 - 1 ~ 1 2 0 - 3 のエントリを指す。同様に、L S M T 1 2 0 - 1 ~ 1 2 0 - 3 のエントリからストレージ 1 3 - 1 ~ 1 3 - 2 の L P M T 1 3 2 a へ向かう一点鎖線の矢印は、L S M T 1 2 0 - 1 ~ 1 2 0 - 3 のエントリの情報によって参照される L P M T 1 3 2 a のエントリを指す。同様に、ストレージ 1 3 - 1 ~ 1 3 - 2 の L P M T 1 3 2 a のエントリからストレージ 1 3 - 1 ~ 1 3 - 2 の物理ボリューム 1 3 0 - 1 ~ 1 3 0 - 3 へ向かう一点鎖線の矢印は、L P M T 1 3 2 a の情報に基づいてアクセスされる物理エクステントを指す。

20

【 0 0 5 4 】

次に、図 7 を参照して、ホストからノード 1 - 1 に対してパスが張られた後に、当該ホストからの、論理ボリューム L V o l = 1 / 論理エクステント L E I D = 1 で指定されるブロックに対するアクセスの流について説明する。

30

【 0 0 5 5 】

まず、図 7 において矢印 B 1 で示されるように、ホストからノード 1 - 1 に対して、論理ボリューム L V o l = 1 内のあるブロックへのアクセス要求が発行されたものとする。ノード 1 - 1 のコントローラ 1 1 - 1 ( に含まれているターゲットインタフェイス 1 1 1 ) は、アクセスされるべきブロックが L E I D = 1 の論理エクステント ( 論理エクステント L E I D = 1 ) に格納されていることを求める。

【 0 0 5 6 】

コントローラ 1 1 - 1 のターゲットインタフェイス 1 1 1 は、論理エクステント L E I D = 1 がどのノードのストレージに格納されているかを、当該コントローラ 1 1 - 1 内のサーバ特定部 1 1 2 に問い合わせる。サーバ特定部 1 1 2 は、自身の L M M T 1 1 2 a 内の L E I D = 1 が保持されているエントリ ( 図 8 の例では、1 行目のエントリ ) を参照する。これによりサーバ特定部 1 1 2 は、論理ボリューム L V o l = 1 / 論理エクステント L E I D = 1 に対応する所在管理サーバがノード 1 - 1 に存在することを認識する。

40

【 0 0 5 7 】

コントローラ 1 1 - 1 のサーバ特定部 1 1 2 は、図 7 において矢印 B 2 で示されるように、ノード 1 - 1 に存在する所在管理サーバ 1 2 - 1 に対して、論理ボリューム L V o l = 1 / 論理エクステント L E I D = 1 がどのノードのストレージ ( 物理ボリューム ) に格納されているかを問い合わせる。

【 0 0 5 8 】

所在管理サーバ 1 2 - 1 は、自身の L S M T 1 2 0 - 1 内で L V o l = 1 / L E I D = 1 が

50

保持されているエントリ（図9の例では、1行目のエントリ）を参照して、論理ボリュームLVol = 1 / 論理エクステントLEID = 1の実体がノード1-1のストレージ（物理ボリューム）に格納されていることを認識する。所在管理サーバ12-1は、ノード1-1のID（SIDx = PVol = 1）を、図7において矢印B3で示されるように、コントローラ11-1のサーバ特定部112に回答する（ステップB3）。これを受けてサーバ特定部112は、ノード1-1のID（SIDx = PVol = 1）をターゲットインタフェース111に通知する。

【0059】

ターゲットインタフェース111（コントローラ11-1のターゲットインタフェース111）は、図7において矢印B4で示されるように、サーバ特定部112によって通知されたノードIDで指定されるストレージ、即ちノード1-1のストレージ13-1に対して、論理ボリュームLVol = 1 / 論理エクステントLEID = 1へのアクセスをIOMネージャ113を介して要求する。

10

【0060】

ノード1-1のストレージ13-1は、自身のLPM T132a内の論理ボリュームLVol = 1 / 論理エクステントLEID = 1が保持されているエントリ（図10の例では、1行目のエントリ）を参照する。これによりストレージ13-1は、論理ボリュームLVol = 1 / 論理エクステントLEID = 1に対応する物理エクステントのPEIDがPEID = 1であることを認識する。ストレージ13-1は、このPEID = 1の物理エクステント（物理エクステントPEID = 1）にアクセスして、図7において矢印B5で示されるように、コントローラ11-1に対して応答する。

20

【0061】

本実施形態において、各ノード1-1～1-3のコントローラ11-1～11-3が有するLMM T112aは同一の内容である。つまり各ノード1-1～1-3（コントローラ11-1～11-3）はLMM T112aを共有する。このためホストは、ストレージクラスタを構成するノード1-1～11-3のいずれのノード（のコントローラ）に対してもパスを張ることができる。例えば、ホストがパスを張っている先のコントローラが異常停止した場合、或いは当該コントローラの負荷が高まったときなどは、ホストは別のコントローラに対してパスを張り替えることができる。

【0062】

LMM T112aなど、全ノード1-1～11-3が共有する情報（共有情報）は、当該ノード1-1～11-3のコントローラ11-1～11-3が有する不揮発性メモリの記憶領域、当該ノード1-1～11-3のストレージ13-1～13-3の記憶領域、或いはストレージ13-1～13-3のうちの任意のストレージの記憶領域に格納される。コントローラ11-1～11-3は、この記憶領域を参照することにより共有情報を管理する。この管理のために、コントローラ11-1～11-3は上記記憶領域のアドレスを共有する。共有情報はLMM T112aの他に、論理ボリューム構成情報を含む。論理ボリューム構成情報は、図1のストレージクラスタがホストに対してどのような論理ボリュームを提供するかを示す。上述した各ノード1-1～11-3が情報を共有する構成により、従来技術と比較して、当該各ノード1-1～11-3が管理しなければならない情報量を減らすことができる。

30

40

【0063】

次に、論理エクステントを構成する物理エクステントのデータを、ストレージをまたがって別の物理エクステントにマイグレート（移動）する場合の動作について、図11を参照して説明する。

【0064】

まず、例えば構成管理部によってデータの移動元及び移動先が決定される。ここでは、図1のストレージクラスタが図10に示す状態にあるときに、データの移動元として、論理エクステントLEID = 4に対応付けられている、ノード1-1のストレージ13-1内の物理ボリュームPVol = 1（130-1）/ 物理エクステントPEID = 2が、移動先として、ノード1-2のストレージ13-2内の物理ボリュームPVol = 2（130-2）/ 物

50

理エクステント P E I D = 1 が、それぞれ決定されたものとする。

【 0 0 6 5 】

この場合、ノード 1 -1 のストレージ 1 3 -1 内の物理ボリューム P V o l = 1 / 物理エクステント P E I D = 2 に格納されているデータが、図 1 1 おいて矢印 C 1 で示されるように、ノード 1 -2 のストレージ 1 3 -2 内の物理ボリューム P V o l = 2 / 物理エクステント P E I D = 1 に移動される。

【 0 0 6 6 】

また、移動元のストレージ 1 3 -1 及び移動先のストレージ 1 3 -2 の各々の L P M T 1 3 2 a が書き換えられる。ここでは、ストレージ 1 3 -1 の L P M T 1 3 2 a 内の P E I D = 2 が保持されているエントリの L V o l , L E I D , M S がそれぞれ「 1 」 「 N U L L 」 , 「 4 」 「 N U L L 」 , 「 1 」 「 N U L L 」 に書き換えられ、ストレージ 1 3 -2 の L P M T 1 3 2 a 内の P E I D = 1 が保持されているエントリの L V o l , L E I D , M S が、図 1 1 おいて矢印 C 2 で示されるように、それぞれ「 N U L L 」 「 1 」 , 「 N U L L 」 「 4 」 , 「 N U L L 」 「 1 」 に書き換えられる。

10

【 0 0 6 7 】

また、論理エクステント L E I D = 4 を参照していた（つまり論理エクステント L E I D = 4 に関する所在管理サービスを担当している）所在管理サーバ 1 2 -1 の L S M T 1 2 0 -1 が書き換えられる。ここでは、L S M T 1 2 0 -1 内の L E I D = 4 が保持されているエントリの P V o l が、図 1 1 おいて矢印 C 3 で示されるように、「 1 」 「 2 」 に書き換えられる。

20

【 0 0 6 8 】

このように本実施形態において、論理エクステント（ L E I D = 4 ）のマイグレーションに係わる通知（更新）は、その論理エクステント（ L E I D = 4 ）を管理する所在管理サーバ（ 1 2 -1 ）に対して行われるだけで良い。つまり、上記の通知はストレージクラスタを構成する全てのノード 1 -1 ~ 1 -3 に対して行われる必要がない。このため、通知コストを小さくすることができる。

【 0 0 6 9 】

次に、論理エクステントを構成する物理エクステントのデータを、ストレージをまたがって別の物理エクステントにコピーする場合の動作について、図 1 2 を参照して説明する。このようなコピー処理は、論理エクステントを構成する物理エクステントを冗長化し、耐障害性を高め、また並列処理によりリードアクセスを高速化するために行われる。

30

【 0 0 7 0 】

まず、例えば構成管理部によってコピー元及びコピー先が決定される。ここでは、図 1 のストレージクラスタが図 1 1 に示す状態にあるときに、データのコピー元として、論理エクステント L E I D = 4 に対応付けられている、ノード 1 -2 のストレージ 1 3 -2 内の物理ボリューム P V o l = 2 ( 1 3 0 -2 ) / 物理エクステント P E I D = 1 が、コピー先として、ノード 1 -3 のストレージ 1 3 -3 内の物理ボリューム P V o l = 1 ( 1 3 0 -3 ) / 物理エクステント P E I D = 2 が、それぞれ決定されたものとする。

【 0 0 7 1 】

この場合、ノード 1 -2 のストレージ 1 3 -2 内の物理ボリューム P V o l = 2 / 物理エクステント P E I D = 1 に格納されている論理ボリューム L V o l = 1 / 論理エクステント L E I D = 4 のデータが、図 1 2 において矢印 D 1 で示されるように、ノード 1 -3 のストレージ 1 3 -3 内の物理ボリューム P V o l = 1 / 物理エクステント P E I D = 2 にコピーされる。このコピーは、例えばストレージ 1 3 -2 及び 1 3 -3 の I O ドライバ 1 3 1 によって行われる。

40

【 0 0 7 2 】

また、コピー先のストレージ 1 3 -3 の L P M T 1 3 2 a が書き換えられる。ここでは、ストレージ 1 3 -2 の L P M T 1 3 2 a 内の P E I D = 1 が保持されているエントリの L V o l = 1 , L E I D = 4 , M S = 1 が、図 1 2 において矢印 D 2 で示されるように、ストレージ 1 3 -3 の L P M T 1 3 2 a 内の P E I D = 2 が保持されているエントリの L V o l

50

、LEID、MSにコピーされる。

【0073】

また、論理エクステントLEID = 4を参照していた所在管理サーバ12-1のLSMT120-1が、例えばストレージ13-1のIODライバ131によって書き換えられる。ここでは、ノード1-1の所在管理サーバ12-1が管理するLSMT120-1内の、LEID = 4が保持されているエントリのPVOLに、図12において矢印D3で示されるように、「3」が追加される。このLSMT120-1内のエントリの書き換えにより、論理エクステントLEID = 4が、ノード1-2の物理ボリュームPVOL = 2に加えて、新たにノード1-3の物理ボリュームPVOL = 3にも対応付けられる。つまり、論理エクステントLEID = 4を構成する物理エクステントが、ノード1-2の物理ボリュームPVOL = 2及びノード1-3の物理ボリュームPVOL = 3に冗長化される。

10

【0074】

論理エクステントを構成する物理エクステントの冗長化構成後に、ホストからコントローラ11-iに対して、当該論理エクステントへのライトアクセスの要求が発生したものとす。この場合、コントローラ11-iは、上記論理エクステントに対応する複数の物理エクステントの全てに対してライトを行う。したがって、論理エクステントLEID = 4へのライトアクセスが発生した場合であれば、ノード1-2の物理ボリュームPVOL = 2 / 物理エクステントPEID = 1及びノード1-3の物理ボリュームPVOL = 3 / 物理エクステントPEID = 2に対してデータライトが行われる。リードの場合はこの限りではない。つまり、物理ボリュームPVOL = 2 / 物理エクステントPEID = 1または物理ボリュームPVOL = 3 / 物理エクステントPEID = 2の一方からのデータリードが行われる。

20

【0075】

次に、論理エクステントに関する所在管理サービスを、ノードをまたがってマイグレート（移動）する場合の動作について、図13を参照して説明する。この所在管理サービスのマイグレーションは、当該所在管理サービスの負荷分散、或いはノードを構成するハードウェアの交換 / 更新などのために行われる。

【0076】

論理エクステントに関する所在管理サービスのマイグレーションに際し、例えば構成管理部により移動元及び移動先が決定される。ここでは、図1のストレージクラスタが図12に示す状態にあるときに、所在管理サービスの移動元として、論理ボリュームLVOL = 1 / 論理エクステントLEID = 4に関するネームサービスを担当する、ノード1-1の所在管理サーバ12-1が決定されたものとする。また、そのネームサービスの移動先として、ノード1-2の所在管理サーバ12-2が決定されたものとする。

30

【0077】

すると、移動元ノード1-1の論理ボリュームLVOL = 1 / 論理エクステントLEID = 4に関するLSMT120-1内エントリの情報が、図13において矢印E1で示されるように、移動先ノード1-2のLSMT120-2にコピーされる。

【0078】

移動元ノード1-1の例えば所在管理サーバ12-1または移動先ノード1-2の例えば所在管理サーバ12-2は、論理ボリュームLVOL = 1 / 論理エクステントLEID = 4に関する所在管理サービスがノード1-2に移動されたことを、全てのコントローラ11-1 ~ 11-3に通知する。

40

【0079】

コントローラ11-1 ~ 11-3の各々は、論理ボリュームLVOL = 1 / 論理エクステントLEID = 4に関するLMMT112a内のエントリのMSを、図13において矢印E2で示されるように、ノード1-1（の所在管理サーバ12-1）のID（MS = 1）からノード1-2（の所在管理サーバ12-2）のID（MS = 2）に書き換える（張り替える）。

【0080】

また、論理ボリュームLVOL = 1 / 論理エクステントLEID = 4を管理するノード

50

1-2及び1-3のストレージ13-2及び13-3は、自身のLPM T132a中の論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関する所在管理サービスの逆引きを、図13において矢印E3で示されるように、ノード1-1(の所在管理サーバ12-1)のID(MS = 1)からノード1-2(の所在管理サーバ12-2)のID(MS = 2)に変更する。ストレージ13-2及び13-3は、この変更を所在管理サービス移動の通知元ノードに  
10

【0081】

ここで、全コントローラ11-1~11-3がLMM T112aを書き換えるには、A)所在管理サーバ(ここでは所在管理サーバ12-1または12-2)によるブロードキャスト(一斉通知)による方法、B)コントローラのサーバ特定部(ここではコントローラ11-1  
10

【0082】

ノード1-1の例えば所在管理サーバ12-1は、全ノード1-1~1-3からの応答を確認して、当該所在管理サーバ12-1の論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関するLSMT120-1内のエントリを削除する。これにより、移動元ノード1-1の論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関するLSMT120-1内エントリのデータが、移動先ノード1-2のLSMT120-2に移動されたことになる。つまり、論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関する移動  
20

【0083】

次に、論理エクステントに関する所在管理サービスを、ノードをまたがってコピーする場合の動作について、図14を参照して説明する。

【0084】

所在管理サービスの多重化がなされていない状態では、当該所在管理サービスを担当するあるノードが故障した場合、ホストからのアクセス要求が実行できなくなる。そこで、所在管理サービスをノードをまたがってコピーして当該所在管理サービスを多重化(冗長化)することにより、耐障害性向上及び負荷分散を図ることが必要となる。  
30

【0085】

所在管理サービスをノードをまたがってコピー(冗長化)するに際し、例えば構成管理部によりコピー元及びコピー先が決定される。ここでは、図1のストレージクラスタが図13に示す状態にあるときに、所在管理サービスのコピー元として、論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関するネームサービスを担当する、ノード1-2の所在管理サーバ12-2が決定されたものとする。また、そのネームサービスのコピー先として、ノード1-3の所在管理サーバ12-3が決定されたものとする。

【0086】

すると、コピー元ノード1-2の論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関するLSMT120-2内エントリの情報が、図14において矢印F1で示されるように、コピー先ノード1-3のLSMT120-3にコピーされる。  
40

【0087】

コピー元ノード1-2の例えば所在管理サーバ12-2またはコピー先ノード1-3の例えば所在管理サーバ12-3は、論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関する所在管理サービスがノード1-2からノード1-3にも機能分担されたことを、全てのコントローラ11-1~11-3に通知する。

【0088】

コントローラ11-1~11-3の各々は、論理ボリュームLVol = 1 / 論理エクステントLEID = 4に関するLMM T112a内のエントリのMS、つまりノード1-2(の所在管理サーバ12-2)のID(MS = 2)が設定されているMSに、図14において矢印F  
50

2で示されるように、ノード1-3(の所在管理サーバ12-3)のID(MS=3)を加える。これにより、論理ボリュームLVol=1/論理エクステントLEID=4に関するLMMT112a内のエントリが、ノード1-2だけでなくノード1-3にも向けられる。

【0089】

また、論理ボリュームLVol=1/論理エクステントLEID=4を新たに管理するノード1-3のストレージ13-3は、自身のLPMT132a中の論理ボリュームLVol=1/論理エクステントLEID=4に関する所在管理サービスの逆引き(MS)を、図14において矢印F3で示されるように、ノード1-2(の所在管理サーバ12-2)だけでなく、ノード1-3(の所在管理サーバ12-3)にも向ける。ストレージ13-3は、この変更を所在管理サービスコピーの通知元ノードに回答する。全コントローラ11-1~11-3がLMMT112aを書き換えるには、前述した所在管理サービスの移動の場合と同様の方法A)、B)またはC)が適用可能である。

10

【0090】

以上、複数のノード(ストレージ)にまたがるエクステント移動/コピー及び所在管理サービス移動/コピーについて説明した。しかし、エクステント移動/コピーは1つのノード(ストレージ)の物理ボリューム内でも行うことができる。この場合、そのノード(ストレージ)のLPMT132aを更新するだけで良い。また、例えば所在管理サーバ12-iのLSMT120-i及びストレージ13-iのLPMT132aを冗長化のためにコピーすることも可能である。

【0091】

20

次に、図1のストレージクラスタに、新たにノードを参加させる(追加する)場合の動作について説明する。この場合、既にストレージクラスタを構成しているノード1-1~1-1のいずれかから、当該ストレージクラスタに新たに追加されるノード(新たなノード)に対して、少なくともLMMT112aをコピーする。これにより、新たなノードはIOを実行することができる。

【0092】

新たなノードにノード1-1~1-1のいずれかの所在管理サーバの処理(所在管理サービス)を分担させるには、例えば構成管理部と連携して、前述の所在管理サービスをマイグレートする処理を行えば良い。ノードの参加で新たに増えたストレージ内の物理ボリューム/物理エクステントも、構成管理部と連携して、適宜分散して利用することができる。

30

【0093】

上述したように本実施形態(第1の実施形態)においては、複数のノード(実施形態では3つのノード)を用いてストレージクラスタを構成したとき、各々のノードが管理しなければならない情報量を減らすことができる。また本実施形態においては、あるノード(のストレージ)の物理ボリューム内で、或いは複数のノード(のストレージ)の物理ボリュームをまたがって、論理エクステントに対応する物理エクステントを移動/コピーする場合に、ストレージクラスタに含まれる全てのコントローラが情報を書き換える必要はないため(所在管理サーバの情報を書き換えるだけで良いため)、コントローラ間での同期更新処理が必要なく、エクステント移動/コピー時の処理を簡単にし、総じて性能を向上させることができる。

40

【0094】

[第2の実施形態]

前記第1の実施形態では、ストレージクラスタを構成する各ノード1-1~1-3は、それぞれコントローラ11-1~1-3、所在管理サーバ12-1~12-3及びストレージ13-1~13-3を備えている。しかし、各ノードが、コントローラ、所在管理サーバ及びストレージの全てを必ずしも備える必要はない。そこで、ストレージクラスタを構成するノードの幾つかは、コントローラ、所在管理サーバ及びストレージの1つまたは2つのみを備える本発明の第2の実施形態について説明する。

【0095】

図15は本発明の第2の実施形態に係るストレージクラスタ構成のストレージシステム

50



の構成を示すブロック図である。図15において、図7と同様の要素には便宜的に同一参照符号を付してある。図15のストレージシステム(ストレージクラスタ)は、前記第1の実施形態と同様に、3つのノード1-1~1-3から構成される。図15のストレージクラスタが第1の実施形態と異なるのは、ノード1-1~1-3のうちノード1-2及び1-3がコントローラ、所在管理サーバ及びストレージの1つまたは2つのみを備える点にある。ここでは、ノード1-1は、第1の実施形態と同様に、コントローラ11-1、所在管理サーバ12-1及びストレージ13-1を備える。これに対し、ノード1-2はコントローラ11-2のみを備え、ノード1-3はコントローラを除く所在管理サーバ12-3及びストレージ13-3を備えている。この図15に示すストレージクラスタの動作自体は、第1の実施形態と同様である。

10

## 【0096】

本実施形態では、ストレージクラスタが3つのノードから構成されているものの、コントローラ、所在管理サーバ及びストレージの数が2つである場合を想定している。しかし、所在管理サーバ及びストレージの数はストレージクラスタ全体で1つであっても構わない。つまり複数のノードから構成されるストレージクラスタにおいて、コントローラの数は少なくとも2つであれば良く、所在管理サーバ及びストレージの数は少なくとも1つであれば良い。この制約のもとで、コントローラ、所在管理サーバ及びストレージの数は適宜増減可能である。つまり本実施形態においては、ストレージクラスタを構成するとき、状況に応じて、必要な機能のみを拡張することができる。

## 【0097】

20

## [第3の実施形態]

次に、本発明の第3の実施形態について説明する。この第3の実施形態の特徴は、実際にホストからある論理エクステントへの初回のアクセスが要求されるまで、ストレージ側で当該論理エクステントに対応する物理エクステントの割り付けを行わず、ディスクの有効利用を図ることにある。

## 【0098】

図16は、本発明の第3の実施形態に係るストレージクラスタ構成のストレージシステムの構成を示すブロック図である。図16において、図7と同様の要素には便宜的に同一参照符号を付してある。図16のストレージシステム(ストレージクラスタ)は、前記第1の実施形態と同様に、3つのノード1-1~1-3から構成される。図16のストレージクラスタが第1の実施形態と異なるのは、ホストから論理エクステントへの初回のアクセス、例えばライトアクセスが要求されるまで、ストレージ13-i側で当該論理エクステントに対応する物理エクステントの割り付けを行わない点にある。

30

## 【0099】

図16の例では、論理ボリュームLV01=1/論理エクステントLEID=1及び論理ボリュームLV01=1/論理エクステントLEID=4に対して、ノード1-1における物理ボリューム130-1の物理エクステントの割り付けが行われていない状態が、ストレージ13-1のLPMT132aによって示されている。同様に、論理ボリュームLV01=1/論理エクステントLEID=2に対して、ノード1-2内の物理ボリューム130-2の物理エクステントの割り付けが行われていない状態が、ストレージ13-2のLPMT132aによって示されている。また、論理ボリュームLV01=1/論理エクステントLEID=3に対して、ノード1-3内の物理ボリューム130-3の物理エクステントの割り付けが行われていない状態が、ストレージ13-3のLPMT132aによって示されている。

40

## 【0100】

ここでは、未割り付けの状態は、LPMT132a内で列項目PEIDに“NA(Not Assignment)”が設定されているエントリにより示される。但し、LPMT132a内で列項目PEIDに“NA”が設定されているエントリの他の列項目LV01、LEID及びMSには、それぞれ論理ボリューム、論理エクステント及びノードを指定する情報が設定されている。このことは、論理ボリューム/論理エクステントに対するストレージ(物

50

理ボリューム)の割り付けは行われていることを示す。換言するならば、論理ボリューム / 論理エクステントに対する L P M T 1 3 2 a 内のエン트리 ( L V o l , L E I D ) で指定される論理ボリューム / 論理エクステントに必要な領域が、当該エントリの列項目 M S で指定されるノード上で予約されていることを示す。なお、図 1 6 では、L P M T 1 3 2 a 内のエン트리における各列項目の並びが便宜的に変更されている点に注意されたい。

【 0 1 0 1 】

L P M T 1 3 2 a 内のエン트리 ( L V o l , L E I D ) で指定される論理ボリューム / 論理エクステントに対して物理エクステントの割り付けが行われると、当該エン트리中の列項目 P E I D が、“ N A ” から当該割り付けられた物理エクステントの I D で書き換えられる。

10

【 0 1 0 2 】

この図 1 6 に示すストレージクラスタの基本的な動作は、第 1 の実施形態と同様である。但し、論理ボリューム作成時と、論理エクステントへの初回のアクセス ( ライトアクセス ) 時の挙動は、第 1 の実施形態と異なる。

【 0 1 0 3 】

そこで、論理ボリューム作成時において第 1 の実施形態と異なる点について説明する。第 1 の実施形態では、前述した ( a ) , ( b ) , ( c ) の手順で論理ボリュームが作成される。これに対して第 3 の実施形態では、論理エクステントへの初回のライトアクセス時に当該論理エクステントに対する物理エクステントの割り付けが行われるため、上記 ( a ) の動作が次の ( a 1 ) のようになる。

20

【 0 1 0 4 】

( a 1 ) 構成管理部は、作成されるべき論理ボリュームの容量を確保するために、当該論理ボリュームを構成する論理エクステントの各々を、どのノード I D ( P V o l ) のストレージ ( 物理ボリューム ) に格納するか決定する。但し、上記 ( a ) とは異なり、各論理エクステントへの物理エクステントの割り付けは行われぬ。構成管理部は、ストレージ 1 3 -1 ~ 1 3 -3 毎に、各論理エクステントに対して物理エクステントが未割り付け状態にあることを示す L P M T 1 3 2 a を生成して、そのストレージに保持する。

【 0 1 0 5 】

次に、図 1 6 のストレージクラスタにおける、ホストから論理エクステントへのライトアクセス時の動作について、図 1 乃至図 3 を援用しながら、図 1 7 のフローチャートを参照して説明する。図 1 7 において、図 4 のフローチャートと同様の部分には同一符号を付してある。

30

【 0 1 0 6 】

図 1 7 のフローチャートが、図 4 のフローチャートと異なるのは、ある論理エクステントへの初回のアクセスがライトアクセスの場合に、当該論理エクステントに物理エクステントが割り付けられていないため、割り付け先の物理エクステント先を決定する処理 ( ステップ S 1 8 ) が追加されている点である。

【 0 1 0 7 】

今、第 1 の実施形態と同様に、ホストからノード 1 -1 のコントローラ 1 1 -1 に含まれているターゲットインタフェース 1 1 1 に対してアクセス ( ライトアクセス ) が要求されたものとする。この場合、図 4 のフローチャートと同様に、まずステップ S 0 ~ S 7 が実行される。但し図 1 7 では、ステップ S 0 ~ S 7 のうちのステップ S 1 ~ S 6 が省略されている。ステップ S 7 では、ノード I D = S I D x で示されるノードのストレージ、例えばノード 1 -1 のストレージ 1 3 -1 に対し、論理エクステント ( L V o l x , L E I D x ) に対応する物理エクステントへのアクセスが要求される。するとストレージ 1 3 -1 では、ステップ S 1 8 が次のように実行される。

40

【 0 1 0 8 】

まず、ストレージ 1 3 -1 の I O ドライバ 1 3 1 は、論理エクステント ( L V o l x , L E I D x ) に対応付けられている L P M T 1 3 2 a 内のエントリの列項目 P E I D を参照する ( ステップ S 1 8 a ) 。ストレージ 1 3 -1 は、この P E I D が非 N A であるか否かに

50

より、論理エクステント (LVOLx, LEIDx) に物理エクステントが割り付けられているか否かを判定する (ステップ S18b)。

【0109】

もし、論理エクステント (LVOLx, LEIDx) に物理エクステントが割り付けられていないならば、I/Oドライバ131は構成管理部と連携して、論理エクステント (LVOLx, LEIDx) にストレージ13-1 (物理ボリューム130-1) 内の適当な物理エクステントPEIDxを割り付ける (ステップ S18c)。このステップ S18cにおいて、I/Oドライバ131は、論理エクステント (LVOLx, LEIDx) に対応付けられているLPMT132a内のエントリの列項目PEIDを、“NA”から物理エクステントPEIDxのID (= PEIDx) に更新する。そしてI/Oドライバ131は、物理エクステントPEIDxに対してアクセスし、コントローラ11-1のI/Oマネージャ113に対して応答する (ステップ S18d)。

10

【0110】

一方、論理エクステント (LVOLx, LEIDx) に物理エクステント (例えば物理エクステントPEIDx) が既に割り付けられているならば (ステップ S18b)、I/Oドライバ131はステップ S18cをスキップしてステップ S18dを実行する。

【0111】

このように第3の実施形態においては、ある論理エクステントに対する初回のアクセスがライトアクセスの場合に、当該論理エクステントに物理エクステントが割り付けられる。これに対し、ある論理エクステントに対する初回のアクセスがリードアクセスの場合、つまりデータがライトされていない論理エクステントに対してリードアクセスが発生した場合には、当該論理エクステントに物理エクステントを割り付けることも、割り付けないことも可能である。割り付ける場合は、I/Oドライバ131によって上記ステップ S18が実行されれば良い。いずれの場合にも、リード要求に対する応答として、全てゼロなどの適切なデータが返されるようにすれば良い。なお、本実施形態 (第3の実施形態) で適用される物理エクステントの割り付け方法は、第2の実施形態にも同様に適用可能である。本実施形態においては、各ストレージにおいて、ホストからアクセスが要求された際の状態に応じて、最適な物理エクステントの割り付けが可能になる。

20

【0112】

[第4の実施形態]

次に本発明の第4の実施形態について説明する。この第4の実施形態の特徴は、実際にホストからある論理エクステントへの初回のアクセスが要求されるまで、当該論理エクステントの領域の予約が行われない点にある。

30

【0113】

図18は、本発明の第4の実施形態に係るストレージクラスタ構成のストレージシステムの構成を示すブロック図である。図18において、図7と同様の要素には便宜的に同一参照符号を付してある。図18のストレージシステム (ストレージクラスタ) は、前記第1の実施形態と同様に、3つのノード1-1~1-3から構成される。図18のストレージクラスタが第1の実施形態と異なるのは、実際にホストからある論理エクステントへの初回のアクセスが要求されるまで、当該論理エクステントへの物理エクステントの割り付けだけでなく、当該論理エクステントへのストレージ (物理ボリューム) の割り付け (予約) も行われない点にある。

40

【0114】

このように、図18のストレージクラスタでは、ある論理エクステントへの初回のアクセスが要求されるまで、当該論理エクステントの領域の予約が行われないため、論理ボリューム作成時には、ストレージ13-1~13-3のLPMT132aのエントリを予約する必要もない。論理エクステントに対するストレージ (物理ボリューム) 及び物理エクステントの割り付けは、コントローラ11-i (i = 1, 2, 3) からの所在管理サーバ12-j に対する問い合わせに応じて当該所在管理サーバ12-j (j = 1, 2, 3) によって行われる。この物理エクステントの割り付け計画は、例えば構成管理部により行われる。

50

## 【 0 1 1 5 】

図 1 8 の例では、論理ボリュームを構成する各論理エクステントに対するストレージ（物理ボリューム）の割り付けが未だ行われていない状態が示されている。つまり図 1 8 の例では、論理ボリューム L V o l = 1 / 論理エクステント L E I D = 1 及び論理ボリューム L V o l = 1 / 論理エクステント L E I D = 4 に対して、ストレージ（物理ボリューム）の割り付けが行われていない状態が、所在管理サーバ 1 2 -1 の L S M T 1 2 0 -1 によって示されている。同様に、論理ボリューム L V o l = 1 / 論理エクステント L E I D = 2 に対して、ストレージ（物理ボリューム）の割り付けが行われていない状態が、所在管理サーバ 1 2 -2 の L S M T 1 2 0 -2 によって示されている。また、論理ボリューム L V o l = 1 / 論理エクステント L E I D = 3 に対して、ストレージ（物理ボリューム）の割り付けが行われていない状態が、所在管理サーバ 1 2 -3 の L S M T 1 2 0 -3 によって示されている。ここでは、ストレージ（物理ボリューム）の未割り付けの状態は、L S M T 1 2 0 -1 ~ 1 2 0 -3 内で列項目 P V o l に “ N A ” が設定されているエントリにより示される。図 1 8 において、L S M T 1 2 0 -1 ~ 1 2 0 -3 のエントリから、ストレージ 1 3 -1 ~ 1 3 -3 の L P M T 1 3 2 a のエントリに向かう矢印が描かれていないことに注意されたい。

10

## 【 0 1 1 6 】

この図 1 8 に示すストレージクラスタの基本的な動作は、第 1 の実施形態と同様である。但し、論理ボリューム作成時と、論理エクステントへの初回のアクセス時の挙動は、第 1 の実施形態と異なる。

## 【 0 1 1 7 】

そこで、論理ボリューム作成時において第 1 の実施形態と異なる点について説明する。第 1 の実施形態では、前述した（ a ）, （ b ）, （ c ）の手順で論理ボリュームが作成される。これに対して第 4 の実施形態では、論理エクステントへの初回アクセス時に当該論理エクステントを格納するストレージ（物理ボリューム）が決定され、この決定されたストレージ（物理ボリューム）内で当該論理エクステントに対して物理エクステントが割り付けられる。このため、上記（ a ）, （ b ）の動作が次の（ a 2 ）, （ b 2 ）のようになる。

20

## 【 0 1 1 8 】

（ a 2 ）構成管理部は、作成されるべき論理ボリュームを構成する論理エクステントの各々を、どのノード I D （ P V o l ）のストレージ（物理ボリューム）に格納するか決定しない。このため構成管理部は、ストレージ 1 3 -1 ~ 1 3 -3 毎に、図 1 8 に示されるように空の状態の L P M T 1 3 2 a を生成して、そのストレージに保持する。

30

## 【 0 1 1 9 】

（ b 2 ）構成管理部は、どのノードの所在管理サーバに論理エクステントの所在管理サーバを担当させるかを決定する。但し、その論理エクステントをどのストレージ（物理ボリューム）に格納するかは決定しない。このため構成管理部は、所在管理サーバ 1 2 -1 ~ 1 2 -3 毎に、その所在管理サーバが担当する論理エクステントについてのエントリ情報であって、列項目 P V o l が “ N A ” のエントリの情報を含む L S M T 1 2 0 -1 ~ 1 2 0 -3 を生成して、その所在管理サーバ 1 2 -1 ~ 1 2 -3 内に設定する。なお、第 2 の実施形態と同様に、ストレージクラスタ内の全てのノード 1 -1 ~ 1 -3 に所在管理サーバ（ L S M T ）が設けられる必要はない。

40

## 【 0 1 2 0 】

次に、図 1 8 のストレージクラスタにおける、ホストから論理エクステントへのライトアクセス時の動作について、図 1 乃至図 3 を援用しながら、図 1 9 のフローチャートを参照して説明する。図 1 9 において、図 4 または図 1 7 のフローチャートと同様の部分には同一符号を付してある。

## 【 0 1 2 1 】

今、第 1 の実施形態と同様に、ホストからノード 1 -1 のコントローラ 1 1 -1 に含まれているターゲットインタフェース 1 1 1 に対してアクセスが要求されたものとする。この場合、図 4 のフローチャートと同様に、まずステップ S 0 ~ S 4 が実行される。図 1 9 では

50

、ステップS 0 ~ S 4のうちのステップS 1 ~ S 3が省略されている。ステップS 4では、例えばコントローラ1 1-1のサーバ特定部1 1 2が例えば所在管理サーバ1 2-1に対して、論理エクステント(ターゲット論理エクステント)(LV o l x , LE I D x)に対応する物理エクステントを格納しているストレージのノードID = S I D xを問い合わせる。すると所在管理サーバ1 2-1では、ステップS 2 5が次のように実行される。

【0 1 2 2】

まず、所在管理サーバ1 2-1は、どのノードID = S I D xのストレージ(物理ボリューム)が、論理エクステント(LV o l x , LE I D x)に対応する物理エクステントを格納しているかを調べるために、当該論理エクステント(LV o l x , LE I D x)に対応付けられているLSMT 1 2 0-1内のエントリの列項目PV o l (= S I D x)を参照する(ステップS 2 5 a)。

10

【0 1 2 3】

所在管理サーバ1 2-1は、このPV o l が非NAであるか否かにより、論理エクステント(LV o l x , LE I D x)にストレージ(物理ボリューム)が割り付けられているか否かを判定する(ステップS 2 5 b)。

【0 1 2 4】

もし、論理エクステント(LV o l x , LE I D x)にストレージ(物理ボリューム)が割り付けられていないならば、所在管理サーバ1 2-1は、ストレージ1 3-1~ 1 3-3(の物理ボリューム1 3 0-1~ 1 3 0-3)における物理エクステントの空きを確認する(ステップS 2 5 c)。この確認処理は、例えば所在管理サーバ1 2-1がストレージ1 3-1~ 1 3-3のI Oドライバ1 3 1と通信を行うことにより実現される。

20

【0 1 2 5】

もし、ストレージ1 3-1~ 1 3-3のいずれにも物理エクステントの空きがない(または要求に応えられない)場合には、I Oエラーとなる(ステップS 2 5 d)。この場合、ホストからのアクセス要求に対する応答として、エラーが返される(ステップS 9)。

【0 1 2 6】

一方、ストレージ1 3-1~ 1 3-3のいずれかに物理エクステントの空きがあるならば、所在管理サーバ1 2-1は物理エクステントに空きのあるストレージ(例えばストレージ1 3-1)を選択して、論理エクステント(LV o l x , LE I D x)に当該ストレージを割り付ける(ステップS 2 5 e)。このステップS 2 5 eにおいて、所在管理サーバ1 2-1は、論理エクステント(LV o l x , LE I D x)へのストレージの割り付けがLSMT 1 2 0-内の対応するエントリに反映されるように、当該エントリを更新する。ここでは、論理エクステント(LV o l x , LE I D x)に対応付けられているLSMT 1 2 0-1内のエントリの列項目PV o l が、“NA”から当該論理エクステント(LV o l x , LE I D x)に割り付けられるストレージ1 3-1のノードID(物理ボリューム1 3 0-1のノードID = PV o l = 1)に更新される。

30

【0 1 2 7】

ステップS 2 5 eにおいてLSMT 1 2 0-1が更新されると、第1の実施形態で図4のフローチャートのステップS 4が実行された場合と同様に、ステップS 5 , S 6が実行される。これにより、所在管理サーバ1 2-1からコントローラ1 1-1に対し、論理エクステント(LV o l x , LE I D x)に割り付けられるストレージ1 3-1のノードID(物理ボリューム1 3 0-1のノードID = PV o l ) = S I D x = 1が問い合わせ元のサーバ特定部1 1 2を介して通知される。なお、論理エクステント(LV o l x , LE I D x)にストレージ(物理ボリューム)が割り付けられている場合には(ステップS 2 5 b)、直ちにステップS 5 , S 6が実行される。

40

【0 1 2 8】

ターゲットインタフェース1 1 1は、サーバ特定部1 1 2を介して通知されたノードID = S I D x = 1を受け取ると、第1及び第3の実施形態と同様に、ノードID = S I D x = 1で示されるノード1-1のストレージ1 3-1に対し、論理エクステント(LV o l x , LE I D x)に対応する物理エクステントへのアクセスをI Oマネージャ1 1 3を介し

50

て要求する（ステップS7）。このステップS7以降の動作は、第3の実施形態と同様であり、ステップS18a～S18dからなるステップS18がストレージ13-1において実行される。これにより、論理エクステント（LVolx, LEIDx）に物理エクステントが割り付けられていない場合には、物理ボリューム130-1内の空き物理エクステントが当該論理エクステント（LVolx, LEIDx）に割り付けられる。ここでは、そして、物理エクステントPEIDxが割り付けられたものとする。この場合、論理エクステント（LVolx, LEIDx）に対応付けられているLPMT132a内のエントリの列項目PEIDを、“NA”から物理エクステントPEIDxのID（=PEIDx）に更新される。

【0129】

本実施形態においては、各ストレージにおいて、ホストからアクセスが要求された際の状態に応じて、最適な物理エクステントの割り付けが可能になるだけでなく、総物理ボリューム容量の上限値を超える容量の論理ボリュームを作成することができる。

【0130】

[第5の実施形態]

次に本発明の第5の実施形態について説明する。この第5の実施形態の特徴は、コントローラ11-i（i=1, 2, 3）のサーバ特定部112から所在管理サーバ12-j（j=1, 2, 3）への問い合わせと当該問い合わせに対する回答（結果）とに関する情報をキャッシュ112bに保持することにある。この情報は、論理ボリュームLVolx/論理エクステントLEIDx、つまり論理エクステント（LVolx, LEIDx）を示す情報（LVolx, LEIDx）と、当該論理エクステント（LVolx, LEIDx）に対応している物理エクステントを格納しているストレージのノードID=SIDxとを含む。このキャッシュ112bを利用することで、再度の問い合わせが必要となった場合に、当該問い合わせに対する回答に相当する情報をキャッシュ112bから速やかに取得することが可能となる。

【0131】

次に、第5の実施形態における、ホストから論理エクステントへのアクセス時の動作について、図1乃至図3及び図7を援用しながら、図20のフローチャートを参照して説明する。

【0132】

今、第1の実施形態と同様に、ホストからノード1-1のコントローラ11-1に含まれているターゲットインタフェイス111に対してアクセスが要求されたものとする。この場合、図4のフローチャートと同様に、まずステップS0～S2が実行される。図20では、ステップS0～S2のうちのステップS0, S1が省略されている。ステップS2では、コントローラ11-1のターゲットインタフェイス111からコントローラ11-1のサーバ特定部112に対して、論理エクステント（LVolx, LEIDx）が、どのノードのストレージに格納されているかが問い合わせられる。

【0133】

するとサーバ特定部112は、第1の実施形態とは異なって、自身が管理しているキャッシュ112bを参照して、論理エクステント（LVolx, LEIDx）に対応している物理エクステントを格納しているストレージのノードID=SIDxを調べる（ステップS31）。

【0134】

もし、目的のノードID=SIDxが（LVolx, LEIDx）と対応付けてキャッシュ112bに格納されていないならば、つまりミスヒットならば、サーバ特定部112は、第1の実施形態におけるターゲットインタフェイス111からの問い合わせ（ステップS2）に対するのと同様に、ステップS3及びS4を実行する。

【0135】

本実施形態では、論理エクステント（LVolx, LEIDx）に対応している物理エクステントが格納されているストレージの所在が、第1の実施形態と同様にノード1-1の

10

20

30

40

50

所在管理サーバ12-1によって管理され、当該ストレージのノードID (= PVOL) が SIDx = 1 であるものとする。このノードIDはLSMT120-1に基づいて求められる。この場合、所在管理サーバ12-1は、論理エクステント (LVOLx, LEIDx) に対応している物理エクステントが格納されているストレージのノードID = SIDx と、LSMT120-1のシリアル番号 (LSMTシリアル番号) SNとを、問い合わせ元のサーバ特定部112に対して通知する (ステップS5a)。第1の実施形態のステップS5との違いは、ノードIDに加えて、当該ノードIDを求めるのに用いられたLSMT120-1のシリアル番号SNが通知される点である。このシリアル番号SNは、LSMT120-1が更新される毎に、所在管理サーバ12-1によって例えばインクリメントされる。このことは、LSMT120-2, 120-3についても同様である。

10

## 【0136】

サーバ特定部112は、所在管理サーバ12-1からの通知を受けて、キャッシュ112bに、{(LVOLx, LEIDx), SIDx, SN}を含む情報 (つまり、問い合わせと回答とに関する情報) を格納する (ステップS5b)。これにより、サーバ特定部112は、論理エクステント (LVOLx, LEIDx) がどのノードのストレージに格納されているかをターゲットインタフェース111によって再び問い合わせられた場合に、当該問い合わせに対する回答に相当する情報を上記ステップS31の処理で速やかに取得することができる。つまり、サーバ特定部112から所在管理サーバへの問い合わせの回数を減らして、アクセス速度を向上させることができる。

20

## 【0137】

但し、所在管理サーバ12-1~12-3上のLSMT120-1~120-3が変更された場合、その変更がキャッシュ112bに反映される必要がある。その理由は、例えば第1の実施形態のように、ストレージ間で物理エクステントのコピー (冗長構成) を管理する場合 (図12参照)、LSMT120-1~120-3の更新がキャッシュ112bに反映されていないと、ライトアクセスによって一方の物理エクステントのみが書き換えられてしまい、冗長構成が崩れてしまうからである。そのため、本実施形態では、以下のようなキャッシュ更新の仕組みが用いられる。

## 【0138】

まず、LSMT120-1~120-3の各々にシリアル番号SNが付与される。このシリアル番号SNは、ストレージ13-1~13-3のI/Oドライバ131によってLSMT120-1~120-3が更新される都度、例えば当該I/Oドライバ131によってインクリメント (更新) される。

30

## 【0139】

また、ステップS5aにおいてノードID = SIDxを求めるのに用いられたLSMT120-j (j = 1, 2, 3) のシリアル番号SNが、(LVOLx, LEIDx), SIDxと組をなしてキャッシュ112bに格納される (ステップS5b)。

## 【0140】

コントローラ11-i (i = 1, 2, 3) は、ストレージ13-jへのアクセスを行うときに、このシリアル番号SNを受け渡す。ストレージ13-jのI/Oドライバ131は、このシリアル番号SNと自身が認識しているLSMT120-jの最新のシリアル番号SNとを比較する。I/Oドライバ131は、両シリアル番号SNが一致していないならば、コントローラ11-iのサーバ特定部112が管理するキャッシュ112bには、LSMT120-jの最新の更新内容が反映されていないと判定する。この場合、ストレージ13-jのI/Oドライバ131はコントローラ11-iに対してキャッシュ112bの該当する情報の無効化を通知する。つまりI/Oドライバ131は、キャッシュ無効化手段として機能する。コントローラ11-iのサーバ特定部112は、この通知に応じて、ステップS32でミスヒットが判定された場合と同様に、所在管理サーバ12-jに対して再度の問い合わせを行う (ステップS3, S4)。これによりサーバ特定部112は、所在管理サーバ12-jからLSMT120-jの最新の情報を取得して、キャッシュ112bに反映させることができる。

40

50

## 【0141】

なお、以下に述べるように、所在管理サーバ12-jがキャッシュ無効化手段として機能しても良い。つまり、所在管理サーバ12-jがコントローラ11-jのサーバ特定部112上のキャッシュ112bを監視して、自身が管理するLSMT120-jとの不整合が生じているかを判別する構成とすることも可能である。もし不整合を判別した場合、所在管理サーバ12-jからコントローラ11-jのサーバ特定部112に、キャッシュ112bの該当するエントリ情報の破棄・更新が必要なことを通知すれば良い。

## 【0142】

なお、本発明は、上記各実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記各実施形態に開示されている複数の構成要素の適宜な組み合わせにより種々の発明を形成できる。例えば、各実施形態に示される全構成要素から幾つかの構成要素を削除してもよい。

## 【図面の簡単な説明】

## 【0143】

【図1】本発明の第1の実施形態に係るストレージクラスタ構成のストレージシステムの概略構成を示すブロック図。

【図2】図1に示されるコントローラの構成を示すブロック図。

【図3】図1に示されるストレージの構成を示すブロック図。

【図4】同第1の実施形態においてホストからアクセス要求が発行された場合の処理の手順を示すフローチャート。

【図5】同第1の実施形態においてホストからアクセス要求が発行された場合の情報の流れを示す図。

【図6】同第1の実施形態においてホストからアクセス要求が発行された場合の情報の流れの変形例を示す図。

【図7】図1のシステムにおいて、論理ボリュームLV01がLV01=1であって、当該論理ボリュームLV01=1を、LEIDがLEID=1~LEID=4の4つの論理エクステントから構成した場合のシステム構成を示す図。

【図8】図7の状態におけるLMMTのデータ構造例を示す図。

【図9】図7の状態におけるLSMTのデータ構造例を示す図。

【図10】図7の状態におけるLPM Tのデータ構造例を示す図。

【図11】同第1の実施形態において、論理エクステントを構成する物理エクステントのデータを、ストレージをまたがって別の物理エクステントにマイグレートする場合の動作を説明するための図。

【図12】同第1の実施形態において、論理エクステントを構成する物理エクステントのデータを、ストレージをまたがって別の物理エクステントにコピーする場合の動作を説明するための図。

【図13】同第1の実施形態において、論理エクステントに関する所在管理サービスを、ノードをまたがってマイグレートする場合の動作を説明するための図。

【図14】同第1の実施形態において、論理エクステントに関する所在管理サービスを、ノードをまたがってコピーする場合の動作を説明するための図。

【図15】本発明の第2の実施形態に係るストレージクラスタ構成のストレージシステムの構成を示すブロック図。

【図16】本発明の第3の実施形態に係るストレージクラスタ構成のストレージシステムの構成を示すブロック図。

【図17】同第3の実施形態においてホストからアクセス要求が発行された場合の処理の手順を示すフローチャート。

【図18】本発明の第4の実施形態に係るストレージクラスタ構成のストレージシステムの構成を示すブロック図。

【図19】同第4の実施形態においてホストからアクセス要求が発行された場合の処理の手順を示すフローチャート。

10

20

30

40

50



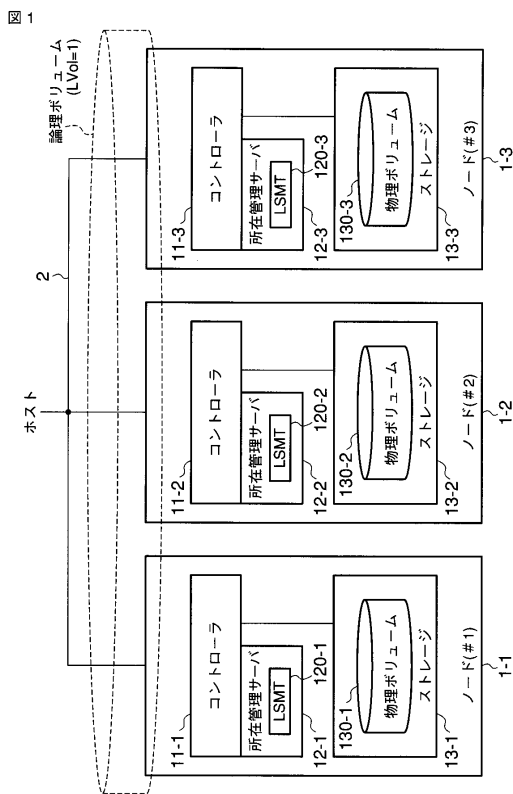
【図20】本発明の第5の実施形態においてホストからアクセス要求が発行された場合の処理の手順を示すフローチャート。

【符号の説明】

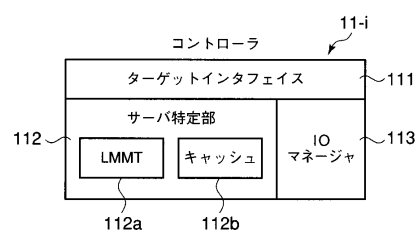
【0144】

1-1~1-3...ノード、11-1~11-3, 11-i...コントローラ、12-1~12-3...所在管理サーバ、13-1~13-3, 13-i...ストレージ、112a...LMMT(論理エクステント-所在管理サーバマッピングテーブル、第1のマッピングテーブル)、112b...キャッシュ、120-1~120-3, 120-i...LSMT(論理エクステント-ストレージマッピングテーブル、第2のマッピングテーブル)、130-1~130-3, 130-i...物理ボリューム、132a...LPMT(論理エクステント-物理エクステントマッピングテーブル、第3のマッピングテーブル)。

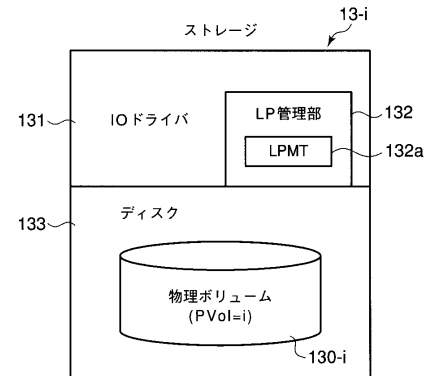
【図1】



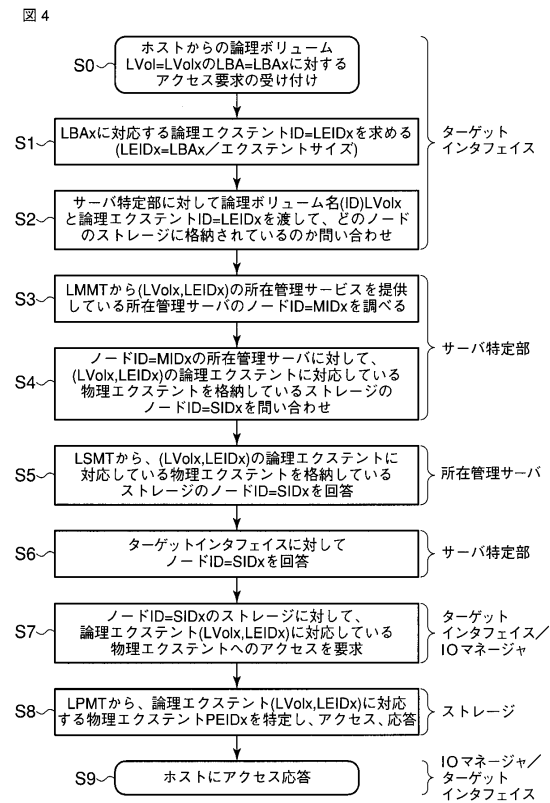
【図2】



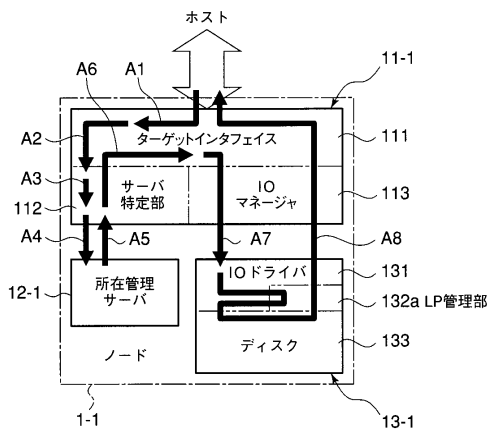
【図3】



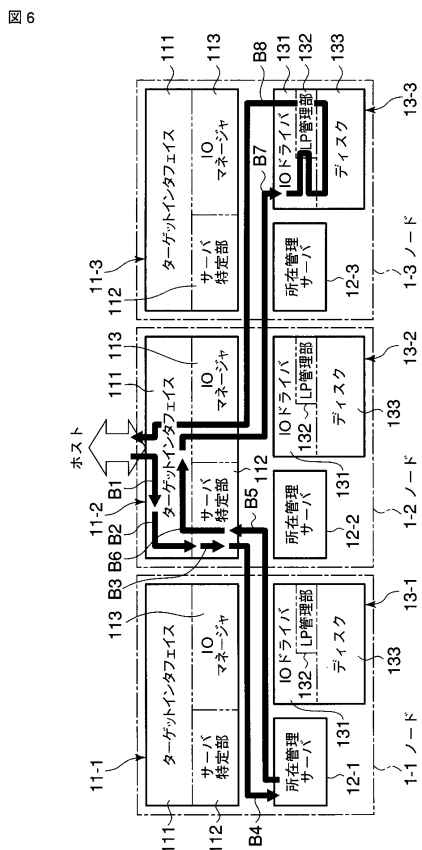
【図4】



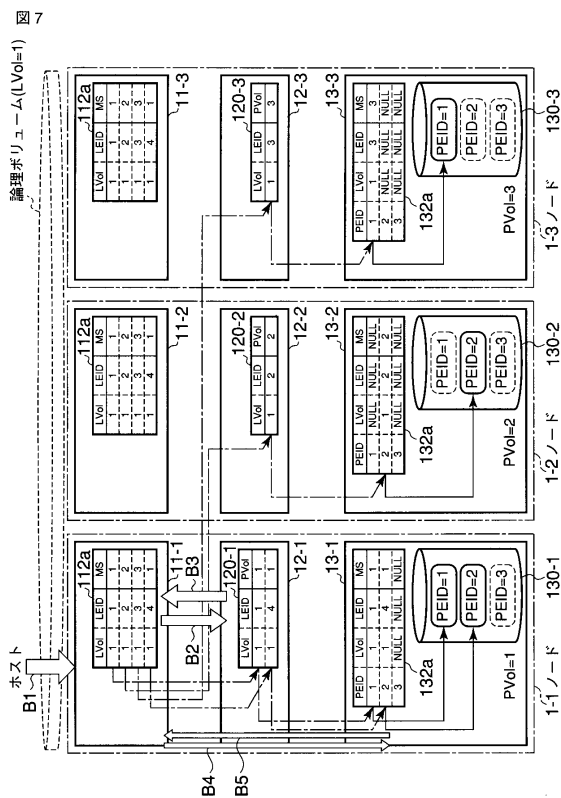
【図5】



【図6】



【図7】



【 8 】

図 8

LMMT		
LVol	LEID	MS
1	1	1
1	2	2
1	3	3
1	4	1

【 9 】

図 9

LSMT		
LVol	LEID	PVol
1	1	1
1	4	1

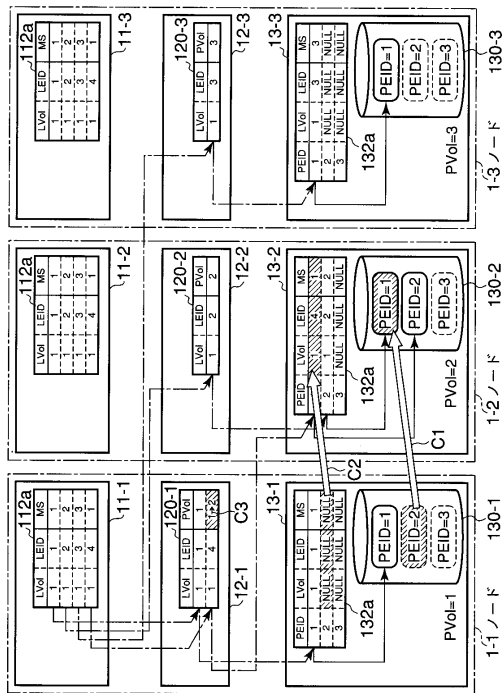
【 10 】

図 10

LPMT			
PEID	VLoI	LEID	MS
1	1	1	1
2	1	4	1
3	NULL	NULL	NULL

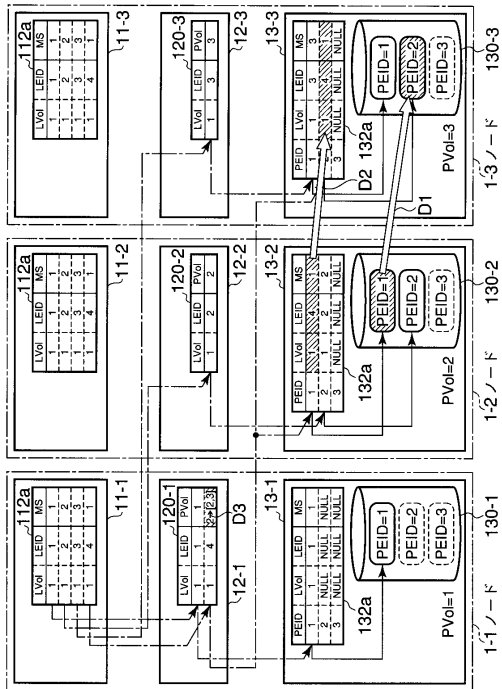
【 11 】

図 11



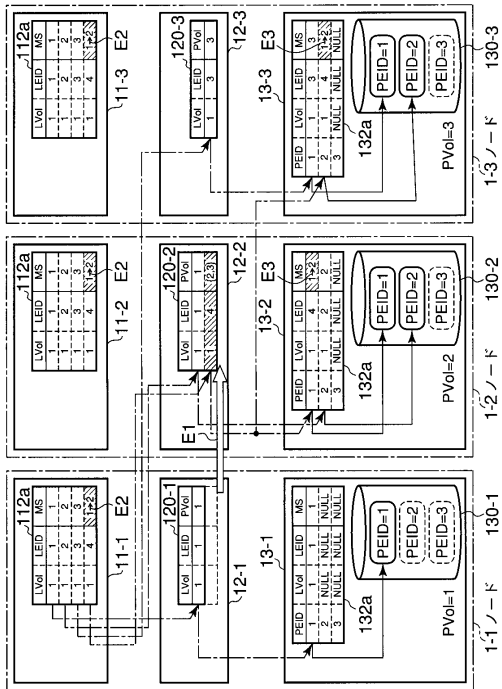
【 12 】

図 12

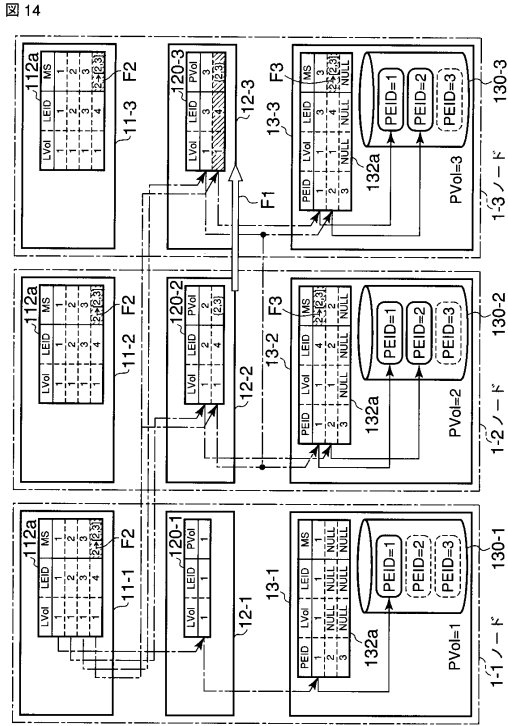


【 13 】

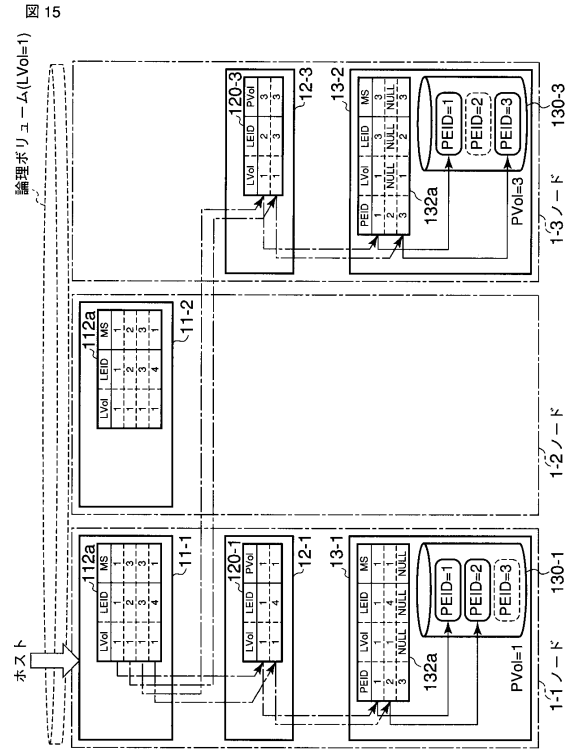
図 13



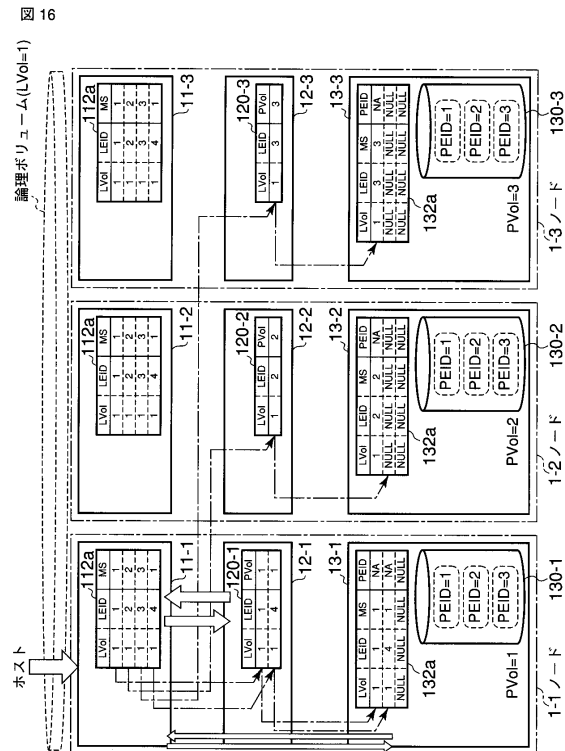
【 図 14 】



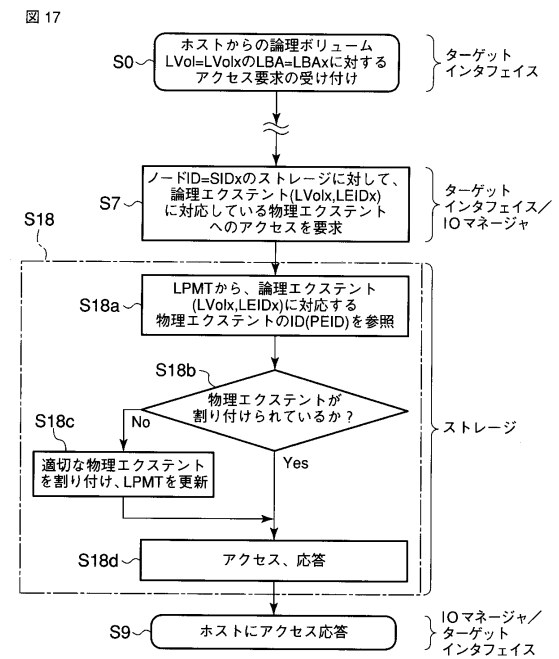
【 図 15 】



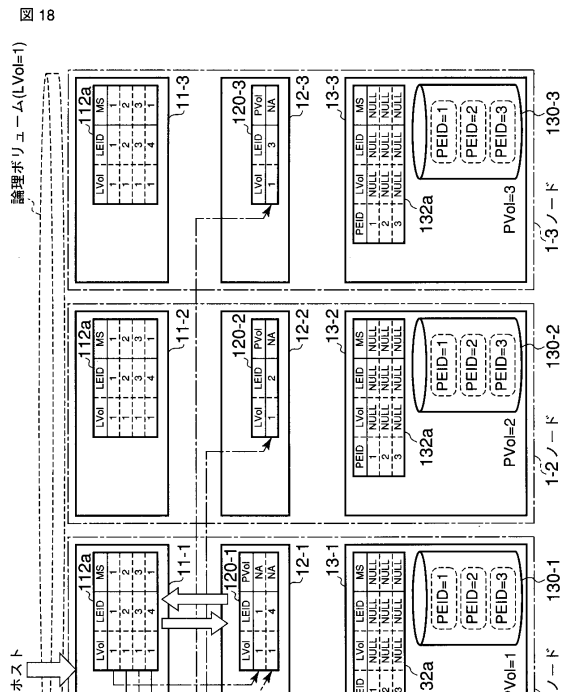
【 図 16 】



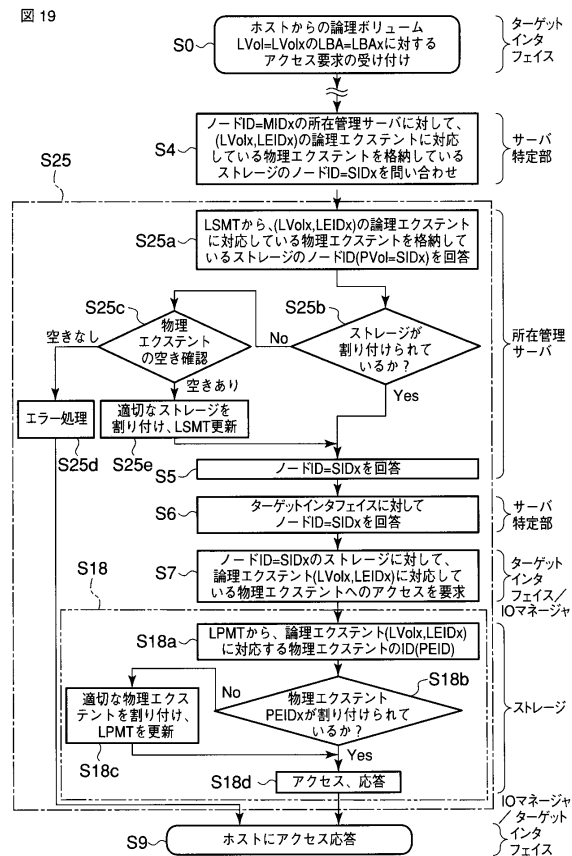
【 図 17 】



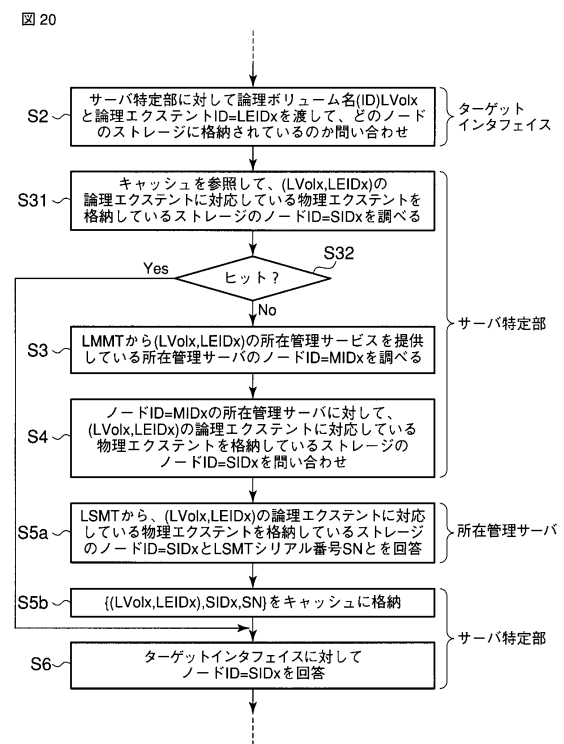
【図18】



【図19】



【図20】



---

フロントページの続き

(74)代理人 100075672

弁理士 峰 隆司

(74)代理人 100109830

弁理士 福原 淑弘

(74)代理人 100084618

弁理士 村松 貞男

(74)代理人 100092196

弁理士 橋本 良郎

(72)発明者 小原 誠

東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

審査官 木村 雅也

(56)参考文献 特開平07-152491(JP,A)

特開2004-078398(JP,A)

特開2001-312372(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06