

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7253007号  
(P7253007)

(45)発行日 令和5年4月5日(2023.4.5)

(24)登録日 令和5年3月28日(2023.3.28)

|                         |         |       |         |  |
|-------------------------|---------|-------|---------|--|
| (51)国際特許分類              | F I     |       |         |  |
| G 0 6 F 3/06 (2006.01)  | G 0 6 F | 3/06  | 3 0 1 Z |  |
| G 0 6 F 11/07 (2006.01) | G 0 6 F | 3/06  | 3 0 1 X |  |
| G 0 6 F 11/34 (2006.01) | G 0 6 F | 11/07 | 1 4 0 M |  |
| G 0 6 F 13/10 (2006.01) | G 0 6 F | 11/07 | 1 9 3   |  |
|                         | G 0 6 F | 11/34 | 1 3 3   |  |
| 請求項の数 15 (全46頁) 最終頁に続く  |         |       |         |  |

|          |                                  |          |                    |
|----------|----------------------------------|----------|--------------------|
| (21)出願番号 | 特願2021-90224(P2021-90224)        | (73)特許権者 | 000005108          |
| (22)出願日  | 令和3年5月28日(2021.5.28)             |          | 株式会社日立製作所          |
| (65)公開番号 | 特開2022-182584(P2022-182584<br>A) | (74)代理人  | 110002365          |
| (43)公開日  | 令和4年12月8日(2022.12.8)             |          | 弁理士法人サンネクスト国際特許事務所 |
| 審査請求日    | 令和3年12月23日(2021.12.23)           | (72)発明者  | 山本 貴大              |
|          |                                  |          | 東京都千代田区丸の内一丁目6番6号  |
|          |                                  | (72)発明者  | 株式会社日立製作所内         |
|          |                                  | (72)発明者  | 坂下 悠貴              |
|          |                                  |          | 東京都千代田区丸の内一丁目6番6号  |
|          |                                  | (72)発明者  | 株式会社日立製作所内         |
|          |                                  | (72)発明者  | 伊藤 晋太郎             |
|          |                                  |          | 東京都千代田区丸の内一丁目6番6号  |
|          |                                  | (72)発明者  | 株式会社日立製作所内         |
|          |                                  | (72)発明者  | 揚妻 匡邦              |
|          |                                  |          | 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】

複数の領域を含むボリュームを1以上のホストに提供するための処理を行うプロセッサを備える複数のノードと、前記プロセッサと接続され、前記ボリュームのデータを記憶する1以上の記憶デバイスとを備えるストレージシステムであって、

前記複数のノードの各々は、自ノードが提供するボリュームの負荷および前記ボリュームの領域を複数に分割した領域の負荷を監視し、

監視している一のボリュームの負荷が閾値以上であると判定した第1のノードは、前記一のボリュームの領域を複数に分割した領域の負荷と負荷分散のポリシーとに応じて、前記一のボリュームに含まれる一部の領域を前記第1のノードとは異なる第2のノードのボリュームに移動する、

ストレージシステム。

【請求項2】

前記第1のノードは、前記一のボリュームの負荷が閾値未満であると判定した場合、前記第2のノードのボリュームから、前記一部の領域を前記一のボリュームに移動する、

請求項1に記載のストレージシステム。

【請求項3】

前記第1のノードは、前記第2のノードとして、前記一部の領域を移動した後のボリュームの負荷が閾値を超えないノードを選択する、

請求項1に記載のストレージシステム。

**【請求項 4】**

前記複数のノードの各々には、前記 1 以上の記憶デバイスの少なくとも 1 つが対応して設けられ、

前記複数のノードの各々は、自ノードに割り当てられている領域のデータを、自ノードに設けられている記憶デバイスに記憶し、

前記第 1 のノードは、前記一のボリュームの容量が提供可能な容量を超えると判定した場合、前記一部の領域を前記第 2 のノードのボリュームに移動する、

請求項 1 に記載のストレージシステム。

**【請求項 5】**

前記複数のノードの各々は、自ノードが提供するボリュームに対するリードの負荷と、前記ボリュームに対するライトの負荷とを監視し、

前記第 1 のノードは、前記一のボリュームに対するリードの負荷が第 1 の閾値以上であると判定した場合、前記一部の領域を前記第 2 のノードのボリュームに移動し、前記一のボリュームに対するライトの負荷が前記第 1 の閾値とは異なる第 2 の閾値以上であると判定した場合、前記一部の領域を前記第 2 のノードのボリュームに移動する、

請求項 1 に記載のストレージシステム。

**【請求項 6】**

前記第 1 のノードは、前記複数のノードに対して前記一のボリュームに含まれる領域を均等に割り振り、前記第 1 のノードとは異なるノードに割り振った領域を、前記ノードのボリュームに移動する、

請求項 1 に記載のストレージシステム。

**【請求項 7】**

前記第 1 のノードは、前記一のボリュームの負荷が前記閾値を下回るまで、前記一のボリュームに含まれる領域を 1 つずつ、前記第 1 のノードとは異なるノードのボリュームに移動する、

請求項 1 に記載のストレージシステム。

**【請求項 8】**

前記一のボリュームが、複数のホストに提供されている場合、前記第 1 のノードは、前記複数のホストの各々がアクセスする領域をまとめてホストごとの移動対象とし、ホストごとの移動対象の領域を、前記第 1 のノードとは異なるノードのボリュームに移動する、

請求項 1 に記載のストレージシステム。

**【請求項 9】**

前記第 1 のノードは、自ノードが提供しているボリュームのうち前記一のボリュームとは異なる他のボリュームを前記第 1 のノードとは異なるノードに移動し、前記第 2 のノードのボリュームから、前記一部の領域を前記一のボリュームに移動する、

請求項 1 に記載のストレージシステム。

**【請求項 10】**

前記一部の領域が移動された前記第 2 のノードのボリュームと、前記一部の領域にアクセスするホストとの間にパスが設定されていない場合、前記第 2 のノードおよび前記ホストは、前記パスを設定する、

請求項 1 に記載のストレージシステム。

**【請求項 11】**

前記複数のノードの各々は、自ノードのボリュームの負荷が特定のホストに偏っていると判定した場合、前記特定のホストとのパスが最適属性であることを前記特定のホストに通知する、

請求項 1 に記載のストレージシステム。

**【請求項 12】**

前記複数のノードの各々は、自ノードのボリュームの負荷が前記特定のホストに偏っていないと判定した場合、前記ボリュームに定義されている全てのパスが最適属性であることを前記パスが設定されている全てのホストに通知する、

10

20

30

40

50

請求項 1 1 に記載のストレージシステム。

【請求項 1 3】

前記 1 以上の記憶デバイスは、前記複数のノードに共通して設けられている、  
請求項 1 に記載のストレージシステム。

【請求項 1 4】

前記第 1 のノードは、前記第 2 のノードに前記一部の領域を移動する際、前記一部の領域を管理するためのデータを更新し、前記 1 以上の記憶デバイスに記憶される前記一部の領域のデータを移動しない、

請求項 1 3 に記載のストレージシステム。

【請求項 1 5】

前記複数のノードが提供するボリュームについて、当該ボリュームの領域の移動先ノードの数に応じたスループットと応答時間とを計算して出力する計算機を備える、

請求項 1 に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、概して、ボリュームの負荷を分散する技術に関する。

【背景技術】

【0002】

特許文献 1 は、複数のサーバをネットワークで接続して、各サーバのローカルストレージをストレージ制御ソフトウェアにより統合して、1つのストレージプールとして提供し、当該ストレージプールからボリュームを提供するストレージシステムについて開示している。上記ストレージシステムは、当該ボリュームにデータを書き込む際に、異なるサーバに格納されるデータを組み合わせ、パリティを計算し、データとは異なるサーバに格納することで、サーバ障害からデータを保護する。上記ストレージシステムでは、サーバの追加により、ストレージの容量と性能とをスケールアウトできることが特徴となっている。

【0003】

また、特許文献 1 は、ボリュームへ書き込まれたデータについてアクセス頻度を採取し、アクセス頻度の高いデータは、ローカルストレージに格納し、アクセス頻度の低いデータは異なるサーバのストレージ（ここでは、リモートストレージと呼ぶ）に格納するようにデータ配置を変更する技術についても開示している。ホストがボリュームのデータにアクセスする際は、ボリュームのローカルストレージを割り当てたサーバにアクセスし、当該サーバにて、データがローカルストレージにあるか、リモートストレージにあるかを判定し、リモートストレージにある場合は、異なるサーバにアクセスを転送し、データへのアクセスを行う。このように、アクセス頻度の高いデータをローカルストレージに格納しておくことで、ネットワークを介することなくデータへアクセスできるため、ホストに対して高速に応答することができる。

【先行技術文献】

【特許文献】

【0004】

【文献】米国特許出願公開第 2016 / 0371145 号明細書

【発明の概要】

【発明が解決しようとする課題】

【0005】

特許文献 1 に記載の技術に基づいてホストにボリュームを提供する場合、ボリュームへのアクセスは、必ずローカルストレージを有するサーバを経由してからアクセスする。このため、1つのボリュームの性能がローカルストレージを有するサーバの性能が上限となる。上記ストレージシステムは、サーバの追加により、性能がスケールアウトすることが特徴である。特許文献 1 で開示されている技術では、システム性能（複数のボリュームの合計性能）は、サーバを追加することで、スケールアウトできるが、1つのボリュームの

10

20

30

40

50

性能は、サーバを追加してもスケールアウトすることができない。

【0006】

本発明は、以上の点を考慮してなされたもので、サーバ（ノード）を追加したときにシステム性能がスケールアウトすると共に、1つのボリュームについても性能をスケールアウトし得るストレージシステム等を提案しようとするものである。

【課題を解決するための手段】

【0007】

かかる課題を解決するため本発明においては、複数の領域を含むボリュームを1以上のホストに提供するための処理を行うプロセッサを備える複数のノードと、前記プロセッサと接続され、前記ボリュームのデータを記憶する1以上の記憶デバイスとを備えるストレージシステムであって、前記複数のノードの各々は、自ノードが提供するボリュームの負荷および前記ボリュームの領域を複数に分割した領域の負荷を監視し、監視している一のボリュームの負荷が閾値以上であると判定した第1のノードは、前記一のボリュームの領域を複数に分割した領域の負荷と負荷分散のポリシとに応じて、前記一のボリュームに含まれる一部の領域を前記第1のノードとは異なる第2のノードのボリュームに移動するようにした。

10

【0008】

上記構成では、一のボリュームの負荷が高まったときに、当該ボリュームの一部の領域が他のノードのボリュームに移動されるので、例えば、ノードを追加した場合、一のボリュームについても性能をスケールアウトすることができるようになる。

20

【発明の効果】

【0009】

本発明によれば、一のボリュームの性能をスケールアウトし得るストレージシステムを実現することができる。

【図面の簡単な説明】

【0010】

【図1】第1の実施の形態によるボリュームのデータの割り当て変更の概要を示すイメージ図である。

【図2】第1の実施の形態によるストレージシステムに係る物理構成の一例を示す図である。

30

【図3】第1の実施の形態によるストレージシステムに係る論理構成の一例を示す図である。

【図4】第1の実施の形態によるメモリ内の情報の一例を示す図である。

【図5】第1の実施の形態によるクラスタ管理テーブルの一例を示す図である。

【図6】第1の実施の形態によるデータ保護セット管理テーブルの一例を示す図である。

【図7】第1の実施の形態によるストレージプール管理テーブルの一例を示す図である。

【図8】第1の実施の形態によるボリューム管理テーブルの一例を示す図である。

【図9】第1の実施の形態によるモニタ情報管理テーブルの一例を示す図である。

【図10】第1の実施の形態によるフロントエンドパス管理テーブルの一例を示す図である。

40

【図11】第1の実施の形態によるリード処理に係るフローチャートの一例を示す図である。

【図12】第1の実施の形態によるライト処理に係るフローチャートの一例を示す図である。

【図13A】第1の実施の形態によるモニタ情報採取処理に係るフローチャートの一例を示す図である。

【図13B】第1の実施の形態によるモニタ情報採取処理に係るフローチャートの一例を示す図である。

【図14】第1の実施の形態によるリバランス要否判定処理に係るフローチャートの一例を示す図である。

50

【図15】第1の実施の形態によるリソース割当決定処理に係るフローチャートの一例を示す図である。

【図16】第1の実施の形態によるリソース割当決定処理に係るフローチャートの一例を示す図である。

【図17】第1の実施の形態によるリソース移動処理に係るフローチャートの一例を示す図である。

【図18】第1の実施の形態によるフロントエンドパス設定処理に係るフローチャートの一例を示す図である。

【図19】第1の実施の形態によるフロントエンドパス設定処理に係るフローチャートの一例を示す図である。

10

【図20】第1の実施の形態によるクラスタ構成変更処理に係るフローチャートの一例を示す図である。

【図21A】第1の実施の形態によるGUIの一例を示す図である。

【図21B】第1の実施の形態によるGUIの一例を示す図である。

【図22】第2の実施の形態によるストレージシステムに係る構成の一例を示す図である。

【図23】第2の実施の形態によるリソース移動処理に係るフローチャートの一例を示す図である。

【発明を実施するための形態】

【0011】

(I) 第1の実施の形態

20

以下、本発明の一実施の形態を詳述する。ただし、本発明は、実施の形態に限定されるものではない。

【0012】

本実施の形態に係るストレージシステムにおいては、ボリュームの領域を複数のスライスと呼ぶ領域に分割し、スライス単位で複数のサーバ計算機に領域を割り当て、ボリュームへのアクセス負荷をモニタリングする。本ストレージシステムでは、主にスライスが割り当てられたサーバ計算機に負荷が生じ、アクセス負荷が低く1つのサーバ計算機でボリュームが要求する性能を提供できる場合は、ボリュームを構成するスライスを1つのサーバ計算機に集約するように割り当てを制御する。また、ストレージシステムでは、アクセス負荷が高く1つのサーバ計算機でボリュームが要求する性能を提供できない場合は、ボリュームを構成するスライスを複数のサーバ計算機に分散して割り当てるように制御する。また、本ストレージシステムでは、ホストがボリュームのデータにアクセスする際は、各サーバ計算機にてアクセス先のスライスがどのサーバ計算機に割り当たっているか判定することで、アクセス時の負荷が特定のサーバ計算機に偏らないようにする。

30

【0013】

これにより、ボリュームのアクセス負荷が1つのサーバ計算機で充足する場合、必ずローカルストレージのデータに対してアクセスできるため、ホストに対して、高速に応答することができる。また、ボリュームのアクセス負荷が1つのサーバ計算機では充足しない場合、複数のサーバ計算機でアクセスを処理することでホストに対して高いスループット(IOPS: Input/Output Per Seconds)を提供できる。また、これらの制御は、ユーザが意識することなく、ストレージシステムが自動で行うため、ユーザは、特許文献1に記載のストレージシステムと変わらない運用負荷で上記の利益を得ることができる。

40

【0014】

本ストレージシステムによれば、1つのボリュームに対してもサーバ計算機の追加に合わせて、容量および性能をスケールアウトし、ボリュームのアクセス負荷に応じて、応答時間とスループットとを自動で好適な状態に変更できる。

【0015】

次に、本発明の実施の形態を図面に基づいて説明する。以下の記載および図面は、本発明を説明するための例示であって、説明の明確化のため、適宜、省略および簡略化がなさ

50

れている。本発明は、他の種々の形態でも実施することが可能である。特に限定しない限り、各構成要素は、単数でも複数でも構わない。

#### 【0016】

なお、以下の説明では、同種の要素を区別しないで説明する場合には、枝番を含む参照符号のうちの共通部分（枝番を除く部分）を使用し、同種の要素を区別して説明する場合は、枝番を含む参照符号を使用することがある。例えば、物理領域を特に区別しないで説明する場合には、「物理領域121」と記載し、個々の領域を区別して説明する場合には、「物理領域121-1」、「物理領域121-2」のように記載することがある。

#### 【0017】

本明細書等における「第1」、「第2」、「第3」等の表記は、構成要素を識別するために付するものであり、必ずしも、数または順序を限定するものではない。また、構成要素の識別のための番号は、文脈毎に用いられ、1つの文脈で用いた番号が、他の文脈で必ずしも同一の構成を示すとは限らない。また、ある番号で識別された構成要素が、他の番号で識別された構成要素の機能を兼ねるものではない。

#### 【0018】

図1は、ストレージシステムにおいて、ボリュームのデータの割り当て変更の概要を示すイメージ図である。ストレージシステム110からストレージシステム120に構成を変更する場合を例に挙げて説明する。

#### 【0019】

ストレージシステム110では、ホスト101にボリューム102が接続されており、ボリューム102内のデータ105A, 105B, 105Cが、ノード100A内のストレージプール103Aへ割り当てられている。ノード100Aは、ボリューム102に割当たったデータ105へのアクセス負荷をモニタリングする。モニタリングの結果、ストレージシステム110は、ノード100Aで提供できる性能を超える負荷を検出した場合、ボリューム102のデータ105B, 105Cをノード100B, 100Cに移動し、負荷を分散し、ストレージシステム120の状態に遷移する。本制御により、ボリューム102の高負荷時は、多数のノード100に処理を分散することで単体のボリューム102について高い性能を提供できるようにする。

#### 【0020】

ストレージシステム120では、ホスト101にボリューム102が接続されている。ボリューム102内のデータ105Aがノード100A内のストレージプール103Aに割り当てられている。データ105Bがノード100B内のストレージプール103Bに割り当てられている。データ105Cがノード100C内のストレージプール103Cに割り当てられている。ノード100A, 100B, 100Cは、ボリューム102に割当たったデータ105A, 105B, 105Cへのアクセス負荷をモニタリングする。モニタリングの結果、ストレージシステム120は、各データ105のアクセス負荷がノード100Aで提供できる性能を超えない負荷であることを検出した場合、ボリューム102のデータ105B, 105Cをノード100Aへ移動し、負荷を集約し、ストレージシステム110の状態に遷移する。本制御により、ボリューム102の低負荷時は、単一ノード100で処理を集約することでネットワークの利用効率を高め、ストレージシステム全体について、高い性能を提供できるようにする。

#### 【0021】

図2は、ストレージシステム200に係る物理構成の一例を示す図である。

#### 【0022】

ストレージシステム200には、1以上のサイト201が設けられてもよい。各サイト201は、ネットワーク202を介して通信可能に接続される。ネットワーク202は、例えば、WAN(Wide Area Network)であるが、WANに限定するものではない。

#### 【0023】

サイト201は、データセンタ等であり、1以上のノード100を含んで構成される。

10

20

30

40

50

## 【0024】

ノード100は、一般的なサーバ計算機の構成を備えてよい。ノード100は、例えば、プロセッサ211、メモリ212等を含む1以上のプロセッサパッケージ213、1以上のドライブ214、1以上のポート215を含んで構成される。各構成要素は、内部バス216を介して接続されている。

## 【0025】

プロセッサ211は、例えば、CPU(Central Processing Unit)であり、各種の処理を行う。

## 【0026】

メモリ212は、ノード100の機能を実現する上で必要な制御用の情報を格納したり、データを格納したりする。また、メモリ212は、例えば、プロセッサ211により実行されるプログラムを格納する。メモリ212は、揮発性のDRAM(Dynamic Random Access Memory)であってもよいし、不揮発のSCM(Storage Class Memory)であってもよいし、その他の記憶デバイスであってもよい。

10

## 【0027】

ドライブ214は、各種のデータ、プログラム等を記憶する。ドライブ214は、SAS(Serial Attached SCSI)またはSATA(Serial Advanced Technology Attachment)接続のHDD(Hard Disk Drive)やSSD(Solid State Drive)、NVMe(Non-Volatile Memory Express)接続のSSDの他、SCM等であってもよく、記憶デバイスの一例である。

20

## 【0028】

ポート215は、ネットワーク220に接続され、サイト201内の他のノード100と通信可能に接続されている。ネットワーク220は、例えば、LAN(Local Area Network)であるが、LANに限定するものではない。

## 【0029】

ストレージシステム200に係る物理構成は、上述の内容に限定されるものではない。例えば、ネットワーク220、202については、冗長化されていてもよい。また、例えば、ネットワーク220は、管理用のネットワークとストレージ用のネットワークとで分離してもよく、接続規格は、Ethernet(登録商標)、Infiniband、無線でもよく、接続トポロジも図2に示す構成に限定しない。

30

## 【0030】

なお、ホスト101は、ノード100と同じ構成要素を備えてもよく、ホスト101の物理構成については、その説明を省略する。

## 【0031】

図3は、ストレージシステム200に係る論理構成の一例を示す図である。ストレージシステム200では、ストレージ仮想化が行われ、複数の物理領域が仮想的に統合され、ストレージプール312として利用される。さらに、ストレージシステム200では、シンプロビジョニングにより、各ホスト101が現在利用している容量だけが割り当てられている。

40

## 【0032】

より具体的には、図3に示すように、ドライブ214は、データ、パリティ等を格納する物理的な領域であるデータ格納領域を有する。データ格納領域のうちの全部または一部の領域であり、連続した領域である論理ドライブ318は、ノード100を跨る複数の論理ドライブ318を組み合わせるパリティグループ317を構築する。

## 【0033】

パリティグループ317は、複数のノード100のドライブ214の論理ドライブ318から構成される。例えば、データ保護ポリシーが2D1Pである場合、異なるノード100のドライブ214から確保した3つの論理ドライブ318でパリティグループ317が

50

構成される。

【0034】

ここで、データ保護ポリシーとしては、例えば、EC ( Erasure Coding ) がある。なお、ECとしては、データローカリティを保持しない第1の方式と、データローカリティを保持する第2の方式 (例えば、国際公開第2016/52665号に記載の方式) とがあるが、ストレージシステム200には、何れの方式も適用可能である。なお、本実施の形態では、第2の方式を適用したケースを例に挙げて主に説明する。

【0035】

付言するならば、例えば、第1の方式の2D1PのECでは、ライト要求のデータを第1のデータと第2のデータとに分け、第1のデータを第1のノード100に格納し、第2のデータを第2のノード100に格納し、第1のデータおよび第2のデータで計算されたパリティを第3のノード100に格納することで冗長化が行われる。また、例えば、第2の方式の2D1PのECでは、ライト要求のデータを第1のデータと第2のデータとに分け、第1のデータおよび第2のデータを第1のノード100 (自ノード100) に格納し、第1のデータのパリティを第2のノード100に格納し、第2のデータのパリティを第3のノード100に格納することで冗長化が行われる。

10

【0036】

パリティグループ317からは、プールボリューム316が切り出される。プールボリューム316は、各ノード100のストレージプール312に容量を割り当てる単位である。1つのパリティグループ317から1つのプールボリューム316が切り出されてもよいし、複数のプールボリューム316が切り出されてよい。

20

【0037】

付言するならば、例えば、データ保護ポリシーが2D1Pである場合、データの格納領域として利用できるのは、パリティグループ317に割り当てられた論理ドライブ318の総量の2/3となり、パリティの格納領域として利用できるのは、パリティグループ317に割り当てられた論理ドライブ318の総量の1/3となる。つまり、プールボリューム316として切り出せる最大の容量は、データ保護ポリシーに応じて異なる。

【0038】

切り出されたプールボリューム316は、ストレージプール312にアタッチされる。ストレージプール312は、1以上のプールボリューム316を含んで構成される。ストレージプール312からは、アプリケーション301により利用される仮想ボリューム313が切り出される。つまり、ストレージプログラム311は、利用者の要求に応じた容量を、ドライブ214に割り当てず、仮想ボリューム313として割り当てる。

30

【0039】

ストレージプール312から仮想ボリューム313を切り出す際、複数のストレージプール312から仮想ボリューム313の領域をスライス314として部分的に切り出し、スライス314を束ねることで仮想ボリューム313を構築する。スライス314は、ストレージプール312に対して、仮想的に割り当てられた領域であり、仮想ボリューム313を作成した時点では、物理的な領域が割り当てられない。スライス314には、1以上のページ315が割り当てられる。例えば、ストレージプログラム311は、アプリケーション301からライト要求を受信した場合は、新規のライトであるときは、仮想ボリューム313のスライス314にページ315 (より詳細には、ページ315に紐づく論理ドライブ318の物理領域) を割り当てる。なお、ページ315には、プールボリューム316のページが対応付けられている。更新のライトであるときは、ストレージプログラム311は、割り当てたページ315に紐づく論理ドライブ318の物理領域を特定してデータを更新する。なお、ライト要求のデータ (または後述の中間データ) は、データの冗長化に係る他のノード100に転送されてパリティが更新される。

40

【0040】

仮想ボリューム313とアプリケーション301とは、フロントエンドパス320 (以降、単にパスとも記述する) で接続される。フロントエンドパス320の接続および設定

50

は、ストレージプログラム 3 1 1 とホスト 1 0 1 上で動作するパス設定プログラム 3 0 2 とにより制御される。なお、図 3 では、スライス 3 1 4 を第 1 のノード 1 0 0 「Node 0」から第 2 のノード 1 0 0 「Node 1」に移動した後、移動先（割当先）の第 2 のノード 1 0 0 「Node 1」にフロントエンドパス 3 2 0 が設定されていない例を示している。この場合、一度、第 1 のノード 1 0 0 「Node 0」を経由してスライス 3 1 4 の割当先の第 2 のノード 1 0 0 「Node 1」に I/O コマンドが転送されて処理される。ただし、後述するように、移動先の第 2 のノード 1 0 0 「Node 1」にフロントエンドパス 3 2 0 が設定され、最適化が行われることが好適である。

#### 【0041】

このように、ストレージプログラム 3 1 1 は、ドライブ 2 1 4 を共有のストレージプール 3 1 2 として管理し、仮想ボリューム 3 1 3 に書き込まれたデータ量に応じてドライブ 2 1 4 に容量を割り当てる。これにより、使用されないドライブ 2 1 4 の無駄をなくし、効率的な運用が行われる。

10

#### 【0042】

以下では、データを更新するにあたり、当該データは、ライト要求を受領したノード 1 0 0 のドライブ 2 1 4（ローカルドライブ）に格納される構成（データローカリティを維持し、リード時のネットワークオーバーヘッドを排除する構成）を例に挙げて主に説明する。

#### 【0043】

なお、データにアクセスするアプリケーション 3 0 1 は、ホスト 1 0 1 に設けられて動作するものであってもよいし、ストレージプログラム 3 1 1 と同一ノード 1 0 0 に設けられて動作するものであってもよいし、別のノード 1 0 0 に設けられて動作するものであってもよい。

20

#### 【0044】

図 4 は、メモリ 2 1 2 内の情報（ドライブ 2 1 4 からメモリ 2 1 2 に読み出される情報）の一例を示す図である。なお、制御情報テーブル 4 1 0、各種のプログラム（ストレージプログラム 3 1 1 等）は、実行中はメモリ 2 1 2 上に展開されるが、停電等に備えてドライブ 2 1 4 等の不揮発な領域に格納されている。

#### 【0045】

制御情報テーブル 4 1 0 には、クラスタ構成管理テーブル 4 1 1、データ保護セット管理テーブル 4 1 2、ストレージプール管理テーブル 4 1 3、ボリューム管理テーブル 4 1 4、モニタ情報管理テーブル 4 1 5、およびフロントエンドパス管理テーブル 4 1 6 が含まれる。各テーブルについては、図 5 ~ 図 1 0 を用いて後述する。

30

#### 【0046】

ストレージプログラム 3 1 1 は、リード処理プログラム 4 2 1、ライト処理プログラム 4 2 2、モニタ情報採取処理プログラム 4 2 3、リバランス要否判定処理プログラム 4 2 4、リソース割当決定処理プログラム 4 2 5、リソース移動処理プログラム 4 2 6、フロントエンドパス設定処理プログラム 4 2 7、およびクラスタ構成変更処理プログラム 4 2 8 を備える。なお、パス設定プログラム 3 0 2 は、フロントエンドパス設定処理プログラム 4 2 7 を備える。

#### 【0047】

ノード 1 0 0 の機能（リード処理プログラム 4 2 1、ライト処理プログラム 4 2 2、モニタ情報採取処理プログラム 4 2 3、リバランス要否判定処理プログラム 4 2 4、リソース割当決定処理プログラム 4 2 5、リソース移動処理プログラム 4 2 6、ストレージプログラム 3 1 1 のフロントエンドパス設定処理プログラム 4 2 7、クラスタ構成変更処理プログラム 4 2 8 等）は、例えば、プロセッサ 2 1 1 がドライブ 2 1 4 に格納されたプログラムをメモリ 2 1 2 に読み出して実行すること（ソフトウェア）により実現されてもよいし、専用の回路等のハードウェアにより実現されてもよいし、ソフトウェアとハードウェアとが組み合わされて実現されてもよい。また、ノード 1 0 0 の機能の一部は、ノード 1 0 0 と通信可能な他のコンピュータにより実現されてもよい。

40

#### 【0048】

50

ホスト101の機能(例えば、パス設定プログラム302のフロントエンドパス設定処理プログラム427)は、例えば、プロセッサ211がドライブ214に格納されたプログラムをメモリ212に読み出して実行すること(ソフトウェア)により実現されてもよいし、専用の回路等のハードウェアにより実現されてもよいし、ソフトウェアとハードウェアとが組み合わされて実現されてもよい。また、ホスト101の機能の一部は、ホスト101と通信可能な他のコンピュータにより実現されてもよい。

【0049】

図5は、クラスタ構成管理テーブル411の一例を示す図である。

【0050】

クラスタ構成管理テーブル411は、サイト201、ノード100、ドライブ214の構成を管理するための情報を格納する。 10

【0051】

クラスタ構成管理テーブル411は、サイト構成管理テーブル510、ノード構成管理テーブル520、およびドライブ構成管理テーブル530を含んで構成される。なお、ストレージシステム200は、サイト構成管理テーブル510を管理し、サイト201は、サイト201内の複数のノード構成管理テーブル520を管理し、ノード100は、ノード100内の複数のドライブ構成管理テーブル530を管理する。

【0052】

サイト構成管理テーブル510は、サイト201に係る構成(サイト201とノード100との関係等)を示す情報を格納する。より具体的には、サイト構成管理テーブル510は、サイトID511と、状態512と、ノードIDリスト513とが対応付けられた情報を格納する。 20

【0053】

サイトID511は、サイト201を識別可能な識別情報である。状態512は、サイト201の状態を示す状態情報(NORMAL、WARNING、FAILURE等)である。ノードIDリスト513は、サイト201に設けられるノード100を識別可能な識別情報のリストである。

【0054】

ノード構成管理テーブル520は、サイト201ごとに設けられ、サイト201に設けられるノード100に係る構成(ノード100とドライブ214との関係等)を示す情報を格納する。より具体的には、ノード構成管理テーブル520は、ノードID521と、状態522と、ドライブIDリスト523とが対応付けられた情報を格納する。 30

【0055】

ノードID521は、ノード100を識別可能な識別情報である。状態522は、ノード100の状態を示す状態情報(NORMAL、WARNING、FAILURE等)である。ドライブIDリスト523は、ノード100に設けられるドライブ214を識別可能な識別情報のリストである。

【0056】

ドライブ構成管理テーブル530は、ノード100ごとに設けられ、ノード100に設けられるドライブ214に係る構成を示す情報を格納する。より具体的には、ドライブ構成管理テーブル530は、ドライブID531と、状態532と、サイズ533とが対応付けられた情報を格納する。 40

【0057】

ドライブID531は、ドライブ214を識別可能な識別情報である。状態532は、ドライブ214の状態を示す状態情報(NORMAL、WARNING、FAILURE等)である。サイズ533は、ドライブ214の容量を示す情報(TB(テラバイト)、GB(ギガバイト)等)である。

【0058】

図6は、データ保護セット管理テーブル412の一例を示す図である。

【0059】

データ保護セット管理テーブル 4 1 2 は、論理ドライブ 3 1 8 を組み合わせて構成したパリティグループ 3 1 7 の構成を管理するための制御情報を格納する。

【 0 0 6 0 】

データ保護セット管理テーブル 4 1 2 は、プールボリューム管理テーブル 6 1 0、パリティグループ管理テーブル 6 2 0、論理ドライブ管理テーブル 6 3 0、およびストライプマッピングテーブル 6 4 0 を含んで構成される。

【 0 0 6 1 】

プールボリューム管理テーブル 6 1 0 は、パリティグループ 3 1 7 から切り出されたプールボリューム 3 1 6 に係る情報を格納する。より具体的には、プールボリューム管理テーブル 6 1 0 は、プールボリューム ID 6 1 1 と、サイズ 6 1 2 と、パリティグループ ID 6 1 3 と、論理ドライブ ID 6 1 4 とが対応付けられた情報を格納する。

10

【 0 0 6 2 】

プールボリューム ID 6 1 1 は、パリティグループ 3 1 7 から切り出されたプールボリューム 3 1 6 を識別可能な識別情報である。サイズ 6 1 2 は、プールボリューム 3 1 6 の容量を示す情報 ( T B ( テラバイト )、 G B ( ギガバイト ) 等 ) である。パリティグループ ID 6 1 3 は、プールボリューム 3 1 6 が属するパリティグループ 3 1 7 を識別可能な識別情報である。論理ドライブ ID 6 1 4 は、プールボリューム 3 1 6 に格納するデータ領域を提供する論理ドライブ 3 1 8 を識別可能な識別情報である。

【 0 0 6 3 】

パリティグループ管理テーブル 6 2 0 は、パリティグループ 3 1 7 に係る情報を格納する。より具体的には、パリティグループ管理テーブル 6 2 0 は、パリティグループ ID 6 2 1 と、冗長化ポリシー 6 2 2 と、論理ドライブ ID リスト 6 2 3 とが対応付けられた情報を格納する。

20

【 0 0 6 4 】

パリティグループ ID 6 2 1 は、パリティグループ 3 1 7 を識別可能な識別情報である。冗長化ポリシー 6 2 2 は、パリティグループ 3 1 7 の冗長化方法に関する設定である。論理ドライブ ID リスト 6 2 3 は、パリティグループ 3 1 7 に割り当てられた論理ドライブ 3 1 8 を識別可能な識別情報のリストである。

【 0 0 6 5 】

論理ドライブ管理テーブル 6 3 0 は、論理ドライブ 3 1 8 に係る情報 ( 開始オフセットからサイズ分だけドライブ 2 1 4 の物理領域を切り出して論理ドライブ 3 1 8 として管理するための情報 ) を格納する。より具体的には、論理ドライブ管理テーブル 6 3 0 は、論理ドライブ ID 6 3 1 と、開始オフセット 6 3 2 と、サイズ 6 3 3 と、ドライブ ID 6 3 4 とが対応付けられた情報を格納する。

30

【 0 0 6 6 】

論理ドライブ ID 6 3 1 は、論理ドライブ 3 1 8 を識別可能な識別情報である。開始オフセット 6 3 2 は、ドライブ 2 1 4 から論理ドライブ 3 1 8 を切り出すときの開始位置を示す情報である。サイズ 6 3 3 は、論理ドライブ 3 1 8 の容量を示す情報 ( ブロックの数 ) である。ここでブロックとは、ドライブ 2 1 4 へのアクセス単位を意味しており、典型的には、1 ブロックのサイズは、5 1 2 B y t e である。ただし、ブロックのサイズは、5 1 2 B y t e に限定せず、4 K B、8 K B 等でもよい。ドライブ ID 6 3 4 は、論理ドライブ 3 1 8 が切り出されている記憶資源を識別可能な識別情報 ( 論理ドライブ 3 1 8 がどのドライブ 2 1 4 から切り出されているかを示す情報 ) である。

40

【 0 0 6 7 】

ストライプマッピングテーブル 6 4 0 は、パリティグループ 3 1 7 に係る情報 ( データおよびパリティの格納先アドレスを計算するための情報 ) を格納する。一例として、ストライプマッピングテーブル 6 4 0 が、 E C ( 2 D 1 P ) のストライプマッピングテーブル 6 4 1、および M i r r o r ( 2 - R e p l i c a t i o n ) のストライプマッピングテーブル 6 4 2 の情報を格納するケースについて説明する。2 D 1 P とは、2 つのデータの組み合わせで 1 つのパリティを算出し、データを保護することを意味する。

50

## 【 0 0 6 8 】

ストライプマッピングテーブル 6 4 1 , 6 4 2 は、あるデータ領域の物理 L B A ( L o g i c a l B l o c k A d d r e s s ) に対して、パリティ領域の物理 L B A ( 冗長化先ノード ) を特定するために使用する。

## 【 0 0 6 9 】

ストライプマッピングテーブル 6 4 1 , 6 4 2 は、表、配列形式等で格納されており、横軸の要素としてノード I D に対応する情報を保持し、縦軸の要素としてアドレスに対応する情報を保持している。横軸の情報は、直接的にノード I D の情報を格納していてもよいし、ストライプマッピングテーブル 6 4 1 , 6 4 2 の横軸の I D とノード I D とを対応づける別のテーブルを介して管理されていてもよい。縦軸の情報は、直接的に L B A の情報を格納している必要はなく、例えば L B A から縦軸の I D へは以下のように変換することができる。

10

## 【 0 0 7 0 】

$$\text{RowID} = \text{LBA} \bmod \text{Rowmax}$$

( Rowmax は、ストライプマッピングテーブル 6 4 1 , 6 4 2 では「 6 」となる )

## 【 0 0 7 1 】

図 7 は、ストレージプール管理テーブル 4 1 3 の一例を示す図である。

## 【 0 0 7 2 】

ストレージプール管理テーブル 4 1 3 は、ストレージプール 3 1 2 の構成を管理するための制御情報を格納する。ストレージプール管理テーブル 4 1 3 は、ストレージプール情報テーブル 7 1 0 を含んで構成される。

20

## 【 0 0 7 3 】

ストレージプール情報テーブル 7 1 0 は、ストレージプール 3 1 2 に係る情報を格納する。より具体的には、ストレージプール情報テーブル 7 1 0 は、ストレージプール I D 7 1 1 と、合計容量 7 1 2 と、使用容量 7 1 3 と、ノード I D 7 1 4 と、プールボリューム I D リスト 7 1 5 とが対応付けられた情報を格納する。

## 【 0 0 7 4 】

ストレージプール I D 7 1 1 は、ストレージプール 3 1 2 を識別可能な識別情報である。合計容量 7 1 2 は、ストレージプール 3 1 2 に割り当てられた合計容量を示す情報 ( T B ( テラバイト ) 、 G B ( ギガバイト ) 等 ) である。使用容量 7 1 3 は、ストレージプール 3 1 2 で使用している容量を示す情報 ( T B ( テラバイト ) 、 G B ( ギガバイト ) 等 ) である。ノード I D 7 1 4 は、ストレージプール 3 1 2 を提供するノード 1 0 0 を識別可能な識別情報である。プールボリューム I D リスト 7 1 5 は、ストレージプール 3 1 2 に割り当てられたプールボリューム 3 1 6 を識別可能な識別情報のリストである。

30

## 【 0 0 7 5 】

図 8 は、ボリューム管理テーブル 4 1 4 の一例を示す図である。

## 【 0 0 7 6 】

ボリューム管理テーブル 4 1 4 は、仮想ボリューム 3 1 3 の構成情報と、ノード 1 0 0 間に割り当てられたスライス 3 1 4 の構成情報と、シンプロビジョニング機能のための制御情報とを格納する。

40

## 【 0 0 7 7 】

ボリューム管理テーブル 4 1 4 は、仮想ボリューム管理テーブル 8 1 0 、スライス管理テーブル 8 2 0 、およびページマッピングテーブル 8 3 0 を含んで構成される。

## 【 0 0 7 8 】

仮想ボリューム管理テーブル 8 1 0 は、仮想ボリューム 3 1 3 に係る情報 ( 仮想ボリューム 3 1 3 と仮想ボリューム 3 1 3 に割り当てられたスライス 3 1 4 との対応関係を示す情報等 ) を格納する。より具体的には、仮想ボリューム管理テーブル 8 1 0 は、仮想ボリューム I D 8 1 1 と、サイズ 8 1 2 と、スライス I D リスト 8 1 3 と、最大分散度 8 1 4 とが対応付けられた情報を格納する。

## 【 0 0 7 9 】

50

仮想ボリュームID 8 1 1は、仮想ボリューム3 1 3を識別可能な識別情報である。サイズ8 1 2は、仮想ボリューム3 1 3の容量を示す情報（TB（テラバイト）、GB（ギガバイト）等）である。スライスIDリスト8 1 3は、仮想ボリューム3 1 3に割り当てられたスライス3 1 4を識別可能な識別情報のリストである。最大分散度8 1 4は、仮想ボリューム3 1 3に割り当てるスライス3 1 4を分散させるノード数の最大値である。これを超えた数のノード1 0 0には、スライス3 1 4が割り当てられないように制御される。

【0 0 8 0】

スライス管理テーブル8 2 0は、仮想ボリューム3 1 3に割り当てたスライス3 1 4に係る情報（スライス3 1 4とスライス3 1 4に対応するストレージプール3 1 2との対応関係を示す情報等）を格納する。より具体的には、スライス管理テーブル8 2 0は、スライスID 8 2 1と、サイズ8 2 2と、ストレージプールID 8 2 3と、状態8 2 4とが対応付けられた情報を格納する。

10

【0 0 8 1】

スライスID 8 2 1は、スライス3 1 4を識別可能な識別情報である。サイズ8 2 2は、スライス3 1 4の容量を示す情報（TB（テラバイト）、GB（ギガバイト）、Logical Block数等）である。ストレージプールID 8 2 3は、スライス3 1 4に対応するストレージプール3 1 2を識別可能な識別情報である。2つのストレージプール3 1 2間でスライス3 1 4を移動中の場合、ストレージプールID 8 2 3は、移動前後のストレージプール3 1 2を識別可能な識別情報を格納する。状態8 2 4は、スライス3 1 4の状態を示す情報である。状態8 2 4には、正常（Normal）、障害状態（Failure）、2つのストレージプール3 1 2間を移動中（Migrating）といった状態がある。

20

【0 0 8 2】

ページマッピングテーブル8 3 0は、仮想ボリューム3 1 3に割り当てたページ3 1 5に係る情報（ページ3 1 5とプールボリューム3 1 6との対応関係を示す情報等）を格納する。より具体的には、ページマッピングテーブル8 3 0は、ページID 8 3 1と、仮想ボリュームID 8 3 2と、仮想ボリュームLBA 8 3 3と、サイズ8 3 4と、プールボリュームID 8 3 5と、プールボリュームLBA 8 3 6とが対応付けられた情報を格納する。

【0 0 8 3】

ページID 8 3 1は、ページ3 1 5を識別可能な識別情報である。仮想ボリュームID 8 3 2は、ページ3 1 5が割り当てられている仮想ボリューム3 1 3を識別可能な識別情報である。仮想ボリュームLBA 8 3 3は、仮想ボリューム3 1 3におけるページ3 1 5の位置を特定可能な情報であり、例えば、仮想ボリューム3 1 3の最初のページ3 1 5から何番目であるかを示す情報である。なお、ページ3 1 5は、ストレージプログラム3 1 1が仮想ボリューム3 1 3にアクセスする単位である。サイズ8 3 4は、ページ3 1 5の容量を示す情報（TB（テラバイト）、GB（ギガバイト）、Logical Block数等）である。プールボリュームID 8 3 5は、ページ3 1 5に対応するプールボリューム3 1 6を識別可能な識別情報である。プールボリュームLBA 8 3 6は、ストレージプール3 1 2におけるプールボリューム3 1 6の位置を特定可能な情報であり、例えば、ストレージプール3 1 2の最初のプールボリューム3 1 6から何番目であるかを示す情報である。

30

40

【0 0 8 4】

なお、サイズ8 3 4は、全てのページ3 1 5で同じであってもよいし、ページ3 1 5ごとに異なってもよい。

【0 0 8 5】

付言するならば、ストレージプログラム3 1 1は、仮想ボリューム3 1 3のアドレスからストレージプール3 1 2のアドレスへの変換を行う際にページマッピングテーブル8 3 0を参照する。また、ストレージプログラム3 1 1は、新規ライトを受領する度に、ページ3 1 5の割当て（ページマッピングテーブル8 3 0へのレコードの追加）を行う。

【0 0 8 6】

50

図 9 は、モニタ情報管理テーブル 4 1 5 の一例を示す図である。

【 0 0 8 7 】

モニタ情報管理テーブル 4 1 5 は、ノード 1 0 0 で動作するプロセスのプロセッサ 2 1 1、ドライブ 2 1 4、およびポート 2 1 5 の使用量と、仮想ボリューム 3 1 3 のスライス 3 1 4、およびフロントエンドパス 3 2 0 に対するアクセス頻度とを管理するためのするための制御情報を格納する。

【 0 0 8 8 】

モニタ情報管理テーブル 4 1 5 は、プロセッサモニタ情報管理テーブル 9 1 0、ドライブモニタ情報管理テーブル 9 2 0、ネットワークモニタ情報管理テーブル 9 3 0、スライスモニタ情報管理テーブル 9 4 0、フロントエンドパスモニタ情報管理テーブル 9 5 0 を含んで構成される。

10

【 0 0 8 9 】

プロセッサモニタ情報管理テーブル 9 1 0 は、プロセッサ 2 1 1 に係る情報（プロセスとプロセッサ 2 1 1 の使用量との関係を示す情報）を格納する。より具体的には、ノード ID 9 1 1 と、プロセッサ ID 9 1 2 と、プロセス ID 9 1 3 と、プロセス名 9 1 4 と、使用率 9 1 5 とが対応付けられた情報を格納する。

【 0 0 9 0 】

ノード ID 9 1 1 は、ノード 1 0 0 を識別可能な識別情報である。プロセッサ ID 9 1 2 は、ノード 1 0 0 内に複数のプロセッサコアが搭載されている場合にプロセッサコアを識別可能な識別情報である。プロセス ID 9 1 3 は、ノード 1 0 0 で動作するプログラムを識別可能な識別情報である。プロセス名 9 1 4 は、ノード 1 0 0 で動作するプログラムを識別可能な文字列情報である。使用率 9 1 5 は、ノード 1 0 0 で動作するプログラムが動作するプロセッサコアの占有率を示す。例えば、ストレージプログラム 3 1 1 の使用率が 5 0 % である場合、ストレージプログラム 3 1 1 が動作するプロセッサコア動作周波数の半分を占有していることを意味する。

20

【 0 0 9 1 】

ドライブモニタ情報管理テーブル 9 2 0 は、ドライブ 2 1 4 に係る情報（ドライブ 2 1 4 の使用量の関係を示す情報）を格納する。より具体的には、ドライブ ID 9 2 1 と、リード I O P S 9 2 2 と、ライト I O P S 9 2 3 と、リード転送量 9 2 4 と、ライト転送量 9 2 5 と、使用率 9 2 6 とが対応付けられた情報を格納する。

30

【 0 0 9 2 】

ドライブ ID 9 2 1 は、ドライブ 2 1 4 を識別可能な識別情報である。リード I O P S 9 2 2 は、当該ドライブ 2 1 4 に対してのリードコマンドの秒間あたりの処理数である。ライト I O P S 9 2 3 は、当該ドライブ 2 1 4 に対してのライトコマンドの秒間あたりの処理数である。リード転送量 9 2 4 は、当該ドライブ 2 1 4 に対してのリードコマンドの秒間あたりのデータ転送量である。ライト転送量 9 2 5 は、当該ドライブ 2 1 4 に対してのライトコマンドの秒間あたりのデータ転送量である。使用率 9 2 6 は、当該ドライブ 2 1 4 の負荷度合いを示し、1 0 0 % となった場合、当該ドライブ 2 1 4 は、それ以上 I / O を処理できず、ドライブ 2 1 4 が受領した I / O 要求は待たされることになる。

【 0 0 9 3 】

ネットワークモニタ情報管理テーブル 9 3 0 は、ネットワーク 2 2 0 に接続されているポート 2 1 5 に係る情報（ポート 2 1 5 の使用量の関係を示す情報）を格納する。より具体的には、ノード ID 9 3 1 と、NIC ( Network Interface Card ) ID 9 3 2 と、送信転送量 9 3 3 と、受信転送量 9 3 4 と、最大転送量 9 3 5 とが対応付けられた情報を格納する。

40

【 0 0 9 4 】

ノード ID 9 3 1 は、ノード 1 0 0 を識別可能な識別情報である。NIC ID 9 3 2 は、ノード 1 0 0 内に複数の NIC ( ポート 2 1 5 ) が搭載されている場合に NIC を識別可能な識別情報である。なお、本実施の形態では、NIC が 1 つのポート 2 1 5 を備える場合を例に挙げて説明する。送信転送量 9 3 3 は、当該 NIC に対しての送信処理の秒

50

間あたりの転送量である。受信転送量 9 3 4 は、当該 N I C に対しての受信処理の秒間あたりの転送量である。最大転送量 9 3 5 は、当該 N I C で処理可能な送受信の秒間あたりの最大転送量である。

【 0 0 9 5 】

スライスモニタ情報管理テーブル 9 4 0 は、スライス 3 1 4 へのアクセス頻度の情報を格納する。より具体的には、スライスモニタ情報管理テーブル 9 4 0 は、スライス I D 9 4 1 と、リードカウンタ 9 4 2 と、ライトカウンタ 9 4 3 と、リード転送量 9 4 4 と、ライト転送量 9 4 5 と、モニタ開始時刻 9 4 6 とが対応付けられた情報を格納する。

【 0 0 9 6 】

スライス I D 9 4 1 は、スライス 3 1 4 を識別可能な識別情報である。リードカウンタ 9 4 2 は、当該スライス 3 1 4 をリードした回数を管理するための情報である。ライトカウンタ 9 4 3 は、当該スライス 3 1 4 に対してライトした回数を管理するための情報である。リード転送量 9 4 4 は、当該スライス 3 1 4 をリードした転送量を管理するための情報である。ライト転送量 9 4 5 は、当該スライス 3 1 4 に対してライトした転送量を管理するための情報である。モニタ開始時刻 9 4 6 は、当該スライス 3 1 4 に対するアクセスの監視が開始された時間を示す情報である。

10

【 0 0 9 7 】

フロントエンドパスモニタ情報管理テーブル 9 5 0 は、フロントエンドパス 3 2 0 へのアクセス頻度の情報を格納する。より具体的には、フロントエンドパスモニタ情報管理テーブル 9 5 0 は、パス I D 9 5 1 と、リード I O P S 9 5 2 と、ライト I O P S 9 5 3 と、リード転送量 9 5 4 と、ライト転送量 9 5 5 とが対応付けられた情報を格納する。

20

【 0 0 9 8 】

パス I D 9 5 1 は、フロントエンドパス 3 2 0 を識別可能な識別情報である。リード I O P S 9 5 2 は、当該フロントエンドパス 3 2 0 に対してのリードコマンドの秒間あたりの処理数である。ライト I O P S 9 5 3 は、当該フロントエンドパス 3 2 0 に対してのライトコマンドの秒間あたりの処理数である。リード転送量 9 5 4 は、当該フロントエンドパス 3 2 0 に対してのリードコマンドの秒間あたりのデータ転送量である。ライト転送量 9 5 5 は、当該フロントエンドパス 3 2 0 に対してのライトコマンドの秒間あたりのデータ転送量である。

【 0 0 9 9 】

図 1 0 は、フロントエンドパス管理テーブル 4 1 6 の一例を示す図である。

30

【 0 1 0 0 】

フロントエンドパス管理テーブル 4 1 6 は、フロントエンドパス 3 2 0 の構成を管理するための制御情報を格納する。フロントエンドパス管理テーブル 4 1 6 は、フロントエンドパス情報テーブル 1 0 1 0 を含んで構成される。

【 0 1 0 1 】

フロントエンドパス情報テーブル 1 0 1 0 は、フロントエンドパス 3 2 0 に係る情報を格納する。より具体的には、フロントエンドパス情報テーブル 1 0 1 0 は、パス I D 1 0 1 1 と、仮想ボリューム I D 1 0 1 2 と、I n i t i a t o r I D 1 0 1 3 と、A L U A 設定 1 0 1 4 と、接続ノード I D 1 0 1 5 とが対応付けられた情報を格納する。

40

【 0 1 0 2 】

パス I D 1 0 1 1 は、フロントエンドパス 3 2 0 を識別可能な識別情報である。仮想ボリューム I D 1 0 1 2 は、フロントエンドパス 3 2 0 が割り当てられた仮想ボリューム 3 1 3 を識別可能な識別情報である。I n i t i a t o r I D 1 0 1 3 は、フロントエンドパス 3 2 0 の接続先であるホスト 1 0 1 を識別可能な識別情報である。A L U A 設定 1 0 1 4 は、ストレージシステム 2 0 0 にとって、対応するフロントエンドパス 3 2 0 が好適であるかどうかの設定を示す情報である。A L U A 設定 1 0 1 4 に基づく情報を、ホスト 1 0 1 に通知することで、ホスト 1 0 1 は、好適なパスへ I / O リクエストを発行することができ、ストレージシステム 2 0 0 の処理効率を向上させることができる。接続ノード I D 1 0 1 5 は、フロントエンドパス 3 2 0 を有するノード I D を識別可能な識別情報で

50

ある。

【0103】

図11は、リード処理に係るフローチャートの一例を示す図である。リード処理では、アプリケーション301からのデータのリード処理要求を受けて、自ノード100のドライブ214からデータが読み出される。なお、リード処理要求では、リード先（例えば、LUN(Logical Unit Number)のような仮想ボリュームID、LBAのようなアドレス等)が指定されている。アクセス先(ドライブ214等)が障害状態である場合、冗長データからリード対象のデータが修復されて応答される。以下、詳細について説明する。

【0104】

ステップS1101では、リード処理プログラム421は、アクセス先LBAからスライスIDを計算する。より具体的には、リード処理プログラム421は、仮想ボリューム管理テーブル810を参照し、スライスIDリスト813の先頭のスライスから仮想ボリューム313のLBAが連続的に割り当てられているとき、リストを順に辿ることでアクセス先のLBAに該当するスライスIDを取得する。

【0105】

ステップS1102では、リード処理プログラム421は、ステップS1101で取得したスライスID(対象スライス)が自身のノード100(自系ノード)に割り当てられているか否かを判定する。より具体的には、リード処理プログラム421は、スライス管理テーブル820を参照し、該当するスライスIDに対応するストレージプールIDを取得する。次に、リード処理プログラム421は、ストレージプール情報テーブル710を参照し、取得したストレージプールIDに対応するノードIDを取得する。リード処理プログラム421は、取得したノードIDと自系ノードのノードIDとを比較し、同じノードIDである場合、対象スライス(アクセス先のスライス314)が自系ノードに割り当たっていると判定する。リード処理プログラム421は、取得したノードIDと自系ノードのノードIDとが異なるノードIDである場合、対象スライスは、他のノード100(他系ノード)に割り当たっていると判定する。リード処理プログラム421は、対象スライスが自系ノードに割り当てられていると判定した場合、ステップS1105に処理を移し、対象スライスが自系ノードに割り当てられていないと判定した場合、ステップS1103に処理を移す。

【0106】

ステップS1103では、リード処理プログラム421は、対象スライスを割り当てた先の他系ノードにリード処理要求を転送する。

【0107】

ステップS1104では、リード処理プログラム421は、ステップS1103にて転送したリード処理要求の実行結果を待ち受け、実行結果を受信し、ステップS1111に処理を移す。

【0108】

ステップS1105では、リード処理プログラム421は、アクセス先の領域に関しての排他制御を取得する。

【0109】

ステップS1106では、リード処理プログラム421は、リード処理要求のデータについて、ストレージプール312にページ315が未割当てであるか否かを判定する。リード処理プログラム421は、未割当てであると判定した場合、ステップS1107に処理を移し、未割当てでないと判定した場合、ステップS1108に処理を移す。

【0110】

ステップS1107では、リード処理プログラム421は、データがないことを示す0データを生成し、ステップS1110に処理を移す。

【0111】

ステップS1108では、リード処理プログラム421は、アクセス先のアドレス(割

10

20

30

40

50

当先アドレス)を取得する。

【0112】

ステップS1109では、リード処理プログラム421は、自系ノードのドライブ214(ローカルドライブ)からデータを読み出す。

【0113】

ステップS1110では、リード処理プログラム421は、取得した排他制御を解放する。

【0114】

ステップS1111では、リード処理プログラム421は、ホスト101にリード処理結果を応答する。

【0115】

ステップS1112では、リード処理プログラム421は、モニタ情報採取処理を実行する。なお、モニタ情報採取処理については、図13Bを用いて後述する。

【0116】

図12は、ライト処理に係るフローチャートの一例を示す図である。ライト処理では、アプリケーション301からのライト処理要求を受けて、自系ノードのドライブ214にデータが書き込まれ、さらに他系ノードのドライブ214に冗長データ(パリティ)が書き込まれる。なお、ライト処理要求では、ライト先(例えば、LUNのような仮想ボリュームID、LBAのようなアドレス等)が指定されている。以下、詳細について説明する。

【0117】

ステップS1201では、ライト処理プログラム422は、アクセス先LBAからスライスIDを計算する。より具体的には、ライト処理プログラム422は、仮想ボリューム管理テーブル810を参照し、スライスIDリスト813の先頭のスライスから仮想ボリューム313のLBAが連続的に割り当てられているとき、リストを順に辿ることでアクセス先のLBAに該当するスライスIDを取得する。

【0118】

ステップS1202では、ライト処理プログラム422は、ステップS1201で取得したスライスID(対象スライス)が自系ノードに割り当てられているか否かを判定する。なお、ライト処理プログラム422は、リード処理プログラム421で説明した方法と同様に判定する。ライト処理プログラム422は、対象スライスが自系ノードに割り当てられていると判定した場合、ステップS1205に処理を移し、対象スライスが自系ノードに割り当てられていないと判定した場合、ステップS1203に処理を移す。

【0119】

ステップS1203では、ライト処理プログラム422は、対象スライスを割り当てた先の他系ノードにライト処理要求を転送する。

【0120】

ステップS1204では、ライト処理プログラム422は、ステップS1203にて転送したライト処理要求の実行結果を待ち受け、実行結果を受信し、ステップS1224に処理を移す。

【0121】

ステップS1205では、ライト処理プログラム422は、アクセス先の領域に関しての排他制御を取得する。

【0122】

ステップS1206では、ライト処理プログラム422は、対象スライスの状態が移動中であるか否かを判定する。より具体的には、ライト処理プログラム422は、スライス管理テーブル820を参照し、アクセス先となるスライス314のスライスIDに対応する状態が、Migratingである場合、移動中であると判定し、Migratingでない場合、移動中でないと判定する。ライト処理プログラム422は、移動中であると判定した場合、ステップS1207に処理を移し、移動中でないと判定した場合、ステップS1209に処理を移す。

10

20

30

40

50

## 【 0 1 2 3 】

ステップ S 1 2 0 7 では、ライト処理プログラム 4 2 2 は、対象スライスの移動先のノード 1 0 0 ( 移動先ノード ) にライト処理要求を転送する。

## 【 0 1 2 4 】

ステップ S 1 2 0 8 では、ライト処理プログラム 4 2 2 は、ステップ S 1 2 0 7 にて転送したライト処理要求の実行結果を待ち受け、実行結果を受信し、ステップ S 1 2 0 9 に処理を移す。

## 【 0 1 2 5 】

ステップ S 1 2 0 9 では、ライト処理プログラム 4 2 2 は、ライト処理要求のデータについて、ストレージプール 3 1 2 にページ 3 1 5 が未割当てであるか否かを判定する。ライト処理プログラム 4 2 2 は、未割当てであると判定した場合、ステップ S 1 2 1 0 に処理を移し、未割当てでないと判定した場合、ステップ S 1 2 1 1 に処理を移す。

10

## 【 0 1 2 6 】

ステップ S 1 2 1 0 では、ライト処理プログラム 4 2 2 は、自系ノードのドライブ 2 1 4 の論理ドライブ 3 1 8 が関連付けられているプールボリューム 3 1 6 ( 自系プールボリューム ) にページ 3 1 5 を割り当てる。

## 【 0 1 2 7 】

ステップ S 1 2 1 1 では、ライト処理プログラム 4 2 2 は、アクセス先のアドレス ( 割当先アドレス ) を取得する。

## 【 0 1 2 8 】

ステップ S 1 2 1 2 では、ライト処理プログラム 4 2 2 は、書込み前のデータ ( 旧データ ) を読み込む。ライト処理プログラム 4 2 2 は、読み込み先のドライブ 2 1 4 またはノード 1 0 0 が障害状態である場合、リード処理プログラム 4 2 1 で説明したようにパリティから読み込み対象のデータを復元して、旧データを読み込む。

20

## 【 0 1 2 9 】

ステップ S 1 2 1 3 では、ライト処理プログラム 4 2 2 は、中間データを生成する。中間データは、データを部分的に更新するとき作成する一時的なデータであり、新旧の差分を示すデータである。例えば、旧データのストライプが「 A 1 - A 2 - A P 」である場合、中間データは、次のように求められる。

## 【 0 1 3 0 】

$$A P ( 旧パリティ ) = A 1 ( 旧データ ) \text{ XOR } A 2 ( 旧データ )$$

$$A 1 ( 新データ ) \text{ XOR } A 1 ( 旧データ ) = M ( 中間データ )$$

なお、新パリティについては、次のように求められる。

$$A P ( 旧パリティ ) \text{ XOR } M ( 中間データ ) = A P ( 新パリティ )$$

30

## 【 0 1 3 1 】

ステップ S 1 2 1 4 では、ライト処理プログラム 4 2 2 は、冗長化先のノード 1 0 0 に中間データ ( パリティ更新要求 ) を送信する。なお、ライト処理プログラム 4 2 2 は、冗長度に応じて ( 冗長度が 2 以上である場合、 2 以上のノード 1 0 0 に ) 中間データを転送する。

## 【 0 1 3 2 】

ステップ S 1 2 1 5 では、ライト処理プログラム 4 2 2 は、自系ノードのドライブ 2 1 4 に新データを書き込む。

40

## 【 0 1 3 3 】

ステップ S 1 2 1 6 では、冗長化先のノード 1 0 0 のライト処理プログラム 4 2 2 は、中間データを受信する。

## 【 0 1 3 4 】

ステップ S 1 2 1 7 では、冗長化先のノード 1 0 0 のライト処理プログラム 4 2 2 は、排他制御を取得する。

## 【 0 1 3 5 】

ステップ S 1 2 1 8 では、冗長化先のノード 1 0 0 のライト処理プログラム 4 2 2 は、

50

自系ノードのドライブ 2 1 4 から旧パリティを読み出す。

【 0 1 3 6 】

ステップ S 1 2 1 9 では、冗長化先のノード 1 0 0 のライト処理プログラム 4 2 2 は、中間データと旧パリティとから新パリティを計算する。

【 0 1 3 7 】

ステップ S 1 2 2 0 では、冗長化先のノード 1 0 0 のライト処理プログラム 4 2 2 は、自系ノードのドライブ 2 1 4 に新パリティを書き込む。

【 0 1 3 8 】

ステップ S 1 2 2 1 では、冗長化先のノード 1 0 0 のライト処理プログラム 4 2 2 は、取得した排他制御を解放し、中間データを転送してきたノード 1 0 0 にパリティ更新結果を応答する。

10

【 0 1 3 9 】

ステップ S 1 2 2 2 では、ライト処理プログラム 4 2 2 は、冗長化先のノード 1 0 0 から書き込み応答を受信する。

【 0 1 4 0 】

ステップ S 1 2 2 3 では、ライト処理プログラム 4 2 2 は、取得した排他制御を解放する。

【 0 1 4 1 】

ステップ S 1 2 2 4 では、ライト処理プログラム 4 2 2 は、ホスト 1 0 1 にライト処理結果を応答する。

20

【 0 1 4 2 】

ステップ S 1 2 2 5 では、ライト処理プログラム 4 2 2 は、モニタ情報採取処理を実行する。なお、モニタ情報採取処理については、図 1 3 B を用いて後述する。

【 0 1 4 3 】

図 1 3 A は、プロセッサ 2 1 1 と、ドライブ 2 1 4 と、およびネットワーク 2 2 0 に関するモニタ情報採取処理に係るフローチャートの一例を示す図である。

【 0 1 4 4 】

ステップ S 1 3 0 1 では、モニタ情報採取処理プログラム 4 2 3 は、プロセッサ 2 1 1 のモニタ情報をテーブルに登録する。より具体的には、プロセッサモニタ情報管理テーブル 9 1 0 にあるように、モニタ情報採取処理プログラム 4 2 3 は、ノード 1 0 0 ごと、プロセスごと、プロセッサコアごとにプロセッサの使用率の情報を収集し、当該テーブルの情報を更新する。図示はしていないが、使用率以外の情報（ I D L E 、 I O W A I T 、 ハイパーバイザ上の仮想マシンとして実行していれば、 S T E A L 等）が取得され、テーブルに加えられてもよい。

30

【 0 1 4 5 】

ステップ S 1 3 0 2 では、モニタ情報採取処理プログラム 4 2 3 は、ドライブ 2 1 4 のモニタ情報をテーブルに登録する。より具体的には、ドライブモニタ情報管理テーブル 9 2 0 にあるように、モニタ情報採取処理プログラム 4 2 3 は、ドライブ 2 1 4 ごとに、リード I O P S と、ライト I O P S と、リード転送量と、ライト転送量との情報を収集し、当該テーブルの情報を更新する。図示はしていないが、これら以外の情報（リード応答時間、ライト応答時間、キューサイズ等）が取得され、テーブルに加えられてもよい。

40

【 0 1 4 6 】

ステップ S 1 3 0 3 では、モニタ情報採取処理プログラム 4 2 3 は、ネットワーク 2 2 0 ( N I C ) のモニタ情報をテーブルに登録する。より具体的には、ネットワークモニタ情報管理テーブル 9 3 0 にあるように、モニタ情報採取処理プログラム 4 2 3 は、ノード 1 0 0 ごとの N I C ごとに、送信転送量と、受信転送量と、最大転送量との情報を収集し、当該テーブルの情報を更新する。図示はしていないが、これら以外の情報（パケットドロップ数、再送パケット数等）が取得され、テーブルに加えられてもよい。

【 0 1 4 7 】

ステップ S 1 3 0 4 では、モニタ情報採取処理プログラム 4 2 3 は、一定時間処理を停

50

止し、その後、ステップ S 1 3 0 1 に処理を移す。つまり、図 1 3 A のモニタ情報採取処理は、周期的に実行される。

【 0 1 4 8 】

図 1 3 B は、スライス 3 1 4 のアクセス頻度、およびフロントエンドパス 3 2 0 のアクセス頻度に関するモニタ情報採取処理に係るフローチャートの一例を示す図である。

【 0 1 4 9 】

ステップ S 1 3 1 1 では、モニタ情報採取処理プログラム 4 2 3 は、I / O を受信したフロントエンドパス 3 2 0 のモニタ情報とアクセス先のスライス 3 1 4 のモニタ情報とを取得する。より具体的には、モニタ情報採取処理プログラム 4 2 3 は、アクセスを受信したフロントエンドパス 3 2 0 に該当するフロントエンドパスモニタ情報管理テーブル 9 5 0 のレコード、およびアクセス先のスライス 3 1 4 に該当するスライスモニタ情報管理テーブル 9 4 0 のレコードを取得する。

10

【 0 1 5 0 】

ステップ S 1 3 1 2 では、モニタ情報採取処理プログラム 4 2 3 は、受信した I / O タイプは、リードであるか否かを判定する。モニタ情報採取処理プログラム 4 2 3 は、リードであると判定した場合、ステップ S 1 3 1 3 に処理を移し、リードでない（ライトである）と判定した場合、ステップ S 1 3 1 5 に処理を移す。

【 0 1 5 1 】

ステップ S 1 3 1 3 では、モニタ情報採取処理プログラム 4 2 3 は、ステップ S 1 3 1 1 で取得したレコードの現行リードカウンタに受領した I / O のカウンタを加算する。ここで、I O P S は、秒間当たりの処理量であるので、モニタ情報採取処理プログラム 4 2 3 は、1 秒毎にカウンタ値を確定させる、つまり、1 秒経過したときに発生したカウンタ値を計算することで I O P S を求め、ステップ S 1 3 1 1 で取得したレコードのリード I O P S に設定する。

20

【 0 1 5 2 】

ステップ S 1 3 1 4 では、モニタ情報採取処理プログラム 4 2 3 は、ステップ S 1 3 1 1 で取得したレコードの現行リード転送量に受領した I / O の転送量を加算する。

【 0 1 5 3 】

ステップ S 1 3 1 5 では、モニタ情報採取処理プログラム 4 2 3 は、ステップ S 1 3 1 1 で取得したレコードの現行ライトカウンタに受領した I / O のカウンタを加算する。また、モニタ情報採取処理プログラム 4 2 3 は、1 秒毎にカウンタ値を確定させる、つまり、1 秒経過したときに発生したカウンタ値を計算することで I O P S を求め、ステップ S 1 3 1 1 で取得したレコードのライト I O P S に設定する。

30

【 0 1 5 4 】

ステップ S 1 3 1 6 では、モニタ情報採取処理プログラム 4 2 3 は、ステップ S 1 3 1 1 で取得したレコードの現行ライト転送量に受領した I / O の転送量を加算する。

【 0 1 5 5 】

図 1 4 は、リバランス要否判定処理に係るフローチャートの一例を示す図である。本処理は、ストレージシステム 2 0 0 により周期的に実行されてもよいし、ユーザ（手動）により任意の契機に実行されてもよいし、リード処理またはライト処理の完了後に実行されてもよいし、後述のクラスタ構成変更処理の実行後に実行されてもよい。

40

【 0 1 5 6 】

ステップ S 1 4 0 1 では、リバランス要否判定処理プログラム 4 2 4 は、ストレージプール 3 1 2 の使用率が上限閾値以上のノード 1 0 0 が存在するか否かを判定する。リバランス要否判定処理プログラム 4 2 4 は、ストレージプール 3 1 2 の使用率が上限閾値以上のノード 1 0 0 が存在すると判定した場合、ステップ S 1 4 0 5 に処理を移し、存在しないと判定した場合、ステップ S 1 4 0 2 に処理を移す。

【 0 1 5 7 】

ステップ S 1 4 0 2 では、リバランス要否判定処理プログラム 4 2 4 は、プロセッサ 2 1 1 の使用率が上限閾値以上のノード 1 0 0 が存在するか否かを判定する。リバランス要

50

否判定処理プログラム424は、プロセッサ211の使用率が上限閾値以上のノード100が存在すると判定した場合、ステップS1405に処理を移し、存在しないと判定した場合、ステップS1403に処理を移す。

【0158】

ステップS1403では、リバランス要否判定処理プログラム424は、ドライブ214の使用率が上限閾値以上のノード100が存在するか否かを判定する。リバランス要否判定処理プログラム424は、ドライブ214の使用率が上限閾値以上のノード100が存在すると判定した場合、ステップS1405に処理を移し、存在しないと判定した場合、ステップS1404に処理を移す。

【0159】

ステップS1404では、リバランス要否判定処理プログラム424は、ネットワーク220(NIC)の使用率が上限閾値以上のノード100が存在するか否かを判定する。リバランス要否判定処理プログラム424は、ネットワーク220の使用率が上限閾値以上のノード100が存在すると判定した場合、ステップS1405に処理を移し、存在しないと判定した場合、ステップS1406に処理を移す。

【0160】

ステップS1405では、リバランス要否判定処理プログラム424は、リソース割当決定処理(分散ポリシ)を実行する。なお、リソース割当決定処理(分散ポリシ)については、図15を用いて後述する。

【0161】

ステップS1406では、リバランス要否判定処理プログラム424は、プロセッサ211の使用率が下限閾値未満のノード100が存在するか否かを判定する。リバランス要否判定処理プログラム424は、プロセッサ211の使用率が下限閾値未満のノード100が存在すると判定した場合、ステップS1409に処理を移し、存在しないと判定した場合、ステップS1407に処理を移す。

【0162】

ステップS1407では、リバランス要否判定処理プログラム424は、ドライブ214の使用率が下限閾値未満のノード100が存在するか否かを判定する。リバランス要否判定処理プログラム424は、ドライブ214の使用率が下限閾値未満のノード100が存在すると判定した場合、ステップS1409に処理を移し、存在しないと判定した場合、ステップS1408に処理を移す。

【0163】

ステップS1408では、リバランス要否判定処理プログラム424は、ネットワーク220(NIC)の使用率が下限閾値未満のノード100が存在するか否かを判定する。リバランス要否判定処理プログラム424は、ネットワーク220の使用率が下限閾値未満のノード100が存在すると判定した場合、ステップS1409に処理を移し、存在しないと判定した場合、処理を終了する。

【0164】

ステップS1409では、リバランス要否判定処理プログラム424は、リソース割当決定処理(集約ポリシ)を実行する。なお、リソース割当決定処理(集約ポリシ)については、図16を用いて後述する。

【0165】

なお、ステップS1402~ステップS1404、および、ステップS1406~ステップS1408は、ノード100の負荷を判定するものであり、ノード100に仮想ボリューム313が1つ設けられている場合は、仮想ボリューム313の負荷を判定するものでもある。

【0166】

また、図14で説明したリソース割当決定処理を実行するか否かを判定するためのメトリクスは、プロセッサ211、ドライブ214、ポート215以外にも仮想ボリューム313に対するIOPSおよび/または転送量を用いてもよい。IOPSまたは転送量をメ

10

20

30

40

50

トリクスとして用いる場合、リードとライトとで異なる閾値を設けて判定してもよい。

【0167】

また、各ノード100は、自系ノードの負荷を低減するために、自系ノードが提供している仮想ボリューム313のうち、領域を移動していない仮想ボリューム313を他系ノードに移動してもよい。これにより、移動している領域がある仮想ボリューム313に、当該領域を戻すことが（当該領域を移動）できる場合がある。

【0168】

図15は、分散ポリシーに基づくリソース割当決定処理に係るフローチャートの一例を示す図である。

【0169】

ステップS1501では、リソース割当決定処理プログラム425は、各メトリクス（プロセッサ211、ドライブ214、ネットワーク220等）が上限閾値以上のノード100を移動元のノード100（移動元ノード）に選択し、移動対象とする仮想ボリューム313を選択する。言い換えると、リソース割当決定処理プログラム425は、プロセッサ211、ドライブ214、ネットワーク220の負荷に余裕のないノード100を選択し、データの移動元（ここでは分散元ともいえる）の仮想ボリューム313を選択する。

【0170】

より具体的には、リソース割当決定処理プログラム425は、リバランス要否判定処理のステップS1402、ステップS1403、またはステップS1404で選択されたノード100に定義された仮想ボリューム313を選択する。例えば、リソース割当決定処理プログラム425は、仮想ボリューム313を選択する際、ストレージプログラム311と同じノード100にアプリケーション301が動作するHCI（Hyper-Converged Infrastructure）構成をとっている場合、同じノード100内のアプリケーション301が使用している仮想ボリューム313を避けて選択する。これは、同じノード100内のアプリケーション301が使用している仮想ボリューム313を選択し、分散させるとアプリケーション301は、ネットワーク220を介してデータへアクセスすることになり、処理効率が低下するためである。

【0171】

付言するならば、ステップS1403、またはステップS1404で選択されたノード100に定義された仮想ボリューム313が1つである場合、リソース割当決定処理プログラム425は、当該仮想ボリューム313を選択する。

【0172】

ステップS1502では、リソース割当決定処理プログラム425は、仮想ボリューム313またはストレージシステム200に設定された分散ポリシーを判定する。リソース割当決定処理プログラム425は、分散ポリシーがボリューム単位分散ポリシーであると判定した場合、ステップS1503に処理を移し、分散ポリシーがスライス単位最大分散ポリシー（スライス単位均等分散ポリシー）であると判定した場合、ステップS1506に処理を移し、分散ポリシーがスライス単位最小分散ポリシーであると判定した場合、ステップS1509に処理を移す。

【0173】

ボリューム単位分散ポリシーは、仮想ボリューム313単位で負荷が分散されるポリシーである。ボリューム単位分散ポリシーでは、仮想ボリューム313単位でまとめてスライス314が移動される。ボリューム単位分散ポリシーでは、仮想ボリューム313単位でスライス314が移動されるため、常にデータの集約が保たれた状態で負荷分散を行うことができる。

【0174】

スライス単位最大分散ポリシーは、スライス314単位で負荷が分散されるポリシーである。スライス単位最大分散ポリシーでは、仮想ボリューム313に設定された最大分散度のノード数だけスライス314が分散される。スライス単位最大分散ポリシーでは、最大分散度でスライス314が分散されるため、高負荷なノード100（仮想ボリューム313）を

10

20

30

40

50

迅速に負荷分散することができる。

【0175】

スライス単位最小分散ポリシーは、スライス314単位で負荷が分散されるポリシーである。スライス単位最小分散ポリシーでは、仮想ボリューム313内のスライス314が1つずつ分散されていく。スライス単位最小分散ポリシーでは、1つずつスライス314が分散されていくことで、データのローカルリティを可能な限り保ちつつ、最小限の負荷だけを高負荷なノード100（仮想ボリューム313）から逃がすことで過負荷状態を回避する。

【0176】

これらの分散ポリシーは、ユーザが仮想ボリューム313に対して、事前に設定してもよいし、ストレージシステム200が状況に応じて分散ポリシーを自動で選択してもよい。ストレージシステム200が自動でポリシーを選択する方法の一例として、基本的には、ストレージシステム200は、ボリューム単位分散ポリシーを適用しておき、仮想ボリューム313が1つのノード100の性能で不足する場合に、スライス単位最大分散ポリシーまたはスライス単位最小分散ポリシーに切り替える。加えて、ストレージシステム200は、仮想ボリューム313の負荷が突発的に高くなった場合は、スライス単位最大分散ポリシーを適用し、仮想ボリューム313の負荷が緩やかに高くなった場合は、スライス単位最小分散ポリシーを適用する。

10

【0177】

ステップS1503、ステップS1504、およびステップS1505では、リソース割当決定処理プログラム425は、仮想ボリューム313単位でスライス314を移動するための前処理を行う。

20

【0178】

ステップS1503では、リソース割当決定処理プログラム425は、選択した仮想ボリューム313内の全てのスライス314をスライスグループとしてグルーピングする。例えば、リソース割当決定処理プログラム425は、スライス管理テーブル820から移動対象のスライスIDをメモリ212上にリストとして格納する。

【0179】

ステップS1504では、リソース割当決定処理プログラム425は、計算したスライスグループを移動対象として設定する。

【0180】

ステップS1505では、リソース割当決定処理プログラム425は、移動先ノード数を「1ノード」に設定する。

30

【0181】

ステップS1506、ステップS1507、およびステップS1508では、リソース割当決定処理プログラム425は、仮想ボリューム313に設定された最大分散度でスライス314を移動するための前処理を行う。

【0182】

ステップS1506では、リソース割当決定処理プログラム425は、選択した仮想ボリューム313内の全てのスライス314を仮想ボリューム313に設定された最大分散度で分割し、スライスグループとしてグルーピングする。例えば、リソース割当決定処理プログラム425は、グルーピングする際、選択した仮想ボリューム313にアクセスするホスト101が複数存在する場合で、かつ、ホスト101ごとにアクセス対象のスライス314に偏り（ローカルリティ）がある場合、ホスト101ごとのアクセス対象のスライス314をグルーピングする。また、例えば、リソース割当決定処理プログラム425は、スライスモニタ情報管理テーブル940を確認し、仮想ボリューム313内の全てのスライス314を最大分散度のノード100で負荷が均等になるようにグルーピングしてもよい。

40

【0183】

ステップS1507では、リソース割当決定処理プログラム425は、計算したスライスグループを移動対象として設定する。

50

## 【 0 1 8 4 】

ステップ S 1 5 0 8 では、リソース割当決定処理プログラム 4 2 5 は、移動先ノード数を最大分散度と同値に設定する。

## 【 0 1 8 5 】

ステップ S 1 5 0 9、ステップ S 1 5 1 0、およびステップ S 1 5 1 1 では、リソース割当決定処理プログラム 4 2 5 は、最小分散度（つまり 1 スライス）でスライス 3 1 4 を移動するための前処理を行う。

## 【 0 1 8 6 】

ステップ S 1 5 0 9 では、リソース割当決定処理プログラム 4 2 5 は、選択した仮想ボリューム 3 1 3 内からスライス 3 1 4 を 1 つ選択する。

10

## 【 0 1 8 7 】

ステップ S 1 5 1 0 では、リソース割当決定処理プログラム 4 2 5 は、選択したスライス 3 1 4 を移動対象として設定する。

## 【 0 1 8 8 】

ステップ S 1 5 1 1 では、リソース割当決定処理プログラム 4 2 5 は、移動先ノード数を「1 ノード」に設定する。

## 【 0 1 8 9 】

ステップ S 1 5 1 2 では、リソース割当決定処理プログラム 4 2 5 は、移動対象を 1 つ選択する。前処理で、スライスグループが作られた場合は、リソース割当決定処理プログラム 4 2 5 は、スライスグループを移動対象として選択し、スライス 3 1 4 がそのまま選択された場合は、スライス 3 1 4 を移動対象として選択する。

20

## 【 0 1 9 0 】

ステップ S 1 5 1 3 では、リソース割当決定処理プログラム 4 2 5 は、移動先とするノード 1 0 0 を 1 つ選択する。リソース割当決定処理プログラム 4 2 5 は、移動先ノードの選択方法の一例として、移動元ノードを除き、各メトリクス（プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等）の負荷に余裕のあるノード 1 0 0 を選択する方法がある。

## 【 0 1 9 1 】

ステップ S 1 5 1 4 では、リソース割当決定処理プログラム 4 2 5 は、選択した移動対象を移動先ノードに移動した場合の移動先ノードの各メトリクス（プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等）の負荷を計算する。

30

## 【 0 1 9 2 】

ステップ S 1 5 1 5 では、リソース割当決定処理プログラム 4 2 5 は、ステップ S 1 5 1 4 で計算した移動先ノードの各メトリクス（プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等）の負荷が閾値を超過していないかを判定する。リソース割当決定処理プログラム 4 2 5 は、1 つでも閾値を超過しているメトリクスがあると判定した場合、ステップ S 1 5 1 3 に処理を移し、全てのメトリクスにおいて閾値を超過していないと判定した場合、ステップ S 1 5 1 6 に処理を移す。なお、各メトリクス（プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等）の上限閾値とステップ S 1 5 1 5 の各メトリクスの閾値とは、同じであってもよいし、異なってもよい。

40

## 【 0 1 9 3 】

ステップ S 1 5 1 6 では、リソース割当決定処理プログラム 4 2 5 は、ステップ S 1 5 1 5 で判定したノード 1 0 0 を移動先ノードとして選択する。

## 【 0 1 9 4 】

ステップ S 1 5 1 7 では、リソース割当決定処理プログラム 4 2 5 は、全ての移動対象に対して判定を完了したか否かを判定する。リソース割当決定処理プログラム 4 2 5 は、全ての移動対象に対して判定が完了している場合、ステップ S 1 5 1 8 に処理を移し、全ての移動対象に対して判定が完了していない場合、ステップ S 1 5 1 2 に処理を移す。

## 【 0 1 9 5 】

ステップ S 1 5 1 8 では、リソース割当決定処理プログラム 4 2 5 は、移動対象として

50

スライスグループまたはスライス 3 1 4 を対象にして、リソース移動処理を実行する。リソース移動処理については、図 1 7 を用いて後述する。

【 0 1 9 6 】

ステップ S 1 5 1 9 では、リソース割当決定処理プログラム 4 2 5 は、閾値を超過していたノード 1 0 0 ( 仮想ボリューム 3 1 3 ) の負荷が閾値未満となったか否かを判定する。リソース割当決定処理プログラム 4 2 5 は、閾値未満でないと判定した場合、ステップ S 1 5 0 1 に処理を移し、閾値未満であると判定した場合、処理を終了する。

【 0 1 9 7 】

図 1 6 は、集約ポリシーに基づくリソース割当決定処理に係るフローチャートの一例を示す図である。

【 0 1 9 8 】

ステップ S 1 6 0 1 では、リソース割当決定処理プログラム 4 2 5 は、各メトリクス ( プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等 ) が下限閾値未満のノードを移動先ノードに選択する。言い換えると、リソース割当決定処理プログラム 4 2 5 は、プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等の負荷に余裕のあるノード 1 0 0 ( 仮想ボリューム 3 1 3 ) を選択し、データの移動先 ( ここでは集約先ともいえる ) のノード 1 0 0 とする。例えば、リソース割当決定処理プログラム 4 2 5 は、リバランス要否判定処理のステップ S 1 4 0 6、ステップ S 1 4 0 7、またはステップ S 1 4 0 8 で選択されたノード 1 0 0 を選択する。

【 0 1 9 9 】

ステップ S 1 6 0 2 では、リソース割当決定処理プログラム 4 2 5 は、ステップ S 1 6 0 1 で選択した移動先ノードに、スライス 3 1 4 が分散した仮想ボリューム 3 1 3 が存在するか否かを判定する。リソース割当決定処理プログラム 4 2 5 は、スライス 3 1 4 が分散した仮想ボリューム 3 1 3 が存在すると判定した場合、ステップ S 1 6 0 4 に処理を移し、スライス 3 1 4 が分散した仮想ボリューム 3 1 3 が存在しないと判定した場合、ステップ S 1 6 0 3 に処理を移す。

【 0 2 0 0 】

ステップ S 1 6 0 3 では、リソース割当決定処理プログラム 4 2 5 は、ステップ S 1 6 0 1 で選択した移動先ノード以外のノード 1 0 0 に、スライス 3 1 4 が分散した仮想ボリューム 3 1 3 が存在するか否かを判定する。リソース割当決定処理プログラム 4 2 5 は、スライス 3 1 4 が分散した仮想ボリューム 3 1 3 が存在すると判定した場合、ステップ S 1 6 0 4 に処理を移し、スライス 3 1 4 が分散した仮想ボリューム 3 1 3 が存在しないと判定した場合、集約対象 ( 移動対象 ) の仮想ボリューム 3 1 3 は存在しないため、処理を終了する。

【 0 2 0 1 】

上述したように、ステップ S 1 6 0 2 とステップ S 1 6 0 3 とでは、リソース割当決定処理プログラム 4 2 5 は、移動対象の仮想ボリューム 3 1 3 を選択する。リソース割当決定処理プログラム 4 2 5 は、選択する際、移動先ノードとして選択したノード 1 0 0 上の仮想ボリューム 3 1 3 を優先的に選択することで、集約時のスライス 3 1 4 の移動量を削減する効果がある。

【 0 2 0 2 】

ステップ S 1 6 0 4 では、リソース割当決定処理プログラム 4 2 5 は、選択した仮想ボリューム 3 1 3 を移動対象に選択する。

【 0 2 0 3 】

付言するならば、リソース割当決定処理プログラム 4 2 5 は、仮想ボリューム 3 1 3 の負荷に余裕のあるノード 1 0 0 があり、当該仮想ボリューム 3 1 3 のスライス 3 1 4 が分散している場合、当該ノード 1 0 0 を移動先ノードとし、当該仮想ボリューム 3 1 3 のスライス 3 1 4 が分散している仮想ボリューム 3 1 3 を移動対象に選択することがある。

【 0 2 0 4 】

ステップ S 1 6 0 5 では、リソース割当決定処理プログラム 4 2 5 は、移動対象ボリュ

10

20

30

40

50

ーム内のスライス 3 1 4 を 1 つ選択する。選択の方法としては、例えば、スライスモニタ情報管理テーブル 9 4 0 を参照し、移動対象ボリューム内のスライス 3 1 4 から処理負荷の高いスライス 3 1 4 を選択する。

【 0 2 0 5 】

ステップ S 1 6 0 6 では、リソース割当決定処理プログラム 4 2 5 は、選択したスライス 3 1 4 を移動先ノードに移動した場合の各メトリクス（プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等）の負荷を計算する。

【 0 2 0 6 】

ステップ S 1 6 0 7 では、リソース割当決定処理プログラム 4 2 5 は、ステップ S 1 6 0 6 で計算したスライス 3 1 4 を移動後の各メトリクス（プロセッサ 2 1 1、ドライブ 2 1 4、ネットワーク 2 2 0 等）の負荷が閾値を超過していないかを判定する。リソース割当決定処理プログラム 4 2 5 は、1 つでも閾値を超過しているメトリクスがあると判定した場合、処理を終了し、全てのメトリクスにおいて閾値を超過していないと判定した場合、ステップ S 1 6 0 8 に処理を移す。

【 0 2 0 7 】

ステップ S 1 6 0 8 では、リソース割当決定処理プログラム 4 2 5 は、移動対象のスライス 3 1 4 を対象にして、移動対象の仮想ボリューム 3 1 3 を提供するノード 1 0 0 にリソース移動処理の実行を要求する。リソース移動処理については、図 1 7 を用いて後述する。

【 0 2 0 8 】

ステップ S 1 6 0 9 では、リソース割当決定処理プログラム 4 2 5 は、移動対象ボリューム内の全てのスライス 3 1 4 に対して、移動するか否かの判定を行ったかを判定する。リソース割当決定処理プログラム 4 2 5 は、全てのスライス 3 1 4 に対して移動するか否かの判定を行った場合、処理を終了し、全てのスライス 3 1 4 に対して移動するか否かの判定を行っていない場合、ステップ S 1 6 0 5 に処理を移し、まだ未判定のスライス 3 1 4 に対して処理を行う。

【 0 2 0 9 】

図 1 7 は、リソース移動処理に係るフローチャートの一例を示す図である。リソース移動処理では、処理対象のスライス 3 1 4 について、現在割当てられているストレージプール 3 1 2 から別ノード 1 0 0 のストレージプール 3 1 2 へ割当先が移動される。リソース移動処理プログラム 4 2 6 は、スライス 3 1 4 の割当先を移動するにあたり、スライス 3 1 4 に書き込まれたデータを移動元のストレージプール 3 1 2 から読み出し、移動先のストレージプール 3 1 2 へ書き出す。

【 0 2 1 0 】

ステップ S 1 7 0 1 では、リソース移動処理プログラム 4 2 6 は、移動対象のスライス 3 1 4 の状態を移動中に更新する。より具体的には、リソース移動処理プログラム 4 2 6 は、スライス管理テーブル 8 2 0 から移動対象のスライス 3 1 4 のレコードを取得し、状態 8 2 4 の値を「M i g g r a t i n g」に更新する。

【 0 2 1 1 】

ステップ S 1 7 0 2 では、リソース移動処理プログラム 4 2 6 は、移動対象のスライス 3 1 4 の先頭オフセットを取得する。より具体的には、リソース移動処理プログラム 4 2 6 は、ページマッピングテーブル 8 3 0 を参照し、移動対象のスライス 3 1 4 のスライス ID に該当するレコードを参照し、プールボリューム ID とプールボリューム L B A とを取得する。次に、リソース移動処理プログラム 4 2 6 は、プールボリューム管理テーブル 6 1 0 を参照し、取得したプールボリューム ID に該当するレコードを参照し、論理ドライブ ID を取得する。次に、リソース移動処理プログラム 4 2 6 は、論理ドライブ管理テーブル 6 3 0 を参照し、取得した論理ドライブ ID に該当するレコードを取得し、ドライブ ID と開始オフセットを取得し、先に取得しているプールボリューム L B A からアクセス先のドライブ 2 1 4 のアドレスを求める。例えば、リソース移動処理プログラム 4 2 6 は、以下のように計算する。

10

20

30

40

50

## 【 0 2 1 2 】

アクセス先アドレス = 開始オフセット + プールボリューム L B A

## 【 0 2 1 3 】

ステップ S 1 7 0 3 では、リソース移動処理プログラム 4 2 6 は、アクセス先の領域の排他制御を取得する。

## 【 0 2 1 4 】

ステップ S 1 7 0 4 では、リソース移動処理プログラム 4 2 6 は、処理対象のオフセットのデータについて、ストレージプール 3 1 2 にページ 3 1 5 が未割当てであるか否かを判定する。リソース移動処理プログラム 4 2 6 は、未割当てであると判定した場合、ステップ S 1 7 0 7 に処理を移し、未割当てでないとして判定した場合、ステップ S 1 7 0 5 に処理を移す。

10

## 【 0 2 1 5 】

ステップ S 1 7 0 5 では、リソース移動処理プログラム 4 2 6 は、取得したオフセットに該当するスライス 3 1 4 の移動元のストレージプール 3 1 2 の領域にリードを発行する。より具体的には、リード処理が実行される。なお、リード処理の動作は、図 1 1 で説明したとおりである。

## 【 0 2 1 6 】

ステップ S 1 7 0 6 では、リソース移動処理プログラム 4 2 6 は、取得したオフセットに該当するスライス 3 1 4 の移動先のストレージプール 3 1 2 の領域にステップ S 1 7 0 5 で読み取ったデータでライトを発行する。ライト先として、移動元のスライス 3 1 4 に対してライトを発行する。移動元のスライス 3 1 4 では、ライト処理プログラム 4 2 2 内にてステップ S 1 2 0 6 の分岐で、Y e s の判定となり、移動先のスライス 3 1 4 へデータの書き込みがなされる。ただし、ステップ S 1 7 0 6 のライト発行先は、移動先のスライス 3 1 4 としてもよい。なお、ライト処理の動作は、図 1 2 で説明したとおりである。

20

## 【 0 2 1 7 】

ステップ S 1 7 0 7 では、リソース移動処理プログラム 4 2 6 は、アクセス先の領域の排他制御を解放する。

## 【 0 2 1 8 】

ステップ S 1 7 0 8 では、リソース移動処理プログラム 4 2 6 は、現在処理対象としているアクセス先のオフセットが移動対象のスライス 3 1 4 の終端オフセットであるか否かを判定する。リソース移動処理プログラム 4 2 6 は、終端オフセットであると判定した場合、ステップ S 1 7 1 0 に処理を移し、終端オフセットでないとして判定した場合、ステップ S 1 7 0 9 に処理を移す。

30

## 【 0 2 1 9 】

ステップ S 1 7 0 9 では、リソース移動処理プログラム 4 2 6 は、移動対象のスライス 3 1 4 の次のオフセットを取得する。

## 【 0 2 2 0 】

ステップ S 1 7 1 0 では、リソース移動処理プログラム 4 2 6 は、移動対象のスライス 3 1 4 の制御情報に関する排他制御を取得する。より具体的には、リソース移動処理プログラム 4 2 6 は、スライス管理テーブル 8 2 0 のアクセス先のスライス 3 1 4 のスライス I D に該当するレコードの排他制御を取得する。

40

## 【 0 2 2 1 】

ステップ S 1 7 1 1 では、リソース移動処理プログラム 4 2 6 は、スライス 3 1 4 を割当てするストレージプール 3 1 2 の情報を更新する。より具体的には、リソース移動処理プログラム 4 2 6 は、スライス管理テーブル 8 2 0 のストレージプール I D 8 2 3 の情報を、移動元のストレージプール I D から移動先のストレージプール I D に更新する。

## 【 0 2 2 2 】

ステップ S 1 7 1 2 では、リソース移動処理プログラム 4 2 6 は、移動対象のスライス 3 1 4 の制御情報に関する排他制御を解放する。より具体的には、リソース移動処理プログラム 4 2 6 は、スライス管理テーブル 8 2 0 のアクセス先のスライス 3 1 4 のスライス

50

IDに該当するレコードの排他制御を解放する。

【0223】

ステップS1713では、リソース移動処理プログラム426は、移動対象のスライス314の状態を正常に更新する。より具体的には、リソース移動処理プログラム426は、スライス管理テーブル820から移動対象のスライス314のレコードを取得し、状態824の値を「Normal」に更新する。

【0224】

ステップS1714では、リソース移動処理プログラム426は、移動対象のスライス314を全て移動したか否かを判定する。リソース移動処理プログラム426は、全ての移動対象のスライス314を移動していない判定した場合、ステップS1701に処理を移し、全ての移動対象のスライス314を移動したと判定した場合、処理を終了する。

10

【0225】

図18は、フロントエンドパス設定処理に係るフローチャートの一例を示す図である。

【0226】

ステップS1801では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311にスライス314の移動があったか否かを問合せ。

【0227】

ステップS1802では、ストレージプログラム311のフロントエンドパス設定処理プログラム427は、スライス314の移動の有無を判定し、パス設定プログラム302に  
20  
に  
応答する。より具体的には、ストレージプログラム311は、現在時刻までのスライス314の移動記録と、パス設定プログラム302から前回問い合わせられたときの時刻情報とをログとして保存しておき、パス設定プログラム302からの前回の問い合わせ時刻以降に、スライス314の移動記録が更新されているか確認する。ストレージプログラム311のフロントエンドパス設定処理プログラム427は、スライス314の移動記録が更新されている場合、移動があったと判定し、更新がなかった場合、移動はなかったと判定し、判定結果をパス設定プログラム302に  
20  
に  
応答する。

20

【0228】

ステップS1803では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311からの応答を受信する。  
30

30

【0229】

ステップS1804では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311からの応答内容に基づき、スライス314の移動があったか否かを判定する。パス設定プログラム302のフロントエンドパス設定処理プログラム427は、スライス314の移動があったと判定した場合、ステップS1805に処理を移し、スライス314の移動がなかったと判定した場合、処理を終了する。

【0230】

ステップS1805では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、移動したスライス314の情報をストレージプログラム311に問い合わせる。  
40

40

【0231】

ステップS1806では、ストレージプログラム311のフロントエンドパス設定処理プログラム427は、移動したスライス314と、当該スライス314を含む仮想ボリューム313との情報をパス設定プログラム302に  
50  
に  
応答する。

50

【0232】

ステップS1807では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311からの応答を受信する。

【0233】

ステップS1808では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311からの応答を受信する。

プログラム 4 2 7 は、移動したスライス 3 1 4 の移動先ノード（移動したスライス 3 1 4 を含む仮想ボリューム 3 1 3）にパス 3 2 0 が設定されているか否かを確認する。パス設定プログラム 3 0 2 のフロントエンドパス設定処理プログラム 4 2 7 は、フロントエンドパス 3 2 0 を設定済みであると判定した場合、処理を終了し、フロントエンドパス 3 2 0 を未設定であると判定した場合、ステップ S 1 8 0 9 に処理を移す。

【 0 2 3 4 】

ステップ S 1 8 0 9 では、フロントエンドパス設定処理プログラム 4 2 7 は、移動したスライス 3 1 4 を含む仮想ボリューム 3 1 3 について、当該スライス 3 1 4 の移動先ノードとのフロントエンドパス 3 2 0 の確立を要求する。パス設定プログラム 3 0 2 とストレージプログラム 3 1 1 と間での具体的なパス確立の手順は、i S C S I ( I n t e r n e t S m a l l C o m p u t e r S y s t e m I n t e r f a c e )、F i b e r C h a n n e l のプロトコル等に基づく。

10

【 0 2 3 5 】

ステップ S 1 8 1 0 では、ストレージプログラム 3 1 1 のフロントエンドパス設定処理プログラム 4 2 7 は、パス設定プログラム 3 0 2 から要求された仮想ボリューム 3 1 3 とホスト 1 0 1 とのパス情報をフロントエンドパス情報テーブル 1 0 1 0 に登録し、パス設定プログラム 3 0 2 に応答する。

【 0 2 3 6 】

ステップ S 1 8 1 1 では、パス設定プログラム 3 0 2 のフロントエンドパス設定処理プログラム 4 2 7 は、ストレージプログラム 3 1 1 からパス設定完了の応答を受信する。

20

【 0 2 3 7 】

図 1 9 は、S C S I ( S m a l l C o m p u t e r S y s t e m I n t e r f a c e ) 規格の A L U A ( A s y m m e t r i c L o g i c a l U n i t A c c e s s ) のメカニズムに基づきフロントエンドパス 3 2 0 の設定を最適化する際のフロントエンドパス設定処理に係るフローチャートの一例を示す図である。

【 0 2 3 8 】

ステップ S 1 9 0 1 では、ストレージプログラム 3 1 1 のフロントエンドパス設定処理プログラム 4 2 7 は、処理対象の仮想ボリューム 3 1 3 を選択する。より具体的には、処理対象の仮想ボリューム 3 1 3 については、定義された全ての仮想ボリューム 3 1 3 に対して周期的に選択してもよいし、リード処理またはライト処理の実行の完了後に、当該処理においてアクセスのあった仮想ボリューム 3 1 3 を選択してもよい。

30

【 0 2 3 9 】

ステップ S 1 9 0 2 では、ストレージプログラム 3 1 1 のフロントエンドパス設定処理プログラム 4 2 7 は、仮想ボリューム 3 1 3 に接続されたホスト 1 0 1 を処理対象として選択する。より具体的には、処理対象のホスト 1 0 1 については、処理対象の仮想ボリューム 3 1 3 に定義された全てのホスト 1 0 1 を選択してもよいし、リード処理またはライト処理の実行の完了後に、当該処理においてアクセスのあったホスト 1 0 1 を選択してもよい。

【 0 2 4 0 】

ステップ S 1 9 0 3 では、ストレージプログラム 3 1 1 のフロントエンドパス設定処理プログラム 4 2 7 は、選択した仮想ボリューム 3 1 3 と選択したホスト 1 0 1 と間のパスの情報を取得する。より具体的には、フロントエンドパス設定処理プログラム 4 2 7 は、フロントエンドパス情報テーブル 1 0 1 0 を参照し、選択した仮想ボリューム 3 1 3 の仮想ボリューム ID と、選択したホスト 1 0 1 の I n i t i a t o r I D とに該当するレコードを取得する。

40

【 0 2 4 1 】

ステップ S 1 9 0 4 では、フロントエンドパス設定処理プログラム 4 2 7 は、ステップ S 1 9 0 3 で取得したフロントエンドパス 3 2 0 への I / O の発行比率を計算する。より具体的には、フロントエンドパス設定処理プログラム 4 2 7 は、フロントエンドパスモニタ情報管理テーブル 9 5 0 のリード I O P S 9 5 2 とライト I O P S 9 5 3 とを参照し、

50

パスAにリードおよびライトのIOPSが合計900IOPS発行されており、パスBにリードおよびライトのIOPSが合計100IOPS発行されている場合、パスAとパスBのI/O発行比率を、パスA：パスB = 9：1と計算する。加えて、フロントエンドパス設定処理プログラム427は、同様にリードおよびライトの転送量についても、比率を計算し、IOPSと転送量とで各パスのI/Oの発行比率の分散が大きい方を最終的な比率として採用してもよいし、IOPSと転送量との比率の平均値を最終的な比率として採用してもよい。

**【0242】**

ステップS1905では、ストレージプログラム311のフロントエンドパス設定処理プログラム427は、I/Oの発行比率が特定のフロントエンドパス320（ホスト101）に偏っているか否かを判定する。例えば、フロントエンドパス設定処理プログラム427は、各フロントエンドパス320へのI/O発行比率の分散を計算し、分散が閾値を上回っている場合に偏りが発生していると判定する。ストレージプログラム311のフロントエンドパス設定処理プログラム427は、特定のフロントエンドパス320にI/Oの発行比率が偏っていると判定した場合、ステップS1906に処理を移し、特定のフロントエンドパス320にI/Oの発行比率が偏っていないと判定した場合、ステップS1907に処理を移す。

10

**【0243】**

ステップS1906では、ストレージプログラム311のフロントエンドパス設定処理プログラム427は、I/Oの発行比率が高いフロントエンドパス320を「Optimize（最適パス）」に設定する。例えば、フロントエンドパス設定処理プログラム427は、フロントエンドパス情報テーブル1010を参照し、I/Oの発行比率が閾値より高いフロントエンドパス320のパスIDと一致するレコードのALUA設定1014を「Optimize」に更新し、I/Oの発行比率が閾値より高くないフロントエンドパス320のパスIDと一致するレコードのALUA設定1014を、「Non-Optimize」に更新する。

20

**【0244】**

ステップS1907では、ストレージプログラム311のフロントエンドパス設定処理プログラム427は、選択した仮想ボリューム313に定義された全てのフロントエンドパス320について「Optimize（最適パス）」に設定（いわゆるラウンドロビンに設定）する。

30

**【0245】**

ステップS1908では、ストレージプログラム311のフロントエンドパス設定処理プログラム427は、パス設定プログラム302に、「Optimize」に設定したフロントエンドパス320の情報（最適化情報）を通知する。図19では、ストレージプログラム311からパス設定プログラム302へ通知を発行する形式で情報を伝達するが、パス設定プログラム302からストレージプログラム311に問い合わせる形式で情報を伝達してもよい。

**【0246】**

ステップS1909では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311からフロントエンドパス320の最適化情報を受領する。

40

**【0247】**

ステップS1910では、パス設定プログラム302のフロントエンドパス設定処理プログラム427は、ストレージプログラム311から受信したフロントエンドパス320を「Optimize（最適パス）」に設定する。以降、アプリケーション301は、最適パスに対して優先的にI/Oを発行するように制御される。

**【0248】**

図20は、クラスタ構成変更処理に係るフローチャートの一例を示す図である。

**【0249】**

50

ステップ S 2 0 0 1 では、クラスタ構成変更処理プログラム 4 2 8 は、ユーザからのクラスタ操作要求を受信する。

【 0 2 5 0 】

ステップ S 2 0 0 2 では、クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求が仮想ボリューム 3 1 3 の作成または削除の要求であるか否かを判定する。クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求が仮想ボリューム 3 1 3 の作成または削除の要求であると判定した場合、ステップ S 2 0 0 7 に処理を移し、受信したクラスタ操作要求が仮想ボリューム 3 1 3 の作成または削除の要求でないとは判定した場合、ステップ S 2 0 0 3 に処理を移す。

【 0 2 5 1 】

ステップ S 2 0 0 3 では、クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求が仮想ボリューム 3 1 3 の拡張または縮小の要求であるか否かを判定する。クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求が仮想ボリューム 3 1 3 の拡張または縮小の要求であると判定した場合、ステップ S 2 0 0 7 に処理を移し、受信したクラスタ操作要求が仮想ボリューム 3 1 3 の拡張または縮小の要求でないとは判定した場合、ステップ S 2 0 0 4 に処理を移す。

【 0 2 5 2 】

ステップ S 2 0 0 4 では、クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求がドライブ 2 1 4 の増設または減設の要求であるか否かを判定する。クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求がドライブ 2 1 4 の増設または減設の要求であると判定した場合、ステップ S 2 0 0 7 に処理を移し、受信したクラスタ操作要求がドライブ 2 1 4 の増設または減設の要求でないとは判定した場合、ステップ S 2 0 0 5 に処理を移す。

【 0 2 5 3 】

ステップ S 2 0 0 5 では、クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求がノード 1 0 0 の増設または減設の要求であるか否かを判定する。クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求がノード 1 0 0 の増設または減設の要求であると判定した場合、ステップ S 2 0 0 7 に処理を移し、受信したクラスタ操作要求がノード 1 0 0 の増設または減設の要求でないとは判定した場合、ステップ S 2 0 0 6 に処理を移す。

【 0 2 5 4 】

ステップ S 2 0 0 6 では、クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求がサイト 2 0 1 の増設または減設の要求であるか否かを判定する。クラスタ構成変更処理プログラム 4 2 8 は、受信したクラスタ操作要求がサイト 2 0 1 の増設または減設の要求であると判定した場合、ステップ S 2 0 0 7 に処理を移し、受信したクラスタ操作要求がサイト 2 0 1 の増設または減設の要求でないとは判定した場合、ステップ S 2 0 0 1 に処理を移す。

【 0 2 5 5 】

ステップ S 2 0 0 7 では、クラスタ構成変更処理プログラム 4 2 8 は、ユーザから要求された操作を実行する。

【 0 2 5 6 】

ステップ S 2 0 0 8 では、クラスタ構成変更処理プログラム 4 2 8 は、リバランス要否判定処理を実行する。なお、リバランス要否判定処理については、図 1 4 を用いて説明したので、その説明を省略する。

【 0 2 5 7 】

図 2 1 A は、ユーザが操作可能なボリュームの設定画面 ( GUI : Graphical User Interface ) の一例を示す図である。ボリューム設定画面 2 1 0 1 は、ストレージシステム 2 0 0 と通信可能に接続された所定の計算機 ( 例えば、管理サーバ等 ) へ出力される。

【 0 2 5 8 】

10

20

30

40

50

ボリューム設定画面 2 1 0 1 は、仮想ボリューム 3 1 3 ごとの設定が行われる画面である。ボリューム設定画面 2 1 0 1 は、ボリューム ID 2 1 0 2 および分散度 2 1 0 3 の情報を設定可能に構成される。

【 0 2 5 9 】

ボリューム ID 2 1 0 2 は、設定対象とする仮想ボリューム 3 1 3 を指定する項目である。分散度 2 1 0 3 は、設定対象の仮想ボリューム 3 1 3 のスライス 3 1 4 を分散するノード数の最大値の情報（最大分散度）を設定可能な項目である。

【 0 2 6 0 】

図 2 1 B は、ストレージシステム 2 0 0 がユーザに提示するボリューム性能予測画面 2 1 1 1（GUI）の一例を示す図である。ボリューム性能予測画面 2 1 1 1 は、ストレージシステム 2 0 0 と通信可能に接続された所定の計算機（例えば、管理サーバ等）に出力される。

10

【 0 2 6 1 】

ボリューム性能予測画面 2 1 1 1 は、ボリューム性能予測情報 2 1 1 2 とメッセージ 2 1 2 0 とを含んで構成される。

【 0 2 6 2 】

ボリューム性能予測情報 2 1 1 2 は、確認対象の仮想ボリューム 3 1 3 について、最大分散度とスループット（IOPS）と応答時間とに関する情報をユーザに提示する。ボリューム性能予測情報 2 1 1 2 は、分散度情報 2 1 1 3、スループット情報 2 1 1 4、応答時間情報 2 1 1 5、予測スループット情報 2 1 1 6、予測応答時間情報 2 1 1 7、現行分散度情報 2 1 1 8、および目標性能値情報 2 1 1 9 を含んで構成される。

20

【 0 2 6 3 】

分散度情報 2 1 1 3 は、最大分散度に関する情報を含む。スループット情報 2 1 1 4 は、IOPSに関する情報を含む。応答時間情報 2 1 1 5 は、応答時間に関する情報を含む。予測スループット情報 2 1 1 6 は、該当する最大分散度におけるスループットの予測値に関する情報を含む。予測応答時間情報 2 1 1 7 は、該当する最大分散度における応答時間の予測値に関する情報を含む。現行分散度情報 2 1 1 8 は、確認対象の仮想ボリューム 3 1 3 に現在設定されている最大分散度に関する情報を含む。目標性能値情報 2 1 1 9 は、確認対象の仮想ボリューム 3 1 3 に現在設定されている目標とする（ユーザが期待する）スループットおよび応答時間の性能値に関する情報を含む。

30

【 0 2 6 4 】

メッセージ 2 1 2 0 は、現行分散度情報 2 1 1 8 に基づく予測応答時間情報 2 1 1 7 が目標性能値情報 2 1 1 9 を超過する場合、または、現行分散度情報 2 1 1 8 に基づく予測スループット情報 2 1 1 6 が目標性能値情報 2 1 1 9 未満となる場合、ストレージシステム 2 0 0 は、ユーザに対して、現在の最大分散度では目標性能値を達成できない旨を提示する。また、メッセージ 2 1 2 0 は、現在の設定値から最大分散度を変更したときに予測性能が目標性能値を満たせるか否かをユーザに提示する。

【 0 2 6 5 】

本実施の形態によれば、1つのボリュームについて容量および性能をスケールアウトすることができる。

40

【 0 2 6 6 】

（II）第2の実施の形態

図 2 2 は、本実施の形態のストレージシステム 2 2 0 0 に係る構成の一例を示す図である。第1の実施の形態と同じ構成については、同じ符号を用いて説明を省略する。

【 0 2 6 7 】

ストレージシステム 2 2 0 0 は、図 2 の構成に加えて、ドライブボックス 2 2 0 1 を含んで構成される。ドライブボックス 2 2 0 1 は、例えば、プロセッサ 2 1 1、メモリ 2 1 2 等を含む 1 以上のプロセッサパッケージ 2 1 3、1 以上のドライブ 2 1 4、1 以上のポート 2 1 5 を含んで構成される。各構成要素は、内部バス 2 1 6 を介して接続されている。

【 0 2 6 8 】

50

ポート 215 は、ネットワーク 2202 に接続され、サイト 201 内のノード 100 と通信可能に接続されている。ネットワーク 2202 は、例えば、LAN、SAN (Storage Area Network)、SAS (Serial Attached SCSI) であるが、これらに限定するものではない。

【0269】

プロセッサ 211 では、ドライブボックス 2201 に対する、I/O 処理プログラムが動作しており、必要に応じて、データの圧縮処理、ドライブボックス 2201 内のドライブ 214 に対する RAID 処理等を実施してもよい。また、ドライブボックス 2201 は、上記に加えて、専用のハードウェアである ASIC (Application Specific Integrated Circuit) が搭載されていてもよく、ASIC は、データ圧縮処理、パリティ演算処理等を備えていてもよい。

10

【0270】

図 23 は、本実施の形態におけるリソース移動処理に係るフローチャートの一例を示す図である。リソース移動処理では、リバランス要否判定処理の結果を受けて、リバランスが必要となった場合にノード 100 間でスライス 314 を移動する。

【0271】

第 1 の実施の形態との違いについて説明する。本実施の形態では、スライス 314 内のデータは、ノード 100 間で共有されたドライブボックス 2201 に配置される。このため、ノード 100 間でスライス 314 を移動する場合も仮想ボリューム 313 とスライス 314 の割り当てとに関する制御情報を更新するだけでよく、データ自体は移動する必要はない。以下、詳細について説明する。

20

【0272】

ステップ S2301 では、リソース移動処理プログラム 426 は、移動対象のスライス 314 の状態を移動中に更新する。より具体的には、リソース移動処理プログラム 426 は、スライス管理テーブル 820 から移動対象のスライス 314 のレコードを取得し、状態 824 の値を「Migrating」に更新する。

【0273】

ステップ S2302 では、リソース移動処理プログラム 426 は、移動対象のスライス 314 の制御情報に関する排他制御を取得する。より具体的には、リソース移動処理プログラム 426 は、スライス管理テーブル 820 のアクセス先のスライス 314 のスライス ID に該当するレコードの排他制御を取得する。

30

【0274】

ステップ S2303 では、リソース移動処理プログラム 426 は、スライス 314 を割り当てるストレージプール 312 の情報を更新する。より具体的には、リソース移動処理プログラム 426 は、スライス管理テーブル 820 のストレージプール ID 823 の情報を、移動元のストレージプール ID から移動先のストレージプール ID に更新する。

【0275】

ステップ S2304 では、リソース移動処理プログラム 426 は、移動対象のスライス 314 の制御情報に関する排他制御を解放する。より具体的には、リソース移動処理プログラム 426 は、スライス管理テーブル 820 のアクセス先のスライス 314 のスライス ID に該当するレコードの排他制御を解放する。

40

【0276】

ステップ S2305 では、リソース移動処理プログラム 426 は、移動対象のスライス 314 の状態を正常に更新する。より具体的には、リソース移動処理プログラム 426 は、スライス管理テーブル 820 から移動対象のスライス 314 のレコードを取得し、状態 824 の値を「Normal」に更新する。

【0277】

本実施の形態によれば、データを移動することなく、1 つのボリュームについて性能をスケールアウトすることができる。

【0278】

50

## ( I I I ) 付記

上述の実施の形態には、例えば、以下のような内容が含まれる。

## 【 0 2 7 9 】

上述の実施の形態においては、本発明をストレージシステムに適用するようにした場合について述べたが、本発明はこれに限らず、この他種々のシステム、装置、方法、プログラムに広く適用することができる。

## 【 0 2 8 0 】

また、上述の実施の形態においては、スライスモニタ情報管理テーブル 9 4 0 は、スライス I D 9 4 1 と、リードカウンタ 9 4 2 と、ライトカウンタ 9 4 3 と、リード転送量 9 4 4 と、ライト転送量 9 4 5 と、モニタ開始時刻 9 4 6 とが対応付けられた情報を格納する場合について述べたが、本発明はこれに限らない。例えば、スライスモニタ情報管理テーブル 9 4 0 は、スライス I D 9 4 1 と、リード I O P S と、ライト I O P S と、リード転送量 9 4 4 と、ライト転送量 9 4 5 とが対応付けられた情報を格納するようにしてもよい。

10

## 【 0 2 8 1 】

また、上述の実施の形態においては、フロントエンドパスモニタ情報管理テーブル 9 5 0 は、パス I D 9 5 1 と、リード I O P S 9 5 2 と、ライト I O P S 9 5 3 と、リード転送量 9 5 4 と、ライト転送量 9 5 5 とが対応付けられた情報を格納する場合について述べたが、本発明はこれに限らない。例えば、フロントエンドパスモニタ情報管理テーブル 9 5 0 は、パス I D 9 5 1 と、リードカウンタと、ライトカウンタと、リード転送量 9 5 4 と、ライト転送量 9 5 5 と、モニタ開始時刻とが対応付けられた情報を格納するようにしてもよい。

20

## 【 0 2 8 2 】

また、上述の実施の形態において、各テーブルの構成は一例であり、1つのテーブルは、2以上のテーブルに分割されてもよいし、2以上のテーブルの全部または一部が1つのテーブルであってもよい。

## 【 0 2 8 3 】

また、上述の実施の形態において、図示および説明した画面は、一例であり、受け付ける情報が同じであるならば、どのようなデザインであってもよい。

## 【 0 2 8 4 】

また、上述の実施の形態において、図示および説明した画面は、一例であり、提示する情報が同じであるならば、どのようなデザインであってもよい。

30

## 【 0 2 8 5 】

また、上述の実施の形態において、統計値として分散および平均値を用いる場合について説明したが、統計値は、分散および平均値に限るものではなく、最大値、最小値、最大値と最小値との差、最頻値、中央値、標準偏差等の他の統計値であってもよい。

## 【 0 2 8 6 】

上述した実施の形態は、例えば、以下の特徴的な構成を有する。

## 【 0 2 8 7 】

## ( 1 )

複数の領域を含むボリューム（例えば、仮想ボリューム 3 1 3）を1以上のホスト（例えば、ホスト 1 0 1）に提供するための処理を行うプロセッサ（例えば、プロセッサ 2 1 1、プロセッサパッケージ 2 1 3）を備える複数のノード（例えば、ノード 1 0 0）と、上記プロセッサと接続され、上記ボリュームのデータを記憶する1以上の記憶デバイス（例えば、ドライブ 2 1 4、ドライブボックス 2 2 0 1）とを備えるストレージシステム（例えば、ストレージシステム 2 0 0、ストレージシステム 2 2 0 0）は、上記複数のノードの各々は、自ノードが提供するボリュームの負荷および上記ボリュームの領域を複数に分割した領域の負荷を監視し、監視している一のボリュームの負荷が閾値以上であると判定した第1のノードは、上記一のボリュームの領域を複数に分割した領域の負荷と負荷分散のポリシー（ボリューム単位分散ポリシー、スライス単位最大分散ポリシー、スライス単位最

40

50

小分散ポリシ等)とに応じて、上記一のボリュームに含まれる一部の領域を上記第1のノードとは異なる第2のノードのボリュームに移動する(例えば、図15参照)。

【0288】

上記ボリュームの負荷は、例えば、自ノードのプロセッサの使用率、自ノードの記憶デバイスへのI/O量、自ノードの記憶デバイスへのI/Oレスポンス、自ノードのネットワークインターフェースへのI/O量、自ノードのネットワークインターフェースへのI/Oレスポンス、ボリュームに対するIOPS、および、ボリュームに対する転送量の少なくとも1つである。

【0289】

上記構成では、一のボリュームの負荷が高まったときに、当該ボリュームの一部の領域が他のノードのボリュームに移動されるので、例えば、ノードを追加した場合、一のボリュームについても性能をスケールアウトすることができるようになる。

10

【0290】

(2)

上記第1のノードは、上記一のボリュームの負荷が閾値未満であると判定した場合(例えば、図14参照)、上記第2のノードのボリュームから、上記一部の領域を上記一のボリュームに移動する(例えば、図16参照)。

【0291】

上記構成では、例えば、領域が移動されているボリュームの負荷が低い場合、当該ボリュームに当該領域が集約されるので、当該ボリュームのスループットを向上させ、当該ボリュームのレイテンシを低下させることができる。

20

【0292】

(3)

上記第1のノードは、上記第2のノードとして、上記一部の領域を移動した後のボリュームの負荷が閾値を超えないノードを選択する(例えば、ステップS1514~ステップS1516参照)。

【0293】

上記構成では、例えば、移動先ノードのボリュームが過負荷になってしまい、更に移動しなければならない事態を回避することができる。

【0294】

(4)

上記複数のノードの各々には、上記1以上の記憶デバイスの少なくとも1つが対応して設けられ(例えば、図2参照)、上記複数のノードの各々は、自ノードに割り当てられている領域のデータを、自ノードに設けられている記憶デバイスに記憶し(例えば、図12参照)、上記第1のノードは、上記一のボリュームの容量が提供可能な容量を超えると判定した場合(例えば、ステップS1401参照)、上記一部の領域を上記第2のノードのボリュームに移動する。

【0295】

上記構成によれば、例えば、各ノードが提供するボリュームは、ストレージシステムが備える記憶デバイス分、容量を利用できるようになる。

40

【0296】

(5)

上記複数のノードの各々は、自ノードが提供するボリュームに対するリードの負荷(例えば、リードIOPS、リード転送量、リードカウンタ等)と、上記ボリュームに対するライトの負荷(例えば、ライトIOPS、ライト転送量、ライトカウンタ等)とを監視し、上記第1のノードは、上記一のボリュームに対するリードの負荷が第1の閾値以上であると判定した場合、上記一部の領域を上記第2のノードのボリュームに移動し、上記一のボリュームに対するライトの負荷が上記第1の閾値とは異なる第2の閾値以上であると判定した場合、上記一部の領域を上記第2のノードのボリュームに移動する。

【0297】

50

ボリュームに対するリードとライトとでは、ノードにかかる負荷が異なるが、上記構成では、それぞれに対して監視を行い、別の閾値を設けることで、例えば、より適切にボリュームの負荷を判定することができる。

【0298】

(6)

上記第1のノードは、上記複数のノードに対して上記一のボリュームに含まれる領域を均等に割り振り（例えば、ステップS1506、ステップS1507）、上記第1のノードとは異なるノードに割り振った領域を、上記ノードのボリュームに移動する。

【0299】

上記構成では、例えば、移動する領域による負荷が均等になるようにボリュームの領域を移動することができる。

10

【0300】

(7)

上記第1のノードは、上記一のボリュームの負荷が上記閾値を下回るまで、上記一のボリュームに含まれる領域を1つずつ（例えば、ステップS1509、ステップS1510）、上記第1のノードとは異なるノードのボリュームに移動する。

【0301】

上記構成では、例えば、データのローカリティを極力保ち、最小限の負荷をボリュームから逃すことができる。

【0302】

(8)

上記一のボリュームが、複数のホストに提供されている場合、上記第1のノードは、上記複数のホストの各々がアクセスする領域をまとめてホストごとの移動対象とし（例えば、ステップS1506、ステップS1507）、ホストごとの移動対象の領域を、上記第1のノードとは異なるノードのボリュームに移動する。

【0303】

上記構成では、例えば、ホストごとにアクセスする領域が違う場合に、ホストがアクセスする領域をまとめて移動することができる。

【0304】

(9)

上記第1のノードは、自ノードが提供しているボリュームのうち上記一のボリュームとは異なる他のボリュームを上記第1のノードとは異なるノードに移動し、上記第2のノードのボリュームから、上記一部の領域を上記一のボリュームに移動する。

【0305】

上記構成では、例えば、第1のノードは、一のボリュームとは異なる他のボリュームをまるごと別のノードに移動することにより、移動した領域を集めることができる場合がある。

【0306】

(10)

上記一部の領域が移動された上記第2のノードのボリュームと、上記一部の領域にアクセスするホストとの間にパスが設定されていない場合、上記第2のノードおよび上記ホストは、上記パスを設定する（例えば、図18参照）。

【0307】

上記構成では、領域が移動された第2のノードのボリュームと当該領域にアクセスするホストとの間にパスが設定されていない場合、パスが設定されるので、例えば、当該領域に対するアクセスがあった際、第1のノードを介することなくデータをやり取りすることができるようになる。

40

【0308】

(11)

上記複数のノードの各々は、自ノードのボリュームの負荷が特定のホストに偏っている

50

と判定した場合、上記特定のホストとのパスが最適属性であることを上記特定のホストに通知する（例えば、ステップ S 1 9 0 4 ~ ステップ S 1 9 0 6、ステップ S 1 9 0 8）。

【 0 3 0 9 】

上記構成によれば、例えば、ボリュームの負荷が特定のホストに偏っている場合、特定のホストに優先的に I / O が発行されるようになる。

【 0 3 1 0 】

( 1 2 )

上記複数のノードの各々は、自ノードのボリュームの負荷が上記特定のホストに偏っていないと判定した場合、上記ボリュームに定義されている全てのパスが最適属性であることを上記パスが設定されている全てのホストに通知する（例えば、ステップ S 1 9 0 4、ステップ S 1 9 0 5、ステップ S 1 9 0 7、ステップ S 1 9 0 8）。

10

【 0 3 1 1 】

上記構成によれば、ボリュームの負荷が特定のホストに偏っていない場合、当該ボリュームにアクセスする全てのホストに均等に I / O が発行されるようになる。

【 0 3 1 2 】

( 1 3 )

上記 1 以上の記憶デバイス（例えば、ドライブボックス 2 2 0 1）は、上記複数のノードに共通して設けられている。

【 0 3 1 3 】

上記構成では、ノードのコンピュータ部分と記憶部分とが分かれたことにより、例えば、プロセッサの使用率が余っているが、記憶デバイスの容量が不足している場合、コンピュータ部分を増やすことなく、記憶部分を増やすことができる。

20

【 0 3 1 4 】

( 1 4 )

上記第 1 のノードは、上記第 2 のノードに上記一部の領域を移動する際、上記一部の領域を管理するためのデータ（例えば、スライス管理テーブル 8 2 0）を更新し、上記 1 以上の記憶デバイスに記憶される上記一部の領域のデータを移動しない（例えば、図 2 3 参照）。

【 0 3 1 5 】

上記構成では、どのノードからも、均等にデータにアクセスできるようになっているので、例えば、ノード間でボリュームの負荷を分散するときは、ボリュームのオーナー権（排他権、メタデータ）を移すだけで、ボリュームの負荷を分散できる。

30

【 0 3 1 6 】

( 1 5 )

上記複数のノードが提供するボリュームについて、当該ボリュームの領域の移動先ノードの数に応じたスループットと応答時間とを計算して出力する計算機（例えば、管理サーバ）を備える。

【 0 3 1 7 】

上記構成によれば、例えば、ユーザは、領域の移動先ノードの数（例えば、最大分散度）を容易に決定することができる。

40

【 0 3 1 8 】

上記複数のノードの各々は、ボリュームに含まれる領域単位で、ボリュームの負荷を監視する（例えば、図 1 3 B）。

【 0 3 1 9 】

上記ストレージシステムは、1つのボリュームが分散可能な最大のノードの数をユーザが指定するための GUI を出力する計算機を備える（図 2 1 A）。

【 0 3 2 0 】

また上述した構成については、本発明の要旨を超えない範囲において、適宜に、変更したり、組み替えたり、組み合わせたり、省略したりしてもよい。

【 0 3 2 1 】

50

「A、B、およびCのうちの少なくとも1つ」という形式におけるリストに含まれる項目は、(A)、(B)、(C)、(AおよびB)、(AおよびC)、(BおよびC)または(A、B、およびC)を意味することができると理解されたい。同様に、「A、B、またはCのうちの少なくとも1つ」の形式においてリストされた項目は、(A)、(B)、(C)、(AおよびB)、(AおよびC)、(BおよびC)または(A、B、およびC)を意味することができる。

【符号の説明】

【0322】

100.....ノード、101.....ホスト、200.....ストレージシステム。

10

20

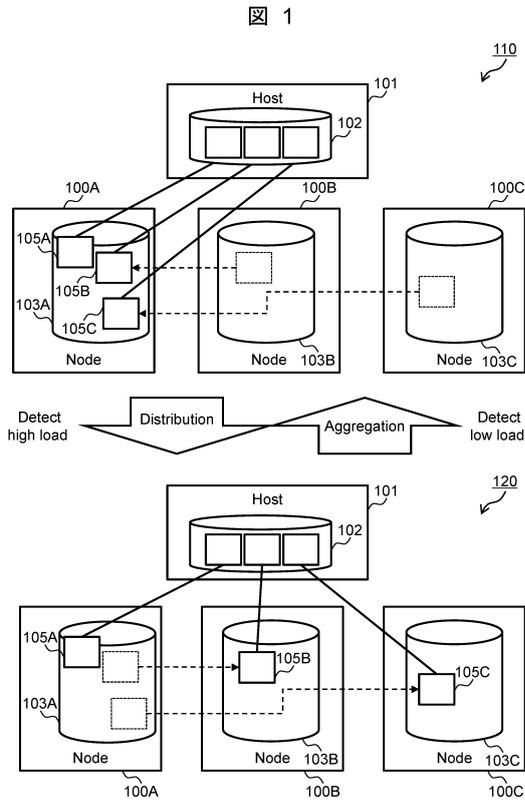
30

40

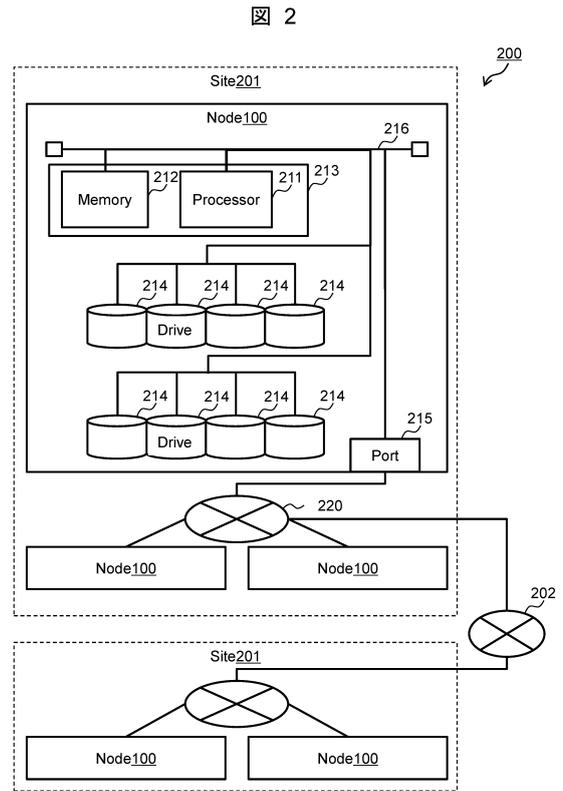
50

【図面】

【図 1】



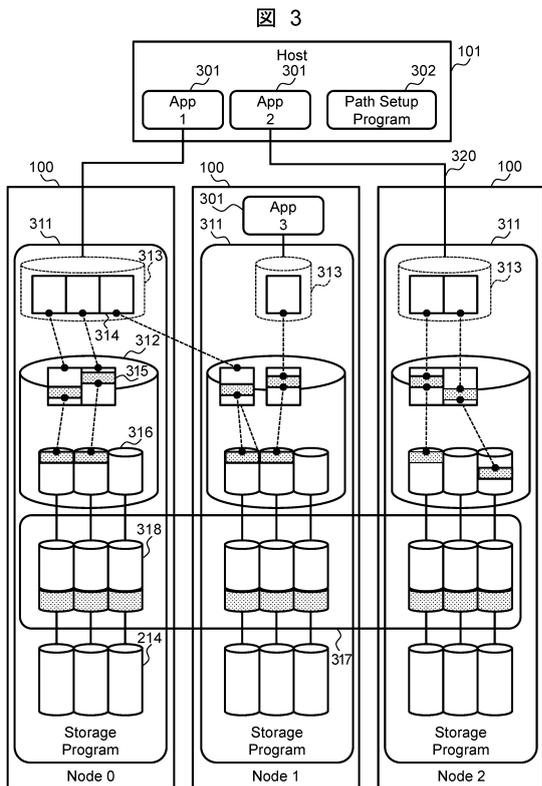
【図 2】



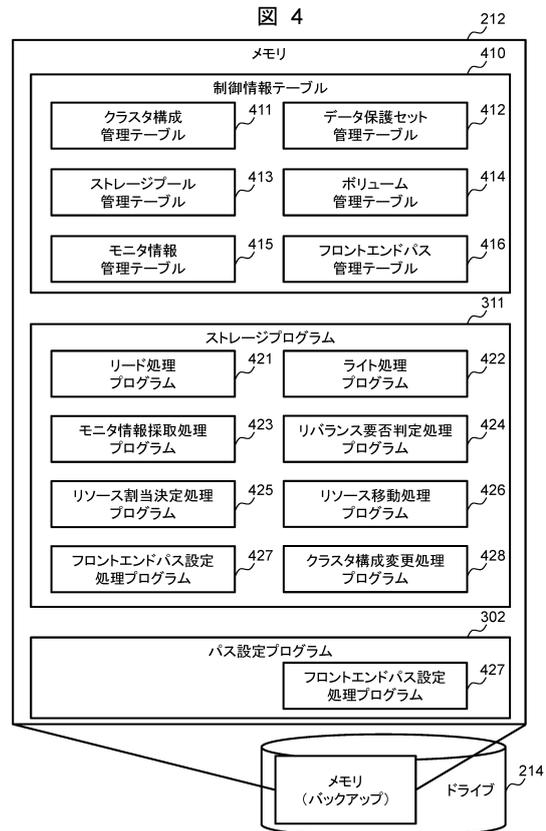
10

20

【図 3】



【図 4】



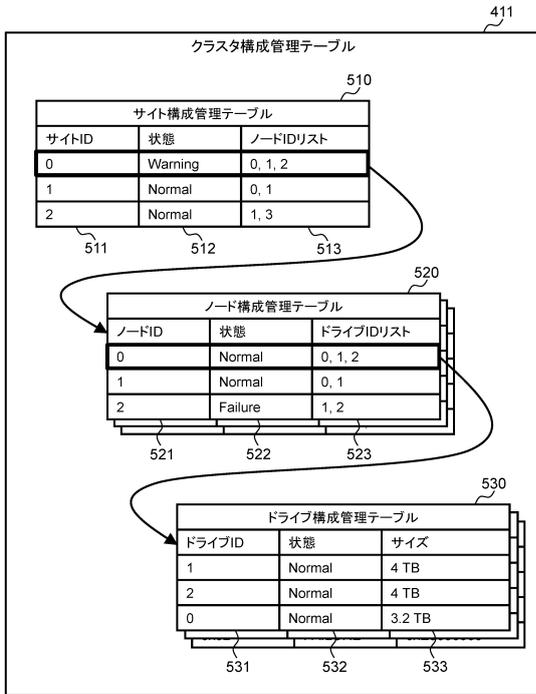
30

40

50

【図5】

図5



【図6】

図6

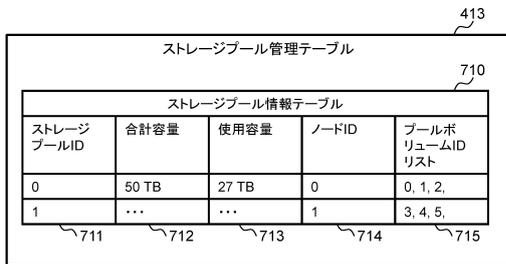


10

20

【図7】

図7



【図8】

図8



30

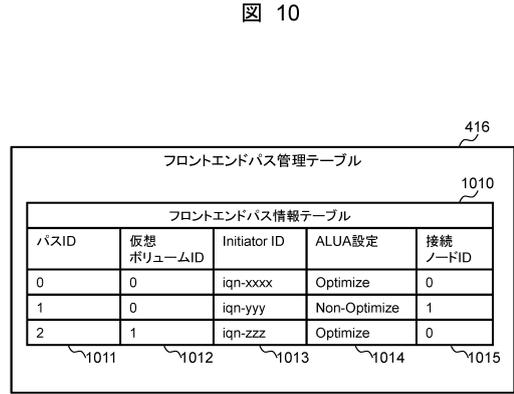
40

50

【 図 9 】



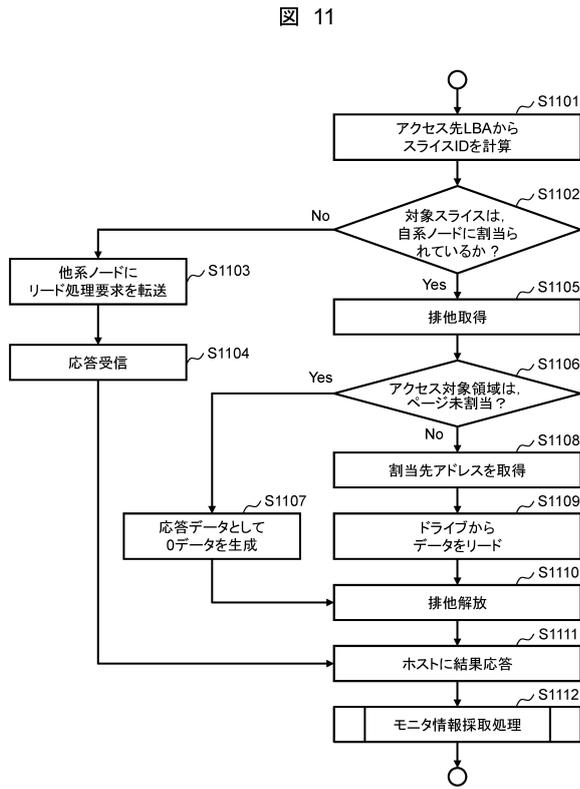
【 図 1 0 】



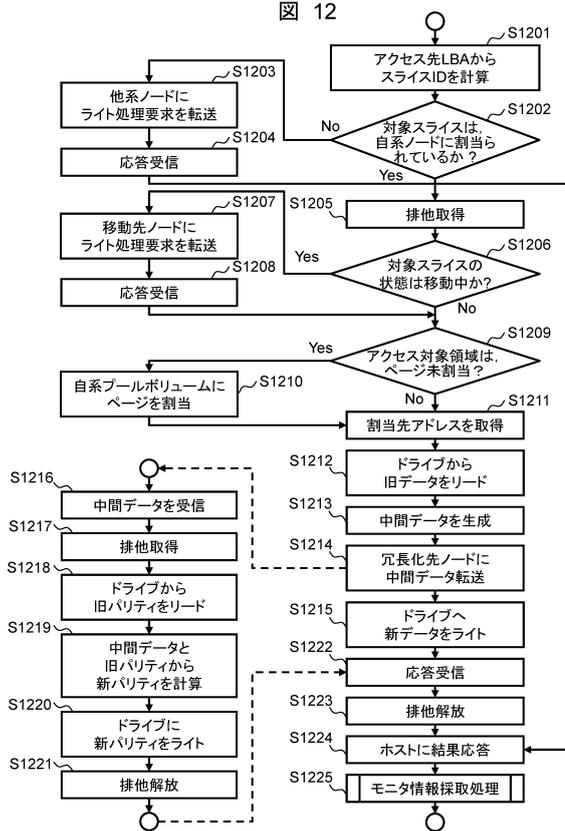
10

20

【 図 1 1 】



【 図 1 2 】



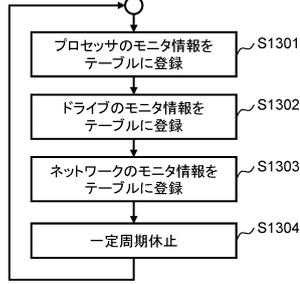
30

40

50

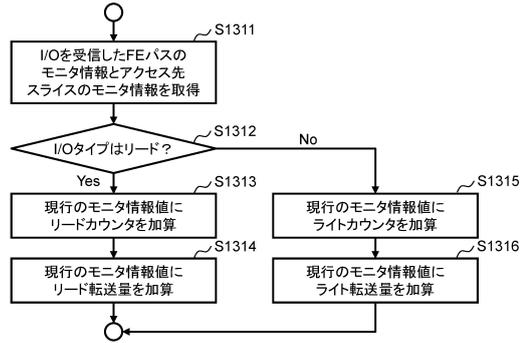
【 図 1 3 A 】

図 13A



【 図 1 3 B 】

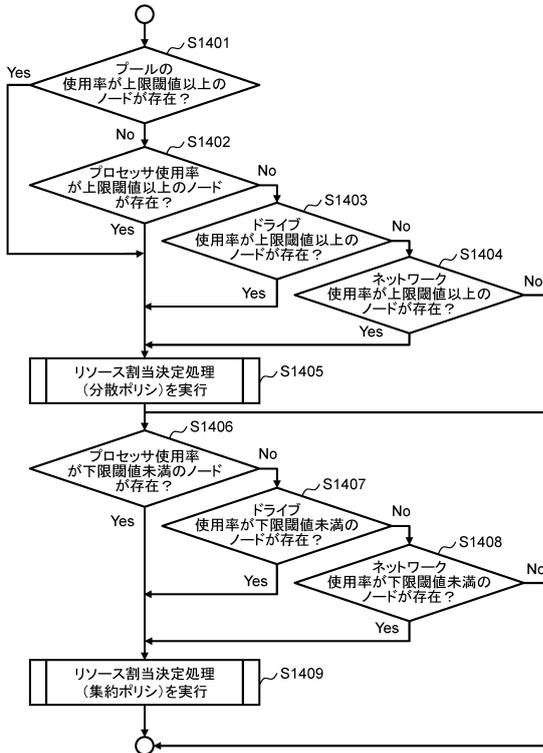
図 13B



10

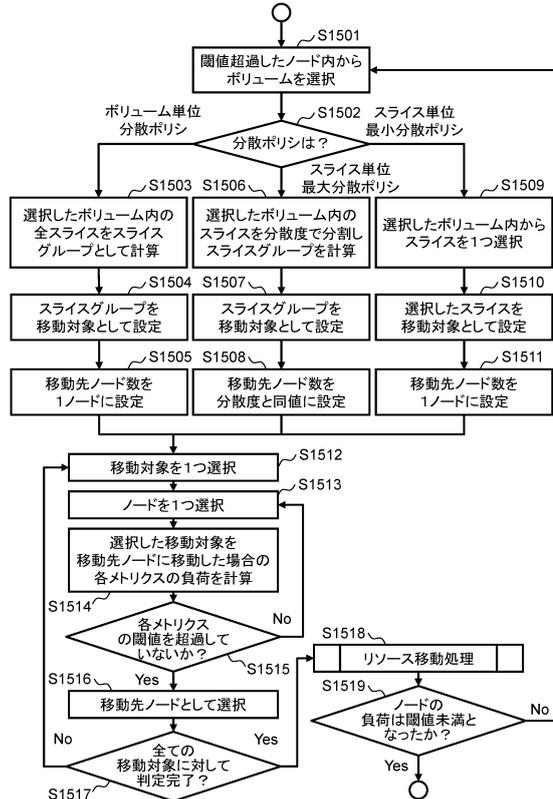
【 図 1 4 】

図 14



【 図 1 5 】

図 15



20

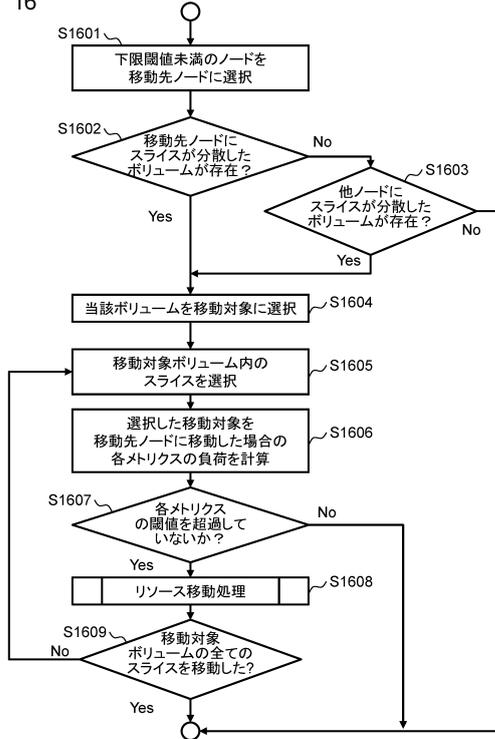
30

40

50

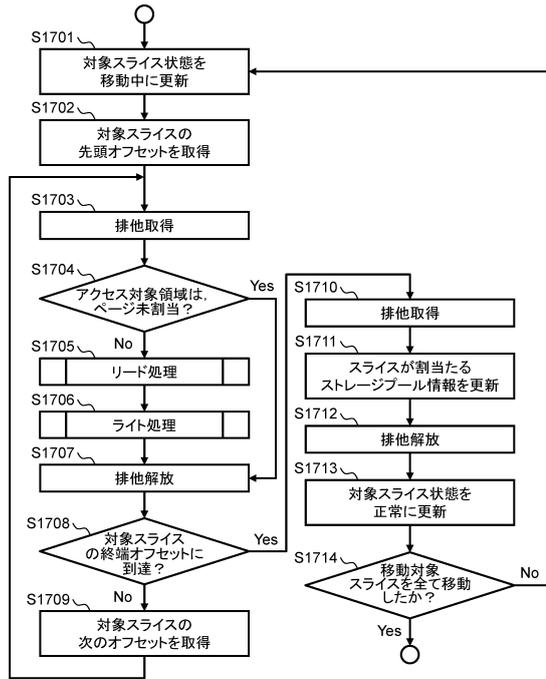
【図16】

図16



【図17】

図17

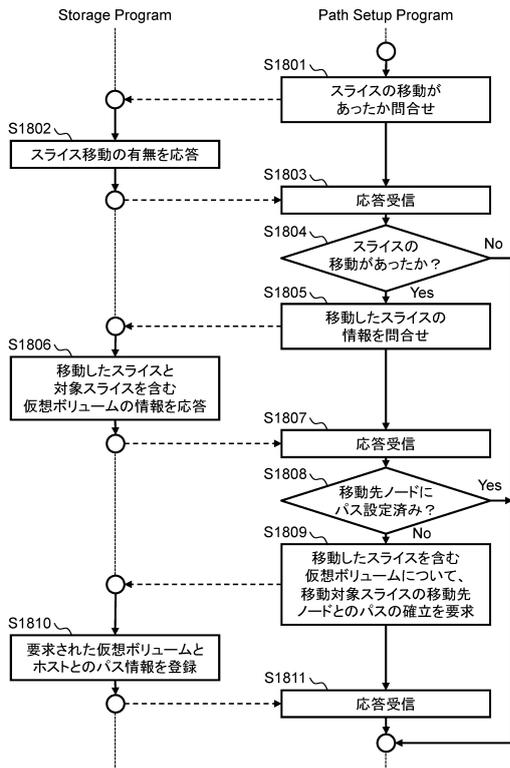


10

20

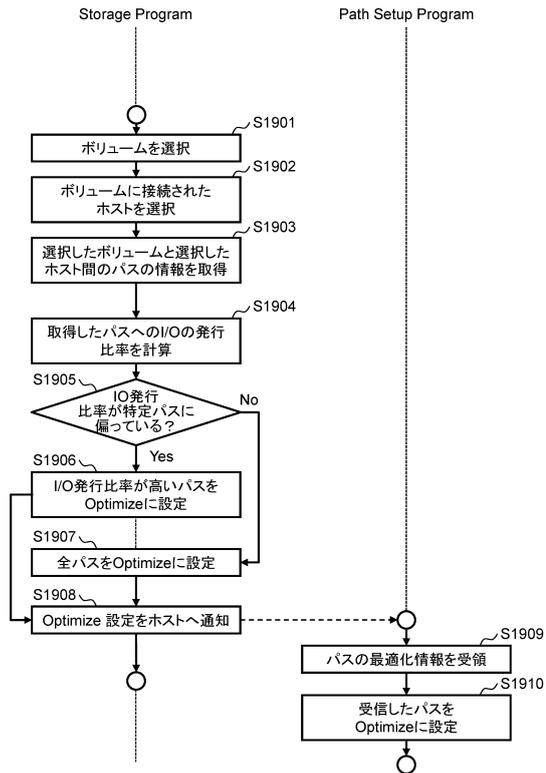
【図18】

図18



【図19】

図19



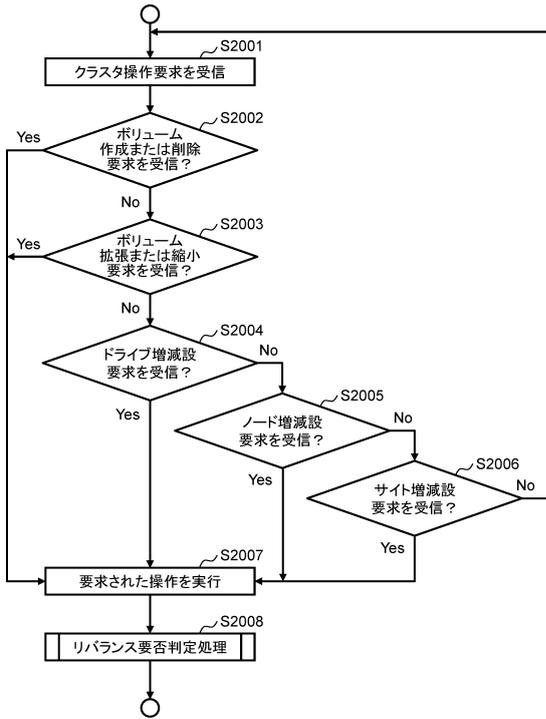
30

40

50

【図 20】

図 20



【図 21 A】

図 21A

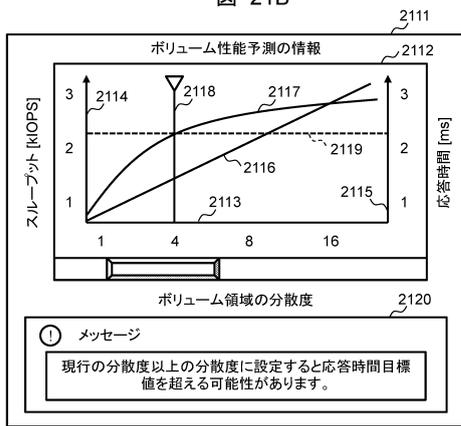


10

20

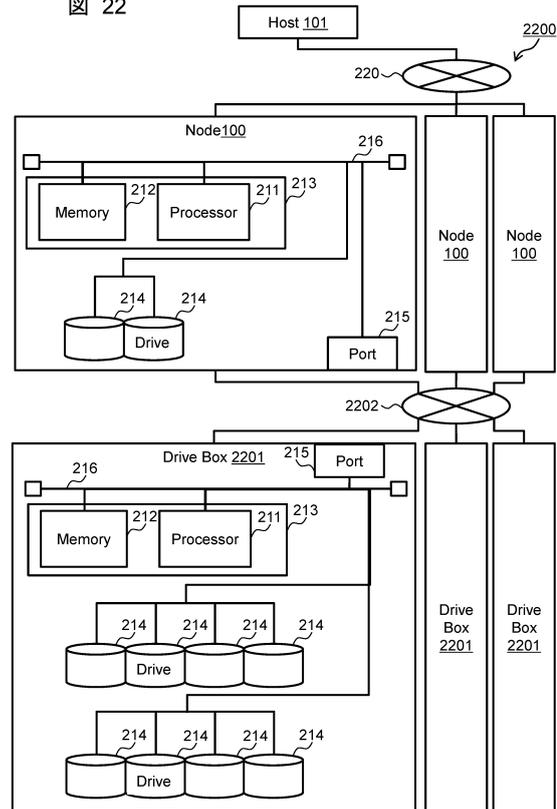
【図 21 B】

図 21B



【図 22】

図 22



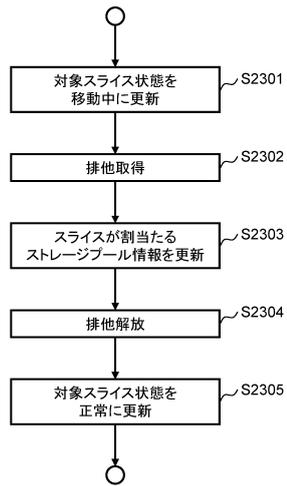
30

40

50

【 図 2 3 】

図 23



10

20

30

40

50

---

フロントページの続き

(51)国際特許分類

F I  
G 0 6 F 13/10 3 4 0 A

東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内

審査官 打出 義尚

(56)参考文献 特開2010-277289(JP,A)  
特開2016-24563(JP,A)  
特開2020-46929(JP,A)  
米国特許出願公開第2019/0265915(US,A1)

(58)調査した分野 (Int.Cl., DB名)  
G 0 6 F 3 / 0 6  
G 0 6 F 1 3 / 1 0  
G 0 6 F 1 1 / 0 7  
G 0 6 F 1 1 / 3 4