

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-272138

(P2010-272138A)

(43) 公開日 平成22年12月2日(2010.12.2)

(51) Int.Cl. F I テーマコード(参考)  
**G06F 3/06 (2006.01)** G06F 3/06 301Z 5B065

審査請求 有 請求項の数 23 O L (全 42 頁)

(21) 出願番号	特願2010-186320 (P2010-186320)	(71) 出願人	506051681 コンベレント・テクノロジーズ アメリカ合衆国ミネソタ州55344, エ デン・プレイリー, ヴァリー・ビュー・ロ ード 12982
(22) 出願日	平成22年8月23日(2010.8.23)	(74) 代理人	100140109 弁理士 小野 新次郎
(62) 分割の表示	特願2006-523434 (P2006-523434) の分割	(74) 代理人	100089705 弁理士 社本 一夫
原出願日	平成16年8月13日(2004.8.13)	(74) 代理人	100075270 弁理士 小林 泰
(31) 優先権主張番号	60/495,204	(74) 代理人	100080137 弁理士 千葉 昭男
(32) 優先日	平成15年8月14日(2003.8.14)	(74) 代理人	100096013 弁理士 富田 博行
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	10/918,329		
(32) 優先日	平成16年8月13日(2004.8.13)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 仮想ディスク・ドライブのシステムおよび方法

(57) 【要約】

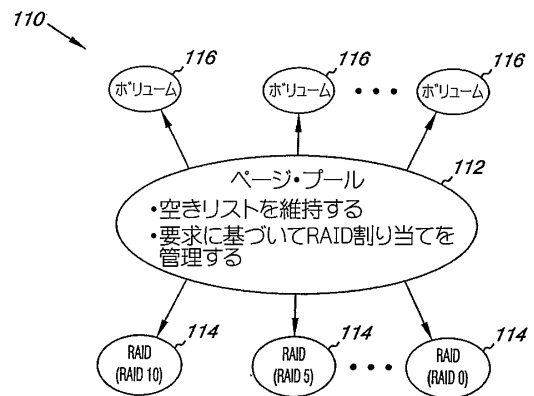
【課題】

動的にデータを割り当てる改善されたシステムを提供する。

【解決手段】

システムは、ストレージのプール(RAIDのフリーリストを維持するストレージのページプール等)又はRAIDのヌルリストを維持するディスクストレージブロックのマトリックスを持つRAIDサブシステムと、ディスク記憶システム制御装置を持つディスクマネージャを含む。サブシステム及びマネージャは、ストレージのプール及びディスクドライブにわたりデータを動的に割り当て、また、ディスクドライブが更に必要かを判定して必要な場合に通知を送る。

【選択図】 図2



**【特許請求の範囲】****【請求項 1】**

ストレージのプールにおいてデータを動的に割り当てる能力を有するディスク・ドライブ・システムであって、

前記ストレージのプールを有する R A I D サブシステムと、

少なくとも 1 つのディスク・ストレージ・システム・コントローラを有するディスク・マネージャと

を含み、

前記 R A I D サブシステムおよび前記ディスク・マネージャは、R A I D - ツー - ディスク・マッピングに基づいて、前記ストレージのプールおよび複数のディスク・ドライブにわたってデータを動的に割り当てる、

システム。

**【請求項 2】**

請求項 1 に記載のシステムであって、前記 R A I D サブシステムおよび前記ディスク・マネージャは、更に別のディスク・ドライブが必要であるかを判定し、前記更に別のディスク・ドライブが必要である場合には通知が送られる、システム。

**【請求項 3】**

請求項 1 に記載のシステムであって、前記ディスク・マネージャは複数のディスク・ストレージ・システム・コントローラを管理する、システム。

**【請求項 4】**

請求項 3 に記載のシステムであって、動作されるディスク・ストレージ・システム・コントローラの障害をカバーするために複数の冗長なディスク・ストレージ・システム・コントローラを更に備える、システム。

**【請求項 5】**

請求項 1 に記載のシステムであって、前記 R A I D サブシステムは、R A I D - 0、R A I D - 1、R A I D - 5、R A I D - 10 などのような R A I D タイプのうちの少なくとも 1 つを含む組合せを更に備える、システム。

**【請求項 6】**

請求項 5 に記載のシステムであって、R A I D - 3、R A I D - 4、R A I D - 6、および R A I D - 7 を含む R A I D タイプを更に備える、システム。

**【請求項 7】**

請求項 1 に記載のシステムであって、前記ストレージのプールは、R A I D のフリー・リストを維持するストレージのページ・プールである、システム。

**【請求項 8】**

請求項 1 に記載のシステムであって、前記ストレージのプールは、R A I D のヌル・リストを維持するディスク・ストレージ・ブロックのマトリックスである、システム。

**【請求項 9】**

動的なデータの割り当ての方法であって、

ストレージのプールを形成する R A I D サブシステムのディスク空間のデフォルトのサイズを定義するステップと、

データを書き込むステップと、

前記データを前記ストレージのプールにおいて割り当てるステップと、

前記 R A I D サブシステムの前記ディスク空間の占有率を、前記 R A I D サブシステムの前記ディスク空間の占有率の履歴に基づいて、決定するステップと、

更に別のディスク・ドライブが必要であるかを判定するステップと、

前記更に別のディスク・ドライブが必要である場合には、通知を前記 R A I D サブシステムへ送るステップと

を備える方法。

**【請求項 10】**

請求項 9 に記載の方法であって、前記通知を電子メールによって送ることを更に含む、

10

20

30

40

50

方法。

【請求項 1 1】

請求項 9 に記載の方法であって、前記ディスク・ストレージ・ブロックのサイズを、デフォルトとして設定し、ユーザによって変更可能であるように設定するステップを更に備える、方法。

【請求項 1 2】

請求項 9 に記載の方法であって、RAID のフリー・リストを維持するステップを更に備え、前記ストレージのプールは、前記 RAID のフリー・リストを維持するストレージのページ・プールである、方法。

【請求項 1 3】

請求項 9 に記載の方法であって、RAID のヌル・リストを維持するステップを更に備え、前記ストレージのプールは、前記 RAID のヌル・リストを維持するディスク・ストレージ・ブロックのマトリックスである、方法。

【請求項 1 4】

データ・インスタント・リプレイの方法であって、  
ストレージのプールを形成する RAID サブシステムのディスク空間のデフォルトのサイズを定義するステップと、

前記ストレージのプールのスナップショットを所定の時間間隔で自動的に生成するステップと、

前記スナップショットまたはデルタのアドレス・インデックスを前記ストレージのプールに保存して、前記ストレージのプールの前記スナップショットまたはデルタが、保存した前記アドレス・インデックスを介して直ちにを見つけ出せるようにするステップと

を備える方法。

【請求項 1 5】

請求項 1 4 に記載の方法であって、前記 RAID サブシステムの前記スナップショットを所定の時間間隔で自動的に生成する前記ステップは、前記 RAID サブシステムのスナップショットをユーザの定義する時間間隔で自動的に生成するステップを含む、方法。

【請求項 1 6】

請求項 1 5 に記載の方法であって、前記時間間隔は数分毎から数時間毎の範囲にある、方法。

【請求項 1 7】

請求項 1 4 に記載の方法であって、前記スナップショットはローカルの RAID サブシステムに保存される、方法。

【請求項 1 8】

請求項 1 4 に記載の方法であって、前記スナップショットはリモートの RAID サブシステムで保存され、重大なシステム・クラッシュが発生した場合に、データ・インテグリティに対する影響を受けないようにし、前記データが直ちにリカバリできるようにする、方法。

【請求項 1 9】

請求項 1 4 に記載の方法であって、前記ストレージのプールは、RAID のフリー・リストを維持するストレージのページ・プールである、方法。

【請求項 2 0】

請求項 1 4 に記載の方法であって、前記ストレージのプールは、RAID のヌル・リストを維持するディスク・ストレージ・ブロックのマトリックスである、方法。

【請求項 2 1】

データ・インスタント・リプレイのシステムであって、  
ストレージのプールを有する RAID サブシステムと、  
少なくとも 1 つのディスク・ストレージ・システム・コントローラを有するディスク・マネージャと  
を含み、

10

20

30

40

50

前記 R A I D サブシステムおよび前記ディスク・マネージャは、R A I D - ツー - ディスク・マッピングに基づいて、データを前記ストレージのプールおよび複数のディスク・ドライブにわたって動的に割り当てる、システム。

【請求項 2 2】

請求項 2 1 に記載のシステムであって、前記ストレージのプールのスナップショットを所定の時間間隔で自動的に生成するため、およびスナップショットまたはデルタのアドレス・インデックスを前記ストレージのプールに保存して、前記ストレージのプールの前記スナップショットまたはデルタが、保存した前記アドレス・インデックスを介して直ちにを見つけ出せるようにするためのディスク・ストレージ・システム・コントローラを更に備えるシステム。

10

【請求項 2 3】

請求項 2 2 に記載のシステムであって、前記ディスク・ストレージ・システム・コントローラは、前記ストレージのプールの前記スナップショットからのデータ使用の頻度を監視し、使用またはアクセスの頻度の低いデータを費用の低い R A I D サブシステムへ移動させるエージング規則を適用する、システム。

【請求項 2 4】

請求項 2 3 に記載のシステムであって、費用の低い R A I D サブシステム内のデータがより頻繁に使用され始めると、前記コントローラは、前記データを、より費用の高い R A I D サブシステムへ移動させ、それにより、前記ディスク・ドライブ・システムのコスト

20

【請求項 2 5】

請求項 2 1 に記載のシステムであって、前記ストレージのプールは、R A I D のフリー・リストを維持するストレージのページ・プールである、システム。

【請求項 2 6】

請求項 2 1 に記載のシステムであって、前記ストレージのプールは、R A I D のヌル・リストを維持するディスク・ストレージ・ブロックのマトリックスである、システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般には、ディスク・ドライブのシステムおよび方法に関し、より詳細には、動的データ割り当てやディスク・ドライブ仮想化などの機能を有するディスク・ドライブ・システムに関する。

30

【背景技術】

【0002】

既存のディスク・ドライブ・システムは、仮想ボリューム (virtual volume) ・データ・ストレージ空間が、データの格納のための特定のサイズおよび場所を有する物理ディスクと静的に関連するように、設計されている。こうしたディスク・ドライブ・システムでは、データを格納するために、データ・ストレージ空間の仮想ボリュームの厳密な場所およびサイズを知ること及び監視 / 制御を行うことが必要である。更に、このシステムでは、しばしば、より大きなデータ・ストレージ空間が必要であり、そのため、R A I D 装置が追加される。しかし、しばしば、こうした追加の R A I D 装置は、高価であり、また、余分なデータ・ストレージ空間が実際に必要となるまでは、必要ではない。

40

【0003】

図 1 4 の A に、データの格納、読み出し / 書き込み、および / またはリカバリ (復旧) のための特定のサイズおよび場所を有する物理ディスクと関連する仮想ボリューム・データ・ストレージ空間を有する従来の既存のディスク・ドライブ・システムを示す。このディスク・ドライブ・システムでは、データの割り当てを、データ・ストレージ空間の仮想ボリュームの特定の場所およびサイズに基づいて、静的に行う。その結果、空のデータ・ストレージ空間は使用されず、このシステムでデータを格納、読み出し / 書き込み、およ

50

び/またはリカバリするために、余分であり時には高価なデータ・ストレージ装置、例えば R A I D 装置が前もって獲得される。こうした余分なデータ・ストレージ空間は、必要でなく、かつ/または、後の時まで使用されない可能性がある。

【発明の概要】

【0004】

従って、改善されたディスク・ドライブのシステムおよび方法が必要とされている。効率のよい、動的なデータ割り当てならびにディスク・ドライブの空間および時間管理のシステムおよび方法が更に必要とされている。

【0005】

本発明は、動的にデータを割り当てる能力を有する改善されたディスク・ドライブのシステムおよび方法を提供する。このディスク・ドライブ・システムは、ディスク・ストレージ・ブロックのマトリックスを有する R A I D サブシステムと、少なくとも1つのディスク・ストレージ・システム・コントローラを有するディスク・マネージャとを含むことができる。R A I D サブシステムおよびディスク・マネージャは、R A I D からディスクへのマッピング (RAID-to-disk mapping、R A I D - ツー - ディスク・マッピング) に基づいて、ディスク・ストレージ・ブロックのマトリックスおよび複数のディスク・ドライブを通して、データを動的に割り当てる。R A I D サブシステムおよびディスク・マネージャは、更に別のディスク・ドライブが必要かどうかを判定し、更に別のディスク・ドライブが必要である場合には、通知を送る。動的なデータの割り当て (アロケーション) により、ユーザは、ディスク・ドライブを、後の時間に、必要となったときに得ることができる。また、動的データ割り当てにより、ディスク・ストレージ・ブロックの仮想ボリューム・マトリックスまたはプールのスナップショット/ポイント・イン・タイム (時) ・コピーの効率のよいデータ格納や、データのバックアップ、リカバリなどのためのインスタント・データ・リプレイおよびデータ・インスタント・フュージョン (fusion) や、リモート・データ・ストレージや、データ・プログレッション (Data Progression) などが、可能になる。また、データ・プログレッションにより、より安価なディスク・ドライブは、後の時間に購入されるため、延期が可能になる。

【0006】

一実施形態では、仮想ボリュームまたはディスク・ストレージ・ブロックのマトリックスまたはプールが、物理ディスクと関連付けるために提供される。仮想ボリュームまたはディスク・ストレージ・ブロックのマトリックスまたはプールの監視/制御は、複数のディスク・ストレージ・システム・コントローラによって動的に行われる。一実施形態では、それぞれの仮想ボリュームのサイズはデフォルトまたはユーザの予め定義したものとすることができ、それぞれの仮想ボリュームの場所はデフォルトでヌル (null) となっている。仮想ボリュームは、データが割り当てられるまでヌルである。データは、マトリックスまたはプールの任意のグリッドに割り当てることができる (例えば、データがそのグリッドに割り当てられた後は、そのグリッド内の「ドット」 (dot) になる)。そのデータが削除された後、その仮想ボリュームは再び使用可能になり「null」として示される。このため、余分なデータ・ストレージ空間や、時には高価なデータ・ストレージ装置、例えば R A I D 装置を、後の時間に、必要となるたびに得ることができる。

【0007】

一実施形態では、1つのディスク・マネージャが複数のディスク・ストレージ・システム・コントローラを管理することができ、複数の冗長ディスク・ストレージ・システム・コントローラを実装して、動作させているディスク・ストレージ・システム・コントローラの障害をカバーすることができる。

【0008】

一実施形態では、R A I D サブシステムは、R A I D - 0、R A I D - 1、R A I D - 5、R A I D - 10 などの R A I D タイプのうち少なくとも1つからなる組合せを含む。代替の R A I D サブシステムでは、R A I D - 3、R A I D - 4、R A I D - 6、R A I D - 7 などのような他の R A I D タイプを使用することも理解されよう。

10

20

30

40

50

## 【 0 0 0 9 】

また、本発明は、動的データ割り当ての方法を提供するものであり、この方法は、論理ブロックまたはディスク・ストレージ・ブロックのデフォルトのサイズを与えて、RAIDサブシステムのディスク空間がディスク・ストレージ・ブロックのマトリックスを形成するようにするステップと、データを書き込み、そのデータをディスク・ストレージ・ブロックのマトリックスに割り当てるステップと、RAIDサブシステムのディスク空間の占有率の履歴に基づいて、RAIDサブシステムのディスク空間の占有率の決定するステップと、更に別のディスク・ドライブが必要かどうかを決定するステップと、更に別のディスク・ドライブが必要である場合に、RAIDサブシステムに通知を送るステップとを含む。一実施形態では、通知は電子メールによって送られる。

10

## 【 0 0 1 0 】

本発明のディスク・ドライブ・システムの利点の1つは、RAIDサブシステムが、仮想的な数のディスクにわたってRAID技法を使用できることである。残りのストレージ空間は自由に使用可能である。RAIDサブシステムのストレージ空間を監視し、そのストレージ空間の占有率を決定することにより、ユーザは、高価であるが購入時点では用途のない大量のドライブを得なくて済む。こうして、ドライブの追加を、ストレージ空間の増大する需要を満たすために実際に必要であるときに行うことにより、ディスク・ドライブの全体のコストが著しく低減される。その一方で、ドライブの使用の効率は著しく改善される。

20

## 【 0 0 1 1 】

本発明の別の利点は、ディスク・ストレージ・システム・コントローラが、特定のコンピュータ・ファイル・システムだけでなく、どのようなコンピュータ・ファイル・システムに対しても普遍的であることである。

## 【 0 0 1 2 】

また、本発明は、データ・インスタント・リプレイ (data instant replay) の方法を提供する。一実施形態では、データ・インスタント・リプレイの方法は、論理ブロックまたはディスク・ストレージ・ブロックのデフォルトのサイズを与えて、RAIDサブシステムのディスク空間がストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスを形成するようにするステップと、ストレージのページ・プールのボリュームのスナップショットまたはディスク・ストレージ・ブロックのマトリックスのスナップショットを、所定の時間間隔で自動的に生成するステップと、スナップショットまたは変分 (delta、デルタ) のアドレス・インデックスを、ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスに保存して、ディスク・ストレージ・ブロックのマトリックスのスナップショットまたはデルタを、その保存したアドレス・インデックスを介して直ちにを見つけ出せるようにするステップとを含む。

30

## 【 0 0 1 3 】

このデータ・インスタント・リプレイの方法では、RAIDサブシステムのスナップショットの自動生成を、ユーザの定義する時間間隔で、または、ユーザのコンフィギュレーションした動的なタイム・スタンプ、例えば、数分または数時間ごとなどに、または、サーバの指示する時刻に、行う。システム障害またはウイルス攻撃が発生した場合、こうしたタイム・スタンプを有する仮想スナップショットにより、データ・インスタント・リプレイおよびデータ・インスタント・リカバリが、約数分または数時間などで可能になる。この技法は、インスタント・リプレイ・フュージョン (instant replay fusion) とも呼ばれる。即ち、クラッシュまたは攻撃のすぐ前のデータのフュージョンが遅れずに行われ、クラッシュまたは攻撃の前に保存したスナップショットを直ちにそれ以降の動作で使用することができる。

40

## 【 0 0 1 4 】

一実施形態では、スナップショットをローカルのRAIDサブシステムまたはリモートのRAIDサブシステムで保存し、重大なシステム・クラッシュが例えばテロリスト攻撃などが原因で生じた場合に、データの完全性 (integrity) に対しての影響を受けないよ

50

うにし、データを直ちにリカバリできるようにすることができる。

【0015】

このデータ・インスタント・リプレイの方法の別の利点は、スナップショットを、システムが動作し続けたままで、検査のために使用できることである。生(live)のデータをリアル・タイムの検査に使用することができる。

【0016】

また、本発明は、RAIDサブシステムと、少なくとも1つのディスク・ストレージ・システム・コントローラを有するディスク・マネージャを含むデータ・インスタント・リプレイのシステムを提供する。一実施形態では、RAIDサブシステムおよびディスク・マネージャは、複数のドライブのディスク空間にわたって、データを、RAID-ツリー・ディスク・マッピングに基づいて動的に割り当てるものであり、ここでRAIDサブシステムのディスク空間は、ディスク・ストレージ・ブロックのマトリックスを形成する。ディスク・ストレージ・システム・コントローラは、ディスク・ストレージ・ブロックのマトリックスのスナップショットを所定の時間間隔で自動的に生成し、スナップショットまたはデルタのアドレス・インデックスをディスク・ストレージ・ブロックのマトリックスに保存して、ディスク・ストレージ・ブロックのマトリックスのスナップショットまたはデルタを、その保存したアドレス・インデックスを介して直ちにを見つけ出せるようにする。

【0017】

一実施形態では、ディスク・ストレージ・システム・コントローラは、ディスク・ストレージ・ブロックのマトリックスのスナップショットからデータ使用の頻度を監視し、そして、使用またはアクセスの頻度の低いデータほど、より安価なRAIDサブシステムへと移動させるような、エイジング(aging、加齢)規則を適用する。同様に、安価なRAIDサブシステム中のデータが頻繁に使用され始めると、コントローラは、そのデータを、より高価なRAIDサブシステムへと移動させる。このため、ユーザは、ユーザ自身のストレージに対する必要性を満たす所望のRAIDサブシステム・ポートフォリオを選ぶことができる。従って、ディスク・ドライブ・システムのコストを著しく低減させること、およびユーザにより動的に制御することが可能となる。

【0018】

本発明のこうしたおよび他の特徴および利点は、本発明を実現するために企図された最良の形態を含む本発明の例示的な実施形態を示し説明する以下の詳細な説明から、当業者に明らかとなろう。本発明は、様々な自明な態様における変更が、すべて本発明の趣旨および範囲から逸脱することなく可能であることが理解されよう。従って、図面および詳細な説明は、本来例示的であって限定的でないと思ふべきである。

【図面の簡単な説明】

【0019】

【図1】図1は、本発明の諸原理に従ったコンピュータ環境におけるディスク・ドライブ・システムの一実施形態を示す。

【図2】図2は、本発明の諸原理に従った、ディスク・ドライブのRAIDサブシステムのためのストレージのページ・プールを有する動的なデータ割り当ての一実施形態を示す。

【図2A】図2Aは、ディスク・ドライブ・システムのRAIDサブシステムにおける従来のデータ割り当てを示す。

【図2B】図2Bは、本発明の諸原理に従った、ディスク・ドライブ・システムのRAIDサブシステムにおけるデータ割り当てを示す。

【図2C】図2Cは、本発明の諸原理に従った動的なデータ割り当ての方法を示す。

【図3】図3のAおよびBは、本発明の諸原理に従った、複数の時間間隔におけるRAIDサブシステムのディスク・ストレージ・ブロックのスナップショットの概略図であり、Cは、本発明の諸原理に従ったデータ・インスタント・リプレイの方法を示す。

【図4】図4は、本発明の諸原理に従った、RAIDサブシステムのディスク・ストレージ

10

20

30

40

50

ジ・ブロックのスナップショットを使用することによるデータ・インスタント・フュージョン機能の概略図である。

【図5】図5は、本発明の諸原理に従った、RAIDサブシステムのディスク・ストレージ・ブロックのスナップショットを使用することによるローカル/リモートのデータの複製およびインスタント・リプレイの機能の概略図である。

【図6】図6は、本発明の諸原理に従った、I/Oを行うために同じRAIDインターフェースを使用し、複数のRAID装置を1つのボリュームへと連結するスナップショットの概略図を示す。

【図7】図7は、本発明の諸原理に従った、スナップショット構造の一実施形態を示す。

【図8】図8は、本発明の諸原理に従った、PITCライフ・サイクルの一実施形態を示す。

【図9】図9は、本発明の諸原理に従った、マルチ・レベル・インデックスを有するPITCテーブル構造の一実施形態を示す。

【図10】図10は、本発明の諸原理に従った、PITCテーブルのリカバリの一実施形態を示す。

【図11】図11は、本発明の諸原理に従った、所有ページ・シーケンスおよび非所有ページ・シーケンスを有する書き込みプロセスの一実施形態を示す。

【図12】図12は、本発明の諸原理に従った、例示的なスナップショット動作を示す。

【図13A】図13Aは、データを静的に割り当てるために特定のサイズおよび場所をもつ物理ディスクと関連付けられた仮想ボリューム・データ・ストレージ空間を有する従来の既存のディスク・ドライブ・システムを示す。

【図13B】図13Bは、図13Aの従来の既存のディスク・ドライブ・システムでのボリューム論理ブロック・マッピングを示す。

【図14】図14のAは、本発明の諸原理に従った、システムにおいてデータを動的に割り当てるための、ディスク・ストレージ・ブロックの仮想ボリューム・マトリックスを有するディスク・ドライブ・システムの一実施形態を示す。図14のBは、図14のAに示すディスク・ストレージ・ブロックの仮想ボリューム・マトリックスにおける動的なデータ割り当ての一実施形態を示す。図14のCは、本発明の諸原理に従った、ストレージの仮想ボリューム・ページ・プールの一実施形態のボリューム - RAIDページ再マッピングの概略図を示す。

【図15】図15は、本発明の諸原理に従った、RAIDサブシステムの複数のディスク・ストレージ・ブロックへマッピングされる3つのディスク・ドライブの例を示す。

【図16】図16は、図15に示す3つのディスク・ドライブへ1つのディスク・ドライブを追加した後のディスク・ドライブ・ストレージ・ブロックの再マッピングの例を示す。

【図17】図17は、本発明の諸原理に従った、データ・プログレッション動作におけるアクセス可能なデータ・ページの一実施形態を示す。

【図18】図18は、本発明の諸原理に従った、データ・プログレッション・プロセスの一実施形態の流れ図を示す。

【図19】図19は、本発明の諸原理に従った、圧縮されたページ・レイアウトの一実施形態を示す。

【図20】図20は、本発明の諸原理に従った、高レベル・ディスク・ドライブ・システムにおけるデータ・プログレッションの一実施形態を示す。

【図21】図21は、本発明の諸原理に従った、サブシステムにおける外部データ・フローの一実施形態を示す。

【図22】図22は、サブシステムにおける内部データ・フローの一実施形態を示す。

【図23】図23は、コヒーレンシを独立して維持する各サブシステムの一実施形態を示す。

【図24】図24は、本発明の諸原理に従った、混合RAIDウォーターフォール・データ・プログレッションの一実施形態を示す。

10

20

30

40

50



【図 2 5】図 2 5 は、本発明の諸原理に従った、ストレージのページ・プールの複数のフリー・リストの一実施形態を示す。

【図 2 6】図 2 6 は、本発明の諸原理に従った、データベースの例の一実施形態を示す。

【図 2 7】図 2 7 は、本発明の諸原理に従った、MRI 画像の例の一実施形態を示す。

【発明を実施するための形態】

【0020】

本発明は、動的にデータを割り当てる能力を有する改善されたディスク・ドライブのシステムおよび方法を提供する。このディスク・ドライブ・システムは、RAID の空きリスト（フリー・リスト、free list）を維持するストレージのページ・プール、または代替例としてはディスク・ストレージ・ブロックのマトリックスを有する RAID サブシステムと、少なくとも 1 つのディスク・ストレージ・システム・コントローラを有するディスク・マネージャとを含むことができる。RAID サブシステムおよびディスク・マネージャは、データを、ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスおよび複数のディスク・ドライブにわたって、RAID - ツー - ディスク・マッピングに基づいて、動的に割り当てる。RAID サブシステムおよびディスク・マネージャは、更に別のディスク・ドライブが必要かどうかを判定し、更に別のディスク・ドライブが必要である場合には通知を送る。動的なデータ割り当てにより、ユーザは、ディスク・ドライブを、後にそれが必要となったときに得ることができる。また、動的なデータ割り当てにより、ディスク・ストレージ・ブロックの仮想ボリュームのマトリックスまたはプールのスナップショット / ポイント・イン・タイム・コピーの効率のよいデータ・ストレージや、データ・バックアップやリカバリなどのためのインスタント・データ・リプレイおよびデータ・インスタント・フュージョン（fusion）や、リモート・データ・ストレージや、データ・プログレッションなどが可能になる。また、データ・プログレッションにより、より安価なディスク・ドライブは後の時間に購入されるため、その延期が可能になる。

10

20

【0021】

図 1 に、本発明の諸原理に従った、コンピュータ環境 102 におけるディスク・ドライブ・システム 100 の一実施形態を示す。図 1 に示すように、ディスク・ドライブ・システム 100 は、RAID サブシステム 104 と、少なくとも 1 つのディスク・ストレージ・システム・コントローラ（図 16）を有するディスク・マネージャ 106 とを含む。RAID サブシステム 104 およびディスク・マネージャ 106 は、データを、複数のディスク・ドライブ 108 のディスク空間にわたって、RAID - ツー - ディスク・マッピングに基づいて、動的に割り当てる。更に、RAID サブシステム 104 およびディスク・マネージャ 106 は、更に別のディスク・ドライブが必要かどうかの判定を、ディスク空間全体にわたるデータ割り当てに基づいて行う能力がある。更に別のディスク・ドライブが必要である場合にはその通知がユーザへ送られて、希望される場合には更に別のディスク空間を追加できる。

30

【0022】

本発明の諸原理に従った、動的なデータ割り当て（または「ディスク・ドライブ仮想化」と呼ばれる）を有するディスク・ストレージ・システム 100 を、図 2 に一つの実施形態の形で示し、図 14 の A および図 14 の B に別の実施形態の形で示している。図 2 に示すように、ディスク・ストレージ・システム 110 は、ストレージのページ・プール 112、即ち、自由にデータを保存できるデータ・ストレージ空間のリストを含むデータ・ストレージのプールを含む。ページ・プール 112 は、RAID 装置 114 のフリー・リストを維持し、読み出し / 書き込みの割り当ての管理を、ユーザの要求に基づいて行う。ユーザの要求したデータ・ストレージ・ボリューム 116 は、ストレージ空間を得るために、ページ・プール 112 へと送られる。それぞれのボリュームは、同じまたは異なる RAID レベル、例えば、RAID 10、RAID 5、RAID 0 などで、同じまたは異なるクラスのストレージ装置を要求することができる。

40

【0023】

50

本発明の動的なデータ割り当ての別の実施形態を、図14のAおよび図14のBに示しており、この例では、複数のディスク・ストレージ・システム・コントローラ1402と、複数のディスク・ストレージ・システム・コントローラ1402により制御されるディスク・ストレージ・ブロックのマトリックス1404とを有するディスク・ストレージ・システム1400は、そのシステムにおいて本発明の諸原理に従ってデータを動的に割り当てる。仮想ボリュームまたはブロックのマトリックス1404は、物理ディスクと関連づけるように提供される。仮想ボリュームまたはブロックのマトリックス1404は、複数のディスク・ストレージ・システム・コントローラ1402によって動的に監視/制御される。一実施形態では、それぞれの仮想ボリューム1404のサイズは、例えば2メガバイトのように、予め定めることができ、それぞれの仮想ボリューム1404の場所はデフォルトでnullとなっている。仮想ボリューム1404のそれぞれは、データが割り当てられるまでnullである。データは、そのマトリックスまたはプールの任意のグリッドに割り当てることができる(例えば、データがそのグリッドに割り当てられた後は、そのグリッド内の「ドット」になる)。そのデータが削除された後は、その仮想ボリューム1404は再び使用可能になり、「null」と示される。このため、余分な、時には高価なデータ・ストレージ装置、例えばRAID装置を、後に必要となったときに得ることができる。

10

20

30

40

50

#### 【0024】

従って、RAIDサブシステムは、仮想的な数のディスクにわたってRAID技法を使用することができる。残りのストレージ空間は自由に使用可能である。RAIDサブシステムのストレージ空間を監視し、そのストレージ空間の占有率を決定することにより、ユーザは、高価であるのに購入時点では用途のない大量のドライブを得なくても済む。このように、ドライブの追加を、ストレージ空間の増大する需要を満たすために実際に必要であるときに行うことにより、ディスク・ドライブの全体のコストが著しく低減される。一方、ドライブの使用の効率も著しく改善される。

#### 【0025】

また、本発明のディスク・ドライブ・システムの動的なデータ割り当てにより、ストレージの仮想ボリューム・ページ・プール、または、ディスク・ストレージ・ブロックの仮想ボリューム・マトリックスの、スナップショット/ポイント・イン・タイム・コピーの効率のよいデータ・ストレージや、データ・リカバリおよびリモート・データ・ストレージのためのインスタント・データ・リプレイおよびデータ・インスタント・フュージョンや、データ・プログレッションが可能になる。

#### 【0026】

ディスク・ドライブ・システム100における動的データ割り当てのシステムおよび方法ならびにその実装形態の結果として得られる上述の特徴および利点を、以下で詳細に論じる。

#### 【0027】

##### 動的なデータ割り当て

図2Aに、ディスク・ドライブ・システムのRAIDサブシステムにおける従来のデータ割り当てを示し、ここでは、空けられたデータ・ストレージ空間は専属(captive)であり、データ・ストレージ用に割り当てることができない。

#### 【0028】

図2Bには、本発明の諸原理に従ったディスク・ドライブ・システムのRAIDサブシステムにおけるデータ割り当てを示し、ここでは、データ・ストレージ用に使用可能である空けられたデータ・ストレージは、一緒にまとめられて、ページ・プール、例えば、本発明の一実施形態では1つのページ・プール、を形成する。

#### 【0029】

図2Cに、本発明の諸原理に従った動的なデータ割り当ての方法200を示す。動的なデータの割り当ての方法200は、論理ブロックまたはディスク・ストレージ・ブロックのデフォルトのサイズを定義して、RAIDサブシステムのディスク空間がディスク・ス

トレージ・ブロックのマトリックスを形成するようにするステップ202と、データを書き込み、そのデータを、マトリックスの、「null」を示すディスク・ストレージ・ブロックに割り当てるステップ204とを含む。この方法は、更に、RAIDサブシステムのディスク空間の占有率を、そのRAIDサブシステムのディスク空間の占有率履歴に基づいて決定するステップ206と、更に別のディスク・ドライブが必要かどうかを判定し、必要な場合に、RAIDサブシステムへ通知を送るステップ208とを含む。一実施形態では、通知は電子メールによって送られる。更に、ディスク・ストレージ・ブロックのサイズはデフォルトに設定し、ユーザにより変更可能とすることができる。

【0030】

一実施形態では、動的なデータ割り当ては、時には「仮想化」または「ディスク空間仮想化」と呼ばれるものであり、数多くの読み出しおよび書き込みの要求を秒単位で効率よく取り扱う。このアーキテクチャでは、キャッシュ・サブシステムを直接に呼び出す割り込みハンドラが必要であることもある。動的なデータの割り当てでは、要求（リクエスト）をキューに入れないので要求の最適化はなされないが、一度に数多くの未解決（pending）の要求を有し得る。

【0031】

また、動的なデータの割り当てでは、データのインテグリティを維持し、データの内容をどのようなコントローラの障害に対しても保護することができる。そうするために、動的なデータの割り当てでは、状態情報を信頼性をもって記憶するためにRAID装置へ書き込む。

【0032】

動的なデータの割り当てでは、更に、読み出しおよび書き込みの要求の順序を維持し、読み出しまたは書き込みの要求を、要求を受け取った順序で完了させることができる。動的なデータの割り当ては、最大のシステム可用性（availability）を提供し、異なる地理的場所へのデータのリモートの複製をサポートする。

【0033】

更に、動的なデータの割り当ては、データ破壊からのリカバリ機能を提供する。スナップショットにより、ユーザは、過去におけるディスクの状態を閲覧することができる。

【0034】

動的なデータの割り当てでは、RAID装置を管理し、大規模な装置を作成し拡張するためのストレージ抽象化を提供する。

【0035】

動的なデータの割り当てでは、サーバに対して仮想ディスク装置を提示するものであり、この装置がボリューム（volume）と呼ばれる。サーバに対して、ボリュームは同じ働きをする。ボリュームは、シリアル番号に関して異なる情報を返し得るが、本質的にはディスク・ドライブのような働きをする。ボリュームは、より大きい動的なボリューム装置を作成するために、複数のRAID装置のストレージ抽象化を提供する。ボリュームは複数のRAID装置を含み、ディスク空間の効率的な使用を可能にする。

【0036】

図21に、従来の既存のボリューム論理ブロック・マッピングを示す。図14のCには、本発明の諸原理に従った、ストレージの仮想ボリュームのページ・プールの一実施形態のボリューム - RAIDページ再マッピングを示す。それぞれのボリュームは1組のページ、例えば、1、2、3などへと分割され、それぞれのRAIDは1組のページへと分割される。ボリュームのページ・サイズと、RAIDのページ・サイズとは、一実施形態では同じとすることができる。従って、本発明のボリューム - RAIDページ再マッピングの一例では、RAID - 2を使用するページ#1がRAIDページ#1へとマッピングされる。

【0037】

動的なデータの割り当てでは、ボリュームのデータ・インテグリティを維持する。データは、ボリュームへと書き込まれ、サーバへの確認が行われる。データ・インテグリティ

は、コントローラの障害を通して、スタンド・アロンや冗長なものを含めての様々なコントローラのコンフィギュレーションに適用される。コントローラの障害としては、電源障害、パワー・サイクル (power cycle)、ソフトウェア例外、およびハード・リセットがある。動的なデータの割り当てでは、一般には、RAIDでカバーされるディスク・ドライブ障害を取り扱わない。

#### 【0038】

動的なデータの割り当ては、コントローラに対する最も高レベルのデータ抽象化を提供する。動的なデータの割り当ては、フロント・エンドから要求を受理し、最終的にはRAID装置を使用してデータをディスクへ書き込む。

#### 【0039】

動的なデータの割り当ては、幾つもの内部サブシステムを含む。

- ・キャッシュ - ボリュームへの読み出しおよび書き込みの動作を、サーバへの迅速な応答時間を提供し、データ・プラグインに対して書き込みをまとめることによって、スムーズにする。

- ・コンフィギュレーション - データ割り当てオブジェクトの作成、削除、取り出し、および変更の方法を含む。より高レベルのシステム・アプリケーションに対するツールボックスを作成するためのコンポーネントを提供する。

- ・データ・プラグイン - ボリュームの読み出しおよび書き込みの要求を、ボリューム・コンフィギュレーションに応じて、様々なサブシステムへ配布する。

- ・RAIDインターフェース - より大きいボリュームを作成するためのRAID装置抽象化を、ユーザおよび他の動的なデータの割り当てサブシステムに提供する。

- ・コピー/ミラー/スワップ - ボリュームのデータを、ローカルおよびリモートのボリュームへと複製する。一実施形態では、サーバの書き込まれたブロックのコピーのみを行うことができる。

- ・スナップショット - データの増分的 (incremental) ボリューム・リカバリを提供する。過去のボリューム状態の閲覧ボリューム (View Volume) を即座に作成する。

- ・プロキシ・ボリューム - リモートの宛先ボリュームへの要求の通信を行えるようにして、リモート複製 (Remote Replication) をサポートする。

- ・請求 (billing) - 割り当てたストレージ、活動、パフォーマンス、およびデータのリカバリに対してユーザに料金請求する機能。

#### 【0040】

また、動的なデータの割り当てでは、エラーおよびコンフィギュレーションにおける顕著な変更のログ記録も行う。

#### 【0041】

図21に、サブシステムにおける外部データ・フローの一実施形態を示す。外部要求はフロント・エンドから来る。要求には、ボリューム情報取得 (get volume information)、読み出し (read)、および書き込み (write) がある。すべての要求はボリュームIDをもつ。ボリューム情報は、ボリューム・コンフィギュレーション・サブシステムによって取り扱われる。読み出しおよび書き込みの要求はLBAを含む。また、書き込み要求はデータを含む。

#### 【0042】

動的なデータの割り当ては、ボリューム・コンフィギュレーションに応じて、要求を幾つもの外部レイヤへ渡す。リモート複製は、要求をフロント・エンドへ渡すが、あて先はリモートの宛先ボリュームである。RAIDインターフェースは要求をRAIDへ渡す。コピー/ミラー/スワップは、要求を、宛先ボリュームに対する動的データ割り当てへと戻す。

#### 【0043】

図22に、サブシステムにおける内部データ・フローの一実施形態を示す。内部データ・フローはキャッシュを行うことから始まる。キャッシュを行うと、書き込み要求がキャッシュに入れられるか、または、その要求が直接にデータ・プラグインへ渡される。キャ

10

20

30

40

50

ッシュは、フロント・エンドのHBA装置からの直接DMAをサポートする。要求は迅速に完了され、応答がサーバに返される。データ・プラグイン・マネージャが、キャッシュより下の要求フローの中心である。それぞれのボリュームに対して、データ・プラグイン・マネージャは、要求ごとに、登録されたサブシステム・オブジェクトを呼び出す。

【0044】

データ・インテグリティに影響を及ぼす動的データ割り当てサブシステムでは、コントローラの一貫性（コヒーレンシ）に対するサポートが必要となることがある。図23に示すように、それぞれのサブシステムは、独立してコヒーレンシを維持する。コヒーレンシの更新により、コヒーレンシ・リンクを通じてのデータ・ブロックのコピーを回避する。キャッシュのコヒーレンシでは、データをピア（peer）コントローラにコピーする必要がある。 10

【0045】

ディスク・ストレージ・システム・コントローラ

図14のAに、本発明の諸原理に従った、システムにおいて動的にデータを割り当てるための、複数のディスク・ストレージ・システム・コントローラ1402と、複数のディスク・ストレージ・システム・コントローラ1402により制御されるディスク・ストレージ・ブロックまたは仮想ボリュームのマトリックス1404とを有するディスク・ストレージ・システム1400を示す。図14のBには、ディスク・ストレージ・ブロックまたは仮想ボリュームの仮想ボリューム・マトリックス1404における動的なデータの割り当ての一実施形態を示す。 20

【0046】

ある動作では、ディスク・ストレージ・システム1400は、ディスク・ストレージ・ブロックまたは仮想ボリュームのマトリックス1404のスナップショットを、所定の時間間隔で時間的に生成し、スナップショットまたはデルタのアドレス・インデックスを、ディスク・ストレージ・ブロックまたは仮想ボリュームのマトリックス1404に保存して、ディスク・ストレージ・ブロックまたは仮想ボリュームのマトリックス1404のスナップショットまたはデルタが、保存したアドレス・インデックスによって直ちに見つけ出せるようにする。

【0047】

更にある動作では、ディスク・ストレージ・システム・コントローラ1402は、ディスク・ストレージ・ブロックのマトリックス1404のスナップショットからデータ使用の頻度を監視し、そして、使用またはアクセスの頻度の低いデータほど、より安価なRAIDサブシステムへと移動させるようなエージング規則を適用する。同様に、より安価なRAIDサブシステム中のデータがより頻繁に使用され始めると、コントローラは、そのデータを、より高価なRAIDサブシステムへと移動させる。従って、ユーザは、ユーザ自身のストレージの必要性を満たすように所望のRAIDサブシステム・ポートフォリオを選ぶことができる。従って、ディスク・ドライブ・システムのコストを著しく低減させ、また、ユーザによる動的な制御を可能とする。 30

【0048】

RAID - ツー - ディスク・マッピング (RAID-to-Disk Mapping)

RAIDサブシステムおよびディスク・マネージャは、データを、複数のディスク・ドライブのディスク空間にわたって、RAID - ツー - ディスク・マッピングに基づいて、動的に割り当てる。一実施形態では、RAIDサブシステムおよびディスク・マネージャは、更に別のディスク・ドライブが必要かどうかを判定し、更に別のディスク・ドライブが必要である場合には通知を送る。 40

【0049】

図15に、本発明の諸原理に従った、RAID - 5サブシステム1500内の複数のディスク・ストレージ・ブロック1502 ~ 1512へとマッピングされる3つのディスク・ドライブ108（図1）の例を示す。

【0050】

10

20

30

40

50

図16は、図15に示す3つのディスク・ドライブ108にディスク・ドライブ1602を追加した後の、ディスク・ドライブ・ストレージ・ブロックの再マッピング1600の例を示す。

【0051】

ディスク・マネージャ

図1に示すディスク・マネージャ106は、一般には、ディスクおよびディスク・アレイの管理を行うものであり、それには、グループ化/資源プーリング、ディスク属性の抽象化、フォーマット、ディスクの追加(addition)/削減(subtraction)、およびディスク・サービス時間およびエラー率の追跡を含む。ディスク・マネージャ106は、ディスクの様々なモデルの間の違いを区別せず、RAIDコンポーネントについて包括的なストレージ装置を提示する。また、ディスク・マネージャ106はグループ化機能を提供し、これは、例えば10000RPMのディスクなどのような特定の特徴をもつRAIDグループの構築を容易にする。

10

【0052】

本発明の一実施形態では、ディスク・マネージャ106は、少なくとも3重のものとなっている。その3重のものとは、抽象化、コンフィギュレーション、およびI/O最適化である。ディスク・マネージャ106は、「ディスク」を上位レイヤに提示し、上位レイヤは、例えば、ローカルまたはリモートに取り付けられた物理ディスク・ドライブや、リモートに取り付けられたディスク・システムなどとすることができる。

20

【0053】

共通の基礎にある特徴は、こうした装置の何れもがI/O動作のターゲットであり得ることである。抽象化サービスは、特にRAIDサブシステムのような上位レイヤに対して、一様なデータ・パス・インターフェースを提供し、そして、管理者に対して、ターゲット装置の管理のための包括的な機構を提供する。

【0054】

また、本発明のディスク・マネージャ106は、ディスクのグループ化機能を提供して管理およびコンフィギュレーションを簡素化する。ディスクには名前を付け、グループに入れることができ、グループにも名前を付けることができる。グループ化は強力な機能であり、これにより、ボリュームを或るディスクのグループから別のグループへと移行させるタスクや、或るグループのディスクを特定の機能専用にするタスクや、或るグループのディスクをスペアに指定するタスクなどのタスクを簡単にする。

30

【0055】

また、ディスク・マネージャは、外部装置の存在を検出する役目をもつSCSI装置サブシステムなどのような装置とインターフェースを行う。SCSI装置サブシステムは、少なくともファイバ・チャネル/SCSIタイプの装置については、ブロック型のターゲット装置である装置のサブセットを判定する能力をもつ。ディスク・マネージャが管理および抽象化を行うのはこうした装置である。

【0056】

更に、ディスク・マネージャは、SCSI装置レイヤからのフロー制御に応答する役目をもつ。ディスク・マネージャは、キューイング機能を有し、これにより、ディスク・ドライブ・システムのスループットを最適化する一方法として、I/O要求を集合させる機会が与えられる。

40

【0057】

更に、本発明のディスク・マネージャは、複数のディスク・ストレージ・システム・コントローラを管理する。また、複数の冗長ディスク・ストレージ・システム・コントローラを実装して、動作させられるディスク・ストレージ・システム・コントローラの障害をカバーすることができる。この冗長ディスク・ストレージ・システム・コントローラもディスク・マネージャによって管理される。

【0058】

ディスク・マネージャと他のサブシステムとの関係

50

ディスク・マネージャは、他の幾つかのサブシステムとやり取りを行う。RAIDサブシステムは、データ・パス活動に関してディスク・マネージャが提供するサービスの重要なクライアントである。RAIDサブシステムは、ディスク・マネージャを、I/Oのためのディスクへの独占的なパスとして使用する。また、RAIDシステムでも、ディスク・マネージャからのイベントについて聴取を行い、ディスクの存在および動作状態の判定を行う。また、RAIDサブシステムはディスク・マネージャと協働して、RAID装置の構築のための範囲(extent)を割り当てる。管理コントロールは、ディスクの存在の情報を得るため、また動作状態の変化について情報を得るために、ディスク・イベントについて聴取を行う。本発明の一実施形態では、RAIDサブシステム104は、RAID-0、RAID-1、RAID-5、RAID-10などのRAIDタイプのうちの少なくとも1つからなる組合せを含むことができる。代替のRAIDサブシステムでは、RAID-3、RAID-4、RAID-6、RAID-7などのような他のRAIDタイプを使用することも理解されよう。

10

**【0059】**

本発明の一実施形態では、ディスク・マネージャは、コンフィギュレーション・アクセスのサービスを使用して、持続的なコンフィギュレーションを保存し、また、統計などのような過渡的な読み出し専用情報をプレゼンテーション・レイヤへ提示する。ディスク・マネージャは、こうしたパラメータへのアクセスのためにコンフィギュレーション・アクセスへのハンドラの登録を行う。

20

**【0060】**

また、ディスク・マネージャは、ブロック・デバイスの存在および動作状態についての情報を得るためにSCSI装置レイヤのサービスを使用するものであり、また、こうしたブロック・デバイスへのI/Oパスをもつ。ディスク・マネージャは、ディスクを一意に識別するための補助的な方法として、装置についてSCSI装置サブシステムに問い合わせを行う。

**【0061】**

データ・インスタント・リプレイおよびデータ・インスタント・フュージョン

また、本発明は、データ・インスタント・リプレイ(data instant replay)およびデータ・インスタント・フュージョン(data instant fusion)の方法を提供する。図3のAおよびBに、本発明の諸原理に従った、複数の時間間隔でのRAIDサブシステムのディスク・ストレージ・ブロックのスナップショットの概略図を示す。図3Cには、データ・インスタント・リプレイの方法300を示し、この方法は、論理ブロックまたはディスク・ストレージ・ブロックのデフォルトのサイズを定義して、RAIDサブシステムのディスク空間がストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスを形成するようにするステップ302と、ページ・プールのボリュームのスナップショットまたはディスク・ストレージ・ブロックのマトリックスのスナップショットを所定の時間間隔で時間的に生成するステップ304と、スナップショットまたはデルタ(delta)のアドレス・インデックスを、ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスに保存して、ディスク・ストレージ・ブロックのマトリックスのスナップショットまたはデルタをその保存したアドレス・インデックスによって直ちに見つけ出せるようにするステップとを含む。

30

40

**【0062】**

図3Bに示すように、所定の各時間間隔、例えば、T1(午後12:00)、T2(午後12:05)、T3(午後12:10)、T4(午後12:15)などのような5分毎に、ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスのスナップショットが自動的に生成される。ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスの中のスナップショットまたはデルタのアドレス・インデックスは、ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスに保存されて、ストレージのページ・プールまたはディスク・ストレージ・ブロックのマトリックスのスナップショットまたはデルタを、その保存したアドレス・

50

インデックスによって直ちに見つけ出せるようにされる。

【0063】

従って、このデータ・インスタント・リプレイの方法では、RAIDサブシステムのスナップショットを、ユーザの定義する時間間隔や、ユーザのコンフィギュレーションする例えば数分や数時間ごとなどのような動的なタイム・スタンプや、サーバの指示する時刻に、時間的に生成する。システム障害やウイルス攻撃にあった場合、こうしたタイム・スタンプを有する仮想スナップショットにより、データ・インスタント・リプレイおよびデータ・インスタント・リカバリが、約数分または数時間などで可能になる。この技法は、インスタント・リプレイ・フュージョンとも呼ばれる。即ち、クラッシュまたは攻撃のすぐ前のデータのフュージョンが遅れずに行われ、クラッシュまたは攻撃の前に保存したスナップショットを直ちにそれ以降の動作で使用することができる。

10

【0064】

図4に、本発明の諸原理に従った、RAIDサブシステムのディスク・ストレージ・ブロックの複数のスナップショットを使用することによるデータ・インスタント・フュージョン機能400の概略図を更に示している。T3で、スナップショットの並列チェーンT3'~T5'が生成され、それにより、フュージョンされたデータT3'によってフュージョンおよび/またはリカバリされたデータを、T4でフュージョンされるデータと置き換えるために使用することができる。同様に、スナップショットの複数の並列チェーンT3''、T4''を生成して、T4'~T5'およびT4''~T5''でフュージョンされるデータと置き換えることができる。一代替実施形態では、T4、T4'~T5'、T5''でのスナップショットは、なおもページ・プールまたはマトリックスに保存することができる。

20

【0065】

スナップショットをローカルのRAIDサブシステムまたはリモートのRAIDサブシステムに保存ことができ、それにより、重大なシステム・クラッシュが、例えば、テロリスト攻撃などが原因で生じた場合にも、データ・インテグリティについて影響を受けず、データを直ちにリカバリできる。図5に、本発明の諸原理に従った、RAIDサブシステムのディスク・ストレージ・ブロックのスナップショットを使用することによる、ローカル・リモートのデータ複製およびインスタント・リプレイ機能500の概略図を示す。

30

【0066】

リモート複製(replication)は、リモート・システムに対してボリューム・データを複製するサービスを行う。リモート複製では、ローカルとリモートとのボリュームを、可能な限り同期させた状態に保つことを試みる。一実施形態では、リモートのボリュームのデータが、ローカルのボリュームのデータの完全なコピーをミラーリングしていないこともあり得る。ネットワークの接続性(コネクティビティ)およびパフォーマンスに起因して、リモートのボリュームがローカルのボリュームと同期しない状態が引き起こされる場合もあり得る。

【0067】

データ・インスタント・リプレイおよびデータ・インスタント・フュージョンの方法の別の特徴は、スナップショットを、システムが動作している状態で検査のために使用できることである。生のデータ(live data)をリアル・タイムの検査に使用することができる。

40

【0068】

スナップショットおよびポイント・イン・タイム・コピー(PITC)

データ・インスタント・リプレイの一例は、本発明の諸原理に従った、RAIDサブシステムのディスク・ストレージ・ブロックのスナップショットを使用することである。スナップショットは、ボリュームへの書き込み動作を記録して、過去のボリュームの内容を見るためにビューを作成できるようにする。また、スナップショットは、ボリュームのそれまでのポイント・イン・タイム・コピー(Point-in-Time-Copy(特定の時にとられたコ

50



ピー)、P I T C) に対するビューを作成することによって、データ・リカバリのサポートを行う。

【0069】

スナップショットのコアは、スナップショットの作成、合体 (coalesce)、管理、および I/O 動作を実現する。スナップショットは、ボリュームへの書き込みを監視し、ボリュームのビューを通してアクセスのためにポイント・イン・タイム・コピー (P I T C) を作成する。スナップショットは、論理ブロック・アドレス (logical block address、L B A) 再マッピング・レイヤを、仮想化レイヤ内のデータ・パスへ追加する。これは、I/Oパス内の仮想 L B A マッピングの別レイヤである。P I T C は、すべてのボリューム情報をコピーしない場合もあり得、単に再マッピングで使用されるテーブルを変更するだけであり得る。

10

【0070】

スナップショットは、ボリューム・データについての変更を追跡し、それまでのポイント・イン・タイムからのボリューム・データを見る機能を提供する。スナップショットは、この機能を、それぞれの P I C T に対するデルタ書き込みのリストを維持することによって、行う。

【0071】

スナップショットは、P I T C プロファイルのための複数の方法を提供するが、これは、アプリケーションにより開始されるもの、および時間により開始されるものを含む。スナップショットは、アプリケーションが P I T C を作成するようにする機能を提供する。アプリケーションは、作成を、サーバ上の A P I を通して制御し、これはスナップショット A P I へと送られる。また、スナップショットは、時間プロファイルを作成する機能を提供する。

20

【0072】

スナップショットは、ジャーナル用システムの実装や、ボリュームへの書き込みのすべてのリカバリを行わない場合もあり得る。スナップショットは、P I T C ウィンドウ内の一つのアドレスへの最後の書き込みを保持するのみであり得る。スナップショットにより、ユーザは、例えば分や時間などの定められた短期間をカバーする P I T C を作成することが可能とされる。スナップショットは、障害を取り扱うために、すべての情報をディスクに書き込む。スナップショットは、デルタ書き込み (delta writes) を含むボリューム・データ・ページ・ポインタを維持する。テーブルはマップをボリューム・データへ提供するものであり、これがなければボリューム・データへはアクセス不能であるため、テーブル情報はコントローラ障害の場合を取り扱わなければならない。

30

【0073】

ビュー・ボリューム機能 (ボリュームを見る機能) により、P I T C へのアクセスが提供される。ビュー・ボリューム機能は、アクティブな P I T C を除き、そのボリューム内の任意の P I T C に付加させることができる。P I T C への付加は、比較的迅速な動作である。ビュー・ボリューム機能の用途は、検査、訓練、バックアップ、およびリカバリを含む。ビュー・ボリューム機能は、書き込み動作を可能にし、その基礎とされている基礎的 P I T C の変更は行わない。

40

【0074】

一実施形態では、スナップショットは、ディスク空間を使用して、パフォーマンスを最適化し、使用を容易にするように設計されており、それは下記のようなものである。

【0075】

・スナップショットにより、ユーザ要求に対する応答時間が速くなる。ユーザ要求には、I/O、P I T C の作成、ビュー・ボリュームの作成/削除が含まれる。これを達成するために、スナップショットは、必要最小量よりも多くのディスク空間を使用してテーブル情報を保存する。I/Oについては、スナップショットは、ボリュームの現在の状態を要約して一つのテーブルにして、読み出しおよび書き込みの要求をすべて一つのテーブルで満たすことができるようにする。スナップショットは、通常の I/O 動作に対する影響

50

をできる限り低減する。第2に、ビュー・ボリューム動作については、スナップショットは、メイン・ボリュームのデータ・パスと同じテーブル機構を使用する。

【0076】

・スナップショットは、コピーされるデータの量を最小化する。このためには、スナップショットは、それぞれのPITCについてのポインタのテーブルを維持する。スナップショットはポインタのコピーおよび移動を行うが、ボリューム上のデータの移動は行わない。

【0077】

・スナップショットは、固定サイズのデータ・ページを用いてボリュームの管理を行う。個々のセクタを追跡するには、一つの手頃なサイズのボリュームでも膨大な量のメモリが必要になる。セクタよりも大きなデータ・ページを使用することにより、ある種のページは、別のページから直接に複製された情報の一部を含むことができる。

【0078】

・スナップショットは、ボリューム上のデータ空間を使用して、データ・ページ・テーブルを保存する。ルックアップ・テーブルは、コントローラ障害の後に複製される。ルックアップ・テーブルは、ページを割り当て、それをさらに分割する。

【0079】

・スナップショットがコントローラ障害を取り扱うが、その際、一実施形態では、スナップショットを用いるボリュームが一つのコントローラ上で動作することが必要とされる。この実施形態ではコピーレンシは必要ない。ボリュームへのすべての変更は、代替のコントローラによるリカバリのために、ディスクまたは信頼できるキャッシュへ記録される。コントローラ障害からのリカバリの際には、一実施形態では、スナップショット情報をディスクから読み取ることが必要である。

【0080】

・スナップショットは、仮想RAIDインターフェースを使用して、ストレージへのアクセスを行う。スナップショットは、複数のRAID装置を単一のデータ空間として使用することができる。

【0081】

・スナップショットは、ボリュームあたり「n」個のPITCおよびボリュームあたり「m」個のビューをサポートする。「n」および「m」に対する制限は、ディスク空間およびコントローラのメモリの関数となる。

【0082】

ボリュームおよびボリュームの割り当て/レイアウト

スナップショットは、LBA再マッピング・レイヤをボリュームに追加する。再マッピングは、I/O要求LBAおよびルックアップ・テーブルを使用して、アドレスをデータ・ページへと変換する。図6に示すように、スナップショットを用いた提示されるボリュームは、スナップショットのないボリュームと同じふるまいをする。それは、線形LBA空間を有し、I/O要求を取り扱う。スナップショットは、RAIDインターフェースを使用してI/Oを行うものであり、複数のRAID装置を1つのボリュームに含める。一実施形態では、スナップショット・ボリュームに対するRAID装置のサイズは、提示されるボリュームのサイズではない。RAID装置により、スナップショットは、ボリューム内でデータ・ページ用の空間を拡張することを可能とされる。

【0083】

新しいボリュームは、初めにスナップショットをイネーブルにしてあると、新しいデータ・ページ用の空間を含めればよいだけである。スナップショットは、ボトム・レベルのPITCに置くページのリストを作成しない。ボトム・レベルPITCは、この場合には空である。割り当て時に、すべてのPITCページはフリー・リスト上にある。初めにスナップショットをイネーブルにしたボリュームを作ることにより、ボリュームが提示するよりも少ない物理空間を割り当て得る。スナップショットは、ボリュームへの書き込みを追跡する。本発明の一実施形態では、ヌル(NULL)・ボリュームは、ページ・ブー

10

20

30

40

50

ルまたはマトリックスへコピーおよび/または保存されず、これにより、ストレージ空間の使用の効率が上がる。

【 0 0 8 4 】

一実施形態では、どちらの割り当て方式でも、P I T Cが、仮想N U L Lボリュームをリストの底部に置く。N U L Lボリュームに対する読み出しは、ゼロのブロックを返す。N U L Lボリュームは、サーバがそれまで書き込んでいないセクタを取り扱う。N U L Lボリュームへの書き込みは起こり得ない。ボリュームは、N U L Lボリュームを、書き込まれていないセクタに対する読み出しのために使用する。

【 0 0 8 5 】

空きページの数、ボリュームのサイズ、P I T Cの数、および予測されるデータ変化レートに依存する。システムは、所与のボリュームに対しての割り当てるページの数を決定する。データ・ページ数は、時間が経つにつれて拡張される可能性がある。拡張により、予測されるよりも急速なデータの変化、より多くのP I T C、またはより大きなボリュームを、サポートすることができる。新しいページはフリー・リストへ追加される。フリー・リストへのページの追加は自動的に行ってよい。

10

【 0 0 8 6 】

スナップショットは、ボリューム空間の管理のためにデータ・ページを使用する。それぞれのデータ・ページは数メガバイトのデータを含み得る。オペレーティング・システムを使用すると、ボリュームの同じ領域へ多くのセクタが書き込まれる傾向がある。メモリ要件からも、スナップショットがボリュームの管理のためにページを使用することが要求される。1テラバイトのボリュームの各セクタについて一つの32ビットポインタを維持すると、8ギガバイトのR A Mが必要となり得る。異なるボリュームは異なるページ・サイズを有し得る。

20

【 0 0 8 7 】

図7に、スナップショット構造の一実施形態を示す。スナップショットは、多くのオブジェクトをボリューム構造に付加する。付加されるオブジェクトとしては、P I T C、アクティブなP I T Cへのポインタ、データ・ページ・フリー・リスト、子のビュー・ボリューム、およびP I T C合体オブジェクトがある。

【 0 0 8 8 】

・アクティブP I T C ( A P ) ポインタはボリュームにより維持される。A Pは、ボリュームに対する読み出しおよび書き込みの要求のマッピングを取り扱う。A Pは、ボリューム内のすべてのデータの現在の場所の要約 ( 一覧、summary ) を含む。

30

【 0 0 8 9 】

・データ・ページ・フリー・リストは、ボリューム上の使用可能なページの追跡を行う。

【 0 0 9 0 】

・オプションの子ビュー・ボリューム ( child view volume ) は、ボリュームP I T Cへのアクセスを提供する。ビュー・ボリュームは、P I T Cへの書き込みを記録するためにそれ自体のA Pを含むが、基礎となるデータの変更は行わない。1つのボリュームが複数の子ビュー・ボリュームをサポートしてもよい。

40

【 0 0 9 1 】

・スナップショット合体オブジェクトは、以前のP I T Cを取り除く目的で2つのP I T Cを一時的にリンクする。P I T Cの合体は、データ・ページの所有権を移動し、データ・ページを解放することを含む。

【 0 0 9 2 】

・P I T Cは、そのP I T Cがアクティブであった間に書き込まれたページに対するテーブルおよびデータ・ページを含む。P I T Cは、P I T Cが書き込み要求の受理を停止した時点の凍結タイム・スタンプを含む。また、P I T Cは、何れの時間にP I T Cが合体することになるかを定めるT i m e - t o - L i v e ( タイム・ツー・ライブ、生きる時間 ) 値を含む。

50

## 【 0 0 9 3 】

また、スナップショットは、ボリューム全体に対するデータ・ページ・ポイントの要約を、P I T C が取られた時点で行って、読み出しおよび書き込みの予測可能なパフォーマンスを提供する。他の解決策は、最新のポイントを見つけるために複数の P I T C を検査するために、読み出しを必要とし得る。こうした解決策では、テーブル・キャッシング・アルゴリズムが必要であるが、パフォーマンスは最低となる場合がある。

## 【 0 0 9 4 】

また、本発明におけるスナップショットを要約することにより、テーブルの最悪のメモリ使用が低減される。その場合、テーブル全体をメモリへロードすることを必要とされ得るが、単一のテーブルだけをロードすることを必要とされ得る。

10

## 【 0 0 9 5 】

要約（一覧）は、現在の P I T C の所有するページを含むが、それまでのすべての P I T C からのページを含むことができる。P I T C がどのページに書き込みを行うかを決定するために、それぞれのデータ・ページについてのページ所有権を追跡する。また、一覧は、合体プロセスに対する所有権の追跡も行う。これを取り扱うために、データ・ページ・ポイントはページ・インデックスを含む。

## 【 0 0 9 6 】

図 8 に、P I T C のライフ・サイクルの一実施形態を示す。それぞれの P I T C は、読み出し専用として記憶（コミット）される前に、多数の下記の状態を経過する。

## 【 0 0 9 7 】

1 . テーブル作成 - 作成を行うと、テーブルが作成される。

20

## 【 0 0 9 8 】

2 . ディスクへのコミット（記憶、commit） - これは、P I T C に対してディスク上にストレージを生成する。この時点でテーブルを書き込むことにより、テーブル情報の保存に必要な空間が、P I T C が取られる前に割り当てられることが保証される。同時に、P I T C オブジェクトもディスクへ記憶（コミット）される。

## 【 0 0 9 9 】

3 . I / O の受け入れ - これは P I T C はアクティブ P I T C ( A P ) になっている。 - ここでは、P I T C は、ボリュームに対する読み出しおよび書き込みの要求を取り扱う。これは、テーブルへの書き込み要求を受理する唯一の状態である。P I T C は、いまアクティブであるというイベントを生成する。

30

## 【 0 1 0 0 】

4 . ディスクへのテーブルの記憶 - P I T C はもはや A P ではなく、もはやそれ以上のページを受理しない。新しい A P による引き継ぎが行われている。この時点より後では、テーブルは、合体動作中に取り除かれない限り、変更されない。これは読み出し専用である。P I T C は、この時点で、凍結され記憶されているというイベントを生成する。何れのサービスもそのイベントを聴取することができる。

## 【 0 1 0 1 】

5 . テーブルのメモリの解放（リリース） - テーブルが必要とするメモリを解放する。また、このステップは、変更がすべてディスクへと書き込まれていることを示すために、ログをクリアする。

40

## 【 0 1 0 2 】

ボリュームまたはビュー・ボリュームのためのトップ・レベル P I T C は、アクティブ P I T C ( A P ) と呼ばれる。A P は、ボリュームへのすべての読み出しおよび書き込みの要求を満たす。A P は、そのボリュームに対しての唯一の、書き込み要求を受理できる P I T C である。A P は、ボリューム全体に対するデータ・ページ・ポイントの一覧を含む。

## 【 0 1 0 3 】

A P は、合体プロセスに関して、宛先（destination）となることはできるが、ソース（source）となることはできない。宛先である場合、A P は、所有するページの数を増加

50

するが、データのビューは変更しない。

【0104】

ボリュームの拡張に関しては、APはボリュームとともに直ちに大きくなる。新しいページは、NULLボリュームを指す。APではないPITCは、ボリュームの拡張のために変更は必要ではない。

【0105】

それぞれのPITCは、入ってくるLBAを、基礎となるボリュームへのデータ・ページ・ポインタへとマッピングするテーブルを維持する。テーブルは、データ・ページへのポインタを含む。テーブルは、提示される論理空間よりも多くの物理ディスク空間をアドレスする必要がある。図9に、マルチ・レベル・インデックスを有するテーブル構造の一実施形態を示す。この構造は、ボリュームLBAを復号化してデータ・ページ・ポインタにする。各レベルでは、図9に示すように、アドレスの増加していくLSB (less significant bit) を復号化する。このテーブルの構造により、高速のルックアップおよびボリュームの拡張の機能が提供される。高速ルックアップのために、マルチ・レベル・インデックス構造は、テーブルを浅い状態に維持し、各レベルで複数のエントリをもつようにしている。このインデックスは、各レベルでアレイのルックアップを行う。ボリュームの拡張をサポートするため、マルチ・レベル・インデックス構造により、拡張をサポートするための別のレイヤを追加することが可能になっている。この場合のボリュームの拡張は、上位レイヤへ提示されるLBAカウント(計数)の拡張であり、そのボリュームへ割り当てられるストレージ空間の実際の量ではない。

10

20

【0106】

マルチ・レベル・インデックスは、ボリューム全体のデータ・ページ再マッピングの一覧(要約)を含む。それぞれのPITCは、記憶されたポイント・イン・タイム(時(特定の時))の、このボリュームに対する完全な再マッピング・リストを含む。

【0107】

マルチ・レベル・インデックス構造は、テーブルのレベルに対して異なるエントリ・タイプを使用する。異なるエントリ・タイプは、情報をディスクから読み出す必要性や、情報をメモリに格納する必要性をサポートする。ボトム・レベルのエントリは、データ・ページ・ポインタのみを含み得る。トップおよび中間のレベルのエントリは2つのアレイを含み、1つは次のレベルのテーブル・エントリのLBAであり、そして、テーブルへのメモリ・ポインタを含む。

30

【0108】

提示されるボリューム・サイズが拡張する際に、以前のPITCテーブルのサイズは増やす必要はなく、また、テーブルを変更する必要もない。テーブルの中の情報は、読み出し専用であるので、変わらず、そして、拡張プロセスは、テーブルを、NULLページ・ポインタを末尾に追加することによって、変更する。スナップショットは、以前のPITCからのテーブルを直接にユーザに提示しない。

【0109】

I/O動作があると、テーブルには、LBAをデータ・ページ・ポインタへマッピングすることが求められる。その場合、I/Oでは、データ・ページ・ポインタにデータ・ページ・サイズを乗算して、基礎となるRAIDのLBAを得る。一実施形態では、データ・ページ・サイズは、2のべき乗である。

40

【0110】

テーブルは、LBAの再マッピング、ページの追加、およびテーブルの合体のためのAPIを提供する。

【0111】

スナップショットは、PITCオブジェクトおよびLBAマッピング・テーブルを保存するためにデータ・ページを使用する。テーブルは、そのテーブル・エントリに対するI/OのためにRAIDインターフェースへ直接にアクセスする。テーブルは、RAID装置に対してそのテーブルの読み出しおよび書き込みを行うときの変更を最小にする。変更

50

がない場合には、テーブルの情報の読み出しおよび書き込みを、テーブル・エントリ構造に対して直接に行うことが可能になる。これにより、I/Oのために必要となるコピーが低減される。スナップショットは、ディスク上のホット・スポットの発生を防ぐために、変化ログを使用することができる。ホット・スポットとは、ボリュームへの更新を追跡するために繰り返し使用される場所のことである。変化ログは、PITCテーブルへの更新、およびボリュームに関するフリー・リストを記録する。リカバリ中には、スナップショットは変化ログを使用して、メモリ内AP (in-memory AP) およびフリー・リスト (free list) の再作成を行う。図10にテーブルのリカバリの一実施形態を示し、これにより、メモリ内APと、ディスク上AP (on-disk AP) と、変化ログとの間の関係を示している。この図は、フリー・リストに関しても同じ関係を示している。メモリ内APテーブルは、ディスク上APテーブルおよびログから再構築することができる。どのようなコントローラ障害の場合でも、APは、ディスク上APを読み出してそれにログ変化を適用することによって、再構築される。変化ログは、システム・コンフィギュレーションに応じて、異なる物理資源を使用する。多重コントローラのシステムでは、変化ログは、ストレージのためにバッテリ・バックアップされたキャッシュ・メモリを頼りにしている。キャッシュ・メモリを使用することにより、スナップショットは、ディスクへのテーブル書き込みの数を減らし且つデータ・インテグリティを維持することができる。変化ログは、リカバリのためにバックアップ・コントローラに対しての複製を行う。単一コントローラのシステムでは、変化ログは、すべての情報をディスクへ書き込む。これには、ディスク上のログの場所にホット・スポットを発生させるという副作用がある。これにより、幾つもの変化を単一のデバイス・ブロックへ書き込むことが可能とされる。

10

20

30

40

50

#### 【0112】

スナップショットは、定期的に、PITCテーブルおよびフリー・リストをディスクへ書き込んで、ログにおけるチェックポイントの作成およびそのクリアを行う。この期間は、PITCへの更新の数に応じて変わり得る。合体プロセスは変化ログを使用しない。

#### 【0113】

スナップショットのデータ・ページI/Oは、要求がデータ・ページ境界内に収まることを必要とし得る。スナップショットは、ページ境界にまたがるI/O要求にであった場合に、その要求を分割する。次いで、スナップショットは、その要求を要求ハンドラへと渡す。書き込みおよび読み出しのセクションは、I/Oがページ境界内に収まることを仮定している。APは、LBA再マッピングを提供してI/O要求を満たす。

#### 【0114】

APはすべての書き込み要求を満たす。スナップショットは、所有ページと非所有ページに対する2つの異なる書き込みシーケンスをサポートする。異なるシーケンスにより、テーブルへのページの追加が可能になる。図11に、所有ページ・シーケンスおよび非所有ページ・シーケンスがある書き込みプロセスの一実施形態を示す。

#### 【0115】

所有ページ・シーケンスでは、プロセスは以下のものを含む。

- 1) テーブル・マッピングを見つける。
- 2) 所有の書き込みをページングする (Page Owned Write) - LBAを再マッピングし、データをRAIDインターフェースへ書き込む。

#### 【0116】

以前に書き込まれたページは、純粋な書き込み要求である。スナップショットは、データをそのページに書き込んで、現在の内容を上書きする。APの所有するデータ・ページだけが書き込まれる。他のPITCの所有するページは、読み出し専用である。

#### 【0117】

非所有ページ・シーケンスでは、プロセスは以下のものを含む。

- 1) テーブル・マッピングを見つける。
- 2) 以前のページを読み出す - そのデータ・ページに対する読み出しを行い、書き込み要求と読み取ったデータとでページ全部が出来上がるようにする。これが、コピー・

オン・ライト (copy on write) ・プロセスの開始となる。

3) データを組み合わせる - データ・ページの読み出しおよび書き込みの要求のペイロードを、単一の連続するブロックに入れる。

4) リスト割り当てを自由にする - 新しいデータ・ページ・ポインタを、フリー・リストから得る。

5) 組み合わせたデータを、新しいデータ・ページへ書き込む。

6) 新しいページ情報をログへ記憶する。

7) テーブルを更新する - テーブル内の L B A 再マッピングを変更して、新しいデータ・ページ・ポインタを反映させる。これでデータ・ページは P I T C により所有される。

10

#### 【0118】

ページを追加するには、そのページがテーブルに追加されるまで、読み出しおよび書き込みの要求を妨げることが必要であり得る。スナップショットは、テーブルの更新をディスクへ書き込み、ログの複数のキャッシュ記憶されたコピーを維持することにより、コントローラのコヒーレンスを達成する。

#### 【0119】

読み出し要求に関して、A P はすべての読み出し要求を満たす。読み出し要求は、A P テーブルを用いて、L B A を、そのデータ・ページの L B A へ再マッピングする。再マッピングされた L B A は R A I D インターフェースへと渡されて、その要求が満たされる。ボリュームは、それまでにそのボリュームへ書き込まれていないデータ・ページに対する読み出し要求を満たすことができる。こうしたページは、P I T C テーブルにおいて N U L L アドレス (すべて 1) でマークされている。このアドレスに対する要求は、N U L L ボリュームにより満たされ、一定のデータ・パターンを返す。異なる P I T C テーブルの所有するページは、ページ境界にまたがる読み出し要求を満たすことができる。

20

#### 【0120】

スナップショットは、N U L L ボリュームを使用して、以前に書き込まれていないデータ・ページに対する読み出し要求を満たす。これは、それぞれのセクタ読み出しに対して、すべてゼロを返す。これは、R A I D 装置も、割り当てられた空間もない。N U L L ボリュームに対する読み出しについてのデータ要件を満たすために、すべてゼロであるブロックをメモリ内に置いておくことが予期される。すべてのボリュームは、読み出し要求を満たすために N U L L ボリュームを共用する。

30

#### 【0121】

一実施形態では、合体プロセスが、P I T C と、その所有ページの一部とをボリュームから取り除く。P I T C を取り除くと、新しい差異を追跡するために使用可能な空間がより多く作成される。合体では、2つの隣接するテーブルを差異があるかどうか比較し、新しい差異だけを維持する。合体は、ユーザによるコンフィギュレーションに従って、定期的にまたは手動で起きる。

#### 【0122】

このプロセスは、2つの P I T C、即ち、ソースと宛先とを含む。一実施形態での、適切なオブジェクトについての規則は、次のようなものである。

40

1) ソースは、宛先に対して、以前の P I T C でなければならない - ソースは宛先より前に作成されなければならない。

2) 宛先は、同時にソースであってはならない。

3) ソースは、複数の P I T C により参照されてはならない。複数の参照は、ビュー・ボリュームが P I T C から作成されたときに、起きる。

4) 宛先は、複数の参照をサポートし得る。

5) A P は宛先であり得るが、ソースではない。

#### 【0123】

合体・プロセスは、すべての変更をディスクへ書き込むものであり、コヒーレンスを必要としない。コントローラが故障した場合に、ボリュームが P I T C 情報をディスクから

50

リカバリし、合体プロセスを再開する。

【0124】

このプロセスは、2つのPITCを合体のためにマークするものであり、以下の諸ステップを含む。

【0125】

1) ソースの状態が合体のソースに設定される - この状態は、コントローラ障害のリカバリのために、ディスクへ記憶される。この時点で、ソースはもはやアクセスされないが、それは、そのデータ・ページが無効であるからである。データ・ページがフリー・リストへ返されるか、または所有権が宛先へ移される。

【0126】

2) 宛先の状態が合体の宛先へ設定される - この状態は、コントローラ障害のリカバリのためにディスクへ記憶される。

【0127】

3) テーブルのロードおよび比較を行う - このプロセスは、データ・ページ・ポインタを移動させる。解放されたデータ・ページは、直ちにフリー・リストに追加される。

【0128】

4) 宛先の状態が正常に設定される - プロセスが完了する。

【0129】

5) リストを調整する - ソースのnext(次)ポインタのprevious(以前)を、宛先へと変更する。これにより、ソースがリストから効率よく取り除かれる。

【0130】

6) ソースを解放する - 制御の情報に使用された何れのデータ・ページをも、フリー・リストへ戻す。

【0131】

上述のプロセスは、2つのPITCの組合せをサポートしている。合体を、単一のパスにおいて複数のPITCの削除と複数のソースの作成とを行うように設計できることが、当業者には理解されている。

【0132】

図2に示すように、ページ・プールは、そのページ・プールに関連するすべてのボリュームが使用するように、データ・ページ・フリー・リストを維持する。フリー・リスト・マネージャは、ページ・プールからのデータ・ページを使用して、フリー・リストを永久的ストレージへ記憶する。フリー・リストの更新を引き起こすものは幾つもあり、書き込みプロセスはページを割り当て、制御ページ・マネージャはページを割り当て、合体用プロセスはページを返す。

【0133】

フリー・リストは、それ自体をある閾値のところで自動的に拡張するためのトリガを維持することができる。このトリガは、ページをページ・プールへ追加するためのページ・プール拡張方法を使用する。自動的な拡張は、ボリューム・ポリシの一機能とすることもできる。重要度の高いデータ・ボリュームには拡張を可能にし、重要度の低いボリュームは強制的に合体させる。

【0134】

ビュー・ボリュームは、以前のポイント・イン・タイムへのアクセスを提供し、正常のボリュームI/O動作をサポートする。PITCは、PITC間の差異を追跡し、ビュー・ボリュームは、PITC内に含まれる情報へユーザがアクセスできるようにする。ビュー・ボリュームはPITCから分岐する。ビュー・ボリュームは、リカバリ、検査、バックアップ動作などをサポートする。ビュー・ボリュームの作成は、データ・コピーを必要としないので、ほぼ瞬間的に生じる。ビュー・ボリュームは、ビュー・ボリュームへの書き込みをサポートするために、それ自体のAPが必要になり得る。

【0135】

ボリュームの現在の状態から取られたビューでは、APを現在のボリュームのAPから

10

20

30

40

50



コピーすることができる。APを用い、ビュー・ボリュームは、ビュー・ボリュームへの書き込み動作を、基礎となるデータを変更せずに行うことを可能にする。OSは、データを使用するために、ファイル・システムまたはファイルの再構築を必要とし得る。ビュー・ボリュームは、APおよび書き込まれたデータ・ページに対して親ボリュームから空間を割り当てる。ビュー・ボリュームは、関連するRAID装置情報を有さない。ビュー・ボリュームを削除すると、空間が解放され親ボリュームへ戻される。

【0136】

図12に、スナップショットを用いてのあるボリュームに関する推移を示す例示的なスナップショット動作を示す。図12は、10ページをもつボリュームを示す。それぞれの状態(ステート)は、ボリュームに対する読み出し要求遂行(Read Request Fulfillment)リストを含む。影を付けたブロックは、所有データ・ページ・ポインタを示す。

10

【0137】

図の左側(即ち、初期状態)から図の中央への推移は、ページ3および8への書き込みを示している。ページ3への書き込みには、PITC I (AP)への変更が必要である。PITC Iは、新ページ書き込み処理に従い、ページ3をテーブルへ追加する。PITC Iは、変更されていない情報をページJから読み出し、ドライブ・ページBを使用してそのページを保存する。このPITC Iのページ3へのこれ以降のすべての書き込みは、ページを移動させずに取り扱われる。ページ8への書き込みでは、ページへの書き込みに関する第2の場合を示す。PITC Iはすでにページ8を含んでいるため、PITC Iは、ページ8のそのデータのその部分へ書き込む。この場合、それはドライブ・ページC

20

【0138】

図の中央から図の右側(即ち、最終状態)への推移は、PITC IIとIIIとの合体を示す。スナップショットの合体は、それぞれ、古いページを削除すること、および両方のPITCにおけるすべての変更を維持することを含む。両方のPITCがページ3および8を含む。プロセスは、PITC IIからの新しい方のページを保ち、PITC IIIからのページを解放し、ページAおよびDをフリー・リストへ返す。

【0139】

スナップショットは、フリー・リストおよびPITCテーブル情報を保存するために、ページ・プールからデータ・ページの割り当てを行う。制御ページの割り当ては、データ・ページを、オブジェクトに必要とされるサイズに合うように、更に細かく割り当てる(sub-allocate)。

30

【0140】

ボリュームは、制御ページ情報の先頭に対するページ・ポインタを含む。このページから、他のすべての情報を読み出すことができる。

【0141】

スナップショットは、特定の時間間隔で、使用中のページの数を追跡する。これにより、スナップショットは、いつユーザがシステムに更なる物理ディスク空間を追加する必要があるかを予測して、スナップショットが不足するのを防ぐ。

【0142】

データ・プログレッション

本発明の一実施形態では、データ・プログレッション(DP、data progression)を使用して、データを、徐々に、適切なコストのストレージ空間へと移動させる。本発明により、ユーザは、ドライブが実際に必要であるときにドライブの追加を行える。これにより、ディスク・ドライブのコスト全体が著しく低減される。

40

【0143】

データ・プログレッションでは、最近アクセスされていないデータおよび履歴のスナップショット・データを、より安価なストレージへ移動させる。最近アクセスされていないデータについては、最近アクセスされていないページがに対してのストレージのコストを徐々に低減する。データは、コストが最も低いストレージへ直ちに移動させない。履歴の

50

スナップショット・データについては、読み出し専用ページを、例えば R A I D 5 などのようなより効率のよいストレージ空間へ移動させ、また、そのページにもはやボリュームからアクセス可能でない場合には、最も安価なストレージへ移動される。

【 0 1 4 4 】

本発明のデータ・プログレッションの他の利点は、現在アクセス中のデータへの高速の I / O アクセスを維持すること、および高速であるが高価なディスク・ドライブを購入する必要性を減らすことを含む。

【 0 1 4 5 】

データ・プログレッションでは、動作時に、ストレージのコストを、物理メディアのコストおよびデータ保護に使用される R A I D 装置の効率を用いて、判定する。また、データ・プログレッションでは、ストレージ効率を判定し、それに従ってデータの移動を行う。例えば、データ・プログレッションは、物理ディスク空間をより効率よく使用するために、R A I D 1 0 装置を R A I D 5 装置へと変換することができる。

【 0 1 4 6 】

データ・プログレッションでは、アクセス可能なデータを、サーバが現在読み出したまたは書き込みできるデータと定義している。データ・プログレッションでは、アクセス可能性 (accessibility) を使用して、ページが使用すべきストレージのクラスを決定する。ページは、履歴 P I T C に属する場合には、読み出し専用である。サーバが、一番最近の P I T C においてそのページを更新していない場合には、そのページはまだアクセス可能である。

【 0 1 4 7 】

図 1 7 に、データ・プログレッション動作におけるアクセス可能なデータ・ページの一実施形態を示す。アクセス可能なデータ・ページは、以下のカテゴリに分かれる。

【 0 1 4 8 】

・アクセス可能な、最近アクセスされたもの - これらは、ボリュームが最もよく使用しているアクティブ・ページである。

【 0 1 4 9 】

・アクセス可能な、最近アクセスされていないもの - 最近使用されていない読み出し / 書き込みページ。

【 0 1 5 0 】

・履歴のアクセス可能なもの - ボリュームによる読み出しが行える読み出し専用ページ - - スナップショット・ボリュームに当てはまる。

【 0 1 5 1 】

・履歴のアクセス不能なもの - ボリュームによって現在アクセスされていない読み出し専用データ・ページ - - スナップショット・ボリュームに当てはまる。スナップショットは、こうしたページをリカバリ目的で維持し、ページは、一般に、可能な限り低コストのストレージに置かれる。

【 0 1 5 2 】

図 1 7 に、スナップショット・ボリューム用の様々な所有ページを有する 3 つの P I T C を示している。動的容量のボリュームは、P I T C C だけで代表させてある。すべてのページは、アクセス可能であり読み出し / 書き込みである。ページは異なるアクセス時間を有し得る。

【 0 1 5 3 】

下記のテーブルは、様々なストレージ装置を、効率の高くなる順または金銭的な支出の減る順に示す。ストレージ装置のリストは、より遅い書き込み I / O アクセスの一般的な順序にも従い得る。データ・プログレッションでは、論理保護空間の効率を R A I D 装置の総物理空間で割ったものを計算する。

【 0 1 5 4 】

10

20

30

40

【表 1】

表 1 : RAIDタイプ

タイプ	サブタイプ	ストレージ 効率	1ブロック書込 I/Oカウント	使用
RAID 10		50%	2	比較的良好な書き込み性能をもつ一次的な読出/書込アクセス可能ストレージ。
RAID5	3ドライブ	66.6%	4 (2読み出し 2書き込み)	RAID10を超えての効率の増加は最小であり、RAID5の書き込みペナルティが生じる。
RAID5	5ドライブ	80%	4 (2読み出し 2書き込み)	読み出し専用の履歴情報用の優れた候補。最近アクセスされていない書き込み可能ページ用の良い候補。
RAID5	9ドライブ	88.8%	4 (2読み出し 2書き込み)	読み出し専用の履歴情報用の優れた候補。
RAID5	17ドライブ	94.1%	4 (2読み出し 2書き込み)	効率の増加は低減されるが RAID装置のフォールト・ドメインは2倍になる。

10

【0155】

RAID5の効率は、ストライプのドライブ数が増えると高くなる。ストライプのディスク数が増えると、フォールト・ドメインが増大する。ストライプのドライブ数の増加により、RAID装置の作成に必要なディスクの最小数も増える。一実施形態では、データ・プログレッションは、9ドライブを超えるRAID5ストライプ・サイズを使用しない。なぜなら、フォールト・ドメイン・サイズが増え、効率の増加が限られるからである。データ・プログレッションは、スナップショット・ページ・サイズの整数倍であるRAID5ストライプ・サイズを使用する。これにより、データ・プログレッションは、ページをRAID5へと移動させるときにフル・ストライプ書き込みを行えるようになり、移動の効率を高める。すべてのRAID5コンフィギュレーションは、データ・プログレッションの目的で、同じ書き込みI/O特性を有する。例えば、2.5インチ(6.5cm)FCディスク上のRAID5は、そうしたディスクの性能を効率よく使用しない可能性がある。この組合せを防ぐために、データ・プログレッションは、あるRAIDタイプが特定のディスクタイプ上で動作しないようにする機能をサポートする必要がある。また、データ・プログレッションのコンフィギュレーションは、システムがRAID10またはRAID5の空間を使用しないようにすることもできる。

20

30

【0156】

ディスクのタイプを次のテーブルに示す。

【0157】

【表 2】

表 2 : ディスク・タイプ

タイプ	速度	費用	論点
2.5インチFC	非常に良い	高い	非常に高価
FC 15 K RPM	良い	中程度	高価
FC 10 K RPM	良い	良い	妥当な価格
SATA	ほどほど	低い	安価/信頼性が低い

40

【0158】

データ・プログレッションは、システム内のドライブに相対的なディスク・ドライブの分類を自動的に行う機能を含む。システムは、ディスクを検査して、そのパフォーマンスをシステム内の他のディスクと比較して決定する。より速いディスクは高評価区分へと分類され、より遅いディスクは低評価区分へと分類される。ディスクがシステムへ追加され

50

ると、システムは、ディスクの評価分類のバランスを自動的にとり直す。このアプローチでは、決して変化することのないシステムと、新しいディスクが追加されると頻繁に変化するシステムとの双方を取り扱う。この自動分類は、複数のドライブ・タイプを同じ評価区分内に入れることがある。ドライブの評価がかなり近接しているものと判定された場合には、その評価は同じになる。

## 【0159】

一実施形態では、システムは次のドライブを含む。即ち、

- 高 - 10K FCドライブ
- 低 - SATAドライブ

## 【0160】

15K FCドライブを追加すると、データ・プログレッションは、ディスクを自動的に再分類し、10K FCドライブを降格させる。この結果として次の分類となる。

- 高 - 15K FCドライブ
- 中 - 10K FCドライブ
- 低 - SATAドライブ

## 【0161】

別の実施形態では、システムには次のドライブ・タイプを有し得る。

- 高 - 25K FCドライブ
- 低 - 15K FCドライブ

## 【0162】

従って、15K FCドライブは評価の低い区分として分類され、これに対して15K FCドライブは評価の高い区分として分類される。

## 【0163】

SATAドライブがシステムに追加された場合、データ・プログレッションは、ディスクを時間的に再分類する。その結果として次の分類となる。

- 高 - 25K FCドライブ
- 中 - 15K FCドライブ
- 低 - SATAドライブ

## 【0164】

データ・プログレッションは、ウォーターフォール(waterfall)・プログレッションを含むことができる。典型的に、ウォーターフォール・プログレッションはデータをより安価な資源へ移動させるが、その移動は、その資源が完全に使用されたときにのみ行う。ウォーターフォール・プログレッションは、最も高価なシステム資源の使用を効率的に最大化する。システムのコストの最小化も行う。安いディスクを最も低いプールに追加すると、ボトム(底部)において、より大きなプールが作成される。

## 【0165】

典型的なウォーターフォール・プログレッションは、RAID10空間を使用し、次いで、RAID5空間などのような、RAID空間の次のものを使用する。このことにより、ウォーターフォールが次のクラスのディスクのうちRAID10へ直接向かうように強制される。代替例として、データ・プログレッションは、図24に示すような混合RAIDウォーターフォール・プログレッションを含んでもよい。この代替のデータ・プログレッション方法により、ディスク空間およびパフォーマンスを最大化する問題が解決され、ストレージは同じディスク・クラスにおいて更に効率のよい形式へと変形できる。また、この代替方法では、RAID10およびRAID5が或るディスク・クラスの資源全体を共用するという要件をサポートする。これには、あるRAIDレベルがあるクラスのディスク用に使用する、ディスク空間の一定の割合をコンフィギュレーションすることが必要であることがある。従って、この代替のデータ・プログレッション方法では、高価なストレージの使用が最大化され、かつ、別のRAIDクラスが共存する余地を残す。

## 【0166】

また、混合RAIDウォーターフォールでは、ストレージが限られているときのみ、ペ

10

20

30

40

50

ージをより安価なストレージへと移動させる。ディスク空間全体のうちの或るパーセンテージなどのような閾値により、特定の R A I D タイプのストレージの量が制限される。これにより、システムにおける最も高価なストレージの使用が最大化される。ストレージがその限度に近づくと、データ・プログレッションは、ページを、コストのより低いストレージへと自動的に移動させる。データ・プログレッションは、書き込みスパイクに対するバッファを提供することができる。

**【 0 1 6 7 】**

上述のウォーターフォール方法では、一部の場合のように、ページをコストが最も低いストレージへと直ちに移動させ得るということが理解され、履歴およびアクセス不能のページをより安価なストレージ上へ適時に移動させる必要性があり得る。また、履歴ページをより安価なストレージへ直ちに移動させることもできる。

10

**【 0 1 6 8 】**

図 1 8 に、データ・プログレッションのプロセス 1 8 0 0 の流れ図を示す。データ・プログレッションは、システムのそれぞれのページを、そのアクセスのパターンおよびストレージのコストに関して絶えず検査して、移動させるべきデータ・ページがあるかどうかを判定する。また、データ・プログレッションは、ストレージがその最大割り当て量に達しているかどうかを判定することもできる。

**【 0 1 6 9 】**

データ・プログレッションのプロセスは、そのページがどのボリュームからもアクセス可能かどうかを判定する。プロセスは、履歴に付加されている各ボリュームについての P I T C を検査して、そのページが参照されているかどうかを判定する。そのページがアクティブに使用中である場合には、そのページは昇格させるに又は降格を遅くするに望ましい。そのページがどのボリュームからもアクセス可能でない場合には、使用可能な最も低コストのストレージへ移動させられる。また、データ・プログレッションは、P I T C が期限切れになる前に時間を計算に入れる。スナップショットで、P I T C がまもなく期限切れになるようにスケジュールしている場合には、ページのプログレッションは行わない。ページ・プールがアグレッシブ・モードで動作中の場合には、ページのプログレッションが行われ得る。

20

**【 0 1 7 0 】**

データ・プログレッションの最近アクセス検出 ( recent access detection、最近されたアクセスの検出 ) では、ページの昇格動作から一塊の活動を取り除く必要がある。データ・プログレッションでは、読み出しと書き込みのアクセス追跡を分離している。これにより、データ・プログレッションは、アクセス可能である R A I D 5 装置上のデータを保持することが可能とされる。ウイルスのスキャンや報告などのような動作は、データを読み出すのみである。データ・プログレッションは、ストレージが残り少なくなっているときに、最近アクセス ( recent access ) の資格を変更する。これにより、データ・プログレッションは、ページの降格をよりアグレッシブに行えるようになる。また、これは、ストレージが残り少なくなっているときに、システムをボトム ( 底 ) の方から詰めていくのに役立つ。

30

**【 0 1 7 1 】**

データ・プログレッションは、システム資源が少なくなるとデータ・ページをアグレッシブ ( 積極的 ) に移動させることができる。より多くのディスクまたはコンフィギュレーションの変更は、こうした場合のすべてでやはり必要である。データ・プログレッションは、時間のない状況においてシステムが動作し得る時間の量を延ばす。データ・プログレッションは、システムをできるだけ長く動作可能に保つように試みる。その時間とは、そのすべてのストレージ・クラスの空間がなくなるときである。

40

**【 0 1 7 2 】**

R A I D 1 0 空間が残り少なくなっており、使用可能なディスク空間全体が残り少なくなっている場合には、データ・プログレッションは、R A I D 1 0 ディスク空間を取り上げて、より効率のよい R A I D 5 へと移動させることができる。これにより、システムの

50

全体容量が、書き込み性能を犠牲にして、増加する。より多くのディスクはやはり必要である。特定のストレージ・クラスが完全に使用された場合、データ・プログレッションは、受け入れ可能でないページを借りてシステムを動作させ続けることができる。例えば、あるボリュームが、RAID10 FCを、そのアクセス可能情報用に使用するようにコンフィギュレーションされている場合に、より多くのRAID10 - FC空間が使用可能になるまで、RAID5 - FCまたはRAID10 - SATAからページを割り当てる。

【0173】

また、データ・プログレッションは、システムの知覚される容量を増加させるための圧縮をサポートする。圧縮が使用されるのは、アクセスされていない履歴ページに対して、またはリカバリ情報の記憶のためにである。圧縮は、ストレージ・コストの最下部近くの別のクラスのストレージとして現れる。

10

【0174】

図25に示すように、ページ・プールは、本質的に、フリー・リストおよび装置情報を含む。ページ・プールは、複数のフリー・リスト、強化されたページ割り当て方式、およびフリー・リストの分類をサポートする必要がある。ページ・プールは、ストレージのクラスごとに別々のフリー・リストを維持する。この割り当て方式により、ページを多くのプールのうちの1つから割り当て、且つ、許可される最小または最大のクラスを設定する。フリー・リストの分類は、装置のコンフィギュレーションに由来する。それぞれのフリー・リストは、統計の収集および表示のためにそれ自体のカウンタを提供する。また、それぞれのフリー・リストは、ストレージ効率の統計の収集のためのRAID装置効率情報を提供する。

20

【0175】

一実施形態では、装置リストには、ストレージ・クラスのコストを追跡する更なる機能が必要であり得る。組合せにより、ストレージのクラスが決定される。これが起きるのは、ユーザが、コンフィギュレーションしたクラスの粒度がより荒いまたは細かいものであることを希望する場合である。

【0176】

図26に、高性能のデータベースの一実施形態を示しており、ここでは、すべてのアクセス可能なデータは、最近アクセスされていない場合でも、2.5FCドライブ上にある。アクセス不能な履歴データは、RAID5ファイバ・チャンネルへ移動される。

30

【0177】

図27は、MRI画像ボリュームの一実施形態を示しており、ここでは、アクセス可能なストレージは、この動的なボリューム用のSATA RAID10およびRAID5である。画像が最近アクセスされていない場合、画像はRAID5へ移動される。その場合、新しい書き込みは、初めはRAID10へ向かう。図19は、圧縮されたページ・レイアウトの一実施形態を示す。データ・プログレッションは、圧縮を、固定サイズのデータ・ページを更に細かく割り当てることによって、実装する。更に細かい割り当て(sub-allocation、サブ・アロケーション)の情報は、ページの空きの部分、およびページの割り当て済みの部分の場所を追跡する。データ・プログレッションは、圧縮の効率を予測することはできないが、そのサブ・アロケーション範囲内で可変サイズのページを取り扱うことができる。

40

【0178】

圧縮されたページは、CPUの性能に著しく影響することがある。書き込みアクセスの場合、圧縮されたページは、ページ全体の圧縮解除および再圧縮が必要となるはずである。従って、アクティブにアクセスされているページは、圧縮されず、非圧縮の状態に戻される。書き込みは、ストレージが極度に限られた状況で必要となり得る。

【0179】

PITC再マッピング・テーブルは、サブ・アロケーション情報を指し、圧縮されているページを示すようにマークされる。圧縮されたページにアクセスするには、非圧縮のページよりも高いI/Oカウントが必要になり得る。アクセスには、実際のデータの場所を

50

取り出すために、サブ・アロケーション情報の読み出しが必要になる。圧縮されたデータは、ディスクから読み出し、プロセッサで圧縮解除することができる。

【0180】

データ・プログレッションでは、圧縮は、ページ全体のうちの一部を圧縮解除できるようにすることが必要とされ得る。これにより、データ・プログレッションの読み出しアクセスは、そのページの小さな部分のみの圧縮解除を行える。読み出しキャッシュの先読み機能は、圧縮の遅延に対して役立ち得る。一つの圧縮解除で、幾つものサーバI/Oを取り扱うことができる。データ・プログレッションは、圧縮の候補として適さないページにマークし、それによって、或るページの圧縮を続けて試みないようにする。

【0181】

図20に、本発明の諸原理に従った、高レベルのディスク・ドライブ・システムにおけるデータ・プログレッションの一実施形態を示す。データ・プログレッションは、ボリュームの外面的な振る舞いやデータ・パスの動作を変更しない。データ・プログレッションは、ページ・プールへの変更を必要とし得る。ページ・プールは、本質的に、フリー・リストおよび装置情報を含む。ページ・プールは、複数のフリー・リスト、強化されたページ割り当て方式、およびフリー・リストの分類をサポートする必要がある。ページ・プールは、ストレージのそれぞれのクラスに対して個々のフリー・リストを維持する。割り当て方式は、許可される最小および最大のクラスを設定しつつ、ページが多くのプールのうちの1つから割り当てられるようにする。フリー・リストの分類は、装置のコンフィギュレーションに由来することができる。それぞれのフリー・リストは、統計の収集および表示のためにそれ自体のカウントを提供する。また、それぞれのフリー・リストは、ストレージの効率の統計の収集のためにRAID装置効率情報を提供する。

【0182】

PITCは、移動の候補を識別し、その移動のときには、アクセス可能なページへのI/Oを妨げる。データ・プログレッションは、候補に関してPITCを継続的に検査する。ページのアクセス可能性は、サーバI/O、新しいスナップショット・ページ更新、およびビュー・ボリュームの作成/削除に起因して、継続的に変化する。また、データ・プログレッションは、ボリュームのコンフィギュレーションの変更を継続的に検査し、ページのクラスおよびカウントの現在のリストをまとめあげる。これにより、データ・プログレッションは、そのまとめ(一覧)を評価し、移動するページがあるかを判定できる。

【0183】

それぞれのPITCは、ストレージの各クラスに使用されるページの数に関するカウントを提示する。データ・プログレッションは、この情報を使用して、閾値に達したときにページを移動させるための良い候補になるPITCを識別する。

【0184】

RAIDは、ディスクのコストに基づいて、1組のディスクから装置を割り当てる。また、RAIDは、装置または潜在的な装置の効率を取り出すAPIを提供する。また、これは、書き込み動作に必要なI/Oの数についての情報を返す必要がある。また、データ・プログレッションは、また、サード・パーティのRAIDコントローラをデータ・プログレッションの一部として使用するために、RAID NULLを必要とし得る。RAID NULLは、ディスク全体を消費し得、単にパス・スルー・レイヤとして働き得る。

【0185】

また、ディスク・マネージャは、ディスク分類を自動的に決定し保存することができる。ディスク分類を自動的に決定するには、SCSIイニシエータへの変更が必要となり得る。

【0186】

上述の説明および図面から、ここに示し説明した特定の諸実施形態は、もっぱら例示を目的としており、本発明の範囲を限定するものではないことが当業者には理解されよう。本発明を他の特定の形式で実施することは、その趣旨または本質的な諸特徴から逸脱することなく行えることが当業者には理解されよう。特定の諸実施形態の詳細への言及は、本

10

20

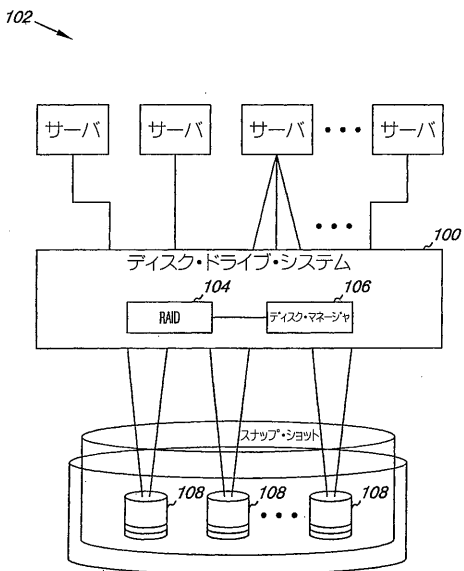
30

40

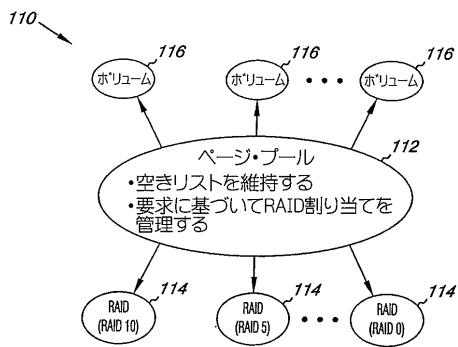
50

発明の範囲を限定するものではない。

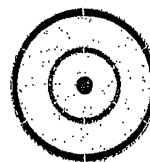
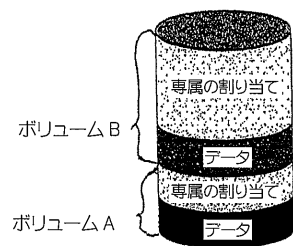
【図1】



【図2】



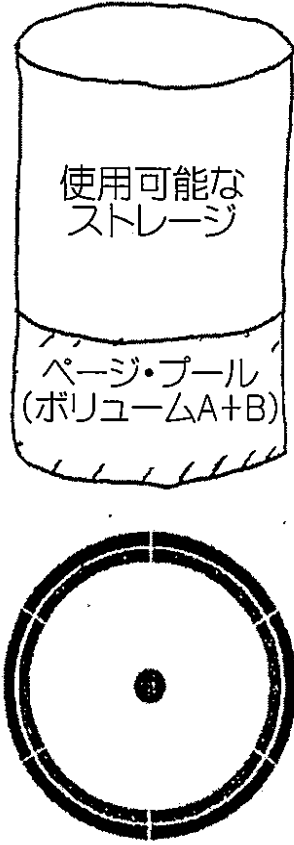
【図2A】



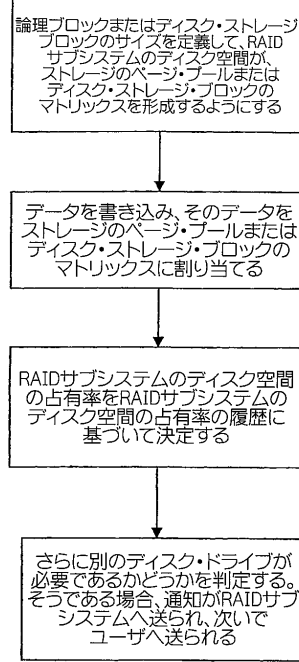
(従来技術)



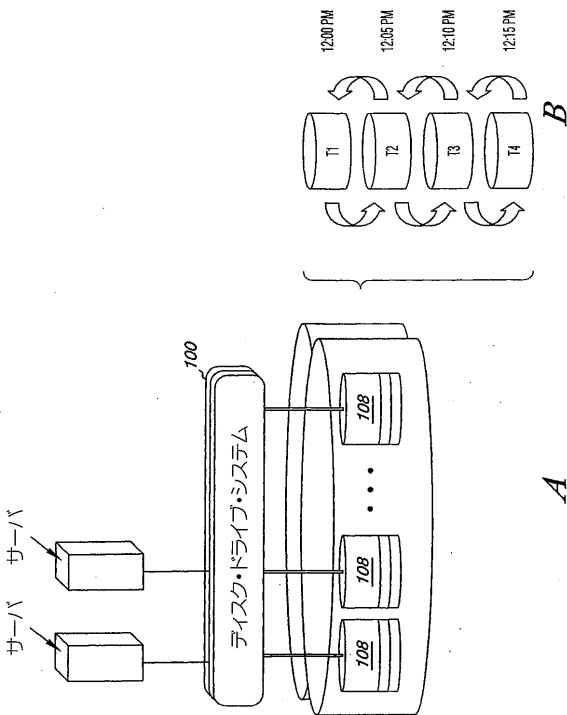
【図2B】



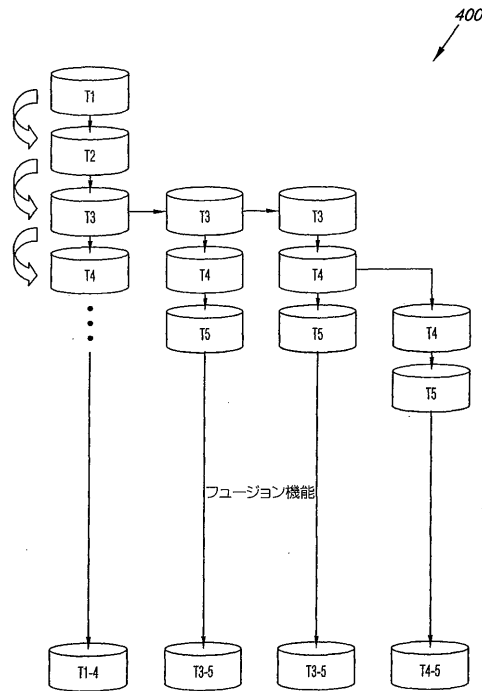
【図2C】



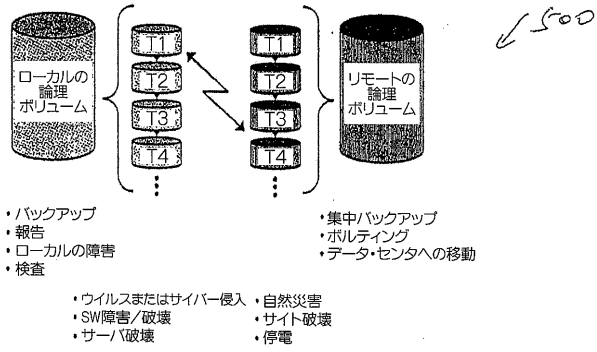
【図3】



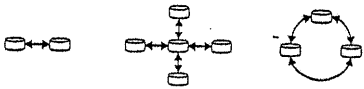
【図4】



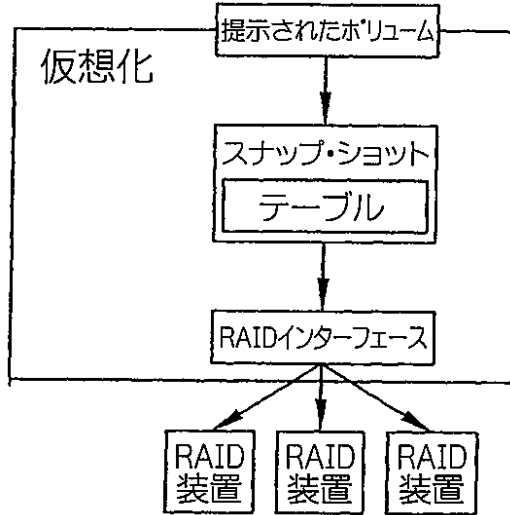
【 図 5 】



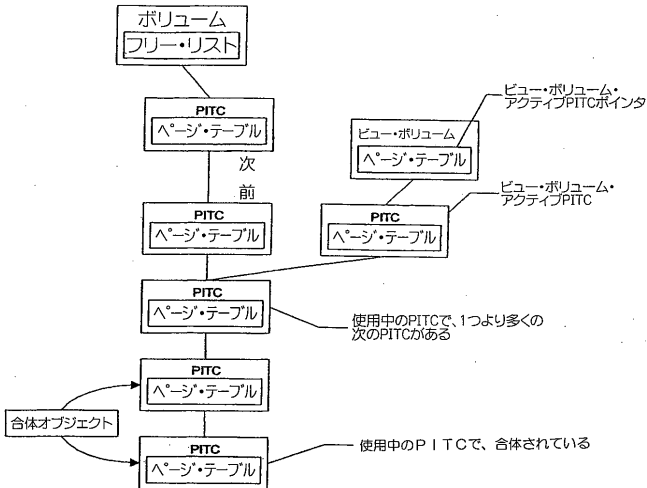
ポイント・ツー・ポイント    ポイント・ツー・マルチポイント    ピア・ツー・ピア



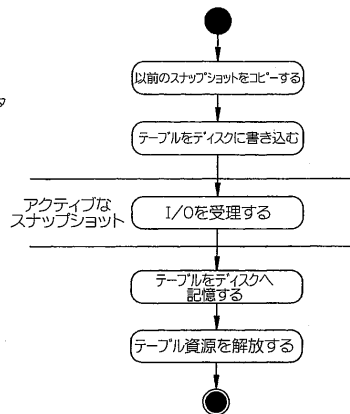
【 図 6 】



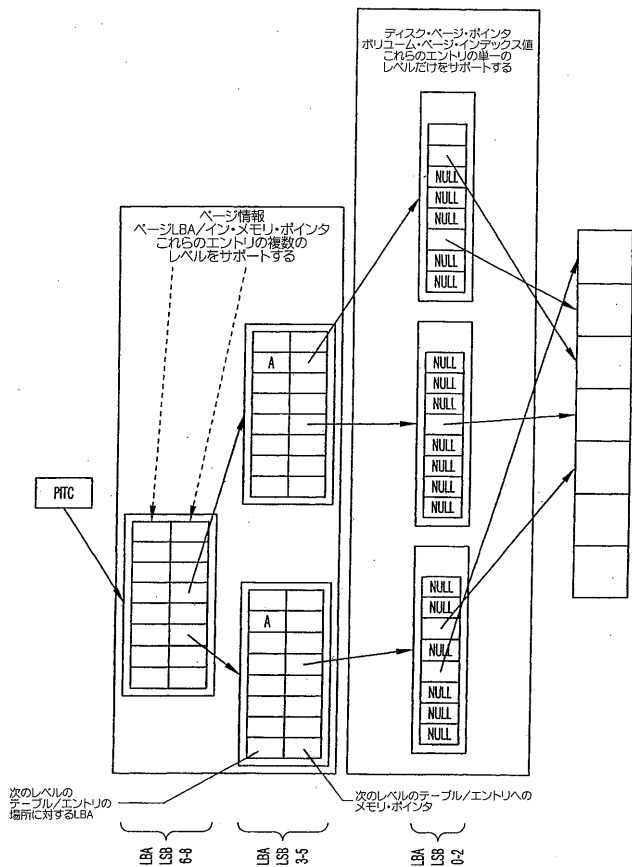
【 図 7 】



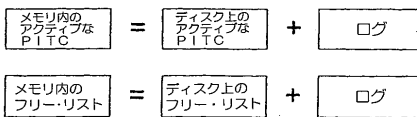
【 図 8 】



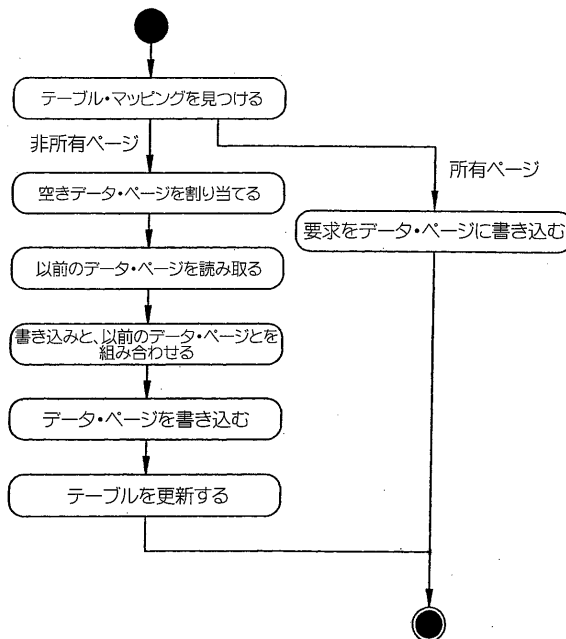
【 図 9 】



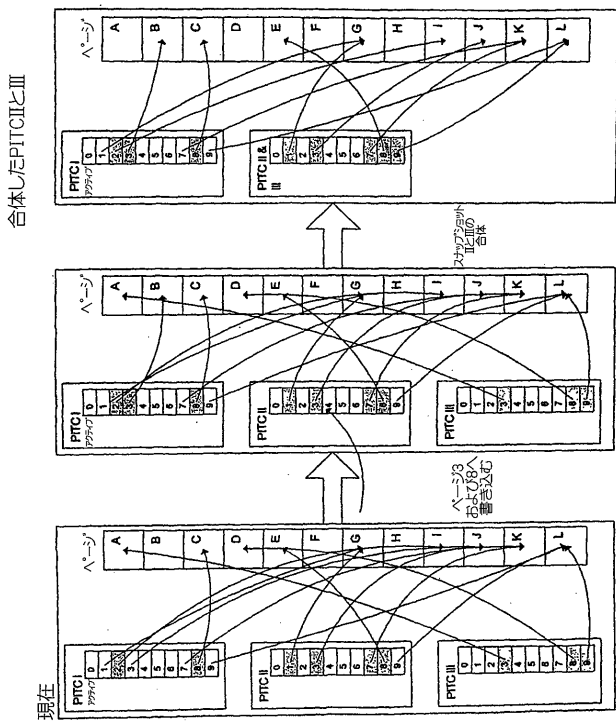
【 図 10 】



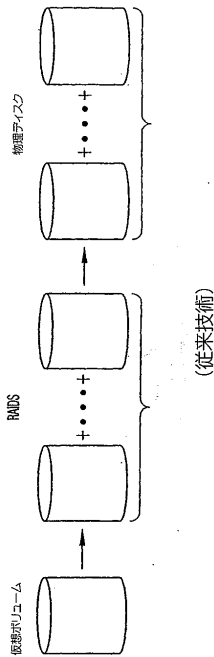
【 図 11 】



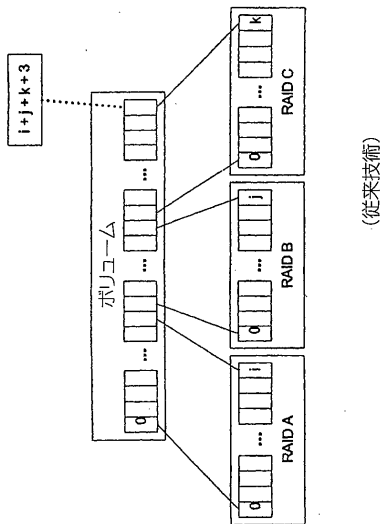
【 図 12 】



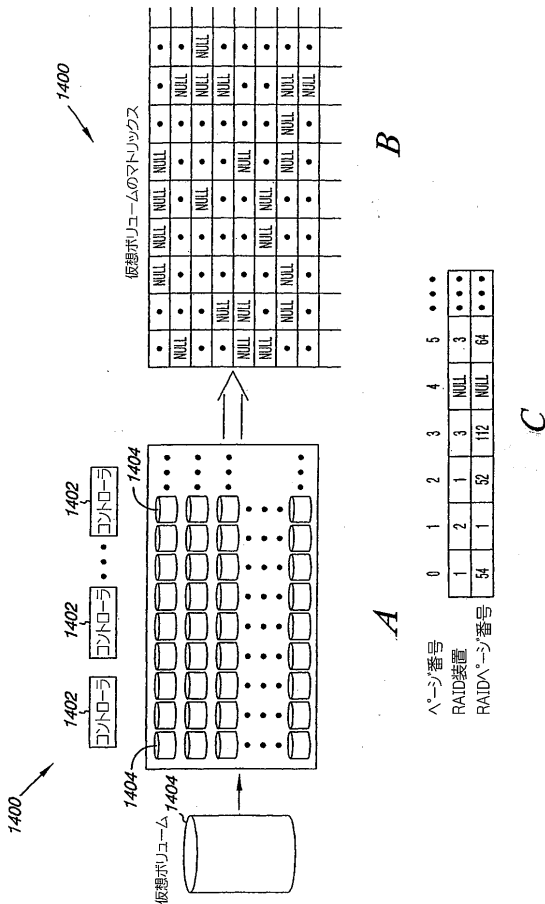
【 図 13 A 】



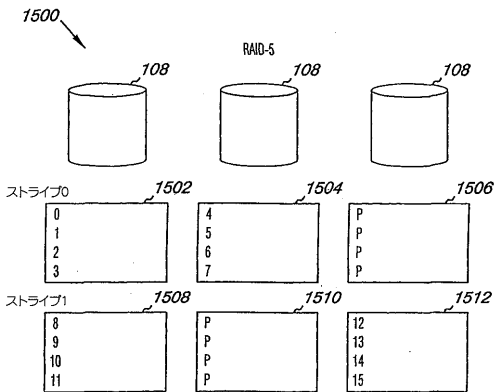
【図 13 B】



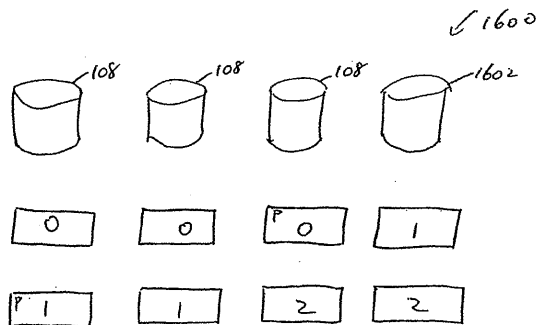
【図 14】



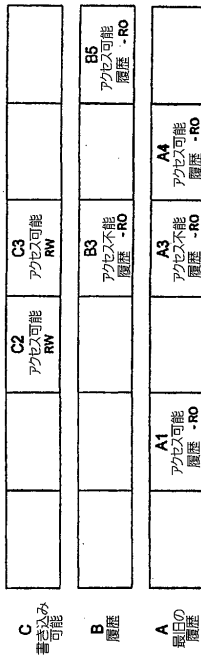
【図 15】



【図 16】



【図 17】



【 図 1 8 】

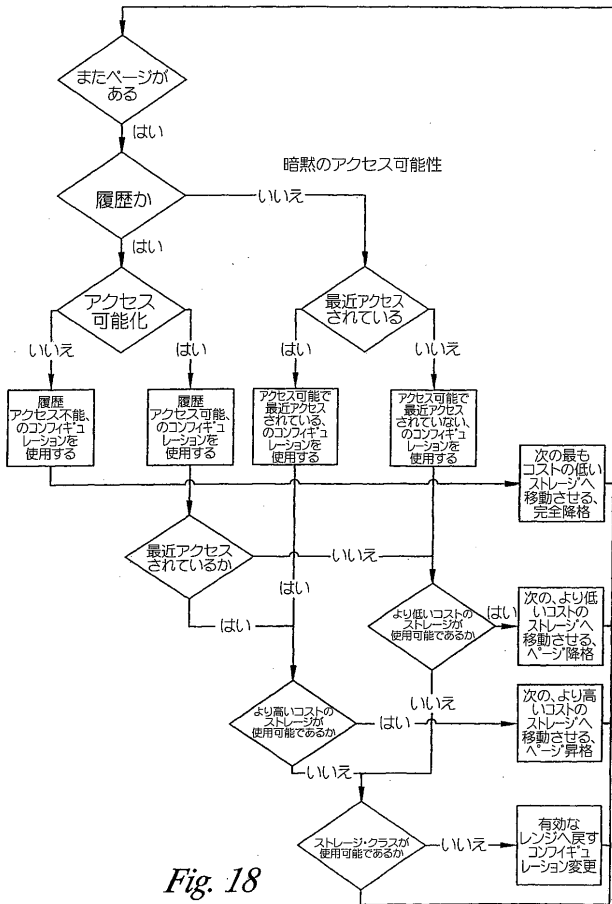
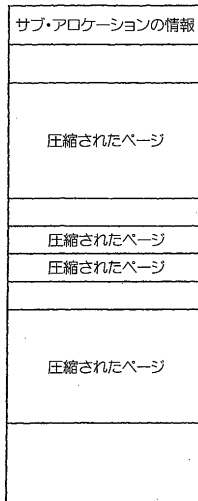
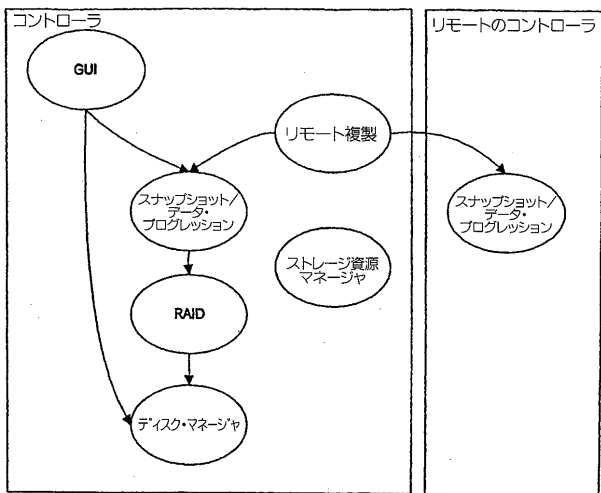


Fig. 18

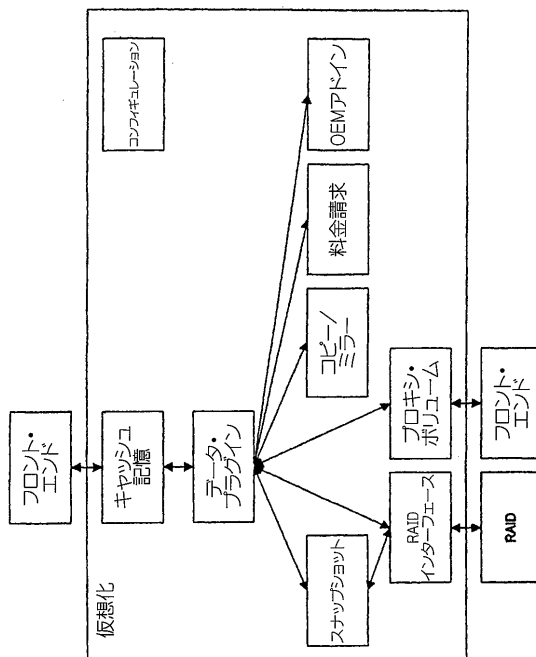
【 図 1 9 】



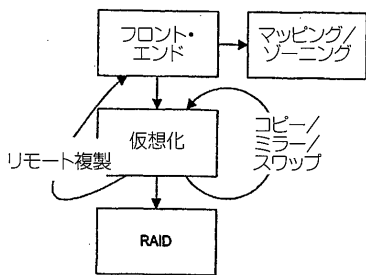
【 図 2 0 】



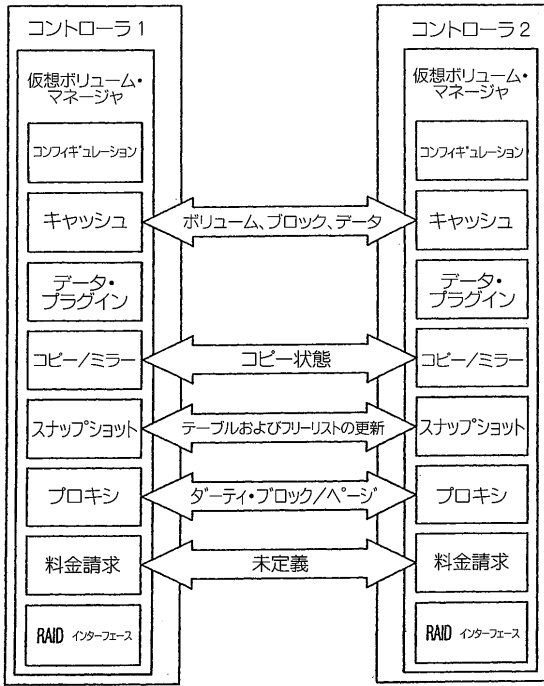
【 図 2 2 】



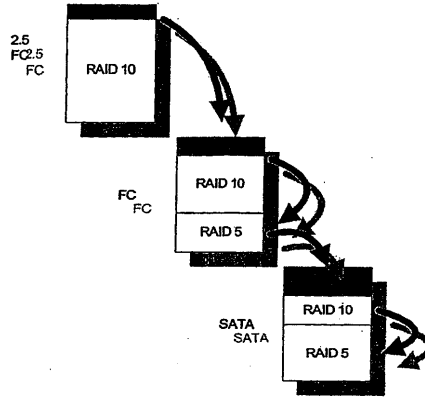
【 図 2 1 】



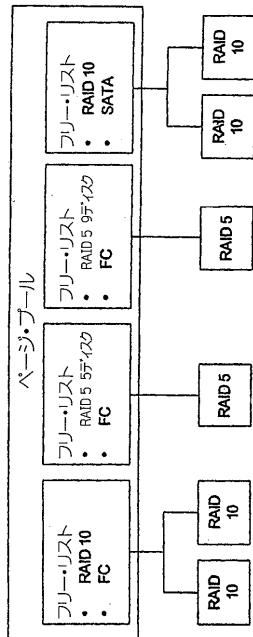
【 図 2 3 】



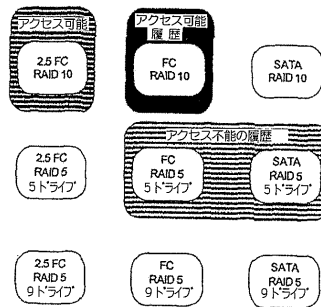
【 図 2 4 】



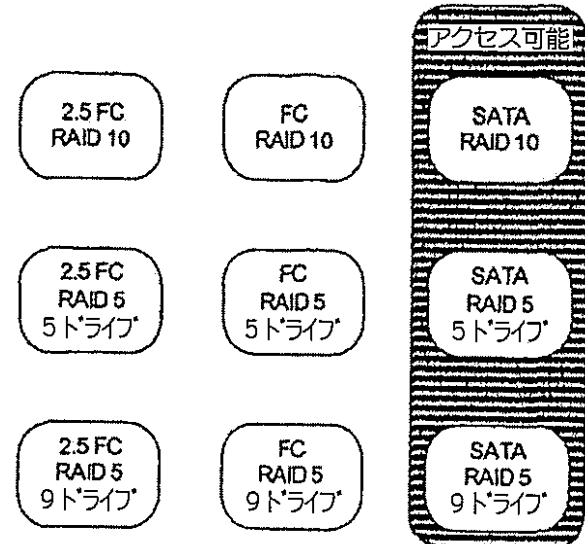
【 図 2 5 】



【 図 2 6 】



【 図 2 7 】



**【手続補正書】**

**【提出日】**平成22年9月22日(2010.9.22)

**【手続補正1】**

**【補正対象書類名】**特許請求の範囲

**【補正対象項目名】**全文

**【補正方法】**変更

**【補正の内容】**

**【特許請求の範囲】**

**【請求項1】**

データ・インスタント・リプレイの方法であって、  
複数のデータ・ストレージ・デバイスからデータ空間ブロックの抽象的なボリュームを作るステップと、

前記抽象的なボリュームのデータのスナップショットを生成するステップと、

前記スナップショットのアドレス・インデックスを記憶するステップと、

記憶された前記アドレス・インデックスを介して前記スナップショットにアクセスするステップと

を備える方法。

**【請求項2】**

請求項1に記載の方法であって、

前記抽象的なボリュームの前記データ空間ブロックを割り当てるステップと、

割り当てられた前記データ空間ブロックへデータを書くステップと

を更に備える方法。

**【請求項3】**

請求項1または2に記載の方法であって、前記スナップショットは所定の時間間隔において生成される、方法。

**【請求項4】**

請求項1ないし3の何れかに記載の方法であって、前記スナップショットはアプリケーションにより開始される、方法。

**【請求項5】**

請求項1ないし4の何れかに記載の方法であって、前記スナップショットを生成するステップは、書き込み動作を前記抽象的なボリュームに記録するステップを備える、方法。

**【請求項6】**

請求項1ないし5の何れかに記載の方法であって、前記スナップショットを生成するステップは、前記抽象的なボリュームに関しての変化を追跡するステップを備える、方法。

**【請求項7】**

請求項6に記載の方法であって、前記変化を追跡するステップは、少なくとも、前記抽象的なボリュームへのデルタ書き込みの部分的リストを維持するステップを備える、方法。

**【請求項8】**

請求項7に記載の方法であって、前記変化を追跡するステップは、少なくとも維持される前記部分的リストは、前のスナップショットの生成以降の、前記抽象的なボリュームの1つのデータ・アドレスへの最も新しい書き込みを含む、方法。

**【請求項9】**

請求項1ないし8の何れかに記載の方法であって、以前の時点での前記抽象的なボリュームのコンテンツのビューを提供する前記スナップショットに基づいて一時的なボリュームを作るステップを更に備える方法。

**【請求項10】**

請求項9に記載の方法であって、前記スナップショットおよび前記抽象的なボリュームにおける基礎となるデータを変更せずに、前記一時的なボリュームを変更するステップを更に備える方法。

**【請求項 1 1】**

請求項 1 0 に記載の方法であって、前記一時的なボリュームは検査に使用され、前記検査は、前記抽象的なボリュームにおける基礎となるデータを変えずに、前記以前の時点にあったデータに対して行われ得る、方法。

**【請求項 1 2】**

請求項 1 0 または 1 1 に記載の方法であって、前記一時的なボリュームは訓練に使用され、前記訓練は、前記抽象的なボリュームにおける基礎となるデータを変えずに、前記以前の時点にあったデータに対して行われ得る、方法。

**【請求項 1 3】**

請求項 1 0 ないし 1 2 の何れかに記載の方法であって、前記一時的なボリュームは、データを複製することにより該データのバックアップに用いられる、方法。

**【請求項 1 4】**

請求項 1 0 ないし 1 3 の何れかに記載の方法であって、前記一時的なボリュームは、データのリカバリに用いられる、方法。

**【請求項 1 5】**

データを再生可能なデスク・ドライブ・システムであって、  
データ空間ブロックの抽象化したものを有するデータ・ストレージ・システムと、  
少なくとも 1 つのデータ・ストレージ・システム・コントローラを有するデータ・マネージャであって、前記データ・ストレージ・システム・コントローラが、自動的に、前記データ空間ブロックの抽象化したもののスナップショットを生成し、前記スナップショットのアドレス・インデックスを記憶し、記憶された前記アドレス・インデックスを介して前記スナップショットにアクセスするものである、データ・マネージャと  
を備えるデスク・ドライブ・システム。

**【請求項 1 6】**

請求項 1 5 に記載のデスク・ドライブ・システムであって、前記システム・コントローラは、自動的に、前記スナップショットを所定の時間間隔において生成する、デスク・ドライブ・システム。

**【請求項 1 7】**

請求項 1 5 または 1 6 に記載のデスク・ドライブ・システムであって、前記システム・コントローラは、自動的に、アプリケーションの制御に基づいて前記スナップショットを生成する、デスク・ドライブ・システム。

**【請求項 1 8】**

請求項 1 5 ないし 1 7 の何れかに記載のデスク・ドライブ・システムであって、前記システム・コントローラは、前記データ空間ブロックの抽象化したものに対する書き込み動作を記憶する、デスク・ドライブ・システム。

**【請求項 1 9】**

請求項 1 5 ないし 1 8 の何れかに記載のデスク・ドライブ・システムであって、前記システム・コントローラは、前記データ空間ブロックの抽象化したものに関しての変化を追跡する、デスク・ドライブ・システム。

**【請求項 2 0】**

請求項 1 9 の何れかに記載のデスク・ドライブ・システムであって、前記システム・コントローラは、少なくとも、前記データ空間ブロックの抽象化したものに対するデルタ書き込みの部分的リストを維持する、デスク・ドライブ・システム。

**【請求項 2 1】**

請求項 2 0 の何れかに記載のデスク・ドライブ・システムであって、少なくとも維持される前記部分的リストは、前のスナップショットの生成以降の、前記データ空間ブロックの抽象化したものの 1 つのデータ・アドレスへの最も新しい書き込みを含む、デスク・ドライブ・システム。

**【請求項 2 2】**

請求項 1 5 ないし 2 1 の何れかに記載のデスク・ドライブ・システムであって、前記シ



システム・コントローラは、以前の時点での前記データ空間ブロックの抽象化したもののコンテンツのビューを提供する前記スナップショットに基づいて一時的なボリュームを作る、デスク・ドライブ・システム。

【請求項 23】

請求項 22 に記載のデスク・ドライブ・システムであって、前記システム・コントローラは、前記一時的なボリュームをデータのリカバリに用いる、デスク・ドライブ・システム。

---

フロントページの続き

- (72)発明者 ソラン, フィリップ・イー  
アメリカ合衆国ミネソタ州5 5 3 4 7, エデン・プレイリー, エイムズバリー・レイン 9 5 0 1
- (72)発明者 ガイダー, ジョン・ピー  
アメリカ合衆国ミネソタ州5 5 1 2 7, ノース・オークス, キャットバード・レイン 7
- (72)発明者 アスズマン, ローレンス・イー  
アメリカ合衆国ミネソタ州5 5 3 7 2, プライアー・レイク, ショア・トレイル・ノースイースト  
5 4 4 5
- (72)発明者 クレム, マイケル・ジェイ  
アメリカ合衆国ミネソタ州5 5 3 0 5, ミネトンカ, リヴェンデル・レイン 2 3 0 1
- Fターム(参考) 5B065 BA06 CA30 CC03 EA02