

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5937599号
(P5937599)

(45) 発行日 平成28年6月22日 (2016. 6. 22)

(24) 登録日 平成28年5月20日 (2016. 5. 20)

(51) Int. Cl.

F I

G 0 6 F 3 / 0 6 (2006. 01)

G 0 6 F 3 / 0 6 3 0 5 C

G 0 6 F 3 / 0 6 5 4 0

請求項の数 7 (全 32 頁)

(21) 出願番号 特願2013-531648 (P2013-531648)
 (86) (22) 出願日 平成23年9月20日 (2011. 9. 20)
 (65) 公表番号 特表2013-539134 (P2013-539134A)
 (43) 公表日 平成25年10月17日 (2013. 10. 17)
 (86) 国際出願番号 PCT/US2011/052276
 (87) 国際公開番号 W02012/044492
 (87) 国際公開日 平成24年4月5日 (2012. 4. 5)
 審査請求日 平成26年9月12日 (2014. 9. 12)
 (31) 優先権主張番号 12/896, 669
 (32) 優先日 平成22年10月1日 (2010. 10. 1)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 513076589
 ビュア・ストレージ・インコーポレイテッド
 アメリカ合衆国・94041・カリフォルニア州・マウンテン ビュー・カストロ
 ストリート・650・スイート・220
 (74) 代理人 100064621
 弁理士 山川 政樹
 (74) 代理人 100098394
 弁理士 山川 茂樹
 (72) 発明者 コルグローヴ, ジョン
 アメリカ合衆国・94024・カリフォルニア州・ロス アルトス・ヴィスタ グラ
 ンデ アヴェニュー・722

最終頁に続く

(54) 【発明の名称】 動的構成のRAIDアレイにおける再構成読み込み

(57) 【特許請求の範囲】

【請求項 1】

独立ディスク冗長アレイ (RAID) における複数のストレージデバイスを備えたデータストレージサブシステムと、ストレージコントローラを備え、

前記ストレージコントローラは、

第1のRAIDストライプを、前記複数のストレージデバイスに書き込むように構成され、前記複数のストレージデバイスの特定のストレージデバイスに、前記特定のストレージデバイスに保存された前記第1のRAIDストライプの第1のRAIDデータを保護するためのデバイス内保護データの第1の量を配置することを含み、

第2のRAIDストライプを書き込むための前記複数のストレージデバイスのストレージデバイスのサブセットを選択するように構成され、前記サブセットのストレージデバイスの量は、前記複数のストレージデバイスのストレージデバイスの量よりも小さく、前記サブセットは、前記特定のストレージデバイスを含み、

第2のRAIDストライプを、前記サブセットに書き込むように構成され、前記特定のストレージデバイスに、前記特定のストレージデバイスに保存された前記第2のRAIDストライプの第2のRAIDデータを保護するためのデバイス内保護データの第2の量を配置することを含み、前記デバイス内保護データの第2の量は、前記デバイス内保護データの第1の量よりも大きい、

コンピュータシステム。

【請求項 2】

10

20

前記複数のストレージデバイスは、ソリッドステートストレージデバイスである、
ことを特徴とする請求項 1 に記載のコンピュータシステム。

【請求項 3】

前記第 1 の RAID ストライプは、 $L + x$ のレイアウトである第 1 の RAID レイアウトを有し、前記第 2 の RAID ストライプは、 $M + y$ のレイアウトである第 2 の RAID レイアウトを有し、 L 、 x 、 M 及び y は整数であり、(1) L は M に等しくない、及び (2) x は y に等しくない、のいずれか又は両方である、
ことを特徴とする請求項 1 に記載のコンピュータシステム。

【請求項 4】

前記複数のストレージデバイスにおける少なくとも 1 つのストレージデバイスは、前記特定のストレージデバイスにおける前記デバイス内保護データの第 1 の量とは異なるデバイス内保護データの量を有する、

10

ことを特徴とする請求項 1 に記載のコンピュータシステム。

【請求項 5】

前記ストレージデバイスのサブセットにおける少なくとも 1 つのストレージデバイスは、前記特定のストレージデバイスにおける前記デバイス内保護データの第 2 の量とは異なるデバイス内保護データの量を有する、

ことを特徴とする請求項 1 に記載のコンピュータシステム。

【請求項 6】

独立ディスク冗長アレイ (RAID) における複数のストレージデバイスを構成するための方法であって、

20

第 1 の RAID ストライプを、前記複数のストレージデバイスに書き込むステップであって、前記複数のストレージデバイスの特定のストレージデバイスに、前記特定のストレージデバイスに保存された前記第 1 の RAID ストライプの第 1 の RAID データを保護するためのデバイス内保護データの第 1 の量を配置することを含むステップと、

第 2 の RAID ストライプを書き込むための前記複数のストレージデバイスのストレージデバイスのサブセットを選択するステップであって、前記サブセットのストレージデバイスの量は、前記複数のストレージデバイスのストレージデバイスの量よりも小さく、前記サブセットは、前記特定のストレージデバイスを含むステップと、

第 2 の RAID ストライプを、前記サブセットに書き込むステップであって、前記特定のストレージデバイスに、前記特定のストレージデバイスに保存された前記第 2 の RAID ストライプの第 2 の RAID データを保護するためのデバイス内保護データの第 2 の量を配置することを含み、前記デバイス内保護データの第 2 の量は、前記デバイス内保護データの第 1 の量よりも大きい、ステップと、
を含むことを特徴とする方法。

30

【請求項 7】

独立ディスク冗長アレイ (RAID) における複数のストレージデバイスを構成ためのプログラム命令を記憶するコンピュータ可読記憶媒体であって、前記プログラム命令は、

第 1 の RAID ストライプを、前記複数のストレージデバイスに書き込むように実行可能であり、前記複数のストレージデバイスの特定のストレージデバイスに、前記特定のストレージデバイスに保存された前記第 1 の RAID ストライプの第 1 の RAID データを保護するためのデバイス内保護データの第 1 の量を配置することを含み、

40

第 2 の RAID ストライプを書き込むための前記複数のストレージデバイスのストレージデバイスのサブセットを選択するように実行可能であり、前記サブセットのストレージデバイスの量は、前記複数のストレージデバイスのストレージデバイスの量よりも小さく、前記サブセットは、前記特定のストレージデバイスを含み、

第 2 の RAID ストライプを、前記サブセットに書き込むように実行可能であり、前記特定のストレージデバイスに、前記特定のストレージデバイスに保存された前記第 2 の RAID ストライプの第 2 の RAID データを保護するためのデバイス内保護データの第 2 の量を配置することを含み、前記デバイス内保護データの第 2 の量は、前記デバイス内保

50

護データの第1の量よりも大きい、
コンピュータ可読記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータネットワークに関し、より詳細には、複数のソリッドステートストレージデバイス間でデータを効率的に分散させる発明に関する。

【背景技術】

【0002】

コンピュータのメモリストレージ及びデータ帯域幅が増すと、企業が日々管理するデータの量及び複雑性も増す。通常、データセンターなどの大規模分散型ストレージシステムは、多くの業務を実行する。分散型ストレージシステムは、1又はそれ以上のネットワークにより相互接続されたクライアントコンピュータに結合することができる。分散型ストレージシステムのいずれかの部分が不良を起こし、又は利用できなくなった場合には、企業活動が損なわれ、又は完全に停止する恐れがある。従って、分散型ストレージシステムは、データの利用可能性及び高パフォーマンス機能のための高い標準を維持すると予想される。本明細書で使用するストレージディスクは、ストレージ技術のタイプによってはディスクを含まないものもあるので、ストレージデバイスと呼ぶことができる。

【0003】

多くの場合、ストレージデバイスは、データ損失からの保護のために、エラー検出機構及びエラー訂正機構を含む。多くの場合、これらの機構は、デバイスにより生成されてデバイス自体に記憶されるエラー訂正符号の形をとる。また、分散型ストレージシステムは、分散アルゴリズムを利用して、一群のストレージデバイス間でデータを分散させることもできる。一般に、これらのアルゴリズムは、中央ディレクトリに依拠せずにデータオブジェクトをストレージデバイスにマッピングする。このようなアルゴリズムの例に、レプリケーション・アンダー・スケーラブル・ハッシング(RUSH)及びコントロールド・レプリケーション・アンダー・スケーラブル・ハッシング(CRUSH)がある。分散型ストレージシステム内の複数のクライアントは、中央ディレクトリを伴わずに複数のサーバ上のデータオブジェクトに同時にアクセスすることができる。また、記憶されているメタデータの量を低減することもできる。しかしながら、容量、入力/出力(I/O)特性及び信頼性問題が異なる複数のストレージディスク間でデータを分散させるという困難なタスクが依然として残る。ストレージデバイス自体と同様に、これらのアルゴリズムも、(RAID5及びRAID6などの)RAIDタイプアルゴリズム又はリードソロモン符号などのエラー検出及び訂正アルゴリズムを含むことができる。

【0004】

動的に追加及び除去できる複数のストレージデバイス間でデータを分散させるために使用する方法は、選択したストレージデバイスに関連する技術及び機構によって決まる。例えば、上述したアルゴリズムは、ハードディスクドライブ(HDD)を利用するシステムに合わせて開発されたものである。HDDは、各々が磁気媒体で被覆された1又はそれ以上の回転ディスクを含む。これらのディスクは、毎日数時間にわたり毎分数千回転の速さで回転する。また、この回転ディスク上への磁気読み込み/書き込み装置の位置付けには、磁気アクチュエータが関与する。これらのアクチュエータは、摩擦、摩耗、振動及び機械的不均衡の影響を受けやすく、結果的に信頼性問題が生じる。上述したデータ分散アルゴリズムは、HDDのこれらの特性及び挙動に基づくものである。

【0005】

別のタイプの記憶ディスクの例に、ソリッドステートディスク(SSD)がある。ソリッドステートディスクは、ソリッドステートドライブと呼ぶこともできる。SSDは、HDDインターフェイスをエミュレートできるが、HDDで見られるような電気機械デバイスではなく固体メモリを利用して永続データを記憶する。例えば、SSDは、一群のフラッシュメモリを含むことができる。可動部品又は機械的遅延がなければ、SSDのアクセ

10

20

30

40

50

ス時間及びレイテンシは、HDDよりも短くなり得る。しかしながら、通常、SSDの書き込みレイテンシはかなり長い。入力/出力(I/O)特性が異なることに加え、SSDの故障モードもHDDとは異なる。従って、記憶のためにSSDを備えたシステムでは、HDDに合わせて開発された分散型データ配置アルゴリズムを利用しながら高パフォーマンス及び高信頼性を実現できない場合がある。

【発明の概要】

【発明が解決しようとする課題】

【0006】

上記に鑑み、複数のソリッドステートストレージデバイス間でデータを効率的に分散し、エラーを検出して訂正するためのシステム及び方法が望まれている。

10

【課題を解決するための手段】

【0007】

複数のソリッドステートストレージデバイス間でデータを効率的に分散して管理するためのコンピュータシステム及び方法の様々な実施形態を開示する。

【0008】

1つの実施形態では、コンピュータシステムが、ネットワークを介して読み込み及び書き込み要求を受け取るように結合された1又はそれ以上のデータストレージアレイにネットワークを介して読み込み及び書き込み要求を伝達するように構成された複数のクライアントコンピュータを備える。複数のストレージデバイス上に複数の記憶位置を有する(単複の)データストレージアレイを企図する。様々な実施形態では、このストレージデバイスが、データを記憶して保護するための独立ドライブ冗長アレイ(RAID)構成で構成される。データストレージデバイスは、フラッシュメモリセルなどの、データ記憶のための固体メモリ技術を含むことができる。データストレージサブシステムは、ストレージコントローラをさらに備え、このストレージコントローラは、ストレージデバイスの第1のサブセットを、第1の冗長データセットを含む第1のRAIDレイアウトで使用するよう

20

30

20

30

【0009】

第1のRAIDレイアウト又は第2のRAIDレイアウトの特定のストレージデバイスを対象とする所与の読み込み要求を受け取ったことに応答して、コントローラが、特定のストレージデバイスが非エラーに関する相対的に遅い読み込み応答を示していると判断したことに応答して、所与の読み込み要求に対応する再構成読み込みを開始するように構成される実施形態も企図する。また、第1のサブセット及び第2のサブセット内の各デバイスは、デバイス内冗長データを記憶するように構成することができ、第1の冗長データセット及び第2の冗長データセットは、いずれもデバイス間冗長データを含む。さらに、様々な実施形態では、第1のRAIDレイアウトが、 $L+x$ のレイアウトであり、第2のRAIDレイアウトが、 $M+y$ のレイアウトであり、 L 、 x 、 M 及び y は整数であり、(1)

40

40

【0010】

以下の説明及び添付図面を検討すると、これらの及びその他の実施形態が明らかになるであろう。

【図面の簡単な説明】

【0011】

【図1】ネットワークアーキテクチャの1つの実施形態を示す汎用ブロック図である。

【図2】動的デバイス内冗長スキームの1つの実施形態の汎用ブロック図である。

【図3】データストレージサブシステムにおいてデバイス内保護を調整する方法の1つの実施形態を示す汎用フロー図である。

50

【図４】ストレージサブシステムの１つの実施形態の汎用ブロック図である。

【図５】デバイスユニットの１つの実施形態の汎用ブロック図である。

【図６】状態テーブルの１つの実施形態を示す汎用ブロック図である。

【図７】フレキシブルなＲＡＩＤデータレイアウトアーキテクチャの１つの実施形態を示す汎用ブロック図である。

【図８】フレキシブルなＲＡＩＤデータレイアウトアーキテクチャの別の実施形態を示す汎用ブロック図である。

【図９】データストレージサブシステム内のレイアウトを動的に決定する方法の１つの実施形態を示す汎用フロー図である。

【図１０】フレキシブルなＲＡＩＤデータレイアウトアーキテクチャのさらに別の実施形態を示す汎用ブロック図である。

【図１１Ａ】デバイスレイアウトの１つの実施形態を示す図である。

【図１１Ｂ】セグメントの１つの実施形態を示す図である。

【図１１Ｃ】異なるページタイプ内のデータストレージ構成の１つの実施形態を示す汎用ブロック図である。

【図１２】ハイブリッドＲＡＩＤデータレイアウトの１つの実施形態を示す汎用ブロック図である。

【図１３】データストレージサブシステム内で代替のＲＡＩＤ構成を選択する方法の１つの実施形態を示す汎用フロー図である。

【発明を実施するための形態】

【００１２】

本発明は様々な修正及び代替形態が可能であるが、図面には特定の実施形態を一例として示し、本明細書ではこれらについて詳細に説明する。しかしながら、図面及びこれらに対する詳細な説明は、開示する特定の形態に本発明を限定することを意図するものではなく、むしろ添付の特許請求の範囲によって定められる本発明の思想及び範囲内にある全ての修正物、同等物及び代替物を含むことを意図するものであると理解されたい。

【００１３】

以下の説明では、本発明を完全に理解できるように数多くの具体的な詳細を示す。しかしながら、当業者であれば、これらの具体的な詳細を伴わずに本発明を実施できると認識すべきである。いくつかの例では、本発明を曖昧にしないように、周知の回路、構造、信号、コンピュータプログラム命令及び技術については詳細に示していない。

【００１４】

図１を参照すると、ネットワークアーキテクチャ１００の１つの実施形態の汎用ブロック図を示している。後述するように、ネットワークアーキテクチャ１００の１つの実施形態は、ネットワーク１８０を介して互いに、及びデータストレージレイ１２０ａ～１２０ｂに相互接続されたクライアントコンピュータシステム１１０ａ～１１０ｂを含む。ネットワーク１８０は、スイッチ１４０を介して第２のネットワーク１９０に結合することができる。このネットワーク１９０を介して、クライアントコンピュータシステム１１０ｃが、クライアントコンピュータシステム１１０ａ～１１０ｂ及びデータストレージレイ１２０ａ～１２０ｂに結合される。また、ネットワーク１９０は、スイッチ１５０を介してインターネット１６０又はその他の外部ネットワークに結合することもできる。

【００１５】

なお、代替の実施形態では、クライアントコンピュータ及びサーバ、スイッチ、ネットワーク、データストレージレイ及びデータストレージデバイスの数及びタイプが、図１に示すものに限定されない。１又はそれ以上のクライアントは、様々な時点でオフライン動作することができる。また、動作中、ユーザがネットワークアーキテクチャ１００への接続、切断及び再接続を行うと、個々のクライアントコンピュータの接続タイプが変化することもある。図１に示す構成要素の各々のさらなる説明を手短に行う。まず、データストレージレイ１２０ａ～１２０ｂにより提供される機能のいくつかの概要について説明する。

【0016】

ネットワークアーキテクチャ100では、データストレージレイ120a~120bの各々を、クライアントコンピュータシステム110a~110cなどの異なるサーバ及びコンピュータ間のデータの共有に使用することができる。また、データストレージレイ120a~120bを、ディスクのミラーリング、バックアップ及び復元、保存データの保管及び検索、並びにストレージデバイス間のデータ移行に使用することもできる。代替の実施形態では、クラスタを形成するために、1又はそれ以上のクライアントコンピュータシステム110a~110cを、高速ローカルエリアネットワーク(LAN)を介して互いにリンクさせることができる。互いにリンクされた1又はそれ以上のノードはクラスタを形成し、これによりデータストレージレイ120a~120bの1つに存在するクラスタ共有ボリュームなどのストレージリソースを共有することができる。

10

【0017】

データストレージレイ120a~120bの各々は、データ記憶のためのストレージサブシステム170を含む。ストレージサブシステム170は、複数のストレージデバイス176a~176mを含むことができる。これらのストレージデバイス176a~176mは、クライアントコンピュータシステム110a~110cにデータ記憶サービスを提供することができる。ストレージデバイス176a~176mの各々は、読み込み及び書き込み要求を受け取るとともに、各々をアレイ内の行及び列としてアドレス指定可能な複数のデータ記憶位置を含むように構成することができる。1つの実施形態では、ストレージデバイス176a~176m内のデータ記憶位置を、論理的で冗長なストレージコンテナ又はRAIDアレイ(低価格/独立ディスク冗長アレイ)内に配置することができる。しかしながら、ストレージデバイス176a~176mは、ディスクを含まないこともある。1つの実施形態では、ストレージデバイス176a~176mの各々が、従来のハードディスクドライブ(HDD)とは異なる技術をデータ記憶に利用することができる。例えば、ストレージデバイス176a~176mの1又はそれ以上は、永続データを記憶するための固体メモリから成るストレージを含み、又はこれにさらに結合することができる。他の実施形態では、ストレージデバイス176a~176mの1又はそれ以上が、スピン注入法、磁気抵抗ランダムアクセスメモリ(MRAM)法、又はその他の記憶技術を利用するストレージを含み、又はこのようなストレージにさらに結合することができる。これらの異なる記憶技術により、ストレージデバイス間で異なる信頼性特性が生じ得る。

20

30

【0018】

ストレージデバイス176a~176mの各々において使用される技術及び機構のタイプにより、データオブジェクトマッピング、並びにエラー検出及び訂正に使用するアルゴリズムを決定することができる。これらのアルゴリズムで使用されるロジックを、基本オペレーティングシステム(OS)116、ファイルシステム140、ストレージサブシステムコントローラ174内の1又はそれ以上のグローバルRAIDエンジン178、及びストレージデバイス176a~176mの各々における制御ロジックのうちの1又はそれ以上に含めることができる。

【0019】

1つの実施形態では、含まれる固体メモリが、ソリッドステートドライブ(SSD)技術を含む。通常、SSD技術は、フラッシュメモリセルを利用する。当業で周知のように、フラッシュメモリセルは、フローティングゲート内に捕捉され蓄積された電子の範囲に基づく二進値を保持する。完全に消去されたフラッシュメモリセルは、フローティングゲート内に電子を全く又は最低数しか蓄積していない。消去されたフラッシュメモリセルには、シングルレベルセル(SLC)フラッシュの二進1などの特定の二進値が関連付けられる。マルチレベルセル(MLC)フラッシュでは、消去されたフラッシュメモリセルに二進値11が関連付けられる。フラッシュメモリセル内の制御ゲートに所与の閾値電圧よりも高い電圧を印加した後、このフラッシュメモリセルは、フローティングゲート内に所与の範囲の電子を捕捉する。従って、プログラムされた(書き込まれた)フラッシュメモリセルには、SLCフラッシュの二進0などの別の特定の二進値が関連付けられる。ML

40

50

Cフラッシュセルでは、制御ゲートに印加された電圧に応じて、プログラムされたメモリセルに複数の二進値の1つを関連付けることができる。

【0020】

一般的に言えば、SSD技術では、読み込みアクセスレイテンシタイムがHDD技術よりも短い。しかしながら、SSDの書き込みパフォーマンスは、SSD内の未使用のプログラマブルブロックの利用可能性によって大きく影響を受ける。SSDの書き込みパフォーマンスは、SSDの読み込みパフォーマンスに比べて大幅に遅いので、同様のレイテンシを予想する一部の機能又は動作に関する問題が生じることがある。また、HDD技術とSSD技術の間の技術及び機構の違いにより、データストレージデバイス176a~176mの信頼性特性に違いが生じることがある。

10

【0021】

様々な実施形態では、SSD内のフラッシュセルに新たなデータが書き込まれる前に、一般にこのフラッシュセルを消去しなければならない。また、様々なフラッシュ技術における消去動作は、ブロック単位で行わなければならない。従って、ブロック(消去セグメント又は消去ブロック)内のフラッシュメモリセルは、全てまとめて消去される。フラッシュ消去ブロックは、複数のページを含むことができる。例えば、1ページのサイズが4キロバイト(KB)であり、1ブロックが64ページ、すなわち256KBを含むことができる。フラッシュデバイスでは、読み込み動作に比べて消去動作のレイテンシの方が相対的に高いことがあり、これにより対応する書き込み動作のレイテンシが増すことがある。フラッシュ技術のプログラミング又は読み込みは、消去ブロックサイズよりも低い粒度レベルで行うことができる。例えば、フラッシュセルには、バイトサイズ、単語サイズ又はその他のサイズでプログラム又は読み込みを行うことができる。

20

【0022】

フラッシュセルには、反復的な消去プログラム動作後に摩耗が生じる。この場合、この摩耗は、MLCフラッシュセルの基板とフローティングゲートの間の誘電酸化物層に注入され捕捉される電荷によるものである。1つの例では、MLCフラッシュセルが、10,000~100,000サイクルなどの、消去及びプログラム動作を受ける回数限界を有することができる。また、SSDには、別のフラッシュセルの消去又はプログラム中に隣接する又は近隣のフラッシュセルに偶発的な状態変化を引き起こすプログラムディスタートエラーが生じることがある。さらに、SSDは、別のフラッシュセルの読み込み中に近隣のフラッシュセルの偶発的な状態変化が生じるリードディスタートエラーも含む。

30

【0023】

1又はそれ以上のストレージデバイス176a~176mの各々の特性が分かると、より効率的なデータオブジェクトマッピング、並びにエラー検出及び訂正を行うことができる。1つの実施形態では、ストレージコントローラ174内のグローバルRAIDエンジン178が、ストレージデバイス176a~176mに関して、I/O要求の応答時間に一貫性がないこと、対応するアクセスに対するデータが誤っていること、エラー率及びアクセス率のうちの少なくとも1つ又はそれ以上を検出することができる。グローバルRAIDエンジン178は、少なくともこれらの特性にตอบสนองして、ストレージデバイス176a~176m内の対応するストレージデバイスグループにいずれのRAIDデータレイアウトアーキテクチャを利用すべきかを判断することができる。また、グローバルRAIDエンジン178は、ストレージデバイス176a~176mの特性に基づいて、デバイス内冗長スキーム及びデバイス間RAIDデータレイアウトを動的に変更することができる。

40

【0024】

図1に、1つの実施形態による、説明した特徴が可能なシステムの例を示す。さらなる詳細については以下で示す。以下、図1を参照しながら、ネットワークアーキテクチャ100の構成要素についてさらに説明する。

【0025】

ネットワークアーキテクチャの構成要素

50

繰り返すが、図示のように、ネットワークアーキテクチャ 100 は、ネットワーク 180 及び 190 を介して互いに及びデータストレージアレイ 120 a ~ 120 b に相互接続されたクライアントコンピュータシステム 110 a ~ 110 c を含む。ネットワーク 180 及び 190 は、無線接続、直接ローカルエリアネットワーク (LAN) 接続、ストレージエリアネットワーク (SAN)、インターネットなどの広域ネットワーク (WAN) 接続及びルータなどを含む様々な技術を含むことができる。ネットワーク 180 及び 190 は、1 又はそれ以上の LAN を含むことができ、これらは無線であってもよい。ネットワーク 180 及び 190 は、リモートダイレクトメモリアクセス (RDMA) ハードウェア及び/又はソフトウェア、伝送制御プロトコル/インターネットプロトコル (TCP/IP) ハードウェア及び/又はソフトウェア、ルータ、リピータ、スイッチ及び/又はグリッドなどをさらに含むことができる。ネットワーク 180 及び 190 内では、イーサネット (登録商標)、ファイバチャネル、ファイバチャネルオーバーイーサネット (FCoE) 及び iSCSI などのプロトコルを使用することができる。スイッチ 140 は、ネットワーク 180 及び 190 の両方に関連するプロトコルを利用することができる。ネットワーク 190 は、伝送制御プロトコル (TCP) 及びインターネットプロトコル (IP)、すなわち TCP/IP などの、インターネット 160 に使用される通信プロトコルの組と整合することができる。スイッチ 150 は、TCP/IP スwitch とすることができる。

【0026】

クライアントコンピュータシステム 110 a ~ 110 c は、デスクトップパソコン (PC)、ワークステーション、ラップトップ、ハンドヘルドコンピュータ、サーバ、サーバファーム、携帯情報端末 (PDA) 及びスマートフォンなどのあらゆる数の固定又はモバイルコンピュータを表す。一般的に言えば、クライアントコンピュータシステム 110 a ~ 110 c は、1 又はそれ以上のプロセッサコアを備えた 1 又はそれ以上のプロセッサを含む。各プロセッサコアは、所定の汎用命令セットに従って命令を実行するための回路を含む。例えば、x86 命令セットアーキテクチャを選択することができる。或いは、Alpha (登録商標)、PowerPC (登録商標)、SPARC (登録商標) 又はその他のいずれの汎用命令セットアーキテクチャを選択してもよい。プロセッサコアは、データ及びコンピュータプログラム命令を求めてキャッシュメモリサブシステムにアクセスすることができる。キャッシュサブシステムは、ランダムアクセスメモリ (RAM) 及びストレージデバイスを含む記憶階層に結合することができる。

【0027】

クライアントコンピュータシステム内の各プロセッサコア及び記憶階層は、ネットワークインターフェイスにさらに接続することができる。クライアントコンピュータシステム 110 a ~ 110 c の各々は、ハードウェア構成要素に加え、記憶階層内に記憶された基本オペレーティングシステム (OS) を含むことができる。この基本 OS は、例えば、MS-DOS (登録商標)、MS-WINDOWS (登録商標)、OS/2 (登録商標)、UNIX (登録商標)、Linux (登録商標)、Solaris (登録商標)、又は別の公知のオペレーティングシステムなどの様々な特定のオペレーティングシステムのいずれかを表すことができる。従って、基本 OS は、エンドユーザに様々なサービスを提供するとともに、様々なプログラムの実行をサポートするソフトウェアフレームワークを提供することができる。また、クライアントコンピュータシステム 110 a ~ 110 c の各々は、高水準バーチャルマシン (VM) をサポートするために使用されるハイパーバイザを含むことができる。当業者には周知のように、デスクトップ及びサーバ内で仮想化を使用して、OS などのソフトウェアをシステムのハードウェアから完全に又は部分的に分離することができる。仮想化により、データストレージアレイ 120 a ~ 120 b の各々におけるストレージデバイス 176 a ~ 176 m に対応する (論理装置番号 (LUN) などの) 論理記憶エンティティなどのように、各々が独自のリソースを有する同じ機械上で複数の OS が実行されているという錯覚をエンドユーザに与えることができる。

【0028】

データストレージアレイ 120 a ~ 120 b の各々は、クライアントコンピュータシス

10

20

30

40

50

テム 110 a ~ 110 c などの異なるサーバ間のデータの共有に使用することができる。データストレージレイ 120 a ~ 120 b の各々は、データを記憶するためのストレージサブシステム 170 を含む。ストレージサブシステム 170 は、複数のストレージデバイス 176 a ~ 176 m を含むことができる。これらのストレージデバイス 176 a ~ 176 m の各々は、SSD とすることができる。コントローラ 174 は、受け取った読み込み / 書き込み要求を処理するためのロジックを含むことができる。例えば、少なくともコントローラ 174 において、手短に上述したアルゴリズムを実行することができる。受け取った書き込み要求などの動作のバッチ処理には、ランダムアクセスメモリ (RAM) 172 を使用することができる。

【0029】

記憶媒体 130 に記憶された基本 OS 132、ファイルシステム 134、いずれかの OS ドライバ (図示せず) 及びその他のソフトウェアは、ファイル及び LUN へのアクセスを可能にする機能を提供し、これらの機能を管理することができる。基本 OS 134 及び OS ドライバは、記憶媒体 130 上に記憶された、受け取った要求に対応する 1 又はそれ以上のメモリアクセス動作をストレージサブシステム 170 内で行うようにプロセッサ 122 により実行可能なプログラム命令を含むことができる。

【0030】

データストレージレイ 120 a ~ 120 b の各々は、ネットワークインターフェイス 124 を使用してネットワーク 180 に接続することができる。1つの実施形態では、クライアントコンピュータシステム 110 a ~ 110 c と同様に、ネットワークインターフェイス 124 の機能をネットワークアダプタカード上に含めることができる。ネットワークインターフェイス 124 の機能は、ハードウェア及びソフトウェアの両方を使用して実装することができる。ネットワークインターフェイス 124 のネットワークカードによる実装上には、ランダムアクセスメモリ (RAM) 及び読み取り専用メモリ (ROM) の両方を含めることができる。1又はそれ以上の特定用途向け集積回路 (ASIC) を使用して、ネットワークインターフェイス 124 の機能を提供することができる。

【0031】

1つの実施形態では、ユーザデータ及び対応するエラー訂正符号 (ECC) 情報のデータレイアウトを最適化しようと努めるデータストレージモデルを作成することができる。1つの実施形態では、このモデルが、ストレージシステム内のストレージデバイスの特性に少なくとも部分的に基づく。例えば、ソリッドステートストレージ技術を利用するストレージシステムでは、特定のデバイスの特性を使用してこのストレージデバイスのためのモデルを作成するとともに、このモデルが、対応するデータストレージ構成アルゴリズムを通知する機能を果たすようにすることもできる。例えば、使用中の特定のストレージデバイスが、時間とともに信頼性の変化を示す場合、データストレージ構成を動的に変更する上でこのような特性を考慮することができる。

【0032】

一般的に言えば、コンピューティングシステムのために開発されるあらゆるモデルは不完全である。多くの場合、現実のシステムで所与のシステムを完全にモデル化するには、考慮すべき変数が単純に多すぎる。場合によっては、完全ではないが価値のあるモデルを開発することが可能な場合もある。以下でより詳細に説明するように、基礎を成すデバイスの特性に基づいてストレージシステムをモデル化する実施形態を説明する。様々な実施形態では、システムがどのように挙動し得るかに関するいくつかの予測に基づいてデータストレージ構成の選択が行われる。デバイスの挙動によっては、デバイスの特性の理解に基づいて、他のデバイスの挙動より予測しやすいものもある。しかしながら、デバイスの動作は時間とともに変化することがあり、これに応じて、選択したデータレイアウトも変化することがある。本明細書で使用するデバイスの特性とは、全体としてのデバイスの特性、チップ又はその他の構成要素などのデバイスの下位部分の特性、消去ブロックの特性、又はデバイスに関する他のあらゆる特性のことを意味することができる。

【0033】

デバイス内冗長性

ここで図2を参照すると、動的デバイス内冗長スキームの1つの実施形態を示す汎用ブロック図を示している。当業者には周知のように、ストレージデバイス内の潜在的なセクタエラーの影響を低減するように複数のデバイス内冗長スキームの1つを選択することができる。通常、「セクタ」という用語は、ディスク上の所与のトラック内のセグメントなどの、HDD上における基本記憶単位を意味する。ここでは、「セクタ」という用語は、SSD上における基本的な割り当て単位（アロケーションユニット）を意味することもできる。

【0034】

SSD内のアロケーションユニットは、SSD内の1又はそれ以上の消去ブロックを含むことができる。図2を参照すると、ユーザデータ210は、エンドユーザが修正及びアクセスすべき記憶データ、並びにデバイス間エラー訂正符号（ECC）データの両方を意味することができる。デバイス間ECCデータは、ユーザデータを保持する他のストレージデバイス上の1又はそれ以上のページから生成されたパリティ情報とすることができる。例えば、デバイス間ECCデータは、RAIDデータレイアウトアーキテクチャで使用されるパリティ情報とすることができる。ユーザデータ210は、ストレージデバイス176a~176kのうちの1又はそれ以上に含まれる1又はそれ以上のページ内に記憶することができる。1つの実施形態では、ストレージデバイス176a~176kの各々がSSDである。

【0035】

SSD内の消去ブロックは、複数のページを含むことができる。上述したように、1つの実施形態では、1ページが4KBのデータ記憶空間を含むことができる。1消去ブロックは、64ページ、すなわち256KBを含むことができる。他の実施形態では、消去ブロックが1メガバイト（MB）もの大きさであり、256ページを含むことができる。アロケーションユニットのサイズは、アロケーションユニットのオーバーヘッドトラッキングを低減するために、十分に大きなサイズの比較的少ないユニットが与えられるように選択することができる。1つの実施形態では、1又はそれ以上の状態テーブルが、アロケーションユニットの状態（割り当て済み、空き、消去済み、エラー）、摩耗レベル、及びアロケーションユニット内で発生した総エラー数（訂正可能及び／又は訂正不能）を維持することができる。様々な実施形態では、アロケーションユニットのサイズを、所与のデバイスが利用できるアロケーションユニットの数とアロケーションユニットを維持するオーバーヘッドとのバランスをとるように選択することができる。例えば、1つの実施形態では、アロケーションユニットのサイズを、SSDの総記憶容量の1/100パーセントとなるように選択することができる。ページ、消去ブロック及びその他の単位構成に関する他の量のデータ記憶空間も可能であり、企図される。

【0036】

ストレージデバイス内の所与のセクタ又はその他の記憶単位にアクセスできない場合、潜在的セクタエラー（LSE）が発生する。この所与のセクタでは、読み込み又は書き込み動作を完了できないことがある。また、訂正不能なエラー訂正符号（ECC）エラーが存在することもある。LSEは、所与のセクタがアクセスを受けるまで検出されないエラーである。従って、この所与のセクタに以前に記憶したあらゆるデータが失われる恐れがある。ストレージデバイスの不具合後のRAID再構成中に1つのLSEが生じた場合、これによりデータ損失に至る恐れがある。SSDでは、デバイスの使用年数、デバイスサイズ、アクセス率、ストレージの圧縮度、及び以前の訂正可能及び訂正不能なエラーの発生といった統計値のうちの少なくとも1つから、別のLSEの発生確率が増加することがある。所与のストレージデバイス内のLSE及びデータ損失を防ぐために、この所与のストレージデバイス内で多様なデバイス内冗長スキームの1つを使用することができる。

【0037】

デバイス内冗長スキームは、所与のストレージデバイス内で、パリティ情報などのECC情報を利用する。このデバイス内冗長スキーム及びそのECC情報は、所与のデバイス

10

20

30

40

50

に対応し、所与のデバイス内に維持することができるが、デバイス自体が内部的に生成して維持できるECCとは異なる。一般的に言えば、デバイスの内部的に生成され維持されるECCは、このデバイスを含むシステムからは見えない。所与のストレージデバイスに含まれるデバイス内ECC情報を使用して、この所与のストレージデバイス内におけるデータ記憶の信頼性を高めることができる。このデバイス内ECC情報は、RAIDデータレイアウトアーキテクチャで利用されるパリティ情報などの、別のストレージデバイスに含めることができる他のECC情報に加えられる。

【0038】

極めて効果的なデバイス内冗長スキームは、所与のRAIDデータレイアウトの信頼性を十分に高めて、パリティ情報を保持するために使用されるデバイスの数を減少させることができる。例えば、各デバイス上のデータを保護するためにデバイス内冗長性が加わった場合、ダブルパリティRAIDレイアウトをシングルパリティRAIDレイアウトに置き換えることができる。一定レベルの記憶効率を得るためには、デバイス内冗長スキームにおける冗長性を高めると、所与のストレージデバイスの信頼性が高まる。しかしながら、このようにして冗長性を高めると、この所与のストレージデバイスの入力/出力(I/O)性能に関する不利点も増える可能性がある。

【0039】

1つの実施形態では、デバイス内冗長スキームが、デバイスを、ユーザデータを記憶するための一群の場所に分割する。例えば、ストライプ250a~250cによって示すような、RAIDレイアウト内のストライプに対応するデバイス内一群の場所に分割することができる。ユーザデータ又はデバイス間RAID冗長性情報は、データ210によって示すように、ストレージデバイス176a~176kの各々の1又はそれ以上のページに記憶することができる。各ストレージデバイス内では、1又はそれ以上のページにデバイス内エラー回復データ220を記憶することができる。本明細書では、デバイス内エラー回復データ220をデバイス内冗長データ220と呼ぶことができる。当業者であれば周知のように、デバイス内冗長データ220は、データ210内の情報の選択した一部に対してある関数を実行することにより取得することができる。XORベースの演算を使用して、デバイス内冗長データ220に記憶すべきパリティ情報を導出することができる。デバイス内冗長スキームの他の例としては、シングルパリティチェック(SPC)、最大距離分離(MDS)消去符号、インタリーブパリティチェック符号(IPC)、ハイブリッドSPC及びMDS符号(MDS+SPC)、及びカラム対角パリティ(CDP)が挙げられる。これらのスキームは、データ220を計算する方法に応じて、もたらされる信頼性及びオーバーヘッドの点で異なる。このシステムは、上述の冗長性情報に加え、デバイス上の領域のチェックサム値を計算するように構成することもできる。例えば、デバイスに情報が書き込まれた時にチェックサムを計算することができる。このチェックサムは、システムによって記憶される。システムは、デバイスから情報が読み戻された時に再びチェックサムを計算し、これを最初に記憶した値と比較することができる。2つのチェックサムが異なる場合、情報が正しく読み込まれておらず、システムは、他のスキームを使用してデータを回復することができる。チェックサム機能の例には、巡回冗長検査(CRC)、MD5及びSHA-1がある。

【0040】

ストライプ250a~250cに示すように、所与のストライプ内にデータ210を記憶するために使用する幅、すなわちページ数は、ストレージデバイス176a~176kの各々において同じものとなり得る。しかしながら、ストライプ250b~250cに示すように、所与のストライプ内にデバイス内冗長データ220を記憶するために使用する幅、すなわちページ数は、ストレージデバイス176a~176kの各々において同じでない場合がある。1つの実施形態では、所与のストレージデバイスの特性又は挙動の変化により、対応するデバイス内冗長データ220を記憶するために使用する幅を少なくとも部分的に決定することができる。例えば、上述したように、フラッシュセルには、ページをプログラムすること又は読み込むことによって近隣のページに支障を来たし、これらの

10

20

30

40

50

近隣のページ内にエラーを引き起こし得るプログラムディスタープエラー及びリードディスタープエラーが生じる。ストレージデバイスが古くなってより多くのエラーが生じている場合、対応するデバイス内冗長データ220の量が増加することがある。例えば、ストライプ250bの書き込み動作前に、ストレージデバイス176a~176kの各々の特性をモニタし、これを使用してエラー率の増加を予測することができる。ストレージデバイス176c及び176jのエラーの増加が予測されることが検出される可能性がある。これにตอบสนองして、ストレージデバイス176c及び176jのデバイス内冗長データ220の量が増加することがある。図2のストライプ250a及び250bの例では、ストライプ250a及び250bのストレージデバイス176c及び176jの記憶されている保護データの量の増加を確認することができる。例えば、この時点で、ストレージデバイス176c及び176jをシングルパリティで保護するのではなく、これらのデバイスをダブルパリティ又はトリプルパリティで保護することができる。なお、デバイス176c及び176jのデバイス内保護の量を増加させても、同じストライプ内の他のデバイスにおける量に対応して増加させる必要はない。むしろ、ストライプのデータは、各デバイスにおいて望むように異なるレベルの保護を有することができる。

10

【0041】

様々な実施形態では、所与のレベルのデータ保護の増減を選択的に行うことができる。例えば、1つの実施形態では、上記の例におけるストレージデバイス176c及び176jなどの、より多くのエラーを生成することが検出されたストレージデバイスに対してしか保護の増加を行わなくてよい。別の実施形態では、ストレージデバイス176c及び176jがより多くのエラーを生成することが検出された場合、ストレージデバイス176a~176kの各々に対して保護の増加を行うことができる。1つの実施形態では、デバイス176kなどのパリティデバイス上におけるデバイス内保護の量を増加させるために、ストライプ内で保護されているデータの量を減少させることが必要となる場合がある。例えば、所与のストライプのパリティデバイス上に記憶されたデバイス内データの量を増加させると、このデバイスが記憶するストライプ内のデータのためのパリティデータの量が必然的に減少する。このパリティデータの量が、ストライプ内のデータを全て保護するために必要な量未満に減少した場合、パリティ保護を継続することが望ましい場合にはストライプ内のデータを減少させなければならない。ストライプ内に記憶されるデータの量の減少させる代替案として、パリティデータを記憶するための異なるデバイスを選択することもできる。様々な選択肢が可能であり、企図される。また、本明細書に記載する図2及びその他の図には、(176kなどの)別個のパリティデバイスを示していることがあるが、様々な実施形態では、パリティを単一のデバイスに記憶するのではなく複数のデバイスに分散させることができる。従って、別個のパリティデバイスの図示は、一般に説明を簡単にするための論理的な表現と見なすことができる。

20

30

【0042】

ここで図3を参照すると、データストレージサブシステムにおけるデバイス内保護を調整する方法300の1つの実施形態を示している。一般に、ネットワークアーキテクチャ100及びデータストレージレイ120a~120b内で具体化される構成要素は、方法300に従って動作することができる。この実施形態のステップを順番に示す。しかしながら、ステップによっては、図示のものとは異なる順序で行なうことができるもの、同時に行うことができるもの、他のステップと組み合わせることができるもの、及び別の実施形態には存在しないものもある。

40

【0043】

ブロック302において、ストレージデバイスにユーザデータを記憶するための第1の空き容量を決定する。このユーザデータは、データ210に関して上述したように、RAIDアーキテクチャで使用されるエンドユーザアプリケーション又はデバイス間パリティ情報内で使用されるデータとすることができる。この第1の空き容量は、上述したようなストレージデバイス内の1又はそれ以上のページを含むことができる。1つの実施形態では、ストレージコントローラ174内のグローバルRAIDエンジン178が、ストレージ

50

ジデバイス 176 a ~ 176 m の各々から挙動統計値を受け取る。グローバル R A I D エンジン 178 は、ストレージデバイス 176 a ~ 176 m の 2 又はそれ以上を含む所与のデバイスグループについて、R A I D データレイアウト及びこれらの 2 又はそれ以上のストレージデバイスの各々に維持すべき最初のデバイス内冗長性の量を決定することができる。ブロック 304 において、R A I D エンジン 178 は、対応するデバイス内保護データをストレージデバイスに記憶するための第 2 の空き容量を決定することができる。この第 2 の空き容量は、ストレージデバイス内の 1 又はそれ以上のページを含むことができる。デバイス内保護データは、上述したデバイス内冗長データ 220 に対応することができる。

【0044】

ブロック 306 において、所与のデバイスグループに含まれる各ストレージデバイス内の第 1 の空き容量にデータを書き込む。1 つの実施形態では、ユーザデータ及びデバイス間パリティ情報が、いずれも所与のデバイスグループに含まれる複数のストレージデバイスにわたる単一の R A I D ストライプとして書き込まれる。再び図 2 を参照して分かるように、対応する書き込まれるデータの幅は、各ストレージデバイスにおいて同じである。ブロック 308 において、E C C アルゴリズム、X O R ベースのアルゴリズム、又はその他のあらゆる適当なアルゴリズムにより、デバイス内保護データを生成する。また、システムは、正しく取り出されなかったデータを識別しやすくするためにチェックサムを生成することができる。ブロック 310 において、生成されたデバイス内保護データを、ストレージデバイス内の第 2 の空き容量に書き込む。

【0045】

ブロック 312 において、R A I D エンジン 178 は、1 又はそれ以上のストレージデバイスの挙動をモニタすることができる。1 つの実施形態では、R A I D エンジン 178 が、対応するストレージデバイスのモデルを含み、このモデルに入力すべき挙動統計値をストレージデバイスから受け取ることができる。このモデルは、ストレージデバイスの既知の特性を利用することにより、ストレージデバイスの挙動を予測することができる。例えば、このモデルは、所与のストレージデバイスのエラー率の増加が近いことを予測することができる。R A I D エンジン 178 は、信頼性に影響を与える所与のストレージデバイスの特性を検出した（条件付きブロック 314）場合、ブロック 316 において、データ及び対応するデバイス内冗長データを記憶するための第 1 及び第 2 の空き容量を調整することができる。例えば、R A I D エンジン 178 は、少なくともデバイスの使用年数、アクセス率及びエラー率などの上述した統計値をモニタすることができる。再び図 2 を参照して分かるように、R A I D エンジン 178 は、ストレージデバイス 176 c 及び 176 j のエラー数が増加したことを検出することができる。或いは、R A I D エンジン 178 は、ストレージデバイス 176 c 及び 176 j のエラー数の増加を予測することができる。従って、R A I D エンジン 178 は、第 2 のストライプ 250 b を書き込む前に、ストレージデバイス 176 a ~ 176 k の各々にデータ 210 及びデータ 220 を記憶するために使用するページ数を調整することができる。同様に、R A I D エンジン 178 は、ストレージデバイス 176 b の信頼性が低下していることを検出することができる。従って、R A I D エンジン 178 は、第 3 のストライプ 250 c を書き込む前に、ストレージデバイス 176 a ~ 176 k の各々にデータ 210 及びデータ 220 を記憶するために使用するページ数を再び調整することができる。

【0046】

ストレージデバイス特性のモニタリング

ここで図 4 を参照すると、ストレージサブシステムの 1 つの実施形態の汎用ブロック図を示している。1 又はそれ以上のデバイスグループ 173 a ~ 173 m の 1 つにおいて、1 又はそれ以上のストレージデバイス 176 a ~ 176 m の各々を分割することができる。同様に、他のデバイスを含む他のデバイスグループも存在し得る。各ストレージデバイスの対応する動作キュー及び状態テーブルを、デバイスユニット 400 a ~ 400 w の 1 つに含めることができる。これらのデバイスユニットを、R A M 172 に記憶することが

10

20

30

40

50

できる。デバイスグループ 173a ~ 173m の各々に関し、対応する RAID エンジン 178a ~ 178m を含めることができる。各 RAID エンジン 178 は、対応するデバイスグループ内のストレージデバイスの各々の統計値を追跡するモニタ 410 を含むことができる。データレイアウトロジック 420 は、対応するストレージデバイス内の、ユーザデータ、デバイス間冗長データ及びデバイス内冗長データに対して割り当てるべき空き容量を決定することができる。ストレージコントローラ 174 は、ウェアレベリング、ガベージコレクション、I/O スケジューリング、重複除外、並びに着信及び発信パケットのプロトコル変換といったタスクのうちの少なくとも 1 つを行うためのその他の制御ロジック 430 を含むことができる。

【0047】

ここで図 5 を参照すると、デバイスユニットの 1 つの実施形態の汎用ブロック図を示している。デバイスユニットは、デバイスキュー 510 及びテーブル 520 を含むことができる。デバイスキュー 510 は、読み込みキュー 512、書き込みキュー 514、及びその他の動作キュー 516 などの 1 又はそれ以上の他のキューを含むことができる。各キューは、1 又はそれ以上の対応する要求 530a ~ 530d を記憶するための複数のエントリを含むことができる。例えば、対応する SSD のデバイスユニットは、少なくとも読み込み要求、書き込み要求、トリム要求及び消去要求などを記憶するためのキューを含むことができる。テーブル 520 は、状態データ又は統計値 530 を記憶するための複数のエントリを各々が含む 1 又はそれ以上の状態テーブル 522a ~ 522b を含むことができる。また、この図及びその他の図には、キュー及びテーブルが特定数のエントリを含むように示しているが、必ずしもエントリ自体が互いに対応するわけではない。さらに、キュー、テーブル及びエントリ数は、図示のものとは異なることがあり、また互いに異なることもある。

【0048】

ここで図 6 を参照すると、所与のデバイスに対応する状態テーブルの 1 つの実施形態を示す汎用ブロック図を示している。1 つの実施形態では、このようなテーブルが、SSD などの所与のストレージデバイスの状態情報、エラー情報及び摩耗率情報に対応するデータを含むことができる。この情報に、対応する RAID エンジンがアクセスすることにより、この RAID エンジン、デバイス間保護及びデバイス内保護の両方に使用するデータ記憶及びスキームのために割り当てられた空間を動的に変更できるようになる。1 つの実施形態では、この情報が、デバイスの使用年数 602、エラー率 604、デバイス 606 上で検出された総エラー数、回復可能なエラー数 608、回復不能なエラー数 610、デバイスのアクセス率 612、記憶されたデータの使用年数 614、及び割り当て空間 616a ~ 616n の 1 又はそれ以上の割り当て状態、のうちの少なくとも 1 つ又はそれ以上を含むことができる。これらの割り当て状態は、使用中、空き、及びエラーなどを含むことができる。

【0049】

フレキシブルな RAID レイアウト

ここで図 7 を参照すると、フレキシブルな RAID データレイアウトアーキテクチャの 1 つの実施形態を示す汎用ブロック図を示している。RAID エンジン、ストレージデバイス 176a ~ 176k に使用すべき保護レベルを決定することができる。例えば、RAID エンジン、ストレージデバイス 176a ~ 176k に RAID ダブルパリティを利用すると決定することができる。デバイス間冗長データ 240 は、対応するユーザデータから生成された RAID ダブルパリティ値を表すことができる。1 つの実施形態では、ストレージデバイス 176j 及び 176k が、ダブルパリティ情報を記憶することができる。他のレベルの RAID パリティ保護も可能であり、企図されると理解されたい。また、他の実施形態では、ダブルパリティ情報を、各 RAID ストラップのストレージデバイス 176j 及び 176k に記憶するのではなく、この記憶をストレージデバイス間で循環させることもできる。ダブルパリティ情報がストレージデバイス 176j 及び 176k に記憶されるように示しているのは、図示及び説明を簡単にするためである。

【 0 0 5 0 】

ここで図 8 を参照すると、フレキシブルな R A I D データレイアウトアーキテクチャの別の 1 つの実施形態を示す汎用ブロック図を示している。図 7 に示す例と同様に、ストレージデバイス 1 7 6 a ~ 1 7 6 k にはダブルパリティを使用することができる。この例では、R A I D ダブルパリティについて説明するが、R A I D データレイアウトアーキテクチャではあらゆる量の冗長性を選択することができる。

【 0 0 5 1 】

動作中、R A I D エンジン 1 7 8 は、ストレージデバイス 1 7 6 a ~ 1 7 6 k の特性をモニタして、デバイスが最初の又はその他の所与の信頼性レベルよりも高い信頼性レベルを示していると判断することができる。これに応答して、R A I D エンジン 1 7 8 は、R A I D 保護を R A I D ダブルパリティから R A I D シングルパリティに変更することができる。他の R A I D データレイアウトアーキテクチャでは、サポートされている冗長性の量を別様に減少させることができる。他の実施形態では、ストレージデバイス 1 7 6 a ~ 1 7 6 k のモニタ及び保護レベルの変更を、ストレージコントローラ 1 7 4 内の他のロジックによって行うことができる。

【 0 0 5 2 】

引き続き上記の例を参照すると、所与の R A I D ストライプ上で実行される後続の書き込み動作に関しては、シングルパリティ情報のみを生成して記憶することができる。例えば、サポートされている冗長性の量を変更した後の書き込み動作に関しては、後続の R A I D ストライプ内でストレージデバイス 1 7 6 k を使用しなくてもよい。また、ストレージデバイス 1 7 6 k に記憶されたデータを無効にし、これによりストレージを解放することもできる。その後、ストレージデバイス 1 7 6 k の解放されたデータに対応するページを、他の用途に再割り当てすることができる。このパリティ保護の量を低減し、パリティ保護データを記憶するために以前に使用されていた空間を解放するプロセスのことを、「パリティシュレディング」と呼ぶことができる。さらに、ストレージデバイス 1 7 6 k が S S D である実施形態では、ストライプ 2 5 0 a 内のページを書き換える前に、ストレージデバイス 1 7 6 k 内で 1 又はそれ以上の消去動作を行うことができる。

【 0 0 5 3 】

上記のパリティシュレディングの例を引き続き参照すると、パリティシュレディング後にストライプ 2 5 0 a 内のストレージデバイス 1 7 6 k の再割り当てされたページに記憶されたデータは、ストライプ 2 5 0 a に対応しない他の R A I D ストライプのユーザデータ又は対応する R A I D シングルパリティ情報を保持することができる。例えば、ストライプ 2 5 0 a 内のストレージデバイス 1 7 6 a ~ 1 7 6 j に記憶されたデータは、パリティシュレディング前に実行される 1 又はそれ以上の書き込み動作に対応することができる。ストライプ 2 5 0 a 内のストレージデバイス 1 7 6 k に記憶されたデータは、パリティシュレディング後に実行される 1 又はそれ以上の書き込み動作に対応することができる。同様に、ストライプ 2 5 0 b 内のストレージデバイス 1 7 6 a ~ 1 7 6 j に記憶されたデータは、パリティシュレディング前に実行される 1 又はそれ以上の書き込み動作に対応することができる。ストライプ 2 5 0 b 内のストレージデバイス 1 7 6 k のページは、解放後に消去して、サポートされている冗長性の量を変更した後に実行される 1 又はそれ以上の書き込み動作に対応するデータに書き換えることができる。なお、冗長性情報が複数のストレージデバイスを循環する場合、このスキームはより一層効果的となり得る。このような実施形態では、シュレディングにより解放される空間も、同様にストレージデバイスにわたって分散される。

【 0 0 5 4 】

再び図 8 を参照すると、ストライプ 2 5 0 c 内のストレージデバイス 1 7 6 k に示す割り当て解除されたページは、パリティシュレディング前に R A I D ダブルパリティ情報を記憶していた可能性のある記憶位置を表す。しかしながら、現時点でこれらのページは無効であり、未だ再割り当てされていない。上記の例では、ストレージデバイス 1 7 6 k 内のページを解放して再割り当てする方法及びタイミングが、S S D の特定の特性によっ

10

20

30

40

50

て決まる。これらの特性の例として、1又はそれ以上のページを再プログラムする（書き換える）前に少なくとも消去ブロック全体を消去することが挙げられる。図8からわかるように、パリティをシュレッドする際に、デバイス全体をシュレッドする必要はない。むしろ、個々のストライプに関して、望むようにパリティをシュレッドすることができる。同様に、ストライプのパリティ保護を増加させることもでき、別のデバイス上に記憶された保護データをストライプに追加することもできる。

【0055】

ここで図9を参照すると、RAIDレイアウトを動的に決定する方法の1つの実施形態を示している。一般に、ネットワークアーキテクチャ100及びデータストレージレイ120a~120b内で具体化される構成要素は、方法900に従って動作することができる。図9には、2つのプロセス910及び920を示している。これらのプロセスの各々は、同時に又は所与の順序で動作することができる。さらに、この実施形態のステップを順番に示す。しかしながら、ステップによっては、図示のものとは異なる順序で行なうことができるもの、同時に行うことができるもの、他のステップと組み合わせることができるもの、及び別の実施形態には存在しないものもある。ブロック910は、ストレージ制御システムが、システム内のストレージデバイスの特性及び挙動をモニタするプロセスを示す（ブロック912）。例えば、図6で説明したような特性を観察及び/又は記録することができる。信頼性の変化などの特定の状態が検出された場合（判定ブロック914）、記憶したデータに使用する保護の量を変更することができる（ブロック916）。例えば、所与のデバイスの使用年数が比較的短い場合には、デバイスの信頼性が分からないことがある（例えば、デバイスが「初期故障」を起こして比較的早い時期に機能しなくなることもある）。従って、RAIDストライプ毎に1又はそれ以上の追加のストレージデバイスを使用して、パリティ情報を記憶することができる。その後の時点で、時間とともにこれらのデバイスが信頼できると判明した時に、この追加の保護を除去することができる。様々な実施形態では、デバイスのエラー率に関する特性を維持することができる。例えば、訂正可能な及び/又は訂正不能なエラーに関する特性を維持し、これらを使用して所与のデバイスの信頼性に関する判断を行うことができる。ストレージコントローラは、この情報に基づいて、デバイス又はストライプの様々な保護レベルを動的に変更することができる。

【0056】

図9のブロック920は、ストレージのストライプ又はその他の部分を割り当てる時点で（判定ブロック922）、データに使用すべきレイアウト及び保護レベルに関する決定を行なう（ブロック924）ことができるプロセスを大まかに示すものである。なお、この時にブロック910のプロセスを実施することもできる。或いは、既にプロセス910によって保護レベルが決定され記憶されている場合もある。この場合、ブロック924の決定は、この記憶されているデータに基づくことができる。1つの実施形態では、所与のレイアウトを決定すると、このレイアウトに使用すべき特定のデバイスをデバイスグループから選択することができる（ブロック925）。例えば、1つの実施形態では、20個のデバイスからなるグループを使用することができる。5+2のレイアウトを決定した場合、この20個のデバイスからなるグループから任意の7つのデバイスを使用するように選択することができる。また、選択した5+2のレイアウトによるその後の書き込みにおいて、この同じ7つのデバイスを使用する必要はない。ストライプのレイアウト、保護レベル及びデバイスを決定した後、このストライプに書き込みを行うことができる（ブロック926）。

【0057】

様々な実施形態では、RUSHアルゴリズムを利用して、所与のストライプのデータ及び冗長性情報がいずれのデバイス上に存在するようになるかを判断することができる。例えば、RUSHアルゴリズムを使用して、ストレージデバイス176a~176kの所与のストライプの8+2のRAIDレイアウトに利用すべき特定のデバイスを選択することができる。一般的に言えば、本明細書では、一般にM+Nのレイアウトは、所与のデータ

10

20

30

40

50

ストライプのM個のデータデバイス及びN個のパリティデバイスを含むレイアウトを表すことができる。また、上述したように、パリティを完全に特定のデバイス内に配置するのではなく、デバイスにわたって分散させることもできる。従って、8 + 2のレイアウトは、10個のデバイスにわたってストライピングされたデータ及びパリティを含むことができ、これらのデバイスのうちの8つがデータを記憶し、2つがパリティを記憶する。その後のある時点で、12 + 2のレイアウトを選択することができる。このように、書き込みを行う（例えば、ストライプを書き込む）時点で、所望のレイアウト及び保護特性を動的に決定することができる。1つの実施形態では、ストレージデバイス176a ~ 176kが、30個、50個、又はそれ以上のストレージデバイスなどの、10個よりも多くのストレージデバイスを含むことができる。しかしながら、8 + 2のレイアウトによるストライプでは、ストレージデバイスのうちの10個しか利用されない。なお、デバイスのうちの任意の10個を選択し、ストライプを記憶する際に使用するこれらの10個のデバイスを選択するために、任意の好適なアルゴリズムを使用することができる。例えば、CRUSHアルゴリズムを使用して、所与の8 + 2のRAIDレイアウトのために、ストレージデバイス176a ~ 176kのうちのいずれの10個を利用すべきかを選択することができる。

10

【0058】

ストレージデバイス176a ~ 176kのために選択される8 + 2のRAIDレイアウトの1つの例では、これらのストレージデバイスのうちの2つを使用して、パリティ情報などのエラー訂正符号（ECC）情報を記憶することができる。この情報を使用して、再構成読み込み要求を行うことができる。再び図8を参照すると、この例では、RAIDダブルパリティ情報を記憶するためにストレージデバイス176j及び176kを選択することができる。この場合も、パリティ情報を常に同じストレージデバイスに記憶するのではなく、RAIDアレイに含まれるストレージデバイス176a ~ 176kの各々の間で循環的に記憶することができる。図示及び説明を簡単にするために、ストレージデバイス176j及び176kがRAIDダブルパリティを記憶しているものとして説明する。

20

【0059】

ブロック926において、書き込み動作の実行中に、RAIDアレイに含まれる複数のストレージデバイスにわたり、メタデータ、ユーザデータ、デバイス内パリティ情報及びデバイス間パリティ情報をRAIDストライプとして書き込むことができる。ブロック912において、RAIDエンジン178は、RAIDアレイ内の1又はそれ以上のストレージデバイスの挙動をモニタすることができる。1つの実施形態では、RAIDエンジン178が、図4に示すようなモニタ410及びデータレイアウトロジック420を含むことができる。RAIDエンジン178は、少なくとも所与のストレージデバイスの使用年数、エラー数及びエラータイプ、最後のデータ割り当て以降に検出された構成変更、所与のデータの使用年数、及びRAIDアレイ内の記憶空間の現在の使用状況などをモニタすることができる。

30

【0060】

RAIDエンジン178によりモニタされたデータは、図4に示すデバイスユニット400a ~ 400wの1つなどのRAM172に記憶することができる。このデータを記憶するために、図5及び図6に示す例のようなテーブルを使用することができる。対応するRAIDエンジンに含まれるロジックは、ストレージデバイスの最新の統計値をモニタすることにより、ストレージデバイスの挙動を検出して予測することができる。例えば、このモデルは、所与のストレージデバイスのエラー率の増加が近いことを予測することができる。

40

【0061】

（単複の）ストレージデバイスの信頼性が増したことが検出された（条件付きブロック908）場合、ブロック910において、RAIDエンジンは、システム内のデータ保護レベルを低下させることができる。例えば、1つの実施形態では、ストレージサブシステムに記憶されているパリティ情報の量を減少させることができる。上記の例に関して、R

50

RAIDエンジンは、対応する $8 + 2$ のRAIDアレイのRAIDダブルパリティをRAIDシングルパリティに低下させて、 $8 + 1$ のRAIDアレイに変換することができる。他の例では、所与のRAIDアレイが、ブロック916の前に、RAIDアーキテクチャにおいてNレベルの量の冗長性又はパリティを利用していることがある。ブロック916において、RAIDエンジンは、 $N > 1$ かつ $1 < m < N$ とする $(N - m)$ レベルの量の冗長性を利用すると決定することができる。従って、その後の所与のRAIDストライプの書き込み動作中には、所与のRAIDストライプ内で、より少ないm個のストレージデバイスに書き込みが行われるようになる。

【0062】

RAIDエンジン（又は別の構成要素）は、システム内のデータ保護レベルを低下させるために、上述したようなパリティシュレディングを行うことができる。その後、ストレージコントローラ174は、シュレディング動作の結果として解放されたページを、その後の書き込み動作で使用されるように再割り当てすることができる。

【0063】

上述したように、ストレージデバイス176a～176kの各々が古くなってデータで満たされると、RAIDアレイから追加のパリティ情報を除去することができる。メタデータ、ユーザデータ、対応するデバイス内冗長性情報、及びデバイス間冗長性情報の一部は残存する。 $8 + 2$ のRAIDアレイを使用した上記の例に関しては、ストレージデバイス176a～176jに記憶された情報が残存する。一方、追加のデバイス間冗長性情報又は追加のパリティ情報は、RAIDアレイから除去することができる。例えば、ストレージデバイス176kに記憶された追加のパリティ情報を、RAIDストライプから除去することができる。

【0064】

上記の例でストレージデバイス176a～176jに記憶された情報などの残存した情報は、適所に残存することができる。上記の例におけるストレージデバイス176k内の対応するページなどの、追加のパリティ情報を記憶する記憶空間は、その後の書き込み動作のために再利用して再割り当てすることができる。1つの実施形態では、個々の新たな割り当てにより、新たな仮想アドレスが受け取られる。個々の新たな割り当ては、所与のサイズ、所与の配置又は構成を有することができ、所与の（仮想又は物理）記憶空間に適應することができる。1つの実施形態では、ストレージデバイス176a～176kの各々及びストレージデバイス内の各割り当てページが、識別情報を含むヘッダを有する。このような識別情報により、所与の構成を変更することなく、解放された追加のパリティ情報のために記憶空間を再利用できるようになる。

【0065】

ストレージデバイス176a～176kの1又はそれ以上がSSDである実施形態では、消去ブロック内の1又はそれ以上のページを再プログラムする前に消去ブロックが消去される。従って、ストレージデバイス176kがSSDである実施形態では、ストレージデバイス176k内の解放されたページを再プログラムする前に、対応する消去ブロックが消去される。元々の $8 + 2$ のRAIDアレイを使用する上記の例に関しては、ページにデータ210を再プログラムする前に、ストライプ250a～250b内のストレージデバイス176kの1又はそれ以上の消去ブロックが消去される。この時、元々の $8 + 2$ のRAIDアレイは $8 + 1$ のRAIDアレイになっており、ストレージデバイス176jは、パリティシュレディング前に書き込まれたRAIDストライプにシングルパリティ情報を提供する。

【0066】

当業者には周知のように、所与のストレージデバイスの読み込み又は書き込みエラー中には、対応するRAIDストライプ内のサポートされているデバイス間パリティ情報からデータを再構成することができる。この再構成されたデータをストレージデバイスに書き込むことができる。しかしながら、この再構成されたデータがストレージデバイスに対して書き込みエラーとなった場合、このストレージデバイス上に記憶されている全てのデー

10

20

30

40

50

タを、対応するパリティ情報から再生することができる。この再生されたデータは、別の場所に再配置ことができる。フラッシュメモリでは、フラッシュ変換層（FTL）が、データの記憶位置をリマップする。また、フラッシュメモリでは、データの再配置が、消去ブロック内の対応するページを再プログラムする前に消去ブロック全体を消去することを含む。マッピングテーブルを消去ブロック対ページの粒度に維持することにより、リマップテーブルをよりコンパクトにすることができる。さらに、再配置中には、パリティシュレディング中に解放された追加のページを使用することができる。

【 0 0 6 7 】

オフセットパリティ

ここで図 10 を参照すると、フレキシブルな RAID データレイアウトアーキテクチャのさらに別の実施形態を示す汎用ブロック図を示している。図 8 に示す汎用ブロック図と同様に、フレキシブルな RAID データレイアウトアーキテクチャを使用することができる。ストレージデバイス 176 a ~ 176 k は、複数のストレージデバイスにわたってレイアウトされた複数の RAID ストライプを含む。ストレージデバイス 176 a ~ 176 k の各々は複数のページを含むが、図を簡単にするために、ページ 1010 及びページ 1020 にしかラベルを付けていない。図示の例では、ストレージデバイス 176 j 及び 176 k がダブルパリティ情報を記憶するダブルパリティ RAID データレイアウトが選択されている。

【 0 0 6 8 】

ストレージデバイス 176 a ~ 176 k のページの各々は、特定のタイプのデータを記憶する。一部のページは、ユーザデータ 210 及び対応する生成されたデバイス間パリティ情報 240 を記憶する。他のページは、対応する生成されたデバイス内パリティ情報 220 を記憶する。さらに他のページは、メタデータ 242 を記憶する。メタデータ 242 は、ページヘッダ情報、RAID ストライプ識別情報、及び 1 又はそれ以上の RAID ストライプのログデータなどを含むことができる。ストレージデバイス 176 a ~ 176 k のページの各々は、デバイス間パリティ保護及びデバイス内パリティ保護に加え、各所与のページに記憶されたチェックサムなどの追加の保護を含むこともできる。様々な実施形態では、各ストライプの最初の単一のメタデータページを、他のストライプヘッダから再生することができる。或いは、データをデバイス間パリティによって保護できるように、このページがパリティシャード内の異なるオフセットに存在することもできる。「シャード」は、デバイスの一部を表す。従って、パリティシャードは、パリティデータを記憶するデバイスの一部を意味する。

【 0 0 6 9 】

物理層

様々な実施形態では、本明細書で説明するシステムが物理層を含み、これを通じてシステムの他の要素がストレージデバイスと通信することができる。例えば、スケジューリングロジック、RAID ロジック及びその他のロジックは、ソフトウェア及び/又はハードウェアのあらゆる好適な組み合わせを含む物理層を介してストレージデバイスと通信することができる。一般に、物理層は、永続ストレージへのアクセスの提供及びデータストレージの完全性に関する機能の実行を含む様々な機能を実行する。

【 0 0 7 0 】

図 11 A に、500 GB のデバイスの仮想デバイスレイアウトの 1 つの実施形態を示す。様々な実施形態では、本明細書で説明するストレージデバイスの最初にパーティションテーブル 1101 をフォーマットし、デバイスの最後にパーティションテーブルのコピーをフォーマットすることができる。また、最初と最後のブロックにデバイスヘッダ 1103 を記憶することができる。例えば、フラッシュベースのストレージデバイスでは、最初と最後の消去ブロックにデバイスヘッダを記憶することができる。上述したように、消去ブロックは、典型的には 256 KB ~ 1 MB のフラッシュ構造である。第 1 の消去ブロック内には、追加の未使用空間を確保することができる（パディング 1105）。各デバイスには、ログ及び診断情報 1107 を書き込むための第 2 の消去ブロックを確保すること

10

20

30

40

50

ができる。これらの間の残りの消去ブロックは、複数の消去ブロックのアロケーションユニット（AU）1109に分割される。AUサイズは、デバイス毎に十分な数のAUが存在して良好な割り当て粒度が得られるように選択することができる。1つの実施形態では、オーバーヘッドを避けるのに十分な多さではあるが、追跡しやすいように多すぎない単位での割り当てを可能にするように、デバイス上に10,000個ほどの範囲内のAUが存在する。AU（割り当て済み/空き/消去済み/不良）の状態の追跡は、AU状態テーブルに維持することができる。AUの摩耗率は、摩耗レベルテーブルに維持することができる。総エラー数は、AUエラーテーブルに維持することができる。

【0071】

様々な実施形態では、物理層が、（異なるノード上に存在し得る）デバイスの組にわたる各デバイス内の1つのセグメントシャードを含むセグメント内で空間を割り当てる。図11Bに、セグメント及びそのセグメントの様々な識別可能部分の1つの実施形態を、1つの考えられるセグメントレイアウトで示す。図示の実施形態では、複数のデバイスに記憶された単一のセグメントを示している。データデバイス：データ0～データN、並びにパリティデバイス：パリティP及びパリティQを示す。1つの実施形態では、各セグメントシャードが、各デバイス上でシャードのサイズが等しくなるように、デバイスに1又はそれ以上のアロケーションユニットを含む。セグメントシャードを示すために、セグメントシャード1123を挙げている。図11Bには、I/O読み込みサイズ1127も示しており、これは1つの実施形態ではページに対応する。I/Oシャードのページパリティの1又はそれ以上のページを含むことができるI/Oパリティチャンク1129も示している。

【0072】

1つの実施形態では、各セグメントが、以下のパラメータのうちの1又はそれ以上を含むことができる独自の構成を有する。

（1）RAIDレベル - セグメント内のデバイス間保護に使用されるRAIDレベル。これにより、ミラーリング、パリティ、又はECC RAID、及びどれだけのセグメントシャードがパリティを含むかを決定することができる。

（2）デバイスレイアウトI/Oシャードサイズ - 書き込み中に各デバイスにわたるストライピングに使用するサイズを表す。このサイズは、典型的には256KB～1MBとなり、恐らくは各デバイス上の消去ブロックサイズの倍数となる。図11Bには、例示目的でI/Oシャードサイズ1125を挙げている。

（3）I/O読み込みサイズ - 論理的読み込みサイズである。各I/Oシャードを一連の論理ページとしてフォーマットすることができる。さらに各ページは、ページ内のデータのヘッダ及びチェックサムを含むことができる。読み込みが発行された場合、この読み込みは、1又はそれ以上の論理ページに対するものとなり、チェックサムを使用して各ページ内のデータを検証することができる。

（4）I/OシャードRAIDレベル - I/Oシャードは、再生中に見つかった潜在的なエラーを処理すべきシャード内パリティを有する。このパラメータは、どのタイプのパリティがシャード内保護に使用されているか、従ってどれだけのシャード内パリティのコピーが維持されるかを決定する。

（5）I/Oパリティチャンク - 様々な実施形態では、ストレージデバイスが、ページ単位でECCを行うことができる。従って、エラーが見つかった場合、物理ページ全体の不具合を示している可能性がある。I/Oパリティチャンクは、セグメント内の各デバイス上の物理ページサイズの最小公倍数であり、I/Oパリティチャンク又はI/O読み込みサイズのうちの大きい方にI/Oシャードをストライピングすることにより、シャード内パリティが計算される。ページパリティの1又はそれ以上のページを含めることができる。様々な実施形態では、チェックサム検証が失敗に終わった場合、このパリティを使用してデータを再生することができる。

【0073】

様々な実施形態では、個々の新たなセグメントに書き込みが行われる際に、そのセグメ

10

20

30

40

50

ントの R A I D 構成が選択される。R A I D 構成の選択は、現在のアクティブなノード及びデバイスの組、並びにセグメント内のデータのタイプなどの因子に基づることができる。例えば、10 個のノード又はデバイスが利用可能な場合、(8 + 2) の R A I D 6 の構成を選択し、2 つのデバイス又はノード故障に耐えるように、これらのノードにわたってセグメントをストライピングすることができる。その後、ノードが故障した場合、次のセグメントを (7 + 2) の R A I D 6 の構成に切り替えることができる。セグメント内では、セグメントシャードの一部がデータを含み、一部が (パリティなどの) E C C を含む。

【 0 0 7 4 】

1 つの実施形態では、5 種類のセグメントが存在する。これらのセグメントのうちの 3 つは、A U 状態テーブル、A U エラーテーブル及び摩耗率テーブルに対応する。いくつかの実施形態では、これらの 3 つのセグメントをミラーリングしてさらに保護することができる。これらの 3 つのセグメントに加え、ミラーリングを通じてさらに保護することもできるメタデータセグメントが存在する。最後に、クライアントブロック及びログ情報を保持するデータセグメントが存在する。ログ情報は、セグメント内のクライアントブロックに関連する最新情報を含む。データセグメントは、パリティ P 及びパリティ Q シャードを使用して、図 1 1 B に示すような R A I D 6 によって保護される可能性が高い。上記に加え、起動時に全てのセグメントシャードのヘッダからの情報を投入されるメモリ内データ構造としてセグメントテーブルが維持される。いくつかの実施形態では、このテーブルを全てのノード上で完全にキャッシュして、あらゆるノードがストレージアクセスを物理アドレスに変換できるようにすることができる。しかしながら、他の実施形態では、論理基準を取ってデータが記憶されているセグメントレイアウトノードを識別できるセグメントテーブルを各ノードが有することができるオブジェクトストレージモデルを使用することができる。その後、ノード上の正確な記憶位置を識別するための要求がノードに伝えられる。図 1 1 B には、セグメント内のかなりの量の空間を占めるあらゆる (ボリューム、スナップショットの) 組み合わせを識別するセグメントテールデータも示す。スナップショットを除去する場合、データスクラバが、このデータに基づいて、ガベージコレクションを行うセグメントの識別を支援することができる。

【 0 0 7 5 】

1 つの実施形態では、基本書き込み単位が、セグメント内のデバイスの各々における 1 つの I / O シャードである `seg io` である。`seg io` 内の各論理ページには、ページのチェックサム (「メディア」チェックサムと呼ぶことができる) を含むページヘッダが、データの実際のページサイズが 1 ページをわずかに下回るようにフォーマットされる。セグメントのパリティシャード内のページについては、データページ内のページチェックサムがパリティページによって保護されるように、ページヘッダが小さくなっている。各 I / O シャードの最終ページは、この場合も小さなヘッダを有して、消去ブロック内の全てのチェックサム及びページデータをページ不具合から保護するパリティページである。ここで言うページサイズとは、1 又はそれ以上の物理フラッシュページとすることができる I / O 読み込みサイズのことである。セグメントによっては、物理ページよりも小さな読み込みサイズを使用することができる。この方法は、検索情報に対する読み込みをインデックスで駆動することができ、より小さなデータ部分を読み込みながら所望のデータを取得できるメタデータに対して行うことができる。このような場合、物理ページの半分を読み込むということは、より少ないデータを I / O パス (及びネットワーク) に結び付けて、より少ないデータを検証する (例えば、チェックサムを行う) ことを意味する。ある実施形態は、物理ページよりも小さい読み込みサイズをサポートするために、消去ブロックの最後に複数のパリティページを含んで、全てのパリティページの全体的サイズがフラッシュページサイズに等しくなるようにすることができる。

【 0 0 7 6 】

消去ブロックの摩耗率が増すにつれ、エラーの可能性は高くなる。摩耗率を追跡することに加え、高いエラー確率が識別される 1 又は複数の消去ブロック上に、エラーが観察される頻度に関するデータを維持することもできる。消去ブロックによっては、単一の R A

10

20

30

40

50

ＩＤ５パリティの代わりに、消去ブロックの最後にダブル又トリプルエラー訂正パリティを保持することを決定できるものもある。この場合、これに応じて `segio` のデータペイロードを低減することができる。全ての消去ブロックではなく、`segio` 内の不良な消去ブロックのみを低減することが必要となり得る。消去ブロック内のページヘッダを使用して、いずれのページがパリティであり、いずれのページがデータであるかを識別することができる。

【００７７】

ストレージからページが読み込まれる時には常に、ページチェックサムを使用してコンテンツを検証することができる。この検証がうまくいかなかった場合、消去ブロックパリティを使用してデータの再生を試みることができる。これがうまくいかなかった場合、セグメントのデバイス間 ECC を使用してデータを再構成することができる。

10

【００７８】

データセグメントでは、ペイロード領域を２つの領域に分割することができる。記憶されたクライアントブロックに関する最新情報を含むことができるログデータとしてフォーマットされたページが存在するようになる。ペイロード領域の残りは、クライアントブロックとしてフォーマットされたページを含むことができる。このクライアントブロックデータは、圧縮された形で記憶することができる。数多くの圧縮アルゴリズムが可能であり、企図される。また、様々な実施形態では、Intel（登録商標）高度暗号化標準命令を使用してチェックサムを生成することができる。また、データと同じページ内に存在して、データの圧縮に使用するアルゴリズムの識別などの、クライアントブロックを読み込むために必要な情報を含むクライアントブロックのヘッダも存在することができる。ガベージコレクションは、`segio` 内のクライアントブロックヘッダ及びログエントリの両方を利用することができる。また、クライアントブロックは、重複除外、及び解凍データの正しさをチェックするために使用される未圧縮データのチェックサムとすることができるデータハッシュを有することができる。

20

【００７９】

いくつかの実施形態では、セグメント及び `segio` が、これらを順序付けるために使用される単調増加する ID 番号を有することができる。`segio` への書き込みの一部として、論理層が、以前のフラッシュに対する依存を記録することができる。起動時には、物理層が、セグメント及び `segios` の順序付きリストを構築することができ、`segio` が別の未完成の `segio` に依存している場合、これをロールバックして書き込まれていないと見なすことができる。

30

【００８０】

摩耗率テーブル

各デバイスのローカルなセグメントには、各デバイスの摩耗率テーブル（WLT）を記憶することができる。この情報を、各セグメントシャードのヘッダに記憶することもできる。１つの実施形態では、摩耗情報が、アロケーションユニットが消去及び再利用された回数を表す整数である。摩耗情報は正確でない場合もあるので、一定量の行動が行われた場合、又は相当期間にわたってシステムがアイドルであった場合に、デバイスに対するテーブルのフラッシュを行なうことができる。WLT は、新たな WLT セグメントを割り当てる際に古い WLT セグメントをクリーニングすることに関与することもできる。追加の保護層を加えるために、古いコピーを解放する前にこれらを維持することができる。例えば、テーブルマネージャが、WLT エントリの以前の消去ブロック及び現在の消去ブロックを常に保持し、新たなセグメントを割り当てる場合には、この新たなセグメントの第２の消去ブロックに書き込みを行うまで古いセグメントを解放しないように保証する。

40

【００８１】

AU 状態テーブル

AU 状態テーブル（AST）は、各 AU の状態を追跡する。この状態には、空き、割り当て済み、消去済み、及び不良がある。この AST を、デバイス上のセグメントに記憶することができる。割り当て済み又は空きへの状態変更は同期更新とすることができ、不良

50

又は消去済みへの状態変更は非同期更新とすることができる。一般に、このテーブルは十分に小さく、N V R A Mに記録できるだけの十分な更新を有することができる。A S Tは、新たなセグメントを割り当てる際に古いW L Tセグメントをクリーニングすることに関与することができる。A S Tは、ドライブ上の各A Uの第1のブロックをスキャンすることによって完全に回復できるので、古いA S Tのコピーを保持する必要はない。

【 0 0 8 2 】

A Uエラーテーブル

A Uエラーテーブル (A E T) を使用して、各A U内の回復可能エラー及び回復不能エラーの数を追跡することができる。A E Tは、デバイス上のセグメントに記憶され、各フィールドは、2 バイトの整数とすることができる。このテーブル全体は、A U当たり4 バイトの比較的小さなものとすることができる。

10

【 0 0 8 3 】

ここで図 1 1 C を参照すると、異なるページタイプ内のデータストレージ構成の1つの実施形態を示す汎用ブロック図を示している。図示の実施形態では3つのページタイプを示しているが、他のタイプも可能であり、企図される。図示のページタイプは、メタデータ 1 1 5 0 を含むページ 1 1 1 0、ユーザデータ 1 1 6 0 を含むページ 1 1 2 0、及びパリティ情報 1 1 7 0 (デバイス間又はデバイス内) を含むページ 1 1 3 0 を含む。ページ 1 1 1 0 ~ 1 1 3 0 の各々は、ヘッダ及び識別情報を含むことができるメタデータ 1 1 4 0 を含む。また、1 1 1 0 ~ 1 1 3 0 ページの各々は、対応するチェックサム又はその他のエラー検出及び / 又は訂正符号などのページ内エラー回復データ 1 1 4 2 を含むことができる。このチェックサム値により、所与のデバイスグループ内のストレージデバイス 1 7 6 a ~ 1 7 6 k に記憶されたデータをさらに保護することができる。

20

【 0 0 8 4 】

さらに、ページ 1 1 3 0 は、ページ間エラー回復データ 1 1 4 4 を含むことができる。このデータ 1 1 4 4 は、他のストレージデバイスに記憶されたページ内データ 1 1 4 2 から導出された E C C 情報とすることができる。例えば、再び図 1 0 を参照すると、デバイス間パリティ情報 2 4 0 を記憶するストレージデバイス 1 7 6 j 内の各ページが、ページ間エラー回復データ 1 1 4 4 を記憶することもできる。データ 1 1 4 4 は、パリティ、チェックサム、又はストレージデバイス 1 7 6 a ~ 1 7 6 i の 1 又はそれ以上に記憶されたページ内エラー回復データ 1 1 4 2 から生成されたその他の値とすることもできる。1つの実施形態では、データ 1 1 4 4 が、他のストレージデバイスに記憶された 1 又はそれ以上の他のチェックサム値 1 1 4 2 から生成されたチェックサム値である。ストレージデバイス 1 7 6 j の所与のページ内のデータ 1 1 4 4 を、ストレージデバイス 1 7 6 a ~ 1 7 6 i の 1 又はそれ以上の対応するページ内のデータ 1 1 4 2 と位置合わせするために、対応するページにパディング 1 1 4 6 を追加することができる。

30

【 0 0 8 5 】

1つの実施形態では、エンドユーザアプリケーションが、H D Dでは5 1 2 バイトのセクタの境界上で I / O 動作を実行する。追加の保護を加えるために、8 バイトのチェックサムを加えて5 2 0 バイトのセクタを形成することができる。様々な実施形態では、フラッシュメモリベースのシステム内で圧縮及びリマッピングを行って、セクタ境界ではなくバイト境界上にユーザデータを配置可能にすることができる。また、ヘッダの後であってユーザデータの前のページ内にチェックサム (8 バイト又は4 バイトなど) を配置して、これを圧縮することができる。ページ 1 1 1 0 ~ 1 1 3 0 の各々には、この配置を示している。

40

【 0 0 8 6 】

エンドユーザアプリケーションが5 1 2 バイトのセクタを読み込む場合、1つの実施形態では2 K B ~ 8 K B のサイズの対応するページが、ページの最初に8 バイトのチェックサムによる追加の保護を有する。様々な実施形態では、2 セクタサイズの非電力のためにこのページをフォーマットしなくてもよい。ページ 1 1 1 0 ~ 1 1 2 0 に示すように、チェックサムをページ内に2、3 バイトオフセットさせることができる。このオフセットに

50

より、ページ 1 1 3 0 などのパリティページが、パリティページをカバーするチェックサム、及び他のページのチェックサムを保護するための ECC を記憶できるようになる。

【 0 0 8 7 】

さらに別の保護レベルでは、チェックサム値を計算する際にデータ位置情報を含めることができる。ページ 1 1 1 0 ~ 1 1 3 0 の各々におけるデータ 1 1 4 2 は、この情報を含むことができる。この情報は、論理アドレス及び物理アドレスの両方を含むことができる。この情報には、セクタ番号、データチャンク及びオフセット番号、トラック番号、平面番号などを含めることもできる。

【 0 0 8 8 】

代替の構成

ここで図 1 2 を参照すると、ハイブリッド RAID データレイアウト 1 2 0 0 の 1 つの実施形態を示す汎用ブロック図を示している。3 つの区分を示しているが、あらゆる数の区分を選択することができる。各区分は、図 1 に示すデバイスグループ 7 1 3 a ~ 1 7 3 b などの別個のデバイスグループに対応することができる。各区分は、複数のストレージデバイスを含む。1 つの実施形態では、CRUSH アルゴリズムなどのアルゴリズムを利用して、データストレージに使用すべき RAID データレイアウトアーキテクチャにおいていずれのデバイスを使用すべきかを選択することができる。

【 0 0 8 9 】

図示の例には、L + 1 の RAID アレイ、M + 1 の RAID アレイ及び N + 1 の RAID アレイを示している。様々な実施形態では、L、M 及び N が全て異なってもよく、同じであってもよく、又はこれらの組み合わせであってもよい。例えば、区分 1 に RAID アレイ 1 2 1 0 を示している。他のストレージデバイス 1 2 1 2 は、区分 1 内の他の RAID アレイの候補である。同様に、RAID アレイ 1 2 2 0 は、区分 2 内の所与の RAID アレイを示す。他のストレージデバイス 1 2 2 2 は、区分 2 内の他の RAID アレイの候補である。RAID アレイ 1 2 3 0 は、区分 3 内の所与の RAID アレイを示す。他のストレージデバイス 1 2 3 2 は、区分 3 内の他の RAID アレイの候補である。

【 0 0 9 0 】

RAID アレイ 1 2 1 0、1 2 2 0 及び 1 2 3 0 の各々では、ストレージデバイス P 1 が、それぞれの RAID アレイ内の RAID シングルパリティ保護を行う。ストレージデバイス D 1 ~ D N は、それぞれの RAID アレイ内のユーザデータを記憶する。この場合も、ストレージデバイス D 1 ~ D N と P 1 の間で、ユーザデータ及び RAID シングルパリティ情報の記憶を循環させることができる。しかしながら、デバイス D 1 ~ D N にユーザデータが記憶されたものとして説明する。同様に、図示及び説明を容易にするために、デバイス P 1 に RAID シングルパリティ情報が記憶されたものとして説明する。

【 0 0 9 1 】

1 又はそれ以上の所与の RAID アレイに対してさらなる量の冗長性をサポートするために、3 つの区分の各々のうちの 1 又はそれ以上のストレージデバイスを選択することができる。例えば、区分 3 内のストレージデバイス Q 1 を、RAID アレイ 1 2 1 0、1 2 2 0 及び 1 2 3 0 の各々と組み合わせることができる。ストレージデバイス Q 1 は、RAID アレイ 1 2 1 0、1 2 2 0 及び 1 2 3 0 の各々に RAID ダブルパリティ情報を提供することができる。この追加のパリティ情報は、アレイ 1 2 1 0、1 2 2 0 又は 1 2 3 0 の 1 つにストライプが書き込まれた時に生成され記憶される。さらに、この追加のパリティ情報は、アレイ 1 2 1 0、1 2 2 0 及び 1 2 3 0 の各々のストライプをカバーすることができる。従って、総ストレージデバイス数に対する RAID パリティ情報を記憶するストレージデバイスの数の比率が小さくなる。例えば、区分の各々が N + 2 の RAID アレイを使用する場合、総ストレージデバイス数に対する RAID パリティ情報を記憶するストレージデバイスの数の比率は、 $3(2) / (3(N + 2))$ 、すなわち $2 / (N + 2)$ である。対照的に、ハイブリッド RAID レイアウト 1 2 0 0 の比率は、 $(3 + 1) / (3(N + 1))$ 、すなわち $4 / (3(N + 1))$ である。

【 0 0 9 2 】

10

20

30

40

50

ユーザデータの記憶に使用するストレージデバイスの数を増やすことにより、上記の比率を低減することが可能である。例えば、ストレージデバイスQ1を利用するのではなく、区分の各々が $3N+2$ のRAIDアレイを利用するようにすることができる。このような場合、総ストレージデバイス数に対するRAIDパリティ情報を記憶するストレージデバイスの数の比率は、 $2/(3N+2)$ である。しかしながら、再構成読み込み動作中には、単一のデバイス故障のための再構成読み込み要求を $(3N+1)$ 個のストレージデバイスが受け取る。対照的に、ハイブリッドRAIDレイアウト1200では、単一のデバイス故障のための再構成読み込み要求をN個のストレージデバイスしか受け取らない。

【0093】

なお、3つの区分の各々は、異なるRAIDデータレイアウトアーキテクチャを利用することができる。所与のRAIDデータレイアウトアーキテクチャの選択は、総ストレージデバイス数に対するRAIDパリティ情報を記憶するストレージデバイスの数の所与の比率に基づくことができる。また、この選択は、再構成中に再構成読み込み要求を受け取ることができる所与のストレージデバイス数に基づくこともできる。例えば、RAIDアレイ1210、1220及び1230は、 $L+a$ 、 $M+b$ 及び $N+c$ などの構成をそれぞれ含むことができる。

【0094】

区分内のRAIDアレイの1又はそれ以上に対してさらなる量の冗長性をサポートするために、上記の又はその他の条件に基づいてストレージデバイスQ1などの1又はそれ以上のストレージデバイスを選択することができる。上記のRAIDアレイ及びこれらのRAIDアレイの各々に対して追加の保護を行うストレージデバイスの数Qを含む3つの区分の例では、総ストレージデバイス数に対するRAIDパリティ情報を記憶するストレージデバイスの数の比率は、 $(a+b+c+Q)/(L+a+M+b+N+c+Q)$ である。単一のデバイス故障の場合、上記の例における区分1～3では、再構成読み込み要求を受け取るストレージデバイスの数は、それぞれL、M及びNである。なお、一般に上記の説明は、図12の3つの異なる区分について行ったものである。このような実施形態では、所与のレイアウトが特定のデバイスグループに制限されるこの種の「ハード」区分が、1つの区分内の再構成読み込みが別の区分の再構成読み込みと衝突しないことを保証することができる。しかしながら、他の実施形態では、区分が、上述したようなハードではないこともある。むしろ、デバイスのプールを考えた場合、これらのデバイスのいずれからもレイアウトを選択することができる。例えば、デバイスを大きなプールとして扱う場合、 $(L+1, M+1, N+1)+1$ などのレイアウトを構成することが可能である。従って、構成が重複し、再構成読み込みが衝突しかねない可能性がある。プールのサイズに対してL、M及びNが小さい場合、通常の読み込みに対する再構成読み込みの割合を低く保つことができる。

【0095】

ここで図13を参照すると、データストレージサブシステムにおいて代替のRAID構成を選択する方法1300の1つの実施形態を示している。一般に、ネットワークアーキテクチャ100及びデータストレージアレイ120a～120b内で具体化される構成要素は、方法1300に従って動作することができる。この実施形態のステップを順番に示す。しかしながら、ステップによっては、図示のものとは異なる順序で行なうことができるもの、同時に行うことができるもの、他のステップと組み合わせることができるもの、及び別の実施形態には存在しないものもある。

【0096】

ブロック1302において、ストレージコントローラ174内のRAIDエンジン178又はその他のロジックが、所与の数のデバイスを使用して、ストレージサブシステムの各区分内のRAIDアレイにユーザデータを記憶すると決定する。その後、RUSH又はその他のアルゴリズムを使用して、いずれのデバイスを使用すべきかを選択することができる。1つの実施形態では、各区分が、同じ数のストレージデバイスを利用する。他の実施形態では、各区分が、異なる固有の数のストレージデバイスを利用してユーザデータを

10

20

30

40

50

記憶することができる。ブロック 1304 において、ストレージコントローラ 174 が、いくつかのストレージデバイスがサブシステムの各区分内に対応するデバイス間エラー回復（パリティ）データを記憶するのをサポートすると決定することができる。この場合も、各区分は、同じ数又は異なる固有の数のストレージデバイスを利用して RAID パリティ情報を記憶することができる。

【0097】

ブロック 1306 において、ストレージコントローラは、Q 個のストレージデバイスが追加のデバイス間エラー回復（パリティ）データをサブシステムの区分にわたって記憶するのをサポートすると決定することができる。ブロック 1308 において、選択したストレージデバイスに、ユーザデータ及び対応する RAID パリティデータを書き込むことができる。再び図 12 を参照すると、区分 1 の RAID アレイ 1210 などの所与の RAID アレイに書き込みが行われた時に、1 又はそれ以上のビットのパリティ情報を生成して、区分 3 のストレージデバイス Q1 に記憶することができる。

【0098】

ストレージコントローラ 174 が、所与の区分内で読み込み再構成を行う状態を検出した（条件付きブロック 1310）場合、及び所与の区分が、RAID パリティ情報を保持しているストレージデバイスを、いくつかの利用できないストレージデバイスを処理するのに十分な数だけ有している（条件付きブロック 1312）場合、ブロック 1314 において、この所与の区分内の 1 又はそれ以上の対応するストレージデバイスを使用して（単複の）再構成読み込み動作を行う。この状態は、デバイス故障によって所与の RAID アレイ内のストレージデバイスを利用できないこと、又はこれらのデバイスが所与のパフォーマンスレベル未満で動作していることを含むことができる。この所与の RAID アレイは、所与の区分内の RAID パリティ情報を記憶している数のストレージデバイスを使用して、最大数の利用できないストレージデバイスを処理することができる。例えば、上記の例の区分 1 の RAID アレイ 1210 が $L + a$ の RAID アレイである場合、この RAID アレイ 1210 は、 $1 \leq k \leq a$ とする k 個のストレージデバイスが利用できない場合に、区分 1 内のストレージデバイスのみを利用して読み込み再構成を行うことができる。

【0099】

所与の区分が、RAID パリティ情報を保持しているストレージデバイスを、いくつかの利用できないストレージデバイスを処理するのに十分な数だけ有していない（条件付きブロック 1312）場合、及びこのいくつかの利用できないストレージデバイスを処理するのに十分な Q 個のストレージデバイスが存在する（条件付きブロック 1316）場合、ブロック 1318 において、1 又はそれ以上の対応する Q 個のストレージデバイスを使用して再構成読み込み動作を行う。読み込み再構成中には、ユーザデータを記憶している他の区分内の 1 又はそれ以上のストレージデバイスにアクセスすることができる。これらのストレージデバイスの選択は、1 又はそれ以上の Q 個のストレージデバイスに記憶されたパリティ情報の導出方法に基づくことができる。例えば、再び図 12 を参照すると、ストレージデバイス Q1 に記憶された対応する RAID パリティ情報を生成するために使用された可能性があるという理由で、読み込み再構成中に区分 2 のストレージデバイス D2 にアクセスすることができる。いくつかの利用できないストレージデバイスを処理するのに十分な数の Q 個のストレージデバイスが存在しない場合（条件付きブロック 1316）、ブロック 1320 において、対応するユーザデータを別のソースから読み込むことができ、又はこれらのユーザデータを失われたものと見なすことができる。

【0100】

なお、上述の実施形態は、ソフトウェアを含むことができる。このような実施形態では、方法及び/又は機構を実装するプログラム命令をコンピュータ可読媒体で搬送し、又はこれに記憶することができる。プログラム命令を記憶するように構成された数多くのタイプの媒体が利用可能であり、これらは、ハードディスク、フロッピー（登録商標）ディスク、CD-ROM、DVD、フラッシュメモリ、プログラマブル ROM（PROM）、ラ

10

20

30

40

50

ンダムアクセスメモリ（ＲＡＭ）及び他の様々な形態の揮発性又は不揮発性ストレージを含む。

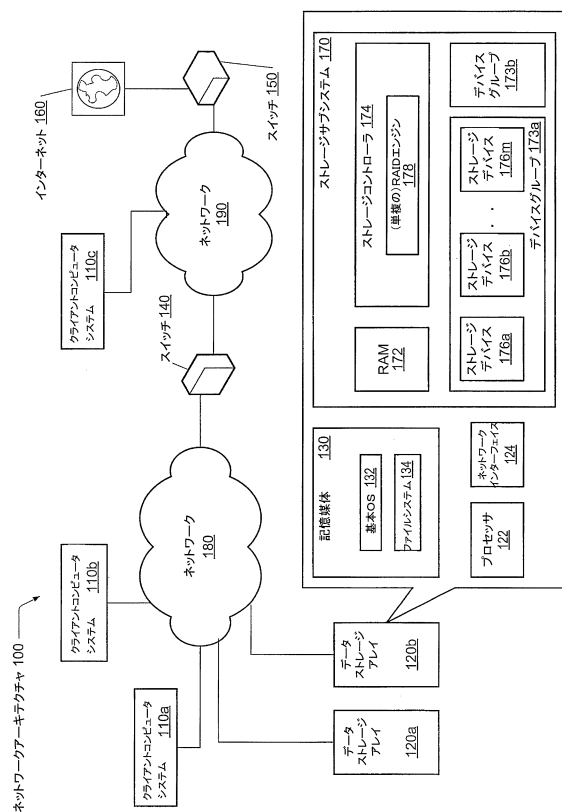
【０１０１】

様々な実施形態では、本明細書で説明した方法及び機構の１又はそれ以上の部分が、クラウドコンピューティング環境の一部を成すことができる。このような実施形態では、１又はそれ以上の様々なモデルにより、インターネットを介してリソースをサービスとして提供することができる。このようなモデルとして、インフラストラクチャ・アズ・ア・サービス（ＩａａＳ）、プラットフォーム・アズ・ア・サービス（ＰａａＳ）、及びソフトウェア・アズ・ア・サービス（ＳａａＳ）を挙げることができる。ＩａａＳでは、コンピュータインフラストラクチャがサービスとして配信される。このような場合、一般にサービスプロバイダがコンピュータ設備を所有し運営する。ＰａａＳモデルでは、開発者がソフトウェアソリューションを開発するために使用するソフトウェアツール及び基本設備をサービスプロバイダがサービスとして供給しホストすることができる。通常、ＳａａＳは、サービスプロバイダのライセンスソフトウェアをサービスオンデマンドとして含む。サービスプロバイダは、このソフトウェアをホストすることができ、又はこのソフトウェアを一定期間にわたって顧客に展開することができる。上記のモデルの数多くの組み合わせが可能であり、企図される。

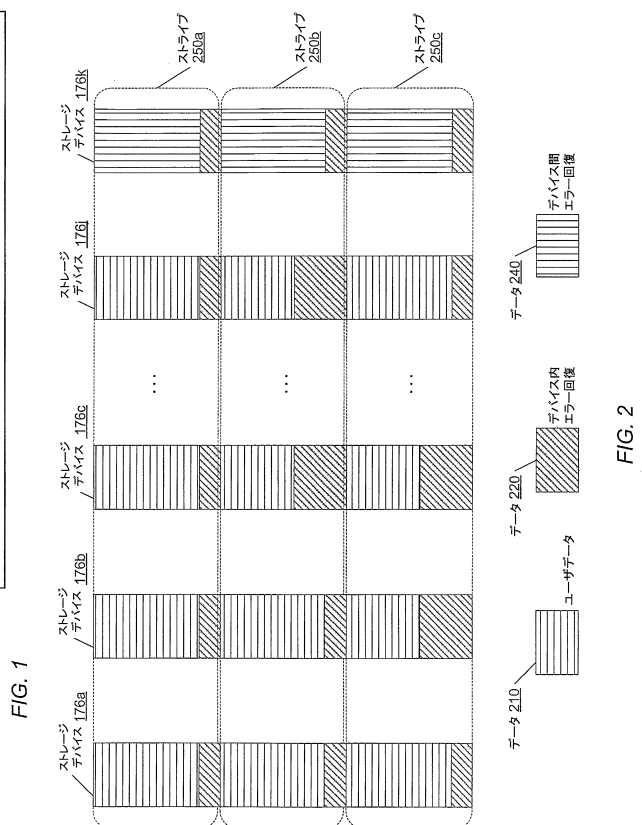
【０１０２】

以上、実施形態についてかなり詳細に説明したが、上記開示を完全に理解すると、当業者には数多くの変形及び修正が明らかになるであろう。以下の特許請求の範囲は、このような変形及び修正を全て含むと解釈すべきであることが意図される。

【図１】



【図２】



【図 3】

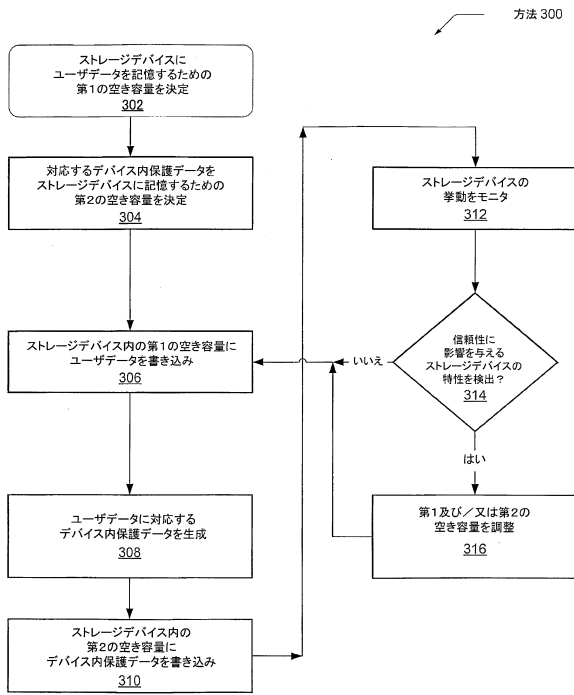


FIG. 3

【図 4】

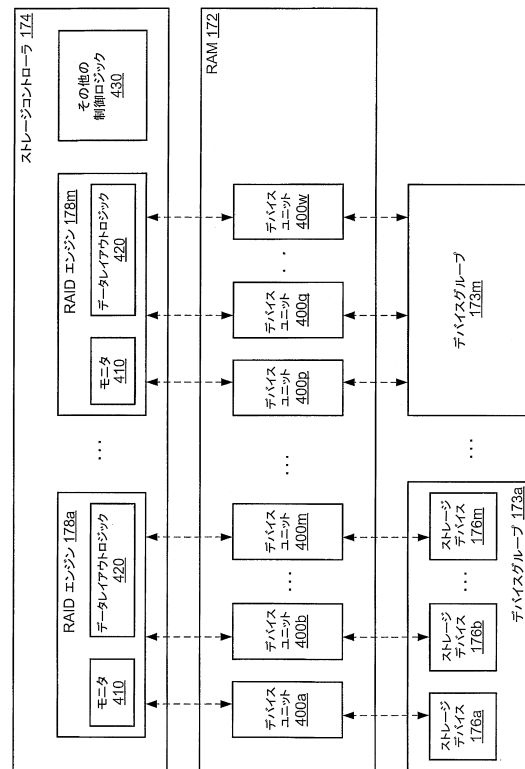


FIG. 4

【図 5】

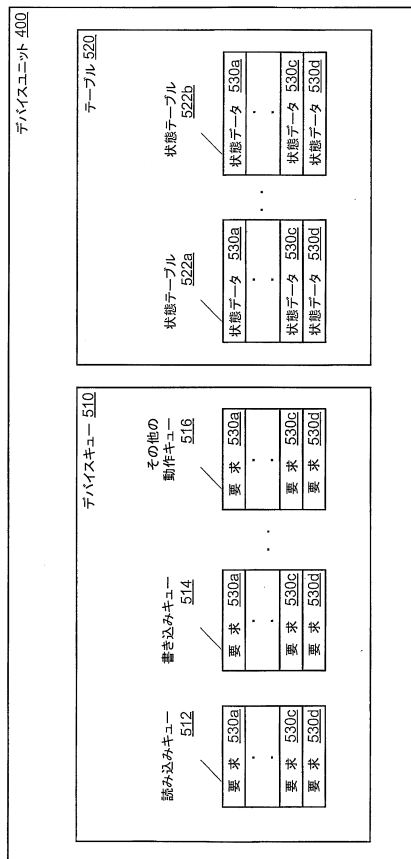


FIG. 5

【図 6】

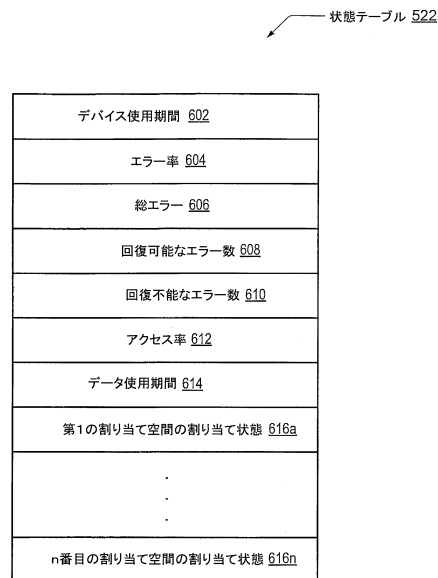
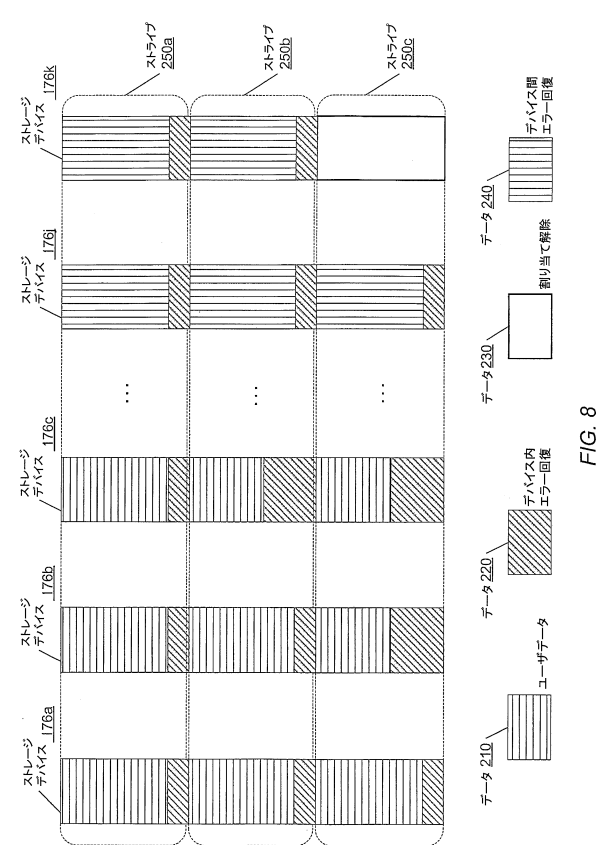


FIG. 6

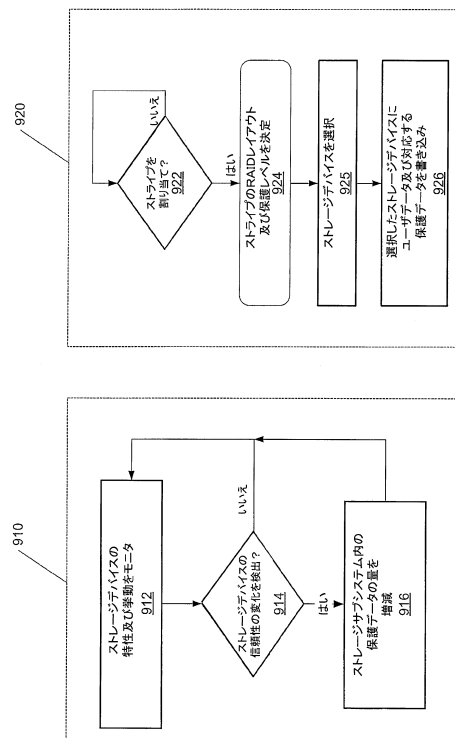
【図 7】



【図 8】



【図 9】



【図 10】

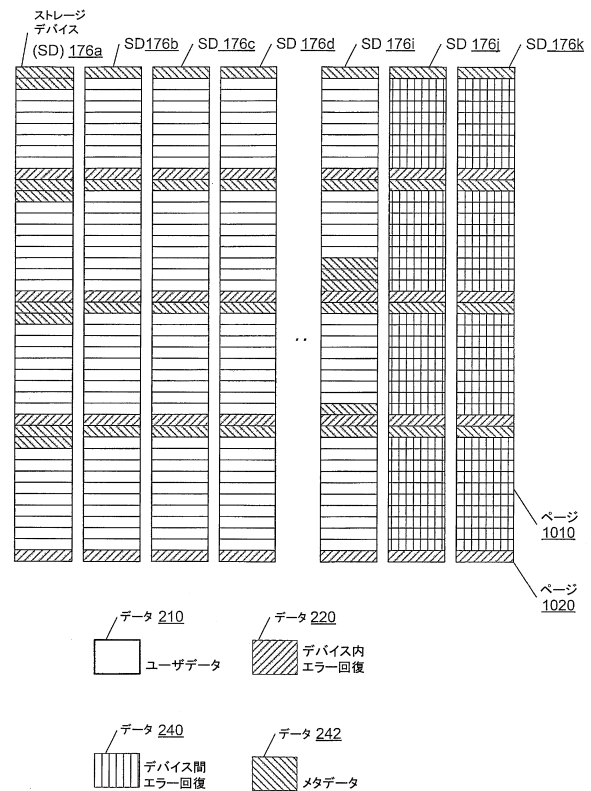


FIG. 10

【図 11A】

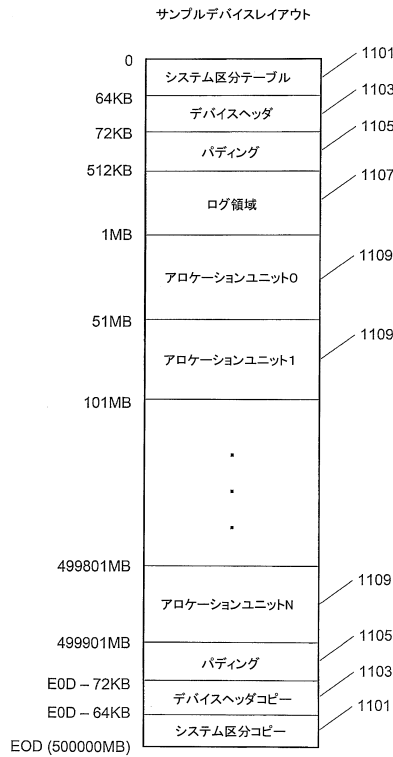


FIG. 11A

【図 11B】

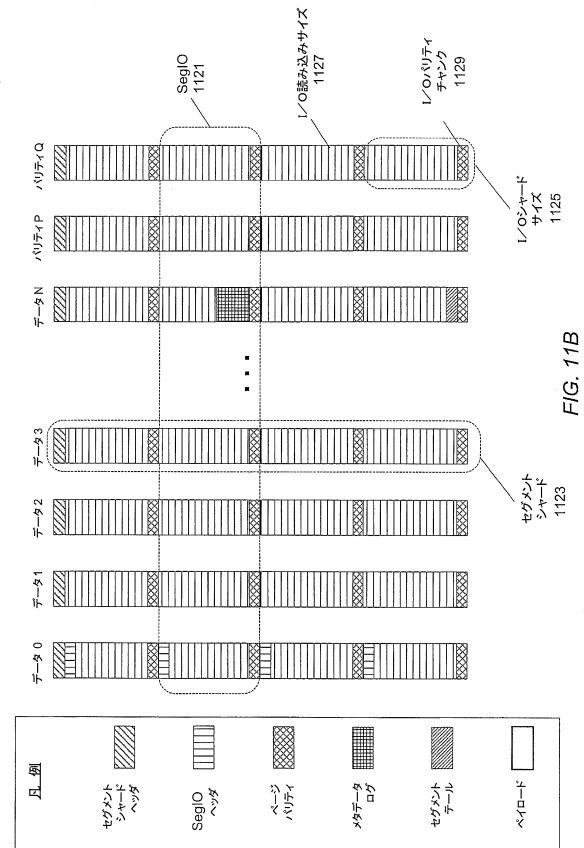


FIG. 11B

【図 11C】

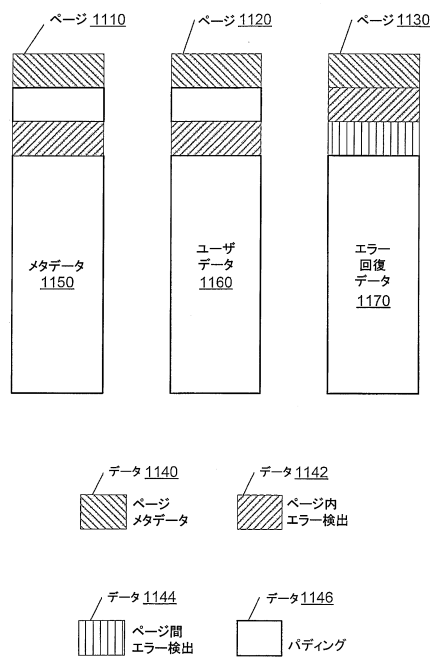


FIG. 11C

【図 12】

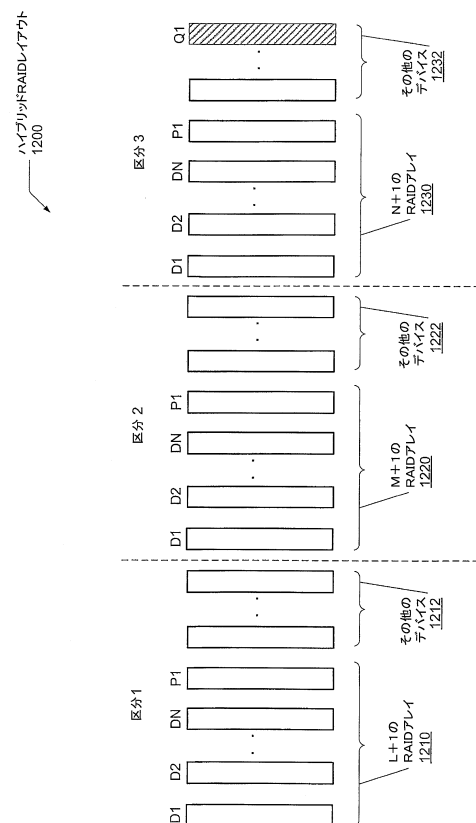


FIG. 12

【図 13】

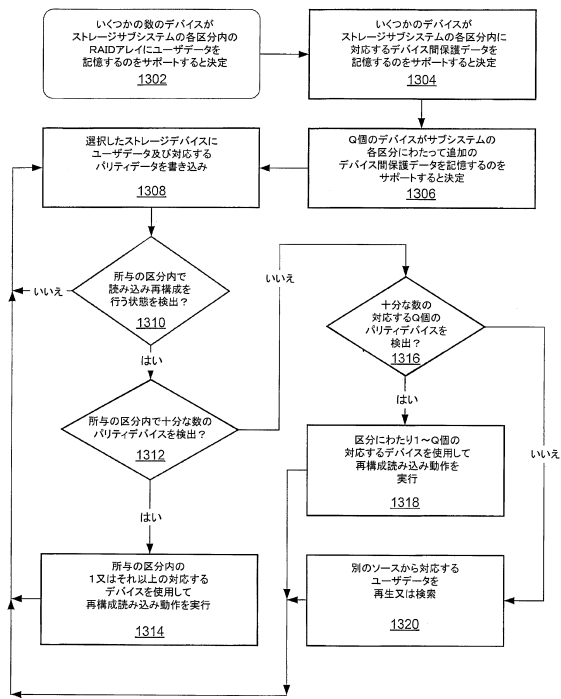


FIG. 13

フロントページの続き

- (72)発明者 ヘイズ, ジョン
アメリカ合衆国・94041・カリフォルニア州・マウンテン ビュー・ハイスクール ウェイ・
800・ナンバー・330
- (72)発明者 ホン, ボー
アメリカ合衆国・94043・カリフォルニア州・マウンテン ビュー・ウェスト ミドルフィー
ルド ロード・1555・アパートメント・95
- (72)発明者 ミラー, イーサン
アメリカ合衆国・95060・カリフォルニア州・サンタクルーズ・カルカー ドライブ・203

審査官 田中 啓介

- (56)参考文献 米国特許出願公開第2009/0210742 (US, A1)
特開平07-200191 (JP, A)
国際公開第2010/106574 (WO, A1)
米国特許第05657439 (US, A)
特開2009-217408 (JP, A)
特開平08-221875 (JP, A)
特開平07-261945 (JP, A)
特開平08-194587 (JP, A)
米国特許出願公開第2008/0229012 (US, A1)

- (58)調査した分野(Int.Cl., DB名)
G06F3/06-3/08
11/08-11/10