

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2010-512568
(P2010-512568A)

(43) 公表日 平成22年4月22日(2010.4.22)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 302A	5B005
G06F 12/08 (2006.01)	G06F 12/08 557	5B065
	G06F 12/08 541Z	
	G06F 3/06 540	

審査請求 未請求 予備審査請求 有 (全 112 頁)

(21) 出願番号 特願2009-540309 (P2009-540309)
 (86) (22) 出願日 平成19年12月6日 (2007.12.6)
 (85) 翻訳文提出日 平成21年7月31日 (2009.7.31)
 (86) 国際出願番号 PCT/US2007/025049
 (87) 国際公開番号 W02008/070173
 (87) 国際公開日 平成20年6月12日 (2008.6.12)
 (31) 優先権主張番号 60/873,111
 (32) 優先日 平成18年12月6日 (2006.12.6)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 60/974,470
 (32) 優先日 平成19年9月22日 (2007.9.22)
 (33) 優先権主張国 米国 (US)

(71) 出願人 509157362
 フリン, デイビッド
 アメリカ合衆国 ユタ州 84093, サ
 ンディ, シェイディメドウッドライブ 88
 56
 (71) 出願人 509157351
 シュトラッサー, ジョン
 アメリカ合衆国 ユタ州 84075, シ
 ラキューズ, サウス 2323
 (71) 出願人 509157339
 サッチャー, ジョナサン
 アメリカ合衆国 ユタ州 84043, リ
 ーハイ, ノース 2080 ウェスト 2
 259

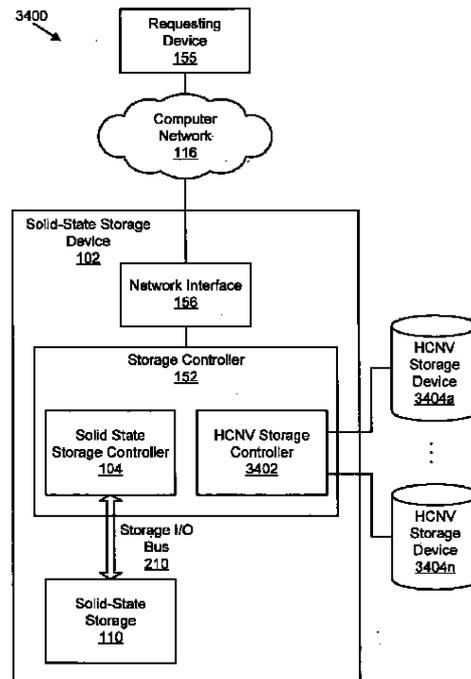
最終頁に続く

(54) 【発明の名称】 高容量不揮発性ストレージ用のキャッシュとしてのソリッドステートストレージのための装置、システム、及び方法

(57) 【要約】

高容量不揮発性ストレージ用のキャッシュとしてのソリッドステートストレージのための装置、システム及び方法が開示されている。その装置、システム及び方法はキャッシュフロントエンドモジュールとキャッシュバックエンドモジュールを具える複数のモジュールで提供される。キャッシュフロントエンドモジュールはストレージ要求に関連するデータ転送を管理する。要求デバイスとソリッドステートストレージとの間のデータ転送は、1又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能し、データ転送はデータ、メタデータ及びメタデータインデックスのうちの1又はそれ以上を含むことができる。ソリッドステートストレージは不揮発性のソリッドステートデータストレージエレメントの阵列を具えることができる。キャッシュバックエンドモジュールはソリッドステートストレージと、1又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理する。

【選択図】 図15



【特許請求の範囲】**【請求項 1】**

1 又はそれ以上の高容量不揮発性（「HCNV」）ストレージデバイスでデータのストレージを管理するための装置であって、当該装置が、

ストレージ要求に関連するデータ転送を管理し、要求デバイスとソリッドステートストレージとの間の前記データ転送が、1 又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能し、前記データ転送がデータ、メタデータ及びメタデータインデックスのうちの1 又はそれ以上を含み、前記ソリッドステートストレージが不揮発性のソリッドステートデータストレージエレメントのアレイを具えるキャッシュフロントエンドモジュールと；

10

前記ソリッドステートストレージと、前記1 又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理するキャッシュバックエンドモジュールと；
を具えることを特徴とする装置。

【請求項 2】

請求項 1 に記載の装置において、前記キャッシュフロントエンドモジュール及び前記キャッシュバックエンドモジュールが、前記ソリッドステートストレージを管理するソリッドステートストレージコントローラと共に同一位置に配置されることを特徴とする装置。

【請求項 3】

請求項 2 に記載の装置において、前記キャッシュフロントエンドモジュール、前記キャッシュバックエンドモジュール、及びソリッドステートストレージコントローラが、前記要求デバイスから自律して動作することを特徴とする装置。

20

【請求項 4】

請求項 1 に記載の装置において、前記ソリッドステートコントローラが、1 又はそれ以上の要求デバイスからオブジェクト要求を提供し、前記ソリッドステートストレージ内の前記オブジェクト要求のオブジェクトを管理するオブジェクトストレージコントローラモジュールを更に具えることを特徴とする装置。

【請求項 5】

請求項 1 に記載の装置が、RAIDレベルの一致する独立ドライブ冗長アレイ（「RAID」）の2 又はそれ以上のHCNVストレージデバイス内の前記ソリッドステートストレージにキャッシュされたデータを記憶するHCNV RAIDモジュールを更に具え、前記データが全体として要求デバイスに見えることを特徴とする装置。

30

【請求項 6】

請求項 1 に記載の装置において、前記ソリッドステートストレージ及び前記1 又はそれ以上のHCNVストレージデバイスが、RAIDグループとして構成されるハイブリッドストレージデバイスセット内にハイブリッドストレージデバイスを具え、前記ソリッドステートストレージにキャッシュされ、HCNVデバイスに後に記憶されるデータセグメントが、ストライプのN個のデータセグメントのうちの1つ、又は、前記ストライプのパリティデータセグメントを含み、前記ハイブリッドストレージデバイスがRAIDストライプのデータセグメントと別個の、1 又はそれ以上のクライアントからストレージ要求を受信することを特徴とする装置。

40

【請求項 7】

請求項 6 に記載の装置において、前記ハイブリッドストレージデバイスが2 又はそれ以上のクライアントから2 又はそれ以上の同時のストレージ要求を受信する、共有されたフロントエンド分散型RAIDグループのストレージデバイスであることを特徴とする装置。

【請求項 8】

請求項 1 に記載の装置において、前記HCNVストレージデバイスがハードディスクドライブ（「HDD」）、光学ドライブ、及びテープストレージのうちの1つであることを特徴とする装置。

【請求項 9】

50

請求項 1 に記載の装置において、前記ソリッドステートストレージ及び前記 1 又はそれ以上の H C N V ストレージデバイスが、ハイブリッドストレージデバイスを含み、当該ハイブリッドストレージデバイスの動作に特有のコードで、1 又はそれ以上の要求デバイスをロードする前に、前記 1 又はそれ以上の要求デバイスに取り付けられた標準デバイスをエミュレートすることによって、前記ハイブリッドストレージデバイスにアクセスを提供する標準デバイスエミュレーションモジュールを更に含み、前記標準デバイスが工業規格 B I O S によって担持されることを特徴とする装置。

【請求項 1 0】

請求項 1 に記載の装置において、前記ソリッドステートストレージデバイスが 2 又はそれ以上の領域に分割可能であり、1 又はそれ以上のパーティションが、前記 H C N V ストレージデバイス用のキャッシュとして機能する前記ソリッドステートストレージと別個の、ソリッドステートストレージとして用いられ得ることを特徴とする装置。

10

【請求項 1 1】

請求項 1 に記載の装置において、1 又はそれ以上のクライアントが、前記ソリッドステートストレージデバイス及び前記 1 又はそれ以上の H C N V ストレージデバイス内に記憶されるファイル又はオブジェクトのうち 1 又はそれ以上の状態を管理するために、前記キャッシュフロントエンドモジュール及び前記キャッシュバックエンドモジュールにキャッシュ制御メッセージを送信することを特徴とする装置。

【請求項 1 2】

請求項 11 に記載の装置において、前記キャッシュ制御メッセージが、
前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージ内のオブジェクト又はファイルの一部を留める (p i n) 制御メッセージ；

20

前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージ内のオブジェクト又はファイルの一部を解放する (u n p i n) 制御メッセージ；

前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージから前記 1 又はそれ以上の H C V N ストレージデバイスまでのオブジェクト又はファイルの一部を消去させる制御メッセージ；

前記キャッシュバックエンドモジュールに、前記 1 又はそれ以上の H C V N ストレージデバイスから前記ソリッドステートストレージに、オブジェクト又はファイルの一部を事前ロードさせる制御メッセージ；及び

30

前記ソリッドステートストレージに、決められたストレージ空き領域量を空けるために、前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージから前記 1 又はそれ以上の H C V N ストレージデバイスまでの 1 又はそれ以上のオブジェクト又はファイルの 1 又はそれ以上の部分をオフロードさせる制御メッセージ；
のうちの 1 又はそれ以上を含むことを特徴とする装置。

【請求項 1 3】

請求項 1 1 に記載の装置において、前記キャッシュ制御メッセージが前記オブジェクト又はファイル用のメタデータ (「キャッシュ制御メタデータ」) を通して通信されることを特徴とする装置。

【請求項 1 4】

請求項 1 3 に記載の装置において、前記キャッシュ制御メタデータが永続的であることを特徴とする装置。

40

【請求項 1 5】

請求項 1 3 に記載の装置において、前記キャッシュ制御メタデータが前記ファイル又はオブジェクトの生成時に設定される属性を通して構築されることを特徴とする装置。

【請求項 1 6】

請求項 1 3 に記載の装置において、前記キャッシュ制御メタデータがファイル又はオブジェクト管理システムから得られることを特徴とする装置。

【請求項 1 7】

請求項 1 に記載の装置が、揮発性キャッシュストレージエレメントを更に含み、前記キ

50

キャッシュフロントエンドモジュール及び前記キャッシュバックエンドモジュールが、前記揮発性キャッシュストレージエレメントにデータを記憶するステップを更に具え、前記ソリッドステートストレージ及び揮発性キャッシュストレージエレメントに記憶されるデータを管理し、前記バックエンドストレージモジュールが、前記揮発性キャッシュストレージエレメントと、前記ソリッドステートストレージと、前記HCNVストレージデバイスとの間のデータ転送を更に管理することを特徴とする装置。

【請求項 18】

請求項 17 に記載の装置において、前記HCNVストレージデバイスに記憶されるオブジェクト及びファイル用のメタデータ及びインデックスメタデータのうちの 1 又はそれ以上が、前記ソリッドステートストレージデバイス及び前記揮発性キャッシュストレージエレメント内に維持されることを特徴とする装置。

10

【請求項 19】

請求項 1 に記載の装置において、前記HCNVストレージデバイスに記憶されるオブジェクト及びファイル用のメタデータ及びインデックスメタデータのうちの 1 又はそれ以上が、前記ソリッドステートストレージデバイス内に維持されることを特徴とする装置。

【請求項 20】

請求項 1 に記載の装置において、前記ソリッドステートストレージ及び前記 1 又はそれ以上のHCNVストレージデバイスが、ストレージデバイスを持って、前記HCNVストレージデバイスが前記ストレージデバイスに接続されたクライアントの表示から隠されるようにすることを特徴とする装置。

20

【請求項 21】

1 又はそれ以上の高容量不揮発性（「HCNV」）ストレージデバイスでデータのストレージを管理するためのシステムであって、当該システムが、

不揮発性のソリッドステートデータストレージエレメントのアレイを具えるソリッドステートストレージと；

1 又はそれ以上のHCNVストレージデバイスと；

ソリッドステートストレージコントローラと；

HCNVストレージデバイスコントローラと；

ストレージ要求に関連するデータ転送を管理し、要求デバイスと前記ソリッドステートストレージとの間の前記データ転送が、前記 1 又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能し、前記データ転送がデータ、メタデータ及びメタデータインデックスのうちの 1 又はそれ以上を含むキャッシュフロントエンドモジュールと；

30

前記ソリッドステートストレージと、前記 1 又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理するキャッシュバックエンドモジュールと；

を具えるストレージコントローラと；

を具えることを特徴とするシステム。

【請求項 22】

請求項 21 に記載のシステムが、前記ストレージコントローラに接続されたネットワークインタフェースを更に具え、当該ネットワークインタフェースがコンピュータネットワークを通して前記要求デバイスと前記ソリッドステートストレージコントローラとの間のデータ転送を促進することを特徴とするシステム。

40

【請求項 23】

請求項 21 に記載のシステムがサーバを更に具え、当該サーバが前記ソリッドステートストレージと、前記 1 又はそれ以上のHCNVストレージデバイスと、前記ストレージコントローラとを具えることを特徴とするシステム。

【請求項 24】

請求項 21 に記載のシステムにおいて、前記 1 又はそれ以上のHCNVストレージデバイスがストレージエリアネットワーク（「SAN」）を通して前記ストレージコントローラに接続されることを特徴とするシステム。

50

【請求項 25】

コンピュータプログラムプロダクトであって、1又はそれ以上の高容量不揮発性（「HCNV」）ストレージデバイスでデータのストレージを管理するための動作を行うのに実行可能な、コンピュータが利用可能なプログラムコードを有するコンピュータ可読媒体を具え、前記コンピュータプログラムプロダクトの動作が、

ストレージ要求に関連するデータ転送を管理するステップであって、要求デバイスとソリッドステートストレージとの間の前記データ転送が、1又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能し、前記データ転送がデータ、メタデータ及びメタデータインデックスのうち1又はそれ以上を含み、前記ソリッドステートストレージが不揮発性のソリッドステートデータストレージエレメントのアレイを具えるステップと；

前記ソリッドステートストレージと、前記1又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理するステップと；

を含むことを特徴とするコンピュータプログラムプロダクト。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互引用

この出願は、David Flynnらによって2006年12月6日出願の“Elemental Blade System”と題された米国仮出願番号60/873,111、及び、David Flynnらによって2007年9月22日出願の“Apparatus, System, and Method for Object-Oriented Solid-State Storage”と題された米国仮出願番号60/974,470の優先権を主張し、上記仮出願は、引用によって本明細書に組み込まれている。

【0002】

本発明は、データを管理することに関し、特に、高容量不揮発性ストレージデバイス用キャッシュとしてソリッドステートストレージを用いることに関する。

【背景技術】

【0003】

一般的にキャッシュは、しばしばアクセスされ、アプリケーション又はオペレーティングシステムの一部として配置されるデータが、ハードディスクドライブ（「HDD」）、光学ドライブ、テープストレージ等のような高容量不揮発性（「HCNV」）ストレージデバイスを通してアクセスされなければならない場合よりも更に早い次のアクセスを伴うキャッシュ内に記憶されるので有利である。キャッシュは一般的にはコンピュータ内に含まれる。

【発明の概要】

【0004】

いくつかのストレージデバイス及びシステムはHCNVストレージデバイスにキャッシュを具えている。いくつかのHCNVストレージデバイスは不揮発性ソリッドステートキャッシュを含み、これらはアクセス時間を低減する利益を提供するが、通常は限定されるHCNVストレージデバイスインタフェースの容量と一致する性能を提供するのみである。マザーボード上に一般的に配置されるいくつかの不揮発性ソリッドステートキャッシュストレージデバイスが存在するが、これらのデバイスはキャッシュコヒーレンスが提供されないようなマルチクライアント環境において用いることができない。HCNVデバイスのいくつかのコントローラは更にキャッシュを具える。冗長HCNVキャッシュコントローラが複数クライアント間で共有される場合、データが破損されていないことを確認するのに高性能なキャッシュコヒーレンシアルゴリズムが要求される。

【0005】

一般的に、キャッシュはDRAM内に実装され、キャッシュ容量をプレミアムにし、性

10

20

30

40

50

能ごとに比較的高い電源を要求する。揮発性キャッシュを支持する電源が失われる場合、キャッシュに記憶されたデータが失われる。一般的にバッテリーバックアップが、電源故障の場合にデータ損失を避けるのに用いられるが、バッテリーバックアップの停止前に不揮発性メモリに対するキャッシュを消去するかなりの可能性を有する。更に、バッテリーバックアップシステムは電力を消費し、冗長性を要求し、信頼性に負の影響を与え、空き領域を消費する。バッテリーは標準ベースで更に提供されなければならない、バッテリーバックアップは比較的高価である。

【 0 0 0 6 】

前述の考察から、キャッシュとしてソリッドステートストレージを用いてデータを管理する装置、システム及び方法に対するニーズが存在することは明らかである。有利には、このような装置、システム及び方法は少ない電力を消費し、かなり大きな容量を有し、キャッシュに記憶されたデータを維持するのにバッテリーバックアップを要求しない不揮発性キャッシュを提供するであろう。

10

【 0 0 0 7 】

本発明は当該技術分野の現在の状態に応じて、特に、データストレージを管理するための現在利用可能なシステムによって未だ十分に解決されない当該技術分野における問題及びニーズに応じて開発されてきた。従って、本発明は、当該技術分野の上述の欠点の大部分又は総てを克服する1又はそれ以上の高容量不揮発性(「HCNV」)ストレージデバイスでデータのストレージを管理するための装置、システム及び方法を提供するように開発されてきた。

20

【 0 0 0 8 】

一実施例においては、装置はキャッシュフロントエンドモジュールと、キャッシュバックエンドモジュールとを具える複数のモジュールで提供される。キャッシュフロントエンドモジュールはストレージ要求に関連するデータ転送を管理する。データ転送は要求デバイスと、ソリッドステートストレージとの間で、1又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能し、データ転送はデータ、メタデータ及びメタデータインデックスのうちの1又はそれ以上を含みうる。ソリッドステートストレージは不揮発性のソリッドステートデータストレージエレメントのアレイを具えることができる。キャッシュバックエンドモジュールはソリッドステートストレージと、1又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理する。

30

【 0 0 0 9 】

装置の一実施例においては、キャッシュフロントエンドモジュール及びキャッシュバックエンドモジュールはソリッドステートストレージを管理するソリッドステートストレージコントローラと共に同一位置に配置される。更なる実施例においては、キャッシュフロントエンドモジュール、キャッシュバックエンドモジュール、及びソリッドステートストレージコントローラは要求デバイスから自律して動作できる。

【 0 0 1 0 】

一実施例においては、装置はRAIDレベルの一致する独立ドライブ冗長アレイ(「RAID」)の2又はそれ以上のHCNVストレージデバイス内のソリッドステートストレージにキャッシュされたデータを記憶するHCNV RAIDモジュールを具える。データは全体として要求デバイスに見えうる。別の実施例においては、ソリッドステートストレージ及び1又はそれ以上のHCNVストレージデバイスは、RAIDグループとして構成されるハイブリッドストレージデバイスセット内にハイブリッドストレージデバイスを具えることができる。ソリッドステートストレージにキャッシュされ、HCNVデバイスで後に記憶されるデータセグメントは、ストライプのN個のデータセグメントのうちの1つ、又はストライプのパリティデータセグメントを含みうる。ハイブリッドストレージデバイスは一般的にRAIDストライプのデータセグメントと別個に、1又はそれ以上のクライアントからストレージ要求を受信する。更なる実施例においては、ハイブリッドストレージデバイスは2又はそれ以上のクライアントから2又はそれ以上の同時のストレージ要求を受信する、共有されたフロントエンド分散型RAIDグループのストレージデバイ

40

50

スにできる。

【0011】

装置の更なる実施例においては、HCNVストレージデバイスはハードディスクドライブ(「HDD」)、光学ドライブ、又はテープストレージにできる。別の実施例においては、ソリッドステートストレージ及び1又はそれ以上のHCNVストレージデバイスは、ハイブリッドストレージデバイスにできる。一実施例においては、装置は、ハイブリッドストレージデバイスの動作に特有のコードで、1又はそれ以上の要求デバイスをロードする前に、1又はそれ以上の要求デバイスに取り付けられた標準デバイスをエミュレートすることによって、ハイブリッドストレージデバイスにアクセスを提供する標準デバイスエミュレーションモジュールを更に具える。標準デバイスは一般的には工業規格BIOSによって担持されうる。

10

【0012】

別の実施例においては、ソリッドステートストレージデバイスは2又はそれ以上の領域に分割可能であり、1又はそれ以上のパーティションが、HCNVストレージデバイス用のキャッシュとして機能するソリッドステートストレージと別個の、ソリッドステートストレージとして用いられ得る。更に別の実施例においては、1又はそれ以上のクライアントは、ソリッドステートストレージデバイス及び1又はそれ以上のHCNVストレージデバイス内に記憶されるファイル又はオブジェクトのうち1又はそれ以上の状態を管理するために、キャッシュフロントエンドモジュール及びキャッシュバックエンドモジュールにキャッシュ制御メッセージを送信する。

20

【0013】

装置の一実施例においては、キャッシュ制御メッセージは1又はそれ以上の制御メッセージを含みうる。制御メッセージの様々な実施例は、キャッシュバックエンドモジュールに、ソリッドステートストレージ内のオブジェクト又はファイルの一部を留める(pin)制御メッセージ、又は、キャッシュバックエンドモジュールに、ソリッドステートストレージ内のオブジェクト又はファイルの一部を解放する(unpin)制御メッセージを含みうる。制御メッセージの他の実施例は、キャッシュバックエンドモジュールに、ソリッドステートストレージから1又はそれ以上のHCNVストレージデバイスまでのオブジェクト又はファイルの一部を消去させる制御メッセージ、又は、キャッシュバックエンドモジュールに、1又はそれ以上のHCNVストレージデバイスからソリッドステートストレージに、オブジェクト又はファイルの一部を事前ロードさせる制御メッセージを含みうる。制御メッセージの更なる別の実施例は、ソリッドステートストレージに、決められたストレージ空き領域量を空けるために、キャッシュバックエンドモジュールに、ソリッドステートストレージから1又はそれ以上のHCNVストレージデバイスまでの1又はそれ以上のオブジェクト又はファイルの1又はそれ以上の部分をオフロードさせる制御メッセージである。一実施例においては、キャッシュ制御メッセージはオブジェクト又はファイル用のメタデータ(「キャッシュ制御メタデータ」)を通して通信される。更なる実施例においては、キャッシュ制御メタデータは永続的にできる。別の実施例においては、キャッシュ制御メタデータはファイル又はオブジェクトの生成時に設定される属性を通して構築されうる。更なる別の実施例においては、キャッシュ制御メタデータはファイル又はオブジェクト管理システムから得うる。

30

40

【0014】

装置の一実施例においては、装置は揮発性キャッシュストレージエレメントを具えることができ、キャッシュフロントエンドモジュール及びキャッシュバックエンドモジュールは、揮発性キャッシュストレージエレメントにデータを記憶し、ソリッドステートストレージ及び揮発性キャッシュストレージエレメントに記憶されるデータを管理できる。バックエンドストレージモジュールは、揮発性キャッシュストレージエレメントと、ソリッドステートストレージと、HCNVストレージデバイスとの間のデータ転送を更に管理できる。更なる実施例においては、HCNVストレージデバイスに記憶されるオブジェクト及びファイル用のメタデータ及び/又はインデックスメタデータは、ソリッドステートスト

50

レージデバイス及び揮発性キャッシュストレージエレメント内に維持されうる。

【0015】

装置の更なる実施例においては、HCVNストレージデバイスに記憶されるオブジェクト及びファイル用のメタデータ及び/又はインデックスメタデータは、ソリッドステートストレージデバイス内に維持されうる。別の実施例においては、ソリッドステートストレージ及び1又はそれ以上のHCNVストレージデバイスが、ストレージデバイスを含み、HCNVストレージデバイスがストレージデバイスに接続されたクライアントの表示から隠されるようにできる。

【0016】

本発明のシステムが更に提供される。システムはモジュール及び装置に関する上述した実施例を実質的に具える。一実施例においては、システムは不揮発性のソリッドステートデータストレージエレメントのレイを具えるソリッドステートストレージを具える。システムは、1又はそれ以上のHCNVストレージデバイスと、ストレージコントローラとを更に具える。一実施例においては、ストレージコントローラは、ソリッドステートストレージコントローラと、HCNVストレージデバイスコントローラとを具えることができる。ストレージコントローラはキャッシュフロントエンドモジュールと、キャッシュバックエンドモジュールとを更に具えることができる。キャッシュフロントエンドモジュールは、ストレージ要求に関連するデータ転送を管理する。データ転送は一般的に、要求デバイスとソリッドステートストレージとの間で1又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能している。データ転送はデータ、メタデータ及びメタデータインデックスのうち1又はそれ以上を含みうる。キャッシュバックエンドモジュールは、ソリッドステートストレージと、1又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理する。

【0017】

一実施例においては、システムはストレージコントローラに接続されたネットワークインタフェースを更に具え、当該ネットワークインタフェースがコンピュータネットワークを通して要求デバイスとソリッドステートストレージコントローラとの間のデータ転送を促進する。別の実施例においては、システムはソリッドステートストレージと、1又はそれ以上のHCNVストレージデバイスと、ストレージコントローラとを具えるサーバを具える。更なる別の実施例においては、1又はそれ以上のHCNVストレージデバイスはストレージエリアネットワーク(「SAN」)を通してストレージコントローラに接続される。

【0018】

本発明の方法は、複数のホスト間でデバイスを共有するために更に提供される。開示された実施例における方法は、説明した装置及びシステムの動作に関する、上に提供した機能を実行するのに必要なステップを実質的に含む。一実施例においては、方法は、ストレージ要求に関連するデータ転送を管理するステップを含み、データ転送は、要求デバイスとソリッドステートストレージとの間で、1又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能している。データ転送はデータ、メタデータ及びメタデータインデックスのうち1又はそれ以上を含みうる。ソリッドステートストレージは、不揮発性のソリッドステートデータストレージエレメントのレイを具えることができる。方法は、ソリッドステートストレージと、1又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理するステップを更に具えることができる。

【0019】

本明細書全体にわたる、特徴、利点または同様の単語の引用は、本発明で実現される特徴および利点の総てが、本発明のいずれかの1つの実施例に存在すべき、または、存在することを意味するものではない。むしろ、特徴及び利点に関する用語は、ある実施例に関連して記載される特定の特徴、利点または特性が、本発明の少なくとも1つの実施例に含まれることを意味すると理解されたい。従って、特徴、利点および同様の用語の記載は、本明細書全体にわたって、必須ではないが、同一の実施例を意味する。

10

20

30

40

50

【0020】

更に、記載された本発明の特徴、利点および特性は、1又はそれ以上の実施例においてあらゆる適宜な方法で組み合わせることができる。関連分野の当業者は、特定の実施例の1又はそれ以上の特定の特徴または利点がなくとも、本発明が実施可能であることが分かるであろう。その他の例においては、更なる特徴および利点は、本発明の総ての実施例において示されていない特定の実施例で分かるであろう。

【0021】

本発明の特徴および利点は、以下の説明および添付の請求の範囲から完全に明らかになり、また、以下に記載する本発明の実施形態によって理解されるであろう。

【図面の簡単な説明】

【0022】

本発明の利点が容易に理解されるように、簡単に上述した本発明のより具体的な説明は、添付の図面に例示した特定の実施例を参照することによってなされるであろう。これらの図面は、本発明の一般的な実施例のみを示し、従って、本発明の範囲を限定するものではないと理解すべきであり、本発明は、添付の図面の利用を通じて、更に具体的かつ詳細に記載および説明されるであろう。

【0023】

【図1A】図1Aは、本発明によるソリッドステートストレージデバイス中のデータ管理のためのシステムの一実施例を示した概略ブロック図である。

【図1B】図1Bは、本発明によるストレージデバイス中のオブジェクト管理のためのシステムの一実施例を示した概略ブロック図である。

【図1C】図1Cは、本発明によるインサーバストレージエリアネットワーク用のシステムの一実施例を示した概略ブロック図である。

【図2A】図2Aは、本発明によるストレージデバイスのオブジェクト管理のための装置の一実施例を示した概略ブロック図である。

【図2B】図2Bは、本発明によるソリッドステートストレージデバイス中のソリッドステートストレージデバイスコントローラの一実施例を示した概略ブロック図である。

【図3】図3は、本発明によるソリッドステートストレージデバイス内の書き込みデータパイプラインおよび読み出しデータパイプラインを有するソリッドステートストレージコントローラを示す概略ブロック図である。

【図4A】図4Aは、本発明によるソリッドステートストレージコントローラ内のバンクインタリーブコントローラの一実施例を示した概略ブロック図である。

【図4B】図4Bは、本発明によるソリッドステートストレージコントローラ内のバンクインタリーブコントローラの代替的な実施例

【図5A】図5Aは、本発明によるデータパイプラインを用いたソリッドステートストレージデバイス内のデータを管理するための方法の一実施例を示した概略フローチャート図である。

【図5B】図5Bは、本発明によるインサーバSANのための方法の一実施例を示した概略フローチャート図である。

【図6】図6は、本発明によるデータパイプラインを用いたソリッドステートストレージデバイス内のデータを管理するための方法の別の実施例を示した概略フローチャート図である。

【図7】図7は、本発明によるバンクインタリーブを用いたソリッドステートストレージデバイス内のデータを管理するための方法の一実施例を示した概略フローチャート図である。

【図8】図8は、本発明によるソリッドステートストレージデバイス内のガベージコレクション用の装置の一実施例を示した概略ブロック図である。

【図9】図9は、本発明によるソリッドステートストレージデバイス内のガベージコレクション用の方法の一実施例を示した概略フローチャート図である。

【図10】図10は、本発明による進行型RAID、フロントエンド分散型RAID、及

10

20

30

40

50

び共有されたフロントエンド分散型 R A I D 用のシステムの一実施例を示した概略ブロック図である。

【図 1 1】図 1 1 は、本発明によるフロントエンド分散型 R A I D 用の装置の一実施例を示した概略ブロック図である。

【図 1 2】図 1 2 は、本発明によるフロントエンド分散型 R A I D 用の方法の一実施例を示した概略フローチャート図である。

【図 1 3】図 1 3 は、本発明による共有されたフロントエンド分散型 R A I D 用の装置の一実施例を示した概略ブロック図である。

【図 1 4】図 1 4 は、本発明による共有されたフロントエンド分散型 R A I D 用の方法の一実施例を示した概略フローチャート図である。

【図 1 5】図 1 5 は、本発明による高容量不揮発性ストレージデバイス用のキャッシュとしてソリッドステートストレージを有するシステムの一実施例を示した概略ブロック図である。

【図 1 6】図 1 6 は、本発明による高容量不揮発性ストレージデバイス用のキャッシュとしてソリッドステートストレージを有する装置の一実施例を示した概略ブロック図である。

【図 1 7】図 1 7 は、本発明による高容量不揮発性ストレージデバイス用のキャッシュとしてソリッドステートストレージを有する方法の一実施例を示した概略フローチャート図である。

【発明を実施するための形態】

【0024】

本明細書に記載された多数の機能ユニットは、その実装独立性をより具体的に強調するために、モジュールとして表示した。例えば、モジュールは、カスタム V L S I 回路またはゲートアレイ、論理チップのような標準規格の半導体、トランジスタまたは他の別個の構成要素を具えるハードウェア回路として実装可能である。更に、モジュールは、フィールドプログラマブルゲートアレイ、プログラマブル配列論理回路、プログラマブル論理デバイスなどのような、プログラム可能なハードウェアデバイスに実装可能である。

【0025】

更に、モジュールは、様々なプロセッサによって実行されるソフトウェアに実装可能である。識別モジュールの実行可能コードは、例えばオブジェクト、手続きまたは関数として構成される、計算機命令の 1 又はそれ以上の物理又は論理ブロックを含む。にもかかわらず、実行可能な識別モジュールは、物理的に一緒に割り当てられる必要はなく、論理的に互いに接続される場合、モジュールを構成し、モジュールのための規定された目的を達成する、異なる位置に記憶される離隔命令を具えることができる。

【0026】

実際に、モジュールの実行可能コードは、単一の命令または多数の命令にでき、複数の異なるコードセグメントにわたり、異なるプログラム間で、および、複数のメモリデバイスにわたってさえ割り当てることができる。同様に、動作データは識別され、本明細書でモジュール内部に例示されており、適宜な形態で組み込まれ、あらゆる適宜な種類のデータ構造内で構成される。動作データは、単一のデータセットとして集められ、異なるストレージデバイスを含む異なる位置にわたって割り当てられ、少なくとも部分的に、システムまたはネットワークの単なる電気信号として存在しうる。モジュールまたはモジュールの一部がソフトウェアに実装される場合、ソフトウェア部分は、1 又はそれ以上のコンピュータ可読メディアに記憶される。

【0027】

本明細書全体における「一実施例 (one embodiment)」、「ある実施例 (an embodiment)」または同様の用語の引用は、実施例に関連して記載される特定の特徵、構造または特性が本発明の少なくとも 1 つの実施例に含まれることを意味する。従って、本明細書全体における「一実施例においては (in one embodiment)」、「ある実施例においては (in an embodiment)」お

10

20

30

40

50

よび同様の用語のフレーズの出現は、必ずではないが、同一の実施例を総て引用している。

【0028】

信号担持媒体については、信号を生成し、信号が生成され、または、デジタル処理装置で機械が読み出し可能な命令のプログラムの実行をさせることが可能なあらゆる形態をとることができる。信号担持媒体は、伝送線路、コンパクトディスク、デジタルビデオディスク、磁気テープ、ベルヌーイドライブ、磁気ディスク、パンチカード、フラッシュメモリ、集積回路または他のデジタル処理装置のメモリデバイスに組み込まれる。

【0029】

更に、記載された本発明の特徴、構造または特性は、1又はそれ以上の実施例において適宜な方法で組み合わせることができる。以下の説明において、例えば、プログラミング、ソフトウェアモジュール、ユーザの選択項目、ネットワークトランザクション、データベースクエリ、データベース構造、ハードウェアモジュール、ハードウェア回路、ハードウェアチップ等のような、多数の特定の詳細な説明が提供されて、本発明の実施例を完全な理解を提供する。しかしながら、関連分野の当業者は、本発明が1又はそれ以上の特定の詳細な説明なしで、あるいは、他の方法、構成、材料などを用いて実施できることが分かるであろう。別の例においては、既知の構造、材料または動作は、本発明の態様があいまいになるのを避けるために、示していないか、または、詳細に記載していない。

【0030】

本明細書に含まれる概略的フローチャートは、論理フローチャート図として一般的に説明されている。このように、示された順番および示されたステップは、本方法の一実施例を示す。他のステップ及び方法は、例示された方法の関数、論理または1又はそれ以上のステップに対する効果、またはその一部において同等であると考えられる。更に、使用したフォーマットおよび符号は、方法の論理的なステップを説明するために提供されたものであり、この方法の範囲を限定するものではないと理解されたい。様々な矢印および線がフローチャート図で用いられているが、対応する方法の範囲を限定するものではないことを理解されたい。実際に、いくつかの矢印または他の結合子は、この方法の論理フローを示すためにのみ用いられている。例えば、矢印は、示された方法の列挙されたステップ間における不特定期間の待機又はモニタリング時間を示す。更に、特定の方法が生じる順番は、示された対応するステップの順番に厳密に準拠してもしなくてもよい。

【0031】

ソリッドステートストレージシステム

【0032】

図1Aは、本発明によるソリッドステートストレージデバイスのデータ管理用のシステム100の一実施例を示した概略ブロック図である。システム100は、以下に示すように、ソリッドステートストレージデバイス102、ソリッドステートストレージコントローラ104、書き込みデータパイプライン106、読み出しデータパイプライン108、ソリッドステートストレージ110、コンピュータ112、クライアント114およびコンピュータネットワーク116を具える。

【0033】

システム100は、少なくとも1のソリッドステートストレージデバイス102を具える。別の実施例においては、システム100は、2又はそれ以上のソリッドステートストレージ102を具える。各ソリッドステートストレージデバイス102は、フラッシュメモリ、ナノランダムアクセスメモリ(「ナノRAMまたはNRAM」)、磁気抵抗RAM(「MRAM」)、ダイナミックRAM(「DRAM」)、相変化RAM(「PRAM」)などのような、不揮発性ソリッドステートストレージ110を具える。ソリッドステートストレージデバイス102は、図2および図3で更に詳細に示す。ソリッドステートストレージデバイス102は、コンピュータネットワーク116を通して、クライアント114に接続したコンピュータ112内に示されている。一実施例においては、ソリッドステートストレージデバイス102は、コンピュータ112の内部にあり、PCI(per

ipheral component interconnect : 周辺装置相互接続規格) エクスプレス(「PCI-e」)バス、シリアルアドバンスドテクノロジーアタッチメント(「シリアルATA」)バスなどのような、システムバスを用いて接続される。別の実施例においては、ソリッドステートストレージデバイス102は、コンピュータ112の外部にあり、ユニバーサルシリアルバス(「USB」)接続、IEEE(Institute of Electrical and Electronics Engineers : 米国電気電子技術者協会)1394バス(「FireWire」)などを用いて接続される。他の実施例においては、ソリッドステートストレージデバイス102は、インフィニバンドまたはPCIエクスプレスアドバンスドスイッチング(「PCIe-AS」)などのような、外部の電気式または光学式バス拡張機能又はバス型ネットワークソリューションを用い、PCI(peripheral component interconnect : 周辺装置相互接続規格)エクスプレスバスを用いたコンピュータ112に接続される。

10

20

30

40

50

【0034】

様々な実施例において、ソリッドステートストレージ102は、デュアルインラインメモリモジュール(「DIMM」)、ドータカードまたはマイクロモジュールの形態でもよい。別の実施例においては、ソリッドステートストレージデバイス102は、ラック取付型ブレード内のエレメントである。別の実施例においては、ソリッドステートストレージデバイス102は、より高位のアセンブリ(例えば、マザーボード、ラップトップ、グラフィックプロセッサ)に直接統合したパッケージ内に含まれる。別の実施例においては、ソリッドステートストレージデバイス102を具える個々の構成は、中位の実装なしにより高位のアセンブリに直接統合される。

【0035】

ソリッドステートストレージデバイス102は、1又はそれ以上のソリッドステートストレージコントローラ104を具え、各々は、書き込みデータパイプライン106および読み出しデータパイプライン108を具え、各々は、ソリッドステートストレージ110を具え、これらは、図2および図3で以下に詳細に示す。

【0036】

システム100は、ソリッドステートストレージデバイス102に接続した1又はそれ以上のコンピュータ112を具える。コンピュータ112は、ホスト、サーバ、ストレージエリアネットワーク(「SAN」)のストレージコントローラ、ワークステーション、パーソナルコンピュータ、ラップトップコンピュータ、ハンドヘルドコンピュータ、スーパーコンピュータ、コンピュータクラスタ、ネットワークスイッチ、ルータまたは機器、データベースまたはストレージ機器、データ取得システムまたはデータキャプチャシステム、診断システム、試験システム、ロボット、携帯型電子機器、ワイヤレスデバイスなどに行うことができる。別の実施例においては、コンピュータ112は、クライアントでもよく、ソリッドステートストレージデバイス102は、コンピュータ112から送信されるデータ要求を提供するために自律動作する。この実施例においては、コンピュータ112およびソリッドステートストレージデバイス102は、コンピュータネットワーク、システムバス、または、コンピュータ112と自律型ソリッドステートストレージデバイス102との間の接続に適切なその他の通信手段を用いて、接続可能である。

【0037】

一実施例においては、システム100は、1又はそれ以上のコンピュータネットワーク116を通して、1又はそれ以上のコンピュータ112に接続された1又はそれ以上のクライアント114を具える。クライアント114は、ホスト、サーバ、SANのストレージコントローラ、ワークステーション、パーソナルコンピュータ、ラップトップコンピュータ、ハンドヘルドコンピュータ、スーパーコンピュータ、コンピュータクラスタ、ネットワークスイッチ、ルータまたは機器、データベースまたはストレージ機器、データ取得またはデータキャプチャシステム、診断システム、試験システム、ロボット、携帯電子機器、ワイヤレスデバイスなどに行うことができる。コンピュータネットワーク116は、インターネッ

ト、ワイドエリアネットワーク（「WAN」）、メトロポリタンエリアネットワーク（「MAN」）、ローカルエリアネットワーク（「LAN」）、トークンリング、ワイヤレスネットワーク、ファイバチャネルネットワーク、SAN、ネットワーク接続ストレージ（「NAS」）、ESCON等、または、いずれかのネットワークの組み合わせを具えることができる。更に、コンピュータネットワーク116は、イーサネット（登録商標）、トークンリング、Wi-Fi、WiMaxなどのような、IEEE802ファミリのネットワーク技術からのネットワークを具えることができる。

【0038】

コンピュータネットワーク116は、サーバ、スイッチ、ルータ、ケーブル、無線、および、コンピュータ112とクライアント114とのネットワーク構築を容易にするのに用いられるその他の設備を具える。一実施例においては、システム100は、コンピュータネットワーク116にわたってピアとして通信する多数のコンピュータ112を具える。別の実施例においては、システム100は、コンピュータ116にわたってピアとして通信する複数のソリッドステートストレージデバイス102を具える。関連分野の当業者は、1又はそれ以上のコンピュータネットワーク116と、1又はそれ以上のクライアント114間の単一または冗長接続を有する関連設備、または、1又はそれ以上のソリッドステートストレージデバイス102を有する他のコンピュータ、または、1又はそれ以上のコンピュータ112に接続された1又はそれ以上のソリッドステートストレージデバイス102とを具える他のコンピュータネットワーク116を認識するであろう。一実施例においては、システム100は、コンピュータ112なしで、コンピュータネットワーク118を通してクライアント116に接続した2又はそれ以上のソリッドステートストレージデバイス102を具える。

10

20

【0039】

ストレージコントローラ管理オブジェクト

【0040】

図1Bは、本発明によるストレージデバイスのオブジェクト管理用のシステム101の一実施例を示す概略ブロック図である。システム101は、各々がストレージコントローラ152および1又はそれ以上のデータストレージデバイス154を有する1又はそれ以上のストレージデバイス150と、1又はそれ以上の要求デバイス155とを具える。ストレージデバイス152は、互いにネットワークで接続され、1又はそれ以上の要求デバイス155に結合される。要求デバイス155は、ストレージデバイス150aにオブジェクト要求を送信する。オブジェクト要求は、オブジェクトを生成する要求、オブジェクトにデータを書き込む要求、オブジェクトからデータを読み出す要求、オブジェクトを削除する要求、オブジェクトを動作確認する要求、オブジェクトをコピーする要求などができる。当分野の当業者は、その他のオブジェクト要求もわかるであろう。

30

【0041】

一実施例においては、ストレージコントローラ152およびデータストレージデバイス154は別個のデバイスである。別の実施例においては、ストレージコントローラ152およびデータストレージデバイス154は、1つのストレージデバイス150に統合される。別の実施例においては、データストレージデバイス154は、ソリッドステートストレージ110であり、ストレージコントローラは、ソリッドステートストレージデバイスコントローラ202である。他の実施例においては、データストレージデバイス154は、ハードディスクドライブ、光学ドライブ、テープストレージ等にできる。別の実施例においては、ストレージデバイス150は、異なる種類の2又はそれ以上のデータストレージデバイス154を具えることができる。

40

【0042】

一実施例においては、データストレージデバイス154は、ソリッドステートストレージ110であり、ソリッドステートストレージエレメント216、218、220のアレイとして配置される。別の実施例においては、ソリッドステートストレージ110は2又はそれ以上のバンク214a-nに配置される。ソリッドステートストレージ110は、

50

図 2 B で以下により詳細に示す。

【 0 0 4 3 】

これらのストレージデバイス 1 5 0 a - n は、互いにネットワークで接続され、分散ストレージデバイスとして機能する。要求デバイス 1 5 5 に接続されたストレージデバイス 1 5 0 a は、分散ストレージデバイスに対するオブジェクト要求を制御する。一実施例においては、ストレージデバイス 1 5 0 および付随するストレージコントローラ 1 5 2 は、オブジェクトを管理し、分散オブジェクトファイルシステムとして要求デバイス 1 5 5 に存在する。この構成において、並行オブジェクトファイルサーバは、分散オブジェクトファイルシステム型の例である。別の実施例においては、ストレージデバイス 1 5 0 および付随のストレージコントローラ 1 5 2 は、オブジェクトを管理し、分散オブジェクトファイルサーバとして要求デバイス 1 5 5 に存在する。この構成において、並行オブジェクトファイルサーバは、分散オブジェクトファイルサーバ型の例である。これらの実施例およびその他の実施例において、要求デバイス 1 5 5 は、オブジェクトを排他的に管理し、または、ストレージデバイス 1 5 0 と併用してオブジェクトを管理することに関与することができ、これは一般的に、ストレージデバイス 1 5 0 の能力を制限せず、他のクライアント 1 1 4 のためのオブジェクトを完全に管理する。縮退型の場合には、各々の分散ストレージデバイス、分散オブジェクトファイルシステムおよび分散オブジェクトファイルサーバは、単一のデバイスとして別個に動作できる。ネットワーク接続されたストレージデバイス 1 5 0 a - n は、分散ストレージデバイス、分散オブジェクトファイルシステム、分散オブジェクトファイルサーバ、および、1 又はそれ以上の要求デバイス 1 5 5 に対して構成された 1 又はそれ以上の上記能力の表現を有するその組み合わせとして、動作できる。例えば、ストレージデバイス 1 5 0 は、第 1 の要求デバイス 1 5 5 a 用の分散ストレージデバイスとして動作するように構成でき、一方、要求デバイス 1 5 5 b 用の分散ストレージデバイスおよび分散オブジェクトファイルシステムとして動作する。システム 1 0 1 が、1 つのストレージデバイス 1 5 0 a を具える場合、ストレージデバイス 1 5 0 a のストレージコントローラ 1 5 2 a は、オブジェクトを管理し、オブジェクトファイルシステムまたはオブジェクトファイルサーバとして要求デバイス 1 5 5 に存在してもよい。

10

20

【 0 0 4 4 】

ストレージデバイス 1 5 0 が、分散ストレージデバイスとして互いにネットワークで接続される一実施例において、ストレージデバイス 1 5 0 は、1 又はそれ以上の分散ストレージコントローラ 1 5 2 によって管理された独立ドライブ冗長アレイ (「 R A I D 」) として動作する。例えば、オブジェクトのデータセグメントを書き込む要求は、R A I D レベルに応じて、データセグメントがパリティストライプを有するデータストレージデバイス 1 5 4 a - n にわたってストライプ化される結果となる。このような配置の利点の 1 つは、単一ストレージデバイス 1 5 0 がストレージコントローラ 1 5 2、データストレージデバイス 1 5 4 またはストレージデバイス 1 5 0 のその他の構成のいずれかに故障を有するときに、このようなオブジェクト管理システムが利用し続けることができることである。

30

【 0 0 4 5 】

冗長ネットワークを用いてストレージデバイス 1 5 0 と要求デバイス 1 5 5 とを相互接続するとき、オブジェクト管理システムは、ネットワークの 1 つが動作している限り、ネットワーク故障の存在下で利用し続けることができる。単一のストレージデバイス 1 5 0 a を有するシステム 1 0 1 は、更に、複数のデータストレージデバイス 1 5 4 a を具えることができ、ストレージデバイス 1 5 0 a のストレージコントローラ 1 5 2 a は、R A I D コントローラとして動作し、ストレージデバイス 1 5 0 a のデータストレージデバイス 1 5 4 a にわたってデータセグメントをストライプ化し、R A I D レベルに応じて、パリティストライプを具えることができる。

40

【 0 0 4 6 】

一実施例においては、1 又はそれ以上のストレージデバイス 1 5 0 a - n は、ソリッドステートストレージデバイスコントローラ 2 0 2 およびソリッドステートストレージ 1 1

50

0を有するソリッドステートストレージデバイス102である場合、ソリッドステートストレージデバイス102は、DIMM構成、ドータカード、マイクロモジュールなどに構成され、コンピュータ112内に常駐できる。コンピュータ112は、サーバ、または、互いにネットワーク接続され、分散RAIDコントローラとして動作するソリッドステートストレージデバイスコントローラ102を有する同様のデバイスでもよい。有利には、ストレージデバイス102は、PCI-e、PCIe-AS、インフィニバンドまたはその他の高性能バス、スイッチングバス、ネットワーク接続バスまたはネットワークを用いて接続されてもよく、ソリッドステートストレージ110a-nにわたってデータセグメントを自律的にストライプ化する単一又は分散ソリッドステートストレージコントローラ202を有する極めて小型の、高性能RAIDストレージシステムを提供することができる。

10

【0047】

一実施例においては、ストレージデバイス150と通信するために要求デバイス155によって用いられる同一ネットワークは、ピアストレージデバイス150b-nと通信しRAID機能を得るために、ピアストレージデバイス150aによって用いられる。別の実施例においては、別個のネットワークは、RAIDするために、ストレージデバイス150間で用いられうる。別の実施例においては、要求デバイス155は、ストレージデバイス150に冗長要求を送信することによってRAID処理に参与することができる。例えば、要求デバイス155は、第1のオブジェクト書き込み要求を第1のストレージデバイス150aに送信し、同一データセグメントを有する第2のオブジェクト書き込み要求

20

【0048】

ストレージデバイス102内でオブジェクト処理するための能力を用いて、ストレージコントローラ152は、1つのRAIDレベルを用いて、1つのデータセグメントまたはオブジェクトを記憶する能力を個別に有し、一方、別のデータセグメントまたはオブジェクトは、異なるRAIDレベルを用いて、または、RAIDストライプをせずに、記憶される。これらの複数のRAIDグループは、ストレージデバイス150内の複数のパーティションに関連していてもよい。RAID0、RAID1、RAID5、RAID6および複合RAIDタイプ10、50、60は、データストレージデバイス154a-nを具

30

【0049】

更に、ストレージコントローラ152は、RAIDコントローラとして自律的に動作するので、RAIDコントローラは、進行型RAIDを行うことができ、1つのRAIDレベルを有するデータストレージデバイス154にわたってストライプ化したオブジェクトまたはオブジェクトの一部を、要求デバイス155が影響を受け、関与し、または、RAIDレベルの変化を検出さえすることなく、別のRAIDレベルに変換可能である。好ましい実施例においては、1つのレベルから別のレベルへとRAID構成を進行させることは、オブジェクトまたはパケットベースでさえも自律的に行われ、ストレージデバイス150の1つまたはストレージコントローラ152において動作する分散RAID制御モジュールによって起動される。一般的に、RAIDの進行は、RAID1などのより高性能でより低効率なストレージ構成から、RAID5などのより低性能でより高い効率のストレージ構成への進行であり、変換は、アクセス頻度に基づいて動的に開始される。なお、RAID5からRAID1へ構成を進行させることも可能であることは分かるであろう。RAID進行を開始するためのその他のプロセスは、ストレージシステム管理サーバ要求など、クライアントまたは外部エージェントから構成または要求されうる。当分野の当業者は、自律的にオブジェクトを管理するストレージコントローラ152を有するストレージデバイス102の特徴および利点分かるであろう。

40

【0050】

50

インサーバSANを用いたソリッドステートストレージデバイス

【0051】

図1Cは、本発明によるインサーバストレージエリアネットワーク（「SAN」）用システム103の一実施例を示した概略ブロック図である。システム103は、通常、サーバ（「サーバ112」）として設定されるコンピュータ112を具える。各サーバ112は、1又はそれ以上のストレージデバイス150を具え、サーバ112およびストレージデバイス150は、共有されたネットワークインタフェース156に接続される。各ストレージデバイス150は、ストレージコントローラ152および関連するデータストレージデバイス154を具える。システム103は、サーバ112の内部または外部にあるクライアント114、114a、114bを具える。クライアント114、114a、114bは、上述したものと実質的に同様の、1又はそれ以上のコンピュータネットワーク116を介して、各サーバ112および各ストレージデバイス150と通信することができる。

10

【0052】

ストレージデバイス150は、DASモジュール158、NASモジュール160、ストレージ通信モジュール162、インサーバSANモジュール164、共通インタフェースモジュール166、プロキシモジュール170、仮想バスモジュール172、フロントエンドRAIDモジュール174およびバックエンドRAIDモジュール176を具え、これらは以下に示す。モジュール158-176はストレージデバイス150内に示される一方、各モジュール158-176の全てまたは一部は、ストレージデバイス150、サーバ112、ストレージコントローラ152またはその他のロケーションにあってもよい。

20

【0053】

インサーバSANと併せて使用されるサーバ112は、サーバとして機能するコンピュータである。このサーバ112は、ファイルサーバ機能など、少なくとも1つのサーバ機能を具えるが、その他のサーバ機能も具えてもよい。サーバ112は、いくつかのファームの一部であってもよく、他のクライアント114に応答する。他の実施例においては、サーバ112は、パーソナルコンピュータ、ワークステーションまたはストレージデバイス150を収容する他のコンピュータでもよい。サーバ112は、ダイレクト接続ストレージ（「DAS」）、SAN接続ストレージ、または、ネットワーク接続ストレージ（「NAS」）としてのサーバ112の1又はそれ以上のストレージデバイス150にアクセスできる。インサーバSANまたはNASに關与するストレージコントローラ150は、サーバ112の内部または外部にあってもよい。

30

【0054】

一実施例においては、インサーバSAN装置は、少なくとも1つのデータストレージデバイス154の少なくとも一部を構成するDASモジュール158を具え、この少なくとも1つのデータストレージデバイス154は、少なくとも1つのクライアント114からサーバ112にストレージ要求を送るように、サーバ112に接続したDASデバイスとしてサーバ112のストレージコントローラ152によって制御される。一実施例においては、第1のデータストレージデバイス154aは、第1のサーバ112に対するDASとして設定され、さらに、第1のサーバ112aに対するインサーバSANストレージデバイスとしても設定される。別の実施例においては、第1のデータストレージデバイス154aは、1つのパーティションは、DASであり、他方が、インサーバSANとなるようにパーティション化される。別の実施例においては、第1のデータストレージデバイス154a内の少なくともストレージスペースの一部は第1のサーバ112aに対するDASとして設定され、第1のデータストレージデバイス154aのストレージスペースの同一部分は、第1のサーバ112aに対するインサーバSANとして設定される。

40

【0055】

別の実施例においては、インサーバSAN装置は、少なくとも1つのクライアント114用のNASデバイスとしてストレージコントローラ152を構成してクライアント11

50

4からのファイル要求に応答する、N A Sモジュール160を具える。さらに、ストレージコントローラ152は、第1のサーバ112a用のインサーバS A Nデバイスとして設定可能である。ストレージデバイス150は、ストレージデバイス150のあるサーバ112から独立して、共有されたネットワークインタフェース155を介して、コンピュータネットワーク116に直接接続できる。

【0056】

一の構成要素形態において、インサーバS A N装置は、第1のサーバ112a内に第1のストレージコントローラ152aを具え、第1のストレージコントローラ152aが少なくとも1つのストレージデバイス154aを制御する。第1のサーバ112aは、第1のサーバ112aおよび第1のストレージコントローラ152aによって共有されるネットワークインタフェース156を具える。インサーバS A N装置は、第1のストレージコントローラ152aと第1のサーバ112aの外部にある少なくとも1つのデバイスとの間の通信を促進するストレージ通信モジュール162を具え、第1のストレージコントローラ152aと前記外部デバイスとの間の通信は、第1のサーバ112aから独立している。ストレージ通信モジュール162によって、第1のストレージコントローラ152aは、外部通信用ネットワークインタフェース156aに独立してアクセス可能になる。一実施例においては、ストレージ通信モジュール162は、ネットワークインタフェース156aのスイッチにアクセスし、第1のストレージコントローラ152aと外部デバイスとの間のネットワークトラフィックを管理(d i r e c t)する。

10

【0057】

さらに、インサーバS A N装置は、ネットワークプロトコルおよびバスプロトコルの一方又は両方を用いてストレージ要求に応答するインサーバS A Nモジュール164を具える。インサーバS A Nモジュール164は、第1のサーバ112aから独立したストレージ要求に応答し、サービス要求は、内部または外部クライアント114、114aから受信する。

20

【0058】

一実施例においては、第1のサーバ112aの外部にあるデバイスは、第2のストレージコントローラ152bである。第2のストレージコントローラは、少なくとも1のデータストレージデバイス154bを制御する。インサーバS A Nモジュール164は、第1のサーバ112aから独立して、第1および第2のストレージコントローラ152a、152b間において、ネットワークインタフェース156aを介した通信を使用して、ストレージ要求に応答する。第2のストレージコントローラ152bは、第2のサーバ112bまたはその他のデバイス内であってもよい。

30

【0059】

別の実施例においては、第1のサーバ112aの外部にあるデバイスは、クライアント114であり、ストレージ要求は、外部クライアント114からきて、第1のストレージコントローラは、S A Nの少なくとも1部として設定され、インサーバS A Nモジュール164は、第1のサーバ112aから独立してネットワークインタフェース156aを介してストレージ要求に応答する。外部クライアント114は、第2のサーバ112b内でもよく、第2サーバ112bの外部内であってもよい。一実施例においては、インサーバS A Nモジュール164は、第1のサーバ112が利用できない場合であっても外部クライアント114からストレージ要求に応答することができる。

40

【0060】

別の実施例においては、ストレージ要求を発するクライアント114aは、第1のサーバ112aの内部にあり、第1のストレージコントローラ152aは、S A Nの少なくとも一部として設定され、インサーバS A Nモジュール164は、1又はそれ以上のネットワークインタフェース156aおよびシステムバスを介してストレージ要求に応答する。

【0061】

従来のS A N設定は、サーバ112から遠隔にあるストレージデバイスが、ダイレクト接続ストレージ(「D A S」)としてサーバ112内にあるようにアクセスでき、ストレ

50

ージデバイスは、ブロックストレージデバイスとしてみえた。一般的に、SANとして接続されたストレージデバイスは、ファイバーチャネル、インターネットスモールコンピュータシステムインタフェース(「iSCSI」)、ハイパーSCSI、ファイバーコネクティビティ(「FICON」)、イーサネット(登録商標)を介したアドバンステクノロジーアタッチメント(「ATA」)など、SANプロトコルを要求する。インサーバSANは、サーバ112内部にストレージコントローラ152を具え、ネットワークプロトコルおよび/またはバスプロトコルを用いて、ストレージコントローラ152aと、リモートストレージコントローラ152bまたは外部クライアント114と、の間のネットワーク接続を可能にしている。

【0062】

一般的に、SANプロトコルは、ネットワークプロトコルのフォームであり、より多くのネットワークプロトコルは、ストレージコントローラ150aおよび関連するデータストレージデバイス154aが、SANとして設定され、外部クライアント114または第2のストレージコントローラ152bと通信可能なInfinibandなどとして存在する。別の実施例においては、第1のストレージコントローラ152aは、イーサネット(登録商標)を使用して、外部クライアント114または第2のストレージコントローラ152bと通信可能である。

【0063】

ストレージコントローラ152は、内部ストレージコントローラ152またはクライアント114aとバスを介して通信できる。例えば、ストレージコントローラ152は、PCIエクスプレス入力/出力ビジュアライゼーション(「PCIe-IOV」)をサポートできるPCI-eを用いてバスを介して通信できる。他の新規(emerging)バスプロトコルによって、システムバスがコンピュータまたはサーバ112の外部に拡張可能となり、ストレージコントローラ152aをSANとして設定可能になる。このようなバスプロトコルの1つは、PCIe-ASである。本発明は、単純なSANプロトコルに限定されないが、ストレージ要求に応答する新規ネットワークおよびバスプロトコルの利点を得ることができる。クライアント114または外部ストレージコントローラ152bの形態の外部デバイスは、拡張システムバスまたはコンピュータネットワーク116を介して通信できる。本明細書で使用されるように、ストレージ要求は、データ書き込み、データ読み取り、データ消去、データクエリなどの要求を具え、オブジェクトデータ、メタデータおよび管理要求ならびにブロックデータ要求を具えることができる。

【0064】

従来サーバ112は、通常、サーバ112内のデバイスへのアクセスを制御するルートコンプレックスを有する。通常、サーバ112のルートコンプレックスは、ネットワークインタフェース156を有し、ネットワークインタフェース156を介したあらゆる通信は、サーバ112によって制御される。しかしながら、インサーバSAN装置の好適な実施例において、ストレージコントローラ152は、独立してネットワークインタフェース156にアクセス可能であり、クライアント114は、SANを形成する第1のサーバ112a内の1又はそれ以上のストレージコントローラ152aと直接通信することができるか、または、1又はそれ以上の第1のストレージコントローラ152aは、第2のストレージコントローラ152bまたはその他のリモートストレージコントローラ152と互いにネットワーク接続されSANを形成する。好適な実施例においては、第1のサーバ112aから遠隔にあるデバイスは、1つの共有されたネットワークアドレスを介して第1のサーバ112aまたは第1のストレージコントローラ152aにアクセスできる。一実施例においては、インサーバSAN装置は、ネットワークインタフェース156、ストレージコントローラ152、およびサーバ112を構成する共通インタフェースモジュール166を具え、サーバ112およびストレージコントローラ152は、共有されたネットワークアドレスを用いてアクセス可能である。

【0065】

他の実施例においては、サーバ112は、2またはそれ以上のネットワークインタフェ

10

20

30

40

50

ース156を具える。例えば、サーバ112は、1つのネットワークインタフェース156を介して通信可能であり、ストレージデバイス150は、別のインタフェースを介して通信できる。別の例においては、サーバ112は、各々がネットワークインタフェースを有する複数のストレージデバイス150を具える。当分野の当業者であれば、1又はそれ以上のストレージデバイス150および1又はそれ以上のネットワークインタフェース156を有するサーバ112のその他の設定であって、1又はそれ以上のストレージデバイス150がサーバ112から独立してネットワークインタフェース156にアクセスすることを理解されたい。当分野の当業者であれば、どのように各種設定が、ネットワーク冗長性をサポートし、有用性を改善するために拡張可能であることを理解されたい。

【0066】

有利なことに、インサーバSAN装置は、複雑性と、従来のSANの費用を大幅に省いている。例えば、一般的なSANは、外部ストレージコントローラ152および関連したデータストレージデバイス154を有するサーバ112を必要とする。このことは、ラックの更なる空間を塞ぎ、ケーブルやスイッチなどを要求する。ケーブル、スイッチング、従来のSANを設定するのに必要なその他のオーバーヘッドは、空間を必要とし、帯域幅を低下させ、費用がかかる。インサーバSAN装置によって、ストレージコントローラ152および関連ストレージ154は、サーバ112フォームファクタにおいて適合し、従って、必要な空間を低減し費用を減少させる。さらに、インサーバSANによって、内部および外部の高速データバスを介して比較的に高速な通信を使用した接続が可能になる。

【0067】

一実施例においては、ストレージデバイス150は、ソリッドステートストレージデバイス102であり、ストレージコントローラ152はソリッドステートストレージコントローラ104であり、データストレージデバイス154はソリッドステートストレージ110である。本明細書に記載するようにソリッドステートストレージデバイス102のスピードにより、この実施例は有利である。さらに、ソリッドステートストレージデバイス102は、サーバ112内で適合可能であり必要とする空間の小さいDIMM内で設定可能である。

【0068】

サーバ112内の1又はそれ以上の内部クライアント114aは、さらに、サーバのネットワークインタフェース156を介してコンピュータネットワーク116に接続可能であり、クライアントの接続は、通常、サーバ112によって制御される。これにはいくつかの利点がある。クライアント114aは、直接ストレージデバイス150に、ローカルまたはリモートでアクセスでき、メモリクライアント114aとストレージデバイス150との間においてローカルまたはリモートダイレクトメモリアクセス(「DMA」、「RDMA」)データ移動を開始できる。

【0069】

別の実施例においては、サーバ112の内部または外部にあるクライアント114、114aは、1又はそれ以上のネットワーク116を介して、クライアント114に対するファイルサーバとして機能し、インサーバSAN、外部SANおよびハイブリッドSANの一部として関与する、ローカルで接続されたソリッドステートストレージデバイス102を使用する。ストレージデバイス150は、DAS、インサーバSAN、SAN、NASなど、および、これらの組み合わせに同時に関与することができる。さらに、各ストレージデバイスをパーティション化することで、第1のパーティションが、DASとしてストレージデバイス150を利用可能にし、第2のパーティションが、インサーバSANの構成要素としてストレージデバイス150を利用可能にし、第3のパーティションがNASとしてストレージパーティション150を利用可能にし、第4のパーティションがSANなどの構成要素としてストレージデバイス150を利用可能にする。同様に、ストレージデバイス150は、安全およびアクセス制御要求に矛盾がないようにパーティション化できる。当分野の当業者であれば、ストレージデバイス、仮想ストレージデバイス、ストレージネットワーク、仮想ストレージネットワーク、プライベートストレージ、共有され

10

20

30

40

50

たストレージ、パラレルファイルシステム、パラレルオブジェクトファイルシステム、ブ
ロックストレージデバイス、オブジェクトストレージデバイス、ストレージアプライア
ンス、ネットワークアプライアンスなどの、多数の組み合わせおよび順列が、構築およびサ
ポートされることを理解されたい。

【 0 0 7 0 】

さらに、コンピュータネットワーク 1 1 6 に直接接続されることによって、ストレージ
デバイス 1 5 0 は、互いに通信可能であり、インサーバ S A N として機能することができ
る。サーバ 1 1 2 内のクライアント 1 1 4 a と、コンピュータネットワーク 1 1 6 を介し
て接続されたクライアント 1 1 4 は、S A N としてのストレージデバイス 1 5 0 にアクセ
スできる。ストレージデバイス 1 5 0 をサーバ 1 1 2 に移動し、S A N としてストレージ
デバイス 1 5 0 を設定する能力を得ることで、サーバ 1 1 2 / ストレージデバイス 1 5 0
の組み合わせは、専用のストレージコントローラ、ファイバチャネルネットワークおよび
他の設備のための従来の S A N の必要性を排除する。インサーバ S A N システム 1 0 3 は
、ストレージデバイス 1 5 0 が電力、冷房、管理および物理的空間などの共通リソースを
クライアント 1 1 4 およびコンピュータ 1 1 2 と共有することができる利点を有する。例
えば、ストレージデバイス 1 5 0 は、サーバ 1 1 2 の空スロットを満たし、S A N または
N A S の実行能力、信頼性および有用性を全てに付与する。当分野の当業者であれば、イ
ンサーバ S A N システム 1 0 3 のその他の特徴および利点を理解されたい。

10

【 0 0 7 1 】

他の設定においては、複数のインサーバ S A N ストレージデバイス 1 5 0 a を 1 つのサ
ーバ 1 1 2 a インフラストラクチャ内に配置することができる。一実施例においては、サ
ーバ 1 1 2 a は、外部ネットワークインタフェース 1 5 6、外部クライアント 1 1 4、1
1 4 b または外部ストレージデバイス 1 5 0 b なしで、P C I エクスプレス I O V を用い
て連結された 1 又はそれ以上の内部ブレード化サーバクライアント 1 1 4 a からなる。

20

【 0 0 7 2 】

さらに、1 又はそれ以上のコンピュータネットワーク 1 1 6 を介して、インサーバ S A
N ストレージデバイス 1 5 0 は、ピア (p e e r) ストレージデバイス 1 5 0 と通信可能
であり、これらのピアストレージデバイス 1 5 0 は、コンピュータ 1 1 2 に位置している
か (図 1 A)、または、S A N およびインサーバ S A N の両方の全ての能力を有するハイ
ブリッド S A N を形成するコンピュータ 1 1 2 なしで、コンピュータネットワーク 1 1 6
に直接接続されている。このフレキシビリティは、各種の可能性のあるソリッドステ
ートストレージネットワークインプリメンテーションの間で拡張性および移動を単純化す
る利点を有する。当分野の当業者であれば、その他の組み合わせ、設定、インプリメン
テーション、およびソリッドステートコントローラ 1 0 4 を配置および相互接続する構
造を理解されたい。

30

【 0 0 7 3 】

ネットワークインタフェース 1 5 6 a は、サーバ 1 1 2 a 内で機能する 1 つのエージェ
ントのみによって制御可能であり、前記エージェント内で機能するリンク設定モジュ
ール 1 6 8 は、ネットワークインタフェース 1 5 6 a から外部ストレージデバイス 1 5 0 b
およびクライアント 1 1 4、1 1 4 b を介して、内部クライアント 1 1 4 a とストレージ
デバイス 1 5 0 a / 第 1 のストレージコントローラ 1 5 2 a との間の通信パスをセットア
ップ可能である。好適な実施例においては、一旦通信パスが構築されると、個々の内
部ストレージデバイス 1 5 0 a および内部クライアント 1 1 4 a は、自身のコマンドキュー
を構築し管理することができ、直接、および、ネットワークインタフェース 1 5 6 a
を制御するプロキシまたはエージェントから独立した R D M A を介して、いずれかの
方向で、ネットワークインタフェース 1 5 6 a から外部ストレージデバイス 1 5 0 b
およびクライアント 1 1 4、1 1 4 b を介し、コマンド及びデータの両方を転送する
ことができる。一実施例においては、リンク設定モジュール 1 6 8 は、ハードウェア
のスタートアップまたは初期化など、初期化プロセス中に通信リンクを構築する。

40

【 0 0 7 4 】

50

別の実施例においては、プロキシモジュール170は、第1のサーバ112aを介してストレージ要求に应答するのに使用されるコマンドの少なくとも1部を指令し、少なくともデータ、および、ストレージ要求と関連したその他の可能性のあるコマンドは、第1のサーバから独立して、第1のコントローラと外部ストレージデバイスとの間で通信される。別の実施例においては、ストレージデバイス150aおよびクライアント114aに変わって、プロキシモジュール170がコマンド又はデータを送る。

【0075】

一実施例においては、第1のサーバ112aは、第1のサーバ112a内の1又はそれ以上のサーバを具え、仮想バスモジュール172を具え、これによって、第1のサーバ112a内の1又はそれ以上のサーバが、別の仮想バスを介して、1又はそれ以上のストレージコントローラ152aに独立してアクセス可能となる。IOVをサポートするネットワークインタフェース156aによって、1又はそれ以上のサーバおよび1又はそれ以上のストレージコントローラは、1又はそれ以上のネットワークインタフェース156aを独立して制御可能になる。

【0076】

各種実施例においては、インサーバSAN装置によって、2又はそれ以上のストレージデバイス150は、RAID内に設定可能である。一実施例においては、インサーバSAN装置は、RAIDとして2又はそれ以上のストレージコントローラ152を構成するフロントエンドRAIDモジュール174を具える。クライアント114、114aからのストレージ要求がデータを記憶する要求を具える場合、フロントエンドRAIDモジュール174は、特定の実装RAIDレベルと一致するRAIDにデータを書き込むことによってストレージ要求に应答する。第2のストレージコントローラ152は、第1のサーバ112a内または第1のサーバ112aの外部に配置可能である。フロントエンドRAIDモジュール174によって、ストレージコントローラ152のRAIDが許容され、ストレージコントローラ152は、ストレージ要求を送るクライアント114、114aに対してビジブル(visible)である。これによって、ストライプおよびパリティ情報は、マスタとして設計されたストレージコントローラ152またはクライアント114、114aによって管理可能となる。

【0077】

別の実施例においては、インサーバSAN装置は、バックエンドRAIDモジュール176を具え、このモジュール176は、RAIDとしてストレージコントローラによって制御される2又はそれ以上のデータストレージデバイス154を設定する。クライアントからのストレージ要求が、データを記憶する要求を具えている場合、バックエンドRAIDモジュール176は、実装されたRAIDレベルと一致するRAIDにデータを書き込むことによってストレージ要求に应答し、RAIDとして設定されたストレージデバイス154は、第1のストレージコントローラ152によって制御される1つのデータストレージデバイス154としてクライアント114、114aによってアクセスされる。このRAID実装によって、RAIDがデータストレージデバイス154にアクセスするあらゆるクライアント114、114aに透明(transparent)であるように、ストレージコントローラ152によって制御されるデータストレージデバイス154のRAIDが可能となる。別の実施例においては、フロントエンドRAIDとバックエンドRAIDの両方が実装され、複数レベルのRAIDを有する。当分野の当業者であれば、本明細書に記載されるように、ソリッドステートストレージコントローラ104と関連するソリッドステートストレージ110と一致するストレージデバイス152をRAIDするその他の方法を理解されたい。

【0078】

ストレージコントローラ管理オブジェクト用装置

【0079】

図2Aは、本発明によるストレージデバイス内のオブジェクト管理のための装置200の実施例を示す概略的ブロック図である。この装置200は、オブジェクト要求受信モジ

ジュール 260、構文解析モジュール 262、コマンド実行モジュール 264、オブジェクトインデックスモジュール 266、オブジェクト要求待ち行列モジュール 268、メッセージモジュール 270 を有するパッケージ 302、および、オブジェクトインデックス復元モジュール 272 を含むストレージコントローラ 152 を具え、これらを以下に示す。

【0080】

ストレージコントローラ 152 は、図 1B のシステム 102 に関連して記載されたストレージコントローラ 152 と実質的に同様で、図 2 と関連して記載されたソリッドステートストレージデバイスコントローラ 202 にできる。この装置 200 は、1 又はそれ以上の要求デバイス 155 からオブジェクト要求を受信するオブジェクト要求受信モジュール 260 を具える。例えば、記憶オブジェクトデータ要求のために、ストレージコントローラ 152 は、ストレージコントローラ 152 に接続したデータストレージデバイス 154 内のデータパケットとしてデータセグメントを記憶する。オブジェクト要求は、一般的に、ストレージコントローラによって管理されたオブジェクト用の 1 又はそれ以上のオブジェクトデータパケットに記憶された、または、記憶されるべきデータセグメントで指令される。オブジェクト要求はオブジェクトを生成し、ストレージコントローラ 152 が、ローカルまたはリモートダイレクトメモリアクセス（「DMA」、「RDMA」）転送を用いることができる、後のオブジェクト要求を通じたデータで後に充填することを要求できる。

10

【0081】

一実施例においては、オブジェクト要求は、既に生成されたオブジェクトにオブジェクトの総てまたは一部を書き込む、書き込み要求である。一例においては、書き込み要求は、オブジェクトのデータセグメント用である。オブジェクトのその他のデータセグメントは、ストレージデバイス 150 またはその他のストレージデバイス 152 に書き込まれてもよい。別の例においては、書き込み要求は、完全なオブジェクト用である。別の例においては、オブジェクト要求は、ストレージコントローラ 152 によって管理されたデータセグメントからのデータを読み出すことである。さらなる別の実施例においては、オブジェクト要求は、データセグメント又はオブジェクトを削除する削除要求である。

20

【0082】

有利には、ストレージコントローラ 152 は、既存のオブジェクトに新規オブジェクトを書き込むか、または、データを加える以上の書き込み要求を受入可能である。例えば、オブジェクト要求受信モジュール 260 が受信する書き込み要求は、ストレージコントローラ 152 によって記録したデータの前にデータを加え、記憶したデータにデータを挿入し、または、データのセグメントを置換する要求を具えることができる。ストレージコントローラ 152 によって維持されるオブジェクトインデックスは、その他のストレージコントローラで利用可能ではないが、サーバや他のコンピュータのファイルシステムのストレージコントローラの外側でのみ現在利用可能である、これらの複雑な書き込み動作に要求されるフレキシビリティを提供する。

30

【0083】

装置 200 は、オブジェクト要求を 1 又はそれ以上のコマンドに構文解析する構文解析モジュール 262 を具える。一般的に、構文解析モジュール 262 は、オブジェクト要求を 1 又はそれ以上のバッファに構文解析する。例えば、オブジェクト要求の 1 又はそれ以上のコマンドは、コマンドバッファに構文解析できる。一般的に、構文解析モジュール 262 は、オブジェクト要求を調整し、オブジェクト要求の情報は、ストレージコントローラ 152 によって理解され実行される。当分野の当業者は、1 又はそれ以上のコマンドにオブジェクト要求を構文解析する構文解析モジュール 262 のその他の機能が分かるであろう。

40

【0084】

装置 200 は、オブジェクト要求から構文解析されたコマンドを実行するコマンド実行モジュール 264 を具える。一実施例においては、コマンド実行モジュール 264 は、1

50

つのコマンドを実行する。別の実施例においては、コマンド実行モジュール264は、複数のコマンドを実行する。一般的に、コマンド実行モジュール264は、書き込みコマンドのようなオブジェクト要求から構文解析されたコマンドを解釈し、次いで、サブコマンドを生成して、待ち行列に入れて実行する。例えば、オブジェクトから構文解析された書き込みコマンドは、ストレージコントローラ152に命令し、複数のデータセグメントを記憶する。更に、オブジェクト要求は、暗号化、圧縮など必要な属性を具備してよい。コマンド実行モジュール264は、ストレージコントローラ152に命令して、データセグメントを圧縮し、データセグメントを暗号化し、1又はそれ以上のデータパケットおよび各データパケット用の関連するヘッダを生成し、メディア暗号化キーでデータパケットを暗号化し、エラー修正コードを加え、そして、特定の位置にデータパケットを記憶することができる。特定の位置およびその他のサブコマンドのデータパケットを記憶することは、その他のより低位のサブコマンドに分解される。当分野の当業者は、コマンド実行モジュール264が、オブジェクト要求から構文解析される1又はそれ以上のコマンドを実行可能であるその他の方法が分かるであろう。

10

20

30

40

50

【0085】

装置200は、オブジェクトインデックスモジュール266を具備し、このモジュール266は、オブジェクトを生成するか、または、オブジェクトのデータセグメントを記憶するストレージコントローラ152に応答して、オブジェクトインデックスのオブジェクトエントリを生成する。一般的に、ストレージコントローラ152は、データセグメントからデータパケットを生成し、データパケットが記憶された位置は、データセグメントが記憶されたときに、割り当てられる。データセグメントとともに、又はオブジェクト要求の一部として受信されたオブジェクトメタデータは、同様の方法で記憶できる。

【0086】

オブジェクトインデックスモジュール266は、データパケットが記憶され、データパケットの物理アドレスが割り当てられるときに、オブジェクトインデックスにオブジェクトエントリを生成する。オブジェクトエントリは、オブジェクトの論理識別子と、ストレージコントローラ152が1又はそれ以上のデータパケットおよびあらゆるオブジェクトメタデータパケットを記憶したところに対応する1又はそれ以上の物理アドレスと、の間におけるマッピングを具備する。別の実施例においては、オブジェクトインデックスのエントリは、オブジェクトのデータパケットが記憶される前に生成される。例えば、データパケットが早期に記憶すべき物理アドレスをストレージコントローラ152が決定する場合、オブジェクトインデックスモジュール266は、オブジェクトインデックス内のエントリを早期に生成できる。

【0087】

一般的には、オブジェクト要求またはオブジェクト要求のグループが、読み出し・変更・書き込み動作中、オブジェクトまたはデータセグメントが変更されることが生じた場合、オブジェクトインデックスモジュール266は、変更したオブジェクトに対応するオブジェクトインデックス内のエントリを更新する。一実施例においては、オブジェクトインデックスは、変更されたオブジェクト用のオブジェクトインデックス中に、新規オブジェクトおよび新規エントリを生成する。一般的に、オブジェクトの一部のみが変更された場合、オブジェクトは、変更されたデータパケットおよび未変化のいくつかのデータパケットを具備する。この場合、新規エントリは、もともと書き込まれていた未変化のデータパケットと、新規位置に書き込まれた変更されたオブジェクトと、に対するマッピングを具備する。

【0088】

別の実施例においては、オブジェクト要求受信モジュール260は、データブロックまたはその他のオブジェクトエレメントを消去するコマンドを具備したオブジェクト要求を受信した場合、ストレージコントローラ152は、オブジェクトに対する基準、オブジェクトとの関係、および、消去したデータブロックのサイズを含む情報を含む消去パケットのような、少なくとも1つのパケットを記憶することができる。更に、それは、更に、消去

オブジェクトエレメントがゼロで充填されていることを示してもよい。従って、消去オブジェクト要求は、消去され、メモリ/ストレージのセル中でゼロを用いて実際に記憶された適宜なメモリ/ストレージの一部を実際に有している、実際のメモリまたはストレージをエミュレートするのに用いることができる。

【0089】

有利には、データセグメントとオブジェクトのメタデータとの間のマッピングを示すエントリを有するオブジェクトインデックスを生成することで、ストレージコントローラ152が、オブジェクトを自律的に処理して管理するのを可能にする。この能力によって、ストレージデバイス150内にデータを記憶する大きなフレキシビリティを可能にする。オブジェクトのインデックスエントリが一度生成されると、オブジェクトに関する次のオブジェクト要求が、ストレージコントローラ152によって効率的に供給される。

10

【0090】

一実施例においては、ストレージコントローラ152は、構文解析モジュール262によって構文解析する前に、オブジェクト要求受信モジュール260によって受信される1又はそれ以上のオブジェクト要求を待ち行列に入れるオブジェクト要求待ち行列モジュールを具える。オブジェクト要求待ち行列モジュール268は、オブジェクト要求が受信されるときと、オブジェクト要求が実行されるときとの間においてフレキシビリティを可能にする。

【0091】

別の実施例においては、ストレージコントローラ152は、1又はそれ以上のデータセグメントから1又はそれ以上のデータパケットを生成するパケットタイザ302を具え、ここで、データパケットがデータストレージデバイス154内のストレージ用にサイズ調整される。パケットタイザ302は、図3でより詳細に以下に示す。パケットタイザ302は、一実施例においては、各パケットのヘッダを生成するメッセージモジュール270を具える。このヘッダは、パケット識別子とパケット長を具える。パケット識別子は、パケットが形成されたオブジェクトにパケットを関連させる。

20

【0092】

一実施例においては、各パケットは、パケット識別子がオブジェクトと、パケット内に含まれるオブジェクトエレメントのオブジェクト内の関係性を識別するのに十分な情報を含むという点で独立型のパケット識別子を含む。しかしながら、より効率的な好ましい実施例は、コンテナにパケットを記憶することである。

30

【0093】

コンテナは、より効率的なパケットのストレージを容易にし、オブジェクトとデータパケットの関係、メタデータパケット、コンテナ内に記憶されるオブジェクトに関連したその他のパケットとを構築するのを助けるデータ構造である。ストレージコントローラ152は、一般的に、同様の方法で、オブジェクトおよびデータセグメントの一部として受信されるオブジェクトメタデータを処理することに留意されたい。一般的に、「パケット」は、データを含むデータパケット、メタデータを含むメタデータパケット、または、別のパケットタイプの別のパケットを意味する。オブジェクトは、1又はそれ以上のコンテナに記憶でき、コンテナは、一般的に、1程度の固有オブジェクト用のパケットを具える。オブジェクトは、複数のコンテナ間に配置される。一般的に、コンテナは、単一の論理消去ブロック(ストレージ部分)内に記憶され、論理消去ブロック間を一般的に分割しない。

40

【0094】

一例においては、コンテナは、2又はそれ以上の論理/仮想ページ間で分割されてもよい。コンテナは、オブジェクトにコンテナを関連させるコンテナラベルによって識別される。コンテナは、多くのパケットに対してゼロを含み、一般的に、コンテナ内のパケットは、1つのオブジェクトからなる。パケットは、オブジェクト属性エレメント、オブジェクトデータエレメント、オブジェクトインデックスエレメントなど、多数のオブジェクトエレメントタイプにできる。2以上のオブジェクトエレメントタイプを含むハイブリッド

50

パケットを生成できる。各パケットは、同一エレメントタイプの多数のエレメントに対してゼロを含む。コンテナ内の各パケットは、一般的に、オブジェクトとの関連性を識別する固有の識別子を含む。

【0095】

各パケットは、1のコンテナに関連している。好ましい実施例においては、コンテナは、消去ブロックに限定されず、各消去ブロックの開始点またはその付近で、コンテナパケットを見つけることができる。このことは、破損したパケットヘッダを有する消去ブロックに対するデータ損失を制限することを支援する。この実施例においては、オブジェクトインデックスが使用不能であり、消去ブロック内のパケットヘッダが破損している場合、パケットヘッダから消去ブロックエンドまでの内容は、次のパケットの位置を決定する信頼性のある機構がない可能性が高いので、損失してしまう。別の実施例においては、より信頼性のある方法は、ページ境界に限定されるコンテナを有することである。この実施例は、より多くのヘッダオーバーヘッドを必要とする。別の実施例においては、コンテナは、ページ及び消去ブロック境界にわたり動くことができる。これは、必要とするヘッダオーバーヘッドは少なくなるが、パケットヘッダが破損している場合、データのより大きな部分が損失される。これらのいくつかの実施例に関して、いくつかのタイプのRAIDが、さらなるデータの完全性を保証するために用いられることが予測される。

10

【0096】

一実施例においては、装置200は、データストレージデバイス154に記憶されるパケットヘッダからの情報を用いて、オブジェクトインデックスのエントリを復元するオブジェクトインデックス復元モジュール272を具える。一実施例においては、オブジェクトインデックス復元モジュール272は、各パケットが属するオブジェクトを判定するヘッダと、オブジェクトにおけるデータまたはメタデータが属する場所を判定する配列情報とを読み出すことによって、オブジェクトインデックスのエントリを復元する。オブジェクトインデックス復元モジュール272は、各パケットおよびタイムスタンプの物理アドレス情報、または、パケットの物理的位置とオブジェクト識別子とデータセグメント配列との間のマッピングを生成する配列情報を用いる。タイムスタンプまたは配列情報はオブジェクトインデックス復元モジュール272によって用いられ、インデックスに対して行われる配列の変化を再生し、これによって、一般的に最新状態を復元する。

20

【0097】

別の実施例においては、オブジェクトインデックス復元モジュール272は、パケットの物理位置と、オブジェクト識別子と、各パケットのシーケンス番号と、を識別して、オブジェクトインデックス内のエントリを復元するコンテナパケット情報とともにパケットヘッダ情報を用いて、パケットを配置する。一実施例においては、消去ブロックは、タイムスタンプされるか、または、パケットが書き込まれるときにシーケンス番号が付与され、タイムスタンプまたは消去ブロックの配列情報は、コンテナヘッダおよびパケットヘッダから集められた情報とともに用いられ、オブジェクトインデックスを復元する。別の実施例においては、消去ブロックが回復されるときに、タイムスタンプまたは配列情報が書き込まれる。

30

【0098】

オブジェクトインデックスが揮発性メモリに記憶されるとき、エラー、電力損失、または、ストレージコントローラ152がオブジェクトインデックスを保存せずにシャットダウンさせる問題が、オブジェクトインデックスが復元できない場合に、問題になるであろう。オブジェクトインデックス復元モジュール272によって、オブジェクトインデックスが揮発性メモリに記憶され、高速アクセスなど揮発性メモリの利点を得ることが可能になる。オブジェクトインデックス復元モジュール272によって、ストレージデバイス150の外部にあるデバイスに依存せず自律的に、オブジェクトインデックスの迅速な復元が可能となる。

40

【0099】

一実施例においては、揮発性メモリのオブジェクトインデックスは、データストレージ

50

デバイス 154 に一時的に記憶される。特定の実施例においては、オブジェクトインデックスまたは「インデックスメタデータ」は、ソリッドステートストレージ 110 に定期的に記憶される。別の実施例においては、インデックスメタデータが、パケットを記憶するソリッドステートストレージ 110 a - 110 n - 1 と別個のソリッドステートストレージ 110 n に記憶される。インデックスメタデータは、データおよび要求デバイス 155 から送信されるオブジェクトメタデータから別個に管理され、ストレージコントローラ 152 / ソリッドステートストレージデバイスコントローラ 202 によって管理される。オブジェクトからの他のデータおよびメタデータと別個にインデックスを管理して記憶することによって、ストレージコントローラ 152 / ソリッドステートストレージデバイスコントローラ 202 がオブジェクトメタデータを不必要に処理することなく、効率的データフローが可能となる。

10

【0100】

一実施例においては、オブジェクト要求受信モジュール 260 によって受信されるオブジェクト要求が書き込み要求を具える場合、ストレージコントローラ 152 は、ローカルまたはリモートダイレクトメモリアクセス（「DMA」、「RDMA」）動作として、要求デバイス 155 のメモリからオブジェクトの 1 又はそれ以上のデータセグメントを受信する。好適な例において、ストレージコントローラ 152 は、1 又はそれ以上の DMA または RDMA 動作において要求デバイス 155 のメモリからデータを引き出す。別の例においては、要求デバイス 155 は、1 又はそれ以上の DMA または RDMA 動作において、データセグメントをストレージコントローラ 152 へ押し出す。別の実施例においては、オブジェクト要求は読み出し要求を具え、ストレージコントローラ 152 は、1 又はそれ以上の DMA または RDMA 動作において、オブジェクトの 1 又はそれ以上のデータセグメントを要求デバイス 155 のメモリに送信する。好適な例においては、ストレージコントローラ 152 は、1 又はそれ以上の DMA または RDMA 動作において、データを要求デバイス 155 のメモリに押し出す。別の実施例においては、要求デバイス 155 は、1 又はそれ以上の DMA または RDMA 動作において、ストレージコントローラ 152 からデータを引き出す。別の例においては、ストレージコントローラ 152 は、1 又はそれ以上の DMA または RDMA 動作において、要求デバイス 155 のメモリからオブジェクトコマンド要求セットを引き出す。別の実施例においては、要求デバイス 155 は、1 又はそれ以上の DMA または RDMA 動作において、ストレージコントローラ 152 にオブジェクトコマンド要求セットを押し出す。

20

30

【0101】

一実施例においては、ストレージコントローラ 152 は、ブロックストレージをエミュレートし、要求デバイス 155 とストレージコントローラ 152 との間で通信されるオブジェクトは、1 又はそれ以上のデータブロックを具える。一実施例においては、要求デバイス 155 は、ドライバを具え、ストレージデバイス 150 はブロックストレージデバイスとして存在する。例えば、要求デバイス 152 は、要求デバイス 155 が記憶されたデータブロックを所望する位置の物理アドレスとともに所定サイズのデータのブロックを送信することができる。ストレージコントローラ 152 は、データブロックを受信し、データブロックで送信される物理ブロックアドレスを用いるか、または、オブジェクト識別子としての物理ブロックアドレスの変換を用いる。次いで、ストレージコントローラ 152 は、データブロックをパケット化してデータブロックを任意に記憶することによって、オブジェクトまたはオブジェクトのデータセグメントとしてデータブロックを記憶する。次いで、オブジェクトインデックスモジュール 266 は、物理的ブロックベースのオブジェクト識別子と、ストレージコントローラ 152 がデータブロックのデータを具えるデータパケットを記憶した実際の物理位置とを用いて、オブジェクトインデックス内にエントリを生成する。

40

【0102】

別の実施例においては、ストレージコントローラ 152 は、ブロックオブジェクトを受け入れることでブロックストレージをエミュレートする。ブロックオブジェクトは、プロ

50

ック構造中に1又はそれ以上のデータブロックを具えることができる。一実施例においては、ストレージコントローラ152は、その他のオブジェクトとしてブロックオブジェクトを処理する。別の実施例においては、オブジェクトは、完全なブロックデバイス、ブロックデバイスのパーティション、又は、トラック、セクタ、チャネルを含むブロックデバイスの他の論理若しくは物理サブエレメントを表す。特筆すべきことは、ブロックデバイスRAIDグループを、進行型RAIDなどの異なるRAID構造をサポートするオブジェクトに再マッピングする能力である。当分野の当業者は、従来または将来のブロックデバイスをオブジェクトへの他のマッピングが分かるであろう。

【0103】

ソリッドステートストレージデバイス

10

【0104】

図2Bは、本発明によるソリッドステートデバイス102内に書き込みデータパイプライン106および読み出しデータパイプライン108を具えるソリッドステートデバイスコントローラ202の一実施例201を示す概略ブロック図である。ソリッドステートストレージデバイスコントローラ202は、各々がソリッドステートストレージ110を制御する多数のソリッドステートストレージコントローラ0-N104a-nを具えることができる。示した実施例においては、2つのソリッドステートコントローラが示され、ソリッドステートコントローラ0104aおよびソリッドステートストレージコントローラN104nは、それぞれ、ソリッドステートストレージ110a-nを制御する。示された実施例においては、ソリッドステートストレージコントローラ0104aは、データチャネルを制御し、接続したソリッドステートストレージ110aはデータを記憶する。ソリッドステートストレージコントローラN104nは、記憶したデータと関連するインデックスメタデータチャネルを制御し、関連するソリッドステートストレージ110nはインデックスメタデータを記憶する。代替の実施例においては、ソリッドステートストレージデバイスコントローラ202は、単一のソリッドステートストレージ110aを有する単一のソリッドステートコントローラ104aを具える。別の実施例においては、複数のソリッドステートストレージコントローラ104a-nおよび関連するソリッドステートストレージ110a-nがある。一実施例においては、付随のソリッドステートストレージ110a-110n-1に接続した1又はそれ以上のソリッドステートコントローラ104a-104n-1は、データを制御し、付随のソリッドステートストレージ110nに接続したソリッドステートストレージコントローラ104nは、インデックスメタデータを制御する。

20

30

【0105】

一実施例においては、少なくとも1つのソリッドステートコントローラ104は、フィールドプログラマブルゲートアレイ(「FPGA」)であり、コントローラ機能は、FPGAにプログラムされる。特定の実施例においては、FPGAはXilinx(登録商標)FPGAである。別の実施例においては、ソリッドステートストレージコントローラ104は、特定用途向け集積回路(「ASIC」)またはカスタム論理回路ソリューションのような特別に設計された構成を具える。各ソリッドステートストレージコントローラ104は、一般的に、書き込みデータパイプライン106および読み出しデータパイプライン108を具えており、これらは、図3で更に説明する。別の実施例においては、少なくとも1つのソリッドステートストレージコントローラ104は、結合型FPGA、ASICおよびカスタム論理回路構成から構成される。

40

【0106】

ソリッドステートストレージ

【0107】

ソリッドステートストレージ110は、バンク214内で構成され、双方向ストレージ入力/出力(「I/O」)バス210を通して並列にアクセスできる、非揮発性ソリッドステートストレージエレメント216、218、220の阵列である。一実施例においては、ストレージI/Oバス210は、どの場合でも単方向通信が可能である。例えば、

50

データがソリッドステートストレージ 110 に書き込まれるとき、データは、ソリッドステートストレージ 110 から読み出すことはできない。別の実施例においては、データは、双方向に同時に流れうる。なお、データバスに関して本明細書で用いられるように、双方向が意味するのは、一度に一方方向のみにしかデータフローを有さないが、双方向データバス上で一方方向のデータフローが停止したときは、データは、双方向データバス上の逆方向にデータをフローすることができるデータ経路である。

【0108】

ソリッドステートストレージエレメント（例えば、SSS 0.0 216a）は、一般的に回路基板上のチップ（1又はそれ以上のダイのパッケージ）またはダイとして構成される。示されているように、これらの複数のエレメントが、チップパッケージ、チップパッケージのスタックまたはその他のパッケージエレメントと一緒にパッケージされている場合であっても、ソリッドステートストレージエレメント（例えば、216a）は、他のソリッドステートストレージエレメント（例えば、218a）から独立または準独立的に動作する。示されているように、ソリッドステートストレージ 216、218、220 の段は、バンク 214 として指定される。示されているように、ソリッドステートストレージ 110 内の $n \times m$ のソリッドステートストレージエレメント 216、218、220 のアレイにおいて、バンクあたり「 n 」バンク 214 a - n および「 m 」ソリッドステートストレージエレメント 216 a - m 、218 a - m 、220 a - m が存在していてもよい。一実施例においては、ソリッドステートストレージ 110 a は、8つのバンク 214 を有するバンク 214 あたり 20 のソリッドステートストレージエレメント 216、218、220 を具え、ソリッドステートストレージ 110 n は、1つのバンク 214 を有するバンク 214 あたり、2つのソリッドステートストレージエレメント 216、218 を具える。一実施例においては、各ソリッドステートストレージエレメント 216、218、220 は、シングルレベルセル（「SLC」）デバイスからなっている。別の実施例においては、各ソリッドステートストレージエレメント 216、218、220 は、マルチレベルセル（「MLC」）デバイスからなっている。

【0109】

一実施例においては、共通ストレージ I/O バス 210 a 列（例えば、216b、218b、220b）を共有する複数のバンク用のソリッドステートストレージエレメントは、一緒にパッケージされる。一実施例においては、ソリッドステートストレージエレメント 216、218、220 は、チップあたり 1 又はそれ以上のダイを有することができ、1 又はそれ以上のチップは、垂直方向にスタックされ、各ダイは別個にアクセス可能である。別の実施例においては、ソリッドステートストレージエレメント（例えば、SSS 0.0 216a）は、ダイあたり 1 又はそれ以上の仮想ダイと、チップあたり 1 又はそれ以上のダイと、垂直方向にスタックされた 1 又はそれ以上のチップと、を有し、各仮想ダイは別個にアクセス可能である。別の実施例においては、ソリッドステートストレージエレメント SSS 0.0 216a は、ダイあたり 1 又はそれ以上の仮想ダイと、チップあたり 1 又はそれ以上のダイと、を有し、1 又はそれ以上のダイの一部または全部が垂直方向にスタックされ、各仮想ダイは別個にアクセス可能である。

【0110】

一実施例においては、2つのダイが、グループあたり 4つのスタックとともに垂直方向にスタックされ、8つのストレージエレメント（例えば、SSS 0.0 - SSS 0.8）216a - 220a を形成し、各々は、別のバンク 214 a - n にある。別の実施例においては、20のストレージエレメント（例えば、SSS 0.0 - SSS 20.0）216 は、仮想バンク 214 a を形成し、8の仮想バンクの各々は、20のストレージエレメント（例えば、SSS 0.0 - SSS 20.8）216、218、220 を有する。データは、ストレージエレメント（SSS 0.0 - SSS 0.8）216a、218a、220a の特定のグループの総てのストレージエレメントへのストレージ I/O バス 210 にわたって、ソリッドステートストレージ 110 に送信される。ストレージ制御バス 212a は、特定のバンク（例えば、バンク - 0 214a）を選択するために用いられ、総ての

10

20

30

40

50

バンク 2 1 4 に接続されたストレージ I / O バス 2 1 0 にわたって受信されるデータは、選択されたバンク 2 1 4 a に書き込まれる。

【 0 1 1 1 】

好ましい実施例においては、ストレージ I / O バス 2 1 0 は、1 又はそれ以上の独立 I / O バスからなり (2 1 0 a . a - m、2 1 0 n . a - m を具える「 I I O B a - m 」)、各列内のソリッドステートストレージエレメントは、各ソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 に並行してアクセスする独立 I / O バスの 1 つを共有するので、総てのバンク 2 1 4 は、同時にアクセスされる。例えば、ストレージ I / O バス 2 1 0 の 1 つのチャンネルは、各バンク 2 1 4 a - n の第 1 のソリッドステートストレージエレメント 2 1 6 a、2 1 8 a、2 2 0 a に同時にアクセスすることができる。ストレージ I / O バス 2 1 0 の第 2 のチャンネルは、各バンク 2 1 4 a - n の第 2 のソリッドステートストレージエレメント 2 1 6 b、2 1 8 b、2 2 0 b に同時にアクセスできる。ソリッドステートストレージ 2 1 6、2 1 8、2 2 0 の各列は、同時にアクセスされる。一実施例においては、ソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 がマルチレベルで (物理的にスタックされている) ある場合、ソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 の総ての物理レベルは、同時にアクセスされる。本明細書で用いられるように、「同時に」は、ほぼ同時のアクセスを含み、デバイスは、スイッチングノイズを回避するために、わずかに異なる間隔でアクセスされる。「同時に」はこのような構成で用いられて、連続または直列アクセスと区別し、コマンドおよび / またはデータは、交互に別個に送信される。

10

20

【 0 1 1 2 】

一般的に、バンク 2 1 4 a - n は、ストレージ制御バス 2 1 2 を用いて別個に選択される。一実施例においては、バンク 2 1 4 は、使用可能チップまたは選択チップを用いて選択される。選択チップおよび使用可能チップの両方が使用可能な場合、ストレージ制御バス 2 1 2 は、マルチレベルのソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 の 1 つのレベルを選択することができる。その他の実施例においては、ストレージ制御バス 2 1 2 によって他のコマンドが用いられて、マルチレベルのソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 の 1 つのレベルを別個に選択する。ソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 は、ストレージ I / O バス 2 1 0 およびストレージ制御バス 2 1 2 上で送信される制御及びアドレス情報との組み合わせを通じて選択されてもよい。

30

【 0 1 1 3 】

一実施例においては、各ソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 は、消去ブロックに分割され、各消去ブロックは、ページに分割される。一般的なページは、2 0 0 0 バイト (「 2 k B 」) である。一例においては、ソリッドステートストレージエレメント (例えば、S S S 0 . 0) は、2 つのレジスタを具え、2 つのページをプログラムして、2 つのレジスタのソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 が、4 k B の容量を有することができる。更に、2 0 のソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 のバンク 2 1 4 は、ストレージ I / O バス 2 1 0 のチャンネルを通る同一アドレスでアクセスされる 8 0 k B のページ容量を有する。

40

【 0 1 1 4 】

8 0 k B のソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 のバンク 2 1 4 における上記ページのグループは、仮想ページと呼ばれる。同様にバンク 2 1 4 a の各ストレージエレメント 2 1 6 a - m の消去ブロックは、仮想消去ブロックを形成するようにグループ化されてもよい。好ましい実施例においては、ソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 内のページの消去ブロックは、消去コマンドがソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 内で受信されたときに消去される。消去ブロック、ページ、プレーンまたはソリッドステートストレージエレメント 2 1 6、2 1 8、2 2 0 内のその他の論理および物理部分のサイズおよび数が、技術の進歩に伴い、時宜にわたって変化することが予想されるが、新規構成に一致する多数の実施例が

50

、可能であり、本明細書の一般的説明と一致すると分かるであろう。

【0115】

一般的に、パケットがソリッドステートストレージエレメント216、218、220内の特定位置に書き込まれ、このパケットが特定バンクの特定エレメントの特定消去ブロックに特有である、特定ページ内の位置に書き込まれるように意図されている場合、物理アドレスは、ストレージI/Oバス210上で送信され、パケットがそれに続く。物理アドレスは、パケットをページ内の指定位置へと配向するための、ソリッドステートストレージエレメント216、218、220の情報を十分に保有している。ストレージエレメント(例えば、SSS0.0 - SSS0.N 216a、218a、220a)の列における総てのストレージエレメントは、ストレージI/Oバス210a . a内の適宜なバスによって同時にアクセスされ、適宜なページに到達し、ストレージエレメント(SSS0.0 - SSS0.N 216a、218a、220a)の列において同様にアドレスされたページに、データパケットを書き込むことを回避するため、データパケットが書き込まれるべき正規のページを有するソリッドステートストレージエレメントSSS0.0 216aを具えるバンク214aは、ストレージ制御バス212によって同時に選択される。

10

【0116】

同様に、ストレージI/Oバス212上を移動する読み出しコマンドは、単一のバンク214aと、このバンク214a内の適宜なページとを選択するために、ストレージ制御バス212上の同時コマンドを要求する。好ましい実施例においては、読み出しコマンドはページ全体を読み出し、バンク214においては並行して複数のソリッドステートストレージエレメント216、218、220があるため仮想ページ全体は読み出しコマンドを用いて読み出される。しかしながら、読み出しコマンドは、バンクインタリーブに関して以下で説明するように、サブコマンドに分割可能である。更に、仮想ページは、書き込み動作でアクセス可能である。

20

【0117】

消去ブロック消去コマンドは、特定の消去ブロックを消去する特定の消去ブロックアドレスを用いて、ストレージI/Oバス210にわたって、消去ブロックを消去するように送信できる。一般的に、消去ブロック消去コマンドは、仮想消去ブロックを消去するために、ストレージI/Oバス210の並列な経路にわたり送信でき、各々が特定の消去ブロックを消去する特定の消去ブロックアドレスを有する。同時に、特定バンク(例えば、バンク-0 214a)がストレージ制御バス212で選択され、総てのバンク(バンク1 - N 214b - n)における同様にアドレスされたブロックが消去されるのを防止する。更に、その他のコマンドは、ストレージI/Oバス210とストレージ制御バス212との組み合わせを用いて特定位置に送信できる。当分野の当業者は、双方向ストレージI/Oバス210およびストレージ制御バス212を用いて、特定のストレージ位置を選択できるその他の方法が分かるであろう。

30

【0118】

一実施例においては、パケットは、ソリッドステートストレージ110に連続に書き込まれる。例えば、パケットは、ストレージエレメント216のバンク214aのストレージ書き込みバッファへとストリーミングされ、バッファが一杯の場合、パケットは、指定の仮想ページにプログラムされる。次いで、パケットは、ストレージ書き込みバッファを補充し、一杯の場合、パケットは、次の仮想ページに書き込まれる。次の仮想ページは、同一バンク214aまたは別のバンク(例えば、214b)内にあってもよい。このプロセスは、一般的に、仮想消去ブロックが満たされるまで、仮想ページごとに繰り返される。別の実施例においては、ストリーミングは、仮想消去ブロック境界にわたって継続し、このプロセスは、仮想消去ブロックごとに継続する。

40

【0119】

読み出し、変更、書き込み動作において、オブジェクトに関連するデータパケットは、配置され、読み出し動作で読み出される。変更されたオブジェクトのデータセグメントは

50

、データセグメントが読み出される位置に書き込まれない。その代わりに、変更されたデータセグメントは、データパケットに再度変換され、次いで、そのとき書き込まれている仮想ページの次の利用可能位置に書き込まれる。それぞれのデータパケットのオブジェクトインデックスエントリは、変更されたデータセグメントを含むパケットに対し指示するように変更される。変更されていない同一オブジェクトに関連したデータパケットのオブジェクトインデックスの1又はそれ以上のエントリは、未変更データパケットの元の位置に対するポインタを具える。従って、例えば、オブジェクトの前バージョンを維持するために、元のオブジェクトを維持している場合、元のオブジェクトは、もともと書き込まれていた総てのデータパケットに対するオブジェクトインデックスのポインタを有する。この新規オブジェクトは、元のデータパケットに対するオブジェクトインデックスのポインタと、その時に書き込まれている仮想ページにおける変更データパケットに対するポインタとを有する。

10

【0120】

コピー動作において、オブジェクトインデックスは、ソリッドステートストレージ110に記憶された多数のパケットにマッピングされた元のオブジェクトのエントリを具える。コピーが生成されるとき、新規オブジェクトが生成され、新規エントリが、元のパケットに新規オブジェクトをマッピングするオブジェクトインデックスに生成される。新規オブジェクトは、ソリッドステートストレージ110に書き込まれ、その位置は、オブジェクトインデックスの新規エントリにマッピングされる。新規オブジェクトパケットは、コピーに伝搬されていない元のオブジェクトで変更が生じ、オブジェクトインデックスが損失または破損した場合に、参照される元のオブジェクト内のパケットを識別するために用いられる。

20

【0121】

有益なことに、パケットを連続に書き込むことは、ソリッドステートストレージ110のより均等な利用を容易にし、ソリッドストレージデバイスコントローラ202がソリッドステートストレージ110内のストレージホットスポットおよび様々な仮想ページの利用レベルをモニタすることを可能にする。更に、パケットを連続に書き込むことは、下記で詳細に示すように、強力で有効なガベージコレクションシステムを容易にする。当分野の当業者は、データパケットの連続ストレージのその他の利点がかかるであろう。

【0122】

ソリッドステートストレージデバイスコントローラ

30

【0123】

様々な実施例においては、ソリッドステートストレージデバイスコントローラ202は、データバス204、ローカルバス206、バッファコントローラ208、バッファ0-N 222 a-n、マスタコントローラ224、ダイレクトメモリアクセス(「DMA」)コントローラ226、メモリコントローラ228、ダイナミックメモリアレイ230、スタティックランダムメモリアレイ232、管理コントローラ234、管理バス236、システムバス240へのブリッジ238、および、種々雑多な論理回路242を更に具え、以下にこれを示す。他の実施例においては、システムバス240は、1又はそれ以上のネットワークインタフェースカード(「NICs」)244に結合し、そのうちのいくつかは、リモートDMA(「RDMA」)コントローラ246、1又はそれ以上の中央処理装置(「CPU」)248、1又はそれ以上の外部メモリコントローラ250および関連する外部メモリアレイ252、1又はそれ以上のストレージコントローラ254、ピアコントローラ256およびアプリケーション特有のプロセッサ258を具え、これらは以下で説明する。システムバス240に接続された構成244-258は、コンピュータ112に配置されてもよく、その他のデバイスであってもよい。

40

【0124】

一般的に、ソリッドステートストレージコントローラ104は、ストレージI/Oバス210を介してソリッドステートストレージ110とデータを通信する。ソリッドステートストレージが、バンク214内に配列され、各々のバンク214が、並列にアクセスさ

50

れる複数のストレージエレメント 216、218、220 を具える一般的な実施例においては、ストレージ I/O バス 210 は、バスのアレイであり、ストレージエレメント 216、218、220 の各列の 1 つは、バンク 214 に及んでいる。本明細書で用いられるように、用語「ストレージ I/O バス」は、1 つのストレージ I/O バス 210、または、データ独立性バス 204 のアレイを意味する。好ましい実施例においては、ストレージエレメントの列（例えば、216 a、218 a、220 a）にアクセスする各ストレージ I/O バス 210 は、ストレージエレメントの列 216 a、218 a、220 a にアクセスされるストレージ部分（例えば、消去ブロック）の論理 - 物理マッピングを具えることができる。このマッピングによって、第 1 のストレージ部分が故障している場合、部分的に故障している場合、アクセス不能な場合、または、その他何らかの問題がある場合に、ストレージ部分の物理アドレスにマッピングされる論理アドレスが、異なるストレージ部分に再マッピングできる。再マッピングは、図 3 の再マッピングモジュール 314 に関連して更に説明する。

10

【0125】

データは、システムバス 240、ブリッジ 238、ローカルバス 206、バッファ 22 を介して、また、最終的にデータバス 204 を介して、要求デバイス 155 からソリッドステートストレージコントローラ 104 と通信可能である。データバス 204 は、一般的に、バッファコントローラ 208 を用いて制御される 1 又はそれ以上のバッファ 222 a - n に接続される。バッファコントローラ 208 は一般的に、ローカルバス 206 からバッファ 222 への、並びに、データバス 204 を通ってパイプライン入力バッファ 306 および出力バッファ 330 への、データの転送を制御する。一般的に、バッファコントローラ 222 は、どのように、要求デバイスからきたデータがバッファ 222 に一時的に記憶可能で、次いで、データバス 204 上を転送されるかを制御でき、あるいは、その逆も同様に制御でき、異なるクロックドメインにし、データの衝突を防ぐこと等ができる。バッファコントローラ 208 は、一般的にマスタコントローラ 224 とともに作動し、データフローに調和させる。データが到着するときは、データは、システムバス 240 の上に到着し、ブリッジ 238 を介してローカルバス 206 に転送される。

20

【0126】

一般的に、データは、マスタコントローラ 224 およびバッファコントローラ 208 によって指示されるように、ローカルバス 206 から 1 又はそれ以上のデータバッファ 222 へと転送される。次いで、データは、ソリッドステートコントローラ 104 を介してバッファ 222 からデータバス 204 へと、及び、NAND フラッシュまたはその他のストレージメディアなどのようなソリッドステートストレージ 110 へと流れる。好ましい実施例においては、データと、このデータとともにくる関連したバンド外のメタデータ（「オブジェクトメタデータ」）は、1 又はそれ以上のソリッドステートストレージコントローラ 104 a - 104 n - 1 および関連するソリッドステートストレージ 110 a - 110 n - 1 を具える 1 又はそれ以上のデータチャンネルを用いて通信され、一方、少なくとも 1 つのチャンネル（ソリッドステートストレージコントローラ 104 n、ソリッドステートストレージ 110 n）は、インデックス情報や、ソリッドステートストレージデバイス 102 に内部生成されるその他のメタデータなどのような、バンド内のメタデータ専用である。

30

40

【0127】

一般的に、ローカルバス 206 は、双方向バスであるか、または、ソリッドステートストレージデバイスコントローラ 202 の内部のデバイス間、および、ソリッドステートストレージデバイス 102 の内部のデバイスと、システムバス 240 に接続したデバイス 244 - 258 との間において、データおよびコマンドの通信を可能にするバスのセットである。ブリッジ 238 は、ローカルバス 206 とシステムバス 240 との間の通信を容易にする。当分野の当業者は、リング構造またはスイッチ式スター構成、ならびに、バス 240、206、204、210 およびブリッジ 238 の機能などのような、その他の実施例が分かるであろう。

50

【0128】

一般的に、システムバス240は、コンピュータ112のバスであり、または、ソリッドステートストレージデバイス102がインストールまたは接続されているその他のデバイスである。一実施例においては、システムバス240は、PCI-eバス、シリアルアドバンステクノロジーアタッチメント(「シリアルATA」)バス、パラレルATAなどでもよい。別の実施例においては、システムバス240は、小型コンピュータシステムインタフェース(「SCSI」)、FireWire、ファイバチャネル、USB、PCIe-A/Sなどのような、外部バスである。ソリッドステートストレージデバイス102は、デバイスの内側に適用するように、または、外側に接続されたデバイスとして、パッケージされてもよい。

10

【0129】

ソリッドステートストレージデバイスコントローラ202は、ソリッドステートストレージ102内のより高位の機能を制御するマスタコントローラ224を具える。マスタコントローラ224は、様々な実施例において、オブジェクト要求およびその他の要求を解釈することによってデータフローを制御し、インデックスの生成を命令し、データに関連したオブジェクト識別子を、関連するデータ、協調するDMA要求等の物理的な位置にマッピングする。本明細書に記載される機能の多くは、マスタコントローラ224によって完全に、又は部分的に制御される。

【0130】

一実施例においては、マスタコントローラ224は、埋込式コントローラを用いる。別の実施例においては、マスタコントローラ224は、ダイナミックメモリアレイ230(ダイナミックランダムアクセスメモリ「DRAM」)、スタティックメモリアレイ323(スタティックランダムアクセスメモリ「SRAM」)など、ローカルメモリを用いる。一実施例においては、ローカルメモリは、マスタコントローラ224を用いて制御される。別の実施例においては、マスタコントローラは、メモリコントローラ228を介してローカルメモリにアクセスする。別の実施例においては、マスタコントローラは、Linuxサーバを運用し、WorldWideWeb、ハイパーテキストマークアップ言語(「HTML」)などのような、様々な共通サーバインタフェースをサポートすることができる。別の実施例においては、マスタコントローラ224は、ナノプロセッサを用いる。マスタコントローラ224は、プログラマブルな若しくは標準の論理回路、または、上記に挙げたコントローラタイプの組み合わせ、を用いて構成されていてもよい。当分野の当業者は、マスタコントローラの多数の実施例が分かるであろう。

20

30

【0131】

一実施例においては、ストレージデバイス152/ソリッドステートストレージデバイスコントローラ202は、複数のデータストレージデバイス/ソリッドステートストレージ110a-nを管理する場合、マスタコントローラ224は、ソリッドステートストレージコントローラ104a-nのような、内部コントローラ中間の作業負荷を分割する。例えば、マスタコントローラ224は、データストレージデバイス(例えば、ソリッドステートストレージ110a-n)に書き込まれるオブジェクトを分割し、オブジェクトの一部は、接続したデータストレージデバイスの各々に記憶できる。この特徴は、オブジェクトへのより迅速なストレージおよびアクセスを可能にする性能向上である。一実施例においては、マスタコントローラ224は、FPGAを用いて実装される。別の実施例においては、マスタコントローラ224内のファームウェアは、管理バス236、NIC244に接続したネットワークを介したシステムバス240、または、システムバス240に接続したその他のデバイスを介して更新可能である。

40

【0132】

一実施例においては、オブジェクトを管理するマスタコントローラ224は、ブロックストレージをエミュレートし、コンピュータ102、または、ストレージデバイス152/ソリッドステートストレージ102に接続したその他のデバイスは、ストレージデバイス152/ソリッドステートストレージデバイス102をブロックストレージデバイスと

50

みなし、ストレージデバイス152/ソリッドステートストレージデバイス102内の特定の物理アドレスにデータを送信する。次いで、マスタコントローラ224は、ブロックを分割し、それをオブジェクトとしてデータブロックに記憶する。更に、マスタコントローラ224は、ブロックと、このブロックとともに送信される物理アドレスとを、マスタコントローラ224によって決定される実際の位置にマッピングする。このマッピングは、オブジェクトインデックスに記憶される。一般的に、ブロックエミュレーションのために、ブロックデバイスアプリケーションプログラムインタフェース(「API」)は、コンピュータ112、クライアント114、または、ブロックストレージデバイスとしてストレージデバイス152/ソリッドステートストレージデバイス102を用いることを望むその他のデバイス内のドライバに提供される。

10

【0133】

別の実施例においては、マスタコントローラ224は、NICコントローラ244および埋込型RDMAコントローラ246で調整し、ジャストインタイムのデータおよびコマンドセットのRDMA転送を配信する。NICコントローラ244は、非透過性のポートの背後に隠れ、カスタムドライバを用いることができる。更に、クライアント114のドライバは、標準スタックAPIを使用し、NIC244とともに機能するI/Oメモリドライバを介してコンピュータネットワーク118にアクセスできる。

【0134】

一実施例においては、マスタコントローラ224も、独立ドライブ冗長アレイ(「RAID」)コントローラである。データストレージデバイス/ソリッドステートストレージデバイス102が、1又はそれ以上の他のデータストレージデバイス/ソリッドステートストレージデバイス102とネットワーク接続される場合、マスタコントローラ224は、単階層型RAID、多階層型RAID、進行型RAID等用のRAIDコントローラでもよい。更に、マスタコントローラ224によって、いくつかのオブジェクトはRAIDアレイに記憶でき、他のオブジェクトはRAIDなしに記憶できる。別の実施例においては、マスタコントローラ224は、分散型RAIDコントローラエレメントでもよい。別の実施例においては、マスタコントローラ224は、多数のRAID、分散型RAID、他に記載されるその他の機能をもつことができる。

20

【0135】

一実施例においては、マスタコントローラ224は、単一または冗長ネットワークマネージャ(例えば、スイッチ)を調整して、ルーティングを構築し、帯域幅使用率、フェイルオーバー等のバランスをとる。別の実施例においては、マスタコントローラ224は、統合型アプリケーション特有の論理回路(ローカルバス206を介して)と、関連するドライバソフトウェアを調整する。別の実施例においては、マスタコントローラ224は、接続されたアプリケーション特有のプロセッサ258または論理回路(外部システムバス240)および関連するドライバソフトウェアを調整する。別の実施例においては、マスタコントローラ224は、リモートのアプリケーション特有の論理回路(コンピュータネットワーク118を介して)および関連するドライバソフトウェアを調整する。別の実施例においては、マスタコントローラ224は、ローカルバス206または外部バス接続ハードディスクドライブ(「HDD」)ストレージコントローラを調整する。

30

40

【0136】

一実施例においては、マスタコントローラ224は、1又はそれ以上のストレージコントローラ254と通信し、ストレージデバイス/ソリッドステートストレージデバイス102は、SCSIバス、インターネットSCSI(「iSCSI」)、ファイバチャネル等を介して接続されたストレージデバイスとして存在してもよい。その間、ストレージデバイス/ソリッドステートストレージデバイス102は、オブジェクトを自律的に管理することができる。オブジェクトファイルシステムまたは分散型オブジェクトファイルシステムとして存在してもよい。マスタコントローラ224は、ピアコントローラ256および/またはアプリケーション特有のプロセッサ258によってアクセスすることができる。

【0137】

50

別の実施例においては、マスタコントローラ 224 は、自律的に統合された管理コントローラを調整し、FPGAコードおよび/またはコントローラソフトウェアを定期的に確認し、(リセット)を実行中にFPGAコードを確認し、および/または、(リセット)出力中にコントローラソフトウェアを確認し、外部リセット要求をサポートし、監視タイムアウトによるリセット要求をサポートし、電圧、電流、電力、温度及びその他の環境測定値をサポートし、閾値の設定が中断する。別の実施例においては、マスタコントローラ 224 は、ガベージコレクションを管理し、再利用のための消去ブロックを含まない。別の実施例においては、マスタコントローラ 224 は、摩耗レベルを管理する。別の実施例においては、マスタコントローラ 224 によって、データストレージデバイス/ソリッドステートストレージデバイス 102 が、複数の仮想デバイスに分割され、パーティションベースメディアの暗号化が可能となる。さらなる別の実施例においては、マスタコントローラ 224 は、有利な複数ビット ECC 修正部を有するソリッドステートストレージコントローラ 104 をサポートする。当分野の当業者は、ストレージコントローラ 152 において、または、より具体的には、ソリッドステートストレージデバイス 102 において、マスタコントローラ 224 のその他の特性および機能が分かるであろう。

10

20

30

40

50

【0138】

一実施例においては、ソリッドステートストレージデバイスコントローラ 202 は、ダイナミックランダムメモリアレイ 230 および/またはスタティックランダムメモリアレイ 232 を制御するメモリコントローラ 228 を具える。上述したように、メモリコントローラ 228 は、独立しているか、または、マスタコントローラ 224 と統合されている。一般的に、メモリコントローラ 228 は、DRAM (ダイナミックランダムメモリアレイ 230) および SRAM (スタティックランダムメモリアレイ 232) などのような、いくつかのタイプの揮発メモリを制御する。他の例においては、メモリコントローラ 228 は、電氣的消去可能なプログラム可能読み出し専用メモリ(「EEPROM」)などのような、その他のメモリタイプを制御できる。他の実施例においては、メモリコントローラ 228 は、2 又はそれ以上のメモリタイプを制御し、メモリコントローラ 228 は、1 以上のコントローラを具えてもよい。一般的に、メモリコントローラ 228 は、SRAM 232 を補助する DRAM 230 による、実行可能な限りの多くの SRAM 232 を制御する。

【0139】

一実施例においては、オブジェクトインデックスは、メモリ 230、232 内に記憶され、次いで、ソリッドステートストレージ 110n のチャンネルまたはその他の非揮発性メモリに定期的にオフロードされる。当分野の当業者は、ダイナミックメモリアレイ 230、メモリコントローラ 228 およびスタティックメモリアレイ 232 のその他の使用および設定が分かるであろう。

【0140】

一実施例においては、ソリッドステートストレージデバイスコントローラ 202 は DMA コントローラ 226 を具え、この DMA コントローラ 226 は、ストレージデバイス/ソリッドステートストレージデバイス 102 および 1 又はそれ以上の外部メモリコントローラ、関連する外部メモリアレイ 252 および CPU 248 間の DMA 動作を制御する。なお、外部メモリコントローラ 250 および関連する外部メモリアレイ 252 は、ストレージデバイス/ソリッドステートストレージデバイス 102 の外部にあるので外部と呼ぶことに留意されたい。更に、DMA コントローラ 226 も、NIC 244 を介した要求デバイスおよび関連する RDMA コントローラ 246 を用いて RDMA 動作を制御することができる。DMA および RDMA は、以下でより詳細に説明する。

【0141】

一実施例においては、ソリッドステートストレージデバイスコントローラ 202 は、管理バス 236 に接続された管理コントローラ 234 を具える。一般的に、管理コントローラ 234 は、環境基準値およびストレージデバイス/ソリッドステートストレージデバイス 102 のステータスを管理する。管理コントローラ 234 は、管理バス 236 を介して

デバイス温度、ファン速度、電力供給設定などをモニタすることができる。管理コントローラは、FPGAコードおよびコントローラソフトウェアのストレージ用の、消去可能なプログラム可能読み出し専用メモリ(「EEPROM」)の読み出しおよびプログラミングをサポートすることができる。一般的に、管理バス236は、ストレージデバイス/ソリッドステートストレージデバイス102内に様々な構成を接続させる。管理コントローラ234は、ローカルバス206を介して、アラート、中断などと通信でき、システムバス240または他のバスに対する別の接続を具えることができる。一実施例においては、管理バス236は、Inter-Integrated Circuit(I2C)バスである。当分野の当業者は、管理バス236によってストレージデバイス/ソリッドステートストレージデバイス102の構成に接続された、管理コントローラ234のその他の関連する機能および利用が分かるであろう。 10

【0142】

一実施例においては、ソリッドステートストレージデバイスコントローラ202は、特定のアプリケーションにカスタマイズ可能な種々雑多な論理回路242を具える。一般的に、ソリッドステートデバイスコントローラ202またはマスタコントローラ224が、FPGAまたはその他の構成可能コントローラを用いて構成される場合、特定のアプリケーション、カスタマ要求、ストレージ要求に基づく、カスタム論理回路が含まれる。

【0143】

データパイプライン

【0144】

図3は、本発明によるソリッドステートストレージデバイス102における書き込みデータパイプライン106および読み出しデータパイプライン108を有する、ソリッドステートストレージコントローラ104の一実施例300を示した概略ブロック図である。実施例300は、データバス204、ローカルバス206およびバッファコントローラ208を具え、これらは、図2のソリッドステートストレージデバイスコントローラ202に関連して記載されたものと実質的に同等である。書き込みデータパイプラインは、パケットタイザ302およびエラー修正コード(「ECC」)ジェネレータ304を具える。他の実施例においては、書き込みデータパイプラインは、入力バッファ306、書き込み同期バッファ308、書き込みプログラムモジュール310、圧縮モジュール312、暗号化モジュール314、ガベージコレクタバイパス316(読み出しデータパイプライン内の一部を有する)、メディア暗号化モジュール318、および、書き込みバッファ320を具える。読み出しデータパイプライン108は、読み出し同期バッファ328、ECC修正モジュール322、デパケットタイザ(depacketizer)324、アライメントモジュール326および出力バッファ330を具える。他の実施例においては、読み出しデータパイプライン108は、メディア復号化モジュール332、ガベージコレクタバイパス316の一部、復号化モジュール334、解凍モジュール336および読み出しプログラムモジュール338を具えることができる。ソリッドステートストレージコントローラ104は、更に、制御およびステータスレジスタ340、制御待ち行列342、バンクインタリーブコントローラ344、同期バッファ346、ストレージバスコントローラ348およびマルチプレクサ(「MUX」)350を具えることができる。ソリッドステートコントローラ104の構成および関連する書き込みデータパイプライン106および読み出しデータパイプライン108は、以下で説明する。他の実施例においては、同期ソリッドステートストレージ110を用いることができ、同期バッファ308、328を省くことができる。 30 40

【0145】

書き込みデータパイプライン

【0146】

書き込みデータパイプライン106は、直接または間接的に別の書き込みデータパイプライン106のステージを介して、ソリッドステートストレージに書き込まれるデータ又はメタデータセグメントを受信する、パケットタイザ302を具え、ソリッドステートスト 50

レンジ 110 用に調整された 1 又はそれ以上のパケットを生成する。データまたはメタデータセグメントは、一般的に、オブジェクトの一部であるが、オブジェクト全体を具えることもできる。別の実施例においては、データセグメントは、データブロックの一部であるが、データブロック全体を具えることもできる。一般的に、オブジェクトは、コンピュータ 112、クライアント 114 またはその他のコンピュータ若しくはデバイスから受信し、ソリッドステートストレージデバイス 102 またはコンピュータ 112 にストリーミングされるデータセグメントでソリッドステートストレージデバイス 102 に送信される。更に、データセグメントは、データ区域など、別の名称として周知であるが、本明細書で規定されるように、オブジェクトまたはデータブロックの総てまたは一部を具える。

【0147】

各オブジェクトは、1 又はそれ以上のパケットとして記憶される。各オブジェクトは、1 又はそれ以上のコンテナパケットを有することができる。各パケットはヘッダを含む。ヘッダは、ヘッダ型フィールドを具えることができる。タイプフィールドは、データ、オブジェクト属性、メタデータ、データセグメントデリミタ（マルチパケット）、オブジェクト構造、オブジェクト結合などを具えることができる。ヘッダは、パケットに含まれるデータのバイト数のような、パケットサイズに関する情報を具えることができる。パケットの長さは、パケットタイプによって構築可能である。ヘッダは、オブジェクトに対するパケットの関係性を構築する情報を具えることができる。一例では、オブジェクト内のデータセグメントの位置を識別するために、データパケットヘッダのオフセットを用いることができる。当分野の当業者は、パケットサイズ 302 によってデータに加えらるるヘッダに含まれるその他の情報、および、データパケットに加えらるるその他の情報が分かるであろう。

【0148】

各パケットは、データまたはデータセグメントからのヘッダおよび可能であればデータを具えることができる。各パケットのヘッダは、パケットが属するオブジェクトにパケットに関連させる、関連情報を具える。例えば、ヘッダは、データパケットを形成していたデータセグメント、オブジェクトまたはデータブロックを示すオブジェクト識別子およびオフセットを具えることができる。更に、ヘッダは、パケットを記憶するための、ストレージバスコントローラによって用いられる論理アドレスを具えることができる。また、ヘッダは、パケットに含まれるバイト数のような、パケットの大きさに関する情報を具えることができる。更に、ヘッダは、データセグメントまたはオブジェクトを復元するとき、データセグメントがオブジェクト内のその他のパケットに対して属しているところを識別するシーケンス番号を具えることができる。ヘッダは、ヘッダ型フィールドを具えることができる。タイプフィールドは、データ、オブジェクト属性、メタデータ、データセグメントデリミタ（マルチパケット）、オブジェクト構造、オブジェクト結合などを具えることができる。当分野当業者は、パケットサイズ 302 によってデータ又はメタデータに加えらるるヘッダに含まれてもよいその他の情報、および、パケットに加えることができる他の情報が分かるであろう。

【0149】

書き込みデータパイプライン 106 は、パケットサイズ 302 から受信する 1 又はそれ以上のパケット用の 1 又はそれ以上のエラー修正コード（「ECC」）を生成する ECC ジェネレータ 304 を具える。ECC ジェネレータ 304 は、一般的に、エラー修正アルゴリズムを用いて、パケットとともに記憶される ECC を生成する。パケットとともに記憶された ECC は、一般的に、送信および記憶を通ったデータに入ったエラーを検出し修正するために用いられる。一実施例においては、パケットは、長さ N の符号化されないブロックとして ECC ジェネレータ 304 内にストリーミングされる。長さ S のシンδροームが算出され、追加され、長さ N + S の符号化ブロックとして出力される。N 値および S 値は、特別な性能、効率およびロバスト性の基準値を得るために選択されるアルゴリズムの特性に依存している。好ましい実施例においては、ECC ブロックとパケットとの間において、固定された関係性はなく、パケットは、1 以上の ECC ブロックを具えることがで

10

20

30

40

50

き、ECCブロックは、1以上のパケットを具えることができ、第1のパケットは、ECCブロック内のどこで終結してもよく、同一のECCブロック内の第1のパケットの末端の後から第2のパケットは開始することができる。好ましい実施例においては、ECCアルゴリズムは、動的に変更されない。好ましい実施例においては、データパケットとともに記憶されたECCは、2より大きいビットにおけるエラーを修正するのに十分にロバストである。

【0150】

有利には、シングルビット以上の修正またはダブルビット修正でさえ可能なロバストなECCアルゴリズムを用いることで、ソリッドステートストレージ110の寿命を延長させることができる。例えば、ソリッドステートストレージ110内のストレージ媒体としてフラッシュメモリが用いられる場合、フラッシュメモリは、消去サイクルにつきエラーなしで、約100,000回書き込み可能である。この利用限界は、ロバストなECCアルゴリズムを用いて拡張することができる。ECCジェネレータ304を有し、ソリッドステートストレージデバイス102に搭載されたECC修正モジュール322を対応させることで、ソリッドステートストレージデバイス102は、エラーを内部で修正ことができ、シングルビット修正のような、より低ロバスト性のECCアルゴリズムを使用した場合よりも長い寿命を有することができる。しかしながら、その他の実施例においては、ECCジェネレータ304は、より低ロバスト性のアルゴリズムを用いることができ、シングルビットまたはダブルビットエラーを修正できる。別の実施例においては、ソリッドステートストレージデバイス110は、容量を増大するためにマルチレベルセル(「MLC」)フラッシュなど信頼性が低いストレージを具えてもよく、ストレージは、よりロバストなECCアルゴリズムがなく、十分な信頼性がなくてもよい。

10

20

【0151】

一実施例においては、書き込みパイプライン106は、入力バッファ306を具え、この入力バッファ306は、ソリッドステートストレージ110に書き込まれるデータセグメントを受信し、パケッタイザ302のような、書き込みデータパイプライン106の次のステージ(またはより複雑な書き込みデータパイプライン106用のその他のステージ)が準備され、次のデータセグメントを処理されるまでに、入ってくるデータセグメントを記憶する。一般的に、入力バッファ306によって、データセグメントが適宜なサイズのデータバッファを用いて書き込みデータパイプライン106によって受信、処理される頻度の差異を許容することができる。更に、入力バッファ306によって、データバス204の動作効率を改善するために、データバス204は、書き込みデータパイプライン106によって維持できるよりも大きな頻度で、データを書き込みデータパイプライン106に転送できる。一般的に、書き込みデータパイプライン106は、入力バッファ306を具えていないときは、バッファリング機能は、ソリッドステートストレージ102内であるが書き込みパイプライン106の外側、ネットワークインタフェースカード(「NIC」)内のようなコンピュータ112内、または、例えば、リモートダイレクトメモリアクセス(「RDAM」)を用いるときの別のデバイスのような、あらゆる場所で行われる。

30

【0152】

別の実施例においては、書き込みデータパイプライン106は、ソリッドステートストレージ110にパケットを書き込む前に、ECCジェネレータ304から受信したパケットをバッファリングする書き込み同期バッファ308を更に具えることができる。書き込み同期バッファ308は、ローカルクロックドメインとソリッドステートストレージクロックドメインとの間の境界に配置され、クロックドメイン差を考慮するためのバッファリングを提供する。他の実施例においては、同期型ソリッドステートストレージ110を用いることができ、同期用バッファ308、328を省くことができる。

40

【0153】

一実施例においては、書き込みデータパイプライン106は、更に、メディア暗号化モジュール318を具えることができ、このモジュール318は、パケッタイザ302から

50

直接または間接的に、1又はそれ以上のパケットを受信して、ECCジェネレータ304にパケットを送信する前に、ソリッドステートストレージデバイス102に固有の暗号化キーを用いて、1又はそれ以上のパケットを暗号化する。一般的に、ヘッダを含む完全なパケットは暗号化される。別の実施例においては、ヘッダは、暗号化されない。本明細書においては、暗号化キーは、ソリッドステートストレージ110を統合する、暗号化保護を必要とする具体例の外部で管理されるシークレット暗号化キーを意味することを理解されたい。メディア暗号化モジュール318および関連するメディア暗号化モジュール332は、ソリッドステートストレージ110に記憶されるデータの安全レベルを提供する。例えば、データが、メディア暗号化モジュール318で暗号化されるとき、ソリッドステートストレージ110が、異なるソリッドステートストレージコントローラ104、ソリッドステートストレージデバイス102またはコンピュータ112に接続される場合、ソリッドステートストレージ110の内容は、一般的に、大きな労力を要せずに、ソリッドステートストレージ110にデータを書き込む間に用いられる同一の暗号化キーを使用せずに読み出すことができない。

10

20

30

40

50

【0154】

一般的な実施例においては、ソリッドステートストレージデバイス102は、不揮発性ストレージの暗号化キーを記憶できず、暗号化キーへのいかなる外部からのアクセスもできない。暗号化キーは、初期化中に、ソリッドステートストレージコントローラ104に提供される。ソリッドステートストレージデバイス102は、暗号化キーとともに用いられる、非シークレット暗号化ナンスを使用および記憶できる。異なるナンスは、パケット毎に記憶できる。データセグメントは、暗号化アルゴリズムによる保護を改善するために、固有のナンスを有する複数のパケット間に分割できる。暗号化キーは、クライアント114、コンピュータ112、キーマネージャ、または、ソリッドステートストレージコントローラ104によって用いられる暗号化キーを管理するその他のデバイスから受信することができる。別の実施例においては、ソリッドステートストレージ110は、2又はそれ以上のパーティションを有し、ソリッドステートストレージコントローラ104は、まるで2又はそれ以上のソリッドステートストレージコントローラ104であるかのように動作し、各々は、ソリッドステートストレージ110内の単一のパーティションで動作する。この実施例においては、固有のメディア暗号化キーは、各パーティションで用いることができる。

【0155】

別の実施例においては、書き込みデータパイプライン106は、更に、パケットサイズ302にデータセグメントを送信する前に、直接または間接的に、入力バッファ306から受信されるデータまたはデータセグメントを暗号化する暗号化モジュール314を具えることができ、データセグメントは、データセグメントとともに受信した暗号化キーを用いて暗号化される。暗号化モジュール314は、データを暗号化する暗号化モジュール318によって用いられる暗号化キーは、ソリッドステートストレージデバイス102内に記憶される総てのデータに共通でなくてもよいが、以下で記載するように、データセグメントを受信するとともに受信されるオブジェクトベースで変化するという点でメディア暗号化モジュール318とは異なる。例えば、暗号化モジュール318によって暗号化されるべきデータセグメント用の暗号化キーは、データセグメントとともに受信され、データセグメントが属するオブジェクトを書き込むコマンドの一部として受信される。ソリッドステートストレージデバイス102は、暗号化キーとともに用いられる各オブジェクトパケットにおいて、非シークレット暗号化ナンスを用いて記憶することができる。異なるナンスは、パケット毎に記憶できる。データセグメントは、暗号化アルゴリズムによる保護を改善するために、固有のナンスを有する複数のパケットに分割できる。一実施例においては、メディア暗号化モジュール318によって用いられるナンスは、暗号化モジュール314によって用いられるものと同ーである。

【0156】

暗号化キーは、クライアント114、コンピュータ112、キーマネージャ、または、

データセグメントを暗号化するのに用いられる暗号化キーを保持したその他デバイスから受信することができる。一実施例においては、ソリッドステートストレージデバイス 102、コンピュータ 112、クライアント 114、または、プライベートおよびパブリックキーを安全に転送して保護する、工業規格の方法を実行する能力を有するその他の外部エージェント、のうちの 1 つから、暗号化キーはソリッドステートストレージコントローラ 104 に転送される。

【0157】

一実施例においては、暗号化モジュール 318 は、第 1 のパケットを暗号化し、第 1 の暗号化キーはこのパケットとともに受信され、また、第 2 のパケットを暗号化し、第 2 の暗号化キーは第 2 のパケットとともに受信される。別の実施例においては、暗号化モジュール 318 は、第 1 のパケットを暗号化し、第 1 の暗号化キーは、このパケットとともに受信され、第 2 のデータパケットを暗号化せず次のステージに通過させる。有利には、ソリッドステートストレージデバイス 102 の書き込みデータパイプライン 106 に含まれる暗号化モジュール 318 によって、オブジェクトごとまたはセグメントごとのデータ暗号化が、単一のファイルシステムまたはその他のファイルシステムなしで、対応するオブジェクトまたはデータセグメントを記憶するのに用いられる異なる暗号化キーのトラックを、保持できる。各要求デバイス 155 または対応するキーマネージャは、要求デバイス 155 によって送信されるオブジェクトまたはデータセグメントのみを暗号化するのに用いられる暗号化キーを別個に管理することができる。

10

【0158】

別の実施例においては、書き込みデータパイプライン 106 は、パケットタイザ 302 にデータセグメントを送信する前に、メタデータセグメント用のデータを圧縮する圧縮モジュール 312 を具える。圧縮モジュール 312 は、一般的に、セグメントのストレージサイズを低減するのに当分野の当業者に周知な圧縮ルーチンを用いて、データまたはメタデータセグメントを圧縮する。例えば、データセグメントは、512 のゼロの記号列を具える場合、圧縮モジュール 312 は、512 のゼロを、512 のゼロで示されるコード又はトークンで置換でき、コードは、512 のゼロによって取られる空き領域よりも更に小型である。

20

【0159】

一実施例においては、圧縮モジュール 312 は、第 1 の圧縮ルーチンを用いて第 1 のセグメントを圧縮し、圧縮せずに第 2 のセグメントを通過させる。別の実施例においては、圧縮モジュール 312 は、第 1 の圧縮ルーチンを用いて第 1 のセグメントを圧縮し、第 2 の圧縮ルーチンを用いて第 2 のセグメントを圧縮する。ソリッドステートストレージデバイス 102 内にフレキシビリティを有することに利点があり、クライアント 114、または、ソリッドステートストレージデバイス 102 にデータを書き込むその他のデバイスは圧縮ルーチンを各々特定でき、あるいは、一方が圧縮ルーチンを特定することができ、他方がいずれの圧縮ルーチンを特定しなくてもよい。圧縮ルーチンの選択は、オブジェクトごとのタイプまたはオブジェクトクラスに基づいたデフォルトの設定に従って選択されてもよい。例えば、特定のオブジェクトの第 1 のオブジェクトは、デフォルトの圧縮ルーチン設定を取り消すことができ、同一オブジェクトクラスおよびオブジェクトタイプの第 2 のオブジェクトは、デフォルトの圧縮ルーチンを用いることができ、同一オブジェクトクラスおよびオブジェクトタイプの第 3 のオブジェクトは、圧縮を利用しなくてもよい。

30

40

【0160】

一実施例においては、書き込みデータパイプライン 106 は、ガベージコレクションシステムのデータパイプの一部として、読み出しデータパイプライン 108 からデータセグメントを受信するガベージコレクションパイプ 316 を具える。ガベージコレクションシステムは、一般的に、パケットが削除用に標識又は変更され、変更されたデータが異なる位置に記憶されるので、もはや有効でないパケットを一般的に標識する。いくつかの場所では、ガベージコレクションシステムは、ストレージの特定のセクションが回復可能

50

かを判定する。この判定は、利用可能なストレージ容量の不足、閾値に達した無効として標識されたデータの割合、有効なデータの統合、閾値に達し、又はデータ配置に基づき性能を改善したストレージセクションのエラー検出率等によってもよい。多くの因子は、ガベージコレクションアルゴリズムによって考慮され、ストレージセクションを回復すべきときを判定する。

【0161】

ストレージセクションが回復のために標識されると、一般的に、そのセクション中の有効なパケットは再配置されなければならない。ガベージコレクタバイパス316によって、パケットは、読み出しデータパイプライン108に読み出され、次いで、ソリッドステートストレージコントローラ104の外にルーティングされずに、書き込みデータパイプライン106に直接転送される。好ましい実施例においては、ガベージコレクタバイパス316は、ソリッドステートストレージデバイス102内で動作する自律型ガベージコレクタシステムの一部である。これによって、ソリッドステートストレージデバイス102は、データを管理することが可能となり、データは、ソリッドステートストレージ110にわたって系統的に分散され、性能、データ信頼性を改善し、ソリッドステートストレージ110のいずれか1の位置または領域の過剰使用または使用不足を回避し、ソリッドステートストレージ110の寿命を延長させる。

10

【0162】

ガベージコレクタバイパス316は、セグメントを書き込みデータパイプライン106に挿入し、その他のセグメントは、クライアント116またはその他のデバイスによって書き込まれるように調整する。示された実施例においては、ガベージコレクタバイパス316は、書き込みデータパイプライン106内のパケットサイズ302の前および読み出しデータパイプライン108内のデパケットサイズ324の後にあるが、読み出しおよび書き込みデータパイプライン106、108内のどこかに配置されていてもよい。ガベージコレクタバイパス316は、書き込みパイプライン106のフラッシュ中に用いて、ソリッドステートストレージ110内のストレージ効率を改善するために仮想ページの残りを満たし、これにより、ガベージコレクションの頻度を低減することができる。

20

【0163】

一実施例においては、書き込みデータパイプライン106は、有効な書き込み動作のデータをバッファリングする書き込みバッファ320を具える。一般的に、書き込みバッファ320は、ソリッドステートストレージ110内の少なくとも1つの仮想ページを満たすのに十分なパケット容量を具える。これによって、書き込み動作は、中断せずに、ソリッドステートストレージ110にデータのページ全体を送信することができる。書き込みデータパイプライン106の書き込みバッファ320と、読み出しデータパイプライン108内のバッファとを、ソリッドステートストレージ110内のストレージ書き込みバッファと同一、または、大きい容量にサイズ変更することによって、書き込みおよび読み出しデータは、複数コマンドの代わりに、単一の書き込みコマンドが、ソリッドステートストレージ110にデータの仮想ページ全体を送信するように生成されるので、より効率的になる。

30

【0164】

書き込みバッファ320が満たされる際に、ソリッドステートストレージ110は、その他の読み出し動作に用いられてもよい。このことが有利であるのは、データがストレージ書き込みバッファに書き込まれ、ストレージ書き込みバッファ内に流れるデータがストールする際に、書き込みバッファがより小さいか、書き込みバッファのないその他のソリッドステートデバイスが、ソリッドステートストレージをタイアップできるからである。読み出し動作は、ストレージ書き込みバッファ全体が満たされ、プログラム化されるまで、ブロックされる。書き込みバッファのないまたは書き込みバッファが小さいシステムに対する別のアプローチは、読み出し可能にするために満たされないストレージ書き込みバッファを消去することである。更に、ページを満たすのに複数の書き込み/プログラムサイクルが必要となるので、これは非効率的である。

40

50

【0165】

仮想ページよりも大きくサイズ変更された書き込みバッファ320を有する、示された実施例について、多数のサブコマンドを具える単一の書き込みコマンドは、次いで、単一のプログラムコマンドが続き、各ソリッドステートストレージエレメント216、218、220におけるストレージ書き込みバッファから、各ソリッドステートストレージエレメント216、218、220内の指定されたページに、データのページを転送する。この技術は、データ信頼性および耐久性を低下させることで周知である部分的なページプログラミングを除去し、バッファが満たされるとき読み出しコマンドおよびその他のコマンドのための指定バンクを解放させる利点がある。

【0166】

一実施例においては、書き込みバッファ320は、ピンポンバッファであり、このバッファの片面が満たされ、次いで、適宜な時間での転送が指定され、ピンポンバッファのもう一方の面が満たされる。別の実施例においては、書き込みバッファ320は、データセグメントの仮想ページよりも大きい容量を有する先入れ先出し(「FIFO」)レジスタを具える。当分野の当業者は、データの仮想ページがデータをソリッドステートストレージ110に書き込む前に記憶させることができるその他の書き込みバッファ320構成が分かるであろう。

【0167】

別の実施例においては、書き込みバッファ320は、仮想ページよりも小さなサイズであり、情報のより少ないページが、ソリッドステートストレージ110内のストレージ書き込みバッファに書き込まれる。その実施例においては、書き込みデータパイプライン106内のストールが読み出し動作を止めるのを防ぐために、データは、1つの位置からガベージコレクションプロセスの一部としての別の位置へと移動する必要があるデータが、ガベージコレクションシステムを用いて待ち行列に入れられる。書き込みデータパイプライン106にけるデータストールの場合、データは、ガベージコレクタバイパス316を介して、書き込みバッファ320、次いで、ソリッドステートストレージ110内のストレージ書き込みバッファへと供給可能であり、データをプログラミングする前に、仮想ページのページを満たす。この方法では、書き込みデータパイプライン106内のデータストールは、ソリッドステートストレージデバイス106から読み出しをストールしない。

【0168】

別の実施例においては、書き込みデータパイプライン106は、書き込みデータパイプライン106内の1又はそれ以上のユーザ定義可能な機能を有する書き込みプログラムモジュール310を具える。書き込みプログラムモジュール310によって、ユーザは書き込みデータパイプライン106をカスタマイズすることができる。ユーザは、特定のデータ要求またはアプリケーションに基づいて書き込みデータパイプライン106をカスタマイズすることができる。ソリッドステートストレージコントローラ104はFPGAである場合、ユーザは、カスタムコマンドおよび機能を有する書き込みデータパイプライン106を比較的容易にプログラムすることができる。更に、ユーザは、書き込みプログラムモジュール310を使用し、ASICを用いたカスタム機能を具えるが、しかしながら、ASICをカスタマイズすることは、FPGAを用いるよりも難しい。書き込みプログラムモジュール310は、バッファおよびバイパス機構を具えることができ、第1のデータセグメントが書き込みプログラムモジュール310で実行され、第2のデータセグメントは、書き込みデータパイプライン106を介して継続できる。別の実施例においては、書き込みプログラムモジュール310は、ソフトウェアを介してプログラム可能なプロセッサコアを具えることができる。

【0169】

書き込みプログラムモジュール310は、入力バッファ306と圧縮モジュール312との間に示されているが、書き込みプログラムモジュール310は、書き込みデータパイプライン106内のあらゆる位置に存在でき、様々なステージ302-320間に配置可能であることに留意されたい。更に、別個にプログラミングおよび動作される様々なステ

10

20

30

40

50

ージ 3 0 2 - 3 2 0 間に配置される、複数の書き込みプログラムモジュール 3 1 0 があってもよい。更に、ステージ 3 0 2 - 3 2 0 の順番は変更可能である。当分野の当業者は、特定のユーザ要求に基づいてステージ 3 0 2 - 3 2 0 の順番に対する作動可能な変更が分かるであろう。

【 0 1 7 0 】

読み出しデータパイプライン

【 0 1 7 1 】

読み出しデータパイプライン 1 0 8 は、要求されたパケットの各 E C C ブロックで記憶される E C C を用いることによって、データエラーが E C C ブロックにある場合に、ソリッドステートストレージ 1 1 0 から受信される要求されたパケットを判定する E C C 修正モジュール 3 2 2 を具える。次いで、E C C 修正モジュール 3 2 2 は、エラーが存在し、エラーが E C C を用いて修正可能である場合、要求されたパケットにおけるエラーを修正する。例えば、E C C が、6 ビットのエラーを検出できるが、3 ビットのエラーしか修正できない場合、E C C 修正モジュール 3 2 2 は、エラー中の 3 ビットまで、要求されたパケットの E C C ブロックを修正する。E C C 修正モジュール 3 2 2 は、エラー中のビットを、正しい 1 またはゼロ状態に変更することで、エラー中のビットを修正し、要求されたデータパケットがソリッドステートストレージ 1 1 0 に書き込まれ、E C C がパケット用に生成された場合と同一になる。

10

【 0 1 7 2 】

E C C が修正できるよりも多く、要求されたパケットがエラー中のビットを含んでいると、E C C 修正モジュール 3 2 2 が判定した場合、E C C 修正モジュール 3 2 2 は、要求されたパケットの破損した E C C ブロック中のエラーを修正できず、中断信号を送信する。一実施例においては、E C C 修正モジュール 3 2 2 は、要求されたパケットがエラー状態にあることを示すメッセージを有する中断信号を送信する。メッセージは、E C C 修正モジュール 3 2 2 がエラーを修正できず、エラーを修正するための E C C 修正モジュール 3 2 2 の能力がないことが示される情報を含む。別の実施例においては、E C C 修正モジュール 3 2 2 は、要求されたパケットの破損 E C C ブロックを、中断信号および/またはメッセージとともに送信する。

20

【 0 1 7 3 】

好ましい実施例においては、E C C 修正モジュール 3 2 2 によって修正不可能な要求されたパケットの、破損した E C C ブロックまたは破損した E C C ブロックの一部は、マスタコントローラ 2 2 4 によって読み出され、修正され、E C C 修正モジュール 3 2 2 に戻されて、読み出しデータパイプライン 1 0 8 によって更に処理される。一実施例においては、要求されたパケットのまたは破損した E C C ブロックの一部は、データを要求しているデバイスに送信される。要求デバイス 1 5 5 は、E C C ブロックを修正することができ、バックアップまたはミラーコピーなど、別のコピーを用いてデータを置換することができ、次いで、要求されたデータパケットの置換データを用いるか、それを読み出しデータパイプライン 1 0 8 に戻すことができる。要求デバイス 1 5 5 は、エラー中の要求されたパケットのヘッダ情報を用いて、破損した要求されたパケットを置換するか、または、パケットが属するオブジェクトを置換するかを必要とするデータを識別する。別の好ましい実施例においては、ソリッドステートストレージコントローラ 1 0 4 は、いくつかのタイプの R A I D を用いてデータを記憶し、破損したデータを回復することができる。別の実施例においては、E C C 修正モジュール 3 2 2 は、中断信号および/またはメッセージを送信し、受信するデバイスは、要求されたデータパケットに関連した読み出し動作ができない。当分野の当業者は、要求されたパケットの 1 又はそれ以上の E C C ブロックが破損され、E C C 修正モジュール 3 2 2 がエラーを修正できないことを判定した E C C 修正モジュール 3 2 2 の結果として、とるべきその他の選択肢や動作が分かるであろう。

30

40

【 0 1 7 4 】

読み出しデータパイプライン 1 0 8 は、デパケッタイザ 3 2 4 を具え、このデパケッタイザ 3 2 4 は、直接または間接的に、E C C 修正モジュール 3 2 2 から要求されたパケッ

50

トのECCブロックを受信し、1又はそれ以上のパケットヘッダをチェック及び除去する。デパケッタイザ324はヘッダ内のパケット識別子、データ長、データ位置等进行检查することによってパケットヘッダを確認する。一実施例においては、ヘッダは、読み出しデータパイプライン108に配信されたパケットが、要求されたパケットであることを確認するために使用可能なハッシュコードを具える。更に、デパケッタイザ324は、パケッタイザ302によって加えられた要求されたパケットからヘッダを除去する。デパケッタイザ324は、所定のパケット上で動作しないように命令できるが、修正されずに先に通過させる。例としては、復元プロセスの途中で要求されるコンテナラベルであってもよく、ヘッダ情報は、オブジェクトインデックス復元モジュール272によって要求される。さらなる例としては、ソリッドステートストレージデバイス102内で用いられるように決められた様々なタイプのパケットの転送を含む。別の実施例においては、デパケッタイザ324の動作は、パケットタイプに応じたものであってもよい。

10

【0175】

読み出しデータパイプライン108は、デパケッタイザ324からデータを受信し、所望されないデータを除去するアライメントモジュール326を具える。一実施例においては、ソリッドステートストレージ110に送信される読み出しコマンドは、パケットのデータを回復させる。データを要求しているデバイスは、回復したパケットの総てのデータを必要としなくてもよく、アライメントモジュール326は、所望されないデータを除去する。回復したページ内の総てのデータが要求されているデータである場合、アライメントモジュール326は、いずれのデータも除去しない。

20

【0176】

アライメントモジュール326は、次のステージにデータセグメントを転送する前に、データセグメントを要求しているデバイスと互換性のある形態で、オブジェクトのデータセグメントとしてデータを再フォーマットする。一般的に、データは、読み出しデータパイプライン108によって処理される場合、データセグメントまたはパケットのサイズは、様々なステージで変更される。アライメントモジュール326は、受信したデータを用いて、データを、要求デバイス155に送信し、応答を形成するのに加えられるのに適したデータセグメントにフォーマットする。例えば、第1のデータパケットの一部からのデータは、第2のデータパケットの一部からのデータと組み合わせることができる。データセグメントが、要求デバイスによって要求されるデータよりも長い場合、アライメントモジュール326は、所望されないデータを廃棄できる。

30

【0177】

一実施例においては、読み出しデータパイプライン108は、読み出しデータパイプライン108によって処理される前に、ソリッドステートストレージ110から読み出される1又はそれ以上の要求されたパケットをバッファリングする読み出し同期バッファ328を具える。読み出し同期バッファ328は、ソリッドステートストレージクロックドメインと、ローカルバスクロックドメインとの間の境界にあり、クロックドメインの差を考慮するバッファリングを提供する。

【0178】

別の実施例においては、読み出しデータパイプライン108は、出力バッファ330を具え、この出力バッファ330は、アライメントモジュール326からの要求されたパケットを受信し、要求デバイスに送信される前にパケットを記憶する。出力バッファ330は、データセグメントが、読み出しデータパイプライン108のステージから受信される時と、データセグメントが、ソリッドステートストレージコントローラ104のその他の部分または要求デバイスに送信される時と、の間の差を考慮している。更に、出力バッファ330によって、データバス204は、データバス204の動作効率を改善するために、読み出しデータパイプライン108によって維持されうるよりも高い頻度で、読み出しデータパイプライン108からのデータを受信できる。

40

【0179】

一実施例においては、読み出しデータパイプライン108は、メディア復号化モジュール

50

ル 3 3 2 を具備、このメディア復号化モジュール 3 3 2 は、E C C 修正モジュール 3 2 2 からの 1 又はそれ以上の暗号化された要求されたパケットを受信し、1 又はそれ以上の要求されたパケットをデパケッタイザ 3 2 4 に送信する前に、ソリッドステートストレージデバイス 1 0 2 に固有の暗号化キーを用いて 1 又はそれ以上の要求されたパケットを復号化する。一般的に、メディア復号化モジュール 3 3 2 によってデータを復号化するのに用いられる暗号化キーは、メディア復号化モジュール 3 1 8 によって用いられる暗号化キーと同一である。別の実施例においては、ソリッドステートストレージ 1 1 0 は、2 又はそれ以上のパーティションを具備していてもよく、ソリッドステートストレージコントローラ 1 0 4 は、2 又はそれ以上のソリッドステートストレージコントローラ 1 0 4 であるかのように動作し、各々は、ソリッドステートストレージ 1 1 0 内の単一のパーティションで動作する。この実施例においては、固有のメディア暗号化キーは、各パーティションで用いられてもよい。

10

【 0 1 8 0 】

別の実施例においては、読み出しデータパイプライン 1 0 8 は、復号化モジュール 3 3 4 を具備、この復号化モジュールは、データセグメントを出力バッファ 3 3 0 に送信する前に、デパケッタイザ 3 2 4 によってフォーマットされるデータセグメントを復号化する。データセグメントは、読み出し要求と共に受信された暗号化キーを用いて復号化され、読み出し要求は、読み出し同期バッファ 3 2 8 によって受信される要求されたパケットの回復を開始させる。復号化モジュール 3 3 4 は、第 1 のパケットの読み出し要求とともに受信される暗号化キーを有する第 1 のパケットを復号化でき、次いで、異なる暗号化キーを有する第 2 のパケットを復号化することができ、または、復号化せずに読み出しデータパイプライン 1 0 8 の次のステージに第 2 のパケットを通過させることができる。一般的に、復号化モジュール 3 3 4 は、メディア復号化モジュール 3 3 2 が要求されたパケットを復号化するのに用いるものとは異なる暗号化キーを用いてデータセグメントを復号化する。パケットが非シークレット暗号化ナンスで記憶されていた場合は、ナンスを、暗号化キーとともに用いて、データパケットを復号化する。暗号化キーは、クライアント 1 1 4 、コンピュータ 1 1 2 、キーマネージャ、または、ソリッドステートストレージコントローラ 1 0 4 によって用いられる暗号化キーを管理するその他のデバイスから、受信できる。

20

【 0 1 8 1 】

別の実施例においては、読み出しデータパイプライン 1 0 8 は、デパケッタイザ 3 2 4 によってフォーマットされたデータセグメントを解凍する解凍モジュール 3 3 6 を具備する。好ましい実施例においては、解凍モジュール 3 3 6 は、パケットヘッダおよびコンテナラベルの 1 又は双方に記憶された圧縮情報を用いて、圧縮モジュール 3 1 2 によってデータを圧縮するのに用いられるのと相補的なルーチンを選択する。別の実施例においては、解凍モジュール 3 3 6 によって用いられる解凍ルーチンは、解凍しようとするデータセグメントを要求しているデバイスによって命令される。別の実施例においては、解凍モジュール 3 3 6 は、オブジェクト毎のタイプまたはオブジェクトクラスベースのデフォルトの設定によって、解凍ルーチンを選択する。第 1 のオブジェクトの第 1 のパケットは、デフォルトの解凍ルーチンを無効にすることができ、同一オブジェクトクラスおよびオブジェクトタイプの第 2 のオブジェクトの第 2 のパケットは、デフォルトの解凍ルーチンを用いることができ、同一オブジェクトクラスおよびオブジェクトタイプの第 3 のオブジェクトの第 3 のパケットは、解凍を利用しなくてもよい。

30

40

【 0 1 8 2 】

別の実施例においては、読み出しデータパイプライン 1 0 8 は、読み出しデータパイプライン 1 0 8 内の 1 又はそれ以上のユーザ定義可能な機能を含む、読み出しプログラムモジュール 3 3 8 を具備する。読み出しプログラムモジュール 3 3 8 は、書き込みプログラムモジュール 3 1 0 と同様の特性を有し、ユーザは、読み出しデータパイプライン 1 0 8 にカスタム機能を提供することができる。読み出しプログラムモジュール 3 3 8 は、図 3 に示されるように配置可能であり、または、読み出しデータパイプライン 1 0 8 内で別のポ

50

ジションに配置可能であり、あるいは、読み出しデータパイプライン108内の複数の位置の複数の部分を具えてもよい。更に、別個に動作する読み出しデータパイプライン108内の複数の位置内に、複数の読み出しプログラムモジュール338があってもよい。当分野の当業者は、読み出しデータパイプライン108内の読み出しプログラムモジュール338のその他の形態が分かるであろう。書き込みデータパイプライン106と同様に、読み出しデータパイプライン108のステージは、再配置可能であり、当分野の当業者は、読み出しデータパイプライン108内のその他の順番が分かるであろう。

【0183】

ソリッドステートストレージコントローラ104は、制御およびステータスレジスタ340と、対応する制御待ち行列342を具える。制御およびステータスレジスタ340および制御待ち行列342は、書き込みおよび読み出しデータパイプライン106、108内で処理されるデータに関連した、制御およびシーケンスコマンドならびにサブコマンドを容易にする。例えば、パケットサイズ302のデータセグメントは、ECCジェネレータに関連した制御待ち行列342において、1又はそれ以上の対応する制御コマンドまたは命令を有することができる。データセグメントが、パケット化される時、命令またはコマンドのいくつかは、パケットサイズ302内で実行されてもよい。その他のコマンドまたは命令は、データセグメントから生成された、新規に形成されたデータパケットが、次のステージへ通過するとき、制御およびステータスレジスタ340を介して次の制御待ち行列342へ通過する。

10

【0184】

コマンドまたは命令は、制御待ち行列342に同時にロードされてもよく、パケットは、書き込みデータパイプライン106に転送され、各パイプラインステージは、各々のパケットがそのステージで実行される時、適宜なコマンドまたは命令を引き出す。同様に、コマンドまたは命令は、制御待ち行列342に同時にロード可能であり、パケットは、読み出しデータパイプライン108から要求されており、各パイプラインステージは、各々のパケットがそのステージで実行される時に、適宜なコマンドまたは命令を引き出す。当分野の当業者は、制御およびステータスレジスタ340、ならびに、制御待ち行列342のその他の特性および機能が分かるであろう。

20

【0185】

ソリッドステートストレージコントローラ104および/またはソリッドステートストレージデバイス102は更に、バンクインタリーブコントローラ344、同期バッファ346、ストレージバスコントローラ348およびマルチプレクサ(「MUX」)350を具えることができ、これらは、図4Aおよび図4Bに関連して記載される。

30

【0186】

バンクインタリーブ

【0187】

図4Aは、本発明によるソリッドステートストレージコントローラ104内のバンクインタリーブコントローラ344の一実施例400を示した概略ブロック図である。バンクインタリーブコントローラ344は、制御およびステータスレジスタ340に接続され、MUX350、ストレージコントローラ348および同期バッファ346を介してストレージI/Oバス210およびストレージ制御バス212に接続され、これらは以下で説明される。バンクインタリーブコントローラは、ソリッドステートストレージ110内のバンク214に対し、読み出しエージェント402、書き込みエージェント404、消去エージェント406、管理エージェント408、読み出し待ち行列410a-n、書き込み待ち行列412a-n、消去待ち行列414a-n、及び管理待ち行列416a-nを、並びに、バンクコントローラ418a-n、バスアービタ420、及びステータスMUX422を具え、これらは以下で説明する。ストレージバスコントローラ348は、再マッピングモジュール430を有するマッピングモジュール424、ステータスキャプチャモジュール426、および、NANDバスコントローラ428を具え、これらは以下で説明する。

40

50

【0188】

バンクインタリーブコントローラ344は、バンクインタリーブコントローラ344において、1又はそれ以上のコマンドを2又はそれ以上の待ち行列に命令し、ソリッドステートストレージ110のバンク214間で、上記待ち行列に記憶されたコマンドの実行を調整し、第1のタイプのコマンドは1つのバンク214aで実行され、第2のタイプのコマンドは第2のバンク214bで実行されるようにする。1又はそれ以上のコマンドは、コマンドタイプによって待ち行列に分けられる。ソリッドステートストレージ110の各バンク214は、バンクインタリーブコントローラ344内の対応する待ち行列のセットを有し、各々の待ち行列のセットは、各コマンドタイプ用の待ち行列を具える。

【0189】

バンクインタリーブコントローラ344は、ソリッドステートストレージ110のバンク214間で、待ち行列に記憶されたコマンドの実行を調整する。例えば、第1のタイプのコマンドは、1つのバンク214aで実行され、第2位のタイプのコマンドは、第2のバンク214bで実行される。一般的に、コマンドタイプおよび待ち行列タイプは、読み出しおよび書き込みコマンド、ならびに、待ち行列410、412を具えるが、更に、ストレージメディア特有のその他のコマンドおよび待ち行列を具えることができる。例えば、図4Aに示される実施例において、消去および管理待ち行列414、416が含まれ、フラッシュメモリ、N RAM、M RAM、D RAM、P RAMなどに好適である。

【0190】

その他のタイプのソリッドステートストレージ110のために、その他のタイプのコマンドおよび対応する待ち行列も、本発明の範囲から逸脱することなく含むことができる。FPGAソリッドステートストレージコントローラ104のフレキシブルな性質は、ストレージメディアにおけるフレキシビリティを与える。フラッシュメモリが、別のソリッドステートストレージタイプに変更された場合、バンクインタリーブコントローラ344、ストレージバスコントローラ348、および、MUX350は、データパイプライン106、108、および、その他のソリッドステートストレージコントローラ104の機能に重大な影響を有さないメディアタイプを収容できるように変更可能である。

【0191】

図4Aに示された実施例においては、バンクインタリーブコントローラ344は、各バンク214用の、ソリッドステートストレージ110からのデータを読み出すための読み出し待ち行列410、ソリッドステートストレージ110に対する書き込みコマンド用の書き込み待ち行列412、ソリッドステートストレージ中の消去ブロックを消去するための消去待ち行列414、管理コマンド用の管理待ち行列416を具える。更に、バンクインタリーブコントローラ344は、対応する読み出し、書き込み、消去および管理エージェント402、404、406、408を具える。別の実施例においては、制御およびステータスレジスタ340ならびに制御待ち行列342、または、同様の構成は、バンクインタリーブコントローラ344なしで、ソリッドステートストレージ110のバンク214に送信されるデータ用コマンドを待ち行列に入れる。

【0192】

エージェント402、404、406、408は、一実施例においては、特定のバンク214a用に定められた適宜なタイプのコマンドを、バンク214a用の正しい待ち行列に命令する。例えば、読み出しエージェント402は、バンク-1 214b用の読み出しコマンドを受信し、読み出しコマンドをバンク-1読み出し待ち行列410bに命令できる。書き込みエージェント404は、ソリッドステートストレージ110のバンク-0 214a内の位置にデータを書き込むための書き込みコマンドを受信でき、次いで、書き込みコマンドをバンク-0書き込み待ち行列412aに送信する。同様に、消去エージェント406は、バンク-1 214bの消去ブロックを消去する消去コマンドを受信でき、次いで、消去コマンドをバンク-1消去待ち行列414bに送信する。管理エージェント408は、一般的に、バンク-0 214aなど、バンク214の構成レジスタを読み出すためのリセットコマンドまたは要求のような、管理コマンド、ステータス要求等を

10

20

30

40

50

受信する。管理エージェント 408 は、管理コマンドをバンク - 0 管理待ち行列 416 a に送信する。

【0193】

更に、エージェント 402、404、406、408 は、一般的に、待ち行列 410、412、414、416 の状態をモニタし、待ち行列 410、412、414、416 が完全、ほぼ完全、機能しない等である場合にステータス、中断またはその他のメッセージを送信する。一実施例においては、エージェント 402、404、406、408 は、コマンドを受信し、対応するサブコマンドを生成する。一実施例においては、エージェント 402、404、406、408 は、制御及びステータスレジスタ 340 を介したコマンドを受信し、待ち行列 410、412、414、416 に転送される、対応するサブコマンドを生成する。当分野の当業者は、エージェント 402、404、406、408 のその他の機能が分かるであろう。

10

【0194】

待ち行列 410、412、414、416 は、一般的に、コマンドを受信し、ソリッドステートストレージバンク 214 に送信されることが要求されるまでコマンドを記憶する。一般的な実施例においては、待ち行列 410、412、414、416 は、先入れ先出し(「FIFO」)レジスタまたは FIFO として動作する同様の構成である。別の実施例においては、待ち行列 410、412、414、416 は、データにマッチする順番、重要度の順番またはその他の基準でコマンドを記憶する。

【0195】

バンクコントローラ 418 は、一般的に、待ち行列 410、412、414、416 からのコマンドを受信し、適宜なサブコマンドを生成する。例えば、バンク - 0 書き込み待ち行列 412 a は、データパケットのページをバンク - 0 214 a に書き込むコマンドを受信することができる。バンク - 0 コントローラ 418 a は、適宜な時間で書き込みコマンドを受信し、1 又はそれ以上の書き込みサブコマンドを生成することができ、書き込みバッファ 320 に記憶される各データパケットは、バンク - 0 214 a に書き込まれる。例えば、バンク - 0 コントローラ 418 a は、バンク 0 214 a およびソリッドステートストレージレイ 216 の状態を有効にするコマンドを生成し、1 又はそれ以上のデータパケットを書き込む適宜な位置を選択し、ソリッドステートストレージメモリアレイ 216 内の入力バッファをクリアし、1 又はそれ以上のデータパケットを入力バッファに転送し、入力バッファを選択した位置にプログラミングし、データが正しくプログラミングされたことを検証し、プログラムに破損がある場合、マスタコントローラの中断するステップ、同一の物理位置への書き込みを再試行するステップ、異なる物理位置への書き込みを再試行するステップのうちの 1 又はそれ以上を行う。更に、例示の書き込みコマンドとともに、ストレージバスコントローラ 348 は、1 又はそれ以上のコマンドを、ストレージ I/O バス 210 a 用の第 1 の物理アドレスにマッピングされ、ストレージ I/O バス 210 b 用の第 2 の物理アドレスにマッピングされ、その他更に以下に述べるようなコマンドの論理アドレスを有する、ストレージ I/O バス 210 a - n の各々に乗算する。

20

30

【0196】

一般的に、バスアービタ 420 は、バンクコントローラ 418 間から選択され、バンクコントローラ 418 内の出力待ち行列からサブコマンドを引き出し、バンク 214 の性能を最適化する順番にストレージバスコントローラ 348 に転送する。別の実施例においては、バスアービタ 420 は、高位中断信号に応答し、標準的な選択基準を変更することができる。別の実施例においては、マスタコントローラ 224 は、制御およびステータスレジスタ 340 を介してバスアービタ 420 を制御することができる。当分野の当業者は、バスアービタ 420 が、バンクコントローラ 418 からソリッドステートストレージ 110 へのコマンドのシーケンスを制御しインタリーブすることができるその他の手段が分かるであろう。

40

【0197】

50

バスアービタ420は、一般的に、適宜なコマンドと、コマンドタイプが必要とされる
ときは、バンクコントローラ418の対応するデータと、の選択を調整し、コマンド及び
データをストレージバスコントローラ348に送信する。一般的に、バスアービタ420
は、更に、適宜なバンク214を選択するのに、ストレージ制御バス212にコマンドを
送信する。フラッシュメモリ、または、非同期式双方向直列ストレージI/Oバス210
を有するその他のソリッドステートストレージ110の場合、1つのコマンド(制御情報
)またはデータセットのみが、同時に送信可能である。例えば、書き込みコマンドまたは
データが、ストレージI/Oバス210のソリッドステートストレージ110に送信され
た場合、読み出しコマンド、読み出されたデータ、消去コマンド、管理コマンドまたはそ
の他のステータスコマンドは、ストレージI/Oバス210に送信できない。例えば、デ
ータがストレージI/Oバス210から読み出された場合、データはソリッドステートス
トレージ110に書き込むことができない。

10

【0198】

例えば、バンク-0の書き込み動作中、バスアービタ420は、バンク-0コントロー
ラ418aを選択し、このコントローラ418aは、ストレージバスコントローラ348
に次のシーケンスを実行させる待ち行列のトップに書き込みコマンドまたは書き込みサブ
コマンドの配列を有することができる。バスアービタ420は、書き込みコマンドをスト
レージバスコントローラ348に転送し、書き込みコマンドを、ストレージ制御バス21
2を介してバンク-0 214aを選択し、バンク-0 214aに関連したソリッドス
テートストレージエレメント110の入力バッファをクリアするコマンドを送信し、バン
ク-0 214aに関連したソリッドステートストレージエレメント216、218、2
20のステータスを有効にするコマンドを送信することによって設定する。次いで、スト
レージバスコントローラ348は、ストレージI/Oバス210に書き込みサブコマンド
を送信し、論理消去ブロックアドレスからマッピングされる時に、個々の物理消去ソリ
ッドステートストレージエレメント216a-m用の論理消去ブロックのアドレスを含む
、物理アドレスを具える。次いで、ストレージバスコントローラ348は、MUX350
を介したストレージI/Oバス210に、書き込み同期バッファ308を介した書き込み
バッファ320を多重化し、書き込みデータを適宜なページにストリーミングする。ペー
ジが完全である場合、次いで、ストレージバスコントローラ348は、バンク-0 21
4aに関連するソリッドステートストレージエレメント216a-mに、入力バッファを
、ソリッドステートストレージエレメント216a-m内のメモセルにプログラミング
させる。最終的に、ストレージバスコントローラ348は、状態を有効にして、ページが
正しくプログラムされていることを確認する。

20

30

【0199】

読み出し動作は、上述の書き込み例と同様である。読み出し動作中、一般的に、バスア
ービタ420またはバンクインタリーブコントローラ344のその他の構成は、データと
、対応するステータス情報を受信し、データを読み出しデータパイプライン108に送信
し、ステータス情報を制御およびステータスレジスタ340に送信する。一般的に、バス
アービタ420からストレージバスコントローラ348に転送される読み出しデータコマ
ンドは、MUX350に、ストレージI/Oバス210の読み出しデータを読み出しデー
タパイプライン108にゲート制御し、ステータスMUX422を介してステータス情報
を適宜な制御及びステータスレジスタ340に送信する。

40

【0200】

バスアービタ420は、様々なコマンドタイプとデータアクセスモードとを調整し、適
宜なコマンドタイプまたは対応するデータのみが、所定時間に、バスに存在している。バ
スアービタ420が、書き込みコマンドを選択しており、書き込みサブコマンドと対応す
るデータとが、ソリッドステートストレージ110に書き込まれている場合、バスアービ
タ420は、ストレージI/Oバス210のその他コマンドタイプを許可しない。有利に
は、バスアービタ420は、バンク214ステータスに関して受信したステータス情報に
沿った、予測したコマンド実行時間のようなタイミング情報を用い、バスのアイドル時間

50

を最小化または除去する目的として、バスの様々なコマンドの実行を調整する。

【0201】

バスアービタ420を介したマスタコントローラ224は、一般的に、ステータス情報とともに、待ち行列410、412、414、416に記憶されたコマンドの予測完了時間を用いて、コマンドに関連したサブコマンドが1つのバンク214aで実行されているとき、他のコマンドの他のサブコマンドが、他のバンク214b-nで実行される。1つのコマンドが、214aで完全に実行されるとき、バスアービタ420は、別のコマンドをバンク214aに命令する。バスアービタ420は更に、待ち行列410、412、414、416に記憶されたコマンドと、待ち行列410、412、414、416に記憶されていない他のコマンドとを調整することができる。

10

【0202】

例えば、消去コマンドは、ソリッドステートストレージ110内の消去ブロックのグループを消去するように送信可能である。消去コマンドは、書き込みまたは読み出しコマンドを実行するよりも10乃至1000倍以上の時間がかかり、プログラムコマンドを実行するよりも10乃至100倍以上の時間がかかる。Nバンク214のために、バンクインタリーブコントローラ344は、消去コマンドをNコマンドに分割することができ、各々は、バンク214aの仮想消去ブロックを消去する。バンク-0 214aが消去コマンドを実行している間、バスアービタ420は、その他のバンク214b-nで実行されるその他のコマンドを選択することができる。更に、バスアービタ420は、ストレージバスコントローラ348、マスタコントローラ224等の他の構成とともに動作することができ、バス間のコマンド実行を調整する。バンクインタリーブコントローラ344のバスアービタ420、バンクコントローラ418、待ち行列410、412、414、416およびエージェント402、404、406、408を用いてコマンド実行を調整することで、バンクインタリーブ機能なしで、その他のソリッドステートストレージシステムを介して、大幅に性能を向上させることができる。

20

【0203】

一実施例においては、ソリッドステートストレージコントローラ104は、ソリッドステートストレージ110のストレージエレメント216、218、220の総てに働く1つのバンクインタリーブコントローラ344を具える。別の実施例においては、ソリッドステートコントローラ104は、ストレージエレメント216a-m、218a-m、220a-mの各列用のバンクインタリーブコントローラ344を具える。例えば、1つのバンクインタリーブコントローラ344は、ストレージエレメントSSS0.0-SSS0.N216a、218a、220aの一行に働き、第2のバンクインタリーブコントローラ344は、ストレージエレメントSSS1.0-SSS1.N216b、218b、220bなどの第2の列に働く。

30

【0204】

図4Bは、本発明によるソリッドステートストレージコントローラ内のバンクインタリーブコントローラの代替的な実施例401を示した概略ブロック図である。図4Bに示された実施例に示された構成210、212、340、346、348、350、402-430は、各バンク214が、単一の待ち行列432a-nを具えて、バンク(例えば、バンク0 214a)用の読み出しコマンド、書き込みコマンド、消去コマンド、管理コマンドが、バンク214aの単一の待ち行列432aに対して命令されることを除いては図4Aに関連して示されたバンクインタリーブ装置400と実質的に同じである。一実施例においては、待ち行列432は、FIFOである。別の実施例においては、待ち行列432は、待ち行列432が記憶されている順番とは別の順番で、待ち行列432から引き出されたコマンドを有することができる。別の代替的な実施例においては(図示せず)、読み出しエージェント402、書き込みエージェント404、消去エージェント406および管理エージェント408は、適宜な待ち行列432a-nに対する単一のエージェント割当コマンドに組み込むことができる。

40

【0205】

50

別の代替的な実施例においては（図示せず）、コマンドは、1つの待ち行列に記憶され、これらのコマンドは、記憶された方法以外の順番で待ち行列から引き出し可能であり、バンクインタリーブコントローラ344は、1つのバンク214aでコマンドを実行でき、他のコマンドは、残りのバンク214b-nで実行される。当分野の当業者であれば、他のコマンドが他のバンク214b-nで実行されている間に、1つのバンク214aでコマンドを実行できる、他の待ち行列構成およびタイプが容易に分かるであろう。

【0206】

ストレージ特有の構成

【0207】

ソリッドストレージコントローラ104は、同期バッファ346を具え、この同期バッファ346は、ソリッドステートストレージ110から送信されて受信されたコマンドおよびステータスメッセージをバッファリングする。同期バッファ346は、ソリッドステートストレージクロックドメインとローカルバスクロックドメインとの間の境界に配置され、クロックドメインの差を考慮するバッファリングを提供する。同期バッファ346、書き込み同期バッファ308および読み出し同期バッファ328は独立にでき、一緒に動作させて、データ、コマンド、ステータスメッセージなどをバッファリングすることができる。好ましい実施例においては、同期バッファ346は、クロックドメインを交差する信号が最小数のところに配置される。当分野の当業者は、クロックドメイン間の同期が、設計の実装のいくつかの態様を最適にするために、ソリッドステートストレージデバイス102内の他の位置に任意に移動可能であることが分かるであろう。

10

20

【0208】

ソリッドステートストレージコントローラ104は、ストレージバスコントローラ348を具え、このストレージバスコントローラ348は、ソリッドステートストレージ110に送信され、読み出されるデータに対するコマンドと、ソリッドステートストレージ110のタイプに基づいてソリッドステートストレージ110から受信したステータスメッセージとを解釈し翻訳する。例えば、ストレージバスコントローラ348は、性能特性の異なるストレージ、製造者の異なるストレージなど、異なるタイプのストレージに対する異なるタイミング要求を具えることができる。更に、ストレージバスコントローラ348は、ストレージ制御バス212に制御コマンドを送信する。

30

【0209】

好ましい実施例においては、ソリッドステートストレージコントローラ104は、MUX350を具え、MUX350は、マルチプレクサ350a-nのアレイを具え、各マルチプレクサは、ソリッドステートストレージアレイ110の列の専用となる。例えば、マルチプレクサ350aは、ソリッドステートストレージエレメント216a、218a、220aに関連している。MUX350は、書き込みデータパイプライン106からのデータと、ストレージI/Oバス210を介したストレージバスコントローラ348からソリッドステートストレージ110へのコマンドをルーティングし、ストレージI/Oバス210を介したソリッドステートストレージ110から、読み出しデータパイプライン106およびストレージバスコントローラ348、同期バッファ346およびバンクインタリーブコントローラ344を介した制御およびステータスレジスタ340へのデータおよびステータスメッセージをルーティングする。

40

【0210】

好ましい実施例においては、ソリッドステートストレージコントローラ104は、ソリッドステートストレージエレメントの各列（例えば、SSS0.1 216a、SSS0.2 218a、SSS0.N 220a）用のMUX350を具える。MUX350は、書き込みデータパイプライン106からのデータと、ストレージI/Oバス210を介してソリッドステートストレージ110へ送信されるコマンドを組合せ、コマンドから読み出しデータパイプライン108によって処理されるようにデータを分割する。書き込みバッファ320に記憶されたパケットは、書き込みバッファ320の外部から、ソリッドステートストレージエレメントの各列（SSSx.0からSSSx.N 216、218、

50

220)用の書き込み同期バッファ308を介し、ソリッドステートストレージエレメントの各列(SSSx.0からSSSx.N216、218、220)用のMUX350までのバスに命令される。コマンドおよび読み出しデータは、ストレージI/Oバス210からMUX350によって受信される。更に、MUX350は、ステータスメッセージをストレージバスコントローラ348に命令する。

【0211】

ストレージバスコントローラ348は、マッピングモジュール424を具える。マッピングモジュール424は、消去ブロックの1又はそれ以上の物理アドレスに消去ブロックの論理アドレスをマッピングする。例えば、ブロック214aにつき20のストレージエレメント(例えば、SSS0.0からSSSM.0216)のアレイを有するソリッドステートストレージ110は、ストレージエレメントにつき1の物理アドレスで、消去ブロックの20の物理アドレスにマッピングされる特定の消去ブロック用の論理アドレスを有することができる。ストレージエレメントが並列にアクセスされるので、ストレージエレメント216a、218a、220aの列の各ストレージエレメントにおける同一位置の消去ブロックは、物理アドレスを共有する。列中の総ての消去ブロック(例えば、ストレージエレメントSSS0.0、0.1、...、0.N216a、218a、220aにおいて)の代わりに、1の消去ブロックを選択するために、1つのバンク(この場合、バンク-0214a)が選択される。

10

【0212】

このような消去ブロック用の論理-物理マッピングは、1つの消去ブロックが損傷またはアクセス不能になった場合、マッピングが別の消去ブロックにマッピングされるように変更可能であるので、有益である。これによって、1つのエレメントの消去ブロックが失敗したとき、仮想消去ブロック全体を失う損失を緩和する。再マッピングモジュール430は、消去ブロックの論理アドレスのマッピングを、仮想消去ブロックの1又はそれ以上の物理アドレス(ストレージエレメントのアレイにわたって配置される)に変更する。例えば、仮想消去ブロック1は、ストレージエレメントSSS0.0216aの消去ブロック1、ストレージエレメントSSS1.0216bの消去ブロック1、...、ストレージエレメントM.0216mにマッピング可能であり、仮想消去ブロック2は、ストレージエレメントSSS0.1218aの消去ブロック2、ストレージエレメントSSS1.1218bの消去ブロック2、...、ストレージエレメントM.1218mにマッピング可能等である。

20

30

【0213】

ストレージエレメントSSS0.0216aの消去ブロック1が損傷し、摩耗などによるエラーが生じるか、何らかの理由で用いられない場合、再マッピングモジュールは、仮想消去ブロック1の消去ブロックをポインティングした論理アドレス用に、論理-物理マッピングを変更することができる。ストレージエレメントSSS0.0216aのスペア消去ブロック(消去ブロック221と呼ぶ)が利用可能であり、その時点でマッピングされていない場合、再マッピングモジュールは、ストレージエレメントSSS0.0216aの消去ブロック221をポインティングするために、仮想消去ブロック1の消去ブロックのマッピングを変更でき、ストレージエレメントSSS1.0216bの消去ブロック1、ストレージエレメントSSS2.0(図示せず)の消去ブロック1、...、ストレージエレメントM.0216mの消去ブロック1のポインティングを継続する。マッピングモジュール424または再マッピングモジュール430は、上述した順番(ストレージエレメントの消去ブロック1に対する仮想消去ブロック1、ストレージエレメントの消去ブロック2に対する仮想消去ブロック2)で消去ブロックをマッピングでき、あるいは、その他の基準に基づいた別の順番でストレージエレメント216、218、220の消去ブロックをマッピングできる。

40

【0214】

一実施例においては、消去ブロックは、アクセス時間によってグループ分け可能である。データを特有の消去ブロックのページにプログラミング(書き込み)するなど、コマン

50

ドを実行する時間を意味する、アクセス時間によるグループ分けは、コマンド完了を平均にすることができ、仮想消去ブロックの消去ブロックにわたって実行されるコマンドは、最も遅い消去ブロックによって制限されることはない。他の実施例においては、消去ブロックは、摩耗レベル、健全などによってグループ分け可能である。当分野の当業者は、消去ブロックをマッピングまたは再マッピングするときに考慮すべきその他のファクタが分かるであろう。

【0215】

一実施例においては、ストレージバスコントローラ348は、ステータスキャプチャモジュール426を具え、このステータスキャプチャモジュール426は、ソリッドステートストレージ110からのステータスメッセージを受信し、ステータスメッセージをステータスMUX422に送信する。別の実施例においては、ソリッドステートストレージ110がフラッシュメモリである場合、ストレージバスコントローラ348は、NANDバスコントローラ428を具える。NANDバスコントローラ428は、読み出しおよび書き込みデータパイプライン106、108からソリッドステートストレージ110内の正しい位置にコマンドを命令し、フラッシュメモリ等の特性に基づいてコマンド実行のタイミングを調整する。ソリッドステートストレージ110が別のタイプのソリッドステートストレージ型である場合、NANDバスコントローラ428は、ストレージタイプに特有のバスコントローラによって置換できる。当分野の当業者は、NANDバスコントローラ428のその他の機能が分かるであろう。

10

【0216】

フローチャート

20

【0217】

図5Aは、本発明によるデータパイプラインを用いてソリッドステートストレージデバイス102におけるデータ管理のための方法500の一実施例を示した概略フローチャート図である。この方法500は開始して(502)、入力バッファ306が、ソリッドステートストレージ110に書き込まれる1又はそれ以上のデータセグメントを受信する(504)。1又はそれ以上のデータセグメントは、一般的に、オブジェクトの少なくとも一部を具えるが、オブジェクト全体であってもよい。パケットサイズ302は、オブジェクトと共に1又はそれ以上のオブジェクト特有のパケットを生成できる。パケットサイズ302は、各パケットにヘッダを加え、一般的に、各パケットは、パケット長およびオブジェクト内のパケット用のシーケンス番号を具える。パケットサイズ302は、入力バッファ306内に記憶される1又はそれ以上のデータまたはメタデータセグメントを受信し(504)、ソリッドステートストレージ110用にサイズを調整された1又はそれ以上のパケットを生成することによって1又はそれ以上のデータまたはメタデータセグメントをパケット化し(506)、各パケットは、1又はそれ以上のセグメントから1つのヘッダおよびデータを具えている。

30

【0218】

一般的に、第1のパケットは、パケットが生成されたオブジェクトを識別するオブジェクト識別子を具える。第2のパケットは、ソリッドステートストレージデバイス102によって用いられる情報を有するヘッダを具え、第2のパケットは、第1のパケット中で識別されたオブジェクトに関連することができ、オフセット情報は、オブジェクト内の第2のパケットおよびデータを配置する。このソリッドステートストレージデバイスコントローラ202は、バンク214と、パケットがストリーミングされる物理領域を管理する。

40

【0219】

ECCジェネレータ304は、パケットサイズ302からのパケットを受信し、データパケット用のECCを生成する(508)。一般的に、パケットとECCブロックとの間には固定した関連性はない。ECCブロックは、1又はそれ以上のパケットを具えることができる。パケットは、1又はそれ以上のECCブロックを具えることができる。パケットは、ECCブロック内のいずれの場所でも開始および終了することができる。パケットは、第1のECCブロック内のいずれかの場所からスタートすることができ、次のECCブ

50

ロックのいずれかの場所で終了することができる。

【0220】

書き込み同期バッファ308は、ECCブロックをソリッドステートストレージ110に書き込む前に、対応するECCブロック内に配置されるように、パケットをバッファリングし(510)、次いで、ソリッドステートストレージコントローラ104は、クロックドメイン差を考慮した適宜の時間でデータを書き込み(512)、方法500は終了する(514)。書き込み同期バッファ308は、ローカルクロックドメインとソリッドステートストレージ110のクロックドメインとの間の境界に配置される。この方法500は、便宜上、1又はそれ以上のデータセグメントを受信するステップと、1又はデータパケットを書き込むステップと、を記載しているが、一般的に、データセグメントのストリームが受信され、グループであることに留意されたい。一般的に、ソリッドステートストレージ110の完全な仮想ページを具える多数のECCブロックが、ソリッドステートストレージ110に書き込まれる。一般的に、パケットサイズ302は、あるサイズのデータセグメントを受信し、別のサイズのパケットを生成する。これは必要に応じて、データ若しくはメタデータセグメント、または、データ若しくはメタデータセグメントの一部が、セグメントのデータ総てがパケットにキャプチャされるように、データパケットを形成するために組み合わせられることを要求する。

10

【0221】

図5Bは本発明によるインサーバSANのための方法の一実施例を示した概略フローチャート図である。方法500は開始し(552)、ストレージ通信モジュール162は、第1のストレージコントローラ152aと、第1のサーバ112aの外部にある少なくとも1のデバイスとの間の通信を促進する(554)。第1のストレージコントローラ152aと外部デバイスとの間の通信は、第1のサーバ112aから独立している。第1のストレージコントローラ112aは第1のサーバ112a内にあり、第1のストレージコントローラ152は少なくとも1のストレージデバイス154aを制御する。第1のサーバ112aは、第1のサーバ112aと第1のストレージコントローラ152aと共に同一位置に配置されるネットワークインタフェース156aを具える。インサーバSANモジュール164はストレージ要求を提供し(556)、方法501は終了する(558)。インサーバSANモジュールはネットワークプロトコル及び/又はバスプロトコルを用いてストレージ要求を提供する(556)。インサーバSANモジュール164は第1のサーバ112aと別個のストレージ要求を提供し(556)、サービス要求はクライアント114、114aから受信される。

20

30

【0222】

図6は、本発明によるデータパイプラインを用いてソリッドステートストレージデバイス102内のデータを管理する方法600の別の実施例を示した概略フローチャート図である。この方法600は開始(602)し、入力バッファ306が、ソリッドステートストレージ110に書き込まれる1又はそれ以上のデータまたはメタデータセグメントを受信する(604)。パケットサイズ302は、ヘッダを各パケットに加え、各パケットは、オブジェクト内のパケット長を一般的に具える。パケットサイズ302は、入力バッファ306に記憶された1又はそれ以上のセグメントを受信し(604)、ソリッドステートストレージ110用にサイズを調整された1又はそれ以上のパケットを生成することによって、1又はそれ以上のセグメントをパケット化し(606)、各パケットは、ヘッダと、1又はそれ以上のセグメントからのデータとを具える。

40

【0223】

ECCジェネレータ304は、パケットサイズ302からのパケットを受信し、パケット用の1又はそれ以上のECCブロックを生成する(608)。書き込み同期バッファ308は、ECCブロックをソリッドステートストレージ110に書き込む前に、対応するECCブロック内に配置されるようにパケットをバッファリングし(610)、次いで、ソリッドステートストレージコントローラ104は、クロックドメイン差を考慮して適宜な時間でデータを書き込む(612)。データが、ソリッドステートストレージ110から

50

要求されるとき、1又はそれ以上のデータパケットを具えるECCは、読み出し同期バッファ328に読み出され、バッファリングされる(614)。パケットのECCブロックは、ストレージI/Oバス210を介して受信される。ストレージI/Oバス210は双方向であるので、データが読み出されるとき、書き込み動作、コマンド動作が中止される。

【0224】

ECC修正モジュール322は、読み出し同期バッファ328で保持される要求されたパケットのECCブロックを受信し、必要に応じて各ECCブロック内のエラーを修正する(616)。ECC修正モジュール322は、ECCブロックに1又はそれ以上のエラーが存在し、そのエラーがECCシンドロームを用いて修正可能であるか判定し、ECC修正モジュール322が、ECCブロックのエラーを修正する(616)。ECC修正モジュール322が、検出されたエラーがECCを用いて修正できないと判定した場合は、ECC修正モジュール322は、中断信号を送信する。

10

【0225】

デパケットタイザ324は、ECC修正モジュール322がエラーを修正した後、要求されたパケットを受信し(618)、各パケットのパケットヘッダをチェックおよび除去することによってパケットを脱パケット化する(618)。アライメントモジュール326は、脱パケット化後、パケットを受信し、所望されないデータを除去し、セグメントまたはオブジェクトを要求するデバイスと互換性のある形態においてオブジェクトのデータまたはメタデータセグメントとしてデータパケットを再フォーマットする(620)。出力バッファ330は、脱パケット化後要求されるパケットを受信し、要求デバイスに送信する前にパケットをバッファリングし(622)、この方法600は終了する(624)。

20

【0226】

図7は、本発明によるバンクインタリーブを用いて、ソリッドステートストレージデバイス102のデータを管理する方法700に関する一実施例を示した概略フローチャート図である。この方法600は開始し(602)、バンクインタリーブコントローラ344は、1又はそれ以上のコマンドを2又はそれ以上の待ち行列410、412、414、416に命令する(604)。一般的に、エージェント402、404、406、408は、コマンドタイプによってコマンドを待ち行列410、412、414、416に命令する。待ち行列410、412、414、416の各セットは、各コマンドタイプの待ち行列を具える。バンクインタリーブコントローラ344は、バンク214間で、待ち行列410、412、414、416に記憶されたコマンドの実行を調整し(606)、第1のタイプのコマンドは、1つのバンク214aで実行され、第2のタイプのコマンドは、第2のバンク214bで実行され、この方法600は終了する(608)。

30

【0227】

ストレージ空き領域回復

【0228】

図8は、本発明によるソリッドステートストレージデバイス102のガベージコレクションのための装置800に関する一実施例を示した概略ブロック図である。この装置800は、シーケンシャルストレージモジュール802、ストレージ部分選択モジュール804、データ回復モジュール806、および、ストレージ部分回復モジュール808を具え、これらは以下で説明される。他の実施例においては、この装置800は、ガベージ標識モジュール810および消去モジュール812を具える。

40

【0229】

装置800は、ストレージ部分内のページにデータパケットを順番に書き込むシーケンシャルストレージモジュール802を具える。パケットは、パケットが新規パケットまたは変更されたパケットであるかにかかわらず順番に記憶される。変更されたパケットは本実施例においては一般的に、既に記憶された位置に戻って書き込まれない。一実施例においては、シーケンシャルストレージモジュール802は、ストレージ部分のページ内の第1の位置にパケットを書き込み、次いで、ページの次の位置にパケットを書き込み、それ

50

から、次、次と、ページが満たされるまで続けられる。次いで、シーケンシャルストレージモジュール 802 は、ストレージ部分内の次のページを満たし始める。これは、ストレージ部分が満たされるまで続く。

【0230】

好ましい実施例においては、シーケンシャルストレージモジュール 802 は、バンク（バンク - 0 214 a）のストレージエレメント（例えば、SSS0.0 乃至 SSSM.0 216）内のストレージ書き込みバッファにパケットを書き込むことを開始する。ストレージ書き込みバッファが完全な場合、ソリッドステートストレージコントローラ 104 は、ストレージ書き込みバッファのデータを、バンク 214 a のストレージエレメント 216 内の指定されたページにプログラミングできる。次いで、別のバンク（例えば、バンク - 1 214 b）が選択され、シーケンシャルストレージモジュール 802 は、バンク 214 b のストレージエレメント 218 のストレージ書き込みバッファにパケットを書き込むことを開始し、第 1 のバンク - 0 は、指定されたページをプログラミングする。このバンク 214 b のストレージ書き込みバッファが完全な場合、ストレージ書き込みバッファの内容は、各ストレージエレメント 218 の別の指定ページにプログラミングされる。1 つのバンク 214 a をページにプログラミングする間に、別のバンク 214 b のストレージ書き込みバッファが満たされうるので、この処理は有効である。

【0231】

ストレージ部分は、ソリッドステートストレージデバイス 102 のソリッドステートストレージ 110 の一部を具える。一般的に、ストレージ部分は消去ブロックである。フラッシュメモリについては、消去ブロックの消去動作は、各セルを充電することによって、1 を消去ブロックの全ビットに書き込む。これは、位置が総て 1 として開始するプログラム動作と比較して、冗長的なプロセスであり、データが書き込まれるとき、いくつかのビットは、ゼロで書き込まれるセルを放電することによってゼロに変更される。しかしながら、ソリッドステートストレージ 110 がフラッシュメモリでないか、あるいは、消去サイクルが、読み出しまたはプログラム等の、その他の動作として同様の時間を費やすフラッシュメモリを有する場合は、ストレージ部分は、消去されることが要求されない。

【0232】

本明細書で用いられるように、ストレージ部分は、消去ブロックに対する領域と同等であるが、消去されても消去されなくてもよい。本明細書で消去ブロックが用いられる場合、消去ブロックは、ストレージエレメント（例えば、SSS0.0 216 a）内の指定されたサイズの特定期間を指し、一般的に、所定量のページを具える。「消去ブロック」は、フラッシュメモリと共に用いられるとき、一般的に、書き込まれる前に消去されるストレージ部分である。「消去ブロック」は、「ソリッドステートストレージ」とともに用いられる場合、消去されても消去されなくてもよい。本明細書で用いられるように、消去ブロックは、1 つの消去ブロック、または、ストレージエレメント（例えば、SSS0.0 乃至 SSSM.0 216 a - n）の各列の消去ブロックを有する消去ブロックのグループ、を具えることができ、消去ブロックは、仮想消去ブロックも意味する。仮想消去ブロックと関連した論理構成を指す場合、消去ブロックは、論理消去ブロック（「LEB」）も意味する。

【0233】

一般的に、パケットは、処理の順番によって、順々に記憶される。一実施例においては、書き込みデータパイプライン 106 が用いられる場合、シーケンシャルストレージモジュール 802 は、パケットが書き込みデータパイプライン 106 からくる順番でパケットを記憶する。有効データは、下記で説明するように、回復動作中にストレージ部分から回復されるので、上記順番は、別のストレージ部分から読み出される有効データのパケットと混合した、要求デバイス 155 から到達するデータセグメントの結果でもよい。回復された有効なデータパケットを書き込みデータパイプライン 106 に再ルーティングすることは、図 3 のソリッドステートストレージコントローラ 104 に関連して上述したように、ガベージコレクタバイパス 316 を具えてもよい。

10

20

30

40

50

【 0 2 3 4 】

装置 8 0 0 は、回復用ストレージ部分を選択するストレージ部分選択モジュール 8 0 4 を具える。回復用ストレージ部分を選択することで、データを書き込むシーケンシャルストレージモジュール 8 0 2 によってストレージ部分を再利用することができ、従って、回復されたストレージ部分をストレージプールに加えるか、または、ストレージ部分がない、信頼性がない、リフレッシュすべき、あるいは、ストレージプールから一時的または永続的にストレージ部分を外すと判定した後、ストレージ部分から有効データを回復する。別の実施例においては、ストレージ部分選択モジュール 8 0 4 は、無効データの多いストレージ部分または消去ブロックを識別することによって、回復用ストレージ部分を選択する。

10

【 0 2 3 5 】

別の実施例においては、ストレージ部分選択モジュール 8 0 4 は、摩耗の少ないストレージ部分または消去ブロックを識別することによって回復用ストレージ部分を選択する。例えば、摩耗の少ないストレージ部分または消去ブロックを識別することは、無効データの少ない、消去サイクルの数が少ない、ビットエラー率の低い、または、プログラムカウンタの低い（回数の低い、パuffa中のデータのページは、ストレージ部分のページに書き込まれ、プログラムカウンタは、デバイスが製造されたときから、ストレージ部分が最後に消去されたときから、その他の任意の事象から、および、これらの組み合わせから、測定されてもよい）ストレージ部分を識別することを含む。このストレージ部分選択モジュール 8 0 4 は、摩耗量の少ないストレージ部分を判定するために、上記パラメータまたは他のパラメータの組み合わせを用いることができる。摩耗量の低いストレージ部分を判定することによる、回復用ストレージ部分の選択は、用いられているストレージ部分を見つけるのに所望され、摩耗レベルなどに合わせて回復可能である。

20

【 0 2 3 6 】

別の実施例においては、ストレージ部分選択モジュール 8 0 4 は、摩耗量の多いストレージ部分または消去ブロックを識別することによって回復用ストレージ部分を選択する。例えば、摩耗量の多いストレージ部分または消去ブロックを識別することは、消去サイクル数の多いストレージ部分、ビットエラー率の高いストレージ部分、回復不能 E C C ブロックを有するストレージ部分、または、プログラムカウンタの高いストレージ部分を識別することを含む。更に、ストレージ部分選択モジュール 8 0 4 は、摩耗量の多いストレージ部分を判定する上記または他のパラメータの組み合わせを用いることができる。摩耗量の多いストレージ部分を判定することによって回復用ストレージ部分を選択することは、過剰に用いられているストレージ部分を見つけるのに所望され、消去サイクルなどを用いてストレージ部分をリフレッシュすることによって回復可能であり、または、利用不能としてサービスからストレージ部分を廃棄することができる。

30

【 0 2 3 7 】

装置 8 0 0 は、データ回復モジュール 8 0 6 を具え、このデータ回復モジュール 8 0 6 は、回復用に選択されたストレージ部分から有効なデータバケットを読み出し、他のデータバケットを有する有効なデータバケットを待ち行列に入れてシーケンシャルストレージモジュール 8 0 2 によって順々に書き込み、シーケンシャルストレージモジュール 8 0 2 によって書き込まれた有効データの新規物理アドレスを有するインデックスを更新する。一般的に、このインデックスは、データオブジェクトから由来するバケットがソリッドステートストレージ 1 1 0 に記憶される物理アドレスに、オブジェクトのデータオブジェクト識別子をマッピングするオブジェクトインデックスである。

40

【 0 2 3 8 】

一実施例においては、装置 8 0 0 は、ストレージ部分回復モジュール 8 0 8 を具え、このストレージ部分回復モジュール 8 0 8 は、使用または再使用のためのストレージ部分を準備し、データ回復モジュール 8 0 6 がストレージ部分から有効データをコピーするのが完了した後、データバケットを順々に書き込むために、利用可能なストレージ部分をシーケンシャルストレージモジュール 8 0 2 に標識する。別の実施例においては、装置 8 0 0

50

は、データを記憶するために利用不可能なとき、回復用に選択されたストレージ部分を標識するストレージ部分回復モジュール 808 を具える。一般的に、これは、ストレージ部分または消去ブロックが、信頼性のあるデータストレージ用に用いられる状況ではないと、摩耗量の多いストレージ部分または消去ブロックを識別するストレージ部分選択モジュール 804 によるものである。

【0239】

一実施例においては、装置 800 は、ソリッドストレージデバイス 102 のソリッドストレージデバイスコントローラ 202 中にある。別の実施例においては、装置 800 は、ソリッドステートストレージデバイスコントローラ 202 を制御する。別の実施例においては、装置 800 の一部が、ソリッドステートストレージデバイスコントローラ 202 中にある。別の実施例においては、データ回復モジュール 806 によって更新されたオブジェクトインデックスは更に、ソリッドステートストレージデバイスコントローラ 202 中に配置される。

10

【0240】

一実施例においては、ストレージ部分は、消去ブロックであり、装置 800 は、データ回復モジュール 806 が、選択された消去ブロックから有効なデータパケットをコピーした後で、ストレージ部分回復モジュール 808 が、利用可能として消去ブロックを標識する前に、回復用に選択された消去ブロックを消去する消去モジュール 810 を具える。読み出しまたは書き込み動作よりも大幅に長い時間を要する消去動作を有するフラッシュメモリおよびその他のソリッドステートストレージのために、新規データを書き込むことを可能にする前にデータブロックを消去することは、効率的な動作として所望される。ソリッドステートストレージ 110 がバンク 214 中に配置される場合、消去モジュール 810 による消去動作は、1つのバンクで実行され、その他のバンクは、読み出し、書き込みまたはその他の動作を実行する。

20

【0241】

一実施例においては、装置 800 は、ガベージ標識モジュール 812 を具え、このモジュール 812 は、データパケットがもはや有効ではないことを示す動作に応じて、ストレージ部分中のデータパケットが無効であると識別する。例えば、データパケットが削除される場合、ガベージ標識モジュール 812 は、データパケットが無効であると識別できる。読み出し - 修正 - 書き込み動作は、データパケットが無効であるとして識別される別の方法である。一実施例においては、ガベージ標識モジュール 812 は、インデックスを更新することによって、データパケットが無効であると識別することができる。別の実施例においては、ガベージ標識モジュール 812 は、無効なデータパケットが削除されたことを示す別のデータパケットを記憶することによって、データパケットが無効であることを識別できる。このことが有利であるのは、ソリッドステートストレージ 110 中に、データパケットが削除されたという情報を記憶することによって、オブジェクトインデックス復元モジュール 262 または同様のモジュールが、無効データパケットが削除されたことを示すエントリを有するオブジェクトインデックスを復元することができるからである。

30

【0242】

一実施例においては、装置 800 は、性能全体を改善するために、消去コマンドに続いて、データの仮想ページの残りの部分を満たすのに使用でき、消去コマンドは、書き込みパイプライン 106 が空になり、総てのパケットが、非揮発性ソリッドステートストレージ 110 内に永続的に書き込まれるまで、書き込みパイプライン 106 にデータが流れるのを中止する。これは、必要とされるガベージコレクションの量、ストレージ部分を消去するのにかかる時間、仮想ページをプログラムするのに必要な時間を低減することができる利点を有する。例えば、ソリッドステートストレージ 100 の仮想ページ内に書き込むために 1つの小さなパケットのみが準備されるときに、消去コマンドが受信される。ほぼ空の仮想ページをプログラミングすることが、浪費された空き領域をすぐに回復する必要を生じさせ、収集された不必要なガベージにして、消去され、回復され、および、シーケンシャルストレージモジュール 802 によって書き込む利用可能な空き領域のプールに戻

40

50

されるストレージ部分内の有効データにする。

【0243】

無効なデータパケットを実際に消去するよりも、データパケットを無効にして標識することが有効であるのは、上述したように、フラッシュメモリおよび他の同様のストレージについては、消去動作が大幅に時間をとるからである。装置800に記載されるように、ガベージコレクションシステムが、ソリッドステートストレージ110内で自律的に動作できることが、読み出し、書き込みおよびその他のより迅速な動作から消去動作を分ける方法を提供することができ、ソリッドステートストレージデバイス102は、他の多くのソリッドステートストレージシステムまたはデータストレージデバイスよりも迅速に動作できる。

10

【0244】

図9は、本発明によるストレージ回復のための方法900の一実施例を示す概略フローチャート図である。方法900は開始し(902)、シーケンシャルストレージモジュール802は、ストレージ部分にデータパケットを順々に書き込む(904)。ストレージ部分は、ソリッドステートストレージデバイス102内のソリッドステートストレージ110の一部である。一般的に、ストレージ部分は消去ブロックである。データパケットは、オブジェクト由来であり、データパケットは、処理の順番によって順々に記憶される。

【0245】

ストレージ部分選択モジュール804は、回復用のストレージ部分を選択し(906)、データ回復モジュール806は、回復用に選択されたストレージ部分から有効なデータパケットを読み出す(908)。一般的に、有効なデータパケットは、消去または削除用に標識されていないデータパケットであるか、標識されているその他の無効データであり、有効または「優良」データとみなされる。データ回復モジュール806は、シーケンシャルストレージモジュール802によって順々に書き込まれるようにスケジューリングされた他のデータパケットを有する有効なデータパケットを待ち行列に入れる(910)。データ回復モジュール806は、シーケンシャルストレージモジュール802によって書き込まれる有効データの新規物理アドレスを有するインデックスを更新する(912)。このインデックスは、オブジェクト識別子に対するデータパケットの物理アドレスのマッピングを具える。データパケットは、ソリッドステートストレージ110中に記憶される際に記憶されるもの、およびデータパケットに対応するオブジェクト識別子である。

20

30

【0246】

データ回復モジュール806は、ストレージ部分からの有効データをコピー完了後、ストレージ部分回復モジュール808は、データパケットを順々に書き込むシーケンシャルストレージモジュール802が利用可能であると、回復用に選択されたストレージ部分を標識し(914)、方法900は終了する(916)。

【0247】

空データセグメント指令

【0248】

一般的に、データがもはや有用できないとき、そのデータは消去可能である。多くのファイルシステムでは、消去コマンドは、ファイルシステム中のディレクトリエントリを削除し、データを含むストレージデバイスの所定位置にあるデータを代わりに残している。一般的に、データストレージデバイスは、この種の消去動作に含まれない。データを消去する別の方法は、ゼロ、1またはその他のヌルデータキャラクタをデータストレージデバイスに書き込み、実際に消去されたファイルに置換することである。しかしながら、これが非効率なのは、データを送信することが上書きされることであるため、有効なバンド幅が用いられるからである。更に、ストレージデバイス中の空き領域は、無効データを上書きするのに用いられるデータによって占有される。

40

【0249】

ここで記載されるソリッドステートストレージデバイス102のようないくつかのストレージデバイスは、ランダムアクセスストレージデバイスではないので、既に記憶された

50

データを更新することは既存のデータを上書きすることではない。上記デバイス上で、1の記号列またはゼロの記号列を有するデータを上書きする試みは、既存のデータを上書きする所望の意図を完了することなく、有用な空き領域を埋める。ソリッドステートストレージデバイス102のような、これらの非ランダムアクセスデバイスについては、クライアント114は、一般的に、データを消去するためにデータを上書きする能力を有しない。

【0250】

反復したキャラクタ列またはキャラクタ列を受信するときに、受信したデータは、高度に圧縮可能であるが、一般的に、圧縮は、ストレージデバイスに送信する前に、ファイルシステムによってなされる。一般的なストレージデバイスは、圧縮データと非圧縮データとを区別できない。更に、ストレージデバイスは、消去されたファイルを読み出すコマンドを受信することができ、ストレージデバイスは、ゼロ、1またはヌルキャラクタの記号列を要求デバイスに送信することができる。更に、バンド幅は、消去されたファイルを示すデータを送信するのに要求される。

10

【0251】

上述した説明から、ストレージデバイスが、空データセグメントまたは反復したキャラクタ若しくはキャラクタ列を有するデータを表すデータセグメントトークンを記憶することができるように、データを消去すべきという指令を受信するストレージデバイスのための装置、システムおよび方法に対するニーズが存在することは明らかである。装置、システムおよび方法は、更に、既存のデータを消去することができ、結果として用いられたストレージ空き領域は、小さなデータセグメントトークンを持つ。従来技術の欠点の一部または全部を克服する装置、システムおよび方法が示される。

20

【0252】

図10は、本発明によるトークン指令を生成する装置を有するシステム1000の一実施例を示した概略ブロック図である。この装置は、トークン指令生成モジュール1002、トークン指令送信モジュール1004、読み出し受信モジュール1006、読み出し要求送信モジュール1008、読み出しトークン指令受信モジュール1010、要求クライアント応答モジュール1012およびデータセグメント再生モジュール1014を具備し、これらは以下で説明する。一実施例においては、この装置は、ストレージコントローラ152およびデータストレージデバイス154を有するストレージデバイス150に接続したサーバ112内にあり、これらは、上述したものと実質的に同等である。

30

【0253】

一実施例においては、この装置は、トークン指令を有するストレージ要求を生成するトークン指令生成モジュール1002を具備する。トークン指令は、ストレージデバイス150上のデータセグメントトークンを記憶する要求を具備する。トークン指令は、ストレージデバイス150に送信されて、データセグメントトークンが所定位置に送信されていなかった場合にデータセグメントとして記憶される、反復する同一キャラクタの配列、又は反復する同一キャラクタ列の配列に置換されるように意図される。一実施例においては、反復する同一キャラクタの配列は、データセグメントが空であることを示唆する。例えば、反復する同一キャラクタの配列は、ゼロまたは1であってもよく、ゼロまたは1で埋められたデータはセグメントは、空として解釈される。

40

【0254】

トークン指令は、少なくともデータセグメント識別子およびデータセグメント長を具備する。データセグメント識別子は、一般的に、オブジェクトID、ファイル名、あるいは、ストレージデバイスに反復する同一キャラクタまたはキャラクタ列を記憶しようとする、ファイルシステム、アプリケーション、サーバ112などに対する既知のその他の識別子である。データセグメント長は、反復する同一キャラクタまたはキャラクタ列の配列によって求められるストレージ空き領域である。データセグメントトークンおよびトークン指令は、一般的に、反復する同一キャラクタの配列のような、データセグメントのデータを具備していない。

50

【 0 2 5 5 】

しかしながら、トークン指令は、少なくとも1インスタンスの反復する同一キャラクタまたはキャラクタ列のような、データセグメントトークンを形成する他の関連情報を含める。更に、トークン指令は、データセグメント位置、ファイルシステムからのアドレス、データセグメントに対応するデータストレージデバイスの位置などのような、メタデータを含めることができる。当分野の当業者は、トークン指令とともに含むことができるその他の情報が分かるであろう。一実施例においては、指令生成モジュール1002は、トークン指令とともにデータセグメントトークンを生成する。

【 0 2 5 6 】

一実施例においては、トークン指令生成モジュール1002は、ストレージデバイス150に既存のデータを上書きする要求に回答して、トークン指令および安全消去コマンドの双方を生成する。既存のデータは、トークン指令におけるデータセグメント識別子と同一のデータセグメント識別子を有するストレージデバイスで識別されたデータを含める。一般的に、データを上書きする要求は、無効またはガベージとしてデータを単に標識するか、データまたは他の一般的な削除動作に対するポイントを削除するのに十分ではないが、データは、回復不能なように上書きされるように要求される場合に送信される。例えば、データを上書きする要求は、データが機密情報と考えられ、安全上の理由から破壊されなければならない場合に要求される。

【 0 2 5 7 】

安全消去コマンドは、ストレージデバイス150が既存のデータを上書きするように命令し、既存のデータは回復不能となる。次いで、ストレージデバイス150は、既存のデータを上書き、回復、消去するとともに、データセグメントトークンを生成する。結果として、既存のデータは、回復不能になり、データセグメントトークンは、ストレージデバイス150に記憶され、一般的に、既存のデータよりも大幅に少ないストレージ空き領域を取る。

【 0 2 5 8 】

さらなる実施例においては、装置は、消去確認モジュール1016を含み、この消去確認モジュール1016は、ストレージデバイスの既存のデータが、キャラクタとともに上書きされ、既存のデータが回復不能であることの確認を受信する。この確認は、要求デバイスまたはクライアント114に転送され、既存のデータが回復不能である条件におかれたことを確認するために用いられる。他の実施例においては、安全消去コマンドは、ストレージデバイス150が、特有のキャラクタまたはキャラクタ列を有する既存のデータを上書きするように命令でき、実行コマンドは、複数回実行することができる。当分野の当業者は、既存のデータが回復不能であることを確認するために、1又はそれ以上の安全消去コマンドを構成する他の方法が分かるであろう。

【 0 2 5 9 】

データは、暗号化可能であり、次いで、ストレージデバイス150に記憶され、データを記憶すると共にストレージデバイス150によって受信される暗号化キーを用いて暗号化がなされる。既存のデータが、記憶される前に受信された暗号化キーで暗号化される場合、別の実施例においては、トークン指令生成モジュール1002は、既存のデータを上書きするための要求を受信するのに応じて、トークン指令を生成するとともに暗号化消去コマンドを生成する。暗号化消去コマンドは、既存のデータを記憶するのに用いられる暗号化キーを消去し、暗号化キーは回復不能になる。

【 0 2 6 0 】

一実施例においては、暗号化キーを消去することは、要求デバイスから暗号化キーを消去することを含める。別の実施例においては、暗号化キーを消去することは、サーバ、キーの保管場所、キーが記憶されるその他の位置から暗号化キーを消去することを含める。暗号化キーを消去することは、暗号化キーは、いずれの方法でも回復することができないように、他のデータまたはキャラクタの配列を用いて暗号化キーを置換することを含むことができる。一般的に、暗号化ルーチンが、既存のデータを復号化する試みを邪魔するの

10

20

30

40

50

に十分ロバストである既存のデータを暗号化するのに用いられた場合、暗号化キーを消去することは、ストレージデバイス150の既存のデータを回復不能にする。既存のデータを上書きする要求は、データを安全上の理由から上書きする安全消去指令、データを消去するためのデータを上書きする要求、または、反復する同一キャラクタまたはキャラクタ列で既存のデータを置換しようとする要求などでもよい。一実施例においては、安全消去指令は、デバイスに、暗号化キーを安全に消去させ、既存のデータを安全に消去させることを両方行う。一実施例においては、暗号化キーの消去によって、ガベージコレクションプロセスがストレージ空き領域回復プロセスの一部としてデータを消去するまで、ストレージデバイスのデータの安全消去を延期することが可能となる。当分野の当業者は、暗号化キーを消去するその他の方法、および、既存のデータを上書きする要求を受信するその他の方法が分かるであろう。

10

【0261】

一実施例においては、トークン指令は、データセグメントトークンを含み、トークン指令送信モジュール1004は、トークン指令とともにデータセグメントトークンを送信する。別の実施例においては、トークン指令は、データセグメントトークンを含みず、データセグメントトークンを生成するためのストレージデバイス150用のコマンドを含み、この実施例においては、トークン指令送信モジュール1004は、データセグメントトークンを生成するコマンドを有するトークン指令を送信し、データセグメントトークンを送信しない。

20

【0262】

この装置は、ストレージデバイス150にトークン指令を送信するトークン指令送信モジュール1004を含み、一般的に、トークン指令送信モジュール1004は、ストレージ要求の一部としてトークン指令を送信する。ストレージ要求は、オブジェクト要求の形態、データ要求、または当分野の当業者に周知の他の形態であってもよい。トークン指令生成モジュール1002が安全消去指令を生成する場合、トークン指令送信モジュール1004は、安全消去指令をストレージデバイス150に送信する。トークン指令生成モジュール1002は、消去暗号化キーコマンドを生成する場合、必要がある場合には、コマンドを実行する別のデバイスに消去暗号化キーコマンドを送信される。

【0263】

一実施例においては、トークン指令送信モジュール1004は、データセグメントトークンなしでトークン指令を送信する。この実施例においては、トークン指令は、データセグメントトークンを生成するための、ストレージデバイス150に対する命令及び情報を含み、別の実施例においては、トークン指令送信モジュール1004は、データセグメントトークンを含み、トークン指令を送信する。この実施例においては、ストレージデバイス150は、トークン指令と共に受信されたデータセグメントが、データセグメントを表し、適宜なアクションをとり、データセグメントトークンを記憶し、データセグメントトークンは、一般的なデータとしてデータセグメントトークンを単に記憶するのではないデータセグメントを表すことを認識できる。

30

【0264】

特定の実施例においては、装置は、ストレージデバイス150からデータセグメントを読み出すストレージ要求を受信する読み出し受信モジュール1006と、ストレージデバイス150にストレージ要求を送信する読み出し要求送信モジュール1008とを含み、一般的に、ストレージ要求は、外部クライアント114のような要求クライアント114、サーバ112で動くアプリケーションまたはファイルサーバのような、サーバ112内部のクライアント114等から受信される。当分野の当業者は、読み出し受信モジュール1006がストレージ要求を受信できる要求クライアント114として機能するその他のデバイスおよびソフトウェアが分かるであろう。

40

【0265】

ストレージ要求は、トークン指令送信モジュール1004によってストレージデバイス150に送信されるトークン指令に記憶されるように要求されたデータセグメントトーク

50

ンに対応する、データセグメントを読み出すための要求を具える。一実施例においては、要求クライアント 114 は、データセグメントがデータセグメントトークンの形態で記憶されたことを認識しない。別の実施例においては、要求デバイスは、データセグメントが、データセグメントトークンとして記憶されたことを認識するが、データセグメントトークンに記憶された情報を認識しない。

【0266】

特定の実施例においては、装置は更に、ストレージデバイスから要求されたデータセグメントトークンに対応するメッセージを受信する読み出しトークン指令受信モジュール 1010 を具え、このメッセージは、少なくともデータセグメント識別子およびデータセグメント長を具える。このメッセージは、一般的に、データセグメントのデータを含まない。更に、このメッセージは、データセグメント位置または反復する同一キャラクタ若しくはキャラクタ列のような、データセグメントトークンに記憶されたその他の情報を具えることができる。特定の実施例においては、この装置は、ストレージデバイス 150 から受信したメッセージから構築された要求クライアント 113 に、応答を送信する要求クライアント応答モジュール 1012 を具える。

10

【0267】

一実施例においては、読み出しトークン指令受信モジュール 1010 も、メッセージにおいて、既存のデータがキャラクタで上書きされたことの確認を受信し、既存のデータは回復不能となり、既存のデータは、ストレージデバイスに以前に記憶されており、メッセージで受信されたデータセグメントトークンからの同一データセグメント識別子として参照される。この確認は、データセグメントを読み出すためにストレージ要求から別個にストレージデバイス 150 から受信されてもよい。

20

【0268】

別の実施例においては、要求クライアント 114 は、データセグメントを要求する場合、この装置は、メッセージ中に含まれる情報を用いて、データセグメントのデータを復元するデータセグメント再生モジュール 1014 を具える。この場合、要求クライアントに送信される応答は、復元されたデータセグメントを含む。別の実施例においては、要求クライアントに送信される応答は、ストレージデバイス 150 から受信したメッセージに含まれる情報を具える。要求クライアント 114 は、次いで、データセグメントを復元し、その他の方法において、この情報を用いることができる。別の実施例においては、メッセージは、データセグメントトークンを具える。データセグメントトークンは、要求クライアント 114 に転送する前にデータセグメントを復元するために、データセグメント再生モジュール 1014 によって用いられ、あるいは、要求クライアント応答モジュール 1012 が、データセグメントトークンを単純に転送することができる。

30

【0269】

一実施例においては、トークン指令を有するストレージ要求は、更に、ストレージデバイス 150 にストレージ空き領域を割り当てる要求を具え、要求された、割り当てられた空き領域は、データセグメント長とほぼ同一のストレージ空き領域量である。別の実施例においては、ストレージ空き領域を割り当てる要求は、データセグメント長とは異なるストレージ空き領域量用である。例えば、ストレージデバイス 150 がソリッドステートストレージデバイス 102 である場合、ソリッドステートストレージデバイス 102 は、ハードドライブまたはその他の低価格の長期ストレージに接続でき、ソリッドステートストレージ 110 は、長期ストレージ用キャッシュとして構成される。ストレージを割り当てる要求は、ソリッドステートストレージデバイス 102 に、ソリッドステートストレージデバイス 102 にデータを書き込む準備をする際、長期ストレージに対するキャッシュの一部を消去させる。当分野の当業者は、所望されるストレージ空き領域を割り当てる要求であるその他の条件が分かるであろう。

40

【0270】

一実施例においては、装置は、読み出し受信モジュール 1006、読み出し要求送信モジュール 1008、読み出しトークン指令受信モジュール 1010 および要求クライアン

50

ト応答モジュール 1012 で提供され、これらは、上述したものと実質的に同等なものにできる。その実施例においては、モジュール 1006 - 1012 は、トークン生成モジュール 1002 またはトークン指令送信モジュール 1004 を具える装置から独立していてもよい。一実施例においては、この装置は、データセグメント再生成モジュール 1014 を具え、これは、上述したデータセグメント再生成モジュール 1014 と実質的に同じである。

【0271】

図 11 は、本発明によるトークン指令を生成し送信する方法 1100 の一実施例を示した概略フローチャート図である。この方法 1100 は開始し (1102)、トークン指令生成モジュール 1002 は、トークン指令を有するストレージ要求を生成し (1104)、トークン指令は、ストレージデバイス 150 にデータセグメントトークンを記憶する要求を具える。トークン指令送信モジュール 1004 は、ストレージデバイス 150 にトークン指令を送信し (1106)、方法 1100 は終了する (1108)。一実施例においては、ストレージ要求は、データセグメントトークンを記憶するトークン指令を具え、ストレージ要求は、データセグメントのデータを実質的に含まない。別の実施例においては、ストレージ要求は、データセグメントからのデータを含む。好ましい実施例においては、ソフトウェアアプリケーションは、データセグメントの生成を回避するためにトークン指令を用いてストレージ要求を生成する。別の実施例においては、ソフトウェアアプリケーションは、トークン指令の生成を要求する。

10

【0272】

図 12 は、本発明によるデータセグメントトークンを読み出す方法 1200 の一実施例を示した概略フローチャート図である。この方法 1200 は、開始し (1202)、読み出し受信モジュール 1006 は、要求クライアント 114 から、ストレージデバイス 150 からのデータセグメントを読み出すストレージ要求を受信する (1204)。読み出し要求送信モジュール 1008 は、ストレージデバイス 150 にストレージ要求を送信する (1206)。

20

【0273】

読み出しトークン指令受信モジュール 1008 は、ストレージデバイス 150 から、要求されたデータセグメントトークンに対応するメッセージを受信し (1208)、このメッセージは、少なくともデータセグメント識別子およびデータセグメント長を具える。このメッセージは、データセグメントのデータを実質的に含まない。要求クライアント応答モジュール 1012 は、ストレージデバイス 150 から受信されたメッセージから構築された要求クライアントに対する応答を送信し (1210)、この方法 1200 は、終了する (1212)。

30

【0274】

図 13 は、本発明によるデータセグメントトークンを管理するための装置を有するシステム 1300 の一実施例を示した概略ブロック図である。このシステム 1300 は、書き込み要求受信モジュール 1302 およびデータセグメントトークンストレージモジュール 1304 を有する装置を具え、様々な実施例においては、トークン指令生成モジュール 1306、読み出し要求受信モジュール 1308、読み出しデータセグメントトークンモジュール 1310、送信データセグメントトークンモジュール 1314 および送信データセグメントモジュール 1316 を有する読み出し要求応答モジュール 1312、復元データセグメントモジュール 1318、消去確認モジュール 1322 を有する安全消去モジュール 1320、ストレージ空き領域割り当てモジュール 1324、を具え、これらは、以下で説明する。システム 1300 は、ストレージコントローラ 152 およびデータストレージデバイス 154 を有するストレージデバイス 150 を具え、これらは、上述したものと実質的に同じである。システム 1300 は、要求デバイス 1326 を具え、これは以下に記載されているように、ストレージデバイス 150 と接続している。

40

【0275】

示された実施例においては、モジュール 1302 - 1324 は、ストレージデバイス 1

50

50 またはストレージコントローラ152に含まれる。別の実施例においては、1又はそれ以上のモジュール1302 - 1324の少なくとも一部分は、ストレージデバイス150の外側に配置される。更なる実施例においては、要求デバイス1326は、ドライバ、ソフトウェアまたは1又はそれ以上のモジュール1302 - 1324のその他の機能の形態において、モジュール1302 - 1324の一部を具える。例えば、トークン生成モジュール1306および復元データセグメントモジュール1318は、要求デバイス1326中に示されている。当分野の当業者は、モジュール1302 - 1324の機能を分散し、実装するその他の方法が分かるであろう。

【0276】

この装置は、要求デバイス1326からのストレージ要求を受信する書き込み要求受信モジュール1302を具え、ストレージ要求は、ストレージデバイス150内にデータセグメントを記憶する要求を具える。データセグメントは、反復する同一キャラクタ又はキャラクタ列の配列を具える。一般的に、反復する同一キャラクタの配列は、データセグメントが空であることを意味する。これは特に、反復する同一キャラクタの配列が1又はゼロである場合に正しい。この装置は、ストレージデバイス150内にデータセグメントトークンを記憶するデータセグメントトークンストレージモジュール1304を具える。データセグメントトークンは、少なくともデータセグメント識別子およびデータセグメント長を具える。データセグメントトークンは、データセグメントからの実際のデータを実質的に含まない。

10

【0277】

データセグメントトークンは、多数の形態で記憶可能である。一実施例においては、インデックス中のエントリを具えており、このインデックスは、ストレージデバイス150に記憶された情報およびデータに対応する。例えば、このインデックスは、図2に示された装置200に関連して上述したように、オブジェクトインデックスであってもよい。更に、このインデックスは、ファイルシステムインデックス、ブロックストレージインデックス、または、当分野の当業者に周知なその他のインデックスであってもよい。別の実施例においては、データセグメントトークンは、ストレージデバイス150に記憶されるメタデータを具え、あるいは、その形態である。別の実施例においては、セグメントトークンが、ストレージデバイスにメタデータとして記憶され、インデックス中のエントリを具える。当分野の当業者は、データセグメントトークンを記憶する他の方法が分かるであろう。

20

30

【0278】

一実施例においては、ストレージ要求は、データセグメントトークンを記憶するトークン指令を具え、このストレージ要求は、データセグメントのデータを実質的に含まない。トークン指令は、データセグメントトークンを生成するデータセグメントトークンまたはコマンドを具えることができる。トークン指令は、データセグメントトークンを具えていない場合、データトークンストレージモジュール1304は、トークン指令中の情報からデータセグメントトークンを生成する。トークン指令がデータセグメントトークンを具える場合は、データセグメントトークンストレージモジュール1304が、トークン指令中のデータセグメント識別子によって識別されるデータセグメントを表すデータ構造としてデータセグメントを認識し、適宜にデータセグメントトークンを記憶する。

40

【0279】

一般的に、データセグメントトークンストレージモジュール1304が、データセグメントトークンを認識する場合、データセグメントトークンは、ストレージデバイス150に記憶される他のデータとは何らかの方法で区別される。例えば、要求デバイス1326は、データを圧縮のみし、圧縮オブジェクト、ファイルまたはセグメントを送信するが、ストレージデバイス150は、他のストレージ要求によって受信された他のデータから圧縮データセグメントを区別できない。

【0280】

データセグメントトークンストレージモジュール1304は、受信されたデータセグメ

50

ントークンがデータセグメントトークンであることを認識する場合、データセグメントトークンストレージモジュール1304は、読み出すとき、データセグメントトークンが、データセグメントトークンではなくデータセグメントとして存在するような方法で、データセグメントトークンを記憶することができる。当分野の当業者は、受信したデータセグメントトークンが、データセグメントではなくデータセグメントトークンであると認識した後に、データセグメントトークンストレージモジュール1304がデータセグメントトークンを記憶することができるような他の方法が分かるであろう。

【0281】

別の実施例においては、ストレージ要求は、データセグメントからのデータを含める。この実施例においては、この装置は、データセグメントからデータセグメントトークンを生成するトークン生成モジュール1306を含み、データセグメントトークンは、データセグメントを記憶するストレージ要求に回答して生成される。更なる実施例においては、トークン生成モジュール1306は、可能なドライバ形態において、要求デバイス1326中に常駐する。

10

【0282】

一実施例においては、この装置は、既存のデータが回復不能であるように、キャラクタを用いて既存のデータを上書きする安全消去モジュール1320を含み、既存のデータは、ストレージ要求で識別されたデータセグメントと同一のデータセグメント識別子で識別された、ストレージデバイスで既に記憶されたデータセグメントのデータを含める。この実施例においては、データセグメントトークンは、データセグメント識別子およびデータセグメント長とともに記憶され、データセグメントトークンに記憶された同一のデータセグメント識別子によって識別される既存のデータは、既存のデータを上書きすることで消去される。一般的に、既存のキャラクタは、ゼロ、1、または、その他のキャラクタ列によって上書きされ、データは破壊され、回復不能となる。

20

【0283】

さらなる実施例においては、安全消去モジュールは、既存のデータが上書きされることを示すメッセージを送信する消去確認モジュール1322を更に含める。一般的に、メッセージは、要求デバイス1326に送信される。消去確認メッセージは、安全消去モジュール1320が既存のデータを上書きした後に送信される。このメッセージは、ストレージ要求として同一のトランザクションにおいて、あるいは、異なるトランザクションにおいて、送信されてもよい。

30

【0284】

別の実施例においては、安全消去モジュール1320は、ストレージ空き領域回復動作中に、既存のデータを上書きする。例えば、ストレージデバイス150がソリッドステートストレージデバイス102である場合、上述したように、ストレージ空き領域回復動作は、図8に示した装置800に関して記載されたガベージコレクションに関連してもよい。しかしながら、既存のデータを上書きする要求を含むストレージ空き領域回復動作は、一般的に促進され、既存のデータが記憶されるストレージ位置が、確認メッセージが消去確認モジュール1322によって送信される前に、必ず回復される。一実施例においては、既存のデータは、安全消去が要求されていることを示すように、標識され、あるいは、識別される。一般的に、確認メッセージは、消去を標識した既存のデータが上書きされて回復不能になるまで、送信されない。別の実施例においては、安全消去モジュール1320は、後のストレージ空き領域回復のために無効であることを既存のデータに標識するのみである。別の実施例においては、安全消去は、既存のデータが無効であることを示すインデックスを更新し、後のストレージ空き領域回復中にデータが上書きされるまで、このデータへのアクセスを阻止する。

40

【0285】

一実施例においては、安全消去モジュール1320は、データセグメントが記憶される毎に、既存のデータを上書きする。別の実施例においては、既存のデータを上書きする要求を特に含み、安全消去モジュール1320は、既存のデータを上書きする要求に回答し

50

て、既存のデータを上書きする。別の実施例においては、安全消去モジュール 1320 は、既存のデータが消去されたという確認に関するメタデータ情報を記憶し、次の読み出しが消去を指示することができる。

【0286】

安全消去が受信されない場合の他の実施例において、既存のデータが削除される。一実施例においては、データを削除することは、インデックスエントリ、アドレスなどを削除することを含める。好ましい実施例においては、データセグメントトークンが記憶される時、対応する既存のデータは、無効あるいはストレージ回復の準備ができていと標識される。このデータは、ストレージ回復動作またはガベージコレクション動作などにおいて後で回復可能である。

10

【0287】

特定の実施例においては、この装置は、データセグメントを読み出すためのストレージ要求を受信する読み出し要求受信モジュール 1308、ストレージ要求によって要求されるデータセグメントに対応するデータセグメントトークンを読み出し読み出しデータセグメントトークンモジュール 1310、および、要求デバイス 1326 に対する応答を送信する読み出し要求応答モジュール 1312 を含める。この応答は、要求されたデータセグメントに対応するデータセグメントトークンを用いて生成される。

【0288】

一実施例において、データセグメントを読み出すストレージ要求は、ストレージ要求に関連しており、ストレージ要求が成功したことを確認するのに役立つ。別の実施例においては、データセグメントを読み出す要求は、ストレージ要求から独立しており、ストレージ要求を生成した要求デバイス 1326 または別の別個の要求デバイス 1326 によって開始されてもよい。

20

【0289】

一実施例においては、要求デバイスが、実際のデータの位置にあるデータセグメントトークンからの情報を受信する場合、読み出し要求応答モジュール 1312 は、要求デバイス 1326 へのメッセージに返信して送信する送信データセグメントトークンモジュール 1314 を含める。このメッセージは、少なくともデータセグメント識別子およびデータセグメント長を含めるが、更に、データセグメント位置、少なくとも 1 インスタンスの回復する同一キャラクタ又はキャラクタ列、あるいは、その他の関連情報を含めてもよい。一般的に、メッセージは、データセグメントトークンに含まれるもの以外に、データセグメントの実際のデータを含んでいない。

30

【0290】

別の実施例においては、要求デバイス 1326 は、データセグメントを受信するとされる場合、この装置は、データセグメントトークンを用いて、データセグメントのデータを復元する復元データセグメントモジュール 1318 を含める。読み出し要求応答モジュール 1312 は、更に、要求デバイス 1326 に復元された、要求されたデータセグメントを送信する、送信データセグメントモジュール 1316 を含める。別の実施例において、復元データセグメントモジュール 1318 は、可能であれば、ドライバの形態において、要求デバイス 1326 に常駐し、送信データセグメントトークンモジュール 1314 は、データセグメントトークン情報を有するメッセージを要求デバイス 1326 に送信する。復元データセグメントモジュール 1318 は、要求デバイス 1326 で、メッセージから要求されたデータセグメントを復元する。

40

【0291】

一実施例においては、システム 1300 は、読み出し要求受信モジュール 1308、読み出しデータセグメントトークンモジュール 1310、読み出し要求応答モジュール 1312 を含む別の装置を含め、これらは、上述したものと実質的に同一である。この装置は、書き込み要求受信モジュール 1302 およびデータセグメントトークンストレージモジュール 1304 を含める装置から独立していてもよい。一実施例においては、読み出し要求応答モジュール 1312 は、送信データセグメントトークンモジュール 1314 および

50

／または送信データセグメントモジュール1316を具え、この装置は、復元データセグメントモジュール1318を具えてもよく、モジュール1314、1316、1318が上述したものと実質的に同じである。

【0292】

図14は、本発明によるデータセグメントトークンを記憶する方法1400の一実施例を示す概略フローチャート図である。この方法1400は、開始し(1402)、書き込み要求受信モジュール1302が、要求デバイス1326からのストレージ要求を受信し(1404)、上記ストレージ要求は、ストレージデバイス150にデータセグメントを記憶する要求を具える。データセグメントは、反復する同一キャラクタ又はキャラクタ列の配列を具える。データセグメントトークンストレージモジュール1304は、ストレージデバイス150にデータセグメントトークンを記憶し(1406)、方法1400は終了する(1408)。データセグメントトークンは、少なくともデータセグメント識別子およびデータセグメント長を具え、データセグメントトークンは、ほとんどの部分に関してデータセグメントからのデータを具えていない。

10

【0293】

図15は、本発明によるデータセグメントトークンを読み出す方法1500の一実施例を示した概略フローチャート図である。この方法1500は開始し(1502)、読み出し要求受信モジュール1308が、ストレージデバイス150からデータセグメントを読み出すストレージ要求を受信する(1504)。データセグメントは、データセグメントトークンによってストレージデバイスに表現され、このデータセグメントは、反復する反復する同一キャラクタ又はキャラクタ列の配列を具える。このデータセグメントトークンは、少なくともデータセグメント識別子およびデータセグメント長を具え、データセグメントトークンは、データセグメントからのデータを具えていない。読み出しデータセグメントトークンモジュール1310は、ストレージ要求において要求されたデータセグメントに対応するデータセグメントを読み出し(1506)、読み出し要求応答モジュール1312は、要求デバイス150に応答を送信し(1508)、この方法1500は終了する(1510)。この応答は、要求されたデータセグメントに対応するデータセグメントトークンを用いて生成される。

20

【0294】

進行型RAID

30

【0295】

独立ドライブ冗長アレイ(「RAID」)は様々な目的を達成するために多くの方法で構成される。本明細書に記載のように、ドライブはデータ用の大容量ストレージデバイスである。ドライブ又はストレージデバイスは、ソリッドステートストレージ110、ハードディスクドライブ(「HDD」)、テープドライブ、光学ドライブ又は当該技術分野の当業者に公知のその他の大容量ストレージデバイスにできる。一実施例においては、ドライブは仮想容量としてアクセスされる大容量ストレージデバイスの一部分を具える。別の実施例においては、ドライブは仮想容量として共にアクセス可能で、かつ、RAIDとして、「ただのディスク/ドライブの束」(「JBOD」)として等でストレージエリアネットワーク(「SAN」)に構成される2又はそれ以上のデータストレージデバイスを具える。一般的には、ドライブはストレージコントローラ152を通して単一ユニット又は仮想容量としてアクセスされる。好ましい実施例においては、ストレージコントローラ152はソリッドステートストレージコントローラ104を具える。当該技術分野の当業者は、RAIDに構成されうる大容量ストレージデバイスの形態における、他のドライブの形態が分かるであろう。本明細書に記載のそれらの実施例においては、ドライブ及びストレージデバイス150は互換的に用いられる。

40

【0296】

従来的には、様々なRAID構成はRAIDレベルと呼ばれる。1の基本RAID構成はストレージデバイス150のミラーコピーを生成するRAIDレベル0である。RAID0の利点は、1又はそれ以上のストレージデバイス150のデータの完全なコピーが1

50

又はそれ以上のストレージデバイス150のミラーコピーで更に利用可能であり、プライマリドライブ又はミラーリングドライブのデータの読み出しが比較的高速であることである。RAID0は更にプライマリストレージデバイス150における故障の場合に、データのバックアップコピーを更に提供する。RAID0の欠点は、書き込みデータが二度書き込まれなければならないため、書き込みが比較的低速であることである。

【0297】

別の従来のRAID構成はRAIDレベル1である。RAID1においては、RAIDに書き込まれるデータがストレージデバイス150のセットにおけるN個のストレージデバイス150に対応するN個のデータセグメントに分割される。N個のデータセグメントは「ストライプ」を形成する。複数のストレージデバイス150にわたってデータをストライピングすることによって、単一のストレージデバイス150がN個のデータセグメントを含むデータを保存できるよりも高速に、N個のデータセグメントを記憶するのに並行してストレージデバイス150が動くため、性能が向上する。しかしながら、データが複数のストレージデバイス150にわたって分散し、複数のストレージデバイス150のアクセス時間が総ての所望のデータを含むストレージデバイス150からデータの読み出しよりも一般的に低速であるため、データを読み出すのは比較的低速である。更には、RAID1は故障の保護を提供しない。

【0298】

普及しているRAID構成は、N個のストレージデバイス150にわたるN個のデータセグメントのストライピングするステップと、N+1ストレージデバイス150でパリティデータセグメントを記憶するステップとを具えるRAIDレベル5である。そのRAIDは、ストレージデバイス150の単一の故障を許容するため、RAID5は故障耐性を提供する。例えば、ストレージデバイス150が故障した場合、ストライプのデータセグメントの欠損は、他の利用可能なデータセグメント及びストライプ用に特に算出されたパリティデータセグメントを用いて生成されうる。更に、RAID5は一般的には、RAIDされたストレージデバイス150のセットの各ストレージデバイス150がデータの完全なコピーではなく、ストライプのデータセグメント又はパリティデータセグメントのみを記憶することを要求するため、RAID0より少ないストレージ空き領域を用いる。RAID1と同様、RAID5はデータの書き込みが比較的高速であるが、データを読み出すのが比較的低速である。しかしながら、パリティデータセグメントはストライプのN個のデータセグメントから各ストライプ用に算出しなければならないため、一般的な従来のRAID5用のデータの書き込みは、RAID1よりも低速である。

【0299】

別の普及しているRAID構成は二重分散パリティを具えるRAIDレベル6である。RAID6においては、2のストレージデバイス150がパリティミラーデバイス(例えば、1602a、1602b)として割り当てられる。ストライプ用の各パリティデータセグメントは別個に算出されて、ストレージデバイスセットのいずれかの2のストレージデバイス150の欠損が残りの利用可能なデータセグメント及び/又はパリティデータセグメントを用いて回復可能である。RAID6はRAID5と同様の性能の利点及び欠点を有する。

【0300】

入れ子構造RAIDは、高信頼性が要求される場合に故障耐性を増大させるのに更に用いられうる。例えば各々がRAID5として構成される2のストレージデバイスセットがRAID0構成でミラーリングされうる。生じた構成はRAID50と呼ばれる。RAID6が各々のミラーリングセットで用いられる場合、その構成はRAID60と呼ばれる。入れ子構造RAID構成は基礎をなすRAIDグループと同様の性能の問題を一般的に有している。

【0301】

前述の記載から、故障耐性の利益と、RAID0、RAID5、RAID6のような従来の故障耐性のあるRAIDレベルよりも高速なデータ書き込み等を提供する一方、RA

10

20

30

40

50

RAID 1、RAID 5、RAID 6 等のような従来のストライピングの RAID レベルよりも高速なデータ読み出しを更に提供する進行型 RAID 用の装置、システム及び方法に対するニーズが存在することは明らかであろう。有益には、このような装置、システム及び方法は、パリティデータセグメントが、ストレージ統合動作前又はストレージ統合動作の一部のように算出されるの要求されるまで、N 個のデータセグメントをパリティミラーストレージデバイス 1602 に書き込み、RAID 0 システムの利点を提供する。

【0302】

図 10 は本発明による進行型 RAID 用のシステム 1600 の一実施例を示す概略ブロック図である。システム 1600 は 1 又はそれ以上のクライアント 114 によりコンピュータネットワーク 116 を通してアクセス可能な N 個のストレージデバイス 150 と、M 個のパリティミラーストレージデバイス 1602 とを具える。N 個のストレージデバイス 150 及びパリティミラーストレージデバイス 1602 は、1 又はそれ以上のサーバ 112 に配置できる。ストレージデバイス 150、サーバ 112、コンピュータネットワーク 116、及びクライアント 114 は実質的に上述のものと同様である。パリティミラーデバイス 1602 は一般的に N 個のストレージデバイス 150 と同様または同一であり、一般的にストライプ用のパリティミラーストレージデバイス 1602 として指定される。

10

【0303】

一実施例においては、N 個のストレージデバイス 150 及び M 個のパリティミラーストレージデバイス 1602 は、1 のサーバ 112 に含まれるか、それを通してアクセス可能であり、システムバス、SAN、又は他のコンピュータネットワーク 116 を用いて互いにネットワーク化できる。別の実施例においては、N 個のストレージデバイス 150 及び M 個のパリティミラーストレージデバイス 1602 は複数のサーバ 112 a - n + m に配置され、それを通してアクセス可能にできる。例えば、ストレージデバイス 150 及びパリティミラーストレージデバイス 1602 は、図 1C のシステム 103 及び図 5B の方法 105 に関連して上述したようなインサーバ SAN の一部にできる。

20

【0304】

一実施例においては、パリティミラーストレージデバイス 1602 は進行型 RAID に記憶されたストライプの総てのパリティデータセグメントを記憶する。別の好ましい実施例においては、進行型 RAID に割り当てられたストレージデバイスセット 1604 のストレージデバイス 150 は、特定のストライプ用のパリティミラーストレージデバイス 1602 であるように割り当てられ、その割り当てはローテーションし、パリティデータセグメントは各ストライプについて、N + M 個のストレージデバイス 150 の間でローテーションする。この実施例は、単一のストレージデバイス 150 を割り当てて、各ストライプ用のパリティミラーストレージデバイス 1602 にするのを介して、性能の利点を提供する。パリティミラーストレージデバイス 1602 をローテーションさせることによって、パリティデータセグメントを算出及び記憶することに関するオーバーヘッドは、分散される。

30

【0305】

一実施例においては、ストレージデバイス 150 は各々が関連するソリッドステートストレージ 110 及びソリッドステートストレージコントローラ 104 を有するソリッドステートストレージデバイス 102 である。別の実施例においては、各ストレージデバイス 150 はソリッドステートストレージコントローラ 104 を具え、関連するソリッドステートストレージ 110 はテープストレージ又はハードディスクドライブのような低価格、低性能のストレージ用のキャッシュとして動作する。別の実施例においては、1 又はそれ以上のサーバ 112 はストレージ要求を進行型 RAID に送信する 1 又はそれ以上のクライアント 114 を具える。当該技術分野の当業者は、進行型 RAID 用に構成されうる N 個のストレージデバイス 150 と、1 又はそれ以上のパリティミラーストレージデバイス 1602 とを有する他のシステム構成が分かるであろう。

40

【0306】

図 17 は、本発明による進行型 RAID 用の装置 1700 の一実施例を示す概略ブロッ

50

ク図である。様々な実施例においては、装置 1700 はストレージ要求受信モジュール 1702、ストライピングモジュール 1704、パリティミラーモジュール 1706、パリティ進行モジュール 1708、パリティ変更モジュール 1710、ミラーリング設定モジュール 1712、更新モジュール 1714、直接クライアント応答モジュール 1718 を有するミラーリング復元モジュール 1716、事前統合モジュール 1720、事後統合モジュール 1722、データ再構築モジュール 1724、及びパリティ再構築モジュール 1726 を具え、以下に述べられている。モジュール 1702 - 1726 はサーバ 112 に示されるが、モジュール 1702 - 1726 の機能の一部又は総ては複数のサーバ 112、ストレージコントローラ 152、ストレージデバイス 150、クライアント 114 に更に分散されうる。

10

【0307】

装置 1700 はデータを記憶する要求を受信するストレージ要求受信モジュール 1702 を具え、データはファイル又はオブジェクトのデータである。一実施例においては、ストレージ要求はオブジェクト要求である。別の実施例においては、ストレージ要求はブロックストレージ要求である。一実施例においては、ストレージ要求はデータを含まず、ストレージデバイス 150 及びパリティミラーストレージデバイス 1602 によって、クライアント 114 又は他のソースからの DMA 又は RDMA データに対して用いられうるコマンドを含む。別の実施例においては、ストレージ要求はストレージ要求の結果として記憶されるべきデータを含む。別の実施例においては、ストレージ要求はストレージデバイスセット 1604 に記憶されるデータを有することが可能な 1 のコマンドを含む。別の実施例においては、ストレージ要求は複数のコマンドを含む。当該技術分野の当業者は進行型 RAID に好適なデータを記憶する他のストレージ要求が分かるであろう。

20

【0308】

データは装置 1700 にアクセス可能な位置に記憶される。一実施例においては、データはクライアント 114 又はサーバによって用いられる RAM のようなランダムアクセスメモリ（「RAM」）で利用可能である。別の実施例においては、データはハードディスクドライブ、テープストレージ、又は他の大容量ストレージデバイスに記憶される。一実施例においては、データはオブジェクトとして、又は、ファイルとして構成される。別の実施例においては、オブジェクト又はファイルの一部であるデータブロックとして構成される。当該技術分野の当業者はストレージ要求の対象であるデータの他の形式及び位置が分かるであろう。

30

【0309】

装置 1700 はデータ用のストライプパターンを算出するストライピングモジュール 1704 を具える。ストライプパターンは 1 又はそれ以上のストライプを含み、各ストライプは N 個のデータセグメントのセットを具える。一般的に、ストライプのデータセグメントの数はどのくらいの数のストレージデバイス 150 が RAID グループに割り当てられるかに依存する。例えば、RAID 5 が用いられた場合、1 のストレージデバイス 150 はパリティミラーストレージデバイス 1602 a として割り当てられて、特定のストライプ用のパリティデータを記憶する。4 の他のストレージデバイス 150 a、150 b、150 c、150 d が RAID グループに割り当てられた場合、ストライプはパリティデータセグメントに加えた、4 のデータセグメントを有する。ストライピングモジュール 1704 は N 個のデータセグメントを、N 個のストレージデバイス 150 a - n に対する N 個のストライプに書き込み、N 個のデータセグメントの各々は、ストライプに割り当てられたストレージデバイス 150 のセット 1604 内の別個のストレージデバイス 150 a、150 b、...、150 n に書き込まれる。当該技術分野の当業者は、特定の RAID レベルの RAID グループに割り当てられうるストレージデバイス 150 の様々な組合せ、及びストライピングパターンを生成し、ストライプごとに N 個のデータセグメントにデータを分割する方法が分かるであろう。

40

【0310】

装置 1700 は、ストライプの N 個のデータセグメントのセットを、ストレージデバイ

50

セット 1604 内の 1 又はそれ以上のパリティミラーストレージデバイス 1602 に書き込むパリティミラーモジュール 1706 を具え、パリティミラーストレージデバイス 1602 が N 個のストレージデバイス 150 に加えて存在する。N 個のデータセグメントはその後パリティデータセグメントの将来の算出のために利用できる。パリティデータセグメントをすぐに算出するのではなく、パリティミラーモジュール 1706 は N 個のデータセグメントのセットをパリティミラーストレージデバイス 1602 にコピーし、一般的に N 個のデータセグメントを記憶するよりも少ない時間を必要とする。N 個のデータセグメントが一度パリティミラーストレージデバイス 1602 に記憶された場合、N 個のデータセグメントは、N 個のストレージデバイス 150 のうちの 1 つが利用不可能になる場合に、データを復元するために読み出し、使用するために利用可能である。N 個のデータセグメントの総てが 1 のストレージデバイス（例えば、1602 a）から一緒に利用できるために、データの読み出しは更に RAID0 構成の利点を有する。1 以上のパリティミラーストレージデバイス（例えば、1602 a、1602 b）のために、パリティミラーモジュール 1706 は N 個のデータセグメントを各パリティミラーストレージデバイス 1602 a、1602 b にコピーする。

10

【0311】

装置 1700 は、ストレージ統合動作に応じてストライプ用の 1 又はそれ以上のパリティデータセグメントを算出するパリティ進行モジュール 1708 を具える。N 個のデータセグメントから算出される 1 又はそれ以上のパリティデータセグメントは、パリティミラーストレージデバイス 1602 に記憶される。パリティ進行モジュール 1708 は、1 又はそれ以上のパリティミラーストレージデバイス 1602 の各々に、パリティデータセグメントを記憶する。ストレージ統合動作は 1 又はそれ以上のパリティミラーストレージデバイス 1602 のうちの少なくとも 1 つの、少なくともストレージ空き領域又はデータ又は両方を回復するように導かれる。例えば、ストレージ統合動作は図 8 及び 9 の装置 800 及び方法 900 に関連して上述したような、ソリッドステートストレージデバイス 102 のデータガベージコレクションにできる。ストレージ統合動作はストレージ空き領域を増加させるためにデータを統合する、ハードディスクドライブ用のデフラグ動作又は他の同様の動作を更に含む。本明細書で用いられるように、ストレージ統合動作はデータを回復するための動作を含み、例えば、ストレージデバイス 150 が利用不可能である場合、エラー、又は、パリティミラーストレージデバイス 1602 からデータを読み出す他の理由から回復させる。別の実施例においては、パリティ生成モジュール 1708 は、パリティミラーストレージデバイス 1602 が混雑が少ない場合には、パリティデータセグメントを単に算出する。

20

30

【0312】

有利には、ストライプのパリティデータセグメントの算出及び記憶を遅らせることにより、パリティミラーストレージデバイス 1602 の N 個のデータセグメントは、データセグメントを読み出すため、データを回復するため、より多くのストレージ空き領域がパリティミラーストレージデバイス 1602 に必要となるまでにデータを再構築するために、又はストレージ統合動作に対する他の理由のために利用可能である。パリティ進行モジュール 1708 はその時に、ストレージ要求受信モジュール 1702、ストライピングモジュール 1704、又はパリティミラーモジュール 1706 から自律して、バックグラウンド動作として実行できる。当該技術分野の当業者は、進行型 RAID 動作の一部としてパリティデータセグメントの算出を遅らせる他の理由が簡単に分かるであろう。

40

【0313】

一実施例においては、データを記憶する要求を受信し、ストライプパターンを算出して N 個のデータセグメントを N 個のストレージデバイスに書き込み、N 個のデータセグメントのセットをパリティミラーストレージデバイスに書き込み、パリティデータセグメントを算出するモジュール 1702 - 1708 の機能の一部又は総てが、ストレージデバイスセット 1604 のストレージデバイス 150、クライアント 114 及びサードパーティの RAID 管理デバイスに生じる。サードパーティの RAID 管理デバイスはサーバ 114

50

又はその他のコンピュータにできる。

【0314】

一実施例においては、装置1700は、各ストライプのために、ストレージデバイスセット1604内のストレージデバイス150のどちらがストライプ用の1又はそれ以上のパリティミラーストレージデバイス1602になるように割り当てられるかを交互に代替するパリティ代替モジュール1710を具える。図10のシステム1600に関連して上述したように、ストライプ用のパリティミラーストレージデバイスのためにどちらのストレージデバイス150が用いられるかをローテーションすることによって、様々なパリティデータセグメントの作業領域算出が、ストレージデバイスセット1604のストレージデバイス150の間で分散される。

10

【0315】

別の実施例においては、ストレージデバイスセット1604は第1のストレージデバイスセットであり、装置1700は第1のストレージセット1604に加えて1又はそれ以上のストレージデバイスセットを生成するミラーリング設定モジュール1712を具えて、1又はそれ以上の追加のストレージデバイスセットの各々が、1又はそれ以上の追加のストレージセットの各々のN個のストレージデバイス150にN個のデータセグメントを書き込む、少なくとも1の付随するストライピングモジュール1704を具えるようにする。関連する実施例においては、1又はそれ以上の追加のストレージデバイスセットの各々は、N個のデータセグメントのセットを記憶するための付随したパリティミラーモジュール1706と、1又はそれ以上のパリティデータセグメントを算出するためのパリティ進行モジュール1708とを具える。ミラーリング設定モジュール1712が1又はそれ以上のミラーリングストレージデバイスセットを生成する場合、RAIDはRAID50のような入れ子構造RAIDにできる。この実施例においては、RAIDレベルはRAID10から進行され、データはRAID50又はRAID60までストライピング及びミラーリングされ、パリティデータセグメントは各ストレージデバイスセット1604用に算出及び記憶される。

20

【0316】

一実施例においては、装置1700は更新モジュール1714を具える。更新モジュール1714はパリティミラーストレージデバイス1602のN個のデータセグメントがパリティデータセグメントに進行されなかった場合に、一般的に用いられる。更新モジュール1714は更新されたデータセグメントを受信し、更新されたデータセグメントはN個のストレージデバイス150に記憶されるN個のデータセグメントの既存のデータセグメントに対応する。更新モジュール1714は更新されたデータセグメントを、既存のデータセグメントが記憶されるストライプのストレージデバイス150に、及び、ストライプの1又はそれ以上のパリティミラーストレージデバイス1602にコピーする。更新モジュール1714はN個のストレージデバイス150a-nのストレージデバイス150に記憶された既存のデータセグメントを、更新されたデータセグメントと置換し、1又はそれ以上のパリティミラーストレージデバイス1602に記憶された対応する既存のデータセグメントを更新されたデータセグメントと置換する。

30

【0317】

一実施例においては、データセグメントを置換するステップは、データセグメントをストレージデバイス150に書き込むステップと、次いで後のガベージコレクションのために、対応するデータセグメントを無効と標識するステップとを含む。この実施例のある例が、図8及び9に関連して上述されたソリッドステートストレージ110及びガベージコレクション装置のために記載されている。別の実施例においては、データセグメントを置換するステップは更新されたデータセグメントで既存のデータセグメントを上書きするステップを含む。

40

【0318】

一実施例においては、ストレージデバイス1604のセットは第1のストレージデバイスセットであり、装置1700は第1のストレージセット1604のストレージデバイス

50

150に記憶されたデータセグメントを回復するミラーリング復元モジュール1716を
具え、第1のストレージセット1604のストレージデバイス150は利用不可能になる
。データセグメントはデータセグメントのコピーを含むミラーストレージデバイスから回
復される。ミラーストレージデバイスはN個のデータセグメントのコピーを記憶する1又
はそれ以上のストレージデバイス150のセットの1つを具える。

【0319】

更なる実施例においては、ミラーリング復元モジュール1716は、データセグメント
を読み出す、クライアント114からの読み出し要求に応じてデータセグメントを回復す
る。別の関連する実施例においては、ミラーリング復元モジュール1716は、ミラー
ストレージデバイスからクライアント114に要求されたデータセグメントを送信する直接
クライアント応答モジュール1718を更に具える。この実施例においては、要求された
データセグメントはクライアント114にコピーされ、クライアント114はデータセグ
メントをクライアント114に送信する前にデータセグメントが回復されるまでクライ
アント114は待つ必要はない。

10

【0320】

一実施例においては、装置1700は、データセグメントを読み出す要求に応じて、ス
トレージセット1604のストレージデバイス150に記憶されるデータセグメントを回
復する事前統合復元モジュールを具える。その実施例においては、ストレージデバイス1
50は利用不可能であり、パリティ進行モジュール1708が1又はそれ以上のパリティ
ミラーストレージデバイス1602に1又はそれ以上のパリティデータセグメントを生成
する前に、データセグメントはパリティミラーストレージデバイス1602から回復され
る。

20

【0321】

別の実施例においては、装置1700はストレージセットのストレージデバイス150
に記憶されるデータセグメントを回復する事後統合復元モジュール1724を具える。そ
の実施例においては、ストレージデバイス150は利用不可能であり、パリティ進行モ
ジュール1708が1又はそれ以上のパリティデータセグメントを生成した後に、パリティ
ミラーストレージデバイス150の1又はそれ以上に記憶される1又はそれ以上のパリ
ティデータセグメントを用いて、データセグメントは回復される。例えば、事後統合復元
モジュール1724はパリティデータセグメント及び利用可能なN個のストレージデバイス
150の利用可能なデータセグメントを用いて、欠損したデータセグメントを再生成する
。

30

【0322】

一実施例においては、装置1700は再構築動作中の置換ストレージデバイスの回復し
たデータセグメントを記憶するデータ再構築モジュール1724を具え、回復したデータ
セグメントは利用不可能なストレージデバイス150に記憶された利用不可能なデータ
セグメントと合致する。利用不可能なストレージデバイス150はストレージデバイスセ
ット1602のN個のストレージデバイス150のうちの一つである。一般的には、再構築
動作は利用不可能なデータセグメントを記憶するストレージデバイス150の故障後に生
じる。再構築動作はデータセグメントを置換ストレージデバイスに復元することであり、
利用不可能なストレージデバイス150に既に記憶されたデータセグメントと合致させる
。

40

【0323】

データセグメントは、いくつかのソースからの再構築動作で回復されうる。例えば、合
致するデータセグメントがパリティミラーストレージデバイス1602に常駐する場合、
進行の前にパリティミラーストレージデバイス1602からデータセグメントが回復され
うる。別の例においては、データセグメントは利用不可能なデータセグメントのコピーを
含むミラーストレージデバイスから回復されうる。一般的には、回復したデータセグ
メントが1又はそれ以上のパリティミラーストレージデバイス1602に常駐しない場合、
データセグメントはミラーストレージデバイスから回復されるが、合致するデータセグメン

50

トがミラーストレージデバイスで利用可能である場合でさえも、ミラーストレージデバイスから回復されうる。

【0324】

別の例においては、回復したデータセグメントがパリティミラーストレージデバイス1604又はミラーストレージデバイスに常駐しない場合、再生成されるデータセグメントは1又はそれ以上のパリティデータセグメント及びN個のデータセグメントの利用可能なデータセグメントから再生成される。一般的には、欠損したデータセグメントは、いくつかの形態で別のストレージデバイス150に存在しない場合のみにのみ再生成される。

【0325】

別の実施例においては、装置1700はパリティ再構築動作中に置換ストレージデバイスで回復したパリティデータセグメントを再構築するパリティ再構築モジュール1726を具え、回復したパリティデータセグメントは、利用不可能なパリティミラーストレージデバイスに記憶された、利用不可能なパリティデータセグメントに合致させる。利用不可能なパリティミラーストレージデバイスは、1又はそれ以上のパリティミラーストレージデバイス1602のうちの一つである。パリティ再構築動作は置換ストレージデバイスにパリティデータセグメントを復元して、利用不可能なパリティミラーストレージデバイスに前に記憶されたパリティデータセグメントに合致させる。

【0326】

再構築動作中に回復したパリティデータセグメントを再生成するために、再構築用に用いられたデータは様々なソースを元にできる。一例においては、回復したパリティデータセグメントは、ストライプのミラーコピーを記憶するストレージデバイス150の第2のセット内の、パリティミラーストレージデバイス1602に記憶されたパリティデータセグメントを用いて回復される。ミラーコピーが利用可能な場合、回復したパリティデータセグメントは再算出する必要がないため、ミラーリングパリティデータセグメントを用いることが所望される。別の例においては、N個のデータセグメントがN個のストレージデバイスで利用可能な場合、回復したパリティデータセグメントは、N個のストレージデバイス150のうちの一つに記憶されたN個のデータセグメントから再生成される。一般的には、単一の故障が、パリティミラーストレージデバイス1602が再構築される際に生じた場合、N個のデータセグメントはN個のストレージデバイス150で利用可能となる。

【0327】

別の例においては、N個のデータセグメントのうちの一つ又はそれ以上が第1のストレージデバイスセット1604のN個のストレージデバイス150から利用不可能であり、合致するパリティデータセグメントがストレージデバイス150の第2のセットで利用可能でない場合、回復したパリティデータセグメントは、N個のデータセグメントのコピーを記憶するストレージデバイス150の第2のセットの一つ又はそれ以上のストレージデバイス150から再生成される。更に別の例においては、回復したパリティデータセグメントは、ストレージデバイス150の一つ又はそれ以上のセット内の位置に拘らず、利用可能なデータセグメント及び不合致のパリティデータセグメントから再生成される。

【0328】

パリティミラーストレージデバイスがストレージデバイスセット1604のストレージデバイス150間で代替された場合、一般的にはデータ再構築モジュール1724及びパリティ再構築モジュール1726は、再構築されるストレージデバイス150でデータセグメント及びパリティデータセグメントを再構築すると共に動作する。第2のパリティミラーストレージデバイス1602bが利用可能である場合、データ再構築モジュール1724及びパリティ再構築モジュール1726は、ストレージデバイスセット1604の2のストレージデバイス150、1602の故障後に、2のストレージデバイスを再構築できる。パリティミラーストレージデバイス1602がパリティミラーデータセグメントを生成するように進行されない場合、パリティミラーストレージデバイス1602が進行され、ストライプ用のパリティデータセグメントが算出されて記憶され、パリティデータセ

10

20

30

40

50

グメントを算出するのに用いられるパリティミラーストレージデバイス1602のN個のデータセグメントが消去された場合よりも、データセグメント又はストレージデバイス150の回復は迅速である。

【0329】

図18は本発明による進行型RAIDを用いてデータセグメントを更新するための装置1800の一実施例を示した概略ブロック図である。一般的には、装置1800はRAIDグループに属し、パリティミラーストレージデバイス1602のうち1又はそれ以上が進行され、パリティデータセグメントを含み、パリティデータセグメントを生成するのに用いられるN個のデータセグメントを含まない。装置1800は更新受信モジュール1802、更新コピーモジュール1804、パリティ更新モジュール1806を具え、以下

10

【0330】

ストライプ、データセグメント、ストレージデバイス150、ストレージデバイスセット1604、パリティデータセグメント、及び1又はそれ以上のパリティミラーストレージデバイス1602は図17の装置1700に関連して上述したようなストライプと実質的に同様である。装置1800は更新されたデータセグメントを受信する更新受信モジュール1802を具え、更新されたデータセグメントは既存のストライプの既存のデータセグメントに対応する。別の実施例においては、更新受信モジュール1802は複数の更新

20

【0331】

装置1800は、対応する既存のデータセグメントが記憶されるストレージデバイス150に、及び、既存のストライプに対応する1又はそれ以上のパリティミラーストレージデバイス1602に、更新されたデータセグメントをコピーする更新コピーモジュール1804を具える。別の実施例においては、更新コピーモジュール1804はパリティミラーストレージデバイス1602に、又は、既存のデータセグメントを記憶するストレージデバイス150に更新されたデータセグメントをコピーし、その後、更新されたデータセグメントのコピーが他のデバイス1602、150に転送されるのを確認する。

【0332】

装置1800はストレージ統合動作に応じて既存のストライプの1又はそれ以上のパリティミラーストレージデバイス用の1又はそれ以上の更新されたパリティデータセグメントを算出するパリティ更新モジュール1806を具える。ストレージ統合動作は図17の装置1700に関連して上述したようなストレージ統合動作と同様である。ストレージ統合動作は1又はそれ以上の更新されたパリティデータセグメントを有する1又はそれ以上のパリティミラーストレージデバイス1602の少なくともストレージ空き領域及び/又はデータを回復するように導かれる。1又はそれ以上のパリティデータセグメントを更新するのを待機することによって、ストレージ空き領域を統合することがより便利に又は必要になるまで、更新が延期できる。

30

【0333】

一実施例においては、更新されたパリティデータセグメントは、既存のパリティデータセグメント、更新されたデータセグメント、及び既存のデータセグメントから算出される。一実施例においては、既存のデータセグメントは更新されたパリティデータセグメントの生成用の既存のデータセグメントを読み出す前に、所定の位置に維持される。この実施例の利点は、パリティミラーストレージデバイス1602、又は、更新されたパリティデータセグメントが生成される他の位置に既存のデータセグメントをコピーすることに関連するオーバーヘッドが、必要になるまで延期されうることである。この実施例の欠点は、既存のデータセグメントを維持するストレージデバイス150が故障した場合に、更新されたパリティデータセグメントが生成されうる前に、既存のデータセグメントが回復されなければならないことである。

40

50

【0334】

別の実施例においては、既存のデータセグメントが記憶されるN個のストレージデバイス150a-nのストレージデバイス150が、更新されたデータセグメントのコピーを受信する場合に、既存のデータセグメントはデータ-ミラーストレージデバイス1602にコピーされる。既存のデータセグメントは次いでストレージ統合動作までに記憶される。別の実施例においては、既存のデータセグメントはN個のストレージデバイス150a-nのストレージデバイス150のストレージ統合動作に応じてデータ-ミラーストレージデバイス1602にコピーされ、更新されたパリティデータセグメントの算出をトリガするストレージ統合動作の前に、ストレージ統合動作が生じる場合に、既存のデータセグメントが記憶される。後者の実施例は、既存のデータセグメントが記憶されるストレージ

10

【0335】

一実施例においては、更新されたパリティデータセグメントは、既存のパリティデータセグメント、更新されたデータセグメント、及びデルタデータセグメントから算出され、デルタデータセグメントは更新されたデータセグメントと、既存のデータセグメントとの間の差として生成される。一般的には、デルタデータセグメントを生成するステップは、パリティデータセグメントを更新する際の部分的解決又は介在ステップである。デルタデータセグメントを生成するステップは、高度に圧縮可能であり、送信前に圧縮されうるために有利である。

20

【0336】

一実施例においては、デルタデータセグメントは、更新されたパリティデータセグメントの生成用のデルタデータセグメントを読み出す前に、既存のデータセグメントを記憶するストレージデバイスに記憶される。別の実施例においては、既存のデータセグメントが記憶されるストレージデバイス150が、更新されたデータセグメントのコピーを受信する場合に、デルタデータセグメントはデータ-ミラーストレージデバイス1602にコピーされる。別の実施例においては、デルタデータセグメントは、既存のデータセグメントが記憶されるストレージデバイス150のストレージ統合動作に応じて、データ-ミラーストレージデバイス1602にコピーされる。既存のデータセグメントをコピーすると同様に、後者の実施例は、既存のデータセグメントを記憶するストレージデバイス150のストレージ統合動作、又は、更新されたパリティデータセグメントの算出をトリガする別のストレージ統合動作のうち早い方まで、デルタデータファイルが移動しないため、有利である。

30

【0337】

様々な実施例においては、モジュール1802、1804、1806の動作部分の総て、すなわち、更新されたデータセグメントを受信するステップ、更新されたデータセグメントをコピーするステップ、及び更新されたパリティデータセグメントを算出するステップは、ストレージデバイスセット1604のストレージデバイス150、クライアント114、又はサードパーティのRAID管理デバイスで生じる。別の実施例においては、ストレージ統合動作は更新受信モジュール1802及び更新コピーモジュール1804の動作から自律して導かれる。

40

【0338】

図19は、本発明による進行型RAIDを用いたデータを管理するための方法1900の一実施例を示した概略フローチャート図である。方法1900は開始し(1902)、ストレージ要求受信モジュール1702はデータを記憶する要求を受信し(1904)、データがファイル又はオブジェクトのデータとなる。ストライピングモジュール1704はデータ用のストライプパターンを算出し、N個のデータセグメントをN個のストレージデバイス150に書き込む(1906)。ストライプパターンは1又はそれ以上のストライプを含む。各ストライプはN個のデータセグメントのセットを含み、N個のデータセグ

50

メントの各々が、ストライプに割り当てられるストレージデバイス 1604 のセット内の別個のストレージデバイス 150 に書き込まれる。

【0339】

パリティミラーモジュール 1706 は、ストレージデバイス 1604 のセット内の 1 又はそれ以上のパリティミラーストレージデバイス 1602 に、ストライプの N 個のデータセグメントのセットを書き込む (1908)。1 又はそれ以上のパリティミラーストレージデバイスは、N 個のストレージデバイス 150 a - n に追加した状態にある。パリティ生成モジュール 1708 は待ち状態のストレージ統合動作があるかどうかを判定する (1910)。パリティ生成モジュール 1708 が待ち状態のストレージ統合動作はないと判定した (1910) 場合、方法 1900 は戻り待ち状態のストレージ統合動作があるかどうかを再度判定する (1910)。別の実施例においては、ストレージ要求受信モジュール 1702、ストライピングモジュール 1704、及びパリティミラーモジュール 1706 はストレージ要求を受信し続け、ストライピングパターンを算出し、データセグメントを記憶する。

10

【0340】

パリティ生成モジュール 1708 が、待ち状態のストレージ統合動作がないと判定した場合、パリティ生成モジュール 1708 はストライプ用のパリティデータセグメントを算出する (1914)。パリティデータセグメントはパリティミラーストレージデバイス 1602 に記憶された N 個のデータセグメントから算出される。パリティ生成モジュール 1708 はパリティミラーストレージデバイス 1602 にパリティデータセグメントを記憶し (1912)、方法 1900 は終了する (1916)。ストレージ統合動作は、N 個のデータセグメントを記憶する要求を受信し (1904)、N 個のデータセグメントを N 個のストレージデバイスに書き込み (1906)、あるいは、N 個のデータセグメントを 1 又はそれ以上のパリティミラーストレージデバイスに書き込むことから自律して導かれる。ストレージ統合動作は少なくともパリティミラーストレージデバイス 1602 の少なくともストレージ空き領域又はデータを回復するように導かれる。

20

【0341】

図 20 は、本発明による進行型 RAID を用いてデータセグメントを更新するための方法 2000 の一実施例を示す概略フローチャート図である。方法 2000 は開始し (2002)、更新受信モジュール 1802 は更新されたデータセグメントを受信し (2004)、更新されたデータセグメントが既存のストライプの既存のデータセグメントに対応する。更新コピーモジュール 1804 は対応する既存のデータセグメントが記憶されたストレージデバイス 150 に、及び既存のストライプに対応する 1 又はそれ以上のパリティミラーストレージデバイス 1602 に、更新されたデータセグメントをコピー (2006) する。

30

【0342】

パリティ更新モジュール 1806 はストレージ統合動作が待ち状態かどうかを判定する (2008)。パリティ更新モジュール 1806 が待ち状態のストレージ統合動作がないと判定した (2008) 場合、パリティ更新モジュール 1806 はストレージ統合動作を待つ。一実施例においては、方法 2000 は戻って、他の更新されたデータセグメントを受信し (2004)、更新されたデータセグメントをコピーする (2006)。パリティ更新モジュール 1806 が待ち状態のストレージ統合動作はないと判定した (2008) 場合、パリティ更新モジュール 1806 は、既存のストライプの 1 又はそれ以上のパリティミラーストレージデバイス用の 1 又はそれ以上の更新されたパリティデータセグメントを算出し (2010)、方法 2000 は終了する (2012)。

40

【0343】

フロントエンド分散型 RAID

【0344】

従来の RAID システムはデータを受信し、データ用のストライピングパターンを算出し、データをデータセグメントに分割し、パリティストライプを算出し、ストレージデバ

50

イスのデータを記憶し、データセグメントを更新する、等をするのに機能する R A I D コントローラで構成される。いくつかの R A I D コントローラによって、いくつかの機能が分散可能となるが、R A I D コントローラにより管理されるストレージデバイスは、R A I D でストライピングされるデータを記憶するために直接にはクライアント 1 1 4 と通信しない。代わりに、ストレージ要求及び R A I D するためのデータはストレージコントローラを通過する。

【 0 3 4 5 】

R A I D に記憶されるべきデータの総てに接触するために R A I D コントローラを要求することは、データフローのボトルネックを生成するため非効率的である。このことは読み出し - 変更 - 書き込み処理中は特に当てはまり、バンド幅及び R A I D グループのドライブの総ての性能が消費される一方、サブセットだけが実際には更新される。更には R A I D コントローラによって管理されるデータ用に指定されたストレージデバイスの領域は、一般的には R A I D グループ専用であり、別個にアクセスできない。クライアントによるストレージデバイス 1 5 0 へのアクセスは、一般的にストレージデバイス 1 5 0 をパーティショニングすることによって得られる。パーティショニングが用いられる場合、汎用ストレージ用のアクセス可能なパーティションは R A I D 用に用いられず、R A I D グループに割り当てられたパーティションは汎用データストレージ用にアクセス可能ではない。包括的に利用を最適化するためにパーティションをオーバスクライブする方式は、複雑であり管理するのがより難しい。更に、ある R A I D グループ用に割り当てられたストレージ空き領域は、1 の R A I D コントローラがマスタとして指定されない限り、1 以上の R A I D コントローラによってアクセスされず、マスタ R A I D コントローラが停止、非機能性等でない限り、他の R A I D コントローラがスレイブとして動作する。

10

20

【 0 3 4 6 】

一般的な R A I D コントローラは R A I D グループのストレージデバイス 1 5 0 の外側にパリティデータセグメントを更に生成する。パリティデータセグメントは一般的に生成されて、その後ストレージのためにストレージデバイス 1 5 0 に送信され、R A I D コントローラの算出能力を要求するためこれは非効率になりうる。更に、パリティデータセグメントの位置及び更新を追跡することは、ストレージデバイス 1 5 0 で自律して行う代わりに、R A I D コントローラで行わなければならない。

【 0 3 4 7 】

別個の R A I D コントローラがオフラインである場合にデータが利用可能のままであることを保証することが必要である場合、R A I D コントローラは一般的には、ドライブに及び互いに相互接続、及び / 又は完全なセットしてミラーリングされ、データ利用可能性を管理するのを高価及び難易にし、ストレージサブシステムの信頼性を劇的に低下させる。

30

【 0 3 4 8 】

必要となるのは、データセグメントごと、オブジェクトごと、ファイルごと、又は同様のベースで R A I D を可能にし、R A I D コントローラ及び、クライアントとストレージデバイスとの間に位置している対の R A I D コントローラの必要性を除去する、フロントエンド分散型 R A I D 用のシステム、装置及び方法である。このようなシステム、装置及び方法において、R A I D グループは、あるデータセグメント、オブジェクト、又はファイルのために生成され、ある R A I D コントローラによってストレージデバイスの 1 のグループ内で管理されうる一方、第 2 の R A I D グループは第 1 の R A I D グループの同一のストレージデバイスの一部の周りにある別のデータセグメント、オブジェクト、又はファイルのために生成されうる。R A I D 制御機能はクライアント 1 1 4、サードパーティの R A I D 管理デバイスの間で、又はストレージデバイス 1 5 0 の間で分散されうる。フロントエンド分散型 R A I D システム、装置、及び方法は R A I D グループのストレージデバイス 1 5 0 にコマンドを更に送信でき、ダイレクトメモリアクセス (「 D M A 」) 又はリモート D M A (「 R D M A 」) を通してデータに直接的にアクセスし、コピーすることをストレージデバイス 1 5 0 に許容できる。

40

50

【0349】

図10は本発明によるフロントエンド分散型RAIDに対しアクセスされうるシステム1600の一実施例を示す概略ブロック図である。進行型RAIDに関連する図10に示した構成に対する上の記載が、フロントエンド分散型RAIDに更に適用可能である。フロントエンド分散型RAIDに関し、ストレージデバイスセット1604はRAIDグループを形成し、自律的にかつ、ネットワーク116又は1又はそれ以上の冗長ネットワーク116上でクライアント114からストレージ要求を別個に受信及び提供できるストレージデバイス150を具える。

【0350】

ストレージデバイスセット1604内のストレージデバイス150の間から、1又はそれ以上がストライプ用のパリティミラーストレージデバイス1602として指定される。一般的には1又はそれ以上のパリティミラーストレージデバイス1602は、他のストレージデバイス150と実質的に同様に機能する。指定されたパリティミラーストレージデバイス1602がストレージデバイスセット1604のストレージデバイス150の間で代替する場合の一般的な構成においては、パリティミラーストレージデバイス1602はノンパリティミラーストレージデバイスとして更に動作しなければならないために、他のストレージデバイス150と同様の特性を実質的に有している。同様の特性はRAIDグループ内の動作、及び、上述したような別個のクライアント114の通信に対する自律動作に関して存在する。様々な実施例においては、ストレージデバイスセット1604のストレージデバイス150は、述べられたRAID環境内で機能することに関連しない他の態様においては異なってもよい。

【0351】

ストレージデバイスセット1604のストレージデバイス150は、1又はそれ以上のサーバ112内にグループ化されたスタンドアロンにでき、各々がサーバ112に常駐でき、1又はそれ以上のサーバ112を通してアクセスされうる等ができる。1又はそれ以上のクライアント114は、1又はそれ以上のストレージデバイス150を具えるサーバ112に常駐でき、別個のサーバ112に常駐でき、1又はそれ以上のコンピュータネットワーク116を通してストレージデバイス150をアクセスするコンピュータ、ワークステーション、ラップトップ等に常駐できる等ができる。

【0352】

一実施例においては、ネットワーク116はシステムバスを具え、ストレージデバイスセット1604のストレージデバイス150、1602の1又はそれ以上は、システムバスを用いて通信する。例えば、システムバスはPCI-eバス、シリアルアドバンスドテクノロジーアタッチメント(「シリアルATA」)バス、パラレルATA等にできる。別の実施例においては、システムバスは小型コンピュータシステムインタフェース(「SCSI」)、FireWire、ファイバチャネル、USB、PCIe-AS、インフィニバンド等のような外部バスである。当該技術分野の当業者は、自律的であり、1又はそれ以上のネットワーク116上のクライアント114からストレージ要求を別個に受信及び提供することができるストレージデバイス150を有する他のシステム1600の構成が分かるであろう。

【0353】

図11は本発明によるフロントエンド分散型RAID用の装置2100の一実施例を示す概略ブロック図である。様々な実施例においては、装置2100はストレージ要求受信モジュール2102、ストライピング連結モジュール2104、パリティミラー連結モジュール2106、ストレージ要求送信モジュール2108、フロントエンドパリティ生成モジュール2110、パリティ代替モジュール2112、データセグメント回復モジュール2114、データ再構築モジュール2116、パリティ再構築モジュール2118、及びピアツーピア通信モジュール2120を具え、以下に述べられている。様々な実施例においては、装置2100はソリッドステートストレージデバイス102のようなストレージデバイス150、ソリッドステートストレージコントローラ104のようなストレージ

デバイスコントローラ 1 5 2、サーバ 1 1 2、サードパーティの R A I D 管理デバイス等に含まれ、1 以上の構成の間で分散されうる。

【 0 3 5 4 】

ストレージデバイスセット 1 6 0 4 にデータを記憶するストレージ要求を受信するストレージ要求受信モジュール 2 1 0 2 を具える。データはファイル又はオブジェクトの一部にでき、あるいは、完全なファイル又はオブジェクトにできる。ファイルは任意の情報のブロック、又は情報を記憶するためのリソースを含み、コンピュータプログラムで利用できる。ファイルはプロセッサによってアクセスされるいずれかのデータ構造を含みうる。ファイルは、データベース、テキストの記号列、コンピュータコード等を含みうる。オブジェクトは一般的にはオブジェクト指向型プログラミング用のデータ構造であり、データを有する又は有しない構造を含みうる。一実施例においては、オブジェクトはファイルのサブセットである。別の実施例においては、オブジェクトはファイルから独立している。いずれのケースにおいても、オブジェクト又はファイルは、完全なデータのセット、データ構造、コンピュータコード、及びストレージデバイスで記憶されうる他の情報を含むように本明細書において規定されている。

10

【 0 3 5 5 】

ストレージデバイスセット 1 6 0 4 は、1 又はそれ以上のネットワーク 1 1 6 上のクライアント 1 1 4 からストレージ要求を独立して受信する R A I D グループを形成する自律型ストレージデバイス 1 5 0 を具える。ストレージデバイスセット 1 6 0 4 内の自律型ストレージデバイス 1 5 0 の 1 又はそれ以上は、ストライプ用のパリティミラーストレージデバイス 1 6 0 2 として指定される。別のクライアント 1 1 4 からの他のストレージ要求は、第 2 のストレージデバイスセットで記憶され、第 2 のストレージデバイスセットは、第 1 のストレージデバイスセット 1 6 0 4 と同一のストレージデバイス 1 5 0 (及び、パリティミラーストレージデバイス 1 6 0 2) の 1 又はそれ以上を具えることができる。双方のストレージデバイスセット 1 6 0 4 に共通のストレージデバイス 1 5 0 は、共通ストレージデバイス 1 5 0 内の重複する割り当てられたストレージ空き領域を有することができる。

20

【 0 3 5 6 】

装置 2 1 0 0 はデータ用のストライプパターンを算出するストライピング連結モジュール 2 1 0 4 を具える。ストライプパターンは 1 又はそれ以上のストライプを含む。各ストライプは N 個のデータセグメントのセットからなる。ストライプの N 個のデータセグメントは 1 又はそれ以上の空データセグメントを更を含みうる。ストライピング連結モジュール 2 1 0 4 は、ストライプに割り当てられたストレージデバイスセット 1 6 0 4 の N 個のストレージデバイス 1 5 0 a - n のうちの 1 つと、N 個のデータセグメントの各々を連結させる。一実施例においては、ストライピング連結モジュール 2 1 0 4 は、ストレージ要求を送信するクライアント 1 1 4 からのデータセグメントに対応するデータを得るためにストレージデバイス 1 5 0 に命令する、ストレージデバイス 1 5 0 に送信すべきストレージ要求とともに、データセグメントをストレージデバイス 1 5 0 と連結させる。

30

【 0 3 5 7 】

別の実施例においては、ストレージ要求はデータセグメントのデータを実質的に含まない。データを実質的に含まないことは、ストレージ要求が一般的に、ストレージ要求の対象となるデータを含まず、データの一部である、キャラクタ、キャラクタ列等を含みうることを意味する。例えば、データがゼロの配列のような、反復する同一キャラクタの配列を含む場合、ストレージ要求は、データ中に含まれる総てのゼロを含まないゼロの配列を含むという指示を含みうる。当該技術分野の当業者は、大部分のデータを送信しないが、ストレージ要求中の特定のキャラクタ又はキャラクタ列の少量又は単一のインスタンスを許容する他の方法が分かるであろう。ストレージ要求は、N 個のストレージデバイス 1 5 0 a - n が D M A 又は R D M A 動作等を用いてデータを検索するのを可能にするコマンドを含む。

40

【 0 3 5 8 】

50

別の実施例においては、ストライピング連結モジュール2104は、ストレージデバイス150へ送信すべきストレージ要求内で、データセグメントのデータを識別することによって、データセグメントをストレージデバイス150と連結させる。データセグメントのデータを識別することは、データセグメント識別子、データセグメント位置若しくはアドレス、データセグメント長、又はどのデータがデータセグメントを含むのかをデータセグメントが認識するのを可能にする他の情報を含みうる。

【0359】

一実施例においては、ストライピング連結モジュール2104はストレージ要求においてデータセグメントをストレージデバイス150と連結させて、クライアント114がデータセグメントを含むデータをブロードキャストにおいて送信できるようにし、各ストレージデバイス150が付随するデータセグメントを記憶し、ストレージデバイス150に割り当てられないデータセグメントに対応するデータを廃棄できるようにする。別の実施例においては、ストライピング連結モジュール2104は、可能であれば各データセグメントをアドレス付けすることによって、ストレージ要求においてデータセグメントをストレージデバイス150と連結させて、クライアント114がデータセグメントを含むデータをマルチキャストにおいて送信できるようにし、各ストレージデバイス150が付随するデータセグメントを記憶し、ストレージデバイス150に割り当てられないデータセグメントに対応するデータを廃棄できるようにする。当該技術分野の当業者は、1又はそれ以上のデータセグメントを1又はそれ以上のストレージデバイスにブロードキャスト、マルチキャスト、ユニキャスト、エニーキャスト等するために、データセグメントをストレージデバイス150と連結させるストライピングモジュール2104用の他の方法が分かるであろう。

【0360】

関連する実施例においては、ストライピング連結モジュール2104はストレージ要求においてデータセグメントをストレージデバイス150と連結させて、クライアント114がストレージ要求をブロードキャスト、マルチキャスト、ユニキャスト等をできるようにし、ストレージデバイス150と連結するデータセグメントに属するクライアント114からストレージ要求の一部を各ストレージデバイス150が受信できるようにし、ストレージデバイス150と連結する1又はそれ以上のデータセグメントに属しないストレージ要求部分を廃棄できるようにする。

【0361】

別の実施例においては、ストレージ要求受信モジュール2102により受信されたストレージ要求は、ストレージ要求の対象となるデータを含み、ストライピング連結モジュール2104は、データセグメントを含むストレージデバイス150のために、ストレージ要求を与えることによって、データセグメントをストレージデバイス150と連結させる。ストライピング連結モジュール2104はクライアント114、サードパーティのRAID管理デバイス、ストレージデバイス150、1602等内で動作可能である。

【0362】

装置2100は、N個のデータセグメントのセットを、ストレージデバイスセット1604の1又はそれ以上のパリティミラーストレージデバイス1602と連結させるパリティミラー連結モジュール2106を具える。1又はそれ以上のパリティミラーストレージデバイス1602はN個のストレージデバイス150a-nに加えて存在する。一実施例においては、パリティミラー連結モジュール2106はN個のデータセグメントのセットを、各パリティミラーストレージデバイス1602に連結させて、各パリティミラーストレージデバイス1602は、ストライプのN個のデータセグメントを受信及び記憶して、パリティデータセグメントを生成できるようにする。別の実施例においては、パリティミラー連結モジュール2106はストライプのデータセグメントを各パリティミラーストレージデバイス1602と連結させて、パリティミラーストレージデバイス1602a-mはN個のストレージデバイス150a-nに記憶されるN個のデータセグメントに対するミラーとして動作するようになる。

10

20

30

40

50

【0363】

様々な実施例においては、パリティミラー連結モジュール2106は、単一のストレージ要求、複数のストレージ要求、又はDMA、RDMA、ブロードキャスト、マルチキャスト用のパリティミラーストレージデバイス1602を設定し、あるいは、ストレージ要求中にN個のデータセグメントを含むストレージ要求のような、ストライピング連結モジュール2104に関連して上述したような他の関連技術を用いて、N個のデータセグメントのセットを1又はそれ以上のパリティミラーストレージデバイス1602と連結させる。パリティミラー連結モジュール2106は、クライアント114、サードパーティのRAID管理デバイス、ストレージデバイス150、1602等の中で動作可能である。

【0364】

装置2100はストレージデバイスセット1604において、1又はそれ以上のストレージ要求を各ストレージデバイス150、1602に送信するストレージ要求送信モジュール2108を具え、各ストレージ要求はストレージデバイス150、1602に、ストレージ要求を受信するストレージデバイス150、1602に連結する1又はそれ以上のデータセグメントを記憶するのに十分である。一実施例においては、各ストレージ要求はストレージ要求の対象となるデータを含まない。更なる実施例においては、各ストレージ要求は、ストレージデバイスセット1604のN個のストレージデバイス150及びパリティミラーストレージデバイス1602が、DMA又はRDMAを用いて付随するデータセグメントのデータをダウンロードするのを可能にする。別の実施例においては、ストレージ要求は、関連するストレージ要求又はクライアント114のブロードキャストからの、付随するデータセグメント用の関連するデータを選定するのに十分な情報を含む。別の実施例においては、ストレージ要求は、付随するデータセグメントのデータを含む。

【0365】

一実施例においては、各ストレージ要求はストライプのストレージデバイスセット1604の一部であるストレージデバイス150、1602を識別する。ストレージデバイスセット1604のストレージデバイス150、1604の識別を含むことによって、マスタとして動作するストレージデバイス150の故障の場合には、別のストレージデバイス150が、RAIDされたデータを管理するためのマスタとして引き継ぐことができる。別の実施例においては、ストレージデバイスセット1604の識別によって、自律型ストレージデバイス150、1602は、ストレージデバイスがオフラインの場合にはデータを回復し、置換ストレージデバイスが、クライアントと別個にストレージデバイスセット1604内に付加される場合にはデータを再構築することを可能にする。別の実施例においては、ストレージデバイスセット1604のストレージデバイス150、1602の識別は、データセグメント又はストレージ要求の送信用のマルチキャストグループを表わしている。この識別はストレージデバイスセット1604内のストレージデバイス150、1602に記憶されたオブジェクト又はファイル用のメタデータと共に、記憶されうる。

【0366】

一実施例においては、パリティミラー連結モジュール2106がN個のデータセグメントのセットを1又はそれ以上のパリティミラーストレージデバイス1602の各々と連結させる場合、装置2100はクライアント114とは別個に、ストライプ用のパリティデータセグメントを算出し、パリティミラーストレージデバイス1602にパリティデータセグメントを記憶するフロントエンドパリティ生成モジュール2110を具える。パリティデータセグメントは、パリティミラーストレージデバイス1602に提供されるN個のデータセグメントのセットから算出される。1以上のパリティミラーストレージデバイス1602がストレージデバイスセット1604に含まれる場合、フロントエンドパリティ生成モジュール2110は一般的に様々なパリティデータセグメントを生成して、ストレージデバイスセット1604内の2又はそれ以上のストレージデバイス150、1602は故障可能になり、パリティデータセグメント情報が、利用不可能なデータセグメント又はパリティデータセグメントの回復を可能にする。

【0367】

別の実施例においては、フロントエンドパリティ生成モジュール 2 1 1 0 は、ストレージデバイスセット 1 6 0 4 のストレージデバイス 1 5 0 及び / 又はサードパーティの R A I D 管理デバイス内で動作する場合、パリティデータセグメントを算出する。例えば、ストレージ要求を送信するクライアント 1 1 4 と別個のサーバ 1 1 2 は、パリティデータセグメントを算出できる。別の実施例においては、フロントエンドパリティ生成モジュール 2 1 1 0 はパリティミラーストレージデバイス内で動作して、パリティデータセグメントを算出する。例えば、パリティミラーストレージデバイス 1 6 0 2 内のストレージコントローラ 1 5 2 は、ストレージデバイスセット 1 6 0 4 によって形成される R A I D グループ用のマスタストレージコントローラとして動作できる。

【 0 3 6 8 】

10

別の実施例においては、フロントエンドパリティ生成モジュール 2 1 1 0 はパリティデータセグメントを算出し、次いで、算出されたパリティデータセグメントを、ミラーを形成する第 2 のストレージデバイスのセットの 1 又はそれ以上の更なるパリティミラーストレージデバイス 1 6 0 4 に送信する。この実施例は、パリティデータセグメントを算出することに関連するオーバーヘッドが、ネットワーク 1 1 6 上のデータトラフィックを低減させる更なる利益を有する各ストレージデバイスセット 1 6 0 4 に対してではなく、一度に消費されるため、有利である。

【 0 3 6 9 】

一実施例においては、ストレージデバイス 1 5 0 が利用不可能であり、要求が、利用不可能なデータセグメントか、利用不可能なデータセグメントを含むデータのいずれかを読み出すように受信される場合、装置 2 1 0 0 はストレージデバイスセット 1 6 0 4 のストレージデバイス 1 5 0 に記憶されたデータセグメントを回復するデータセグメント回復モジュール 2 1 1 2 を具える。データセグメントは、ストレージデバイスセット 1 6 0 4 の利用可能なストレージデバイス 1 5 0 のデータセグメント、又は、ストレージデバイスセット 1 6 0 4 の利用可能なストレージデバイス 1 5 0、1 6 0 2 のパリティデータセグメント及びデータセグメントの組合せを用いて、あるいは、データセグメントのコピーを含むミラーストレージデバイスから回復される。一般的には、ミラーストレージデバイスは N 個のデータセグメントのコピーを記憶するミラーストレージデバイスセットの 1 のストレージデバイス 1 5 0 である。データセグメント回復モジュール 2 1 1 2 は動作して、ストレージデバイス 1 5 0、パリティミラーストレージデバイス 1 6 0 2、サードパーティの R A I D 管理デバイス、ミラーストレージデバイス等から利用不可能なデータセグメントを回復できる。

20

30

【 0 3 7 0 】

別の実施例においては、装置 2 1 0 0 は再構築動作中に、置換ストレージデバイス 1 5 0 に回復したデータセグメントを記憶するデータ再構築モジュール 2 1 1 4 を具える。例えば、ストレージデバイス 1 5 0 が故障、同期の損失等のために利用不可能になった場合、データ再構築モジュール 2 1 1 4 はストレージデバイス 1 5 0 を再構築して、利用不可能なストレージデバイス 1 5 0 を置換できる。一実施例においては、再構築されるストレージデバイス 1 5 0 は、利用可能にされた元のストレージデバイス 1 5 0 である。

【 0 3 7 1 】

40

回復したデータセグメントは、ストレージデバイスセット 1 6 0 4 の利用不可能なストレージデバイス 1 5 0 に記憶された利用不可能なデータセグメントに合致させる。再構築動作は一般的に、データセグメント及びパリティデータセグメントの 1 又はそれ以上を置換ストレージデバイス 1 5 0 に復元して、利用不可能なストレージデバイス 1 5 0 に前に記憶されたデータセグメント及びパリティデータセグメントに合致させる。

【 0 3 7 2 】

一実施例においては、回復したデータセグメントはストレージデバイスセット 1 6 0 2 の利用可能なストレージデバイス 1 5 0 で利用可能なデータセグメントを用いて、再構築動作によって回復される。別の実施例においては、回復したデータセグメントは、パリティミラーストレージデバイス 1 6 0 2 の 1 又はそれ以上からのパリティデータセグメント

50

と、ストレージデバイスセット 1604 の利用可能なストレージデバイス 150 の利用可能なデータセグメントとの組合せを用いて、再構築動作によって回復される。別の実施例においては、回復したデータセグメントは、パリティミラーストレージデバイス 1602 から読み出された、合致するデータセグメントを用いて、再構築動作により回復される。更なる別の実施例においては、回復したデータセグメントは、ミラーストレージデバイスからの合致するデータセグメントを用いて、再構築動作により回復される。データ再構築モジュール 2114 は動作し、クライアント 114、サードパーティの RAID 管理デバイス、ストレージデバイス 150、1602、ミラーストレージデバイス等からの受信したデータセグメントを記憶するのである。

【0373】

別の実施例においては、装置 2100 は再構築動作で、置換ストレージデバイス 1602 に回復したパリティデータセグメントを再構築するパリティ再構築モジュール 2116 を具える。再構築動作は、データ再構築モジュール 2114 に関連して記載した再構築動作と実質的に同様である。パリティ再構築モジュール 2116 が回復したパリティデータセグメントを再構築することを除いては、パリティ再構築モジュール 2116 は、データ再構築モジュール 2114 と同様に動作する。回復したパリティデータセグメントは、ストライプに割り当てられた利用不可能なパリティミラーストレージデバイス 1602 に記憶される利用不可能なパリティデータセグメントに合致される。

【0374】

様々な実施例において、ミラーリングストレージデバイスセットのパリティミラーストレージデバイス 1602 に記憶されたパリティデータセグメントをコピーすることによって、ストレージデバイスセット 1604 のパリティミラーストレージデバイス 1602 からパリティデータセグメントをコピーすることによって（利用不可能なパリティデータセグメントと同一の場合）、ストレージデバイスセット 1604 の利用可能なストレージデバイス 150、1602、及びデータセグメントのコピーを含むミラーストレージデバイスに記憶された N 個のデータセグメント及びパリティデータセグメントの 1 又はそれ以上を用いてパリティデータセグメントを生成することによって、等によって回復される。データ再構築モジュール 2116 は動作して、回復したデータセグメントを記憶する一方、クライアント 114、サードパーティの RAID 管理デバイス、ストレージデバイス 150、ミラーストレージデバイス等に常駐する。

【0375】

有利には、装置 2100 は本明細書に記載のフロントエンド分散型 RAID 動作専用のパーティションに、ストレージデバイス 150、1602 のデータを記憶することに限定されない。代わりに、自律型ストレージデバイス（例えば、150a）は、クライアント 114 からストレージ要求を別個に受信して、ストライピング連結モジュール 2104、パリティミラー連結モジュール 2106 及びフロントエンドパリティ生成モジュール 2110 によってデータを記憶するのに更に利用可能なストレージデバイス 150a の 1 又はそれ以上の領域の RAID された、又は RAID されないデータを記憶できる。

【0376】

一実施例においては、ストレージ要求受信モジュール 2102 によって受信され、又は、ストレージ要求送信モジュール 2108 によって送信される 1 又はそれ以上のストレージ要求が、ストライプのストレージデバイスセット 1604 を具えるストレージデバイス 150 を識別する。有利には、ストレージ要求でストレージデバイスセット 1604 のストレージデバイス 150 を識別することは、バックアップ RAID コントローラを促進して、マスタコントローラが非機能性である場合に動作する。例えば、ストレージデバイスセット 1604 のストレージデバイス 150 がストレージ要求で識別され、マスタコントローラがパリティミラーストレージデバイス 1602 にあり、利用不可能である場合、別のパリティミラーストレージデバイス 1602、又は、N 個のストレージデバイス 150 a - n のうちの別のものが、マスタコントローラになりうる。

【0377】

10

20

30

40

50

一実施例においては、装置 2100 は各ストライプに対して、ストレージデバイスセット 1604 内のどのストレージデバイス 150 がストライプ用のパリティミラーストレージデバイス 1602 として指定されるかを代替するパリティ代替モジュール 2118 を具える。パリティ代替モジュール 2118 の利点は上述のとおりである。別の実施例においては、ストレージデバイスセット 1604 のストレージデバイス 150 はピアのグループを形成し、装置 2100 は、ストレージデバイスセット 1604 のストレージデバイス 150、1602 内のストレージ要求を送信及び受信するピアツーピア通信モジュール 2120 を具える。ピアツーピア通信モジュール 2120 はピアデバイスをストレージデバイスセット 1604 の外側に有して、ストレージ要求を更に送信及び受信できる。

【0378】

好ましい実施例においては、ストレージ要求は装置 2100 のモジュール 2102 - 2120 を用いてストレージデバイスセット 1604 のストレージデバイス 150、1602 にわたって、オブジェクトのデータをストライピングすることによって、オブジェクトを記憶するオブジェクト要求である。別の実施例においては、ストレージデバイスセット 1604 の自律型ストレージデバイス 150、1602 の 1 又はそれ以上は、少なくとも一部の第 1 のオブジェクト又はファイル用の第 1 の RAID グループ内に割り当てられ、少なくとも一部の第 2 のオブジェクト又はファイル用の第 2 の RAID グループ内に割り当てられる。例えば、1 のストレージデバイス 150 a は 1 又はそれ以上のストライプ用のストレージデバイスセット 1604 のためのマスター RAID コントローラにでき、第 2 のストレージデバイス 150 b はストレージデバイスセットのストレージデバイス 150 の一部又は総てを含む、RAID グループ用のマスター RAID コントローラにできる。有利には、装置 2100 はストレージデバイス 150、1602 をグループ化して、様々なクライアント 114 のために RAID グループを形成する際にフレキシビリティを与える。

【0379】

図 12 は、本発明によるフロントエンド分散型 RAID 用の方法の一実施例を示す概略フローチャート図である。方法 2200 は開始し (2202)、ストレージ要求受信モジュール 2102 はストレージ要求を受信して (2204)、ストレージデバイスセット 1604 の N 個のストレージデバイス 150 a - n のデータを記憶する。ストライピング連結モジュール 2104 は、データ用のストライプパターンを算出し (2206)、N 個のデータセグメントのそれぞれを、N 個のストレージデバイス 150 a - n の 1 つと連結させる (2208)。

【0380】

パリティミラー連結モジュール 2106 は、N 個のデータセグメントのセットを 1 又はそれ以上のパリティミラーストレージデバイス 1602 と連結させる (2210)。ストレージ要求送信モジュール 2108 は 1 又はそれ以上のストレージ要求をストレージデバイスセット 1604 の各ストレージデバイス 150、1602 に送信する (2212)。各ストレージ要求は、ストレージデバイス 150 に、ストレージ要求を受信するストレージデバイス 150 に関連する 1 又はそれ以上のデータセグメントを記憶するのに十分である。データのデータセグメントは次いで、ストレージ要求によって命令されるように、DMA、RDMA、ブロードキャスト、マルチキャストを用いてストレージデバイスセット 1604 のストレージデバイス 150、1602 に転送される。選択的に、フロントエンドパリティ生成モジュール 2110 はストライプ用のパリティデータセグメントを算出し (2214)、方法 2200 は終了する (2216)。

【0381】

共有されたフロントエンド分散型 RAID

【0382】

従来の RAID は RAID 専用の各ディスクの少なくとも一部で、ディスク又は他のストレージデバイスのアレイを利用する。RAID コントローラは RAID グループに対するストレージ要求を管理する。冗長システムについては、マスター RAID コントローラは

10

20

30

40

50

故障又は利用不可能になった場合に、引き継ぐようなバックアップRAIDコントローラを有している。RAIDに記憶された同一データにアクセスするようになる複数クライアントからのストレージ要求は、一般的に出現順に順々に実行される。

【0383】

図10、11、及び12にそれぞれ示されたシステム1600、装置2100、及び方法2200に関して上述したように、フロントエンド分散型RAIDシステムは、分散型RAIDコントローラとして機能できるストレージコントローラ152をそれぞれ具えることができる自律型ストレージデバイス150を具え、ストレージデバイス150はそれぞれ、複数のクライアント114を供する複数の重複するRAIDグループで構成される。場合により、2のクライアント114は同一データにアクセスするようである。ストレージ要求が最初に到達して実行する場合、一般的にはデータの不一致がない。他方、同一データに対する2又はそれ以上のストレージ要求が同時、又は実質的に同時に到達した場合、データは破損するであろう。

10

【0384】

例えば、データがRAIDグループの4のストレージデバイス150に記憶され、ストレージデバイス150の1つがパリティミラーストレージデバイス1602として割り当てられ、第1のクライアント114がストレージ要求をRAIDコントローラとして動作する第1のストレージコントローラ150aに送信し、第2のクライアント114が第2のストレージ要求を第2のRAIDコントローラとして動作する第2のストレージデバイス150aに送信し、双方のストレージ要求が同一データにアクセスする場合、第1のストレージデバイス150aは、第1のストレージデバイス150a、次いで、RAIDグループの他のストレージデバイス150b-nのストレージ要求の実行を開始できる。同時に、第2のストレージデバイス150bの第2のRAIDコントローラは別のストレージデバイス(例えば、150b)、次いで、RAIDグループの残りのストレージデバイス150a、c-nの第2のストレージ要求の実行を開始できる。この実行における不整合は、ストレージデバイス150の間の物理的な距離、実行時間の差異等によって生じうる。結果として破損したデータとなる。

20

【0385】

必要となるのは、同一データのアクセスを求める同時のストレージ要求を処理する共有されたフロントエンド分散型RAIDのためのシステム、装置及び方法である。有益には、システム、装置及び方法はデータへのアクセスを制御し、1のストレージ要求は第2のストレージ要求が実行される前に、実行される。

30

【0386】

図10は、進行型RAID及びフロントエンド分散型RAIDに加えて、本発明による共有されたフロントエンド分散型RAIDに対するシステム1600として働く一実施例を示す概略ブロック図である。進行型RAID及びフロントエンド分散型RAIDに関連する図10に示した構成に対する上の記載が、共有されたフロントエンド分散型RAIDに更に適用可能である。フロントエンド分散型RAIDと同様、ストレージデバイスセット1604はRAIDグループを形成し、自律的にかつ、ネットワーク116上でクライアント114からストレージ要求を別個に受信及び提供できるストレージデバイス150を具える。

40

【0387】

共有されたフロントエンド分散型RAIDに関し、システム1600は2又はそれ以上のクライアント114のそれぞれが同一データに関するストレージ要求を送信するような、2又はそれ以上のクライアント114を具える。1のストレージ要求が別のストレージ要求の別の到達前に完了しないようにストレージ要求が到達する点で、ストレージ要求は同時である。ストレージデバイスセット1604内のストレージデバイス150の間から、1又はそれ以上がストライプ用のパリティミラーストレージデバイス1602として指定される。一般的には、1又はそれ以上のパリティミラーストレージデバイス1602は他のストレージデバイス150と実質的に同様に機能する。

50

【0388】

指定されたパリティミラーストレージデバイス1602がストレージデバイスセット1604のストレージデバイス150の間で代替する場合の一般的な構成においては、パリティミラーストレージデバイス1602はノンパリティミラーストレージデバイスとして更に動作しなければならないために、他のストレージデバイス150と同様の特性を実質的に有している。同様の特性はRAIDグループ内の動作、及び、上述したような別個のクライアント114の通信に対する自律動作に関して存在する。様々な実施例においては、ストレージデバイスセット1604のストレージデバイス150は、述べられたRAID環境内で機能することに関連しない他の態様においては異なってもよい。

【0389】

ストレージデバイスセット1604のストレージデバイス150は、1又はそれ以上のサーバ112内にグループ化されたスタンドアロンにでき、各々がサーバ112に常駐でき、1又はそれ以上のサーバ112を通してアクセスされうる等ができる。1又はそれ以上のクライアント114は、1又はそれ以上のストレージデバイス150を具えるサーバ112に常駐でき、別個のサーバ112に常駐でき、コンピュータネットワーク116を通してストレージデバイス150をアクセスするコンピュータ、ワークステーション、ラップトップ等に常駐できる等ができる。

【0390】

一実施例においては、ネットワーク116はシステムバスを具え、ストレージデバイスセット1604のストレージデバイス150、1602の1又はそれ以上は、システムバスを用いて通信する。例えば、システムバスはPCI-eバス、シリアルアドバンスドテクノロジーアタッチメント(「シリアルATA」)バス、パラレルATA等にできる。別の実施例においては、システムバスは小型コンピュータシステムインタフェース(「SCSI」)、FireWire、ファイバチャネル、USB、PCIe-AS、インフィニバンド等のような外部バスである。当該技術分野の当業者は、自律的であり、ネットワーク116上のクライアント114からストレージ要求を別個に受信及び提供することができるストレージデバイス150を有する他のシステム1600の構成が分かるであろう。

【0391】

図13は、本発明による共有されたフロントエンド分散型RAID用の装置2300の一実施例を示す概略ブロック図である。装置2300は、多重ストレージ要求受信モジュール2302、ストライピングモジュール2304、パリティミラーモジュール2306、シーケンサモジュール2308、マスタ確認モジュール2310、マスタ判定モジュール2312、マスタエラーモジュール2314、パリティ生成モジュール2316、及びパリティ代替モジュール2318を具え、以下に述べられている。

【0392】

装置2300は、少なくとも2のクライアント114から少なくとも2のストレージ要求を受信する多重ストレージ要求受信モジュール2302を具えて、ストレージデバイスセット1602のストレージデバイス150にデータを記憶する。データはファイル又はオブジェクトのデータを含む。装置に関連するストレージ要求は、各々が共通のデータの少なくとも一部を有し、更に、1のストレージ要求がストレージ要求の到達前に完了されないように到達することによる同時のストレージ要求である。これらの同時のストレージ要求はフロントエンド分散型RAIDシステム1600中の共通データを破損する危険を冒す。一実施例においては、同時のストレージ要求は1のクライアント114からにできる。別の実施例においては、同時のストレージ要求は2又はそれ以上のクライアント114からになる。

【0393】

複数のストレージ要求がストレージデバイスセット1602のストレージデバイス150に記憶された1又はそれ以上のデータセグメントを更新でき、既に記憶されたデータが、ストレージデバイスセット1602のストレージデバイス150に記憶されたデータセグメント内にストライピングモジュール2304によってストライピングされる。一実施

10

20

30

40

50

例においては、ストレージ要求は初めてのデータをRAIDグループに書き込む。このケースにおいては、データは一般的に他の場所に存在し、1又はそれ以上のクライアント114によってアクセスし、あるストレージ要求はデータをRAIDグループにコピーする一方、別のストレージ要求はデータに同時アクセスする。

【0394】

複数のストレージ要求はストレージデバイスセット1602のストレージデバイス150に記憶される1又はそれ以上のデータセグメントを更新する1の要求と、共通にデータの少なくとも一部分を標的にする1又はそれ以上の読み出し要求を含みうる。更新要求が完全ではない場合に、次いでストレージデバイスセット1602のストレージデバイス150からの読み出し要求応答は、データを破損する既存の更新データの組合せを構成しうる。

10

【0395】

装置2300は同時のストレージ要求の各々に対して、データ用のストライプパターンを算出し、ストライプのN個のデータセグメントをストレージデバイスセット1604内のN個のストレージデバイス150a-nに書き込むストライピングモジュール2304を具える。ストライプパターンは1又はそれ以上のストライプを含み、各ストライプはN個のデータセグメントのセットを含む。N個のデータセグメントの各々はストレージデバイスセット1604内の別個のストレージデバイス150に書き込まれ、ストライプに割り当てられる。装置2300は同時のストレージ要求の各々に対して、ストライプのN個のデータセグメントをパリティミラーストレージデバイス1602として指定されたストレージデバイスセット1604内のストレージデバイス150に書き込むパリティミラーモジュール2306を具える。パリティミラーストレージデバイス1602はN個のストレージデバイス150a-nに加えられて存在する。

20

【0396】

ストライピングモジュール2304はファイル又はオブジェクトの一部である1又はそれ以上のデータセグメントを読み出すことから、1又はそれ以上のストレージデバイス150a-nの同一性を算出するのに更に用いられる。

【0397】

装置2300は、第2のクライアント114からの第2のストレージ要求を実行する前に、第1のクライアント114からの第1のストレージ要求の完了を確認するシーケンサモジュール2308を具え、少なくとも2の同時のストレージ要求は第1及び第2のストレージ要求を含む。他の実施例においては、シーケンサモジュール2308は2又はそれ以上の他の同時のストレージ要求を実行する前に、第1のストレージ要求の完了を確認する。有益には、シーケンサモジュール2308は同時のストレージ要求の順番の実行を促進して、データの損失を防ぐ。一実施例においては、シーケンサモジュール2308は、総てのストレージ要求がデータ用にアクセスされなければならないマスタコントローラを用いることによって、ロッキングシステム、2相コミットを用いることによって、又は当該技術分野に公知の他の手段を用いることによって、同時のストレージ要求の実行を調整する。シーケンサモジュール2308によって用いられる方法の一部は以下に述べられている。

30

40

【0398】

一実施例においては、シーケンサモジュール2308は、第2のストレージ要求の実行前の第1のストレージ要求とともに、ストレージ要求を受信したストレージデバイスセット1602のストレージデバイス150の各々から受取り文字を受信することによって、同時のストレージ要求の実行前に第1のストレージ要求の完了を確認する。一般的には、受取り文字はストレージ要求の完了を通知する。一実施例においては、ストレージ要求によって影響が与えられるストレージデバイス150の各々は、シーケンサモジュール2308が第2のストレージ要求の実行を開始する前に、ストレージデバイス150の各々から受信された受取り文字に書き込まれる。

【0399】

50

一実施例においては、ストレージ要求の完了は、同一のストレージデバイス150aの待ち状態の第2のストレージ要求の一部の実行前に、単一のストレージデバイス（例えば、150a）に命令された第1のストレージ要求の一部の完了を含みうる。シーケンサモジュール2308は、ストレージデバイス150のストレージ要求部分の完了を別個に確認できる。この実施例においては、第1のストレージ要求に関連するデータセグメントを書き込むことは、第1のストレージ要求の総てのデータセグメントが完了されるまで遅らせる必要はない。シーケンサモジュール2308はストレージデバイスセット1604のストレージデバイス150に生じる様々な実行を調整して、データが破損されないことを確認できる。

【0400】

一実施例においては、ストレージ要求の完了についての受取り文字は、ストライピングモジュール2304及びパリティミラーモジュール2306がそれぞれ、ストレージ要求に関連するデータセグメントをストレージデバイスセット1604のストレージデバイス150に書き込んだ後に受信される。別の実施例においては、ストレージ要求の完了についての受取り文字は、ストライピングモジュール2304及びパリティミラーモジュール2306がそれぞれ、ストレージ要求に関連するデータセグメントをストレージデバイスセット1604のストレージデバイス150に書き込み、ストレージデバイス150、1602の各々がデータセグメントが書き込まれたことを確認した後に受信される。

【0401】

一実施例においては、シーケンサモジュール2308は、最初に到達した同時の要求の中から、ストレージ要求を選択することによって、実行用の第1のストレージ要求を選択する。別の実施例においては、シーケンサモジュール2308は、最も早いタイムスタンプを有するストレージ要求を選択することによって、実行用の第1のストレージ要求を選択する。別の実施例においては、シーケンサモジュール2308は、いくつかの選択基準を用いてストレージ要求を選択することによって、実行用の第1のストレージ要求を選択する。例えば、シーケンサモジュール2308は、高優先度として、要求クライアント114によっていくつかの方法で標識されたストレージ要求を選択、好まれるクライアント114からのストレージ要求を選択等を行うことができる。当該技術分野の当業者は、シーケンサモジュール2308がいくつかの選択基準を用いて第1のストレージ要求を選択できる他の方法が分かるであろう。

【0402】

一実施例においては、ストレージ要求受信モジュール2302、ストライピングモジュール2304、パリティミラーモジュール2306、及びシーケンサモジュール2308は、同時のストレージ要求を制御し提供するマスタコントローラ（図示せず）の一部である。マスタコントローラの総て又は一部は、クライアント114、サードパーティのRAID管理デバイス、ストレージデバイスセット1604のストレージデバイス150、又はストレージデバイス150のストレージコントローラ152で常駐及び動作可能である。マスタコントローラを用いてデータ用のサービス要求を実行することによって、シーケンサモジュール2302はデータに向けられたストレージ要求を認識でき、次いで、同時のストレージ要求を認識でき、ストレージデバイスセットのストレージデバイス150に記憶されたデータが破損されない方法で同時のストレージ要求を配列できる。当該技術分野の当業者はデータに向けられたストレージ要求の提供を制御するマスタコントローラの他の実装が分かるであろう。

【0403】

別の実施例においては、マスタコントローラは、1又はそれ以上のクライアント114から同時のストレージ要求を提供できる2又はそれ以上のマスタコントローラのグループの一部であり、ストレージ要求はストレージデバイスセット1604のストレージデバイス150に記憶されたデータに向けられる。例えば、マスタコントローラは第1のクライアント114用のストレージ要求を提供でき、第2のマスタコントローラは第2のクライアント114用のストレージ要求を提供できる。第1及び第2のクライアント114は双

10

20

30

40

50

方ともストレージデバイスセット 1604 のストレージデバイス 150 に記憶されたデータにアクセスでき、このようにして同時のストレージを可能にする。あるマスタコントローラは 1 のストレージデバイス 150 a の一部でできるが、他のマスタコントローラは第 2 のストレージデバイス 150 b の一部にできる。別の実施例においては、第 1 のマスタコントローラは第 1 のストレージデバイスセット 1604 a の一部にでき、第 2 のマスタコントローラはミラーリングストレージデバイスセット 1604 b の一部にできる。

【0404】

マスタコントローラが、ストレージデバイスセット 1604 のストレージデバイス 150 にアクセスするマスタコントローラのグループの一部である場合、装置 2300 は、受信されたストレージ要求を提供するマスタコントローラが、受信されたストレージ要求の
10
実行前の 1 又はそれ以上の同時のストレージ要求の実行より先に、ストレージ要求の実行を制御していることを確認するマスタ確認モジュール 2310 を具えることができる。この実施例においては、同時のストレージ要求が他のマスタコントローラによって受信され、サービス要求は他のマスタコントローラによって受信される同時のストレージ要求と共通の少なくともデータの一部分を有する。

【0405】

例えば、マスタコントローラはストレージ要求を受信でき、マスタ確認モジュール 2310 は次いでストレージ要求の実行前に他のマスタコントローラをポーリングして、そのマスタコントローラが更にストレージ要求のデータ用のマスタコントローラであることを確認できる。確認の一部は、マスタコントローラが相互に通信する確認を含み、指定されたマスタコントローラがストレージ要求の実行前に確認できるようにする。これは、ある
20
フロントエンド RAID がマスタとして指定され、別のものがバックアップである場合に有用にできる。別の例においては、マスタコントローラはストレージ要求を受信して、ファイル又はオブジェクトからデータセグメントを読み出すストレージ要求を受信でき、マスタ確認モジュール 2310 は次いで他のマスタコントローラをポーリングして、更新がファイル又はオブジェクトで進行中ではないことを確認できる。別の例においては、マスタコントローラはマスタ確認モジュールを用いて、ストレージ要求用のデータの制御を取得できる。

【0406】

マスタコントローラがストレージ要求の実行用のマスタであると確認する 1 の方法は、スリーウェイポーリング方式を用いることであり、コントローラが進行するストレージ要求用のマスタであるとみなすために、2 のデバイス/コントローラは利用可能でなければならない。その方式はマスタにすべきと争うコントローラに対しサードパーティとなるデバイス（図示せず）を用いて、どちらのコントローラがマスタにすべく割り当てられるかの記録を維持する。このマスタ確認デバイスは別のコントローラにし、サーバにあるクライアントにする等を行うことができ、マスタコントローラとして動作できるグループでコントローラと通信できる。マスタ確認モジュール 2310 の一部はその後各コントローラに属するマスタ確認モジュール 2310 の一部を伴い、マスタ確認デバイスに常駐できる
30
。

【0407】

一例においては、システム 1600 はその各々をマスタにできる、第 1 のフロントエンド分散型 RAID コントローラ（「第 1 のコントローラ」）及び第 2 のフロントエンド分散型 RAID コントローラ（「第 2 のコントローラ」）と、別個のマスタ確認デバイスとを具える。第 1 及び第 2 のコントローラ及びマスタ確認デバイスは総て互いに通信する。マスタ確認デバイス 2310 は第 1 のコントローラをマスタコントローラとして、及び第 2 のコントローラをストレージデバイスセット 1604 のストレージデバイス 150 に記憶されたデータ用のバックアップとして指定でき、マスタ確認デバイス 2310 はコントローラ及びマスタ確認デバイスの、このマスタ情報を記憶できる。通信が第 1 のコントローラと、第 2 のコントローラと、マスタ確認デバイスとの間で維持される限り、マスタ確認モジュール 2310 は第 1 のコントローラがマスタであると確認できる。
40
50

【0408】

第1のマスターコントローラがストレージ要求を受信し、第2のバックアップコントローラが利用不可能になり、又は、第1のコントローラとマスター確認デバイスで通信が損失した場合、マスター確認モジュール2310はマスター確認デバイスと第1のコントローラとの間の通信を通して第1のコントローラが未だマスターであることを確認し、第1のコントローラとマスター確認デバイスの双方が、第1のコントローラが実際のマスターであることを確認するため、マスター確認モジュール2310によりストレージ要求を進行することが可能となる。第2のバックアップコントローラにより受信されたストレージ要求はマスター確認モジュール2310を通して、第2のコントローラがマスターではないと認識するため、進行しない。

10

【0409】

一方、第1のマスターコントローラが利用不可能であり、又は、第2のバックアップコントローラ及びマスター確認デバイスと通信できず、第2のバックアップコントローラがストレージ要求を受信する場合、マスター確認モジュール2310は、第2のコントローラとマスター確認モジュールの双方が第1のコントローラと通信できないことを認識し、マスター確認モジュール2310は第2のバックアップコントローラをマスターにすべきと指定し、ストレージ要求は進行できる。マスター指定の変更は、その際第2のコントローラに記録される。

【0410】

第1のコントローラが動作し、第2のコントローラ及びマスター確認デバイスで損失した通信のみがある場合、第1のコントローラによって受信されるデータ用のどのストレージ要求も実行されない。通信が復元される場合、第2のコントローラとマスター確認デバイスの双方が、第2のコントローラをマスターと認識するため、第1のコントローラはストレージ要求をこれ以上実行しないであろう。当然にして、このマスター指定はリセットされる。当該技術分野の当業者はマスターコントローラのうちの1つにマスター指定を割り当て及び再割り当てする、様々な静的及び動的手段が分かるであろう。

20

【0411】

マスター確認デバイスが利用不可能であり、第1のストレージコントローラがストレージ要求を受信する場合、第1及び第2のコントローラで動作するマスター確認モジュール2310の一部は、第1のコントローラがマスターであり、ストレージ要求が進行できることを確認できる。第2のコントローラがストレージ要求を受信する場合、第1及び第2のコントローラで動作するマスター確認モジュール2310の一部は、第1のコントローラがマスターであり、ストレージ要求が進行できないことを確認できる。他の実施例においては、2以上のコントローラがポーリング方式の一部である。当該技術分野の当業者はマスター確認モジュール2310がストレージ要求の実行の前にコントローラがマスターであることを確認できる他の方法が分かるであろう。

30

【0412】

別の実施例においては、装置2300はマスター判定モジュール2312を具える。ストレージ要求を送信する前に、マスター判定モジュール2312はマスターコントローラのグループにマスター判定要求を送信する。マスターコントローラのグループは次いで、どちらのコントローラがストレージ要求用のマスターとして指定されるかを識別し、マスターコントローラを識別する応答をマスター判定モジュール2312に送信する。マスター判定モジュール2312はストレージ要求用のマスターコントローラの識別表示を受信し、要求デバイスに命令して、ストレージ要求を指定されたマスターコントローラに送信する。一実施例においては、マスター判定モジュール2312はクライアント114で常駐し動作する。別の実施例においては、マスター判定モジュール2312はサードパーティのRAID管理デバイスで常駐し実行される。別の実施例においては、マスター判定モジュール2312はストレージデバイス150に常駐する。別の実施例においては、マスター判定モジュールは2以上のストレージデバイス150間で分散される。

40

【0413】

50

更なる実施例においては、装置 2300 はエラー表示を返すマスタエラーモジュール 2314 を具える。一実施例においては、マスタコントローラによって制御される多重ストレージ要求受信モジュール 2302 が、マスタコントローラによって制御されないストレージ要求を受信する場合、マスタエラーモジュール 2314 はエラー表示を返す。

【0414】

別の実施例においては、マスタコントローラがストレージ要求の実行完了の時に、判定されたマスタではもはやないと、マスタ判定モジュール 2312 又はマスタ確認モジュール 2310 が判定する場合に、マスタエラーモジュール 2314 はエラー表示を返す。この実施例は一般的にマスタコントローラがストレージ要求を実行し始め、グループの他のマスタコントローラで通信を損失し、又は、ポーリング方式の場合、他のマスタコントローラ及びマスタ確認デバイスで通信を損失する場合に生じる。別の実施例においては、マスタコントローラによって制御される多重ストレージ要求受信モジュール 2302 が、マスタコントローラによって制御されないストレージ要求を受信する場合、マスタエラーモジュール 2312 はエラー表示を返す。

10

【0415】

別の実施例においては、マスタコントローラは 1 又はそれ以上のセカンダリマスタコントローラに対するストレージ要求を制御する。セカンダリマスタコントローラはそれぞれ、ストレージデバイスセット 1604 のストレージデバイス 150、1602 に記憶されたデータ用のストレージ要求を制御する。別の実施例においては、セカンダリマスタコントローラを制御するマスタコントローラは更に、ストレージデバイスセット 1604 のストレージデバイス 150、1602 に記憶されたデータに向けられたストレージ要求用のセカンダリマスタコントローラである。

20

【0416】

別の実施例においては、マスタコントローラは 1 又はそれ以上のセカンダリマスタコントローラに対するストレージ要求を制御し、各セカンダリマスタコントローラはセカンダリマスタコントローラに固有の、ストレージデバイスセットのストレージデバイス 150 に記憶されるデータ用のストレージ要求を制御する。装置 2300 はフレキシブルであり、いずれのマスタコントローラも、セカンダリマスタコントローラとして動作する他のコントローラに対してマスタにできる。いくつかのセカンダリマスタコントローラはストレージデバイスセット 1604 を共有でき、その他のものは異なるストレージデバイスセットを制御できる。他の実施例においては、マスタコントローラはパリティミラーストレージデバイス 1602、又は、N 個のストレージデバイス 150 a - n のうちの 1 つにできる。

30

【0417】

別の実施例においては、マスタコントローラがオフラインであるか、指定されたマスタであると判定できない場合、セカンダリマスタコントローラはマスタコントローラにできる。当該技術分野の当業者は、1 又はそれ以上のセカンダリマスタコントローラの間で、マスタ指定を割り当て及び再割り当てするための様々な静的及び動的手段が分かるであろう。

【0418】

好ましい実施例においては、装置 2300 はストライプ用のパリティデータセグメントを算出し、パリティミラーストレージデバイス 1602 にパリティデータセグメントを記憶するパリティ生成モジュール 2316 を具える。パリティストライプはパリティミラーストレージデバイス 1602 の N 個のデータセグメントのセットから算出される。この実施例は一般的な RAID 5、RAID 6、又はその他の RAID レベルであるが、一般的には RAID 0、RAID 1、RAID 10 等には含まれない。

40

【0419】

別の好ましい実施例においては、装置 2300 は各ストライプに対して、ストレージデバイスセット 1604 内のストレージデバイス 150 のどちらが、ストライプ用の 1 又はそれ以上のパリティミラーストレージデバイス 1602 であると割り当てられるかを代替

50

するパリティ代替モジュール 2318 を具える。ストライプごとにパリティデータセグメントをローテーションすることは、性能を改善する。パリティ代替モジュール 2318 はストライピングモジュール 2304 とともに用いられ、ファイル又はオブジェクトの一部である 1 又はそれ以上のデータセグメントを読み出し、書き込み、又は更新することから、1 又はそれ以上のストレージデバイス 150a - n の同一性を計算できる。

【0420】

様々なモジュール 2302 - 2318 の機能は、単一のマスタコントローラと一緒に常駐でき、又は、1 又はそれ以上のクライアント 114、サードパーティの RAID 管理デバイス、及び 1 又はそれ以上のストレージデバイス 150、1602 の間で分散されうる。当該技術分野の当業者は本明細書に記載の機能が分散される様々な実施例が分かるであろう。

10

【0421】

図 14 は、本発明による共有されたフロントエンド分散型 RAID 用の方法 2400 の実施例を示す概略フローチャート図である。

方法 2400 は開始し (2402)、多重ストレージ要求受信モジュール 2302 は少なくとも 2 のクライアント 114 から少なくとも 2 のストレージ要求を受信して (2404)、ストレージデバイスセット 1604 の 1 又はそれ以上のストレージデバイス 150 のデータを読み出すか記憶する。データは、ファイル又はオブジェクトからのものであり、ストレージ要求はそれぞれ、共通にデータの少なくとも一部分を有し、到達によって同時になり、別の少なくとも 2 のストレージ要求の到達前にストレージ要求が完了しない。ストライピングモジュール 2304 はデータ用のストライプパターンを算出し (2406)、ストライプパターンは 1 又はそれ以上のストライプを含み、各ストライプは N 個のデータセグメントのセットを含む。ストライピングモジュール 2304 は更に、ストレージデバイスセット 1604 内の N 個のストレージデバイス 150a - n にストライプの N 個のデータセグメントを読み出し又は書き込み (2408)、N 個のデータセグメントの各々は、別個のストレージデバイス 150 に書き込まれ又はそこから読み出される。

20

【0422】

ストレージ要求が書き込み動作である場合、パリティミラーモジュール 2306 はストライプの N 個のデータセグメントのセットをストレージデバイスセット 1604 内の 1 又はそれ以上のパリティミラーストレージデバイス 1602 に書き込み (2410)、パリティミラーストレージデバイス 1602 は N 個のストレージデバイス 150a - n に加えた状態になる。パリティミラーモジュール 2306 は更に、パリティミラーデバイス 1602 に記憶されたデータセグメント、又はパリティデータセグメントを読み出すことができる (2410)。シーケンサモジュール 2308 は第 2 のクライアント 114 からの第 2 のストレージ要求を実行する前に、第 1 のクライアント 114 からの第 1 のストレージ要求の完了を確認し (2412)、方法 2400 は終了する (2416)。第 1 及び第 2 のストレージ要求は同時のストレージ要求である。

30

【0423】

本発明は、本発明の精神または本質的特徴から逸脱することなく、その他の特定の形態で実施され得る。示された実施例は、例示のみであって、限定ではないと総ての点で見なされるべきである。本発明の範囲は、従って、上述した説明ではなく、添付の請求の範囲によって規定される。請求項と同等の意味および範囲にある総ての変更は、本発明の範囲に包含されるべきである。

40

【0424】

高容量不揮発性ストレージ用のキャッシュとしてのソリッドステートストレージ

【0425】

一般的にキャッシュは、しばしばアクセスされ、アプリケーション又はオペレーティングシステムの一部として配置されるデータが、ハードディスクドライブ (「HDD」)、光学ドライブ、テープストレージ等のような高容量不揮発性 (「HCNV」) ストレージデバイスを通してアクセスされなければならない場合よりも更に早い次のアクセスを伴う

50

キャッシュ内に記憶されるので有利である。キャッシュは一般的にはコンピュータ内に含まれる。

【0426】

いくつかのストレージデバイス及びシステムはHCNVストレージデバイスにキャッシュを具備している。いくつかのHCNVストレージデバイスは不揮発性ソリッドステートキャッシュを含み、これらはアクセス時間を低減する利益を提供するが、通常は限定されるHCNVストレージデバイスインタフェースの容量と一致する性能を提供するのみである。マザーボード上に一般的に配置されるいくつかの不揮発性ソリッドステートキャッシュストレージデバイスが存在するが、これらのデバイスはキャッシュコヒーレンスが提供されないようなマルチクライアント環境において用いることができない。HCNVデバイスのいくつかのコントローラは更にキャッシュを具備する。冗長HCNVキャッシュコントローラが複数クライアント間で共有される場合、データが破損されていないことを確認するのに高性能なキャッシュコヒーレンシアルゴリズムが要求される。

10

【0427】

一般的に、キャッシュはDRAM内に実装され、キャッシュ容量をプレミアムにし、性能ごとに比較的高い電源を要求する。揮発性キャッシュを支持する電源が失われる場合、キャッシュに記憶されたデータが失われる。一般的にバッテリーバックアップが、電源故障の場合にデータ損失を避けるのに用いられるが、バッテリーバックアップの停止前に不揮発性メモリに対するキャッシュを消去するかなりの可能性を有する。更に、バッテリーバックアップシステムは電力を消費し、冗長性を要求し、信頼性に負の影響を与え、空き領域を消費する。バッテリーは標準ベースで更に提供されなければならない、バッテリーバックアップは比較的高価である。

20

【0428】

前述の考察から、キャッシュとしてソリッドステートストレージを用いてデータを管理する装置、システム及び方法に対するニーズが存在することは明らかである。有利には、このような装置、システム及び方法は少ない電力を消費し、かなり大きな容量を有し、キャッシュに記憶されたデータを維持するのにバッテリーバックアップを要求しない不揮発性キャッシュを提供するであろう。

【0429】

図15は、本発明による大容量不揮発性ストレージデバイス用のキャッシュとしてソリッドステートストレージ110を有するシステム3400の一実施例を示す概略ブロック図である。システム3400はソリッドステートストレージコントローラ104及びHCLVコントローラ3402と、ソリッドステートストレージ110と、ネットワークインタフェース156とを具備するストレージコントローラ152を有するソリッドステートストレージデバイス102を具備する。システム3400はコンピュータネットワーク116を通してソリッドステートストレージデバイス102に接続された要求デバイス155と、1又はそれ以上のHCNVストレージデバイス3404a-nとを具備する。当該技術分野の当業者は図15に示されたシステム3400が単なる一実施例であって、ソリッドステートストレージ110をストレージデバイス用のキャッシュにできる多くのその他の構成が可能であることが分かるであろう。

30

40

【0430】

システム3400はネットワークインタフェース156とストレージコントローラ152とを有するソリッドステートストレージデバイス102を具備する。別の実施例においては、ネットワークインタフェース156はソリッドステートストレージデバイス102の外側にある。例えば、ネットワークインタフェース156はソリッドステートストレージデバイス102を具備する又は具備しないサーバ112内にできる。

【0431】

示された実施例においては、ソリッドステートストレージデバイス102は、ソリッドステートストレージコントローラ104と、大容量不揮発性(「HCNV」)ストレージコントローラ3402とを具備するストレージコントローラ152を具備する。別の実施例に

50

おいては、ソリッドステートストレージデバイス 102 はストレージコントローラ 152 内にない、ソリッドステートストレージコントローラ 104 及び HCNV ストレージコントローラ 3402 を具える。他の実施例においては、ソリッドステートストレージデバイス 102 は HCNV ストレージコントローラ 3402 を含むソリッドステートストレージコントローラ 104 を具え、その逆もできる。

【0432】

示された実施例においては、システム 3400 は統合型ソリッドステートストレージ 110 を有し、外部 HCNV ストレージデバイス 3404 a - n を有するソリッドステートストレージデバイス 102 を具える。別の実施例においては、ストレージコントローラ 152、104、3402 はソリッドステートストレージ 110 から分離できる。別の実施例においては、コントローラ 152、104、3402 及びソリッドステートストレージ 110 は、HCNV ストレージデバイス 3404 に含まれる。HCNV ストレージデバイス 3404 はネットワークインタフェース 156 を更に具えることができる。当該技術分野の当業者は、多くのその他の構成が可能であることが分かるであろう。ソリッドステートストレージデバイス 102、ソリッドステートストレージコントローラ 104、ソリッドステートストレージ 110、ストレージ I/O バス 210、ネットワークインタフェース 156、コンピュータネットワーク 116、及び要求デバイス 155 は上述のデバイス及びバスの他の実施例と実質的に同様である。

10

【0433】

一実施例においては、要求デバイス 155 はシステムバスを通してソリッドステートストレージデバイス 102、ストレージコントローラ 152、ソリッドステートストレージコントローラ 104 等に接続される。要求デバイス 155 と、ソリッドステートストレージ 110 との間のデータ転送は、システムバスで生じうる。

20

【0434】

HCNV ストレージデバイス 3404 は一般的には不揮発性ストレージを提供する高容量ストレージデバイスであり、一般的にはソリッドステートストレージ 110 よりもデータを読み出すこと及び書き出すことについては遅い。HCNV ストレージデバイス 3404 はソリッドステートストレージ 110 よりも単位あたりのストレージ容量が安くなる。HCNV ストレージデバイス 3404 はハードディスクドライブ(「HDD」)、光学ドライブ、テープストレージ等にできる。HCNV ストレージデバイス 3404 用のキャッシュとしてソリッドステートストレージ 110 を提供することは、一般的にデータアクセス及びストレージの速度を増加させる。当該技術分野の当業者は、HCNV ストレージデバイス 3404 用のキャッシュとしてのソリッドステートストレージ 110 の利点があるであろう。

30

【0435】

一実施例においては、HCNV ストレージデバイス 3404 はストレージエリアネットワーク(「SAN」)を通してストレージコントローラ 152 に接続される。一実施例においては、別個の SAN コントローラが HCNV ストレージデバイス 3404 をストレージコントローラ 152 に接続する。別の実施例においては、HCNV ストレージコントローラ 3402 又はストレージコントローラ 152 は SAN コントローラとして作用する。当該技術分野の当業者は HCNV ストレージデバイス 3404 が SAN に接続されうる他の方法があるであろう。

40

【0436】

図 16 は、本発明による高容量不揮発性ストレージデバイス用のキャッシュとしてのソリッドステートストレージを有する装置 3500 の一実施例を示した概略ブロック図である。装置 3500 はキャッシュフロントエンドモジュール 3502、キャッシュバックエンドモジュール 3504、オブジェクトストレージコントローラ 3506、HCNV モジュール 3508、及び標準デバイスエミュレーションモジュール 3510 を具え、以下に説明されている。装置 3500 のモジュール 3502 - 3510 は、ソリッドステートストレージコントローラ 104 と、HCNV ストレージコントローラ 3402 とを有するス

50

トレージコントローラ 152 に示されているが、各モジュール 3502 - 3510 の一部及び総てがソリッドステートストレージコントローラ 104、HCNVストレージコントローラ 3402、サーバ 112、HCNVストレージデバイス 3404、又はその他の場所に含まれている。

【0437】

装置 3500 はストレージ要求と関連するデータ転送を管理するキャッシュフロントエンドモジュール 3502 を具え、データ転送は要求デバイス 155 とソリッドステートストレージ 110 との間で、1 又はそれ以上の HCNV ストレージデバイス 3404 a - n 用のキャッシュとして機能している。装置 3500 はソリッドステートストレージ 110 と、HCNV ストレージデバイス 3404 a - n との間でデータ転送を管理するキャッシュバックエンドモジュール 3504 を更に具える。データ転送はデータ、メタデータ、及び/又はメタデータインデックスを含みうる。上述のように、ソリッドステートストレージ 110 は不揮発性のソリッドステートデータストレージエレメントのアレイ 216、218、220 であり、一般的にバンク 214 に配置される。様々な実施例においては、ソリッドステートストレージ 110 はフラッシュメモリ、ナノランダムアクセスメモリ(「ナノRAM」又は「NRAM」)、磁気抵抗RAM(「MRAM」)、ダイナミックRAM(「DRAM」)、相変化RAM(「PRAM」)等にできる。

10

【0438】

一般的に、キャッシュフロントエンドモジュール 3502、キャッシュバックエンドモジュール 3504、及びソリッドステートストレージコントローラ 104 は要求デバイス 155 から自律して動作する。例えば、要求デバイス 155 は単一のストレージデバイスとして付随するストレージ 110、3404 a - n とともに、ソリッドステートストレージコントローラ 104 及び HCNV ストレージコントローラ 3402 を有するストレージコントローラ 152 を表示する。別の例においては、要求デバイス 155 は HCNV ストレージデバイス 3404 a - n を表示でき、ソリッドステートストレージ 110 は非表示にできる。

20

【0439】

一実施例においては、ソリッドステートストレージコントローラ 104 は 1 又はそれ以上の要求デバイス 155 からオブジェクト要求を提供し、ソリッドステートストレージ 110 内のオブジェクト要求のオブジェクトを管理するオブジェクトストレージコントローラモジュール 3506 を具える。その実施例においては、ソリッドステートコントローラ 104 はオブジェクトストレージコントローラモジュール 3506 とともに、上述のような、特に図 2A に示した装置に関連して記載されるような、オブジェクト要求を管理する。

30

【0440】

一実施例においては、装置 3500 は RAID レベルの一致する独立ドライブ冗長アレイ(「RAID」)の 2 又はそれ以上の HCNV ストレージデバイス 3404 a - n 内のソリッドステートストレージ 110 にキャッシュされたデータを記憶する HCNV RAID モジュール 3508 を具える。その実施例においては、データは全体として要求デバイス 155 に見え、RAID することは要求デバイス 155 から隠されるようになる。例えば、キャッシュフロントエンドモジュール 3502 はソリッドステートストレージ 110 に要求デバイス 155 からデータをキャッシュでき、キャッシュバックエンドモジュール 3504 は、HCNV RAID モジュール 3508 と提携して、データをストライプし、RAID レベルの一致する HCNV ストレージデバイス 3404 a - n にデータセグメント及びパリティデータセグメントを記憶できる。当該技術分野の当業者は、要求デバイス 155 からのデータが HCNV ストレージデバイス 3404 a - n に RAID できる他の方法が分かるであろう。

40

【0441】

別の実施例においては、ソリッドステートストレージ 110 及び HCNV ストレージデバイスは、RAID グループとして構成されたハイブリッドストレージデバイスセット内

50

のハイブリッドストレージデバイスを用意する。例えば、ハイブリッドストレージデバイスセットは、フロントエンド分散型RAIDストレージデバイスセット1604にでき、ハイブリッドストレージデバイスは図10、11、及び12にそれぞれ示されたシステム1600、装置2100、及び方法2200に関連して上述のように設定されたストレージデバイスのストレージデバイス150、1602にできる。その実施例においては、ソリッドステートストレージ110にキャッシュされ、HCNVデバイス3404で後に記憶されるデータセグメントは、ストライプのN個のデータセグメントのうちの1つ、又はストライプのパリティデータセグメントである。フロントエンドRAIDと同様、ハイブリッドストレージデバイスはRAIDストライプのデータセグメントと別個に、1又はそれ以上のクライアント114からストレージ要求を受信する。

10

【0442】

更なる実施例においては、ハイブリッドストレージデバイスは、2又はそれ以上のクライアント114から2又はそれ以上の同時のストレージ要求を受信する共有されたフロントエンド分散型RAIDグループ1604のストレージデバイス150、1602であり、図10、13、及び14にそれぞれ示されたシステム1600、装置2300、及び方法2400に関連して述べたような共有されたフロントエンドRAIDに関連して上述される。有益には、この実施例は、共有された冗長キャッシュが、更なる複雑なコヒーレンスアルゴリズム及びプロトコルを有することなく、コヒーレンスを維持することを保証する。

20

【0443】

別の実施例においては、ソリッドステートストレージ110及びHCNVストレージデバイス3404a-nはハイブリッドストレージデバイスを用意し、装置3500はハイブリッドストレージデバイスの動作に特有のコードで、要求デバイス155をロードする前に、1又はそれ以上の要求デバイス155に取り付けられた標準デバイスをエミュレートすることによって、ハイブリッドストレージデバイスにアクセスを提供する標準デバイスエミュレーションモジュール3510を用意する。その実施例においては、標準デバイスは工業規格BIOSによって担持される。このブートストラップ動作によって、ハイブリッドデバイスが、ソリッドステートストレージコントローラ104、HCNVストレージコントローラ3404、及び装置3500のできる限りの他のモジュール3502-3510に特有のドライバが、要求デバイス155にロードできるまで、限られた機能性を有する要求デバイス155によって認知、アクセスされることが可能である。

30

【0444】

一実施例においては、ソリッドステートストレージデバイス110は2又はそれ以上の領域に分割可能であり、1又はそれ以上のパーティションが、HCNVストレージデバイス3404用のキャッシュとして機能するソリッドステートストレージと別個の、ソリッドステートストレージとして用いられる。例えば、ソリッドステートストレージ110のいくつかのパーティションは汎用データストレージ用のクライアント114によってアクセスできる一方、1又はそれ以上のパーティションはHCNVストレージデバイス3404用のキャッシュとして割り当てられる。

40

【0445】

一実施例においては、1以上のクライアント114（又は要求デバイス155）はキャッシュフロントエンドモジュール3502及びキャッシュバックエンドモジュール3504にキャッシュ制御メッセージを送信して、ソリッドステートストレージデバイス110及び1又はそれ以上のHCNVストレージデバイス3404内に記憶されるファイル又はオブジェクトのうちの1又はそれ以上の状態を管理できる。有益には、ファイルごと、オブジェクトごと、又はデータセグメントベースごとのキャッシュを管理するクライアント114/要求デバイス155の能力は、ソリッドステートストレージデバイス110を共有するようかなりのフレキシビリティを提供する。

【0446】

多数のキャッシュ制御メッセージが許容可能であり、実現可能である。例えば、キャッ

50

シュ制御メッセージは、キャッシュバックエンドモジュール3504に、ソリッドステートストレージ110内のオブジェクト又はファイルの一部を留める(p i n)制御メッセージを含みうる。別のキャッシュ制御メッセージは、キャッシュバックエンドモジュール3504に、ソリッドステートストレージ110内のオブジェクト又はファイルの一部を解放する(u p p i n)制御メッセージを含みうる。別のキャッシュ制御メッセージは、キャッシュバックエンドモジュール3404に、ソリッドステートストレージ110から1又はそれ以上のH C V Nストレージデバイス3404までのオブジェクト又はファイルの一部を消去させる制御メッセージを含みうる。別のキャッシュ制御メッセージは、キャッシュバックエンドモジュール3404に、1又はそれ以上のH C V Nストレージデバイス3404からソリッドステートストレージ110に、オブジェクト又はファイルの一部を事前ロードさせる制御メッセージを含みうる。別のキャッシュ制御メッセージは、ソリッドステートストレージ110に、決められたストレージ空き領域量を空けるために、キャッシュバックエンドモジュール3504に、ソリッドステートストレージから1又はそれ以上のH C V Nストレージデバイス3404までの1又はそれ以上のオブジェクト又はファイルの1又はそれ以上の部分をオフロードさせる制御メッセージを含みうる。当該技術分野の当業者は、他の可能なキャッシュ制御メッセージが分かるであろう。

10

【0447】

一実施例においては、キャッシュ制御メッセージはオブジェクト又はファイル用のメタデータ(「キャッシュ制御メタデータ」)を通して通信される。一実施例においては、キャッシュ制御メタデータは永続的である。別の実施例においては、キャッシュ制御メタデータはファイル又はオブジェクトの生成時に設定される属性を通して構築される。その実施例においては、属性は特定のオブジェクトクラス、デフォルト、特定のファイル型の特徴等に対する関係を通して継承されうる。当該技術分野の当業者は、キャッシュ制御メッセージがメタデータを通して通信されうる他の方法が分かるであろう。

20

【0448】

一実施例においては、システム3400は揮発性キャッシュストレージエレメントを具える。例えば、ソリッドステートストレージ110に加えて、システム3400は揮発性のいくつかのタイプのランダムアクセスメモリ(「RAM」)を更に具えることができる。この実施例においては、キャッシュフロントエンドモジュール3502及びキャッシュバックエンドモジュール3504は揮発性キャッシュストレージエレメントにいくつかのデータを記憶し、ソリッドステートストレージ110及び揮発性キャッシュストレージエレメントに記憶されたデータを管理し、バックエンドストレージモジュール3504は、揮発性キャッシュストレージエレメントと、ソリッドステートストレージと、H C V Nストレージデバイスとの間のデータ転送を更に管理する。例えば、重要ではない、又は、別のソースから簡単に回復されうるデータは揮発性キャッシュに記憶できる一方、その他のデータは、キャッシュとして機能するソリッドステートストレージ110に記憶できる。

30

【0449】

更なる実施例においては、H C V Nストレージデバイス3404に記憶されたオブジェクト及びファイル用のメタデータ及び/又はインデックスメタデータは、ソリッドステートストレージデバイス110内及び揮発性キャッシュストレージエレメントに維持される。図2Aに示された装置200に関連して上述したように、特定のメタデータは揮発性キャッシュストレージエレメントに記憶でき、揮発性キャッシュストレージエレメント中のデータが損失した場合に、インデックスを再構築するのに用いられうる。一実施例においては、メタデータ及びインデックスメタデータは、ソリッドステートストレージ110に記憶され、揮発性キャッシュストレージエレメントは含まれない。当該技術分野の当業者は、キャッシュとして機能するソリッドステートストレージ110と共に揮発性キャッシュストレージエレメントを用いる他の利点及び方法が分かるであろう。

40

【0450】

図17は、本発明による高容量不揮発性ストレージデバイス用のキャッシュとしてソリ

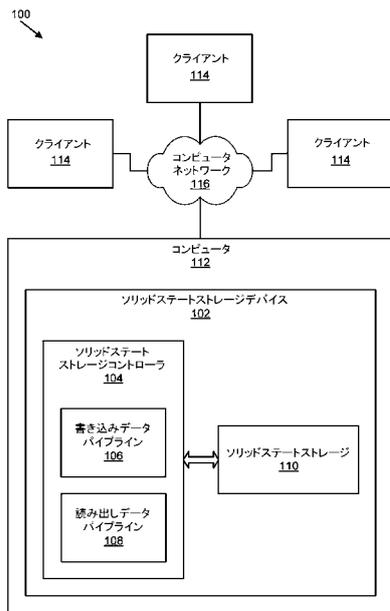
50

ッドステートストレージを有する方法 3600 の一実施例を示した概略フローチャート図である。方法 3600 は開始し (3602)、キャッシュフロントエンドモジュール 3502 はストレージ要求と関連するデータ転送を管理し (3604)、データ転送は、要求デバイス 155 とソリッドステートストレージ 110 との間で、1 又はそれ以上の HCNV ストレージデバイス 3404 a - n 用のキャッシュとして機能している。キャッシュバックエンドモジュール 3504 は、ソリッドステートストレージ 110 と、1 又はそれ以上の HCNV ストレージデバイス 110 との間のデータ転送を管理し (3606)、方法 3600 は終了する (3608)。方法 3600 は図 10 の装置 3500 に関連して上述したのと実質的に同様に動作する。

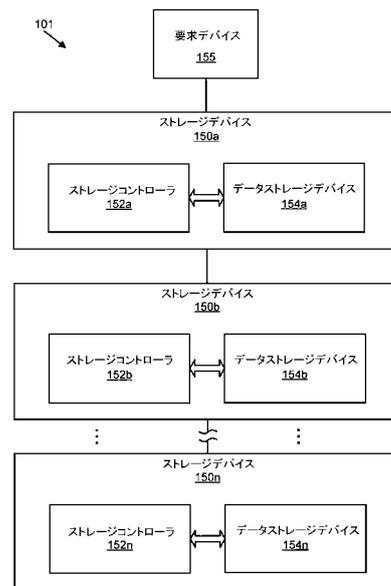
【0451】

本発明は、本発明の精神または本質的特徴から逸脱することなく、その他の特定の形態で実施され得る。示された実施例は、例示のみであって、限定ではないと総ての点で見なされるべきである。本発明の範囲は、従って、上述した説明ではなく、添付の請求の範囲によって規定される。請求項と同等の意味および範囲にある総ての変更は、本発明の範囲に包含されるべきである。

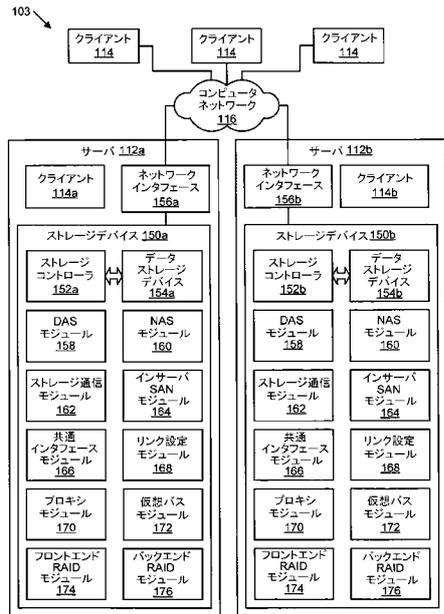
【図 1 A】



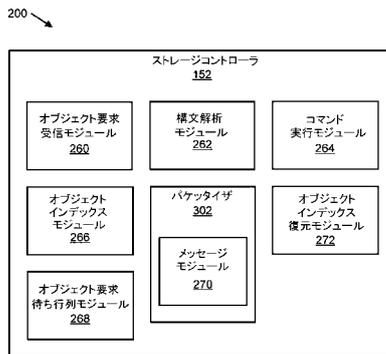
【図 1 B】



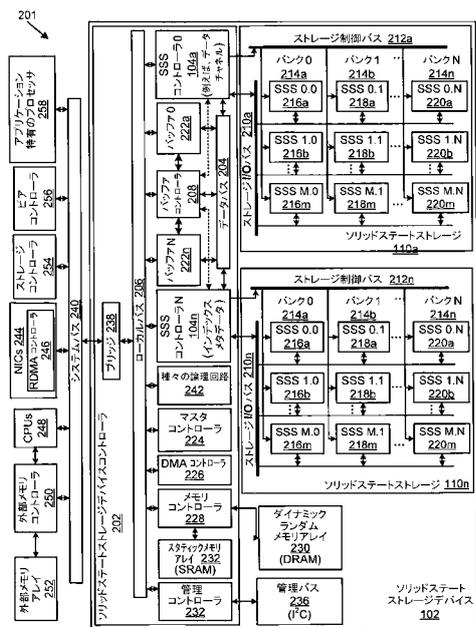
【図1C】



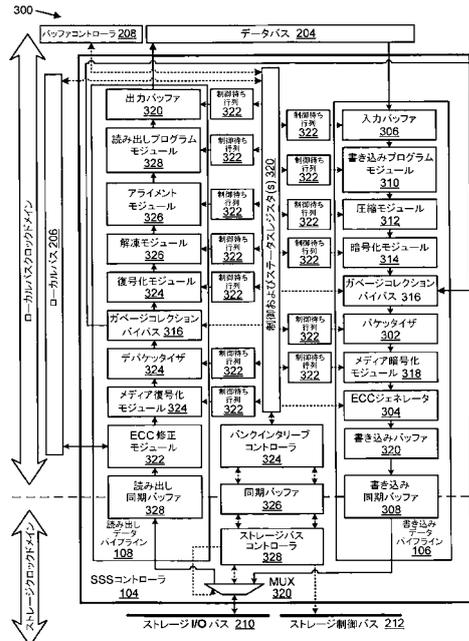
【図2A】



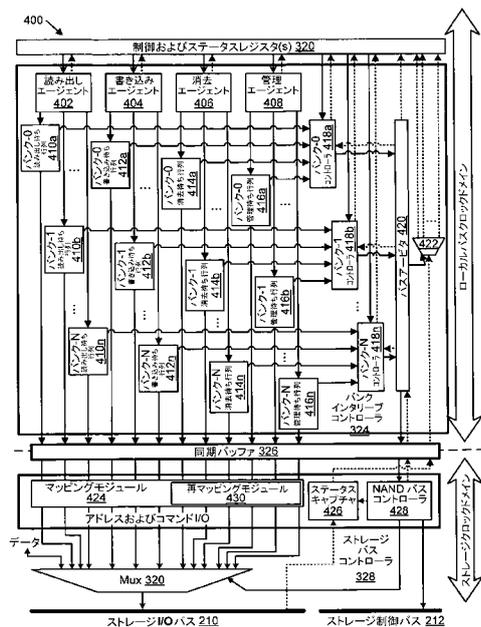
【図2B】



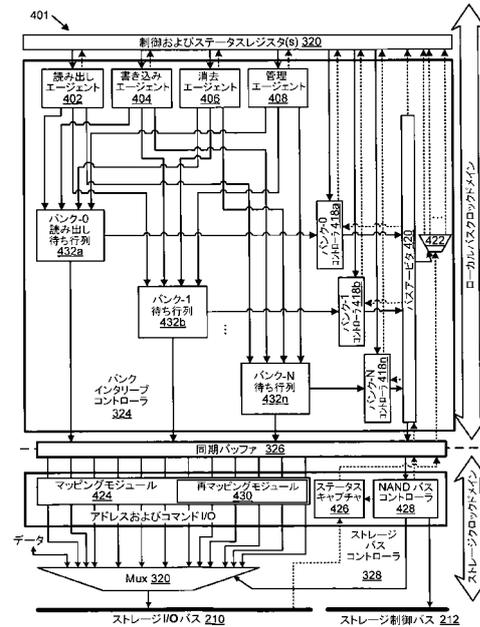
【図3】



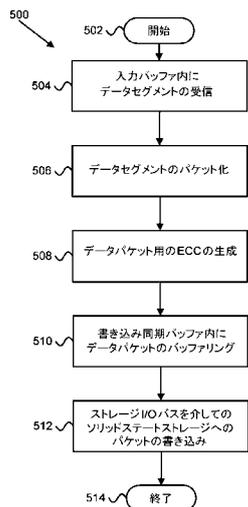
【図 4 A】



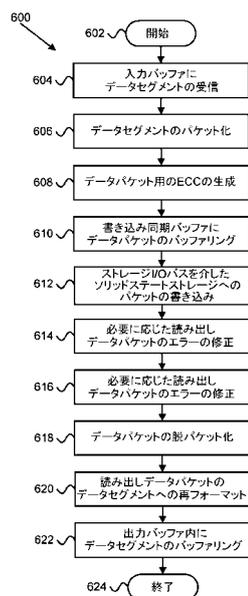
【図 4 B】



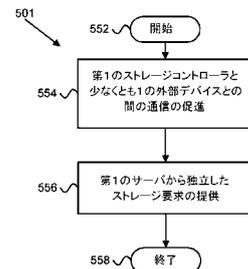
【図 5 A】



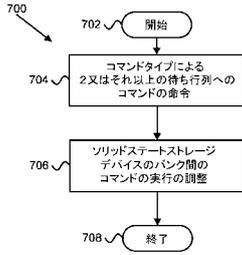
【図 6】



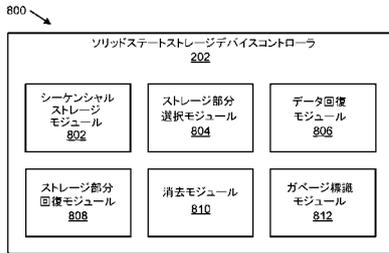
【図 5 B】



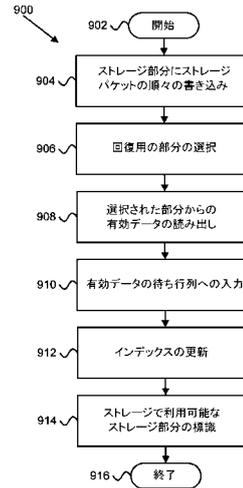
【 図 7 】



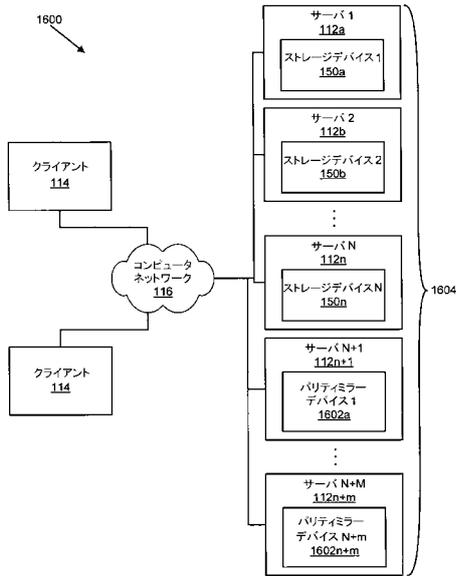
【 図 8 】



【 図 9 】



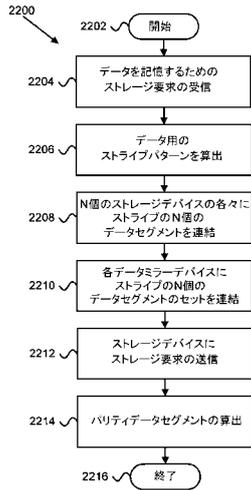
【 図 10 】



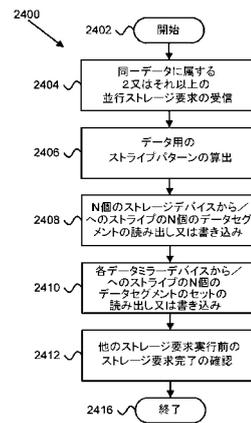
【 図 11 】



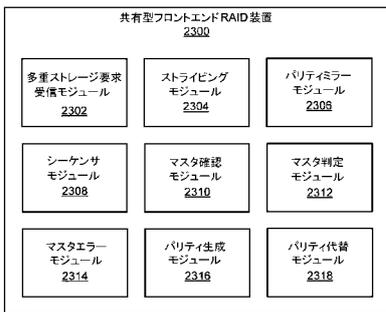
【 図 1 2 】



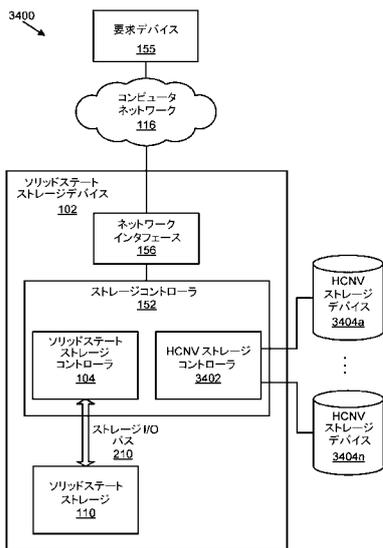
【 図 1 4 】



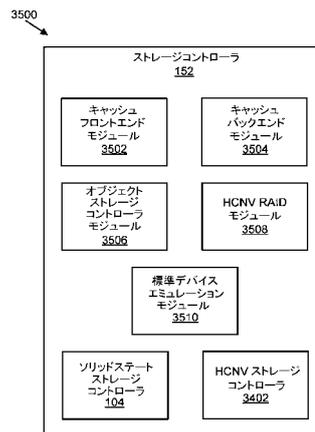
【 図 1 3 】



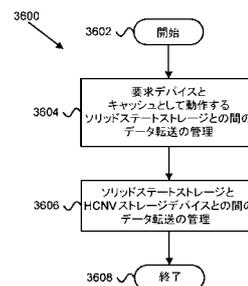
【 図 1 5 】



【 図 1 6 】



【 図 1 7 】



【手続補正書】

【提出日】平成20年10月9日(2008.10.9)

【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】

1又はそれ以上の高容量不揮発性(「HCNV」)ストレージデバイスでデータのストレージを管理するための装置であって、当該装置が、

ストレージ要求に関連するデータ転送を管理し、要求デバイスとソリッドステートストレージとの間の前記データ転送が、1又はそれ以上のHCNVストレージデバイス用のキャッシュとして機能し、前記データ転送がデータ、メタデータ及びメタデータインデックスのうち1又はそれ以上を含み、前記ソリッドステートストレージが不揮発性のソリッドステートデータストレージエレメントのレイを具えるキャッシュフロントエンドモジュールと；

前記ソリッドステートストレージと、前記1又はそれ以上のHCNVストレージデバイスとの間のデータ転送を管理するキャッシュバックエンドモジュールと；

を具え、ソリッドステートコントローラが、1又はそれ以上の要求デバイスからオブジェクト要求を提供し、前記ソリッドステートストレージ内の前記オブジェクト要求のオブジェクトを管理するオブジェクトストレージコントローラモジュールを更に具えることを特徴とする装置。

【請求項2】

請求項1に記載の装置において、前記キャッシュフロントエンドモジュール及び前記キャッシュバックエンドモジュールが、前記ソリッドステートストレージを管理するソリッドステートストレージコントローラと共に同一位置に配置されることを特徴とする装置。

【請求項3】

請求項1に記載の装置において、前記キャッシュフロントエンドモジュール、前記キャッシュバックエンドモジュール、及びソリッドステートストレージコントローラが、前記要求デバイスから自律して動作することを特徴とする装置。

【請求項4】

請求項1に記載の装置において、前記ソリッドステートコントローラが、現在のデータの直前に記憶されたデータに追加された、前記現在のデータを記憶するログストレージデバイスを具えることを特徴とする装置。

【請求項5】

請求項1に記載の装置が、RAIDレベルの一致する独立ドライブ冗長アレイ(「RAID」)の2又はそれ以上のHCNVストレージデバイス内の前記ソリッドステートストレージにキャッシュされたデータを記憶するHCNV RAIDモジュールを更に具え、前記データが全体として要求デバイスに見えることを特徴とする装置。

【請求項6】

請求項1に記載の装置において、前記ソリッドステートストレージ及び前記1又はそれ以上のHCNVストレージデバイスが、RAIDグループとして構成されるハイブリッドストレージデバイスセット内にハイブリッドストレージデバイスを具え、前記ソリッドステートストレージにキャッシュされ、HCNVデバイスに後に記憶されるデータセグメントが、ストライプのN個のデータセグメントのうち1つ、又は、前記ストライプのパーティータセグメントを含み、前記ハイブリッドストレージデバイスがRAIDストライプのデータセグメントと別個の、1又はそれ以上のクライアントからストレージ要求を受信することを特徴とする装置。

【請求項7】

請求項 6 に記載の装置において、前記ハイブリッドストレージデバイスが 2 又はそれ以上のクライアントから 2 又はそれ以上の同時のストレージ要求を受信する、共有されたフロントエンド分散型 R A I D グループのストレージデバイスであることを特徴とする装置。

【請求項 8】

請求項 1 に記載の装置において、前記 H C N V ストレージデバイスがハードディスクドライブ（「H D D」）、光学ドライブ、及びテープストレージのうちの 1 つであることを特徴とする装置。

【請求項 9】

請求項 1 に記載の装置において、前記ソリッドステートストレージ及び前記 1 又はそれ以上の H C N V ストレージデバイスが、ハイブリッドストレージデバイスを具え、当該ハイブリッドストレージデバイスの動作に特有のコードで、1 又はそれ以上の要求デバイスをロードする前に、前記 1 又はそれ以上の要求デバイスに取り付けられた標準デバイスをエミュレートすることによって、前記ハイブリッドストレージデバイスにアクセスを提供する標準デバイスエミュレーションモジュールを更に具え、前記標準デバイスが工業規格 B I O S によって担持されることを特徴とする装置。

【請求項 10】

請求項 1 に記載の装置において、前記ソリッドステートストレージデバイスが 2 又はそれ以上の領域に分割可能であり、1 又はそれ以上のパーティションが、前記 H C N V ストレージデバイス用のキャッシュとして機能する前記ソリッドステートストレージと別個の、ソリッドステートストレージとして用いられ得ることを特徴とする装置。

【請求項 11】

請求項 1 に記載の装置において、1 又はそれ以上のクライアントが、前記ソリッドステートストレージデバイス及び前記 1 又はそれ以上の H C N V ストレージデバイス内に記憶されるファイル又はオブジェクトのうちの 1 又はそれ以上の状態を管理するために、前記キャッシュフロントエンドモジュール及び前記キャッシュバックエンドモジュールにキャッシュ制御メッセージを送信することを特徴とする装置。

【請求項 12】

請求項 11 に記載の装置において、前記キャッシュ制御メッセージが、

前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージ内のオブジェクト又はファイルの一部を留める（p i n）制御メッセージ；

前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージ内のオブジェクト又はファイルの一部を解放する（u n p i n）制御メッセージ；

前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージから前記 1 又はそれ以上の H C V N ストレージデバイスまでのオブジェクト又はファイルの一部を消去させる制御メッセージ；

前記キャッシュバックエンドモジュールに、前記 1 又はそれ以上の H C V N ストレージデバイスから前記ソリッドステートストレージに、オブジェクト又はファイルの一部を事前ロードさせる制御メッセージ；及び

前記ソリッドステートストレージに、決められたストレージ空き領域量を空けるために、前記キャッシュバックエンドモジュールに、前記ソリッドステートストレージから前記 1 又はそれ以上の H C V N ストレージデバイスまでの 1 又はそれ以上のオブジェクト又はファイルの 1 又はそれ以上の部分をオフロードさせる制御メッセージ；

のうちの 1 又はそれ以上を含むことを特徴とする装置。

【請求項 13】

請求項 11 に記載の装置において、前記キャッシュ制御メッセージが前記オブジェクト又はファイル用のメタデータ（「キャッシュ制御メタデータ」）を通して通信されることを特徴とする装置。

【請求項 14】

請求項 13 に記載の装置において、前記キャッシュ制御メタデータが永続的であること

を特徴とする装置。

【請求項 15】

請求項 13 に記載の装置において、前記キャッシュ制御メタデータが前記ファイル又はオブジェクトの生成時に設定される属性を通して構築されることを特徴とする装置。

【請求項 16】

請求項 13 に記載の装置において、前記キャッシュ制御メタデータがファイル又はオブジェクト管理システムから得られることを特徴とする装置。

【請求項 17】

請求項 1 に記載の装置が、揮発性キャッシュストレージエレメントを更に具え、前記キャッシュフロントエンドモジュール及び前記キャッシュバックエンドモジュールが、前記揮発性キャッシュストレージエレメントにデータを記憶するステップを更に具え、前記ソリッドステートストレージ及び揮発性キャッシュストレージエレメントに記憶されるデータを管理し、前記バックエンドストレージモジュールが、前記揮発性キャッシュストレージエレメントと、前記ソリッドステートストレージと、前記 H C V N ストレージデバイスとの間のデータ転送を更に管理することを特徴とする装置。

【請求項 18】

請求項 17 に記載の装置において、前記 H C V N ストレージデバイスに記憶されるオブジェクト及びファイル用のメタデータ及びインデックスメタデータのうちの 1 又はそれ以上が、前記ソリッドステートストレージデバイス及び前記揮発性キャッシュストレージエレメント内に維持されることを特徴とする装置。

【請求項 19】

請求項 1 に記載の装置において、前記 H C V N ストレージデバイスに記憶されるオブジェクト及びファイル用のメタデータ及びインデックスメタデータのうちの 1 又はそれ以上が、前記ソリッドステートストレージデバイス内に維持されることを特徴とする装置。

【請求項 20】

請求項 1 に記載の装置において、前記ソリッドステートストレージ及び前記 1 又はそれ以上の H C N V ストレージデバイスが、ストレージデバイスを持って、前記 H C N V ストレージデバイスが前記ストレージデバイスに接続されたクライアントの表示から隠され、前記要求デバイスに単一のストレージデバイスとして見えるようにすることを特徴とする装置。

【請求項 21】

1 又はそれ以上の高容量不揮発性（「H C N V」）ストレージデバイスでデータのストレージを管理するためのシステムであって、当該システムが、

不揮発性のソリッドステートデータストレージエレメントのアレイを具えるソリッドステートストレージと；

1 又はそれ以上の H C N V ストレージデバイスと；

ソリッドステートストレージコントローラと；

H C N V ストレージデバイスコントローラと；

ストレージ要求に関連するデータ転送を管理し、要求デバイスと前記ソリッドステートストレージとの間の前記データ転送が、前記 1 又はそれ以上の H C N V ストレージデバイス用のキャッシュとして機能し、前記データ転送がデータ、メタデータ及びメタデータインデックスのうちの 1 又はそれ以上を含むキャッシュフロントエンドモジュールと；

前記ソリッドステートストレージと、前記 1 又はそれ以上の H C N V ストレージデバイスとの間のデータ転送を管理するキャッシュバックエンドモジュールと；

を具えるストレージコントローラと；

を具え、前記ソリッドステートコントローラが、1 又はそれ以上の要求デバイスからオブジェクト要求を提供し、前記ソリッドステートストレージ内の前記オブジェクト要求のオブジェクトを管理するオブジェクトストレージコントローラモジュールを更に具えることを特徴とするシステム。

【請求項 2 2】

請求項 2 1 に記載のシステムが、前記ストレージコントローラに接続されたネットワークインタフェースを更に具備し、当該ネットワークインタフェースがコンピュータネットワークを通して前記要求デバイスと前記ソリッドステートストレージコントローラとの間のデータ転送を促進することを特徴とするシステム。

【請求項 2 3】

請求項 2 1 に記載のシステムがサーバを更に具備し、当該サーバが前記ソリッドステートストレージと、前記 1 又はそれ以上の H C N V ストレージデバイスと、前記ストレージコントローラとを具備することを特徴とするシステム。

【請求項 2 4】

請求項 2 1 に記載のシステムにおいて、前記 1 又はそれ以上の H C N V ストレージデバイスがストレージエリアネットワーク（「S A N」）を通して前記ストレージコントローラに接続されることを特徴とするシステム。

【請求項 2 5】

コンピュータプログラムプロダクトであって、1 又はそれ以上の高容量不揮発性（「H C N V」）ストレージデバイスでデータのストレージを管理するための動作を行うのに実行可能な、コンピュータが利用可能なプログラムコードを有するコンピュータ可読媒体を具備し、前記コンピュータプログラムプロダクトの動作が、

キャッシュフロントエンドモジュールを用いて、ストレージ要求に関連するデータ転送を管理するステップであって、要求デバイスとソリッドステートストレージとの間の前記データ転送が、1 又はそれ以上の H C N V ストレージデバイス用のキャッシュとして機能し、前記データ転送がデータ、メタデータ及びメタデータインデックスのうちの 1 又はそれ以上を含み、前記ソリッドステートストレージが不揮発性のソリッドステートデータストレージエレメントのアレイを具備するステップと；

キャッシュバックエンドモジュールを用いて、前記ソリッドステートストレージと、前記 1 又はそれ以上の H C N V ストレージデバイスとの間のデータ転送を管理するステップと；

を含み、ソリッドステートコントローラが、1 又はそれ以上の要求デバイスからオブジェクト要求を提供し、前記ソリッドステートストレージ内の前記オブジェクト要求のオブジェクトを管理するオブジェクトストレージコントローラモジュールを更に具備することを特徴とするコンピュータプログラムプロダクト。

【請求項 2 6】

請求項 2 0 に記載の装置において、前記単一のストレージデバイスが少なくともブロックストレージデバイス、オブジェクトストレージデバイス、及びファイルストレージデバイスのうちのいずれかとして前記要求デバイスに見えるように構成可能であることを特徴とする装置。

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

		International application No PCT/US2007/025049
A. CLASSIFICATION OF SUBJECT MATTER INV. G06F12/08 G06F3/06		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06F		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practical, search terms used) EPO-Internal		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	WO 02/01365 A (INTEL CORP [US]; COULSON RICHARD [US]) 3 January 2002 (2002-01-03) abstract page 3, line 1 - page 4, line 23 page 7, line 3 - page 8, line 26 page 10, line 28 - page 11, line 2 page 14, line 13 - line 15 page 14, line 22 - line 27 page 15, line 12 - line 24 page 16, line 2 figures 1,7,8 ----- -/--	1-9, 11-18, 21-23,25
<input checked="" type="checkbox"/>	Further documents are listed in the continuation of Box C.	<input checked="" type="checkbox"/> See patent family annex.
* Special categories of cited documents :		
A document defining the general state of the art which is not considered to be of particular relevance *E* earlier document but published on or after the international filing date *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) *O* document referring to an oral disclosure, use, exhibition or other means *P* document published prior to the international filing date but later than the priority date claimed		*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art *&* document member of the same patent family
Date of the actual completion of the international search 17 April 2008		Date of mailing of the international search report 14/05/2008
Name and mailing address of the ISA/ European Patent Office, P.B. 5618 Patentlaan 2 NL - 2280 HW Rijswijk Tel. (+31-70) 340-2040, Tx. 31 651 epo nl, Fax: (+31-70) 340-3016		Authorized officer Knutsson, Frédéric

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2007/025049

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>US 2005/177687 A1 (RAO RAGHAVENDRA J [US]) 11 August 2005 (2005-08-11) abstract paragraph [0008] paragraph [0016] paragraph [0019] - paragraph [0025] figures 1,2</p>	<p>1-8, 11-25</p>
X	<p>"Windows PC Accelerators"[Online] 5 December 2006 (2006-12-05), pages 1-16, XP002476842 Retrieved from the Internet: URL:http://download.microsoft.com/download/9/c/5/9c5b2167-8017-4bae-9fde-d599bac8184a/perfacce1.doc> [retrieved on 2008-04-16] abstract page 3, line 1 - page 4, line 6 page 4, line 28 - page 5, line 12 page 5, line 20 - page 7, line 15 page 7, line 42 - line 46 page 8, line 38 - page 13, line 2 page 14, line 1 - line 9</p>	<p>1-4, 8-16,19, 21-23,25</p>

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2007/025049

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 0201365 A	03-01-2002	AU 7514701 A	08-01-2002
		CN 1527973 A	08-09-2004
		DE 10196380 T0	16-10-2003
		GB 2379538 A	12-03-2003
		JP 3951918 B2	01-08-2007
		JP 2004506256 T	26-02-2004
US 2005177687 A1	11-08-2005	EP 1723528 A1	22-11-2006
		WO 2005078588 A1	25-08-2005

フロントページの続き

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

(特許庁注：以下のものは登録商標)

1. Linux

(71)出願人 509157063

ザッペ, マイケル

アメリカ合衆国 コロラド州 80033, ホイートリッジ, シムズストリート 4615

(71)出願人 309039325

フュージョン マルチシステムズ, インク. (ディービエイ フュージョン - アイオー)

アメリカ合衆国 ユタ州 84121, ソルトレイクシティ, 6番フロア, サウス 3000 イースト 6350

(74)代理人 100096024

弁理士 柏原 三枝子

(74)代理人 100125520

弁理士 高橋 剛一

(74)代理人 100155310

弁理士 柴田 雅仁

(74)代理人 100156339

弁理士 米村 道子

(72)発明者 フリン, デイビッド

アメリカ合衆国 ユタ州 84093, サンディ, シェイディメドウッドライブ 8856

(72)発明者 シュトラッサー, ジョン

アメリカ合衆国 ユタ州 84075, シラキユース, サウス 2323

(72)発明者 サッチャー, ジョナサン

アメリカ合衆国 ユタ州 84043, リーハイ, ノース 2080 ウェスト 2259

(72)発明者 ザッペ, マイケル

アメリカ合衆国 コロラド州 80033, ホイートリッジ, シムズストリート 4615

Fターム(参考) 5B005 JJ11 KK02 KK15 MM12

5B065 BA05 CC08 CE12 CE26 CH01