

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4464279号
(P4464279)

(45) 発行日 平成22年5月19日 (2010.5.19)

(24) 登録日 平成22年2月26日 (2010.2.26)

(51) Int. Cl.	F I
G06F 12/00 (2006.01)	G06F 12/00 545Z
G06F 3/06 (2006.01)	G06F 12/00 514Z
	G06F 12/00 520Z
	G06F 12/00 531Z
	G06F 3/06 301J

請求項の数 10 (全 61 頁)

(21) 出願番号	特願2004-553816 (P2004-553816)	(73) 特許権者	504041295
(86) (22) 出願日	平成15年11月14日 (2003.11.14)		アイシロン・システムズ・インコーポレーテッド
(65) 公表番号	特表2006-506741 (P2006-506741A)		アメリカ合衆国・ワシントン・98119
(43) 公表日	平成18年2月23日 (2006.2.23)		・シアトル・ウエスト・マーサー・ストリート・220・#501
(86) 国際出願番号	PCT/US2003/036699	(74) 代理人	100064908
(87) 国際公開番号	W02004/046971		弁理士 志賀 正武
(87) 国際公開日	平成16年6月3日 (2004.6.3)	(74) 代理人	100089037
審査請求日	平成18年7月28日 (2006.7.28)		弁理士 渡邊 隆
(31) 優先権主張番号	60/426,464	(74) 代理人	100108453
(32) 優先日	平成14年11月14日 (2002.11.14)		弁理士 村山 靖彦
(33) 優先権主張国	米国 (US)	(74) 代理人	100110364
			弁理士 実広 信哉

最終頁に続く

(54) 【発明の名称】 分散ファイルシステムにおけるファイルの再ストライピングのためのシステム及び方法

(57) 【特許請求の範囲】

【請求項1】

相互に通信するように構成され、それぞれが、少なくとも一つの記憶装置と、少なくとも一つのプロセッサと、前記プロセッサによって実行される少なくとも一つの実行可能なソフトウェアモジュールとを有する複数の記憶ユニットを備え、

少なくともファイルの第1の部分が前記複数の記憶ユニットのうちの第1の記憶ユニットに格納されると共に少なくとも前記ファイルの第2の部分が前記複数の記憶ユニットのうちの第2の記憶ユニットに格納されるようにして、前記ファイルが当該分散ファイル記憶システムに格納され、

前記ファイルに関連するエラー訂正データが、前記ファイルを再構成するために利用可能なファイルデータと共に使用され得るようにして、前記複数の記憶ユニットのうちの1又は2以上の記憶ユニットに格納され、

当該分散ファイル記憶システムに格納された前記ファイルの各部分のアドレスロケーションと前記エラー訂正データのアドレスロケーションとを含むメタデータが、前記複数の記憶ユニットのうちの第3の記憶ユニットに格納され、

前記各記憶ユニットの少なくとも一つの実行可能なソフトウェアモジュールは、前記複数の記憶ユニットのうちの前記第3の記憶ユニット上の前記メタデータをアクセスし、

前記複数の記憶ユニットのうちの2以上の記憶ユニットの間で前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部を再ストライピングし、

10

20

前記複数の記憶ユニットのうちの前記第3の記憶ユニット上の前記メタデータを更新して、再ストライピングされた前記ファイルの各部についての新たなアドレスロケーションと、再ストライピングされた前記エラー訂正データのそれぞれについての新たなアドレスロケーションとを前記メタデータに含めるように構成され、

前記新たなアドレスロケーションは、別の記憶ユニット上のデータロケーションを含む分散ファイル記憶システム。

【請求項2】

前記再ストライピングは、前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部の再ストライピング中のデータの移動を最小化することを含む1又は2以上のプリフェレンスまたは制約に従う請求項1記載の分散ファイル記憶システム。

10

【請求項3】

前記少なくとも一つの実行可能なソフトウェアモジュールは、第1のプロテクションスキームから別の第2のプロテクションスキームへ前記ファイルのプロテクションスキームを変更するために、指示にตอบสนองして前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部を再ストライピングするように更に構成された請求項1記載の分散ファイル記憶システム。

【請求項4】

前記少なくとも一つの実行可能なソフトウェアモジュールは、当該分散ファイル記憶システムへの記憶ユニットの追加に対して、前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部を再ストライピングするように更に構成された請求項1記載の分散ファイル記憶システム。

20

【請求項5】

前記少なくとも一つの実行可能なソフトウェアモジュールは、当該分散ファイル記憶システムにおける記憶ユニットの故障に対して、前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部を再ストライピングするように更に構成された請求項1記載の分散ファイル記憶システム。

【請求項6】

前記少なくとも一つの実行可能なソフトウェアモジュールは、当該分散ファイル記憶システムからの記憶ユニットの除去に対して、前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部を再ストライピングするように更に構成された請求項1記載の分散ファイル記憶システム。

30

【請求項7】

前記少なくとも一つの実行可能なソフトウェアモジュールは、前記エラー訂正データの少なくとも一部または前記ファイルの少なくとも一部が無事に再ストライピングされたかどうかを判断するように更に構成された請求項1記載の分散ファイル記憶システム。

【請求項8】

アドレスロケーションは、前記ファイルの一部が格納された装置のロケーションと、前記ファイルの一部についてのメモリブロックアドレスとを含む請求項1記載の分散ファイル記憶システム。

【請求項9】

前記エラー訂正データは、パリティデータブロック及びミラーされたデータブロックのうちの一つを含む請求項1記載の分散ファイル記憶システム。

40

【請求項10】

前記複数の記憶ユニットのうち各記憶ユニットの前記少なくとも一つの実行可能なソフトウェアモジュールは、前記複数の記憶ユニットのうちの前記第1の記憶ユニットに格納された前記ファイルの少なくとも第1の部分と、前記複数の記憶ユニットのうちの前記第2の記憶ユニットに格納された前記ファイルの少なくとも第2の部分とに対するリード要求にตอบสนองするように更に構成された請求項1記載の分散ファイル記憶システム。

【発明の詳細な説明】

【技術分野】

50

【 0 0 0 1 】

本発明のシステム及び方法は、概して、分散ファイルストレージの分野に関し、特にインテリジェント分散ファイル管理システムに格納されたファイルの再ストライピングに関する。

【 0 0 0 2 】

本出願は、2002年11月14日に出願された米国仮出願60/426464の利益を主張するものであり、その出願は、参照することによりそっくりそのまま本明細書に組み込まれる。

【背景技術】

【 0 0 0 3 】

インターネットの爆発的な成長は、情報が常時交換され且つアクセスされる新規な領域において先導役を務めてきた。この成長に対応して、共有(share)されているデータのサイズが増加している。ユーザは、標準HTMLドキュメント以上のものを要求しており、オーディオデータ、ビデオデータ、イメージデータおよびプログラミングデータのような種々のデータにアクセスすることを望んでいる。従って、大きなデータのセットに対応できると共に、そのデータに対する高速かつ信頼性のあるアクセスを同時に提供するデータストレージに対する要請がある。

【 0 0 0 4 】

一つの回答は、大量のデータを格納できるが高スループットレートの提供が困難な単一の記憶装置(storage device)を利用することであった。データ容量が増加するにつれ、データをアクセスするのに要する時間が同様に増加する。処理速度とパワーについては改善されたが、ディスクI/O(入力/出力)オペレーション性能は、特に大きなデータファイルについては、I/Oオペレーションを非効率にするのと同じ進度では改善されていない。

【 0 0 0 5 】

他の回答は、ストレージエリアネットワークソリューション(SAN; Storage Area Network solutions)のようなアーキテクチャを用いて複数のサーバが共有ディスクにアクセスすることを可能にしたが、このようなシステムは高価であり、しかもデータの完全性(integrity)を制御し且つセットアップするための複雑な技術を必要とする。更に、大量のデータ要求を取り扱うためには高速アダプタを必要とする。

【 0 0 0 6 】

従来のアプローチに関する一つの問題は、それらが拡張性(scalability)において制約されていることである。従って、データ量が増加するにつれてシステムが大きくなる必要があるが、拡大にはコストを要し、極めて秩序を乱す(disruptive)。

【 0 0 0 7 】

従来のアプローチに関する他の共通問題は、それらが、柔軟性(flexibility)において制約されていることである。システムは、しばしば、所定のエラー訂正制御を使用するように構成される。例えば、RAIDシステムは、データが格納されるべき場所または使用されるべき冗長パラメータのタイプを決定することにおいて管理者に柔軟性をほとんど与えないか全く与えない物理ディスクレベルで、データファイルのミラーリング(mirroring)及び冗長性(redundancy)を提供するために使用される。

【発明の開示】

【課題を解決するための手段】

【 0 0 0 8 】

インテリジェント分散システムは、単一ファイルシステムとしてアクセスされるスマート記憶ユニット(smart storage unit)のセット間でのファイルデータの格納を有利に可能にする。このインテリジェント分散ファイルシステムは、メタデータ構造を利用して各ファイルについての詳細な情報を管理(manage)すると共に追跡(track)し、例えば、ファイルのデータブロックのブロックロケーション及び装置を備え、単一ファイルシステム内の冗長性(redundancy)及び/又は複製(replication)の異なるレベルを可能とし、冗長パラ

10

20

30

40

50

メータの変更を容易化し、メタデータのための高レベルのプロテクションを提供し、実時間でデータを移動すると共に複製する。等々。

【 0 0 0 9 】

また、インテリジェント分散ファイルシステムは、スマート記憶ユニットのセット間に分散されたファイルを再ストライピングするためのシステム及び方法を有利に備え、ここで、特定のストライプについてのデータブロックは一般に記憶ユニット上の何れかの特定のロケーションに位置される必要がなく、データは最小データ移動でスマート記憶ユニット間で再ストライピングされ、データは一般にシステム故障が再ストライピング処理中に発生しても回復可能であり且つ保護される。

【 0 0 1 0 】

本発明の一態様は、複数のインテリジェント記憶装置と通信する分散ファイルシステムに関し、この分散ファイルシステムは、メッセージングシステム(messaging system)を備え、且つメタデータが、メタデータのデータブロックのロケーション、コンテンツデータブロック、およびパリティデータブロックを含むように、分散ファイルシステム上に格納されたディレクトリおよびファイルについてのメタデータを格納し且つ管理するように構成され、且つ、この分散ファイルシステムは、更にデータファイルを再ストライピングするように構成される。

【 0 0 1 1 】

本発明の他の態様は、インテリジェント分散ファイルシステムにおける複数のスマート記憶装置に格納されたデータの再ストライピングの方法に関する。この方法は、既存のファイルからのデータのクラスタ(cluster)が格納される少なくとも一つのスマート記憶装置を識別(identify)するステップと、前記クラスタを前記識別されたスマート記憶装置上に格納するステップと、前記クラスタが無事に格納されたことを判断するステップと、前記データのクラスタに関連するメタデータを更新して新たな割り当て(assignment)を反映させるステップとを備える。

【 0 0 1 2 】

本発明の追加の態様は、データを再ストライピングするためのシステムに関する。このシステムは、複数のスマート装置(smart devices)であって該スマート装置間に分散されたストライプにデータブロックを格納するように構成されたスマート装置と、データブロックをスマート装置に割り当てるように構成された割り当てモジュール(assignment module)と、前記データブロックが割り当てられた後に前記データブロックを格納するためのスマート装置に命令を送信するように構成された要求モジュール(request module)と、前記データブロックが格納された後に前記データブロックに関連するメタデータを更新するように構成された更新モジュール(update module)とを備える。

【 0 0 1 3 】

本発明の他の態様は、データブロックからなるデータを再ストライピングする方法に関する。この方法は、データブロックを記憶ユニットに割り当てるステップと、前記データブロックが既に前記記憶ユニットに格納されているかどうかを判断するステップと、それが前記記憶ユニットにまだ存在していなければ、前記データブロックを前記記憶ユニットに格納するステップと、前記データブロックが格納されていれば、前記データブロックに関連するメタデータを更新するステップとを備える。

【 0 0 1 4 】

本発明の追加の態様は、記憶ユニットにデータブロックを割り当てる方法に関する。この方法は、利用可能な記憶ユニットを識別するステップと、各データブロックのための利用可能な記憶ユニットを選択するステップと、選択された各記憶ユニットについて関連アドレスロケーション(related address locations)を判断するステップとを備える。

【 0 0 1 5 】

本発明の他の態様は、プロテクショングループ(protection group)におけるデータのブロックを格納するための記憶ユニットのセットから記憶ユニットを選択する方法に関する。この方法は、前記記憶ユニットのセットからの何れの記憶ユニットが、プロテクション

10

20

30

40

50

グループに関して最適な記憶ユニットであるかを識別するステップと、前記最適な記憶ユニットを選択するステップとを備える。

【0016】

本発明の追加の態様は、データのブロックのための記憶ユニットを選択する方法に関する。この方法は、前記ブロックデータが現在格納されている記憶ユニットを識別するステップと、第1記憶ユニット上の最適なロケーションを識別するステップと、最適な記憶ユニットのセットを判断するステップと、前記第1記憶ユニットが、前記データのブロックが現在格納されている前記記憶ユニットのうちの一つであれば、前記第1記憶ユニットを選択するステップと、前記第1の記憶ユニットが、前記データのブロックが現在格納されている前記記憶ユニットのうちの一つでなければ、前記最適な記憶ユニットのセットのうちの一つを選択するステップとを備える。

10

【0017】

本発明の他の態様は、命令を含む記憶装置を備えたシステムに関し、この命令は、実行された場合に、既存ファイルからのデータのクラスタが格納される記憶装置を識別するステップと、前記クラスタを前記識別された記憶装置に格納するステップと、前記クラスタが無事に格納されたことを判断するステップと、前記データのクラスタに関連するメタデータを更新するステップの方法を前記システムに実行させる。

【0018】

要約すると、本発明の或る態様、利点、及び新規な特徴が本明細書で述べられている。必ずしも、全てのこのような利点は本発明の何れかの特定の実施形態に従って達成されるものではないことが理解されるべきである。従って、例えば、当業者であれば、本明細書において示唆または教示されるような他の利点を必ずしも達成することなく、本明細書で教示されるような一つの利点または利点のグループを達成する方法で、本発明が実行または具現されることが分かる。

20

【0019】

図1は、本発明の一実施形態のハイレベルブロック図である。

【0020】

図2は、図1に図解される構成要素間でのデータの流れの一例を示す。

【0021】

図3は、スマート記憶ユニットの一例のハイレベルブロック図である。

30

【0022】

図4は、ファイルディレクトリの一例を図解する。

【0023】

図5は、メタデータ構造の維持を図解する。

【0024】

図6Aは、データロケーションテーブル構造の一実施形態を図解する。

【0025】

図6Bは、データロケーションテーブル構造の追加の実施形態を図解する。

【0026】

図6Cは、データロケーションテーブル構造の追加の実施形態を図解する。

40

【0027】

図6Dは、データロケーションテーブル構造の追加の実施形態を図解する。

【0028】

図7Aは、ディレクトリのためのメタデータ構造の一実施形態を図解する。

【0029】

図7Bは、ファイルのためのメタデータ構造の一実施形態を図解する。

【0030】

図8Aは、データロケーションテーブルの一実施形態を図解する。

【0031】

図8Bは、データロケーションテーブルの追加の実施形態を図解する。

50

【 0 0 3 2 】

図 8 C は、データロケーションテーブルの追加の実施形態を図解する。

【 0 0 3 3 】

図 9 は、対応するサンプルデータを有するファイルのメタデータの一例を図解する。

【 0 0 3 4 】

図 1 0 は、ファイルを読み出すためのフローチャートの一実施形態を図解する。

【 0 0 3 5 】

図 1 1 は、名前解決(name resolution)を実施するためのフローチャートの一実施形態を図解する。

【 0 0 3 6 】

図 1 2 は、ファイルを読み出すためのフローチャートの一実施形態を図解する。

【 0 0 3 7 】

図 1 3 は、パリティ情報を生成するためのフローチャートの一実施形態を図解する。

【 0 0 3 8 】

図 1 4 は、エラー訂正を実施するためのフローチャートの一実施形態を図解する。

【 0 0 3 9 】

図 1 5 は、インテリジェント分散ファイルシステムにおいてデータを再ストライピングするためのフローチャートの一実施形態を図解する。

【 0 0 4 0 】

図 1 6 は、スマート記憶ユニットにデータを割り当てるためのフローチャートの一実施形態を図解する。

【 0 0 4 1 】

図 1 7 は、利用可能な記憶ユニットのセットの中から選択するためのフローチャートの一実施形態を図解する。

【 発明を実施するための最良の形態 】

【 0 0 4 2 】

本発明の種々の実施形態および適用例を表すシステムおよび方法が図面を参照して説明されるであろう。また他の実施形態を表すシステム及び方法に対する変形も述べられるであろう。

【 0 0 4 3 】

説明のため、いくつかの実施形態が、コンテンツ配信およびウェブホスティングとの関連において述べられるであろう。本発明者は、本システムおよび方法が使用される環境の形式に本発明が制限されず、本システムおよび方法が、例えばインターネット、ワールドワイドウェブ、病院の私設ネットワーク、政府機関のブロードキャストネットワーク、法人事業の内部ネットワーク、イントラネット、ローカルエリアネットワーク、ワイドエリアネットワーク、等々のような他の環境において使用してもよいことを意図している。しかしながら、図面および説明は、環境がインターネットコンテンツ配信およびウェブホスティングのものであるところの本発明の実施形態に関する。他の実施形態において、本システムおよび方法は、単一モジュールとして実施されてもよく、及び/又は、他のモジュールやその同等物と併せて実施されてもよい。さらに、本明細書で述べられる特定の実施は説明のためのものであって本発明を限定するものではない。本発明の範囲は添付の特許請求の範囲によって定められる。

【 0 0 4 4 】

これら及び他の特徴は上述の図面を参照して説明される。図面及び関連する説明は、本発明の実施形態を説明するために提供されるものであって、本発明の範囲を限定するものではない。図面の全体を通して、参照番号は、参照される要素間の対応を示すために再使用される。加えて、各参照番号の最初の桁は、概して要素が最初に現れる図面を示す。

【 0 0 4 5 】

I. 概要

本発明のシステム及び方法は、インテリジェント分散ファイルシステムを提供し、この

10

20

30

40

50

インテリジェント分散ファイルシステムは、単一のファイルシステムとしてアクセスされるスマート記憶ユニットのセット間でのデータの格納を可能にする。本インテリジェント分散ファイルシステムは、各ファイルについての詳細なメタデータを追跡し且つ管理する。メタデータは、例えば、装置およびブロックロケーション情報の双方を含むファイルのデータブロックのロケーション、メタデータ及び/又はデータブロック(もしあれば)の冗長なコピーのロケーション、エラー訂正情報、アクセス情報、ファイルのネーム、ファイルのサイズ、ファイルのタイプ、ファイルのデータ及びプロテクション情報が格納されるスマート記憶ユニット、等々のようなファイルに関し及び/又は記述する任意のデータを含んでもよい。加えて、インテリジェント分散ファイルシステムは、本ファイルシステムによって管理されるところの異なるファイル及び/又はデータブロックについて、異なるレベルの複製及び/又は冗長を可能にし、システムがアクティブ中に冗長パラメータの変更を容易化し、システムがアクティブ中に欠損データの回復を容易化し、そしてメタデータ及びデータの実時間の複製及び移動を可能にする。更に、各スマート記憶ユニットは、スマート記憶ユニットからのファイルのデータを配置(locating)及び収集(collecting)することによりファイル要求に応答してもよい。

10

【0046】

インテリジェント分散ファイルシステムは、有利には、特にライト(WRITE)要求の数に比例する多くのリード(READ)要求が存在する状況におけるデータへのアクセスを提供する。これは、スマート記憶ユニットのグループをロックング(locking)すること、またはライト要求に対する整合性を確保するためにスマート記憶ユニット上でジャーナル処理(journaling)を行うことの追加的複雑性のためである。さらに、インテリジェント分散ファイルシステムは、有利には、大きなデータのブロックに対する要求が一般的であるブロックトランザクションを取り扱う。

20

【0047】

いくつかの実施形態のひとつの利点は、ファイル及びディレクトリのためのメタデータがインテリジェント分散ファイルシステムによって管理され且つアクセスされることである。このメタデータは、ディレクトリまたはファイルのためのメタデータが配置された場所、コンテンツデータが格納された場所、同様にシステムに関連する他のエラー訂正情報またはパリティが格納された場所を表してもよい。データロケーション情報は、例えば装置およびブロックロケーション情報を用いて格納されてもよい。従って、このインテリジェント分散ファイルシステムは、その両方がスマート記憶ユニット間に分散されて格納されるメタデータを用いて、要求されたコンテンツデータを配置し且つ読み出してもよい。加えて、インテリジェント分散ファイルシステムはメタデータへのアクセスを有するので、インテリジェント分散ファイルシステムは、データが格納されるべき場所を選択するために使用されてもよく、そしてスマート記憶ユニットのセットを混乱させることなくデータを移動(move)し、複製(replicate)し、回復(recover)し、及び/又は変更(change)するために使用されてもよい。

30

【0048】

いくつかの実施形態の他の利点は、各ファイルについてのデータは、多くのスマート記憶ユニットにわたって格納されてもよく、且つタイミングを得た方法でアクセスされてもよい。各ファイルについてのデータブロックは、データアクセスタイムが減少されるようなスマート記憶ユニットのサブセット間で分散されてもよい。さらに、異なるファイルは、スマート記憶ユニットの異なるセットにわたると同様に、異なる数のスマート記憶ユニットにわたって格納されてもよい。このアーキテクチャは、インテリジェント分散ファイルシステムが、利用可能な記憶容量、CPU稼働率、および各スマート記憶ユニットのネットワーク稼働率と同様に、ファイルのサイズ、重要性、予期されたアクセスレートに基づいてデータブロックを知的(intelligently)に格納することを可能にする。

40

【0049】

いくつかの実施形態の追加的利点は、本システムおよび方法が、スマート記憶ユニット間に格納される異なるブロック又はファイルが異なるタイプのプロテクションを備えても

50

よいようなブロック又はファイルベースで、エラー訂正、冗長、およびミラーリングのような種々のプロテクションスキームを提供するために使用されてもよいということである。例えば、いくつかのディレクトリまたはファイルはミラーされてもよく、他は、種々のエラー又は欠損訂正スキームを用いてエラー及び/又は欠損訂正データで保護されてもよく、重要度の低い他のものは何らプロテクションスキームを使用しなくてもよい。

【0050】

いくつかの実施形態の更なる利点は、本システムおよび方法が、進行中のデータ要求を混乱または中断することなく、スマート記憶ユニットの実時間の付加(addition)、削除(delete)、及び/又は修正(modification)を可能にするということである。従って、より多くの記憶が要求されるほど、追加のスマート記憶ユニットがスマート記憶ユニットのセットに加えられ、そしてファイル要求を中断することなく、または既存のスマート記憶ユニットをオフラインにする必要なく、実時間でインテリジェント分散ファイルシステムに組み込まれてもよい。既存のファイル又は新たなファイルのデータブロックが、今や新たなスマート記憶ユニットを備えるスマート記憶ユニットにわたってインテリジェント分散ファイルシステムにより分散されているので、既存のスマート記憶ユニットはファイルの要求を処理してもよい。

10

【0051】

いくつかの実施形態の他の利点は、本システムおよび方法は、1または2以上のスマート記憶ユニット上のそれらのブロックを複製し従って任意の個々のデータブロックのためのアクセスの複数のポイントを生成することにより、データブロックの記憶に対する実時間の修正を実施してもよい。この複製は、CPU稼働率、および、頻繁なアクセスパターンが観られるファイル又はファイルのグループのための個々のスマート記憶ユニットに対するネットワーク資源必要条件を低減させるのに役立つ。これらのアクセスパターンは、スマート記憶ユニットによりモニタされ、そしてインテリジェント分散ファイルシステムは、インテリジェント分散ファイルシステムの稼動中にこのようなデータ複製を作成するための柔軟性をスマート記憶ユニットに与える。

20

【0052】

いくつかの実施形態の更なる利点は、本システムおよび方法が、1または2以上のスマート記憶ユニット上でファイルの再ストライピングを通じてデータファイルの実時間の再配置(relocation)を可能にする。いくつかの実施形態において、この再ストライピングは、有利には、システム故障が再ストライピング処理中に発生してもデータの完全性を保ち且つ回復可能性を確保する最小限のデータブロックの移動で実施される。この再ストライピングシステムは、有利には、特定のデータブロックがスマート記憶ユニット上の何れの特定のロケーションにも配置されることを必要とせず、格納されたデータブロックに対するプロテクションを提供することを継続し、たとえ1または2以上のスマート記憶ユニットが故障してもデータの回復可能性を可能にする。

30

【0053】

II. 動作例

説明のため、インテリジェント分散ファイルシステムが動作において使用される場所のシナリオ例が述べられるであろう。このシナリオ例においては、インテリジェント分散ファイルシステムは、インターネットウェブサイトを紹介した映画のダウンロードを提供する会社によって使用される。この会社は、ウェブサイトを紹介して顧客によってアクセスされる映画の予告編、広告、および顧客情報のみならず、ダウンロード可能な映画のコピーを格納し且つ管理するためにインテリジェント分散ファイルシステムを使用してもよい。このデータは、種々のレベルのプロテクションを用いて格納されてもよく、且つ高速アクセスのために複数のスマート記憶ユニットにわたって格納されてもよい。

40

【0054】

例えば、電子メールに高速にアクセスするために、会社は、インテリジェント分散ファイルシステムにおける多数のスマート記憶ユニットにわたって顧客調査(customer survey)電子メールを格納することを望んでもよい。しかしながら、会社は、全ての電子メール

50

のバックアップテープを保持してもよく、そして顧客調査の迅速な回復を可能にすることは重要ではないように感じるかもしれない。会社は、インテリジェント分散ファイルシステムに、顧客調査電子メールにエラー訂正またはミラーリングプロテクションを使用しないように指示してもよい。従って、もし1または2以上のスマート記憶ユニットがアクセス不能になれば、会社は、その電子メールがバックアップテープから回復できるまでは、それらのスマート記憶ユニット上の顧客調査電子メールが遅延されることを受け入れることができると感じるかもしれない。

【0055】

広告については、会社は、もし1または2以上のスマート記憶ユニットが故障すれば、インテリジェント分散ファイルシステムが、広告の表示を中断することなくデータを回復できるような高度のエラー訂正パラメータを使用するように、インテリジェント分散ファイルシステムに指示してもよい。例えば、会社は、種々の障害許容性測定(fault tolerance measurement)を信頼して、どの程度のプロテクションが特定のファイルに与えられるべきかを判断することを支援してもよい。重要な情報については、会社はXの障害許容性レベルを確保し、且つさほど重要でない情報については、会社はYの障害許容性を確保することを望んでもよい。ここで $X > Y$ である。従って、会社は、その広告主に、たとえ1または2以上のスマート記憶ユニットが故障したとしても、広告が信頼できるベースで利用可能であることを保証してもよい。障害許容性に加えて又は障害許容性に代えて他の測定が使用されてもよいことが理解され、そして障害許容性は信頼性のひとつの測定を説明するために使用される。

【0056】

人気のある映画のダウンロードについて、会社は、有利には、インテリジェント分散ファイルシステムをセットアップして自動的に複数の映画データのコピーを格納し、より多くのユーザがそのデータにアクセスすることを可能とし、もし1または2以上のスマート記憶ユニットが故障すれば、欠損したデータが再生成または他のロケーションから読み出される。更にまた、もし要求の数が増加し、またはもし1または2以上のスマート記憶ユニットがスマート記憶ユニット上に存在するデータの要求で氾濫し始めていれば、人気のある映画ダウンロードの追加コピーは、生成され且つスマート記憶ユニット間に格納されてもよい。

【0057】

会社は、さほど人気のない他の映画を提供してもよく、そして少ない需要のために少ないコピーを格納するようにインテリジェント分散ファイルシステムに対し指示してもよい。更に、“人気のあるダウンロード映画”の人気がなくなるにつれ、会社は、有利には、映画が格納されたスマート記憶ユニットから映画の余分なコピーを削除し、そして“人気の少ない”映画を低速性能を有するスマート記憶ユニット(例えば、少ない利用可能ディスクスペースを有するそれらのスマート記憶ユニット)に移動するように、インテリジェント分散ファイルシステムをセットアップしてもよい。インテリジェント分散ファイルシステムは、自動的に、スマート記憶ユニットを用いてこれらのタスクを処理してもよい。

【0058】

加えて、会社がより多くの映画を獲得するにつれ、会社は、追加のスマート記憶ユニットをインテリジェント分散ファイルシステムに加えてもよい。そして、会社は、新たなスマート記憶ユニットを使用して、より多くの映画を格納し、より多くの既存の映画のコピーを格納し、応答時間を改善するために既存の映画データを再分散し、及び/又は1または2以上のファイルのプロテクションスキームを変更してもよい。たとえインテリジェント分散ファイルシステムが複数の記憶ユニット間にデータを格納し且つ管理するとしても、インテリジェント分散ファイルシステムが単一のファイルシステムとして見えるように、追加のスマート記憶ユニットがインテリジェント分散ファイルシステムに組み込まれても良い。

【0059】

もし、スマート記憶ユニットのひとつが故障すれば、インテリジェント分散ファイルシ

10

20

30

40

50

システムは、故障したユニット上に格納された任意のファイルデータを再構築し、その情報を稼動中のユニットに格納し、そして適切なメタデータのデータ構造を更新してもよい。従って、もしユーザが、故障したユニットに部分的に格納されたファイルを要求すれば、そのユーザは、それでもユニットのひとつがダウンしていることを知らずにそのファイルのデータを受信してもよい。

【 0 0 6 0 】

この例では、インテリジェント分散ファイルシステムは、首位映画ダウンロードに対する確実且つ高速なアクセス、人気の少ない映画に対する高速アクセス、および顧客調査電子メールに対するアクセスを提供するための能力を会社に提供する。各ファイルについて、会社は、エラー及び/又は欠損訂正パラメータを設定してもよく、そしてどれだけ多くのファイルの追加的コピーが格納されるべきかを選択してもよい。いくつかの状況において、会社は、手動で、どれだけ多くのデータのコピーが格納されるべきかを選択し、そのデータを格納するための場所を決定してもよい。他の状況においては、会社は、どれだけ多くのデータのコピーが格納されるべきか、(もしあれば)使用されるべきエラー及び/又は欠損訂正スキームを選択し、及び/又はデータが格納されるべき場所を選択するために、本インテリジェント分散ファイルシステムの特徴を当てにしてもよい。従って、会社はその記憶スペースを効率的に使用してユーザの要求に一層良好に対処できる。記憶スペースは、まばらに要求されたファイルについては浪費されず、そして重要でないファイルについては、エラー訂正情報は生成されず且つ格納されない。

【 0 0 6 1 】

上述の例は、ダウンロード用の映画を提供する会社を含むが、この例は、インテリジェント分散ファイルシステムの一実施形態の特徴を例示するためだけに使用される。さらに、このインテリジェント分散ファイルシステムは、他の環境で使用されもよく、且つ、他とパ、サウンドファイル、オーディオファイル、グラフィックファイル、マルチメディアファイル、デジタル写真、実行ファイル等々を含む他のタイプのデータ及び/又はデータの組み合わせと共に使用してもよい。

【 0 0 6 2 】

III. インテリジェント分散ファイルシステム

図1は、ネットワークサーバ120と通信してリモートファイルアクセスを提供するインテリジェント分散ファイルシステムの一実施形態を示す。このインテリジェント分散ファイルシステム110は、例えばNFSまたはCIFSのような種々のプロトコルを用いてネットワークサーバ120と通信してもよい。ユーザ130は、インターネット145のような通信媒体140を介してネットワークサーバ120と情報のやりとりを行って、インテリジェント分散ファイルシステム110によって管理されるファイルを要求する。代表的なインテリジェント分散ファイルシステム110は、スマート記憶ユニット114のセットおよびネットワークサーバ120と通信するスイッチコンポーネント125を使用する。このインテリジェント分散ファイルシステム110は、個々のファイルのブロックが複数のスマート記憶ユニット114にわたって分散されることを可能にする。このデータへのアクセスが、このデータが単一の装置上に格納された場合よりも高いスループットを提供するように、このデータが格納される。加えて、インテリジェント分散ファイルシステム110は、種々のプロテクションスキームを用いて格納された種々のデータファイルを格納するために使用されてもよい。

【 0 0 6 3 】

代表的なインテリジェント分散ファイルシステム110は、スマート記憶ユニット114のセット間にデータを格納する。このスマート記憶ユニット114についての詳細な説明に関しては、以下の“スマート記憶ユニット”と題された節を参照されたい。

【 0 0 6 4 】

代表的なインテリジェント分散ファイルシステムは、要求されているデータのタイプを処理することができるアプリケーションサーバに要求を方向付ける負荷バランシングスイッチのようなスイッチコンポーネント125を使用する。入来する要求は、遅延を最小化

し且つデータの完全性を確保するために、高速技術を用いて適切なアプリケーションサーバに転送される。

【 0 0 6 5 】

標準イーサネット（登録商標）スイッチまたは他の負荷バランシングスイッチと同様に、例えば 1 0 0 0 Base-T（銅）ギガビット負荷バランシングイーサネット（登録商標）スイッチ、エクストリームネットワークサミット(Extreme Networks Summit) 7 I、ファウンドリファストアイアン(Foundry Fast Iron) I I、ノーテルネットワークアルテオン A C E スイッチ(Nortel Networks Alteon ACEswitch) 1 8 0、F 5 B i g - I p のような種々の負荷バランシングスイッチ 1 2 5 が使用されてもよいことが理解される。インテリジェント分散ファイルシステムは、例えば“ジャンボ”イーサネット（登録商標）フレームのような大きなフレームサイズをサポートするスイッチを使用する。加えて、負荷バランシングスイッチ 1 2 5 は、他の市販製品及び/又は有標製品のみならず、ファウンドリネットワークスのサーバアイアン(ServerIron)スイッチ、アサンテ(Asante)のインストラスイッチ(InstraSwith) 6 2 0 0 スイッチ、アサンテのホットスタック(HotStack)、シスコ(Cisco)のカタリスト(Catalyst)スイッチを用いて実施されてもよい。しかしながら、当業者であれば、多様なスイッチコンポーネント 1 2 5 が使用され、または他の技術が使用されてもよいことが分かる。さらに、種々のネットワークフレームサイズを伝送するようにスイッチコンポーネント 1 2 5 を構成してもよいことが分かる。

10

【 0 0 6 6 】

ディスク、マザーボード、CPU、オペレーティングシステム、または 1 または 2 以上のスマート記憶ユニットへのアクセスを妨害する他のハードウェアまたはソフトウェア故障の場合には、重要性の高いファイルは、データに高い回復率を提供する高いエラー訂正パラメータを用いて格納されてもよい。もし、何れかのデータが失われ又は損なわれていれば、スマート記憶ユニット 1 1 4 は、他のロケーションからデータを取得し又はデータを再生成するために、メタデータにおける冗長情報またはミラーリング情報を使用してもよい。高い需要があるファイルは、なお一層高いスループットレートを提供するために、追加のスマート記憶ユニット 1 1 4 にわたって実時間でミラーされてもよい。

20

【 0 0 6 7 】

インテリジェント分散ファイルシステム 1 1 0 の一実施形態において、メタデータ構造は、少なくとも、メタデータ構造に対応するディレクトリに由来するもの(descendant)を含んで参照するデータと同じプロテクションを有する。メタデータ構造なしでデータを読み出すことは困難であるので、メタデータ構造におけるデータの欠損(loss)は、インテリジェント分散ファイルシステム 1 1 0 を害する。インテリジェント分散ファイルシステム 1 1 0 では、メタデータ構造の代わりにコピーは、要求されたプロテクションを提供するために、必要に応じた多くのロケーションにおいてミラーされてもよい。従って、パリティプロテクションを有するファイルは、少なくとも同一又はより大きなパリティプロテクションで格納されたそのメタデータ構造を備えてもよく、2 回ミラーされたファイルは、少なくとも 2 つのロケーションでミラーされたそのメタデータ構造を備えてもよい。

30

【 0 0 6 8 】

図 1 はインテリジェント分散ファイルシステム 1 1 0 の一実施形態を示すが、他の実施形態が用いられてもよいことが理解される。例えば、アプリケーションサーバのような追加のサーバがスイッチコンポーネント 1 2 5 と通信してもよい。これらアプリケーションサーバは、例えば、オーディオストリーミングサーバ、ビデオストリーミングサーバ、画像処理サーバ、データベースサーバなどを含んでも良い。さらに、スイッチコンポーネント 1 2 5 と通信するワークステーションのような追加の装置が存在していてもよい。加えて、図 1 は、4 つのスマート記憶ユニット 1 1 4 と共に動作するインテリジェント分散ファイルシステム 1 1 0 を示すが、このインテリジェント分散ファイルシステム 1 1 0 は、異なる数のスマート記憶ユニット 1 1 4 と共に動作してもよいことが理解される。

40

【 0 0 6 9 】

用語“リモート”は、装置、コンポーネント、及び/又はローカルに格納されないモジ

50

ジュールを含んでも良く、それは、ローカルバスを介してアクセスできない。従って、リモート装置は、物理的に同一の部屋に設置され且つスイッチ又はローカルエリアネットワークのような装置を介して接続された装置を含んでも良い。他の状況において、リモート装置は、また、例えば異なったロケーション、土地(country)などのような離れた地理上の区域に配置されてもよい。

【 0 0 7 0 】

また、インテリジェント分散ファイルシステム 1 1 0 を用いて種々のタイプのデータが格納されてもよいことも理解される。例えば、インテリジェント分散ファイルシステム 1 1 0 は、例えば、ビデオオンデマンド(video-on-demand)、オンラインミュージックシステム、ウェブサイトミラーリング、大規模データベース、大規模グラフィックファイル、CAD/CAM設計、ソフトウェア更新、会社のプレゼンテーション、保険金請求ファイル、医療画像ファイル、会社の文書記憶などのような大きなファイルのアプリケーションで使用されてもよい。

【 0 0 7 1 】

図 2 は、ウェブサイトユーザ 1 3 0 がオンデマンドデジタルビデオを観るための要求を送信した環境例を示す。イベント A において、ユーザ 1 3 0 は、インターネット 1 4 5 を介してウェブサイトに要求を送信し、映画movie.movieのコピーを観ることを要求する。この要求は、ウェブサイトサーバ 1 2 0 によって受信され、そしてサーバ 1 2 0 は、このファイルが、\movies\comedy\mymovie.movieに配置されていることを判断する。イベント B において、インテリジェント分散ファイルシステムのスイッチコンポーネント 1 2 5 は、この要求を調べてインテリジェント分散ファイルシステム 1 1 0 に接続し、そして標準負荷バランシング技術を用いて、この要求を、スマート記憶ユニット 0 のような利用可能なスマート記憶ユニット 1 1 4 に転送する。イベント C において、スマート記憶ユニット 0 は、ファイル/DFSR/movies/comedy/mymovie.movieに対する要求を受信し、そしてそのルートメタデータのデータ構造(ルートディレクトリ/DFSRのための)から、サブディレクトリmoviesのメタデータのデータ構造がスマート記憶ユニット 2 で格納されていることを判断する。イベント D において、スマート記憶ユニット 0 は、サブディレクトリcomedyのためのメタデータ構造のロケーションを要求するスマート記憶ユニット 2 に要求を送信する。イベント E において、スマート記憶ユニット 0 は、サブディレクトリcomedyのメタデータ構造がスマート記憶ユニット 3 で格納されているという情報を受信する。イベント F において、スマート記憶ユニット 0 は、ファイルmymovie.movieのメタデータ構造のロケーションを要求しているスマート記憶ユニット 3 に要求を送信する。イベント G において、スマート記憶ユニット 0 は、ファイルmymovie.movieのメタデータのデータ構造がスマート記憶ユニット 0 で格納されているという情報を受信する。そして、スマート記憶ユニット 0 は、イベント H においてファイルmymovie.movieのメタデータのデータ構造をローカル記憶から読み取る。このメタデータのデータ構造から、スマート記憶ユニット 0 は、ファイルにおけるデータの各ブロックのロケーションを格納しているmymovie.movieのデータロケーションテーブルを読み取る。そして、スマート記憶ユニット 0 は、データロケーションテーブル情報を使用して、ローカルに格納されたブロックの読み取りと、他のスマート記憶ユニットで格納されたデータに対する要求の送信とを開始する。

【 0 0 7 2 】

ファイルのデータ又はデータの一部が読み取られた後に、ファイルデータが、要求元のサーバ 1 2 0 に送信されて、要求元のユーザ 1 3 0 に転送される。一例において、ファイルデータは、データがいつどのようにユーザ 1 3 0 に送信されるかを定めるビデオストリーミングサーバにルーティングされてもよい。いくつかの実施形態において、要求の遅れ(latency)を低減するために、要求されたものよりも多くのデータを読み出す先見読み出し技術(read ahead techniques)を利用することが有利であるが理解される。

【 0 0 7 3 】

IV. インテリジェントファイルシステム構造

テーブル I は、ファイルシステムレイヤのセット例の一実施形態を示し、そのファイル

10

20

30

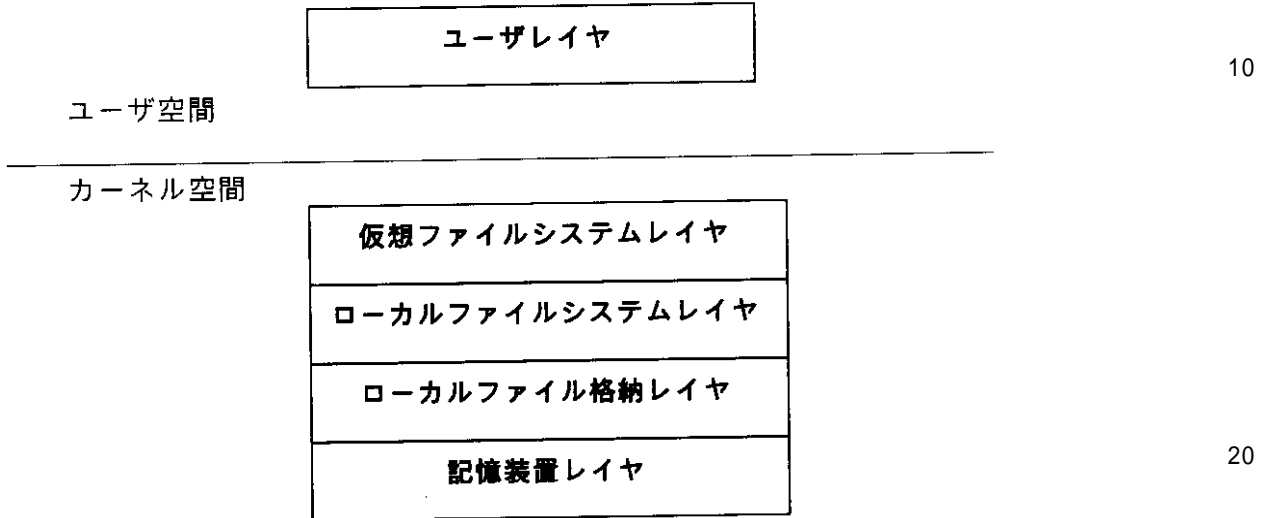
40

50

システムレイヤを通して、物理的記憶装置をアクセスするためにファイル要求が処理される。この代表的ファイルシステムレイヤは、ユーザレイヤ(User Layer)、仮想ファイルレイヤ(Virtual File System Layer)、ローカルファイルシステムレイヤ(Local File System Layer)、ローカルファイルストアレイヤ(Local File Store Layer)、および記憶装置レイヤ(Storage Device Layer)である。

【 0 0 7 4 】

【表 1】



【 0 0 7 5 】

ファイル要求の一つのタイプにおいて、例えば、マイクロソフトウィンドウズ（登録商標）ファイルシェアリングサーバプロトコルのバージョンを実施するユニックス上で使用されるSMBD、FTPD、HTTDP（アパッチウェブサーバ）のようなファイルシェアリングのためのユーザレベルプロトコルアプリケーションを介して、要求が受信される。ユーザレベルプロトコルアプリケーションは、カーネルレベルのオープン(open)、リード(read)、シーク(seek)、ライト(write)、または、例えば、Cランタイムライブラリlibcへのファンクションコールを行うことによるようなクローズシステムコール(close system call)を実行する。

【 0 0 7 6 】

システムコールは、バッファキャッシュを保持する仮想ファイルシステムレイヤ（“VFS”）に受け渡される。このバッファキャッシュは、例えば、最も使用頻度が低いバッファのキャッシュ（“LRU”）であってもよく、このバッファは、下層のファイルシステムレイヤから受信されたメタデータ構造またはデータを格納するために使用される。

【 0 0 7 7 】

次のレイヤは、ローカルファイルシステムレイヤであり、このローカルファイルシステムレイヤは、ファイルの階層的ネーミングシステムを保持すると共にディレクトリおよびファイル名前要求を、下のレイヤであるローカルファイルストアレイヤに送信する。ローカルファイルシステムレイヤは、メタデータのデータ構造ルックアップ及びマネジメント(management)を処理する。例えば、いくつかのシステムにおいて、ユニックススペースのファイルシステムのように、メタデータのデータ構造は、ファイルアクセス許可、データブロックロケーション、およびレファレンスカウント(reference counts)についての情報を備えるファイルアブストラクション(file abstraction)である。一旦、その名前によってファイルが開かれると、他のファイルオペレーションは、特定のファイルについてメタデータ構造を識別するユニークな識別子を介してファイルを参照する。このアプローチの利点は、単一のファイルが多くの異なる名前を有してもよく、単一のファイルが異なる複数のパスを介してアクセスされてもよく、そして、例えば‘mv’コマンドのような標準UNIX（登録商標）ユーザレベルユーティリティを介して実際のデータファイル

を上書きすることなく、V F S 名前スペースにおける古いファイル上に新規なファイルがコピーされてもよい。現在のコンテンツサービスを混乱させることなく適切にコンテンツが更新されるので、これらの利点は、コンテンツ配信およびウェブホスティングのような環境において一層有利である。メタデータのデータ構造内のレファレンスカウントは、一旦、全てのオープンファイルハンドルが閉じられると、システムがデータブロックを無効とすることのみを可能とする。

【 0 0 7 8 】

4 番目のレイヤは、ローカルファイルストアレイヤであり、このローカルファイルストアレイヤは、“バッファ要求からブロック要求へ”の変換(translation)およびデータバッファ要求マネージメントを処理する。例えば、ローカルファイルストアレイヤは、ブロック割り当てスキーム(block allocation scheme)を使用して、読み出し(reading)のためのブロック読み出しスキーム(block retrieval scheme)と同様に、ライトおよびリードについてのスループットを改善し且つ最小化する。

10

【 0 0 7 9 】

最後のレイヤは記憶装置レイヤであり、この記憶装置レイヤは、ファイルシステムによって使用される特定の一つのディスクハードウェア用の装置ドライバをホスト(host)する。例えば、もし物理的な記憶装置が A T A ディスクであれば、記憶装置レイヤは A T A ディスクドライバをホストする。

【 0 0 8 0 】

テーブル 1 はファイルシステムの一つのセット例を示すが、H は、ファイルシステムレイヤが異なって定義されてもよく又は全く存在しなくてもよいような異なる環境におけるのと同様に種々の異なるレイヤでこのシステム及び方法が動作してもよいと認識される。

20

【 0 0 8 1 】

V. スマート記憶ユニット

一実施形態において、スマート記憶ユニット 1 1 4 は、プラグアンドプレイ(plug-and-play)の、高密度で、ラックマウント可能な装置であり、それは高スループットデータ配信に最適化されている。このスマート記憶ユニットは、単一の仮想ファイルシステムを提供するために、種々の他のスマート記憶ユニットと通信するように構成されてもよい。より多くの記憶スペースが必要とされるので、または、もし 1 または 2 以上のスマート記憶ユニットが故障すれば、追加のスマート記憶ユニットがシステム全体をダウンさせ又はサービスの中断を生じる必要なくインストールされてもよい。

30

【 0 0 8 2 】

本明細書で使用されるように、モジュールなる用語は、ハードウェアまたはファームウェアで実現されたロジックを称し、またはソフトウェア命令の集まりを称し、できれば入口点(entry point)または出口点(exit point)を有し、例えば C または C + + のようなプログラミング言語で記述される。ソフトウェアモジュールは、コンパイルされると共に実行プログラムにリンクされて、動的リンクライブラリにインストールされてもよく、または、B A S I C、P e r l、P y t h o n のような翻訳プログラミング言語(interpreted programming language)で記述されてもよい。ソフトウェアモジュールは、他のモジュールまたは自身から読み出し可能であってもよく、及び/又はイベントまたはインタラプトに応答して呼び起こされてもよい。また、ソフトウェア命令は、E P R O M のようなファームウェアに組み込まれてもよいことが理解される。さらに、ハードウェアモジュールが、ゲートおよびフリップフロップのようなロジックユニットから構成されてもよく、及び/又はプログラマブルゲートアレイまたはプロセッサのようなプログラム可能なユニットから構成されてもよい。本明細書で述べられるモジュールは、好ましくは、ソフトウェアモジュールとして実施されるが、ハードウェア又はファームウェアで表現されてもよい。

40

【 0 0 8 3 】

図 3 は、マネージメントモジュール 3 2 0、処理モジュール 3 3 0、キャッシュ 3 4 0、スタック 3 5 0、および記憶装置 3 6 0 を備えるスマート記憶ユニット 1 1 4 の一実施形態を示す。代表的なスマート記憶ユニット 1 1 4 は、図 1 に示されるように、要求を送

50

信および受信するために、スイッチコンポーネント 1 2 5 と通信するように構成される。

【 0 0 8 4 】

A . マネージメントモジュール

一実施形態において、スマート記憶ユニットは、例えば、インストール(installation)、パラメータ設定、インテリジェント分散ファイルシステムのモニタ、インテリジェント分散ファイルシステム 1 1 0 で発生するイベントのロギング(logging)、更新などの管理タスクを実行するためのマネージメントモジュール 3 2 0 を備える。

【 0 0 8 5 】

B . 処理モジュール

代表的な処理モジュール 3 3 0 は、データファイルについての要求を受信し、要求されたデータファイルについてのローカル及び/又はリモートに格納されたメタデータを読み出し、そしてこの要求されたデータファイルのローカル及び/又はリモートに格納されたデータブロックを読み出すように構成されてもよい。加えて、処理モジュール 3 3 0 は、また、1または2以上の要求されたデータブロックが崩壊または失われたイベントにおいてエラー訂正およびデータ回復を実施してもよい。

10

【 0 0 8 6 】

一実施形態において、処理モジュール 3 3 0 は、ファイル要求に応答するための5つのモジュール、即ちブロック割り当てマネージャモジュール 3 3 1、ブロックキャッシュモジュール 3 3 3、ローカルブロックマネージャモジュール 3 3 5、リモートブロックマネージャモジュール 3 3 7、およびブロック装置モジュール 3 3 9 を備える。

20

【 0 0 8 7 】

1 . ブロック割り当てマネージャモジュール

ブロック割り当てマネージャモジュール 3 3 1 は、リード要求に응答して、ブロックをどこに割り当て、このブロックをどこに配置するかを決定し、そして装置に故障回復を指示する。ブロックをどこに割り当てるかについての情報は、グラフィカルユーザインターフェイス又はシェルインターフェイス、またはこれらの組み合わせのようなツールを介してシステム管理者によって設定されたポリシー、デフォルトパラメータとして設定されたポリシーによって決定されてもよい。一実施形態において、ブロック割り当てマネージャモジュール 3 3 1 は、少なくともローカルファイルシステムレイヤに存在し、そしてパークレイソフトウェアデザインユニバーサルファイルシステム(“BSD UFS”)に代

30

【 0 0 8 8 】

代表的なブロック割り当てマネージャモジュール 3 3 1 は、3つのサブモジュール、即ちブロック要求変換モジュール(block request translator module)、フォワード割り当てモジュール(forward allocator module)、および故障回復モジュール(failure recovery module)を備える。

【 0 0 8 9 】

a . ブロック要求変換モジュール

ブロック要求変換モジュールは、到来するリード要求を受信し、ネームルックアップを実施し、適切な装置を探し、そしてこの装置からデータを引き出して要求を満足させる。もしこのデータが直ちに利用可能であれば、ブロック要求変換モジュールは、データのブロックがローカル記憶装置に格納されているか又は他のスマート記憶ユニットの記憶装置に格納されているかによって、データ要求をローカルブロックマネージャモジュールまたはリモートブロックマネージャモジュールに送信する。

40

【 0 0 9 0 】

一実施形態において、ブロック要求変換モジュールは、“インテリジェント分散ファイルシステム処理 - ネームルックアップ処理”と題された節において述べられる以下のネームルックアップ処理を含む。

【 0 0 9 1 】

50

ブロック要求変換モジュールは、また、装置故障に応答してもよい。例えば、もし装置がダウンしていれば、ブロック要求変換モジュールは、使用中のデータ(data using)、例えばパリティ情報を再構成するために使用されるローカル及びリモートデータブロックを要求してもよい。従って、このデータは、たとえリードが実施されなくても生成される。加えて、ブロック要求変換モジュールは、故障回復モジュールがパリティ又は他のエラー又は欠損訂正データを用いてデータを再生成し、インテリジェント分散ファイルシステムにおける自由スペースにわたって欠損訂正データ(loss correction data)を再ストライピングするように、故障回復モジュールと通信してもよい。他の実施形態において、ブロック要求変換モジュールは、崩壊したデータ又は欠落したデータの鮮明なコピー(clean copy)を要求してもよい。

10

【0092】

b. フォワード割り当てモジュール

フォワード割り当てモジュールは、例えば冗長性、スペース、および性能のような要素に基づき、どの装置のブロックがライト要求に使用されるべきかを決定する。これらのパラメータは、システム管理者によって設定されてもよく、インテリジェント分散ファイルシステム110に組み込まれた情報から導き出されてもよく、インテリジェント分散ファイルシステム110にロジックとして組み込まれてもよく、これらの組み合わせによってもよい。フォワード割り当てモジュール110は、インテリジェント分散ファイルシステムを使用する他のスマート記憶ユニットから統計データ(statistics)を受信し、そしてこれらの統計データを使用して、新たに到来するデータを配置すべき最良のロケーションはどこなのかを決定する。集められたこの統計データは、例えば、CPU稼働率、ネットワーク稼働率、およびディスク稼働率を含む。

20

【0093】

フォワード割り当てモジュールは、また、リモートのスマート記憶ユニットの応答時間に基づき、このリモートのスマート記憶ユニットからの遅延情報(latency information)を受信してもよい。もし、中間装置の遅延が他のスマート記憶ユニットに比較して高いレベルに達すれば、割り当てスキームは、できれば冗長設定に基づき、遅いスマート記憶ユニットを十分に活用しない他のスマート記憶ユニットを助けるように調整されてもよい。一つの好都合な例において、インテリジェント分散ファイルシステムは、データのブロックを或るスマート記憶ユニットから他のスマート記憶ユニットに移動させ、対応するメタデータ構造を適宜更新してもよい。遅延状況(latency condition)は、ロギングシステムを通して記録されてもよく、そしてシステム管理者に報告されてもよい。遅いリンク状況(slow link condition)の原因は、例えば、不良のネットワークカード、不適切な双方向情報交換(duplex negotiation)、または比較的頻繁にリードまたはライトされる装置のデータであってもよい。

30

【0094】

種々の方策が、データをどこに格納するかを決定するために使用される。これらの方策は、パラメータとの整合性、またはシステムの管理者により設定されるプリフェレンス(reference)、選択された冗長性レベルの一致、及び/又は性能改善のようなシステムの目標に基づき調整されてもよい。以下では、データを格納するためにフォワード割り当てモジュールによって採用される方策の例を少し提供する。以下に述べられることに加えて、またはそのことと共に、広範な種々の方策が使用されてもよいことが理解される。

40

【0095】

フォワード割り当てモジュールは、複数のスマート記憶ユニットにわたってデータをストライピングするための割り当てスキームを備えてもよい。データのストライピングは、一般にはハイエンドのRAID記憶装置で使用される公知技術であるが、複数のディスクを備える単一ユーザーワークステーションマシンに採用されてもよい。データのストライピングは、単に、ファイルのデータの異なる部分が、存続(live)し及び/又は異なる記憶装置またはディスクに格納されることを意味する。データのストライピングに対する利点は、リード要求が複数のディスクに配置されたブロックに及ぶときに、各ディスクがデータ

50

読み出しの総スループットに關与することである。代表的なシステムでは、データのストライピングはソフトウェア装置レイヤで行われる。即ち、ファイルシステムは、データのストライピングに関する情報を持たない。ファイルシステム下にあるソフトウェアレイヤのみがこの構造を理解する。いくつかの特化したハードウェアにおいては、このストライピングは、実際のハードウェアレイヤでソフトウェア装置レイヤの下でさえ行われる。インテリジェント分散ファイルシステム 1 1 0 においては、ファイルシステム自体は、データのストライピングを処理する。この実施は、より大きな柔軟性をストライピング構成に与える。例として、一般的な R A I D 技術は、全ディスクが同一のサイズでなければならず、且つ同一の性能特性を有していなければならないという制約がある。これらの制約は、データが異なる装置に均一に分散されることを保証するために必要である。R A I D についての更なる詳細については、Paul Massiglia による “The RAID Book” 第 6 版 (1 9 9 7) を参照されたく、それは参照することにより本明細書に組み込まれる。

10

【 0 0 9 6 】

加えて、フォワード割り当てモジュールは、ユーザ要求 (例えば、新しいスマート記憶ユニットのプロテクションスキームの付加など) に応答して、及び / 又はシステム (例えば、故障したスマート記憶ユニットの検出) に応答してデータのストライピングを実施してもよい。しかしながら、再ストライピングは、ブロック割り当て管理モジュールにおける他のモジュールにより又は他のモジュールと共に実施されてもよく、またはスマート記憶ユニットの他の部分により又はスマート記憶ユニットの他の部分と共に実施されてもよい。

20

【 0 0 9 7 】

インテリジェント分散ファイルシステム 1 1 0 については、ディスクおよびディスクサイズを異なえることは、種々のスマート記憶ユニット 1 1 4 で使用され、そしてファイルストライピングに關与する。フォワード割り当てモジュールは、ディスク装置情報のためのルートメタデータ構造を調べ (look up)、そして、性能メトリックまたは前もってセットされた規則を用いて、ファイルデータが拡散されるべきスマート記憶ユニットの数を計算する。そして、フォワード割り当てモジュールは、ファイルのブロックをスマート記憶ユニットのセットに割り当ててもよい。

【 0 0 9 8 】

また、フォワード割り当てモジュールは、パリティ又は他のエラー又は欠損訂正プロテクションのための割り当てスキームを備える。ほとんどの R A I D システムにおいては、ファイルストライピングが使用されるとき、ひとつを除いて全てのディスクがデータ記憶のために使用されるようにパリティプロテクションもまた使用される。最後のディスクは、純粹にパリティ情報のために使用される。このパリティ情報は、一般には、データディスクの全てにわたって各データのブロックのビット単位の排他的論理和 (“ X O R ”) をとることにより計算される。このパリティ情報は、ディスク故障が発生したときにデータ回復を実施するために使用される。失われたデータは、残りのディスクのデータブロックおよびパリティ情報のビット単位の X O R をとることで再計算される。一般的な R A I D システムでは、データは、失われたデータを再編するために交換用ディスクが配列 (array) に挿入されるまで回復不能である。

30

40

【 0 0 9 9 】

インテリジェント分散ファイルシステム 1 1 0 では、パリティプロテクションがソフトウェア装置レイヤの代わりにファイルシステムレイヤで行われるので、失われたデータは、再計算されて、残りのスマート記憶ユニットの他の部分の自由スペースに再書き込みされてもよい。もし、そのデータを再書き込みするための十分な自由スペースが残っていなければ、パリティデータは、再計算されたデータで上書きされてもよく、そして冗長性がオリジナルのレベル以下に低下したという事実が記録され及び / 又はシステム管理者に報告されてもよい。

【 0 1 0 0 】

また、フォワード割り当てモジュールは、データのミラーリングのための割り当てスキ

50

ームを備えてもよく、それは異なるスマート記憶ユニット上で利用可能なデータの複数のコピーを作成している。フォワード割り当てモジュールは、記憶スペース、ネットワーク稼働率、及び/又はCPU稼働率の観点から使用される度合いが最も少ないスマート記憶ユニットを用いて、それらスマート記憶ユニットにわたってデータのブロックのロケーションを負荷バランスするための割り当てスキームを使用してもよい。ミラーリングは、さらなる性能および更なる耐故障性を提供してもよい。もし、ミラーリングが、あるいくつかのコンテンツに対して要求されれば、フォワード割り当てモジュールは、ミラーされたデータと同様にオリジナルのデータにスペースを割り当てて、もし、1よりも大きな耐故障性レベルが要求されれば、フォワード割り当ては、耐故障性カウント(fault tolerance count)によって、スマート記憶ユニット、またはスマート記憶ユニットのサブセットを論理的に分割し、そしてストライピングされたデータを生成してもよい。例えば、もし、10のスマート記憶ユニットがインテリジェント分散ファイルシステム110に存在し、且つ2の耐故障性が要求されれば、フォワード割り当ては、インテリジェント分散ファイルシステムを、それぞれが5つのスマート記憶ユニットの2つのセクションに分け、各セクションにおける4つのスマート記憶ユニットにわたってデータをストライピングし、そして各セクションから5番目のスマート記憶ユニットをパリティディスクとして使用してもよい。このスマート記憶ユニットの分割は、配列ミラー分割(array mirror split)と称される。配列ミラー分割の種々の実施は、記憶時にデータが重複(overlap)され又は歪曲(skew)される実施を含んで使用されてもよい。

10

【0101】

20

c. 故障回復モジュール

故障回復モジュールは、実時間でインテリジェント分散ファイルシステム110を再構成し、装置故障のためにはや利用できないデータを回復する。故障回復モジュールは、性能を維持しながらサービスを中断することなく再構成を実施してもよく、そしてそのデータを短時間で所望の冗長レベルに戻してもよい。

【0102】

上述したように、リモートブロック管理モジュール337は、故障を検出して、このような故障の通知を故障回復モジュールに受け渡す。初期故障については、故障回復モジュールは、システム管理者により要求され又はインテリジェント分散ファイルシステム110によって設定されるような冗長パラメータに合致しない何れかのデータを配置する。

30

【0103】

最初に、パリティ情報から再生成され得るデータが再生成され、そして新たなデータにスペースを割り当てて、要求がフォワード割り当てモジュールに送られる。フォワード割り当ては、CPUとネットワークの稼働率をモニタし、そしてCPUおよびネットワーク稼働率が所定の目標に達するまで積極的に動作を始める。この所定の目標は、システム管理者によって設定されてもよく、あるいは、例えばコンピュータプロセッサのような要素によって予め設定されてもよい。一旦、その目標に達すると、故障回復モジュールは、スマート記憶ユニットの性能に与える影響を低減するために、その目標の時点で到達したレートで再計算を有利に実施してもよい。

【0104】

40

もし、最近故障した装置がオンラインに復帰すれば、故障回復モジュールはその回復した装置のリモートブロック管理モジュール337と通信して、データの完全性を確かめ、そして不一致があれば修正する。

【0105】

また、インテリジェント分散ファイルシステム110は、ホットスタンバイ装置の具備を支援してもよい。このホットスタンバイ装置は、現在のところは何のデータも処理していないが装置故障のときに使用される待機状態(idle)の記憶装置である。このような状況において、故障回復モジュールは、ホットスタンバイ装置のリモートブロック管理モジュール337と通信することにより、ホットスタンバイ装置を用いて、失われたデータを再構築してもよい。

50

【 0 1 0 6 】

2 . ブロックキャッシュモジュール

ブロックキャッシュモジュール 3 3 3 は、データブロックのキャッシング(caching)、ネームルックアップ、およびメタデータ構造を管理する。一実施形態において、ブロックキャッシュモジュール 3 3 3 は、BSD 仮想ファイルシステム(BSD Virtual File System)のバッファキャッシュの代わりに、またはBSD 仮想ファイルシステムのバッファキャッシュと共に動作する。

【 0 1 0 7 】

ブロックキャッシュモジュール 3 3 3 は、例えば頻度キャッシング(frequency caching)のような種々のキャッシングアルゴリズムを使用できるが、使用頻度最低(Least Recently Used)キャッシングアルゴリズムを用いてデータブロックとメタデータのデータブロックを格納(cache)してもよい。ブロックキャッシュモジュール 3 3 3 は、どちらが最善を遂行するかに応じて、どのブロックキャッシングアルゴリズムを使用するかを決定してもよく、他の実施形態では、アルゴリズムはデフォルトとして設定されてもよい。

10

【 0 1 0 8 】

使用頻度最低キャッシング(“LRU”)は、ほとんどのシステムで使用されている代表的なキャッシングスキームである。LRUは、一旦データがアクセスされると、そのデータは再びアクセスされる可能性が高いという原理に基づいている。従って、データは、最も長い時間アクセスされていないデータが捨てられるように、その最近の使用の順に格納される。

20

【 0 1 0 9 】

頻度キャッシングは、最も頻繁にアクセスされているデータを格納する。ディスクライト(disk write)は、比較的集約的な動作であるため、追加的性能は、メタデータのデータ構造でアクセス頻度を追跡すること、およびアクセス頻度に基づきキャッシングすることにより得られる。

【 0 1 1 0 】

加えて、ブロックキャッシュモジュール 3 3 3 は、より多くのデータが必要以上に要求される“先見読み出し(read ahead)”または“オンデマンド(on demand)”プロトコルを利用してよい。ブロックキャッシュモジュール 3 3 3 は、データのセットに対する要求を送信してもよく、またデータのセットに先んじてその若干のデータを要求してもよい。例えば、ブロックキャッシュモジュール 3 3 3 は、1 パケット先見読み出し、2 パケット先見読み出し、10 パケット先見読み出し、12 パケット先見読み出し等々のような先見読み出しを実施してもよい。他の実施形態において、ブロックキャッシュモジュール 3 3 3 は、要求の遅延(latency)に基づいて先見読み出しを利用してよい。例えばブロックキャッシュモジュール 3 3 3 は、V パケット先見読み出しを実施してもよく、ここでVはリンクの遅延(latency)および読み出しレート(read rate)を用いて計算される。ブロックキャッシュモジュール 3 3 3 は、また、CPUおよびネットワーク稼働率に基づいた他のアルゴリズムを使用して先見読み出しデータのサイズを決定してもよい。さらにまた、ブロックキャッシュモジュール 3 3 3 は、セットキャッシングプロトコルを使用してよく、またはシステムの性能レベルに応じてキャッシングプロトコルを変えてもよい。

30

40

【 0 1 1 1 】

キャッシュ 3 4 0 は、一般のマルチユーザのオペレーティングシステムに備わるデフォルトサイズを用いて実施されてもよく、またはシステム性能に大幅に影響を与えることのない範囲でキャッシュブロックサイズを異なる量に増加するように修正されてもよい。このような修正は、例えば、格納されるデータのタイプ、処理速度、インテリジェント分散ファイルシステムにおけるスマート記憶ユニットの数、及び使用されているプロテクションスキームのような要素に依存する種々の性能テストによって決定されてもよい。

【 0 1 1 2 】

3 . ローカルブロックマネージャモジュール

ローカルブロックマネージャモジュール 3 3 5 は、記憶装置 3 6 0 にローカルに格納さ

50

れたデータブロックの読み出し(retrieval)、記憶(storage)、割り当て(allocation)を管理する。ローカルブロックマネージャモジュール335は、ディスクから、例えばネットワークカードのような記憶装置360の他の部分にデータを移動させるためにゼロコピーファイル読み出し(zero copy file read)を実施し、これにより性能を改善してもよい。ローカルブロックマネージャモジュール335は、また、性能を向上させるために、使用されている記憶装置360に基づく修正を実施してもよい。一実施形態において、ローカルブロックマネージャモジュール335は、ローカルファイル格納レイヤ(Local File Store layer)に存在してもよく、FreeBSD高速ファイルシステム(FreeBSD Fast File System)の代わりに又はFreeBSD高速ファイルシステムと共に稼動してもよい。

【0113】

一実施形態において、ローカルブロックマネージャモジュール335は、データを記憶装置360に格納するための要求を処理してもよい。一実施形態において、ローカルブロック管理モジュール335は、データが記憶装置360上のどこに格納されるかを決定する。例えば、ローカルブロック管理モジュールは、既に格納されているデータと関連するデータを受信したときに、新たなデータが、可能な限りその関連データの近くに格納されるように、関連データを隣接させて格納するように試みてもよい。しかしながら、種々の記憶プリフェレンス(storage preference)が使用され、且つ各スマート記憶ユニットが1または2以上の異なる記憶プリフェレンスを使用してもよいことが理解される。一実施形態において、インテリジェント分散ファイルシステムにおけるスマート記憶ユニットは同一の記憶プリフェレンスを使用してもよい。

【0114】

4. リモートブロックマネージャモジュール

リモートブロックマネージャモジュール337は、例えば、ブロック要求、ブロック応答、およびリモート装置故障の検出を含む中間装置通信を管理する。一実施形態において、リモートブロックマネージャモジュール337は、ローカルファイルシステムレイヤ(Local File System layer)に存在する。

【0115】

一実施形態において、スマート記憶ユニット114は、リモートブロックマネージャモジュール337を介してインテリジェント分散ファイルシステム110における他のスマート記憶ユニットと通信及び/又は接続してもよい。

【0116】

リモートブロックマネージャモジュール337は、スマート記憶ユニット114がTCPのようなコネクションを介して互いにやり取り(talk)することを可能にしてもよい。一実施形態においては、各スマート記憶ユニットについて少なくとも2つのTCPコネクションが存在し、ひとつはファイルデータトランスポートのためのものであり、もうひとつはコントロールメッセージトランスポートのためのものである。この二重チャンネル(dual channel)TCP通信アーキテクチャの利点は、データがシステムの或る部分から他の部分にコピーされる必要なく、データブロックが複数のページサイズで送信される限りデータがDMA転送を介して直接的にネットワークインターフェイスカードからシステムメモリに送信されてもよく、且つDMA転送を介してシステムメモリからシステムの他の部分(おそらくは再びネットワークインターフェイスカード)に送信されてもよいことである。これは、データパケットが非データヘッダ(non-data header)を含まないときにCPUがそのデータパケットを解析し、または情報がコントロールチャンネル上に伝送されるときにこの情報をCPUが識別することに関与する必要がないためである。高性能のサーバ及びオペレーションシステムにおいては、システムメモリの或る部分から他の部分へのこれらのメモリコピーは、システム性能に関する厳しい制限になる。

【0117】

一実施形態において、リモートブロックマネージャモジュール337は、例えば、データブロックアクセスメッセージ(例えば、READ, READ_RESPONSE, WRITE, WRITE_RESPONSE)、メタデータアクセスメッセージ(例えば、GET_INODE, GET_INODE_RESPONSE, SET_ADD

10

20

30

40

50

RESS, GET_ADDRESS, INVALIDATE_INODE)、ディレクトリメッセージ(例えば、ADD_DIR, REMOVE_DIR)、ステータスメッセージ、同様の種々の他のタイプのメッセージのようなメッセージを利用するメッセージング通信を用いて通信する。

【0118】

ここでは二重チャンネルプロトコルについて述べたが、スマート記憶ユニット114間の通信を可能とするために他の通信プロトコルを使用してもよいことが理解される。

【0119】

5. ブロック装置モジュール

ブロック装置モジュール339は、ファイルシステムによって使用される特定の一部のディスクハードウェアのための装置ドライバをホスト(host)する。例えば、もし物理的記憶装置がATAディスクであれば、ブロック装置モジュール339はATAディスクドライバをホストする。

【0120】

C. キャッシュ

キャッシュメモリ又はキャッシュ340は、例えば、1GRAMキャッシュのような本技術分野において公知の種々の製品を用いて実施されてもよい。図3に示されるキャッシュ340は、最近アクセスされたデータのブロック、または一定時間内にアクセスされるべきデータのブロックを格納してもよい。キャッシュ340は、スタティックRAM装置、ダイナミックRAM装置、内部キャッシュ、ディスクキャッシュ、同様の種々の他のタイプの装置のような高速記憶装置を用いて実施されてもよい。通常、データは、不揮発性記憶装置をアクセスするために要する時間よりも速くキャッシュ340からアクセスされる。もしスマート記憶ユニット114が記憶装置360からデータをアクセスする必要があるならば、そのデータが既に読み出されたかを調べるためにキャッシュ340が最初にチェックされるようにキャッシュ340がデータを格納する。従って、キャッシュ340の使用は、データブロックの読み出しにおいてスマート記憶ユニットの性能を改善する。

【0121】

D. ネットワークスタック

一実施形態において、また、スマート記憶ユニット310は、例えばTCP/IPのようなプロトコルを用いて、入ってくるメッセージトラフィックと、出て行くメッセージトラフィックとを処理するネットワークスタック350を備える。しかしながら、スタック350を実施するためには他のプロトコルまたはデータ構造が使用されてもよいことが理解される。

【0122】

E. 記憶装置

記憶装置360は、1又は2以上の不揮発性メモリ装置のセットであり、それはデータブロックを格納するために使用されてもよい。記憶装置360は、例えば、41.25GB ATA100装置、SCSI装置等のような公知の種々の製品を用いて実施されてもよい。加えて、記憶装置360のサイズは、インテリジェント分散ファイルシステム110における全てのスマート記憶ユニット114について同一であってもよく、或いはそれは異なるスマート記憶ユニット114についてサイズを変えてもよい。

【0123】

F. システム情報

一実施形態において、スマート記憶ユニット114は、スマート記憶ユニット114が他のスマート記憶ユニット114と通信することを可能とするコンピュータ上で稼動する。このコンピュータは、例えば、ペンティアム(登録商標)プロセッサ、ペンティアム(登録商標)IIプロセッサ、ペンティアム(登録商標)Proプロセッサ、ペンティアム(登録商標)IVプロセッサ、xx86プロセッサ、8051プロセッサ、MIPSプロセッサ、パワーPC(登録商標)プロセッサ、SPARC(登録商標)プロセッサ、アルファプロセッサ等のような1又は2以上マイクロプロセッサを用いた汎用のコンピュータであってもよい。

10

20

30

40

50

【 0 1 2 4 】

一実施形態において、プロセッサユニットは、オープンソースのFreeBSDオペレーティングシステムを動作させ、そしてファイルのオープン、読み取り、書き込み、およびクローズのような標準オペレーティングシステムの機能を実施する。例えば、マイクロソフト（登録商標）ウィンドウズ（登録商標）3.X、マイクロソフト（登録商標）ウィンドウズ（登録商標）98、マイクロソフト（登録商標）ウィンドウズ（登録商標）2000、マイクロソフト（登録商標）ウィンドウズ（登録商標）NT、マイクロソフト（登録商標）ウィンドウズ（登録商標）CE、マイクロソフト（登録商標）ウィンドウズ（登録商標）ME、パームパイロットOS、アップル（登録商標）Mac OS（登録商標）、ディスクオペレーティングシステム（DOS）、UNIX（登録商標）、IRIX（登録商標）、Solaris（登録商標）、Sun（登録商標）OS、FreeBSD（登録商標）、Linux（登録商標）、IBM（登録商標）OS/2（登録商標）オペレーティングシステムのような他のオペレーティングシステムが使用されてもよいことが理解される。

10

【 0 1 2 5 】

一実施形態において、コンピュータは、例えば、イーサネット（登録商標）（IEEE 802.3）、トークンリング（IEEE 802.5）、FDDI（Fiber Distributed Data Link Interface）、またはATM（Asynchronous Transfer Mode）のような従来のネットワーク接続性（network connectivity）を備える。さらに、コンピュータは、例えば、NFS v2/v3 over UDP/TCP、マイクロソフト（登録商標）CIFS、HTTP 1.0、HTTP 1.1、DAFS、FTP等のような種々のネットワークプロトコルを支援するように構成されてもよい。

20

【 0 1 2 6 】

一実施形態において、スマート記憶装置114は、ジャンボ9Kイーサネット（登録商標）フレームを支援する1000/100ネットワークインターフェイスカードと同様に、シングル又はデュアルCPU 2Uラック搭載可能な構成を備える。しかしながら、種々の構成が使用されてもよいことが理解される。

【 0 1 2 7 】

上述したように、異なるモジュールがスマート記憶ユニットに関して述べられたが、そのタスクは異なるモジュールによって実行されてもよいことが理解される。加えて、1又は2以上のモジュールが組み合わされることができ、及び/又は、ある特定のモジュールが特定のタスクを実施することを要求されないように、1又は2以上の新たなモジュールが付加されてもよい。

30

【 0 1 2 8 】

VI. インテリジェント分散ファイルシステムデータ構造

図4は、インテリジェント分散ファイルシステムで使用されるディレクトリ構造例を示す。この例において、ルートディレクトリは、“DFSR”と名づけられ、そしてサブディレクトリIMPORTANT、TEMP、USERを備える。サブディレクトリIMPORTANTは、サブディレクトリPASSWORDSおよびCREDITCARDを備える。ファイルUSER.TXTおよびADMIN.TXTは、サブディレクトリPASSWORDSに格納される。従ってファイルUSER.TXTについてのアドレスは、

/DFSR/IMPORTANT/PASSWORDS/USER.TXT

40

である。ディレクトリおよびファイルについての情報またはメタデータは、インテリジェント分散ファイルシステム110により格納され且つ保持される。

【 0 1 2 9 】

A. メタデータのデータ構造

図5は、メタデータを格納するためのデータ構造例を示す。代表的なデータ構造510は、次の情報を格納する。

【 0 1 3 0 】

【表 2】

フィールド	説明
モード	ファイルのモード（例えば、レギュラーファイル、ブロックスPECIAL、キャラクターSPECIAL、ディレクトリ、シンボリックリンク、fifo、ソケット、ホワイトアウト、アンノウン）
オーナー	ファイルのオーナーシップを有するスマート記憶ユニットに関するアカウント
タイムスタンプ	最終変更ファイルのタイムスタンプ
サイズ	メタデータファイルのサイズ
パリティカウント	使用されたパリティ装置の数
ミラーカウント	使用されたミラー装置の数
バージョン	メタデータ構造のバージョン
タイプ	データロケーションテーブルのタイプ（例えば、タイプ0、タイプ1、タイプ2、タイプ3）
データロケーションテーブル	データロケーションテーブルまたは実際のデータロケーションテーブル情報のアドレス
リファレンスカウント	これを参照するメタデータ構造の数
フラグ	ファイルパーミッション（例えば標準UNIXのパーミッション）
パリティマップポインタ	パリティブロック情報に対するポインタ

10

20

【0131】

データ構造510の例は、メタデータを格納するためのデータ構造510の実施形態を示し、且つ種々の実施が本発明に従って実施されてもよいことが理解される。例えば、データ構造510は異なるフィールドを備えてもよく、このフィールドは異なるタイプのものであってもよく、このフィールドはグループ化されて別々に格納されてもよい。等々。

【0132】

図6A、6B、6C、6Dは、データロケーションテーブルのタイプ、即ちタイプ0、タイプ2、タイプ3についてのデータロケーションテーブル構造例をそれぞれ示す。図6Aにおいて、タイプ0のデータロケーションテーブルは、24の直接ブロックエントリを備え、これは、データが格納されるロケーションを示す装置/ブロック番号対をデータロケーションテーブルにおけるエントリが備えることを意味する。図6Bにおいて、タイプ1のデータロケーションテーブルは、15の直接ブロックエントリと、3つのシングル間接エントリと、3つのダブル間接エントリと、3つのトリプル間接エントリとを備える。シングル間接エントリについてのエントリは、直接エントリの追加的データロケーションテーブルが格納されるロケーションを示す。ダブル間接エントリについてのエントリは、データロケーションテーブルが格納されるロケーションを示し、このデータロケーションテーブルはシングル間接エントリを備える。トリプル間接エントリについてのエントリは、データロケーションテーブルが格納されるロケーションを示し、このデータロケーションテーブルはダブル間接エントリを備える。

30

40

【0133】

任意数の装置にわたって任意のブロックがミラーされてもよいので、メタデータのデータ構造510は複数のロケーションを用いてブロックを表現するのに十分柔軟性があり、且つ一定のスペース内での直接指標付け(direct indexing)に由来する高速アクセスを依然として提供する。従って、使用されるべきデータロケーションテーブルのタイプを指示するために、タイプは、メタデータ構造510と有利に関連づけられてもよい。メタデータのデータ構造510の一実施形態において、例えば24のポインタのような24のデー

50

タエントリのための場所 (room) があってもよい。

【 0 1 3 4 】

タイプ 0 は、データファイルが小さいとき：データロケーションアドレスが直接エントリとして格納されるときに使用されてもよい。従って、タイプ 0 のメタデータのデータ構造は 2 4 の直接エントリを備える。タイプ 1 は、より大きなファイルと 2 回までのミラー（ファイルの 3 つのコピー）を支援するために使用されてもよい。タイプ 1 は、1 5 の直接エントリと、3 つのシングル間接エントリと、3 つのダブル間接エントリと、3 つのトリプル間接エントリとを使用する。タイプ 2 は、7 回までのミラーリング（ファイルの 8 つのコピー）を支援するために使用されてもよく、そして 8 つのシングル間接エントリ、8 つのダブル間接エントリ、8 つのトリプル間接エントリを備えてもよい。タイプ 3 のデータロケーションテーブルは、ディスクアドレスの全てがトリプル間接エントリとして格納されるので更なるミラーリングを可能する。結果として、2 4 までの完全なファイルコピーが格納される。

10

【 0 1 3 5 】

種々のデータロケーションテーブルが使用されてもよいことが理解され、図 6 A , 6 B , 6 C , 6 D は実施例を示す。他の実施形態において、例えば、データロケーションテーブルは、直接および間接エントリの異なる混合体を備えてもよい。さらに、他の実施形態において、データロケーションテーブルは、テーブルの各エントリについてエントリのタイプを指定するエントリフィールドを備えてもよい。タイプは、例えば、上述したもの（例えば、直接、シングル間接、ダブル間接、トリプル間接）を備えてもよく、他（例えば、クアドラブル間接など）も同様である。加えて、データロケーションテーブルは、X レベルまでのデータロケーションテーブルのより深いネスティング (nesting) を備えてもよく、ここで X は整数である。

20

【 0 1 3 6 】

また、メタデータのデータ構造は、どのスマート記憶ユニットがファイルのコンテンツとプロテクションデータを備えるかについての情報を備える。加えて、メタデータのデータ構造は、ファイルのコンテンツデータのために使用される最後のブロックアドレスとファイルのプロテクションデータのために使用される最後のブロックアドレスとを追跡 (tracking) する各スマート記憶ユニットのための情報を格納してもよい。例えば、メタデータのデータ構造は、MYFILE.TXT が、装置 0、装置 2、装置 3、および装置 5 に格納されたそのデータを有することを記録してもよい。メタデータのデータ構造は、また、次のことを記録してもよい。

30

【 0 1 3 7 】

【表 3】

	コンテンツのための最後の ブロックアドレス	パリティのための最後の ブロックアドレス
装置 0	3 0 0	0 0 1
装置 2	3 0 7	2 0 3
装置 3	2 0 0	3 0 3
装置 5	1 0 3	5 0 1

40

【 0 1 3 8 】

1 . ディレクトリメタデータ

図 7 A は、ディレクトリ PASSWORDS についてのメタデータのセット例を示す。図 7 A において、データ構造はディレクトリ PASSWORDS についての情報を格納する。ディレクトリは 2 回ミラーされる（合計 3 つのコピー）。ディレクトリ構造は比較的小さいので（例えばブロックに収まる）、使用される 3 つの直接ポイントのみが存在し、ひとは各コピーのためのものである。メタデータのセット例は、使用されないブロックエントリ 7 3 0 のセットと同様に装置 / ブロック番号対を用いてデータブロックのロケーションを示す直接

50

エントリ 7 2 0 を有するデータロケーションテーブル 7 1 0 を備える。

【 0 1 3 9 】

2 . ファイルメタデータ

図 7 B は、ファイル USER.TXT についてのメタデータのセット例を示す。図 7 B において、データ構造は、ファイル USER.TXT についての情報を格納する。ファイルデータ USER.TXT についてデータブロックのそれぞれの一つのコピーが存在し、そしてそのデータは 3 + 1 パリティスキームを用いて保護される。USER.TXT についてのコンテンツデータは、サイズが 4 5 K であり、ブロックサイズが 8 K であり、従って 6 番目のデータブロックが全く使用されない 6 つのデータブロックが存在する。データロケーションテーブル 7 1 0 は、6 つのデータブロックのそれぞれが格納されるロケーションを示し、ここでデータブロックは、装置番号とブロック番号とによって参照され、最初のエントリが最初のデータブロックに対応する。さらに、コンテンツデータのためのパリティ情報のロケーションはパリティマップ 7 4 0 に格納され、そのパリティマップのロケーションは、“パリティマップポインタ”としてのデータ構造の最後のロケーションによって指定される。ファイル USER.TXT は、3 + 1 パリティスキームを用いて格納され、従って 3 つのデータブロックに対しパリティデータのブロックが格納される。この 3 + 1 パリティスキームには 6 つのブロックが存在するので、パリティデータの 2 つのブロックが存在する（6 を 3 で割って、最も近い整数に丸める）。パリティマップは、両方のパリティデータのブロックが格納されたロケーションを示し、ここでパリティデータのブロックは、装置番号とブロック番号とによって参照され、最初のエントリがパリティデータの最初のブロックに対応する。

10

20

【 0 1 4 0 】

B . データロケーションテーブルデータ構造

インテリジェント分散ファイルシステム 1 1 0 は、どのようにファイルが格納されるかについての柔軟性のみならず、多種多様なファイルのための格納を提供する。ファイルデータのミラーリングおよび冗長は、インテリジェント分散ファイルシステムが、異なるファイルについての变化する冗長パラメータを支援することを可能とするファイルシステムレベルで実施される。例えば、いくつかのディレクトリはミラーされ、パリティが保護され、または全く保護されない。

【 0 1 4 1 】

図 8 A , 8 B , 8 C は、变化するプロテクションタイプとレベルのデータファイルのためのデータロケーション情報を格納するために使用されるデータロケーションテーブル例を示す。図 8 A , 8 B , 8 C は、さまざまなデータロケーションテーブルを図解するためのものであり、種々の異なるフォーマット及び / 又は構造が使用されてもよい。

30

【 0 1 4 2 】

図 8 A は、対応するファイルの各データブロックが格納される場所を示すデータロケーションテーブル 8 1 0 の例を図解する。図 7 B におけるように、ファイルについての対応するメタデータは示されていないが、データロケーションテーブル 8 1 0 がメタデータのセットに対応してもよいことが理解される。代表的なデータロケーションテーブル 8 1 0 は、直接エントリおよび間接エントリの両方を備える。

【 0 1 4 3 】

直接エントリは装置 ID / ブロック対を備える。装置 ID は、データが格納されたスマート記憶ユニットを示し、オフセットまたはブロックアドレスは、データが格納された記憶装置上のロケーションを示す。データロケーションテーブルにおける一つのエントリ例は、

40

エントリ	装置	ブロック
1	7	1 2 7

であり、これは、データのブロック 1 がブロック 1 2 7 の装置番号 7 上に格納されていることを示す。

【 0 1 4 4 】

データロケーションテーブル 8 1 0 の例は、また、データロケーションテーブルがより

50

大きなデータのセットのためのデータロケーションを追跡することを可能とする追加のデータロケーションテーブルを指し示す間接エントリを備える。間接エントリのレベルは理論的に制限されないが、そのレベルは、スループットレートを改善するために有利に制限されてもよい。例えば、データロケーションテーブルは、せいぜいダブル間接エントリ、またはせいぜいトリプル間接エントリのみを許容するように制限されてもよい。代表的なデータロケーションテーブル 810 は 2 つのレベルの間接エントリを示す。

【0145】

さらに、データロケーションテーブルの最後のエントリは、パリティマップのアドレス（もしあれば）を格納するために予約されてもよい。他の実施形態において、パリティマップのアドレスは、例えば、メタデータのデータ構造におけるエントリのように、他のロケーションに格納されてもよい。もし、データのセットがパリティプロテクションを備えていなければ、アドレス値は NULL のような標準的な値に設定されてもよい。

10

【0146】

図 8 B は、2 つの追加的なロケーションにおいてミラーされたデータののためのデータロケーションテーブルを示す。このデータロケーションテーブルは、データの各コピーについてのブロック又はオフセットアドレスおよび装置 ID を備える。代表的なデータロケーションテーブルにおいて、ミラーされたロケーションは、ブロック単位ベース (block-by-block basis) で選択されている。例えば、特定のスマート記憶ユニットをミラーするために 1 又は 2 以上の記憶ユニットを選択することのような他のスキームが使用されてもよいことが理解される。図 8 B におけるデータロケーションテーブルは直接エントリのみを備えているが、間接エントリが使用されてもよいことが理解される。

20

【0147】

一実施形態において、ファイルのためのミラーリング情報は、ファイルの対応するメタデータ構造に格納されてもよい。この情報は、例えば、各コピーのためのデータロケーションテーブルのロケーションのみならず、データのコピー数を備えてもよい。データロケーションテーブルは、シングルデータ構造として格納されてもよく、及び/又はデータロケーションテーブルの別個のコピーは異なるロケーションに格納されてもよいことが理解される。

【0148】

ミラーされたデータを有する図 8 B のデータロケーションテーブル例は、パリティプロテクションを備えないが、このデータロケーションテーブルは、パリティ情報を備えてもよいことが理解される。

30

【0149】

図 8 C は、パリティマップを有するデータロケーションテーブルを示す。代表的なデータロケーションテーブルにおいて、データは、3 + 1 パリティスキームを用いて保護され、即ちパリティデータのセットは 3 つのデータブロックの全てから生成されている。例えば、ビット単位 (bit-by-bit)、バイト単位 (byte-by-byte)、ブロック単位 (block-by-block) ベースでデータブロックの XOR を演算してパリティブロックを生成することによるような、データ生成のための技術分野における公知技術が使用されてもよい。

【0150】

代表的なデータロケーションテーブルは、21 個のデータブロック（ブロック 0 からブロック 20）からなるデータファイルに関する情報を提供する。パリティスキームが 3 + 1 であるので、パリティブロックは、3 つのデータブロックの各セットについて生成される。テーブル 2 は、図 8 C に示されるいくつかのパリティブロックといくつかのデータブロックとの対応を示す。

40

【0151】

【表 4】

テーブル 2

データブロック			パリティブロック
0 装置 5 ブロック 100	1 装置 9 ブロック 200	2 装置 7 ブロック 306	0 装置 0 ブロック 001
3 装置 5 ブロック 103	4 装置 9 ブロック 203	5 装置 7 ブロック 303	1 装置 8 ブロック 001

10

【0152】

データロケーションテーブルの例は、パリティマップまたはパリティロケーションテーブルを備える。この代表的なパリティマップにおいて、データを生成するために使用されるブロックエントリのセットとパリティマップの間には1対1のマッピングが存在する。他の実施形態において、パリティマッピングは、また、装置故障によりその直接ロケーション(direct locations)の何れにおいても利用可能でないイベントにおいて、何れのブロックがデータを再生成するために一緒にXORがとられたパリティであるかを、装置およびブロック番号により特定する可変サイズのエントリを備える。他の実施形態において、パリティ生成スキームは、パリティデータの対応(correspondence)およびロケーション

20

【0153】

一実施形態において、パリティマップは、例えば、メタデータのデータ構造に備えられるよりは、むしろ、メタデータのデータ構造の最後のエントリにおけるようなメタデータのデータ構造によって指し示されてもよい。このマップは、故障したスマート記憶ユニット114のまれな場合においてその使用が要求されるだけなので、メタデータ構造に直接的に備えられる代わりに指し示されてもよい。また、スマート記憶ユニット114が、データを再構築しながら一回でパリティマップをトラバース(traverse)し且つトラバースされるときにパリティマップを解析することを可能とするパリティ再結合ブロックを表わすために、パリティマップは可変サイズのエントリを使用してもよい。或る状況では、エントリを取得して解析するための計算およびI/O時間はパリティ計算時間に比較して無視できる。

30

【0154】

パリティロケーション情報を有する図8Cのデータロケーションテーブル810の例は、ミラーリング情報または間接エントリを備えないが、一方または両方がパリティロケーション情報と共に使用されてもよいことが理解される。さらに、他のデータ構造が使用されてもよく、データロケーションテーブルデータ構造は本発明の一実施形態を説明することのみを意図したものであることが理解される。

40

【0155】

C. データ例

図9は、データロケーションテーブル910およびパリティマップ920の例と、そのデータが格納される対応装置を示す。図9の例は、装置上の変化するロケーションにデータがどのように格納されるかを示し、データの“ストライプ(stripe)”が各装置上の異なるオフセットアドレスにわたって格納され、そしてパリティデータが、同一ファイルからのデータについてさえ種々の装置に格納されてもよいことを示す。他の実施形態においては、データは各装置上の同一オフセットアドレスに格納されてもよい。

【0156】

例えば、最初のストライプのためのパリティデータは、ロケーション400で装置3上

50

に格納され、ロケーション 100 で装置 0 に格納されたデータブロック 0 と、ロケーション 200 で装置 1 上に格納されたデータブロック 1 と、ロケーション 300 で装置 2 上に格納されたデータブロック 2 とに関連する。2 番目のストライプのためのパリティデータは、ロケーション 600 で装置 2 上に格納され、そしてロケーション 300 で装置 0 上に格納されたデータブロック 3 と、ロケーション 800 で装置 4 上に格納されたデータブロック 4 と、ロケーション 700 で装置 1 上に格納されたデータブロック 5 とに関連する。

【0157】

いくつかの実施形態において、個々のスマート記憶ユニットは、どこで及び/又はどのようにロケーションをディスク上の実際のロケーションにマッピングするかを決定する。例えば、もし装置 0 が 4 つの物理的なハードディスクを備え、且つ各ハードディスクが 1000 ブロックの記憶容量を備えていれば、装置 0 は、ロケーション 0 からロケーション 399 への格納を可能にする。ロケーションをディスク上のブロックにどのようにマッピングするかを決定するために使用されるガイドラインのセット例は、次のようである。

ディスク番号 = (ディスクごとのブロックの番号 / ロケーション) のフロア

ディスク上のブロック = ロケーション MOD ディスクごとのブロックの番号

【0158】

MOD は、除算の余りをとるモジュラス(modulus)演算子である。上述のガイドラインは、ロケーションをディスク及びディスクブロックにマッピングさせるために使用されるガイドラインの単なる一例を表し、多くの他のガイドラインまたはスキームを使用できることが理解される。例えば、一実施形態は、各ディスクを表すブロック範囲の関連したリスト(linked list)を利用し、そしてリストトラバーサル(list traversal)を実施してもよい。関連したリストは、複数のサイズ化されたディスクを可能にする利点を有する。

【0159】

データ及びパリティ情報の記憶の柔軟性のために、新たなスマート記憶ユニットが付け加えられるに従って、システムを混乱させることなく、既存のデータが新たなスマート記憶ユニットに移動されてもよく(例えば、既存ユニット上のデータを削除する前にコピーを作成することによって)、及び/又は、新たなデータが新たなスマート記憶ユニット上に格納されてもよい。加えて、データブロックまたはファイル全体が、大量の要求、ディスク故障、冗長またはパリティパラメータにおける変更などに応じて実時間でコピー(copy)され又は移動(move)されてもよい。

【0160】

VII. インテリジェント分散ファイルシステム処理

A. データの読み取り

図 10 は、データの読み取り(“データ読み取り処理”)のためのフローチャートの一実施形態を示す。例えばディレクトリメタデータ、ファイルメタデータ、コンテンツデータなどのような種々のデータタイプが読み取られる。

【0161】

スタート状態から始めると、データ読み取り処理は、データが格納されているロケーションを受信する(ステート 1010)。一実施形態において、ロケーションは、スマート記憶ユニット ID およびオフセットまたはブロックアドレスを用いて指定されてもよい。他の実施形態において、記憶装置の ID が使用されてもよいが、他の実施形態において、ID を他の ID にマッピングするためにテーブルが使用されてもよい。等々。

【0162】

次に、データ読み取り処理は、データがローカルに格納されているかどうかを決定する(ステート 1020)。もし、データがローカルに格納されていれば、データ読み取り処理は、そのデータをローカル記憶装置から読み取る(ステート 1030)。一実施形態において、データ読み取り処理は、最初にキャッシュをチェックし、そしてもしデータが存在しなければ、記憶装置をチェックしてもよい。他の実施形態において、データ読み取り処理は記憶装置のみをチェックしてもよい。

【0163】

10

20

30

40

50

もし、データがローカルに格納されていなければ、データ読み取り処理は、データの要求を、そのデータが格納されたスマート記憶ユニットに送信する(ステート1040)。一実施形態において、この要求は、図1に示されるスイッチコンポーネント125を介して送信される。それから、データ読み取り処理は、要求されたデータを受信する(ステート1050)。

【0164】

データ読み取り処理は、要求されているデータを収集し、そしてそのデータを戻す(ステート1060)。いくつかの実施形態において、データは、データのセット全体が収集された後に戻される。他の実施形態において、データの一部またはセットは、そのデータがローカル記憶から読み出され又は他のスマート記憶ユニットから受信されるときに戻される。その一部は、ファイルロケーションテーブルに従う順番(sequential order)で戻されてもよく、またはそれらは読み出され又は受信されるときに戻されてもよい。データが戻された後、データ読み出し処理はエンドステートに進む。

【0165】

図10はデータ読み出し処理の一実施形態を示しており、そして他の実施形態が使用されてもよいことが理解される。他の例において、例えば、並列処理、パイプライン、または非同期I/Oのような技術または技術の組み合わせを用いて、複数のデータ読み出し処理によりデータが並列に読み出されるように2以上のデータ読み出し処理が同時に使用されてもよい。

【0166】

B. ネームルックアップ処理

図11は、ネームルックアップのための処理(“ネームルックアップ処理”)の一実施形態を示す。スタートステートから始めると、ネームルックアップ処理は、ファイルネームを受信し(ステート1110)、ルートディレクトリのメタデータを読み出し、そしてルートメタデータのロケーションをCURRENTとして設定する(ステート1120)。一実施形態において、ルートディレクトリのデータは、図5のデータ構造のようなデータ構造に格納されてもよいが、種々のデータ構造が、ルートディレクトリのメタデータを格納するために使用されてもよいことが理解される。さらにまた、いくつかの実施形態において、ルートディレクトリのメタデータは、各スマート記憶ユニット114がルートディレクトリのメタデータのコピーと同一または類似のコピーを備えるように各スマート記憶ユニット114で格納されてもよい。他の実施形態において、ルートディレクトリのメタデータは、インテリジェント分散ファイルシステム110における他のロケーションに格納され、またはファイル要求と共にスマート記憶ユニット114に送信されてもよい。例えば、ミューテックス(mutexes)及び/又はセマフォ(semaphores)などによってロックするような、複数のデータの複製の完全性を確保するための公知技術が使用されてもよい。

【0167】

そして、ネームルックアップ処理は、ファイルネームの一部である次のトークン(token)を読み取ってもよい(ステート1130)。そして、ネームルックアップ処理は、CURRENTのデータを格納するスマート記憶ユニット114からトークンのメタデータのロケーションのアドレスを要求する(ステート1140)。この要求は、ローカルまたはリモートであってもよい。それから、ネームルックアップ処理は、戻されたアドレスをCURRENTとして設定し(ステート1150)、そして他のトークンが存在しないかどうかを判断してもよい(ステート1160)。ここでトークンはディレクトリ階層における単一のレベルを表す。もし、他のトークンが存在すれば、ネームルックアップ処理はブロック1130に戻る。もし、それ以上トークンが存在しなければ、ネームルックアップ処理は、参照値をCURRENTに戻し(ステート1170)、そしてエンドステートに進む。

【0168】

ネームルックアップ処理の他の実施が使用され得ることが理解される。例えば、ネームルックアップ処理はファイルのメタデータを読み出してよい。加えて、一旦、要求されたデータのロケーションが発見されると、ネームルックアップ処理は、データがローカル

10

20

30

40

50

または他のスマート記憶ユニットで格納されているかどうかを判断してもよい。もし、データがローカルに格納されていれば、ネームルックアップ処理はリード要求をスマート記憶ユニット 114 のローカルブロックマネージャモジュール 335 に送信してもよく、もしデータが他のスマート記憶ユニットに格納されていれば、ネームルックアップ処理は、リード要求をリモートスマート記憶ユニット 114 のリモートブロックマネージャモジュール 337 に送信してもよい。

【0169】

C. ファイル要求処理

図 12 は、ファイル要求のための処理（“ファイル要求処理”）のフローチャートの一実施形態を示す。スタート状態から始めると、ファイル要求処理は、ファイルを読み出すための要求を受信する（状態 1210）。一実施形態において、ファイルは、ロケーションおよびファイルネームを備えるファイルの完全なパスネームを用いて指定される。他の実施形態において、パスは、相対パスであってもよく、及び/又は、ファイルのアドレスについての情報が格納されるテーブルのような他のデータ構造であってもよい。次に、ファイル要求処理は、図 11 に示されるようなネームルックアップ処理を実施して、ファイルのメタデータのデータ構造のロケーションを決定する。

【0170】

そして、ファイル要求処理は、図 10 に示され且つ上述したようなファイル読み出し処理を用いてファイルのメタデータを読み出すが（状態 1230）、他のファイル読み出し処理が使用されてもよい。一実施形態において、ファイルのメタデータは、インテリジェント分散ファイルシステムの全体にわたってファイルにおける各データのブロックが格納されているロケーションへのアクセスを提供するデータロケーションテーブルを備えてもよい。

【0171】

そして、ファイルにおける各データのブロックに対し（状態 1240, 1270）、ファイル要求処理は、ファイルのメタデータにおいてそれを探し出すことによりデータブロックのロケーションを取得し（状態 1250）、そして図 10 に示され且つ上述したようなファイル読み出し処理を用いることにより、そのデータブロックを読み出す（状態 1260）。ただし、他のファイル読み出し処理が使用されてもよい。

【0172】

それから、ファイル要求処理はファイルのデータを戻し（状態 1280）、そしてエンド状態に進む。いくつかの実施形態において、ファイルは、データのセット全体が収集された後に戻される。他の実施形態において、1 又は 2 以上のデータのブロックは、データが読み出されるときに戻されてもよい。一部は、ファイルロケーションテーブルに従う順番で戻されてもよく、またはそれらは、それらが読み出され又は受信されるときに戻されてもよい。一実施形態において、ファイル要求処理は、順番にデータブロックを置いてよく、及び/又はストリーミングサーバのような他のモジュールがデータブロックを順序づけてもよい。データが戻された後に、データ読み出し処理はエンド状態に進む。

【0173】

図 12 はファイル要求処理の一実施形態を示し、そして他の実施形態が使用されてもよいことが理解される。例えば、ファイル要求処理は、図 11 に示されるものとは異なるネームルックアップ処理を用いてファイルのロケーションを決定してもよい。他の例において、データブロックを読み出すために 2 以上のデータ読み出し処理が同時に使用されてもよく、例えば、並列処理、パイプライン、または非同期 I/O のような技術または技術の組み合わせを用いて、複数のデータ読み出し処理により並列にデータを読み出すことを可能にしてもよい。

【0174】

D. パリティ生成処理

図 13 は、パリティ情報の生成（“パリティ生成処理”）のためのフローチャートの一

10

20

30

40

50

実施形態を示す。スタートステートから始めると、パリティ生成処理は、データのセットに関連するパリティスキーム情報を受信する(ステート1310)。データのセットは、ファイルデータ、ファイルメタデータ、ディレクトリメタデータ、ファイルデータのサブセットなどを表してもよい。パリティ生成処理は、データのセットに関連するデータロケーション情報を読み出す(ステート1320)。次に、各パリティのセットに対し(ステート1330, 1370)、パリティ生成処理はデータのセットを読み出す(ステート1340)。例えば、もし、パリティが3+1であれば、パリティ生成処理は、図10に示されるようなデータ読み出し処理を用いて最初の3つのデータブロックを読み出す。次に、パリティ生成処理は、ビット単位、バイト単位、またはブロック単位ベースでデータのXOR演算を実施するようにして、データのセットについてのパリティデータを生成する(ステート1350)。それから、パリティ生成処理は、データをバッファに格納し、そしてデータのセットについてのパリティ情報が生成されるまでブロック1330に戻る。パリティ情報が生成された後に、パリティ生成処理は、そのパリティデータを格納すべき場所を決定する(ステート1380)。パリティ生成処理は、輪番方式のパリティスキーム(rotating parity scheme)を使用してもよい。ここで、ファイルデータの各連続したストライプについての各パリティブロックは輪番で次の装置上に格納される。パリティ生成処理は、パリティ情報がデータ情報と同時に失われることのない装置故障のイベントにおいて確保すべき現在のストライプのためのデータを保持しているどの装置とも異なる装置にパリティブロックを割り当てる。また、パリティ生成処理は、パリティ記憶のために考慮されることからいくつかの装置をふるい落とすために、記憶容量、CPU稼働率、およびネットワーク稼働率のような他の要素を考慮に入れてもよい。そして、パリティ生成処理は、バッファされたデータを割り当てられたスペースに格納し(ステート1390)、パリティマップにパリティデータのロケーションを記録し(ステート1395)、そしてエンドステートに戻る。

【0175】

図13はパリティ生成処理の一実施形態を示し、他の実施形態が使用されてもよいことが理解される。例えば、パリティ生成は、並列にデータのブロックを読み出し、そして並列にパリティ情報を生成してもよく、または公知のパイプラインあるいは非同期I/O技術を用いてもよい。さらに、パリティ生成処理は、一時的なバッファに書き込むことなくパリティ情報およびパリティ情報のロケーションを格納してもよく、または、パリティ生成処理は、パリティデータまたはパリティデータへのポインタを戻してもよい。

【0176】

E. データ回復処理

図14は、欠損または崩壊したデータを回復するための処理(“データ回復処理”)の一実施形態を示す。スタートステートから始めると、データ回復処理は、使用されたパリティスキームに関する情報を受信する(ステート1410)。そして、データ回復処理は、故障または崩壊したディスクまたはデータについての情報を受信する(ステート1420)。次に、データ回復処理は、欠損または崩壊したデータが割り当てられたパリティブロックグループについてのアドレス情報を受信する(ステート1430)。そして、データ回復処理は、利用可能なスマート記憶ユニットからデータブロックを受信する(ステート1440)。データは、図10のもののようなデータ読み取り処理を用いて読み取られてもよい。データ回復処理は、パリティスキームに従ってブロックのXORをとることのような誤り訂正を実施し(ステート1450)、そしてその結果をバッファに格納する(ステート1460)。バッファのデータは欠損したデータ(missing data)を表す。そして、データ回復処理は、バッファにあるデータを戻し(ステート1470)、そしてエンドステートに進む。

【0177】

図14はデータ回復処理の一実施形態を示し、そして他の実施形態が使用されてもよいことが理解される。例えば、データ回復処理は回復されたデータを格納することなく、それを戻してもよい。

10

20

30

40

50

【 0 1 7 8 】

VIII. 分散ファイルシステムにおけるファイルの再ストライピング

いくつかの実施形態において、インテリジェント分散ファイルシステムは、スマート記憶ユニット間に分散されたファイルを再ストライピングするためのシステム及び方法を備える。既に分散されてインテリジェント分散ファイルシステムに格納されたファイルは、ファイルに対するユーザのアクセスを中断することなく且つシステムをオフラインにすることなく、再分散されてシステムに再格納されてもよい。加えて、データは、最小限のデータ移動でスマート記憶ユニット間に再ストライピングされてもよく、再ストライピング処理中にシステム故障が発生したとしても、通常は保護され且つ回復可能である。

【 0 1 7 9 】

再ストライピング処理は、例えば、欠損データが再生成されてシステムに格納されるような或るタイプの故障にスマート記憶ユニットの一つが陥ったときに使用されてもよい。再ストライピング処理は、また、1又は2以上のスマート記憶ユニットが、データが新たなスマート記憶ユニットに付加され又は他のスマート記憶ユニットに再分散されるようなインテリジェント分散ファイルシステムに付加され又はこのインテリジェント分散ファイルシステムから削除されるときに使用されてもよい。加えて、再ストライピング処理は、ファイルのプロテクションスキームが変更されるときに使用されてもよい。例えば、もし、ファイルが3 + 1パリティプロテクションから4 + 1パリティプロテクションに移れば、再ストライピング処理は、新しいパリティプロテクションに合ったレイアウトでスマート記憶ユニットにデータを移動してもよく、新しいレイアウトが完成するまで、古いパリティスキームによってデータが保護されるような古いレイアウトの下でユーザがファイルにアクセスすることを可能とし続ける。

【 0 1 8 0 】

一実施形態において、再ストライピング処理は、ブロック割り当てマネージャによって実施されるが、他の実施形態においては、再ストライピング処理は、インテリジェント分散ファイルシステムの他の部分によって実施されてもよい。

【 0 1 8 1 】

上述の方法論を用いるスマート記憶ユニットは、如何なる特定のストライピングもスマート記憶ユニット内の如何なる特定のロケーションに存在することを必要としないという利点を提供する。従って、“データストライプ”の抽出は、複数のスマート記憶ユニットにわたってブロックの何れの特定のセットにも関与する必要はないが、異なるユニットからの何れかの利用可能なブロックを有利に含む。

【 0 1 8 2 】

A. 再ストライピング処理

図15は、インテリジェント分散ファイルシステム内のデータの再ストライピング(“再ストライピング処理”)のための方法の一実施形態を示す。ファイルは、ファイルの所望のプロテクションスキームに基づき決定されるプロテクショングループのセットとして論理的に表される。例えば、もし、ファイルの所望のプロテクションスキームが3 + 1パリティであれば、ファイルは、4つのクラスタまたはブロック、即ち3つのコンテンツブロックと1つのパリティブロックを有するプロテクショングループに分割される。もし、ファイルの所望のプロテクションスキームが、3 × ミラーされたプロテクションスキームであれば、ファイルは、3つのクラスタ、即ち3つの同一のコンテンツデータブロックを有するプロテクショングループに分割される。

【 0 1 8 3 】

本明細書で述べられる再ストライピング処理は、システム内のデータを移動させるための方法論を有利に表すが、スマート記憶ユニット上の特定のロケーションにデータが存在することを必要としない。再ストライピング処理は、また、再ストライピング処理中のデータ回復を可能とし、且つデータのストライピングにおける実質的な変化を可能にする。加えて、プロテクションスキームの制約が満たされる。例えば、パリティデータおよびその関連するコンテンツデータは、異なるスマート記憶ユニットにそれぞれ格納される。ミ

10

20

30

40

50

ラーされたデータについては、各データのコピーは異なるスマート記憶ユニット上に格納される。

【0184】

－実施形態において、再ストライピング処理は、スマート記憶ユニット間にデータをどのように格納するかを決定するとき他の目的を可能とするためのプリフェレンス(preference)を用いて有利に再ストライピングを行ってもよい。

【0185】

－例において、上記のプリフェレンス(preferance)とは、最小限の移動が、もしブロックが既にスマート記憶ユニットに格納されていれば、上記ブロックが異なるスマート記憶ユニット上にリストアされる必要がないように同一のスマート記憶ユニットがブロックに割り当てられるような優先事項(priority)であるということである。このプリフェレンスは、スマート記憶ユニットが故障した後にファイルを修復するとき使用できる。

【0186】

他のプリフェレンスは、最適なスマート記憶ユニットがレイアウトのために使用されることである。もし、ブロックが既にスマート記憶ユニット上に配置され且つそのスマート記憶ユニットが“最適な”ユニットの一つであれば、同一のスマート記憶ユニットがそのブロックに割り当てられる。従って、ブロックのかなりの移動が回避されるが、ブロックは、スマート記憶ユニットのセットにわたってファイルをバランスさせることが必要となるときに移動され得る。このようなプリフェレンスは、例えば、プロテクションの設定が変更され、または新しいスマート記憶ユニットがインテリジェント分散ファイルシステムに付加されたときに使用され得る。

【0187】

さらにプリフェレンスは、データが既に格納された場所とは無関係に最適なレイアウトのためのものである。従って、再ストライピングは、既存のブロックのロケーションに何ら構うことなく行われる。代わりに、既存のブロックがファイルを最適にレイアウトするために移動されてもよい。しかしながら、いくつかのデータブロックは、適切なスマート記憶ユニットに既に格納されているので移動の必要がないことが理解される。いくつかの実施形態において、インテリジェント分散ファイルシステムは、記憶装置上のフラグメンテーション(fragmentation)を修復するため、最適なスマート記憶ユニット上にたまたま存在していたとしても、ブロックを再配置することを望んでも良い。このプリフェレンスは、例えばリード性能のための最適化又はデフラグメンテーション(defragmentation)のような動作を調整するために使用されてもよく、そして“必要に応じた(as-needed)”基準で、又は通常ファイルシステム動作中のバックグラウンドで使用されるかもしれない。

【0188】

“最適な”レイアウトの目標は、1又は2以上の要素に基づいている。例えば、その要素は、ファイルサイズ、リード性能、ライト性能、予想されるアクセス頻度、システムスループット、ネットワーク速度、利用可能な記憶装置スペースなどを含む。

【0189】

－実施形態において、再ストライピング処理は、もし再ストライピング処理中に1または2以上のスマート記憶ユニットが故障すれば、プロテクションデータを用いてファイルが再生成されることを可能とするステートにおいてファイルを残す(leave)ことを試みる。プリフェレンスができれば使用されるかもしれないが、プリフェレンスが適合しない状況が存在するかもしれないことが理解される。例えば、もしブロックA、B、Cが3×ミラープロテクションスキームの一部であり、ブロックA、Bが装置0上にあり、且つブロックCが装置1上にあれば、たとえ装置0上の残留ブロックBが最小限のデータ移動のプリフェレンスを満足するとしても、プロテクションスキームの制約を満たすために、ブロックAまたはBの何れかが他の利用可能なスマート記憶ユニット上にリストア(restore)されなければならないであろう。このアルゴリズムの実行中に、プロテクションスキームの制約およびプリフェレンスが互いに対峙している周期的なインスタンス(instance)が存在すれば、いくつかの実施形態において、システムは、余分のデータ移動の追加費用でプ

10

20

30

40

50

ロテクションを維持することを選択する。他の実施形態において、1又は2以上のプリフェレンスは、特定の実施によってはプロテクションスキームに優先してもよい。

【0190】

ここで、再ストライピング処理の維持が、図15を参照して述べられる。スタートステートから始めると、再ストライピング処理は次のステートに進み、ここで、各プロテクショングループに対し(ステート1520, 1540)、再ストライピング処理は、スマート記憶ユニットをプロテクショングループのブロックに割り当てる(ステート1530)。割り当て処理の一実施形態は、図16を参照して以下に述べられるが、他の割り当て処理が使用されてもよい。

【0191】

ブロックがスマート記憶ユニットに割り当てられた後、もしブロックが新たなスマート記憶ユニットに割り当てられれば(即ち、それは、割り当てられたスマート記憶ユニットに既に格納されていない)、再ストライピング処理は、ブロックを格納するため、割り当てられたスマート記憶ユニットに要求を送信する(ステート1550)。再ストライピング処理は、プロテクショングループのセットが割り当てられた後に、割り当てられたスマート記憶ユニットに要求を送信し、そして次のステートに進む前にプロテクショングループの全てが割り当てられる必要がないというような他の実施形態が使用されてもよいことが理解される。例えば、もし1又は2以上のスマート記憶ユニットがシステムに加えられるときに再ストライピング処理が使用されれば、または、もしスマート記憶ユニットの一つが故障すれば、再ストライピング処理は、各プロテクショングループが割り当てられた後、割り当てられたスマート記憶ユニットに要求を送信してもよい。

【0192】

次に、再ストライピング処理は、データが無事に格納されたことを確かめる(ステート1560)。代表的な実施形態において、割り当てられたスマート記憶ユニットに既に格納されたデータはこの問い合わせ(query)を満足する。もしデータが無事に格納されていれば、再ストライピング処理は、ファイルに関連するメタデータのデータ構造を更新してエンドステートに進む(ステート1590)。一旦、更新が完了すると、ファイルの新たな情報は、システムによって使用される。先に使用された何れのメモリも解放され、そしてメタデータはファイルの新しいステートを反映する。

【0193】

メタデータの更新は、例えば、コンテンツブロックの新しいアドレスロケーション、プロテクションブロックの新しいアドレスロケーション、新しいプロテクションスキーム、ファイルのデータを格納するために使用されるスマート記憶ユニットの新しいリスト、およびファイルのプロテクション情報を格納するために使用されるスマート記憶ユニットの新しいリストを含む。再ストライピング処理は、また、各スマート記憶ユニットについてファイルデータおよびパリティデータのための新たな“最後のブロックアドレス”を用いてメタデータのデータ構造を更新してもよい。プロテクショングループのセットが格納された後に、データが無事に格納されたことを確かめ(ステート1560)、そして次のステップに進む前にブロックの全てが格納される必要がないというような他の実施形態が使用されてもよいことが理解される。このような実施形態において、再ストライピング処理は、新しいデータのプロテクション値および古いデータのプロテクション値の最小公倍数(the least common multiple)に基づきメタデータのデータ構造を更新する前に無事に格納されなければならないプロテクショングループの数を判断してもよい。このようなインスタンスにおいて、ミラーされたデータについてのプロテクション値は1である(例えば、3xミラーされたプロテクションスキームについて、プロテクション値は1である)。パリティ保護されたデータについてのプロテクション値は、各パリティグループにおけるコンテンツブロックの数である(例えば、4+1パリティスキームの下でのデータについて、プロテクション値は4である)。

【0194】

もし、データが無事に格納されなかったならば(例えば、1又は2以上のスマート記憶

10

20

30

40

50

ユニットがデータを格納できなかつたならば)、再ストライピング処理は、エラーを戻し(ステート1580)、そしてエンドステートに進む。従って、もし1又は2以上のスマート記憶ユニットがデータを格納できなければ、元のデータが維持され、そしてファイルがその元のプロテクションスキームによって依然として保護される。新たに格納された何れのデータも解放される。

【0195】

再ストライピング処理の他の実施形態が使用されてもよいことが理解される。例えば、割り当て処理は、アドレスおよびプロテクション設定に基づき最適なスマート記憶ユニットが推奨され、且つ各メタデータのデータ構造が調べられるようにブロック単位ベースで割り当てる。

10

【0196】

B. 格納処理

一実施形態において、割り当てられたスマート記憶ユニットがデータを格納するための要求を受信したとき、割り当てられた各スマート記憶ユニットは、データが記憶装置上のどこに格納されるかを決定する。一実施形態において、割り当てられたスマート記憶ユニットは、新しいデータを受信したときにスマート記憶ユニットがファイルの関連データの最後のブロックアドレスを用いてその新しいデータを最初の利用可能なロケーションに格納するように、関連データを隣接させて格納することを試みる。しかしながら、種々の記憶のプリフェレンスが使用されてもよく、且つ各スマート記憶ユニットが異なるセットの記憶のプリフェレンスを使用してもよいことが理解される。他の実施形態において、イン

20

【0197】

C. 割り当て処理

図16はデータブロックおよびプロテクションブロックを利用可能なスマート記憶ユニットに割り当てるための方法(“割り当て処理”)の一実施形態を示す。割り当て処理は、再ストライピング中にブロックがどこに割り当てられるべきかを決定する。代表的な実施形態において、割り当ては、単一のプロテクショングループにおいて実施される。従って、割り当ては、データのストライプのために実施され、ここでストライプはデータブロックおよび任意の関連プロテクションデータを含む(例えば、パリティデータ、ミラーされたデータのコピー)。割り当て処理は、最も起こり得ないブロック再配置(relocation)を含むレイアウトを構築してもよく、また一方では、ブロックについて最適なロケーションを決定し、そして任意の指定されたレイアウト目標を満足する。

30

【0198】

スタートステートから始めると、割り当て処理は次のステートに進み、どのスマート記憶ユニットが利用可能であるかを識別する(ステート1620)。一実施形態において、利用可能なユニットは、利用可能な自由スペースによって格納された残りのスマート記憶

40

【0199】

次に、割り当て処理は、割り当てを実施するための利用可能なスマート記憶ユニットが

50

十分存在するかどうかを決定する(ステート1630)。スマート記憶ユニットの好ましい数は、例えば、プロテクショングループにおけるデータブロックの数、使用されているプロテクションスキーム、システムにおけるスマート記憶ユニットの最大数、ファイルの所望のスループットレートなどのような種々の要素に依存してもよい。例えば、3+1パリティプロテクションを有するファイルは、少なくとも4つのスマート記憶ユニットを必要とする。もし、利用可能なスマート記憶ユニットが十分になれば、割り当て処理は、エラーを戻し(ステート1640)、そしてエンドステートに進む(ステート1670)。

【0200】

もし、利用可能なスマート記憶ユニットが十分あれば、割り当て処理は、プロテクショングループにおける各ブロックについてスマート記憶ユニットを選択する(ステート1650)。スマート記憶ユニットを選択するための処理の更なる詳細な説明は、図17を参照して以下に説明されるが、他の処理が使用されてもよい。

【0201】

いくつかの実施形態において、スマート記憶ユニットが選択された後、選択された各スマート記憶ユニットについて、割り当て処理は、現在のファイルに関連するデータブロックの最後のロケーション、及び/又は、現在のデータに関連するプロテクションブロックの最後のロケーションを決定し(ステート1660)、そしてエンドステートに進む。これは、新しいデータブロック及び/又は新しいプロテクションブロックがファイルからの他のデータブロック及びプロテクションブロックの近くに格納されることを可能にする。他の実施形態において、他の記憶のプリフェレンスが使用されてもよいことが理解される。例えば、データは、隣接して又は隣接しないで格納されてもよく、データは、先頭よりはむしろ最後尾に格納されてもよい。等々。

【0202】

D. 選択処理

図17は、単一プロテクショングループのクラスタを格納するスマート記憶ユニットを選択するための方法(“選択処理”)の一実施形態を示す。代表的な実施形態において、選択処理は、装置をファイルのプロテクショングループに割り当てるために使用される。上述のように、プロテクションスキームの制約は、プロテクショングループにおけるブロックが異なるスマート記憶ユニット上にそれぞれ格納されることを要求する。従って、ミラーされたファイル、ブロックの各コピーは異なるスマート記憶ユニット上に格納されるべきであり、且つパリティ保護されたファイルについては、コンテンツデータの各ブロック及びその関連パリティデータは異なるスマート記憶ユニット上に格納されるべきである。代表的な選択処理は、スマート記憶ユニットを単一のプロテクショングループに割り当てるために使用されるが、他の実施形態において、選択処理は、スマート記憶ユニットを、より小さな及び/又はより大きなデータのセットに割り当てるために使用されてもよい。例えば、選択処理は、ブロック上、プロテクショングループのセット上等で、単一ブロックと連動してもよい。

【0203】

1. 実施例

ここでは装置として参照されるスマート記憶ユニットのセット全体、インテリジェント分散ファイルシステムにおいて利用可能な装置は、Tとして表現される。動作している装置のセットは、インテリジェント分散ファイルシステムにおいて稼働している装置の全てのセットであり、Gとして表現される。Gは、或るタイプの故障(例えば、システム故障、ネットワーク故障、記憶装置故障など)により“ダウン”していたインテリジェント分散ファイルシステムにおける装置を除く。加えて、Wは、ファイルのコンテンツデータ及びプロテクションデータが現在格納されている装置のセット(即ち、ファイルによって占められている装置)を表現する。もし、ファイルがインテリジェント分散ファイルシステムに決して格納されることがないとすれば、Wは空(empty)である。

【0204】

ステートステートから始めると、選択処理は次のステートに進み、そして好ましい装置のセットを識別する(ステート1715)。

【0205】

代表的な実施形態において、ファイルが分散される装置の数は、選択されたプロテクションスキームの制約を満たすのに十分に大きくあるべきである。例えば、 $m+n$ のパリティプロテクションを達成するために、少なくとも $m+n$ の装置が必要とされる。 k 回のミラーリングを達成するためには、 k 個の装置が必要とされる。

【0206】

加えて、ファイルが分散される装置の数は、ファイルが損傷を受けやすくしないように十分に小さくあるべきである。特定の装置が故障すればファイルが影響を受ける可能性が増加するため、ファイルは、より多くの装置に拡散されるときに損傷を受けやすくなるかもしれない。従って、システムは、ファイルのデータが占められる装置の最大数 max として上限を選択してもよい。

【0207】

従って、ファイルの幅に関する制約は、次のように表現される。

$m+n$ を用いたパリティ保護: $\{m+n\} \leq \text{ファイルの幅} \leq max$

k 回のミラーリング: $\{k\} \leq \text{ファイルの幅} \leq max$

【0208】

好ましい装置 P のセットを選択するため、システムは、 max を使用して P のサイズを設定する。 G および W の共通集合(intersection)からの装置は P に加えられ、そしてもし P が依然として最大のサイズでなければ、 G 及び W の共通集合からの装置は、 P が最大サイズに到達するまで P に加えられる。 P に加えられない G 及び W の共通集合からの残りの装置は S に加えられる(ステート1720)。

【0209】

従って、もし $|max| = |W \& G|$ であれば、 P はデータが既に存在する UP 装置のみを含み、且つ S は、データが既に存在しない UP 装置のみを含むであろう。もし $|max| < |W \& G|$ であれば、 P は、データが既に存在する UP 装置のいくつかを含み、且つ S は、データが既に存在しない UP 装置のいくつかと同様のデータが既に存在する UP 装置のいくつかを含むであろう。もし $|max| > |W \& G|$ であれば、 P は、データが既に存在する UP 装置と、データが既に存在しない UP 装置のいくつかを含むであろう。 S は、データが既に存在しない UP 装置のいくつかを含むであろう。

【0210】

好ましい装置を選択する他の方法もまた使用されてもよいことが理解される。

【0211】

次に、選択処理は、順序づけられた最適な装置のセットを識別し、それは現在のプロテクショングループについての“最適なプロテクショングループ”である(ステート1725)。最適なプロテクショングループは、 0 として表現されてもよい。代表的な実施形態において、最適なファイルレイアウトは、初回に好ましい装置のセット全体にファイルを書き込む場合に使用されるレイアウトとして定義される。 P における最後の装置が使用されるとき、使用される次の装置は、ブロックが P における装置によって“包み込まれる(wrapped around)”ことを可能にする P における最初の装置である。従って、最初のプロテクショングループの最初のブロックは、 P における最初の装置に格納され、最初のプロテクショングループの2番目のブロックは、 P における2番目の装置に格納され、最初のプロテクショングループにおける3番目のブロックは、 P における3番目の装置に格納され、最初のプロテクショングループにおける残りのブロック等について同様である。そして、2番目のプロテクショングループの最初のブロックは、先のプロテクショングループの最後のブロックが格納された装置に格納される。2番目のプロテクショングループの2番目のブロックは次の装置に格納される。等々。この実施形態は、 P における装置間でプロテクショングループにおける一つのブロックの重複(overlap)を可能とするが、ゼロを含む他の重複サイズが使用されてもよいことが理解される。この実施形態において、“最適

なプロテクショングループ”は、初回に好ましい装置にファイルを書き込む場合に特定のプロテクショングループが格納される装置の順序づけされたセットとして定義され得る。“最適なファイルレイアウト”および“最適なプロテクショングループ”を選択するための他の定義が使用されてもよい。

【0212】

次に、選択処理は非最適装置のセットを識別し、それはNとして表現される(ステート1730)。この実施形態において、Nは、もし初回にファイルが書き込まれればプロテクショングループが格納されない好ましい装置内の装置のセットである。

【0213】

次に、選択処理は、プロテクショングループにおけるブロックのそれぞれが現在格納される装置の配列(array)または行列(matrix)を生成し(ステート1735)、それは、Cとして参照される。Cは、新しいプロテクショングループに存在するブロックと同数の列(column)と、現在格納されるようなファイルによって使用されるミラーの数に対応する各列における行(row)の数とを有し、行の最小数は1とされる。Cにおける各列のエントリは、ブロックが既に格納された種々の装置を表す。例えば、もしファイルが現在3×ファイルとして格納されており、それが5×ファイルに変更されていれば、新しいプロテクショングループはサイズ5であり、そして現在格納されるようなファイルによって使用されるミラーの数は3である。従って、Cは5列を有し、且つ各列は3行を有する。もしファイルが、現在、3+1パリティ保護され、そして5+1パリティプロテクションに変化していれば、Cは6列を有し、そして各列は1行を有する。もしファイルがまだシステムに格納されていないとすれば、Cにおける各エントリが、現在システムにブロックが格納されていないことを表すゼロであるように、Cは新しいプロテクショングループにおける各ブロックのための列を有すると共に1行を有する。ゼロエントリは、ブロックが現在システムに格納されていないことを表し、例えば、パリティが前に使用されていないとすればパリティブロックがシステムにまだ格納されていないように、パリティプロテクションがファイルに加えられているときに使用されてもよい。

【0214】

そして、選択処理は、このファイルについて何れのプリフェレンスが選択されたかを決定する(ステート1740)。代表的な実施形態において、3つのプリフェレンスが存在する。1番目のプリフェレンス(“修復(repair)”)は、データが既に存在するスマート記憶ユニットにブロックを割り当てることによってデータの移動を最小化することを支持する。2番目のプリフェレンス(“再バランス(rebalance)”)は、“最適ナリスト”における装置上に既に存在するブロックがその装置に残るようにブロックを割り当てることを支持し、そして他のブロックが“最適ナリスト”における残りの装置に移動される。3番目のプリフェレンス(“再チューン(retune)”)は、“最適ナリスト”における順序づけられた値に全てのブロックを割り当てることを支持する。プリフェレンスは、選択されたプロテクションスキームに抵触すれば、必ずしも満足されとは限らないことが理解される。

【0215】

上述したプリフェレンスは代表的なものであり、種々のプリフェレンスが使用されてもよいことが理解される。さらにまた、1又は2以上のプリフェレンスがシステムの目標に最も適合するように組み合わせられてもよい。加えて、いくつかの一実施形態は、プリフェレンスを使用しなくてもよいが、装置を選択するためにプロテクションスキームの制約を使用することだけはよい。他の実施形態において、プリフェレンスは、プロテクションスキームの制約に優先してもよい。

【0216】

もし修復(REPAIR)が上記プリフェレンスであれば、選択処理は、Cにおける各列を通してトラバース(traverse)し、列における装置の一つが0における装置の一つと一致するかどうかを調べるためにチェックする(ステート1745)。もし一致すれば、選択処理は、最後の割り当てリスティング(final assignment listing)において装置IDを記録する

10

20

30

40

50

ことにより、一致した装置をそのブロックに割り当て、Oからその装置を削除し、そしてCにおける次の列に移動する。もし一致がなければ、選択処理は、Cにおける次の列に移動する。一旦、Cにおける各列がトラバースされると、選択処理は次のステートに進む。

【0217】

このステートでは、割り当てられていない全てのブロックについて、選択処理がCにおける対応列を通してトラバースし、そして列における装置の一つがNからの装置の一つと一致するかどうかを調べるためにチェックする(ステート1750)。もし一致すれば、選択処理は、最後の割り当てリスティングにおいて装置IDを記録することにより、一致した装置をそのブロックに割り当て、Nからその装置を削除し、そして割り当てられていない次のブロックに移動する。もし一致しなければ、選択処理は、割り当てられていない次のブロックに移動する。一旦、割り当てられていないブロックがトラバースされると、選択処理は次のステートに進む。

10

【0218】

このステートでは、割り当てられていない全てのブロックについて、選択処理は、Cにおける対応列を通してトラバースし、そして列における装置の一つがSからの装置の一つと一致するかどうかを調べるためにチェックする(ステート1755)。もし一致すれば、選択処理は、最後の割り当てリスティングにおいて装置IDを記録することにより、一致した装置をそのブロックに割り当て、Sからその装置を削除し、そして割り当てられていない次のブロックに移動する。もし一致しなければ、選択処理は、割り当てられていない次のブロックに移動する。一旦、割り当てられていない全てのブロックがトラバースされると、選択処理は次のステートに移動する。

20

【0219】

このステートでは、まだ割り当てられていない全てのブロックについて、選択処理はOからの装置を割り当て(ステート1760)。一旦、各ブロックが割り当てられると、選択処理はエンドステートに進む。

【0220】

もし再バランス(REBALANCE)がプリフェレンスであれば、選択処理は、Cにおける各列を通してトラバースし、そして列における装置の一つがOにおける装置の一つと一致するかどうかを調べるためにチェックする(ステート17650)。もし一致すれば、選択処理は、最後の割り当てリスティングにおいて装置IDを記録することにより、一致した装置をそのブロックに割り当て、Oからその装置を削除し、そしてCにおける次の列に移動する。もし一致しなければ、選択処理は、Cにおける次の列に移動する。一旦、Cにおける各列がトラバースされると、選択処理は次のステートに進む。

30

【0221】

このステートでは、まだ割り当てられていない全てのブロックについて、選択処理は、Oからの装置を割り当て(ステート1770)、そして最後の割り当てリスティングにおいてその割り当てを記録する。一旦、各ブロックが割り当てられると、選択処理は、エンドステートに進む。

【0222】

もし再チューン(RETUNE)がプリフェレンスであれば、選択処理は、最初のブロックから始めて、Oにおける最初の装置を割り当て、2番目のブロックは、Oにおける2番目の装置に割り当てられ、装置は順序づけられたOリストを用いて割り当てられる。等々。一旦、各ブロックが割り当てられると、選択処理はエンドステートに進む。

40

【0223】

従って、選択処理は、プロテクショングループにおける各ブロックについて装置を選択する。

【0224】

図17は、選択処理の一実施形態を示し、他の実施形態が使用されてもよいことが理解される。例えば、ファイルを最初に書き出すために使用される方法は異なってもよい。上

50

述したように、代表的な実施形態は、好ましい装置間にブロックを分散させるために、スキュー(skew)、および一つのブロックによるプロテクショングループ重複(Protection Group Overlap)のようなものを使用する。加えて、システムはスキュー値の使用を排除できる。このようなケースでは、各プロテクショングループは、プロテクショングループへのユニットの割り当ての後またはその度に、同一のユニットを用いてそのデータのその割り当てを開始することができ、好ましいユニットのリストは、種々のパラメータに基づき再順序づけされてもよい。さらにまた、他の実施形態において、選択処理は、Cにおける次の列に移動する前に、セットO、N、Sのそれぞれを通してトラバースしてもよい。

【0225】

加えて、好ましい装置及び最適な装置のセットが、追加の及び/又は他の基準を用いて選択されてもよい。このような基準は、例えば、どのスマート記憶ユニットが動作しているか、スマート記憶ユニットの記憶装置上の利用可能なスペースの量、スマート記憶ユニットのスループットレート、プロテクショングループからのデータが既に格納されたスマート記憶ユニットなどを含んでもよい。異なるシステムが異なる目標を有してもよいことが理解される。

10

【0226】

2. 基礎の概要

以下では、上述した数学的構成のいくつかの概要を提供する。

【0227】

a. オペレータ

以下のものは、図17に関して説明される代表的な実施形態のための基礎を提供するために使用される演算子である。

20

& = 交差部(Intersection)

| = ユニオン(Union)

! = 反対(Inverse)

- = 減算(Subtraction)

【0228】

b. セット/リスト

以下のものは、図17に関して説明される代表的な実施形態についての基礎を提供するために使用されるセット及び順序づけされたリストである。

30

T = インテリジェント分散ファイルシステムにおける利用可能な装置のセット

G = インテリジェント分散ファイルシステムにおける動作中の装置の全てのセット

W = ファイルによって占められた装置のセット

S = 予備装置のセット

P = 好ましい装置の順序づけされたセット

O[size] = プロテクショングループにおける各ブロックへの装置の最適な割り当てであり、“size”はプロテクショングループのサイズである。

N = 非最適セット = P - O

F[size] = プロテクショングループにおける各ブロックについての装置の最後の割り当て

40

C[size,k] = 各ブロックに割り当てられた現在の装置であり、“k”は現在格納されるようなファイルによって使用されるミラーの数である。

【0229】

3. 例

装置の選択例が説明される。しかしながら、この例は、本発明の範囲を制限することを意図したものではなく、種々の実施形態についての詳細を提供するものに過ぎない。

【0230】

a. 修復(Repair)

我々は、ちょうどノード3上の単一の装置故障にであった4ノード配列上の3xファイル修復することを試みると仮定する。

50


```

W == (1,2,3)
G == (1,2,4)
dev(X) 3 ; widthpolicy(x) == 8
R == Repair
group_index = 0

```

【 0 2 3 1 】

ここで、 $|G| \geq dev(X)$ とすれば、

```
P = G & W == (1,2)
```

```
S = G & !W == (4)
```

を得る。

10

【 0 2 3 2 】

ここで、我々は、S から引いて P に追加すると、

```
P = (1,2,4)
```

```
S = ( )
```

となる。

【 0 2 3 3 】

cycle_offset は、 $group_index \% |P|$ 、又は $0 \% 3 = 0$ である。これから、

```
O = (1,2,4)
```

```
N = ( )
```

を得る。

20

【 0 2 3 4 】

F における全ての値にわたって繰り返す。

【 0 2 3 5 】

F[0] について、C[0] & O は (1 , 2) である。これら二つのミラーは可能な選択だからである。装置 1 および 2 の両方は F[0] のミラーを備える。F[0] = 1 に選択する。O から 1 を削除すると、O は (2 , 4) になる。

【 0 2 3 6 】

さて、F[1] について、C[1] & O は (2) である。F[1] = 2 に選択する。O から 2 を削除する。O は (4) となる。

30

【 0 2 3 7 】

F[2] について、C[2] & O は空(empty)のセットである。我々は、最後のステップに進み、O の残りの値を F[2] に割り当てる。そこで F[2] = 4 となる。最終的な F は、

```
F = (1,2,4)
```

となる。

【 0 2 3 8 】

修復の目標が達成された。一つのクラスタのみが再構成されなければならない。他の二つのクラスタは残される。

【 0 2 3 9 】

b . 再バランス(Rebalance)

3 ノード配列上に 2 + 1 ファイルが与えられると、再バランス目標は、4 番目のノードが配列に加えられるときに 3 + 1 にファイルを再ストライピングするために使用されてもよい。

40

```
W == (1,2,3)
```

```
G == (1,2,3,4)
```

```
dev(X) == 4; widthpolicy(X) == 8
```

```
R == Rebalance
```

```
group_index = 30
```

【 0 2 4 0 】

ここで $|G| \geq dev(X)$ であり、

50

P = (1,2,3)

S = (4)

を得る。

【 0 2 4 1 】

S から引いて P に追加すると、

P = (1,2,3,4)

S = ()

を得る。

【 0 2 4 2 】

group_index は 30 であり、cycle_offset は $30 \% 4 = 2$ である、従って、

O = (3,4,1,2)

N = ()

である。

【 0 2 4 3 】

ファイルはミラーされず、従って各ソークラスタのただ一つのデータコピーしか存在しないので、全ての C[i] はせいぜい単一要素を持つことに注意されたい。この例のために、つぎの事柄を仮定する。

C = ((3), (4), (2), (3))

従って、

C[0] = (3), C[1] = (4), 等

【 0 2 4 4 】

我々は、F[0] から始める。このケースでは C[i] & O は 3 を提供し、従って F[0] = 3 である。同様に、F[1] = 4 であり、且つ F[2] = 2 である。各ステップで、我々は、O が (1) のみを含んで残されるように、O から削除する。

【 0 2 4 5 】

F[3] について、C[3] & O は空(empty)である。我々は、O の最後のメンバーを F[3] に割り当てる最終ステップに進み、そして従って F[3] = 1 である。最終的な F は、

F = (3,4,2,1)

である。

【 0 2 4 6 】

ここで、二つのブロックが移動され、そして F の全てのメンバーは O にある。

【 0 2 4 7 】

c . 再チューン(Retune)

代表的なインテリジェント分散ファイルシステムは、8つの装置 A, B, C, D, E, F, G, H を備え、そしてファイルは、12個のコンテンツブロック b1, b2, b3, b4, b5, b6, b7, b8, b9, b10, b11, b12 から構成される。現在、ファイルは装置 A, B, C, D に格納されており、そして2回ミラーされているが、パリティプロテクションは有していない。以下に、ファイルの各ブロックが現在格納されている装置を示す。

【 0 2 4 8 】

【表 5】

b1	b2	b3	b4	b5	b6	b7	b8	b9	b10	b11	b12
A	B	C	D	A	B	C	D	A	B	C	D
B	C	D	A	B	C	D	A	B	C	D	A

【 0 2 4 9 】

もし新しいデータレイアウトが 3 + 1 パリティプロテクションを備えるが、ミラーされたデータを備えていなければ、ファイルは4つのプロテクショングループを有し、それぞ

10

20

30

40

50

れは、4つのブロック、即ち3つのコンテンツブロックと1つのパリティブロックを有する。ファイルが分散される装置の最大数 max は4に設定される。

【0250】

従って、

$T = \{A, B, C, D, E, F, G, H\}$

$W = \{A, B, C, D\}$

である。

【0251】

次に、装置 F, G, H がダウンしたとする。このことは、

$G = \{A, B, C, D, E\}$

を意味する。

【0252】

P を構築するため、セット W および G の共通集合からの装置が P に加えられる。もし P がまだ装置の " max " 数を有していなければ、セット G および W の共通集合からの装置が、" max " 数に到達するまで P に加えられる。セット G および W の共通集合における残りの装置は S に配置される。一実施形態において、 G および W の共通集合における装置は、前部 (front) により多くの自由スペースを持つもので順序づけされる。

$P = \{A, B, C, D\}$

$S = \{F\}$

【0253】

上述したスキューを用いて、最適なファイルレイアウトは、次のようになる。

【0254】

【表6】

	A	B	C	D	E
パリティグループ1	b1	b2	b3	p1	--
パリティグループ2	b6	p2	--	b4	b5
パリティグループ3	--	b7	b8	b9	p3
パリティグループ4	b11	b12	p4	--	b10

【0255】

従って、パリティグループ1:

$O = \{A, B, C, D\}$

$N = \{E\}$

$C = \{\{A, B\}, \{B, C\}, \{C, D\}, \{0, 0\}\}$

【0256】

パリティグループ2について:

$O = \{D, E, A, B\}$

$N = \{C\}$

$C = \{\{D, A\}, \{A, B\}, \{B, C\}, \{0, 0\}\}$

【0257】

パリティグループ3について:

$O = \{B, C, D, E\}$

$N = \{A\}$

$C = \{\{C, D\}, \{D, A\}, \{A, B\}, \{0, 0\}\}$

【0258】

パリティグループ4について:

$O = \{E, A, B, C\}$

$N = \{D\}$

$C = \{\{B,C\}, \{C,D\}, \{D,A\}, \{0,0\}\}$

【 0 2 5 9 】

ここで、パリティグループ 1 は、リターン(RETURN)プリフェレンスを用いて装置を選択するための例として使用される。

$O = \{A,B,C,D\}$

を呼び出す。

【 0 2 6 0 】

従って、Fにおける最初のブロックはAに割り当てられ、Fにおける2番目のブロックはBに割り当てられ、Fにおける3番目のブロックはCに割り当てられ、そしてFにおける4番目のブロックはDに割り当てられる。

10

【 0 2 6 1 】

選択処理は、プロテクショングループのそれぞれに呼び出され、つぎの割り当てが得られる。

$F_1 = \{A,B,C,D\}$

$F_2 = \{D,E,A,B\}$

$F_3 = \{B,C,D,E\}$

$F_4 = \{E,A,B,C\}$

【 0 2 6 2 】

従って、ブロックは、ファイルがあれば最初に書き込まれたであろう場所と同じロケーション - 上述した最適なファイルレイアウトと同様のレイアウト - に格納される。

20

【 0 2 6 3 】

4 . 擬似コード実例 (Sample Pseudocode)

以下に、選択処理の実施例についての擬似コードの例を提供する。しかしながら、この例は、本発明の範囲を制限することを意図したのではなく、特定の実施形態のための詳細を提供するものに過ぎない。

【 0 2 6 4 】

[著作権 2 0 0 3 アイシロン・システムズ・インコーポレーテッド]

プロテクショングループについての装置選択の実施

30

‘ & ’ は共通集合を取ることを(intersect)を意味し、‘ | ’ はユニオン(union)を意味し、‘ ! ’ は反転(inverse)を意味する。

本アルゴリズムは、プロテクショングループについての装置選択のセットであるFを提供することを試みる。Fは配列であり、F[i]のようにインデックスが付けられる。F E C保護された(パリティ保護された)ファイルについて、Fは、ファイルの全体を通してオフセット順にプロテクショングループにおけるクラスタのセットのためのデヴィッド(devid)選択を記述する。ミラーされたファイルについては、Fは、プロテクショングループにおける全てのミラーのためのデヴィッド(devid)選択を記述する。

本アルゴリズムは多くの入力を考慮する：

group_indexは、プロテクショングループのインデックスである。どのファイルも多く、互いに素(disjoint)のプロテクショングループに分割され、インデックスが割り当てられる。

40

Gは、全ての‘アップ(up)’の装置のセットであり、即ち利用可能であり且つ故障状態に陥っていない装置である。

Wは、既にファイルによって占められている装置のセットである。それは、順序づけられたリスト(an ordered list)である。

Xは、ファイルがどのようにレイアウトされるかを表すポリシーである。どのXについても、Xによってファイルを保護するために必要とされる装置の数を示す値dev(X)が存在する。これは、また、プロテクショングループにおけるエレメントの数であり、従って、 $|F|=dev(X)$ である。Xは、また、ファイルの所望の幅widthpolicy(X)を提供し、それは、

50

もし $|G| \geq \text{widthpolicy}(X)$ であれば、理想条件下でファイルが占めるべき装置の数である。

C は、既存のレイアウトを表す。それは、もしあれば、各 $F[i]$ に必要とされるデータのコピーが発見される装置のセットを各 $F[i]$ について記述する。このような装置のセットを $C[i]$ とする。 $|C| == |F|$ である。

R は、本アルゴリズムがめざす目的物(goal)である。それは、アルゴリズムが、次のようであるかどうかを表す。

R == 修復(repair) : F を取得すると、ブロック移動を最小化することを試みる。

R == 再バランス(rebalance) : F を取得しながら $|W| == \text{widthpolicy}(X)$ を達成することを試みる。

1 . G が、サブミットされた X に適合するかを検証することから始める。もし、 $|G| < \text{dev}(X)$ であれば、アルゴリズムは継続できず、エラーを発呼側(caller)に戻し、別の X を用いて再試行することを可能にする。

我々は、好ましい装置 P の順序づけされた配列を確立することを望み、それから F の選択がなされる。次のように、順序づけされた W のサブセットとして P を確立する。

P == G & W

最初に、自由スペースの量によって順序付けられた、W にない G の全てのメンバーのように、予備配列 S をセットアップする。

S == G & !W

P および S は、両方とも順序づけされる。 $|P| == \text{wp}$ になるまで、S からメンバーを P に追加する :

```
while(|P| < widthpolicy(X) and |S| > 0) {
    x=S.pop() // take from front
    P.append(x) // add to end
}
```

$|P|$ は、ファイルによる占有のための装置の理想的なリストを表す。それは、G が与えられると、 $\text{widthpolicy}(X)$ に近いサイズとされるべきであるが、もちろん $|G| < \text{widthpolicy}(X)$ である。

2 . ミラー及びパリティのローテーション(rotation)を可能にするためにプロテクショングループインデックスに基づきサイクルオフセットを計算する。P におけるサイクルオフセットから初めて、装置をカウントし、そしてそれらを最適なセット 0 に加える :

$\text{cycle_offset} = \text{group_index} \% |P|$

for (i=0; i<dev(X); i++)

0.append(P[(cycle_offset + i) % |P|])

N を、あまり最適でない装置(the less-than-optimal device)として定義する :

$N = P - 0$

注記 :

-N, 0, S は互いに素(disjoint)である

-N|0 == P

3 . 既存の現在の装置にわたって繰り返し、R が達成されるように各 $F[i]$ を選択することを試みる。R == Repair について、我々は、ブロック移動を最小化したい。R == Rebalance について、我々は、F を 0 に近づけたい。

R == Repair について、選択セットの順序づけされた配列(array)を次のように定義する :

A = (0, N, S)

R == Rebalance について、A を次のように簡単に定義する :

A = 0

R による A の適切な選択の後、我々は、次のように F を完結する :

for SET in A {

for (i=0; i<|F|; i++) {

10

20

30

40

50

```

x = C[i] & SET // 既存の装置と所望のものとの共通集合をとる
if (|x|>0) {
    y = x.first() // これらのうちの一つを取得
    SET = SET - y // 再考されることを防止
    F[i] = y // Fに入れる
}
}
}

```

4 . 我々は、依然として、C [i]からの適切な選択が存在しない任意のF [i]についてFにおけるデヴィッド(devids)を選択しなければならない。

10

```

for (i=0; i<|F|; i++) {
    if (F[i] unassigned) {
        F = 0.pop();
    }
}

```

【 0 2 6 5 】

E . 擬似コードの例(Example Pseudocode)

以下では、再ストライピング処理の実施例についての擬似コードの例を提供する。しかしながら、この例は、本発明の範囲を制限することを意図したのではなく、一実施形態の詳細を提供するものに過ぎない。

20

【 0 2 6 6 】

< 著作権 2 0 0 3 アイシロン・システムズ・インコーポレーテッド >

このセクションでは、(BAMにおいて)どのようにブロックが装置レベルで割り当てられるかを決定するファイルシステムのコンポーネントを述べる。ディスクおよびブロックレベルでの割り当ては、LBM割り当てコードによって処理される。以下に述べられるモジュールは、書き込み及び読み出し処理の間に新しいブロックが割り当てられるべき場所を見つけ出すために、主として、BAMの書き込み及び読み出しコンポーネントによって使用されるであろう。しかしながら、APIは、通常、必要に応じて他のコンポーネントによって十分に使用され得る。

ファンクションはbam_layout.cに存在し、そして最初は新たな書き込み及び読み出しコードによって使用される。DFMのような他のモジュールは、必要に応じて将来使用するために修正されてもよい：

30

```

int bam_layout_protection_group
(
    const struct gmp_group_info *gi,
    long *free_cluster_counts,
    struct inode *ip,
    enum layout_goal_t goal,
    struct protection_level *protection,
    int width_policy,
    ifs_lbn_t start_lbn,
    int num_clusters,
    int curr_devid_depth,
    ifs_devid_t **curr_devids,
    u_int8_t *alloc_counts,
    ifs_devid_t *result_devids,
    ifs_baddr_t *result_pbas);

```

40

< 概要 >

このファンクションは、特定のプロテクション設定下で単一のプロテクショングループのためのレイアウト情報を計算する。このファンクションは、(上記特定のプロテクション設定に基づき)完全なプロテクショングループのために呼び出される。レイアウトリコ

50

メンテーション(Layout recommendation)は、グループにおける各クラスタのための単一の装置idの形式で、アウトパラメータ(result_devids)で戻されるであろう。他のアウトパラメータ(result_pbas)は、割り当てにおいて使用される各クラスタについての前のブロックアドレスを備えるであろう。

<パラメータ>

Gi : 現在のクラスタグループ情報。

free_cluster_counts : グループにおける各ノードについてのフリークラスタカウント(free cluster count)。これは、その現在の自由スペースに基づき、どの装置を使用するかをレイアウトエンジンが決定することを可能にする。この配列は、まさに、長さが 'ds et_size(&gi->group)' であるべきである。

10

lp : inodeである。これは、主として、width_device情報と前のブロックアドレス情報とをアクセスするために使用される。

Goal : このレイアウト動作の目標を特定するフラグ。このフラグは、ブロックのレイアウトの仕方を決定する時を参照するための明確な目標(objective)をレイアウトエンジンに与える。現在3つの目標(goal)が存在する :

LAYOUT_REPAIR : これは、レイアウトエンジンに、最小のブロック移動が最優先であることを知らせる。レイアウトエンジンは、新たなデヴィッドを、我々のプロテクション制約(ミラーは異なる装置上にあるべきである、など)を維持する必要がある場所のみ割り当てる。既存のブロックは、絶対必要なときにのみ移動される。これは、装置またはノード故障の後にファイルを修復するときで使用され得る目標である。割り当てられないクラスタが、このシナリオの下に最適に依然として完全に配置されているので、これは、書き込みパス(bam_new_write)によって使用され得る目標である。それは、少なくとも、他の二つの目標と同様に全てのシナリオにおいて高速であり、多くの場合、'とても(much)' 高速である。

20

LAYOUT_REBALANCE : これは、レイアウトエンジンに、最適な装置がレイアウトのために使用されるべきであることを知らせる。しかしながら、それらの装置は、既存のブロックの移動を避けるためにシャッフルされてもよい。この方法では、システムは、ブロックの移動を極力避けるが、正確な数のノードにわたってファイルをバランスする必要があるときにブロックを移動するであろう。多くの場合、これは、LAYOUT_REPAIRと同様に高速であるが、しかしながら、より最適且つより良好なバランスされたファイルレイアウトをもたらす。この目標は、新たなノードの付加によりプロテクション設定または再バランスングを変更するときに、再ストライパー(restriper)によって使用され得る。

30

LAYOUT_RETUNE - これは、レイアウトエンジンに、最適なレイアウトが最優先であることを知らせる。既存のブロックには何の注意も払われず、そして任意の又は既存のブロックは、ファイルを完全にレイアウトするために移動されてもよい。もし我々が、相当に良好なレイアウトファイルを手がければ、ブロックは、適切な装置上に図らずも既に存在するので、移動される必要がないが、しかし、これを成し遂げるための試みはなされていない。実際、発呼側は、たとえ、ディスクにおけるフラグメンテーションを修復するために図らずも最適な装置に位置していても、ブロックを再割り当てすることを望んでも良い。この目標は、大部分のシナリオにおいて明らかに最も遅く、且つ決して、LAYOUT_REPAIRまたはLAYOUT_REBALANCEよりも速くはない。発呼側は、この目標が、リード性能についての最適化またはデフラグメンテーションのような'チューニング'動作であることを選択し得る。それらの処理は、通常のファイルシステム動作中にバックグラウンドまたは必要に応じて稼動してもよいであろう。

40

Protection : 所望のプロテクションの設定。これは、inodeにおけるプロテクション設定と一致しない(例えば、もし発呼側が2xでのパリティファイルを考慮していれば、パリティは可能ではない)。レイアウトコードは、ブロックを置く場所を決定するときに、このプロテクション設定を使用する。それは、inode自体におけるプロテクション設定を無視するであろう。

width_policy : ファイルについての所望の幅ポリシー。これは、ファイル全体が配置さ

50

れるべきノードの目標数である。これは、多くの状況において達成できなくてもよいが、しかしながら、レイアウトエンジンは、このポリシーを満足するように努める。

start_lbn : 領域(region)における最初のlbn。このファンクションは操作すべきプロテクショングループの全数を必要とするので、これはプロテクショングループの始まりにある。Start_lbnは、スキューを計算するための容易な手法を提供し、クラスタの周辺のパリティ情報およびデータのローテーションを可能にする。これらのlbnは、また、結果として得られた構造における特定のブロックを識別するために使用される。

num_clusters : 配置されるべきプロテクショングループにおけるクラスタの全数。例えば、3 + 1プロテクション設定においては、各プロテクショングループは、4つのクラスタを備える一方、3 x プロテクショングループは3つのクラスタを備える。Num_clustersは、まさにプロテクショングループと一致する。もしそうでなければ、EINVALが戻される。num_clusterパラメータは、3つの配列パラメータ、即ちcurr_devids, alloc_counts, result_devidsの長さを規定することに注意されたい。

curr_devids : 各クラスタにおいて最初に割り当てられたブロックのデヴィッド(devid)。これは、まさに長さが ' num_clusters ' であるべき2次元配列である。第2の次元は、発呼側が、配列されるべき各クラスタについて代替の既存のロケーションを提供することを可能にする。このパラメータの更なる詳細な説明については、本文書の ' Curr_devids ' のセクションを参照されたい。この配列の多くのエレメントは0に設定されてもよく、このことは、その特定のクラスタには現在ブロックが割り当てられていないことを示す。もしレイアウト制約がそのように行うことを妨げないのであれば、レイアウトエンジンは、既存のブロックが割り当てられた場所に新しいブロックを置くことを試みる。他の配列パラメータと同様にパリティクラスタがこの最後にリストされるべきであることに注意されたい。

alloc_counts : 各クラスタについて、発呼側が割り当てようとするブロックの数。これは、レイアウトエンジンが、完全に割り当てられたクラスタをそのままにしておき、且つ、部分的に割り当てられているが現在の装置上で完結されていないクラスタを移動することを可能にする(なぜなら、その装置は、自由スペースがなくなるので)。このパラメータは、1または2以上の装置が容量に近づいているときに参照されるのみである。このパラメータは、まさに、長さが ' num_clusters ' である。

result_devids : このアウトパラメータは、発呼側によって割り当てられ、まさに長さが ' num_clusters ' である。それは、プロテクショングループにおける各クラスタについて推奨されたデヴィッド(devid)を備える。その上、将来において、レイアウトエンジンは推奨された特定のドライブに拡大されてもよい。発呼側は、各クラスタについて推奨されたデヴィッドを調べる。もしそのデヴィッドが、そのブロックに既に割り当てられたブロックのデヴィッドに等しくなければ、これらの既存のブロックは、レイアウト制約を満たすために移動される。新たなブロックが、推奨されたデヴィッドに割り当てられる。再び、任意のパリティクラスタは、これと他の配列パラメータの最後にリストされ、そして、ミラーは昇順にリストされる。

result_pbas : このアウトパラメータは、発呼側によって割り当てられ、そして長さが ' num_clusters ' のサイズである。それは、各クラスタについての前のブロックアドレス (p b a) を備え、最適な連続した割り当てについてのブロック割り当て要求と共に送信される。もしそのノード上に既に割り当てられたパリティブロック又は前のデータが存在しなければ、これら p b a のいくつかはゼロであってもよい。他の配列パラメータのように、パリティクラスタはこの配列において最下位(last)であり、ミラーは昇順にリストされる。

<Curr_devids フォーマット>

上述したcurr_devidsの説明は、発呼側が、プロテクショングループにおいて現在割り当てられたブロックのロケーションについてのレイアウトエンジン情報に伝えることを可能にする。レイアウトエンジンは、この情報を参照して、要求された ' 目標 ' を依然として満足しながら、可能な限り少ないブロックを移動するように試みる。

他の配列パラメータのように、curr_devids配列は、長さが 'num_clusters' のサイズ (まさに1つのプロテクショングループのサイズ) である。パリティの場合には、パリティクラスタは、プロテクショングループにおいて最下位にリストされ、そしてミラーは、ミラーインデックスの昇順にリストされる。

このパラメータが他の配列と異なる点は、これが2次元配列であることである。この配列は、新たなプロテクショングループにおいて各クラスタについて1列(column)のデバイスを含むであろう(それは、長さが 'num_clusters' である)。各列は、既存のデータブロックのミラーを備える(従って、この配列の行の数は、古いレイアウトのデータミラーカウントにほとんど常に等しい)。このようにして、発呼側は、レイアウトエンジンに、既存のクラスタが装置 A, B, C 上にミラーを有することを知らせることができる。そして、レイアウトエンジンは、もしそれが、そのクラスタについて、A、BまたはCを選択すれば、発呼側が最近割り当てられたどのブロックも移動しないかもしれないことを知る。

10

少しの例がこのことを更に明確にするであろう。システムが3xから3+1にファイルを再ストライピングすると仮定する。システムは、レイアウトエンジンに、1つの3+1プロテクショングループ(ファイルにおける最初の3つのクラスタを含む)を配置することを要求する。システムは、3行4列の配列を割り当てることにより、curr_devidsを構成する。4列については、新たなプロテクショングループが4つのクラスタを有するためであり、3行については、各既存のブロックが3つのミラーを有するためである：

```
curr_devids[4][3] :
    [1][2][3][0]
    [2][3][4][0]
    [3][4][5][0]
```

20

4番目の列は、パリティブロックが現在のところ割り当てられていないので、空(empty)であるが、各データクラスタは3つの既存のミラーを有している。これは、レイアウトエンジンに、それが現在のファイルレイアウトについて知る必要のあることの全てを知らせる。

他の例。システムが3+1プロテクションから3xプロテクションでファイルを再ストライピングしていると仮定する。システムは、長さが3クラスタ(我々の新たな3xプロテクショングループが3のクラスタを有するため)であり、且つ深さが1クラスタ(古いデータはミラーされていないため)であるcurr_devidsリストを作成する：

30

```
curr_devids[3][1] :
    [1][2][3]
```

このことは、現在のクラスタがミラーされれば余剰の行が必要とされるだけであり、且つ、システムが、既存のミラーの場所をレイアウトエンジンに知らせることを欲することを示している。前述したように、2番目の次元は、ほとんど常に、ファイルの古いデータミラーカウント(old data mirror count)に設定されるであろう。もし、古いファイルがパリティプロテクションであれば、これは1であろう。

この構造の実際の能力(power)は、あるミラー設定から他の設定に再ストライピングするときに現れる。例えば、システムが5xから2xに再ストライピングしていることを想定する。システムは、2の長さ(我々の新たなプロテクショングループは2クラスタをそなえるため)を有すると共に5の深さ(各既存のデータクラスタは5ミラーを備えるため)を有するcurr_devids配列を生成するであろう。レイアウトエンジンは、既存のミラーのロケーションを知る必要があり、従ってそれはベスト2を選択して保存する：

40

```
curr_devids[2][5] :
    [1][1]
    [2][2]
    [3][3]
    [4][4]
    [5][5]
```

50

このケースでは、情報のいくつかは冗長である。なぜなら、システムは、1 データクラスタについて語るだけであるから（それは、5 回ミラーされる）。しかしながら、この構造は、レイアウトエンジンが容易に理解する。エンジンが気をつけることの全ては、それが最初のミラーについて最初の列の任意のメンバーを選択できることであり、そして、結果として得られるレイアウトは全くデータブロックの移動を必要としないであろう。このフォーマットは、レイアウトエンジンに、それが効率的なレイアウトの決定をすることを必要としているという情報を与える。

大体的場合、通常のファイルを書き込むとき、この配列は極めてシンプルである。余分な複雑性は、再ストライピング、または装置が利用不能のときに回復された書き込みを実施するときに必要なとされるのみである。例えば、もし3番目のクラスタを3 + 1ファイル

10

```
curr_devids[4][1]:
    [1][2][0][4]
```

同様に、3 x ファイルの3番目のクラスタを書き込む場合、配列は次のように単純にできる：

```
curr_devids[3][1]:
    [1][2][0]
```

これらのケースでは、既存のプロテクション設定は、まさに新たな設定と同一であり、従ってシステムは、フラットな配列で既存のレイアウト情報を伝えることができる。

<リターン値>

20

0：成功(success)。アウトパラメータは、割り当て又は移動を必要とするブロックについての装置IDを含み、且つ移動の必要のないブロックについての結果は含まない。割り当てを必要とするクラスタは、どれも、実際の割り当てコールで用いる関連の前のブロックアドレス(pba)を有するであろう。

EROF S：レイアウト制約に違反することなく、要求されたブロックを配置するための利用可能な装置が十分に存在しない。発呼側は、このエラーをその発呼側に戻すか、または他の所望のプロテクション設定を用いてbam_layout_protection_group()を呼び出すことの何れかを行うことができる。

E I N V A L：無効パラメータ(Invalid parameter)。もし特定された'num_clusters'が、ひとつの完全なプロテクショングループと厳密に等しくなければ、このエラーは戻されるであろう。デバッグビルド(debug builds)において、また、不良の発呼側を見つけて出すを手助けすることは、レイアウトモジュールにおけるアサーション故障(assertion failure)を伴う。このエラーコードは、通常の動作では予期されず、且つコーディングエラーを示さない。

30

<擬似コード>

この文書では、BAMレベルレイアウトモジュールの基本的実施について述べる。

フェーズ1：グループ情報(group info)、幅要件(width requirement)、および装置ごとの自由ブロック統計に基づき、使用するのに十分な自由スペースを有する利用可能な装置のリストが形成される。システムのスペースが尽きるか、または多くの装置がダウンして所望のプロテクションレベルで書き込みできないケースについてここでは理解され、そしてエラーは発呼側に戻されるであろう。

40

フェーズ2：別のデヴィッドが装置のリストからプロテクショングループにおける各クラスタに割り当てられる。もし特定の装置に割り当てられたブロックが既に存在すれば、システムは、既存のブロックが移動される必要のないようにデヴィッドを割り当てることを試みる。これは、既存のブロックレイアウトおよび現在利用可能な装置によっては、可能であるかもしれないし、可能でないかもしれない。目標のパラメータが、ここでは大きな役割を演じるであろうことに注意されたい。もし目標がLAYOUT_REPAIRであれば、システムはブロックを移動することを避けるが、もしファイルを適切にバランスすることが必要とされていれば、そのようにする。最終に、もし目標がLAYOUT_REBALANCEであれば、既存のブロックがある場所とは無関係に、最適な配置が計算される。

50

フェーズ3：どのクラスタにもデヴィッドを割り当てると、システムは、割り当て目的で各クラスタに p b a を割り当てる。これは、各装置上に割り当てられた最後のブロックのブロックアドレスを与える struct inode における p b a リストを参照することにより行われる。ここでの唯一の希薄な点は、システムが p b a を使用するためのデータブロックと、パリティ p b a を使用するためのパリティブロックとを欲することである。それらは、struct inode に別々に格納され、従ってリストが参照されるべきである。

<問題>

1) ディスクスペースの不足：このアルゴリズムひとつの未解決の問題は、いくつかの又は全てのノードが容量に近づいたときにすべきことである。もし可能であれば、システムは、割り当てに関して E N O S P C エラーをもたらすであろう発呼側に装置を推奨することは避けたい。それらのエラーは、もし配列ステートが変更されていなければ、シスコールリスタート (syscall restart) を起動せず、したがってユーザにまで遡って伝搬される。この問題は、我々が実際に装置を結果リストに割り当てるときに、フェーズ2において最もよく取り扱われる。そのときに、我々は、その装置上の利用可能なディスクスペースを調べることができる（我々がアクセスすることができない予約スペースを捕らえることを記憶している）。情報が各シナリオにどのように影響を及ぼすかは次のようである：

LAYOUT_RETUNE：振る舞いに変化がない。このシナリオは、どんなことがあっても最適なレイアウトを選ぶ。もしシステムが最適でない配置を与えれば、チャンスは、それが既に利用可能であることと同然であることである。

LAYOUT_REPAIRE/LAYOUT_REBALANCE：もしシステムが、廃された current_devid と符合するデヴィッドを割り当てていれば、システムは、自由スペースについて心配する必要はない。システムは、既に、そこで割り当てられたブロックを割り当てており、従ってシステムは、その何れも移動していないか、または少しばかりを割り当てているかのどちらかである。システムは、これらのシナリオでは、可能な限り少ないブロックを移動することを欲し、従ってそれらの装置をそれまでのままにしておく。

システムが空のロットにデヴィッドをまさに割り当てようとしているときに（そのクラスタについては current_devid が存在しないことを意味する）、その装置について自由スペースを参照する。もしそれがほとんど満杯であれば、より多くの自由スペースを有する他の好ましい装置を探すことを試み、そしてそれを代わりに使用する。もし好ましい装置の全てがほとんど満杯であれば、好ましい装置と予備装置のリストから最大量の自由スペースを有する装置を割り当てる。

E N O S P C は、レイアウトエンジン A P I から戻されないことに注意されたい。自由スペース情報は、それらのカウントをリモート装置から読み出すのにある程度の時間を要するので、少しばかり期限を過ぎている。スペースが厳しいとき、システムは、既存のカウントに基づき最も見込みのある装置を推奨する。最悪の場合、発呼側は、割り当てが試みられるときに E N O S P C を取得するであろう。

<アルゴリズム擬似コード：>

```
int
bam_layout_protection_group(gi, free_clusters, ip, goal, protection,
    width_policy, start_lbn, num_clusters, curr_devid_depth, curr_devids,
    alloc_counts, *result_devids, *result_pbas)
{
    /* dinodeから現在幅装置リスト(current width device list)を取得する。*/
    /* リストから n のダウン装置を削除する。*/
    /* max_width(MAX_WIDTH - n)を計算する。*/
    /* リストに残りの全てのアップ装置を加える。*/
    while (width device list size < UP device list size) {
        /* アップ装置に最も自由なクラスタを加える。*/
    }
    /*
```

10

20

30

40

50

```

* 長さがmax_widthにリストを切り詰める。このことは、dinode width device
* limitをオーバーフローすることを防止する。
*/
/*
* width_policy及び所望プロテクションパラメータを用いて、preferred_width
* 及びmin_widthを計算する。
*/
/* 我々は利用可能な装置を十分に備えていることを検証する。*/
if (not)
    return EROFS;
/*
* この点で、我々は、おそらく我々が使用することができるオンライン装置の全
* てのリストを有する。リストの始まりには、このファイルによって現在使用
* される装置があり、この後には、利用可能な自由スペースによってソート
* (sort)された残りの全ての装置が続く。
*/
/*
* リストを2つの部分に分割する。第1のサブリストはまさに長さにおいて好ま
* しい幅(preferred_width)であり、我々のマスターリストの先頭から始まる。
* 第2のサブリストはマスターリストの残りを含む。我々は、第1のリストを
我々の ' 好ましい装置(preferred_devices) ' と呼び、そして、第2のリスト
を ' 我々の ' 予備の装置(spare_devices) ' と呼ぶ。第2のリストは空
(empty)であ
* ってもよい。
*/
/* ここで、プロテクショングループのレイアウトを開始する。*/
/* num_clusters及びalloc_countsの正常チェック(sanity check)。*/
if (insane)
/*
* 我々は、全数のプロテクショングループを必要とし、そして、alloc_counts
* は0と16の間である。
Return EINVAL;
/* start_lbnを計算する(プロテクショングループにおける最初のl b m)。*/
/* start_lbnを用いてgroup_numberを計算する。*/
/* ダウン(DOWN)であるcurrent_devidsのエレメントをクリアする。*/
/* current_devidsにおける複製装置(duplicate devices)をクリアする。*/
/*
* preferred_deviceにおけるstart_offsetを計算する(恐らく(group_number %
* list size)にちょうどなる)。これはスキューを受け持つ。
*/
/* preferred_listにおけるstart_offsetから始めて、我々が割り付ける各クラ
/* スタについてのリストを通したステップを進め、そして各デヴィッドを新た
な /* リストにコピーする(必要であれば、好ましいリストの先頭に包み込む)
。こ /* の新規なサブリストは最適装置リスト(optimal_device list)と
呼ばれる。
*/
/* 発呼側特定の目標に基づきデヴィッドを割り当てる。*/
switch (goal) {
    case LAYOUT_REPAIR:
        /* ブロック移動を避ける。*/
        error = bam_layout_repair();

```

10

20

30

40

50

```

        break;
    case LAYOUT_REBALANCE:
        /* バランスされたレイアウトを確保する。 */
        error = bam_layout_rebalance();
        break;
    case LAYOUT_RETUNE:
        /* 最適なレイアウト。 */
        error = bam_layout_retune();
        break;
    default:
        ASSERT(0);
};
/*
 * もうすぐ終わり。全クラスタはデヴィッドを割り当てられるべき。
 */
/* ここで、我々は、前のブロックアドレスを選択しなければならない・・・ */
for (each cluster in the protection group) {
    /* p b a をzero_baddrに初期化する。 */
    /* これがパリティであるのかデータクラスタであるのかを決定する。 */
    /* struct inodeから適切な p b a リストを取得する。 */
    /* 使用可能な p b a を探す。 */
    for (each item in the pba list) {
        /*
         * もし、デヴィッドがこのクラスタについての我々の推奨デヴ
         * イッドと一致すれば、この p b a を使用し、そしてブレーク
         * する。
         */
    }
}
/* ここで、本当に終わり！ */
}
int
bam_layout_repair()
{
    /* 可能な限り既存のデヴィッドを割り当てる。 */
    /* 最適ナリストからのデヴィッドで空(empty)のロットを満たす。 */
}
int
bam_layout_rebalance()
{
    /* 最適ナリストに存在しない既存のデヴィッドはない。 */
    /*
     * もし既存のブロックの移動を最小化することが必要であれば、最適ナリストを
     * 再シャッフルし、そして出力に割り当てる。
     */
}
int
bam_layout_retune()
{
    /* 最適ナリストを出力に割り当て、何も質問なし。 */

```

【 0 2 6 7 】

IX. 著作権情報

この特許文書の開示の一部は、著作権保護の対象となっている。著作権所有者は、特許商標局の特許ファイルまたは記録にあるので、特許文書または特許開示のいずれによるファクシミリ複製に対しても異議をもたないが、そうでなければ、どんなものであれ、すべての著作権を留保する。

【 0 2 6 8 】

X. 結論

本発明のいくつかの実施形態を説明したが、これらの実施形態は単なる一例を表したものであり、本発明の範囲を制限することを意図したものではない。従って、本発明の広がり
10
と範囲は、添付の特許請求の範囲及びそれらの均等物によって定められるべきである。

【 図面の簡単な説明 】

【 0 2 6 9 】

【 図 1 】 本発明の一実施形態のハイレベルブロック図である。

【 図 2 】 図 1 に図解される構成要素間でのデータの流れの一例を示す図である。

【 図 3 】 スマート記憶ユニットの一例のハイレベルブロック図である。

【 図 4 】 ファイルディレクトリの一例を示す図である。

【 図 5 】 メタデータ構造の維持を説明するための図である。

【 図 6 A 】 データロケーションテーブル構造の一実施形態を示す図である。

【 図 6 B 】 データロケーションテーブル構造の追加の実施形態を示す図である。
20

【 図 6 C 】 データロケーションテーブル構造の追加の実施形態を示す図である。

【 図 6 D 】 データロケーションテーブル構造の追加の実施形態を示す図である。

【 図 7 A 】 ディレクトリのためのメタデータ構造の一実施形態を示す図である。

【 図 7 B 】 ファイルのためのメタデータ構造の一実施形態を示す図である。

【 図 8 A 】 データロケーションテーブルの一実施形態を示す図である。

【 図 8 B 】 データロケーションテーブルの追加の実施形態を示す図である。

【 図 8 C 】 データロケーションテーブルの追加の実施形態を示す図である。

【 図 9 】 対応するサンプルデータを有するファイルのメタデータの一例を示す図である。

【 図 1 0 】 ファイルを読み出すためのフローチャートの一実施形態を示す図である。

【 図 1 1 】 名前解決(name resolution)を実施するためのフローチャートの一実施形態を示す図である。
30

【 図 1 2 】 ファイルを読み出すためのフローチャートの一実施形態を示す図である。

【 図 1 3 】 パリティ情報を生成するためのフローチャートの一実施形態を示す図である。

【 図 1 4 】 エラー訂正を実施するためのフローチャートの一実施形態を示す図である。

【 図 1 5 】 インテリジェント分散ファイルシステムにおいてデータを再ストライピングするためのフローチャートの一実施形態を示す図である。

【 図 1 6 】 スマート記憶ユニットにデータを割り当てるためのフローチャートの一実施形態を示す図である。

【 図 1 7 】 利用可能な記憶ユニットのセットの中から選択するためのフローチャートの一実施形態を示す図である。
40

【 符号の説明 】

【 0 2 7 0 】

1 1 0 インテリジェント分散ファイルシステム

1 1 4 スマート記憶ユニット

1 2 0 サーバ

1 2 5 スイッチ

1 4 0 通信媒体

1 4 5 インターネット

1 3 0 ユーザ

【図1】

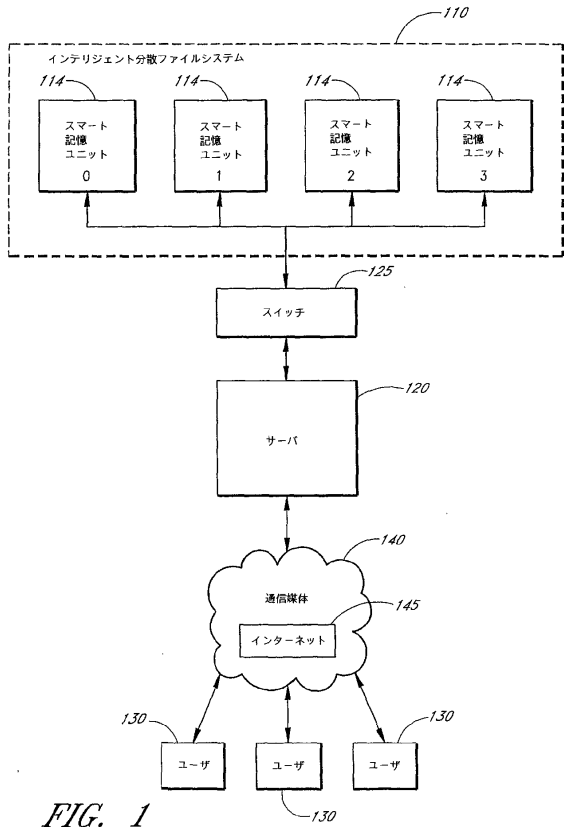


FIG. 1

【図2】

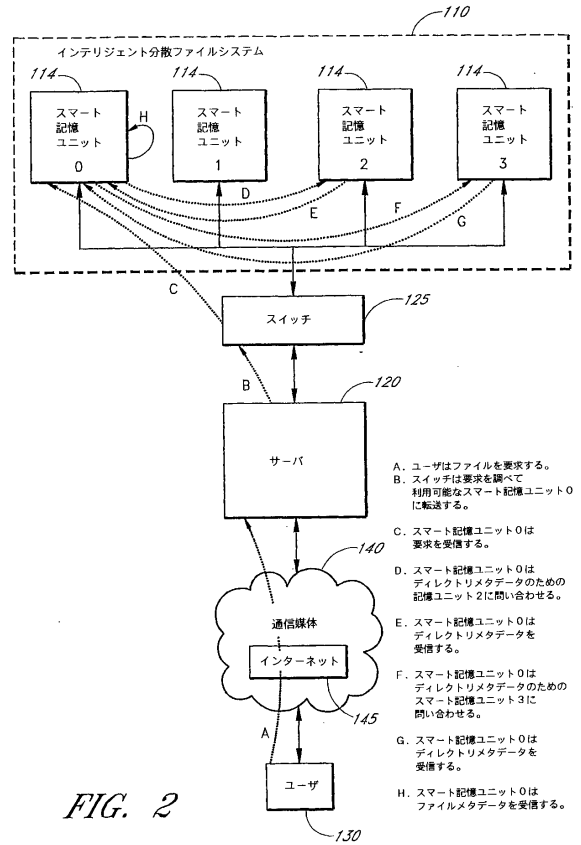


FIG. 2

【図3】

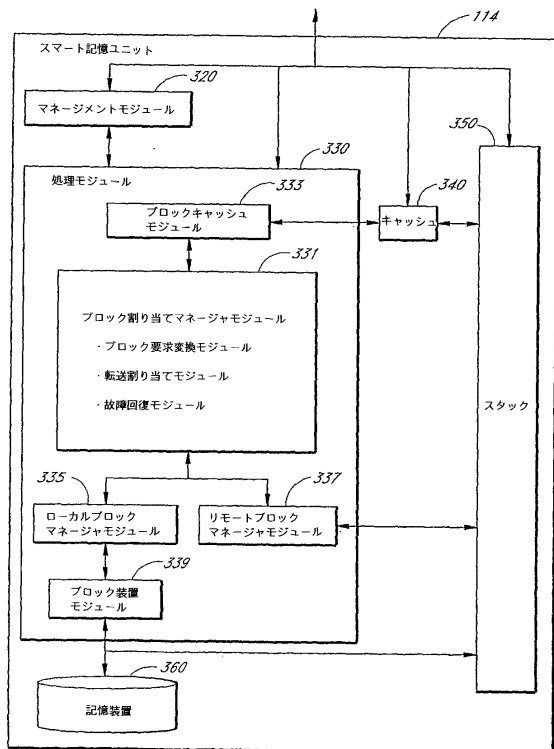


FIG. 3

【図4】

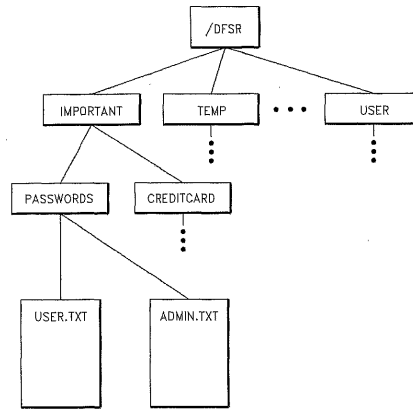


FIG. 4

【図5】

510

モード
オーナー
タイムスタンプ
サイズ
パリティカウント
ミラーカウント
バージョン
タイプ
0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
リファレンスカウント
フラグ
パリティマップポイント

データ
ロケーション
テーブル

FIG. 5

【図6A】

タイプ0

装置	ブロック
0	D0
1	D1
2	D2
3	D3
4	D4
5	D5
6	D6
7	D7
8	D8
9	D9
10	D10
11	D11
12	D12
13	D13
14	D14
15	D15
16	D16
17	D17
18	D18
19	D19
20	D20
21	D21
22	D22
23	D23

FIG. 6A

【図6B】

タイプ1

装置	ブロック
0	D0
1	D1
2	D2
3	D3
4	D4
5	D5
6	D6
7	D7
8	D8
9	D9
10	D10
11	D11
12	D12
13	D13
14	D14
15	SI0
16	D10
17	TI0
18	SI1
19	DI1
20	TI1
21	SI2
22	DI2
23	TI2

FIG. 6B

【図6C】

タイプ2

装置	ブロック
0	SI0
1	DI0
2	TI0
3	SI1
4	DI1
5	TI1
6	SI2
7	DI2
8	TI2
9	SI3
10	DI3
11	TI3
12	SI4
13	DI4
14	TI4
15	SI5
16	DI5
17	TI5
18	SI6
19	DI6
20	TI6
21	SI7
22	DI7
23	TI7

FIG. 6C

【図 6 D】

タイプ3		
装置	ブロック	
0		TI0
1		TI1
2		TI2
3		TI3
4		TI4
5		TI5
6		TI6
7		TI7
8		TI8
9		TI9
10		TI10
11		TI11
12		TI12
13		TI13
14		TI14
15		TI15
16		TI16
17		TI17
18		TI18
19		TI19
20		TI20
21		TI21
22		TI22
23		TI23

FIG. 6D

【図 7 A】

モード	ディレクトリ
オーナー	ルート
タイムスタンプ	65536
サイズ	345
パリティカウント	0
ミラーカウント	3
バージョン	1
タイプ	1

装置1	ブロック11
装置2	ブロック21
装置3	ブロック31

リファレンスカウント	1
フラグ	777
パリティマップポインタ	ゼロ(NULL)

FIG. 7A

【図 7 B】

モード	レギュラーファイル
オーナー	ルート
タイムスタンプ	65892
サイズ	45897
パリティカウント	3+1
ミラーカウント	1
バージョン	1
タイプ	1

装置1	ブロック11
装置2	ブロック21
装置3	ブロック31
装置2	ブロック100
装置3	ブロック343
装置4	ブロック34

リファレンスカウント	1
フラグ	777
パリティマップポインタ	0x12

装置4	ブロック45
装置1	ブロック87

FIG. 7B

【図 8 A】

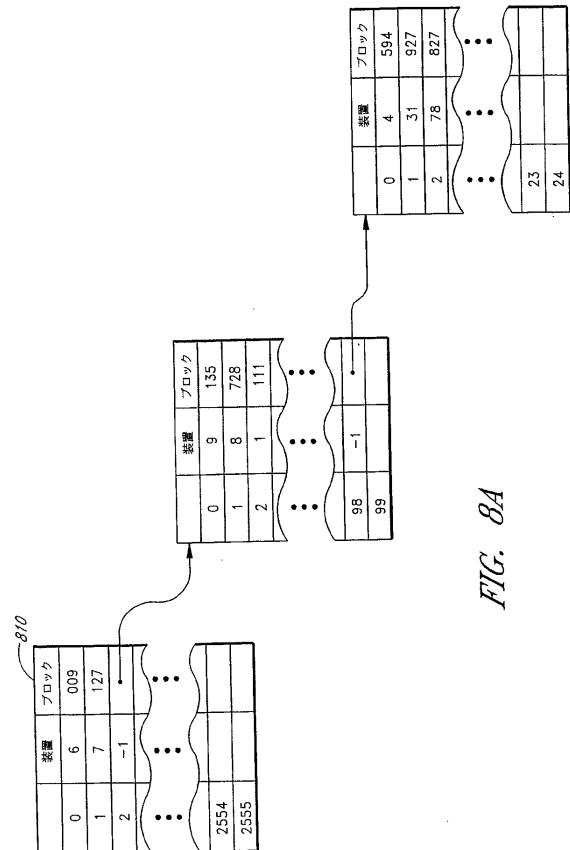


FIG. 8A

【図 8 B】

装置	ブロック
0	11
1	12
2	13
3	14
4	15
5	11
6	12
7	13
8	14
9	15
10	11
11	12

装置	ブロック
0	6
1	7
2	8
3	9
4	10
5	6
6	7
7	8
8	9
9	10
10	6
11	7

装置	ブロック
0	1
1	2
2	3
3	4
4	5
5	1
6	2
7	3
8	4
9	5
10	1
11	2

FIG. 8B

【図 8 C】

装置	ブロック
0	5
1	9
2	7
3	5
4	9
5	7
6	5
7	9
8	7
9	5
18	5
19	9
20	7

装置	ブロック
0	6
1	8
2	10
3	6
4	8
5	10
6	6

FIG. 8C

【図 9】

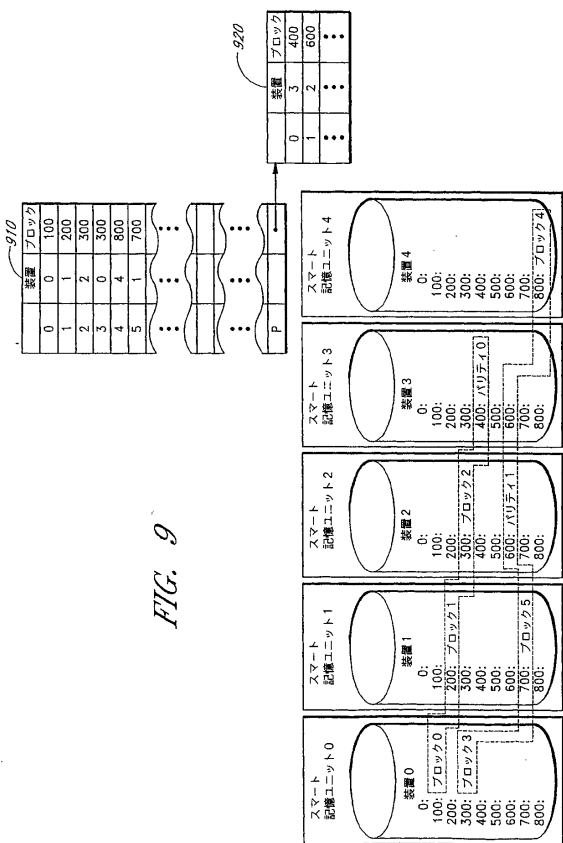


FIG. 9

【図 10】

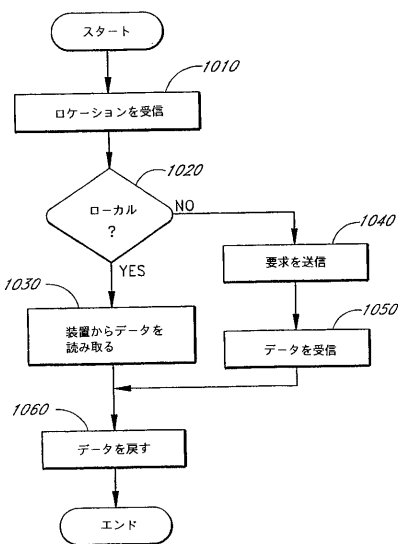


FIG. 10

【図 11】

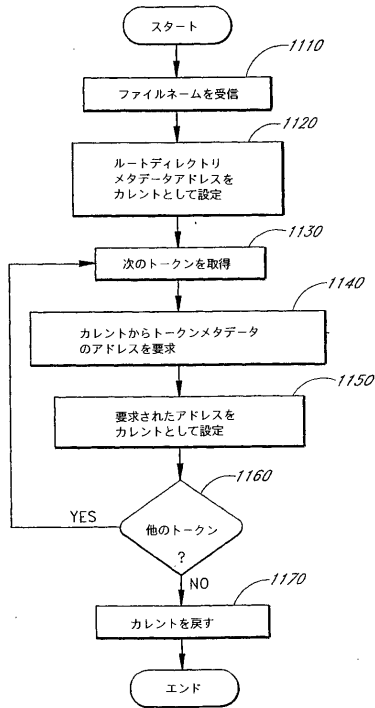


FIG. 11

【図 12】

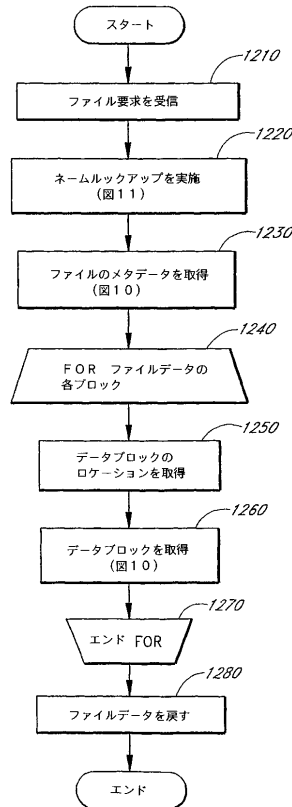


FIG. 12

【図 13】

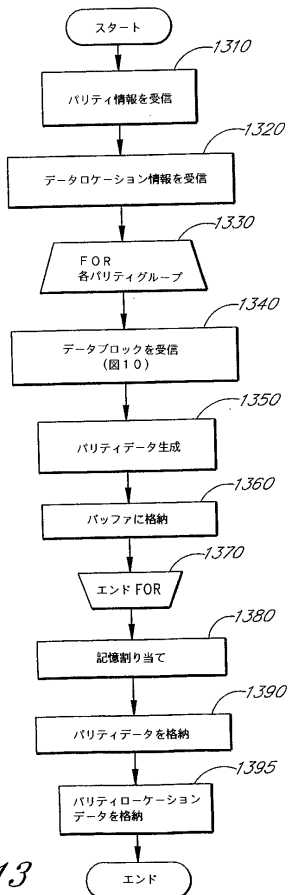


FIG. 13

【図 14】

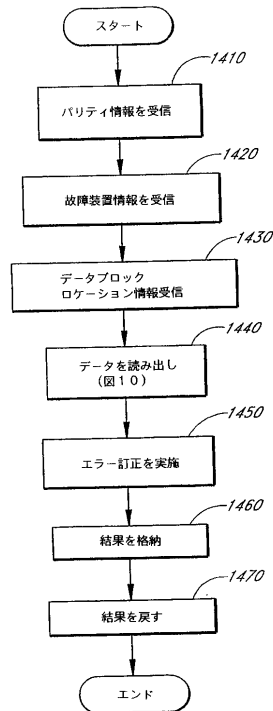


FIG. 14

【図15】

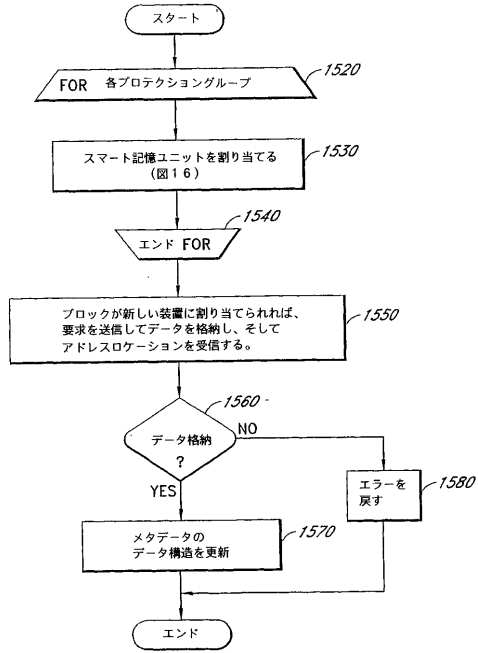


FIG. 15

【図16】

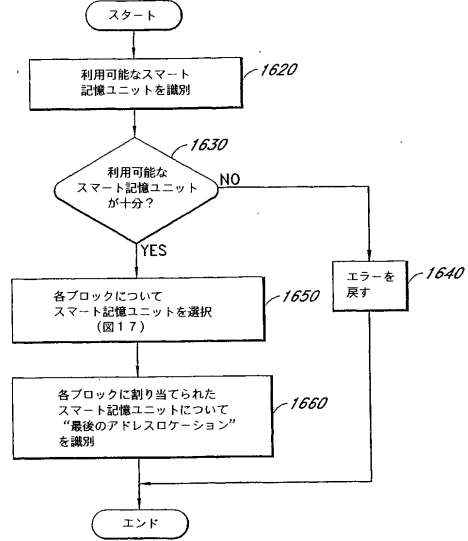


FIG. 16

【図17】

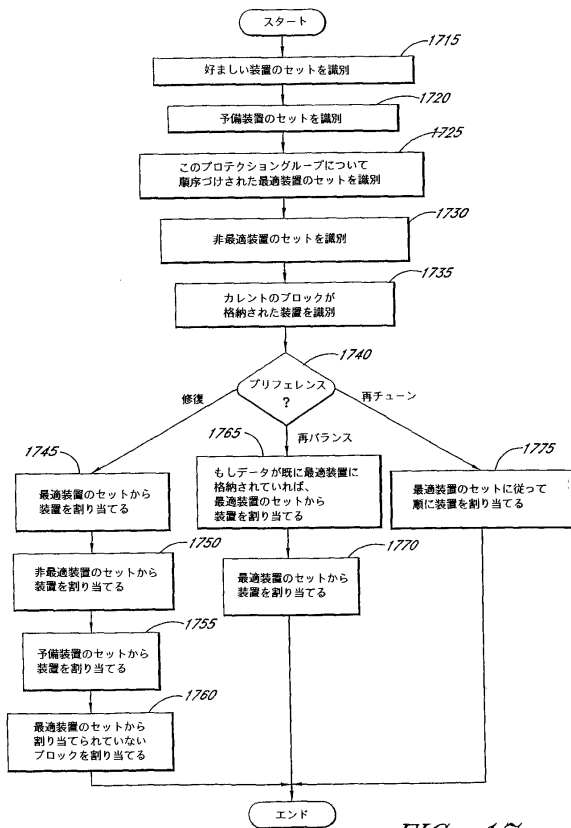


FIG. 17

フロントページの続き

- (72)発明者 ポール・エー・マイクセル
アメリカ合衆国・ワシントン・98121・シアトル・フォース・アヴェニュー・#171・240
0
- (72)発明者 ロブ・アンダーソン
アメリカ合衆国・ワシントン・98106・シアトル・フォーティーンズ・アヴェニュー・ノース
ーイスト・6330
- (72)発明者 アーロン・ジェイ・パッセイ
アメリカ合衆国・ワシントン・98107・シアトル・セヴェンティーンズ・アヴェニュー・ノース
ーイスト・5402
- (72)発明者 ピーター・ジョン・ゴッドマン
アメリカ合衆国・ワシントン・98116・シアトル・エイキンス・アヴェニュー・サウスウエスト
・4221
- (72)発明者 ハッサン・エフ・カーン
アメリカ合衆国・ワシントン・98121・シアトル・フォース・アヴェニュー・2514・アパー
トメント・801
- (72)発明者 ダレン・ピー・シャック
アメリカ合衆国・ワシントン・98119・シアトル・ウェスト・レパブリック・ストリート・#
440・415

審査官 工藤 嘉晃

(56)参考文献 特開2000-99282(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

G06F 3/06

JSTPlus(JDreamII)