

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5102776号
(P5102776)

(45) 発行日 平成24年12月19日(2012.12.19)

(24) 登録日 平成24年10月5日(2012.10.5)

(51) Int.Cl.	F I
G06F 3/06 (2006.01)	G06F 3/06 305C
G06F 12/16 (2006.01)	G06F 3/06 540
	G06F 12/16 320L

請求項の数 34 (全 38 頁)

(21) 出願番号	特願2008-545792 (P2008-545792)	(73) 特許権者	303039534
(86) (22) 出願日	平成18年12月14日(2006.12.14)		ネットアップ、インコーポレイテッド
(65) 公表番号	特表2009-524124 (P2009-524124A)		アメリカ合衆国 カリフォルニア 940
(43) 公表日	平成21年6月25日(2009.6.25)		89, サニーヴェール, イースト ジ
(86) 国際出願番号	PCT/US2006/047647		ャバ ドライブ 495
(87) 国際公開番号	W02007/078803	(74) 代理人	100087642
(87) 国際公開日	平成19年7月12日(2007.7.12)		弁理士 古谷 聡
審査請求日	平成21年9月24日(2009.9.24)	(74) 代理人	100076680
(31) 優先権主張番号	11/304,369		弁理士 溝部 孝彦
(32) 優先日	平成17年12月15日(2005.12.15)	(74) 代理人	100121061
(33) 優先権主張国	米国 (US)		弁理士 西山 清春
		(72) 発明者	コルベット, ピーター, エフ.
			アメリカ合衆国マサチューセッツ州024
			20, レキシントン, サマー・ストリート
			・33

最終頁に続く

(54) 【発明の名称】 ストレージアレイにおける三重故障からの効率的な復旧を可能にする三重パリティ技術

(57) 【特許請求の範囲】

【請求項1】

ストレージアレイにおける記憶装置の3以下の同時故障からの復旧を可能にする方法であって、

データ、及び行パリティを格納するように構成された複数の第1の装置、対角パリティを格納するように構成された1つの対角パリティ装置、並びに反対角パリティを格納するように構成された反対角パリティ装置を含む所定数の記憶装置を備え、前記記憶装置の所定数 n が $p + 2$ であり、 p が素数であるアレイを用意するステップと、

各装置を複数のブロックに分割するステップと、

前記ブロックを各装置上の同数のブロックを含む複数のストライプに編成するステップであって、各ストライプが $n - 3$ 行のブロックを含むように前記ブロックを編成するステップと、

前記複数の第1の装置にわたって規定される対角パリティ集合に沿って対角パリティを定義するステップであって、前記対角パリティ集合が、 $n - 3$ 行のグループの中で循環し、あるストライプの対角パリティ集合に属するブロックが全て、そのストライプに格納されるように、対角パリティを定義するステップと、

1つを除く全ての対角パリティ集合について、対角パリティを計算し、前記対角パリティ装置に格納するステップと、

前記複数の第1の装置にわたって規定される反対角パリティ集合に沿って反対角パリティを定義するステップであって、前記反対角パリティ集合が、 $n - 3$ 行のグループの中

10

20

で循環し、あるストライプの反対角パリティ集合に属するブロックが全て、そのストライプに格納されるように、反対角パリティを定義するステップと、

1つを除く全ての反対角パリティ集合について、反対角パリティを計算し、前記反対角パリティ装置に格納するステップと
からなる方法。

【請求項2】

あるストライプ中の行パリティブロックが全て、単一の装置に格納される、請求項1に記載の方法。

【請求項3】

前記第1の装置は複数のデータ装置を含み、前記データ装置は全て存在するのではなく、不在のデータ装置はゼロ値のデータを格納するものとして処理される、請求項1に記載の方法。

10

【請求項4】

異なる装置内において、パリティブロックの位置が、装置ごとにシフトされる、請求項1に記載の方法。

【請求項5】

前記ストライプのサイズは、2の乗数個のビットである、請求項1に記載の方法。

【請求項6】

ストレージアレイにおける記憶装置の3以下の同時故障からの復旧を可能にするように構成されたシステムであって、

20

データ、及びパリティを格納するように構成された複数の第1の装置、対角パリティを格納するように構成された1つの対角パリティ装置、並びに反対角パリティを格納するように構成された反対角パリティ装置を含む所定数の記憶装置を備え、前記記憶装置の所定数 n が $p + 2$ であり、 p が素数である前記ストレージアレイと、

(i) 前記複数の第1の装置にわたって規定される対角パリティ集合に沿って対角パリティを計算し、(ii) 1つを除く全ての対角パリティ集合について対角パリティを前記対角パリティ装置に格納し、(iii) 前記複数の第1の装置にわたって規定される反対角パリティ集合に沿って反対角パリティを計算し、(iv) 1つを除く全ての反対角パリティ集合について反対角パリティを前記反対角パリティ装置に格納する、三重パリティ(TP)技術を実施するように構成されたデバイスストレージ層を含むストレージオペレーティングシステムと、

30

前記ストレージオペレーティングシステムを実行することにより、前記TPパリティ技術にしたがって前記ストレージアレイに対する双方向のストレージアクセス操作を実施するように構成された処理要素と

からなるシステム。

【請求項7】

あるストライプ中の行パリティブロックが全て、単一の装置に格納される、請求項6に記載のシステム。

【請求項8】

前記デバイスストレージ層はRAIDシステムであり、前記記憶装置はディスクである、請求項6に記載のシステム。

40

【請求項9】

前記RAIDシステムは、各ディスクをさらに複数のブロックに分割し、該ブロックを複数のストライプに編成する、請求項8に記載のシステム。

【請求項10】

$n = p + 2$ として、各ストライプが、 $n - 3$ 行のブロックを含み、各行が、各ディスクから1つのブロックを有する、請求項9に記載のシステム。

【請求項11】

前記装置は、ビデオテープ、磁気テープ、光学媒体、DVD、バブルメモリ、磁気ディスク、電氣的ランダムアクセスメモリ、及びMEMSデバイスのうちのいずれかである、

50

請求項 6 に記載のシステム。

【請求項 1 2】

前記第 1 の装置は複数のデータ装置を含み、前記データ装置は全て存在するのではなく、不在のデータ装置はゼロ値のデータを格納するものとして処理される、請求項 6 に記載のシステム。

【請求項 1 3】

ストレージアレイにおける記憶装置の 3 以下の同時故障からの復旧を可能にする装置であって、

データ、及び行パリティを格納するように構成された複数の第 1 の装置、対角パリティを格納するように構成された 1 つの対角パリティ装置、並びに反対角パリティを格納するように構成された反対角パリティ装置を含む所定数の記憶装置を備え、前記記憶装置の所定数 n が $p + 2$ であり、 p が素数であるアレイを設ける手段と、

各装置を複数のブロックに分割する手段と、

前記複数の第 1 の装置にわたって規定される対角パリティ集合に沿って対角パリティを定義する手段であって、前記対角パリティ集合が、 $n - 3$ 行のグループの中で循環し、あるストライプの対角パリティ集合に属するブロックが全て、そのストライプに格納されるように、対角パリティを定義する手段と、

1 つを除く全ての対角パリティ集合について、対角パリティを計算し、前記対角パリティ装置に格納する手段と、

前記複数の第 1 の装置にわたって規定される反対角パリティ集合に沿って反対角パリティを定義する手段であって、前記反対角パリティ集合が、 $n - 3$ 行のグループの中で循環し、あるストライプの反対角パリティ集合に属するブロックが全て、そのストライプに格納されるように、反対角パリティを定義する手段と、

1 つを除く全ての反対角パリティ集合について、反対角パリティを計算し、前記反対角パリティ装置に格納する手段と

からなる装置。

【請求項 1 4】

あるストライプ中の行パリティブロックが全て、単一の装置に格納される、請求項 1 3 に記載の装置。

【請求項 1 5】

各装置を複数のブロックに分割する手段と、

前記ブロックを複数のストライプに編成する手段と

をさらに含む、請求項 1 3 に記載の装置。

【請求項 1 6】

前記格納する手段は、1 つを除き、ストライプの対角パリティ集合のそれぞれについて、対角パリティブロックを前記対角パリティ装置に格納する手段を含む、請求項 1 5 に記載の装置。

【請求項 1 7】

前記格納する手段は、1 つを除き、ストライプの反対角パリティ集合のそれぞれについて、反対角パリティブロックを前記反対角パリティ装置に格納する手段を含む、請求項 1 5 に記載の装置。

【請求項 1 8】

前記第 1 の装置は複数のデータ装置を含み、前記データ装置は全て存在するのではなく、不在のデータ装置はゼロ値のデータを格納するものとして処理される、請求項 1 3 に記載の装置。

【請求項 1 9】

前記ストライプのサイズは、2 の乗数個のビットである、請求項 1 5 に記載の装置。

【請求項 2 0】

各ストライプ中のブロックの数に第 1 のブロックサイズを乗じたものは、ファイルシステムがストレージアレイにアクセスするために使用する第 2 のブロックサイズに等しい、

10

20

30

40

50

請求項 1 5 に記載の装置。

【請求項 2 1】

ストレージアレイにおける 2 つの記憶装置の 2 以下の同時故障からの復旧を可能にする実行可能プログラム命令が格納されたコンピュータ読み取り可能な記憶媒体であって、前記実行可能プログラム命令が、

データ、及び行パリティを格納するように構成された複数の第 1 の装置、対角パリティを格納するように構成された 1 つの対角パリティ装置、並びに反対角パリティを格納するように構成された反対角パリティ装置を含む所定数の記憶装置を備え、前期記憶装置の所定数 n が $p + 2$ であり、 p が素数であるアレイを形成し、

各装置を複数のブロックに分割し、

各ストライプが $n - 3$ 行のブロックを含むように、前記ブロックを各装置上の同数のブロックを含む複数のストライプに編成し、

前記複数の第 1 の装置にわたって規定される対角パリティ集合に沿って対角パリティを定義し、前記対角パリティ集合が、 $n - 3$ 行のグループの中で循環し、あるストライプの対角パリティ集合に属するブロックが全て、そのストライプに格納されるようにし、

1 つを除く全ての対角パリティ集合について、対角パリティを計算し、前記対角パリティ装置に格納し、

前記複数の第 1 の装置にわたって規定される反対角パリティ集合に沿って反対角パリティを定義し、前記反対角パリティ集合が、 $n - 3$ 行のグループの中で循環し、あるストライプの反対角パリティ集合に属するブロックが全て、そのストライプに格納されるようにし、

1 つを除く全ての反対角パリティ集合について、反対角パリティを計算し、前記反対角パリティ装置に格納するためのプログラム命令からなる、コンピュータ読み取り可能な記憶媒体。

【請求項 2 2】

あるストライプ中の行パリティブロックが全て、単一の装置に格納される、請求項 2 1 に記載のコンピュータ読み取り可能な記憶媒体。

【請求項 2 3】

ストレージアレイにおける 3 つの記憶装置の 3 以下の同時故障からの復旧を可能にする方法であって、

データ、及び行パリティを格納するように構成された複数の第 1 の装置、対角パリティを格納するように構成された対角パリティ装置、並びに反対角パリティを格納するように構成された反対角パリティ装置を含む所定数の記憶装置を備えたアレイを用意するステップと、

前記複数の第 1 の装置にわたって規定される対角パリティ集合に沿って対角パリティを計算するステップと、

1 つを除く全て前記対角パリティ集合について、対角パリティを前記対角パリティ装置に格納するステップと、

前記複数の第 1 の装置にわたって規定される反対角パリティ集合に沿って反対角パリティを計算するステップと、

1 つを除く全ての前記反対角パリティ集合について、反対角パリティを前記反対角パリティ装置に格納するステップと

からなる方法。

【請求項 2 4】

各装置を複数のブロックに分割するステップと、

前記ブロックを複数のストライプに編成するステップと

をさらに含む、請求項 2 3 に記載の方法。

【請求項 2 5】

前記第 1 の装置は複数のデータ装置を含み、前記データ装置は全て存在するのではなく、不在のデータ装置はゼロ値のデータを格納するものとして処理される、請求項 2 3 に記

10

20

30

40

50

載の方法。

【請求項 2 6】

前記記憶装置の所定数は n であり、 $n = p + 2$ であり、 p は素数である、請求項 2 3 に記載の方法。

【請求項 2 7】

ストレージアレイにおける記憶装置の 3 つの同時故障からの復旧を可能にする方法であって、

失われた対角パリティ、及び反対角パリティを計算するステップと、
行、対角、及び反対角のそれぞれに沿って、一組の故障した記憶装置のそれぞれに対し、代数計算を実施するステップと、
故障した中間記憶装置に対して 4 タプル一組の和を計算するステップと
からなる方法。

10

【請求項 2 8】

前記失われた対角パリティを計算するステップは、対角パリティ記憶装置上の複数のブロックを計算することからなる、請求項 2 7 に記載の方法。

【請求項 2 9】

前記失われた反対角パリティを計算するステップは、反対角パリティ記憶装置上の複数のブロックを計算することからなる、請求項 2 7 に記載の方法。

【請求項 3 0】

前記中間記憶装置に対して 4 タプル一組の和を計算するステップは、
故障した記憶装置に順位を付けるステップと、
前記記憶装置の各行について
(a) 一行の記憶装置を選択し、
(b) 前記選択された行に対応する失われた記憶装置上のブロックの行パリティ和を読み出し、
(c) 第 1 の故障した記憶装置について、前記選択された行にあるブロックの対角パリティを読み出し、
(d) 第 2 の故障した記憶装置について、前記選択された行にあるブロックの反対角パリティを読み出し、
(e) 前記反対角、及び対角の最後の行に対応する失われた記憶装置のブロックの行パリティ和を読み出すステップと、
前記ステップ(b)、(c)、(d)、及び(e)の結果に対して排他的論理和 (XOR) 演算を実施することにより、4 タプル和を形成するステップと
をさらに含む、請求項 2 7 に記載の方法。

20

30

【請求項 3 1】

前記 4 タプル和を前記中間記憶装置上での二つ一組の和に低減するステップと、
前記中間記憶装置を復旧するステップと
をさらに含む、請求項 2 7 に記載の方法。

【請求項 3 2】

行 - 対角パリティ復元技術を実施し、残りの故障した記憶装置を復元するステップをさらに含む、請求項 3 1 に記載の方法。

40

【請求項 3 3】

前記記憶装置はディスクである、請求項 2 7 に記載の方法。

【請求項 3 4】

前記代数演算は、排他的論理和 (XOR) からなる、請求項 2 7 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

[発明の分野]

本発明は、ストレージシステムのアレイに関し、より具体的には、ストレージアレイの

50

任意の3つの記憶装置の故障を効率的に復元する技術に関する。

【背景技術】

【0002】

[発明の背景]

ストレージシステムは通常、要求に応じてデータを出し入れすることが可能な1以上の記憶装置を含む。ストレージシステムは、種々のストレージアーキテクチャにしたがって実施され、限定はしないが、例えば、ネットワーク・アタッチド・ストレージ環境、ストレージ・エリア・ネットワーク、あるいは、クライアント、若しくはホストコンピュータに直接取り付けられたディスクアセンブリのような種々のストレージアーキテクチャにしたがって実施される。記憶装置は通常、ディスクドライブであり、ここで言う「ディスク」という用語は一般に、内蔵型の回転式磁気媒体記憶装置を意味する。この文脈において「ディスク」という用語とは、ハードディスクドライブ(HDD)やダイレクト・アクセス・ストレージ・デバイス(DASD)と同義である。

10

【0003】

ストレージシステム内のディスクは一般に1以上のグループに編成され、各グループが、RAID (Redundant Array of Independent (Inexpensive) Disks) として運用される。大半のRAID実施形態は、RAIDグループを成す所与の数の物理ディスクにわたってデータ「ストライプ」を冗長書き込みし、そのストライピングされたデータに対する冗長情報を適切に記憶することによって、データ記憶の信頼性/完全性を向上させる。冗長情報により、記憶装置が故障したときに失われたデータの復元が可能になる。

20

【0004】

ディスクアレイの動作中に、ディスクは故障することがある。高い処理能力を持つストレージシステムの目標は、平均データ紛失時間(MTTL)を出来る限り長くすることであり、MTTLは、システムの期待サービス寿命よりも遥かに長いことが好ましい。1以上のディスクが故障すると、データは失われることがあり、装置からデータを復旧できなくなることがある。データ紛失を回避する典型的な手法としては、ミラーリング、バックアップ、及びパリティ保護などがある。ディスクのようなストレージリソースの消費の観点からすれば、ミラーリングは高価な解決策である。バックアップは、バックアップ作成後に変更が加えられたデータを保護することができない。パリティ技術が一般的である理由は、システムに一台のディスクドライブを追加するだけで、単一消去(一台のディスクの紛失)を許容するデータの冗長符号化が可能となるからである。

30

【0005】

パリティ保護は、コンピュータシステムにおいて、ディスクのような記憶装置上のデータの紛失に対する保護を提供するために使用される。パリティ値は、異なるデータを有する複数の同様のディスクにわたって特定ワードサイズ(通常は1ビット)のデータを加算(通常はモジュロ2)することによって計算され、その結果が、さらに別の同様のディスクに格納される。すなわち、パリティは、各ディスク上の対応する位置にあるビット幅から構成される1ビット幅の種々のベクトルに対して計算される。パリティが1ビット幅のベクトルに対して計算される場合、パリティは、和として計算される場合もあれば、その補数として計算される場合もある。それらのパリティはそれぞれ、偶数(EVEN)パリティ、及び奇数(ODD)パリティと呼ばれる。1ビットベクトルに対する加算、及び減算はいずれも、排他的論理和(XOR)演算と等価である。次に、いずれかの1つのディスクの紛失、またはいずれか1つのディスクにおける任意部分のデータの紛失から、データは保護される。パリティを格納するディスクが失われた場合、そのパリティは、データから復元することができる。データディスクの1つが失われた場合、そのデータは、生き残っているデータディスクの中身を加算し、その結果を格納されたパリティから減算することによって再現することができる。

40

【0006】

通常、ディスクは幾つかのパリティグループに分割され、各パリティグループは、1以上のデータディスクと、1つのパリティディスクとを含む。パリティ集合は、幾つかのデ

50

ータブロックと、1つのパリティブロックとを含むブロックの集合であり、パリティブロックは、全てのデータブロックをXOR演算したものである。パリティグループは、1以上のパリティ集合を選出する元になる一組のディスクである。ディスク空間は幾つかのストライプに分割され、各ストライプは、各ディスクから1つのブロックを格納する。ストライプを形成する幾つかのブロックは通常、パリティグループ内の各ディスク上の同じ位置にある。ストライプ内で、1つを除く全てのブロックは、データを格納するブロック（「データブロック」）であり、1つのブロックだけは、全てのデータのXORを取ることによって計算されたパリティを格納するブロック（「パリティブロック」）である。パリティブロックが全て1つのディスクに格納され、それによってパリティ情報を全て含む（且つ、パリティ情報しか持たない）単一のディスクが形成される場合、RAID - 4実施形態が提供される。各ストライプにおいて異なるディスク上に（通常は、巡回パターンを成すようにして）パリティブロックが格納される場合、実施形態はRAID - 5である。RAIDという用語、及びその種々の実施形態については広く知られており、1998年6月に、D.A. Gibson、及びR.H. Katzにより、「A Case for Redundant Arrays of Inexpensive Disks (RAID)」と題するデータ管理に関する国際会議論文（Proceedings of the International Conference on Management of Data）に開示されている。

10

【0007】

本明細書において、「符号化（encoding）」という用語は、データブロックの所定のサブセットに対する冗長値の計算を指し、「復号（decoding）」という用語は、データブロック、及びパリティブロックのサブセットを利用した、データブロック、又はパリティブロックの復元を意味する。パリティグループ内の1つのディスクが故障した場合、そのディスクの中身は、残りのデータブロックの中身を全て加算し、その結果をパリティブロックから差し引くことにより、予備ディスク（複数の場合もあり）上に復号（復元）することができる。1ビット幅での2の補数による加算と減算はいずれも、XOR演算と等価であるから、復元は、生き残ったデータブロックとパリティブロックを全てXOR演算することからなる。同様に、パリティディスクが失われた場合も、同様の仕方で、生き残ったデータからそれを再計算することが出来る。

20

【0008】

パリティ技術は通常、パリティグループ内の単一のディスク故障に対する保護を提供する。各故障が異なるパリティグループで発生する限り、パリティ技術は、複数のディスク故障に対する保護を提供することも可能である。ただし、1つのパリティグループ内で2つのディスクが同時に故障した場合、復元不能なデータ紛失を被る。1つのパリティグループ内で2つのディスクが同時に故障することは、極めて一般的に起こりうる。特に、ディスクの「磨耗」やディスクの動作に関する環境要因が原因で発生する。この文脈において、1つのパリティグループ内における2つのディスクの同時故障は、「二重故障」と呼ばれる。

30

【0009】

二重故障は一般に、一台のディスクが故障した後、その最初の故障からの復元を試みている間に、別のディスクが続けて故障する結果として発生する。復元時間、すなわち復旧時間は、ストレージシステムの活動レベルに応じて変わる。すなわち、故障したディスクを復元している間も、ストレージシステムは「オンライン」状態にあり、（クライアント、又はユーザからの）データアクセス（すなわち、読み出し、及び/又は書き込み）の要求に対してサービスを提供し続けることができる。ストレージシステムが要求に対する応答に忙しい場合、復元のための経過時間は増大する。また、失われたデータを復元するためには、生き残ったディスクを全て読み出さなければならないため、ストレージシステム内のディスクのサイズや数が増えるほど、復元処理時間も増大する。さらに、二重故障率は、パリティグループ内のディスク数の二乗に比例する。しかしながら、パリティグループを小さくすると、各パリティグループにつき、一台のディスク全体を冗長データの記憶のために専用に使わなければならないため、費用がかかる。

40

【0010】

50

ディスクのさらに別の故障形態は、ディスクの単一のブロック、又はセクタが読み出せなくなるメディア読み出しエラーである。ストレージアレイにパリティが保持されていれば、読み出せないデータを復元できることがある。ただし、アレイ内の一台のディスクが既に故障しているときに、さらに別のディスク上にメディア読み出しエラーが発生した場合、データは失われる。これが、二重故障の第2の形態である。

【0011】

二重故障の訂正に必要な冗長情報の最小量が2単位であることは、簡単に示すことができる。したがって、データディスクに追加することが可能なパリティディスクの最小数は、2である。これは、複数のディスクにわたってパリティが分散配置されるか、追加された2台のディスク上にパリティが集中配置されるかに関わらず、常に当てはまる。

10

【0012】

二重故障を訂正する既知のパリティ技術として、失われた(故障した)ディスクの逐次復元が可能な EVENODD XOR を利用する技術がある。EVENODD パリティは、ちょうどディスク2台分の冗長データを必要とし、この量が最適である。このパリティ技術によれば、ディスクブロックは全て2つのパリティ集合に属する。一方は、全てのデータディスクにわたる通常の RAID-4 スタイルでの XOR 演算によって計算され、他方は、斜めに隣り合う幾つかのディスクブロックの集合から計算される。対角パリティ集合は、1つを除く全てのデータディスクからブロックを含む。n 個のデータディスクに対し、1つのストライプ中には、n - 1 行のブロックが存在する。各ブロックは、一本の対角上にあり、n - 1 ブロック分の長さをそれぞれ有する n 本の対角が存在する。なお、EVENODD 技術は、n が素数でないと動作しない点に注意して欲しい。EVENODD 技術は、1995年2月に「EVENODD: An Efficient Scheme for Tolerating Double Disk Failures in RAID Architectures」と題する IEEE Transactions on Computers, Vol. 44, No.2 の記事に Blaum 他により開示されている。EVENODD の変形は、1996年11月26日に、「METHOD AND MEANS FOR ENCODING AND REBUILDING THE DATA CONTENTS OF UP TO TWO UNAVAILABLE DASDS IN A DASD ARRAY USING SIMPLE NON-RECURSIVE DIAGONAL AND ROW PARITY」と題する米国特許第5,579,475号に開示されている。上記の記事、及び特許は、参照により、本明細書の中で完全に説明されたものとして本明細書に援用される。

20

【0013】

EVENODD 技術は、p を素数として、総数 p + 2 個のディスクを使用する。うち、p 個のディスクはデータを格納し、残り2つのディスクはパリティ情報を格納する。一方のパリティディスクは、行パリティブロックを格納する。行パリティは、各データディスク上の同じ位置にあるデータブロック全ての XOR として計算される。他方のパリティディスクは、対角パリティブロックを格納する。対角パリティは、複数のデータディスクにわたって対角パターンを成して配置される p - 1 個のデータブロックから構成される。ブロックは、p - 1 行のストライプにグループ化される。これは、行パリティ集合へのデータブロックの割り当てには影響を及ぼさない。ただし、対角は、ある対角内のブロックが全て、同じストライプに属するようなパターンを成すように構成される。これは、対角が複数のディスクにわたって規定されるときに、大抵の対角は、ストライプ中で「循環する」ことを意味する。

30

40

【0014】

具体的には、n × (n - 1) 個のデータブロックからなるアレイの場合、アレイの端部で対角が「循環」されるとすれば、n - 1 の長さをそれぞれ有する対角が、ちょうど n 個存在する。EVENODD パリティ構成の復元において重要な点は、各対角パリティ集合が、データディスクの1つからは、何も情報を有しないことにある。ただし、パリティを格納するブロックを有する対角の他に、さらにもう1つ対角が存在する。すなわち、EVENODD パリティ構成では、ある1つの対角パリティ集合については、独立したパリティブロックを持たない。この余分な「抜けている」パリティブロックを許容するために、EVENODD 構成では、ある特定の対角に対するパリティ計算の結果は、残りの対角の

50

それぞれに関するパリティブロックにXOR演算される。

【0015】

図1は、従来のEVENODDパリティ構成に従って構成された従来のディスクアレイ100を示す略ブロック図である。各データブロック D_a, b は、パリティ集合a及びbに属する。ただし、各パリティ集合のパリティブロックは、 P_a と表記する。また、1つの特別な対角(X)については、対応するパリティブロックが格納されない。ここに、EVENODDの特徴が現れる。2つの故障からの復元を可能にするために、各データディスクは、少なくとも1つの対角パリティ集合に属してはならない。 $n \times (n - 1)$ 個のデータブロックからなる矩形アレイを採用する場合、対角パリティ集合は、 $n - 1$ 個のデータブロック要素を有する。また、上記のように、そのような構成は、全ての対角についてパリティブロックを格納する位置を持つのではない。したがって、余分な(抜けている)対角パリティブロック(X)のパリティは、その対角パリティを残りの対角パリティブロックにそれぞれXOR演算することによって記録される。具体的には、抜けているパリティ集合のパリティは、対角パリティブロック $P_4 \sim P_7$ のそれぞれにXOR演算され、それらのブロックが $P_4X \sim P_7X$ と表記される。

【0016】

2つのデータディスクの故障を復旧させる場合、まず、全てのパリティブロックのXORを求めることによって、対角パリティを持たない対角のパリティを計算する。例えば、全ての行パリティの和は、全てのデータブロックの和に等しい。全ての対角パリティの和は、全てのデータブロックの和から、抜けている対角パリティブロックの和を差し引いたものに等しい。したがって、全てのパリティブロックのXORは、全てのブロックの和(行パリティ和)から、抜けている対角を除く全てのブロックの和を差し引いたものに等しい。ここで、抜けている対角とは、要するに、抜けている対角のパリティである。実際には、各対角パリティブロックについて1つ、抜けている対角パリティの $n - 1$ 個のコピーが、その結果に加算される。 n は2よりも大きな素数であるから、 $n - 1$ は偶数であり、したがって、あるブロックをそれ自体と偶数回だけXOR演算した結果は、ゼロブロックになる。したがって、抜けているパリティを各対角パリティブロックに加えた後の対角パリティブロックの和は、その追加の対角パリティ以外の対角パリティブロックの和に等しい。

【0017】

次に、対角パリティブロックのそれぞれから、抜けている対角パリティを差し引く。2つのデータディスクが故障した後、1ブロックだけ抜けているブロックを含む対角パリティ集合が少なくとも2つ存在する。それらのパリティ集合のそれぞれにおいて抜けているブロックは、それらの対角パリティ集合のうち的一方が、たとえパリティブロックを持たない対角であったとしても、復元することができる。それらのブロックが復元されると、2つの行パリティ集合の1つを除く全ての要素が、利用可能になる。その結果、それらの行の抜けている要素の復元が可能になる。この復元は、他の対角に対しても行われ、それらの対角上の最後の抜けている幾つかのブロックを復元するための十分な情報を提供する。行パリティと対角パリティを交互に使用したこの復元このパターンは、抜けているブロックの復元が完了するまで、繰り返される。

【0018】

n は素数であるから、復元中、全ての対角に遭遇する前に、すなわち、抜けているデータブロックの復元が全て完了する前に、循環が形成されることはない。もし n が素数でなければ、これは成り立たないこともある。パリティディスクが両方とも失われた場合、データからのパリティの単純な復元を実施することができる。データディスクと対角パリティディスクが失われた場合、行パリティを使用して単純なRAID-4スタイルでデータディスクの復元が実施された後、続いて、対角パリティディスクの復元が実施される。データディスクと行パリティディスクが失われた場合、1つの対角パリティを計算できる場合がある。対角は全て同じパリティを有しているので、続いて、各対角上の抜けているブロックを計算することが出来る。

【 0 0 1 9 】

各データブロックは、いずれかの対角パリティ集合の要素であるから、2つのデータディスクが失われた場合（二重故障）、1要素しか失われてない対角パリティ集合が2つ存在する。各ディスクは、そのディスク上に表現されない対角パリティ集合を1つ有する。したがって、二重故障の場合、復元可能な対角パリティ集合は2つ存在する。EVENODDも、両方のパリティディスクの故障からの復旧、及び1つのデータディスクと1つのパリティディスクの任意の組み合わせからの復旧が可能である。この技術も、任意の単一の故障からの復旧が可能である。

【 0 0 2 0 】

EVENODD技術は、パリティ情報の量の点では最適であるが、エンコードとデコードの両方に要する計算量の点では、最適とは言えない。その理由は、抜けている対角パリティを対角パリティブロックのそれぞれに加算するために、余分な計算が必要になるからである。つまり、1ストライプ中の $p - 1$ 個のブロックは、 p 個の対角から生成される p 個のパリティブロックを保持するのに十分ではない。これを克服するために、EVENODD技術では、いずれか1つの対角のパリティを残りの全ての対角のパリティブロックにXOR演算しなければならず、その結果、計算オーバーヘッドが増大する。

10

【 0 0 2 1 】

一般に、直接的なパリティブロックを持たない対角上のデータブロックに対する小さな書き込み処理については常に、全ての対角パリティブロックを更新しなければならない。大きな書き込み処理の場合は、余分な計算が更に必要となる。本明細書では、「大きな書き込み」処理とは、1ストライプ中の全てのブロックの書き換えが必要な処理を言い、「小さな書き込み」処理とは、少なくとも1つのデータブロック、及びその関連パリティの変更が必要な処理を言う。

20

【 0 0 2 2 】

ストレージ環境によっては、データをテープその他の長期保管システムにバックアップするまでの短期記憶装置として、相当な数の、例えばニア・ライン記憶システムのような低品質ディスクドライブを使用することがよくある。しかしながら、アレイ内のディスク数が増えるにつれて、多重故障が発生する確率も増大する。この確率は、安価な記憶装置のMTTF（故障までの平均時間）が短くなるほど悪化する。つまり、ストレージシステムは三重故障、すなわち、ストレージアレイ内の3つの装置の同時故障を被る可能性がある。また、SAS（Serial Attached SCSI）、ファイバ・チャネル（FC）などのような膨大な数のストレージプロトコルが存在する結果として、ディスクシェルフのアーキテクチャは益々複雑になっており、それに伴い、ディスクシェルフが受ける故障の数も増加しており、その結果、故障したディスクシェルフに接続された各ディスクへのアクセスは失われることがある。

30

【 0 0 2 3 】

三重故障を訂正する一つの技術として、STAR技術と呼ばれるEVENODD技術の拡張がある。この技術は、2005年8月、Cheng Wang著の「Efficient and Effective Schemes for Streaming Media Delivery」に記載されており、この文献は参照により本明細書に援用される。

40

【 0 0 2 4 】

STAR技術は、 p を素数として $p + 3$ 個のディスクからなるアレイを使用し、EVENODD符号化方式を使用して、対角パリティ集合と反対角パリティ集合の両方をエンコードし、 p 個のデータディスク、1つの行パリティディスク、1つの対角パリティディスク、及び1つの反対角パリティディスクを作成する。反対角パリティ集合は、対角パリティ集合と同様に計算されるが、対角パリティ集合では傾き1を使用するのに対し、反対角パリティ集合は、傾き -1 を使用して計算される点が異なる。

【 0 0 2 5 】

STAR技術の顕著な欠点は、EVENODD符号化技術を使用する点にある。実際、STAR技術は、EVENODDを使用してアレイを符号化し、行パリティ、及び対角パ

50

リティを生成した後、EVENODD符号化を使用して第2のパリティ計算を実施し、反対角パリティを生成する。EVENODD技術の欠点は、データブロックに書き込みを行うときに、対角パリティ集合と反対角パリティ集合の両方を更新しなければならないことによって悪化する。

【0026】

STAR技術が有する更に別の顕著な欠点は、大きな素数に関わるアレイ、及び/又は最悪故障条件に関わるアレイの場合、復元処理の複雑度が非常に高くなることである。この非常に高い複雑度は、復元処理を実施するために必要となる計算回数を増加させる。

【発明の開示】

【課題を解決するための手段】

10

【0027】

[発明の概要]

本発明は、アレイ内の最大3つまでの記憶装置の同時故障からの効率的復旧が可能であるように構成されたストレージアレイに対し、パリティ計算のオーバーヘッドを低減する三重パリティ(TP)技術を含む。このTP技術は、好ましくは、複数のデータディスク、1つの行パリティディスク、1つの対角パリティディスク、及び1つの反対角パリティディスクを含むn個のディスクのような記憶装置を含むアレイにおいて実施される。ただし、pを素数として、数 $n = p + 2$ である。ディスクは、複数のブロックに分割され、ブロックは複数のストライプに編成される。ただし、各ストライプは、 $n - 3$ (又は、 $p - 1$) 行を含む。

20

【0028】

1つのストライプを形成するように選択された複数行のブロックは通常、各ディスク上で連続しているが、それは本発明の必須条件ではない。対角パリティディスクは、アレイの複数の対角パリティ集合(「対角」)に沿って計算されたパリティ情報を格納している。1つのストライプを形成するブロックは、 $n - 2$ 本の対角に編成され、各対角は、データディスク、及び行パリティディスクから $n - 3$ 個のブロックを含み、1つを除く全ての対角は、自分のパリティを対角パリティディスク上の1つのブロックに格納する。同様に、反対角パリティディスクは、アレイの複数の反対角パリティ集合(「反対角」)に沿って計算されたパリティ情報を格納している。特に、反対角は、対角に対して直交する傾きを有する。その結果、新規なTP技術は、一様なストライプ深さ、及び3台のディスク分に相当するパリティ情報の量を提供する。この量は、任意の3台のディスク故障からの復旧に必要な最小量である。

30

【0029】

本発明によれば、TP技術は、アレイ内のデータディスクの各行に沿った行パリティの計算を含み、以後、(反)対角パリティディスクに格納される対角パリティ、及び反対角パリティを計算するときには、行パリティブロックとデータブロックを区別しない。すなわち、(反)対角パリティは、全データディスク、及び行パリティディスクにわたって規定される幾つかの(反)対角に沿って計算される。また、1つを除く全ての(反)対角について、(反)対角パリティディスクにパリティが格納される。換言すれば、(反)対角パリティディスクは、1つのストライプの1つを除く各(反)対角のパリティを格納する。ただし、(反)対角パリティのうちの1つについては、パリティが計算も格納もされないが、本発明の技術によれば、アレイ内の任意の3つの同時ディスク故障から復旧するために十分なだけのパリティ情報が提供される。

40

【0030】

アレイ内の1以上の記憶装置故障に応答し、本発明は例えば、適切な復元技術を決定するために、ストレージ・オペレーティング・システムのディスクストレージ層(RAIDシステム)においてマルチステップ・ステートマシンを実施する。発明の目的のために、データディスクと行パリティディスクは、まとめて「RAID4ディスク」と呼ばれる。一台のディスクが故障した場合、故障したディスクから失われたブロックは、従来の行パリティ技術を使用して復元される。(反)対角パリティディスクがm故障した場合、デー

50

タディスクと行パリティディスクを使用して、適当な（反）対角パリティが再計算される。二重RAID4ディスク故障が発生した場合、対角パリティと反対角パリティのいずれかを使用し、行対角（R-D）パリティ復元技術にしたがって、データは復元される。

【0031】

三重ディスク故障が発生した場合、どのディスクが故障したかに関する判定がなされる。一台のRAID4ディスクに加え、対角パリティディスク、及び反対角パリティディスクが故障した場合、失われた対角パリティ、及び反対角パリティを再計算する前に、従来の行パリティディスクを使用して、故障したRAID4ディスクが復元される。2台のRAID4ディスク、及び1つの（反）対角パリティディスクが故障した場合、失われた対角パリティ、又は反対角パリティを計算する前に、R-Dパリティ技術を使用して、失われたRAID4ディスクが復元される。

10

【0032】

ただし、3台のRAID4ディスクが故障した場合、「抜けている」対角パリティ、及び反対角パリティ、すなわち、以前に格納されなかった対角パリティ、及び反対角パリティをまず計算することによって、三重パリティ復元技術が実施される。次にRAIDシステムは、多数の十字を生成することにより、故障したディスクの中間（すなわち、真ん中）にあるディスクに沿って、合計p個の四要素和を計算する。その後、この四要素和は、中間ディスク上の合計p-1個の二つ組の要素和にまで減らされる。生成された二つ組の要素和は、解放可能な一組の一次方程式を形成する。例えば、方程式の系を解くことによって、中間ディスク上のデータがまず復元される。中間ディスク上の第1のブロックが求まると、中間ディスク上のブロックが全て復元されるまで、その解は、他の式にも代入される。中間ディスクが復元された後、システムは、R-Dパリティ技術を使用して、残りの2つのディスクも復元する。

20

【0033】

有利なことに、本発明の技術によれば、故障なし条件の下でアレイに記憶されるパリティを計算する際の計算負荷を最小限に抑えることができる。また、本発明の技術は、パリティ計算のオーバーヘッドを低減し、STARのような従来の手法に比べて、所与の数のデータディスクに対する計算量も少なく済む。さらに、本発明は、行パリティブロックが全て同じディスク上に格納される集中パリティ方式を使用して実施することができ、既存のパリティ情報を再フォーマットしたり、再計算したりすることなく、データディスクを少しづつアレイに追加することができる。アレイへのデータディスクの追加に関する唯一の制限は、アレイ内で使用可能なディスクの最大数を、前述のようにして（前もって）決定しなければならないことだけである。この制限は、（反）対角の使用に起因するものであり、（反）対角の長さは、ストライプ深さによって決まる。

30

【0034】

存在する実際のディスクの数と、アレイ内のディスクの最大数との差は、例えば、全てゼロ値のデータを有する「仮想」ディスクを使用して埋められる。

【0035】

本発明の上記の利点、及びその他の利点は、添付の図面と併せて下記の説明を読めれば、より深く理解できるであろう。図中、同じ参照符号は、同一の要素、又は機能的に同じ要素を意味している。

40

【発明を実施するための最良の形態】

【0036】

[例示的实施形態の詳細な説明]

A. ストレージシステム環境

図2は、本発明とともに有利に使用されるストレージシステム220を含む環境200を示す略ブロック図である。本明細書に記載する本発明の技術は、ストレージシステム200として実施され、又はストレージシステム200を含む形で実施されるスタンドアロンのコンピュータ、又はその一部を含む、いかなるタイプのコンピュータにも適用することができ、特殊な用途のコンピュータ（例えばファイルサーバ、又はファイラ）にも、汎

50

用コンピュータにも適用することができる。また、本明細書の教示は、限定はしないが、ネットワーク・アタッチド・ストレージ環境、ストレージ・エリア・ネットワーク、あるいは、クライアント若しくはホストコンピュータに直接取り付けられたディスクアセンブリを含む、種々のストレージシステムアーキテクチャに適合する。したがって、「ストレージシステム」という用語は、記憶機能を実施するように構成され、他の装置、又はシステムに関連する何らかのサブシステムだけでなく、そうした構成も含むものとして広い意味で解釈しなければならない。

【 0 0 3 7 】

図示の実施形態において、ストレージシステム 220 は、システムバス 232 によって相互接続されたプロセッサ 222、メモリ 224、ネットワークアダプタ 225、及びストレージアダプタ 228 を含む。メモリ 224 は、本発明に関連するソフトウェアプログラムコード、及びデータ構造を格納するために、プロセッサ、及びアダプタによってアドレス指定可能な複数の記憶場所を有する。そして、プロセッサ、及びアダプタは、そのソフトウェアコードを実行し、データ構造を操作するように構成された処理要素、及び/又はロジック回路を含む。ストレージオペレーティングシステム 300 は、その一部が通常、メモリに常駐し、処理要素によって実行され、とりわけ、ストレージシステムによって実行される記憶機能を実施することによって、ストレージシステム 220 を機能的に編成する。当業者には明らかなように、本明細書に記載する本発明の技術に関連するプログラム命令の格納、及び実行には、種々のコンピュータ読取可能媒体を含む他のプロセッサや記憶手段を使用してもよい。

【 0 0 3 8 】

ネットワークアダプタ 225 は、ポイント・ツー・ポイントリンク、ワイド・エリア・ネットワーク、公共ネットワーク（インターネット）上で実施される仮想私設ネットワーク、あるいは、共有ローカルエリアネットワークを介して、ストレージシステム 220 を 1 以上のクライアント 210 に接続するように構成された複数のポートを有する。したがって、ネットワークアダプタ 225 は、ノードをネットワークに接続するために必要とされる機械的、電氣的、及び信号回路を含む。たとえば、ネットワーク 205 は、イーサネット(R)ネットワーク、またはファイバチャネル（FC）ネットワークとして実施される場合がある。各クライアント 210 は、TCP/IP のような所定のプロトコルにしたがって、ネットワーク 205 を介して個々のデータフレームやデータパケットを交換することにより、ストレージシステム 220 と通信する。

【 0 0 3 9 】

ストレージアダプタ 228 は、ストレージシステム 220 上で実行されているストレージオペレーティングシステム 300 と協働し、ユーザ（又は、クライアント）から要求された情報をアクセスする。情報は、ビデオテープ、光学、DVD、磁気テープ、バブルメモリ、電氣的ランダムアクセスメモリ、MEMS のような書換型記憶装置媒体の任意のタイプのアタッチド・アレイ、及びデータやパリティ情報のような情報を記憶するように構成された任意の他の同様の媒体に格納される。ただし、本明細書に例示的に記載されるように、情報は、好ましくは、アレイ 240 の HDD、及び/又は DASD のようなディスク 250 に格納される。ストレージアダプタは、従来の高性能 FC シリアルリンクトポロジのような I/O 相互接続構成を介してディスクを接続するための入出力（I/O）インタフェース回路を含む。

【 0 0 4 0 】

アレイ 240 への情報の格納は、ディスク空間の全体的論理構成を規定する、一群の物理記憶ディスク 250 を含む 1 以上の「ボリューム」として実施されることが好ましい。もちろん必須ではないが、各ボリュームは一般に、独自のファイルシステムに関連する。ボリューム/ファイルシステム内のディスクは通常、1 以上のグループに編成され、各グループが、RAID（Redundant Array of Independent Disks）として運用される。大半の RAID 実施形態は、所与の数の物理ディスクにわたってデータ「ストライプ」を冗長書き込みし、そのストライピングされたデータに関する適当なパリティ情報を格納するこ

とによって、データ記憶の信頼性／完全性を向上させる。

【 0 0 4 1 】

本発明は、ディスクアレイ上の行パリティ、対角パリティ、及び反対角パリティを使用した三重故障パリティ訂正復旧を提供する「三重」パリティ（TP）技術を含む。本発明の技術は、好ましくは、ストレージオペレーティングシステム 30 のディスクドライバ層（図 3 に符号 340 で示す）によって実施され、ストレージシステム内のディスクのような複数の記憶装置にわたって規定される複数のストライプにパリティを形成する方法、及びシステムを提供する。アレイ内の 3 台のディスクはパリティ専用で使用され、残りのディスクがデータを保持する。データディスク上のデータは、「クリアな状態の」で格納される。「クリアな状態」とは、格納の際にそれ以上エンコードされないことを意味する。10
任意の 1 台、2 台、又は 3 台の同時故障の後でも、データを失うことなく、アレイの内容を完全に復元することができる。本発明は、従来の方法に比べて、必要とされるパリティ情報の計算量を減らすことができるだけでなく、3 台のディスク故障からの復旧に必要とされる計算量も低減する。さらに、本発明は、均一なストライプ深さ（各ディスクが、1 ストライプあたり同数のブロックを有する）を提供し、また、任意の 3 台のディスク故障からの復旧を可能にするために必要となる最少量である、ディスク 3 台分の量のパリティ情報を提供する。

【 0 0 4 2 】

概して言えば、本発明は、 n 台の記憶装置を含む。ただし、 p は素数であり、 $n = p + 2$ である。記憶装置は、複数の同じサイズのブロックに分割される。全ての記憶装置にわたって、各記憶装置の中から $n - 3$ 個のブロックが自由に選択され、グループ化され、ストライプが形成される。ストライプ中で、1 台の記憶装置は、他の記憶装置からブロックを入力として選択することによって形成されたパリティを保持するように指定される。後で詳しく説明する単純化された形成技術の結果として、この記憶装置は、「対角パリティ装置」と呼ばれ、対角パリティ装置に保持されるパリティは「対角パリティ」と呼ばれる。同様に、各ストライプにおいて、1 台の記憶装置は、「反対角パリティ装置」としてパリティを保持するように指定され、反対角パリティ装置に保持されるパリティは「反対角パリティ」と呼ばれる。20

例えば、対角パリティと反対角パリティはアレイ上で直交する傾きを有し、例えば、対角パリティが傾き 1 を有するのに対し、反対角パリティは傾き - 1 を有する。各ストライプにおいて、そのストライプ中の（反）対角パリティ装置以外の記憶装置からそれぞれ、1 つのブロックが選択される。このブロックの集合は「行」と呼ばれる。行内の 1 つのブロックは、その行のパリティ（行パリティ）を保持するために選択され、残りのブロックはデータを保持する。行の形成は、ストライプ中の対角パリティ装置、又は反対角パリティ装置上にないブロックが全て、ちょうど 1 行に割り当てられるまで継続される。全部で $n - 3$ 行になる。30

【 0 0 4 3 】

各ストライプ中で、そのストライプ中の対角パリティ装置でも反対角パリティ装置でもない記憶装置のうち、一つを除く全ての記憶装置から、1 つのブロックが選択される。ただし、その際、選択されたブロックのうちの 2 つが同じ行に属することがないようにするという条件が課される。これは、「対角パリティ集合」、又は「対角」と呼ばれる。対角の形成は、例えば、データディスク、行パリティディスク、及び対角パリティディスクに 0 から $n - 2$ の番号を付け、行に 0 から $n - 3$ まで番号を付けた後、装置 i の行 j にあるブロックを対角 $(i + j) \bmod (n - 2)$ に割り当てることによってなされる。対角の形成は、ストライプ内の対角装置、及び反対角装置上にないブロックが全て、対角に割り当てられるまで継続される。ただし、その際、同じ装置から選択されたブロックを全く有しない対角が 2 つ存在しないようにするという条件が更に課される。 $n - 2$ 本の対角が存在し、その $n - 2$ 本の対角から、 $n - 3$ 本の対角が選択される。これらの対角上のブロックは、データを保持しているかパリティを保持しているかに関わらず、対角パリティプロ 40
50

ックを形成するために結合される $n - 3$ 個の対角パリティブロックは、ストライプ中の対角パリティを保持する装置上にある、ストライプ中の $n - 3$ 個のブロックに任意の順序で格納される。同様の技術は、パリティやブロックを「反対角パリティ集合」、すなわち「反対角」に割り当てるときにも使用される。反対角の形成は、対角パリティ装置上にも反対角パリティ装置上にもないブロックが全て、反対角に割り当てられるまで継続される。ただし、その際、同じ装置から選択されたブロックを全く有しない反対角が2つ存在しないようにするという条件が更に課される。反対角の形成は、例えば、データデバイス、行パリティデバイス、反対角パリティデバイスに0から $n - 2$ まで番号を付け、行に0から $n - 3$ の番号を付けた後、デバイス i の行 j にあるブロックを対角 $(n - 3 - i + j) \bmod (n - 2)$ に割り当てることによってなされる。

10

【0044】

本発明は、各ディスク上の同じ位置にあるブロックを含む行を選択し、 $n - 3$ 行の連続したグループを選択して複数のストライプを形成し、さらにストライプ内のブロックを選択し、各（反）対角上のブロックにより循環（反）対角パターンが形成されるようにすることで、単純に実施されることがある。さらに、本発明は、ストライプ中の全ての行パリティブロックを同じ装置に記憶することによって実施されることがある。好ましい実施形態として、本発明は、ストライプごとに、行パリティ装置、（反）対角パリティ装置、及びデータ装置のような装置の使用を同様に維持することによって実施されることがある。あるいは、本発明の他の好ましい実施形態では、行パリティ装置、（反）対角パリティ装置、及びデータ装置のような装置がストライプごとに異なるように、装置の使用を循環、その他の方法で異ならせる場合がある。

20

【0045】

パリティブロックを形成する際、パリティは一般に、データブロックの排他的論理和（XOR）として計算される。XOR演算は一般に、各入力ブロック中の同じ1ビットフィールドに対して実施され、対応する1ビットの出力を生成する。上記のように、XOR演算は、2つの1ビットフィールドにおける2の補数による加算、又は減算に等しい。また、冗長パリティ情報は、全ての入力における同じサイズの多数ビットフィールド（例えば8、16、32、64、128ビットなど）の和として計算される場合もある。例えば、パリティに相当するものは、32ビットフィールドに対して2の補数加算を使用してデータを加算することによって計算され、それぞれ32ビットの冗長情報を生成する場合がある。あるブロックをそれ自体とXOR演算したものはゼロになることから、これは、同じ入力を2回、あるブロックに対してXOR演算すれば、そのブロックの元の内容が得られるという事が信頼できないと想定される場合だけである。

30

【0046】

当業者には明らかなように、ブロック（パリティ計算のための）は、ファイルブロック、データブロック、ディスクセクタ、又は何らかの他の便利なサイズの単位に対応する場合もあれば、対応しない場合もある。パリティ計算に使用されるブロックサイズが、システム内で使用される他の何らかのブロックサイズと何らかの関係を有している必要はない。しかしながら、1以上の整数個のパリティブロックは、1以上の整数個のディスクセクタとして規定される単位にぴったりと収まるものであることが期待される。多くの場合、幾つかのブロックが、ファイルシステム、又は幾つかのデータベースブロックに対応し、通常は、 $4k$ （4096）バイト、又は2バイトのそれより大きな次数の乗数（例えば、 $8k$ 、 $16k$ 、 $32k$ 、 $64k$ 、 $128k$ 、 $256k$ ）のサイズを有する。

40

【0047】

本明細書に記載するシステムは、好ましくは、フルストライプ書き込み処理を実施する。具体的には、一般に $4k$ バイト、又は $8k$ バイトである個々のファイルブロックは、パリティ計算のときにしか使用されないより小さな複数のブロックに分割され、全ストライプの例えば、 $4k$ バイトサイズのブロックが、ディスクアレイに書き込まれる。全ストライプをディスクに書き込むとき、パリティ計算は全てメモリ上で実施され、その後、その結果がディスクに書き込まれる。したがって、ディスク上でのパリティの計算、及び更新

50

に関する負担が軽減される。

【0048】

B．ストレージオペレーティングシステム

ディスクに対するアクセスを容易にするために、ストレージオペレーティングシステム300は、仮想化モジュールと協働するwrite-anywhereファイルシステムを実施し、ディスクによって提供される記憶空間を「仮想化」する。ファイルシステムは、情報を名前付きディレクトリ、及びファイルオブジェクト（以後、「ディレクトリ」、及び「ファイル」）の階層構造としてディスク上に論理編成する。「ディスク上」の各ファイルは、データのような情報を格納するように構成されたディスクブロックの集合として実施される一方、ディレクトリは、特殊フォーマットのファイルとして実施され、その中に、名前や、他のファイル、及びディレクトリへのリンクが格納される。仮想化システムによれば、ファイルシステムは、情報を名前付きdiskの階層構造としてディスク上にさらに論理編成することが可能となり、それによって、NASシステムとSANシステムの統合アプローチ提供し、ファイルやディレクトリに対するアクセスにはファイルベースのアクセス（NAS）を可能にする一方、ファイルベースのストレージプラットフォーム上のdiskに対するアクセスには、ブロックベースのアクセス（SAN）を可能にする。

10

【0049】

例示的实施形態として、ストレージオペレーティングシステムは、カリフォルニア州サンバベルにあるネットワーク・アプライアンス・インコーポレイテッドから販売されているNetApp Data ONTAPオペレーティングシステムであることが好ましい。このオペレーティングシステムは、Write Anywhere File Layout（WAFL）ファイルシステムを実施する。ただし、当然ながら、write-in-placeファイルシステムのような任意の他のストレージオペレーティングシステムを、本明細書に記載する本発明の原理にしたがって使用されるように拡張してもよい。したがって、「ONTAP」という用語を使用した場合であっても、この用語は、本発明の教示に適合させることが可能な任意のストレージオペレーティングシステムを指すものとして広い意味で捉えなければならない。

20

【0050】

本明細書では、「ストレージオペレーティングシステム」とは、コンピュータ上で実行可能な、データアクセスを管理するためのコンピュータ実行可能コードを言い、ストレージシステムの場合、マイクロカーネルとして実施されるData ONTAPストレージオペレーティングシステムのように、データアクセスセマンティックを実施する場合がある。また、ストレージオペレーティングシステムは、UNIXやWindows NTのような汎用コンピュータ上で動作するアプリケーションプログラムとして実施してもよいし、あるいは、本明細書に記載するようなストレージアプリケーションのために構成された構成変更機能を備えた汎用オペレーティングシステムとして実施してもよい。

30

【0051】

また、当業者には明かなように、本明細書に記載する本発明の技術は、いかなるタイプの特殊目的のコンピュータ（例えば、ストレージを提供するアプライアンス）にも、汎用コンピュータにも適用することができ、ストレージシステムとして実施され、又はストレージシステムを含む形で実施されるスタンドアロンのコンピュータ、又はその一部にも適用することができる。さらに、本発明の教示は、種々のストレージシステムアーキテクチャに適合させることができ、限定はしないが例えば、ネットワーク・アタッチド・ストレージ環境、ストレージ・エリア・ネットワーク、及びクライアントやホストコンピュータに直接取り付けられるディスクアセンブリにも適合させることができる。したがって、「ストレージシステム」という用語は、ストレージ機能を実施するように構成され、他の装置、又はシステムに関連する任意のサブシステムだけでなく、それらの構成も含むものとして広い意味で解釈しなければならない。

40

【0052】

50

図 3 は、本発明とともに有利に使用されるストレージオペレーティングシステム 300 を示す略ブロック図である。ストレージオペレーティングシステムは、統合ネットワークプロトコルスタック、すなわち、より一般的には、マルチプロトコルストレージシステム上に格納された情報をクライアントがブロックアクセスプロトコルやファイルアクセスプロトコルを使用してアクセスするためのデータバスを提供するマルチプロトコルエンジンを形成するように編成された一連のソフトウェア層を含む。プロトコルスタックは、IP 層 312、並びに、その支持搬送機構である TCP 層 314、及びユーザデータグラムプロトコル (UDP) 層 316 といったネットワークプロトコル層へのインタフェースを提供するネットワークドライバ (例えば、ギガビットイーサネットドライバ) のメディアアクセス層 310 を含む。ファイルシステムプロトコル層は、マルチプロトコルファイルアクセスを提供し、その目的のために、DAFS プロトコル 318、NFS プロトコル 320、CIFS プロトコル 322、及びハイパーテキストトランスファプロトコル (HTTP) プロトコル 324 をサポートする。VI 層 326 は、VI アーキテクチャを実施し、DAFS プロトコル 318 に必要とされる RDMA のようなダイレクトアクセスポート (DAT) 機能を提供する。

【0053】

iSCSI ドライバ層 328 は、TCP/IP ネットワークプロトコル層を介したブロックプロトコルアクセスを可能にする一方、FC ドライバ層 330 は、ネットワークアダプタと協働し、ストレージシステムに対するブロックアクセス要求、及び応答の送受信を行う。FC ドライバ、及び iSCSI ドライバは、LUN (vdisk) に対する FC 固有の、及び iSCSI 固有のアクセス制御を提供し、したがって、マルチプロトコルストレージシステム上の単一の vdisk をアクセスするときに、iSCSI と FC のどちらか一方、あるいは両方への vdisk のエクスポートを管理する。さらに、ストレージオペレーティングシステムは、RAID プロトコルやディスクドライバ層 350 のようなディスクストレージプロトコルを実施する RAID システムのようなディスクストレージ層 340 を含み、ディスクドライバ層 350 は、例えば SCSI プロトコルのようなディスクアクセスプロトコルを実施する。

【0054】

本発明の例示的实施形態として、ディスクストレージ層 (例えば RAID システム 340) は、新規な TP 技術を実施する。例えば、書き込み処理の際に、RAID システム 340 は、データを以下に説明する符号化技術に従ってデータをエンコードし、記憶装置の 1 以上の故障の検出に回答して、後で詳しく説明される新規な復元技術を実施する。なお、代替実施形態では、この新規な TP 技術は、RAID システム 340 以外のストレージオペレーティングシステムのモジュールによって実施される場合もある。したがって、新規な TP 技術を実施する RAID システム 340 の説明は、単なる例として捉えなければならない。

【0055】

ディスクソフトウェア層を統合ネットワークプロトコルスタック層に橋渡しするのは、仮想化システム 355 である。仮想化システム 355 はファイルシステム 365 によって実施され、ファイルシステム 365 は、例えば vdisk モジュール 370、及び SCSI ターゲットモジュール 360 として実施される仮想化モジュールと対話する。なお、vdisk モジュール 370、ファイルシステム 365、及び SCSI ターゲットモジュール 360 は、ソフトウェアで実施しても、ハードウェアで実施しても、ファームウェアで実施しても、それらの組み合わせにより実施してもよい。vdisk モジュール 370 は、ファイルシステム 365 と対話し、システム管理者がマルチプロトコルストレージシステム 220 に対して発行したコマンドに回答して、管理者インタフェースを使用したアクセスを可能にする。実際、vdisk モジュール 370 は、とりわけ、システム管理者がユーザインタフェースを介して発行した vdisk (LUN) コマンドの複雑な組み合わせを実施することにより、SAN デプロイメントを管理する。こうした vdisk コマンドは、vdisk を実施するためのファイルシステム 365 や SCSI ターゲットモジュ

ール360と対話する原始的なファイルシステムオペレーション(「プリミティブ」)に変換される。

【0056】

次に、SCSIターゲットモジュール360は、LUNを特殊なvdiskファイルタイプに変換するマッピング手順を提供することにより、ディスク、又はLUNのエミュレーションを開始する。SCSIターゲットモジュールは、例えば、FCドライバ、iSCSIドライバ330、328と、ファイルシステム365との間に配置され、SANブロック(LUN)空間と、ファイルシステム空間(LUNがvdiskとして表現される)との間に、仮想化システム355の変換層を提供する。ファイルシステム365の上にSAN仮想化を「配置」することにより、マルチプロトコルストレージシステムは、従来のシステムによって行われるアプローチの逆を行うことができ、それによって、実質的に全てのアクセスプロトコルに対して単一の統一されたストレージプラットフォームを提供することができる。

10

【0057】

ファイルシステム365は、例えば、メッセージベースのシステムである。したがって、SCSIターゲットモジュール360は、SCSI要求を、ファイルシステムに対する操作を表わすメッセージに変換する。例えば、SCSIターゲットモジュールによって生成されるメッセージは、操作のタイプ(たとえば、読み出し、書き込み)だけでなく、ファイルシステム上に表現されるvdiskオブジェクトのパス名(例えば、パス記述子)、及びファイル名(例えば、特殊なファイル名)を含む場合がある。SCSIターゲットモジュール360は、そのメッセージを例えばファンクションコールとして、操作が実施される場所であるファイルシステム365に渡す。

20

【0058】

ファイルシステム365は例えば、例えば4キロバイト(KB)ブロックを使用し、inodeを使用してファイルを表現するブロックベースのオンディスクフォーマット表現を備えたWAFSファイルシステムを実施する。WAFSファイルシステムは、ファイルを使用して、自己のファイルシステムのレイアウトを表わすメタデータを格納する。そうしたメタデータには、とりわけ、inodeファイルがある。ディスクからinodeを読み出すために、ファイルハンドル、すなわちinode番号を含む識別子が使用される。オンディスクinode、及びinodeファイルを含む、ファイルシステムの構造の説明については、「METHOD FOR MAINTAINING CONSISTENT STATES OF A FILE SYSTEM AND FOR CREATING USER-ACCESSIBLE READ-ONLY COPIES OF A FILE SYSTEM」と題する米国特許第5,819,292号に、David Hitz他が記載している。

30

【0059】

動作的に、クライアント210からの要求は、コンピュータネットワーク205を介してパケットとしてストレージシステム220へと転送され、そこで要求はネットワークアダプタ225によって受信される。ネットワークドライバは、そのパケットを処理し、必要であれば、それをネットワークプロトコル層やファイルシステム層に渡して更なる処理を施した後、それをwrite-anywhereファイルシステム365に転送する。ここで、要求されたデータが「コア内」になれば、すなわち、メモリ224上になれば、ファイルシステムは、要求されたデータをディスク250からロードする(読み出す)ための処理を生成する。情報がメモリ上になれば、ファイルシステム365は、inode番号を使用してinodeファイル内を検索し、適当なエントリにアクセスし、論理ボリュームブロック番号(vbn)を読み出す。次にファイルシステムは、その論理vbnを含むメッセージをRAIDシステム340に渡す。論理vbnは、ディスク識別子、及びディスクブロック番号(disk、dbn)にマッピングされ、ディスクドライバシステム350の適当なドライバ(例えば、SCSI)に送られる。ディスクドライバは、指定されたディスク250からそのdbnにアクセスし、要求されたデータブロック(複数の場合もあり)をメモリ上にロードし、ストレージシステムによって処理する。要求の処理が完了すると、ストレージシステム(及び、オペレーティングシステム)は、ネッ

40

50

トワーク 205 を介してクライアント 210 に返答を返す。

【0060】

なお、ストレージシステムで受信されたクライアント要求に対し、データストレージアクセスを実施するために必要とされる、ストレージオペレーティングシステム層を通る上記のソフトウェア「パス」は、代替として、ハードウェアとして実施してもよい。すなわち、本発明の代替実施形態において、ストレージアクセス要求データパスは、フィールド・プログラマブル・ゲート・アレイ (FPGA) や特定用途向け集積回路 (ASIC) の中に論理回路として実施される場合がある。この手のハードウェア実施形態によれば、クライアント 210 によって発行される要求に応答してストレージシステム 220 が提供するストレージサービスの性能を向上させることができる。また、本発明の更に他の実施形態において、アダプタ 225、228 の処理要素は、プロセッサ 222 から、パケット処理やストレージアクセス処理の負荷の一部、又は全部を取り除くように構成され、それによって、システムによって提供されるストレージサービスの性能を向上させる場合もある。本明細書に記載する種々の処理、アーキテクチャ、及び手順は、ハードウェアで実施しても、ファームウェアで実施しても、ソフトウェアで実施してもよいものと考えられる。

【0061】

本明細書では、「ストレージオペレーティングシステム」は通常、ストレージシステムにおけるストレージ機能、例えば、データアクセスを管理するコンピュータ実行可能コードを意味し、場合によっては、ファイルシステムセマンティックを実施する場合もある。その意味で、ONTAPソフトウェアは、マイクロカーネルとして実施され、WAFLファイルシステムセマンティックを実施し、データアクセスを管理するためのWAFL層を含むストレージオペレーティングシステムの一例である。ストレージオペレーティングシステムは、UNIXやWindows NTのような汎用オペレーティングシステム上で動作するアプリケーションとして実施してもよいし、あるいは、本明細書に記載するストレージアプリケーションに合わせて構成された構成変更機能を備えた汎用オペレーティングシステムとして実施してもよい。

【0062】

さらに、当業者には明らかなように、本明細書に記載する本発明の教示は、いかなるタイプの特殊目的のコンピュータ (例えば、ファイルサーバ、ファイラ、又はストレージシステム) にも、汎用コンピュータにも適用することができ、ストレージシステム 220 として実施され、又はストレージシステム 220 を含む形で実施されるスタンドアロンのコンピュータ、又はその一部にも適用することができる。本発明とともに有利に使用されるストレージシステムの一例は、2002年8月8日に出版されたBrian Pawlowski他による「MULTI-PROTOCOL STORAGE APPLIANCE THAT PROVIDES INTEGRATED SUPPORT FOR FILE AND BLOCK ACCESS PROTOCOLS」と題する米国特許出願第10/215,917号に記載されている。また、本発明の教示は、種々のストレージシステムアーキテクチャに適合させることができ、限定はしないが例えば、ネットワーク・アタッチド・ストレージ環境、ストレージエリアネットワーク、及びクライアントやホストコンピュータに直接取り付けられたディスクアセンブリに適合させることができる。したがって、「ストレージシステム」という用語は、ストレージ機能を実施するように構成され、他の装置、又はシステムに関連する任意のサブシステムだけでなく、そうした構成も含むものとして広い意味で捉えなければならない。

【0063】

C. 三重パリティ符号化

本発明は、アレイ内の最大3つまでの記憶装置の同時故障からの効率的な復旧が可能となるように構成された、ストレージアレイに関するパリティ計算のオーバーヘッドを低減する三重パリティ (TP) 技術を含む。TP技術は、好ましくは、 p を素数として、 $p = n + 2$ 個のディスクのような記憶装置を含むアレイにおいて実施され、例えば、複数のデータディスク、1つの行パリティディスク、1つの対角パリティディスク、及び1つの反対角パリティディスクからなるアレイにおいて実施される。ディスクは複数のブロックに分

10

20

30

40

50

割され、ブロックは複数のストライプに編成され、各ストライプが、 $n - 3$ （すなわち $p - 1$ ）行を含む。1つのストライプを形成するように選択された複数行のブロックは通常、各ディスク上で連続しているが、これは、本発明にとって必須ではない。対角パリティディスクは、アレイの対角パリティ集合（「対角」）に沿って計算されたパリティ情報を記憶する。1ストライプ中のブロックは、 $n - 2$ 本の対角に編成され、各対角は、データディスクと行パリティディスクから $n - 3$ 個のブロックを含み、1つを除く全ての対角が、自己のパリティをブロックとして対角パリティディスク上に格納する。同様に、反対角パリティディスクは、アレイの反対角パリティ集合（「反対角」）に沿って計算されたパリティ情報を記憶する。特に、反対角は、対角に対して直交する傾きを有する。そのため、新規なTP技術によれば、一様なストライプ深さ、及びディスク3台分に相当するパリティ情報の量が得られる。ディスク3台分は、任意の3つのディスク故障からの復旧に必要となる最小量である。

【0064】

本発明によれば、TP技術は、アレイ内のデータディスクの各行における行パリティの計算を必要とし、その後、（反）対角パリティディスク上に格納される（反）対角パリティを計算するときに、行パリティブロックとデータブロックを区別しない。すなわち、（反）対角パリティは、全データディスク、及び行パリティディスクにわたって規定される幾つかの（反）対角に沿って計算される。また、1つを除く全ての（反）対角について、パリティが（反）対角パリティディスク上に格納される。換言すれば、（反）対角パリティディスクは、ストライプ中の1つを除く全ての（反）対角について、パリティブロックを有する。さらに、（反）対角のうちの1つについては、パリティが計算も格納もされないが、本発明の技術によれば、アレイ内の3台のディスクの同時故障からの復旧に十分なだけのパリティ情報が提供される。

【0065】

図4は、本発明の一実施形態によるTP技術を実施する手順400のステップを示すフロー図である。手順400はステップ406から始まり、ステップ410へ進み、そこでまず、素数 p に等しい数のディスクのような記憶装置でアレイを構成する。 p 個のディスクは、幾つかのデータディスク、及び1つの行パリティディスクに相当する。ステップ415において、さらにもう1つ別の対角パリティディスクをアレイに含め、アレイ全体が $p + 1$ 個のディスクから構成されるようにする。本明細書に記載されるように、対角パリティディスクは、アレイ内の全データディスク、及び行パリティディスクにわたって規定される幾つかの対角に沿って計算された対角パリティを記憶する。したがって、この時点で、アレイは、 $p - 1$ 個のデータディスク、1つの行パリティディスク、及び1つの対角パリティディスクを含む。ステップ420において、反対角パリティをアレイに追加する。対角パリティディスクと同様に、反対角パリティディスクも、アレイ内の全データディスク、及び行パリティディスクにわたって規定される幾つかの反対角に沿って計算された反対角パリティを記憶する。特に、対角と反対角は互いに直交し、例えば、傾き ± 1 を有する。したがってアレイは、 $p - 1$ 個のデータディスク、1つの行パリティディスク、1つの対角パリティディスク、及び1つの反対角パリティディスクを含み、総数 $n = p + 2$ 個のディスクを含む。ステップ425において、これらのディスクは複数のブロックに分割され、ステップ430において、ブロックはストライプに編成され、各ストライプは、 $n - 3$ 行のブロックを含む（ただし、 $n = p + 2$ ）。ステップ435では、ある行の各データブロックは、各データディスク上の同じ位置にあるデータブロックを全てXOR演算したものを保持する、その行の行パリティブロックにXOR演算される。

【0066】

次に、ステップ440において、全データブロック、及び行パリティブロックが、対角に割り当てられる。 p 個のディスクを含むアレイの場合、対角は、 $p - 1$ 行のブロックのグループに収容される。ちょうど p 本の対角が存在し、各対角は、ちょうど $p - 1$ 個のデータブロック、及び/又は行パリティブロックをXOR演算したものを含む1つの対角パリティブロックを含む。 p 個の対角集合がそれぞれ、ちょうど1つのディスクを含まない

ようにして、対角はアレイの端部で循環される。各対角は異なる 1 つのディスクを含んではならない。p - 1 行の集合内で、あらゆるディスクブロックは、p 個の対角のうちちょうど 1 つ上にある。表 1 は、0 から 4 まで番号の付いた対角を有する、p = 5 の場合のアレイの一実施形態を示している。表中の番号は、各ブロックが属する対角パリティ集合を示している。

【 0 0 6 7 】

【表 1】

0 1 2 3 4

1 2 3 4 0

2 3 4 0 1

3 4 0 1 2

10

【 0 0 6 8 】

なお、1 行中の 2 つのブロックが同じ対角パリティ集合に属することがなく、任意の 2 つのディスク故障からアレイを復元できるという性質が変わらない限り、列の順序は変更してもよく、各列における要素の位置も変更してよい。一般性を失うことなく、ブロックを対角パリティ集合に割り当てる方法は、表 1 に実質的に従うものと仮定してよい。さらに、行内のブロックの要素は、順序変更してもよい。

20

【 0 0 6 9 】

上記のように、対角パリティ集合のパリティは、対角パリティディスク上に格納される。本発明によれば、TP パリティ技術は、対角パリティディスク上に格納される対角パリティを計算するときに、行パリティブロックとデータブロックを区別しない。換言すれば、元のアレイの全てのディスクは等しく扱われ、ディスクの 1 つに格納された情報が、行パリティ集合中の他の全てのディスクの XOR から復元できるような形で扱われる。したがって、対角パリティディスクは、アレイ内の全データディスク、及び行パリティディスクにわたって規定される幾つかの対角パリティ集合に沿って計算された対角パリティを格納する（ステップ 445）。なお、RAID 5 スタイルの分散パリティ実施形態が可能となるように、データディスク、行パリティディスク、又は対角パリティディスクのようなディスクの役割は、ストライプごとに違っていてもよい。

30

【 0 0 7 0 】

しかしながら、p - 1 本の行に対して規定される p 本の対角に関するパリティ情報を全て保持するだけの十分な空間が、対角パリティディスク上には無い。具体的には、対角パリティディスク上には、p - 1 ブロック分の対角パリティを入れる空間しか存在しない。データディスク、及び行パリティディスクはそれぞれ、多くとも 1 ブロックしか対角に貢献せず、また、データブロックであるか行パリティブロックであるかに関わらず、1 つの行が、同じ対角の要素である 2 つのブロックを有することはない。1 ストライプ中にはちょうど p 本の対角が存在するが、対角パリティディスク上には p - 1 個の対角パリティブロックしか存在しない。

40

【 0 0 7 1 】

これを克服するために、対角パリティ集合のうちの 1 つについては、対角パリティを、対角パリティディスクに格納しない（ステップ 450）。すなわち、対角パリティディスクは、ストライプ中の幾つかの対角のうちの 1 つを除くそれぞれについて、パリティブロックを保持する。どの対角パリティブロックを格納しないかは、自由である。そのパリティは格納されないで、計算もされない。対角のうちの 1 つについてはパリティが格納されないが、本発明の技術によれば、反対角パリティをさらに使用することにより、アレイ内の任意の 3 つの同時ディスク故障から復旧するのに十分なだけのパリティ情報が提供される。つまり、本発明の一復旧態様によれば、アレイ内の任意の 3 台のディスクが失われたときでも、ストライプの中身を完全に復元することができる。

50

【 0 0 7 2 】

対角パリティを計算し、格納した後、R A I Dシステムは、ステップ4 5 5 ~ 4 6 5 の処理を実施し（ステップ4 4 0 ~ 4 5 0 の処理と同様に）、反対角パリティを計算し、格納する。したがって、ステップ4 5 5 において、データブロック、及び行パリティブロックは全て、反対角に割り当てられる。上記のように、反対角の傾きは - 1 であり、すなわち、対角の傾きに対して直交している。表 2 は、0 から 4 まで番号を付けた反対角を有する、 $p = 5$ の場合のアレイの一実施形態を示している。表中の数字は、各ブロックが属する反対角パリティ集合を示している。

【 0 0 7 3 】

【表 2】

10

4 3 2 1 0

0 4 3 2 1

1 0 4 3 2

2 1 0 4 3

【 0 0 7 4 】

次に、ステップ4 6 0 において、全データディスク、及び行パリティディスクにわたって規定される幾つかの反対角に沿って反対角パリティを計算し、ステップ4 6 5 において、1 つを除く全ての反対角の反対角パリティを反対角パリティディスクに格納する。そして、手順4 0 0 はステップ4 7 0 で終了する。

20

【 0 0 7 5 】

図 5 は、本発明の新規な T P 技術にしたがって編成されたディスクアレイ 5 0 0 を示すブロック図である。 n をアレイ内のディスク数とし、 $n = p + 2$ であるものと仮定する。最初の $n - 3$ 台のディスク ($D_0 \sim D_3$) はデータを保持し、ディスク $n - 2$ ($R P$) は、データディスク $D_0 \sim D_3$ に対する行パリティを保持し、ディスク $n - 1$ ($D P$) は対角パリティを保持し、ディスク n ($A D P$) は反対角パリティを保持している。この実施形態の場合、アレイ内のディスク数 n は 7 ($p = 5$) である。ディスクはブロックに分割され、ブロックはストライプにグループ化され、各ストライプは、 $n - 3$ (例えば 4) 行に相当する。また、1 つの対角あたり、 $n - 2$ (例えば 5) 個の対角が存在する。

30

【 0 0 7 6 】

各行において、各ブロックが 1 つの対角パリティ集合に属し、且つ各ブロックが異なる対角パリティ集合に属するように、データブロックと行パリティブロックに番号が付けられている。 $D_{a, b, c}$ 、及び $P_{a, b, c}$ という記述は、特定の行 (a)、対角パリティ (b)、及び反対角パリティ (c) の計算に対するデータブロック (D)、及びパリティブロック (P) それぞれの貢献を意味する。すなわち、 $D_{a, b, c}$ という記述は、そのデータブロックが、行パリティ a 、対角パリティ b 、及び反対角パリティ c の計算に使用される行、又は対角に属することを意味し、 $P_{a, b, c}$ は、そのデータブロックが、行パリティ集合 a のパリティを格納し、且つ対角パリティ集合 b 、及び反対角パリティ c に貢献することを意味する。例えば、下記のようなものである。

40

【 0 0 7 7 】

【数 1】

$$P_{0,8,10} = D_{0,4,9} \oplus D_{0,5,13} \oplus D_{0,6,12} \oplus D_{0,7,11}$$

【 0 0 7 8 】

また、特定の対角の対角パリティの計算に使用される行パリティブロックを含む記述もある。例えば、下記のようなものである。

【 0 0 7 9 】

【数 2】

$$P_4 = D_{0,4,9} \oplus D_{3,4,10} \oplus D_{2,4,13} \oplus P_{1,4,11}$$

【0080】

なお、対角パリティディスクに格納される対角パリティブロックはそれぞれ、アレイ内の他のディスク（行パリティディスクは含むが、反対角ディスクは含まない）のうちの1つを除く全てのディスクからの貢献を含む。例えば、対角パリティブロック P_4 は、 D_0 ($D_{0,4,9}$)、 D_2 ($D_{3,4,10}$)、 D_3 ($D_{2,4,13}$)、及び RP ($P_{1,4,11}$) からの貢献はあるが、 D_1 からの貢献はない。また、対角 8 (P_8) は、計算されず、対角パリティディスク DP にも格納されない。

10

【0081】

図6、及び図7は、対角、及び反対角へのブロックの割り当てをそれぞれ示すアレイの概略図である。図6は、対角へのブロックの割り当てを示すアレイ600を示し、各ブロックには、そのブロックが属する対角に対応する番号が付けられている。同図にさらに（破線で）示されているのは、アレイを符号化するときには格納されなかった抜けている対角である。同様に、図7は、反対角へのブロックの割り当てを示すアレイ700を示し、アレイに格納されなかった抜けている反対角を（破線で）示している。表1、及び表2に関して上で述べたように、単一のディスクが、同じ（反）対角のブロックを2つ有することがないようにさえすれば、（反）対角へのブロックの割り当ては、自由に変更してよい。

【0082】

20

D．ディスク故障と復元

図8は、本発明の新規なTP技術を使用するときには実施される適当な復元手順を判定する手順800のステップの詳細を示すフロー図である。上記のように、説明の都合上、「RAID4ディスク」は、データディスク、及び行パリティを意味するものとする。なお、データディスクと行パリティディスクは、RAID4構成以外の構成、例えばRAID5で構成される場合もある。手順800はステップ805から始まり、ステップ810へ進み、そこで、1以上のディスク故障が発生する。故障の原因には例えば、ディスクの完全な故障や、ディスクの一部に対するメディアエラーがある。ステップ815において、RAIDシステム340は、故障したディスクが1つであるか、2つであるか、それとも3つであるかを判定し、適当な復旧技術を使用して、故障したディスクの復旧を行う。故障したディスクが1つである場合、手順はステップ820へ分岐する。そこでシステムは、故障したディスクがRAID4ディスクであれば、従来の行パリティを使用してその失われたブロックを計算することにより、あるいは故障したディスクが（反）対角パリティディスクであれば、（反）対角パリティを計算することにより、その単一のディスクを復旧する。

30

【0083】

2台のディスクが故障した場合、手順はステップ900へ分岐し、そこで、R-Dパリティ技術を実施し、二重ディスク故障からの復旧を行う。R-Dパリティ復元技術は、行パリティと対角パリティ、あるいは行パリティと反対角パリティを復旧に使用する場合がある。4以上のディスクが故障した場合、手順800は、ステップ840で完了する前に、エラー条件により、ステップ835で終了する。

40

【0084】

一方、3台のディスクが故障した場合、ステップ845において、1つのRAID4ディスク、対角パリティディスク、及び反対角パリティディスクが故障したか否かが判定される。そうであれば、まずステップ850において、従来の行パリティ技術を使用して失われたRAID4ディスクを復元し、その後ステップ855において、対角パリティと反対角パリティを再計算することによって、アレイは復元される。そうでなければ、手順はステップ860へ進み、そこで、RAIDシステムは、2つのRAID4ディスクと、1つの（反）対角パリティディスクが故障したか否かを判定する。そうであれば、ステップ900において、システムは、R-Dパリティ復元技術を使用して、故障したRAID4

50

ディスクを復元する。この復元は、良好な（反）対角パリティを使用して実施される。すなわち、対角パリティディスクが故障した場合、R - Dパリティ復元技術は反対角パリティを使用するが、反対角パリティディスクが故障した場合、R - Dパリティ復元技術は対角パリティを使用する。ステップ900におけるR - Dパリティ復元が完了した後、ステップ870において、システムは次に、失われた（反）対角パリティを再計算する。一方、ステップ860において、3台のRAID4ディスクが故障したものと判定された場合、手順はステップ1000へと分岐し、そこでRAIDシステムは、新規な三重RAID4故障手順を実施する。

【0085】

E．行 - （反）対角復元

10

ディスクを対角に割り当てるときに、行パリティディスクとデータディスクの間に区別はないので、（反）対角パリティ集合からの復元を処理するとき、行パリティディスクとデータディスクの違いは無視することができる。例えば、任意の2台のデータディスク、あるいは、任意の1つのデータディスクと行パリティディスクが故障したものと仮定する。あるパリティ集合中の失われたブロックを復元できるのは、そのパリティ集合を構成する残りのブロックが全て、利用可能である場合だけである。XORパリティのアーチファクトは、最初にデータを保持していたかパリティを保持していたかに関わらず、全てのブロックが数学的に等価である点にある。例えば、

【0086】

【数3】

20

$$a \oplus b \oplus c = d$$

【0087】

というパリティ構成を考える。ただし、

【0088】

【数4】

$$[\oplus]$$

【0089】

は、XOR演算を表わす。式の両辺にdをXOR演算すると、

30

【数5】

$$a \oplus b \oplus c \oplus d = 0$$

【0090】

となる。したがって、復元の際に、データディスクと行パリティディスクは全て、同様に処理することができる。

【0091】

これらのディスクのそれぞれにおいて、ちょうど1つの（反）対角だけは現れない。したがって、復元は、その（反）対角の要素を含まない他のディスクから開始することができる。2つのディスクが故障しているので、大抵の場合は、（反）対角パリティデータから直ぐに復元可能なブロックが2つ存在する。これは、1ブロックだけを失った（反）対角の一方が、パリティを持たない（反）対角でない限り成り立つ。しかしながら、そのパリティ集合について多くとも1つのディスクしかデータを失っていないので、直ぐに復元可能なブロックが少なくとも1つ存在する。1つ、又は2つのブロックを（反）対角パリティから復元した後、次に、その行、又はそれらの行にある残りの失われたブロックを行パリティから復元することができる。なぜなら、この時点で、（反）対角パリティ（（反）対角パリティブロックを含まない）を使用して復元されたブロックを有する行パリティ集合から失われているのは、1ブロックだけだからである。それらのブロックを復元した後、行ブロックと同じ（反）対角（複数の場合もあり）上にある1、又は2以上のブロックを復元することができる。

40

50

【 0 0 9 2 】

このように、復元は、一連の（反）対角「移動」、及び水平「移動」によって進められる。p が素数であるから、一連の水平移動、及び（反）対角移動は全て、同じ行に 2 回遭遇するより前に、ストライプのあらゆる行に「遭遇」する。ただし、（反）対角上で（反）対角移動が全くできない（反）対角が 1 つ存在する。なぜなら、その（反）対角については、パリティが格納されていないからである。一般性を損なわずに言えば、（反）対角に 0 から $p - 1$ の番号を付した場合、パリティは、（反）対角 0 を除く全ての（反）対角について計算される。アレイ内でディスクが互いに所定距離だけ離れている場合、（反）対角 0 で復元を完了することが可能な（反）対角の固定シーケンスが常に存在する。ディスクに 0 から $p - 1$ まで番号を付し、ディスク $p - 1$ （行パリティディスク）がディスク 0 の隣にくるようなディスクの循環を想定すると、 $(p - 1)$ 個のシーケンスが考えられる。各シーケンスは、その距離だけ隔てられた任意の対を成すディスクの復元に対応する。表 3 は、例えば、 $p = 3$ の場合のシーケンスを示している。

【 0 0 9 3 】

【表 3】

Disks 1 apart: 1 2 3 4 5 6 7 8 9 10 11 12 0
 Disks 2 apart: 2 4 6 8 10 12 1 3 5 7 9 11 0
 Disks 3 apart: 3 6 9 12 2 5 8 11 1 4 7 10 0
 Disks 4 apart: 4 8 12 3 7 11 2 6 10 1 5 9 0
 Disks 5 apart: 5 10 2 7 12 4 9 1 6 11 3 8 0
 Disks 6 apart: 6 12 5 11 4 10 3 9 2 8 1 7 0
 Disks 7 apart: 7 1 8 2 9 3 10 4 11 5 12 6 0
 Disks 8 apart: 8 3 11 6 1 9 4 12 7 2 10 5 0
 Disks 9 apart: 9 5 1 10 6 2 11 7 3 12 8 4 0
 Disks 10 apart: 10 7 4 1 11 8 5 2 12 9 6 3 0
 Disks 11 apart: 11 9 7 5 3 1 12 10 8 6 4 2 0
 Disks 12 apart: 12 11 10 9 8 7 6 5 4 3 2 1 0

【 0 0 9 4 】

なお、ディスク k 離れている場合のシーケンスは、常に対角 k から開始され、その対角を毎回モジュロ p で k だけインクリメントすることによって継続され、 $p \bmod p = 0$ で終了する。また、ディスク k 離れている場合のシーケンスの最初の $p - 1$ 項は、ディスク $p - k$ 離れている場合のシーケンスの最初の $p - 1$ 項を逆にしたものである。

【 0 0 9 5 】

シーケンス上の開始位置は、どのディスク対が故障したかによって変わる。前述のようにディスクと対角に番号を付した場合、すなわち、ディスクに 0 から $n - 2$ まで順番に番号を付し、行に 0 から $n - 3$ まで順番に番号を付し、故障した各ディスク j において、ディスク j のブロック i が、対角パリティ集合 $(i + j + 1) \bmod (n - 2)$ に属している場合、失われる対角は常に、対角 j である。したがって、 k だけ離れた一対のディスクの場合、修復を開始することが可能な 2 つの対角は、 j と、 $(j + k) \bmod (n - 2)$ である。なお、ディスク k 離れている場合、復旧シーケンスにおいて、それら 2 つの対角は常に隣り合うものとなる。反対角についても、同様の計算を行うことが出来る。復旧は、2 つのシーケンス上の開始点から右へ移動することによって決定される対角のシーケンスに従って進められ、 $k < p / 2$ とした場合、ディスク k 離れている場合は記号 $(j + k) \bmod (n - 2)$ から開始され、ディスク $p - k$ 離れている場合は記号 j から開始される。したがって、2 つのデータディスクのどのような組み合わせが故障しても、また、1 つのデータディスクと行パリティパリティディスクのどのような組み合わせが故障した場

合でも、常に完全な復旧が可能である。対角パリティディスクと、もう1つ他のディスクが故障した場合は、それがデータであるか行パリティであるかに関わらず、格納されている行パリティから他のディスクを復元した後、対角パリティディスクを復元することは簡単なことである。

【0096】

なお、全てのデータブロックが、パリティの計算される（反）対角に属するとは限らない。実際、（反）対角パリティは、データブロック、及び行パリティブロックのうちの $(p-1)/p$ に対してしか計算されない。単一ブロックの書き換えは、そのブロックの行パリティを更新しなければならないだけでなく、そのブロックの（反）対角パリティも再計算しなければならないため、高くつく。また、そのブロックの行パリティを更新するときに、さらに、変化分をその行パリティブロックの（反）対角パリティブロックに加算しなければならない。ただし、1つのストライプが「1ブロック」幅であり、パリティ計算にしか使用されない幾つかのサブブロックから構成されるシステムの場合は、この計算を簡略化出来ることもある。ここで、計算されたパリティ更新は、行パリティに加算される。同じパリティ更新ブロックの幾つかの部分は、そのストライプの（反）対角パリティブロックの幾つかの部分にも直接加算される。

【0097】

ディスク（ADP）DP上の（反）対角パリティブロックは、自己のXOR計算に行パリティブロックを含める。換言すれば、ディスク（ADP）DPに格納された（反）対角パリティは、データディスクの内容に従って計算されるだけでなく、行パリティディスクの内容にも従って計算される。アレイ500に示されているように（反）対角パリティを符号化することによって、システムは、抜けている（反）対角パリティ（例えば、対角パリティの場合、対角番号8）を除く任意の2つの同時ディスク故障からの復旧が可能となる。なぜなら、行パリティブロックが、（反）対角パリティディスクDP/ADPに格納される（反）対角パリティブロックの計算に要素として含まれるからである。これに対し、従来のEVENODD技術は、行パリティブロックを対角パリティ集合の計算に要素として含めない。むしろ、従来のEVENODDアプローチは、抜けている対角パリティブロックを自己の対角パリティディスクに格納される他の対角パリティブロックのそれぞれに要素として含める。

【0098】

動作として、対角パリティディスクと、何らかのデータディスクが故障した場合、まず、そのデータディスクを行パリティディスクに基づいて復元し（例えば、従来のRAID-4復元技術にしたがって）、次に（反）対角パリティディスクを復元することによって、復旧は達成される。同様に、2つのパリティディスクが故障した場合、まず、データディスクから行パリティディスクを復元し、次に（反）対角パリティディスクを復元することによって、復旧は達成される。一方、任意の一对のデータディスクが故障した場合は、少なくとも1つの、大抵は2つの（反）対角パリティ集合から、一方のブロックを直ちに復元することができる。そしてシステムは、失われたデータブロックのうちの残りの一方を復元することができる。なお、行パリティディスクとデータディスクを失うことは、2つのデータディスクが失われた場合と全く同じであり、その復旧も、同様の仕方で達成される。

【0099】

図9は、行-（反）対角パリティのための復旧手順（復元プロセス）900に必要なとされる一連のステップを示すフロー図である。図示のように、RAIDシステム340が、二重故障が発生したものと判断すると、手順900が開始される。あるいは、手順900は、三重故障においてディスクの1つが復元された後、二重故障が残っているときに実施される場合がある。なお、手順900は、対角を使用して実施してもよいし、反対角を使用して実施してもよい。一つを除く全ての（反）対角について、DP/ADPディスク上に（反）対角パリティブロックが格納される。したがって、手順900はステップ905から開始され、ステップ910へ進み、そこで、（反）対角パリティを使用して、失われ

たブロックのうちの少なくとも一方、大抵は2つの復元を開始する。

【0100】

一方の失われたブロックを復元した後、行パリティを使用して、その行にある他方の失われたブロックを復元することにより、行の復元は完了する（ステップ915）。他方のブロックを復元するときに、ステップ920において、そのブロックが、パリティを有する（反）対角に属しているか否かの判定がなされる。そのブロックが、パリティを有する（反）対角に属している場合、（反）対角パリティを使用して、その（反）対角パリティ上にある他のディスクから、（反）対角上の他方の失われたブロックを復元することができる（ステップ925）。つまり、抜けている（反）対角を除く全ての（反）対角について、その（反）対角上にある一方のブロックが復元されれば、他方のブロックも復元することができる。次に、手順はステップ915へ戻り、そこで、行パリティ集合中の他方の失われたブロックが復元される。一方、そのブロックが、パリティを有しない（反）対角（すなわち、抜けている（反）対角）に属している場合、ステップ930において、全てのブロックが復元されたか否かに関する判定がなされる。否であれば、手順はステップ910へ戻り、そこで、まず、（反）対角パリティを利用した復元を行い、次に行パリティを利用した復元を行うというパターンが、抜けている（反）対角パリティ集合の計算に使用される最後のデータブロックに達するまで続けられる。全てのブロックの復元が完了すれば、ステップ935において手順は終了する。一連の復元手順は常に、抜けている（反）対角パリティ集合で終了する。

10

【0101】

20

要するに、復元手順は、復元される最初の（反）対角から始まって、抜けている（反）対角パリティで終わる、復元可能な（反）対角を順番に列記することによって表される。互いに k だけ離れた2つのディスク j 、及びディスク $j+k$ が故障した場合、復元可能な（反）対角の一方のシーケンスは毎回 k だけインクリメントされ、他方のシーケンスは毎回 k だけデクリメントされる。なぜなら、行復元はディスク k 個分右（又は左）へ移動し、且つ更にモジュロ p でそこから k だけ高い（又は低い）（反）対角へ移動するからである。大抵の場合、複数（例えば、少なくとも2つ）の並列復元スレッドが存在する。例外は、「抜けている」（反）対角のブロックを有しないディスクであるディスク D_0 と、（反）対角パリティディスク DP/ADP 以外のいずれかの他のディスクとが失われた場合である。その場合、復元されるブロックの流れは、失われた他のディスク上にある抜けている（反）対角パリティ集合に属するブロックで終了する単一の流れしか存在しない。

30

【0102】

F. 三重パリティ復元

アレイ内の1以上の記憶装置故障に回答し、適当な復元技術を判定するために、本発明は例えば、ストレージオペレーティングシステムのディスクストレージ層（RAIDシステム）において、マルチステップ・ステートマシンを実施する。具体的には、3つのRAID4ディスクが故障した場合、まず、「抜けている」対角パリティ、及び反対角パリティ、すなわち、先に格納されなかった対角パリティブロック、及び反対角パリティブロックを計算することにより、三重復元技術を実施する。次にRAIDシステムは、多数のクロスを生成することにより、故障したディスクのうちの中間の（すなわち、中間）ディスクに沿って、多数の4タプル和を計算する。その後、それらの4タプル和は、中間ディスク上の多数の二つ一組の和にまで低減される。生成された二つ一組の和は、解法可能な一組の一次方程式を形成する。中間ディスク上のデータは、例えば、この一組の方程式を解くことによって復元される。中間ディスク上の最初のブロックが復元された後、その解は他の方程式に代入され、これは、中間ディスク上の全てのブロックが復元されるまで継続される。中間ディスクが復元された後、次にシステムは、R-Dパリティ技術を実施し、残り2つのディスクを復旧する。

40

【0103】

図10は、本発明の一実施形態による、3つの故障したRAID4ディスクを復元するための手順1000のステップの詳細を示すフロー図である。手順1000は、ステップ

50

1 0 0 5 から開始され、ステップ 1 0 1 0 へ進み、そこで、R A I D システム 3 4 0 は、抜けている対角パリティ、及び反対角パリティを計算する。上記のように、抜けている対角、及び反対角パリティは、ディスクアレイに格納されていない（反）対角に関係する。例えば、図 6 のアレイ 6 0 0 では、第 4 の対角については、パリティが格納されていない。同様に、図 7 のアレイ 7 0 0 では、第 0 の反対角については、パリティが格納されていない。この抜けている（反）対角の計算は、比較的簡単である。抜けている（反）対角のパリティは、例えば、（反）対角パリティディスク上のブロックの和として計算される場合がある。すなわち、下記のように計算される。

【 0 1 0 4 】

【 数 6 】

10

$$\sum \oplus \text{Diagonal_Parity_Blocks} =$$

$$\sum \oplus \text{Data_Blocks} \oplus \sum \oplus \text{Data_Blocks_on_Dropped_Diagonal} \oplus$$

$$\sum \oplus \text{Row_Parity_Blocks} \oplus$$

$$\sum \oplus \text{Row_Parity_Blocks_on_Dropped_Diagonal}$$

【 0 1 0 5 】

ただし、

20

【 0 1 0 6 】

【 数 7 】

$$\sum \oplus$$

【 0 1 0 7 】

は、指定されたブロックの X O R 演算による和を表わす。

【 0 1 0 8 】

【 数 8 】

$$\sum \oplus \text{Row_Parity_Blocks} \equiv \sum \oplus \text{Data_Blocks}$$

30

【 0 1 0 9 】

であり、2つの同一の対象を X O R 演算した結果が 0 になることに留意すると、この式は、

【 0 1 1 0 】

【 数 9 】

$$\sum \oplus \text{Data_Blocks_on_Dropped_Diagonal} \oplus$$

$$\sum \oplus \text{Row_Parity_Blocks_on_Dropped_Diagonal}$$

$$= \sum \oplus \text{Blocks_on_Dropped_Diagonal}$$

40

【 0 1 1 1 】

のように低減することができる。

【 0 1 1 2 】

次に R A I D システムは、故障したディスクをアレイ内の 3 つの故障したディスクの索引値に等しい値を有する X、Y、及び Z として識別し、ディスクに 0 - p から始まるラベルを付ける。つまり、例えば、アレイ 5 0 0 のうちディスク D 0、D 1、及び D 3 が故障した場合は、X = 0、Y = 1、Z = 3 となる。次に、ステップ 1 0 1 2 においてシステムは、代数演算を行い、例えば、行、対角、及び反対角の 3 つの次元のそれぞれについて、

50

3つの故障したディスク上の失われたブロックのXORを計算する。例えば、この時点で和を計算をしておくことは、中間ディスクの復元が終わった後、残りの2つのディスクを復元するために必要となるXOR演算の回数を低減するのに役立つ。ステップ1010における抜けている/失われた(反)対角パリティの復元の結果、行パリティ集合、対角パリティ集合、及び反対角パリティ集合のそれぞれに沿って、3つのXORを計算することが可能となる。

【0113】

例えば、失われたブロック D_{00} 、 D_{01} 、及び D_{03} (最初の行にあるブロック)の行パリティ和は、下記のように計算することができる。

【0114】

【数10】

$$D_{00} \oplus D_{01} \oplus D_{03} = RP_0 \oplus D_{02}$$

【0115】

同様に、失われたブロック D_{30} 、 D_{31} 、及び D_{33} (第4の行にあるブロック)の行パリティ和は、下記のように計算することができる。

【0116】

【数11】

$$D_{30} \oplus D_{31} \oplus D_{33} = RP_3 \oplus D_{32}$$

【0117】

失われたブロック D_{00} 、 D_{11} 、及び D_{33} (図7を参照すると、これらのブロックは、反対角4上のブロックである)の反対角パリティ和は、下記のように計算することができる。

【0118】

【数12】

$$D_{00} \oplus D_{11} \oplus D_{33} = ADP_0 \oplus D_{22}$$

【0119】

失われたブロック D_{30} 、 D_{21} 、及び D_{03} (図6を参照すると、これらは、対角3上にあるブロックである)の対角パリティ和は、下記のように計算することができる。

【0120】

【数13】

$$D_{30} \oplus D_{21} \oplus D_{03} = DP_3 \oplus D_{12}$$

【0121】

次にRAIDシステムは、故障した中間ディスク上のp個の4タプル和の総計を計算する(ステップ1013~1018)。ステップ1013では、故障したディスクの行間距離を下記のように定義することにより、ディスクに順番を付ける。

$$g = Y - X$$

$$h = Z - Y$$

【0122】

したがって、 $X = 0$ 、 $Y = 1$ 、及び $Z = 3$ である上記の例を使用すると、 $g = 1 - 0 = 1$ となり、 $h = 3 - 1 = 2$ となる。この定義により、ディスクYが中間ディスクとなる。

【0123】

ステップ1014において、システムは次に、行kを選択する。例えば、 $k = 0$ であるものと仮定する。そして、システムは、その選択された行に対応する、失われたディスク上のブロックの行パリティの和を読み出す(ステップ1015)。この例では、行 $k = 0$ の場合の和は、下記のように既に計算されている。

【0124】

10

20

30

40

50

【数 1 4】

$$D_{00} \oplus D_{01} \oplus D_{03} = RP_0 \oplus D_{02}$$

【0 1 2 5】

次に、ステップ 1 0 1 6 において、システムは、ディスク Z 上の行 k にあるブロックの対角を読み出す。例えば、k = 0 であるものと仮定すると、この対角上にある失われたブロックの和は、

【0 1 2 6】

【数 1 5】

$$D_{30} \oplus D_{21} \oplus D_{03} = DP_3 \oplus D_{12}$$

10

【0 1 2 7】

となる。次に、ステップ 1 0 1 7 において、RAIDシステムは、ディスク X 上の行 k にあるブロックの反対角パリティを読み出す。例えばこれは、下記のようにになる。

【0 1 2 8】

【数 1 6】

$$D_{00} \oplus D_{11} \oplus D_{33} = ADP_0 \oplus D_{22}$$

【0 1 2 9】

(反)対角の最後の行を行 q と呼ぶものとする。次に、ステップ 1 0 1 8 において、RAIDシステムは、行 q に対応する失われたディスク上のブロックの行パリティ和を読み出す。これは例えば、下記のようにになる。

20

【0 1 3 0】

【数 1 7】

$$D_{30} \oplus D_{31} \oplus D_{33} = RP_3 \oplus D_{32}$$

【0 1 3 1】

図示の実施形態では、各ステップ 1 0 1 5、1 0 1 6、1 0 1 7 において、読み出した和が、前回の和と XOR 演算される。例えば、行 k = 0 である場合、総計は、下記のようにになる。

30

【0 1 3 2】

【数 1 8】

$$D_{00} \oplus D_{01} \oplus D_{03} \oplus D_{30} \oplus D_{31} \oplus D_{33} \oplus D_{00} \oplus D_{11} \oplus D_{33} \oplus D_{30} \oplus D_{21} \oplus D_{03} =$$

$$RP_0 \oplus D_{02} \oplus RP_3 \oplus D_{32} \oplus ADP_0 \oplus D_{22} \oplus DP_3 \oplus D_{12}$$

【0 1 3 3】

この式は、下記のように低減される。

【0 1 3 4】

【数 1 9】

40

$$D_{01} \oplus D_{11} \oplus D_{21} \oplus D_{31} = RP_0 \oplus D_{02} \oplus RP_3 \oplus D_{32} \oplus ADP_0 \oplus D_{22} \oplus DP_3 \oplus D_{12}$$

【0 1 3 5】

この式の右辺は分かっているから、この式には、中間ディスク上の 4 つの不明なものが残っている。より一般的には、各クロスにおける一番上の行と一番下の行にある重複項を削除すれば、中間ディスク上の多くとも 4 つのブロックの和が得られる。失われたデータを求めるために、異なるストライプからクロスを開始することにより、p 個のそのような和が計算される。アレイは、p - 1 行しか有していないので、第 p 番目の 4 タプル和は、ディスク Z とディスク X のそれぞれに対応する抜けている対角ディスク、及び反対角ディスクを使用し、クロスを形成することによって形成される。例えば、図 1 1 において、第

50

p 番目の 4 タプル和は、下記 4 つの X O R 和を使用して形成される。

【 0 1 3 6 】

【 数 2 0 】

$$\text{(対角)} \quad D_{40} \oplus D_{01} \oplus D_{23}$$

$$\text{(反対角)} \quad D_{43} \oplus D_{11} \oplus D_{20}$$

$$\text{(行)} \quad D_{40} \oplus D_{41} \oplus D_{43}$$

$$D_{20} \oplus D_{21} \oplus D_{23}$$

10

【 0 1 3 7 】

したがって、下記のような 4 タプル和が得られる。

【 0 1 3 8 】

【 数 2 1 】

$$D_{41} \oplus D_{01} \oplus D_{11} \oplus D_{21}$$

【 0 1 3 9 】

20

この例において、第 1 のクロス (行 0 に基づく) の結果、[0 , 1 , 2 , 3] からなるタプルが得られる。このタプルは、和の計算対象となる中間ディスク上のブロックを表わす。p 個のクロスを生成することによって、行 k に対応するタプルは、[k , k + g , k + h , k + h + g] で表わされる。ただし、加算は全てモジュロ p で実施される。

【 0 1 4 0 】

なお、4 タプル和を計算するステップは、故障したディスクの任意の順序の和として実施される。上記の例では、この順序が、X = 0、Y = 1、Z = 2 のように選択された。あるいは、別の順序 X = 0、Y = 2、及び Z = 1 を選択してもよい。この場合、中間ディスクは Y = 2 になる。その場合の値は、g = Y - X = 2、h = Z - Y = - 1 となる。3 つのディスク故障の場合、全部で 6 つの順序が可能であることが、簡単に見て取れる。各順序付けの結果、g と h に対し、異なる値の集合が生成されるため、4 タプル和を二つ一組の和に減らすために必要となるタプル数は異なる。したがって、必要となる X O R 演算の回数を最小限に抑えるために、4 タプル和を二つ一組の和に減らすために必要になるタプル集合の数が結果として最小になる順序が選択される。

30

【 0 1 4 1 】

また、このステップでは、削除列間の距離が重要となる。当然ながら、g = h であれば、中間ディスクに更に 2 つのブロックを追加することで、4 タプルから 2 タプルへの変換は不要となる。したがって、等距離 / 対称故障の場合、次のステップは不要となる。また、行「0」に対応するタプルを [0 , g , h , h + g] と表した場合、等距離故障は、条件 g = h MOD (p)、すなわち [(g - h) MOD p = 0] のように一般化される。この条件は、4 タプル中の第 2 のブロックと第 3 のブロックが同一のものであるため、削除されることを意味する。

40

【 0 1 4 2 】

4 タプル和を計算した後、ステップ 1 0 2 0 において、R A I D システムは、4 タプル和を中間ディスク上の二つ一組の和に低減する。二つ一組の和を形成するために、システムは、2 ブロックだけ残して共通のブロックを削除するように、一部の式を選択する。一部について和を計算すると、中間ディスク上に二つ一組のブロックの和が得られる。

【 0 1 4 3 】

タプルを二つ一組の和に減らすために一部のタプルを選択する方法は、多数存在する。一実施形態では、行 k に対応するタプルから開始して、最後からオフセット g (又は h)

50

の位置までにある後続のタプルを選択することによって、一部のタプルを選択する。各ステップにおいて共通のブロックは削除され、手順は、2つの不明なブロックだけが残るまで続けられる。その結果、二つ一組の和が得られる。

【0144】

例えば、行0に対応する4タプル和、すなわち $[0, g, h, g + h]$ から開始する場合、オフセット g にある次のタプルを選択すると、2つの新たなブロックを追加するとともに2つのブロックを削除することが可能となり、その結果、不明なブロックの総数をそのまま維持できることが、簡単に分かる。なぜなら、行 g に対応するタプルは $[g, 2g, h + g, 2g + h]$ であり、ブロック g とブロック $h + g$ は両方のタプルに存在するため、共通ブロック g と $h + g$ は削除できるからである（ただし、加算、及び乗算は全て、モジュロ p で行われるものと仮定する）。したがって、行0に対応する4タプルから開始する場合（これを第1ステップと呼ぶことにする）、ステップ m の結果、オフセット g にある連続したタプルが選択され、ブロック $[0, (m * g) \text{ MOD } p, h, (m * g + h) \text{ MOD } p]$ の和が得られる。

【0145】

p が素数であり、条件 $\{g, h < p\}$ が変わらないものと仮定すれば、 $[(m * g + h) \text{ MOD } p] = 0$ が成り立つような m ($0 < m < p$) が、常に見つかる。同様に、 $[(m * g - h) \text{ MOD } p] = 0$ が成り立つような m が、常に見つかる。したがって、 $[(m * g + h) \text{ MOD } p] = 0$ が成り立つような m を選択すれば、第 m のステップの後、得られる結果 $[0, (m * g) \text{ MOD } p, h, (m * g + h) \text{ MOD } p]$ において、第1のブロックと第4のブロックは削除することが出来る。あるいは、 $[(m * g - h) \text{ MOD } p] = 0$ が成り立つような m を選択すれば、第 m のステップの後、第2のブロックと第3のブロックは削除することが出来る。不明なブロックは2つしか残っていないので、タプルを選択するプロセスは、このステップで終了する。 $p - 1$ 行のそれぞれにおいて4タプル和から開始し、上記のステップを繰り返すことにより、 $p - 1$ 個の二つ一組の和が得られる。

【0146】

ステップ1025では、代数演算を実施することにより、式の一つから、既知の大きさに等しい単一の未知の値が得られる。そして、その値は、前述の式に代入され、中間ディスク上の全ての未知の値が解明され、それによって中間ディスクが復旧される。

【0147】

具体的には、アレイの形成には、 $p - 1$ 行しか使用されないので、ディスク Y 上の第 p のブロックはゼロであるものと仮定してよい。したがって、第 p のブロックと対を成してXOR演算されるブロックの値は、二つ一組の和の生成の完了時に分かっている。すなわち、式は未知の値を1つしか有しない。解を復元されたブロックに代入し、他の二つ一組の和を使用することにより、中間ディスク上の残りのブロックを復元することができる。この時点で、RAIDシステムは、故障したディスクを2つだけ残して、中間ディスクの復元を完了する。この問題は、行 - 対角パリティを使用して既に解決されているので、システムは、そのような $R - D$ パリティを実施し、失われた2つのディスクを復元する。したがって、中間ディスクの復旧が完了した後、ステップ900において、RAIDシステムは、 $R - D$ パリティを使用し、残りの2つのディスクを復元し、その後ステップ1035において終了する。

【0148】

アレイの構築に特に良好な幾つかの素数が存在する。それは、2の乗数に1を加算したものであって、且つディスクアクセスに使用されるブロックサイズよりも小さく、さらに、アレイ上に想定されるディスク数と同じ、又はそれよりも大きな素数である。2の乗数に1を加えたものである最初の幾つかの素数は、5、17、及び257である。そのうち5は、ディスクアレイに、多くても4つのデータディスクまで含めることができないので、多くの場合、小さすぎる。一方、17と257はいずれも、良い選択肢である。なぜなら、大半のストレージシステムは、ディスクストレージを通常、4k (4096) バイト

10

20

30

40

50

、 $8k(8192)$ バイト、又は他の同様の2の乗数のサイズの複数のブロックに分割するからである。最大で15個、又は255個のデータディスクを有するアレイにおいてそれぞれ、 $p=17$ 、又は $p=257$ にすることにより、16行、又は256行のグループの中で、対角パリティが計算される。これらはいずれも妥当な選択肢である。なぜなら、例えば $4k$ バイトサイズのデータブロックを $4k/16=256$ バイト、又は $4k/256=16$ バイトのサブブロックに均等に分割することができるからである。(反)対角パリティは、サブブロックに沿って(反)対角を定義することによって計算される。行パリティは、例えば $4k$ ブロック全体に対してパリティを計算することによって計算され、対角パリティを持たないRAID4、若しくはRAID5のアレイにおける計算と全く同じ方法で計算することができる。

10

【0149】

(反)対角パリティ集合の計算のために、各 $4k$ ディスクブロックを16個、又は256個のサブブロックに均等に分割することも可能であるが、例示的实施形態として、このアプローチのソフトウェア、又はハードウェアによる実施形態は、各 $4k$ ブロックのうちの、 $4k$ 対角パリティブロックと重ならない1つ、又は2つの連続した非重複領域を、1以上のサブブロックから構成される各領域にXOR演算しなければならない。データブロックの中身は、シフトパターンを成すようにして(反)対角パリティブロックにXOR演算され、抜けている(反)対角パリティ集合に属する各データブロックのサブブロックが、(反)対角パリティブロックに貢献しないようにする。(反)対角パリティデータを(反)対角パリティブロックに加算するための総計算時間は、データブロックを行パリティディスクに加算するための計算時間に匹敵することがある。

20

【0150】

有利なことに、本発明によれば、例えば、パリティ情報が全てディスクのような3つの装置に格納される、RAID4スタイルの集中パリティフォーマットの三重ディスク故障保護が提供される。したがって、本発明のパリティ技術によれば、既存のパリティ情報を再形成、すなわち再計算することなく、データディスクをディスクアレイに徐々に追加してゆくことが可能となる。本発明は、最少量の冗長ディスク空間、すなわち、アレイ一つあたりちょうど3つのディスクしか使用しない。また、本発明によれば、所与の数のデータディスクに対し、EVENODDやSTARといった従来技術のアプローチに比べて、パリティ計算のオーバーヘッドも低減される。パリティ計算オーバーヘッドは、本明細書に記載するTP技術の場合に最適なものとなる。

30

【0151】

なお、本発明のTP技術は、冗長データストリームに依存する他の用途における三重故障からの復旧にも使用される場合がある。例えば、TP技術は、データ通信アプリケーションにおいて使用されることがある。データ通信アプリケーションでは、最大で3つまでの失われた、及び/又は壊れたパケットを復元するために、追加のデータを伝送し、再送への依存性を低減する。また、さらに別の実施形態として、XOR演算以外の他の代数演算が使用される場合もある。

【0152】

上記の説明は、本発明の特定の幾つかの実施形態に関するものである。しかしながら、当業者には明らかなように、それらの利点の一部、又はは全部を獲得しつつも、記載した実施形態に対して他の変更、または修正を施すことも可能である。したがって、添付の特許請求の範囲の目的は、そうした変形や修正も、本発明の真の思想、及び範囲に入るものとしてカバーすることにある。

40

【図面の簡単な説明】

【0153】

【図1】従来のEVENODDパリティ構成に従って構成された従来のディスクアレイを示す略ブロック図である。

【図2】本発明の一実施形態によるストレージシステムを含む環境を示す略ブロック図である。

50

【図 3】本発明の一実施形態による、図 2 のストレージシステム上で使用される例示的ストレージオペレーティングシステムを示す略ブロック図である。

【図 4】本発明の一実施形態による三重パリティを符号化する手順のステップの詳細を示すフロー図である。

【図 5】本発明の一実施形態により編成されたディスクアレイを示すブロック図である。

【図 6】本発明の一実施形態による対角パリティストライプを示すディスクアレイの略ブロック図である。

【図 7】本発明の一実施形態による反対角パリティストライプを示すディスクアレイの略ブロック図である。

【図 8】本発明の一実施形態による、復元を実施する手順のステップの詳細を示すフロー図である。

10

【図 9】本発明の一実施形態による、行 - 対角 (R - D) パリティ復元を実施する手順のステップの詳細を示すフロー図である。

【図 10】本発明の一実施形態による、三重パリティ復元を実施する手順のステップの詳細を示すフロー図である。

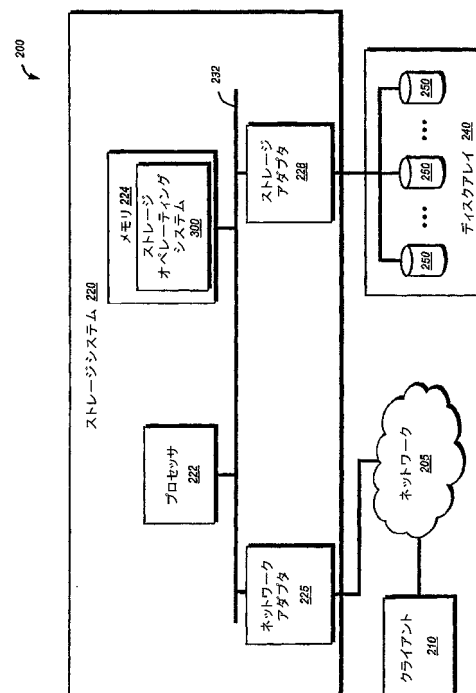
【図 11】本発明の一実施形態による、ディスク識別子を示すディスクアレイの略ブロック図である。

【図 1】

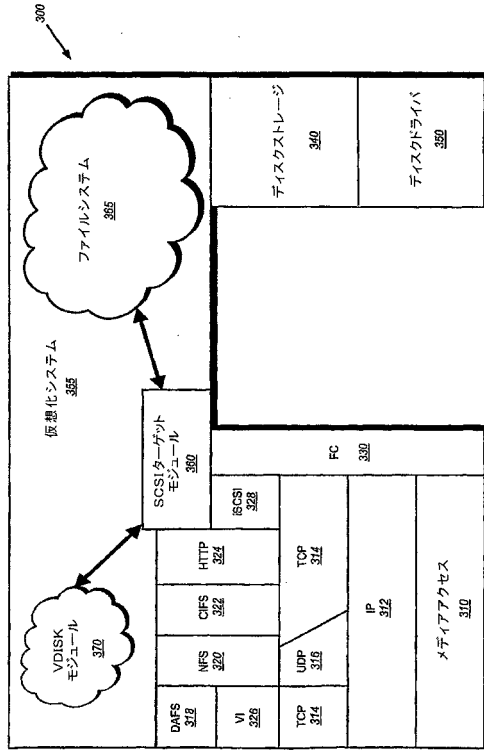
データディスク0	データディスク1	データディスク2	データディスク3	データディスク4	行パリティディスク	対角パリティディスク
D04	D05	D06	D07	D0X	P0	P0X
D15	D16	D17	D1X	D14	P1	P6X
D26	D27	D2X	D24	D25	P2	P6X
D37	D3X	D34	D35	D36	P3	P7X

(従来技術)

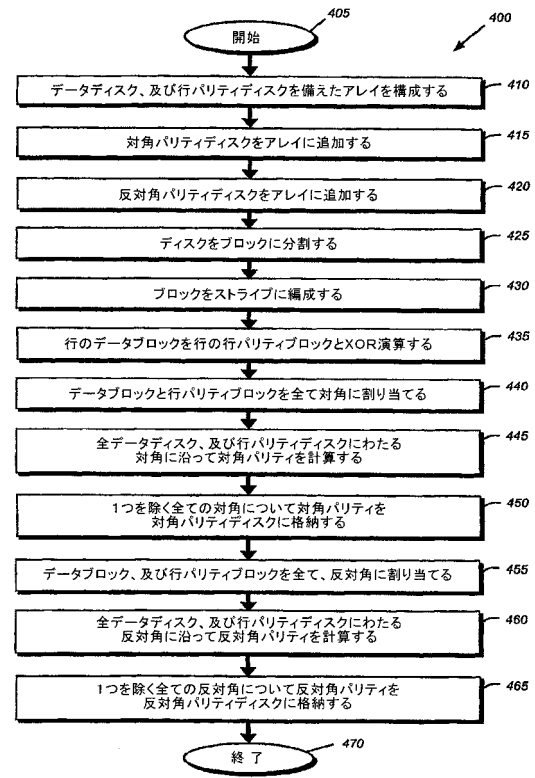
【図 2】



【図 3】



【図 4】



【図 5】

		N3				N2			
		D0	D1	D2	D3	RP	DP	ADP	
	D _{0,0}	D _{0,1}	D _{0,2}	D _{0,3}	P _{0,0}	P _{0,1}	P _{0,2}	P _{0,3}	
	D _{1,0}	D _{1,1}	D _{1,2}	D _{1,3}	P _{1,0}	P _{1,1}	P _{1,2}	P _{1,3}	
	D _{2,0}	D _{2,1}	D _{2,2}	D _{2,3}	P _{2,0}	P _{2,1}	P _{2,2}	P _{2,3}	
	D _{3,0}	D _{3,1}	D _{3,2}	D _{3,3}	P _{3,0}	P _{3,1}	P _{3,2}	P _{3,3}	

FIG. 5

【図 6】

		N3				N2			
		D0	D1	D2	D3	RP	DP	ADP	
	0	1	2	3	4	0	1	2	3
	1	2	3	4	0	1	2	3	4
	2	3	4	0	1	2	3	4	0
	3	4	0	1	2	3	4	0	1
	4	0	1	2	3	4	0	1	2

FIG. 6

【圖 7】

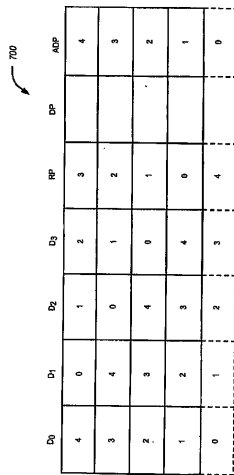
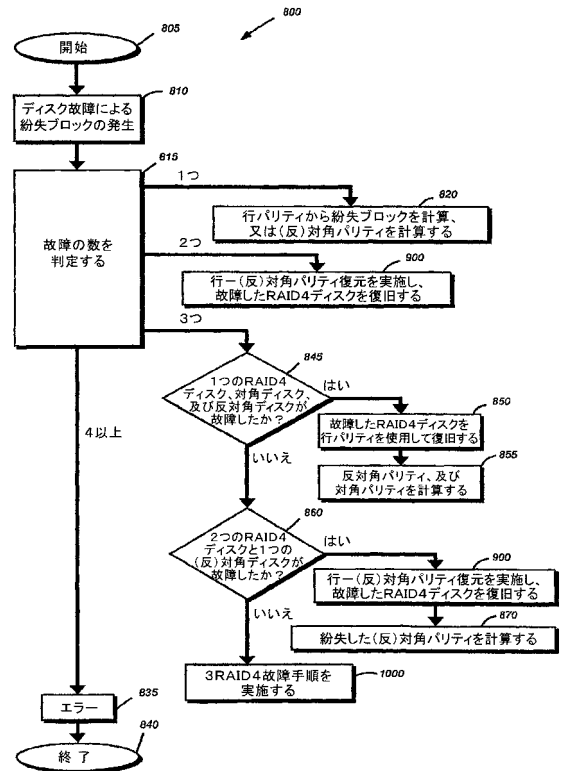
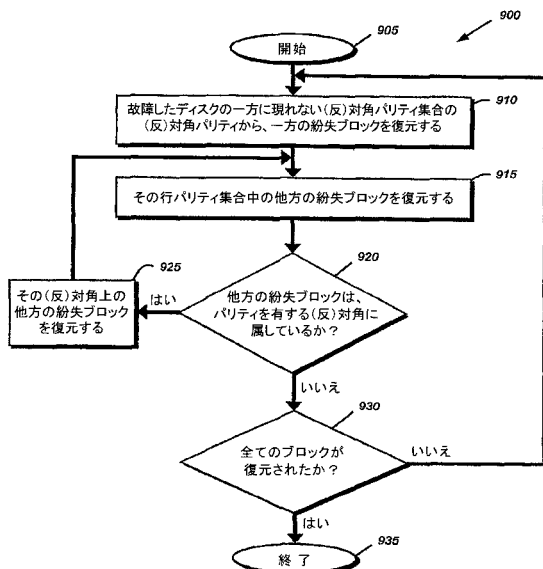


FIG. 7

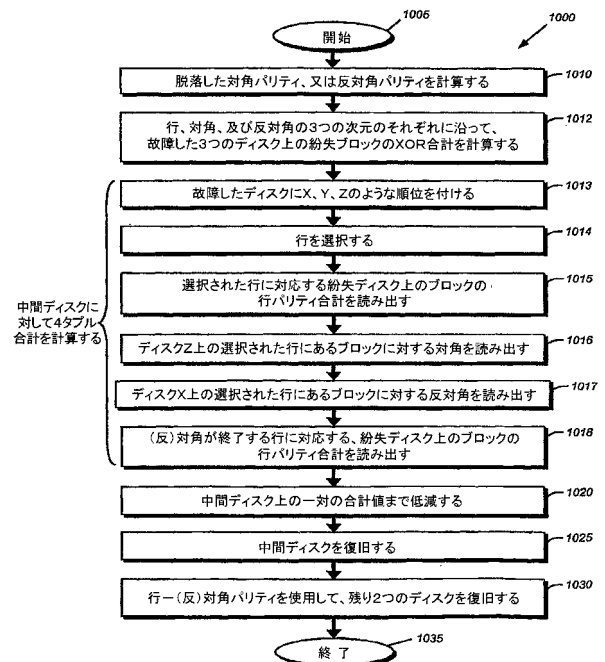
【 図 8 】



【圖 9】



【 図 1 0 】



1150

D0	D1	D2	D3	R0	DP	ADP
D00	D01	D02	D03	R00	DP0	ADP0
D04	D05	D06	D07	R01	DP1	ADP1
D08	D09	D10	D11	R02	DP2	ADP2
D12	D13	D14	D15	R03	DP3	ADP3
D16	D17	D18	D19	R04	DP4	ADP4

FIG. 11

フロントページの続き

(72)発明者 ゴエル, アトゥル

アメリカ合衆国カリフォルニア州 9 4 4 0 4 , フォスター・シティ, フォスター・シティ・ブルバード・ 1 1 1 5 , アpartment・ナンバー 4

審査官 菅原 浩二

(56)参考文献 特開 2 0 0 3 - 2 3 3 4 6 8 (J P , A)

特開平 1 1 - 0 3 9 1 0 4 (J P , A)

特開 2 0 0 5 - 1 6 6 0 1 6 (J P , A)

特表 2 0 0 6 - 5 0 5 0 3 5 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G06F 3/06

G06F 12/16