

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5124792号
(P5124792)

(45) 発行日 平成25年1月23日 (2013. 1. 23)

(24) 登録日 平成24年11月9日 (2012. 11. 9)

(51) Int. Cl.

F I

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

G 0 6 F 3 / 0 6 3 0 2 Z

G 0 6 F 3 / 0 6 5 4 0

請求項の数 22 (全 18 頁)

(21) 出願番号 特願2009-521760 (P2009-521760)
 (86) (22) 出願日 平成19年7月17日 (2007. 7. 17)
 (65) 公表番号 特表2009-545062 (P2009-545062A)
 (43) 公表日 平成21年12月17日 (2009. 12. 17)
 (86) 国際出願番号 PCT/US2007/016159
 (87) 国際公開番号 W02008/013695
 (87) 国際公開日 平成20年1月31日 (2008. 1. 31)
 審査請求日 平成22年7月9日 (2010. 7. 9)
 (31) 優先権主張番号 60/820, 180
 (32) 優先日 平成18年7月24日 (2006. 7. 24)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 11/591, 325
 (32) 優先日 平成18年11月1日 (2006. 11. 1)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 502188642
 マーベル ワールド トレード リミテッ
 ド
 バルバドス国 ビービー 1 4 0 2 7, セン
 トマイケル、ブリトンズ ヒル、ガンサイ
 トロード、エル ホライズン
 (74) 代理人 100104156
 弁理士 龍華 明裕
 (74) 代理人 100118005
 弁理士 飯山 和俊
 (74) 代理人 100143502
 弁理士 明石 英也
 (74) 代理人 100138128
 弁理士 東山 忠義

最終頁に続く

(54) 【発明の名称】 RAID (Redundant Array of Independent Disks) システム用
 のファイルサーバ

(57) 【特許請求の範囲】

【請求項 1】

ネットワークインターフェースと、
 N個のストレージアレイと、
 データ処理モジュールと

前記データ処理モジュールと前記N個のストレージアレイとの間、および、前記N個の
 ストレージアレイの間でデータを切り替えるスイッチモジュールと

を備え、

前記N個のストレージアレイのそれぞれは、

ターゲット処理モジュールと、

1個からM個のハードディスクドライブと

を有し、

MおよびNは1よりも大きい整数であり、

前記データ処理モジュールは、第1のデータブロックに対する第1のデータ格納要求を
 、前記第1のデータブロックに対するエラー検出訂正 (ECC) データ処理のために、複
 数の前記ターゲット処理モジュールのうちの1つのターゲット処理モジュールに割り当て

、

前記複数のターゲット処理モジュールのうちの前記1つのターゲット処理モジュールは
 、前記第1のデータブロックの一部を格納し、前記第1のデータブロックの第1の部分お
 よび前記第1のデータブロックに対応付けられているECCデータを、前記複数のターゲ

ット処理モジュールのうちの別の1つのターゲット処理モジュールに送信し、

前記別の1つのターゲット処理モジュールは、前記第1のデータブロックの前記第1の部分及び前記ECCデータを受信および格納し、

前記データ処理モジュールおよび前記スイッチモジュールは前記N個のストレージレイと前記ネットワークインターフェースとの間に接続され、

前記N個のストレージレイは前記データ処理モジュールおよび前記スイッチモジュールを介して前記ネットワークインターフェースと通信するRAIDシステム。

【請求項2】

前記複数のターゲット処理モジュールのうち前記別の1つのターゲット処理モジュールが受信する前記ECCデータは、前記第1の部分に対応する請求項1に記載のRAIDシステム。

10

【請求項3】

前記ネットワークインターフェースは、前記第1のデータブロックを受信して、前記第1のデータブロックを前記データ処理モジュールに転送する請求項1または2に記載のRAIDシステム。

【請求項4】

前記ネットワークインターフェースは、ギガビットイーサネットネットワークインターフェース、およびデータバスのうち少なくとも1つを有する請求項3に記載のRAIDシステム。

【請求項5】

前記スイッチモジュールは、マルチポート高速スイッチを有する

請求項1から4の何れか1項に記載のRAIDシステム。

20

【請求項6】

前記データ処理モジュールは、第2のデータブロックに対する第2のデータ格納要求を、前記第2のデータブロックに対するECCデータ処理のために、第2のターゲット処理モジュールに割り当て、前記第2のターゲット処理モジュールは、前記第2のデータブロックのうち第1の部分および前記第2のデータブロックに対応付けられているECCデータを前記複数のターゲット処理モジュールのうち第3のターゲット処理モジュールに送信する請求項1から5の何れか1項に記載のRAIDシステム。

【請求項7】

前記第1および第2のデータブロックはそれぞれ、前記複数のターゲット処理モジュールのうち前記1つのターゲット処理モジュールおよび前記第2のターゲット処理モジュールにおいて並列に処理される請求項6に記載のRAIDシステム。

30

【請求項8】

前記データ処理モジュールは、前記第1のデータブロックに対してファイルシステム(FS)プロトコル機能を実行し、

前記FSプロトコルは、ネットワークファイルサーバ(NFS)および共通インターネットファイルサーバ(CIFS)のうち少なくとも1つを含む請求項1から7の何れか1項に記載のRAIDシステム。

【請求項9】

前記データ処理モジュールは、前記第1のデータブロックに適用されるRAIDストレージレベルを決定する請求項1から8の何れか1項に記載のRAIDシステム。

40

【請求項10】

前記データ処理モジュールは、前記N個のストレージレイのうち選択される複数のストレージレイに前記第1のデータブロックをマッピングして、前記N個のストレージレイに対するストレージマップを更新する請求項1から9の何れか1項に記載のRAIDシステム。

【請求項11】

前記データ処理モジュールは、データ取得要求を受信すると、前記データ取得要求を前記複数のターゲット処理モジュールのうち第1のターゲット処理モジュールに割り当て、

50

前記複数のターゲット処理モジュールのうち前記第1のターゲット処理モジュールは、前記複数のターゲット処理モジュールのうちの他の複数のターゲット処理モジュールから前記データ取得要求に対応するデータを取得すると共に、前記データのうちエラーを含む部分に関連するECCデータを取得し、

前記複数のターゲット処理モジュールのうち前記第1のターゲット処理モジュールは、前記ECCデータを用いて前記部分に対してデータ復元を実行する請求項1から10の何れか1項に記載のRAIDシステム。

【請求項12】

前記データ処理モジュールは、データ取得要求を受信すると、前記データ取得要求に対応するデータを含む前記複数のターゲット処理モジュールにデータ取得メッセージを送信し、

10

前記複数のターゲット処理モジュールは、前記データ取得要求に対応する前記データおよび前記データのうちエラーを含む部分に関連するECCデータを取得し、

前記複数のターゲット処理モジュールは、前記データ取得要求に対応する取得された前記データおよび前記ECCデータを、前記データ処理モジュールに送信し、前記データ処理モジュールは前記ECCデータを用いて前記部分に対してデータ復元を実行する請求項1から11の何れか1項に記載のRAIDシステム。

【請求項13】

前記データ処理モジュールは、データ取得要求を受信すると、前記データ取得要求に対応するデータを含む前記複数のターゲット処理モジュールにデータ取得メッセージを送信し、

20

前記複数のターゲット処理モジュールは、前記データ取得要求に対応する前記データおよび前記データのうちエラーを含む部分に関連するECCデータを取得し、

前記複数のターゲット処理モジュールは、前記ECCデータを用いて前記部分に対してデータ復元を実行する請求項1から12の何れか1項に記載のRAIDシステム。

【請求項14】

前記データ処理モジュールは、前記スイッチモジュールと前記ネットワークインターフェースとの間に接続される請求項1から13の何れか1項に記載のRAIDシステム。

【請求項15】

前記ターゲット処理モジュールとは別物である前記データ処理モジュールは、前記RAIDシステムとは別物のホストデバイスから前記第1のデータブロックを受信する請求項1から14の何れか1項に記載のRAIDシステム。

30

【請求項16】

前記データ処理モジュールは、オペレーティングシステム、ファイルシステムプロトコルファンクションのうち少なくとも1つを用いて前記第1のデータブロックを処理し、選択された前記ターゲット処理モジュールの1つにRAID冗長処理および復元処理を割り当てる請求項1から15の何れか1項に記載のRAIDシステム。

【請求項17】

前記データ処理モジュールは、前記N個のストレージアレイに格納されたデータブロックのマップを管理し、前記ターゲット処理モジュールに前記マップを割り当てる請求項1から16の何れか1項に記載のRAIDシステム。

40

【請求項18】

前記データ処理モジュールは、前記ターゲット処理モジュールのうちの1つを選択し、選択された前記ターゲット処理モジュールのうちの1つにRAID冗長処理および復元処理を割り当てる請求項1から17の何れか1項に記載のRAIDシステム。

【請求項19】

前記データ処理モジュールは、エラー訂正コード処理を実行し、前記N個のストレージアレイへの格納後、前記第1のデータブロックがリカバーされる請求項1から18の何れか1項に記載のRAIDシステム。

【請求項20】

50

前記 R A I D ストレージのレベルは、前記 R A I D システムが分散化、ディスクミラーリング、およびパリティストレージを実行するか否かを識別する請求項 9 に記載の R A I D システム。

【請求項 2 1】

前記データ処理モジュールは、要求されたデータに対するデータ要求信号を生成する、請求項 1 に記載の R A I D システム。

【請求項 2 2】

前記 1 つのターゲット処理モジュールが、前記第 1 のデータブロックに対して E C C データを生成する間、前記データ処理モジュールは、第 2 のデータブロックを受信し、前記第 2 のデータブロックを他の前記ターゲット処理モジュールに割り当て、

10

前記他のターゲット処理モジュールは、前記第 2 のデータブロックに対して E C C データを生成する

請求項 1 から 2 1 のいずれか 1 項に記載の R A I D システム。

【発明の詳細な説明】

【関連出願】

【0 0 0 1】

本願は、米国仮特許出願第 6 0 / 8 2 0 , 1 8 0 号 (出願日 : 2 0 0 6 年 7 月 2 4 日) 、米国特許出願第 1 1 / 7 2 4 , 5 4 9 号 (出願日 : 2 0 0 7 年 3 月 1 5 日) 、および米国特許出願第 1 1 / 5 9 1 , 3 2 5 号 (出願日 : 2 0 0 6 年 1 1 月 1 日) による恩恵を主張する。上記出願の開示内容はすべて、参照により本願に組み込まれる。

20

【技術分野】

【0 0 0 2】

本開示は、R A I D (R e d u n d a n t A r r a y o f I n d e p e n d e n t D i s k s) システムに関する。

【背景技術】

【0 0 0 3】

本明細書における「背景技術」の説明は、本開示内容がどのような状況の中で考案されているのかを概して説明するためになされる。現時点において名前が挙げられている発明者による研究は、この「背景技術」部分で説明されている範囲内において、説明されていなければ出願時において先行技術としての基準を満たすものでない本明細書の側面と同様に、本開示内容に対する先行技術として、明示または黙示を問わず、認められない。

30

【0 0 0 4】

R A I D (R e d u n d a n t A r r a y o f I n d e p e n d e n t D i s k s) システムは、複数のハードディスクにデータを冗長に格納する。一部の R A I D レベルでは、データブロックを分割して複数の異なるディスクに格納して、データ格納量および検索レイテンシを低減している。また、複数のディスクを用いる場合、平均故障間隔 (M T B F) が大きくなると共にフォールトトレランスが高くなる傾向がある。

【0 0 0 5】

R A I D システムは、アクセスしているデバイスまたはホストデバイスにとっては、単一の論理ハードディスクドライブに見える。R A I D システムは、各ドライブの格納スペースを複数のユニットに分割することを含むディスクストライピングを採用するとしてもよい。ユニットサイズは、アプリケーションに応じて、セクタ (5 1 2 バイト) から数メガバイトまでの間で、選択される。すべてのディスクのストライプは通常、インターリーブされて順にアドレスが割り当てられる。

40

【0 0 0 6】

R A I D システムには、冗長性を持たないアレイ (R A I D - 0) に加えて、多くの種類が存在する。R A I D - 0 では、データに冗長性を持たせずにストライピングを採用したものである。性能は最も良好であるが、フォールトトレランスがない。R A I D - 1 では、ストライピングは用いられずにディスクミラーリングが利用され、データの格納時に複製を可能とするべく少なくとも 2 つのドライブが必要となる。各ディスクを同時に読み

50

出すことができるので、読出性能は向上する。書込性能は、単一のディスクストレージの場合と変わらない。RAID - 1は、マルチユーザシステムにおいて、性能およびフォールトトレランスが最も高い。

【0007】

RAID - 2では、複数のディスクにわたってストライピングが利用される。一部のディスクは、エラー検出訂正(ECC)情報を格納している。RAID - 3では、ストライピングが用いられて、1つのドライブがパリティ情報の格納専用とされる。埋め込まれたエラー検出訂正(ECC)情報を用いてエラーを検出する。データ復元は、残りのドライブに格納されている情報の論理的排他和(XOR)を算出することによって実行される。I/O動作ではすべてのドライブについて同時にアドレスを指定するので、RAID - 3はI/Oを重複して行うことができない。このため、RAID - 3は、記録内容が長いアプリケーションの単一ユーザシステムでの利用が最適である。

10

【0008】

RAID 4では、大きいストライプが利用される。記録内容は、どの1つのドライブからも読み出すことができる。このため、読出動作においてはI/Oを重複させて実行することができる。書込動作ではパリティドライブを更新するので、I/Oを重複させることはできない。RAID - 5では、RAID - 4の書込動作に関する制約に対処するべく、回転式パリティアレイを用いている。このため、読出動作および書込動作を共に重複させることができる。RAID - 5では、パリティ情報を格納するが、冗長データは利用しない。しかし、パリティ情報を用いてデータを再構成することができる。RAID - 5は、アレイについて、少なくとも3つ、通常は5つのディスクを必要とする。RAID - 5は、性能が重要視されないマルチユーザシステムまたは書込動作がほとんど行われないマルチユーザシステムに最適である。

20

【0009】

RAID - 6は、RAID - 5と同様であるが、複数の異なるドライブにわたって分散される第2のパリティが用いられる。RAID - 6は、フォールトトレランスおよびドライブ故障トレランスが高い。RAID - 7では、リアルタイムの埋め込みオペレーティングシステムおよびコントローラが用いられる。RAID - 7は、高速バスを介したキャッシングおよびスタンドアローンのコンピュータのほかの特性を利用する。

【0010】

30

RAID - 10は、RAID - 0とRAID - 1とを組み合わせたものである。RAID - 10はさらに2つに分類することができ、そのうち一方のRAID - 0 + 1では、データはストライプとして複数のディスクに渡って編成され、ストライプ状のディスク群をミラーリングさせる。もう一方のRAID - 1 + 0では、データをミラーリングさせて、ミラーリングの結果にストライピングを施す。

【0011】

RAID - 50(またはRAID - 5 + 0)では、一連のRAID - 5群が用いられる。RAID - 5群を、RAID - 0方式に従ってストライピングして、データ保護レベルを下げることなくRAID - 5の性能を改善する。RAID - 53(またはRAID - 5 + 3)では、ストライピング(RAID - 0方式)をRAID - 3の仮想ディスクブロックに用いる。この結果、RAID - 3に比べて性能も高くなるが、コストも同様に大きくなる。

40

【0012】

ホストデバイスがストレージに対してデータブロックを送信すると、選択されているRAID方式についてRAID処理が実施される。このようなRAID処理は、選択されているRAIDレベルおよび/またはその他の処理に対する冗長処理および復元処理(例えば、エラー検出訂正(ECC))を含むとしてもよい。

【0013】

1つの方法を挙げると、単一の中央演算処理装置(CPU)が別のデバイスからデータブロックを受信する。CPUはECCを含むすべてのRAID処理を実行する。当該方法

50

によれば、E C C 関連処理は一定でなく時間がかかり得るので、C P U によってデータストレージの速度が制限されることが多い。つまり、C P U での処理によってボトルネックが生じると共にレイテンシが大きくなってしまう場合がある。C P U が 1 つの場合、1 つのデータブロックの R A I D 再構成は、後続のデータブロックの処理の前に完了しなければならない。

【発明の概要】

【課題を解決するための手段】

【0014】

R A I D (R e d u n d a n t A r r a y o f I n d e p e n d e n t D i s k s) システムは、N 個のストレージレイを備え、N 個のストレージレイのそれぞれは、ターゲット処理モジュールと、1 個から M 個のハードディスクドライブとを有し、M および N は 1 よりも大きい整数である。データ処理モジュールは、第 1 のデータブロックに対する第 1 のデータ格納要求を、第 1 のデータブロックに対するエラー検出訂正 (E C C) データ処理のために、複数のターゲット処理モジュールのうちの 1 つのターゲット処理モジュールに割り当てる。複数のターゲット処理モジュールのうちの 1 つのターゲット処理モジュールは、第 1 のデータブロックの第 1 の部分および第 1 のデータブロックに対応付けられている E C C データを、複数のターゲット処理モジュールのうちの別の 1 つのターゲット処理モジュールに送信する。

10

【0015】

別の特徴によると、複数のターゲット処理モジュールのうち別の 1 つのターゲット処理モジュールが受信する E C C データは、第 1 の部分に対応する。インターフェースは、第 1 のデータブロックを受信して、第 1 のデータブロックをデータ処理モジュールに転送する。インターフェースは、ネットワークインターフェース、ギガビットイーサネットネットワークインターフェース、およびデータバスのうち少なくとも 1 つを有する。スイッチモジュールは、データ処理モジュールと N 個のストレージレイとの間、および、N 個のストレージレイの間でデータを切り替える。スイッチモジュールは、マルチポート高速スイッチを有する。データ処理モジュールは、第 2 のデータブロックに対する第 2 のデータ格納要求を、第 2 のデータブロックに対する E C C データ処理のために、第 2 のターゲット処理モジュールに割り当て、第 2 のターゲット処理モジュールは、第 2 のデータブロックのうち第 1 の部分および第 2 のデータブロックに対応付けられている E C C データを複数のターゲット処理モジュールのうち第 3 のターゲット処理モジュールに送信する。第 1 および第 2 のデータブロックはそれぞれ、複数のターゲット処理モジュールのうち 1 つのターゲット処理モジュールおよび第 2 のターゲット処理モジュールにおいて、重複して処理される。

20

30

【0016】

別の特徴によると、データ処理モジュールは、インターフェース、メモリ、および少なくとも 1 つのプロセッサを有する。データ処理モジュールは、第 1 のデータブロックに対してファイルシステム (F S) プロトコル機能を実行する。F S プロトコルは、ネットワークファイルサーバ (N F S) および共通インターネットファイルサーバ (C I F S) のうち少なくとも 1 つを含む。データ処理モジュールは、第 1 のデータブロックに適用される R A I D ストレージレベルを決定する。データ処理モジュールは、N 個のストレージレイのうち選択される複数のストレージレイに第 1 のデータブロックをマッピングして、N 個のストレージレイに対するストレージマップを更新する。

40

【0017】

別の特徴によると、データ処理モジュールは、データ取得要求を受信すると、データ取得要求を複数のターゲット処理モジュールのうち第 1 のターゲット処理モジュールに割り当てる。複数のターゲット処理モジュールのうち第 1 のターゲット処理モジュールは、複数のターゲット処理モジュールのほかの複数のターゲット処理モジュールからデータ取得要求に対応するデータを取得すると共に、データのうちのエラーを含む部分に関連する E C C データを取得する。

50

【 0 0 1 8 】

別の特徴によると、複数のターゲット処理モジュールのうち第1のターゲット処理モジュールは、ECCデータを用いて部分に対してデータ復元を実行する。データ処理モジュールは、データ取得要求を受信すると、データ取得要求に対応するデータを含む複数のターゲット処理モジュールにデータ取得メッセージを送信する。複数のターゲット処理モジュールは、データ取得要求に対応するデータおよびデータのうちのエラーを含む部分に関連するECCデータを取得する。複数のターゲット処理モジュールは、データ取得要求に対応する取得されたデータおよびECCデータを、データ処理モジュールに送信し、データ処理モジュールはECCデータを用いて部分に対してデータ復元を実行する。複数のターゲット処理モジュールは、ECCデータを用いて部分に対してデータ復元を実行する。

10

【 0 0 1 9 】

RAID (Redundant Array of Independent Disks) システムは、N個のストレージレイを備え、N個のストレージレイのうちそれぞれは、ターゲット処理モジュールと、1個からM個のハードディスクドライブとを有し、MおよびNは1よりも大きい整数である。データ処理モジュールは、複数のデータブロックに対するエラー検出訂正 (ECC) 処理を複数のターゲット処理モジュールのうち選択される複数のターゲット処理モジュールに、重複しないように、選択的に割り当てる。スイッチモジュールは、データ処理モジュールとN個のストレージレイとの間、および、N個のストレージレイのそれぞれとN個のストレージレイのほかの複数のストレージレイとの間において、通信経路を提供する。

20

【 0 0 2 0 】

別の特徴によると、データ処理モジュールは、第1のデータブロックに対するデータ格納要求を、第1のデータブロックに対するECCデータ処理のために、複数のターゲット処理モジュールのうち1つのターゲット処理モジュールに割り当て、複数のターゲット処理モジュールのうち1つのターゲット処理モジュールは、第1のデータブロックの第1の部分および第1のデータブロックに対応付けられているECCデータを、複数のターゲット処理モジュールのうち別の1つのターゲット処理モジュールに送信する。複数のターゲット処理モジュールのうち別の1つのターゲット処理モジュールが受信するECCデータは、第1の部分に対応する。インターフェースは、複数のデータブロックを受信して、複数のデータブロックをデータ処理モジュールに転送する。インターフェースは、ネットワークインターフェース、ギガビットイーサネットネットワークインターフェース、およびデータパスのうち少なくとも1つを有する。

30

【 0 0 2 1 】

別の特徴によると、スイッチモジュールは、マルチポート高速スイッチを有する。スイッチモジュールは、毎秒1ギガビット以上の速度で動作するマルチポートスイッチを有する。スイッチモジュールは、マルチポートギガビットイーサネットスイッチを有する。データ処理モジュールは、第2のデータブロックに対する第2のデータ格納要求を、第2のデータブロックに対するECCデータ処理のために、第2のターゲット処理モジュールに割り当てる。第2のターゲット処理モジュールは、第2のデータブロックのうち第1の部分および第2のデータブロックに対応付けられているECCデータを複数のターゲット処理モジュールのうち第3のターゲット処理モジュールに送信する。

40

【 0 0 2 2 】

別の特徴によると、第1および第2のデータブロックはそれぞれ、複数のターゲット処理モジュールのうち1つのターゲット処理モジュールおよび第2のターゲット処理モジュールにおいて、重複して処理される。データ処理モジュールは、インターフェース、メモリ、および少なくとも1つのプロセッサを有し、データ処理モジュールは、第1のデータブロックに対してファイルシステム (FS) プロトコル機能を実行する。FSプロトコルは、ネットワークファイルサーバ (NFS) および共通インターネットファイルサーバ (CIFS) のうち少なくとも1つを含む。データ処理モジュールは、複数のデータブロックに適用されるRAIDストレージレベルを決定する。データ処理モジュールは、N個の

50

ストレージレイのうち選択される複数のストレージレイに複数のデータブロックをマッピングして、N個のストレージレイに対するストレージマップを更新する。

【0023】

別の特徴によると、データ処理モジュールは、データ取得要求を受信すると、データ取得要求を複数のターゲット処理モジュールのうち第1のターゲット処理モジュールに割り当てる。複数のターゲット処理モジュールのうち第1のターゲット処理モジュールは、複数のターゲット処理モジュールのほかのターゲット処理モジュールに対してデータ取得要求に対応するデータを要求すると共にデータのうちエラーを含む部分に関連するECCデータを要求する。複数のターゲット処理モジュールのうち第1のターゲット処理モジュールは、ECCデータを用いて部分に対してデータ復元を実行する。

10

【0024】

別の特徴によると、データ処理モジュールは、データ取得要求を受信すると、データ取得要求に対応するデータを含む複数のターゲット処理モジュールにデータ取得メッセージを送信する。複数のターゲット処理モジュールは、データ取得要求に対応するデータおよびデータのうちエラーを含む部分に関連するECCデータを取得する。複数のターゲット処理モジュールは、データ取得要求に対応する取得されたデータおよびECCデータを、データ処理モジュールに送信し、データ処理モジュールはECCデータを用いて部分に対してデータ復元を実行する。複数のターゲット処理モジュールは、ECCデータを用いて部分に対してデータ復元を実行する。

【0025】

20

本開示内容はさらに異なる分野でも応用が可能であり、そのような分野は以下の詳細な説明から明らかとなる。本開示内容の好ましい実施形態を説明するためのものとして、詳細な記載および具体的な例を挙げるが、これらは説明を目的としたものに過ぎず本開示内容の範囲を限定するものではないと理解されたい。

【図面の簡単な説明】

【0026】

以下の詳細な説明および添付図面を参照して本開示内容をより詳しく説明する。添付図面は以下の通りである。

【0027】

【図1】本開示に係るRAIDシステムを示す機能ブロック図である。

30

【0028】

【図2A】データ処理モジュールを示す機能ブロック図である。

【0029】

【図2B】ターゲット処理モジュールを示す機能ブロック図である。

【0030】

【図2C】ターゲット処理モジュールの一例を示す、より詳細な機能ブロック図である。

【0031】

【図3】データブロックの処理を示す機能ブロック図である。

【0032】

【図4】データブロックの処理を示す機能ブロック図である。

40

【0033】

【図5】データブロックの処理を示す機能ブロック図である。

【0034】

【図6】データブロックの処理のタイミングを示す図である。

【0035】

【図7】格納されるべきデータブロックの処理方法を示すフローチャートである。

【0036】

【図8A】データ処理モジュールによってECC処理が実施されるデータブロックを取得する方法の一例を示すフローチャートである。

【図8B】データ処理モジュールによってECC処理が実施されるデータブロックを取得

50

する方法の一例を示すフローチャートである。

【 0 0 3 7 】

【図 9 A】ターゲット処理モジュールのそれぞれによって E C C 処理が実施されるデータブロックを取得する方法の一例を示すフローチャートである。

【図 9 B】ターゲット処理モジュールのそれぞれによって E C C 処理が実施されるデータブロックを取得する方法の一例を示すフローチャートである。

【 0 0 3 8 】

【図 1 0】ターゲット処理モジュールのうち選択された 1 つのターゲット処理モジュールによって E C C 処理が実施されるデータブロックを取得する方法の一例を示すフローチャートである。

10

【発明を実施するための形態】

【 0 0 3 9 】

以下の説明は、本質的に例示に過ぎず、本開示内容、その用途または利用を限定するものではない。本開示内容を明確に説明するべく、類似の構成要素を図中で指定する際には複数の図面にわたって同一の参照番号を使用する。本明細書で使用する場合、モジュール、回路および/またはデバイスという用語は、1 以上のソフトウェアまたはファームウェアプログラムを実行する特定用途集積回路 (A S I C)、電子回路、プロセッサ (共有、専用または群) およびメモリ、組み合わせ論理回路、ならびに/または本明細書で記載する機能を提供する上記以外に適切な構成要素を指すものとする。本明細書で言及される場合、「 A、B および C のうち少なくとも 1 つ」という表現は、論理演算 (A または B または C)、非排他的論理 O R を意味すると解釈されたい。尚、方法が含むステップは、本開示内容の原則を変更することなく、別の順序に従って実行され得る。

20

【 0 0 4 0 】

図 1 は、R A I D (R e d u n d a n t A r r a y o f I n d e p e n d e n t D i s k s) システム 1 0 0 を示す。インターフェース 1 0 4 は、R A I D システム 1 0 0 に格納されるべきデータブロックを受信する。例えば、インターフェース 1 0 4 は、ギガビットイーサネット (登録商標) ネットワークインターフェース、データバス等の高速インターフェースであってもよいが、任意のほかの種類のインターフェースを用いてもよい。データ処理モジュール 1 0 8 は、R A I D 処理の一部を実行する。つまり、データ処理モジュール 1 0 8 は、インターフェース 1 0 4 からデータブロックを受信して、当該データに対してオペレーティングシステム (O S) およびファイルシステム (F S) プロトコル機能を実行する。例えば、F S プロトコルは、ネットワークファイルサーバ (N F S)、共通インターネットファイルサーバ (C I F S)、および/または、その他の適切なプロトコルを含むとしてもよい。データ処理モジュール 1 0 8 は、後述するように、他のターゲット処理デバイスに、冗長処理および復元処理 (例えば、エラー検出訂正 (E C C)) を分配する。

30

【 0 0 4 1 】

データ処理モジュール 1 0 8 は、スイッチモジュール 1 1 2 と通信する。一例に過ぎないが、スイッチモジュール 1 1 2 は、クロスバースイッチ、ギガビットスイッチまたはギガビットイーサネットスイッチといったマルチポート高速スイッチであってもよい。スイッチモジュール 1 1 2 は、データパケットとして編成されたデータを切り替えるとしてもよい。想到されることであるが、スイッチモジュール 1 1 2 によれば、ハードワイヤード接続に比べ、拡張性および柔軟性が実現される。

40

【 0 0 4 2 】

スイッチモジュール 1 1 2 は、2 つ以上のストレージアレイ 1 2 0 - 1、1 2 0 - 2、
・ ・ ・ および 1 2 0 - X (ストレージアレイ 1 2 0 と総称する) と通信する。ここで、X
は 1 よりも大きい整数である。ストレージアレイ 1 2 0 はそれぞれ、ターゲット処理モ
ジュール 1 2 2 - 1、1 2 2 - 2、
・ ・ ・ および 1 2 2 - X (ターゲット処理モジュール 1
2 2 と総称する) と、1 以上のハードディスクドライブ (H D D) 1 2 4 - 1 1、1 2 4
- 1 2、
・ ・ ・ および 1 2 4 - X Y (H D D 1 2 4 と総称する) とを有する。ここで、Y

50

は0よりも大きい整数である。想到されることであるが、ストレージレイ120の数および各ストレージレイ120が有するHDD124の数を変化させてスケールリングを可能とするとしてもよい。

【0043】

図2Aは、データ処理モジュール108の一例をさらに詳細に示す図である。データ処理モジュール108は、インターフェース104を介して、データ格納されるべきデータブロックを受信する。データ処理モジュール108は、インターフェース150、メモリ154、および1以上のプロセッサ156を有するとしてもよい。

【0044】

データ処理モジュール108は、適用すべきRAIDストレージレベルを決定して、FS関連処理を実行して、ストレージレイに対してデータブロックをマッピングして、RAID冗長処理および復元処理（例えば、エラー検出訂正（ECC））を、選択されるターゲット処理モジュールに割り当てて、ストレージマップを更新する等としてもよい。

【0045】

RAID冗長処理および復元処理を実行するように割り当てられているターゲット処理モジュール122は、データ処理モジュール108から命令を受信する。選択ターゲット処理モジュール122は、割り当てられたデータブロックについてエラー検出訂正（ECC）を生成する。完了すると、当該ターゲット処理モジュール122は、データ処理モジュール108が与えたRAID命令に基づいて、他のレイに格納させるべく、他のターゲット処理モジュールに、データブロックおよび/またはECCデータの一部分を選択的に送信することによって、データ拡散処理を実行する。一部のデータおよびECCデータは、ローカルに格納されるとしてもよい。

【0046】

同時に、他のターゲット処理モジュール122には、他のデータブロックについてのRAID冗長処理および復元処理が割り当てられるとしてもよい。他のターゲット処理モジュール122は、重複して、他のデータブロックに対するECCを処理する。データ処理モジュール108は、格納されるべきデータブロックのいずれについてもECCを処理しないので、ボトルネックを生じさせることはない。データ処理モジュール108に対応付けられているメモリ154は、ストレージレイ120のデータのグローバルドライブマップ158を格納および更新するとしてもよい。

【0047】

図2Bを参照しつつ説明すると、各ターゲット処理モジュール122は、RAID構成モジュール170およびRAID取得モジュール168を有するとしてもよい。RAID構成モジュール170はECCを処理する。RAID取得モジュール168は、後述するように、RAID取得要求を処理する。

【0048】

RAID構成モジュール170は、選択ターゲット処理モジュール122に対応付けられているローカルドライブ124に格納されるべきデータブロックの一部分についてのECCを処理する。さらに、RAID構成モジュール170は、リモートストレージレイ120に対応付けられているリモートドライブについてのECCを処理する。RAID命令モジュール172は、別のターゲット処理についてのRAID命令を生成し、他のターゲット処理モジュール122から受信するRAID命令を処理するとしてもよい。RAID命令モジュール172は、RAID構成モジュール170と一体化されるとしてもよい。

【0049】

選択ターゲット処理モジュール以外のターゲット処理モジュールに対応付けられているリモートストレージレイ120は、選択ターゲット処理モジュールから受信するデータおよび/またはECCデータを格納する。リモートストレージレイ120は単純に、選択ターゲット処理モジュール122から送信されるRAID命令に従うとしてもよい。

【0050】

想到されるように、リモートストレージレイ 120 が実行する処理の量は、選択ターゲット処理モジュール 122 が実行する RAID 構成処理よりもはるかに少ない。このため、リモートストレージレイ 120 のターゲット処理モジュール 122 は、重複して他のデータブロックの RAID 構成を処理するべく利用され得る。

【0051】

図 2C は、ターゲット処理モジュール 122 の一例をさらに詳細に示す図である。ターゲット処理モジュール 122 は、データ処理モジュールからの RAID 構成実行要求および/またはスイッチモジュール 112 を介してリモートターゲット処理モジュールが送信する RAID 命令を受信する。ターゲット処理モジュール 122 は、インターフェース 178、メモリ 182、および 1 以上のプロセッサ 184 を有する。

10

【0052】

図 3 を参照しつつ利用について説明すると、第 1 のデータブロック 200 - 1 はインターフェース 104 を介してデータ処理モジュール 108 で受信される。データ処理モジュール 108 は、当該データブロックに対して OS および FS プロトコル機能を実行する。データ処理モジュール 108 は、当該データブロックを、ストレージレイ 120 のうちの 1 つと対応付けられているターゲット処理モジュール 122 に割り当てる。さらに、当該データ処理モジュール 108 は、適用されるべき RAID ストレージレベルを決定し、ストレージレイに対してデータブロックをマッピングして、ストレージマップを更新する等してもよい。

【0053】

20

例えば、第 1 のデータブロック 200 - 1 は第 1 のストレージレイ 120 - 1 のターゲット処理モジュール 122 - 1 に割り当てられるとしてもよい。選択ターゲット処理モジュール 122 - 1 は、当該データブロックに対して ECC を生成する。ストレージレイ 120 - 1 が第 1 のデータブロックについて ECC を生成している間、データ処理モジュール 108 はインターフェース 104 を介して第 2 のデータブロック 200 - 2 を受信する。データ処理モジュール 108 は、ECC 生成のために、ストレージレイ 120 - 2 に対応付けられているターゲット処理モジュール 122 - 2 に第 2 のデータブロックを割り当てる。

【0054】

このようにデータブロックについての RAID 構成処理を重複して実行させるが、これは、すべてのターゲット処理モジュールがデータブロックを処理するようになるまで、さらにデータブロック 200 - P まで継続することができる。このため、他の方法に比べると、スループットを大幅に向上させ得る。

30

【0055】

図 4 および図 5 は、データブロック 200 - 1 の処理をさらに詳細に示す図である。処理が終わると、データ処理モジュール 108 は、ストレージレイ 120 - 1 のターゲット処理モジュール 122 - 1 にデータブロック 200 - 1 を送信する。データ処理モジュール 108 はまた、ドライブマップを更新するとしてもよい。ターゲット処理モジュール 122 - 1 は、当該データブロックについて ECC を処理する。ターゲット処理モジュール 122 - 1 は、データブロック 200 - 1 に対応付けられているデータのうち一部を、ストレージレイ 120 - 1 に対応付けられているローカルドライブ 124 に格納するとしてもよい。さらに、ターゲット処理モジュール 122 - 1 は、RAID 命令、データ、および/または ECC データを、他のストレージレイに対応付けられているターゲット処理モジュール 122 - 2、・・・および 122 - X に送信するとしてもよい。リモートストレージレイ 120 - 2、・・・および 120 - X が有する他のターゲット処理モジュール 122 - 2、・・・および 122 - X は、RAID 命令に単純に従っており、処理負荷は制限されている。このため、リモートストレージレイ 120 - 2、・・・および 120 - X が有するターゲット処理モジュール 122 - 2、・・・および 122 - X は、他のデータブロックについて ECC を処理することができる。

40

【0056】

50

図5において、ターゲット処理モジュール122-1が第1のデータブロック200-1についてECCを処理している間に、データ処理モジュール108は第2のデータブロック200-2を受信する。データ処理モジュール108は、第2のデータブロック200-2を、ストレージレイ120-2に対応付けられているターゲット処理モジュール122-2に割り当てる。追加のデータブロック200-Pは、その他のストレージレイ120のターゲット処理モジュール122に割り当てられるとしてもよい。

【0057】

図6は、データブロックのRAID処理250の一例を概して示す図である。この種のRAID処理はボトルネックを生じさせることがあり、データアクセス時間およびデータ取得時間が短くなる。本開示の一部の実施形態に係るデータ処理を252として示す。それぞれのデータブロックに対するRAID構成に係る時間は、一定でないとしてもよい。本開示に係るRAIDシステムは、格納要求のうちの1つを処理する時間が大幅に長くなっても、データブロックの処理を継続することができる。

【0058】

図7は、データ格納要求時におけるRAIDシステム動作方法を示す図である。ステップ300から開始される。ステップ302では、格納されるべきデータブロックがデータ処理モジュール108で受信されたか否かを判断する。ステップ302が真であれば、ステップ304において、データ処理モジュール108は、当該データブロックについてのECC処理を、複数のターゲット処理モジュール122のうち1つに割り当てる。データ処理モジュール108はさらに、グローバルドライブマップを更新して、上述したほかの機能を実行するとしてもよい。ステップ306において、選択ターゲット処理モジュールは、データブロックについてECCを処理する。選択ターゲット処理モジュールは、リモートストレージレイに対応付けられるリモートターゲット処理モジュールに、RAID命令、データ、および/またはECCデータを送信するとしてもよい。ステップ310で終了する。

【0059】

図8Aから図10は、データ取得方法についてさまざまな例を示す。データ取得について説明すると、取得時にエラーが検出されるとECC処理が実行されるとしてもよい。エラーは、当該エラーに対応付けられるサブブロックを格納しているハードディスクドライブによって検出されるとしてもよい。エラーが検出されると、ECC復元は、同種の複数のターゲット処理モジュール、ターゲット処理モジュールの中から選択される1つのターゲット処理モジュール、および/または、データ処理モジュールによって、ローカルに実行されるとしてもよい。

【0060】

図8Aおよび図8Bは、データブロックを取得する方法の一例を示すフローチャートである。本実施形態によると、ECCエラーを含むデータに対するECC処理は、データ処理モジュールによって実行される。図8Aに示す方法はステップ320から開始され、ステップ322に進んで、データ処理モジュールがデータ取得要求を受信したか否か判断する。ステップ322が真であれば、ステップ324において、データ処理モジュールは、データ取得要求に対応付けられているデータを有するターゲット処理モジュールすべてに対してブロードキャストメッセージを送信する。これに代えて、データ処理モジュールは、マップを用いて各ターゲット処理モジュールに個別に個別メッセージを送信してもよい。

【0061】

ステップ326において、データ処理モジュールは、エラーを有するターゲット処理モジュールからデータブロック（およびエラーを有するデータに対する対応するECCデータ）を受信したか否かを判断する。ステップ326が真であれば、データ処理モジュールはECCデータを用いてデータを復元する。ステップ326および327はステップ328に進み、データ処理モジュールは修正されたデータを要求元に送信する。エラーが修正できない場合、データ処理モジュールはエラーメッセージを送信するとしてもよいし、お

10

20

30

40

50

よび／または、取得を再試行してもよい。ステップ 3 2 9 で終了する。

【 0 0 6 2 】

図 8 B に示す方法は、ステップ 3 3 0 から開始され、ステップ 3 3 2 に進む。ステップ 3 3 2 において、ターゲット処理モジュールは、データ処理モジュールからデータ取得要求を受信したか否かを判断する。ステップ 3 3 4 において、ターゲット処理モジュールは、取得要求に関するデータを取得してデータ処理モジュールに送信する。ステップ 3 3 6 において、ターゲット処理モジュールはサブブロックにエラーが検出されるか否かを判断する。ステップ 3 3 6 が真であれば、ターゲット処理モジュールは、サブブロックに関する E C C データを、データ処理モジュールに送信する。ステップ 3 3 6 および 3 3 7 はステップ 3 3 8 に進み、データ取得要求に関するデータがすべて送信されたか否かを判断する。送信されていない場合、ステップ 3 3 4 に戻る。ステップ 3 3 8 が真の場合は、ステップ 3 3 9 で終了する。

10

【 0 0 6 3 】

図 9 A および図 9 B は、データブロックを取得する方法の一例を示すフローチャートである。本実施形態によると、E C C 処理は、データを格納しているターゲット処理モジュールそれぞれによって実行される。図 9 A に示す方法は、ステップ 3 4 0 から開始される。ステップ 3 4 2 において、データ処理モジュールはデータ取得要求を受信したか否かを判断する。ステップ 3 4 2 が真であれば、ステップ 3 4 4 において、データ処理モジュールはすべてのターゲット処理モジュールに対してブロードキャストメッセージを送信する。これに代えて、データ処理モジュールは、マップに基づいて、ターゲット処理モジュールに対して個別メッセージを送信するとしてもよい。ステップ 3 4 8 において、データ処理モジュールはデータを受信して要求元に転送する。ステップ 3 4 9 で終了する。

20

【 0 0 6 4 】

図 9 B に示す方法は、ステップ 3 5 0 から開始される。ステップ 3 5 2 において、ターゲット処理モジュールはデータ取得要求を受信したか否かを判断する。ステップ 3 5 2 が真であれば、ステップ 3 5 4 において、ターゲット処理モジュールは、データ処理モジュールに対して、サブブロックで取得要求に関するデータを送信する。ステップ 3 5 6 において、サブブロックにおいてエラーが検出されたか否かを判断する。ステップ 3 5 6 が真であれば、ステップ 3 5 7 に進み、E C C を処理してデータを復元して、復元されたデータを送信する。データを復元できない場合には、エラーメッセージを送信するとしてもよいし、および／または、再試行してもよい。ステップ 3 5 6 および 3 5 7 からステップ 3 5 8 に進み、データ取得要求に対応付けられているサブブロックすべてが送信されたか否かを判断する。送信されていない場合、ステップ 3 5 4 に戻る。送信されている場合には、ステップ 3 5 9 で終了する。

30

【 0 0 6 5 】

図 1 0 は、E C C 処理が実施されるデータブロックを取得する方法の一例を示すフローチャートである。本実施形態によると、データ復元は、ターゲット処理モジュールのうち選択される 1 つのターゲット処理モジュールによって実行されるとしてもよい。ステップ 3 6 0 から開始され、ステップ 3 6 1 に進む。ステップ 3 6 1 において、データ処理モジュールはデータ取得要求を受信したか否かを判断する。ステップ 3 6 1 が真であれば、ステップ 3 6 2 において、データ処理モジュールは、当該データ取得をターゲット処理モジュールのうちの 1 つに割り当てる。選択ターゲット処理モジュールおよび／またはデータ処理モジュールは、リモートターゲット処理モジュールに対してデータを要求する。

40

【 0 0 6 6 】

ステップ 3 6 4 において、リモートターゲット処理モジュールは、選択ターゲット処理モジュールに対して、取得要求に関連するデータサブブロックを送信する。同様に、選択ターゲット処理モジュールは、ローカルドライブから、取得要求に関連するデータを取得する。これに代えて、リモートターゲット処理モジュールは、データにエラーがなければ、データを直接データ処理モジュールに送信するとしてもよい。エラーがあれば、リモートターゲット処理モジュールは、選択ターゲット処理モジュールにデータを送信してデー

50

タ復元を行うとしてもよい。

【 0 0 6 7 】

リモートターゲット処理モジュールのそれぞれについて説明すると、ステップ 3 6 6 において、リモートターゲット処理モジュールは、データサブブロックのうちのいずれかでエラーが検出されたか否かを判断する。ステップ 3 6 6 が真である場合、リモートターゲット処理モジュールは、エラーがあるサブブロックに対応付けられている E C C データを、選択ターゲット処理モジュールに送信する。エラーがないデータサブブロックは、ターゲット処理モジュールまたはデータ処理モジュールに送信されるとしてもよい。

【 0 0 6 8 】

ステップ 3 6 6 および 3 6 7 からステップ 3 6 8 に進み、リモートターゲット処理モジュールに関連する制御において、データサブブロックのすべてが送信されたか否かを判断する。ステップ 3 7 0 において、選択ターゲット処理モジュールは、データ復元に E C C データを用いる、つまり、エラーを修正する。選択ターゲット処理モジュールは、データ処理モジュールにデータを転送する。ステップ 3 7 2 において、データ処理モジュールは、復元されたデータを要求元のデバイスに転送する。

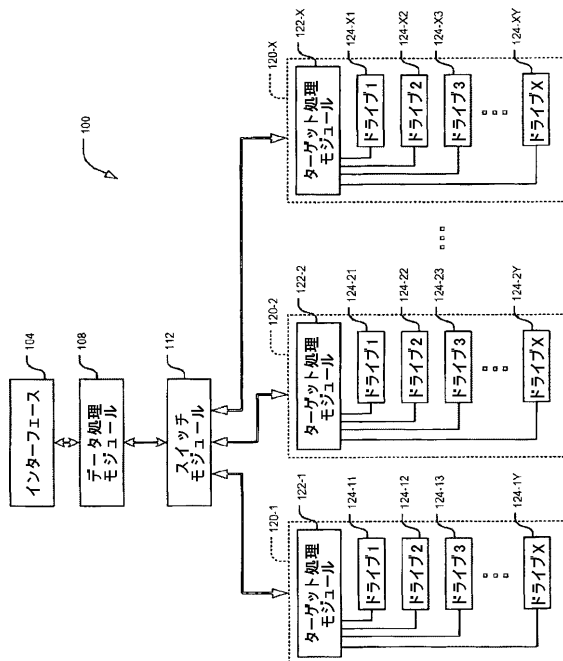
【 0 0 6 9 】

当業者であれば、上述の説明に基づき、本開示の広範囲な教示内容がさまざまな形態で実施され得ることに想到することができる。このため、本開示では具体的な例を紹介したが、本願の図面、明細書および特許請求の範囲を参照することによって当業者はほかの変形例を提案できることが明らかであるので、本開示の真の範囲は記載された具体例に限定

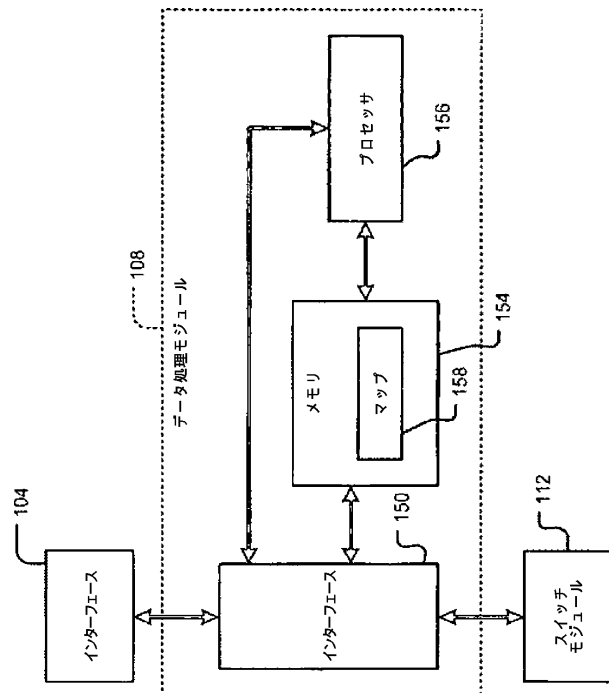
10

20

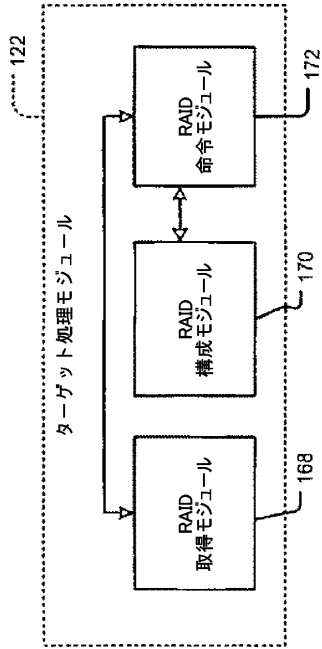
【 図 1 】



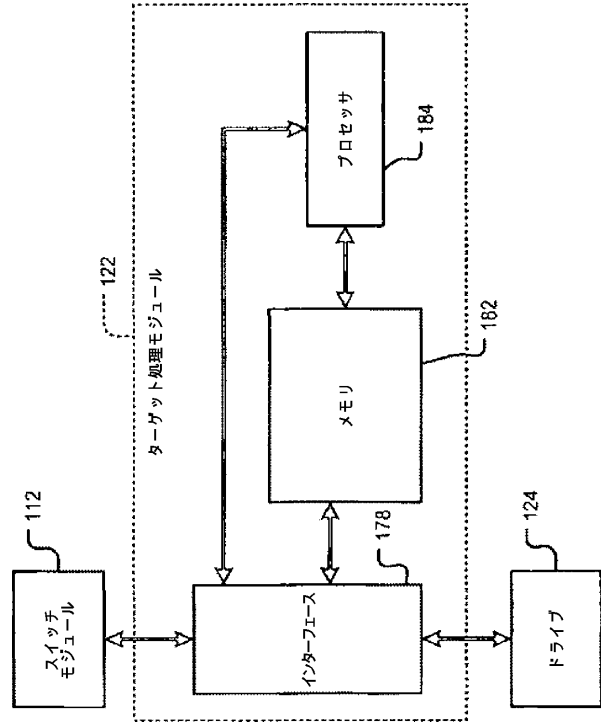
【 図 2 A 】



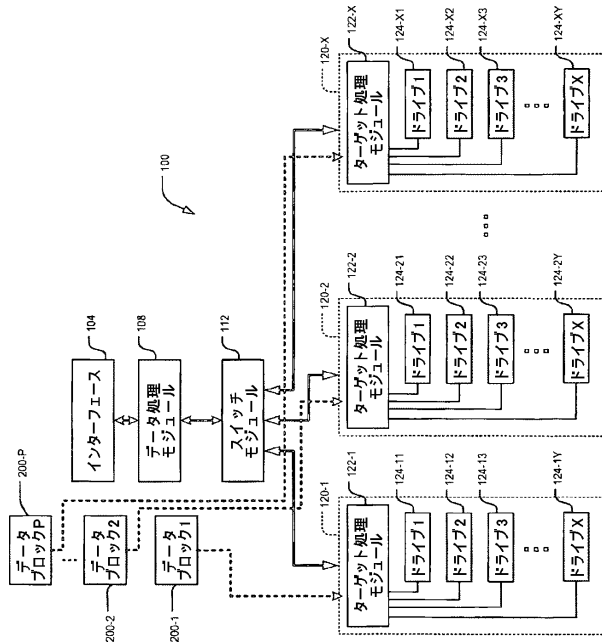
【図 2 B】



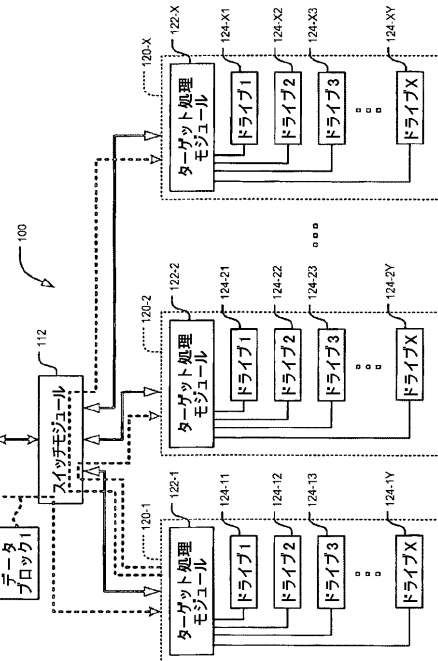
【図 2 C】



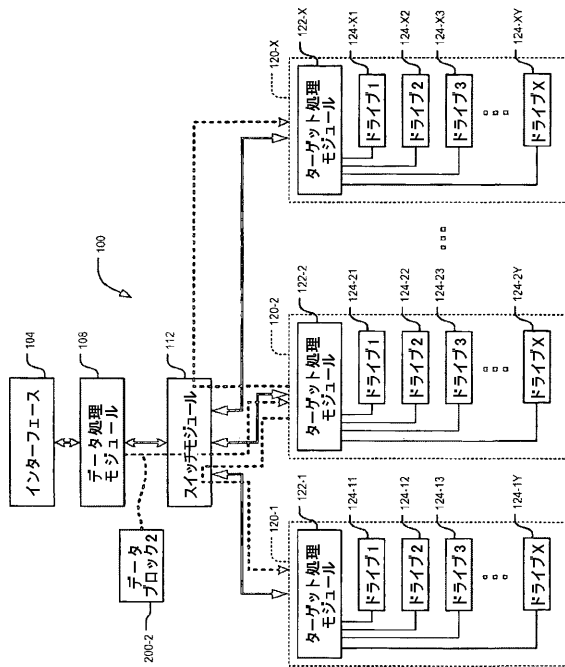
【図 3】



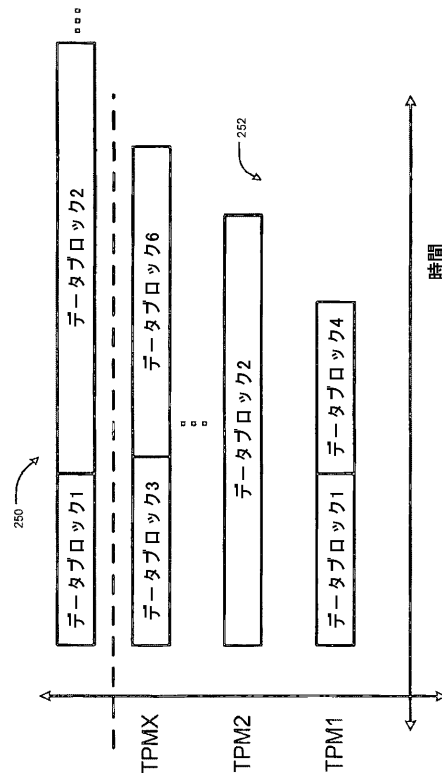
【図 4】



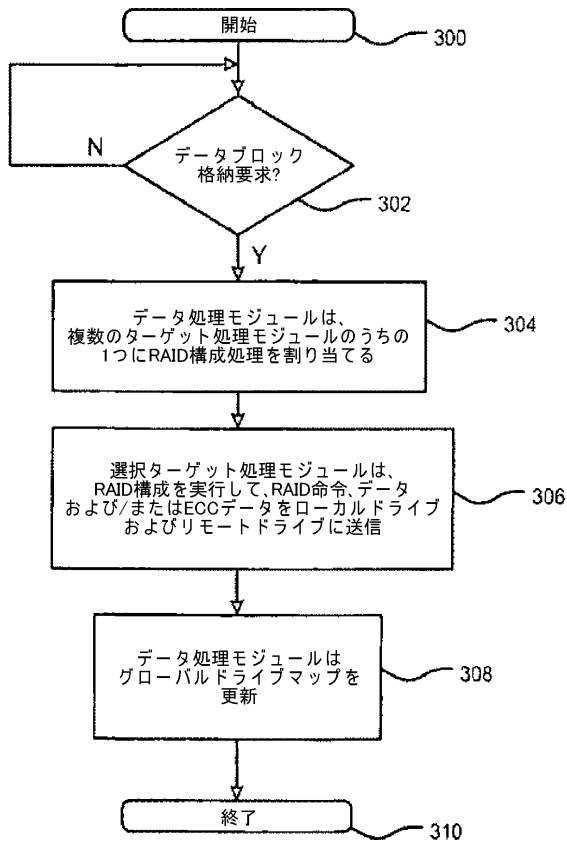
【図 5】



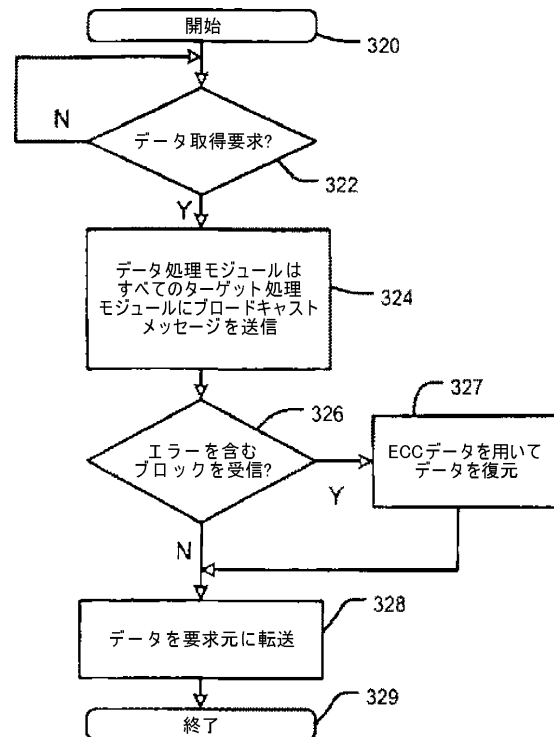
【図 6】



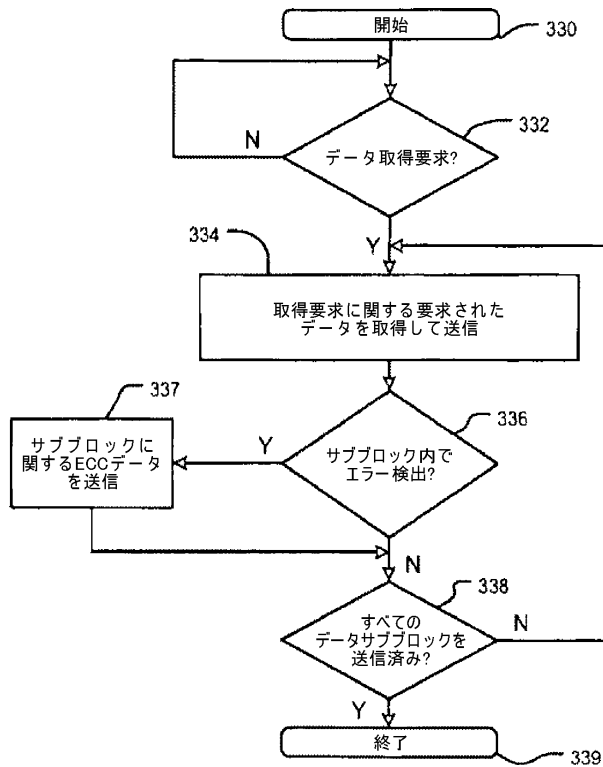
【図 7】



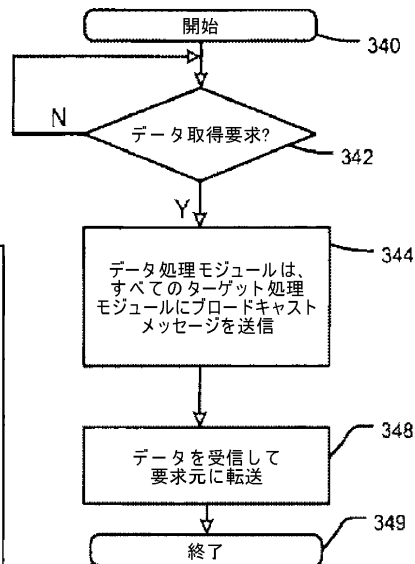
【図 8 A】



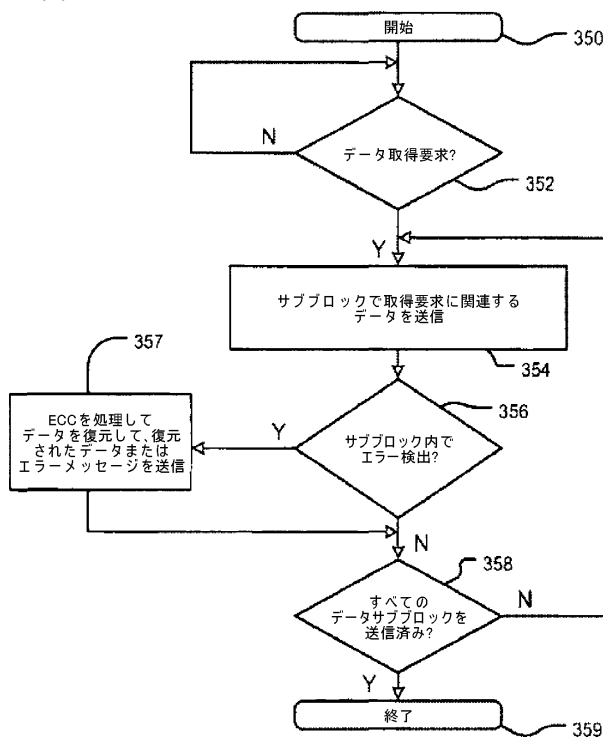
【図 8 B】



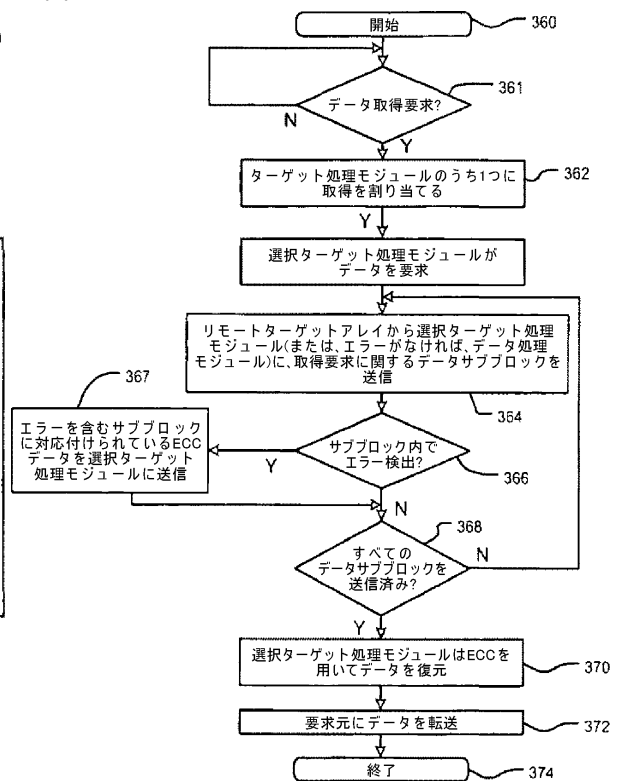
【図 9 A】



【図 9 B】



【図 10】



フロントページの続き

(31)優先権主張番号 11/724,549

(32)優先日 平成19年3月15日(2007.3.15)

(33)優先権主張国 米国(US)

(74)代理人 100112520

弁理士 林 茂則

(74)代理人 100156591

弁理士 高田 学

(72)発明者 スタルジャ、パントス

アメリカ合衆国、95054 カリフォルニア州、サンタ クララ、マーベル レーン 5488
マーベル セミコンダクター インコーポレイテッド内

審査官 坂東 博司

(56)参考文献 特開平08-263224(JP,A)

特開平07-200187(JP,A)

特開平06-067814(JP,A)

特開2004-192483(JP,A)

特開2005-275829(JP,A)

特開2004-334706(JP,A)

欧州特許出願公開第00730229(EP,A1)

特表平05-504431(JP,A)

特表2004-514968(JP,A)

米国特許出願公開第2005/0022052(US,A1)

米国特許第04989206(US,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06