

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5399345号
(P5399345)

(45) 発行日 平成26年1月29日(2014.1.29)

(24) 登録日 平成25年11月1日(2013.11.1)

(51) Int.Cl. F I
G06F 3/06 (2006.01) G O 6 F 3/06 3 O 4 Z
G06F 12/00 (2006.01) G O 6 F 12/00 5 1 4 A

請求項の数 13 外国語出願 (全 25 頁)

(21) 出願番号	特願2010-192834 (P2010-192834)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成22年8月30日(2010.8.30)	(74) 代理人	100093861 弁理士 大賀 真司
(65) 公開番号	特開2011-221981 (P2011-221981A)	(74) 代理人	100129218 弁理士 百本 宏之
(43) 公開日	平成23年11月4日(2011.11.4)	(72) 発明者	川口 智大 アメリカ合衆国 カリフォルニア州 95 014 クパチーノ プルネリッジ・アベ ニュー#9304 19500
審査請求日	平成24年7月24日(2012.7.24)	(72) 発明者	山本 彰 神奈川県横浜市戸塚区吉田町292番地 株式会社日立製作所 システム開発研究所 内
(31) 優先権主張番号	12/756,475		
(32) 優先日	平成22年4月8日(2010.4.8)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 外部ストレージシステムに結合されたストレージシステムのエラーコード管理方法及び装置

(57) 【特許請求の範囲】

【請求項1】

ホストコンピュータから入出力オペレーションを受信する第1ポートと、
 第1プロセッサ及び第1メモリを含む第1ストレージコントローラと、
 前記ホストコンピュータから受信したデータを格納する複数の第1ストレージ装置と
 を含む、第1ストレージシステムと、
 前記第1ストレージコントローラを介して前記ホストコンピュータから入出力オペレー
 ションを受信する第2ポートと、
 第2プロセッサ及び第2メモリを含む第2ストレージコントローラと、
 前記ホストコンピュータから受信したデータを格納する複数の第2ストレージ装置と
 を含む、第2ストレージシステムとを備えるシステムであって、
 前記第1ストレージコントローラは、
 前記複数の第2ストレージ装置に基づき前記ホストコンピュータに第1仮想ボリューム
 を提供し、

前記第1仮想ボリュームに対するライトデータを前記第1メモリに一時的に格納し、一
 時的に格納した前記ライトデータを前記第1メモリから前記複数の第2ストレージ装置に
 デスティングする際、前記ライトデータに対応するエラーチェックコードを生成し、生
 成したエラーチェックコードを前記第1メモリに格納し、

前記ホストコンピュータからのデータ転送要求に応じて、データ転送要求の対象となる
 リードデータが前記第1メモリに格納されているか否かを確認し、

10

20

前記リードデータが前記第1メモリに格納されている場合、前記第1メモリに格納されている前記リードデータを前記ホストコンピュータに転送し、

前記リードデータが前記第1メモリに格納されていない場合、前記リードデータを前記第2ストレージ装置から前記第1メモリにステージングする際、前記リードデータに対応するエラーチェックコードを計算し、かつ、前記第1メモリに格納されている前記リードデータに対応するエラーチェックコードと比較し、前記比較結果が一致した場合、前記第2ストレージシステムに格納されている前記リードデータを前記ホストコンピュータに転送し、前記比較結果が一致しなかった場合、前記第2ストレージシステムに前記リードデータを回復するように要求し、前記第2ストレージシステムに格納された前記リードデータを前記第2ストレージコントローラが回復した場合、回復した前記リードデータに対応するエラーチェックコードを計算し、かつ、前記第1メモリに格納されている前記リードデータに対応するエラーチェックコードと比較して、比較結果が一致することを確認した後、回復した前記リードデータを前記ホストコンピュータに転送する、システム。

10

【請求項2】

前記複数の第1ストレージ装置は、
前記複数の第2ストレージ装置よりも高い信頼性を有し、
前記エラーチェックコードは、
第1ストレージコントローラによって計算されて、前記第1メモリに格納される、請求項1に記載のシステム。

【請求項3】

前記第2ストレージシステムにデータを書き込むための前記ホストコンピュータからの書込みI/Oオペレーションにตอบสนองして、前記データは前記第1ポートを介して前記第1メモリに受信され、前記エラーチェックコードは前記受信データのために前記第1ストレージコントローラによって生成され、その後、前記データは前記第2ポートを介して前記複数の第2ストレージ装置に格納され、
前記エラーチェックコードはハッシュ機能を用いて生成される、請求項1に記載のシステム。

20

【請求項4】

前記第2ストレージシステムに格納されたデータはRAIDレベル5によって格納され、前記データは、対応するストライプセットのパリティビットを用いて回復される、請求項1に記載のシステム。

30

【請求項5】

前記比較結果が一致しない場合、前記第2ストレージシステムに格納されたデータは、前記第2ストレージシステムに格納されたデータ及びパリティを用いて、前記第1ストレージコントローラによって回復される、請求項1に記載のシステム。

【請求項6】

ホストコンピュータからの入出力オペレーションを受信する第3ポートと、
第3プロセッサ及び第3メモリを含む第3ストレージコントローラと、
前記ホストコンピュータから受信したデータを格納する複数の第3ストレージ装置とを含む、第3ストレージシステムをさらに備えるシステムであって、
前記複数の第3ストレージ装置に格納されたデータに対応するエラーチェックコードは前記第1ストレージシステムに格納される、請求項1に記載のシステム。

40

【請求項7】

前記複数の第3ストレージ装置は、
前記複数の第2ストレージ装置の重複データを格納し、
前記第2ストレージシステムに格納されたデータを読み出すための前記ホストコンピュータからの入出力オペレーションにตอบสนองして、前記格納されたデータのチェックコードが計算され、かつ、前記第1ストレージシステムに格納された対応するエラーチェックコードと比較され、
前記比較結果が一致しない場合、前記第3ストレージシステムに格納された前記データ

50

は、前記第3及び第1ポートを介して前記ホストコンピュータに転送される、請求項6に記載のシステム。

【請求項8】

前記第2ストレージシステムにデータを書き込むための前記ホストコンピュータからの入出力オペレーションにตอบสนองして、前記データは前記第1ポートを介して前記第1メモリに受信され、前記受信データが前記複数の第2及び第3ストレージ装置にデステージングされる際、前記第1ストレージコントローラによってエラーチェックコードが生成され、その後、前記エラーチェックコードは前記第1ストレージ装置に格納される一方で、前記データは前記複数の第2及び第3ストレージ装置に格納され、

前記エラーチェックコードはハッシュ機能を用いて生成される、請求項7に記載のシステム。

【請求項9】

第1ストレージシステムに結合された外部ストレージシステムを制御する方法であって、

前記第1ストレージシステムの第1ストレージコントローラが、
前記外部ストレージシステムに基づいてホストコンピュータに第1仮想ボリュームを提供することと、

前記第1仮想ボリュームに対するライトデータを前記ホストコンピュータから受信すると、受信した前記ライトデータを前記第1ストレージシステムの第1メモリに一時的に格納し、一時的に格納した前記ライトデータを前記第1メモリから前記外部ストレージシステムにデステージングする際、前記ライトデータに対応するエラーチェックコードを生成し、生成したエラーチェックコードを前記第1メモリに格納することと、

前記ホストコンピュータからのデータ転送要求に応じて、データ転送要求の対象となるリードデータが前記第1メモリに格納されているか否かを確認することと、

前記リードデータが前記第1メモリに格納されている場合、前記第1メモリに格納されている前記リードデータを前記ホストコンピュータに転送することと、

前記第1メモリに前記リードデータが格納されていない場合、前記リードデータを前記外部ストレージシステムから前記第1メモリにステージングする際、前記リードデータに対応するエラーチェックコードを計算し、かつ、前記第1メモリに格納されている前記リードデータに対応するエラーチェックコードと比較し、前記比較結果が一致した場合、前記外部ストレージシステムに格納された前記リードデータを前記ホストコンピュータに転送し、前記比較結果が一致しなかった場合、前記外部ストレージシステムに前記リードデータを回復するように要求し、前記外部ストレージシステムに格納された前記リードデータを前記外部ストレージシステムが回復した場合、回復した前記リードデータに対応するエラーチェックコードを計算し、かつ、前記第1メモリに格納されている前記リードデータに対応するエラーチェックコードと比較して、比較結果が一致することを確認した後、回復した前記リードデータを前記ホストコンピュータに転送すること

を含む、方法。

【請求項10】

前記比較結果が一致しない場合、前記外部ストレージシステムへ回復データを要求し、かつ、前記第1ストレージコントローラによって前記回復データに基づいて前記リードデータを計算することをさらに備える、請求項9に記載の方法。

【請求項11】

前記第1ストレージシステムの対応するキャッシュスロットに前記リードデータを書き込むことをさらに備える、請求項10に記載の方法。

【請求項12】

前記リードデータのハッシュ値を計算して、前記リードデータの前記ハッシュ値と、前記第1ストレージシステムに格納された前記対応するエラーチェックコードとを比較することと、

前記リードデータの前記ハッシュ値と、前記対応するエラーチェックコードとが一致し

10

20

30

40

50

た場合、前記リードデータを前記ホストコンピュータに転送することとをさらに備える、請求項9に記載の方法。

【請求項13】

前記比較結果が一致しない場合、第2ストレージシステムからデータを転送することをさらに備え、

前記第2ストレージシステムは、前記外部ストレージシステムに書き込まれた重複データを格納し、かつ、前記第1ストレージシステムに結合される、請求項9に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

0001 本発明は、高階層ストレージシステムに結合された低階層ストレージシステムに格納されたデータの高い信頼性を管理する方法及び装置に関する。

【背景技術】

【0002】

0002 多階層ストレージシステムでは、システムは、異なる容量及び性能信頼性を有する複数のストレージシステムから構成され得る。ユーザは、データの予算、負荷及び重要性によって、そのデータを格納する階層を決定することになる。データ信頼性を高めるためには、日本特許第2000-347815号に開示されているように、データ訂正コードをデータに加えてもよい。しかし、データがデータ訂正コードを含む場合、低階層ストレージは、追加されたエラーコードをサポートできないかもしれない。このことは、システム内のデータ移行を考慮すると、システムの総合的信頼性をもたらし得る。

【発明の概要】

【0003】

0003 本発明の例示的实施形態は、複数レベルの信頼性のストレージ装置を用いる複数のストレージシステムを備えるシステムを提供する。比較的高い信頼性のストレージシステムの中に比較的低い信頼性のストレージ装置のためのエラーコードを維持することによって、システム全体としての信頼性は高められる。エラーコードは、ハッシュ機能を用いて計算され、この値は、比較的低い信頼性のストレージディスクから読み出されたデータのハッシュ値と比較するために使用される。

【0004】

0004 一実施形態では、比較的高い信頼性のストレージシステムは、比較的低い信頼性のストレージシステムから要求された該当データを取得することによって、訂正データを計算する。他の実施形態では、比較的高い信頼性のストレージシステムが、比較的低い信頼性のストレージシステムに対して、訂正データを生成するように要求する。

【図面の簡単な説明】

【0005】

【図1】0005 図1は、本発明の構成の概要の一例を示す。

【図2】0006 図2は、図1のストレージサブシステム100のメモリの一例を示す。

【図3】0007 図3は、図2のメモリのRAIDグループ管理テーブルの一例を示す。

【図4】0008 図4は、図3のメモリのボリューム管理テーブルの一例を示す。

【図5】0009 図5は、図3のメモリの外部ボリュームエラーコードテーブルの一例を示す。

【図6】0010 図6は、図3のメモリのキャッシュ管理テーブルの一例を示す。

【図7】0011 図7は、図1のメモリのキャッシュ管理テーブルの割当の一例を示す。

【図8】0012 図8は、図1のストレージサブシステム100の書込みI/O制御シーケンスの一例を示す。

【図9】0013 図9は、図1のストレージサブシステム100の読出しI/O制御

10

20

30

40

50

シーケンスの一例を示す。

【図10】0014 図10は、図1のストレージサブシステム100のステージング制御シーケンスの一例を示す。

【図11】0015 図11は、図1のストレージサブシステム100のデステージング制御シーケンスの一例を示す。

【図12】0016 図12は、図1のストレージサブシステム100のフラッシュ制御シーケンスの一例を示す。

【図13】0017 図13は、図1のストレージサブシステム100の外部ボリューム実装制御シーケンスの一例を示す。

【図14】0018 図14は、図1のストレージサブシステム400のメモリの一例を示す。 10

【図15】0019 図15は、図1のストレージサブシステム400のステージング制御シーケンスの一例を示す。

【図16】0020 図16は、図1のストレージサブシステム400のデステージング制御シーケンスの一例を示す。

【図17】0021 図17は、図1のストレージサブシステム400の回復データ転送制御シーケンスの一例を示す。

【図18】0022 図18は、図1のシステムの書き込みI/Oプロセスフローの一例を示す流れ図である。

【図19】0023 図19は、図1のシステムの読出しI/Oプロセスフローの一例を示す流れ図である。 20

【図20】0024 図20は、図1のストレージサブシステム100のステージング制御シーケンスの一例を示す。

【図21】0025 図21は、図1のストレージサブシステム400のメモリの一例を示す。

【図22】0026 図22は、図1のストレージサブシステム100のデータ回復制御シーケンスの一例を示す。

【図23】0027 図23は、図1のシステムの読出しI/Oプロセスフローの一例を示す流れ図である。

【図24】0028 図24は、本発明の構成の概要の一例を示す。 30

【図25】0029 図25は、図1のストレージサブシステム100のメモリの一例を示す。

【図26】0030 図26は、図25のメモリのRAIDグループ管理テーブルの一例を示す。

【図27】0031 図27は、図1のストレージサブシステム100のステージング制御シーケンスの一例を示す。

【図28】0032 図28は、図24のシステムの読出しI/Oプロセスフローの一例を説明する流れ図の一例を示す。

【図29】0033 図29は、図2のメモリのストレージリストの一例を示す。

【発明を実施するための形態】 40

【0006】

0034 第1の実施形態

【0007】

0035 図1は、本発明の方法及び装置が適用され得るシステムのハードウェア構成を示す。ストレージサブシステム100は、SAN(ストレージエリアネットワーク)200を介してホストコンピュータ300に接続される。ストレージサブシステム400は、ファイバ・チャネル(FC)を介してストレージサブシステム100に接続される。ストレージサブシステム100は、ホストコンピュータ200からI/Oコマンドを受信して、両方のストレージサブシステム100、400のストレージ装置121、421を用いて、ホストコンピュータ200にストレージボリュームを提供する。ストレージサブ 50

システム 100 は、ストレージサブシステム 400 より高い信頼性を有する。

【0008】

0036 ストレージサブシステム 100 は、CPU 111 と、メモリ 112 と、ストレージインタフェース 113、114 と、ディスクインタフェース 115 とを含む、ストレージコントローラ 110 を有する。CPU 111 は、ストレージサブシステム 100 を制御し、メモリ 112 からプログラム及びテーブルを読み出す。メモリ 112 はプログラム及びテーブルを格納する。ストレージインタフェース 113 は、ストレージネットワーク 200 を介してホストコンピュータ 300 に接続する。ストレージインタフェース 114 は、ストレージサブシステム 400 のストレージインタフェース 413 に接続する。ディスクインタフェース 115 は複数のストレージ装置 121 に接続し、これら複数のストレージ装置 121 はディスクユニット 120 に格納される。ストレージ装置 121 は、データを格納するためのソリッドステートデバイス（例えば、フラッシュメモリ及び/又はハードディスクドライブ（HDD）など）から構成される。ストレージネットワーク 200 は、ストレージサブシステム 100 及びホストコンピュータ 300 に接続する。ホストコンピュータ 300 は、ストレージネットワーク 200 を介してストレージサブシステム 100 に I/O 要求を送信し、ストレージネットワーク 200 を介してストレージサブシステム 100 との間でデータを送受信する。ストレージサブシステム 400 は、CPU 411 と、メモリ 412 と、ストレージインタフェース 413 と、ディスクインタフェース 115 とを含む、ストレージコントローラ 410 を有する。ストレージサブシステム 200 は、ストレージサブシステム 100 の外部ストレージであって、ストレージサブシステム 100 及びストレージネットワーク 200 を介してホストコンピュータ 300 との間でデータを送受信することになる。CPU 411 はストレージサブシステム 400 を制御し、メモリ 412 からプログラム及びテーブルを読み出す。メモリ 412 はプログラム及びテーブルを格納する。ディスクインタフェース 415 は複数のストレージ装置 421 に接続し、これらのストレージ装置 421 はディスクユニット 420 に格納される。ストレージ装置 421 は、データを格納するためのソリッドステートデバイス（例えば、フラッシュメモリ及び/又はハードディスクドライブ（HDD）など）から構成される。二つのストレージサブシステムを比較すると、ストレージサブシステム 100 は、ストレージサブシステム 400 より比較的高い信頼性を有する。この例では、内部ボリュームに使用されるディスクユニット 120 は、外部ボリュームに使用されるディスクユニット 420 に比べて、より高いグレードのストレージ装置から構成され、例えば、ディスクユニット 120 では SLC（シングルレベルセル）フラッシュメモリが使用され、ディスクユニット 420 では、MLC（マルチレベルセル）フラッシュメモリ又は比較的安価な SATA（シリアル ATA）HDD が使用される。ストレージコントローラ 110 の CPU 内のプロセッサの数又はグレード、又はメモリの容量は、ストレージコントローラ 410 よりも大きいものであり得る。ストレージサブシステム 400 で使用されているものに比べて、ストレージサブシステム 100 において比較的高いグレードのプロセッサを使用することによって、ストレージコントローラ 100 によるデータ処理のより高い信頼性は、ストレージサブシステム 400 に格納されたデータの信頼性を高めることになるだろう。

【0009】

0037 図 2 は、図 1 のストレージサブシステム 100 のメモリ 112 の一例を示す。メモリ 112 は、RAID グループ管理テーブル 112 - 11 - 1 と、ボリューム管理テーブル 112 - 11 と、外部ボリュームエラーチェックコードテーブル 112 - 11 - 3 と、高信頼性ストレージリスト 112 - 11 - 4 とを含む、ストレージ管理テーブル 112 - 11 を含む。管理テーブル 112 - 11。RAID グループ管理テーブル 112 - 11 - 1 は、ストレージ装置 121、外部ボリューム及びこれらのグループの物理構造管理を提供する。ボリューム管理テーブル 112 - 11 - 2 は論理ボリューム構成を提供する。外部ボリュームエラーチェックコードテーブル 112 - 11 - 3 は、外部ボリュームのいくつかの領域についてエラーチェックコードを格納する。ある領域のエラーチェックコードの値は、当該領域に格納されたデータからハッシュ計算によって計算される。高

10

20

30

40

50

信頼性ストレージリスト 112 - 11 - 4 は、高ストレージ製品名又は製品 ID を格納し、これは、ストレージが比較的低い信頼性のものであるか否かを判断するために使用される。ストレージシステムのために使用されるストレージ製品がリストに記憶されていない場合、そのストレージは比較的低い信頼性のものとして扱われ、比較的高い信頼性のストレージシステムにエラーコードが格納される。キャッシュデータ領域 112 - 30 の管理及び LRU / MRU 管理のために、キャッシュ管理テーブル 112 - 14 が設けられる。ボリューム I / O 制御 112 - 21 は、書込み I / O 要件によって実行され、ライトデータを受信し、キャッシュデータ領域 112 に格納する書込み I / O 制御 112 - 21 - 1 (図 8) と、読出し I / O 要件によって実行され、キャッシュデータ領域 112 からリードデータを送信する読出し I / O 制御 112 - 21 - 2 (図 9) とを含む。ディスク制御 112 - 22 は、ディスク 121 からキャッシュデータ領域 112 へデータを転送するステージング制御 112 - 22 - 1 (図 10) と、キャッシュデータ領域 112 からディスク 121 へデータを転送するデステージング制御 112 - 22 - 2 (図 11) とを含む。メモリ 112 は、キャッシュデータ領域からディスク 121 へと定期的にダーティデータをフラッシュするフラッシュ制御 112 - 23 (図 12) と、キャッシュデータ領域にキャッシュされたデータを見つけ、かつキャッシュデータ領域に新しいキャッシュ領域を割り当てるキャッシュ制御 112 - 24 とをさらに含む。メモリ 112 は、リード及びライトキャッシュデータを格納するキャッシュデータ領域 112 - 30 を含む。この領域は、複数のキャッシュスロット用に分割される。各キャッシュスロットはデータストライプのために割り当てられる。メモリ 112 は、プログラム実行スケジュールを制御し、マルチタスク環境をサポートするカーネル 112 - 40 を含む。プログラムが ACK (確認) を待っている場合、CPU 111 は別のタスクを実行するために変更する (例えば、ディスク 121 からキャッシュデータ領域 112 - 30 へのデータ転送待ち)。メモリ 112 は、外部ボリューム実装の実装を制御する外部ボリューム実装制御 112 - 26 (図 13) を含む。

【 0010 】

0038 図 3 は、図 2 のメモリ 112 の RAID グループ管理テーブル 112 - 11 - 1 の一例を示す。RAID グループ管理テーブル 112 - 11 - 1 は、RAID グループの ID としての RAID グループ番号 112 - 11 - 1 - 1 欄と、RAID グループの構造を表す RAID レベル 112 - 11 - 1 - 2 欄とを含む。例えば、「5」は「RAID レベルが 5 である」ことを意味する。「NULL」は、RAID グループが存在しないことを意味する。「Ext / 1」は、RAID グループが内部ボリュームの外にある外部ボリュームとして存在し、かつ RAID レベルが 1 であることを意味する。RAID グループ管理テーブル 112 - 11 - 1 は、内部ボリュームの場合には、RAID グループに属する HDD の ID リストを表す HDD 番号の欄、外部ボリュームの場合には WWN (ワールドワイドネーム) の欄 112 - 11 - 1 - 3 を含む。RAID グループ管理テーブル 112 - 11 - 1 はさらに、重複領域を除く RAID グループの全容量を表す RAID グループ容量 112 - 11 - 1 - 4 を含む。RAID グループ管理テーブル 112 - 11 - 1 さらに、ストレージ装置の信頼性を表す信頼性 112 - 11 - 1 - 5 を含む。ストレージ装置の信頼性は、管理サーバによって手動で設定されてもよく、又は製品が図 29 のような高信頼性ストレージリスト 112 - 11 - 4 に含まれているか否かチェックすることによって判断してもよい。ストレージ装置の製品タイプ ID が、リストの製品タイプ ID 112 - 11 - 4 - 1 のうちの一つに一致した場合、当該製品は比較的高い信頼性を有すると判断され、いずれの製品タイプ ID にも一致しない場合、当該製品は比較的低い信頼性を有すると判断される。製品タイプ ID は、管理サーバによって追加又は削除され得る。リストには比較的低い信頼性の製品を記載してもよく、逆の場合も同様に判断され得る。

【 0011 】

0039 図 4 は、図 2 のメモリ 112 のボリュームグループ管理テーブル 112 - 11 - 2 の一例を示す。ボリュームグループ管理テーブル 112 - 11 - 2 は、ボリュー

10

20

30

40

50

ムのIDとしてのボリューム番号112-11-2-1欄と、ボリュームの容量を表す容量112-11-2-1欄とを含む。「N/A」は、ボリュームが実際に存在せず、よって、そのボリュームについての関連情報がないことを意味する。ボリュームグループ管理テーブル112-11-2はさらに、ボリュームによって使用されるRAIDグループ番号112-11-1-1を表すRAIDグループ番号112-11-2-3と、ボリュームについて使用されるアドレス範囲を示すアドレス範囲112-11-2-5とを含む。ボリュームグループ管理テーブル112-11-2はさらに、それによってボリュームにアクセスすることができるポート番号を表すポート番号112-11-2-6と、ポートを介して認識されるボリュームのIDを表すLUN112-11-2-7とを含む。

【0012】

0040 図5は、図2のメモリ112の外部ボリュームエラーチェックコードテーブル112-11-3の一例を示す。外部ボリュームエラーチェックコードテーブル112-11-3は、仮想ボリュームのIDとしての仮想ボリューム番号112-11-3-1欄と、スロットのIDを表すスロット番号112-11-3-2欄とを含む。外部ボリュームエラーチェックコードテーブル112-11-3はさらに、外部ボリュームのエラーチェックコードを表すエラーチェックコード112-11-3-3を含み、当該エラーチェックコードは、スロット内のデータの計算されたハッシュ値である。

【0013】

0041 図6は、図2のメモリ112のキャッシュ管理テーブル112-14の一例を示す。キャッシュ管理テーブル112-14は、キャッシュデータ領域112-30におけるキャッシュスロットのIDとしてのインデックス112-14-1欄と、対応するデータを格納するキャッシュスロットのディスク121のIDを表すディスク番号112-14-2欄とを含む。キャッシュ管理テーブル112-14はさらに、対応するデータを格納するディスクの論理ブロックアドレスを表すLBA112-14-3と、キュー管理のための次キャッシュスロット番号を表す次112-14-4とを含む。「NULL」とは、陰に連続するキューがないこと、及びそのキューは当該スロットで終了することを意味する。キャッシュ管理テーブル112-14はさらに、キャッシュスロットキューの種類(タイプ)を表すキューの種類112-14-5と、処理されるべき次スロットであるキュースロットキューのトップスロットIDを表すキューインデックスポインタ112-14-6とを含む。「フリー」スロットキューは、未使用キャッシュスロットを有するキューであり、これは新しいライトデータを割り当てるために使用されることになる。「クリーン」スロットキューは、キューがディスクスロットの中の同じデータを格納するキャッシュスロットを有し、かつ当該データがディスクにフラッシュアウトされている。「ダーティ」スロットキューは、ディスクにデータをフラッシュアウトしていないキューである。キャッシュスロットは、対応するディスクスロットから異なるデータを格納するため、ストレージコントローラ110は、将来的にフラッシュ制御112-23を使用してキャッシュスロット中のデータをディスクスロットにフラッシュする必要がある。「ダーティ」スロットがディスクにフラッシュされた後、スロットのスロット状態は「クリーン」に変更することになる。

【0014】

0042 図7は、図1のストレージシステム100の論理構造の一例を示す。点線は、ポインタがオブジェクトを参照することを表す。実線は、オブジェクトが計算によって参照されることを表す。図2のキャッシュデータ112-30は、複数のキャッシュスロット112-30-1に分割される。キャッシュスロットのサイズは、容量プールストライプ121-3及び仮想ボリュームスロット141-3のサイズと同じである。キャッシュ管理テーブル112-14とキャッシュスロット112-30-1は互いに対応し、一対一の関係である。キャッシュ管理テーブル112-14は仮想ボリュームスロット141-3及び容量プールストライプ121-3を参照する。

【0015】

0043 図8は、図2のメモリ112の書込みI/O制御112-21-1のプロ

10

20

30

40

50

セスフローの一例を示す。プログラムは112-21-1-1で開始する。ステップ112-21-1-2では、プログラムは、キャッシュ制御112-24を呼び出して、キャッシュスロット112-30-1を検索する。ステップ112-21-1-3では、プログラムは、ホストコンピュータ300から書込みI/Oデータを受信し、当該データを上述のキャッシュスロット112-30-1に格納する。プログラムは112-21-1-4で終了する。

【0016】

0044 図9は、図2のメモリ112の読出しI/O制御112-21-2のプロセスフローの一例を示す。プログラムは、112-21-2-1で開始する。ステップ112-21-2-2では、プログラムは、キャッシュ制御112-24を呼び出して、キャッシュスロット112-30-1を検索する。ステップ112-21-2-3では、プログラムは、上述のキャッシュスロット112-30-1の状態をチェックし、データが既にそこに格納されているか否かを判断する。データがキャッシュスロット112-30-1に格納されていない場合、プログラムは、ステップ112-21-2-4でステージング制御112-22-1を呼び出す。ステップ112-21-2-5で、プログラムは、キャッシュスロット112-30-1の中のデータをホストコンピュータ300に転送する。プログラムは、112-21-2-6で終了する。

【0017】

0045 図10は、図2のメモリ112のステージング制御112-22-1のプロセスフローの一例を示す。プログラムは、112-22-1-1で開始する。ステップ112-22-1-2では、プログラムは、ボリューム管理テーブル112-11-2及びRAIDグループ管理テーブル112-11-1を参照して、物理ディスク及びデータのアドレスを判断する。ステップ112-22-1-3で、プログラムは、ディスク121、141の-slotからデータを読み出すことを要求し、当該データをバッファに格納する。ステップ112-22-1-4では、プログラムは、RAIDグループ管理テーブル112-11-1を用いて比較的低い信頼性のストレージディスクによって割り当てられた外部ボリュームにデータが格納されているか否かをチェックする。当該データが比較的低い信頼性のストレージディスクに格納されている場合、プログラムは、ステップ112-22-1-5において、バッファ内のデータからハッシュ値を計算し、計算したハッシュ値と外部ボリュームエラーコードテーブル112-11-3に格納されたエラーコードとを比較する。データが比較的低い信頼性のストレージディスクに格納されていない場合、プログラムはステップ112-22-1-9に進む。ステップ112-22-1-6で、プログラムは、比較された値が一致するか否かをチェックすることで、比較的低い信頼性のストレージディスクに格納されたデータエラーを検出することができる。比較された値が一致しない場合、プログラムは、ステップ112-22-1-7で外部ボリュームに対して回復データを転送するように要求する。よって、外部ボリュームがRAID5である場合、正しいデータを計算するためにストライプ列の-slotの重複データを要求することになる。そして、ステップ112-22-1-8では、プログラムは、送信された回復データから正しいデータを生成し、回復された-slotに対してパーティ属性を設定する。正しいデータはバッファに格納されることになる。外部ボリュームがRAID5である場合、正しいデータを生成するためにパリティ計算を実行する。比較的低い信頼性のストレージディスクに格納されたデータがデータエラーを含まず、かつ比較された値が一致する場合、プログラムはステップ112-22-1-9に進む。ステップ112-22-1-9で、プログラムは、バッファからキャッシュスロット112-30へ-slotデータを転送することによって、フラッシュ制御112-23及びデステージング制御112-22-2によって、比較的低い信頼性のストレージシステムの中のディスク及びキャッシュに、訂正されたデータを最終的に置き換える。プログラムは、112-22-1-10で終了する。

【0018】

0046 図11は、図2のメモリ112のデステージング制御112-22-2の

プロセスフローの一例を示す。プログラムは、112-22-2-1で開始する。ステップ112-22-2-2で、プログラムはボリューム管理テーブル112-11-2及びRAIDグループ管理テーブル112-11-1を参照して、物理ディスク及びデータのアドレスを判断する。ステップ112-22-2-3で、プログラムはステージング制御112-22-1を呼び出し、最新スロット領域をステージする。ステップ112-22-2-4では、プログラムは、キャッシュスロット112-30の書き込みされていない領域を送信されたデータで満たす。ステップ112-22-2-5で、プログラムは、RAIDグループ管理テーブル112-11-1を用いて比較的低い信頼性のストレージディスクによって割り当てられた外部ボリュームにデータが格納されているか否かをチェックする。データが比較的低い信頼性のストレージディスクに格納されている場合、プログラムは、ステップ112-22-2-6において、キャッシュスロット内のデータからハッシュ値を計算し、計算したチェックコードを外部ボリュームエラーコードテーブル112-11-3に格納する。データが比較的低い信頼性のストレージディスクに格納されていない場合、プログラムはステップ112-22-2-7に進む。ステップ112-22-2-7で、プログラムはキャッシュデータ領域112-30内のスロットからデータを読み出し、かつ内部又は外部ボリュームに格納する。プログラムは、112-22-2-8で終了する。

10

【0019】

0047 図12は、図2のメモリ112のフラッシュ制御112-23のプロセスフローの一例を示す。プログラムは、112-23-1で開始する。ステップ112-23-2で、プログラムは、キャッシュ管理テーブル112-14の「ダーティキュー」を読み取る。ダーティキャッシュ領域が見つかった場合、プログラムは、ステップ112-23-3で、見つかったダーティキャッシュスロット112-30-1のためにデステージング制御112-22-2を呼び出す。プログラムは、112-23-4で終了する。

20

【0020】

0048 図13は、図2のメモリ112の外部ボリューム実装制御112-25-1のプロセスフローの一例を示す。プログラムは、112-25-1-1で開始する。ステップ112-25-1-2で、プログラムは、使用されるストレージ装置のRAIDレベル、構造、製品名及び外部ボリュームの信頼性情報を含む構成情報を要求する。信頼性情報は、RAIDグループ管理テーブル112-11-1の信頼性112-11-1-5欄に格納される。外部ストレージの製品名が高信頼性ストレージリスト112-11-4にリストされている場合、又は外部ストレージが比較的高い信頼性を有すると報告した場合、RAIDグループ信頼性112-11-1-5に「高い」と格納する。上述の通りではない場合、RAIDグループ信頼性112-11-1-5に「低い」と格納する。プログラムは、112-25-1-3で終了する。

30

【0021】

0049 図14は、図1のストレージサブシステム400のメモリ412の一例を示す。メモリ412は、RAIDグループ管理テーブル112-11-1及びボリューム管理テーブル112-11（これらは112-11のテーブルと同一である）を含むストレージ管理テーブル412-11を含む。しかし、ストレージ管理テーブル412-11は、メモリ112の場合のような外部ボリュームエラーチェックコードテーブル112-11-3及び高信頼性ストレージリスト112-11-4を含まない。メモリ112の場合のようにキャッシュ管理テーブル112-14は、キャッシュデータ領域112の管理及びLRU/MRU管理のために設けられる。ボリュームI/O制御112-21は、メモリ112の場合のように、書き込みI/O要件によって実行し、ライトデータを受信し、かつキャッシュデータ領域112に格納する書き込みI/O制御112-21-1（図8）と、読み出しI/O要件によって実行し、かつキャッシュデータ領域412-30からリードデータを送信する読み出しI/O制御112-21-2（図9）とを含む。ディスク制御412-22は、ディスク421からキャッシュデータ領域412-30へデータを転送するステージング制御412-22-1（図15）と、キャッシュデータ領域412-3

40

50

0 からディスク 4 2 1 へデータを転送するデステージング制御 4 1 2 - 2 2 - 2 (図 1 6) と、正しいデータを生成するために指定領域のパリティビットを含む重複データを転送する回復データ転送制御 4 1 2 - 2 2 - 3 (図 1 7) とを含む。メモリ 1 1 2 は、メモリ 1 1 2 の場合のように、キャッシュデータ領域からディスク 4 2 1 へと定期的にデータデータをフラッシュするフラッシュ制御 1 1 2 - 2 3 と、キャッシュデータ領域にキャッシュされたデータを見つけ、かつキャッシュデータ領域に新しいキャッシュ領域を割り当てるキャッシュ制御 1 1 2 - 2 4 とをさらに含む。メモリ 4 1 2 は、リード及びライトキャッシュデータを格納するキャッシュデータ領域 4 1 2 - 3 0 を含む。この領域は、複数のキャッシュスロット用に分割される。各キャッシュスロットはデータストライプのために割り当てられる。メモリ 4 1 2 は、メモリ 1 1 2 の場合のように、プログラム実行スケジュールを制御し、マルチタスク環境をサポートするカーネル 1 1 2 - 4 0 を含む。

10

【 0 0 2 2 】

0 0 5 0 図 1 5 は、図 1 4 のメモリ 4 1 2 のデステージング制御 4 1 2 - 2 2 - 1 のプロセスフローの一例を示す。プログラムは、4 1 2 - 2 2 - 1 - 1 で開始する。ステップ 4 1 2 - 2 2 - 1 - 2 では、プログラムはボリューム管理テーブル 1 1 2 - 1 1 - 2 及び RAID グループ管理テーブル 1 1 2 - 1 1 - 1 を参照して、物理ディスク及びデータのアドレスを判断する。ステップ 4 1 2 - 2 2 - 1 - 3 で、プログラムは、ディスク 4 2 1 からのデータの読み取りを要求し、当該データをキャッシュデータ領域 4 1 2 - 3 0 に格納する。ステップ 4 1 2 - 2 2 - 1 - 4 では、プログラムはデータ転送の終了を待つ。メモリ 4 1 2 のカーネル 1 1 2 - 4 0 は、文脈切替えを行うために命令を発行する。プログラムは、4 1 2 - 2 2 - 1 - 5 で終了する。

20

【 0 0 2 3 】

0 0 5 1 図 1 6 は、図 1 4 のメモリ 4 1 2 のデステージング制御 4 1 2 - 2 2 - 2 のプロセスフローの一例を示す。プログラムは、4 1 2 - 2 2 - 2 - 1 で開始する。ステップ 4 1 2 - 2 2 - 2 - 2 で、プログラムは、ボリューム管理テーブル 1 1 2 - 1 1 - 2 及び RAID グループ管理テーブル 1 1 2 - 1 1 - 1 を参照して、物理ディスク及びデータのアドレスを判断する。ステップ 4 1 2 - 2 2 - 2 - 3 では、プログラムは、キャッシュデータ領域 4 1 2 - 3 0 からのデータの読み取りを要求し、当該データをディスク 4 2 1 に格納する。ステップ 4 1 2 - 2 2 - 2 - 4 で、プログラムはデータ転送の終了を待つ。メモリ 4 1 2 のカーネル 1 1 2 - 4 0 は、文脈切替えを行うために命令を発行する。プログラムは、4 1 2 - 2 2 - 2 - 5 で終了する。

30

【 0 0 2 4 】

0 0 5 2 図 1 7 は、図 1 4 のメモリ 4 1 2 の回復データ転送制御 4 1 2 - 2 1 - 3 のプロセスフローの一例を示す。プログラムは、4 1 2 - 2 1 - 3 - 1 で開始する。ステップ 4 1 2 - 2 1 - 3 - 2 で、プログラムは、ボリューム管理テーブル 1 1 2 - 1 1 - 2 及び RAID グループ管理テーブル 1 1 2 - 1 1 - 1 を参照して、物理ディスク及びデータのアドレスを判断する。ステップ 4 1 2 - 2 1 - 3 - 3 では、プログラムは、キャッシュ制御 1 1 2 - 2 4 を呼び出して、対応するキャッシュスロット 4 1 2 - 3 0 - 1 を検索する。ステップ 4 1 2 - 2 1 - 3 - 4 では、プログラムは、前記キャッシュスロット 4 1 2 - 3 0 - 1 の状態をチェックする。データがまだキャッシュに格納されていない場合、ステップ 4 1 2 - 2 1 - 3 - 5 で、プログラムはデステージング制御 4 1 2 - 2 1 - 1 を呼び出す。データが既にキャッシュに格納されている場合、プログラムはステップ 4 1 2 - 2 1 - 3 - 6 に移動する。ステップ 4 1 2 - 2 1 - 3 - 6 では、プログラムは、キャッシュスロット 1 1 2 - 3 0 - 1 データをインシエータに転送する。よって、メモリ 1 1 2 のデステージング制御 1 1 2 - 2 2 - 1 がプログラムを呼び出した場合、データはストレージコントローラ 1 1 0 に転送されることになり、従って比較的高いストレージシステム 1 0 0 に正しいデータを生成することができる。プログラムは、4 1 2 - 2 1 - 3 - 7 で終了する。

40

【 0 0 2 5 】

0 0 5 3 図 1 8 は、図 1 のシステムで行われる書込みオペレーションの一例を示す

50

。ホストコンピュータ300は、高信頼性ストレージサブシステム100に書き込まれるべきデータと共に書込みI/O要求を送信する(W1001)。高信頼性ストレージサブシステム100のCPU111は、書込みI/O要求を受信し、このデータを高信頼性ストレージサブシステム100のキャッシュスロット112-30-1に格納する(W1002)。キャッシュ領域112-30は書込みI/Oデータを受信する(W1003)。CPU111は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、かつデステージング制御112-22-2を実行して、エラーチェックコードを生成する(W1004)。キャッシュ領域112-30は、ダーティスロットデータを外部ボリュームに転送する(W1005)。低信頼性ストレージサブシステム400のCPU411は書込みI/O要求を受信し、かつデータを低信頼性ストレージサブシステム400のキャッシュスロット412-30-1に格納する(W1006)。キャッシュ領域412-30は書込みI/Oデータを受信する(W1007)。CPU111は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、かつデステージング制御112-22-2を実行する(W1008)。キャッシュ領域412-30は、ダーティスロットデータをディスク421に転送する(W1009)。ディスク421は、当該データを受信し格納する(W1010)。

【0026】

0054 図19は、図1のシステムで行われる読出しオペレーションの一例を示す。ホスト300は、高信頼性ストレージサブシステム100に読出しI/O要求を送信する(R1001)。高信頼性ストレージサブシステム100のCPU111は、読出しI/O要求を受信し、ステージング制御112-22-1を呼び出して、読出しI/Oデータをキャッシュスロット112-30-1に格納する。ステージング制御112-22-1は、データエラーが存在するか否かチェックし、データエラーが存在する場合には、データを回復した後、データを転送する(R1002)。キャッシュ領域112-30は、外部ボリュームデータの読出しを要求し、データをホスト300に転送する(R1003)。低信頼性ストレージサブシステム100のCPU411は、読出しI/O要求を受信し、ステージング制御412-22-1を呼び出して、読出しI/Oデータをキャッシュスロット412-30-1に格納する(R1004)。キャッシュ領域412-30は、ディスク421からディスクデータを読み出すことを要求する(R1005)。ディスク421は、当該要求に従ってデータを送信する(R1006)。CPU111は、エラーチェックコードを計算し、かつ外部ボリュームエラーチェックコード112-11-3のエラーチェックコードと比較することによって、データのエラーを検出する(R1007)。キャッシュ領域112-30は、回復データの読出しを要求し、データを転送する(R1008)。低信頼性ストレージサブシステム100のCPU411は、回復データ読出し要求を受信し、ステージング制御112-22-1を呼び出して、読出しI/Oデータをキャッシュスロット112-30-1に格納する(R1009)。データが破損した場合、ステップW1004~W1010に示されるように、正しいデータを低信頼性ストレージサブシステム400のキャッシュ及びディスクに書き込まなければならない。CPU111は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、デステージング制御112-22-2を実行し、これによってエラーチェックコードを生成する(W1004)。キャッシュ領域412-30は、ダーティスロットデータを外部ボリュームに転送する(W1005)。低信頼性ストレージサブシステム400のCPU411は、書込みI/O要求を受信し、データを低信頼性ストレージサブシステム400のキャッシュスロット412-30-1に格納する(W1006)。キャッシュ領域412-30は書込みI/Oデータを受信する(W1007)。CPU411は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、デステージング制御112-22-2を実行する(W1008)。キャッシュ領域412-30は、ダーティスロットデータをディスク421に転送する(W1009)。ディスク421は、当該データを受信し格納する(W1010)。

【0027】

10

20

30

40

50

0055 第2の実施形態

【0028】

0056 第1の実施形態では、高信頼性ストレージサブシステム100のストレージコントローラ110が、低信頼性ストレージサブシステム400から正しいデータを生成するために必要なデータを検索したが、第2の実施形態は、低信頼性ストレージサブシステム400のストレージコントローラ410によって正しいデータを生成するための方法を示す。第1の実施形態との差異についてのみ、図20～図22を用いて説明する。

【0029】

0057 図20は、図2のメモリ112のステージング制御112-22-1のプロセスフローの一例を示す。プログラムは、112-22-1-1で開始する。ステップ112-22-1-2で、プログラムは、ボリューム管理テーブル112-11-2及びRAIDグループ管理テーブル112-11-1を参照して、物理ディスク及びデータのアドレスを判断する。ステップ112-22-1-3で、プログラムは、ディスク121の-slotからのデータの読出しを要求し、当該データをバッファに格納する。ステップ112-22-1-4で、プログラムは、RAIDグループ管理テーブル112-11-1を用いて比較的低い信頼性のストレージディスクによって割り当てられた外部ボリュームにデータが格納されているか否かをチェックする。データが比較的低い信頼性のストレージディスクに格納されている場合、プログラムは、ステップ112-22-1-5で、バッファ内のデータからハッシュ値を計算し、計算したハッシュ値と、外部ボリュームエラーコードテーブル112-11-3に格納されたエラーコードとを比較する。データが比較的低い信頼性のストレージディスクに格納されていない場合、プログラムは、ステップ112-22-1-9に進む。ステップ112-22-1-6で、プログラムは、比較された値が一致するか否かをチェックし、これによって、比較的低い信頼性のストレージディスクに格納されたデータエラーを検出することができる。比較された値が一致しない場合、プログラムは、ステップ112-22-1-7'において、外部ボリュームに当該データを回復するように要求する。その後、ステップ112-22-1-8'で、プログラムは、低信頼性ストレージサブシステム400自体が正しいデータを生成する代わりに、正しいデータを転送するのを待ち、ステップ112-22-1-3に進んで、破損したデータと置換された回復データが訂正されているかをチェックする。データが比較的低い信頼性のストレージディスクに格納されていて、かつ比較された値が一致する場合、プログラムは、ステップ112-22-1-9に進む。ステップ112-22-1-9で、プログラムは、バッファからキャッシュスロット112-30へスロットデータを転送し、それによってフラッシュ制御112-23及びデステージング制御112-22-2によって、訂正されたデータは、最終的に比較的低い信頼性のストレージシステムのディスク及びキャッシュに置き換えられることになる。プログラムは、112-22-1-10で終了する。

【0030】

0058 図21は、図1のストレージサブシステム400のメモリ412の一例を示す。図14のメモリ412との違いは、ディスク制御412-22にデータ回復制御412-22-4(図22)を含むということである。データ回復制御412-22-4は、重複データを使用することによって、指定領域のデータを回復する。

【0031】

0059 図22は、図21のメモリ412のデータ回復制御412-22-4のプロセスフローの一例を示す。プログラムは、412-22-4-1で開始する。ステップ412-22-4-2で、プログラムは、ボリューム管理テーブル112-11-2及びRAIDグループ管理テーブル112-11-1を参照して、物理ディスク及びデータのアドレスを判断する。ステップ412-22-4-3で、プログラムは、重複データを使用することによってデータを回復する。プログラムは、412-22-4-4で終了する。

【0032】

10

20

30

40

50

0060 図23は、図1のシステムで行われる読出しオペレーションの一例を示す。ホスト300は、高信頼性ストレージサブシステム100に読出しI/O要求を送信する(R1001)。高信頼性ストレージサブシステム100のCPU111は、読出しI/O要求を受信し、ステージング制御112-22-1を呼び出して、読出しI/Oデータをキャッシュスロット112-30-1に格納する。ステージング制御112-22-1は、データエラーが存在するか否かをチェックし、もしもデータエラーが存在する場合には、低信頼性ストレージサブシステム400に回復を要求した後、低信頼性ストレージサブシステム400によって受信されたデータ、正しいデータを転送する(R2002)。キャッシュ領域112-30は、外部ボリュームデータの読出しを要求し、データをホスト300に転送する(R1003)。低信頼性ストレージサブシステム100のCPU411は、読出しI/O要求を受信し、ステージング制御412-22-1を呼び出して、読出しI/Oデータをキャッシュスロット412-30-1に格納する(R1004)。キャッシュ領域412-30は、ディスク421からディスクデータを読み出すことを要求する(R1005)。ディスク421は、当該要求に従ってデータを送信する(R1006)。CPU411は、データ回復要求を受信し、データ回復制御412-22-4を呼び出して、データを回復する(R2007)。キャッシュ領域412-30は、外部ボリューム回復データの読出しを要求し、回復を実行する(R2008)。その後、ステップR1003~R1006が繰り返されて、回復されたデータが正しいか否かをチェックする。データが破損した場合、ステップW1008~W1010に示されるように、正しいデータを低信頼性ストレージサブシステム400のキャッシュ及びディスクに書き込まなければならない。CPU411は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、デステージング制御112-22-2を実行する(W1008)。キャッシュ領域412-30は、ダーティスロットデータをディスク421に転送する(W1009)。ディスク421は、当該データを受信し格納する(W1010)。

【0033】

0061 第2の実施形態では、回復プロセスが比較的低い信頼性のストレージサブシステムによって処理される。これによって、ストレージサブシステム100のより高い処理容量を可能にするが、これは負荷がストレージサブシステム400にシフトされるからである。しかし、正しいデータを計算するためのデータ処理は、ストレージコントローラ410によって行われるため、計算の正確度は、ストレージコントローラ110によって処理される場合よりも低くなるかもしれない。よって、本実施形態では、計算された正確なデータのハッシュ値は、高い信頼性を維持するために実際に使用される前に、メモリ112に格納されているエラーチェックコードに一致させられる。

【0034】

0062 第3の実施形態

【0035】

0063 本実施形態では、ストレージシステムは、二つ以上の低信頼性ストレージサブシステム400を有し、ここに重複データが格納される。よって、低信頼性ストレージサブシステム400のうちの一つから読み出されたデータが破損している場合、データは、もう一方の低信頼性ストレージサブシステム400から読み出される。第1の実施形態との差異についてのみ、図24~図28を用いて説明する。

【0036】

0064 図24は、本発明の方法及び装置が適用され得るシステムのハードウェア構成を示す。ストレージサブシステム100は、SAN(ストレージエリアネットワーク)200を介してホストコンピュータ300に接続される。ストレージサブシステム400は、ファイバ・チャネル(FC)を介してストレージサブシステム100に接続される。ストレージサブシステム100はホストコンピュータ200からI/Oコマンドを受信し、両方のストレージサブシステム100、400のストレージ装置121、421を使用してストレージボリュームをホストコンピュータ200に提供する。ストレージサブシステム100は、ストレージサブシステム400より高いデータ信頼性を有する。例えば

、ストレージサブシステム 100 に使用されるストレージ装置（例えば SAS）は、ストレージサブシステム 100 で使用されるもの（例えば SATA）に比べてより高い信頼性を有し、又は異なる RAID ランクが適用され得る。

【0037】

0065 ストレージサブシステム 100 は、CPU 111 と、メモリ 112 と、ストレージインタフェース 113、114 と、ディスクインタフェース 115 とを含むストレージコントローラ 110 を有する。CPU 111 は、ストレージサブシステム 100 を制御し、メモリ 112 からプログラム及びテーブルを読み出す。メモリ 112 はプログラム及びテーブルを格納する。ストレージインタフェース 113 は、ストレージネットワーク 200 を介してホストコンピュータ 300 に接続する。ストレージインタフェース 114 は、ストレージサブシステム 400 a、b のストレージインタフェースに接続する。ディスクインタフェース 115 は複数のストレージ装置 121 に接続し、これらはディスクユニット 120 に格納される。ストレージ装置 121 は、データを格納するためのソリッドステートデバイス（例えばフラッシュメモリ及び/又はハードディスクドライブ（HDD））から構成される。ストレージネットワーク 200 は、ストレージサブシステム 100 及びホストコンピュータ 300 に接続する。ホストコンピュータ 300 は、ストレージネットワーク 200 を介してストレージサブシステム 100 に I/O 要求を送信し、ストレージネットワーク 200 を介してストレージサブシステム 100 との間でデータを送受信する。ストレージサブシステム 400 a、b は基本的に、図 1 のストレージサブシステム 400 における構造と同じ構造を有する。

【0038】

0066 図 25 は、図 24 のストレージサブシステム 100 のメモリ 112 の一例を示す。メモリ 112 は、RAID グループ管理テーブル 112 - 11 - 1' と、ボリューム管理テーブル 112 - 11 と、外部ボリュームエラーチェックコードテーブル 112 - 11 - 3 とを含む、ストレージ管理テーブル 112 - 11 を含む。RAID グループ管理テーブル 112 - 11 - 1' は、ストレージ装置 121、外部ボリューム及びこれらグループの物理構造管理を提供し、かつ二つの外部ボリューム 441 間の重複構造を管理する。ボリューム管理テーブル 112 - 11 - 2 は論理ボリューム構成を提供する。外部ボリュームエラーチェックコードテーブル 112 - 11 - 3 は、外部ボリュームのいくつかの領域のためのエラーチェックコードを格納する。ある領域のエラーチェックコードの値は、ハッシュ計算によって当該領域に格納されたデータから計算される。キャッシュ管理テーブル 112 - 14 は、キャッシュデータ領域 112 - 30 の管理及び LRU/MRU 管理のために設けられる。ボリューム I/O 制御 112 - 21 は、書込み I/O 要件によって実行され、ライトデータを受信し、キャッシュデータ領域 112 に格納する書込み I/O 制御 112 - 21 - 1（図 8）と、読出し I/O 要件によって実行され、キャッシュデータ領域 112 からリードデータを送信する読出し I/O 制御 112 - 21 - 2（図 9）とを含む。ディスク制御 112 - 22 は、ディスク 121 からキャッシュデータ領域 112 へデータを転送するステージング制御 112 - 22 - 1（図 10）と、キャッシュデータ領域 112 からディスク 121 へデータを転送するデステージング制御 112 - 22 - 2（図 11）とを含む。メモリ 112 は、キャッシュデータ領域からディスク 121 へと定期的にダーティデータをフラッシュするフラッシュ制御 112 - 23（図 12）と、キャッシュデータ領域にキャッシュされたデータを見つけ、かつキャッシュデータ領域に新しいキャッシュ領域を割り当てるキャッシュ制御 112 - 24 とをさらに含む。メモリ 112 は、リード及びライトキャッシュデータを格納するキャッシュデータ領域 112 - 30 を含む。この領域は、複数のキャッシュスロット用に分割される。各キャッシュスロットはデータストライプのために割り当てられる。メモリ 112 は、プログラム実行スケジュールを制御し、マルチタスク環境をサポートするカーネル 112 - 40 を含む。プログラムが ACK（確認）を待っている場合、CPU 111 は別のタスクを実行するために変更する（例えば、ディスク 121 からキャッシュデータ領域 112 - 30 へのデータ転送待ち）。

10

20

30

40

50

【 0 0 3 9 】

0 0 6 7 図 2 6 は、図 2 のメモリ 1 1 2 の R A I D グループ管理テーブル 1 1 2 - 1 1 - 1 ' の一例を示す。R A I D グループ管理テーブル 1 1 2 - 1 1 - 1 ' は、R A I D グループの I D としての R A I D グループ番号 1 1 2 - 1 1 - 1 - 1 欄と、R A I D グループの構造を表す R A I D レベル 1 1 2 - 1 1 - 1 - 2 欄とを含む。例えば、数字は、R A I D レベルが当該数字であること（「5」は「R A I D レベルが 5 である」こと）を意味する。「N U L L」は、R A I D グループが存在しないことを意味する。「E x t」は、R A I D グループが内部ボリュームの外にある外部ボリュームとして存在することを意味する。R A I D グループ管理テーブル 1 1 2 - 1 1 - 1 は、内部ボリュームの場合には、R A I D グループに属する H D D の I D リストを表す H D D 番号の欄、外部ボリュームの場合には W W N の欄 1 1 2 - 1 1 - 1 - 3 を含む。R A I D グループが二つの外部ボリュームから構成される場合、この欄は、2 セットの W W N を含む。なぜならば、外部ボリュームは重複データを格納することになるからである。R A I D グループ管理テーブル 1 1 2 - 1 1 - 1 ' はさらに、重複領域を除く R A I D グループの全容量を表す R A I D グループ容量 1 1 2 - 1 1 - 1 - 4 を含む。

10

【 0 0 4 0 】

0 0 6 8 図 2 7 は、図 2 5 のメモリ 1 1 2 のステージング制御 1 1 2 - 2 2 - 1 のプロセスフローの一例を示す。プログラムは、1 1 2 - 2 2 - 1 - 1 で開始する。ステップ 1 1 2 - 2 2 - 1 - 2 では、プログラムは、ボリューム管理テーブル 1 1 2 - 1 1 - 2 及び R A I D グループ管理テーブル 1 1 2 - 1 1 - 1 ' を参照して、物理ディスク及びデータのアドレスを判断する。ステップ 1 1 2 - 2 2 - 1 - 3 で、プログラムは、ディスク 1 2 1 のスロットからデータを読み出すことを要求し、当該データをバッファに格納する。ステップ 1 1 2 - 2 2 - 1 - 4 では、プログラムは、データが外部ボリュームに格納されているか否かをチェックする。当該データが外部ボリュームに格納されている場合、プログラムは、ステップ 1 1 2 - 2 2 - 1 - 5 において、バッファ内のデータからハッシュ値を計算し、計算したハッシュ値と外部ボリュームエラーコードテーブル 1 1 2 - 1 1 - 3 に格納されたエラーコードとを比較する。データが比較的低い信頼性のストレージディスクに格納されていない場合、プログラムはステップ 1 1 2 - 2 2 - 1 - 9 に進む。ステップ 1 1 2 - 2 2 - 1 - 6 で、プログラムは、比較された値が一致するか否かをチェックすることで、比較的低い信頼性のストレージディスクに格納されたデータエラーを検出することができる。比較された値が一致しない場合、プログラムは、ステップ 1 1 2 - 2 2 - 1 - 7 ' ' で他方の外部ボリュームから回復データを読み出す。そして、ステップ 1 1 2 - 2 2 - 1 - 8 ' ' では、プログラムは、回復されたスロットに対してダーティ属性を設定する。正しいデータはバッファに格納されることになる。外部ボリュームは重複データを格納するので、正しいデータを生成する必要はない。当該データが外部ボリュームに格納されており、かつ比較された値が一致する場合、プログラムはステップ 1 1 2 - 2 2 - 1 - 9 に進む。ステップ 1 1 2 - 2 2 - 1 - 9 で、プログラムは、バッファからキャッシュスロット 1 1 2 - 3 0 へスロットデータを転送することによって、訂正されたデータは、フラッシュ制御 1 1 2 - 2 3 及びデステージング制御 1 1 2 - 2 2 - 2 によって、比較的低い信頼性のストレージシステム（ハッシュ値が一致しなかったデータを含む）の中のディスク及びキャッシュに最終的に置き換えられる。プログラムは、1 1 2 - 2 2 - 1 - 1 0 で終了する。

20

30

40

【 0 0 4 1 】

0 0 6 9 図 2 8 は、図 2 4 のシステムで行われる読出しオペレーションの一例を示す。ホスト 3 0 0 は、高信頼性ストレージサブシステム 1 0 0 に読出し I / O 要求を送信する（R 1 0 0 1）。高信頼性ストレージサブシステム 1 0 0 の C P U 1 1 1 は、読出し I / O 要求を受信し、ステージング制御 1 1 2 - 2 2 - 1 を呼び出して、読出し I / O データをキャッシュスロット 1 1 2 - 3 0 - 1 に格納する。ステージング制御 1 1 2 - 2 2 - 1 は、データエラーが存在するか否かをチェックし、もしもデータエラーが存在する場合には、他方の外部ボリュームからデータを読み出し、その後、当該データをホスト 3 0 0

50

に転送する（R3002）。キャッシュ領域112-30は、外部ボリュームデータの読出しを要求する（R1003）。低信頼性ストレージサブシステム100のCPU411は、読出しI/O要求を受信し、ステージング制御412-22-1を呼び出して、読出しI/Oデータをキャッシュスロット412-30-1に格納する（R1004）。キャッシュ領域412-30は、ディスク421からディスクデータを読み出すことを要求する（R1005）。ディスク421は、当該要求に従ってデータを送信する（R1006）。低信頼性ストレージサブシステム400aに格納されたデータが破損していた場合、ステップW1004~W1010に示されるように、低信頼性ストレージサブシステム400bによって取得された正しいデータが、低信頼性ストレージサブシステム400aのキャッシュ及びディスクに書き込まなければならない。CPU111は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、デステージング制御112-22-2を実行し、これによってエラーチェックコードを生成する（W1004）。キャッシュ領域412-30は、ダーティスロットデータを外部ボリュームに転送する（W1005）。低信頼性ストレージサブシステム400のCPU411は、書込みI/O要求を受信し、データを低信頼性ストレージサブシステム400のキャッシュスロット412-30-1に格納する（W1006）。キャッシュ領域412-30は書込みI/Oデータを受信する（W1007）。CPU411は、フラッシュ制御112-23によってダーティキャッシュスロットを見つけ、デステージング制御112-22-2を実行する（W1008）。キャッシュ領域412-30は、ダーティスロットデータをディスク421に転送する（W1009）。ディスク421は、当該データを受信し格納する（W1010）。

10

20

【0042】

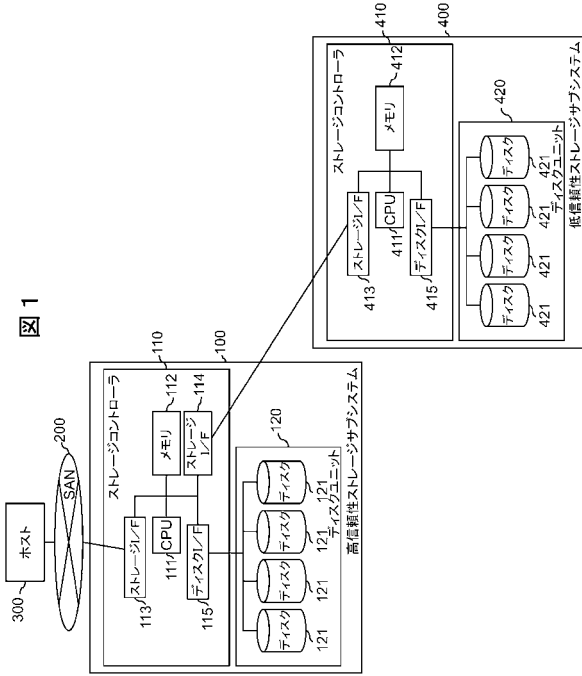
0070 第3の実施形態では、ストレージサブシステム100又は400のいずれによっても回復プロセスは必要とされない。このことにより、ストレージサブシステム100、400のより高い処理容量を可能にする。但し、データは二つの外部ストレージシステムに書き込む必要がある。

【0043】

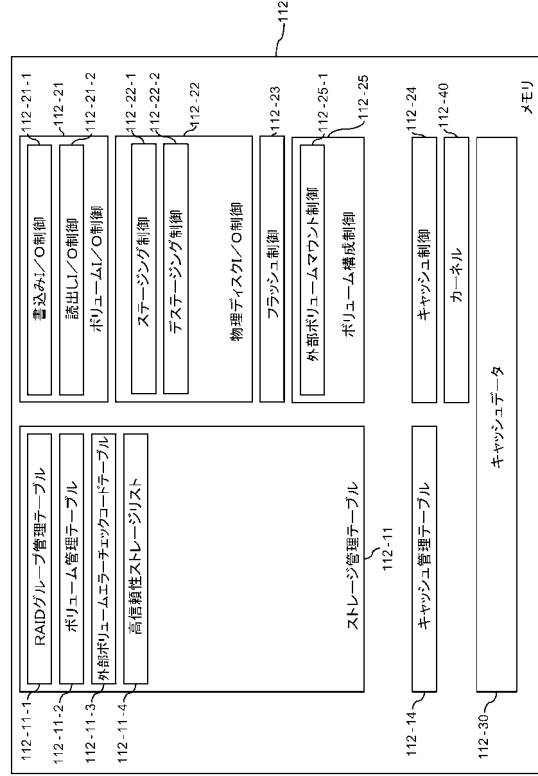
0071 本発明は、比較的低い信頼性のストレージディスクを使用することによってコストを節減することができ、比較的高い信頼性のストレージシステムにおいて比較的低い信頼性のストレージディスクのためのエラーコードを維持することによって、システム全体としての高い信頼性を維持するストレージシステムを提供する。但し、本発明について上述したが、添付の請求の範囲によって定義される本発明の精神から逸脱することなしに、本発明への多くの修正が、本発明に関わる分野の当業者にとって明らかになるであろう。

30

【図1】



【図2】



【図3】

図3

RAIDグループ番号	RAIDレベル	ディスク番号	容量	信頼性
0	5	0-3	900[GB]	高
1	5	4-7	3000[GB]	高
2	5	8-11	3000 [GB]	高
3	Ext/5	12:3 4:56 :78:9 A:B C/0	0[GB]	低
4	Ext/1	12:3 4:56 :78:9 A:B C/1	0[GB]	低
5	NULL	NULL	0[GB]	NULL
6	Ext/5	11:2 2:33 :44:5 5:6 6/0	1500 [GB]	高
7	10	68-72	1500 [GB]	高

RAIDグループ管理テーブル

【図4】

図4

ボリューム番号	容量	RAIDグループ番号	アドレス範囲	ポート番号	LUN
0	10[GB]	1	0x00000000 - 0x0FFFFFFF	0	0
1	30[GB]	0	0x00000000 - 0x2FFFFFFF	0	1
2	20[GB]	1	0x10000000 - 0x2FFFFFFF	0	2
3	60[GB]	7	0x00000000 - 0x5FFFFFFF	0	3
4	N/A	N/A	N/A	N/A	N/A
5	60[GB]	2	0x00000000 - 0x5FFFFFFF	1	0
6	N/A	N/A	N/A	N/A	N/A
7	N/A	N/A	N/A	N/A	N/A

ボリューム管理テーブル

【 図 5 】

112-11-3-1	112-11-3-2	112-11-3-3
仮想ボリューム番号	スロット番号	エラーチェックコード
0	0	48c959b9b900656bb883 d6f91ab1cdec12a41de2
0	1	6c272d271202e4d48936e 5e395a7ea90e0ff98566
0	2	10ea371937de976c0622 06f3564c064e795c0b13
0	3	36a0b960d7236803c026 7459517297360e67c1cc
0	4	191779748647cba5c123 bca2d001bcaead13cc25
0	5	0aa42567b1810aa0c764 ebe715773d2b57db153a
1	0	1f9155a6b9e94b792a12c dafa3a10172b95cef2e

外部ボリュームエラーチェックコードテーブル
112-11-3

【 図 6 】

図 6

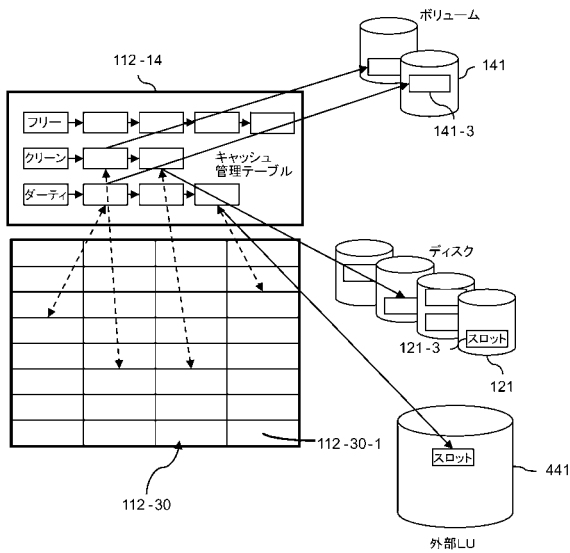
112-14-1	112-14-2	112-14-3	112-14-4
インデックス	ディスク番号	LBA	次
0	2	0xA00	1
1	1	0x7E000	2
2	1	0x9700	3
3	0	0x0000	NULL
4	2	0xC500	5
5	1	0x1100	6
6	1	0xFF00	NULL

112-14-5	112-14-6
キューの種類	ポインタ
フリー	2
クリーン	1
ダーティ	4

キャッシュ管理テーブル
112-14

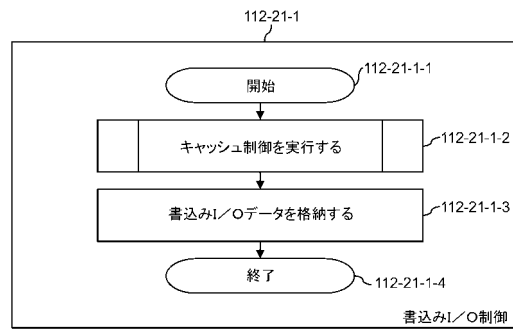
【 図 7 】

図 7

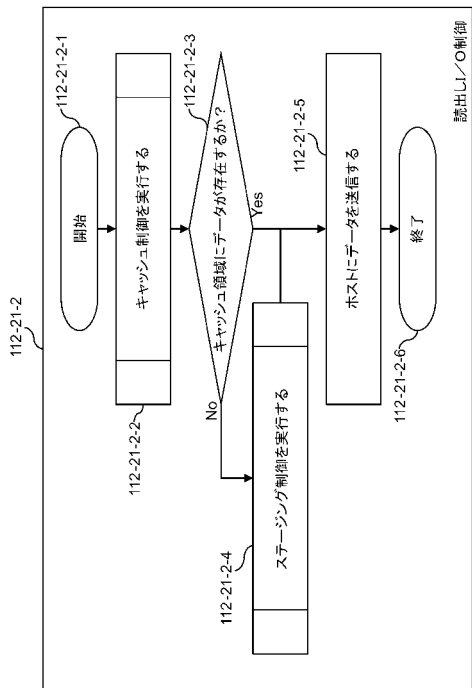


【 図 8 】

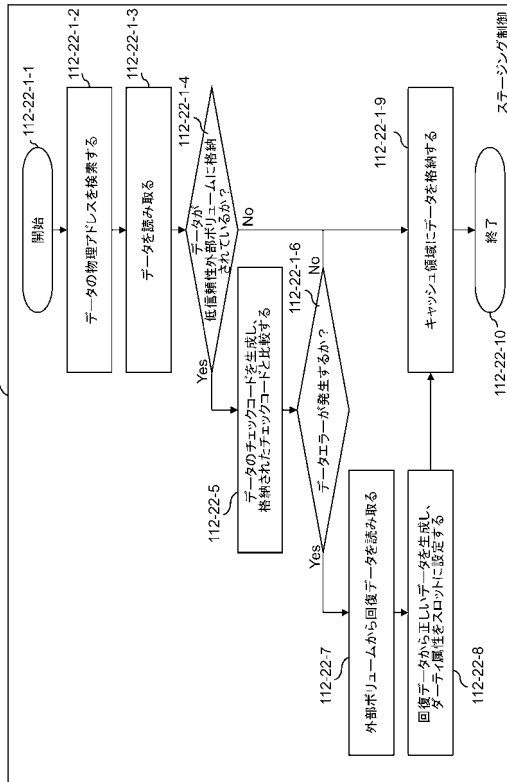
図 8



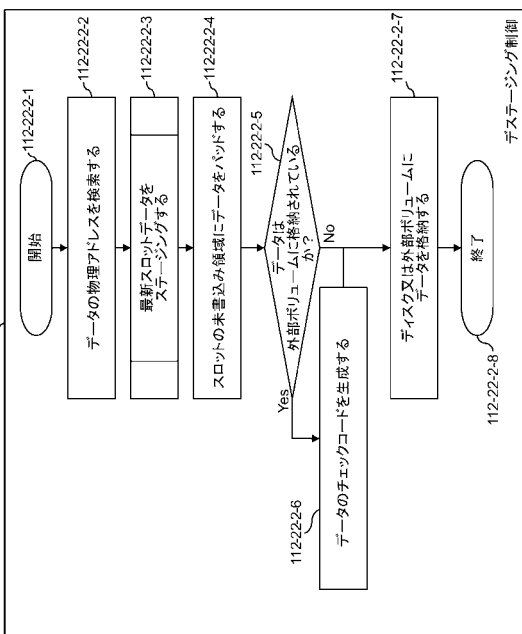
【 図 9 】



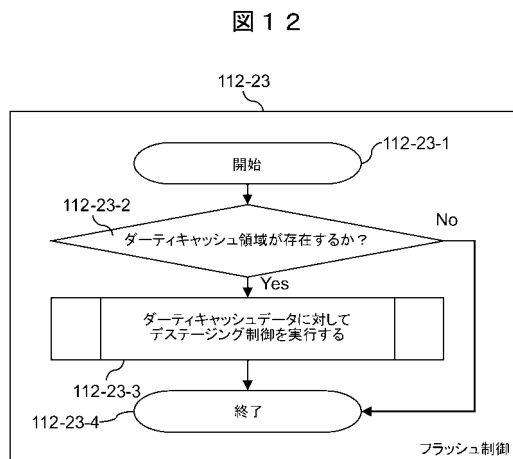
【 図 10 】



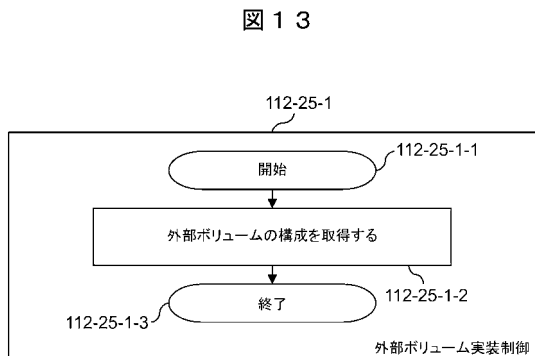
【 図 11 】



【 図 12 】

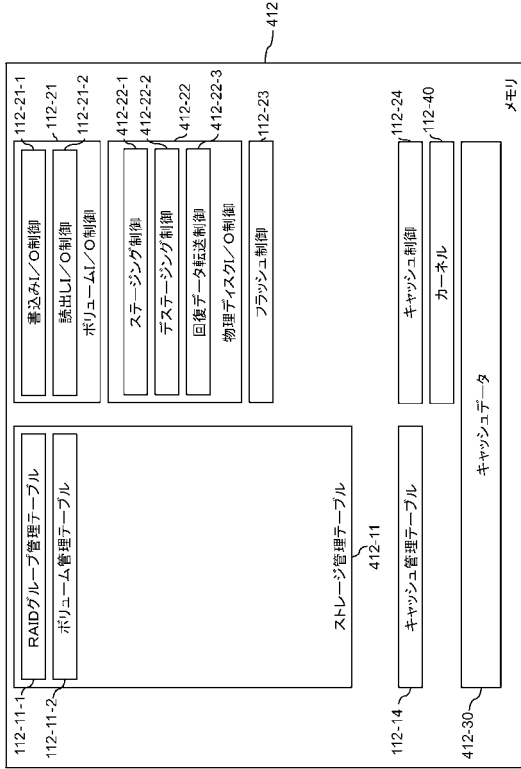


【 図 13 】



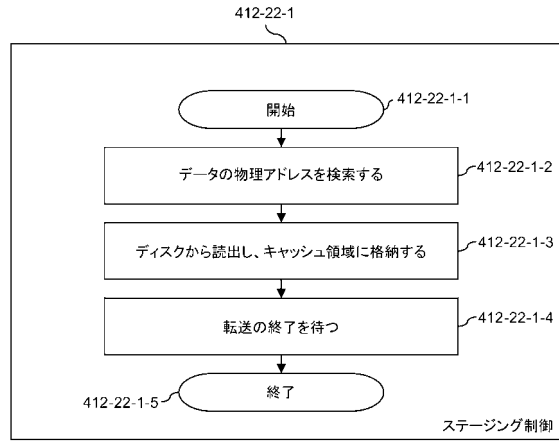
【 図 1 4 】

図 1 4



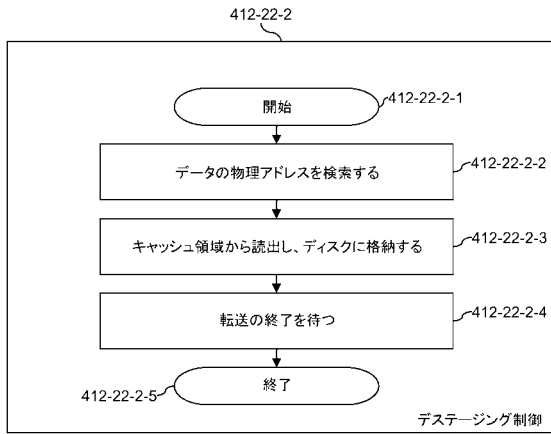
【 図 1 5 】

図 1 5



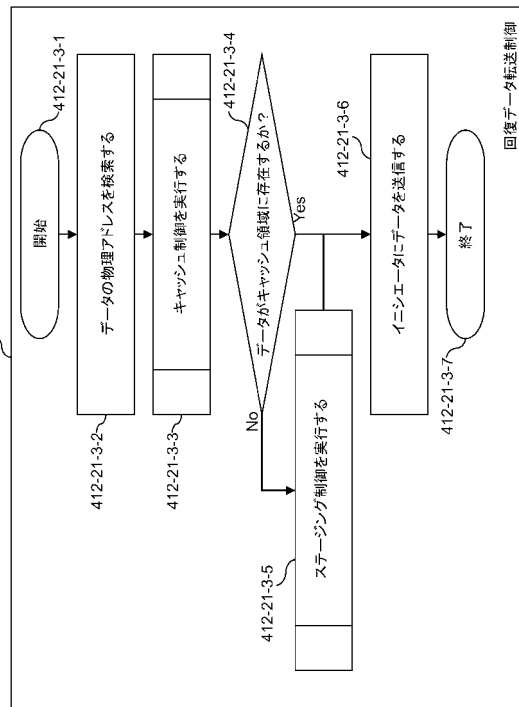
【 図 1 6 】

図 1 6



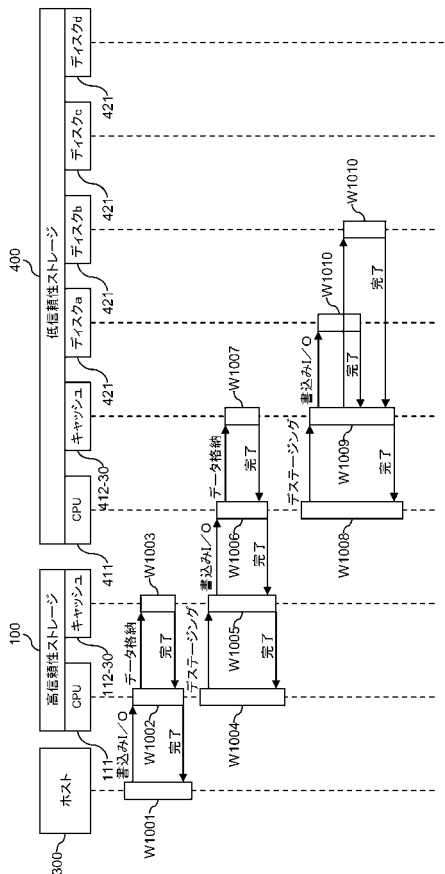
【 図 1 7 】

図 1 7



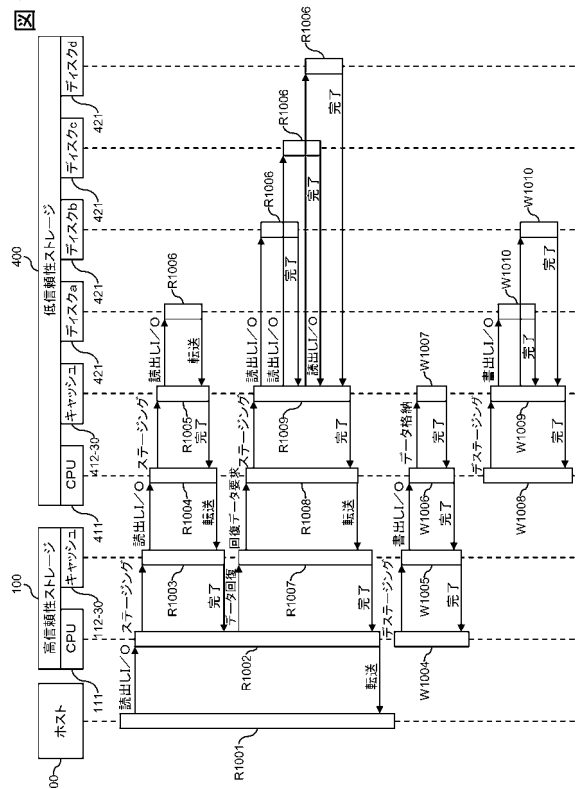
【図18】

図18



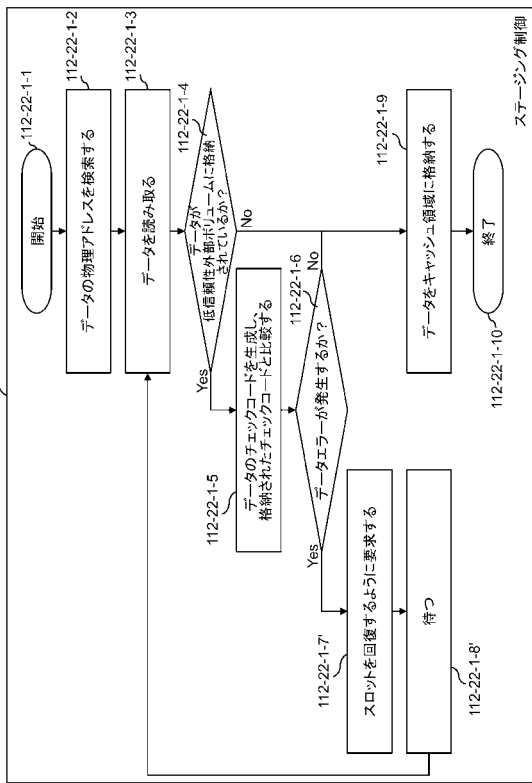
【図19】

図19



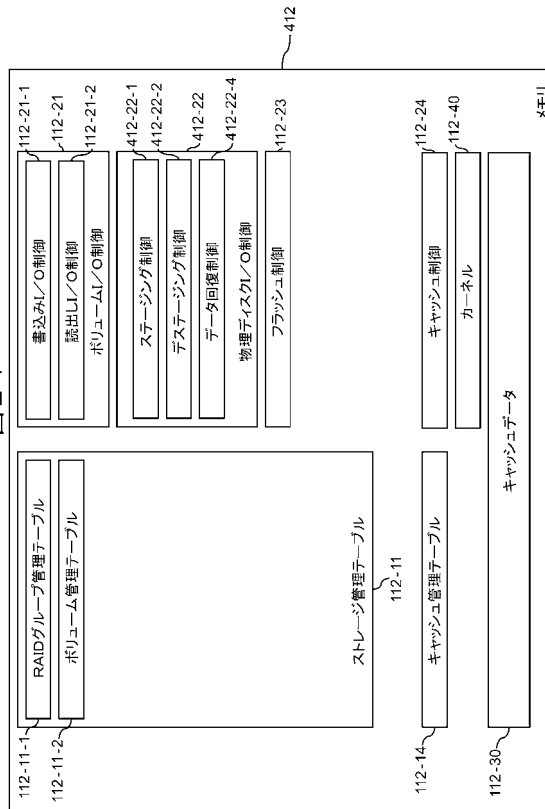
【図20】

図20

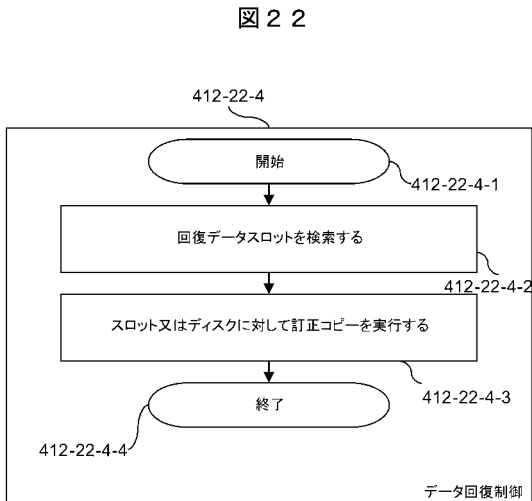


【図21】

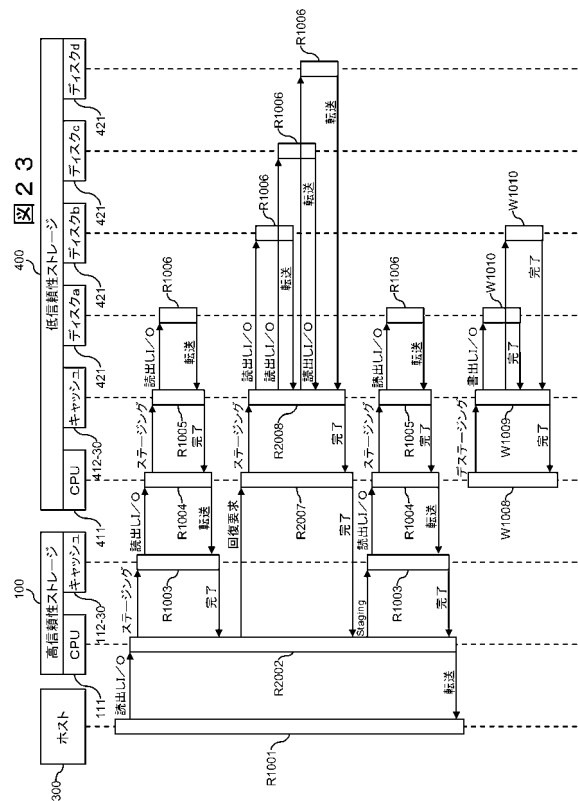
図21



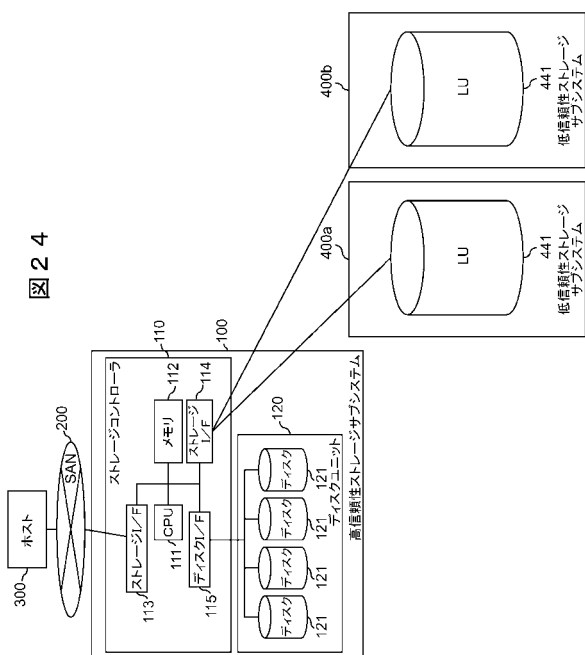
【図22】



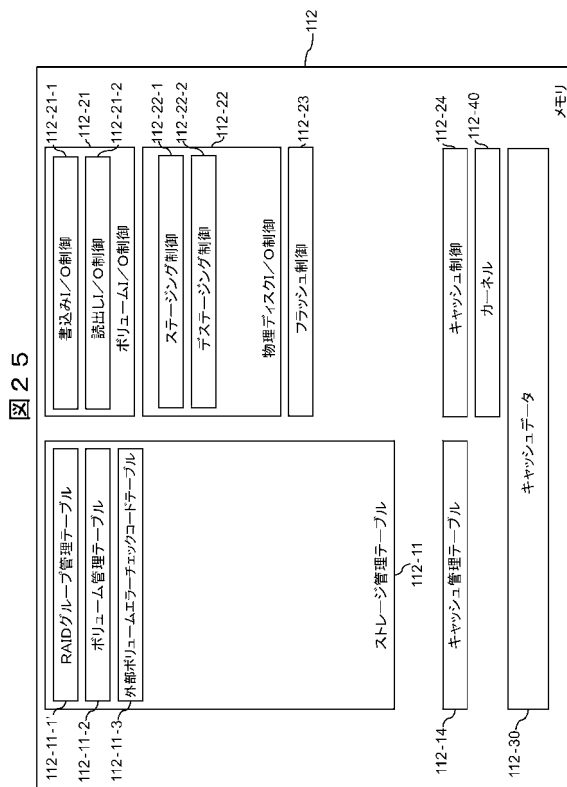
【図23】



【図24】



【図25】

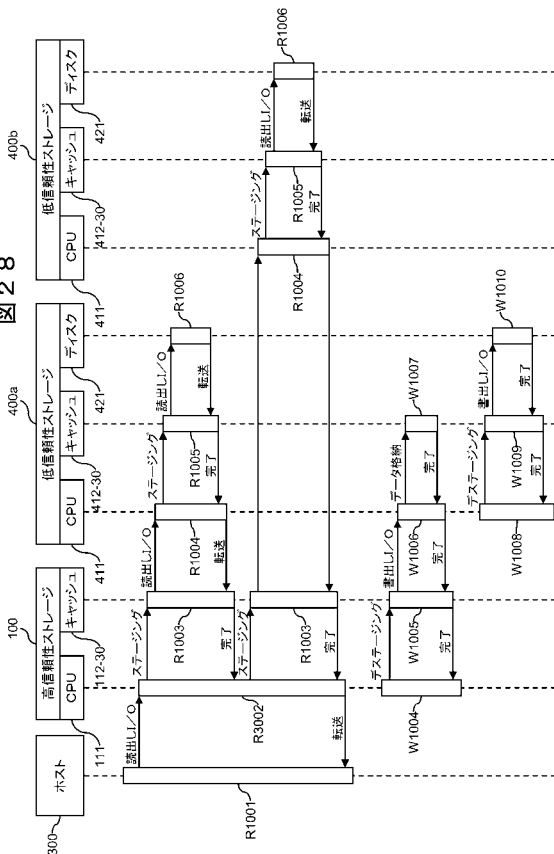


【 図 26 】

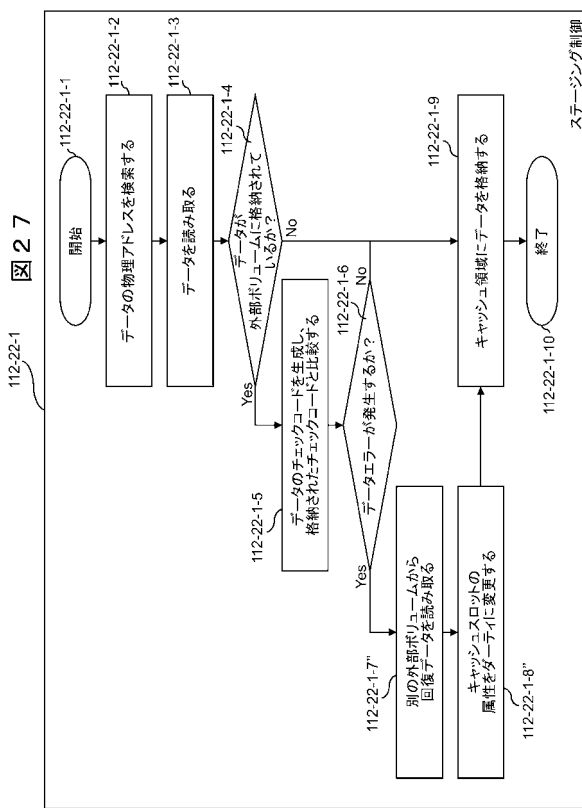
RAIDグループ番号	RAIDレベル	ディスク番号	容量
0	5	0-3	900[GB]
1	5	4-7	3000[GB]
2	5	8-11	3000[GB]
3	Ext	12,34,56,78,9A,BC/0, 12,34,56,78,9A,BC/1	0[GB]
4	Ext	12,34,56,78,9A,BC/2, 12,34,56,78,9A,BC/3	0[GB]
5	NULL	NULL	0[GB]
6	10	64-67	1500[GB]
7	10	68-72	1500[GB]

RAIDグループ管理テーブル

【 図 28 】



【 図 27 】



【 図 29 】

製品タイプID
AAA0001112223
AAA0001112224
AAA0001112225
BASKLF000
BASKLF001
BASKLF002
CA0000000000007
CA0000000000008

高信頼性ストレージリスト

フロントページの続き

審査官 稲葉 崇

- (56)参考文献 特開2005-025683(JP,A)
特開2008-117253(JP,A)
特開2010-026873(JP,A)
特開2009-104420(JP,A)
特開2008-293350(JP,A)
特開2007-026453(JP,A)
特開2005-165619(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08

G06F 12/00

G06F 12/16