

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2011-221981

(P2011-221981A)

(43) 公開日 平成23年11月4日(2011.11.4)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 304Z	5B065
G06F 12/00 (2006.01)	G06F 12/00 514A	5B082

審査請求 未請求 請求項の数 18 O L 外国語出願 (全 26 頁)

(21) 出願番号	特願2010-192834 (P2010-192834)	(71) 出願人	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
(31) 優先権主張番号	12/756, 475	(74) 代理人	100129218 弁理士 百本 宏之
(32) 優先日	平成22年4月8日 (2010. 4. 8)	(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
(33) 優先権主張国	米国 (US)	(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 Fターム(参考) 5B065 BA01 CC08 EA03 ZA08 5B082 DB00 FA12

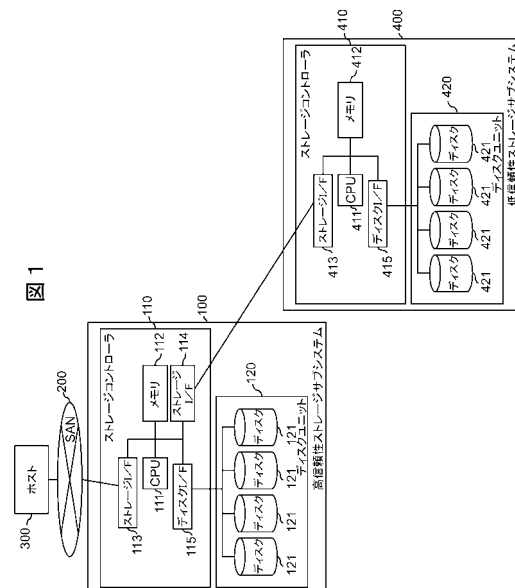
(54) 【発明の名称】 外部ストレージシステムに結合されたストレージシステムのエラーコード管理方法及び装置

(57) 【要約】 (修正有)

【課題】 システム全体としての信頼性を高める。

【解決手段】 複数レベルの信頼性を有するストレージ装置を用いた複数のストレージシステムを備えるシステムに於いて、比較的高い信頼性のストレージシステムに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信頼性が高まる。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。

【選択図】 図1



【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、
前記複数の第 2 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、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 オペレーションに回答して、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回答して、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た場合、前記第 2 ストレージシステムに格納された前記データは
、前記ホストコンピュータへ転送される、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前記デー
タは、前記第 2 ストレージコントローラによって回復された後、前記第 2 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 4 に記載のシステム。

【請求項 6】

前記第 2 ストレージシステムに格納されたデータは R A I D 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 5 に記載のシステム。

40

【請求項 7】

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 4 に記載のシステム。

【請求項 8】

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、

50

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。

【請求項9】

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8に記載のシステム。

10

【請求項10】

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た前記データは、前記ホストコンピュータへ転送される、請求項9に記載のシステム。

【請求項11】

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、前記受信データのために前記第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、

20

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10に記載のシステム。

【請求項12】

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、

ホストコンピュータから前記第1ストレージシステムを介して前記外部ストレージシステムへと書込みI/Oオペレーションを受信することと、

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によって、前記書込みI/Oオペレーションによって書き込まれるデータのエラーコードを計算することと、

30

前記エラーコードを前記第1ストレージシステムに格納することと、

前記書込みオペレーションによって書き込まれる前記データを前記外部ストレージシステムのボリュームに格納することと

を含む、方法。

【請求項13】

ホストコンピュータから前記第1ストレージシステムを介して前記外部ストレージシステムへと読出しI/Oオペレーションを受信することと、

前記第1ストレージシステムによって読出しオペレーションに従ってデータのエラーチェックコードを計算することと、

前記計算されたエラーチェックコードと、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を比較することと、

40

前記比較結果が一致した場合、読出しオペレーションに従った前記データを、前記外部ストレージシステムから前記ホストコンピュータへ転送することと

をさらに備える、請求項12に記載の方法。

【請求項14】

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の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正しいデータを計算することをさらに備える、請求項13に記載の方法。

【請求項15】

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正しいデータを書き

50

込むことをさらに備える、請求項 14 に記載の方法。

【請求項 16】

前記比較結果が一致しない場合、前記第 2 ストレージコントローラによって正しいデータを計算することをさらに備える、請求項 13 に記載の方法。

【請求項 17】

前記正しいデータのハッシュ値を計算して、前記正しい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、

前記正しい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た場合、前記正しいデータを前記ホストコンピュータに転送することと

10

をさらに備える、請求項 16 に記載の方法。

【請求項 18】

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 13 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。

20

【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。

30

【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。

40

【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。

50

- 【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。
- 【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。
- 【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。
- 【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】
- 【0006】

10

20

30

40

50

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・チャンネル（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。

10

【0008】

0036 ストレージサブシステム100は、CPU111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。CPU111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、CPU411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。CPU411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。

20

30

40

【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

50

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。

10

【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。

20

【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。40
「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。

30

40

【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

50

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。

10

【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。

20

【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スロット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、回復されたスロット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

30

40

50

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 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、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、1 1 2 - 2 2 - 1 - 1 0 で終了する。

【0018】

0046 図 1 1 は、図 2 のメモリ 1 1 2 のデステージング制御 1 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、1 1 2 - 2 2 - 2 - 1 で開始する。ステップ 1 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 1 1 2 - 2 2 - 2 - 3 で、プログラムはステージング制御 1 1 2 - 2 2 - 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 1 1 2 - 2 2 - 2 - 4 では、プログラムは、キャッシュスロット 1 1 2 - 3 0 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 1 1 2 - 2 2 - 2 - 5 で、プログラムは、RAID グループ管理テーブル 1 1 2 - 1 1 - 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 1 1 2 - 2 2 - 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 1 1 2 - 1 1 - 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 1 1 2 - 2 2 - 2 - 7 に進む。ステップ 1 1 2 - 2 2 - 2 - 7 で、プログラムはキャッシュデータ領域 1 1 2 - 3 0 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、1 1 2 - 2 2 - 2 - 8 で終了する。

10

20

【0019】

0047 図 1 2 は、図 2 のメモリ 1 1 2 のフラッシュ制御 1 1 2 - 2 3 のプロセスフローの一例を示す。プログラムは、1 1 2 - 2 3 - 1 で開始する。ステップ 1 1 2 - 2 3 - 2 で、プログラムは、キャッシュ管理テーブル 1 1 2 - 1 4 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 1 1 2 - 2 3 - 3 で、見つかったダーティキャッシュスロット 1 1 2 - 3 0 - 1 のためにデステージング制御 1 1 2 - 2 2 - 2 を呼び出す。プログラムは、1 1 2 - 2 3 - 4 で終了する。

30

【0020】

0048 図 1 3 は、図 2 のメモリ 1 1 2 の外部ボリューム実装制御 1 1 2 - 2 5 - 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 5 - 1 - 1 で開始する。ステップ 1 1 2 - 2 5 - 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 1 1 2 - 1 1 - 1 の信頼性 1 1 2 - 1 1 - 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 1 1 2 - 1 1 - 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 1 1 2 - 1 1 - 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 1 1 2 - 1 1 - 1 - 5 に「低い」と格納する。プログラムは、1 1 2 - 2 5 - 1 - 3 で終了する。

40

【0021】

0049 図 1 4 は、図 1 のストレージサブシステム 4 0 0 のメモリ 4 1 2 の一例を示す。メモリ 4 1 2 は、RAID グループ管理テーブル 1 1 2 - 1 1 - 1 及びボリューム管理テーブル 1 1 2 - 1 1 (これらは 1 1 2 - 1 1 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 4 1 2 - 1 1 を含む。しかし、ストレージ管理テーブル 4 1 2 - 1 1 は、メモリ 1 1 2 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 1 1 2 -

50

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412-30からディスク421へデータを転送するデステージング制御412-22-2(図16)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412-22-3(図17)とを含む。メモリ112は、メモリ112の場合のように、キャッシュデータ領域からディスク421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412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412-30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412は、メモリ112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112-40を含む。

10

20

【0022】

0050 図15は、図14のメモリ412のステージング制御412-22-1のプロセスフローの一例を示す。プログラムは、412-22-1-1で開始する。ステップ412-22-1-2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412-22-1-3で、プログラムは、ディスク421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412-30に格納する。ステップ412-22-1-4では、プログラムはデータ転送の終了を待つ。メモリ412のカーネル112-40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、412-22-1-5で終了する。

30

【0023】

0051 図16は、図14のメモリ412のデステージング制御412-22-2のプロセスフローの一例を示す。プログラムは、412-22-2-1で開始する。ステップ412-22-2-2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412-22-2-3では、プログラムは、キャッシュデータ領域412-30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421に格納する。ステップ412-22-2-4で、プログラムはデータ転送の終了を待つ。メモリ412のカーネル112-40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、412-22-2-5で終了する。

40

【0024】

0052 図17は、図14のメモリ412の回復データ転送制御412-21-3のプロセスフローの一例を示す。プログラムは、412-21-3-1で開始する。ステップ412-21-3-2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412-21-3-3では、プログラムは、キャッシュ制御112-24を呼び出して、対応するキャッシュスロット412-30-1を検索する。ステップ412-21-3-4では、プログラムは、前記キャッシュスロット412-30-1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412-21-3-5で、プログラムはステージング制御412-21-1を呼

50

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412-21-3-6に移動する。ステップ412-21-3-6では、プログラムは、キャッシュスロット112-30-1データをイニシエータに転送する。よって、メモリ112のステージング制御112-22-1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110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100に正しいデータを生成することができる。プログラムは、412-21-3-7で終了する。

【0025】

0053 図18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。10
 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4
 12-30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。30
 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
 40
 50

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】

0055 第2の実施形態

10

【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のスロット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。

20

30

40

【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。

50

【 0 0 3 1 】

0 0 5 9 図 2 2 は、図 2 1 のメモリ 4 1 2 のデータ回復制御 4 1 2 - 2 2 - 4 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 - 4 - 1 で開始する。ステップ 4 1 2 - 2 2 - 4 - 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 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 4 1 2 - 2 2 - 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、4 1 2 - 2 2 - 4 - 4 で終了する。

【 0 0 3 2 】

0 0 6 0 図 2 3 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。10
 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、低信頼性ストレージサブシステム 4 0 0 に回復を要求した後、低信頼性ストレージサブシステム 4 0 0 によって受信されたデータ、正しいデータを転送する (R 2 0 0 2)。キャッシュ領域 1 1 2 - 3 0 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 3 0 0 に転送する (R 1 0 0 3)。低信頼性ストレージサブシステム 1 0 0 の C P U 4 1 1 は、読出し I / O 要求を受信し、ステージング制御 4 1 2 - 2 2 - 1 を呼び出して 20
 、読出し I / O データをキャッシュロット 4 1 2 - 3 0 - 1 に格納する (R 1 0 0 4)。キャッシュ領域 4 1 2 - 3 0 は、ディスク 4 2 1 からディスクデータを読み出すことを要求する (R 1 0 0 5)。ディスク 4 2 1 は、当該要求に従ってデータを送信する (R 1 0 0 6)。C P U 4 1 1 は、データ回復要求を受信し、データ回復制御 4 1 2 - 2 2 - 4 を呼び出して、データを回復する (R 2 0 0 7)。キャッシュ領域 4 1 2 - 3 0 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 (R 2 0 0 8)。その後、ステップ R 1 0 0 3 ~ R 1 0 0 6 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 W 1 0 0 8 ~ W 1 0 1 0 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 4 0 0 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。C P U 4 1 1 は、フラッシュ制御 1 1 2 - 2 3 によってダーティキ 30
 ャッシュロットを見つけ、デステージング制御 1 1 2 - 2 2 - 2 を実行する (W 1 0 0 8)。キャッシュ領域 4 1 2 - 3 0 は、ダーティロットデータをディスク 4 2 1 に転送する (W 1 0 0 9)。ディスク 4 2 1 は、当該データを受信し格納する (W 1 0 1 0)。

【 0 0 3 3 】

0 0 6 1 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 1 0 0 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 4 0 0 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 4 1 0 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 1 1 0 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 1 1 2 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。

【 0 0 3 4 】

0 0 6 2 第 3 の実施形態

【 0 0 3 5 】

0 0 6 3 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 4 0 0 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 4 0 0 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 4 0 0 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 2 4 ~ 図 2 8 を用いて説明する。

10

20

30

40

50

【 0 0 3 6 】

0 0 6 4 図 2 4 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 1 0 0 は、S A N (ストレージエリアネットワーク) 2 0 0 を介してホストコンピュータ 3 0 0 に接続される。ストレージサブシステム 4 0 0 は、ファイバ・チャネル (F C) を介してストレージサブシステム 1 0 0 に接続される。ストレージサブシステム 1 0 0 はホストコンピュータ 2 0 0 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 1 0 0、4 0 0 のストレージ装置 1 2 1、4 2 1 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 2 0 0 に提供する。ストレージサブシステム 1 0 0 は、ストレージサブシステム 4 0 0 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 1 0 0 に使用されるストレージ装置 (例えば S A S) は、ストレージサブシステム 1 0 0 で使用されるもの (例えば S A T A) に比べてより高い信頼性を有し、又は異なる R A I D ランクが適用され得る。

10

【 0 0 3 7 】

0 0 6 5 ストレージサブシステム 1 0 0 は、C P U 1 1 1 と、メモリ 1 1 2 と、ストレージインタフェース 1 1 3、1 1 4 と、ディスクインタフェース 1 1 5 とを含むストレージコントローラ 1 1 0 を有する。C P U 1 1 1 は、ストレージサブシステム 1 0 0 を制御し、メモリ 1 1 2 からプログラム及びテーブルを読み出す。メモリ 1 1 2 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 1 1 3 は、ストレージネットワーク 2 0 0 を介してホストコンピュータ 3 0 0 に接続する。ストレージインタフェース 1 1 4 は、ストレージサブシステム 4 0 0 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 1 1 5 は複数のストレージ装置 1 2 1 に接続し、これらはディスクユニット 1 2 0 に格納される。ストレージ装置 1 2 1 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 (H D D)) から構成される。ストレージネットワーク 2 0 0 は、ストレージサブシステム 1 0 0 及びホストコンピュータ 3 0 0 に接続する。ホストコンピュータ 3 0 0 は、ストレージネットワーク 2 0 0 を介してストレージサブシステム 1 0 0 に I / O 要求を送信し、ストレージネットワーク 2 0 0 を介してストレージサブシステム 1 0 0 との間でデータを送受信する。ストレージサブシステム 4 0 0 a、b は基本的に、図 1 のストレージサブシステム 4 0 0 における構造と同じ構造を有する。

20

【 0 0 3 8 】

0 0 6 6 図 2 5 は、図 2 4 のストレージサブシステム 1 0 0 のメモリ 1 1 2 の一例を示す。メモリ 1 1 2 は、R A I D グループ管理テーブル 1 1 2 - 1 1 - 1' と、ボリューム管理テーブル 1 1 2 - 1 1 と、外部ボリュームエラーチェックコードテーブル 1 1 2 - 1 1 - 3 とを含む、ストレージ管理テーブル 1 1 2 - 1 1 を含む。R A I D グループ管理テーブル 1 1 2 - 1 1 - 1' は、ストレージ装置 1 2 1、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 4 4 1 間の重複構造を管理する。ボリューム管理テーブル 1 1 2 - 1 1 - 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 1 1 2 - 1 1 - 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 1 1 2 - 1 4 は、キャッシュデータ領域 1 1 2 - 3 0 の管理及び L R U / M R U 管理のために設けられる。ボリューム I / O 制御 1 1 2 - 2 1 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 1 1 2 に格納する書込み I / O 制御 1 1 2 - 2 1 - 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 1 1 2 からリードデータを送信する読出し I / O 制御 1 1 2 - 2 1 - 2 (図 9) とを含む。ディスク制御 1 1 2 - 2 2 は、ディスク 1 2 1 からキャッシュデータ領域 1 1 2 へデータを転送するステージング制御 1 1 2 - 2 2 - 1 (図 1 0) と、キャッシュデータ領域 1 1 2 からディスク 1 2 1 へデータを転送するデステージング制御 1 1 2 - 2 2 - 2 (図 1 1) とを含む。メモリ 1 1 2 は、キャッシュデータ領域からディスク 1 2 1 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 1 1 2 - 2 3 (図 1 2) と、

30

40

50

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 1 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 1 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 1 1 2 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 1 1 2 - 4 0 を含む。プログラムが A C K (確認) を待っている場合、C P U 1 1 1 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 1 2 1 からキャッシュデータ領域 1 1 2 - 3 0 へのデータ転送待ち) 。

【 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。

【 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

10

20

30

40

50

- 10で終了する。

【0041】

0069 図28は、図24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、他方の外部ボリュームからデータを読出し、その後、当該データをホスト300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

30

【0042】

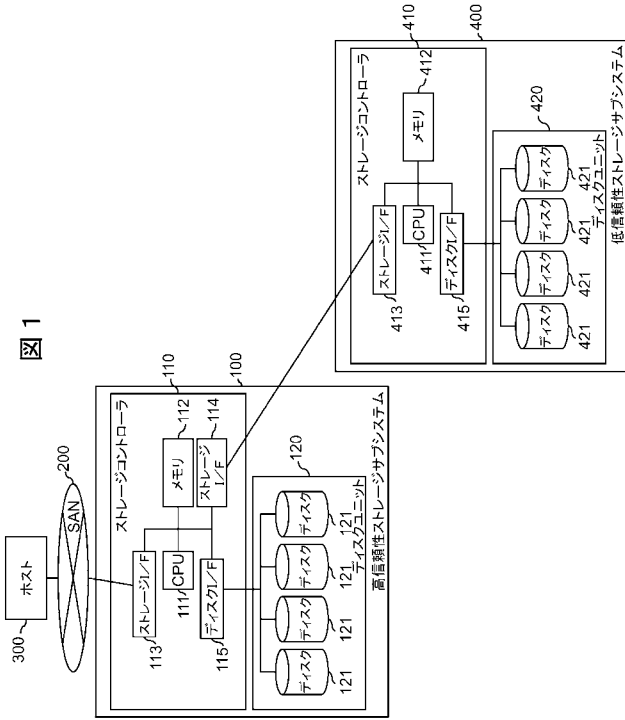
0070 第3の実施形態では、ストレージサブシステム100又は400のいずれによっても回復プロセスは必要とされない。このことにより、ストレージサブシステム100、400のより高い処理容量を可能にする。但し、データは二つの外部ストレージシステムに書き込む必要がある。

【0043】

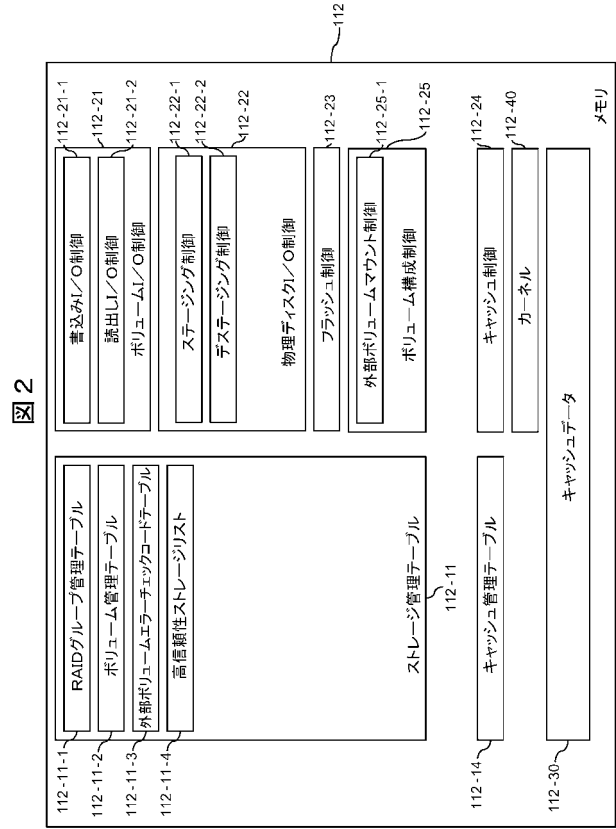
0071 本発明は、比較的低い信頼性のストレージディスクを使用することによってコストを削減することができ、比較的高い信頼性のストレージシステムにおいて比較的低い信頼性のストレージディスクのためのエラーコードを維持することによって、システム全体としての高い信頼性を維持するストレージシステムを提供する。但し、本発明について上述したが、添付の請求の範囲によって定義される本発明の精神から逸脱することなしに、本発明への多くの修正が、本発明に関わる分野の当業者にとって明らかになるであろう。

40

【 図 1 】



【 図 2 】



【 図 3 】

RAIDグループ管理テーブル

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

【 図 4 】

ボリューム管理テーブル

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

【 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	48c959b8b900566b883 dbf91ab1ddecf2a41de2
0	1	6d272d71202ed48936e 5e395a7a90e0f798566
0	2	10ea371937de976c0622 06f6584d064e795d0b13
0	3	36a0b960d7238803c026 7459517297360e67c1cc
0	4	1917797748647cba5c123 bca2d001bcaead13c25
0	5	0aa42567b1810aa0c764 e8e215773d2b57db153a
1	0	1f9155acbce94b792a42c data3a10172b95cef2e

外部ボリュームエラーチェックコードテーブル

図 5

【 図 6 】

図 6

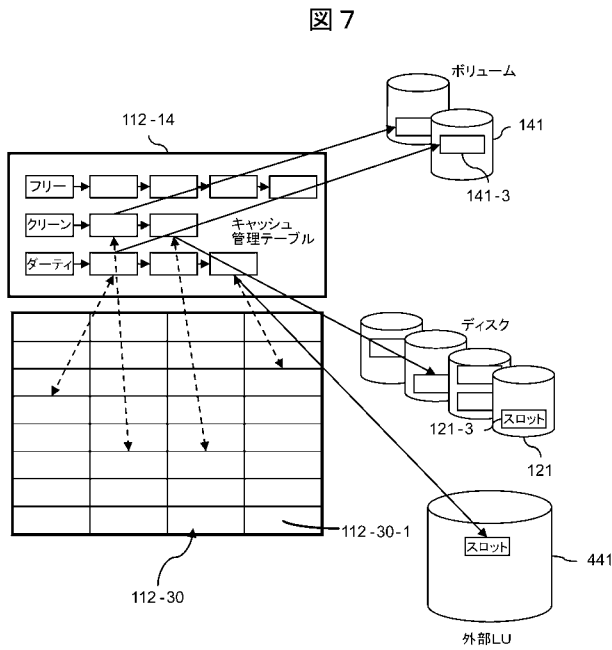
112-14-1	112-14-2	112-14-3	112-14-4
インデックス	ディスク番号	LBA	次
0	2	0xA00	1
1	1	0x7E000	2
2	1	0x9700	3
3	0	0x0000	NULL
4	2	0xC500	5
5	1	0x1100	6
6	1	0xFF00	NULL

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

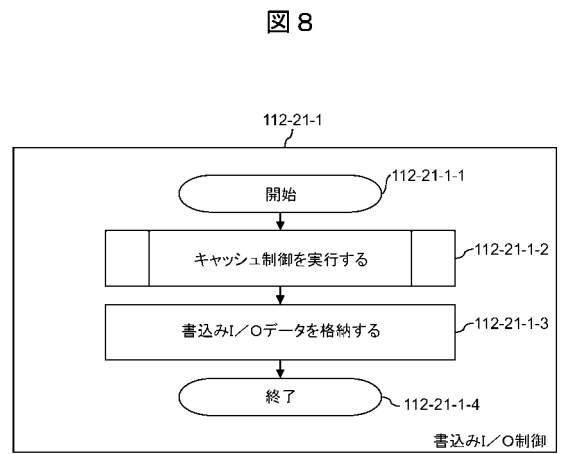
キャッシュ管理テーブル

112-14

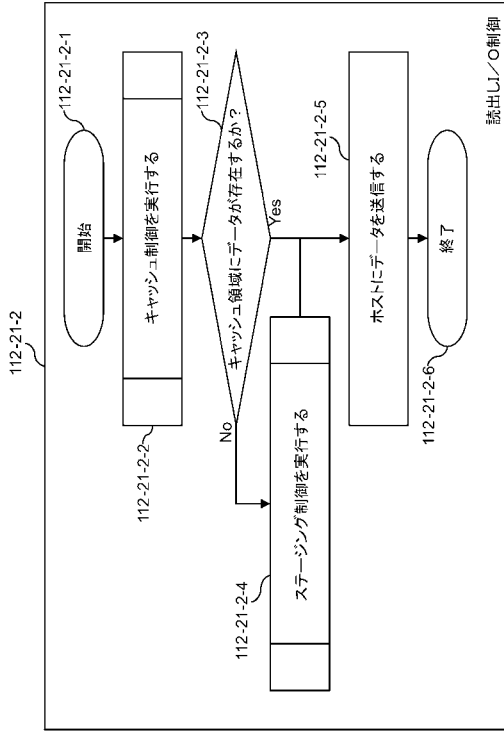
【 図 7 】



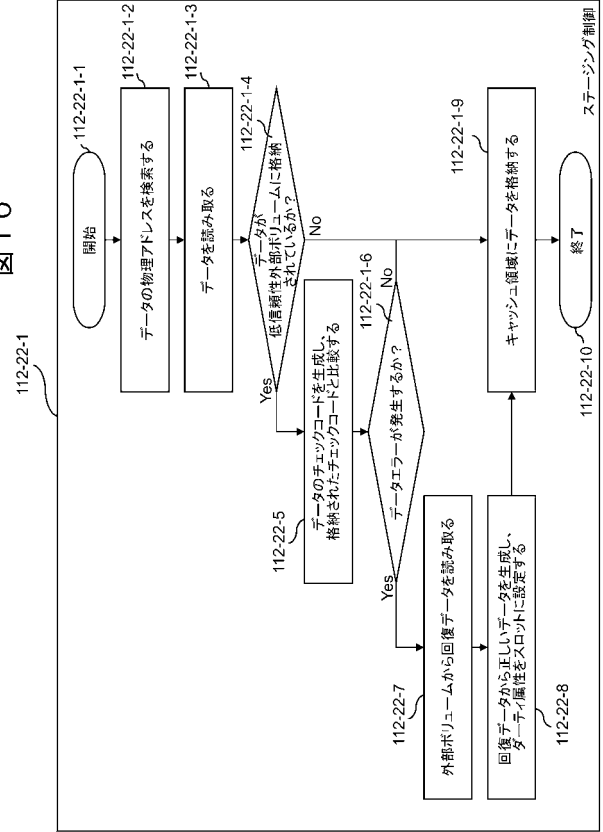
【 図 8 】



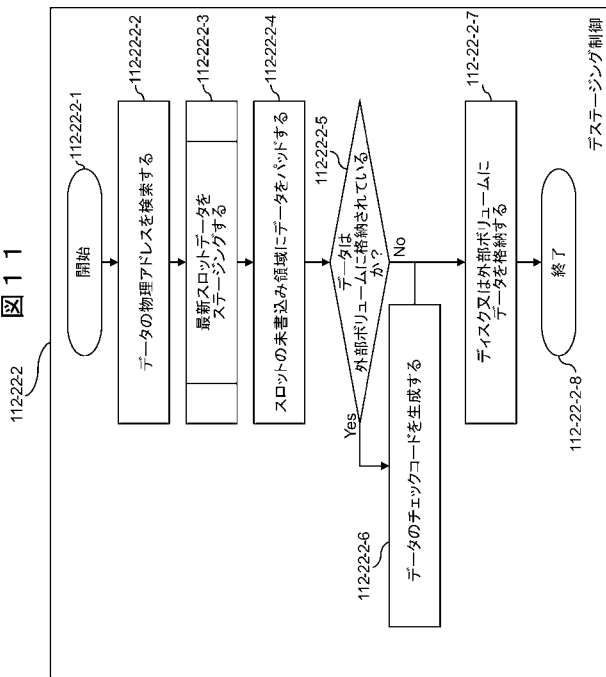
【 図 9 】



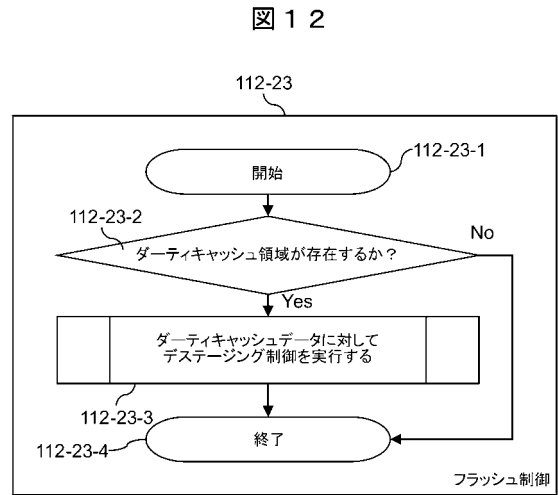
【 図 10 】



【 図 11 】

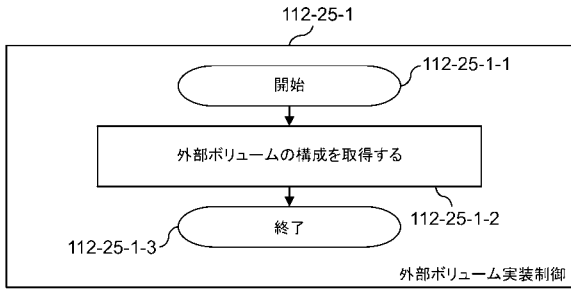


【 図 12 】



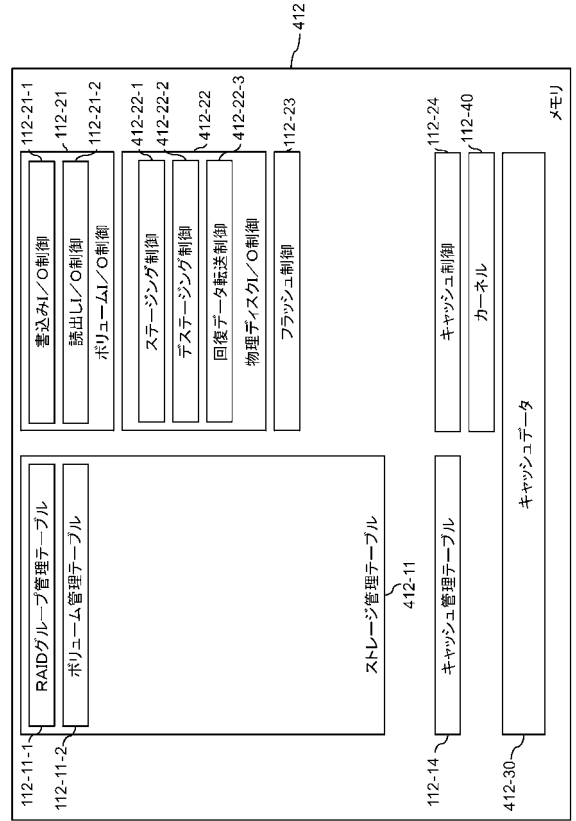
【図 13】

図 13



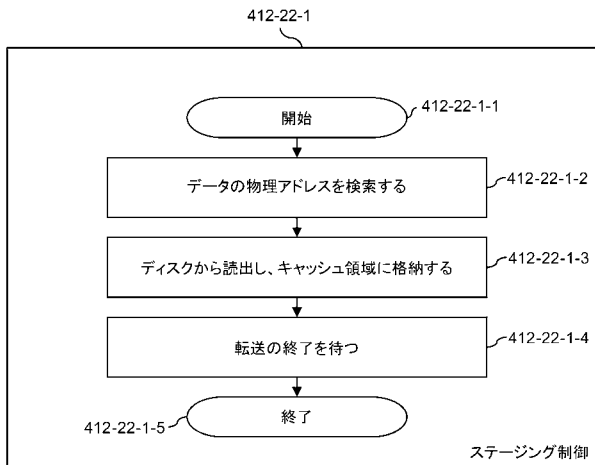
【図 14】

図 14



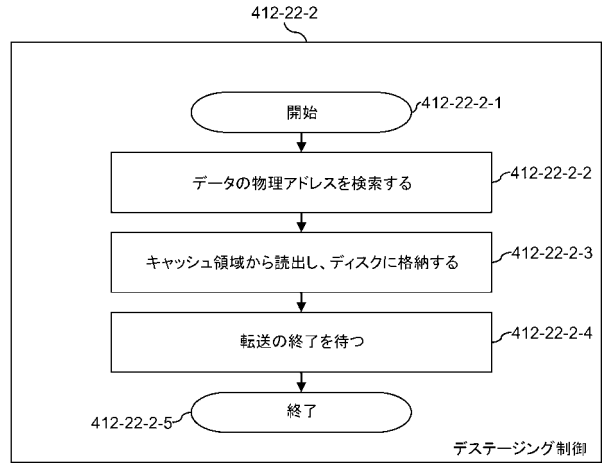
【図 15】

図 15

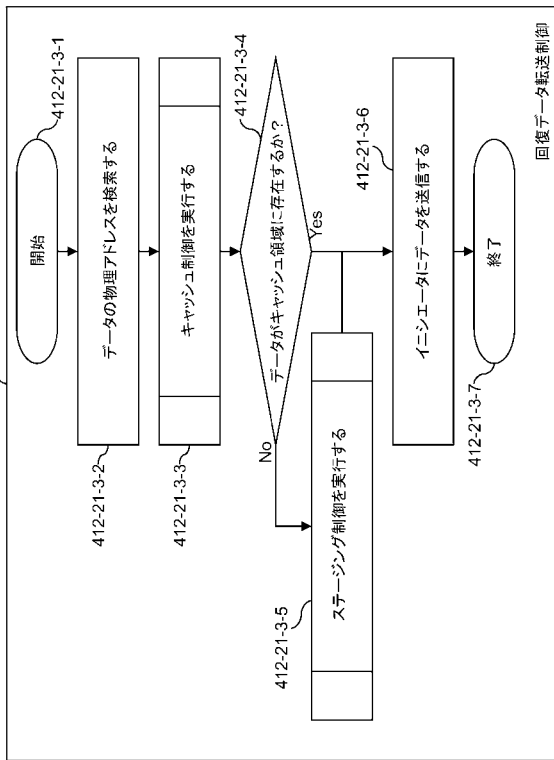


【図 16】

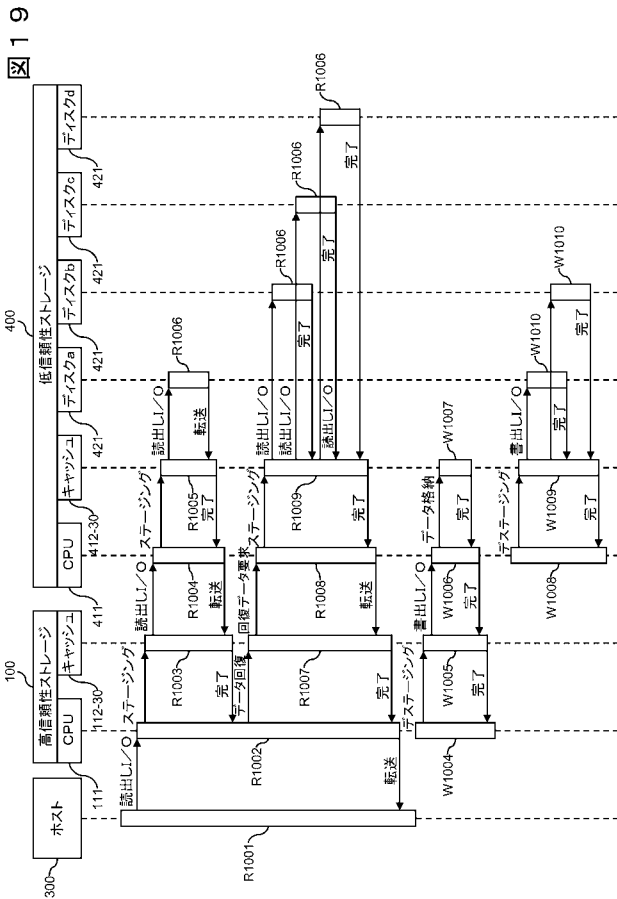
図 16



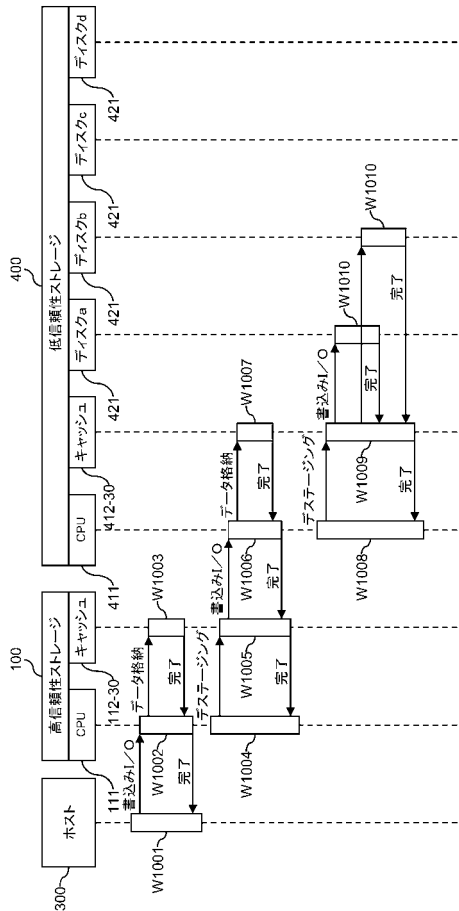
【 図 17 】



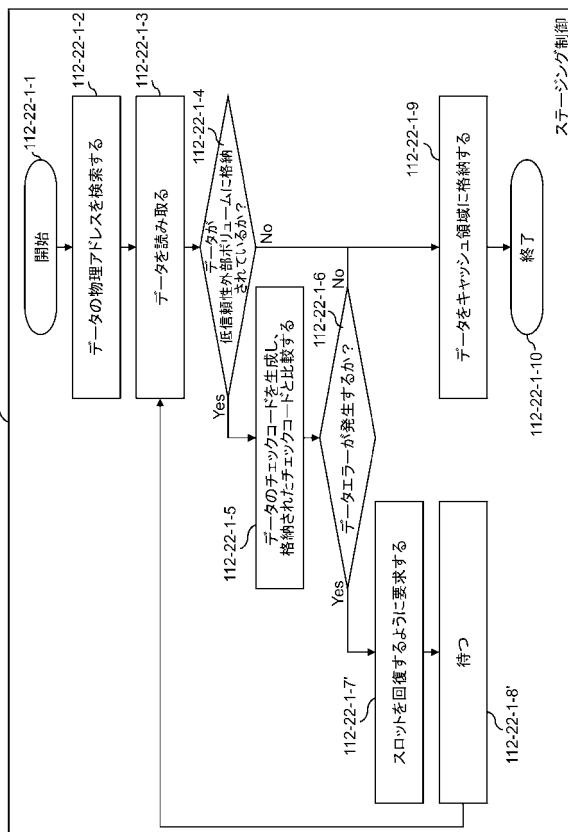
【 図 19 】



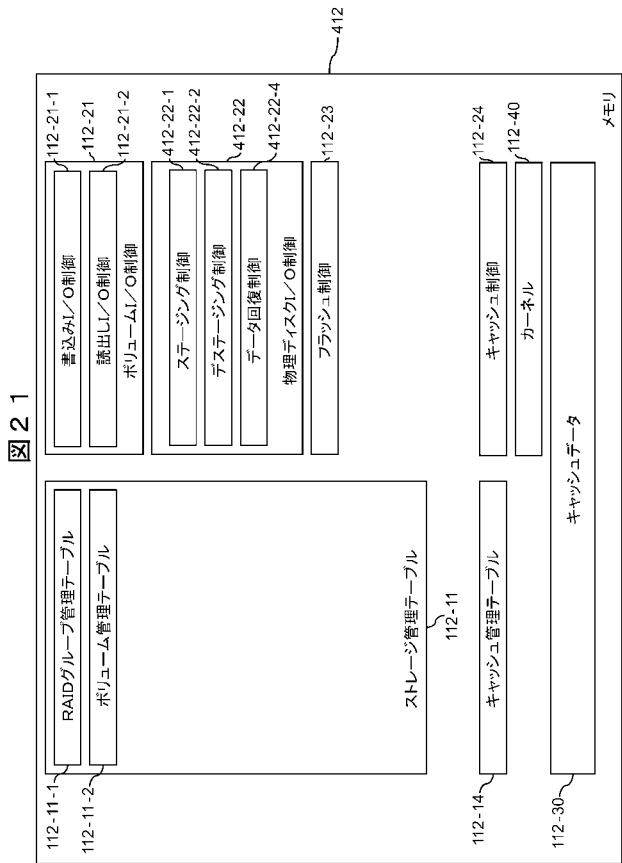
【 図 18 】



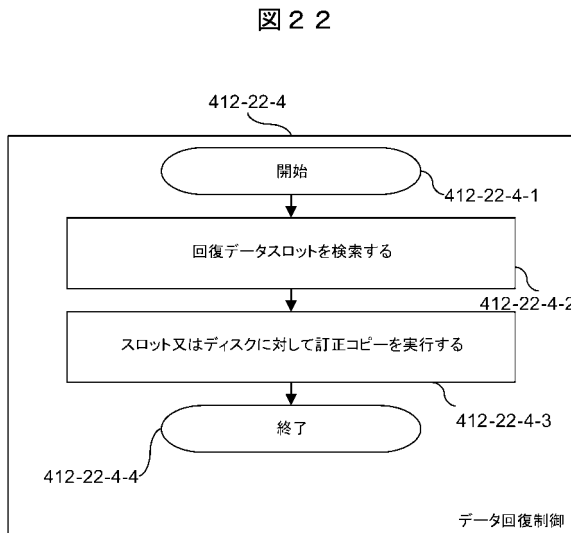
【 図 20 】



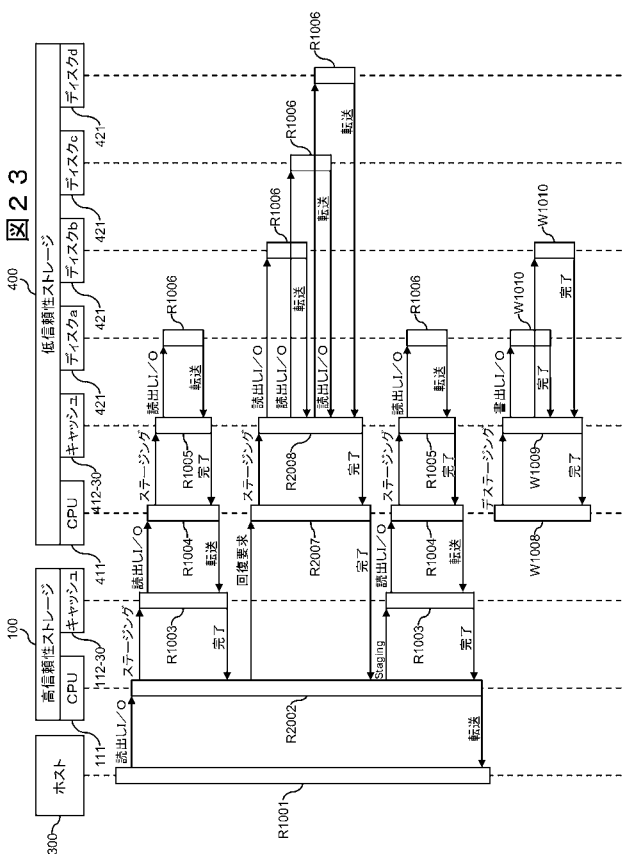
【 図 2 1 】



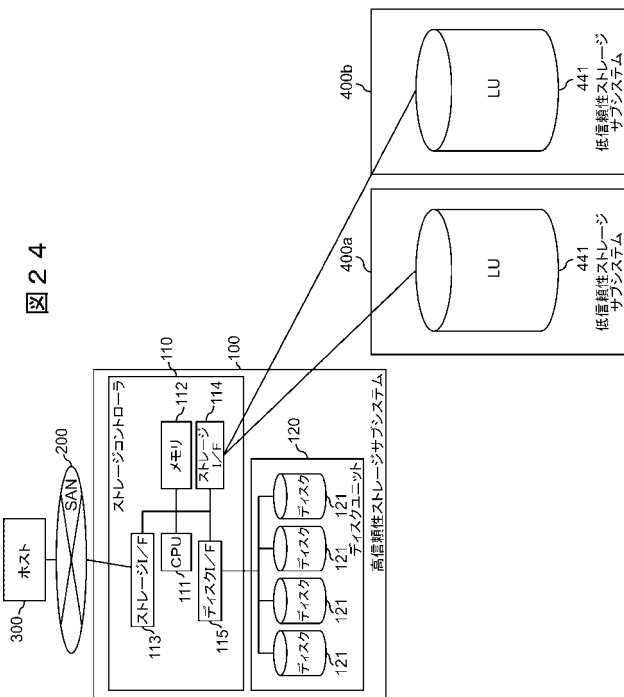
【 図 2 2 】



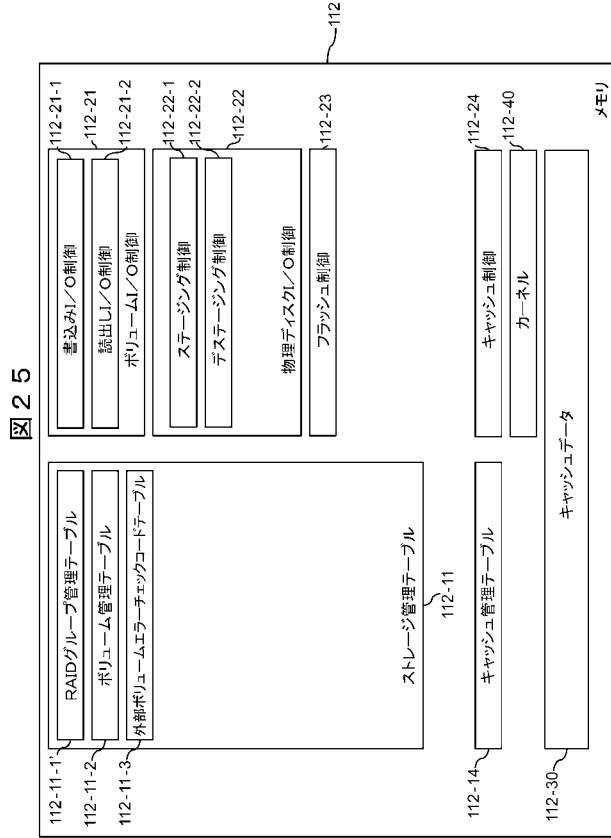
【 図 2 3 】



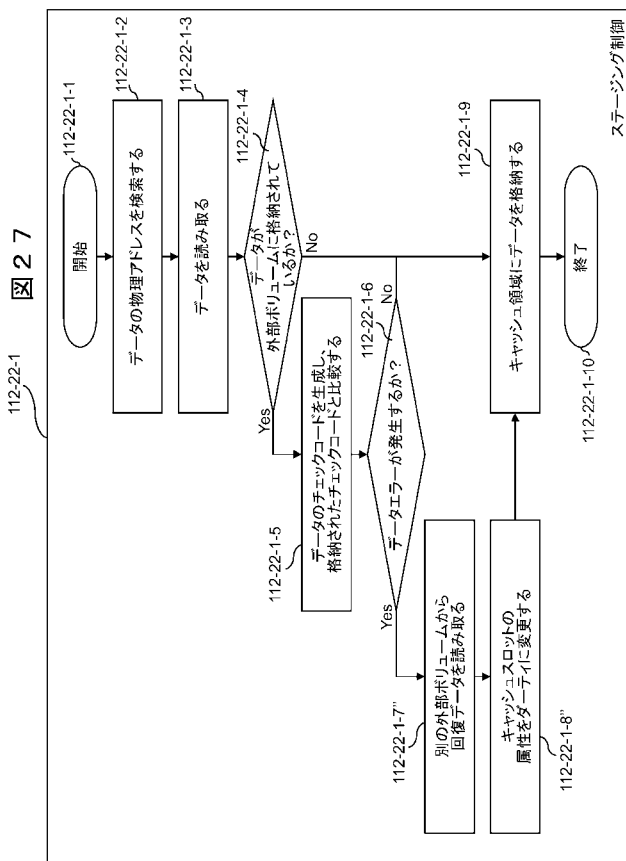
【 図 2 4 】



【 図 2 5 】



【 図 2 7 】

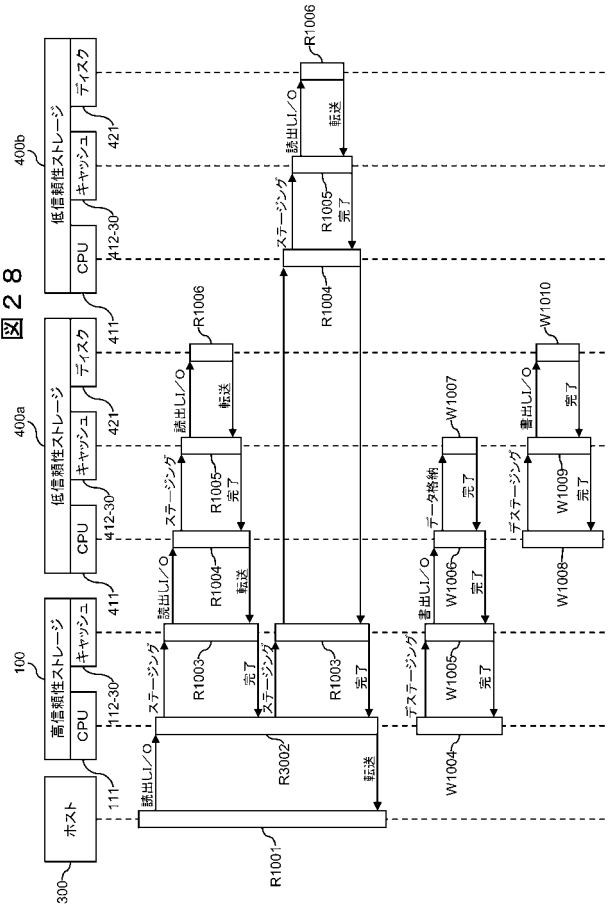


【 図 2 6 】

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 2 8 】



【 図 2 9 】

図 2 9

112-11-4

112-11-4-1

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

【外国語明細書】

2011221981000001.pdf