

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7316242号
(P7316242)

(45)発行日 令和5年7月27日(2023.7.27)

(24)登録日 令和5年7月19日(2023.7.19)

(51)国際特許分類	F I
G 0 6 F 3/06 (2006.01)	G 0 6 F 3/06 3 0 4 E
G 0 6 F 13/10 (2006.01)	G 0 6 F 3/06 3 0 2 A
G 0 6 F 11/16 (2006.01)	G 0 6 F 13/10 3 4 0 A
	G 0 6 F 11/16 6 6 6

請求項の数 11 (全48頁)

(21)出願番号	特願2020-47216(P2020-47216)	(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	令和2年3月18日(2020.3.18)	(74)代理人	110001689 青稜弁理士法人
(65)公開番号	特開2021-149366(P2021-149366 A)	(72)発明者	山本 彰 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(43)公開日	令和3年9月27日(2021.9.27)	(72)発明者	達見 良介 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
審査請求日	令和3年2月8日(2021.2.8)	(72)発明者	大平 良徳 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
		(72)発明者	小川 純司

最終頁に続く

(54)【発明の名称】 ストレージシステムおよびデータ転送方法

(57)【特許請求の範囲】

【請求項1】

キャッシュと、ストレージコントローラと、データを格納するストレージ装置を含む一つ以上のストレージシステムにおいて、

前記ストレージ装置は、記憶媒体を含む一つ以上のストレージボックスを含み、

前記ストレージ装置は、サーバからのライト要求にかかるデータを、前記キャッシュに格納してから前記ストレージボックスに格納する第1のライト処理と、前記キャッシュに格納せずに前記ストレージボックスに格納する第2のライト処理と、を実行可能であり、

前記ストレージコントローラは、前記第2のライト処理を行う場合に、サーバから受信したライト要求に基づき、前記ライト要求にかかるデータを格納するアドレスを設定して前記ストレージボックスに通知し、

前記ストレージボックスは、

サーバから受信したライトデータを受信し、受信したライトデータから冗長データを生成し、

前記ライトデータ及び前記冗長データを、前記ストレージコントローラから通知されたアドレスに従って前記記憶媒体に書き込み、

前記記憶媒体への書き込みの際に異常が発生したとき、前記ライトデータを、前記ストレージコントローラに送信し、

前記ストレージコントローラは、

受信した前記ライトデータを前記キャッシュに書き込むことを特徴とするストレージシ

ステム。

【請求項 2】

請求項 1 に記載のストレージシステムにおいて、
前記ストレージボックスは、
前記ライトデータの中で、前記ストレージコントローラに送ることができないデータが存在する場合、その旨とそのデータを書き込むよう指定された領域のアドレスを、前記ストレージコントローラに伝え、
前記ストレージコントローラは、
前記データを書き込むよう指定された領域のアドレスに対応するビットマップをオンにすることを特徴とするストレージシステム。

10

【請求項 3】

請求項 2 に記載のストレージシステムにおいて、
前記ストレージコントローラは、
前記ビットマップにおいてオンになっているビットに対応する領域に対し、前記サーバからリード要求を受けつけたとき、前記リード要求を、異常終了させることを特徴とするストレージシステム。

【請求項 4】

請求項 1 に記載のストレージシステムにおいて、
前記ストレージボックスは、前記冗長データを生成できない場合に、前記ライトデータを前記ストレージコントローラに送り、
前記ストレージコントローラは、前記ストレージボックスから受信した前記ライトデータを前記キャッシュに、二重に書き込むことを特徴とするストレージシステム。

20

【請求項 5】

請求項 1 に記載のストレージシステムにおいて、
前記ストレージボックスは、
前記ライトデータ及び前記冗長データを、前記記憶媒体の指定された領域に書き込めないとき前記ライトデータを前記ストレージコントローラに送ることを特徴とするストレージシステム。

【請求項 6】

請求項 1 に記載のストレージシステムにおいて、
前記ストレージ装置は、
前記サーバから、書込み領域を指定したライト要求を受けつけ、
前記指定された書込み領域がシーケンシャルになっているかを判別し、
シーケンシャルになっている場合、前記ストレージボックスに、前記サーバからライトデータを受け付けるよう指示することを特徴とするストレージシステム。

30

【請求項 7】

請求項 1 に記載のストレージシステムにおいて、
前記ストレージ装置は、
前記サーバから、書込み領域を指定したライト要求を受け、
前記ライト要求で書き込まれるデータを前記記憶媒体にログストラクチャ形式で書き込むよう制御し、
前記ストレージボックスに、前記サーバからのライト要求で書き込まれるデータを受け取るよう指示することを特徴とするストレージシステム。

40

【請求項 8】

請求項 7 に記載のストレージシステムにおいて、
前記ストレージボックスは、
受け取ったライトデータのハッシュ値を計算し、前記ハッシュ値と対応するライトデータの識別子を、前記ストレージコントローラに送り、
前記ストレージコントローラは、前記ハッシュ値と前記対応するライトデータの識別子

50

を記憶することを特徴とするストレージシステム。

【請求項 9】

請求項 8 に記載のストレージシステムにおいて、

前記記憶媒体への書き込みに際して異常が発生した場合に、

前記ストレージボックスは、

ハッシュ値が一致するものがある場合、指定された領域のデータを読み出し、ライトデータの値と一致するかどうかをチェックし、

ハッシュ値が一致するものがなかったライトデータと、一致するものがあっても読み出した領域のデータが一致しなかったライトデータから、冗長データを生成して、前記ライトデータと前記冗長データとを、前記ストレージコントローラに送信し、

前記ストレージコントローラは、前記ストレージボックスから受信した前記ライトデータと前記冗長データを、前記キャッシュに記憶し、

前記ストレージボックスは、

前記読み出した領域データとライトデータが一致した場合、その旨と前記ライトデータの識別子と読み出した領域を、前記ストレージコントローラに送信し、

前記ストレージコントローラは、

受け取ったライトデータの識別子と読み出した領域を記憶することを特徴とするストレージシステム。

【請求項 10】

請求項 1 に記載のストレージシステムにおいて、

前記ストレージコントローラは、前記第 2 のライト処理を行う場合に、前記ストレージボックスが前記ライトデータを受信した報告を受領したら、前記サーバに完了報告を送信し、

前記完了報告の送信後に、前記ストレージボックスは、前記冗長データの生成と、前記ライトデータ及び前記冗長データの記憶媒体への書き込みと、を行うことを特徴とするストレージシステム。

【請求項 11】

キャッシュと、ストレージコントローラと、データを格納するストレージ装置とを含む一つ以上のストレージシステムにおけるデータ転送方法において、

前記ストレージ装置は、記憶媒体を含む一つ以上のストレージボックスを含み、

前記ストレージ装置は、サーバからのライト要求にかかるデータを、前記キャッシュに格納してから前記ストレージボックスに格納する第 1 のライト処理と、前記キャッシュに格納せずに前記ストレージボックスに格納する第 2 のライト処理と、を実行可能であり、

前記ストレージコントローラは、前記第 2 のライト処理を行う場合に、サーバから受信したライト要求に基づき、前記ライト要求にかかるデータを格納するアドレスを設定して前記ストレージボックスに通知し、

前記ストレージボックスは、

サーバから受信したライトデータから冗長データを生成し、

前記ライトデータ及び前記冗長データを前記、前記ストレージコントローラから通知されたアドレスに従って前記記憶媒体に書き込み、

前記記憶媒体への書き込みに際して異常が発生したとき、前記ライトデータを、前記ストレージコントローラに送信し、

前記ストレージコントローラは、

受信した前記ライトデータを前記キャッシュに格納することを特徴とするデータ転送方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、フラッシュメモリや磁気ディスクをストレージ（記憶媒体）として用いるスケールアウト型のストレージシステムライト処理の高性能化、高信頼化に関する。

10

20

30

40

50

【背景技術】**【0002】**

近年、SSDをはじめ、不揮発性メモリを使用したフラッシュストレージのために最適化された通信プロトコルであるNVMe (Non-Volatile Memory Express) を利用することで、ストレージシステムの大幅な処理高速化が図られている。このようなストレージシステムでは、ストレージコントローラが性能のボトルネックとなることが考えられる。そのため、ストレージシステムに接続されるドライブボックスにホストコンピュータをファブリックネットワーク (NVMeoF:NVMe over Fabric) により直接接続し、ホストコンピュータとドライブボックス間でストレージシステムを介さず直接データ転送をすることによりストレージコントローラの性能ボトルネックを解消し、高速なデータ転送を実現することが考えられる。

10

【0003】

(直接データ転送実現の課題)

直接データ転送を実現するにあたり、以下の2つの課題がある。

【0004】

(課題1) ストレージシステムの提供する論理ボリュームについて、ホストコンピュータから見たアドレス空間と、FBOF内のドライブのアドレス空間が異なり、ホストコンピュータは所望のデータがFBOF内のドライブのどのアドレスに格納されているか識別できない。

【0005】

(課題2) ストレージシステムのキャッシュを使うことでデータアクセスの高性能化を行う場合、キャッシュに新データがある時は、ストレージのキャッシュからデータを読み込む必要があるが、ホストコンピュータは新データの有無を判断できない。

20

【0006】

このような課題に対して、ホストコンピュータで動作するAgentソフトウェアが、ホストコンピュータのアクセス先データに対応するFBOF内のドライブとそのアドレスをストレージコントローラに問い合わせ、得られた情報を元に直接FBOF内のドライブにアクセスする発明が開示されている。

【0007】

特許文献1に開示された発明では、ホストコンピュータが直接FBOFのドライブにアクセスできる一方で、AgentソフトウェアでRAIDなどのデータ保護のための計算をしなくてはならず、ホストコンピュータ側に高信頼な処理を行うための計算負荷がかかる問題がある。

30

【先行技術文献】**【特許文献】****【0008】**

【文献】US9,800,661号明細書

【発明の概要】**【発明が解決しようとする課題】****【0009】**

特許文献1に開示された発明では、ホストコンピュータが直接FBOFのドライブにアクセスできる一方で、AgentソフトウェアでRAIDなどのデータ保護のための計算をしなくてはならず、ホストコンピュータ側に高信頼な処理を行うための計算負荷がかかる問題がある。

40

【0010】

また、サーバからストレージボックスに、データを直接転送するライト処理を実行する場合、以下の三つの課題を解決する必要がある。

【0011】

第1は、従来のライト処理は、ストレージ制御装置が有するキャッシュに、サーバから受け取ったデータを書き込んだ段階で、サーバに対する終了報告を行っていた。本発明に

50

おいても、サーバから見て、従来のライト処理の同等の時間で、ライト処理を同程度の時間で完了させる必要がある。

【0012】

第2は、冗長データの作成が課題となる。従来のストレージ制御装置は、データを受け取った後、受け取ったデータから冗長データを生成して、受け取ったデータ、冗長データを記憶装置に格納した。記憶装置が故障しても、他の記憶装置のデータと冗長データから、故障した記憶装置のデータを復元する機能を有していた。受け取ったデータ、冗長データを記憶装置に格納するまでの間、ストレージ制御装置が、受け取ったデータを喪失しないように、キャッシュにはバッテリなどが装備され、データもキャッシュに2重化して格納していた。本発明では、データをストレージ制御装置が受け取らないので、データの喪失を回避しつつ、冗長データをどのように作成するかが課題となる。

10

【0013】

第3は、障害が発生したときの対応である。サーバに対して、完了を報告したデータについては、ストレージ制御装置は、責任を持って保管する必要があった。このため、障害などが発生しても、データの喪失等を極力抑える必要があるため、ストレージ制御装置は、これらを防ぐため、膨大な論理を有している。同様の論理をストレージボックスで開発しようとする、開発工数が多大となり現実的でない。

【0014】

本発明は、複数のサーバ、複数のストレージシステム、フラッシュメモリや磁気ディスクなどのストレージ装置を格納するストレージボックスをネットワークで共有する構成において、ライト処理を、サーバからストレージボックスに、データを直接転送する際の課題を解決する複合的システムおよびデータ転送方法を提供することを目的とする。

20

【課題を解決するための手段】

【0015】

上記課題を解決する本発明の一側面は、キャッシュとストレージコントローラを含む一つ以上のストレージシステムと、記憶媒体を含む一つ以上のストレージボックスと、を含む複合的システムにおいて、前記ストレージボックスは、サーバから受信したライトデータから冗長データを生成し、前記ライトデータ及び前記冗長データを前記記憶媒体に書き込み、前記ストレージボックスは、前記冗長データを生成できない、または前記ライトデータ及び前記冗長データを前記記憶媒体に書き込みできないとき、前記ライトデータを、前記ストレージシステムに送信し、前記ストレージシステムは、受信した前記ライトデータを前記キャッシュに格納する。

30

【発明の効果】

【0016】

本発明により、複数のサーバ、複数のストレージシステム、フラッシュメモリや磁気ディスクなどのストレージ装置を格納するストレージボックスをネットワークで共有する構成において、ライト処理を実行する際に、ストレージ制御装置の指示により、サーバから送られたデータを、直接ストレージボックスが受け取ることで、ネットワーク負荷の低減を図ることができる。

【図面の簡単な説明】

40

【0017】

【図1】実施例1における情報システムの構成を示す図である。

【図2】実施例1におけるサーバポート情報の構成を示す図である。

【図3】実施例1における実ストレージシステムの構成を示す図である。

【図4】実施例1におけるキャッシュの構成を示す図である。

【図5】実施例1におけるストレージシステムの共有メモリに格納された情報を示す図である。

【図6】実施例1における仮想ストレージシステム情報を示す図である。

【図7】実施例1における他ストレージシステムの情報の形式を示す図である。

【図8】実施例1における仮想論理ボリューム情報の形式を示す図である。

50

- 【図 9】実施例 1 における論理ボリューム情報の形式を示す図である。
- 【図 10】実施例 1 におけるキャッシュ管理情報の形式を示す図である。
- 【図 11】実施例 1 における空きキャッシュ管理情報キューの形式を示す図である。
- 【図 12】実施例 1 におけるストレージボックス情報の形式を示す図である。
- 【図 13】実施例 1 におけるストレージグループ情報の形式を示す図である。
- 【図 14】実施例 1 におけるストレージ装置情報の形式を示す図である。
- 【図 15】実施例 1 におけるストレージコントローラが実行するプログラムの構成を示す図である。
- 【図 16 A】実施例 1 におけるライト要求受付部の処理フローを示す図である。
- 【図 16 B】実施例 1 におけるライト要求受付部の処理フローを示す図である。 10
- 【図 17】実施例 1 におけるライト異常処理対応処理部の処理フローを示す図である。
- 【図 18】実施例 1 におけるリード処理実行部の処理フローを示す図である。
- 【図 19】実施例 1 におけるストレージボックスが実行するプログラムの構成を示す図である。
- 【図 20】実施例 1 におけるライトデータ受領部の処理フローを示す図である。
- 【図 21】実施例 1 におけるライトデータ書込み部の処理フローを示す図である。
- 【図 22】実施例 1 における一時書込み領域転送部の処理フローを示す図である。
- 【図 23】実施例 1 におけるリードデータ直接転送部の処理フローを示す図である。
- 【図 24】実施例 2 における情報システムの構成を示す図である。
- 【図 25】実施例 2 におけるストレージシステムの共有メモリに格納された情報を示す図 20
- 【図 26】実施例 2 における論理ボリューム情報の形式を示す図である。
- 【図 27】実施例 2 におけるストレージグループ情報の形式を示す図である。
- 【図 28】実施例 2 における実ページ情報の形式を示す図である。
- 【図 29】実施例 2 における空きページ管理情報キューの構成を示す図である。
- 【図 30】実施例 2 におけるストレージコントローラが実行するプログラムの構成を示す図である。
- 【図 31 A】実施例 1 におけるライト要求受付部の処理フローを示す図である。
- 【図 31 B】実施例 1 におけるライト要求受付部の処理フローを示す図である。
- 【図 32】実施例 2 におけるライト異常処理対応処理部の処理フローを示す図である。 30
- 【図 33】実施例 2 におけるリード処理実行部の処理フローを示す図である。
- 【図 34】実施例 2 における重複排除処理実行部の処理フローを示す図である。
- 【図 35】実施例 2 におけるストレージボックスが実行するプログラムの構成を示す図である。
- 【図 36】実施例 2 におけるデータ移動部の処理フローを示す図である。
- 【図 37】実施例 2 におけるライトデータ受領部の処理フローを示す図である。
- 【図 38】実施例 2 におけるライトデータ書込み部の処理フローを示す図である。
- 【図 39】実施例 2 における一時書込み領域転送部の処理フローを示す図である。
- 【図 40】実施例 2 におけるリードデータ直接転送部の処理フローを示す図である。
- 【図 41】実施例 2 におけるボックス間データ転送部の処理フローを示す図である。 40
- 【発明を実施するための形態】
- 【実施例 1】
- 【0018】
- < 発明の概要 >
- サーバからストレージボックスに、データを直接転送するライト処理を実行する場合、以下の三つの課題を解決する必要がある。
- 【0019】
- 第 1 は、従来のライト処理は、ストレージ制御装置が有するキャッシュに、サーバから受け取ったデータを書き込んだ段階で、サーバに対する終了報告を行っていた。本発明においても、サーバから見て、従来のライト処理の同等の時間で、ライト処理を同程度の時

間で完了させる必要がある。

【0020】

第2は、冗長データの作成が課題となる。従来のストレージ制御装置は、データを受け取った後、受け取ったデータから冗長データを生成して、受け取ったデータ、冗長データを記憶装置に格納した。記憶装置が故障しても、他の記憶装置のデータと冗長データから、故障した記憶装置のデータを復元する機能を有していた。受け取ったデータ、冗長データを記憶装置に格納するまでの間、ストレージ制御装置が、受け取ったデータを喪失しないように、キャッシュにはバッテリーなどが装備され、データもキャッシュに2重化して格納していた。本発明では、データをストレージ制御装置が受け取らないので、データの喪失を回避しつつ、冗長データをどのように作成するかが課題となる。

10

【0021】

第3は、障害が発生したときの対応である。サーバに対して、完了を報告したデータについては、ストレージ制御装置は、責任を持って保管する必要があるため、障害などが発生しても、データの喪失等を極力抑える必要があるため、ストレージ制御装置は、これらを防ぐため、膨大な論理を有している。同様の論理をストレージボックスで開発しようとする、開発工数が多大となり現実的でない。

【0022】

本発明では、基本的には、ストレージ制御装置とストレージボックスが、機能分担を行い、これらの課題を解決する。

【0023】

第1と第2の課題を解決するための方法を説明する。なお、本発明では、ストレージ制御装置のように、データはストレージ制御装置を通過しないので、ストレージボックスが、冗長データを生成する。冗長データは、主に、データを2重化する方法と、複数のデータから、冗長データを生成し、失われたデータを、他のデータと冗長データから、復元する方法がある。後者を詳細に説明する。後者において、冗長データを生成する方法を説明する。冗長データを生成する、複数のデータの組をパリティグループと呼ぶ。第1の方法は、パリティグループ全体のデータの書込みが行われたとき、全体のデータから冗長データを生成する方法である。もう一つの方法は、比較的少量のデータが書き込まれたときに適用される方法である。本発明は、ネットワークの負荷を低減することを目的としているので、前者を前提とし、ストレージボックスにデータを直接転送するようにする。つまり、このデータの更新前の値と更新後の値と更新前の冗長データとから、更新後のパリティを生成する方法である。

20

30

【0024】

パリティグループ全体のデータを更新するアクセスパターンは、データをアドレス順に更新するシーケンシャルライトである。また、1回のシーケンシャルライトでパリティグループ全体が更新されるわけではない。パリティグループ全体のデータが揃うまで、冗長データは生成できないので、この間、ストレージボックスが受け取ったデータを喪失しないように、本発明では、ストレージボックスがボックス内の記憶装置に、受け取ったデータを一時的に2重に書き込んでおくようにする。ストレージボックスは、記憶装置に2重にデータを書き終えた契機で、ライト処理を完了させる。これにより、第1の課題である、サーバから見たライト処理を、従来と同等の時間で完了させることができる。

40

【0025】

記憶装置にパリティグループ全体のデータが2重に書き込まれると、ストレージボックスが冗長データを生成し、パリティグループ全体のデータと冗長データを、ストレージボックスの中の本来書き込むべき領域に、書き込む。その後、一時的に2重にデータを書き込んでおいた領域を開放する。パリティグループ全体のデータと冗長データを書き込むまで、別の領域に、データを2重書き込まれているので、データの喪失を回避しつつ、冗長データを生成するという第2の課題を解決できる。

【0026】

第3の課題の解決方法を説明する。ストレージ制御装置では、キャッシュにライトデー

50

データを2重化して保持した状態で、障害等が発生した場合の処理を実行していた。このため、本発明では、パリティグループ全体のデータと冗長データを、書き込む前に、何か障害が発生すると、ストレージボックスは、それまで受け取っていたデータをストレージ制御装置に転送し、ストレージ制御装置は、このデータを、キャッシュに2重に格納する。この状態にすると、従来ストレージ制御装置がもっていた障害に対応する論理をそのまま適用することが可能となるので、第3の課題を解決できる。

【0027】

図1は、実施例1における情報システムの構成を示す。情報システムは、一つ以上のストレージボックス130、一つ以上の実ストレージシステム100、一つ以上のサーバ110と、実ストレージシステム100、サーバ110、ストレージボックス130とを接続するネットワーク120と、ストレージ管理サーバ150とを有する。

10

【0028】

サーバ110は、一つ以上のサーバポート195によって、実ストレージシステム100は、一つ以上のストレージポート(ストレージパス)197によって、ストレージボックス130は、ボックスポート180によって、ネットワーク120に接続される。

【0029】

サーバ110は、内部に、サーバポート情報198をもつ。サーバ110は、ユーザアプリケーションが動作するシステムで、実ストレージシステム100に対して、ネットワーク120経由で、データの読み書き要求を発行する。また、実ストレージシステム100は、ネットワーク120経由で、サーバ110から受け取った読み書き要求で指定されたデータを格納したストレージボックス130に要求を伝える。

20

【0030】

実ストレージシステム100同士も、互いに、ネットワーク120経由で、データの送受信を行う。ネットワーク120は例えば、NVM-e over Ethernet等のプロトコルを用いることができる。実施例1では、一つ以上の実ストレージシステム100から構成される仮想ストレージシステム190が存在する。サーバ110からは、仮想ストレージシステム190が、一台のストレージシステムに見える。ただし、本実施の形態は、仮想ストレージシステム190が存在しない場合も有効である。仮想ストレージシステム190が存在しない場合、実ストレージシステム100が、サーバ110からは、ストレージシステムに見える。

30

【0031】

ストレージボックス130は、HDD(Hard Disk Drive)、フラッシュメモリを記憶媒体とするフラッシュストレージなどのストレージ装置160(共有ストレージ装置175)とネットワーク120に接続するための一つ以上のボックスポート180、実ストレージシステム100の指示にしたがって、ストレージ装置160との間で転送を行うひとつ以上のボックスプロセッサ170、ボックスプロセッサ170が使用するボックスメモリ181などが含まれる。また、一つのストレージ装置160は、一つ以上のストレージパス197に接続されているものとする。

【0032】

また、フラッシュストレージも、いくつかの種類があり、高価格、高性能、消去可能回数の多いSLCと、これに対し、低価格、低性能、消去可能回数の少ないMLCとがある。さらに、相変化メモリなどの新しい記憶媒体が含まれていてもよい。

40

【0033】

実施例1においては、ストレージボックス130の中のストレージ装置160は、複数の実ストレージシステム100によって共有される。したがって、実施例1では、ストレージボックス130の中のストレージ装置160を、共有ストレージ装置175と呼ぶ。

【0034】

実ストレージシステム100とストレージボックス130は、ネットワーク120経由で、仮想ストレージシステム190の中の一つ以上の実ストレージシステム100に接続される。ストレージボックス130は、必ずしも、仮想ストレージシステム190の中の

50

すべての実ストレージシステム 100 と接続されている必要はない。同様に、ある実ストレージシステム 100 と別の実ストレージシステム 100 に接続されているストレージボックス 130 の集合は、まったく同じである必要はない。ストレージボックス 130 は、ネットワーク 120 経由で、実ストレージシステム 100 のいくつかによって、共有されている。

【0035】

ストレージ管理サーバ 150 は、ネットワーク 120 などを経由して、実ストレージシステム 100 と接続されていて、実ストレージシステム 100、ストレージボックス 130 の管理を行うため、ストレージ管理者が利用する装置である。また、実施例 1 では、実ストレージシステム 100 は、容量仮想化機能をもたないものとする。ただし、実施例 1 は、実ストレージシステム 100 が、容量仮想化機能をもつ場合でも、有効である。

10

【0036】

図 2 は、サーバポート情報 198 のフォーマットである。サーバポート情報 198 は、サーバポート 195 ごとに持つ情報である。サーバ 110 は、実ストレージシステム 100 に読み書き要求を発行する際、ストレージシステムの識別子、論理ボリュームの識別子、ストレージポート（ストレージパス 197）の識別子を設定する。このため、サーバポート情報 198 は、サーバポート識別子 24000 と、当該サーバポート 195 からアクセスする一つ以上の論理ボリューム識別子 24001、この論理ボリュームを含むストレージシステム識別子 24002、ストレージポート識別子 24003 が含まれる。当該論理ボリュームが複数のストレージパス 197 に接続されている場合、複数のストレージポート識別子（ストレージパス識別子）24003 が設定される。なお、実施例 1 では、ストレージシステム識別子 24002 には、仮想ストレージシステム 190 の識別子が設定される。ただし、本発明は、仮想ストレージシステム 190 が存在しない場合も有効で、その場合には、ストレージシステム識別子 24002 には、実ストレージシステム 100 の識別子が設定される。ストレージポート識別子 24003 には、ストレージパス 197 の識別子が設定される。

20

【0037】

実施例 1 では、論理ボリューム識別子 24001 には、仮想論理ボリュームの識別子が設定される。仮想論理ボリュームの識別子は、仮想ストレージシステム 190 内で、ユニークな値である。一方、それぞれの実ストレージシステム 100 も論理ボリュームをもつ。論理ボリュームの識別子は、実ストレージシステム 100 内では、ユニークである。サーバ 110 の読み書き要求は、仮想ストレージシステムの識別子、仮想論理ボリュームの識別子、ストレージパス 197 の識別子を含む。ストレージパス 197 の識別子は、リアルな値であるため、この要求を受け取る、実ストレージシステム 100 が決定される。なお、本発明は、仮想ストレージシステム 190 が存在しない場合も有効で、その場合には、論理ボリューム識別子 24001 には、実ストレージシステム 100 の論理ボリュームの識別子が設定される。また、実施例 1 では、仮想論理ボリュームを読み書きする実ストレージシステム 100 を変更するが、その際にも、仮想論理ボリュームの識別子は、変化せず、接続されたストレージパス 197 を変更する。

30

【0038】

図 3 は、実ストレージシステム 100 の構成である。実ストレージシステム 100 は、一つ以上のストレージコントローラ 200、キャッシュ 210、共有メモリ 220、ストレージ装置 160（内部ストレージ装置 230）と、これらの構成要素を接続する一つ以上の接続装置 250、ネットワーク 120 とのインターフェイスであるストレージポート 197 とを有する。

40

【0039】

実ストレージシステム 100 内に含まれる内部ストレージ装置 230 も、ストレージボックス 130 内に含まれるストレージ装置 160 も、データを格納する装置である。特に、実施例 1 では、実ストレージシステム 100 に含まれるストレージ装置 160 を内部ストレージ装置 230、ストレージボックス 130 内のストレージ装置 160 を共有ストレ

50

ージ175と呼ぶ。ただし、本発明の対象は、ストレージボックス130内のストレージ装置であるため、基本的には、ストレージ装置160は、共有ストレージ175を指すものとする。また、本発明においては、ストレージシステム100は、必ずしも、内部ストレージ装置230をもつ必要はない。また、本発明は、複数の実ストレージシステム100によって共有される一つ以上のストレージボックス130によって構成される複合ストレージシステムに関するものなので、内部ストレージ装置230に対し、ストレージコントローラ200が実行する処理、共有メモリ220に保有する情報などについては、説明を省略する。

【0040】

ストレージコントローラ200は、サーバ110から発行されたリード・ライト要求を処理するプロセッサ260、プログラムや情報を保管するメモリ270、バッファ275と、から構成される。特に、本発明は、バッファ275は、(1)後述する冗長データを生成する際、生成に必要な情報や、生成した冗長データを格納するため、(2)実ストレージシステム100のキャッシュ領域に格納されたデータを、恒久的に格納するストレージ装置160に書き込む際の一時的な格納領域とし使用される。

10

【0041】

接続装置250は、実ストレージシステム100内の各構成要素を接続する機構である。本発明の特徴は、接続装置250経由で、一つ以上のストレージボックス130が接続されている点である。これにより、ストレージコントローラ200は、ストレージボックス130内のストレージ装置160に、読み書きできる。また、実施例1では、ストレージボックス130には、実ストレージシステム100内の一つ以上のストレージコントローラ200に接続されているものとする。

20

【0042】

キャッシュ210、共有メモリ220は、通常DRAMなどの揮発メモリで構成されるが、バッテリーなどにより不揮発化されているものとする。また、実施例1では、高信頼化のため、それぞれが2重化されているものとする。ただ、本発明は、キャッシュ210、共有メモリ220が不揮発化されていなくても、2重化されていなくとも、有効である。また、実施例1では、キャッシュ210、共有メモリ220は、ストレージコントローラ200の外側にあるが、ストレージコントローラ200の内部のメモリ270上に、構成されていても良い。

30

【0043】

キャッシュ210には、内部ストレージ装置230、ストレージボックス130の中のストレージ装置160に格納されたデータの中で、ストレージコントローラ200からよくアクセスされるデータが格納される。

【0044】

ストレージコントローラ200は、サーバ110から、リード・ライト要求を受け付ける。通常は、ストレージ装置160に書き込むよう受け取ったデータを、キャッシュ210に書き込んで、該当するライト要求を完了させた後から、このデータを、ストレージボックス130のストレージ装置160に書き込むことが多い。しかし、この方法では、データは2回、ネットワーク120を通過することになるので、ネットワーク120に対する負荷が大きい。

40

【0045】

本発明では、書き込むデータを、サーバ110から直接、ストレージボックス130に転送することにより、ネットワーク120の負荷を軽減する。

【0046】

図4は、キャッシュ210の構成を表しているものとする。キャッシュ210は、固定長のスロット21100に分割されている。スロット21100が、リード・ライトデータの割り当て単位となる。

【0047】

なお、実施例1では、実ストレージシステム100のストレージ装置160の中の一

50

の装置が故障しても、その装置のデータを回復できるRedundancy Array Independent Device (RAID) 機能をもっているものとする。RAID機能もった場合、複数の同一種類の記憶装置が、一つのRAID構成をとる。これをストレージグループ280(図3参照)と呼ぶ。

【0048】

実施例1では、RAID構成は、一つのストレージボックス130内の共有ストレージ175の集合、あるいは、一つの実ストレージシステム100内の内部ストレージ装置230の集合から構成されるものとする。RAIDには、いくつかのタイプがある。例えば、RAID1は、データを2重に書く。RAID4やRAID5は、複数のデータから、パリティ(以下、パリティ)を生成する。RAIDのタイプは、ストレージ装置グループごとに決まっている。なお、実施例1では、RAIDのタイプは、RAID5とする。もちろん、本発明は、RAID5以外のRAIDタイプでも有効である。

10

【0049】

本発明では、ストレージボックス130がサーバ110からデータを直接受け取り、データはストレージ制御装置を通過しないので、ストレージボックス130が、パリティを生成する。

【0050】

RAID5において、パリティを生成する方法を説明する。パリティを生成する、複数のデータの組をパリティグループと呼ぶ。

【0051】

第1の方法は、パリティグループ全体のデータの書込みが行われたとき、全体のデータから冗長データを生成する方法である。もう一つの方法は、比較的少量のデータが書き込まれたときに適用される方法である。本発明は、ネットワークの負荷を低減することを目的としているので、前者を前提とし、ストレージボックスにデータを直接転送するようにする。

20

【0052】

パリティグループ全体のデータを更新するアクセスパターンは、データをアドレス順に更新するシーケンシャルライトである。また、1回のシーケンシャルライトでパリティグループ全体が更新されるわけではない。パリティグループ全体のデータが揃うまで、パリティは生成できないので、この間、ストレージボックス130が受け取ったデータを喪失しないように、実施例1では、ストレージボックス130がボックス内のストレージ装置160に、受け取ったデータを一時的に2重に書き込んでおくようにする。ただし、本発明は、ストレージボックス130が、ストレージコントローラ200と同様に、ストレージ装置160以外に、不揮発化されたキャッシュを持ち、このキャッシュに、受け取ったデータを2重に書き込んで有効である。キャッシュの記憶媒体は、バッテリーバックアップされたDRAMでも、新しい不揮発の半導体メモリでも、本発明は有効である。ストレージボックス130は、ストレージ装置160に2重にデータを書き終えた契機で、ライト処理を完了させる。

30

【0053】

ストレージ装置160にパリティグループ全体のデータが2重に書きこまれると、ストレージボックス130がパリティを生成し、パリティグループ全体のデータとパリティを、ストレージボックス130の中の本来書き込むべき領域(サーバが書込みを指示した領域)に、書き込む。この後、一時的に2重にデータを書き込んでおいた領域を開放する。

40

【0054】

パリティグループ全体のデータとパリティを書き込むまで、データは、別の領域に2重書き込まれているので、データの喪失を回避しつつ、パリティを生成することができる。ただし、2重に書き込む一方を、本来書き込むべき領域(サーバが書込みを指示した領域)に書き込んでおき、パリティグループ全体のデータが2重に書きこまれると、ストレージボックス130がパリティを生成し、パリティを書き込むようにしても、本発明は有効である。

50

【 0 0 5 5 】

ストレージコントローラ 200 では、キャッシュ 210 にライトデータを 2 重化して保持した状態で、障害等が発生した場合の処理を実行していた。このため、本発明では、パリティグループ全体のデータとパリティを、書き込む前に、何か障害が発生すると、ストレージボックス 130 は、それまで受け取っていたデータをストレージコントローラ 200 に転送し、ストレージコントローラ 200 は、このデータを、キャッシュ 210 に 2 重に格納する。この状態にすると、ストレージコントローラ 200 が従来からもっていた障害に対応する論理をそのまま適用することが可能となる。

【 0 0 5 6 】

図 5 は、実施例 1 における実ストレージシステム 100 の共有メモリ 220 の中の実施例 1 に関する情報を示している。共有メモリ 220 には、ストレージシステム情報 2060、他ストレージシステム情報 2070、仮想論理ボリューム情報 2085、論理ボリューム情報 2000、ストレージボックス情報 2050、ストレージグループ情報 2300、ストレージ装置情報 2500、キャッシュ管理情報 2750、空きキャッシュ管理情報ポインタ（キュー）2650、によって構成される。

10

【 0 0 5 7 】

この中で、ストレージシステム情報 2060 は、図 6 に示すように、実ストレージシステム 100 に関する情報で、実施例 1 では、仮想ストレージシステム識別子 2061、実ストレージシステム識別子 2062 から構成される。仮想ストレージシステム識別子 2061 は、当該実ストレージシステム 100 が含まれる仮想ストレージシステム 190 の識別子である。実ストレージシステム識別子 2062 は、当該実ストレージシステム 100 の識別子である。

20

【 0 0 5 8 】

図 7 は、他ストレージシステム情報 2070 に関する情報を示した図であり、仮想ストレージシステム識別子 2071、他実ストレージシステム識別子 2072 から構成される。仮想ストレージシステム識別子 2071 は、図 6 に含まれる仮想ストレージシステム識別子 2051 と同じで、当該実ストレージシステム 100 が含まれる仮想ストレージシステム 190 の識別子である。他実ストレージシステム識別子 2072 は、当該実ストレージシステム 100 を含む仮想ストレージシステム 190 に含まれる他の実ストレージシステムである。

30

【 0 0 5 9 】

図 8 は、仮想論理ボリューム情報 2085 に関する情報を示した図である。論理ボリューム対応に作成される仮想論理ボリューム情報 2085 は、仮想論理ボリューム識別子 2086、制御権情報 2087、実ストレージシステム情報 2088、ストレージポート識別子（ストレージパス識別子）2089、サーバポート識別子 2091、論理ボリューム識別子 2090 等から構成される。

【 0 0 6 0 】

仮想論理ボリューム識別子 2086 は、当該仮想論理ボリュームの識別子である。実施例 1 では、仮想論理ボリュームは、どれか一つの実ストレージシステム 100 が読み書きを行う権限をもつ。制御権情報 2087 は、当該実ストレージシステム 100 が制御権をもっているか、そうでないかを表す。制御権をもっていない場合、実ストレージシステム情報 2088、ストレージパス識別子 2089 に、どの実ストレージシステム 100 が制御権をもっているか、その仮想論理ボリュームが接続されている一つ以上のストレージパス 197 が示される。

40

【 0 0 6 1 】

論理ボリューム識別子 2090 は、当該実ストレージシステム 100 内の対応する論理ボリュームの識別子が格納され、そうでない場合、制御権をもつ実ストレージシステム 100 内の対応する論理ボリュームの識別子が格納される。

【 0 0 6 2 】

図 9 は、本発明に関する論理ボリューム情報 2000 に関する情報を示した図である。

50

実施例1では、サーバ110がデータをリード・ライトするストレージ装置160は、仮想論理ボリュームである。また、サーバ110は、仮想論理ボリュームのID、仮想論理ボリューム内のアドレス、読み書きしたいデータの長さを指定して、リード要求やライト要求を発行する。実ストレージシステム100は、サーバ110から読み書き要求を受け取ると、仮想論理ボリューム情報2085により、対応する論理ボリュームの識別子を認識する。実施例1では、論理ボリュームの識別子は、実ストレージシステム100内でユニークな情報である。論理ボリューム情報2000は、論理ボリュームごとに存在する情報である。この情報は、論理ボリューム識別子2001、論理容量2002、論理ボリュームタイプ2005、論理ボリュームRAIDタイプ2003、割り当て範囲2026、第1キャッシュ管理情報ポインタ2022、第2キャッシュ管理情報ポインタ2023、第1ストレージポインタ2024と第2ストレージポインタ2025等より構成される。

10

【0063】

論理ボリューム識別子2001は、対応する論理ボリュームのIDを示す。論理容量2002は、この論理ボリュームの容量である。論理ボリュームタイプ2005は、論理ボリュームのタイプを表す。実施例1では、当該論理ボリュームが、内部ストレージ装置230に格納されているか、共有ストレージ装置175に格納されているかを示す。論理ボリュームRAIDタイプ2003は、該当する論理ボリュームのRAIDタイプ、RAID0、RAID1などを指定する。RAID5のように、N台の容量に対し、1台の容量の冗長データを格納する場合、Nの具体的な数値を指定するものとする。ただし、任意のRAIDタイプが指定できるわけではなく、少なくとも一つストレージグループがもつRAIDタイプである必要がある。割り当て範囲2026は、当該論理ボリュームに割り当てたストレージグループの識別子と、最も小さい番号のセグメントの番号を示している。セグメントビットマップは、最も小さい番号のセグメントを最初1ビットとして、以降の各セグメントに、当該論理ボリュームに割り当てられたか、そうでないかを表す。

20

【0064】

第1キャッシュ管理情報ポインタ2022は、当該論理ボリュームを、スロット21100に相当する容量で分割したそれぞれの領域が、スロット21100が割り当てられている(キャッシュ210に格納されている)かを表す。割り当てられている場合、対応するキャッシュ管理情報2750がポイントされる。割り当てられていない場合、ヌル状態になる。実施例1では、ストレージ装置160に未書込みデータは、高信頼化のために、2重化してキャッシュ210に格納する。このため、第2キャッシュ管理情報ポインタ2023は、未書込みデータを2重化して格納するキャッシュ管理情報2750へのポインタである。

30

【0065】

第1ストレージポインタ2024と第2ストレージポインタ2025は、当該論理ボリュームを、パリティグループに相当する容量で分割したそれぞれの領域が(ストレージボックス130に、ライトデータを直接転送している場合)、一時的にストレージ装置160に、二重に、データを書き込むために、確保した領域のアドレスを示す。アドレスには、ストレージボックス130の識別子、ストレージ装置160の識別子、ストレージ装置160のアドレスが示される。

40

【0066】

図10は、キャッシュ管理情報2750の構成である。キャッシュ管理情報2750は、スロット21100対応に存在する情報である。次キャッシュ管理情報ポインタ2751、割り当て論理ボリュームアドレス2752、ブロックビットマップ2753、更新ビットマップ2754、アクセス不可ビットマップ2755、最終アクセスアドレスフラグ2756、直接転送フラグ2757、直接転送中断フラグ2758、旧データ破壊フラグ2759から構成される。

【0067】

次キャッシュ管理情報ポインタ2751は、データを格納していない状態のあるスロットに対応するキャッシュ管理情報2750をポイントし、当該スロットは次のデータを格

50

納する。

【0068】

割り当て論理ボリュームアドレス2752は、対応するスロットにデータを格納したとき、どの論理ボリュームのどのアドレスから割当を開始する領域のデータを格納したかを示す。ブロックビットマップ2753は、この領域の中で、キャッシュ210、あるいは、ストレージ装置160のバッファに格納されたブロック（読み書きの最小単位）を示す。格納時にはONとなる。

【0069】

更新ビットマップ2754は、サーバ110から書き込み要求を受け、まだ、ストレージ装置160に書き込んでいないブロックを示す。未書きこみは、ONとなる。

10

【0070】

アクセス不可ビットマップ2755は、障害などが発生していて、リード要求がきた場合、エラーを返すブロックを示す。ライト要求がきた場合、データを受付、正常終了すれば、読んでもよいデータとなるため、アクセス不可ビットマップ2755の対応するビットをオフにする。実施例1では、シーケンシャルライトに対して、データを直接、ストレージボックス130に転送するようにする。シーケンシャルライトは、アドレス順に、データが書き込まれる。

【0071】

最終アクセスアドレスフラグ2756は、前回の要求で、アクセスされたアドレスを示しており、これにより、アクセスがアドレス順になっているかをチェックできる。また、ストレージボックス130に直接データの送信、あるいは送信の中止することができる。ビットマップ情報や最終アクセスアドレスフラグ2756については、直接データ転送を行わない従来のキャッシュに書き込む処理を実行していても、正しい値が格納されているものとする。

20

【0072】

直接転送フラグ2757は、現在、このスロットに対応する領域は、サーバ110から、ストレージボックス130に直接データ転送が行われていることを示す。直接転送中止フラグ2758は、なんらかの理由により、直接転送が中止されたことを示すフラグである。旧データ破壊フラグ2759は、ストレージボックス130がパリティグループ全体のデータを受領後、パリティを生成し、パリティグループ全体のデータとパリティをストレージ装置160に書き込んでいる際に、障害が発生した場合、すでに、旧データを破壊してしまっている可能性がある際、当該パリティグループのすべてのキャッシュ管理情報2750の当該フラグをオンにする。

30

【0073】

図11は、空きキャッシュ管理情報ポインタ2650によって管理されるデータを格納されていないスロットに対応するキャッシュ管理情報2750の管理方式である。空きキャッシュ管理情報ポインタ2650が、データが格納されていないスロットに対応する先頭のキャッシュ管理情報2750をポイントし、次のキャッシュ管理情報2750を、次ぎキャッシュ管理情報ポインタ2751がポイントする。

【0074】

図12は、ストレージボックス情報2050の構成である。ストレージボックス情報2050は、ストレージボックス130ごとに設ける情報である。

40

【0075】

ストレージボックス情報2050は、ストレージボックス識別子7000、接続情報7001、ストレージ台数7002、接続ストレージ台数7003、ボックスポート数7004、ボックスポート識別子7005、ボックスプロセッサ数7006、ボックスプロセッサ識別子7007、パス数7008、パス識別子7009、ストレージグループ数7010、ストレージグループ識別子7011から構成される。

【0076】

ストレージボックス識別子7000は、対応するストレージボックス130の識別子で

50

ある。接続情報 7001 は、対応するストレージボックス 130 が、それぞれの実ストレージシステム 100 に接続されているか、いないかを示す情報である。

【0077】

ストレージ台数 7002 は、当該ストレージボックス 130 に接続可能なストレージ装置 160 の台数である。接続ストレージ台数 7003 は、実際に接続されているストレージ装置 160 の台数である。ボックスポート数 7004 は、ネットワーク 120 に接続された当該ストレージボックス 130 のボックスポート 180 の数で、ボックスポート識別子 7005 は、それぞれのボックスポート 180 ごとの識別子である。

【0078】

ボックスプロセッサ数 7006 は、当該ストレージボックス 130 内のボックスプロセッサ 170 の数である。ボックスプロセッサ識別子 7007 は、それぞれのボックスプロセッサ 170 との識別子である。パス数 7008 は、当該、ストレージボックス 130 内のストレージ装置 160 とボックスプロセッサ 170 の間のパスの数である。パス識別子 7009 は、それぞれのストレージパス 197 ごとの識別子である。ストレージグループ数 7010 は、当該ストレージボックス 130 に含まれるストレージグループの数である。ストレージグループ識別子 7011 は、当該ストレージボックス 130 に含まれるストレージグループの識別子である。実施例 1 では、ストレージグループを構成するストレージ装置 160 は、一つのストレージボックス 130 に含まれるとする。ただし、本発明は、ストレージグループが複数のストレージボックス 130 のストレージ装置 160 で構成されても有効である。

【0079】

図 13 は、ストレージグループ情報 2300 の形式を示す。ストレージグループ情報 2300 は、ストレージグループ ID 2301、ストレージグループ RAID タイプ 2302、セグメント数 2303、割り当て権限セグメント数 2309、空きセグメント数 2304、割り当て権限セグメントビットマップ 2308、空きセグメントビットマップ 2307、ストレージ装置ポインタ 2305 から構成される。

【0080】

ストレージグループ ID 2301 は、当該ストレージグループ情報 2300 の識別子である。この識別子は、ストレージボックス 130 の識別子も含む。ストレージグループ RAID タイプ 2302 は、当該ストレージグループの RAID タイプである。実施例 1 における RAID タイプは、論理ボリューム RAID タイプ 2003 を説明したときに述べたとおりである。実施例 1 では、ストレージコントローラ 200 は、容量仮想化機能をもたないので、論理ボリュームを定義したときに、その容量分の領域が確保される。実施例 1 では、ストレージグループ情報 2300 の容量は、セグメントという単位に分割されているとする。論理ボリュームの容量が定義されると、その容量以上で最小の数のセグメントが、確保されることになる。セグメント数 2303 は、当該ストレージグループのセグメントの数を示す。割り当て権限セグメント数 2309 は、その中で、当該実ストレージシステム 100 が、割り当て権限をもっているセグメントの数を示す。ストレージボックス 130 の中のストレージグループの場合、複数の実ストレージシステムにシェアされるので、それぞれ、実ストレージシステム 100 ごとに、割り当て権限をもつセグメントの集合を決めておく。割り当て権限ビットマップ 2308 は、それぞれのセグメントについて、当該実ストレージシステム 100 が、割り当て権限をもっているか、いないかを表す。空きセグメント数 2304 は、該ストレージグループが割り当て権限をもっているセグメントの中で、空いた状態のセグメントの数を示す。空きセグメントビットマップ 2307 は、該ストレージグループが割り当て権限をもっているセグメントの中で、それぞれのセグメントが、空いた状態か、割り当て済みであることを示す情報である。ストレージ装置ポインタ 2305 の数は、当該ストレージグループに属するストレージ装置 160 の数であるが、これは、ストレージグループ RAID タイプ 2302 によって決まる値である。ストレージ装置ポインタ 2305 は、当該ストレージグループに属するストレージ装置 160 の識別子である。

【0081】

10

20

30

40

50

図14は、ストレージ装置情報2500のフォーマットである。ストレージ装置情報2500のフォーマットは、ストレージ装置識別子2510、接続パス数2501、接続パス識別子(接続パス)2502、ストレージタイプ2503、容量2504、一時書込み領域数2505、一時書込み領域アドレス2506、使用中フラグ2507である。

【0082】

ストレージ装置識別子2510は、当該ストレージ装置160の識別子である。接続パス数2501は、当該ストレージ装置160とボックスプロセッサ170との間のストレージパス197の数である。接続パス識別子2502は、接続されているパスの識別子を示す。ストレージタイプ2503は、当該ストレージ装置160が、HDD、フラッシュメモリ、などのどんな記憶媒体なのかを示す。容量2504は、当該ストレージ装置の容量である。なお、実施例1では、ストレージグループを構成するストレージ装置160のストレージタイプ2503と容量2504は、等しいとする。一時書込み領域数2505は、ストレージボックス130にデータを直接書き込んだときに、パリティグループ全体のデータを受領するため一時的に書き込んでおくために確保している領域の数である。実施例1においては、領域の大きさは、一つのパリティグループの容量である。一時書込み領域アドレス2506は、一時書込み領域のアドレスである。実施例1では、二つのストレージ装置の一時書込み領域を選択し、それぞれの領域にパリティグループのデータを書き込む。ただし、本発明は、一つの一時書込み領域を複数のn台のストレージ装置160から構成しておき、二つの書込み領域を選択し、2n台のストレージ装置に書き込んででも有効である。なお、一時書込み領域は、サーバ110からの書込みが、直接指定されない領域であるものとする。使用中フラグ2507は、一時書込み領域対応に存在するフラグで、当該一時書込み領域が、使用中であることを示す。

【0083】

次に、上記に説明した管理情報を用いて、ストレージコントローラ200が実行する動作の説明を行う。まず、ストレージコントローラ200の動作を説明する。ストレージコントローラ200の動作は、ストレージコントローラ200内のプロセッサ260が実行し、そのプログラムは、メモリ270に格納されている。図15は、メモリ270内に、格納された実施例1に関するプログラムが示されている。実施例1に関するプログラムは、ライト要求受付部4100、ライト異常終了対応処理部4200、リード処理実行部4300である。

【0084】

図16A、Bは、ライト要求受付部4100の処理フローである。ライト要求受付部4100は、ストレージコントローラ200が、サーバ110からライト要求を受け付けたときに実行される。

【0085】

図16Aから説明する。

ステップS6000：プロセッサ260は、受け取ったライト要求で指定された仮想論理ボリュームを、仮想論理ボリューム情報2055によって、論理ボリュームに変換する。

【0086】

ステップS6001a：受け取ったライト要求のアドレスが、パリティグループの先頭のアドレスかを判断する。そうでなければ、ステップS6007へジャンプする。

ステップS6001b：キャッシュ管理情報を割り当てる。

【0087】

ステップS6002：先頭の場合、一つ前の領域にキャッシュ管理情報2750が割り当てられているかをチェックする。割り当てられていない場合、直接転送を行わないので、キャッシュ210にライトデータを転送する従来の処理を実行するので、ステップS6027へジャンプする。

【0088】

ステップS6003：当該パリティグループ全体のキャッシュ管理情報2750がミスしているかをチェックする。割り当てられているものが一つでも存在する場合、実施例1

10

20

30

40

50

では直接転送に入らないので、ステップS 6 0 2 7へジャンプする。なお、割り当てられているものがあつた場合にも、直接転送を実行しても、本発明は有効である。

【0089】

ステップS 6 0 0 4：一つ前のキャッシュ管理情報2 7 5 0の最終アクセスアドレスをチェックし、今回のアクセスが、シーケンシャルアクセスになっているかチェックする。なっていない場合、直接転送を行わないので、キャッシュ2 1 0にライトデータを転送する従来の処理を実行するので、ステップS 6 0 2 7へジャンプする。

【0090】

ステップS 6 0 0 5：一つ前のキャッシュ管理情報2 7 5 0のすべての更新ビットマップ2 7 5 4、アクセス不可ビットマップ2 7 5 5がオフかチェックする。すべて、オフであれば、このキャッシュ管理情報2 7 5 0を、空きキャッシュ管理情報ポインタ(キュー)2 6 5 0に戻す。

10

【0091】

ステップS 6 0 0 6：ここで、ストレージボックス1 3 0に直接ライトデータを転送する準備処理に入る。ここでは、空きキャッシュ管理情報ポインタ2 6 5 0から、キャッシュ管理情報2 7 5 0を、パリティグループに相当する領域の、第1キャッシュ管理情報ポインタ2 0 2 2と、第2キャッシュ管理情報ポインタ2 0 2 3に割り当てる。そして、それぞれのキャッシュ管理情報2 7 5 0の直接転送フラグ2 7 5 7をオンにする。また、ブロックビットマップ2 7 5 3、更新ビットマップ2 7 5 4、アクセス不可ビットマップ2 7 5 5は、オールクリアする。次に、第1ストレージポインタ2 0 2 4と第2ストレージポインタ2 0 2 5に、当該パリティグループを格納するストレージボックス1 3 0のストレージ装置1 6 0の一時書込み領域アドレス2 5 0 6の中で、使用中フラグ2 5 0 7がオフのものの中から、異なるストレージ装置1 6 0の一時書込み領域アドレス2 5 0 6のふたつを選んで、設定する。選んだら、使用中フラグ2 5 0 7をオンにする。(なお、空き状態にあるスロット管理情報や未使用の一時書込み領域が不足した場合、割り当てをやめ、直接転送は実行せず、ステップS 6 0 2 7へジャンプし、従来の処理を実行する。)この後、ステップS 6 0 1 5にジャンプする。

20

【0092】

ステップS 6 0 0 7：このステップは、パリティグループの先頭以外のアドレスのライト要求を受け付けた場合に実行される。まず、対応する領域に、キャッシュ管理情報2 7 5 0が割り当てられているかチェックする。割り当てられていない場合、直接データ転送は実行しないので、ステップS 6 0 1 2を判定し、ステップS 6 0 2 7へジャンプする。

30

【0093】

ステップS 6 0 0 8：以下では、アクセスがシーケンシャルになっているかチェックする。本ステップでは、アクセスが対応するスロットの先頭かチェックする。先頭でない場合、ステップS 6 0 2 7へジャンプする。

【0094】

ステップS 6 0 0 9：先頭のアドレスの場合、一つ前の領域に、キャッシュ管理情報2 7 5 0が割り当てられているかチェックする。割り当てられていない場合、直接データ転送は実行しないので、ステップS 6 0 2 7へジャンプする。

40

【0095】

ステップS 6 0 1 0：割り当てられている場合、一つ前のキャッシュ管理情報2 7 5 0の直接転送フラグ2 7 5 7をチェックする。オフの場合、直接データ転送は実行しないので、ステップS 6 0 2 7へジャンプする。

【0096】

ステップS 6 0 1 1：次に、アクセスがシーケンシャルかを確認するため、一つ前のキャッシュ管理情報2 7 5 0の最終アクセスアドレスフラグ2 7 5 6を、チェックする。シーケンシャルになっていない場合、直接転送を中止するため、ステップS 6 0 2 7にジャンプする。加えて、直接転送中止フラグ2 7 5 8がオンになっているかチェックする。オンの場合、直接転送を中止するため、ステップS 6 0 2 7にジャンプする。オフの場合

50

、直接転送を実行するため、ステップ S 6 0 1 5 へジャンプする。

【 0 0 9 7 】

ステップ S 6 0 1 2 : 当該キャッシュ管理情報 2 7 5 0 の直接転送フラグ 2 7 5 7 をチェックする。オフの場合、直接データ転送は実行しないので、ステップ S 6 0 2 7 へジャンプする。

【 0 0 9 8 】

ステップ S 6 0 1 3 : 次に、アクセスがシーケンシャルかを確認するため、当該キャッシュ管理情報 2 7 5 0 の最終アクセスアドレスフラグ 2 7 5 6 を、チェックする。シーケンシャルになっていない場合、直接転送を中止するため、ステップ S 6 0 2 7 にジャンプする。

【 0 0 9 9 】

ステップ S 6 0 1 4 : 次に、直接転送中止フラグ 2 7 5 8 がオンになっているかをチェックする。オンの場合、直接転送を中止するため、ステップ S 6 0 2 7 にジャンプする。オフの場合、以下のステップで直接転送に入る。

【 0 1 0 0 】

ステップ S 6 0 1 5 : 本ステップでは、ストレージボックス 1 3 0 に、直接転送を行うように指示する。この際、ライトデータを、本来書き込むアドレスと二つの一時的に書き込むアドレスを指示する。また、本来書き込むアドレスをサーバ 1 9 8 に送る。

【 0 1 0 1 】

図 1 6 B の説明を行う。

ステップ S 6 0 1 6 : 完了をまつ。

ステップ S 6 0 1 7 : 正常終了か異常終了かを判別する。異常終了の場合、まず、サーバ 1 1 0 に異常終了報告を行う。また、パリティグループの先頭以外のアクセスの場合、ステップ S 6 0 2 6 へジャンプする。(パリティグループの先頭のアクセスの場合、まだ、正常終了報告を返したライトデータはないので、後処理を行う必要がない。)なお、異常終了の場合、ある時間たっても、ストレージボックス 1 3 0 からの完了報告がなく、ストレージコントローラ 2 0 0 の時間監視で、タイムオーバを起こした場合も含めるとする。

【 0 1 0 2 】

ステップ S 6 0 1 8 : 一時領域に格納したデータに対応する、ブロックビットマップ 2 7 5 3 , 更新ビットマップ 2 7 5 4 をセットする。

【 0 1 0 3 】

ステップ S 6 0 1 9 : 当該ライト要求のアクセスアドレスを、最終アクセスアドレスに記憶する。

ステップ S 6 0 2 0 : サーバ 1 1 0 に正常終了を報告する。

ステップ S 6 0 2 1 : パリティグループ全体の書込みが完了したかをチェックする。完了していない場合、処理を完了する。

ステップ S 6 0 2 2 : 完了した場合、ストレージボックス 1 3 0 に、パリティを生成して、ライトデータとパリティを書き込むべき領域に書き込むよう指示する。

ステップ S 6 0 2 3 : 完了をまつ。

【 0 1 0 4 】

ステップ S 6 0 2 4 : 正常終了か、異常終了かを判別する。異常終了の場合、直接転送を中止するための処理を実行するため、ステップ S 6 0 2 6 へジャンプする。なお、異常終了の場合、ある時間たっても、ストレージボックス 1 3 0 からの完了報告がなく、ストレージコントローラ 2 0 0 の時間監視で、タイムオーバを起こした場合も含めるとする。また、本ステップで異常終了した場合、パリティグループ全体のデータとパリティをストレージ装置 1 6 0 に書き込んでいる際に、障害が発生した場合、すでに、旧データを破壊してしまっている可能性があるため、当該パリティグループのすべてのキャッシュ管理情報 2 7 5 0 の旧データ破壊フラグ 2 7 5 9 をオンにする。

【 0 1 0 5 】

ステップ S 6 0 2 5 : ここでは、当該パリティグループに割り当てたキャッシュ管理情

10

20

30

40

50

報 2750 と一時書込み領域を開放する。具体的には、最後のキャッシュ管理情報 2750 以外のキャッシュ管理情報 2750 は、空きキャッシュ管理情報キュー 2650 に戻す。また、キューに戻すキャッシュ管理情報 2750 のブロックビットマップ等のビットマップ情報や最終アクセスアドレス、フラグ情報は、すべてオフにする。加えて、最後のキャッシュ管理情報 2750 は、最終アドレスアドレスフラグ 2756 以外の情報は、すべてオフにする。この後、処理を終了する。

【0106】

ステップ S6026：異常処理が発生したので、直接転送の後処理を実行するため、直接転送異常処理をコールする。この後、処理を終了する。

【0107】

ステップ S6027：キャッシュにライトデータを受け取る処理を実行する。この際、当該領域に、キャッシュ管理情報 2750 が割り当てられていない場合、キャッシュ管理情報 2750 を割り当てることになる。また、今回アクセスされたアドレスを最終アクセスアドレスに記憶する。加えて、ライト要求が正常終了した場合、当該要求でライトされた領域に対応するアクセス不可ビットマップ 2755 がオンの場合、オフする。さらに、当該キャッシュ管理情報 2750 のすべてのアクセス付加ビットマップがオフになったら、直接転送中止フラグ 2758 をオフする。

【0108】

図 17 は、ライト異常終了対応処理部 4200 の処理フローである。本処理も、プロセッサ 260 が、適宜実行する処理である。

【0109】

ステップ S7000：ストレージボックス 130 が、無応答で異常終了したのかをチェックする。無応答で処理を終了していない場合、ステップ S7004 へジャンプする。

【0110】

ストレージボックス 130 からストレージコントローラ 200 へ、障害通知があった場合は、このステップで、「NO」となり、ステップ S7004 にジャンプする。障害通知としては、ストレージボックス 130 で、冗長データを生成することができないとき、ライトデータをストレージコントローラ 200 に送れる場合が含まれる。

【0111】

また、ライトデータの中で、障害により、ストレージボックス 130 がストレージコントローラ 200 に送ることができないデータが存在する場合、その旨とそのデータを書き込むよう指定された領域のアドレスを、ストレージコントローラ 200 に伝える場合も含まれる。

【0112】

また、ライトデータ或いは冗長データを、障害により、記憶媒体に書き込むよう指定された領域に書き込めないとき、ストレージコントローラ 200 に送る場合も含まれる。

【0113】

ステップ S7001：無応答で処理を終了した場合、ストレージボックス 130 が電源ダウン等で、通信できない可能性がある。本ステップでは、一度、ストレージボックスにリセット要求を発行する。

【0114】

ステップ S7002：リセット要求が完了した場合、ステップ S7004 へジャンプする。

【0115】

ステップ S7003：ストレージボックス 130 がダウンしている等の原因があるので、ストレージ管理者、保守員等に、連絡して、ストレージボックス 130 が立ち上がるのをまつ。立ち上がった後、ステップ S7004 から実行を開始する。

【0116】

ステップ S7004：当該パリティグループに対応したキャッシュ管理情報 2750 のブロックビットマップ 2753，更新ビットマップ 2754 がオンのデータを格納した、

10

20

30

40

50

一時書込み領域のアドレスを算出する。また、当該パリティグループに対応したキャッシュ管理情報 2750 の直接転送中止フラグ 2758 をオンにする。

【0117】

ステップ S7005 : ストレージボックス 130 に、このアドレス (ステップ S7004 で計算したアドレス) を通知して、これらのデータを、ストレージコントローラ 200 に送るよう指示する。なお、実施例 1 では、ストレージボックス 130 には、サーバ 110 から受け取ったデータを転送するよう、指示しているが、ストレージボックス 130 で、すでにパリティを生成している場合、パリティもストレージボックス 130 から受け取って、キャッシュ 210 に格納するようにしてもよい。

【0118】

ステップ S7006 : 完了をまつ。ある時間待って、完了が返って、こなかった場合、ステップ S7000 へジャンプする。

【0119】

ステップ S7007 : 正常終了か異常終了かをチェックする。ストレージボックス 130 との間で、一部のデータが転送に失敗した場合や、ストレージボックス 130 から、一時書込み領域から読み出せなかったデータがあった旨の報告を受けた場合も、ステップ S7009 へジャンプする。

【0120】

ステップ S7008 : ストレージボックスとの間で、すべてのデータの転送が正常終了して、かつ、ストレージボックスから一時書込み領域のすべてのデータを読み出せた旨を受けとった場合、受け取ったデータを対応するスロット管理情報 2705 に対応したキャッシュ 210 の領域に 2 重に書き込む。この後、処理を終了する。

【0121】

ステップ S7009 : 転送に成功して、かつ、一時書込み領域から正常に読み出せたデータのみ、キャッシュ管理情報 2750 に対応した領域に、2 重に書き込む。

【0122】

ステップ S7010 : 転送に失敗したデータあるいは、一時書込み領域から読み出せなかったデータに対応するアクセス不可ビットマップ 2755 をオンにする。また、ライトデータ、データを書き込むよう指定された領域を記憶する。

【0123】

当該パリティグループに対応するアクセス不可ビットマップ 2755 アクセス不可ビットマップ 2755 にオンがない場合、(すべてのデータがキャッシュ 210 に正常に格納できた場合)、ストレージコントローラ 200 がパリティを生成して、データとパリティを、ストレージ装置 160 に書き込むことができる。アクセス不可ビットマップ 2755 がオンのビットをもつスロットの場合、このアクセス不可ビットマップ 2755 がオンに対応する領域に、新しいライト要求が入り、キャッシュ 210 にデータを正常に格納できれば、対応するアクセス不可ビットマップ 2755 をオフにできる。当該スロットに対応するキャッシュ管理情報 2750 にすべてのアクセス不可ビットマップ 2755 がオフになった場合、旧データ破壊フラグ 2759 がオフであれば、ストレージコントローラ 200 が持つ、キャッシュ 210 に旧データや旧パリティを格納して、これらのデータから、新パリティを生成して、データとパリティをストレージ装置 160 に格納する論理を適用できる。旧データ破壊フラグ 2759 がオンの場合、パリティの生成に旧データが使用できないので、当該パリティグループに対応するすべてのキャッシュ管理情報 2750 のすべてのアクセス不可ビットマップ 2755 がオフになった場合、ストレージコントローラ 200 が、残りのパリティグループのデータをキャッシュ 210 に格納し、パリティグループ全体のデータから、パリティを生成して、更新データとパリティを、ストレージ装置 160 に書き込むことができる。また、旧データ破壊 2759 がオフの場合にも、残りのパリティグループのデータをキャッシュ 210 に格納し、パリティグループ全体のデータから、パリティを生成してもよい。

【0124】

10

20

30

40

50

図18は、リード処理実行部4300の処理フローである。実施例1ではストレージボックス130からサーバに、要求されたデータを直接転送するようにする。もちろん、本発明は、リード要求で指定されたデータを、ストレージボックス130からストレージコントローラ200を経由して、サーバに送っても、本発明は有効である。

【0125】

ステップS5000：プロセッサ260は、受け取ったリード要求で指定された仮想論理ボリュームを、仮想論理ボリューム情報2055によって、論理ボリュームに変換し、対応する論理ボリューム情報2000を獲得する。

【0126】

ステップS5001：受け取ったリード要求のアドレスから、当該領域にキャッシュ管理情報2750が割り当てられているかをチェックする。割り当てられている場合、ステップS5006へジャンプする。

10

【0127】

ステップS5002：ストレージボックス130に、指定されたデータを、サーバ110に送るよう指示する。

ステップS5003：完了をまつ。

ステップS5004：サーバ110に完了報告を行い、処理を終了する。

ステップS5006：要求されたデータに対応するアクセス不可ビットマップ2755のビットがオンになっていないかをチェックする。オフであれば、ステップS5007へジャンプする。

20

ステップS5007：サーバ110に対し、当該データがアクセスできないという異常報告を行う。この後処理を終了する。

ステップS5008：直接転送フラグ2757がオンで、直接転送中止フラグ2758がオフかをチェックし、条件が満足した場合、ステップS5012へジャンプする。ストレージコントローラ200が、記憶した領域に対し、サーバ110からリード要求を受けつけたとき、リード要求を、異常終了させる。

ステップS5008：対応するブロックビットマップ2753がオンかをチェックする。オフの場合、ステップS5012へジャンプする。

ステップS5009：キャッシュから、要求されたデータを、サーバに転送する。

ステップS5010：転送が完了するのを待つ。

30

ステップS5011：サーバ110に対し、完了報告を行う。

ステップS5012：ストレージボックス130に、一時書込み領域の中の対応する領域から、サーバ110に送るよう指示する。

ステップS5013：完了をまつ。

ステップS5014：サーバ110に完了報告を行い、処理を終了する。

【0128】

次に、ストレージコントローラ200の指示にしたがって、ストレージボックス130のボックスプロセッサ170が実行する動作の説明を行う。実行するプログラムは、ボックスメモリ181に格納されている。図19は、ボックスメモリ181内に、格納された実施例1に関するプログラムが示されている。実施例1に関するプログラムは、ライトデータ受領部4400、ライトデータ書込み部4500、一時書込み領域転送部4600、リードデータ直接転送部4700である。

40

【0129】

図20は、ライトデータ受領部4400の処理フローである。ライトデータ受領部4400は、ストレージコントローラ200の指示にしたがって、サーバ110から、ライトデータを受け取り、一時書込み領域にデータを書き込む処理である。

ステップS8000：サーバ110に、データを送るよう指示する。

ステップS8001：転送の完了を待つ。

ステップS8002：正常終了であれば、ステップS8004へジャンプする。

ステップS8003：異常終了報告を、ストレージコントローラに返し、処理を終了す

50

る。

ステップ S 8 0 0 4 : 受け取ったライトデータを、ストレージコントローラ 2 0 0 から指定された、二つの一時書込み領域のアドレスに書き込む。

ステップ S 8 0 0 5 : 書込みが完了するのをまつ。

ステップ S 8 0 0 6 : 正常終了であれば、ステップ S 8 0 0 8 へジャンプする。

ステップ S 8 0 0 7 : 異常終了報告を、ストレージコントローラに返し、処理を終了する。

ステップ S 8 0 0 8 : 正常終了報告を、ストレージコントローラに返し、処理を終了する。

【 0 1 3 0 】

図 2 1 は、ライトデータ書込み部 4 5 0 0 の処理フローである。ライトデータ書込み部 4 5 0 0 は、ストレージコントローラ 2 0 0 の指示にしたがって、一時書込み領域からライトデータを読み出し、パリティを生成し、ライトデータとパリティを、ストレージコントローラ 2 0 0 から、指示されたストレージ装置 1 6 0 の領域に書き込む。

ステップ S 9 0 0 0 : ストレージコントローラ 2 0 0 から指定された二つの一時書込み領域の一つから、データを読み出す。

ステップ S 9 0 0 1 : 処理の完了をまつ。

ステップ S 9 0 0 2 : すべてのデータの読み出しが正常終了したかをチェックする。正常終了していれば、ステップ S 9 0 0 7 へジャンプする。

ステップ S 9 0 0 3 : もう一つの一時書込み領域から、異常終了したデータだけを読み出す。

ステップ S 9 0 0 4 : 完了をまつ。

ステップ S 9 0 0 5 : 指定したすべてのデータの読み出しが正常終了したかをチェックし、正常終了した場合、ステップ S 9 0 0 7 へジャンプする。

ステップ S 9 0 0 6 : ストレージコントローラ 2 0 0 に、異常終了を報告し、処理を完了する。

ステップ S 9 0 0 7 : 読み出したライトデータから、パリティを生成する。

ステップ S 9 0 0 8 : ライトデータとパリティを、ストレージコントローラ 2 0 0 が指定したストレージ装置 1 6 0 の領域に書き込む。

ステップ S 9 0 0 9 : 完了をまつ。

ステップ S 9 0 1 0 : すべての書込みが正常終了した場合、ストレージコントローラ 2 0 0 に、正常終了報告を行い、一つでも、異常終了があった場合、異常終了報告をストレージコントローラ 2 0 0 に行う。

【 0 1 3 1 】

図 2 2 は、一時書込み領域転送部 4 6 0 0 の処理である。一時書込み領域転送部 4 6 0 0 は、ストレージコントローラ 2 0 0 の指示にしたがって、一時書込み領域からライトデータを読み出し、ストレージコントローラ 2 0 0 に送る。一時書込み領域から、データを読み出す処理は、ライトデータ書込み処理部のステップ S 9 0 0 0 からステップ S 9 0 0 5 までの処理と同様である。読み出したデータをストレージコントローラ 2 0 0 に送る以下の処理のみが異なる。

ステップ S 1 0 0 0 0 : 正常に読み出せなかったデータが存在する処理である。正常に読み出せたデータ(ライトデータと冗長データ)とそのアドレスをすべてストレージコントローラ 2 0 0 に送る。また、読み出せなかったデータについては、読み出せなかった旨とそのアドレスをストレージコントローラ 2 0 0 に通知する。

【 0 1 3 2 】

その後、処理を完了する。

ステップ S 1 0 0 0 1 : すべてのデータを読み出せたことになるので、データとそのアドレスを、ストレージコントローラ 2 0 0 に送る。この後、処理を終了する。

【 0 1 3 3 】

以上の通り、障害により、ストレージボックス 1 3 0 で、冗長データを生成することが

10

20

30

40

50

できないとき、ライトデータをストレージコントローラ 200 に送る。

【0134】

また、ライトデータの中で、障害により、ストレージボックス 130 がストレージコントローラ 200 に送ることができないデータが存在する場合、その旨とそのデータを書き込むよう指定された領域のアドレスを、ストレージコントローラ 200 に伝える。

【0135】

また、ライトデータ或いは冗長データを、障害により、記憶媒体に書き込むよう指定された領域に書き込めないとき、ストレージコントローラ 200 に送る。

【0136】

図 23 は、リードデータ直接転送部 4700 の処理フローである。リードデータ直接転送部 4700 は、ストレージコントローラ 200 からの指示によって、指定されたデータをサーバ 110 に送る処理である。実施例 1 では、一時書込み領域の管理を、ストレージコントローラ 200 が行っているため、サーバ 110 が読み書きするデータを恒久的に格納した領域か、パリティを作成するための一時的な書込み領域の区別を、ストレージボックス 130 は、認識していない。このため、サーバ 110 が読み書きするデータを恒久的に格納した領域か、パリティを作成するための一時的な書込み領域に係わず、ストレージボックス 130 は、指定された領域のデータを送るだけである。以上から、サーバ 110 が読み書きするデータを恒久的に格納した領域のデータを送る処理と、パリティを作成するための一時的な書込み領域のデータを送る処理の処理フローは、同一となり、いずれも、リードデータ直接転送部 4700 が実行する。

ステップ S11000：指定された領域のデータをサーバ 110 に送る。

ステップ S11001：処理が完了するのを待つ。

ステップ S11002：正常終了した場合、正常終了報告を、異常終了した場合、異常終了報告を、ストレージコントローラ 200 に返す。この後、処理を終了する。

【実施例 2】

【0137】

次に、実施例 2 について説明する。実施例 2 では、ストレージコントローラ 200 が容量仮想化機能 (Thin Provisioning) をもつとする。また、ストレージコントローラ 200 とストレージボックス 130 が連携して、圧縮・重複排除機能を実現する。圧縮・重複排除機能を実現する場合、ライト処理は、ログストラクチャ方式を採る。この場合、ライトデータは、追記型となり、新しいアドレスに書き込まれる。

【0138】

図 24 は、実施例 2 における情報システムの構成を示す。実施例 2 における情報システムと実施例 1 の情報システムの相違は、ストレージボックス 130 の構成が異なる点である。実施例 2 におけるストレージボックス 130 は、ストレージコントローラ 200 と連携して、実施例 1 におけるストレージボックス 130 の構成要素に、複数の圧縮・伸長回路 1950、ハッシュ回路 1960 を備える点である。ただし、ボックスプロセッサ 170 で、これらの処理を実行しても、本発明は有効である。また、実施例 2 では、圧縮・伸長回路 1950、ハッシュ回路 1960 は、ボックスプロセッサ 170 対応に存在するものとする。ただし、圧縮・伸長回路 1950、ハッシュ回路 1960 は、ボックスプロセッサ 170 対応に存在しない場合にも、本発明は有効である。

【0139】

実施例 2 のサーバポート情報 198 のフォーマットは、第 1 のサーバポート情報と同じである。第 2 の実ストレージシステム 100 の構成は、実施例 1 と同じである。

【0140】

図 25 は、実施例 2 における実ストレージシステム 100 の共有メモリ 220 の中の実施例 2 に関する情報を示しており、ストレージシステム情報 2060、仮想論理ボリューム情報 2085、他ストレージシステム情報 2070、論理ボリューム情報 2000、ストレージボックス情報 2050、ストレージグループ情報 2300、ストレージ装置情報 2500、キャッシュ管理情報 2750、空きキャッシュ管理情報ポインタ 2650、に

10

20

30

40

50

よって構成される。また、実施例 2 では、新たに、実ページ情報 2100、空きページ管理情報ポインタ 2200、仮想ページ容量 2600 が新しく加わる。これらの情報は、容量仮想化機能をサポートしたために必要となる情報である。この中で、ストレージシステム情報 2060 は、実施例 1 と同様である。他ストレージシステム情報 2070 の形式も、実施例 1 と同様である。仮想論理ボリューム情報 2085 の形式も、実施例 1 と同様である。

【0141】

図 26 は、実施例 2 における論理ボリューム情報 2000 の形式を示したものである。実施例 2 でも、サーバ 110 がデータをリード・ライトする記憶装置は、仮想論理ボリュームである。また、サーバ 110 は、仮想論理ボリュームの ID、仮想論理ボリューム内のアドレス、読み書きしたいデータの長さを指定して、リード要求やライト要求を発行する。実ストレージシステム 100 は、サーバ 110 から読み書き要求を受け取ると、仮想ボリューム情報 2085 により、対応する論理ボリュームの識別子を認識する。実施例 2 でも、論理ボリュームの識別子は、実ストレージシステム 100 内でユニークな情報である。

10

【0142】

論理ボリューム情報 2000 は、論理ボリュームごとに存在する情報である。この情報は、論理ボリューム識別子 2001、論理容量 2002、論理ボリュームタイプ 2005、論理ボリューム RAID タイプ 2003、第 1 キャッシュ管理情報ポインタ 2022、第 2 キャッシュ管理情報ポインタ 2023、マッピング情報（マッピングアドレス）2004、圧縮後データ長 2018、ハッシュ値 2019、ログストラクチャポインタ 2005、ログストラクチャ書込みポインタ 2006、第 1 圧縮前データ格納領域ポインタ 2007、第 2 圧縮前データ格納領域ポインタ 2008、第 1 圧縮後データ格納領域ポインタ 2009、第 2 圧縮後データ格納領域ポインタ 2010、第 1 無効圧縮前ビットマップ 2020、第 2 圧縮前無効ビットマップ 2021、第 1 無効圧縮後ビットマップ 2011、第 2 圧縮後無効ビットマップ 2012、第 1 圧縮前書込みポインタ 2013、第 2 圧縮前書込みポインタ 2014、第 1 圧縮後書込みポインタ 2015、第 2 圧縮後書込みポインタ 2016、論理ボリューム直接転送中止フラグ 2030 より構成される。

20

【0143】

論理ボリューム識別子 2001、論理容量 2002、論理ボリュームタイプ 2005、論理ボリューム RAID タイプ 2003、第 1 キャッシュ管理情報ポインタ 2022、第 2 キャッシュ管理情報ポインタ 2023 までの情報は、第 1 の実施例と同様である。

30

【0144】

第 2 の実施例では、ストレージコントローラ 200 は、圧縮・重複排除機能をもつ。圧縮・重複排除機能をもつ場合、ライトデータは、ログストラクチャストリームで、書込み順に追記される。実施例 2 では、論理ボリューム対応に、ログストラクチャストリームをもつ。実施例 2 では、論理ボリュームごとにログストラクチャストリームは一つであるが、高速化のために、論理ボリュームごとに複数もっても、本発明は有効である。また、本発明は、ログストラクチャストリームが、論理ボリューム対応でなくとも有効である。マッピング情報（マッピングアドレス）2004 は、論理ボリュームの圧縮・重複排除単位が、どの実ページのどの相対アドレスに割り当てられているかを示す情報である。また、本発明では、ストレージボックス 130 にライトデータが直接転送され、パリティが生成されるまでは、一時的な領域に二重に格納されている。その場合、そのアドレスが二つ示される。割り当てられていない場合、ヌル状態となる。

40

【0145】

また、実施例 2 では、圧縮後データ長 2018 はその管理単位の圧縮後のデータ長、ハッシュ値 2019 はその管理単位のハッシュ値を示す。

【0146】

実施例 2 では、ストレージコントローラ 200 は、容量仮想化機能をもつ。この場合、データはページに書き込まれる。ログストラクチャ書込みポインタ 2006 には、ログストラクチャストリームに割り当てられた実ページ情報 2100 が、キュー状に接続されて

50

いるものである。ログストラクチャ書込みポインタ2006は、次に、データを書き込むページのアドレスを示すものである。実施例2では、ライトデータは、直接、ストレージボックス130に転送される。この際、ストレージボックス130では、データが二重化され、パリティグループのデータが揃うまで、パリティが生成されない。また、実施例2では、ストレージコントローラ200とストレージボックス130が、連携して圧縮・重複排除を行うため、実際に記憶装置に書き込むデータの量は、受け取ったライトデータに比べ小さくなる。

【0147】

また、データはすべて2重化するのので、実施例2では、圧縮重複・排除前のデータを格納するための領域を $m \times 2$ 面で $2m$ 個割り当てているものとする。第1圧縮前データ格納領域ポインタ2007、第2圧縮前データ格納領域ポインタ2008は、二重に圧縮前データを格納する領域で、それぞれ m 個の領域に対応し、それぞれが m 個存在する。第1圧縮後データ格納領域ポインタ2009、第2圧縮後データ格納領域ポインタ2010は、二重に圧縮前データを格納する領域で、それぞれ2個の領域に対応し、それぞれが2個存在する。

10

【0148】

また、実施例2では、ログストラクチャストリームで、ライトデータは新しい領域に書き込まれるので、それまでのデータは無効になる。これは、まだ、第1圧縮前データ格納領域ポインタ2007、第2圧縮前データ格納領域ポインタ2008、第1圧縮後データ格納領域ポインタ2009、第2圧縮後データ格納領域ポインタ2010が示す領域に存在する間にも発生する可能性がある。この場合、対応する領域の第1無効圧縮前ビットマップ2020、第2圧縮前無効ビットマップ2021、第1無効圧縮後ビットマップ2011、第2圧縮後無効ビットマップ2012のビットなどをオンにすることになる。

20

【0149】

また、次にサーバ110から受け取ったライトデータを書き込むべきアドレスが、第1圧縮前書込みポインタ2013、第2圧縮前書込みポインタ2014である。また、次に圧縮・重複排除後のデータを書き込むべきデータのアドレスが、第1圧縮後書込みポインタ2015、第2圧縮後書込みポインタ2016である。実施例2では、ライトデータは論理ボリューム単位に、ログストラクチャストリームに書き込まれ、基本的にはストレージボックス130に直接転送される。このため、障害が発生すると、論理ボリューム単位で、直接転送が中止になる。直接転送中止フラグ2758は、当該論理ボリュームのライトデータの直接転送が中止になったことを示す。

30

【0150】

キャッシュ管理情報2750は、基本的には、実施例1と同様である。ただし、最終アクセスアドレスフラグ2756はなくともよい。

【0151】

空きキャッシュ管理情報ポインタ2650の形式は、実施例1と同様である。

【0152】

実施例2の特徴は、実ストレージシステム100は、容量仮想化機能をサポートしている点である。通常、容量仮想化機能において、記憶領域の割り当て単位は、ページと呼ばれる。また、論理ボリュームは、通常、サーバ110が読み書きをする論理的な記憶装置である。ただし、本発明では、キャッシングのために使用する記憶装置の領域を論理ボリュームとして定義するのが特徴である。そして、この論理ボリュームにも、容量仮想化機能を適用して、ページを割り当てることで、実記憶領域を確保する。

40

【0153】

なお、実施例2では、論理ボリュームの空間は、仮想ページという単位で、分割されているものとし、実際の記憶装置グループは、実ページという単位で分割されているものとする。容量仮想化においては、論理ボリュームの記憶容量を、実際の記憶媒体の容量よりも大きく見せる。このため、仮想ページの数のほうが、実ページの数より大きいのが、一般的である。容量仮想化機能を実現した場合、ストレージコントローラ200がサーバ1

50

10からのライト不要で書込みを指示されたアドレスを含む仮想ページに実ページが割り当てていないとき、実ページを割り当てる。仮想ページ容量2600は、仮想ページの容量である。

【0154】

しかし、実施例2では、仮想ページ容量2600と実ページの容量は等しいというわけではない。というのは、実ページの容量は、RAIDのタイプにより異なってくる冗長データを含むためである。したがって、実ページの容量は、その実ページが割り当てられた記憶装置グループのRAIDタイプにより決まる。たとえば、RAID1のようにデータを2重に書き込む場合、実ページの容量は、仮想ページ容量2600の2倍になる。RAID5のように、N台の記憶装置の容量に対し、1台分の記憶装置の容量の冗長データを格納する場合、
10
仮想ページ容量2600の $(N+1)/N$ の容量が確保される。当然、RAID0のように、冗長性がない場合、仮想ページ容量2600と等しい容量が実ページの容量ということになる。

【0155】

なお、実施例2においては、仮想ページ容量2600はストレージシステム100の中で共通であるが、ストレージシステム100に仮想ページ容量2600に異なったものがあったとしても、本発明は有効である。なお、実施例2では、それぞれの記憶装置グループは、RAID5で構成されているものとする。もちろん、本発明は、記憶装置グループが任意のRAIDグループで構成されていても有効である。

【0156】

ストレージボックス情報2050の構成は、実施例1と同様である。

【0157】

図27に、実施例2のストレージグループ情報2300の形式を示す。第3のストレージグループ情報2300は、実施例1では含まれるセグメント数2303、割り当て権限セグメント数3209、空きセグメント数2304、割り当て権限セグメントビットマップ3208、空きセグメントビットマップ2307を含まないかわりに、空きページ情報ポインタ2310を含む。

【0158】

図28は、実ページ情報2100の形式である。実ページ情報2100は、実ページごとに存在する該当する実ページの管理情報である。実ページ情報2100は、ストレージグループ識別子2101、実ページアドレス2102、空きページ情報ポインタ2310、ページ無効ビットマップ2104、第1キャッシュ管理情報ポインタ2105、第2キャッシュ管理情報ポインタ2106から構成される。

【0159】

実施例2では、ログストラクチャストリームで、ライトデータは新しい領域に書き込まれるので、それまでのデータは無効になる。この場合、対応する領域のページ無効ビットマップ2104のビットをオンにする。実施例2では、当該ログストリームに書き込むデータを、パリティを生成し、ストレージ装置160に書き込むまでの間、キャッシュに2重に格納する。第1キャッシュ管理情報ポインタ2105、第2キャッシュ管理情報ポインタ2106は、このデータを格納したキャッシュ管理情報である。
40

【0160】

また、ここでは、生成したパリティもキャッシュ210に格納するため、このためのキャッシュ管理情報2750もポイントされる。第1キャッシュ管理情報ポインタ2105、第2キャッシュ管理情報ポインタ2106の数は、実ページ容量((パリティも含む)ノートのキャッシュ管理2750で管理するキャッシュ容量)である。

【0161】

ストレージグループ識別子2101は、該当する実ページが、どのストレージグループ280に割り当てられているかを示す。なお、実施例2でも、ストレージグループ280を構成するストレージ装置160は、一つのストレージボックス130に含まれるものとする。なお、本発明は、ストレージグループ280を構成するストレージ装置160
50

が、複数のストレージボックス 130 に含まれていても有効である。実ページアドレスは、対応ストレージグループ 280 の中で、どの相対的なアドレスに割り当てられているかを示す情報である。空きページポインタ 2103 は、この実ページに、仮想ページが割り当てられていない場合、有効な値となる。この場合、その値は、対応するストレージグループ 280 の中で、次の仮想ページが割り当てられていない空きページ情報をさす。

【0162】

空きページ情報ポインタ 2310 は、ストレージグループごとに設けられる情報である。図 29 は、空きページ情報ポインタ 2310 によって管理される空き実ページの集合を表している。空き実ページとは、仮想ページに割り当てられていない実ページを意味する。また、空き実ページに対応した実ページ情報 2100 を空き実ページ情報 2100 と呼ぶ。全体の構造を、空き実ページ情報キュー 2900 と呼ぶ。空きページ管理情報ポインタ 2200 (図 25 参照) は、先頭の空き実ページ情報 2100 のアドレスを指す。次に、先頭の空き実ページ情報 2100 の中の空きページポインタ 2103 が次の空き実ページ情報 2100 を指す。図 29 では、最後の空き実ページ情報 2100 の空きページポインタ 2103 は、空きページ管理情報ポインタ 2200 を示しているが、ヌル値でもよい。ストレージコントローラ 200 は、実ページを割り当てていない仮想ページに書込み要求を受け付けると、論理ボリューム RAID タイプ 2003 と割り当て範囲 2026 (図 9 参照) に該当する、ストレージグループのいずれか、例えば、該当する記憶装置グループの中の空き実ページ数の最も多いストレージグループに対応する空きページ管理情報ポインタ 2200 から、空き実ページを探し、仮想ページに割り当てる。仮想ページ容量 2600 は、仮想ページの容量である。

10

20

【0163】

ストレージ装置情報 2500 のフォーマットは、実施例 1 と同様である。ただし、一時書込み領域アドレス 2506 が示す領域は、実ページを割り当てていない領域とする。

【0164】

次に、上記に説明した管理情報を用いて、ストレージコントローラ 200 が実行する動作の説明を行う。まず、ストレージコントローラ 200 の動作を説明する。ストレージコントローラ 200 の動作は、ストレージコントローラ 200 内のプロセッサ 260 が実行し、そのプログラムは、メモリ 270 に格納されている。図 30 は、メモリ 270 内に、格納された実施例 2 に関するプログラムが示されている。実施例 2 に関するプログラムは、実施例 1 に加えて、重複排前処理部 4800 が加わる。

30

【0165】

次に、上記に説明した管理情報を用いて、ストレージコントローラ 200 が実行する動作の説明を行う。

【0166】

図 31 A、B は、実施例 2 のライト要求受付部 4100 の処理フローである。ライト要求受付部 4100 は、ストレージコントローラ 200 が、サーバ 110 からライト要求を受け付けたときに実行される。

ステップ S12000 : プロセッサ 260 は、受け取ったライト要求で指定された仮想論理ボリュームを、仮想論理ボリューム情報 2085 によって、論理ボリュームに変換する。

40

ステップ S12001 : 論理ボリューム直接転送中止フラグ 2030 がオンになっているかをチェックする。この場合、直接転送を中止するので、ステップ S12026 へジャンプする。

ステップ S12002 : 指定されたライトを行う領域に、キャッシュ管理情報 2750 が、割り当てられているかをチェックする。割り当てられている場合、ステップ S12004 へジャンプする。

【0167】

ステップ S12003 : 割り当てられていない場合、当該領域に、キャッシュ管理情報 2750 を割り当てる。ここでは、空きキャッシュ管理情報ポインタ 2650 から、キャ

50

ッシュ管理情報 2750 を、パリティグループに相当する領域の、第 1 キャッシュ管理情報ポインタ 2022 と、第 2 キャッシュ管理情報ポインタ 2023 に割り当てる。そして、それぞれのキャッシュ管理情報 2750 の直接転送フラグ 2757 をオンにする。また、ブロックビットマップ 2753、更新ビットマップ 2754、アクセス不可ビットマップ 2755 は、オールクリアする。なお、空き状態にあるスロット管理情報や未使用の一時書込み領域が不足した場合、割り当てをやめ、直接転送は実行せず、従来の処理を実行する。この後、ステップ S 12005 にジャンプする。

【0168】

ステップ S 12004：マッピング情報（マッピングアドレス）2004 を参照して、ライトを行う領域が、現在は、どこに格納されているかをチェックする。実ページに格納されている場合、対応する実ページ情報 2100 の対応するページ無効ビットマップ 2104 をオンにする。圧縮前格納領域存在する場合、第 1 圧縮前無効ビットマップ 2020、第 2 圧縮前無効ビットマップ 2021 の対応する無効ビットをオンにする。圧縮後格納領域存在する場合、第 1 無効圧縮後ビットマップ 2011、第 2 圧縮後無効ビットマップ 2012 の対応する無効ビットをオンにする。マッピング情報（マッピングアドレス）2004 がヌルの場合、特に何もしない。

10

【0169】

ステップ S 12005：ここでは、ライトを指定されたアドレスに対応する更新ビットマップ 2754、ブロックビットマップ 2753 をオンにする。また、マッピング情報（マッピングアドレス）2004 には、第 1 圧縮前書込みポインタ 2013、第 2 圧縮前書込みポインタ 2014 が示すアドレスを設定する。

20

【0170】

ステップ S 12006：本ステップでは、ストレージボックス 130 に、直接転送を行うように指示する。この際、ライトデータを書き込むアドレスとして、第 1 圧縮前書込みポインタ 2013、第 2 圧縮前書込みポインタ 2014 が示すアドレスと二つのアドレスを、一時的に書き込むアドレスとして指示する。また、第 2 の実施例では、ログストラクチャストリームでデータが書き込まれるので、論理ボリュームのバラバラなアドレスのデータが圧縮前格納領域に書き込まれることになる。実施例 2 では、ストレージボックス 130 に、論理ボリュームの識別子と書込みが発生した相対アドレス（実ページ内のオフセット）を送るものとする。もちろん、ストレージボックス 130 に、論理ボリュームの識別子と書込みが発生した相対アドレスをストレージボックス 130 にこれらの情報を送らなくとも、本発明は有効である。

30

ステップ S 12007：完了をまつ。

ステップ S 12008：正常終了か異常終了かを判別する。異常終了の場合、直接転送の中止を行うため、ステップ S 12025 へジャンプする。なお、異常終了の場合、ある時間たっても、ストレージボックス 130 からの完了報告がなく、ストレージコントローラ 200 の時間監視で、タイムオーバを起こした場合も含める。

【0171】

ステップ S 12009：ストレージボックス 130 から受け取ったハッシュ値を、対応するハッシュ値 2019 に設定する。また、第 1 圧縮前書込みポインタ 2013、第 2 圧縮前書込みポインタ 2014 を更新する。

40

ステップ S 12010：サーバ 110 に正常終了を報告する。

ステップ S 12011：第 1 圧縮前格納領域全体の書込みが完了したかをチェックする。完了していない場合、処理を完了する。

ステップ S 12012：完了した場合、ログストラクチャ書込みポインタ 2006 が示す実ページの相対アドレスに対応するキャッシュ管理情報 2750 が割り当てられているかをチェックする。割り当てられていれば、ステップ S 12014 へジャンプする。

ステップ S 12013：割り当てられていなければ、キャッシュ管理情報を割り当てる。割り当てた場合、また、ブロックビットマップ 2753、更新ビットマップ 2754、アクセス不可ビットマップ 2755 をオールクリアする。

50

ステップ S 1 2 0 1 4 : 実施例 2 では、第 1 圧縮前格納領域全体の書込みが完了した場合、ストレージボックス 1 3 0 から、第 1 圧縮前格納領域に格納されていたデータの論理ボリュームの相対アドレスが送られてくるものとする。ここでは、第 1 圧縮前格納領域に格納されていたデータの論理ボリュームの相対アドレスをパラメータにして、重複排除前処理部 4 8 0 0 をコールする。

【 0 1 7 2 】

図 3 1 B の説明を行う。

ステップ S 1 2 0 1 5 : 処理終了後以下の処理を実行する。ストレージボックス 1 3 0 に、それぞれの重複排除単位の論理ボリュームの相対アドレスと、その単位が、重複排除できないか、できる可能性があるかを示す情報をおくる。加えて、できる可能性がある場合、データを読み出し、チェックする必要のあるアドレスを通知する。さらに、重複排除できなかったデータを、圧縮し、第 1 圧縮後書込みポインタ 2 0 1 5、第 2 圧縮後書込みポインタ 2 0 1 6 が示す、アドレスから格納していくように指示する。

ステップ S 1 2 0 1 6 : 完了をまつ。

【 0 1 7 3 】

ステップ S 1 2 0 1 7 : 正常終了か、異常終了かを判別する。異常終了の場合、直接転送を中止するための処理を実行するため、ステップ S 1 2 0 2 5 へジャンプする。なお、異常終了の場合、ある時間たっても、ストレージボックス 1 3 0 から の完了報告がなく、ストレージコントローラ 2 0 0 の時間監視で、タイムオーバを起こした場合も含める。

【 0 1 7 4 】

ステップ S 1 2 0 1 8 : 対応する圧縮後データ長 2 0 1 8 に、ストレージボックス 1 3 0 から受け取ったデータ長を設定する。重複排除できたデータに対応するマッピング情報 (マッピングアドレス) 2 0 0 4 に、同じデータを格納していたアドレスを設定する。また、重複排除できなかったデータについては、対応するマッピング情報 (マッピングアドレス) 2 0 0 4 に、ストレージボックス 1 3 0 から受け取った、このデータを格納した第 1 圧縮後データ格納領域ポインタ 2 0 0 9 と第 2 圧縮後データ格納領域 2 0 1 2 のアドレスを設定する。また、圧縮を行ったデータに対応する論理ボリュームの相対アドレスに割り当てたキャッシュ管理情報の更新ビットマップ 2 7 5 4、ブロックビットマップ 2 7 5 3 をオフする。キャッシュ管理情報のすべての更新ビットマップ 2 7 5 4、ブロックビットマップ 2 7 5 3 がオフになった場合、キャッシュ管理情報 2 7 5 0 を、空きキャッシュ管理情報ポインタ 2 6 5 0 に戻す。

【 0 1 7 5 】

ステップ S 1 2 0 1 9 : 圧縮後データに対応する実ページに割り当てたキャッシュ管理情報 2 7 5 0 の対応するブロックビットマップ 2 7 5 3、更新ビットマップ 2 7 5 4 をオンにする。

ステップ S 1 2 0 2 0 : 第 1 圧縮後データ格納領域ポインタ 2 0 0 9 が満杯になった (格納したデータ長が、パリティグループ長になった) かを、チェックする。なっていないければ、処理を終了する。

ステップ S 1 2 0 2 1 : 満杯になっている場合、ストレージボックス 1 3 0 に、パリティを生成し、ログストラクチャ書込みポインタ 2 0 0 6 が示す領域からデータと生成したパリティを記憶装置に書き込むよう、指示する。

ステップ S 1 2 0 2 2 : 完了をまつ。

【 0 1 7 6 】

ステップ S 1 2 0 2 3 : 正常終了か、異常終了かを判別する。異常終了の場合、直接転送を中止するための処理を実行するため、ステップ S 1 2 0 2 5 へジャンプする。なお、異常終了の場合、ある時間たっても、ストレージボックス 1 3 0 から の完了報告がなく、ストレージコントローラ 2 0 0 の時間監視で、タイムオーバを起こした場合も含める。

【 0 1 7 7 】

ステップ S 1 2 0 2 4 : 圧縮後データを格納した実ページと相対アドレスを、対応するマッピング情報 2 0 0 4 に格納する。また、ここでは、当該パリティグループに割り当て

10

20

30

40

50

たキャッシュ管理情報 2750 を開放する。具体的には、最後のキャッシュ管理情報 2750 以外のキャッシュ管理情報 2750 は、空きキャッシュ管理情報キュー 2650 に戻す。また、キューに戻すキャッシュ管理情報 2750 のブロックビットマップ 2753 等のビットマップ情報やフラグ情報は、すべてオフにする。この後、処理を終了する。

【0178】

ステップ S12025：異常処理が発生したので、直接転送の後処理を実行するため、直接転送異常処理をコールする。この後、処理を終了する。

ステップ S12026：キャッシュにライトデータを受け取る処理に入る。この処理は、公知の技術であるため、詳細を説明しない。終了後、処理を完了させる。

【0179】

図 32 は、実施例 2 のライト異常終了対応処理部 4200 の処理フローである。本処理も、プロセッサ 260 が、適宜実行する処理である。

【0180】

ステップ S13000：ストレージボックス 130 が、無応答で異常終了したのかをチェックする。無応答で処理を終了していない場合、ステップ S13004 へジャンプする。

【0181】

ステップ S13001：無応答で処理を終了した場合、ストレージボックス 130 が電源ダウン等で、通信できない可能性がある。本ステップでは、一度、ストレージボックスにリセット要求を発行する。

ステップ S13002：リセット要求が完了した場合、ステップ S13004 へジャンプする。

ステップ S13003：ストレージボックス 130 がダウンしている等の原因があるので、ストレージ管理者、保守員等に、連絡して、ストレージボックス 130 が立ち上がるのをまつ。立ち上がった後、ステップ S13004 から実行を開始する。

ステップ S13004：当該ログストラクチャストリームに対応する論理ボリューム直接転送中止フラグ 2030 をオンにする。

【0182】

ステップ S13005：ストレージボックス 130 に、m 個の第 1 圧縮前データ格納領域ポインタ 2007、m 個の第 2 圧縮前データ格納領域ポインタ 2008、2 個の第 1 圧縮後データ格納領域ポインタ 2009、2 個の第 2 圧縮後データ格納領域ポインタ 2010 が示すアドレスと、それぞれの格納領域に対応した第 1 無効圧縮前ビットマップ 2020、第 2 圧縮前無効ビットマップ 2021、第 1 無効圧縮後ビットマップ 2011、第 2 圧縮後無効ビットマップ 2012 を通知して、無効ビットがオフの領域に格納されたデータとその論理ボリュームのアドレスを、ストレージコントローラ 200 に送るよう指示する。

【0183】

なお、実施例 2 では、ストレージボックス 130 には、サーバ 110 から受け取ったデータを転送するよう、指示しているが、ストレージボックス 130 で、すでにパリティを生成している場合、パリティもストレージボックス 130 から受け取って、キャッシュ 210 に格納するようにしてもよい。

【0184】

ステップ S13006：完了をまつ。ある時間待って、完了が返って、こなかった場合、ステップ S13000 へジャンプする。

ステップ S13007：一部のデータが転送に失敗した場合、ステップ S13008 へジャンプする。

【0185】

ステップ S13008：すべてのデータの転送が正常終了した場合、以下の処理を実行する。第 1 圧縮前データ格納領域、あるいは、第 2 圧縮前データ格納領域 から受け取ったデータは、対応する論理ボリュームの相対アドレスから論理ボリュームの空間に割り当てた対応するスロット管理情報 2705 に対応したキャッシュ 210 の領域に 2 重に書き

10

20

30

40

50

込む。また、第1圧縮後データ格納領域、あるいは、第2圧縮後データ格納領域 から受け取ったデータは、実ページに空間に割り当てた対応するスロット管理情報2705に対応したキャッシュ210の領域に2重に書き込む。この後、処理を終了する。

【0186】

ステップS13009：転送に成功したデータのみ、キャッシュ管理情報2750に対応した領域に、2重に書き込む。

ステップS1310：転送に失敗したデータに対応するアクセス不可ビットマップ2755をオンにする。

【0187】

当該パリティグループに対応するアクセス不可ビットマップ2755にオンがない場合、（すべてのデータがキャッシュ210に正常に格納できた場合）、ストレージコントローラ200がパリティを生成して、データとパリティを、ストレージ装置160に書き込むことができる。

10

【0188】

アクセス不可ビットマップ2755がオンのビットをもつスロットの場合、このアクセス不可ビットマップ2755がオンに対応する領域に、新しいライト要求が入り、キャッシュ210にデータを正常に格納できれば、対応するアクセス不可ビットマップ2755をオフにできる。当該スロットに対応するキャッシュ管理情報2750にすべてのアクセス不可ビットマップ2755がオフになった場合、旧データ破壊フラグ2759がオフであれば、ストレージコントローラ200が、キャッシュ210に旧データや旧パリティを読み込み、これらから、パリティを生成して、データとパリティをストレージ装置160に格納する論理を適用できる。

20

【0189】

旧データ破壊フラグ2759がオンの場合、パリティの生成に旧データが使用できないので、当該パリティグループに対応するすべてのキャッシュ管理情報2750のすべてのアクセス不可ビットマップ2755がオフになった場合、ストレージコントローラ200が、残りのパリティグループのデータをキャッシュ210に読み込み、パリティを生成して、パリティグループ全体のデータとパリティを、ストレージ装置160に書き込むことができる。

【0190】

図33は、実施例2のリード処理実行部4300の処理フローである。実施例2では、ストレージボックス130からサーバに、要求されたデータを直接転送するようにする。もちろん、本発明は、リード要求で指定されたデータを、ストレージボックス130からストレージコントローラ200を経由して、サーバに送っても、本発明は有効である。

30

【0191】

ステップS14000：プロセッサ260は、受け取ったリード要求で指定された仮想論理ボリュームを、仮想論理ボリューム情報2085によって、論理ボリュームに変換し、対応する論理ボリューム情報2000を獲得する。

【0192】

ステップS14001：受け取ったリード要求のアドレスから、マッピング情報（マッピングアドレス）2004，キャッシュ管理情報ポインタ2022，実ページ情報2100、キャッシュ管理情報2750のブロックビットマップ2753などから、当該要求で、指定されたデータが、キャッシュ210にヒットしているかをチェックする。ミスであれば、ステップS14002へジャンプする。ヒットであれば、ステップS14011へジャンプする。

40

【0193】

ステップS14002：マッピング情報（マッピングアドレス）2004を参照して、当該データが、圧縮前データ格納領域、あるいは、圧縮後データ格納領域に格納されているかをチェックする。そうであれば、ステップS14006へジャンプする。

【0194】

50

ステップ S 1 4 0 0 3 : マッピング情報 (マッピングアドレス) 2 0 0 4 に格納された実ページ上のアドレスを指定して、ストレージボックス 1 3 0 に、格納されたデータを、サーバ 1 1 0 に送るよう指示する。

ステップ S 1 4 0 0 4 : 完了をまつ。

ステップ S 1 4 0 0 5 : サーバ 1 1 0 に完了報告を行い、処理を終了する。

ステップ S 1 4 0 0 6 : 要求されたデータに対応するアクセス不可ビットマップ 2 7 5 5 のビットがオンになっていないかをチェックする。オフであれば、ステップ S 1 4 0 0 8 へジャンプする。

ステップ S 1 4 0 0 7 : サーバ 1 1 0 に対し、当該データがアクセスできないという異常報告を行う。この後処理を終了する。

ステップ S 1 4 0 0 8 : ストレージボックス 1 3 0 に、マッピング情報に示された第 1 圧縮前データ格納領域と第 2 圧縮前データ格納領域の二つのアドレス、あるいは、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域の二つのアドレスから、サーバ 1 1 0 に送るよう指示する。

ステップ S 1 4 0 0 9 : 完了をまつ。

ステップ S 1 4 0 1 0 : サーバ 1 1 0 に完了報告を行い、処理を終了する。

ステップ S 1 4 0 1 1 : キャッシュから、要求されたデータを、サーバに転送する。

ステップ S 1 4 0 1 2 : 転送が完了するのを待つ。

ステップ S 1 4 0 1 3 : サーバ 1 1 0 に対し、完了報告を行い、処理を終了する。

【 0 1 9 5 】

図 3 4 は、重複排除前処理部 4 8 0 0 の処理フローである。重複排除処理とは、まったく同じ内容のデータを格納しないことにより、容量を削減する機能である。効率よく行うため、以下のステップを踏む。(1) データのハッシュ値をとり、ハッシュ値が同じ値になるデータの集合を見つける。(2) ハッシュ値が同じ値になるデータに関しては、データそのものを比較し、一致しているものがみつかった場合、そのデータを格納しない。実施例 2 では、(1) を、ストレージコントローラで実現し、(2) をストレージボックスで実現する。重複排除前処理部 4 8 0 0 は、(1) の処理を実行する。この処理自体は、公知であるため、簡単に説明する。

【 0 1 9 6 】

ステップ S 1 5 0 0 0 : マッピングテーブルを参照して、満杯になった第 1 圧縮前データ格納領域に格納されたデータのハッシュ値と同じハッシュ値をもつものを探す。それぞれのデータごとに、同じハッシュ値をもつ場合、そのデータが格納されているアドレスの集合を獲得する。以上を実行して、処理を完了する。

【 0 1 9 7 】

次に、ストレージコントローラ 2 0 0 の指示にしたがって、ストレージボックス 1 3 0 のボックスプロセッサ 1 7 0 が実行する動作の説明を行う。実行するプログラムは、ボックスメモリ 1 8 1 に格納されている。

【 0 1 9 8 】

図 3 5 は、ボックスメモリ 1 8 1 内に、格納された実施例 2 に関するプログラムが示されている。実施例 2 に関するプログラムは、実施例 1 に加えて、データ移動部 4 9 0 0、ボックス間データ転送部 4 9 5 0 が含まれる。

【 0 1 9 9 】

図 3 6 は、実施例 2 におけるライトデータ受領部 4 4 0 0 の処理フローである。ライトデータ受領部 4 4 0 0 は、ストレージコントローラ 2 0 0 の指示にしたがって、サーバ 1 1 0 から、ライトデータを受け取り、一時書込み領域にデータを書き込む処理である。

ステップ S 1 6 0 0 0 : サーバ 1 1 0 に、データを送るよう指示する。

ステップ S 1 6 0 0 1 : 転送の完了を待つ。

ステップ S 1 6 0 0 2 : 正常終了であれば、ステップ S 1 6 0 0 4 へジャンプする。

ステップ S 1 6 0 0 3 : 異常終了報告を、ストレージコントローラに返し、処理を終了する。

10

20

30

40

50

ステップ S 1 6 0 0 4 : 受け取ったライトデータのハッシュ値を計算する。

ステップ S 1 6 0 0 5 : ストレージコントローラ 2 0 0 から指定された、二つの一時書込み領域のアドレスに書き込む。

ステップ S 1 6 0 0 6 : 書込みが完了するのをまつ。

ステップ S 1 6 0 0 7 : 正常終了であれば、ステップ S 1 6 0 0 9 へジャンプする。

ステップ S 1 6 0 0 8 : 異常終了報告を、ストレージコントローラに返し、処理を終了する。

ステップ S 1 6 0 0 9 : ハッシュ値と正常終了報告を、ストレージコントローラに返し、処理を終了する。

【 0 2 0 0 】

図 3 7 は、実施例 2 におけるデータ移動部 4 9 0 0 の処理フローである。データ移動部は、まず、重複排除処理の(2)の部分に関する処理を実行する。具体的には、ストレージコントローラから受け取った、同じハッシュ値をもつデータのアドレスのデータを読み出し内容が一致するかどうかをチェックする。一致したものがない場合、圧縮処理を行い、圧縮後のデータを、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域の指定されたアドレスから格納していく。この圧縮後のデータの長さは、ストレージコントローラに返す。また、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域が満杯になった場合、次の第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域の先頭からデータを格納する。加えて、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域が満杯になった場合、この旨をストレージコントローラに報告する。

【 0 2 0 1 】

ステップ S 1 7 0 0 0 : 最初のデータを実行対象とする。

ステップ S 1 7 0 0 1 : 対象としたデータが、重複排除できる可能性のあるデータかどうかをチェックする。具体的には、同じハッシュ値をもつデータがあるかどうかをチェックする。同じハッシュ値をもつデータがない場合、ステップ S 1 7 0 0 6 へジャンプする。

【 0 2 0 2 】

ステップ S 1 7 0 0 2 : 指定されたアドレスに該当するデータを読み出す要求を発行する。この場合、複数のアドレスが指定された場合、複数の読み出し要求を並行して発行する。また、指定されたアドレスが当該ストレージボックス 1 3 0 うちストレージ装置 1 6 0 である場合、そのストレージに読み出す領域のアドレスを指定して、要求を発行する。また、当該ストレージボックス 1 3 0 のストレージ装置 1 6 0 でない場合、そのストレージ装置を含むストレージボックス 1 3 0 に要求を発行する。

【 0 2 0 3 】

ステップ S 1 7 0 0 3 : 要求が完了するのをまつ。

ステップ S 1 7 0 0 4 : 読み出したデータの中に、当該データと一致するものがないかチェックする。一致するものがない場合、ステップ S 1 7 0 0 6 にジャンプする。

ステップ S 1 7 0 0 5 : 当該データが重複排除できたことを記憶しておく。この後、ステップ S 1 7 0 0 8 へジャンプする。

ステップ S 1 7 0 0 6 : 扱っているデータを圧縮し、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域に格納する。

ステップ S 1 7 0 0 7 : 第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域が満杯になったかチェックする。満杯でなければ、ステップ S 1 7 0 0 9 にジャンプする。

ステップ S 1 7 0 0 8 満杯になった場合、これを記憶し、次の第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域を使用するようにする。

【 0 2 0 4 】

ステップ S 1 7 0 0 9 : 指定された第 1 圧縮前データ格納領域と第 2 圧縮前データ格納領域のすべてのデータの処理が実行されたかチェックする。実行されていない場合、次の第 1 圧縮前データ格納領域と第 2 圧縮前データ格納領域のデータを実行対象として、ステップ S 1 7 0 0 1 へジャンプする。

【 0 2 0 5 】

10

20

30

40

50

ステップ S 1 7 0 1 0 : コントローラに完了報告を通知する。この際、重複排除できたデータのアドレスを通知する。また、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域が満杯になった場合、この旨を報告する。

【 0 2 0 6 】

図 3 8 は、実施例 2 におけるライトデータ書込み部 4 5 0 0 の処理フローである。ライトデータ書込み部 4 5 0 0 は、ストレージコントローラ 2 0 0 の指示にしたがって、第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域からライトデータを読み出し、パリティを生成し、ライトデータとパリティを、ストレージコントローラ 2 0 0 から、指示されたストレージ装置 1 6 0 の領域に書き込む。

【 0 2 0 7 】

ステップ S 1 8 0 0 0 : ストレージコントローラ 2 0 0 から指定された二つの第 1 圧縮後データ格納領域と第 2 圧縮後データ格納領域の一つから、データを読み出す。

ステップ S 1 8 0 0 1 : 処理の完了をまつ。

ステップ S 1 8 0 0 2 : すべてのデータの読み出しが正常終了したかをチェックする。正常終了していれば、ステップ S 1 8 0 0 5 へジャンプする。

ステップ S 1 8 0 0 3 : もう一つの圧縮後データ格納領域から、異常終了したデータだけを読み出す。

ステップ S 1 8 0 0 4 : 完了をまつ。

ステップ S 1 8 0 0 5 : 指定したすべてのデータの読み出しが正常終了したかをチェックし、正常終了した場合、ステップ S 1 8 0 0 7 へジャンプする。

ステップ S 1 8 0 0 6 : ストレージコントローラ 2 0 0 に、異常終了を報告し、処理を完了する。

ステップ S 1 8 0 0 7 : 読み出したライトデータから、パリティを生成する。

ステップ S 1 8 0 0 8 : ライトデータとパリティを、ストレージコントローラ 2 0 0 が指定したストレージ装置 1 6 0 の領域に書き込む。

ステップ S 1 8 0 0 9 : 完了をまつ。

ステップ S 1 8 0 1 0 : すべての書込みが正常終了した場合、ストレージコントローラ 2 0 0 に、正常終了報告を行い、一つでも、異常終了があった場合、異常終了報告をストレージコントローラ 2 0 0 に行う。

【 0 2 0 8 】

図 3 9 は、実施例 2 における一時書込み領域転送部 4 6 0 0 の処理フローである。一時書込み領域転送部 4 6 0 0 は、ストレージコントローラ 2 0 0 の指示にしたがって、第 1 圧縮前データ格納領域、第 2 圧縮前データ格納領域、第 1 圧縮後データ格納領域、第 2 圧縮後データ格納領域から、データを読み出し、ストレージコントローラ 2 にデータを転送する。

【 0 2 0 9 】

ステップ S 1 9 0 0 0 : ストレージコントローラ 2 が指定した、第 1 圧縮前データ格納領域、第 1 圧縮後データ格納領域のデータを読み出すため、対応する記憶装置に、データ読み出す要求を発行する。

ステップ S 1 9 0 0 1 : 要求が完了するのをまつ。

ステップ S 1 9 0 0 2 : すべてのデータが正常に読みだせたかをチェックする。読み出せていれば、ステップ S 1 9 0 0 7 へジャンプする。

ステップ S 1 9 0 0 3 : 読み出せなかったデータと同じデータを格納した第 2 圧縮前データ格納領域、第 2 圧縮後データ格納領域からデータを読み出すため、対応する記憶装置から、データを読み出す要求を発行する。

ステップ S 1 9 0 0 4 : 要求が完了するのをまつ

ステップ S 1 9 0 0 5 : すべてのデータが正常に読みだせたかをチェックする。読み出せていれば、ステップ S 1 9 0 0 7 へジャンプする。

【 0 2 1 0 】

ステップ S 1 9 0 0 6 : 正常に読み出せたデータとそのアドレスをすべてストレージコ

10

20

30

40

50

ントローラ 200 に送る。また、読み出せなかったデータについては、読み出せなかった旨とそのアドレスをストレージコントローラ 200 に通知する。その後、処理を完了する。

ステップ S 19007：すべてのデータを読み出せたことになるので、データとそのアドレスを、ストレージコントローラ 200 に送る。この後、処理を終了する。

【0211】

図 40 は、実施例 2 におけるリードデータ直接転送部 4700 の処理フローである。リードデータ直接転送部 4700 は、ストレージコントローラ 200 からの指示によって、指定されたデータをサーバ 110 に送る処理である。実施例 2 では、圧縮前データ格納領域と圧縮後データ格納領域の管理を、ストレージコントローラ 200 が行っているため、サーバ 110 が読み書きするデータを恒久的に格納した領域か、パリティを作成するための一時的な書込み領域の区別を、ストレージボックス 130 は、認識していない。このため、サーバ 110 が読み書きするデータを恒久的に格納した領域か、パリティを作成するための一時的な書込み領域に係わらず、ストレージボックス 130 は、指定された領域のデータを送るだけである。以上から、サーバ 110 が読み書きするデータを恒久的に格納した領域のデータを送る処理と、パリティを作成するための一時的な書込み領域のデータを送る処理の処理フローは、同一となり、いずれも、リードデータ直接転送部 4700 が実行する。

10

【0212】

ステップ S 20000：指定された領域のデータをサーバ 110 に送る。

ステップ S 20001：処理が完了するのを待つ。

20

ステップ S 20002：正常終了した場合、正常終了報告を、異常終了した場合、異常終了報告を、ストレージコントローラ 200 に返す。この後、処理を終了する。

【0213】

図 41 は、実施例 2 におけるボックス間データ転送部 4950 の処理フローである。ボックス間データ転送部 4950 は、他のストレージボックス 130 から、要求されたデータを送る際に、実行される処理フローである。

【0214】

ステップ S 21000：指定されたストレージ装置 160 のアドレスから、データを読み出すよう、ストレージ装置 160 に要求を発行する。

ステップ S 21001：完了をまつ。

30

ステップ S 22001：読み出し要求を受け取ったストレージボックス 130 に読み出したデータを送信する。完了後、処理を終了する。

【0215】

本発明は、サーバ、ストレージコントローラとストレージボックスがネットワークを介して、接続された環境で、サーバからストレージボックスにデータを書き込んだ場合、障害が発生した場合、書き込んだデータをストレージコントローラのキャッシュに格納することにより、ストレージコントローラのもつ障害に対応した各種の処理を適応することができる。

【符号の説明】

【0216】

40

100 実ストレージシステム

110 サーバ

120 ネットワーク

130 ストレージボックス

150 ストレージ管理サーバ

160 ストレージ装置

175 共有ストレージ装置

170 ボックスプロセッサ

180 ボックスポート

181 ボックスメモリ

50

1 9 0	仮想ストレージシステム	
1 9 5	サーバポート	
1 9 7	ストレージバス	
1 9 8	サーバポート情報	
2 0 0	ストレージコントローラ	
2 1 0	キャッシュ	
2 2 0	共有メモリ	
2 3 0	内部ストレージ装置	
2 5 0	接続装置	
2 6 0	プロセッサ	10
2 7 0	メモリ	
2 7 5	バッファ	
2 0 5 0	ストレージボックス情報	
2 0 5 5	仮想論理ボリューム情報	
2 0 7 0	他ストレージシステム情報	
2 0 8 4	仮想論理ボリューム情報	
2 0 0 0	論理ボリューム情報	
2 0 6 0	ストレージシステム情報	
2 3 0 0	ストレージグループ情報	
2 3 5 0	ストレージ装置情報	20
2 7 5 0	キャッシュ管理情報	
2 6 5 0	空きキャッシュ管理情報ポインタ	
4 1 0 0	ライト要求受付部	
4 2 0 0	ライト異常終了対応処理部	
4 3 0 0	リード処理実行部	
4 4 0 0	ライトデータ受領部	
4 5 0 0	ライトデータ書込み部	
4 6 0 0	一時書込み領域転送部	
4 7 0 0	リードデータ直接転送部	
1 9 6	圧縮/伸長回路	30
1 9 7	ハッシュ回路	
2 1 0 0	実ページ管理情報	
2 2 0 0	空きページ管理情報ポインタ	
2 6 0 0	仮想ページ容量	
4 8 0 0	重複排除前処理部	
4 9 0 0	データ移動部	
4 9 5 0	ボックス間データ転送部	

【 図面 】

【 図 1 】

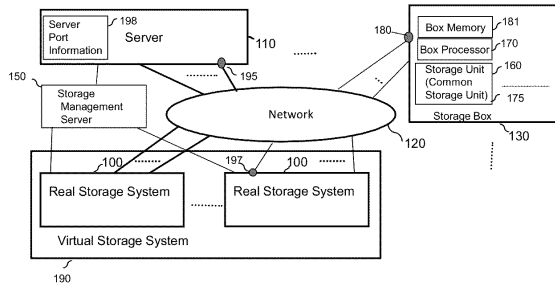
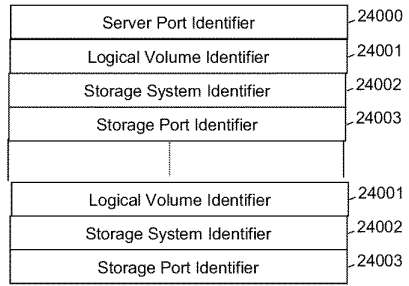


Fig. 1

【 図 2 】



Sever port Information 198

Fig. 2

10

【 図 3 】

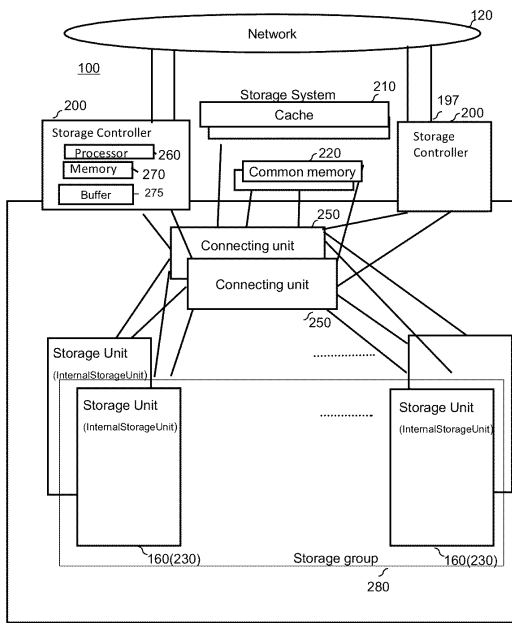


Fig. 3

【 図 4 】

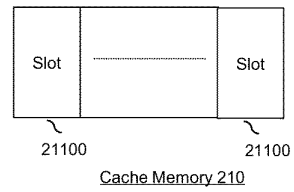


Fig.4

20

30

40

50

【 図 5 】

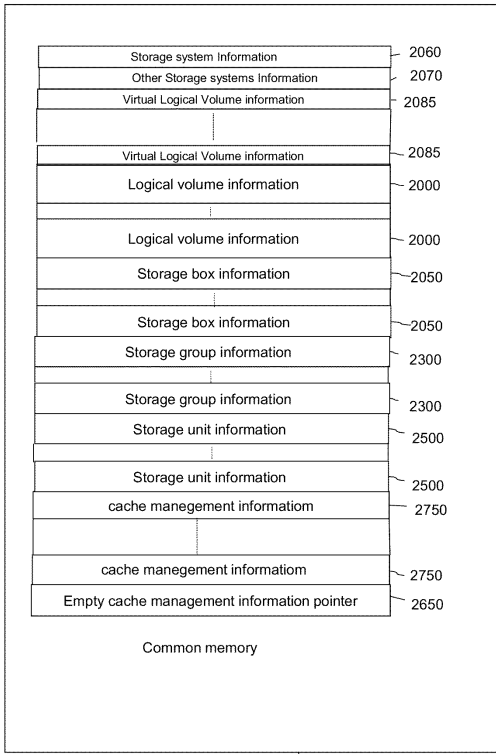


Fig. 5

【 図 6 】

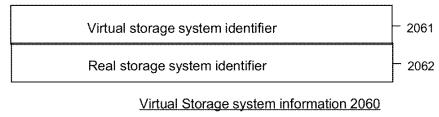


Fig.6

10

20

【 図 7 】

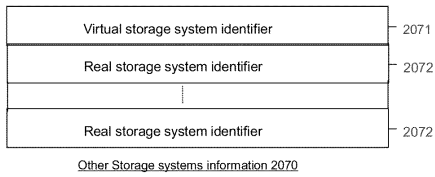
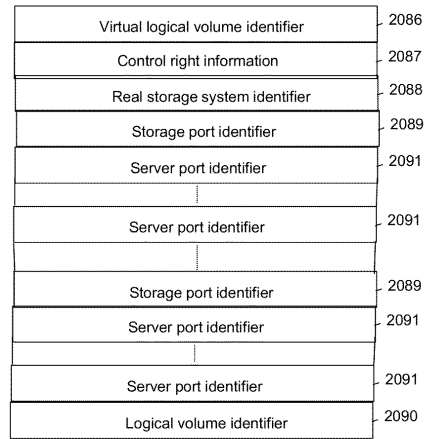


Fig.7

【 図 8 】



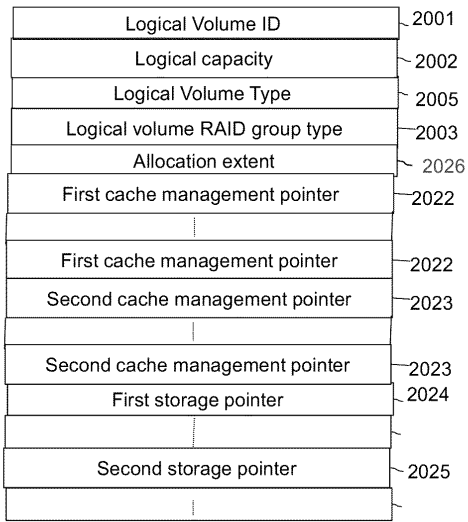
Virtual logical volume information 2085

Fig.8

30

40

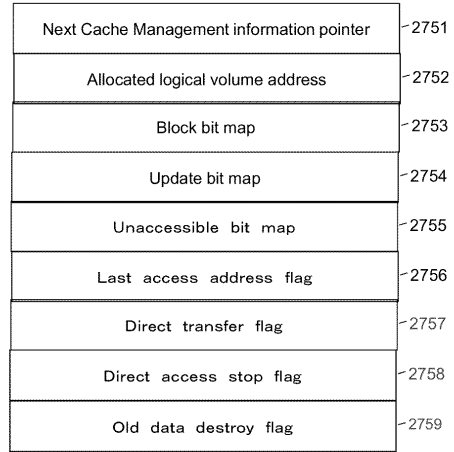
【 図 9 】



Logical volume information 2000

Fig. 9

【 図 1 0 】



Cache Management information 2750

Fig.10

10

20

【 図 1 1 】

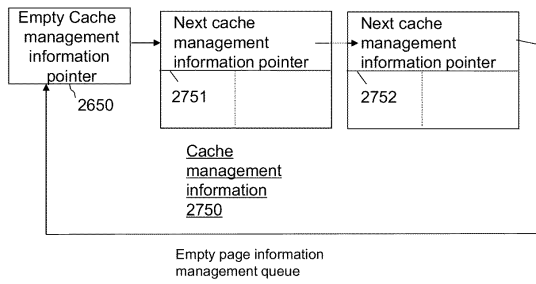
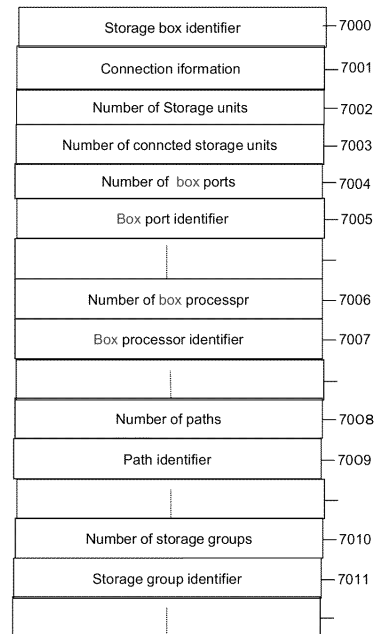


Fig. 11

【 図 1 2 】



Storage_box information 2050

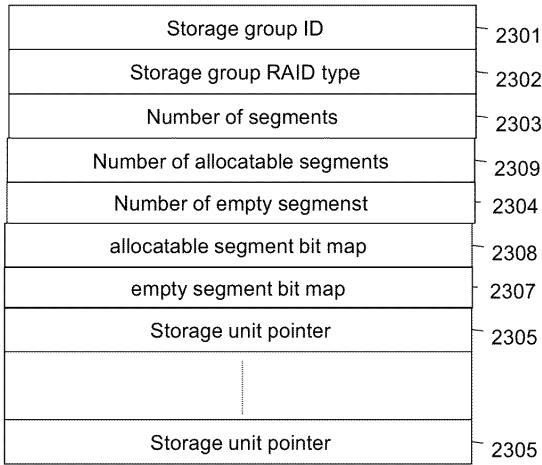
Fig.12

30

40

50

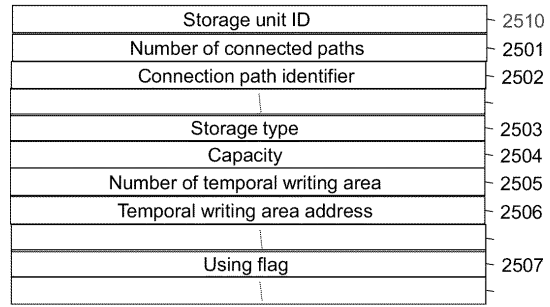
【 1 3 】



Storage group information 2300

Fig. 13

【 1 4 】



Storage unit information 2500

Fig. 14

10

20

【 1 5 】

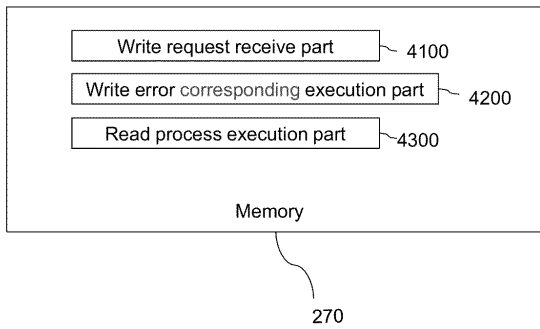


Fig. 15

【 1 6 A 】

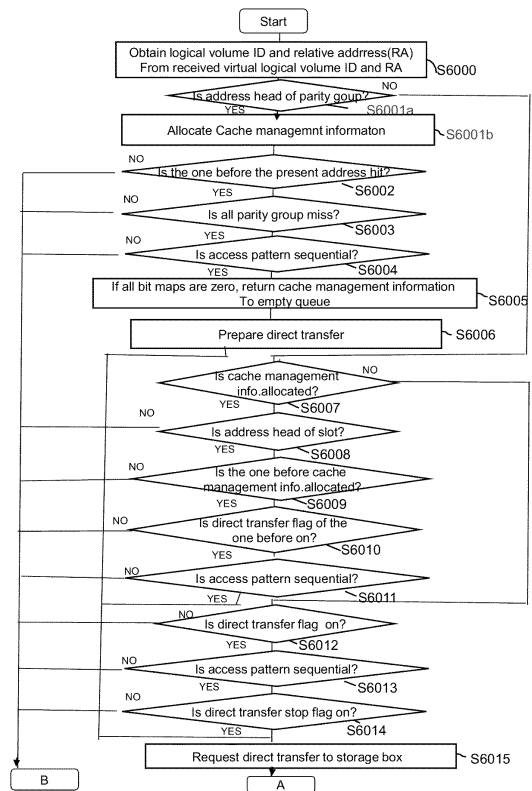


Fig. 16A

30

40

50

【 図 1 6 B 】

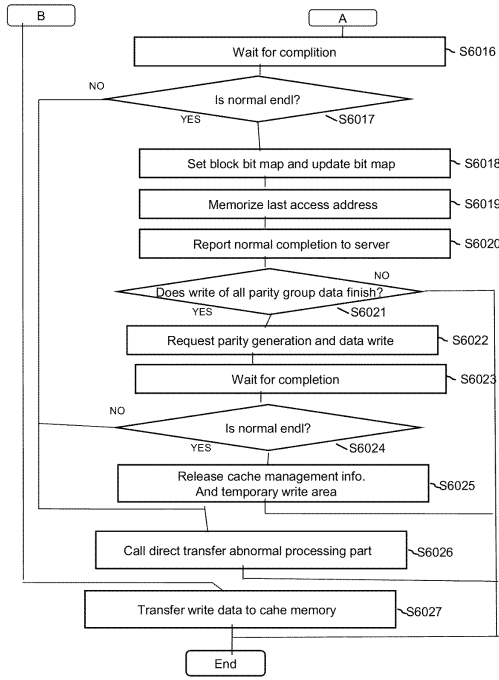


Fig. 16B

【 図 1 7 】

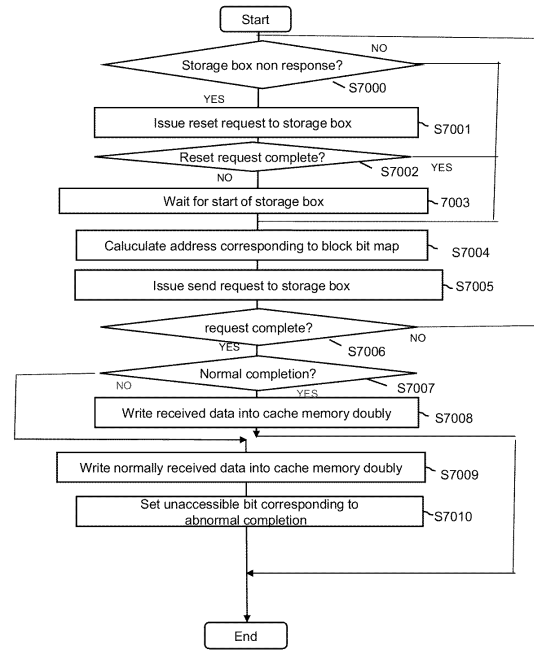


Fig. 17

【 図 1 8 】

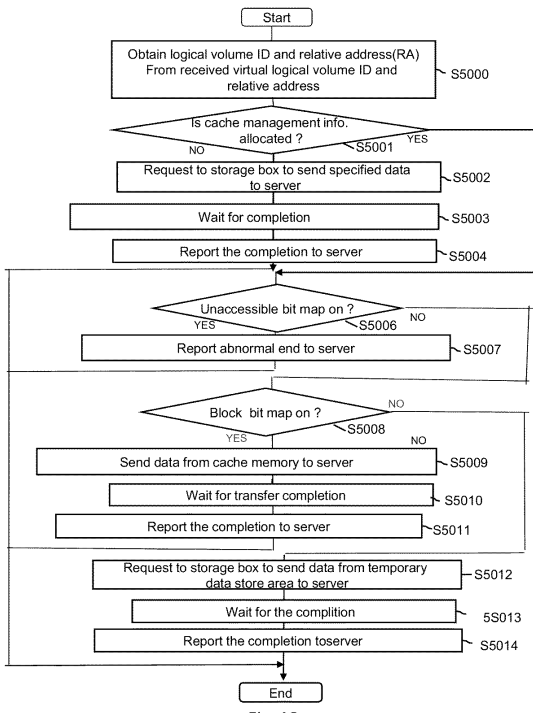


Fig. 18

【 図 1 9 】

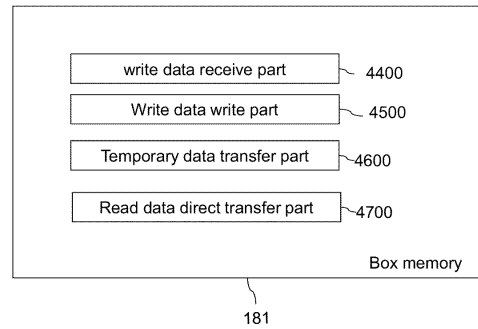


Fig. 19

10

20

30

40

50

【 2 0 】

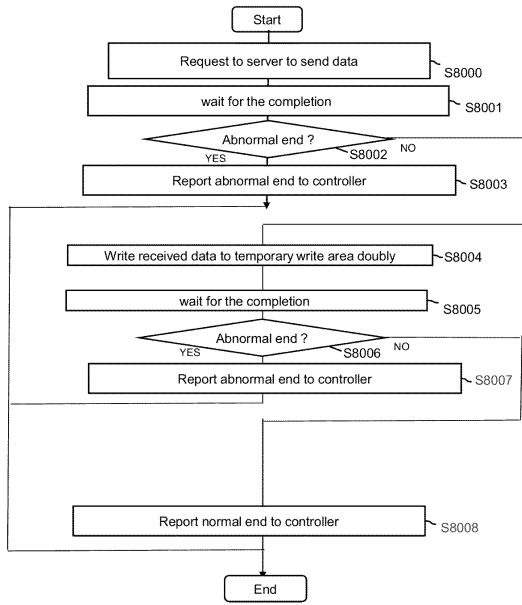


Fig. 20

【 2 1 】

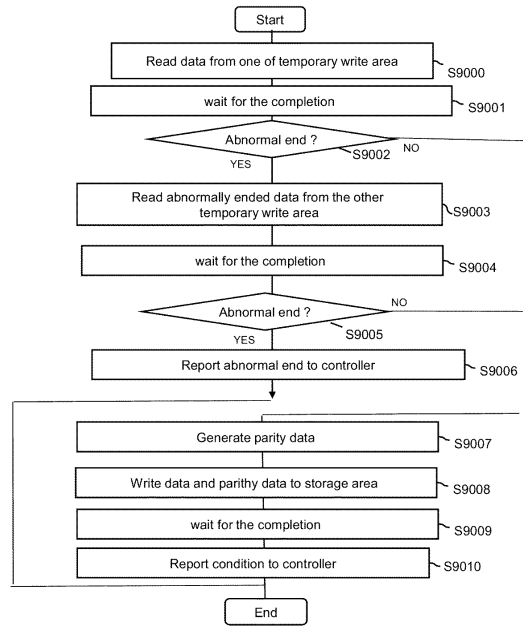


Fig. 21

【 2 2 】

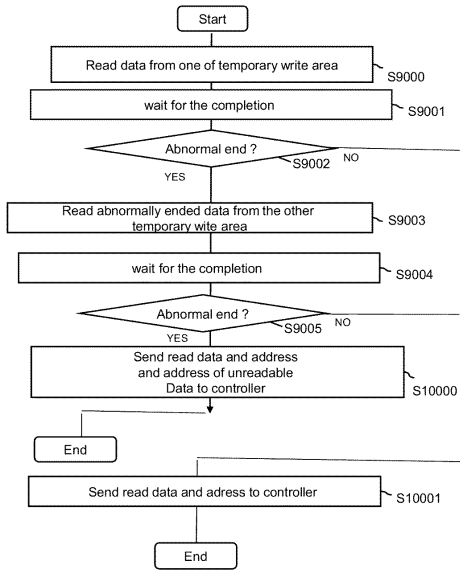


Fig. 22

【 2 3 】

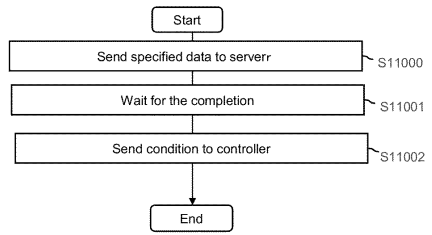


Fig. 23

10

20

30

40

50

【 2 4 】

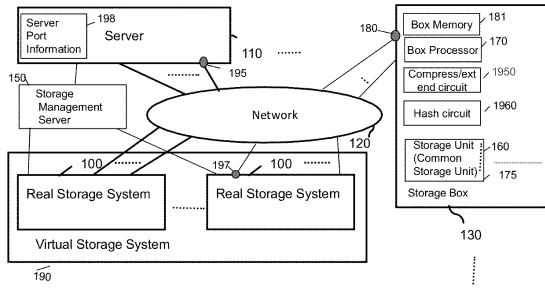


Fig. 24

【 2 5 】

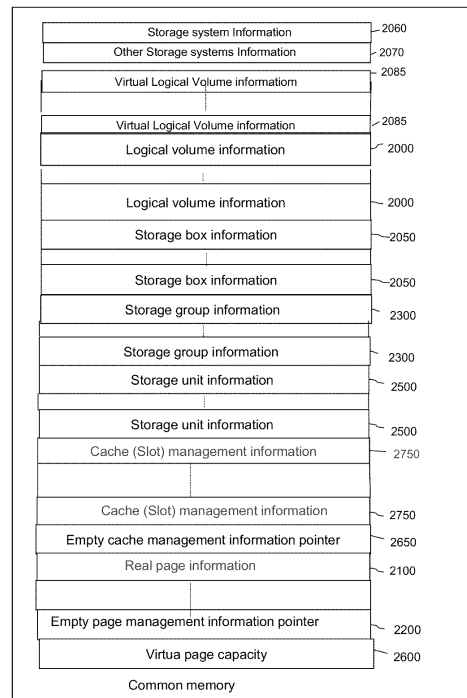
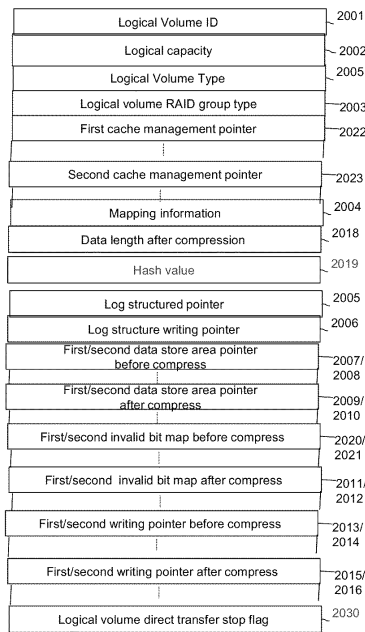


Fig.25 220

10

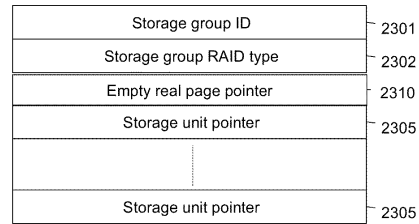
20

【 2 6 】



Logical volume information 2000
Fig. 26

【 2 7 】



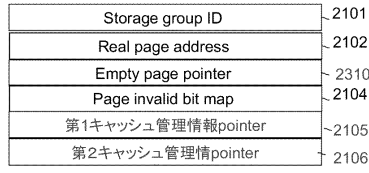
Storage group information 2300
Fig. 27

30

40

50

【 図 2 8 】



Real page information 2100

Fig. 28

【 図 2 9 】

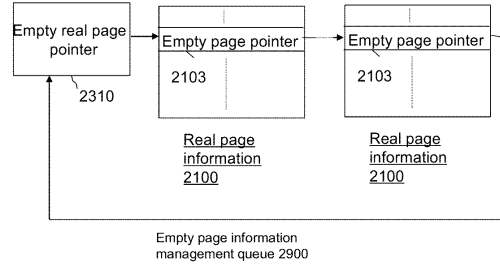


Fig.29

10

【 図 3 0 】

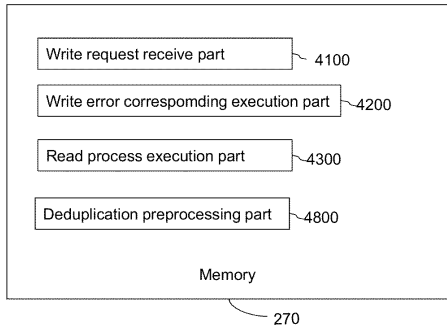


Fig. 30

【 図 3 1 A 】

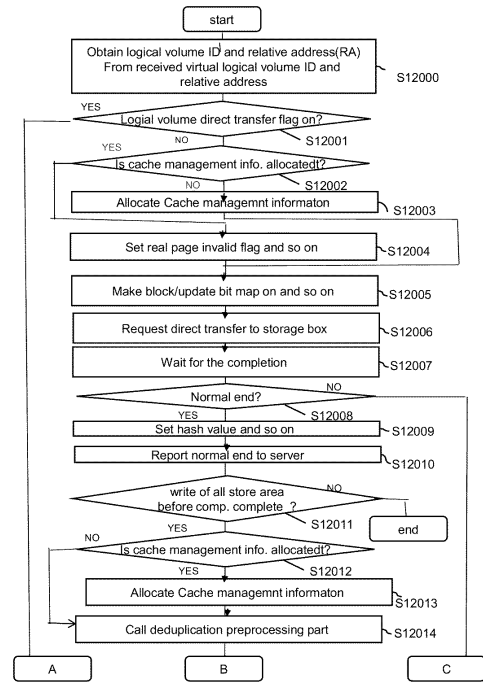


Fig. 31A

20

30

40

50

【 3 1 B 】

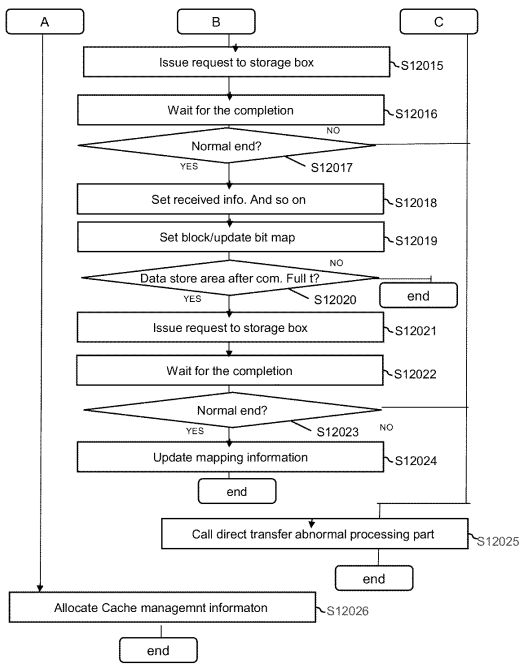


Fig. 31B

【 3 2 】

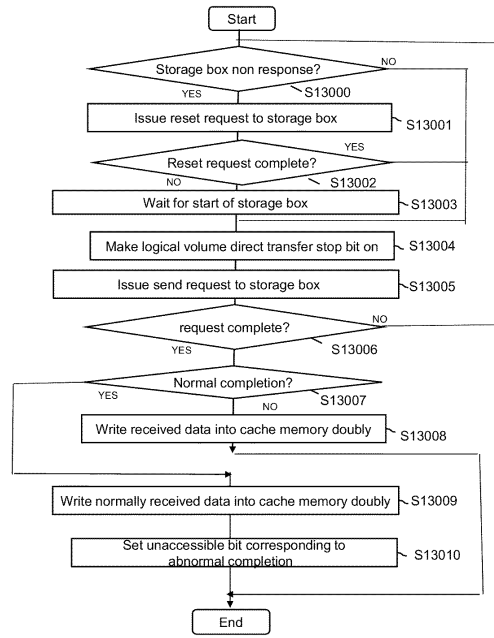


Fig.32

【 3 3 】

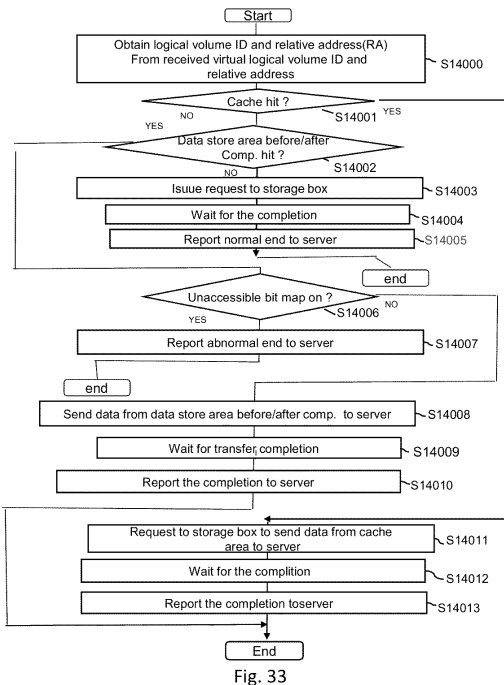


Fig. 33

【 3 4 】

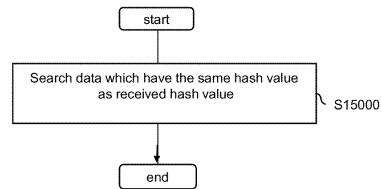


Fig.34

10

20

30

40

50

【 3 5 】

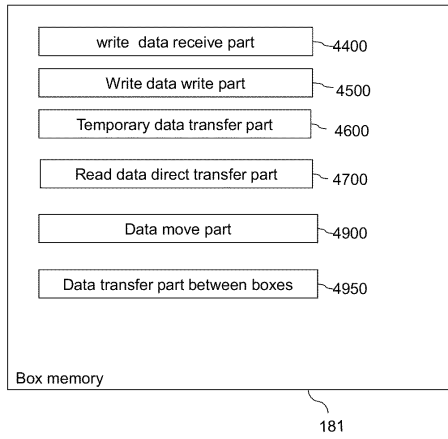


Fig. 35

【 3 6 】

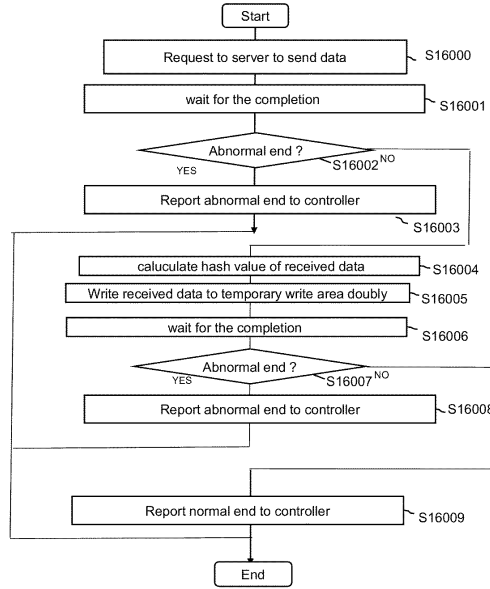


Fig. 36

【 3 7 】

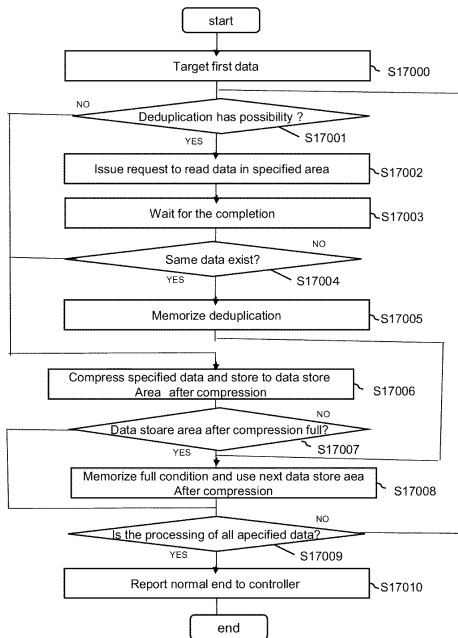


Fig. 37

【 3 8 】

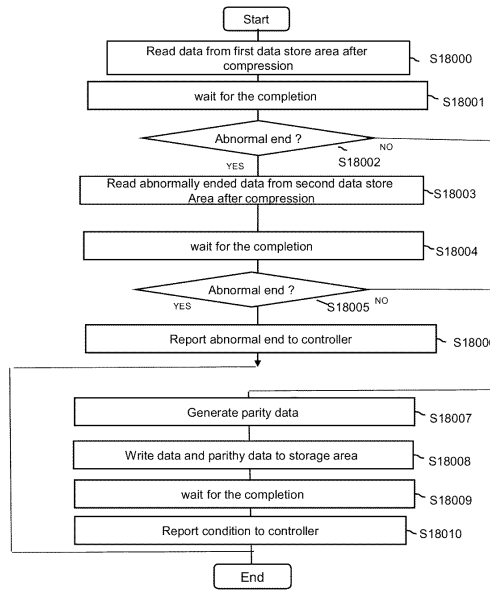


Fig. 38

10

20

30

40

50

【 3 9 】

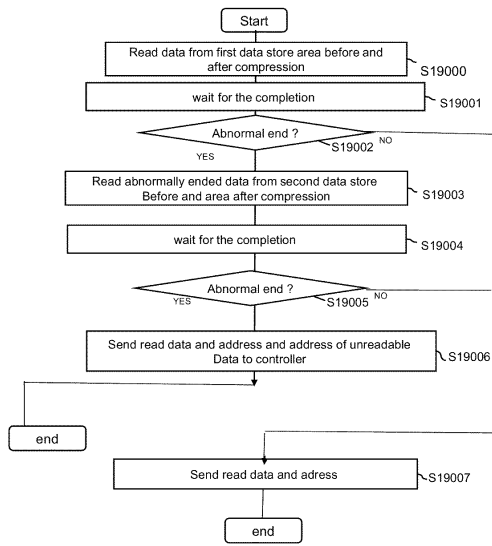


Fig. 39

【 4 0 】

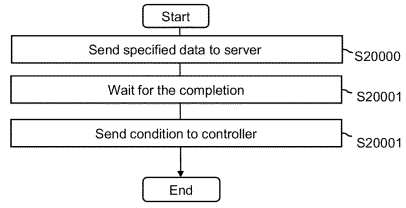


Fig. 40

10

【 4 1 】

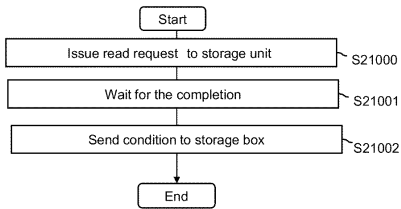


Fig. 41

20

30

40

50

フロントページの続き

東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内

審査官 田名網 忠雄

- (56)参考文献 特開2008-102967(JP,A)
特開2009-093225(JP,A)
米国特許出願公開第2002/0087751(US,A1)
米国特許出願公開第2015/0286438(US,A1)
再公表特許第2015/052798(JP,A1)
- (58)調査した分野 (Int.Cl., DB名)
G06F 3/06 - 3/08
G06F 13/10 - 13/14
G06F 11/16