

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7201775号
(P7201775)

(45)発行日 令和5年1月10日(2023.1.10)

(24)登録日 令和4年12月26日(2022.12.26)

(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 B
G 0 6 F	13/10 (2006.01)	G 0 6 F	3/06	3 0 2 A
G 0 6 F	11/20 (2006.01)	G 0 6 F	13/10	3 4 0 B
G 0 6 F	12/0866(2016.01)	G 0 6 F	11/20	6 8 9
G 0 6 F	12/0868(2016.01)	G 0 6 F	12/0866	1 0 0

請求項の数 6 (全23頁) 最終頁に続く

(21)出願番号 特願2021-182809(P2021-182809)
 (22)出願日 令和3年11月9日(2021.11.9)
 (62)分割の表示 特願2018-210691(P2018-210691)
)の分割
 原出願日 平成30年11月8日(2018.11.8)
 (65)公開番号 特開2022-10181(P2022-10181A)
 (43)公開日 令和4年1月14日(2022.1.14)
 審査請求日 令和3年11月9日(2021.11.9)

(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
 (74)代理人 110000279
 弁理士法人ウィルフォート国際特許事務所
 (72)発明者 鴨生 悠冬
 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 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
 (72)発明者 吉原 朋宏
 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

最終頁に続く

(54)【発明の名称】 ストレージシステム、データ管理方法、及びデータ管理プログラム

(57)【特許請求の範囲】

【請求項1】

複数のコントローラと、データを格納可能な記憶デバイスユニットとを有するストレージシステムであって、

前記コントローラは、

プロセッサ部と、

メモリと、を有し、

ライト要求にかかる新データが第1の前記コントローラの第1の前記メモリに格納された場合に、前記ライト要求にかかる新データにかかる第2のコントローラの第2のメモリ及び第3のコントローラの第3のメモリの旧データの状態がダーティである場合に、前記第1のコントローラは、ライト要求に対応する新データが格納された第1のメモリから、前記新データを第2のコントローラの第2のメモリに対して転送し、前記第2のメモリへの転送が完了した後に、前記新データを第3のコントローラの第3のメモリに対して転送することにより、前記第2のメモリと前記第3のメモリに別々に転送を行って前記旧データを上書きし、

前記ライト要求にかかる新データにかかる前記第2のメモリ及び前記第3のメモリの旧データの状態がダーティではない場合に、前記新データを前記第2のメモリと前記第3のメモリに並行して転送を行って前記旧データを上書きする

ストレージシステム。

【請求項2】

前記別々に転送を行う場合には、前記第 2 のメモリへの転送の成功を確認してから前記第 3 のメモリへの転送を行い、

前記転送時に障害が発生した場合には、前記障害により損傷が発生していない前記新データまたは旧データを用いて処理を行う

請求項 1 に記載のストレージシステム。

【請求項 3】

前記転送時に障害が発生した場合に、

前記第 2 のメモリまたは前記第 3 のメモリのいずれかの転送が正常に行われた場合には、正常に転送された新データを前記記憶デバイスユニットにデステージし、

前記第 2 のメモリまたは前記第 3 のメモリのいずれかにも転送が正常に行われていない場合には、前記第 2 のメモリまたは前記第 3 のメモリに格納済みの旧データを前記記憶デバイスユニットにデステージして、前記第 2 のメモリまたは前記第 3 のメモリのデータをダーティでなくする

請求項 2 に記載のストレージシステム。

【請求項 4】

前記第 1 のコントローラは、前記第 2 のコントローラの第 2 のメモリ及び前記第 3 のコントローラの第 3 のメモリに直接アクセス可能である

請求項 1 に記載のストレージシステム。

【請求項 5】

前記転送は、前記第 1 のコントローラの DMA (Direct Memory Access) 部が行う

請求項 4 に記載のストレージシステム。

【請求項 6】

複数のコントローラと、データを格納可能な記憶デバイスユニットとを有するストレージシステムによるデータ管理方法であって、

前記コントローラは、

プロセッサ部と、

メモリと、を有し、

ライト要求にかかる新データが第 1 の前記コントローラの第 1 の前記メモリに格納された場合に、前記ライト要求にかかる新データにかかる第 2 のコントローラの第 2 のメモリ及び第 3 のコントローラの第 3 のメモリの旧データの状態がダーティである場合に、前記第 1 のコントローラは、ライト要求に対応する新データが格納された第 1 のメモリから、前記新データを第 2 のコントローラの第 2 のメモリに対して転送し、前記第 2 のメモリへの転送が完了した後に、前記新データを第 3 のコントローラの第 3 のメモリに対して転送することにより、前記第 2 のメモリと前記第 3 のメモリに別々に転送を行って前記旧データを上書きし、

前記ライト要求にかかる新データにかかる前記第 2 のメモリ及び前記第 3 のメモリの旧データの状態がダーティではない場合に、前記新データを前記第 2 のメモリと前記第 3 のメモリに並行して転送を行って前記旧データを上書きする

データ管理方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージシステムにおけるデータを管理する技術に関する。

【背景技術】

【0002】

ストレージシステムでは、電源喪失などの障害からデータを保護するために、複数のストレージコントローラ間でデータを多重化（一般には二重化）している。また、ストレージシステムでは、複数の専用回路を用い、ライトデータをキャッシュ領域に同時に多重化することも行われている。

10

20

30

40

50

【 0 0 0 3 】

例えば、特許文献 1 には、第一キャッシュと F I F O バッファとにデータを格納した時点でホストにライト完了を送信し、その後、F I F O バッファから第二キャッシュにデータを送付することで、キャッシュを二重化するストレージシステムでの書込みを高速化する技術が開示されている。

【 0 0 0 4 】

一方、専用回路の開発コスト削減を目的として、特許文献 2 には、専用回路の処理を汎用コントローラでエミュレーションするストレージシステムにおいて、データの一貫性を保証する技術が開示されている。この技術では、コントローラ外部から受領するデータをバッファ領域に格納し、このコントローラがバッファ領域からキャッシュ領域に転送することで、I / O 処理中に障害が発生しても、データの破壊を防ぐことができるようにしている。

10

【 先行技術文献 】

【 特許文献 】

【 0 0 0 5 】

【 文献 】 特開 2 0 0 5 - 4 4 0 1 0 号公報
国際公開第 2 0 1 5 / 0 5 2 7 9 8 号

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

ストレージシステムの性能向上のために、汎用コントローラを多数搭載したストレージシステムが登場している。

20

【 0 0 0 7 】

このようなストレージシステムにおいて、メモリ容量の増加を抑える観点では、データの多重化を二多重とすることが望ましい。

【 0 0 0 8 】

例えば、このようなストレージシステムにおいてデータを二重化する場合においては、ホストからデータを受領するコントローラと、データの二重化先のコントローラとが全て異なる場合がある。このような場合に、特許文献 2 に記載の技術では、二重化先のコントローラのそれぞれにおいて、バッファ領域からキャッシュ領域へデータを転送する必要があり、コントローラのプロセッサへの処理負荷が掛かり、コントローラの性能が低下してしまう虞がある。

30

【 0 0 0 9 】

本発明は、上記事情に鑑みなされたものであり、その目的は、コントローラのプロセッサへの処理負荷を抑制しつつ、適切にデータの一貫性を確保することのできる技術を提供することにある。

【 課題を解決するための手段 】

【 0 0 1 0 】

上記目的を達成するため、一観点に係るストレージシステムは、複数のコントローラと、データを格納可能な記憶デバイスユニットとを有するストレージシステムであって、メモリに直接アクセス可能であるとともに、他のコントローラとの通信が可能な D M A (D i r e c t M e m o r y A c c e s s) 部を有し、コントローラは、プロセッサ部と、データを一時的に格納するバッファ領域と、データをキャッシュするキャッシュ領域とを有するメモリと、を有し、ライト要求にかかる新データがバッファ領域に格納された場合に、コントローラのプロセッサ部は、D M A 部を用いて、ライト要求に対応する新データが格納されたバッファ領域から、複数のコントローラのキャッシュ領域に対して、他のバッファ領域を介さずに順次転送させる。

40

【 発明の効果 】

【 0 0 1 1 】

本発明によれば、コントローラのプロセッサへの処理負荷を抑制しつつ、適切にデータ

50

の一貫性を確保することができる。

【図面の簡単な説明】

【0012】

【図1】図1は、実施例1に係るストレージシステムのライト処理の概要を説明する図である。

【図2】図2は、実施例1に係るストレージシステムの障害発生時のライト処理の概要を説明する図である。

【図3】図3は、実施例1に係る計算機システムの構成図である。

【図4】図4は、実施例1に係るコントローラ状態管理情報のデータ構造の一例を示す図である。

【図5】図5は、実施例1に係るキャッシュ状態管理情報のデータ構造の一例を示す図である。

【図6】図6は、実施例1に係る転送管理情報のデータ構造の一例を示す図である。

【図7】図7は、実施例1に係る転送状態管理情報のデータ構造の一例を示す図である。

【図8】図8は、実施例1に係る逐次転送依頼処理のフローチャートである。

【図9】図9は、実施例1に係る逐次転送完了待ち処理のフローチャートである。

【図10】図10は、実施例1に係る逐次転送処理のフローチャートである。

【図11】図11は、実施例1に係る障害対応処理のフローチャートである。

【図12】図12は、実施例2に係る逐次転送完了待ち処理のフローチャートである。

【図13】図13は、実施例2に係る障害対応処理のフローチャートである。

【発明を実施するための形態】

【0013】

以下、本発明の実施例を、図面を用いて説明する。ただし、本発明は以下に示す実施の形態の記載内容に限定して解釈されるものではない。本発明の思想ないし趣旨から逸脱しない範囲で、その具体的構成を変更し得ることは当業者であれば容易に理解される。

【0014】

以下に説明する発明の構成において、同一又は類似する構成又は機能には同一の符号を付し、重複する説明は省略する。

【0015】

本明細書等における「第1」、「第2」、「第3」等の表記は、構成要素を識別するために付するものであり、必ずしも、数又は順序を限定するものではない。

【0016】

図面等において示す各構成の位置、大きさ、形状、及び範囲等は、発明の理解を容易にするため、実際の位置、大きさ、形状、及び範囲等を表していない場合がある。したがって、本発明では、図面等に開示された位置、大きさ、形状、及び範囲等に限定されない。

【0017】

また、以下の説明における用語の意味は、下記の通りである。

(*)「PDEV」は、不揮発性の物理的な記憶デバイスの略である。複数のPDEVで複数のRAIDグループが構成されてよい。「RAID」は、Redundant Array of Independent (or Inexpensive) Disksの略である。RAIDグループはパリティグループと呼ばれてもよい。

(*)HCA(Host Channel Adaptor)は、CPUに指示され、コントローラ間の通信を行うデバイスである。HCAは、例えば、DMA(Direct Memory Access)部の一例であり、メモリに直接アクセスすることができる。

(*)プロセッサ部は、1以上のプロセッサを含む。少なくとも1つのプロセッサは、典型的には、CPU(Central Processing Unit)のようなマイクロプロセッサである。1以上のプロセッサの各々は、シングルコアでもよいしマルチコアでもよい。プロセッサは、処理の一部または全部を行うハードウェア回路を含んでもよい。

【実施例1】

【0018】

10

20

30

40

50

まず、実施例 1 に係る計算機システムについて説明する。

【0019】

図 1 は、実施例 1 に係るストレージシステムのライト処理の概要を説明する図である。図 1 は、二重化先のコントローラ 22 (# 1、# 2) それぞれのキャッシュ領域 243 (# 1、# 2) に対して逐次にデータを転送するライト処理の流れを示している。

【0020】

本実施例に係る計算機システム 100 のストレージシステム 2 は、複数のコントローラ 22 (コントローラ # 0、# 1、# 2) を備えている。複数のコントローラ 22 は、相互に接続されている。コントローラ 22 は、例えば、ストレージシステム専用のコントローラではなくて、汎用のコントローラである。コントローラ 22 は、FE - I / F 210 と、プロセッサ部の一例としての CPU 230 と、メモリ 240 とを有する。メモリ 240 は、バッファ領域 242 及びキャッシュ領域 243 を有するとともに、転送状態管理情報 247 を格納している。

10

【0021】

ストレージシステム 2 においては、各コントローラ 22 がホスト計算機 (ホストともいう) 1 からの I / O 要求を並列に処理できるよう、I / O 処理対象の空間 (例えば、論理ユニット : LU) ごとに処理担当 (この処理担当である権利を、オーナー権という) のコントローラ 22 を定めている。例えば、コントローラ # 1 が、LUN # 0 の LU に対するオーナー権を持っているとき、LUN # 0 の LU に対する I / O 要求は、このコントローラ # 1 の制御により処理される。

20

【0022】

ホスト計算機 1 (ホスト計算機 # 0、# 1、# 2) は、通信ネットワーク 11 を介して、ストレージシステム 2 と接続されている。ホスト計算機 1 は、例えば、それぞれ 1 つのコントローラ 22 と接続する。

【0023】

コントローラ 22 は、ホスト計算機 1 からのライト要求に従うライトデータを、記憶デバイスユニット 20 に書き込まず、複数のコントローラ 22 内のキャッシュ領域 243 に二重化して格納した後に、ホスト計算機 1 に対してライト処理の完了を通知する。これにより、高速なライト処理を実現することができる。

【0024】

コントローラ 22 は、ライト要求とは非同期にキャッシュ領域 243 内のライトデータを記憶デバイスユニット 20 に書き込む。既にキャッシュ領域 243 に二重化されたライトデータが格納されていて、まだ記憶デバイスユニット 20 に書き込まれていない状態 (ダーティ状態といい、このライトデータをダーティデータという) において、ホスト計算機 1 から同じ書き込み先への新たなライトデータを受領するときには、コントローラ 22 は、キャッシュ領域 243 内のライトデータの破壊を避けるために、新たなライトデータをバッファ領域 242 に格納する。その後、コントローラ 22 は、バッファ領域 242 内に格納した新たなライトデータを複数 (二重化の場合には 2 つ) のコントローラ内のそれぞれのキャッシュ領域 243 に逐次に転送することでライトデータの一貫性を保持する。転送状態管理情報 247 は、二重化先のコントローラ 22 のそれぞれのキャッシュ領域 243 に対してライトデータを逐次に転送する際の進捗状況 (転送状態) を管理する情報である。

30

40

【0025】

ここで、ライト要求を受領したコントローラ 22 以外のコントローラ 22 のバッファ領域 242 を介さずに、一のコントローラ 22 のキャッシュ領域 243 にライトデータを転送し、転送が完了した後、他のコントローラ 22 のキャッシュ領域 243 にライトデータを転送することを「逐次転送」という。

【0026】

ここで、コントローラ # 0 が、コントローラ # 1 がオーナー権を有する LU に対するライト要求をホスト計算機 # 0 から受け取った場合のライト処理について説明する。

50

【 0 0 2 7 】

コントローラ # 0 は、ホスト計算機 # 0 からライト要求を受信した場合、コントローラ # 0 の CPU # 0 は、ライト要求の対象となる LU のオーナー権を有するコントローラ # 1 の CPU # 1 にライト要求を転送する。

【 0 0 2 8 】

CPU # 1 は、バッファ領域 # 0 上にライトデータを格納する領域を確保させ、ライトデータに対応するキャッシュ領域 2 4 3 (本例では、キャッシュ領域 # 1、# 2) 上に格納されているデータの状態を確認する。本実施形態では、キャッシュ領域 2 4 3 のデータは、ダーティ状態であることとする。キャッシュ領域 2 4 3 のデータがダーティ状態 (ダーティデータ) であるので、CPU # 1 は、逐次転送が必要と判断する。

10

【 0 0 2 9 】

その後、コントローラ # 0 の CPU # 0 は、FE - IF # 0 を介してバッファ領域 # 0 に確保された領域にライトデータを格納する (ステップ S 1)。

【 0 0 3 0 】

次いで、CPU # 1 は、HCA # 0 に、バッファ領域 # 0 からキャッシュ領域 # 2 へライトデータをコピー (転送) し、その後、キャッシュ領域 # 1 にライトデータをコピー (転送) すること (逐次転送) を依頼する (ステップ S 2)。

【 0 0 3 1 】

HCA # 0 は、HCA # 2 を介して、バッファ領域 # 0 からキャッシュ領域 # 2 へライトデータをコピーする (以降、第一転送という) (ステップ S 3)。この際、HCA # 0 は、データコピー時にデータに付与された保証コードを確認する。保証コードは、データの格納位置を示す情報 (VOL 番号や VOL のアドレス等) やデータの一貫性を確認する情報 (CRC (Cyclic Redundancy Check) 等) から構成されてもよい。

20

【 0 0 3 2 】

次に、HCA # 0 は、HCA # 2 を介して、転送状態管理情報 # 2 に、キャッシュ領域 # 2 へのライトデータの転送受領を格納させ、HCA # 2 に、転送状態管理情報 # 1 にキャッシュ領域 # 2 へのライトデータの転送完了を格納させることを依頼する (ステップ S 4)。依頼を受けた HCA # 2 は、HCA # 1 を介して、転送状態管理情報 # 1 にキャッシュ領域 # 2 へのライトデータの転送完了を格納させる (ステップ S 5)。

【 0 0 3 3 】

次に、HCA # 0 は、HCA # 1 を介して、バッファ領域 # 0 からキャッシュ領域 # 1 へライトデータをコピーする (以降、第二転送という) (ステップ S 6)。

30

【 0 0 3 4 】

次に、HCA # 0 は、HCA # 1 を介して、転送状態管理情報 # 1 にキャッシュ領域 # 1 へのライトデータの転送完了を格納させる (ステップ S 7)。

【 0 0 3 5 】

CPU # 1 は、転送状態管理情報 # 1 を参照し、ライトデータの二重化完了を確認する (ステップ S 8)。次に、CPU # 1 は、CPU # 0 及び FE - I / F # 0 を介して、ホスト計算機 # 0 にライト要求完了を報告する (ステップ S 9)。これにより、ホスト計算機 # 1 からのライトデータは、キャッシュ領域 # 1 と、キャッシュ領域 # 2 とに二重化して格納される。

40

【 0 0 3 6 】

なお、上記例では、キャッシュ領域 # 2、# 1 の順でライトデータを順次転送させていたが、キャッシュ領域 # 1、# 2 の順としてもよい。

【 0 0 3 7 】

図 2 は、実施例 1 に係るストレージシステムの障害発生時のライト処理の概要を説明する図である。図 2 は、図 1 に示すライト処理の途中に障害が発生した時のライト処理の概要を示している。

【 0 0 3 8 】

HCA # 0 が、CPU # 1 から逐次転送の依頼を受け、逐次転送を実施中に HCA # 2 5

50

0やHCA250を繋ぐネットワーク(図3のHCAネットワーク23)のパス等に障害が発生すると(ステップS11)、キャッシュ領域#1または#2内のダーティ状態のライトデータ(ダーティデータ)を破壊してしまう虞がある(ステップS12)。つまり、ダーティデータの一部分だけが新たなライトデータの一部分によって上書きされた別のデータとなってしまう虞がある。

【0039】

そこで、ライトデータを管理するCPU#1は、転送状態管理情報247を参照し、正常なダーティデータを保持するキャッシュ領域243を特定する(ステップS13)。その後、CPU#1は、特定したキャッシュ領域243内のダーティデータをデステージ(すなわち、記憶デバイスユニット20に転送)する(ステップS14)。さらに、CPU#1は、デステージを完了後に、キャッシュ領域#1及び#2内のダーティデータを破棄する。なお、以降の説明では、特に記載しない場合には、デステージの完了後に、そのデステージしたデータに対応するデータが格納されていた複数のキャッシュ領域243のデータを破棄するものとする。

10

【0040】

以上の処理により、正常なダーティデータを選択して記憶デバイスユニット20に書き込むことができ、ライトデータの一貫性を保証できる。

【0041】

次に、本実施例に係る計算機システムについて詳細に説明する。

【0042】

図3は、実施例1に係る計算機システムの構成図である。

20

【0043】

計算機システム100は、1以上のホスト計算機1と、ストレージシステム2とを備える。ホスト計算機1と、ストレージシステム2とは、ネットワーク11を介して接続されている。ネットワーク11は、例えば、SAN(Storage Area Network)である。

【0044】

ストレージシステム2は、複数(例えば、3台以上)のコントローラ22(コントローラ22#0, ..., #N)と、記憶デバイスユニット20とを有する。複数のコントローラ22は、HCAネットワーク23を介して相互に接続されている。ストレージシステム2の可用性を向上させるため、コントローラ22毎に専用の電源を用意し、それぞれのコントローラ22に対して、その専用の電源を用いて給電するようにしてもよい。

30

【0045】

コントローラ22は、通信インタフェースと、記憶デバイスと、それらに接続されたプロセッサとを有する。通信インタフェースは、例えば、FE-I/F(Front End Inter/Face)210、BE-I/F(Back End Inter/Face)220、及びHCA250である。記憶デバイスは、例えば、メモリ240である。プロセッサは、例えば、CPU(Central Processing Unit)230である。なお、図3においては、コントローラ22は、1つのメモリ240を備えている構成としているが、メモリ240を複数備えてもよい。

【0046】

FE-I/F210は、ホスト計算機1等のフロントエンドに存在する外部デバイスと通信するためのインタフェースデバイスである。BE-I/F220は、コントローラ22が記憶デバイスユニット20と通信するためのインタフェースデバイスである。HCA250は、各コントローラ22のメモリ240を操作するために他のHCA250と通信するためのインタフェースデバイスである。

40

【0047】

メモリ240は、例えば、RAM(Random Access Memory)であり、バッファ領域242と、キャッシュ領域243とを含む。また、メモリ240は、制御モジュール241、コントローラ状態管理情報244、キャッシュ状態管理情報245、転送管理情報246、及び転送状態管理情報247を記憶する。なお、メモリ240は、不揮発性メモリ

50

であっても、揮発性メモリであってもよい。

【 0 0 4 8 】

制御モジュール 2 4 1 は、CPU 2 3 0 に実行されることにより、ストレージシステム 2 全体を制御するためのモジュール（プログラム）である。より具体的には、制御モジュール 2 4 1 は、CPU 2 3 0 に実行されることにより、I/O 処理の制御等を行う。

【 0 0 4 9 】

バッファ領域 2 4 2 は、ホスト計算機 1 から受領したライトデータを一時的に格納する領域である。

【 0 0 5 0 】

キャッシュ領域 2 4 3 は、ホスト計算機 1 から記憶デバイスユニット 2 0 へ送信されるライトデータをキャッシュする領域である。キャッシュ領域 2 4 3 は、ダーティデータを格納することもあるので、バックアップ電源等により不揮発化されていてもよい。

10

【 0 0 5 1 】

コントローラ状態管理情報 2 4 4 は、コントローラ 2 2 が正常状態か、故障状態かを管理するための情報である。キャッシュ状態管理情報 2 4 5 は、二重化に使用されているキャッシュ領域 2 4 3 を有するコントローラ 2 2 とキャッシュの状態を管理するための情報である。転送管理情報 2 4 6 は、逐次転送で転送するライトデータを受信したコントローラ 2 2 と、転送状態管理情報 2 4 7 のエントリのアドレスを管理するための情報である。転送状態管理情報 2 4 7 は、逐次転送の進捗状況（転送状態）を管理するための情報である。コントローラ状態管理情報 2 4 4、キャッシュ状態管理情報 2 4 5、転送管理情報 2 4 6、及び転送状態管理情報 2 4 7 の詳細は、図 4 乃至図 7 を参照して後述する。

20

【 0 0 5 2 】

記憶デバイスユニット 2 0 は、複数の PDEV 2 0 0 を有する。PDEV 2 0 0 は、HDD (Hard Disk Drive) でよいが、他種の記憶デバイス（不揮発性の記憶デバイス）、例えば、SSD (Solid State Drive) のような FM (Flash Memory) デバイスでもよい。記憶デバイスユニット 2 0 は、異なる種類の PDEV 2 0 0 を有してよい。また、複数の同種の PDEV 2 0 0 で RAID グループが構成されてよい。RAID グループには、所定の RAID レベルに従いデータが格納される。コントローラ 2 2 がホスト計算機 1 から受信したライトデータに対しては、FE - I/F 2 1 0 によって保証コードが付与される。この保証コードが付与されたデータは、RAID グループに格納される。

30

【 0 0 5 3 】

HCA 2 5 0 は、CPU 2 3 0 から指示を受け、自コントローラ 2 2 のメモリ 2 4 0 に対する操作や、HCA ネットワーク 2 3 を経由して、他コントローラ 2 2 のメモリ 2 4 0 に対する操作を行う。

【 0 0 5 4 】

次に、コントローラ状態管理情報 2 4 4 を詳細に説明する。

【 0 0 5 5 】

図 4 は、実施例 1 に係るコントローラ状態管理情報のデータ構造の一例を示す図である。

【 0 0 5 6 】

コントローラ状態管理情報 2 4 4 は、コントローラ 2 2 ごとのエントリを格納する。コントローラ状態管理情報 2 3 3 のエントリは、コントローラ ID 4 0 1 及び状態 4 0 2 のフィールドを含む。コントローラ ID 4 0 1 には、エントリに対応するコントローラ 2 2 の識別子（コントローラ ID）が格納される。状態 4 0 2 には、エントリに対応するコントローラ 2 2 の動作状態が格納される。動作状態としては、正常、故障等がある。

40

【 0 0 5 7 】

次に、キャッシュ状態管理情報 2 4 5 を詳細に説明する。

【 0 0 5 8 】

図 5 は、実施例 1 に係るキャッシュ状態管理情報のデータ構造の一例を示す図である。

【 0 0 5 9 】

キャッシュ状態管理情報 2 4 5 は、データアドレス毎のエントリを格納する。キャッシュ

50

ユ状態管理情報 2 4 5 のエントリは、データアドレス 5 0 1、第一転送先コントローラ ID 5 0 2、第二転送先コントローラ ID 5 0 3、及びキャッシュ状態 5 0 4 のフィールドを含む。

【 0 0 6 0 】

データアドレス 5 0 1 には、エントリに対応するストレージシステム 2 内のユーザデータの格納位置を示す値（データアドレス）が格納される。

【 0 0 6 1 】

第一転送先コントローラ ID 5 0 2 には、エントリに対応するデータアドレスのデータが二重化されてキャッシュされている、第一転送の転送先のキャッシュ領域 2 4 3 を有するコントローラ 2 2（転送先コントローラの一例）の識別子（コントローラ ID：第一転送先コントローラ ID）が格納される。

10

【 0 0 6 2 】

第二転送先コントローラ ID 5 0 3 には、エントリに対応するデータアドレスのデータが二重化されてキャッシュされている、第二転送の転送先のキャッシュ領域 2 4 3 を有するコントローラ 2 2（担当コントローラの一例）の識別子（コントローラ ID：第二転送先コントローラ ID）が格納される。本実施形態では、第二転送先コントローラ ID 5 0 3 には、エントリに対応するデータアドレスのデータが属する論理ユニットのオーナー権を有するコントローラ（オーナーコントローラ）2 2 のコントローラ ID が格納される。

【 0 0 6 3 】

キャッシュ状態 5 0 4 には、エントリに対するデータアドレスのデータのキャッシュの状態を示す情報が格納される。キャッシュの状態としては、記憶デバイスユニット 2 0 にデステージされていないことを示すダーティと、デステージされていることを示すクリーンとがある。

20

【 0 0 6 4 】

次に、転送管理情報 2 4 6 を詳細に説明する。

【 0 0 6 5 】

図 6 は、実施例 1 に係る転送管理情報のデータ構造の一例を示す図である。

【 0 0 6 6 】

転送管理情報 2 4 6 は、データアドレス毎のエントリを格納する。転送管理情報 2 4 6 のエントリは、データアドレス 6 0 1、コントローラ ID 6 0 2、及び転送状態管理情報アドレス 6 0 3 のフィールドを含む。データアドレス 6 0 1 には、エントリに対応するストレージシステム 2 内のユーザデータの格納位置（記憶空間）を示す値（データアドレス）が格納される。コントローラ ID 6 0 2 には、エントリに対応するデータアドレスのライトデータをホスト 1 から受信したコントローラ（受信コントローラ）2 2 の識別情報（コントローラ ID）が格納される。転送状態管理情報アドレス 6 0 3 には、エントリに対応するデータアドレスの転送状態管理情報 2 4 7 における対応するエントリの格納場所を示す値（アドレス）が格納される。

30

【 0 0 6 7 】

次に、転送状態管理情報 2 4 7 を詳細に説明する。

【 0 0 6 8 】

図 7 は、実施例 1 に係る転送状態管理情報のデータ構造の一例を示す図である。

40

【 0 0 6 9 】

転送状態管理情報 2 4 7 は、データアドレス毎のエントリを格納する。転送状態管理情報 2 4 7 のエントリは、データアドレス 7 0 1、第一転送データ受領済フラグ 7 0 2、第一転送完了フラグ 7 0 3、及び第二転送完了フラグ 7 0 4 のフィールドを含む。

【 0 0 7 0 】

データアドレス 7 0 1 は、エントリに対応するストレージシステム 2 内のユーザデータの格納位置を示す値（データアドレス）が格納される。第一転送データ受領済フラグ 7 0 2 には、H C A 2 5 0 によって、データアドレス 7 0 1 のデータアドレスに対応するライトデータの第一転送のデータが受領されたか否かを示す値（受領済みフラグ）が格納され

50

る。受領済みフラグは、受領された場合には、「1」が設定され、受領されていない場合には、「0」が設定される。第一転送完了フラグ703には、HCA250によって、データアドレス701に対応するデータアドレスのライトデータの第一転送が完了したか否かを示す値（第一転送完了フラグ）が格納される。第一転送完了フラグは、第一転送が完了された場合には、「1」が設定され、第一転送が完了されていない場合には、「0」が設定される。第二転送完了フラグ704には、HCA250によって、データアドレス701に対応するデータアドレスのライトデータの第二転送が完了したか否かを示す値（第二転送完了フラグ）が格納される。第二転送完了フラグは、第二転送が完了された場合には、「1」が設定され、第二転送が完了されていない場合には、「0」が設定される。

【0071】

次に、実施例1に係る計算機システムによる処理動作について説明する。

【0072】

まず、逐次転送依頼処理について説明する。

【0073】

図8は、実施例1に係る逐次転送依頼処理のフローチャートである。

【0074】

逐次転送依頼処理は、ライト要求に対応するライトデータ（新データ）が対象とする記憶デバイスユニット20における論理ユニット（記憶空間）のオーナー権を有するコントローラ22（オーナーコントローラ22という。：担当コントローラの一例）がライト要求を受信した場合に実行される。ここで、オーナーコントローラ22にライト要求が送信される場合としては、ホスト計算機1から直接オーナーコントローラ22に送られる場合と、ライト要求に対応するライトデータに対応するキャッシュ領域243を有さず、FE-I/F210を介してホスト計算機1からライトデータを受領したコントローラ22（FEコントローラ22という。受信コントローラの一例）からオーナーコントローラ22に転送される場合と、がある。

【0075】

本例では、FEコントローラ22からオーナーコントローラ22にライト要求が転送された場合を例に説明する。

【0076】

オーナーコントローラ22は、ライト要求を受信する（ステップS101）。次いで、オーナーコントローラ22は、キャッシュ状態管理情報245を参照し、ライト要求のデータアドレスに対応するエントリのキャッシュ状態504からキャッシュ状態を取得し（ステップS102）、キャッシュ状態がダーティであるか否かを判定する（ステップS103）。

【0077】

この結果、ダーティでないと判定された場合（ステップS103：NO）には、キャッシュ領域243のデータ（旧データ）が既に記憶デバイスユニット20に格納されていることを示すので、オーナーコントローラ22は、ライトデータを2つのコントローラ22のキャッシュ領域243に同時に（並行して）転送し、処理を終了する（S106）。

【0078】

一方、ダーティであると判定された場合（ステップS103：YES）には、オーナーコントローラ22は、ライト要求のデータアドレスに対応する転送状態管理情報247のエントリの格納先を示す値（転送状態管理情報アドレス）を取得し、転送管理情報246にエントリを追加する。オーナーコントローラ22は、追加したエントリのデータアドレス601、コントローラID602、及び転送状態管理情報アドレス603に、それぞれ、ライトデータのデータアドレス、FEコントローラ22のコントローラID、及び転送状態管理情報247のエントリの転送状態管理情報アドレスを設定する（ステップS104）。

【0079】

次に、オーナーコントローラ22は、ライトデータの逐次転送をFEコントローラ22内のHCA250に依頼し（ステップS105）、次の処理（図9の逐次転送完了待ち処理

10

20

30

40

50

)を実行する(L0)。なお、FEコントローラ22のHCA250への依頼は、自コントローラ22のHCA250を経由して通知してもよい。

【0080】

次に、逐次転送完了待ち処理について説明する。

【0081】

図9は、実施例1に係る逐次転送完了待ち処理のフローチャートである。

【0082】

オーナコントローラ22は、逐次転送が完了しているか否かを判定する(ステップS201)。すなわち、オーナコントローラ22は、転送状態管理情報247を参照し、ライトデータのデータアドレスに対応するエントリ、すなわち、データアドレス701の値がライトデータのデータアドレスであるエントリにおける第一転送完了フラグ703及び第二転送完了フラグ704のフラグが立っているか否か、すなわち、フラグの値が“1”であるか否かを判定する。なお、本ステップの処理は、一定の周期で行ってもよい。

10

【0083】

この結果、逐次転送が完了していると判定した場合(ステップS201: YES)には、オーナコントローラ22は、FEコントローラ22を経由してホスト1にライト処理が終了したことを意味するGood応答を送信し(ステップS202)、処理を終了する。一方、逐次転送が完了していないと判定した場合(ステップS201: NO)には、オーナコントローラ22は、コントローラ状態管理情報244から他のコントローラ22の状態を取得し、状態が故障であるコントローラID(故障コントローラID)を特定する(ステップS203)。

20

【0084】

次いで、オーナコントローラ22は、第一転送の転送先コントローラ22(第一転送先コントローラ22)が故障しているか否かを判定する(ステップS204)。具体的には、オーナコントローラ22は、ステップS203で特定した故障コントローラIDに、データアドレスに対応するキャッシュ状態管理情報245のエントリにおけるデータアドレス501に格納された第一転送先コントローラID502の第一転送先コントローラIDと一致するものが存在するか否かにより、第一転送先コントローラ22が故障しているか否かを判定する。

【0085】

この結果、第一転送先コントローラ22が故障していると判定した場合(ステップS204: YES)には、オーナコントローラ22は、第二転送が完了しているか否かを判定する(ステップS205)。すなわち、オーナコントローラ22は、転送状態管理情報247を参照し、データアドレスに対応するエントリの第二転送完了フラグ704のフラグが立っているか否かを判定する。

30

【0086】

この結果、第二転送が完了していると判定した場合(ステップS205: YES)には、第二転送により、第二転送先コントローラ(オーナコントローラ22)のキャッシュ領域243に対してライトデータが格納されていることを意味しているため、オーナコントローラ22は、オーナコントローラ22のキャッシュ領域243に格納されているライトデータ(保証データ)をデステージ(記憶デバイスユニット20に転送)する(ステップS207)。次に、オーナコントローラ22は、FEコントローラ22を経由してホスト1に失敗応答を送信し、処理を終了する(ステップS211)。

40

【0087】

ここで、第二転送が完了している場合(ステップS205: YES)には、第二転送の転送先のコントローラ22(第二転送先コントローラ22、オーナコントローラ)のキャッシュ領域243のライトデータは壊れていないことを示しているため、キャッシュ領域243のライトデータをデステージすることで、データの一貫性を保証できる。

【0088】

一方、第二転送が完了していないと判定した場合(ステップS205: NO)には、オ

50

ーナコントローラ 2 2 は、第一転送が完了しているか否かを判定する（ステップ S 2 0 6）。すなわち、オーナコントローラ 2 2 は、転送状態管理情報 2 4 7 を参照し、データアドレスに対応するエントリの第一転送完了フラグ 7 0 3 のフラグが立っているか否かを判定する。

【 0 0 8 9 】

この結果、第一転送が完了していると判定した場合（ステップ S 2 0 6 : Y E S）には、オーナコントローラ 2 2 は、処理をステップ S 2 0 1 に進め、第二転送の完了を待つ。

【 0 0 9 0 】

一方、第一転送が完了していないと判定した場合（ステップ S 2 0 6 : N O）、オーナコントローラ 2 2 は、処理をステップ S 2 0 7 に進める。

【 0 0 9 1 】

ここで、第一転送が完了していない場合、第二転送の開始前に第一転送先コントローラ 2 2 が故障していることを意味しているため、第二転送先コントローラであるオーナコントローラ 2 2 のキャッシュ領域 2 4 3 のパーティデータは更新されておらず、このキャッシュ領域 2 4 3 のパーティデータ（保証データ）をデステージすることで、データの一貫性を保証できる。

【 0 0 9 2 】

一方、ステップ S 2 0 4 で、第一転送先コントローラ 2 2 が故障していないと判定した場合（S 2 0 4 : N O）には、オーナコントローラ 2 2 は、F E コントローラ 2 2 が故障しているか否かを判定する（ステップ S 2 0 8）。すなわち、オーナコントローラ 2 2 は、転送管理情報 2 4 6 を参照し、データアドレスに対応するエントリのコントローラ I D 6 0 2 のコントローラ I D を取得し、このコントローラ I D と一致するものがステップ S 2 0 3 で特定した故障コントローラ I D に存在するか否かにより、F E コントローラ 2 2 が故障しているか否かを判定する。

【 0 0 9 3 】

この結果、F E コントローラ 2 2 が故障していないと判定した場合（ステップ S 2 0 8 : N O）には、オーナコントローラ 2 2 は、処理をステップ S 2 0 1 に戻し、逐次転送の完了を待つ。

【 0 0 9 4 】

一方、F E コントローラ 2 2 が故障していると判定した場合（ステップ S 2 0 8 : Y E S）には、オーナコントローラ 2 2 は、第一転送が完了しているか否かを判定する（ステップ S 2 0 9）。すなわち、オーナコントローラ 2 2 は、転送状態管理情報 2 4 7 を参照し、データアドレスに対応するエントリの第一転送完了フラグ 7 0 3 のフラグが立っているか否かを判定する。なお、第一転送完了フラグ 7 0 3 のフラグで判定する代わりに、第一転送データ受領済フラグ 7 0 2 のフラグが立っているか否かを判定してもよい。また、ステップ S 2 0 9 の前に、第二転送が完了しているか否かを判定し、第二転送が完了している場合、ホスト 1 に G o o d 応答を送信し、処理を終了してもよい。

【 0 0 9 5 】

この結果、第一転送が完了していると判定した場合（ステップ S 2 0 9 : Y E S）、オーナコントローラ 2 2 は、第一転送先のコントローラ 2 2 にそのコントローラ 2 2 のキャッシュ領域 2 4 3 からのデータ（保証データ）のデステージを依頼し（ステップ S 2 1 0）、処理をステップ 2 1 1 に進める。なお、第一転送先のコントローラ 2 2 は、依頼に対応して、キャッシュ領域 2 4 3 のデータをデステージすることとなる。ここで、第一転送が完了している場合、第一転送先コントローラ 2 2 のキャッシュ領域 2 4 3 のライトデータは壊れていないため、このキャッシュ領域 2 4 3 のライトデータがデステージされることによりデータの一貫性を保証できる。

【 0 0 9 6 】

一方、第一転送が完了していないと判定した場合（ステップ S 2 0 9 : N O）には、オーナコントローラ 2 2 は、処理をステップ S 2 0 7 に進め、自コントローラ 2 2 のキャッシュ領域 2 4 3 のライトデータ（保証データ）をデステージする。ここで、第一転送が完

10

20

30

40

50

了していない場合、第二転送の開始前にFEコントローラ22が故障しているため、オーナーコントローラ22のキャッシュ領域243のパーティデータは更新されておらず、このキャッシュ領域243のパーティデータをデステージすることでデータの一貫性を保証できる。

【0097】

次に、逐次転送処理について説明する。

【0098】

図10は、実施例1に係る逐次転送処理のフローチャートである。

【0099】

FEコントローラ22（具体的には、FEコントローラ22のHCA250）は、オーナーコントローラ22から送信された逐次転送依頼を受信し、逐次転送依頼からライトデータのデータアドレスを取得する（ステップS301）。次に、FEコントローラ22は、キャッシュ状態管理情報245を参照し、データアドレスに対応するエントリの第一転送先コントローラID502及び第二転送先コントローラID503から、第一転送先コントローラID及び第二転送先コントローラIDを取得する（ステップS302）。

10

【0100】

次に、FEコントローラ22のHCA250は、第一転送を実行する（ステップS303）。具体的には、FEコントローラ22のHCA250は、ライトデータをバッファ領域242から取り出し、第一転送先コントローラ22のHCA250を介して、ライトデータをキャッシュ領域243に転送する（ステップS303）。この際、ライトデータは、第一転送先コントローラ22のバッファ領域240を経由することなく、また、第一転送先コントローラ22のCPU230の関与なしに、キャッシュ領域243に転送される。

20

【0101】

次に、FEコントローラ22のHCA250は、ライトデータの転送が成功したか否かを判定する（ステップS304）。

【0102】

この結果、転送が失敗したと判定した場合（ステップS304：NO）には、FEコントローラ22のHCA250は、逐次転送処理を終了する。

【0103】

一方、転送が成功したと判定した場合（ステップS304：YES）には、FEコントローラ22のHCA250は、第一転送先コントローラ22のメモリ240内に存在する転送状態管理情報247のデータアドレスに対応するエントリの第一転送データ受領済フラグ702のフラグを立てる、すなわち、フラグを1に設定する（ステップS305）。

30

【0104】

次に、FEコントローラ22のHCA250は、第一転送先コントローラ22のHCA250に第二転送先コントローラ22のメモリ240内に存在する転送状態管理情報247のデータアドレスに対応するエントリの第一転送完了フラグ703のフラグを立てることを指示する（ステップS306）。

【0105】

次に、FEコントローラ22のHCA250は、第二転送を実行する（ステップS307）。具体的には、FEコントローラ22のHCA250は、ライトデータをバッファ領域242から取り出し、第二転送先コントローラ22のHCA250を介して、ライトデータをキャッシュ領域243に転送する（ステップS307）。

40

【0106】

次に、FEコントローラ22のHCA250は、ライトデータの転送が成功したか否かを判定する（ステップS308）。

【0107】

この結果、転送が失敗したと判定した場合（ステップS308：NO）には、FEコントローラ22は、逐次転送処理を終了する。

【0108】

50

一方、転送が成功したと判定した場合（ステップS308：YES）には、FEコントローラ22のHCA250は、第二転送先コントローラ22のメモリ240内に存在する転送状態管理情報247のデータアドレスに対応するエントリの第二転送完了フラグ704のフラグを立てる、すなわち、フラグを1に設定し（ステップS309）、処理を終了する。

【0109】

次に、障害対応処理について説明する。

【0110】

図11は、実施例1に係る障害対応処理のフローチャートである。障害対応処理は、オーナーコントローラ22以外のコントローラ22により実行される処理である。障害対応処理は、一定時間ごとに一度実行されてもよく、コントローラ22により障害が検知された場合に実行されてもよい。

10

【0111】

コントローラ22は、コントローラ状態管理情報244から他のコントローラ22の状態を取得し、状態が故障であるコントローラID401（故障コントローラID）を特定する（ステップS401）。

【0112】

次に、コントローラ22は、オーナーコントローラ22（第二転送先コントローラ22）が故障しているか否かを判定する（ステップS402）。すなわち、オーナーコントローラ22は、障害コントローラIDに、データアドレスに対応するキャッシュ状態管理情報245のエントリにおける第二転送先コントローラID503のコントローラIDと一致するものが存在しているか否かにより、オーナーコントローラ22が故障しているか否かを判定する。

20

【0113】

この結果、オーナーコントローラ22が故障していないと判定した場合（ステップS402：NO）には、コントローラ22は、処理をステップS401に戻す。一方、オーナーコントローラ22が故障していると判定した場合（ステップS402：YES）には、コントローラ22は、自身（自コントローラ）が第一転送先コントローラであるか否かを判定する（ステップS403）。すなわち、コントローラ22は、自身のコントローラID（自コントローラID）と、データアドレスに対応するキャッシュ状態管理情報245のエントリにおける第一転送先コントローラID503のコントローラIDとが同一であるか判定する。なお、上記処理の代わりに、自コントローラがFEコントローラであるか否かを判定し、すなわち、自コントローラIDと、転送管理情報246のデータアドレスに対応するエントリのコントローラID602のコントローラIDとが同一であるか否かを判定し、自コントローラがFEコントローラである場合に、以降の処理を行ってもよい。

30

【0114】

ステップS403の判定の結果、自コントローラが第一転送先コントローラでないと判定した場合（ステップS403：NO）には、コントローラ22は、処理をステップS401に進める。

【0115】

一方、自コントローラが第一転送先コントローラであると判定した場合（ステップS403：YES）には、コントローラ22は、第一転送が完了しているか否かを判定する（ステップS404）。すなわち、コントローラ22は、転送管理情報246のデータアドレスに対応するエントリの転送状態管理情報アドレス603のアドレスを用いて、転送状態管理情報247のエントリを参照し、このエントリの第一転送データ受領済フラグ702のフラグが立っているか否かを判定する。

40

【0116】

この判定結果、第一転送が完了していないと判定した場合（ステップS404：NO）には、コントローラ22は、処理をステップS401に進め、第一転送の完了を待つ。

【0117】

50

一方、第一転送が完了していると判定した場合（ステップ S 4 0 4 : Y E S ）には、コントローラ 2 2 は、キャッシュ領域 2 4 3 のライトデータ（保証データ）をデステージし（ステップ S 4 0 5 ）、F E コントローラ 2 2 を経由してホスト 1 に失敗応答を送信し（ステップ S 4 0 6 ）、処理を終了する。ここで、第一転送が完了している場合（ステップ S 4 0 4 : Y E S ）、第一転送先コントローラ 2 2 のキャッシュ領域 2 4 3 のライトデータは壊れていないため、キャッシュ領域 2 4 3 のライトデータをデステージすることでデータの一貫性を保証できる。

【 0 1 1 8 】

以上説明したように、上記実施例に係る計算機システムでは、ライトデータの二重化の処理の進捗に合わせて、障害発生時に記憶デバイスユニット 2 0 に書き込むキャッシュ領域 2 4 3 を使い分けることで、ライトデータの一貫性を保証できる。

10

【実施例 2】

【 0 1 1 9 】

次に、実施例 2 に係る計算機システムについて説明する。

【 0 1 2 0 】

実施例 2 に係る計算機システムは、図 3 に示す実施例 1 に係る計算機システムにおいて、論理ユニットを担当するコントローラ 2 2 を特定のコントローラ 2 2 に限定しない、すなわち、論理ユニットのオーナー権を設定しないようにしたシステムである。この計算機システムにおいては、例えば、ホスト 1 からのライト要求を受信したコントローラ（受信コントローラ）が担当コントローラとなる。

20

【 0 1 2 1 】

実施例 2 に係る計算機システムでは、図 8 に示す逐次転送依頼処理を、ホスト 1 からライト要求を受信したコントローラ 2 2（F E コントローラ 2 2）が実行する。

【 0 1 2 2 】

次に、逐次転送完了待ち処理について説明する。

【 0 1 2 3 】

図 1 2 は、実施例 2 に係る逐次転送完了待ち処理のフローチャートである。

【 0 1 2 4 】

F E コントローラ 2 2 は、逐次転送が完了しているか否かを判定する（ステップ S 5 0 1 ）。すなわち、F E コントローラ 2 2 は、転送状態管理情報 2 4 7 を参照し、ライトデータのデータアドレスに対応するエン트리、すなわち、データアドレス 7 0 1 の値がライトデータのデータアドレスであるエントリにおける第一転送完了フラグ 7 0 3 及び第二転送完了フラグ 7 0 4 のフラグが立っているか否か、すなわち、フラグの値が“ 1 ”であるか否かを判定する。

30

【 0 1 2 5 】

この結果、逐次転送が完了していると判定した場合（ステップ S 5 0 1 : Y E S ）には、F E コントローラ 2 2 は、ホスト 1 に G o o d 応答を送信し（ステップ S 5 0 2 ）、処理を終了する。一方、逐次転送が完了していないと判定した場合（ステップ S 5 0 1 : N O ）には、F E コントローラ 2 2 は、コントローラ状態管理情報 2 4 4 から他のコントローラ 2 2 の状態を取得し、状態が故障であるコントローラ I D（故障コントローラ I D）を特定する（ステップ S 5 0 3 ）。

40

【 0 1 2 6 】

次いで、オーナーコントローラ 2 2 は、第一転送の転送先コントローラ 2 2（第一転送先コントローラ 2 2）が故障しているか否かを判定する（ステップ S 5 0 4 ）。

【 0 1 2 7 】

この結果、第一転送先コントローラ 2 2 が故障していると判定した場合（ステップ S 5 0 4 : Y E S ）には、F E コントローラ 2 2 は、第二転送が完了しているか否かを判定する（ステップ S 5 0 5 ）。すなわち、F E コントローラ 2 2 は、転送状態管理情報 2 4 7 を参照し、データアドレスに対応するエントリの第二転送完了フラグ 7 0 4 のフラグが立っているか否かを判定する。

50

【 0 1 2 8 】

この結果、第二転送が完了していると判定した場合（ステップ S 5 0 5 : Y E S ）には、F E コントローラ 2 2 は、第二転送先のコントローラ 2 2 のキャッシュ領域 2 4 3 に格納されているライトデータ（保証データ）のデステージを依頼する（ステップ S 5 0 7 ）。次に、F E コントローラ 2 2 は、ホスト 1 に失敗応答を送信し、処理を終了する（ステップ S 5 1 1 ）。

【 0 1 2 9 】

ここで、第二転送が完了している場合（ステップ S 5 0 5 : Y E S ）には、第二転送の転送先のコントローラ 2 2 （第二転送先コントローラ 2 2 ）のキャッシュ領域 2 4 3 のライトデータは壊れていないことを示しているため、キャッシュ領域 2 4 3 のライトデータをデステージすることで、データの一貫性を保証できる。

10

【 0 1 3 0 】

一方、第二転送が完了していないと判定した場合（ステップ S 5 0 5 : N O ）には、F E コントローラ 2 2 は、第一転送が完了しているか否かを判定する（ステップ S 5 0 6 ）。

【 0 1 3 1 】

この結果、第一転送が完了していると判定した場合（ステップ S 5 0 6 : Y E S ）には、F E コントローラ 2 2 は、処理をステップ S 5 0 1 に進め、第二転送の完了を待つ。

【 0 1 3 2 】

一方、第一転送が完了していないと判定した場合（ステップ S 5 0 6 : N O ）、F E コントローラ 2 2 は、処理をステップ S 5 0 7 に進める。

20

【 0 1 3 3 】

ここで、第一転送が完了していない場合、第二転送の開始前に第一転送先コントローラ 2 2 が故障していることを意味しているので、第二転送先コントローラのキャッシュ領域 2 4 3 のダーティデータは更新されておらず、このキャッシュ領域 2 4 3 のダーティデータをデステージすることで、データの一貫性を保証できる。

【 0 1 3 4 】

一方、ステップ S 5 0 4 で、第一転送先コントローラ 2 2 が故障していないと判定した場合（S 5 0 4 : N O ）には、F E コントローラ 2 2 は、第二転送先コントローラ 2 2 が故障しているか否かを判定する（ステップ S 5 0 8 ）。すなわち、F E コントローラ 2 2 は、キャッシュ状態管理情報 2 4 5 を参照し、データアドレスに対応するエントリの第二転送先コントローラ I D 5 0 3 のコントローラ I D を取得し、このコントローラ I D と一致するものがステップ S 5 0 3 で特定した故障コントローラ I D に存在するか否かにより、第二転送先コントローラ 2 2 が故障しているか否かを判定する。

30

【 0 1 3 5 】

この結果、第二転送先コントローラ 2 2 が故障していないと判定した場合（ステップ S 5 0 8 : N O ）には、F E コントローラ 2 2 は処理をステップ S 5 0 1 に戻し、逐次転送の完了を待つ。

【 0 1 3 6 】

一方、第二転送先コントローラ 2 2 が故障していると判定した場合（ステップ S 5 0 8 : Y E S ）には、F E コントローラ 2 2 は、第一転送が完了しているか否かを判定する（ステップ S 5 0 9 ）。

40

【 0 1 3 7 】

この結果、第一転送が完了していないと判定した場合（ステップ S 5 0 9 : N O ）には、F E コントローラ 2 2 は、処理をステップ S 5 0 1 に進め、第一転送が終わるのを待つ。

【 0 1 3 8 】

一方、第一転送が完了していると判定した場合（ステップ S 5 0 9 : Y E S ）、F E コントローラ 2 2 は、第一転送先のコントローラ 2 2 にそのコントローラ 2 2 のキャッシュ領域 2 4 3 からのデータ（保証データ）のデステージを依頼し（ステップ S 5 1 0 ）、処理をステップ S 5 1 1 に進める。ここで、第一転送が完了している場合、第一転送先コントローラ 2 2 のキャッシュ領域 2 4 3 のライトデータは壊れていないため、このキャッシ

50

ユ領域 2 4 3 のライトデータをデステージすることでデータの一貫性を保証できる。

【 0 1 3 9 】

次に、逐次転送処理について説明する。

【 0 1 4 0 】

実施例 2 に係る計算機システムの逐次転送処理は、図 1 0 に示す逐次転送処理とは、ステップ S 3 0 6 とステップ S 3 0 9 における処理内容が異なる。

【 0 1 4 1 】

実施例 2 に係る計算機システムにおいては、ステップ S 3 0 6 では、コントローラ 2 2 の H C A 2 5 0 は、自コントローラ 2 2 のメモリ 2 4 0 内に存在する、データアドレスに対応する転送状態管理情報 2 4 7 のエントリの第一転送完了フラグ 7 0 3 のフラグを立てる。また、ステップ S 3 0 9 では、コントローラ 2 2 の H C A 2 5 0 は、自コントローラ 2 2 のメモリ 2 4 0 内に存在する、データアドレスに対応する転送状態管理情報 2 4 7 のエントリの第二転送完了フラグ 7 0 4 のフラグを立てる。

10

【 0 1 4 2 】

次に、障害対応処理について説明する。

【 0 1 4 3 】

図 1 3 は、実施例 2 に係る障害対応処理のフローチャートである。

【 0 1 4 4 】

障害対応処理は、F E コントローラ 2 2 以外のコントローラ（他コントローラ）が実行する処理である。障害対応処理は、一定時間に一度実施してもよいし、他コントローラ 2 2 の障害検知時に実施してもよい。

20

【 0 1 4 5 】

コントローラ 2 2 は、コントローラ状態管理情報 2 4 4 から他のコントローラ 2 2 の状態を取得し、状態が故障であるコントローラ I D（故障コントローラ I D）を特定する（ステップ S 6 0 1）。次に、コントローラ 2 2 は、F E コントローラ 2 2 が故障しているか否かを判定する（ステップ S 6 0 2）。

【 0 1 4 6 】

この結果、F E コントローラ 2 2 が故障していないと判定した場合（ステップ S 6 0 2：N O）には、コントローラ 2 2 は、処理をステップ S 6 0 1 に戻す。一方、F E コントローラ 2 2 が故障していると判定した場合（ステップ S 6 0 2：Y E S）には、コントローラ 2 2 は、自身（自コントローラ）が第一転送先コントローラであるか否かを判定する（ステップ S 6 0 3）。なお、上記処理の代わりに、自コントローラが第二転送先コントローラであるか否かを判定し、以降の処理を行ってもよい。

30

【 0 1 4 7 】

ステップ S 6 0 3 の判定の結果、自コントローラ 2 2 が第一転送先コントローラでないと判定した場合（ステップ S 6 0 3：N O）には、コントローラ 2 2 は、処理をステップ S 6 0 1 に進める。

【 0 1 4 8 】

一方、自コントローラが第一転送先コントローラであると判定した場合（S 6 0 3：Y E S）には、コントローラ 2 2 は、第一転送が完了しているか否かを判定する（ステップ S 6 0 4）。

40

【 0 1 4 9 】

この判定結果、第一転送が完了していると判定した場合（ステップ S 6 0 4：Y E S）には、コントローラ 2 2 は、キャッシュ領域 2 4 3 のライトデータ（保証データ）をデステージし（ステップ S 6 0 6）、F E コントローラ 2 2 を経由してホスト 1 に失敗応答を送信し（ステップ S 6 0 7）、処理を終了する。ここで、第一転送が完了している場合（ステップ S 6 0 4：Y E S）、第一転送先コントローラ 2 2 のキャッシュ領域 2 4 3 のライトデータは壊れていないため、キャッシュ領域 2 4 3 のライトデータをデステージすることでデータの一貫性を保証できる。

【 0 1 5 0 】

50

一方、第一転送が完了していないと判定した場合（ステップS604：NO）には、コントローラ22は、第二転送先のコントローラ22にキャッシュ領域243に格納されているライトデータ（保証データ）のデステージを依頼し（ステップS605）、処理をステップS607に進める。ここで、第一転送が完了していない場合（ステップS604：NO）、第二転送の開始前に第一転送先コントローラ22が故障しているため、第二転送先コントローラ22のキャッシュ領域243のダーティデータは更新されておらず、このキャッシュ領域243のダーティデータをデステージすることでデータの一貫性を保証できる。

【0151】

以上説明したように、上記実施例に係る計算機システムでは、ライトデータの二重化の処理の進捗に合わせて、障害発生時に記憶デバイスユニット20に書き込むキャッシュ領域243を使い分けることで、ライトデータの一貫性を保証できる。

【0152】

なお、本発明は上記した実施例に限定されるものではなく、様々な変形例が含まれる。また、例えば、上記した実施例は本発明を分かりやすく説明するために構成を詳細に説明したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。また、各実施例の構成の一部について、他の構成に追加、削除、置換することが可能である。

【0153】

また、上記の各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。また、本発明は、実施例の機能を実現するソフトウェア（データ管理プログラム）のプログラムコードによって実現してもよい。この場合、プログラムコードを記録した記憶媒体をコンピュータに提供し、そのコンピュータが備えるプロセッサが記憶媒体に格納されたプログラムコードを読み出す。この場合、記憶媒体から読み出されたプログラムコード自体が前述した実施例の機能を実現することになり、そのプログラムコード自体、及びそれを記憶した記憶媒体は本発明を構成することになる。このようなプログラムコードを供給するための記憶媒体としては、例えば、フレキシブルディスク、CD-ROM、DVD-ROM、ハードディスク、SSD（Solid State Drive）、光ディスク、光磁気ディスク、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどがある。

【0154】

また、本実施例に記載の機能を実現するプログラムコードは、例えば、アセンブラ、C/C++、perl、Shell、PHP、Java（登録商標）等の広範囲のプログラム又はスクリプト言語で実装してもよい。

【0155】

さらに、実施例の機能を実現するソフトウェアのプログラムコードを、ネットワークを介して配信することによって、それをコンピュータのハードディスクやメモリ等の記憶手段又はCD-RW、CD-R等の記憶媒体に格納し、コンピュータが備えるプロセッサが当該記憶手段や当該記憶媒体に格納されたプログラムコードを読み出して実行するようにしてもよい。

【0156】

上記実施例において、制御線や情報線は、説明上必要と考えられるものを示しており、製品上必ずしも全ての制御線や情報線を示しているとは限らない。全ての構成が相互に接続されていてもよい。

【0157】

また、上記実施例では、複数のコントローラ22のキャッシュ領域243上での二重化ができない場合に、正常なデータを記憶デバイスユニット20にデステージすることにより、データの一貫性を保証できるようにしていたが、本発明はこれに限られず、例えば、複数のコントローラ22のキャッシュ領域243上での二重化ができない場合に、正常なデータを、正常な動作が可能なコントローラ22のキャッシュ領域243にコピーして、

10

20

30

40

50

複数のコントローラ 2 2 のキャッシュ領域 2 4 3 上で二重化させるようにしてもよい。

【 0 1 5 8 】

また、上記実施形態では、ライトデータを複数のコントローラ 2 2 のキャッシュ領域 2 4 3 上で二重化をさせるようにしていたが、本発明はこれに限られず、3 以上の多重化をさせるようにしてもよく。

【 符号の説明 】

【 0 1 5 9 】

1 ... ホスト計算機、2 ... ストレージシステム、1 1 ... ネットワーク、2 0 ... 記憶デバイスユニット、2 2 ... コントローラ、2 3 ... H C A ネットワーク、1 0 0 ... 計算機システム、2 0 0 ... P D E V、2 3 0 ... C P U、2 4 0 ... メモリ、2 4 3 ... キャッシュ領域、2 4 4 ... コントローラ状態管理情報、2 4 5 ... キャッシュ状態管理情報、2 4 6 ... 転送管理情報、2 4 7 ... 転送状態管理情報、2 5 0 ... H C A

10

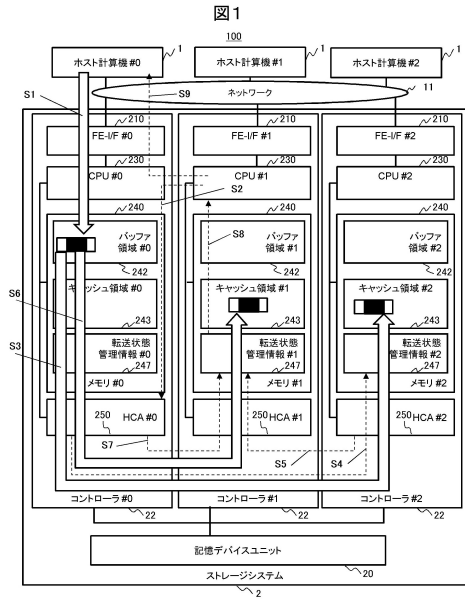
20

30

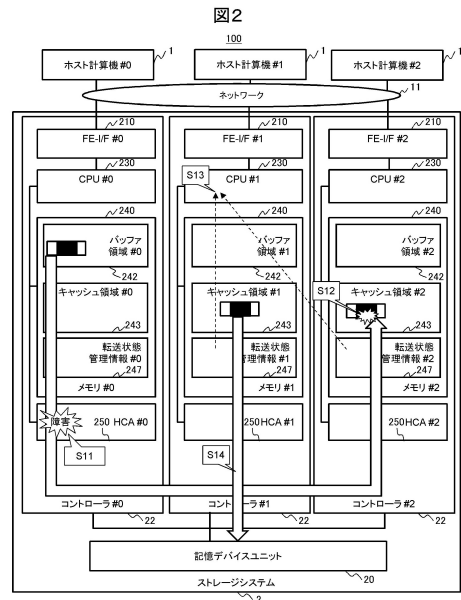
40

50

【図面】
【図 1】



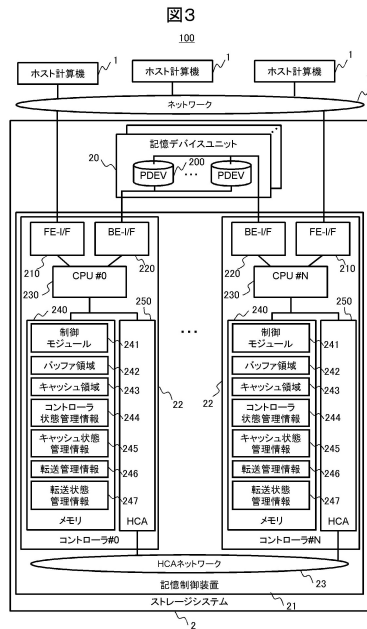
【図 2】



10

20

【図 3】



【図 4】

図 4

401 コントローラID	402 状態
0	正常
1	故障
2	正常

30

40

50

【 図 5 】

図5

データアドレス	第一転送先 コントローラID	第二転送先 コントローラID	キャッシュ 状態
0x1000	2	1	ダーティ
0x2000	0	1	クリーン
0x3000	2	0	ダーティ

【 図 6 】

図6

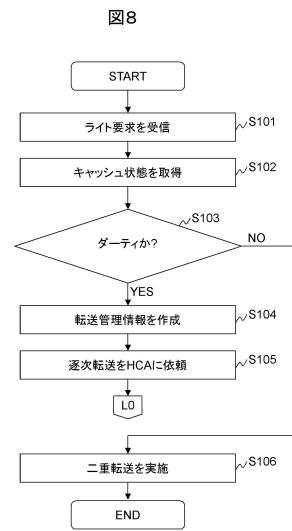
データアドレス	コントローラID	転送状態管理 情報アドレス
0x1000	0	0x10000
0x2000	2	0x20000
0x3000	1	0x30000

【 図 7 】

図7

データアドレス	第一転送データ受領済フラグ	第一転送完了フラグ	第二転送完了フラグ
0x1000	1	1	0
0x2000	1	1	1
0x3000	1	0	0

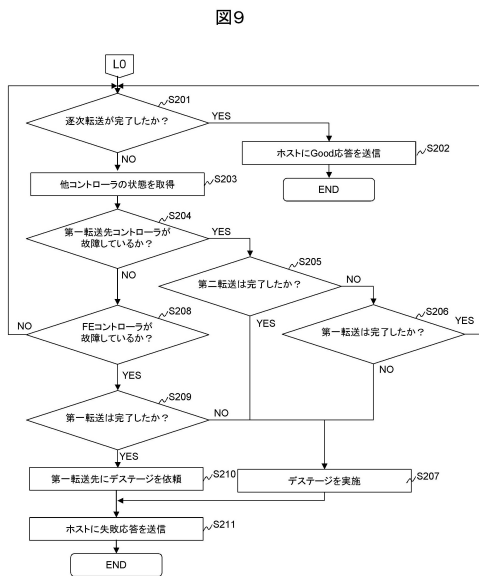
【 図 8 】



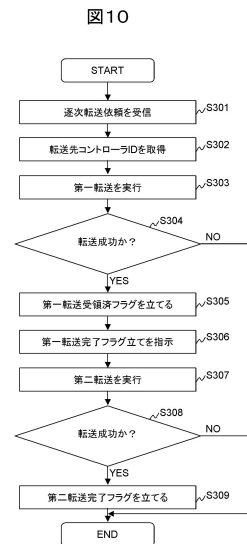
10

20

【 図 9 】



【 図 10 】



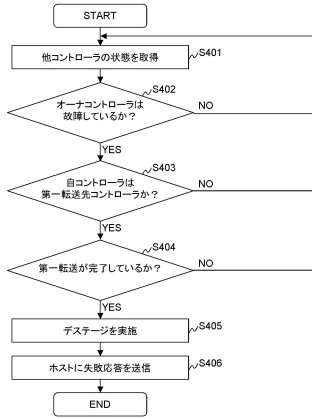
30

40

50

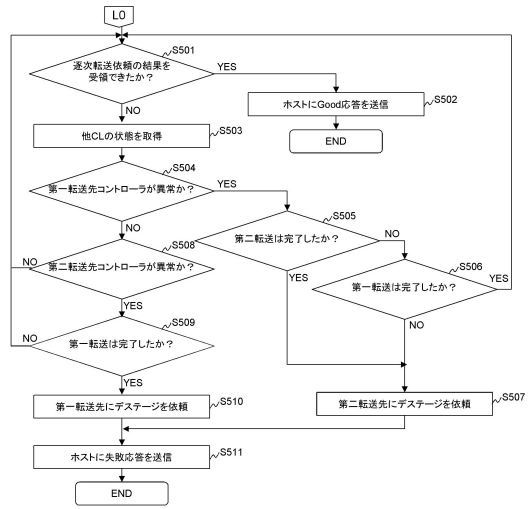
【 図 1 1 】

図 11



【 図 1 2 】

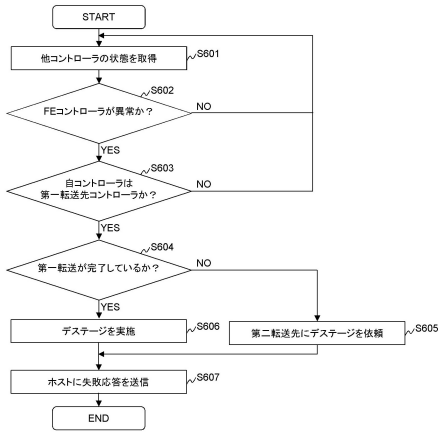
図 12



10

【 図 1 3 】

図 13



20

30

40

50

フロントページの続き

(51)国際特許分類

F I
G 0 6 F 12/0868

(72)発明者 長尾 尚

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 国際公開第2017/056219(WO, A1)

国際公開第2015/052798(WO, A1)

(58)調査した分野 (Int.Cl., DB名)

G 0 6 F 3 / 0 6 - 3 / 0 8

G 0 6 F 1 3 / 1 0 - 1 3 / 1 4

G 0 6 F 1 1 / 2 0

G 0 6 F 1 2 / 0 8 6 6 - 1 2 / 0 8 7 3