

【公報種別】特許法第17条の2の規定による補正の掲載

【部門区分】第6部門第3区分

【発行日】平成18年1月5日(2006.1.5)

【公表番号】特表2005-508551(P2005-508551A)

【公表日】平成17年3月31日(2005.3.31)

【年通号数】公開・登録公報2005-013

【出願番号】特願2003-542490(P2003-542490)

【国際特許分類】

**G 06 F 3/06 (2006.01)**

**G 06 F 12/00 (2006.01)**

**G 06 F 13/28 (2006.01)**

【F I】

G 06 F 3/06 3 0 4 B

G 06 F 3/06 5 4 0

G 06 F 12/00 5 1 4 E

G 06 F 12/00 5 3 1 D

G 06 F 12/00 5 4 5 A

G 06 F 13/28 3 1 0 A

【手続補正書】

【提出日】平成17年9月13日(2005.9.13)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

記憶アレイを含む記憶システムにおいてデータをミラー化するための方法であつて、

第1のダイレクトメモリアクセスエンジンを含む第1の制御装置管理モジュールであつて、前記記憶アレイに関する読み出し／書き込み動作を制御するための前記第1の制御装置管理モジュールを提供するステップ、

第2のダイレクトメモリアクセスエンジンを含む第2の制御装置管理モジュールであつて、前記記憶アレイに関する読み出し／書き込み動作を制御し、前記第1のダイレクトメモリアクセスエンジンが前記第2の制御装置管理モジュールにデータをミラー化するのに用いられ、及び前記第2のダイレクトメモリアクセスエンジンが前記第1の制御装置管理モジュールにデータをミラー化するのに用いられる前記第2の制御装置管理モジュールを提供するステップ、及び

前記第1のダイレクトメモリアクセスエンジンを用いて、前記第1の制御装置管理モジュールから前記第2の制御装置管理モジュールへ第1のデータをミラー化するステップ、を含む方法。

【請求項2】

前記第1の制御装置管理モジュールが第1のプロセッサを含み、前記ミラー化するステップが、前記第1のプロセッサを用いて前記ミラー化するステップが行われるかどうかを判定するステップを含む請求項1に記載の方法。

【請求項3】

前記第2の制御装置管理モジュールが第2のプロセッサを含み、前記ミラー化するステップが、前記第2のプロセッサから独立して前記第1のデータをミラー化するステップを含む請求項1に記載の方法。

**【請求項 4】**

前記第2の制御装置管理モジュールが第2のプロセッサを含み、前記第2のプロセッサが、前記記憶アレイに関する前記読み出し／書き込み動作の前記制御するステップに用いられ、前記第2のプロセッサの割り込みを回避しながら前記ミラー化するステップが行われる請求項1に記載の方法。

**【請求項 5】**

前記第2の制御装置管理モジュールが第2のプロセッサを含み、前記ミラー化するステップの間、前記第2のプロセッサが、第1の読み出し動作及び第1の書き込み動作のうち少なくとも1つを制御するのに用いられるように動作可能である請求項1に記載の方法。

**【請求項 6】**

前記第2の制御装置管理モジュールが第2のプロセッサを含み、前記第2のダイレクトメモリアクセスエンジン及び前記第2のプロセッサを用いない間、前記ミラー化するステップが、前記第1のダイレクトメモリアクセスエンジンを用いて実行される請求項1に記載の方法。

**【請求項 7】**

前記第2の制御装置管理モジュールが第2のプロセッサを含み、前記ミラー化するステップが、前記第2のプロセッサを用いることなく前記第2の制御装置管理モジュールの不揮発性メモリに前記第1のデータを記憶するステップを含む請求項1に記載の方法。

**【請求項 8】**

前記第2の制御装置管理モジュールが不揮発性メモリを含み、前記ミラー化するステップが、

第1に、前記第1のデータを受信する前記不揮発性メモリの内容の部分を無効として記録し、第1のDMAトランザクションにおいて前記不揮発性メモリに前記第1のデータを転送するステップ、及び

第2に、第2のDMAトランザクションにおいて前記不揮発性メモリの内容の前記部分を有効として記録するステップと、を含む請求項1に記載の方法。

**【請求項 9】**

前記第2の制御装置管理モジュールが不揮発性メモリを含み、前記ミラー化するステップが、

第1の文字列を前記不揮発性メモリの第1の記憶領域に記憶し、前記第1のデータを前記不揮発性メモリに転送し、前記不揮発性メモリの第2の記憶領域に前記第1の文字列を記憶するステップを含む請求項1に記載の方法。

**【請求項 10】**

前記ミラー化するステップが単一のDMAトランザクションを用いて実行される請求項9に記載の方法。

**【請求項 11】**

前記第1の制御装置管理モジュールが、第1のプロセッサを含み、更に、前記第1のプロセッサを用い、かつ、パリティを判定するためにXORエンジンを用いて、前記記憶アレイに情報を記憶するステップを含む、請求項1に記載の方法。

**【請求項 12】**

記憶アレイを含む記憶システムにおいてデータをミラー化するための装置であって、

第1の制御装置管理モジュールは第1のプロセッサ及び第1のダイレクトメモリアクセスエンジンを含み、前記第1のプロセッサは記憶アレイに関する読み出し動作及び書き込み動作を制御するのに使用され、前記第1のダイレクトメモリアクセスエンジンは前記第1の制御装置管理モジュールによって受信されるデータを記憶するのに使用され、

第2の制御装置管理モジュールが第2のプロセッサ及び第2のダイレクトメモリアクセスエンジンを含み、前記第2のプロセッサは記憶アレイに関する読み出し動作及び書き込み動作を制御するのに使用され、前記第2のダイレクトメモリアクセスエンジンは前記第2の制御装置管理モジュールによって受信されるデータを記憶するのに使用され、

第1のデータは前記第1の制御装置管理モジュールによってホストから受信され、前記

第1のデータは前記第2のプロセッサの割り込みを回避しながら、前記第1のダイレクトメモリアクセスエンジンを用いて、前記第1の制御装置管理モジュールから前記第2の制御装置管理モジュールへミラー化される装置。

【請求項13】

前記第1の制御装置管理モジュールが不揮発性メモリを含み、前記第1のデータが前記不揮発性メモリに記憶される請求項12に記載の装置。

【請求項14】

前記第1のダイレクトメモリアクセスエンジンが前記第1のプロセッサから分離しているが通信し、前記第1のプロセッサは前記第1のダイレクトメモリアクセスエンジンを用いて前記第1のデータのミラー化を開始する請求項12に記載の装置。

【請求項15】

前記第1の制御装置管理モジュールがフィールドプログラマブルゲートアレイを有し、前記第1のダイレクトメモリアクセスエンジンが少なくともその部分と通信する請求項12に記載の装置

【請求項16】

前記装置が、更に、

第1の共有経路を有する第1のチャネルインターフェースモジュールを有し、前記第1のチャネルインターフェースモジュールは前記第1の制御装置管理モジュールと通信し、前記第1の共有経路が前記第1の制御装置管理モジュールと前記第2の制御装置管理モジュールとの間で前記第1のデータを転送するのに用いられる請求項12に記載の装置。

【請求項17】

前記装置が、更に、

前記第1のチャネルインターフェースモジュール及び前記第1の制御装置管理モジュールを相互接続する受動パックプレーンを有する請求項16に記載の装置。

【請求項18】

前記第2のプロセッサが、前記第1のデータが前記第2の制御装置管理モジュールにミラー化されている間、前記第2の制御装置管理モジュールに関連する動作を制御する請求項12に記載の装置。

【請求項19】

前記第1のデータが、前記第2のダイレクトメモリアクセスエンジンから独立して前記第2の制御装置管理モジュールにミラー化される請求項12に記載の装置。

【請求項20】

前記第2の制御装置管理モジュールが不揮発性メモリを有し、前記第1のデータが前記第2のプロセッサから独立して前記不揮発性メモリに記憶される請求項12に記載の装置。

【請求項21】

前記第2の制御装置管理モジュールが不揮発性メモリを有し、前記第1のダイレクトメモリアクセスエンジンは、前記第1のデータを受信するべき前記不揮発性メモリの少なくとも部分に対して、前記部分が無効であるという指示を提供するのに用いられ、前記第1のデータが不揮発性メモリによって受信された後、前記第1のダイレクトメモリアクセスエンジンが前記部分を有効として記録するのに用いられる請求項12に記載の装置。

【請求項22】

前記第2の制御装置管理モジュールが少なくとも第1の記憶領域及び第2の記憶領域を有する不揮発性メモリを有し、

前記第1のダイレクトメモリアクセスエンジンは、前記第1のデータが前記不揮発性メモリによって受信される前に、前記第1の記憶領域での記憶のために第1の文字列を提供すること、及び、前記第1のデータが前記不揮発性メモリによって受信された後に、前記第2の記憶領域での記憶のために前記第1の文字列を提供することに用いられる請求項12に記載の装置。

【請求項23】

前記第1のダイレクトメモリアクセスエンジンが、単一のダイレクトメモリアクセストランザクションにおいて、前記第1の文字列を提供し、前記第1のデータを転送し、及び、前記第1の文字列を提供するよう、動作可能である請求項21に記載の装置。

【手続補正2】

【補正対象書類名】明細書

【補正対象項目名】0006

【補正方法】変更

【補正の内容】

【0006】

制御装置30、34は、入出力モジュール-1 42と入出力モジュール-2 46の2つの入出力モジュールに接続されるファイバチャネルバス38に接続される。各制御装置30、34は、C P Uサブシステム50、2倍のデータ転送速度の(DDR:Double Data Rate)メモリ54、制御論理58、2つのホストポート62a、62bのデュアル(dual)ポートファイバチャネル接続、及び2つのディスクポート66a、66bのデュアルポートファイバチャネル接続を含む。C P Uサブシステム50は、データの分散化と、読み取り及び書き込みコマンドの開始及び実行とを含む、ディスクアレイへのデータの記憶に必要なタスクを実行する。DDRメモリ54は、データ及びその他の情報のための不揮発性記憶領域である。制御論理58は、C P Uサブシステム50と、DDRメモリ54と、ホストポート62a、62b及びディスクポート66a、66bとの仲介等のいくつかの機能を実行する。制御論理58は、排他的論理和(XOR:Exclusive OR)エンジン等のパリティ生成機能を含むその他の機能をも有していても良い。ホストポート62a、62b及びディスクポート66a、66bは、ファイバチャネルバックプレーン38との通信を提供する。入出力モジュール42、46は、ポートバイパス回路(port bypass circuit)としても良く知られている、各ホストポート14、18及び各ディスクポート22、26を各制御装置30、34に接続するために機能するリンク復元回路(LRC:Link Resiliency Circuit)70を含む。これによって、両制御装置30、34は、両ホストポート14、18及び両ディスクポート22、26にアクセス可能となる。

【手続補正3】

【補正対象書類名】明細書

【補正対象項目名】0016

【補正方法】変更

【補正の内容】

【0016】

本発明によれば、記憶アレイを有する記憶システムにおいてデータをミラー化するための方法及び装置が提供される。前記装置は、第1のプロセッサ及び第1のダイレクトメモリアクセスエンジンを含む第1の制御装置管理モジュールを含む。前記第1のプロセッサは、記憶アレイに関わる読み出し動作及び書き込み動作を制御するのに用いられる。前記第1のダイレクトメモリアクセスエンジンは、前記第1の制御装置管理モジュールによって受信されるデータを記憶するのに用いられる。前記装置はまた、第2のプロセッサ及び第2のダイレクトメモリアクセスエンジンを含む第2の制御装置管理モジュールを有する。前記第2のプロセッサは、記憶アレイに関する読み出し動作及び書き込み動作を制御するのに用いられる。前記第2のダイレクトメモリアクセスエンジンは、前記第2の制御装置管理モジュールから前記第1の制御装置管理モジュールにデータを転送するのに用いられることがある。前記データは、前記第2のプロセッサの割り込みを回避しつつ、前記第1のダイレクトメモリアクセスエンジンを用いて前記第1の制御装置管理モジュールから前記第2の制御装置管理モジュールへミラー化される。前記第1のダイレクトメモリアクセスエンジンは前記第1のプロセッサから分離しているが通信し、前記第1のプロセッサは前記第1のダイレクトメモリアクセスエンジンを用いて前記データのミラー化を制御する。ある実施の形態において、前記第1の制御装置管理モジュールは、フィールドプログラマブルゲートアレイ(field programmable gate array)を有する。前記第1のダイレ

クトメモリアクセスエンジンは、フィールドプログラマブルゲートアレイの少なくとも一部と通信し、前記第1のダイレクトメモリアクセスエンジンは、フィールドプログラマブルゲートアレイの1部であることもできる。

【手続補正4】

【補正対象書類名】明細書

【補正対象項目名】0017

【補正方法】変更

【補正の内容】

【0017】

ある実施の形態において、前記装置は、第1の共有経路を有する第1のチャネルインタフェースモジュールを含む。前記第1のチャネルインタフェースモジュールは前記第1の制御装置管理モジュールと通信し、前記第1の共有経路は、前記第1の制御装置管理モジュールと前記第2の制御装置管理モジュールとの間でデータを転送するのに用いられる。受動バックプレーンは、前記第1のチャネルインタフェースモジュール及び前記第1の制御装置管理モジュールを相互接続する。前記第2のプロセッサは、前記データが前記第2の制御装置管理モジュールにミラー化されている間、前記第2の制御装置管理モジュールに関わる動作を制御する。前記データは、前記第2のダイレクトメモリアクセスエンジンから独立して前記第2の制御装置管理モジュールにミラー化される。前記第2の制御装置管理モジュール内には不揮発性メモリがあり、前記データは前記第2のプロセッサから独立して前記不揮発性メモリ内に記憶され得る。前記第1のダイレクトメモリアクセスエンジンは、前記データが記憶される前記不揮発性メモリの部分を無効として記録し、前記不揮発性メモリに前記データを転送する。前記データが記憶された前記不揮発性メモリの前記部分は、次に有効として記録される。

【手続補正5】

【補正対象書類名】明細書

【補正対象項目名】0021

【補正方法】変更

【補正の内容】

【0021】

ネットワーク記憶装置100は、1つ以上のチャネルインタフェースモジュール(CIM: Channel Interface Module)を有する。図2に示される実施の形態において、CIM-1 136とCIM-2 140の、2つのCIMがあるが、この数はネットワーク記憶装置100が使用される構成及びアプリケーションに応じて変化しても良いことは理解されるであろう。各CIM136、140は、2つのCIMバスインターフェースポート144a、144bを有する。各CIM136、140上では、1つのCIMバスインターフェースポート144aがCIMバス接続146を介してCMM-A104に接続される1つのバスと接続し、1つのCIMバスインターフェースポート144bがCIMバス接続146を介してCMM-B108に接続される1つのバスと接続する。図2に示される実施の形態において、CIM-1 136が第1のデータバス120及び第3のデータバス128に接続し、CIM-2 140が第2のデータバス124及び第4のデータバス132に接続する。各CIM136、140は、ホストコンピュータ(図示せず)に接続されるホストチャネル152に接続する2つのホストポート148を有する。また、各CIM136、140は、1つ以上の記憶装置(図示せず)に接続されるディスクチャネル158に接続する2つのディスクポート156をも有する。記憶装置は、RAIDアレイ等の記憶アレイであっても良い。下記に詳述するように、別の実施の形態において、CIMは、要求されるアプリケーション及びチャネルインタフェースに応じて、複数のホストポート或いは複数のディスクポートを含んでも良い。

【手続補正6】

【補正対象書類名】明細書

【補正対象項目名】0026

【補正方法】変更

【補正の内容】

【0026】

メモリインタフェース176は、メモリ164とインタフェースFPGA168間のインタフェースとしての役割を果たす。XORエンジン180は、書き込まれるデータに関するパリティ情報を得るために、記憶されるデータに関するXOR動作を実行するように働く。XORエンジン180は、ディスクアレイ内の故障ドライブからデータを復元するためにパリティ情報の使用が必要とされる状況においても使用される。XORエンジン180は、PCIインターフェース172を通じてCPUサブシステム160に接続する。複数のデータ FIFO192はメモリインタフェース176及びブリッジコア184に接続し、順次バックプレーンインターフェース112に接続する。複数のデータ FIFOは、読み取り及び書き込み動作を処理するためにCMM104によって使用される待ち行列として働く。下記に詳述されるように、DMAエンジン188は、CMMが冗長性を提供するように動作している時、別のCMMにDMAデータを提供するように働く。ある実施の形態において、DMAエンジン188はまた、XOR動作を実行するために、XORエンジン180と関連して用いられ、メモリ164内の2つの領域からデータを読み出し、XORエンジン180にデータを提供し、メモリ164内の第3の領域にXORエンジンの出力を書き込む。

【手続補正7】

【補正対象書類名】明細書

【補正対象項目名】0028

【補正方法】変更

【補正の内容】

【0028】

チャネルインターフェース204は、交換PCI X経路208、212をホストポート148及びディスクポート156に接続する。チャネルインターフェース204は、交換PCI X経路208、212上で送信されるデータを監視し、データがホストポート148或いはディスクポート156のどちらに経路付けされるかを判定するように動作可能である。この監視及び経路付けは、データを適切なディスク又はホスト位置に通過させ、ミラーハードデータをホスト又はディスクポート148、156にまで通過させない。チャネルインターフェース204は、アプリケーションにとって最適なチャネル媒体上での通信を可能にする。例えば、ホストチャネル152及びディスクチャネル158がファイバチャネルを使用する場合、チャネルインターフェース204は交換PCI X経路208、212とファイバチャネルとの間のインターフェースの役割を果たすであろう。同様に、ホストチャネル152及びディスクチャネル158がSCSIチャネルを使用する場合、チャネルインターフェース204は交換PCI X経路208、212とSCSIチャネルとの間のインターフェースの役割を果たすであろう。ホストチャネル152及びディスクチャネル158の両方が同じチャネル媒体を使用する場合、CIM136は、ホストポート148及びディスクポート156を使用してホストチャネル152及びディスクチャネル158の両方との通信のため用いられることができる。

【手続補正8】

【補正対象書類名】明細書

【補正対象項目名】0030

【補正方法】変更

【補正の内容】

【0030】

図4に示される実施の形態において、バスインターフェースポート144及びCIMバス接続146を通じて、第1の交換PCI X経路208が第1のデータバス120と通信し、第2の交換PCI X経路212が第3のデータバス128と通信する。下記に詳述するように、PCI Xブリッジ200は、一方のCMMが他方のCMMと通信するための通信

経路として使用されても良い。ネットワーク制御装置上に存在する残りの C I M に対しても、同様の構成が使用されると理解されるであろう。例えば、図 2 に示される実施の形態において、C I M - 2 1 4 0 は第 2 のデータバス 1 2 4 及び第 4 のデータバス 1 3 2 に接続され、従って、C I M - 2 1 4 0 は、第 2 のデータバス 1 2 4 及び第 4 のデータバス 1 3 2 とそれぞれ通信する交換 P C I X 経路 2 0 8、2 1 2 を有するであろう。同様に、2 つを超える C I M が存在する場合、アプリケーションによる要求に応じて受動バックプレーン 1 1 6 上の適切なバスと通信するよう構成されるであろう。

#### 【手続補正 9】

【補正対象書類名】明細書

【補正対象項目名】0 0 3 3

【補正方法】変更

【補正の内容】

#### 【0 0 3 3】

当業者によって理解されるように、冗長制御装置は、記憶サブシステムに取り付けられた 2 つの制御装置間のデータのミラー化を必要とする。これは、制御装置がホストコンピュータからデータを受信し、データをキャッシュし、データが書き込まれたというメッセージをホストコンピュータに送信する、書き戻しキャッシュ(write back cache)の使用による。従って、データが実際に制御装置に格納され、ディスクアレイ内のドライブに書き込まれるためにそこで待機している時、ホストコンピュータは、データが書き込まれたと判定する。故障の際に、このデータが失われないよう確保するのを助けるため、冗長制御装置はこのデータを他の制御装置にミラー化し、従って、他の制御装置にデータの別のコピーができる。このことは、キャッシュの一貫性として既知である。ある実施の形態において、C M M 1 0 4、1 0 8 は、ネットワーク記憶装置 1 0 0 a にキャッシュの一貫性を提供するためにデータをミラー化する。これは、C M M - A 1 0 4 と C M M - B 1 0 8 との間に D M A 経路を設けることによって実現できる。これは、図 3 に関して上述したように、インターフェース F P G A 1 6 8 内の D M A エンジン 1 8 8 を提供することによって達成されることが出来、共有経路は、図 4 に関して上述したように、P C I X ブリッジ 2 0 0 を利用する。各 C M M 1 0 4、1 0 8 は、他の C M M にデータを送信するためにこの D M A 経路を用いる。D M A 経路を利用することによって、2 つの C M M 1 0 4、1 0 8 は、ホストチャネル 1 5 2 又はディスクチャネル 1 5 8 を使用する必要なくデータをミラー化でき、従って、ディスクチャネル 1 5 8 又はホストチャネル 1 5 2 のチャネル帯域幅はデータをミラー化している C M M 1 0 4、1 0 8 によって消費されない。

#### 【手続補正 1 0】

【補正対象書類名】明細書

【補正対象項目名】0 0 3 6

【補正方法】変更

【補正の内容】

#### 【0 0 3 6】

他の実施の形態において、図 1 3 に示された 2 つの D M A トランザクションは、単一の順序付けられた D M A トランザクションに組み合わせられる。この実施の形態において、各 C M M 1 0 4、1 0 8 内の D D R メモリ 1 6 4 は、ユーザデータに関連するメタデータを記憶するための 2 つの記憶領域を有する。D M A 転送を開始する場合、第 1 のユニークな文字列が第 1 の記憶領域に記憶され、その後に転送されるデータが続く。D M A 転送の最後に、第 1 のユニークな文字列は第 2 の記憶領域に記憶される。C M M 1 0 4、1 0 8 が故障から回復されなければならない際には、第 1 の及び第 2 の記憶領域内に記憶された文字列が比較される。比較は C M M の故障が起こった際にのみ実行され、従って、通常の動作の間、パフォーマンスを低下させない。文字列が一致した場合、ミラー化されたデータが有効であるということを示す。