

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5358736号  
(P5358736)

(45) 発行日 平成25年12月4日(2013.12.4)

(24) 登録日 平成25年9月6日(2013.9.6)

(51) Int.Cl. F I  
G O 6 F 3/06 (2006.01) G O 6 F 3/06 3 O 1 C

請求項の数 13 (全 30 頁)

(21) 出願番号	特願2012-516982 (P2012-516982)	(73) 特許権者	000005108
(86) (22) 出願日	平成21年11月10日(2009.11.10)		株式会社日立製作所
(65) 公表番号	特表2012-531656 (P2012-531656A)		東京都千代田区丸の内一丁目6番6号
(43) 公表日	平成24年12月10日(2012.12.10)	(74) 代理人	110000279
(86) 国際出願番号	PCT/JP2009/005995		特許業務法人ウィルフォート国際特許事務所
(87) 国際公開番号	W02011/058598	(72) 発明者	中村 崇仁
(87) 国際公開日	平成23年5月19日(2011.5.19)		神奈川県横浜市戸塚区吉田町292番地
審査請求日	平成23年12月27日(2011.12.27)	(72) 発明者	新井 政弘
			神奈川県横浜市戸塚区吉田町292番地
			株式会社日立製作所 横浜研究所内

最終頁に続く

(54) 【発明の名称】 複数のコントローラを備えたストレージシステム

(57) 【特許請求の範囲】

【請求項1】

第1のコントローラと、

前記第1のコントローラに第1パスを介して接続された第2のコントローラとを備え、

ホスト装置が発行したI/O (Input/Output) コマンドを前記第1及び第2のコントローラのいずれかが受け付けた場合、そのI/Oコマンドを受けた第1又は第2のコントローラが、そのI/Oコマンドに従う処理であるI/O処理を行い、前記I/O処理において、そのI/Oコマンドに従うデータのI/Oを記憶デバイスに対して行い、

前記第1のコントローラが、

データ転送を制御する回路である第1の中継回路と、

前記第1の中継回路に第1の第2パスを介して接続された第1のプロセッサとを有し、

前記第2のコントローラが、

データ転送を制御する回路であり、前記第1パスを介して前記第1の中継回路に接続された第2の中継回路と、

前記第2の中継回路に第2の第2パスを介して接続された第2のプロセッサとを有し、

前記第1のプロセッサが、前記第1の中継回路を介することなく第1の第3パスを介して前記第2の中継回路に接続されており、前記第1のコントローラが前記I/O処理を行

っている場合、その I / O 処理において、前記第 1 の第 3 パスを介して前記第 2 の中継回路にアクセスし、

前記第 2 のプロセッサが、前記第 2 の中継回路を介することなく第 2 の第 3 パスを介して前記第 1 の中継回路に接続されており、前記第 2 のコントローラが前記 I / O 処理を行っている場合、その I / O 処理において、前記第 2 の第 3 パスを介して前記第 1 の中継回路にアクセスし、

前記第 1 の第 3 パスは、前記第 1 のプロセッサが前記第 2 の中継回路を直接利用するダイレクトパスであり、前記第 2 の第 3 パスは、前記第 2 のプロセッサが前記第 1 の中継回路を直接利用するダイレクトパスである、

記憶制御装置。

10

【請求項 2】

請求項 1 記載の記憶制御装置であって、

前記第 1 の中継回路が、データの転送を行う回路である第 1 のデータ転送回路を有し、  
前記第 2 の中継回路が、データの転送を行う回路である第 2 のデータ転送回路を有し、  
前記第 1 のコントローラが、前記第 1 の中継回路に接続された第 1 のメモリと、プロセッサ毎に設けられておりデータ転送のパラメータである転送パラメータが蓄積される第 1 キューとを有し、

各第 1 キューは、前記第 1 のデータ転送回路についてのキューであり、

前記第 2 のコントローラが、前記第 2 の中継回路に接続された第 2 のメモリと、プロセッサ毎に設けられており転送パラメータが蓄積される第 2 キューとを有し、

20

各第 2 キューは、前記第 2 のデータ転送回路についてのキューであり、

前記転送パラメータは、前記第 1 及び第 2 プロセッサによって生成され、データの転送元の記憶領域のアドレスである転送元アドレスと、そのデータの転送先のアドレスである転送先アドレスとを含み、

前記第 1 のデータ転送回路が、複数の第 1 キューに蓄積されている転送パラメータの数をそれぞれ記憶する記憶領域である複数の第 1 インデックスと、前記複数の第 1 インデックスから一つの第 1 インデックスを選択する第 1 セレクタと、前記第 1 セレクタによって選択された第 1 インデックスに対応した第 1 キューから転送パラメータを取得しその転送パラメータが有する転送元アドレス及び転送先アドレスを設定する第 1 パラメータ取得回路と、前記設定された転送元アドレスが表す記憶領域内のデータを前記設定された転送先アドレスが表す記憶領域に転送する第 1 転送制御回路とを有し、

30

前記第 2 のデータ転送回路が、複数の第 2 キューに蓄積されている転送パラメータの数をそれぞれ記憶する記憶領域である複数の第 2 インデックスと、前記複数の第 2 インデックスから一つの第 2 インデックスを選択する第 2 セレクタと、前記第 2 セレクタによって選択された第 2 インデックスに対応した第 2 キューから転送パラメータを取得しその転送パラメータが有する転送元アドレス及び転送先アドレスを設定する第 2 パラメータ取得回路と、前記設定された転送元アドレスが表す記憶領域内のデータを前記設定された転送先アドレスが表す記憶領域に転送する第 2 転送制御回路とを有し、

前記第 1 及び第 2 プロセッサのうちの I / O コマンドを受領したプロセッサである対象プロセッサが、下記の (A) 及び (B) :

40

(A) データの転送先の記憶領域が、前記対象プロセッサを有するコントローラである対象コントローラに存在するか否か；

(B) 前記対象プロセッサに対応した第 1 キューについての未完了の転送パラメータの数である第 1 の数と、前記対象プロセッサに対応した第 2 キューについての未完了の転送パラメータの数である第 2 の数、

に基づいて、前記第 1 及び第 2 のデータ転送回路のうちのいずれかのデータ転送回路を選択し、選択したデータ転送回路を有するコントローラ内の、前記対象プロセッサに対応したキューに、転送パラメータを格納する、

記憶制御装置。

【請求項 3】

50

請求項 2 記載の記憶制御装置であって、

前記対象プロセッサが、下記 ( X ) 及び ( Y ) の条件に適合するデータ転送回路を優先的に選択する、

( X ) 前記転送先記憶領域が前記対象コントローラ内にある場合、前記対象コントローラに存在する；

( Y ) 前記第 1 の数と前記第 2 の数とのうちの少ない方に対応する、  
記憶制御装置。

【請求項 4】

請求項 3 記載の記憶制御装置であって、

前記第 1 のメモリが、第 1 のキャッシュメモリ領域を有し、

前記第 2 のメモリが、第 2 のキャッシュメモリ領域を有し、

前記第 1 のプロセッサが、受け付けた I / O コマンドに従うデータの転送先アドレスとして前記第 2 のキャッシュメモリ領域の任意の位置のアドレスを含んだ転送パラメータを、前記第 1 の第 3 パスを介して、前記第 1 のプロセッサに対応した第 2 キューに格納し、

前記第 2 のプロセッサが、受け付けた I / O コマンドに従うデータの転送先アドレスとして前記第 1 のキャッシュメモリ領域の任意の位置のアドレスを含んだ転送パラメータを、前記第 2 の第 3 パスを介して、前記第 1 のプロセッサに対応した第 1 キューに格納する、

記憶制御装置。

【請求項 5】

請求項 1 記載の記憶制御装置であって、

プロセッサ毎に設けられておりデータ転送のパラメータである転送パラメータが蓄積される第 1 パラメータ領域と、

プロセッサ毎に設けられており転送パラメータが蓄積される第 2 パラメータ領域とを有し、

前記第 1 の中継回路が、データの転送を行う回路である第 1 のデータ転送回路を有し、

前記第 2 の中継回路が、データの転送を行う回路である第 2 のデータ転送回路を有し、

前記第 1 パラメータ領域は、前記第 1 のデータ転送回路についての記憶領域であり、

前記第 2 パラメータ領域は、前記第 2 のデータ転送回路についての記憶領域であり、

前記第 1 のプロセッサが、前記第 1 のデータ転送回路を利用する場合、前記第 1 のプロセッサに対応した第 1 パラメータ領域に転送パラメータを格納し、前記第 2 のデータ転送回路を利用する場合、前記第 1 のプロセッサに対応した第 2 パラメータ領域に転送パラメータを格納し、

前記第 2 のプロセッサが、前記第 2 のデータ転送回路を利用する場合、前記第 2 のプロセッサに対応した第 2 パラメータ領域に転送パラメータを格納し、前記第 1 のデータ転送回路を利用する場合、前記第 2 のプロセッサに対応した第 1 パラメータ領域に転送パラメータを格納し、

前記第 1 のデータ転送回路が、いずれかの第 1 パラメータ領域から転送パラメータを取得し、その転送パラメータに従って、データ転送を実行し、

前記第 2 のデータ転送回路が、いずれかの第 2 パラメータ領域から転送パラメータを取得し、その転送パラメータに従って、データ転送を実行する、  
記憶制御装置。

【請求項 6】

請求項 5 記載の記憶制御装置であって、

前記第 1 及び第 2 プロセッサのうちの I / O コマンドを受領したプロセッサである対象プロセッサが、データの転送先の記憶領域が、前記対象プロセッサを有するコントローラである対象コントローラに存在する場合、前記対象コントローラ内のデータ転送回路を選択し、選択したデータ転送回路を有するコントローラ内の、前記対象プロセッサに対応したパラメータ領域に、転送パラメータを格納する、

記憶制御装置。

10

20

30

40

50

## 【請求項 7】

請求項 5 記載の記憶制御装置であって、

前記第 1 及び第 2 プロセッサのうちの I / O コマンドを受領したプロセッサである対象プロセッサが、前記対象プロセッサに対応した第 1 パラメータ領域内の未処理の転送パラメータの数と、前記対象プロセッサに対応した第 2 パラメータ領域内の未処理の転送パラメータの数とのうち少ない方に対応したデータ転送回路を選択し、選択したデータ転送回路を有するコントローラ内の、前記対象プロセッサに対応したパラメータ領域に、転送パラメータを格納する、

記憶制御装置。

## 【請求項 8】

請求項 5 記載の記憶制御装置であって、

前記第 1 のコントローラが、前記第 1 の中継装置に接続された第 1 のキャッシュメモリ領域を有し、

前記第 2 のコントローラが、前記第 2 の中継装置に接続された第 2 のキャッシュメモリ領域を有し、

前記第 1 のプロセッサが、前記第 2 のキャッシュメモリ領域内の任意のアドレスを転送先アドレスとした転送パラメータを、前記第 1 のプロセッサに対応した第 2 パラメータ領域に格納し、

前記第 2 のプロセッサが、前記第 1 のキャッシュメモリ領域内の任意のアドレスを転送先アドレスとした転送パラメータを、前記第 2 のプロセッサに対応した第 1 パラメータ領域に格納する、

記憶制御装置。

## 【請求項 9】

請求項 5 記載の記憶制御装置であって、

前記第 1 のパラメータ領域は、I / O コマンドがライトコマンドである場合の転送パラメータの格納先とされる第 1 のパラメータライト領域と、I / O コマンドがリードコマンドである場合の転送パラメータの格納先とされる第 1 のパラメータリード領域とを有し、

前記第 2 のパラメータ領域は、I / O コマンドがライトコマンドである場合の転送パラメータの格納先とされる第 2 のパラメータライト領域と、I / O コマンドがリードコマンドである場合の転送パラメータの格納先とされる第 2 のパラメータリード領域とを有し、

前記第 1 のデータ転送回路は、前記第 1 のパラメータライト領域よりも前記第 1 のパラメータリード領域内の未処理の転送パラメータを優先的に取得し、

前記第 2 のデータ転送回路は、前記第 2 のパラメータライト領域よりも前記第 2 のパラメータリード領域内の未処理の転送パラメータを優先的に取得する、

記憶制御装置。

## 【請求項 10】

請求項 5 記載の記憶制御装置であって、

前記第 1 のデータ転送回路が、複数の第 1 パラメータ領域に蓄積されている未処理の転送パラメータの数をそれぞれ記憶する記憶領域である複数の第 1 インデックスと、前記複数の第 1 インデックスから一つの第 1 インデックスを選択する第 1 セレクタと、前記第 1 セレクタによって選択された第 1 インデックスに対応した第 1 パラメータ領域から転送パラメータを取得しその転送パラメータが有する転送元アドレス及び転送先アドレスを設定する第 1 パラメータ取得回路と、前記設定された転送元アドレスが表す記憶領域内のデータを前記設定された転送先アドレスが表す記憶領域に転送する第 1 転送制御回路とを有し、

前記第 2 のデータ転送回路が、複数の第 2 パラメータ領域に蓄積されている未処理の転送パラメータの数をそれぞれ記憶する記憶領域である複数の第 2 インデックスと、前記複数の第 2 インデックスから一つの第 2 インデックスを選択する第 2 セレクタと、前記第 2 セレクタによって選択された第 2 インデックスに対応した第 2 パラメータ領域から転送パラメータを取得しその転送パラメータが有する転送元アドレス及び転送先アドレスを設定

10

20

30

40

50

する第 2 パラメータ取得回路と、前記設定された転送元アドレスが表す記憶領域内のデータを前記設定された転送先アドレスが表す記憶領域に転送する第 2 転送制御回路とを有する、

記憶制御装置。

【請求項 1 1】

請求項 1 記載の記憶制御装置であって、

障害の発生を監視する障害監視ユニットを備え、

前記障害監視ユニットは、前記第 1 のプロセッサの障害を検出した場合、前記第 2 のプロセッサに、前記第 1 のプロセッサの障害を通知する、

記憶制御装置。

10

【請求項 1 2】

請求項 1 1 記載の記憶制御装置であって、

前記障害監視ユニットは、前記第 1 の中継装置が有する、

記憶制御装置。

【請求項 1 3】

請求項 1 記載の記憶制御装置であって、

第 1 のスイッチ装置と、第 2 のスイッチ装置とを有し、

前記第 1 パスは、前記第 1 及び第 2 の中継装置を含んだ複数の中継装置に接続された前記第 1 のスイッチ装置で実現され、

前記第 1 の第 3 パス及び前記第 2 の第 3 パスは、前記第 1 及び第 2 のプロセッサを含んだ複数のプロセッサに接続された前記第 2 のスイッチ装置で実現される、

記憶制御装置。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数のコントローラを備えたストレージシステムに関する。

【背景技術】

【0002】

近年、低コストでありながら高性能かつ高機能のストレージシステムが市場から要求されている。ストレージシステムは、一般に、コントローラと物理記憶装置群とに大別することができるが、ストレージシステムを高性能かつ高機能にする方法としては、コントローラを高性能かつ高機能にする方法がある。具体的には例えば、プロセッサの周波数を上げる、コントローラ内でのデータ転送を行う L S I (Large Scale Integration) の性能を高くするなどの方法が考えられる。

30

【0003】

しかし、コントローラを高性能かつ高機能にすると、通常、ストレージシステム自体のコストは上昇する。そのため、コントローラを高性能かつ高機能にすることに代えて、ストレージシステム内での処理を効率的にする必要があると考えられる。

【0004】

ストレージシステム内での処理の効率化に関する技術が、例えば、特許文献 1 及び 2 に開示されている。特許文献 1 によれば、第 1 のコントローラと第 2 のコントローラとの間に、高レスポンスに適した第 1 種のパスと高スループットに適した第 2 種のパスとが設けられ、それら二種類のパスが使い分けられる。特許文献 2 によれば、第 1 種のリソースと第 2 種のリソースのうち一方の種類のリソースの負荷が高く他方の種類のリソースの負荷に余裕がある場合には、他方の種類のリソースの負荷がより高くなるような処理が実行される。

40

【先行技術文献】

【特許文献】

【0005】

50

【特許文献1】特開2001-43026号公報

【特許文献2】特開2008-186108号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

第1及び第2のコントローラを有するストレージシステムでは、通常、第1のコントローラの電源と第2のコントローラの電源は異なっており、第1のコントローラに障害が生じた場合、第2のコントローラが、第1のコントローラに代わって稼働する。このため、第1のコントローラが稼働している間、第2のコントローラが有する資源が利用されないようになっており、第1のコントローラの負荷を低減することができない。

10

【0007】

具体的には、例えば、第1のコントローラは、ホスト装置から受けたライト要求に付随するデータ(ライトデータ)を、ライトデータが失われないようにするために、第1のコントローラが有するキャッシュメモリ(第1のキャッシュメモリ)だけでなく、コントローラ間のバスを介して、第2のコントローラが有するキャッシュメモリ(第2のキャッシュメモリ)にも書き込むことができる。

【0008】

しかし、それは、ライトデータのミラーリングにすぎず、第1のコントローラの負荷を低減することにはならない。

【0009】

そこで、本発明の目的は、第1のコントローラの負荷を低減できるように第2のコントローラの資源を利用可能にすることである。

20

【課題を解決するための手段】

【0010】

第1のコントローラと、前記第1のコントローラに第1バスを介して接続された第2のコントローラとを備える。第1のコントローラが、データ転送を制御する回路である第1の中継回路と、第1の中継回路に第1の第2バスを介して接続された第1のプロセッサとを有する。第2のコントローラが、データ転送を制御する回路であり、第1バスを介して第1の中継回路に接続された第2の中継回路と、第2の中継回路に第2の第2バスを介して接続された第2のプロセッサとを有する。第1のプロセッサが、第1の中継回路を介することなく第1の第3バスを介して第2の中継回路に接続されており、I/O処理において、第1の第3バスを介して前記第2の中継回路にアクセスする。第2のプロセッサが、第2の中継回路を介することなく第2の第3バスを介して第1の中継回路に接続されており、I/O処理において、第2の第3バスを介して第1の中継回路にアクセスする。

30

【0011】

記憶制御装置は、物理記憶デバイスを有したストレージシステムであっても良いし、ストレージシステムに備えられる装置であっても良いし、ホスト装置とストレージシステムとの間の通信を中継する中継装置(例えばスイッチ装置)であっても良い。

【発明の効果】

【0012】

第1のコントローラの負荷を低減できるように第2のコントローラの資源を利用することができる。

40

【図面の簡単な説明】

【0013】

【図1】図1は、本発明の第1の実施形態にかかるストレージシステムを備えた計算機システムの構成の一例を示すブロック図である。

【図2】図2は、ストレージコントローラ21のブロック図である。

【図3】図3は、SPU25aの詳細ブロック図である。なお、他のSPU25bも同様の構成を有する。

【図4】図4は、MP11の詳細ブロック図である。

50

【図5】図5は、キャッシュメモリ24a及び25bの内部のデータ領域を示したブロック図である。

【図6】図6は、ストレージシステム20がホスト装置10からリード要求を受領した場合の処理（I/O処理のうちのリード処理）の流れを示すラダーチャートである。

【図7】図7は、ストレージコントローラ21がホスト装置10からライト要求を受領した場合の処理（I/O処理のうちのライト処理）の流れを示したラダーチャートである。

【図8】図8は、DMAC251の内部の構造、及び、DMAC251とパケットユニット252及びDMAC制御エリア2431との関係を示すブロック図である。

【図9】図9は、転送パラメータ30およびステータスメッセージ40の詳細を示す。

【図10】図10は、DMAC選択処理のフローチャートである。

10

【図11A】図11Aは、キャッシュデータエリアが固定的に区分されている場合のリード時のキャッシュ利用の一例を示す。

【図11B】図11Bは、キャッシュデータエリアが固定的に区分されている場合のライト時のキャッシュ利用の一例を示す。

【図12A】図12Aは、本発明の第1の実施形態でのリード時のキャッシュ利用の一例を示す。

【図12B】図12Bは、本発明の第1の実施形態でのライト時のキャッシュ利用の一例を示す。

【図13】図13は、本発明の第2の実施形態にかかるストレージシステムのストレージコントローラの内部構成を示すブロック図である。

20

【図14】図14は、本発明の第3の実施形態におけるインデックスを示す。

【図15】図15は、本発明の第4の実施形態におけるインデックス及びパラメータキューの構成を示す。

【発明を実施するための形態】

【0014】

以下、図面を参照して、本発明の幾つかの実施形態を説明する。

【実施例1】

【0015】

図1は、本発明の第1の実施形態にかかるストレージシステムを備えた計算機システムの構成の一例を示すブロック図である。

30

【0016】

計算機システムは、ストレージシステム20と、SANスイッチ30と、複数のホスト装置10とを有する。ストレージシステム20が、SAN(Storage Area Network)を構成するSANスイッチ30を介して、複数のホスト装置10に接続される。例えば、SANスイッチ30は、ストレージシステム20にホストチャンネル40を介して接続されるとともに、複数のホスト装置10にそれぞれチャンネル11を介して接続される。

【0017】

ストレージシステム20は、ストレージコントローラ21と、記憶媒体としての複数のHDD(Hard Disk Drive)23とを有する。ストレージコントローラ21に、HDDチャンネル22(Hard Disk Drive Channel)を介して、複数のHDD23が接続される。HDDチャンネル22は冗長構成にしてもよい。記憶媒体は、HDD23に代えて、フラッシュメモリやDRAM(Dynamic Random Access Memory)などを用いたSSD(Solid State Drive)、或いは、他の記憶媒体でもよい。

40

【0018】

ストレージコントローラ21は、ホスト装置10からのライト要求に従い、ライト要求に付随するライト対象のデータをHDD23に格納する。ストレージコントローラ21は、また、ホスト装置10からのリード要求に従い、リード対象のデータをHDD23から読み出し、そのリード対象のデータをホスト装置10へ送信する。

50

## 【0019】

図2は、ストレージコントローラ21のブロック図である。

## 【0020】

本実施形態では、ストレージコントローラ21は障害に備え冗長構成となっている。例えば、図2に示すように、ストレージコントローラ21は、第1のコントローラ部21a及び第2のコントローラ部21bの二つのコントローラ部を備える。各コントローラ部21a、21bを点線で示す。図示の例では、第1のコントローラ部21aと第2のコントローラ部21bは、同じ構成となっているが、異なった構成となってもよい。本実施形態では、後に説明する通り、稼働中において、第1のコントローラ部21a及び第2のコントローラ部21bのいずれも、自分のコントローラ部内の資源だけでなく相手のコントローラ部内の資源を活用することができるため、第1のコントローラ部21a及び第2のコントローラ部21bをともに稼働させることができる。つまり、ホスト装置が一方のコントローラのみに接続されている場合においても、接続された一方が稼働している間に他方が待機している必要は無く、両方が稼働することができる。

10

## 【0021】

第1のコントローラ部21aには第1のPS(Power Supply)50aから電力が供給され、第2のコントローラ部21bには第2のPS50bから電力が供給される。

## 【0022】

コントローラ部21a(21b)は、SPU(Storage Processing Unit)25a(25b)と、MP(Micro Processor)26a(26b)と、HPC(Host Channel Protocol)212a(212b)と、DPC(Disc Channel Protocol Tip)211a(211b)と、キャッシュメモリ24a(24b)と、を備える。

20

## 【0023】

以下、第1のコントローラ部21aを例にとり、コントローラ部21a及び21bを詳細に説明する。なお、以下の説明では、適宜、第1のコントローラ部21aの要素の名称の前に「第1の」と付し、第2のコントローラ部21bの要素の名称の前に「第2の」と付すことにする。

## 【0024】

第1のSPU25aは、ハードウェア回路、例えばLSI(Large Scale Integration)である。第1のSPU25aは、第1のキャッシュメモリ24aに接続され、これを制御する。第1のSPU25aは、第1のHPC212aを介してホストチャネル40に接続される。第1のHPC212aは、ホストチャネル40の下位レベルのプロトコル処理を行って、第1のSPU25aに接続可能な伝送方式に変換する。第1のSPU25aは、また、第1のDPC211aを介してHDDチャネル22に接続される。第1のDPC211aは、HDDチャネル22の下位レベルのプロトコル処理を行って、第1のSPU25aに接続可能な伝送方式に変換する。

30

## 【0025】

第1のコントローラ部21aは第2のコントローラ部21bに接続される。具体的には、第1のSPU25aと第2のSPU25bとが、2本のSPU間パスC1、C2により接続される。SPU間パスの数は、2より多くても少なくともよい。SPU間パスC1、C2は、キャッシュメモリ24a、24bと後述するDMACとの間のアクセスなどに用いられる。

40

## 【0026】

本実施形態では、第1のMP26aとして、MP11及び12があり、第2のMP26bとして、MP21及び22がある。

## 【0027】

第1のMP11(12)が、SPU-MPパスN11(N12)により第1のSPU2

50



5 aに接続されている。同様に、第2のMP 2 1 ( 2 2 )が、SPU - MPパスN 2 1 ( N 2 2 )により第2のSPU 2 5 bに接続されている。以下、MPとそのMPを有するコントローラ部内のSPUとを結ぶSPU - MPパスを、「ノーマルパス」と呼ぶ。

【 0 0 2 8 】

本実施形態では、一方のコントローラ部内のMPが他方のコントローラ部内のSPUにおける資源を直接利用できるようにするための一つの工夫として、コントローラ部をまたぐSPU - MPパスが用意される。以下、そのSPU - MPパスを「ダイレクトパス」と呼ぶ。具体的には、第1のSPU 2 5 aに、第2のMP 2 1 ( 2 2 )が、ダイレクトパスD 2 1 ( 2 2 )により接続される。同様に、第2のSPU 2 5 bに、第1のMP 1 1 ( 1 2 )が、ダイレクトパスD 1 1 ( 1 2 )により接続される。

10

【 0 0 2 9 】

すなわち、本実施形態によれば、第1のSPU 2 5 aには、第1のMP 1 1 ( 1 2 )がノーマルパスN 1 1 ( N 1 2 )で接続されるとともに、第2のMP 2 1 ( 2 2 )がダイレクトパスD 2 1 ( D 2 2 )で接続されることとなる。同様に、第2のSPU 2 5 bには、第2のMP 2 1 ( 2 2 )がノーマルパスN 2 1 ( N 2 2 )で接続されるとともに、第1のMP 1 1 ( 1 2 )がダイレクトパスD 1 1 ( D 1 2 )で接続されることとなる。

【 0 0 3 0 】

この構成において、第1のMP 1 1 ( 1 2 )は、I/O処理(例えば、ライト要求に従う処理、及び、リード要求に従う処理)において、第1のSPU 2 5 a内の資源(例えば後述のDMAC)又は第1のキャッシュメモリ 2 4 aには、ノーマルパスN 1 1 ( N 1 2 )を介してアクセスし、第2のSPU 2 5 b内の資源(例えば後述のDMAC)又は第2のキャッシュメモリ 2 4 bには、ダイレクトパスD 1 1 ( D 1 2 )を介してアクセスする。また、第1のMP 1 1 ( 1 2 )は、ダイレクトパスD 1 1 ( D 1 2 )の障害が検出された場合、第1のSPU 2 5 a及びSPU間パスC 1又はC 2を介して、第2のSPU 2 5 b内の資源(例えば後述のDMAC)又は第2のキャッシュメモリ 2 4 bにアクセスすることができる。

20

【 0 0 3 1 】

一方、第2のMP 2 1 ( 2 2 )は、I/O処理において、第2のSPU 2 5 b内の資源(例えば後述のDMAC)又は第2のキャッシュメモリ 2 4 bに、ノーマルパスN 2 1 ( N 2 2 )を介してアクセスし、第1のSPU 2 5 a内の資源(例えば後述のDMAC)又は第1のキャッシュメモリ 2 4 aに、ダイレクトパスD 2 1 ( D 2 2 )を介してアクセスする。また、第2のMP 2 1 ( 2 2 )は、ダイレクトパスD 2 1 ( D 2 2 )の障害が検出された場合、第2のSPU 2 5 b及びSPU間パスC 1又はC 2を介して、第1のSPU 2 5 a内の資源(例えば後述のDMAC)又は第1のキャッシュメモリ 2 4 aにアクセスすることができる。

30

【 0 0 3 2 】

図3は、SPU 2 5 aの詳細ブロック図である。なお、他のSPU 2 5 bも同様の構成を有する。

【 0 0 3 3 】

SPU 2 5 aは、DMAC ( Direct Memory Access Controller ) 2 5 1と、パケットユニット 2 5 2と、バッファ 2 5 3と、障害管理ユニット 2 5 7と、メモリコントローラ ( MC ) 2 5 8と、複数のポートと、を備える。

40

【 0 0 3 4 】

各ポートは、プロトコルチップポート ( P - P ) 2 5 4、SPU間パスポート ( I - P ) 2 5 5、及び、MPポート ( M - P ) 2 5 6のいずれかである。I - P 2 5 5は、SPU間パスC 1又はC 2が接続されたポートである。M - P 2 5 6は、ノーマルパス又はダイレクトパスが接続されたポートである。P - P 2 5 4は、プロトコルチップ ( HPC又はDPC ) が接続されたポートである。MC 2 5 8は、キャッシュメモリ 2 4 a、2 4 bへのアクセスの制御を行う回路である。

【 0 0 3 5 】

50

パケットユニット 252 は、S P U 25 が受信したデータについて、宛先アドレスを識別し、適切なポート若しくはコンポーネントへのパケット転送を行う。ホスト装置 10 若しくは H D D 23 から、H P C 212 若しくは D P C 211 がデータを受信した場合には、受信したデータを含んだパケットは、パケットユニット 252 を介してバッファ 253 に格納される（ただしデータ以外の要求などを含んだパケットは M P に転送される）。また、バッファ 253 に格納されたデータがホスト装置 10 若しくは H D D 23 へ転送される場合には、バッファ 253 に格納されたデータを含んだパケットは、パケットユニット 252 を介して H P C 212 若しくは D P C 211 に転送され、その後、H P C 212 若しくは D P C 211 からホスト装置 10 若しくは H D D 23 へ転送される。

【0036】

本実施形態においてバッファ 253 は S P U 25 内に備えているが、構成によってはキャッシュメモリ 24 の一部や、後述する M P 11 ( 12 , 21 , 22 ) 内のメモリ 263 の一部を利用して配しても良い。

【0037】

D M A C 251 は、転送パラメータによる M P 26 a 又は 26 b からの指示に基づいて、バッファ 253 に格納されたデータの、キャッシュメモリ 24 a、24 b への転送を制御する。D M A C 251 は、また、転送パラメータによる M P 26 a、26 b からの指示に基づいて、キャッシュメモリ 24 a、24 b に格納されたデータの、バッファ 253 へ転送を制御する。転送パラメータについては後述する。

【0038】

障害管理ユニット 257 は、S P U 25 内の各資源を有効に使用するために、S P U 25 内の各資源の障害発生を監視するユニットである。障害管理ユニット 257 は、S P U 25 内の各資源（コンポーネント）に接続されている。S P U 25 内の或る資源で障害が発生した場合には、障害管理ユニット 257 は、障害の発生を検知し、各 M P 26 a、26 b に、障害が生じた資源を報告する。例えば、或る M - P 256 に障害が発生した場合、障害が発生した M - P 256 に接続している M P 26 a 又は 26 b は、ノーマルパス又はダイレクトパス経由では、S P U 25 内の資源（例えば D M A C 251）にアクセスすることができなくなってしまう。このため、障害管理ユニット 257 は、障害が発生した M - P 256 に接続されている M P 26 a 又は 26 b が行う処理を、障害が発生していない M - P 256 に接続されている M P 26 a 又は 26 b に引き継がせる。

【0039】

図 4 は、M P 11 の詳細ブロック図である。なお、他の M P 12、21 及び 22 もそれぞれ同様の構成を有する。

【0040】

M P 11 は、C P U 261 と、周辺ユニット 262 と、メモリ 263 と、を備える。

【0041】

C P U 261 は、メモリ 263 からプログラムを読み出し、そのプログラムを実行することで、プログラムに基づく処理を行う。

【0042】

周辺ユニット 262 は、C P U 261 とメモリ 263 に接続されているインターフェース回路である。周辺ユニット 262 には、S P U 25 a に接続されるノーマルパス N11 と、S P U 25 b に接続されるダイレクトパス D11 とが接続される。周辺ユニット 262 は、C P U 261 とメモリ 263 との間の通信や、C P U 261 とノーマルパス N11 又はダイレクトパス D11 を介しての S P U 25 a 又は 25 b との通信を制御する。なお C P U 261 の種類によっては、メモリ 263 が周辺ユニット 262 を介さずに C P U 261 に直接接続される形態や、周辺ユニット 262 が C P U 261 と一体となっている形態であっても構わない。

【0043】

メモリ 263 は、C P U 261（または、S P U 25 a 又は 25 b）からアクセスされる。メモリ 263 は、プログラムエリア 263 と、ワークエリア 2632 と、D M A C エ

10

20

30

40

50

リア 2 6 3 3 と、メールボックス 2 6 3 4 とを備える。

【 0 0 4 4 】

プログラムエリア 2 6 3 には、CPU 2 6 1 が実行するプログラムが格納される。ワークエリア 2 6 3 2 は、CPU 2 6 1 がプログラムを実行するために確保された作業領域である。DMA C エリア 2 6 3 3 には、DMA C 2 5 1 からの転送ステータスが格納される。メールボックス 2 6 3 4 は、他の MP 1 2、2 1 又は 2 2 との通信のために用いられる。

【 0 0 4 5 】

図 5 は、キャッシュメモリ 2 4 a 及び 2 4 b の内部のデータ領域を示したブロック図である。

10

【 0 0 4 6 】

キャッシュメモリ 2 4 a ( 2 4 b ) は、制御情報エリア 2 4 2 a ( 2 4 2 b ) と、転送パラメータエリア 2 4 3 a ( 2 4 3 b ) と、キャッシュデータエリア 2 4 1 a ( 2 4 1 b ) と、を備える。

【 0 0 4 7 】

制御情報エリア 2 4 2 a ( 2 4 2 b ) は、制御情報を記憶する。制御情報は、例えば、下記の ( 1 ) ~ ( 5 ) の情報：

( 1 ) データを格納するスロット 2 4 1 1 a ( 2 4 1 1 b ) が指定されるディレクトリ情報；

( 2 ) スロット 2 4 1 1 a ( 2 4 1 1 b ) の使用状況を示す情報；

20

( 3 ) 使用中のスロット 2 4 1 1 a ( 2 4 1 1 b ) について、どのスロット 2 4 1 1 a ( 2 4 1 1 b ) にどのデータが格納されているかを示すデータ格納情報；

( 4 ) 複数の HDD 2 3 を仮想的に 1 つのボリュームとして提供する RAID 制御の設定情報；

( 5 ) バックアップ機能、スナップショット機能及びリモートコピー機能など、機能に関する情報、を含む。

【 0 0 4 8 】

転送パラメータエリア 2 4 3 a ( 2 4 3 b ) には、MP 2 6 a 又は 2 6 b がセットした転送パラメータが格納される。転送パラメータエリア 2 4 3 に格納された転送パラメータは、DMA C 2 5 1 にフェッチされて実行される。転送パラメータについては後述する。

30

【 0 0 4 9 】

本実施形態では、転送パラメータエリア 2 4 3 a ( 2 4 3 b ) はキャッシュメモリ 2 4 a ( 2 4 b ) に設けられている。しかし、例えば、転送パラメータエリア 2 4 3 a 及び / 又は 2 4 3 b は、他の記憶資源、例えば、バッファ 2 5 3 及び / 又はメモリ 2 6 3 ( MP 2 6 a ( 2 6 b ) 内の記憶資源 ) に設けられてもよい。メモリ 2 6 3 が転送パラメータエリア 2 4 3 a ( 2 4 3 b ) を有する場合、例えば、DMA C エリア 2 6 3 3 に転送パラメータエリア 2 4 3 a ( 2 4 3 b ) が設けられる。

【 0 0 5 0 】

キャッシュデータエリア 2 4 1 a ( 2 4 1 b ) は、複数のスロット 2 4 1 1 a ( 2 4 1 1 b ) を有する。スロット 2 4 1 1 a ( 2 4 1 1 b ) には、ホスト装置 1 0 からのライトコマンドに従うライトデータ又はリードデータが一時的に格納される。図では、ライトデータは、アルファベット「W」と、数字との組合せであらわされている。リードデータは、アルファベット「R」と、数字との組合せであらわされている。

40

【 0 0 5 1 】

キャッシュデータエリア 2 4 1 a ( 2 4 1 b ) には、複数のスロットが設けられており、ライトデータの書き込みエリアとリードデータの書き込みエリアといった複数のエリアが固定的に設けられていない。このため、ライトデータ ( W ) 及びリードデータ ( R ) が、キャッシュデータエリア 2 4 1 a ( 2 4 1 b ) の任意のスロット 2 4 1 1 に格納される。特に、ライトデータ ( W ) は、キャッシュデータエリア 2 4 1 a 及び 2 4 1 b の両方に

50

格納される。つまり、データのいわゆるダブルライト、言いかえれば、データのキャッシュミラーリングが行われる。

【 0 0 5 2 】

図 9 は、転送パラメータ 3 0 およびステータスメッセージ 4 0 の詳細を示す。

【 0 0 5 3 】

まず、転送パラメータ 3 0 について以下に説明する。

【 0 0 5 4 】

転送パラメータ 3 0 は、データの転送に関するパラメータである。転送パラメータ 3 0 は、ID フィールド 3 0 1 と、オペレーションフィールド 3 0 2 と、バッファアドレスフィールド 3 0 3 と、キャッシュメモリアドレスフィールド 3 0 4 と、サイズフィールド 3 0 5 とを備える。

10

【 0 0 5 5 】

ID フィールド 3 0 1 は、MP 2 6 a ( 2 6 b ) に付与された ID が設定される領域である。

【 0 0 5 6 】

オペレーションフィールド 3 0 2 は、オペレーションの種類を表す情報が設定される領域である。オペレーションの種類には、例えば、キャッシュメモリ 2 4 a ( 又は 2 4 b ) からデータを転送する ( 読み出す ) 「リード」、単一のキャッシュメモリ 2 4 a ( 又は 2 4 b ) にデータを転送する ( 書き込む ) 「ライト」、及び、キャッシュメモリ 2 4 a 及び 2 4 b の両方にデータを転送する ( 書き込む ) 「ダブルライト」がある。

20

【 0 0 5 7 】

バッファアドレスフィールド 3 0 3 は、バッファ 2 5 3 のアドレスを表す値 ( バッファアドレス値 ) が設定される領域である。

【 0 0 5 8 】

キャッシュアドレスフィールド 3 0 4 は、キャッシュメモリ 2 4 a 及び / 又は 2 4 b のアドレスを表す値 ( キャッシュアドレス値 ) が設定される領域である。キャッシュメモリ値は、例えば、キャッシュメモリ 2 4 a 及び / 又は 2 4 b の識別情報と、スロット 2 4 1 1 の位置情報とを含んでよい。オペレーションの種類が「ダブルライト」のときには、当該アドレスフィールド 3 0 4 には、キャッシュメモリ 2 4 a 及び 2 4 b の両方のキャッシュアドレス値が設定される。

30

【 0 0 5 9 】

サイズフィールド 3 0 5 は、転送対象のデータのデータサイズ ( データ長 ) を表す情報 ( データサイズ情報 ) が設定されるフィールドである。

【 0 0 6 0 】

なお、ストレージコントローラ 2 1 は、ライトの時に、バッファ 2 5 3 内の連続するライトデータを、複数のキャッシュメモリ 2 4 に分割して格納するスキップ機能や、リードの時に、複数のキャッシュメモリ 2 4 に分割されていたリードデータを、一つのバッファ内に連続して格納するギャザー機能を備えてもよい。これらスキップ機能若しくはギャザー機能を用いる場合には、転送パラメータ 3 0 に、キャッシュアドレスフィールド 3 0 4 及びサイズフィールド 3 0 5 は複数組み合わせられることとなる。

40

【 0 0 6 1 】

次に、ステータスメッセージ 4 0 について以下に説明する。

【 0 0 6 2 】

ステータスメッセージ 4 0 は、転送パラメータ 3 0 に対応するメッセージであり、転送パラメータ 3 0 の実行の結果のステータスを表す。ステータスメッセージ 4 0 は、ID フィールド 4 0 1 と、ステータスフィールド 4 0 2 とを備える。

【 0 0 6 3 】

ID フィールド 4 0 1 は、対応する転送パラメータ 3 0 に設定されていた ID と同じ ID が設定される領域である。従って、ID フィールド 4 0 1 内の ID を用いて、その ID を有するステータスメッセージ 4 0 に対応する転送パラメータ 3 0 を特定することができ

50

る。

【 0 0 6 4 】

ステータスフィールド 4 0 2 は、転送パラメータ 3 0 の実行に関するステータスを表す情報が設定される領域である。ステータスは、例えば、「正常終了」、「転送中キャッシュメモリ障害検出」、「転送中バス障害検出」及び「転送パラメータフォーマット不正」がある。

【 0 0 6 5 】

ステータスメッセージ 4 0 は、DMAC 2 5 1 が転送パラメータ 3 0 の実行が完了したことを検知したときに MP 2 6 a、2 6 b に送信するメッセージである。ステータスメッセージ 4 0 を受信することにより、MP 2 6 a、2 6 b は転送状況を知ることができる。

10

【 0 0 6 6 】

すなわち、MP 2 6 a、2 6 b が、受信していないステータスメッセージ 4 0 に対応する転送パラメータ 3 0 をカウントすることで、未完了の転送パラメータ 3 0 の数を知ることができる。例えば、DMAC 2 5 1 の負荷が高ければ、結果として転送パラメータ 3 0 の処理速度は遅くなる。具体的には、例えば、転送パラメータ 3 0 の実行頻度よりも転送パラメータ 3 0 が発行頻度の方が高ければ、未完了の転送パラメータ 3 0 の数が増えていくことになる。つまり、未完了の転送パラメータ 3 0 の数を知ること、各 MP 2 6 a、2 6 b は、他の MP 2 6 a、2 6 b と通信せずとも DMAC 2 5 1 の負荷状況を知ることができる。

【 0 0 6 7 】

20

図 8 は、DMAC 2 5 1 の内部の構造、及び、DMAC 2 5 1 とパケットユニット 2 5 2 及び DMAC 制御エリア 2 4 3 1 との関係を示すブロック図である。以下の説明では、図 8 に示す DMAC 2 5 1 及びパケットユニット 2 5 2 は、SPU 2 5 a 内の要素であるとする。

【 0 0 6 8 】

まず、DMAC 制御エリア 2 4 3 1 の内部構造について説明する。

【 0 0 6 9 】

DMAC 制御エリア 2 4 3 1 は、キャッシュメモリ 2 4 a における転送パラメータエリア 2 4 3 a の内部に設けられた領域である。DMAC 制御エリア 2 4 3 1 は、複数のパラメータキュー 2 4 3 1 1 を備える。DMAC 制御エリア 2 4 3 1 内の複数のパラメータキュー 2 4 3 1 1 は、ストレージコントローラ 2 1 内の MP に一対一で対応づけられる。つまり、各 DMAC 2 5 1 について、MP 毎にパラメータキュー 2 4 3 1 1 が設けられている。各パラメータキュー 2 4 3 1 1 には、例えば、対応する MP の ID ( 対応 MP - ID ) が付与される。

30

【 0 0 7 0 】

次に、DMAC 2 5 1 の内部構造を説明する。

【 0 0 7 1 】

DMAC 2 5 1 は、インデックス 2 5 1 1 と、セクタ 2 5 1 2 と、パラメータフェッチユニット 2 5 1 3 と、アドレスレジスタ 2 5 1 4 と、カウントレジスタ 2 5 1 5 と、転送制御ユニット 2 5 1 6 とを備える。

40

【 0 0 7 2 】

インデックス 2 5 1 1 は、その DMAC 2 5 1 についてのパラメータキュー 2 4 3 1 1 と同様に、ストレージコントローラ 2 1 内の MP に一対一で対応づけられる。つまり、DMAC 2 5 1 には、MP 毎にインデックス 2 5 1 1 が設けられており、各インデックス 2 5 1 1 は各パラメータキュー 2 4 3 1 1 に対応している。各インデックス 2 5 1 1 には、例えば、対応する MP の ID ( 対応 MP - ID ) が付与される。すなわち、MP 2 6 a、2 6 b が個々に有する ID は、各パラメータキュー 2 4 3 1 1 に付与された対応 MP - ID であるとともに、各インデックス 2 5 1 1 に付与された対応 MP - ID でもある。

【 0 0 7 3 】

パラメータキュー 2 4 3 1 1 には、そのキュー 2 4 3 1 1 に対応した MP 2 6 a 又は 2

50

6 b から転送パラメータ 3 0 が格納される。パラメータキュー 2 4 3 1 1 からは、先に格納された転送パラメータ 3 0 から先に、パラメータフェッチユニット 2 5 1 3 によってフェッチされる。つまり、パラメータキュー 2 4 3 1 1 内の未処理の転送パラメータ 3 0 は、格納された順に処理される。

**【 0 0 7 4 】**

インデックス 2 5 1 1 には、そのインデックス 2 5 1 1 に対応したパラメータキュー 2 4 3 1 1 についての未完了の転送パラメータ 3 0 の数（つまりキュー 2 4 3 1 1 に蓄積されている転送パラメータ 3 0 の数）が記憶されている。例えば、MP 2 6 a 又は 2 6 b により、その MP に対応するパラメータキュー 2 4 3 1 1 に転送パラメータ 3 0 が一つ格納された場合、その MP 対応するインデックス 2 5 1 1 が記憶する値に 1 がインクリメントされる。また、パラメータフェッチユニット 2 5 1 3 により、パラメータキュー 2 4 3 1 1 から一つの転送パラメータ 3 0 がフェッチされた場合、そのパラメータキュー 2 4 3 1 1 に対応するインデックス 2 5 1 1 が記憶する値から 1 がデクリメントされる。

10

**【 0 0 7 5 】**

セクタ 2 5 1 2 は、インクリメントのあったインデックス 2 5 1 1（以下、対象インデックスという）を選択し、選択した対象インデックスの対応 ID をパラメータフェッチユニット 2 5 1 3 に送る。セクタ 2 5 1 2 による対象インデックスの選択は、例えば、パラメータフェッチユニット 2 5 1 3 が、使用可能な状況（「READY」という状況）であった場合に行われる。対象インデックスが複数あった場合、セクタ 2 5 1 2 が所定のルール（例えばラウンドロビン）で一つの対象インデックスを選択する。

20

**【 0 0 7 6 】**

パラメータフェッチユニット 2 5 1 3 は、セクタ 2 5 1 2 から受け取った対応 ID が付与されている対象インデックス 2 5 1 1 から転送パラメータ 3 0 をフェッチする。パラメータフェッチユニット 2 5 1 3 は、フェッチした転送パラメータ 3 0 内のバッファアドレス値及びキャッシュアドレス値を、アドレスレジスタ 2 5 1 4 にセットし、且つ、その転送パラメータ 3 0 内のデータサイズ値をカウントレジスタ 2 5 1 5 にセットする。また、転送制御ユニット 2 5 1 6 が使用可能な状況であった場合には、パラメータフェッチユニット 2 5 1 3 は、転送制御ユニット 2 5 1 6 に対して、例えば、転送パラメータ 3 0 のオペレーション種類及び上記受け取った対応 ID を伝えるとともに、そのオペレーション種類に応じた転送開始のトリガーをかける。トリガーとしては、リードのトリガーとライトのトリガーとがある。

30

**【 0 0 7 7 】**

転送制御ユニット 2 5 1 6 は、パラメータフェッチユニット 2 5 1 3 から転送開始のトリガーを受けた場合に、キャッシュメモリとバッファとの間の転送制御を行う。

**【 0 0 7 8 】**

例えば、リードのトリガーがかけられた場合には、転送制御ユニット 2 5 1 6 は、レジスタ 2 5 1 4 に設定されているキャッシュメモリ値が表すスロットに格納されているデータ（リードデータ）を、そのスロットから、レジスタ 2 5 1 4 に設定されているバッファアドレス値が表すバッファ領域（バッファ 2 5 3 内の領域）に転送する。また、例えば、ライトのトリガーがかけられた場合には、転送制御ユニット 2 5 1 6 は、レジスタ 2 5 1 4 に設定されているバッファアドレス値が表すバッファ領域に格納されているデータ（ライトデータ）を、そのバッファ領域から、レジスタ 2 5 1 4 に設定されているキャッシュメモリ値が表すスロットに転送する。このデータ転送が行われる都度に、転送制御ユニット 2 5 1 6 によって、レジスタ 2 5 1 5 に記憶されているデータサイズ値から、転送されたデータのサイズ分の値がデクリメントされる。つまり、このデータ転送は、レジスタ 2 5 1 5 に記憶されているデータサイズ値分のデータが転送されるまで繰り返される。レジスタ 2 5 1 5 に記憶されているデータサイズ値分のデータの転送が完了したときに（例えばレジスタ 2 5 1 5 に記憶されているデータサイズ値がゼロになったときに）、転送制御ユニット 2 5 1 6 が、処理した転送パラメータ 3 0 に対応するステータスメッセージ 4 0 を生成し、そのメッセージ 4 0 を、MP に送信する。具体的には、例えば、転送制御ユニ

40

50

ット2516が、パラメータフェッチユニット2513から受け取った対応IDを含んだステータスメッセージ40を作成し、そのメッセージ40を、その対応IDを有するMPに送信する。このデータ転送では、下記の三パターン：

(P1) 転送元の記憶領域が、図示のDMAC251を有するコントローラ部に存在し、転送先の記憶領域も、図示のDMAC251を有するコントローラ部に存在し、それ故、SPU間でのデータ転送が生じないパターン；

(P2) 転送元の記憶領域が、図示のDMAC251を有するコントローラ部に存在し、転送先の記憶領域が、図示のDMAC251を有するコントローラ部とは別のコントローラ部に存在するパターン；

(P3) 転送元の記憶領域が、図示のDMAC251を有するコントローラ部とは別のコントローラ部に存在し、転送先の記憶領域が、図示のDMAC251を有するコントローラ部に存在するパターン、

がある。転送元と転送先の一方向の記憶領域は、キャッシュデータエリア241a又は241b内のスロットであり、他方の記憶領域は、バッファ領域(又はMP内のメモリなどの他の記憶資源)である。

【0079】

以下、本実施形態で行われる処理の流れを説明する。

【0080】

図6は、ストレージシステム20がホスト装置10からリード要求を受領した場合の処理(I/O処理のうちのリード処理)の流れを示すラダーチャートである。以下の説明では、HPC212aがリード要求を受領するとする。また、図6及び図7の説明では、どの要素がどちらのコントローラ部に存在するのかを区別し易くするため、第1のコントローラ部21aに存在する要素の参照番号の末尾に「a」を付し、第2のコントローラ部21bに存在する要素の参照番号の末尾に「b」を付すことにする。

【0081】

HPC212aは、リード要求をホスト装置10から受領する(s1000)。以下、リード要求に従うデータを「該当リードデータ」と言う。

【0082】

HPC212aは、ホスト装置10から受領したリード要求を、MP26aのメールボックス2634aへ送付する(s1001)。

【0083】

メールボックス2634a内のリード要求を確認したCPU261aは、制御情報エリア242a及び242b内の制御情報を参照し、キャッシュヒットしたか否かを判定する。キャッシュヒットしたか否かとはい、ホスト装置からのI/O要求で指定されている論理アドレス(例えば、LUN(Logical Unit Number)及びLBA(Logical Block Address))に従う場所に格納されるデータがあるか否かである。キャッシュヒットした場合(ここでは、リード要求が指定する論理アドレスとスロット2411a又は2411bとの対応が管理されている場合)、S1009に処理が移る。

【0084】

一方、キャッシュヒットしなかった場合、CPU261aは、該当リードデータをバッファ253a(又はメモリ263a等の他の記憶資源)に転送するようDPC211aに指示する(s1002)。その際、CPU261aは、例えば、リード要求が指定する論理アドレスを基に特定した物理アドレスを、DPC211aに通知する。

【0085】

DPC211aは、HDD23(例えば通知された物理アドレス)にアクセスして該当リードデータを受け取る(s1003)。

【0086】

s1004で、DPC211aは、受け取った該当リードデータをバッファ253a(又はメモリ263a等の他の記憶資源)に格納する。そして、CPU261aは、キャッシュメモリ24a及び24bにおける複数の確保可能スロットから任意のスロット241

10

20

30

40

50

1 a又は2 4 1 1 bを確保する。確保可能スロットとは、例えば、フリー又はクリーンのスロットとして管理されているスロットである。フリーのスロットとは、空きのスロットである。クリーンのスロットとは、HDD 2 3に格納済みのデータを記憶しているスロットである。図6の説明では、キャッシュメモリ2 4 bからスロット2 4 1 1 bが確保されたとする。

【0087】

DPC 2 1 1 aは、CPU 2 6 1 aに対して、該当リードデータのバッファ2 5 3 aへの格納が終了したことを知らせる終了通知を送信する(s 1 0 0 5)。

【0088】

s 1 0 0 6で、以下の処理(s 1 0 0 6 - 1) ~ (s 1 0 0 6 - 4)が行われる。

10

【0089】

(s 1 0 0 6 - 1) 終了通知を受けたCPU 2 6 1 aが、DMAC選択処理を行う。すなわち、CPU 2 6 1 aは、SPU 2 5 a及び2 5 bに存在する複数のDMAC 2 5 1 a及び2 5 1 bから一つのDMAC 2 5 1 a又は2 5 1 bを選択する。DMAC選択処理では、後述するように、データの転送先の記憶領域を有するコントローラ部内のDMACが、データの転送先の記憶領域を有しないコントローラ部内のDMACよりも優先的に選択される。そのため、ここでは、DMAC 2 5 1 bが選択されたとする。

【0090】

(s 1 0 0 6 - 2) CPU 2 6 1 aは、下記(R 0 1) ~ (R 0 5)を含んだ転送パラメータ30：

20

(R 0 1) ID；

(R 0 2) オペレーション種類「ライト」(キャッシュメモリにデータを書くため)、

(R 0 3) 確保したスロットのアドレス値、

(R 0 4) 該当リードデータを格納しているバッファ領域のアドレス値；

(R 0 5) 該当リードデータのデータサイズ値、

を作成する。

【0091】

(s 1 0 0 6 - 3) CPU 2 6 1 aは、その転送パラメータ30を、s 6 0 0 0で選択したDMAC 2 5 1 bについての複数のパラメータキュー2 4 3 1 1 bのうちの、CPU 2 6 1 aを有するMP 2 6 aに対応したパラメータキュー2 4 3 1 1 bにセットする。

30

【0092】

(s 1 0 0 6 - 4) CPU 2 6 1 aは、上記選択したDMAC 2 5 1 bを起動する。

【0093】

起動したDMAC 2 5 1 bは、DMAC制御エリア2 4 3 1 b内のパラメータキュー2 4 3 1 1 bから転送パラメータ30をフェッチし、その転送パラメータ30に基づき、転送元のバッファ2 5 3 aから転送先のスロット2 4 1 1 bに該当リードデータを転送する(s 1 0 0 7)。

【0094】

該当リードデータがスロット2 4 1 1 bに格納された後、DMAC 2 5 1 bは、転送が終了した旨のステータスメッセージ40を、MP 2 6 a内のDMACエリア2 6 3 3に、例えばダイレクトパスD 1 1又はD 1 2を介して送信する(s 1 0 0 8)。

40

【0095】

ステップS 1 0 0 9は、ステップS 1 0 0 8でDMACエリア2 6 3 3 aに格納されたステータスメッセージ40をCPU 2 6 1 aが確認したとき、若しくは、ステップS 1 0 0 2においてキャッシュヒットしたときに、実行される。s 1 0 0 9で、以下の処理(s 1 0 0 9 - 1) ~ (s 1 0 0 9 - 4)が行われる。

【0096】

(s 1 0 0 9 - 1) CPU 2 6 1 aが、DMAC選択処理を行う。ここでは、例えば、DMAC 2 5 1 aが選択されたとする。

【0097】

50



( s 1 0 0 9 - 2 ) C P U 2 6 1 a は、下記 ( R 1 1 ) ~ ( R 1 5 ) を含んだ転送パラメータ 3 0 :

( R 1 1 ) I D ;

( R 1 2 ) オペレーション種類「リード」(キャッシュメモリからデータを読むため)、

( R 1 3 ) 該当リードデータを格納しているスロット 2 4 1 1 b のアドレス値、

( R 1 4 ) バッファ領域のアドレス値 ;

( R 1 5 ) 該当リードデータのデータサイズ値、

を作成する。

【 0 0 9 8 】

( s 1 0 0 9 - 3 ) C P U 2 6 1 a は、その転送パラメータ 3 0 を、 D M A C 2 5 1 a 10  
についての複数のパラメータキュー 2 4 3 1 1 a のうちの、 C P U 2 6 1 a を有する M P  
2 6 a に対応したパラメータキュー 2 4 3 1 1 a にセットする。

【 0 0 9 9 】

( s 1 0 0 9 - 4 ) C P U 2 6 1 a は、上記選択された D M A C 2 5 1 a を起動する。

【 0 1 0 0 】

起動した D M A C 2 5 1 a は、 D M A C 制御エリア 2 4 3 1 a 内のパラメータキュー 2  
4 3 1 1 a から転送パラメータ 3 0 をフェッチし、その転送パラメータ 3 0 に基づき、ス  
ロット 2 4 1 1 b からバッファ領域へ該当リードデータを転送する ( s 1 0 1 0 ) 。

【 0 1 0 1 】

転送終了後、 D M A C 2 5 1 a は、転送が終了した旨のステータスメッセージ 4 0 を M 20  
P 2 6 a の D M A C エリア 2 6 3 3 a に送信する ( s 1 0 1 1 ) 。

【 0 1 0 2 】

ステップ s 1 0 1 1 で D M A C エリア 2 6 3 3 a に格納されたステータスメッセージ 4  
0 を確認した C P U 2 6 1 a は、バッファ 2 5 3 に格納されている該当リードデータをホ  
スト装置 1 0 に転送するよう H P C 2 1 2 に指示する ( s 1 0 1 2 ) 。

【 0 1 0 3 】

指示を受けた H P C 2 1 2 は、バッファ 2 5 3 より該当リードデータを読み出す ( s 1  
0 1 3 ) 。

【 0 1 0 4 】

H P C 2 1 2 は、バッファ 2 5 3 から読みだした該当リードデータをホスト装置 1 0 へ 30  
転送する ( s 1 0 1 4 ) 。

【 0 1 0 5 】

全ての該当リードデータの転送が正常に終了した場合は、 H P C 2 1 2 は、その旨 ( 正  
常終了 ) をホスト装置 1 0 に応答する ( s 1 0 1 5 ) 。

【 0 1 0 6 】

さらに、 H P C 2 1 2 は、ホスト装置 1 0 に応答したことを M P 2 6 a に通知する ( s  
1 0 1 6 ) 。これにより、 M P 2 6 a は、リード要求の処理が終了したことを認識する。

【 0 1 0 7 】

図 7 は、ストレージコントローラ 2 1 がホスト装置 1 0 からライト要求を受領した場合  
の処理 ( I / O 処理のうちのライト処理 ) の流れを示したラダーチャートである。以下の 40  
説明では、 H P C 2 1 2 a がライト要求を受領するとする。

【 0 1 0 8 】

H P C 2 1 2 a は、ライト要求をホスト装置 1 0 から受領する ( s 2 0 0 0 ) 。以下、  
ライト要求に従うデータを「該当ライトデータ」と言う。

【 0 1 0 9 】

H P C 2 1 2 a は、ホスト装置 1 0 から受領したライト要求を、 M P 2 6 a のメールボ  
ックス 2 6 3 4 a へ送付する ( s 2 0 0 1 ) 。

【 0 1 1 0 】

メールボックス 2 6 3 4 a 内のライト要求を確認した C P U 2 6 1 a は、制御情報エリ  
ア 2 4 2 a 及び 2 4 2 b の制御情報を参照し、キャッシュヒットしたか否かを判定する。 50

## 【 0 1 1 1 】

ここで、該当ライトデータは、キャッシュメモリ 2 4 a 及び 2 4 b の両方に格納されるべきデータ、つまり、ダブルライトされるべきデータである。

## 【 0 1 1 2 】

そのため、キャッシュヒット判定では、ライト要求が指定する論理アドレスと対応するスロット 2 4 1 1 a 及び 2 4 1 1 b があるか否かが判定される。キャッシュヒットした場合、それら対応するスロット 2 4 1 1 a 及び 2 4 1 1 b が確保される。キャッシュヒットしなかった場合、CPU 2 6 1 a は、キャッシュデータエリア 2 4 1 a 及び 2 4 1 b からスロット 2 4 1 1 a 及び 2 4 1 1 b を確保する。

## 【 0 1 1 3 】

スロット 2 4 1 1 a 及び 2 4 1 1 b を確保できた場合は、CPU 2 6 1 a は、該当ライトデータの受け入れ準備ができた旨をホスト装置 1 0 に伝えるよう HPC 2 1 2 a に指示する ( s 2 0 0 2 ) 。

## 【 0 1 1 4 】

CPU 2 6 1 a からの指示を受けた HPC 2 1 2 a は、該当ライトデータの受け入れ準備完了した旨のメッセージをホスト装置 1 0 に送信する ( s 2 0 0 3 ) 。

## 【 0 1 1 5 】

ホスト装置 1 0 は、該当ライトデータを HPC 2 1 2 a へ送信する ( s 2 0 0 4 ) 。

## 【 0 1 1 6 】

該当ライトデータを受領した HPC 2 1 2 a は、その該当ライトデータをバッファ 2 5 3 a に格納する ( s 2 0 0 5 ) 。

## 【 0 1 1 7 】

該当ライトデータをバッファ 2 5 3 a に格納した後、HPC 2 1 2 a は、CPU 2 6 1 a に対して、該当ライトデータのバッファ 2 5 3 a への格納が終了したことを知らせる終了通知を送信する ( s 2 0 0 6 ) 。

## 【 0 1 1 8 】

s 2 0 0 7 で、以下の処理 ( s 2 0 0 7 - 1 ) ~ ( s 2 0 0 7 - 4 ) が行われる。

## 【 0 1 1 9 】

( s 2 0 0 7 - 1 ) 終了通知を受けた CPU 2 6 1 a が、DMAC 選択処理を行う。ここでは、DMAC 2 5 1 a が選択されたとする。

## 【 0 1 2 0 】

( s 2 0 0 7 - 2 ) CPU 2 6 1 a は、下記 ( W 0 1 ) ~ ( W 0 5 ) を含んだ転送パラメータ 3 0 :

( W 0 1 ) ID ;

( W 0 2 ) オペレーション種類「ダブルライト」(ライト要求に従うデータのキャッシュメモリへのデータ転送であるため)、

( W 0 3 ) 確保したスロット 2 4 1 1 a 及び 2 4 1 1 b のアドレス値、

( W 0 4 ) 該当ライトデータを格納しているバッファ領域のアドレス値 ;

( W 0 5 ) 該当ライトデータのデータサイズ値、

を作成する。

## 【 0 1 2 1 】

( s 2 0 0 7 - 3 ) CPU 2 6 1 a は、その転送パラメータ 3 0 を、s 6 0 0 0 で選択した DMAC 2 5 1 a についての複数のパラメータキュー 2 4 3 1 1 a のうちの、CPU 2 6 1 a を有する MP 2 6 a に対応したパラメータキュー 2 4 3 1 1 a にセットする。

## 【 0 1 2 2 】

( s 2 0 0 7 - 4 ) CPU 2 6 1 a は、上記選択した DMAC 2 5 1 a を起動する。

## 【 0 1 2 3 】

起動した DMAC 2 5 1 a は、DMAC 制御エリア 2 4 3 1 a 内のパラメータキュー 2 4 3 1 1 a から転送パラメータ 3 0 をフェッチし、その転送パラメータ 3 0 に基づき、バッファ 2 5 3 からスロット 2 4 1 1 a 及び 2 4 1 1 b に該当ライトデータを転送する ( s

10

20

30

40

50

2008)。

【0124】

両データ転送が終了した後、DMAC251aは、転送が終了した旨のステータスメッセージ40をMP26aのDMACエリア2633aに送信する(s2009)。

【0125】

ステップS2009でDMACエリア2633aに格納されたステータスメッセージ40を確認したCPU261aが、HPC212aにライト終了を通知する(s2010)。

【0126】

ライト終了の通知を受けたHPC212aは、ライト終了をホスト装置10に応答する(s2011)。

【0127】

その後、s2012で、以下の処理(s2012-1)～(s2012-4)が行われる。

【0128】

(s2012-1)CPU261aは、該当ライトデータが格納されているスロット2411a及び2411bのうちのいずれかを選択する。ここではスロット2411aが選択されたとする。

【0129】

(s2012-2)CPU261aは、DMAC選択処理を行う。ここでは、例えば、DMAC251aが選択されたとする。

【0130】

(s2012-3)CPU261aは、下記(W11)～(W15)を含んだ転送パラメータ30：

(W11)ID；

(W12)オペレーション種類「リード」(キャッシュメモリからデータを読むため)、

(W13)該当ライトデータを格納しているスロット2411aのアドレス値、

(W14)バッファ領域のアドレス値；

(W15)該当ライトデータのデータサイズ値、

を作成する。

【0131】

(s2012-4)CPU261aは、その転送パラメータ30を、DMAC251aについての複数のパラメータキュー24311aのうちの、CPU261aを有するMP26aに対応したパラメータキュー24311aにセットする。

【0132】

(s2012-5)CPU261aは、上記選択されたDMAC251aを起動する。

【0133】

起動したDMAC251aは、DMAC制御エリア2431a内のパラメータキュー24311aから転送パラメータ30をフェッチし、その転送パラメータ30に基づき、スロット2411aからバッファ領域へ該当ライトデータを転送する(s2013)。

【0134】

転送終了後、DMAC251aは、転送が終了した旨のステータスメッセージ40を、MP26aのDMACエリア2633aへ送信する(s2014)。

【0135】

s2014でDMACエリア2633aに格納されたステータスメッセージ40を確認したCPU261aは、バッファ領域に格納された該当ライトデータをHDD23へ転送するようDPC211aに指示する(s2015)。

【0136】

CPU261aからの指示を受けたDPC211aは、バッファ領域より該当ライトデータを読み出す(s2016)。

10

20

30

40

50

## 【 0 1 3 7 】

D P C 2 1 1 a は、該当ライトデータを H D D 2 3 に格納する ( s 2 0 1 7 ) 。

## 【 0 1 3 8 】

全てのデータ転送が正常に終了した場合は、D P C 2 1 1 a はその旨を C P U 2 6 1 a に終了通知を送付する ( s 2 0 1 8 ) 。

## 【 0 1 3 9 】

終了通知を受けた C P U 2 6 1 a は、該当ライトデータが格納されていたスロット 2 4 1 1 a を開放する。解放されたスロット 2 4 1 1 a は、クリーンのスロットとして管理される。

## 【 0 1 4 0 】

図 1 0 は、D M A C 選択処理のフローチャートである。なお、以下の説明では、M P 2 6 a 及び 2 6 b を「M P 2 6」と総称する。

## 【 0 1 4 1 】

s 3 0 0 1 で、M P 2 6 が、D M A C 2 5 1 毎の D M A C スコアをすべてリセットする。

## 【 0 1 4 2 】

s 3 0 0 2 で、M P 2 6 は、転送先 ( 又は転送元 ) のバッファ 2 5 3 の位置に応じて D M A C スコアの変更を行う。具体的には、例えば、M P 2 6 a は、転送先 ( 又は転送元 ) のバッファ 2 5 3 を含んだ S P U 2 5 内の D M A C 2 5 1 の D M A C スコアを変更せず、転送先のバッファ 2 5 3 を含んでいない S P U 2 5 内の D M A C 2 5 1 の D M A C スコアに第 1 の所定値を加算する。

## 【 0 1 4 3 】

s 3 0 0 3 で、M P 2 6 は、転送先 ( 又は転送元 ) のキャッシュメモリの位置に応じて D M A C スコアの変更を行う。具体的には、例えば、M P 2 6 は、転送先 ( 又は転送元 ) のキャッシュメモリが接続されている S P U 2 5 内の D M A C 2 5 1 の D M A C スコアを変更せず、転送先のキャッシュメモリが接続されていない S P U 2 5 内の D M A C 2 5 1 の D M A C スコアに第 2 の所定値を加算する。第 2 の所定値は第 1 の所定値と同じでも異なっても良い。

## 【 0 1 4 4 】

次に、M P 2 6 は、各パラメータキュー 2 4 3 1 1 における未完了の転送パラメータ 3 0 の数をチェックする ( s 3 0 0 4 ) 。ここでは、M P 2 6 は、各 D M A C について、その D M A C についての未完了の転送パラメータの数と第 3 の所定値 ( 0 より大きい値 ) との乗数を、その D M A C の D M A C スコアに加算する。

## 【 0 1 4 5 】

最後に、M P 2 6 は、すべての D M A C スコアを比較し、D M A C 2 5 1 を選択する ( s 3 0 0 5 ) 。具体的には、例えば、D M A C スコアが最少の D M A C 2 5 1 が選択される。D M A C スコアが最少の D M A C 2 5 1 は、他の D M A C 2 5 1 よりも転送先記憶領域に近く、及び / 又は、他の D M A C 2 5 1 よりも負荷が低い。そのため、D M A C スコアが最少の D M A C 2 5 1 が、D M A C 選択処理を行っている時点において、データ転送を行うのに最適な D M A C 2 5 1 であると考えられる。

## 【 0 1 4 6 】

さらに、D M A C 2 5 1 に障害が発生した場合に、障害が発生していない資源に処理を引き継がせることなどにより、上記第 1 ~ 第 3 の所定値の少なくとも一つがいずれかの M P 2 6 によってより小さい値にされても良い。すなわち、S P U 2 5 内の正常な D M A C 2 5 1 の数に不均衡が発生しても、スコアの重み付けを変えることで、引き続き S P U 2 5 内の資源を有効に使うことが可能となる。

## 【 0 1 4 7 】

本実施形態では、コントローラ部 2 1 a 及び 2 1 b の一方が他方のコントローラ部内のハードウェア資源を利用している。これにより、ストレージコントローラ 2 1 の負荷を低減することができる。

10

20

30

40

50

## 【 0 1 4 8 】

また、本実施形態では、第1のMP 1 1 ( 1 2 ) が、ダイレクトバスD 1 1 ( 1 2 ) を介して第2のSPU 2 5 bに接続され、第2のMP 2 1 ( 2 2 ) が、ダイレクトバスD 2 1 ( 2 2 ) を介して接続される。このため、コントローラ部 2 1 a 及び 2 1 b の一方が他方のコントローラ部内のハードウェア資源を利用する際に、SPU間バスC 1 及びC 2 の経路を不要にすることができる。

## 【 0 1 4 9 】

また、本実施形態では、一つのDMACが複数のMPから使用されることにより生じ得る下記課題を、ソフトウェア制御ではなくハードウェア制御によって解決される。

## 【 0 1 5 0 】

例えば、或るMP ( x ) が転送パラメータをパラメータキューにセットした後、インデックスが更新される前に、他のMP ( y ) が転送パラメータを上書きしてしまうことがあり得る。この場合は、MP x のセットした転送パラメータは実行されない。その上、転送パラメータが実行されなかったことをMP ( x ) は検知できない。このため、バッファやキャッシュメモリに新たなデータが反映されないまま、ホスト装置に対して処理の完了を通知してしまうことになる。つまり、データ消失や誤ったデータの返送といった致命的な不具合がおこってしまう可能性がある。

## 【 0 1 5 1 】

この課題を、ソフトウェアでの制御によって解決することも考え得る。ここで、ソフトウェアの制御によって解決するとは、例えば、MP間であらかじめ調停を行って権利を取得したMPが転送パラメータをセットするなどである。しかしながら、ソフトウェアの制御によって上記不具合を解決した場合には、頻繁なアクセスが必要なDMACなどの資源に関しては、相対的にオーバーヘッドが大きく、資源の有効活用ができなくなってしまうという更なる課題が生じるおそれがある。

## 【 0 1 5 2 】

本実施形態では、パラメータキュー 2 4 3 1 1 及びインデックス 2 5 1 1 がMP 2 6 毎に設けられていることで、ハードウェアの制御によって、上記のようなデータ消失や誤ったデータの返送といった不具合を防止することができる。

## 【 0 1 5 3 】

本実施形態にあっては、各SPU 2 5 a、2 5 b内に障害管理ユニット 2 5 7 がある。このため、障害管理ユニット 2 5 7 が、SPU 2 5 a、2 5 b内で障害が発生した資源に対して、対策を施すことが可能となる。例えば、障害管理ユニット 2 5 7 が、障害が発生したMPポート 2 5 6 に接続されたMP 2 6 a、2 6 bが行う処理を、障害が発生していないMPポート 2 5 6 に接続されたMP 2 6 a、2 6 bに引き継がせることが可能となる。これにより、SPU 2 5 中の資源を有効に使うことができる。

## 【 0 1 5 4 】

本実施形態にあっては、キャッシュメモリ 2 4 a、2 4 bのキャッシュデータエリア 2 4 1 は、下記に示すように、ライトデータの書き込みエリアとリードデータの書き込みエリアを予め固定的に用意しておく必要はない。

## 【 0 1 5 5 】

図 1 1 A 及び図 1 1 B によれば、キャッシュデータエリア ( A ) に、リードキャッシュエリア ( A ) と、ライトキャッシュエリア ( A ) と、ミラーエリア ( A ) とが予め用意されている。同様に、キャッシュデータエリア ( B ) に、リードキャッシュエリア ( B ) と、ライトキャッシュエリア ( B ) と、ミラーエリア ( B ) とが予め用意されている。ライトキャッシュエリア ( A ) とミラーエリア ( B ) はペアになっており、同様に、ライトキャッシュエリア ( B ) とミラーエリア ( A ) はペアになっている。

## 【 0 1 5 6 】

リードキャッシュエリアは、リードデータが格納される専用エリアである。図 1 1 A に示すように、リードキャッシュエリア ( A ) には、リードキャッシュエリア ( A ) と同じコントローラ部内にあるMP ( A ) によって取得されたリードデータ ( R 1 ) ~ ( R 1 0

10

20

30

40

50

)のみが格納され、ライトデータや、リードキャッシュエリア(A)と別のコントローラ部内にあるMP(B)によって取得されたリードデータは格納されない。

【0157】

ライトキャッシュエリアは、ライトデータが格納される専用エリアであり、それとペアになったミラーエリアは、ライトデータのミラーが格納される専用エリアである。図11Bに示すように、ライトキャッシュエリア(A)には、ライトキャッシュエリア(A)と同じコントローラ部内にあるMP(A)によって取得されたライトデータ(W1)~(W5)のみが格納され、リードデータや、他のMP(B)が受けたライトデータや、そのライトデータのミラーデータは格納されない。また、図11Bに示すように、ミラーエリア(B)には、MP(A)によって取得されたライトデータのミラー(W1)~(W5)のみが格納され、リードデータや、MP(A)及び(B)が受けたライトデータ(オリジナルデータ)は格納されない。

10

【0158】

しかし、本実施形態では、図11A及び図11Bを参照して説明したリードキャッシュエリア、ライトキャッシュエリア及びミラーエリアのいずれも、キャッシュデータエリアに用意されていない。本実施形態では、データの種類に限らず、キャッシュデータエリアの任意のスロットに格納することができる。そのため、図12Aに示すように、例えば、第1のMP24aが取得したリードデータ(R1)~(R40)の一部を第1のキャッシュメモリ24aに格納し残りを第2のキャッシュメモリ24bに格納することができる。また、図12Bに示すように、キャッシュメモリ24a及び24bにおけるそれぞれのキャッシュデータエリアを、第1のMP24aが取得したライトデータ(W1)~(W20)で満杯にすることもできる。

20

【実施例2】

【0159】

図13は、本発明の第2の実施形態にかかるストレージシステムのストレージコントローラの内部構成を示すブロック図である。

【0160】

本実施形態では、ノーマルパス及びダイレクトパスがMPスイッチ(MP Switch)270で実現されている。MPスイッチ270は、複数のMP26と複数のSPU25が接続されたスイッチ装置である。各MP26は、MPスイッチ270を介して所望のSPU25にアクセスすることができる。

30

【0161】

また、本実施形態では、SPU間パスがSPUスイッチ(SPU Switch)280で実現されている。複数のSPU25のうちの一つと複数のSPU25のうちの一つとの間の通信は、SPUスイッチ280を介して行われる。

【0162】

SPU25の数、MP26の数、キャッシュメモリ24の数を、スケラブルにする場合には、本実施形態のような構成をとることができる。この場合、MPスイッチ270とSPUスイッチ280は論理的な構成としてもよく、物理的には同一のスイッチで構成しても良い。例えばPCI-Expressのvirtual channel技術が用いられてもよい。またMPスイッチ270およびSPUスイッチ280は、物理的には複数のスイッチをカスケード接続した形態でも良い。

40

【0163】

また、さらにキャッシュメモリのスケラビリティを向上させるため、キャッシュメモリを独立させてもよい。この場合、例えば、キャッシュメモリ機能のみを持ったSPUがSPUスイッチに接続されてもよい。

【実施例3】

【0164】

図14は、本発明の第3の実施形態におけるインデックスを示す。

【0165】

50

インデックス2511が、コンシューマインデックス25111と、プロデューサインデックス25112とで構成されている。

【0166】

コンシューマインデックス25111は、パラメータキューに蓄積されている一以上の転送パラメータのうちどの転送パラメータまでフェッチしたかを表す値を記憶する。

【0167】

プロデューサインデックス25112は、パラメータキューに蓄積されている最後尾の転送パラメータを表す値である。

【0168】

このようなインデックス2511を用いて、パラメータキューのどこから転送パラメータを読み出すべきかを特定することができる。

10

【実施例4】

【0169】

図15は、本発明の第4の実施形態におけるインデックス及びパラメータキューの構成を示す。

【0170】

MP毎に、リード用パラメータキュー26331Rとライト用パラメータキュー26331Wとがある。また、MP毎に、リード用パラメータキュー26331Rに対応したリード用インデックス2511Rと、ライト用パラメータキュー26331Wに対応したライト用インデックス2511Wとがある。

20

【0171】

MPは、リード要求での処理において生成した転送パラメータを、リード用パラメータキュー26331Rに格納する。リード用インデックス2511Rは、リード用パラメータキュー26331Rでの未処理の転送パラメータの数を表す値を記憶する。

【0172】

MPは、ライト要求での処理において生成した転送パラメータを、ライト用パラメータキュー26331Wに格納する。ライト用インデックス2511Wは、ライト用パラメータキュー26331Wでの未処理の転送パラメータの数を表す値を記憶する。

【0173】

例えば、ライト要求の処理では、該当ライトデータがキャッシュメモリに格納されれば、ホスト装置に対してライトの完了とすることができるが、リード要求の処理では、該当リードデータがキャッシュメモリに格納されるだけでなくホスト装置に転送されてリードの完了となる。

30

【0174】

そこで、本実施形態では、DMACは、ライト用パラメータキュー26331Wに格納されている転送パラメータよりも、リード用パラメータキュー26331Rに格納されている転送パラメータを優先的にフェッチする。これにより、リード性能を高めることができる。

【0175】

以上、本発明の幾つかの実施形態を説明したが、本発明は、これらの実施形態に限定されるものでなく、その要旨を逸脱しない範囲で種々変更可能であることはいうまでもない。

40

【0176】

上述した実施形態に従う記憶制御装置を別の観点から抽象的に表現すれば、例えば、下記の通りになる。

<表現1>

ホスト装置から発行されたI/Oコマンドに従うデータの転送を制御する回路である中継回路と、

前記中継回路に接続された複数のプロセッサと

プロセッサ毎に設けられておりデータ転送のパラメータである転送パラメータが蓄積さ

50

れるパラメータ領域と  
を有し、

前記中継回路が、データの転送を行う回路であるデータ転送回路を有し、  
前記パラメータ領域は、前記データ転送回路についての記憶領域であり、  
各プロセッサが、そのプロセッサに対応したパラメータ領域に転送パラメータを格納し

、  
前記データ転送回路が、いずれかのパラメータ領域から転送パラメータを取得し、その  
転送パラメータに従って、データ転送を実行する、  
記憶制御装置。

< 表現 2 >

10

表現 1 記載の記憶制御装置であって、

前記データ転送回路が複数あり、

各データ転送回路について、プロセッサ毎のパラメータ領域があり、

前記複数のプロセッサのうちの I / O コマンドを受領したプロセッサである対象プロセ  
ッサが、前記対象プロセッサに対応した複数のパラメータ領域のうち未完了の転送パラ  
メータの数が最も少ないパラメータ領域に対応したデータ転送回路を選択し、選択したデ  
ータ転送回路についての、前記対象プロセッサに対応したパラメータ領域に、転送パラメ  
ータを格納する、

記憶制御装置。

< 表現 3 >

20

表現 1 又は 2 記載の記憶制御装置であって、

前記複数のプロセッサにそれぞれ対応した複数のキャッシュメモリ領域を有し、

前記複数のプロセッサのうちの或るプロセッサが、前記複数のプロセッサのうちの別の  
プロセッサに対応したキャッシュメモリ領域内の任意のアドレスを転送先アドレスとした  
転送パラメータを、前記或るプロセッサに対応したパラメータ領域に格納する、

記憶制御装置。

< 表現 4 >

表現 1 乃至 3 のうちのいずれかに記載の記憶制御装置であって、

各パラメータ領域は、I / O コマンドがライトコマンドである場合の転送パラメータの  
格納先とされるパラメータライト領域と、I / O コマンドがリードコマンドである場合の  
転送パラメータの格納先とされるパラメータリード領域とを有し、

30

前記データ転送回路は、前記パラメータライト領域よりも前記パラメータリード領域内  
の未完了の転送パラメータを優先的に取得する、

記憶制御装置。

< 表現 5 >

表現 1 乃至 4 のうちのいずれかに記載の記憶制御装置であって、

前記データ転送回路が、複数のパラメータ領域に蓄積されている未完了の転送パラメ  
ータの数をそれぞれ記憶する記憶領域である複数のインデックスと、前記複数のインデッ  
クスから一つのインデックスを選択するセレクトと、前記セレクトによって選択されたイン  
デックスに対応したパラメータ領域から転送パラメータを取得しその転送パラメータが有  
する転送元アドレス及び転送先アドレスを設定するパラメータ取得回路と、前記設定され  
た転送元アドレスが表す記憶領域内のデータを前記設定された転送先アドレスが表す記憶  
領域に転送する転送制御回路とを有する、

40

記憶制御装置。

【符号の説明】

【 0 1 7 7 】

2 0 : ストレージシステム



【 図 1 】

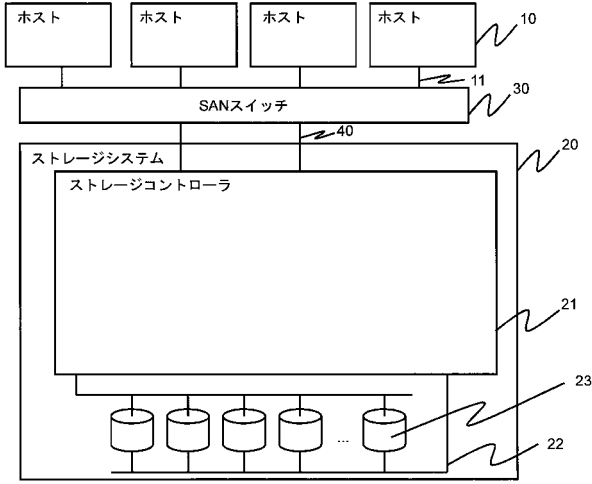


Fig.1

【 図 2 】

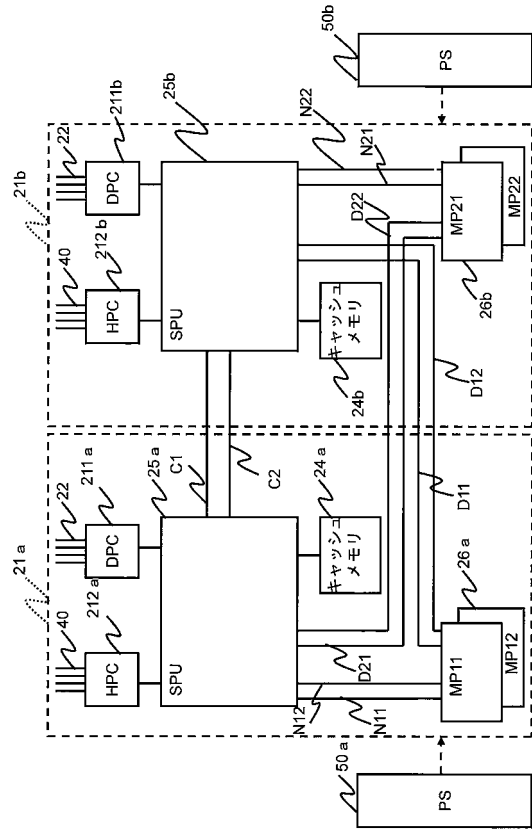


Fig.2

【 図 3 】

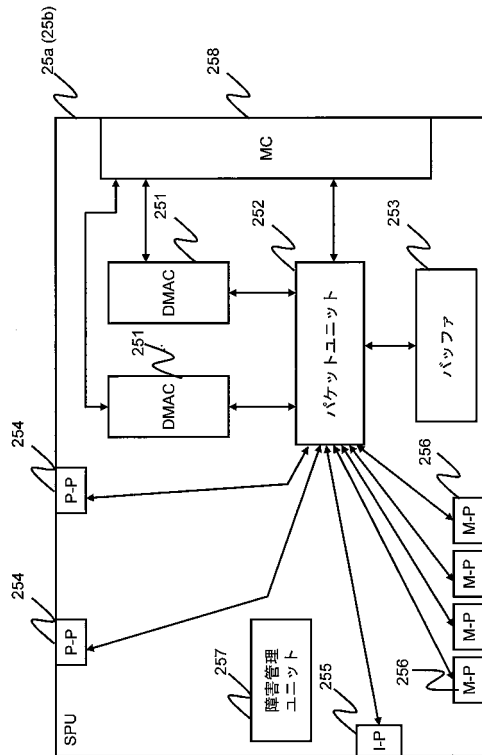


Fig.3

【 図 4 】

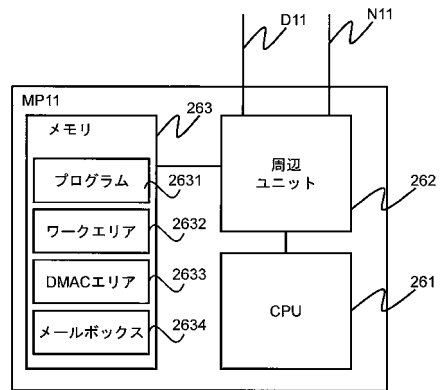


Fig.4

【 図 5 】

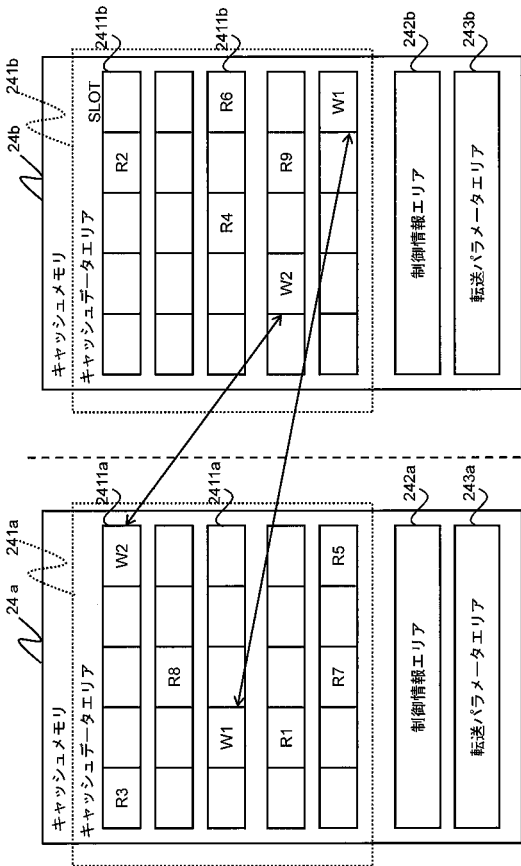


Fig.5

【 図 7 】

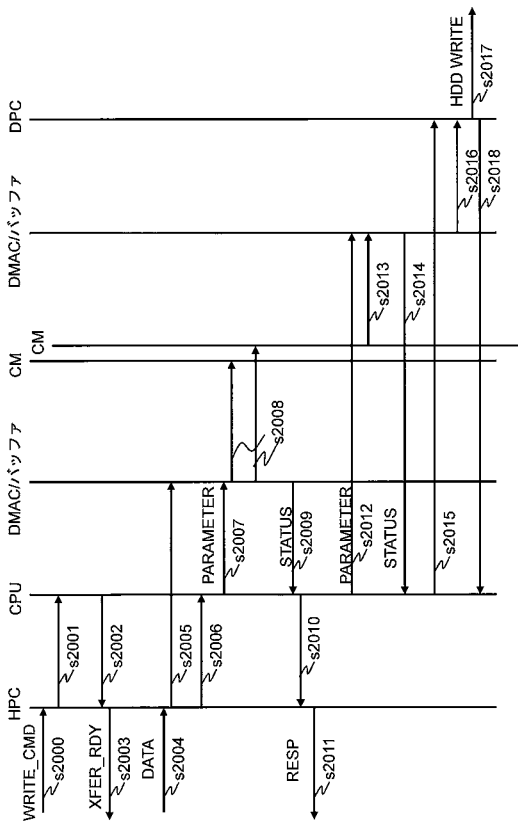


Fig.7

【 図 6 】

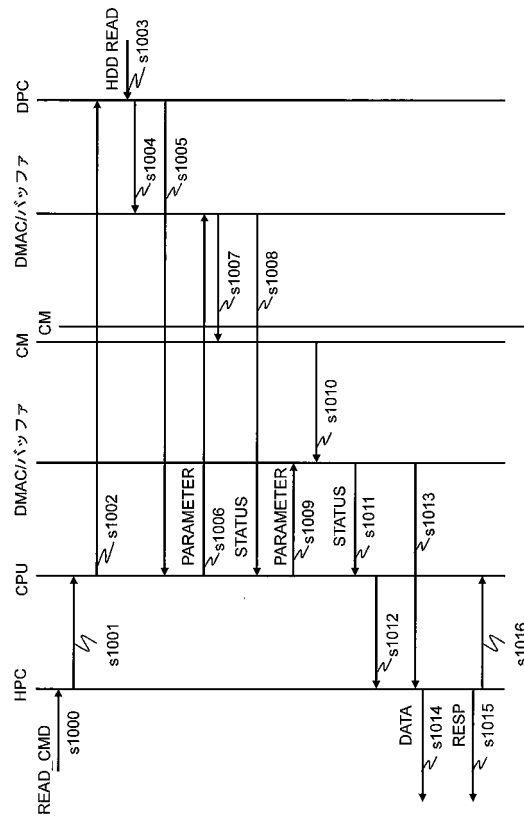


Fig.6

【 図 8 】

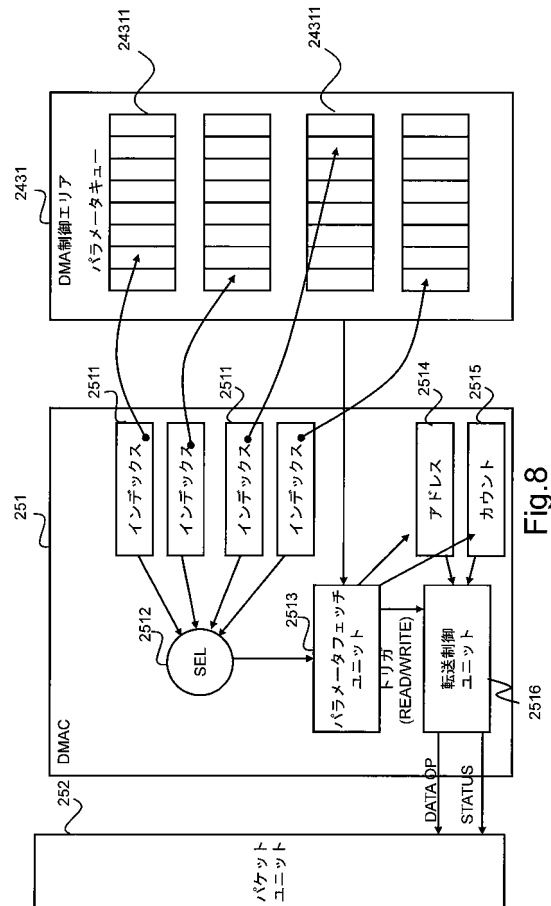
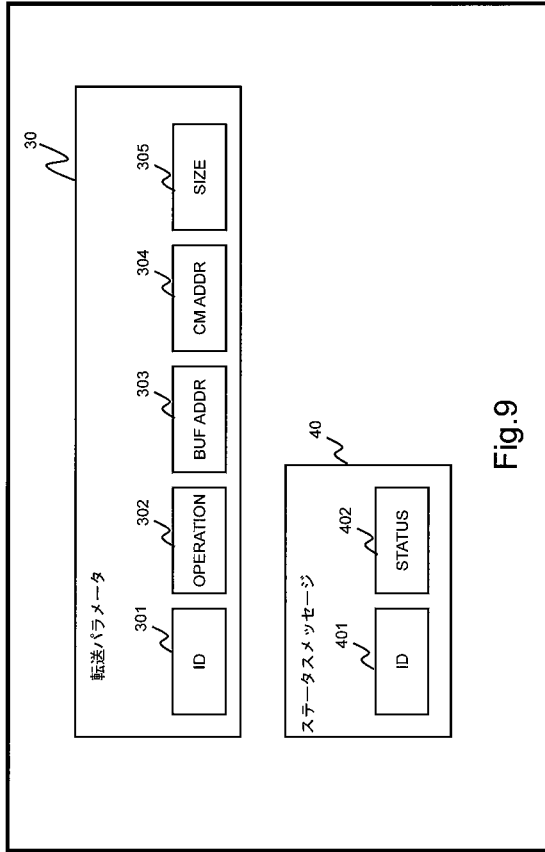
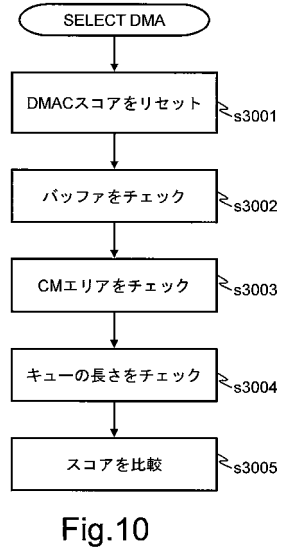


Fig.8

【図9】



【図10】



【図11A】

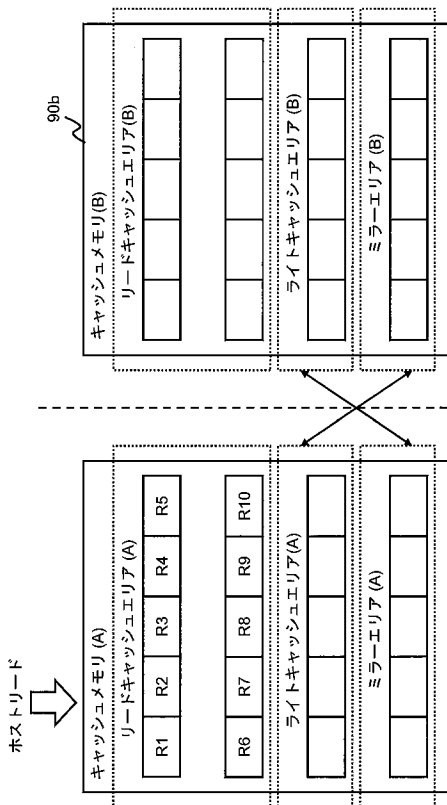


Fig.11A

【図11B】

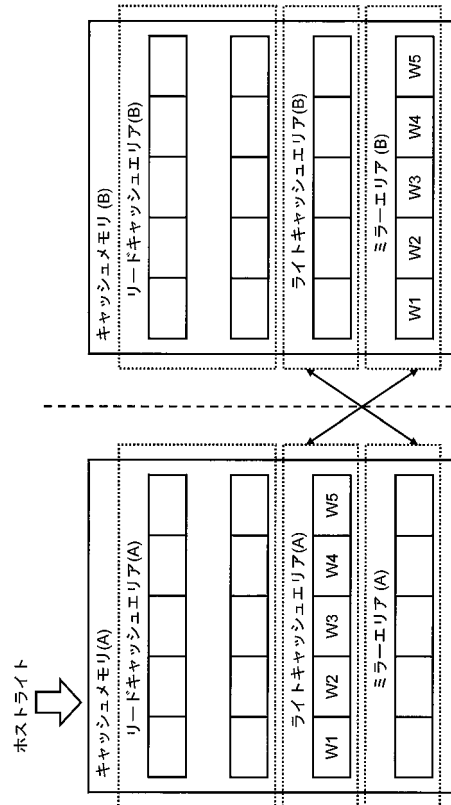


Fig.11B

【図12A】

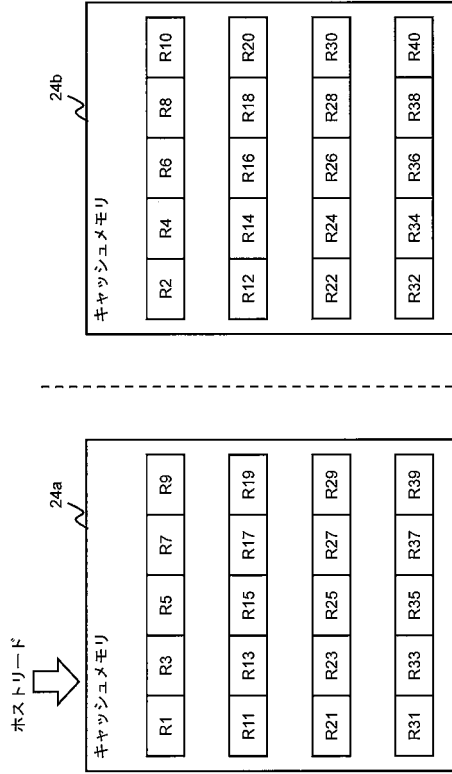


Fig.12A

【図12B】

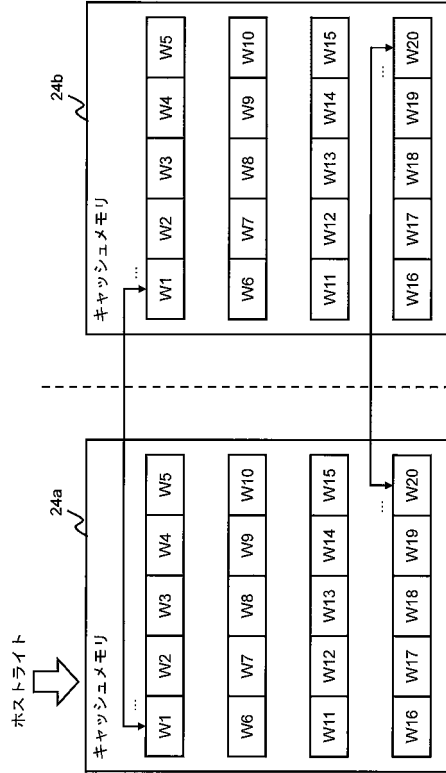


Fig.12B

【図13】

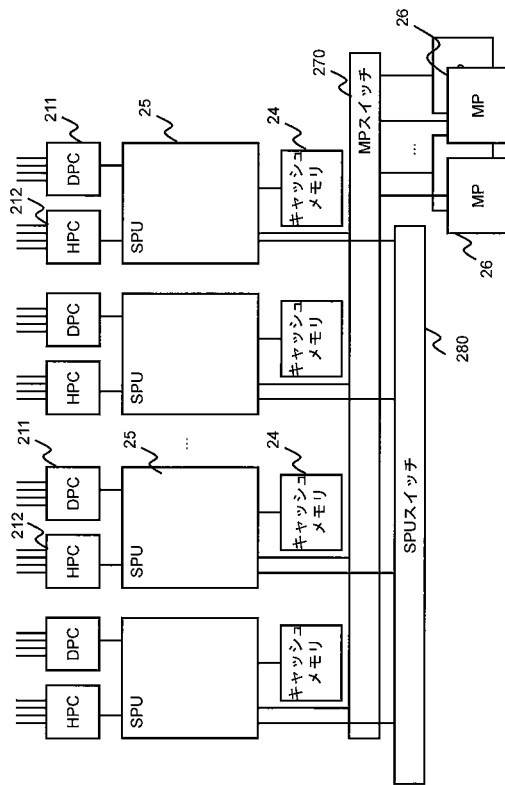


Fig.13

【図14】

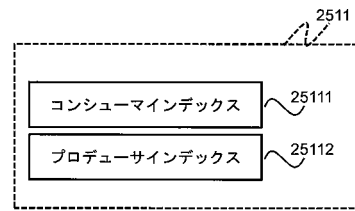


Fig.14

【 図 15 】

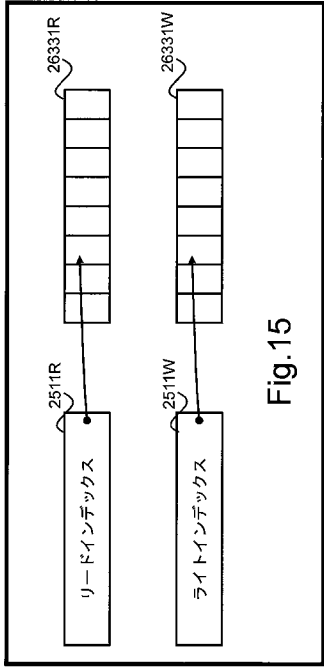


Fig.15

---

フロントページの続き

- (72)発明者 福田 秀明  
神奈川県小田原市中里322番2号 株式会社日立製作所 RAIDシステム事業部内
- (72)発明者 箕輪 信幸  
神奈川県小田原市中里322番2号 株式会社日立製作所 RAIDシステム事業部内

審査官 木村 貴俊

- (56)参考文献 特開平09-128305(JP,A)  
特開平08-335144(JP,A)  
特開2008-134776(JP,A)  
特開2001-318904(JP,A)

- (58)調査した分野(Int.Cl., DB名)
- |      |               |
|------|---------------|
| G06F | 3/06 - 3/08   |
| G06F | 12/00 - 12/16 |
| G06F | 13/00         |