

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5008845号  
(P5008845)

(45) 発行日 平成24年8月22日(2012.8.22)

(24) 登録日 平成24年6月8日(2012.6.8)

(51) Int.Cl. F 1  
**G 0 6 F 3/06 (2006.01)**  
 G 0 6 F 3/06 3 0 2 A  
 G 0 6 F 3/06 5 4 0

請求項の数 3 (全 22 頁)

(21) 出願番号	特願2005-252989 (P2005-252989)	(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	平成17年9月1日(2005.9.1)	(74) 代理人	110000062 特許業務法人第一国際特許事務所
(65) 公開番号	特開2007-66129 (P2007-66129A)	(72) 発明者	中村 崇仁 神奈川県川崎市麻生区王禅寺1099番地 株式会社 日立製作所 システム開発研 究所内
(43) 公開日	平成19年3月15日(2007.3.15)	(72) 発明者	藤本 和久 神奈川県川崎市麻生区王禅寺1099番地 株式会社 日立製作所 システム開発研 究所内
審査請求日	平成20年1月9日(2008.1.9)		

最終頁に続く

(54) 【発明の名称】 ストレージシステムとストレージ装置及びその制御方法

(57) 【特許請求の範囲】

【請求項1】

ホストコンピュータからのデータを保存する1以上のメディアと、前記メディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備えるストレージ装置を具備するストレージシステムにおいて、

前記メディアは、HDDと、フラッシュメモリとからなり、

前記ホストコンピュータからのリードの要求を受けた際に、格納先が前記HDDか前記フラッシュメモリかを判別し、格納先が前記HDDの場合には前記HDDから読み出したデータを前記キャッシュメモリ部に格納した後に前記ホストコンピュータに应答し、格納先が前記フラッシュメモリの場合には前記フラッシュメモリから読み出したデータを前記キャッシュメモリ部には格納せずに前記ホストコンピュータに应答し、

前記ストレージ装置は、前記フラッシュメモリに対する書き込み回数を平均化するように、デステージするデータを選択し、各々のフラッシュメモリに対するデステージした回数を記録し、デステージするデータを選択する際に、各々の前記フラッシュメモリに対するデステージした回数を比較し、デステージした回数が少ないものから優先的にデステージするデータを決定することを特徴とするストレージシステム。

【請求項2】

ホストコンピュータからのデータを保存する1以上のメディアと、前記メディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御

10

20

部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備えるストレージ装置を具備するストレージシステムにおいて、

前記メディアは、HDDと、フラッシュメモリとからなり、

前記ホストコンピュータからのリードの要求を受けた際に、格納先が前記HDDか前記フラッシュメモリかを判別し、格納先が前記HDDの場合には前記HDDから読み出したデータを前記キャッシュメモリ部に格納した後に前記ホストコンピュータに応答し、格納先が前記フラッシュメモリの場合には前記フラッシュメモリから読み出したデータを前記キャッシュメモリ部には格納せずに前記ホストコンピュータに応答し、

前記ストレージ装置は、前記フラッシュメモリに対する書き込み回数を平均化するように、デステージするデータを選択し、デステージするデータを選択する際に、各々の前記フラッシュメモリに対するデステージした回数およびアクセス回数を基に評価関数を算出し、評価関数が少ないものから優先的にデステージするデータを決定することを特徴とするストレージシステム。

【請求項3】

ホストコンピュータからのデータを保存するメディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備えるストレージ装置を制御する方法において、

前記メディアは、HDDと、フラッシュメモリとからなり、

前記ホストコンピュータからのリードの要求を受けた際に、格納先が前記HDDか前記フラッシュメモリかを判別し、格納先が前記HDDの場合には前記HDDから読み出したデータを前記キャッシュメモリ部に格納した後に前記ホストコンピュータに応答し、格納先が前記フラッシュメモリの場合には前記フラッシュメモリから読み出したデータを前記キャッシュメモリ部には格納せずに前記ホストコンピュータに応答し、

前記フラッシュメモリに対する書き込み回数を平均化するように、デステージするデータを選択することを特徴とするストレージ装置の制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージ装置に関し、また、その制御の方法に関するものである。

【背景技術】

【0002】

近年、データセンタなどの情報ビジネスの現場において、ストレージシステムの総所有コスト(TCO)の削減がますます重要になってきている。その一方で、データを長期的かつ確実に記録する要求も高まってきている。このことに関する例としては、金融機関および医療機関などの文書データは、消去せずに蓄積されることが法律によって義務付けられていることなどが挙げられる。こうした背景において、高信頼で大容量のストレージシステムが必要とされている。しかしながら、ハードディスクドライブ(以降では「HDD」という)を用いた大規模なストレージシステムでは、一般に、ストレージ容量に比例して電力消費量が増大する。つまり、大容量のストレージシステムを所有することは、電気料金が含まれる総所有コストが上昇してしまうことを示している。こうした状況を鑑み、キャッシュ管理アルゴリズムにより、HDDの消費する電力を低減させる技術が提案されている(非特許文献1)。また、問題は、電気料金に限られたことではない。一般に、ストレージシステムの容量が大きくなるほど、設置するための床面積が増大する。このことも総所有コストの上昇につながっている。

【0003】

ところで、近年、不揮発メディアとして、フラッシュメモリが注目されている。フラッシュメモリは、一般にHDDと比較し、数十倍以上低消費電力であり、高速に読み出しが可能である。また、HDDのように機械的な駆動部分が不要なため小型である。

【0004】

10

20

30

40

50

しかしながら、フラッシュメモリには、情報を保持するセルの物理的な理由により、書き込み回数の制限がある。こうした制限に対し、上位装置に示すアドレスとセル位置との対応をもち、各セルに書き込まれる回数を均等化するように制御を行う、ウェアレベリングと呼ばれる技術などにより、フラッシュメモリの書き込み可能回数の向上が図られている。なお、以降では、情報を保持するための素子を単に「フラッシュメモリ」といい、上記のウェアレベリングや、上位装置に対するプロトコル処理などを行う機構を含めたものを「フラッシュメモリデバイス」という。こういった技術により、フラッシュメモリデバイスとしての書き込み回数制限に対する効率化が図られているものの、フラッシュメモリデバイスの書き込み回数制限は依然存在している。また、それに加えて、フラッシュメモリに書き込む際に、消去と呼ばれる操作が必要になる場合はHDD同等程度の速度となっ

10

【0005】

このようなフラッシュメモリを用いたストレージシステムを構成する技術としては、例えば特許文献1がある。特許文献1は、RAID構成などにおいて頻繁にアクセスされるパリティデータをフラッシュメモリなどの半導体ディスクにおき、ストレージシステムの性能を向上させる技術を示している。しかし、ストレージシステムとしての書き込み回数制限を回避する手段は示されていない。また、HDDとフラッシュメモリを混在して1つのRAIDグループ、すなわち1つの仮想デバイスを構成するので、それぞれのメディアの特徴を考慮した仮想デバイスの制御はできていない。

【特許文献1】特開平6-324815号公報

20

【非特許文献1】ZHU, Q., DAVID, F., ZHOU, Y., DEVARAJ, C., AND CAO, P., "Reducing Energy Consumption of Disk Storage Using Power-Aware Cache Management". In Proc. of the 10th Intl. Symp. on High Performance Computer Architecture (HPCA-10) (Feb. 2004).

【発明の開示】

【発明が解決しようとする課題】

【0006】

以上のような背景から、低消費電力であり、設置面積が小さく、かつ、大容量を得ることのできる大規模まで構成可能なストレージシステムを提供することが課題となっている。

30

【0007】

また、データを格納するメディアに応じた、高いシステム性能を提供することも課題となっている。

【0008】

また、ストレージシステムとしての信頼性、可用性の向上も課題である。書き込み回数制限を有するメディアに対してはストレージシステムとして緩和を行わなくてはならない。

【課題を解決するための手段】

【0009】

本発明は、ホストコンピュータからのデータを保存するメディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備え、前記メディアは、HDDと、フラッシュメモリとからなり、前記ホストコンピュータからのリードの要求を受けた際に、格納先が前記HDDか前記フラッシュメモリかを判別し、格納先が前記HDDの場合には前記HDDから読み出したデータを前記キャッシュメモリ部に格納した後に前記ホストコンピュータに回答し、格納先が前記フラッシュメモリの場合には前記フラッシュメモリから読み出したデータを前記キャッシュメモリ部には格納せずに前記ホストコンピュータに回答し、前記フラッシュメモリに対する書き込み回数を平均化するように、デステージするデータを選択することを特徴とする。

40

【0010】

50

すなわち、本発明は、ホストコンピュータからのデータを保存する1以上のメディアと、前記メディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備えるストレージ装置を具備するストレージシステムにおいて、前記メディアは、HDDと、フラッシュメモリとからなり、前記ホストコンピュータからのリードの要求を受けた際に、格納先が前記HDDか前記フラッシュメモリかを判別し、格納先が前記HDDの場合には前記HDDから読み出したデータを前記キャッシュメモリ部に格納した後に前記ホストコンピュータに应答し、格納先が前記フラッシュメモリの場合には前記フラッシュメモリから読み出したデータを前記キャッシュメモリ部には格納せずに前記ホストコンピュータに应答し、前記ストレージ装置は、前記フラッシュメモリに対する書き込み回数を平均化するように、デステージするデータを選択し、各々のフラッシュメモリに対するデステージした回数を記録し、デステージするデータを選択する際に、各々の前記フラッシュメモリに対するデステージした回数を比較し、デステージした回数が少ないものから優先的にデステージするデータを決定することを特徴とするストレージシステムである。

10

#### 【発明の効果】

##### 【0011】

本発明のストレージシステムでは、低消費電力であり、設置面積が小さく、かつ、大規模まで構成可能であり、また、データを格納するメディアに応じた、高いシステム性能を提供することもできる。また、各メディアに対する書き込み回数の低減を図るので、書き込み回数制限を有するメディアに対してもストレージシステムとして、信頼性、可用性の向上可能になるといった利点がある。

20

#### 【発明を実施するための最良の形態】

##### 【0012】

本発明を実施するための最良の形態を説明する。

以下に本発明のストレージシステムとストレージ装置及びその制御方法の実施例について、図面に基づいて説明する。

##### 【実施例1】

##### 【0013】

実施例1を説明する。図1は、本発明の第1の実施の形態のストレージシステムの構成のブロック図である。ストレージシステムは、ストレージ制御装置1およびハードディスクドライブ(HDD)50、フラッシュメモリデバイス(本図ではFM制御部16の内部にフラッシュメモリデバイスを備える例を示している)から構成されている。ストレージ制御装置1は、チャンネル4を通じて、SANスイッチ3などで構成されるSAN(Storage Area Network)を経て、ホストコンピュータ2に接続される。また、ディスク側チャンネル60を通じてデータを格納するメディアである複数のHDD50と接続されている。ストレージ制御装置1は、複数のチャンネル制御部11と複数のキャッシュメモリ13と制御メモリ部17と複数のディスク制御部14および複数のFM制御部16と、これらを、内部バス15を介して接続する内部スイッチ12から成る。チャンネル制御部11は、チャンネル4を通じてホストコンピュータ2からの入出力要求を受け取り、この入出力要求の種類(例えば、リード要求、ライト要求)や対象アドレスなどを解釈し、図7以降で述べるような処理を行う。キャッシュメモリ部はHDDやフラッシュメモリに格納されるべきデータやホストコンピュータ2に返すべきデータを一時的に格納する。制御メモリ部17は、キャッシュメモリ部13のディレクトリ情報や、ストレージシステムの構成情報を格納している。ディスク制御部14は、チャンネル制御部11などの要求に基づき、ディスク側チャンネル60を通じてHDD50の制御を行い、ホストコンピュータ2から要求されたデータを取り出しや格納を行う。この際、ディスク制御部14がHDD50に対してRAID制御を行いストレージシステムの信頼性、可用性および性能を向上させてもよい。FM制御部16は、フラッシュメモリまたはフラッシュメモリデバイスの制御を行う。チャンネル制御部11などの要求に基づき、フラッシュメモリまたはフラッシュメモリデバイス

30

40

50

にホストコンピュータ2から要求されたデータを取り出しや格納を行う。この際、FM制御部16がフラッシュメモリデバイスに対してRAID制御を行いストレージシステムの信頼性、可用性および性能を向上させてもよい。なお、本実施例では、ストレージシステムはHDD50と接続されているが、HDD50ならびにディスク制御部14を持たない構成でもよい。また、制御メモリ部17に格納される情報は、物理的にキャッシュメモリ部13と同一のメモリ上に配置してもよい。

#### 【0014】

図2はチャンネル制御部11の詳細の構成のブロック図である。チャンネル制御部11は、複数のプロセッサ111、メモリモジュール112、周辺処理部113、及び、複数のチャンネルプロトコル処理部114と内部ネットワークインタフェース部117から成る。プロセッサ111はバス等で周辺処理部113に接続される。周辺処理部113は、メモリモジュール112に接続され、メモリモジュールの制御を行う。また、制御系バス115を介してチャンネルプロトコル処理部114および内部ネットワークインタフェース部117にも接続される。周辺処理部113は、接続されるプロセッサ111及びチャンネルプロトコル処理部114及び内部ネットワークインタフェース部117からのパケットを受け、パケットの示す転送先アドレスがメモリモジュール112上ならばその処理を行い、必要ならばデータを返す。また転送先アドレスがそれ以外ならば、適切なフォワーディングを行う。また、周辺処理部113は、その他のプロセッサ111がこの周辺処理部113に接続されるプロセッサ111と通信を行うためのメールボックス1131を持つ。プロセッサ111は、周辺処理部113を通してメモリモジュール112にアクセスし、メモリモジュール112に格納された制御プログラム1121に基づいて処理を行う。また、メモリモジュール112には、チャンネルプロトコル処理部114がDMA(Direct Memory Access)を行うための転送リスト1123も格納されている。チャンネルプロトコル処理部114は、チャンネル4上のプロトコル制御を行い、ストレージシステム1内部で処理ができるようなプロトコル方式に変換する。チャンネルプロトコル処理部114は、チャンネル4を通じてホストコンピュータ2からの入出力要求を受けると、その入出力要求のホストコンピュータ番号やLUN(Logical Unit Number)やアクセス先アドレスなどをプロセッサ111に通知する。プロセッサ111は、その通知に基づき、ディレクトリ情報1323にアクセスし、入出力要求のデータを格納すべきアドレスまたは入出力要求のデータが存在する場合はメモリモジュール112上に転送リスト1123を作成し、それに基づきチャンネルプロトコル処理部114に転送を行わせる。ホストコンピュータ2がリード要求したデータがキャッシュメモリ13上になく、HDD50に格納されているならば、ディスク制御部14に、HDD50に格納されている要求データをキャッシュメモリ13に格納する(この動作を「ステージング」という)ように指示を与えた後に転送リスト1123により転送させる。フラッシュメモリ上に格納されているならば、フラッシュメモリのアドレスを転送リストにセットする。転送リストは、キャッシュメモリ13もしくはフラッシュメモリ上のアドレスのリストになっており、入出力要求がライトならば、データ転送系バス115を介して接続された内部ネットワークインタフェース部117を通じ、ホストコンピュータからのデータをリストに記載されたアドレスに書き込んでいく。またリードならば、同様にリストに記載されたアドレスからデータを読み込み、それをホストコンピュータに返す。これらの動作の詳細に関しては図7以降を用いて説明する。内部ネットワークインタフェース部117は、チャンネル制御部11内と他のストレージシステム1内部と内部バス15を経て通信を行う際の、インタフェースとなる部位である。

#### 【0015】

なお、ディスク制御部14もチャンネル制御部11とほぼ同じ構造を持つ。ただし、制御プログラム1121の内容が異なり、また、チャンネルプロトコル処理部114はHDD50との通信を行う(チャンネル4とディスク側チャンネル60のプロトコルは異なってもよい)。しかし、ディスク側チャンネル60上のプロトコル処理を行い、ストレージシステム1内部で処理が出来るように変換する意味ではチャンネル制御部11のチャンネルプロトコル処理

10

20

30

40

50

部 1 1 4 と同様)。プロセッサ 1 1 1 はチャンネル制御部 1 1 からの要求や一定時間間隔で、キャッシュメモリ 1 3 上のデータをハードディスクドライブ 5 0 に書き込みを行う。また、キャッシュメモリ 1 3 上にホストコンピュータから要求されたデータが存在していない場合は、チャンネル制御部 1 1 からの指示を受け、HDD 5 0 からデータを読み込み、キャッシュメモリ 1 3 にそのデータを書き込む。これらの際には、プロセッサ 1 1 1 は、制御メモリ部 1 7 に格納されたディレクトリ情報にアクセスし、ホストコンピュータ 2 の要求するデータを読み出すべきまたは格納すべきキャッシュメモリのアドレスを調査する。そして、要求されたデータがキャッシュメモリ 1 3 上に存在しない場合でキャッシュメモリ 1 3 に空き領域がないときは、要求されたデータを格納するため空き領域を作るべく、すでにあるデータを HDD 5 0 に格納する（この動作をデステージと呼ぶ）。これら HDD 5 0 の操作においては、ディスク制御部 1 4 がディスク側チャンネル 6 0 を通じて HDD 5 0 を制御する。この際、HDD 5 0 全体としての可用性および性能を向上させるため、ディスク制御部 1 4 は HDD 5 0 群に対して RAID 制御を行う。

10

## 【 0 0 1 6 】

図 3 は、FM 制御部 1 6 の詳細の構成のブロック図であり、フラッシュメモリを一体化している。FM 制御部 1 6 は、内部ネットワークインタフェース部 1 6 1、DMA コントローラ 1 6 2、揮発メモリであるメモリモジュール 1 6 4 およびその制御を行うメモリコントローラ 1 6 3、ならびにフラッシュメモリ 1 6 6（図中では FM）とそれを制御する FM コントローラ 1 6 5 が備えられている。内部ネットワークインタフェース部 1 6 1 は、FM 制御部 1 6 内と他のストレージ制御装置 1 内部と内部バス 1 5 を経て通信を行う際の、インタフェースとなる部位である。FM 制御部 1 6 内の DMA コントローラ 1 6 2 は、ホストコンピュータ 2 からのライトの要求を処理する際キャッシュメモリに空き領域を作るなどの場合に、チャンネル制御部 1 1 のプロセッサ 1 1 1 がセットする転送リスト 1 6 4 1 により、キャッシュメモリ部 1 3 からフラッシュメモリ 1 6 6 へのデータの転送を行う。FM コントローラ 1 6 5 は、内部ネットワークを通してなされたチャンネル制御部 1 1 からの読み込み要求や、DMA コントローラ 1 6 2 の書き込み要求により、フラッシュメモリ 1 6 6 を制御しデータのやり取りを行う。図 3 において、フラッシュメモリ 1 6 6 の実装形態としては、基板に直接配置することができる。この際は図 4 のコネクタや FM プロトコル処理部、図 5 の FM 側チャンネルなどの部品が不要になるので、よりコンパクトにストレージシステムを実現することが可能になる。また、FM コントローラ 1 6 5 にてそれぞれのフラッシュメモリ 1 6 6 に対するウェアレベリングなどを行ってもよい。

20

30

## 【 0 0 1 7 】

図 4 は、FM 制御部 1 6 の別の詳細の構成のブロック図である。ここでは、記憶素子としてフラッシュメモリデバイス 1 6 9 を用いており、フラッシュメモリデバイスを別体としている。フラッシュメモリデバイス 1 6 9 は、コネクタ 1 6 8 を介して FM 制御部 1 6 と接続されているため、脱着が可能となっている。そのため、フラッシュメモリデバイス 1 6 9 が故障した際には、交換することも可能となる（これを行うためには、チャンネル制御部 1 1 のプロセッサ 1 1 1 が、予めフラッシュメモリデバイス 1 6 9 間で冗長構成をとるように、転送リスト 1 6 4 1 をセットしていればよい）。また、フラッシュメモリデバイス 1 6 9 自体をより容量の大きなものに交換することも可能となる。このフラッシュメモリデバイス 1 6 9 は、内部においてウェアレベリングなどの技術により信頼性、性能の向上が図られているものであり、外部とのやり取りは、専用プロトコルで行われる。そのため、FM プロトコル処理部 1 6 7 にて、ストレージ制御装置 1 内部で処理できる形式に変換を行っている。

40

## 【 0 0 1 8 】

図 5 は、FM 制御部 1 6 の別の詳細の構成のブロック図である。ここでは、フラッシュメモリデバイス 1 6 9 を FM 側チャンネル 1 6 1 0 により接続している。この構成をとることにより、図 4 の特徴に加え、より多数のフラッシュメモリデバイス 1 6 9 を接続でき、大容量のストレージシステムを実現することができる。また、あるフラッシュメモリデバイス 1 6 9 の一部の領域は、後に説明する緊急デステージ領域 1 6 9 0 としてもよい。

50

## 【 0 0 1 9 】

図6は、内部スイッチ部12の詳細の構成のブロック図である。内部スイッチ部12は、内部ネットワークインタフェース部121と複数のセレクタ122からなっている。セレクタ122は、ストレージ制御装置1の内部のチャンネル制御部11など各部から送られてきた要求の要求先を解析し、その要求先に接続される内部バス15を制御する内部ネットワークインタフェース部121に要求を転送する。その際、各セレクタ間で要求転送先の内部ネットワークインタフェース部121競合を行う。この内部スイッチ部12により、チャンネル制御部11からは、キャッシュメモリ部13、制御メモリ部17およびFM制御部16と直接のやり取りを行うことができる。FM制御部16はチャンネル制御部11、キャッシュメモリ部13および制御メモリ部17とやり取りができる。また、ディスク制御部14からは、キャッシュメモリ部13と制御メモリ部17に直接のやり取りができる。FM制御部16とディスク制御部14の接続に関する相違は、この内部スイッチ部12はチャンネル制御部FM制御部間接続123を備えていることにより、FM制御部16はチャンネル制御部11と直接のやり取りが行える点である。

10

## 【 0 0 2 0 】

図7は、ホストコンピュータ2から、HDD50領域へのリードの要求が行われた場合の処理の流れを表す図である。ホストコンピュータ2からリード要求を、チャンネル4を通じてチャンネル制御部11が受信する(s701)。チャンネル制御部11のプロセッサ111は、受信した要求を解析し、LUNや対象論理ブロックアドレスを得る。ここでは、該当データがHDD50に格納される領域のものであることを知る(s702)。さらにチャンネル制御部11のプロセッサ111は、制御メモリ部17に格納された、ライトキャッシュ領域ならびにリードキャッシュ領域のディレクトリ情報にアクセスし、キャッシュメモリ部13に該当データの格納の有無を確認する(s703、s704。図中では1回のアクセスになっているが、実際には複数回のアクセスになる場合もある。以下も同様)。該当データがキャッシュメモリ部13に既に存在していればs715以降の処理によりホストコンピュータ2に応答するが、ここでは、該当データがキャッシュメモリ部13に存在しなかったとする。この場合、ディスク制御部に該当データをキャッシュメモリ部13に転送させる(ステージング)が、キャッシュメモリ部13に空き領域がない場合は、ステージングの前にデータを格納すべきキャッシュ領域を作らなければならない。どの領域に空き領域にすべきかを決定し、その操作をおこなうステップがs705、s706である。領域が確保された後で、ステージングの要求を、制御メモリ部17の通信領域173にメッセージを書き込むことを介してディスク制御部14に通知する(s707)。ディスク制御部14は定期的もしくは一連の処理の終了毎に、制御メモリ部17の通信領域173を読み込むことによりチャンネル制御部11からの要求があることを知る(s708、s709)。なお、こうした方式でチャンネル制御部11とディスク制御部14が連携を行うのは、HDD50からデータを得られる時間が不定で、その他の処理時間と比較して長いためであり、この方式によりバックグラウンドで他の要求の処理などを行うことができる。ステージングの要求を認知したディスク制御部14は、該当データが得られるようにディスク側チャンネル60を通じてHDD50を制御する(s710)。HDD50からのデータが得られると、s705、s706により確保した領域に、該当データを書き込む(s711。ステージング)。また、ステージングが終了したことを、通信領域173を用いてチャンネル制御部11に通知する(s712)。チャンネル制御部11は、ディスク制御部14同様に通信領域173を読み込んでおり、ステージング終了のメッセージの存在を知る(s713、s714)。その後、チャンネル制御部11のプロセッサ111は、転送リスト1123をセットし(あらかじめステージング要求直後にセットしておいてもよい)、転送の指示を行う。チャンネルプロトコル処理部114は、キャッシュメモリ部13から該当データを読み込み、ホストコンピュータ2に転送する(s715、s716、s717)。以上が、HDD50に格納されたデータのリードの処理の流れである。HDD50の速度が、不定であり、遅いため、キャッシュメモリ部13を介したデータのやり取りが必要となる。

20

30

40

50

## 【 0 0 2 1 】

図 8 は、ホストコンピュータ 2 から、フラッシュメモリ領域へのリードの要求が行われた場合の処理の流れを表す図である。ホストコンピュータ 2 からリード要求を、チャンネル 4 を通じてチャンネル制御部 1 1 が受信する ( s 8 0 1 )。チャンネル制御部 1 1 のプロセッサ 1 1 1 は、受信した要求を解析し、LUN や対象論理ブロックアドレスを得る。ここでは、該当データがフラッシュメモリに格納される領域のものであることを知る ( s 8 0 2 )。さらにチャンネル制御部 1 1 のプロセッサ 1 1 1 は、制御メモリ部 1 7 に格納された、ライトキャッシュ領域のディレクトリ情報にアクセスし、キャッシュメモリ部 1 3 のライトキャッシュ領域 1 3 2 に該当データの格納の有無を確認する ( s 8 0 3、s 8 0 4 )。図 7 の場合と違い、ライトキャッシュ領域 1 3 2 に対する調査のみでよい。ここでは、該当データがキャッシュメモリ部 1 3 に存在しなかったとする。チャンネル制御部 1 1 のプロセッサ 1 1 1 は、転送リスト 1 1 2 3 をセットし、チャンネルプロトコル処理部 1 1 4 に転送の指示を行う。チャンネルプロトコル処理部 1 1 4 は、内部スイッチ部 1 2 のチャンネル制御部 F M 制御部間接続 1 2 3 を介して、F M 制御部 1 6 にデータの要求を行う ( s 8 0 5 )。F M 制御部 1 6 は、フラッシュメモリ 1 6 6 もしくはフラッシュメモリデバイス 1 6 9 からデータを取り出し ( s 8 0 6 )、チャンネル制御部 1 1 にデータを返す ( s 8 0 7 )。こうして得たデータをチャンネルプロトコル処理部はホストコンピュータ 2 に転送する ( s 8 0 8 )。以上が、フラッシュメモリ 1 6 6 もしくはフラッシュメモリデバイス 1 6 9 に格納されたデータのリード手順であった。フラッシュメモリの読み出し速度が一定で、高速のため ( s 8 0 6 )、キャッシュメモリ部 1 3 を介さずにデータのやり取りが可能になる。また、直接転送を行うため、リードに関してはキャッシュ処理を行わず、リードキャッシュ領域 1 3 1 の調査 ( s 7 0 3、s 7 0 4 ) やキャッシュ領域の確保 ( s 7 0 5、s 7 0 6 ) の処理が不要になり、メディアからの読み出し高速化のみならず、付随する処理の高速化も可能になる。また、フラッシュメモリに対しては、リードキャッシュ領域 1 3 1 が不要なため、キャッシュメモリ部 1 3 やそのリードキャッシュディレクトリ情報 1 7 1 1 を格納する制御メモリ部 1 7 の容量を小さくすることが可能になる効果もある。

## 【 0 0 2 2 】

続いて、フラッシュメモリ 1 6 6 またはフラッシュメモリデバイス 1 6 9 に好適なキャッシュメモリ部 1 3 の制御方法を説明するが、まずその前に、ホストコンピュータ 2 からライト要求が行われた場合の処理の流れと、キャッシュメモリ部 1 3 に関わる情報の説明を行う。

## 【 0 0 2 3 】

図 1 1 A は、ホストコンピュータ 2 からライトの要求が行われた場合で、すでにキャッシュメモリ部 1 3 上に該当アドレスのデータが存在する場合の処理の流れを表す図である。ライト要求については、対象が HDD 5 0 領域、フラッシュメモリ領域いずれの場合もデータをキャッシュメモリ部 1 3 に一時的に格納する。まず、ホストコンピュータ 2 からライト要求を、チャンネル 4 を通じてチャンネル制御部 1 1 が受信する ( s 1 1 0 1 )。チャンネル制御部 1 1 のプロセッサ 1 1 1 は、受信した要求を解析し、LUN や対象論理ブロックアドレスを得る ( s 1 1 0 2 )。さらにチャンネル制御部 1 1 のプロセッサ 1 1 1 は、制御メモリ部 1 7 に格納された、ライトキャッシュ領域ならびにリードキャッシュ領域のディレクトリ情報にアクセスし、キャッシュメモリ部 1 3 に該当データの格納の有無を確認する ( s 1 1 0 3、s 1 1 0 4 )。なお、リードキャッシュ領域 1 3 1 に存在した場合は、後に述べるような手順で該当スロットを無効化し、新たにライトキャッシュディレクトリ情報 1 7 1 2 に加える)。ここでは、該当アドレスのデータがライトキャッシュ領域 1 3 2 に既に存在していたとする。この場合、同じキャッシュメモリのスロット ( キャッシュの制御単位。後で詳述 ) にデータを格納する。チャンネル制御部 1 1 は、ホストコンピュータ 2 にデータの送信要求を行い ( s 1 1 0 5 ) つつ、データを該当するスロットへ格納するように転送リスト 1 1 2 3 をセットする。チャンネル制御部 1 1 のチャンネルプロトコル処理部 1 1 4 は、応答したホストコンピュータ 2 からのデータを受信し ( s 1 1 0 6 )、それを転送リスト 1 1 2 3 に基づきキャッシュメモリ部 1 3 に格納する ( s 1 1 0 7 )。



このような場合、フラッシュメモリに対するライト要求であったとしても、フラッシュメモリに書き込みをしなくともよく、書き込み回数の低減が図れる。このため、ライト要求の場合は、リード要求の場合と異なり、データの格納メディアがフラッシュメモリ 166 もしくはフラッシュメモリデバイス 169 の場合でも、キャッシュメモリ部 13 を用いた処理を行う。

#### 【0024】

次に図 11B を用いて、ホストコンピュータ 2 から、ライトの要求が行われた場合で、キャッシュメモリ部 13 上に該当アドレスのデータが存在しないが既に空いているスロットがない場合の処理の流れを説明する。s 1124 までは、図 11A の s 1104 と同様であるが、この図においては、該当アドレスのデータがライトキャッシュ領域 132 に存在しておらず、空きのスロットはない場合とする。そこで、まず、既にキャッシュメモリ部 13 上に存在するデータを、格納されるべき HDD 50、フラッシュメモリに格納し（デステージング）、今回のホストコンピュータ 2 からのデータを格納するための領域を作る。まず、LRU (Least Recently Used) アルゴリズムなどで、デステージするデータを決定する (s 1125)。デステージするデータが決まったら、チャンネル制御部 11 のプロセッサ 111 は、FM 制御部 16 に、キャッシュメモリ部 13 上のアドレスと格納先のフラッシュメモリ 166 もしくはフラッシュメモリデバイス 169 のアドレスの対応が記された転送リスト 1641 をセットする (s 1127)。なお、対象が HDD 50 の場合は、図 7 同様に通信領域 173 を用いてディスク制御部 14 にメッセージを送信することにより、デステージの要求が行われる。続いて、転送リスト 1641 に基づき、FM 制御部 16 の DMA コントローラ 162 はキャッシュメモリ部 13 から該当データを読み込み (s 1127、s 1128)、フラッシュメモリ 166 もしくはフラッシュメモリデバイス 169 に書き込む (s 1129)。一連の処理が終了すると DMA コントローラ 162 はデステージが完了した旨をチャンネル制御部 11 のプロセッサ 111 に通知する (s 1130)。プロセッサ 111 は、デステージしたスロットに今回のホストコンピュータ 2 からのデータを格納するために、ディレクトリ情報を更新し (s 1131、s 1132)、ホストコンピュータ 2 にデータの送信要求を行い (s 1133) つつ、データを該当するスロットへ格納するように転送リスト 1123 をセットする。チャンネル制御部 11 のチャンネルプロトコル処理部 114 は、応答したホストコンピュータ 2 からのデータを受信し (s 1134)、それを転送リスト 1123 に基づきキャッシュメモリ部 13 に格納する (s 1135)。

#### 【0025】

次に、これまで述べてきたキャッシュメモリ部 13 に関する情報を説明する。図 9 は、キャッシュメモリ部 13 ならびに制御メモリ部 17 に格納されるデータの詳細を示したブロック図である。キャッシュメモリ部 13 は、リードキャッシュ領域 131 とライトキャッシュ領域 132 を持つ。リードキャッシュ領域 131 は、ホストコンピュータ 2 からリード要求されたデータを一時的に格納する。これにより、格納されているデータに対して再びホストコンピュータ 2 からリード要求を受け付けた場合に、HDD 50 から再び読み込むことなくリードキャッシュ領域 131 のデータを応答することによりストレージシステムとしての処理の高速化を図る。本実施例では、HDD 50 の読み出しまでの時間が不定であり、低速であるためリードキャッシュ領域 131 を設けているが、フラッシュメモリに格納されているデータに対してはリードキャッシュ領域 131 を用いない。ライトキャッシュ領域 132 は、ホストコンピュータ 2 からライト要求されたデータを一時的に格納し、格納されているデータと同じアドレスに対して再びライト要求を受け付けた場合には、ライトキャッシュ領域 132 上の該当データを上書きする。HDD 50 領域に対しては、高速なキャッシュメモリ部 13 に一時的に書き込んでおくことで性能の向上も期待できる。また、キャッシュメモリ部 13 上で上書きさせるため、フラッシュメモリに書き込む回数を低減することができ、ストレージシステムの信頼性、可用性を向上させることができる。なお、これらの管理している単位をここではスロットという。

#### 【0026】

次に制御メモリ部17について説明する。制御メモリ部17は、ディレクトリ情報171、構成情報172ならびに通信領域173を格納している。構成情報172は、ストレージシステムに関する情報である。例えば、チャンネル制御部11がいくつ存在していて、それぞれ接続されているチャンネル4にどういったLU(Logical Unit)を提供し、そのLUは、どのHDD50もしくはフラッシュメモリデバイス166を用いてデータを格納するか、またそれはどういった形で仮想的に提供するか、といった情報である。通信領域173は、チャンネル制御部11とディスク制御部14とが互いにメッセージを書き込みまたは読み込みを行い協調動作するための領域である。また、ディレクトリ情報171は、キャッシュメモリ部13にどのようなデータが格納されているかを示す情報であり、どのLUNや論理ブロックアドレス(LBA、Logical Block Address)などで示されるホストコンピュータ2からのアドレスのデータがキャッシュメモリ部13上のどこに格納されているかを示すリードキャッシュディレクトリ情報1711、ライトキャッシュディレクトリ情報1712、および、新しいデータをキャッシュメモリ部13に格納するために、どの元からあるデータを無効化するかを決定するための基準となるリードおよびライトのアクセス順序リスト1713、1714を持つ。

#### 【0027】

なお、キャッシュメモリ部13、制御メモリ部17は、高速で書き込み回数制限のない揮発性のメモリで構成される。しかし、停電などに備え、一定時間は内容を保持できるバッテリーバックアップの機能を備えている。しかし、そのバッテリーの保つことの出来る時間を越えると考えられる更に大規模な停電の際には、図11Bのs1125からs1129の処理を用いて、全ての一時的に格納されたキャッシュメモリ部13上のデータをメディアに書き込み、バックアップバッテリーが切れた際にもデータを失わないようにする。もしくは、s1125からs1129を用いた処理はプロセッサ111の計算量が必要となってしまうのでこの方法を用いず、メディア上に予め用意された緊急デスチージ領域に、キャッシュメモリ部13および制御メモリ部17のイメージをそのまま書き込む。これは、転送リスト1641をキャッシュメモリ部13および制御メモリ部17の全域を獲得するようにセットすることによって実行しても、DMAコントローラ162に、キャッシュメモリ部13および制御メモリ部17のイメージを取得する機能を設けることによって実行してもよい。特にフラッシュメモリは、低消費電力であるので、バックアップバッテリーの限られた時間と電力で実行せねばならないこの処理のメディアとして適している。つまり、この緊急デスチージ領域をフラッシュメモリ166もしくはフラッシュメモリデバイス169上に設けることにより大きな効果が得られる。例えば、図5ではあるフラッシュメモリデバイス169の一部の領域に緊急デスチージ領域1690を割り当てている。

#### 【0028】

図10Aは、リードキャッシュディレクトリ情報1711、ライトキャッシュディレクトリ情報1712の詳細を示したブロック図である。ディレクトリ情報は、LUN欄1001、LBA欄1002、最終的にそのデータが格納されるメディアを示すメディア欄1003および、それに続く、ホストコンピュータ2から指定されるアドレス(LBA)とそのデータが格納されているキャッシュメモリ部13上のアドレスの対応がリスト形式で格納されているホスト-キャッシュアドレス対応リスト1004から成る。キャッシュメモリ部13上のホストコンピュータ2の要求したデータがあるかどうかを調査する場合(図7のs703、s704や図8のs803、s804の操作に相当)、まず受領した要求のLUNと一致するものをLUN欄1001から見つける。次に対応するポインタが指し示すテーブルのLBA欄1002と、受領した要求のLBAの一致するものを調査する。LBA欄1002はLBAの範囲が示されており、一致するLBA範囲のデータに関する情報が、その欄に対応するポインタが指す先の、ホスト-キャッシュアドレス対応リスト1004に格納されている。これには、(範囲ではなく)LBAの値と、そのデータが格納されているキャッシュのアドレス(スロット番号)が記されている。このホスト-キャッシュアドレス対応リスト1004の中に対応するLBAがなければキャッシュメモリ部13には該当データが格納されていないことになる。

10

20

30

40

50

## 【 0 0 2 9 】

図 1 0 B は、従来のアクセス順序リスト 1 7 1 3、1 7 1 4 の詳細を示したブロック図である。これは、どのスロットのデータが何番目に最近アクセスされたのかを示す情報である。先頭に格納されているものが最も最近にアクセスされたスロット、最後尾に格納されているものが最も昔にアクセスされたスロットである。リードキャッシュ領域 1 3 1、ライトキャッシュ領域 1 3 2 のそれぞれに対応する情報がリスト構造で格納されている。例えば、ホストコンピュータ 2 からキャッシュメモリ部 1 3 に格納されていないデータに対してのアクセス要求を受領した場合など、キャッシュメモリ部 1 3 に空きスロットを作る場合、一般的な L R U アルゴリズムを用いるとすると、キャッシュに新たな空きスロットを作る場合、最も昔にアクセスされたもの、すなわち、リストの最後尾に示されているスロットを無効化する。その後、新しいデータを該当スロットに格納し、該当スロットは最も最近にアクセスされたことになるので、アクセス順序リスト 1 7 1 3、1 7 1 4 の先頭に該当スロットをおく。このリスト操作の手順としては、まず、最後尾にあった該当スロットの情報を delete 操作により外し、先頭に該当スロットの情報を insert 操作により追加する。なお、この際、キャッシュディレクトリ情報 1 7 1 1 または 1 7 1 2 も更新する。該当スロットに元に格納されていたものを示す情報をホスト - キャッシュアドレス対応リスト 1 0 0 4 から delete 操作により外し、新たな L B A に対応するものをこの該当スロット番号とともに insert 操作により追加する。また、ホストコンピュータ 2 からの要求を受けたが、キャッシュメモリ部 1 3 に、すでに格納されていた場合は、上記の手順によりアクセス順序リスト 1 7 1 3、1 7 1 4 のみを更新し、キャッシュディレクトリ情報 1 7 1 1、1 7 1 2 は更新しない。

10

20

## 【 0 0 3 0 】

以上に示した、図 1 0 B のアクセス順序リストをフラッシュメモリ 1 6 6 およびフラッシュメモリデバイス 1 6 9 に用いると次のような不都合が起こる。アクセス順序によりデステージするスロットを決定しているため、あるメディアにばかりデステージすなわち書き込みが起こってしまうことがある。例えば、スロット数が 3 のキャッシュメモリが備わっている場合に、F M 1 に対するアドレス A、B に対するデータと、F M 2 に対するアドレス C、D に対するデータを、A、B、C、A、B、D のパターンの繰り返しでライト要求された場合は、F M 2 に集中してデステージが起こってしまう。こういった状況では、メディアの書き込み制限回数に早々に到達してしまい、ストレージシステムとして十分な信頼性、可用性を提供できない恐れがある。なお、本実施例では、L R U アルゴリズムを用いて説明しているが、アクセス頻度を用いるなどの L R U 以外のアルゴリズムを用いた場合も同様のことが起こりうる。

30

## 【 0 0 3 1 】

以下では、こうした不都合を回避するキャッシュ制御方法を説明する。図 1 0 C は、好適なキャッシュ制御を行うためのアクセス順序リスト 1 7 1 3、1 7 1 4 の詳細を示したブロック図である。図 1 0 B と異なり、書き込み回数制限があるフラッシュメモリに関しては、各メディア毎（各フラッシュメモリデバイス 1 6 9 毎）にリスト 1 0 6 1 が設けられている。H D D 5 0 に関しては、各々の H D D 5 0 毎に設けても、または H D D 5 0 全体として 1 つのリスト 1 0 6 2 としてもよい。さらに、各リストにデステージ回数欄 1 0 6 3、アクセス回数欄 1 0 6 4、容量欄 1 0 6 5 を持つ表を備える。デステージ回数欄 1 0 6 3 は、そのメディアにデステージを行った総回数を示す。アクセス回数欄 1 0 6 4 は、そのメディアに領域に対して起こったリード、ライトも含めたアクセス回数を示す。これはキャッシュメモリ部 1 3 上にデータが存在した場合も回数に含めている。また容量欄 1 0 6 5 は、そのメディアの容量を示している。同じ頻度でフラッシュメモリデバイス 1 6 9 にデステージしたとしても、容量が大きいほど、フラッシュメモリ内部の 1 つのセルに書き込まれる頻度は少なくなることを考慮する場合に用いる。さらに、各メディアにリストを持つ。リストの内容は図 1 0 B の場合とほぼ同様で、L R U アルゴリズムを用いるとすると、どのスロットが何番目に最近アクセスされたのかを示す情報が格納されている。ただし、この情報は、このメディアに格納されるべきデータに関するもののみである。

40

50

なお、HDD50に関するリストにはデステージ制限欄1066がある。この欄が「なし」の場合は、フラッシュメモリとHDD50で異なるキャッシュ制御を行い、次に説明する図12Aのs1201でメディア種類をHDDと認識する。この場合、ライトキャッシュ領域132は、予めHDD50領域とフラッシュメモリ領域に分けられている。こうした分割をなくキャッシュメモリをより柔軟に用いたい場合は、この欄を「あり」とし、フラッシュメモリと同等のキャッシュ制御を行う。この場合は、メディア種類がHDDであっても、図12Aのs1201でフラッシュメモリとして処理が選択される。また、デステージ回数欄1063は、一定時間間隔でゼロにリセットしてもよい。リセットを行う際は全メディアに関してのデステージ回数欄1063をゼロにする。

#### 【0032】

次に、図10Cのアクセス順序リスト1713、1714を用いた、デステージするスロットを決定する処理の説明を、図12Aを用いて行う。これは、チャンネル制御部11のプロセッサ111により実行されるもので、デステージを行わねばならないことが判明した後の処理である図11Bのs1125以降の処理に相当する。まず、ホストコンピュータ2から受領したライト要求が、HDDに対するものかフラッシュメモリに対するものか、メディア種類を判別する(s1201)。これは、ディレクトリ情報のメディア欄1003を用いて知ることができる。メディア種類がHDDならば、HDDのアクセス順序リスト1062を用いて最も昔にアクセスされたスロットを探し、そのスロットをデステージさせるメッセージを送り、s1207に処理を移す(s1202)。メディア種類がフラッシュメモリ、もしくはHDDであってもデステージ制限欄1066が「あり」の場合は、各メディアのデステージ回数欄1063を比較して、デステージ回数が最小のメディアを探す(s1203)。そのメディアのリストの要素数を調べることで、使用しているスロットがあるかを調査する(s1204)。ないならば、次にデステージ回数の少ないメディアを探して再びs1204に処理を移す(s1205)。あるならば、そのアクセス順序リスト1061を用いて最も昔にアクセスされたスロットを探し、そのスロットをデステージするように指示する(s1206)。次に該当メディアのデステージ回数欄1063に1を足す(s1207)。さらに、該当スロットをディレクトリ情報およびアクセス順序リストから外し(s1208)、さらにステージするデータに関してディレクトリ情報およびアクセス順序リストの最新に追加する(s1209)。その後、該当スロットに今回受領したデータを書き込むことを行う(s1210)。なお、s1203でデステージ回数の最小のメディアを探すが、容量欄1065の説明で述べたとおりセル自体の書き込み回数を考慮するならば、その評価関数として、デステージ回数を容量で除算したものをを用いてもよい。

#### 【0033】

また、デステージ回数のみでなく、アクセス回数も考慮した場合の同処理を、図12Bを用いて説明する。まず、各メディアのデステージ回数欄1063とアクセス回数欄1064について調査し、 $\text{デステージ回数} \times \text{アクセス回数}$ を定数とし、 $\text{デステージ回数} \times \text{アクセス回数}$ をアクセス回数で除したものの和を評価関数とし、これが最少となるメディアを探す(s1241)。そのメディアのリストの要素数を調べることで、使用しているスロットがあるかを調査する(s1242)。ないならば、次に少ないメディアを探して再びs1242に処理を移す(s1243)。あるならば、そのアクセス順序リスト1061、1062を用いて最も昔にアクセスされたスロットを探し、そのスロットをデステージするように指示する(s1244)。以降の処理は、図12Aのs1207と同様である。この処理を用いた場合は、フラッシュメモリ領域とHDD領域に分けておく必要もなく、またメディアがHDDの場合のライト要求後に同じデータに対するリード要求が多発するなどの、キャッシュメモリを用いた性能向上が可能なパターンなどにも、柔軟に対応することができる。なお、メディアをインデックス順で選択することも可能である。

#### 【実施例2】

#### 【0034】

実施例2を説明する。図13に本発明の第2の実施の形態のストレージシステムのプロ

10

20

30

40

50

ック図を示す。この形態では、FM制御部がキャッシュメモリ部13や制御メモリ部15と同等の位置に実装されている。一般にキャッシュメモリ部13や制御メモリ部15は、チャンネル制御部11やディスク制御部14と異なり、チャンネル4やディスク側チャンネル60が接続できる位置に実装しなくてもよい。例えば、チャンネル4接続する作業が容易なように装置の前面に配置するといったことや、ディスク側チャンネル60の伝送路的制限からHDD50の付近に配置せねばならないといった実装上の制限はない。そのため、この位置に配置することにより、さらにコンパクトに実装できる。この場合は、図3あるいは図4のFM制御部16の形式が好適である。

【0035】

さらに、図14にFM制御部16を高機能にした高機能FM制御部160のブロック図を示す。さらに大規模な構成をとる場合は、高機能FM制御部160の使用がより好ましい。ここでは、FM制御部16もチャンネル制御部11同様に複数のプロセッサ111を備え、制御プログラム1121に処理方法を設定することにより高度な信頼性、可用性を得るための制御や管理性向上に向けた制御が行えるようになっている。

10

【0036】

以上実施例で説明したが、本発明の他の実施形態1は、前記ストレージ装置は、前記第1のメディアを制御する第1のメディア制御部から前記チャンネル制御部に直接に転送するパスを備えるストレージシステムである。

【0037】

本発明の他の実施形態2は、前記ストレージ装置は、前記キャッシュメモリ部から、前記第1のメディアを制御する第1のメディア制御部に直接に転送するパスを備えるストレージシステムである。

20

【0038】

本発明の他の実施形態3は、前記メディアは、書き込み回数に制限を有する第1のメディアと、書き込み回数に制限のない第2のメディアとからなるストレージシステムである。

【0039】

本発明の他の実施形態4は、前記第2のメディアは、前記第1のメディアよりも読み込み速度が遅く、消費電力が大きい、第1のメディアよりも書き込み可能回数がかかるに多いメディアであるストレージシステムである。

30

【0040】

本発明の他の実施形態5は、前記ストレージ装置は、前記第2のメディアを制御する第2のメディア制御部から、該第2のメディアに格納されたデータを一時的に格納するキャッシュメモリ部に転送するパスを備えるストレージシステムである。

【0041】

本発明の他の実施形態6は、前記チャンネル制御部は、前記ホストコンピュータからのリード要求を受領し、該リード要求の対象データが前記キャッシュメモリ部に格納されていないが前記第1のメディアに格納されている場合、該第1のメディアを制御する第1のメディア制御部に対し、該当データを直接に転送するストレージシステムである。

【0042】

本発明の他の実施形態7は、前記ストレージ装置は、前記第1のメディアに対する書き込み回数を平均化するように、デステージするデータを選択するストレージシステムである。

40

【0043】

本発明の他の実施形態8は、前記ストレージ装置は、各々の第1のメディアに対するデステージした回数を記録するストレージシステムである。

【0044】

本発明の他の実施形態9は、前記ストレージ装置は、デステージするデータを選択する際に、各々の前記第1のメディアに対するデステージした回数を比較し、デステージした回数が少ないものから優先的にデステージするデータを決定するストレージシステムであ

50

る。

【 0 0 4 5 】

本発明の他の実施形態 1 0 は、前記ストレージ装置は、各々の第 1 のメディアに格納されるべきデータに対するアクセスの時間または頻度の情報を、各々の第 1 のメディアに関連して記録するストレージシステムである。

【 0 0 4 6 】

本発明の他の実施形態 1 1 は、前記ストレージ装置は、デステージするデータを選択する際に、各々の前記第 1 のメディアに対するデステージした回数およびアクセス回数を基に評価関数を算出し、評価関数が少ないものから優先的にデステージするデータを決定するストレージシステムである。

10

【 0 0 4 7 】

本発明の他の実施形態 1 2 は、前記ストレージ装置は、キャッシュメモリ部をバックアップするバッテリーを備えるストレージシステムである。

【 0 0 4 8 】

本発明の他の実施形態 1 3 は、前記第 1 のメディア制御部は、前記チャンネル制御部の機能の一部を行う高機能メディア制御部であるストレージシステムである。

【 0 0 4 9 】

本発明の他の実施形態 1 4 は、ホストコンピュータからのデータを保存するメディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備えるストレージ制御装置において、前記メディアのうちの書き込み回数に制限を有する第 1 のメディアを制御する第 1 のメディア制御部から前記チャンネル制御部に直接に転送するパスを備えるストレージ装置である。

20

【 0 0 5 0 】

本発明の他の実施形態 1 5 は、前記キャッシュメモリ部から、前記第 1 のメディアを制御する第 1 のメディア制御部に直接に転送するパスを備えるストレージ装置である。

【 0 0 5 1 】

本発明の他の実施形態 1 6 は、前記第 1 のメディアよりも読み込み速度が遅く、消費電力が大きい第 1 のメディアよりも書き込み可能回数が多い第 2 のメディアを制御する第 2 のメディア制御部から、該第 2 のメディアに格納されたデータを一時的に格納するキャッシュメモリ部に転送するパスを備えるストレージ装置である。

30

【 0 0 5 2 】

本発明の他の実施形態 1 7 は、前記チャンネル制御部は、前記ホストコンピュータからのリード要求を受領し、該リード要求の対象データが前記キャッシュメモリ部に格納されていないが前記第 1 のメディアに格納されている場合、該第 1 のメディアを制御する第 1 のメディア制御部に対し、該当データを直接に転送するよう指示するストレージ装置である。

【 0 0 5 3 】

本発明の他の実施形態 1 8 は、ホストコンピュータからのデータを保存するメディアを制御するメディア制御部、チャンネルを介して前記ホストコンピュータに接続するチャンネル制御部、及び前記ホストコンピュータからのデータを一時的に保存する揮発メモリであるキャッシュメモリ部を備えるストレージ制御装置を制御する方法において、前記チャンネル制御部が前記ホストコンピュータからのリード要求を受領し、該要求の対象データが前記キャッシュメモリ部に格納されていないが書き込み回数に制限を有する第 1 のメディアに格納されている場合、該第 1 のメディアを制御する第 1 のメディア制御部に対し、該当データを直接転送するよう指示するストレージ装置の制御方法である。

40

【 0 0 5 4 】

本発明の他の実施形態 1 9 は、前記第 1 のメディアに対する書き込み回数を平均化するように、デステージするデータを選択するストレージ装置の制御方法である。

【 図面の簡単な説明 】

50

## 【 0 0 5 5 】

【図 1】実施例 1 のストレージシステムの構成のブロック図。

【図 2】チャンネル制御部 1 1 の詳細の構成のブロック図。

【図 3】F M 制御部 1 6 の詳細の構成のブロック図。

【図 4】F M 制御部 1 6 の別の詳細の構成のブロック図。

【図 5】F M 制御部 1 6 の別の詳細の構成のブロック図。

【図 6】内部スイッチ部 1 2 の詳細の構成のブロック図。

【図 7】ホストコンピュータ 2 から、H D D 5 0 領域へのリードの要求が来た場合の処理の流れを表す図。

【図 8】ホストコンピュータ 2 から、フラッシュメモリ領域へのリードの要求が来た場合の処理の流れを表す図。

10

【図 9】キャッシュメモリ部 1 3 ならびに制御メモリ部 1 7 に格納されるデータの詳細を示したブロック図。

【図 1 0 A】リードキャッシュディレクトリ情報 1 7 1 1、ライトキャッシュディレクトリ情報 1 1 7 2 の詳細を示したブロック図。

【図 1 0 B】従来のアクセス順序リスト 1 7 1 3、1 7 1 4 の詳細を示したブロック図。

【図 1 0 C】好適なキャッシュ制御を行うためのアクセス順序リスト 1 7 1 3、1 7 1 4 の詳細を示したブロック図。

【図 1 1 A】ホストコンピュータ 2 から、ライトの要求が来た場合で、すでにキャッシュメモリ部 1 3 上に該当アドレスのデータが存在する場合の処理の流れを表す図。

20

【図 1 1 B】ホストコンピュータ 2 から、ライトの要求が来た場合で、キャッシュメモリ部 1 3 上に該当アドレスのデータが存在しないが既に空いているスロットがない場合の処理の流れを表す図。

【図 1 2 A】デステージするスロットを決定する処理を表す図。

【図 1 2 B】ステージ回数のみでなく、アクセス回数も考慮した場合のスロット決定の処理を表す図。

【図 1 3】実施例 2 のストレージシステムのブロック図。

【図 1 4】高機能 F M 制御部 1 6 0 のブロック図。

【符号の説明】

## 【 0 0 5 6 】

30

1 ストレージ制御装置

2 ホストコンピュータ

3 S A N スイッチ

4 チャンネル

1 1 チャンネル制御部

1 2 内部スイッチ

1 3 キャッシュメモリ部

1 4 ディスク制御部

1 5 内部バス

1 6 F M 制御部

40

1 7 制御メモリ部

5 0 ハードディスクドライブ

6 0 ディスク側チャンネル

1 1 1 プロセッサ

1 1 2、1 6 4 メモリモジュール

1 1 3 周辺処理部

1 1 4 チャンネルプロトコル処理部

1 1 5 データ転送系バス

1 1 6 制御系バス

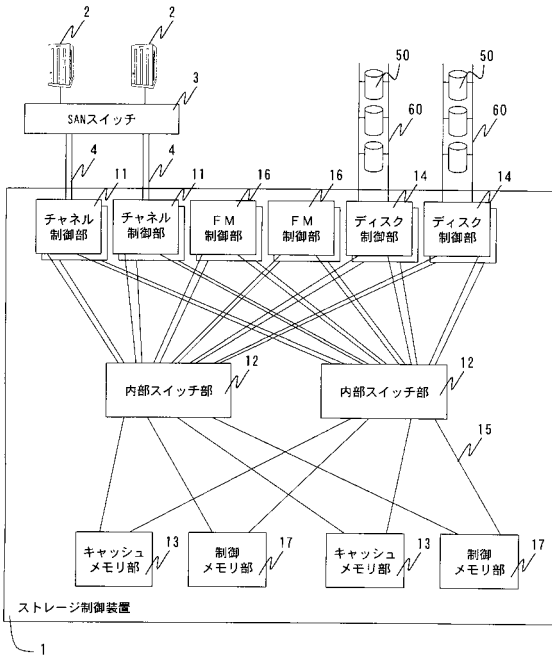
1 1 7、1 2 1、1 6 1 内部ネットワークインタフェース部

50

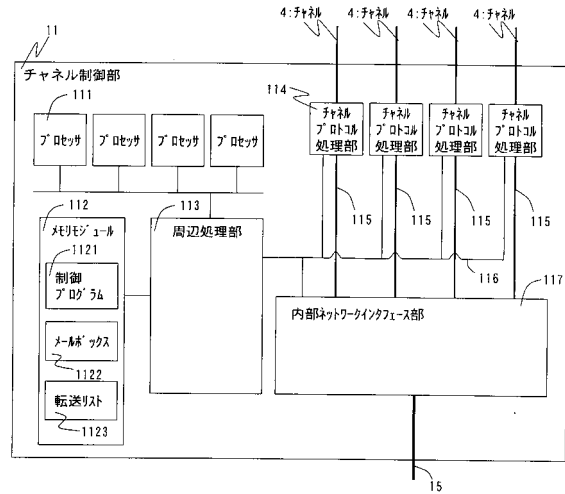
1 2 2	セレクタ	
1 2 3	チャンネル制御部 F M 制御部間接続	
1 3 1	リードキャッシュ領域	
1 3 2	ライトキャッシュ領域	
1 6 0	高機能 F M 制御部	
1 6 2	D M A コントローラ	
1 6 3	メモリコントローラ	
1 6 5	F M コントローラ	
1 6 6	フラッシュメモリ	
1 6 7	F M プロトコル処理部	10
1 6 8	コネクタ	
1 6 9	フラッシュメモリデバイス	
1 7 2	構成情報	
1 7 3	通信領域	
1 0 0 1	L U N 欄	
1 0 0 2	L B A 欄	
1 0 0 3	メディア欄	
1 0 0 4	ホスト - キャッシュアドレス対応リスト	
1 0 6 1、1 0 6 2、1 7 1 3、1 7 1 4	アクセス順序リスト	
1 0 6 3	デステージ回数欄	20
1 0 6 4	アクセス回数欄	
1 0 6 5	容量欄	
1 0 6 6	デステージ制限欄	
1 1 2 1	制御プログラム	
1 1 2 2	メールボックス	
1 1 2 3、1 6 4 1	転送リスト	
1 6 1 0	F M 側チャンネル	
1 6 9 0	緊急デステージ領域	
1 7 1 1	リードキャッシュディレクトリ情報	
1 7 1 2	ライトキャッシュディレクトリ情報	30



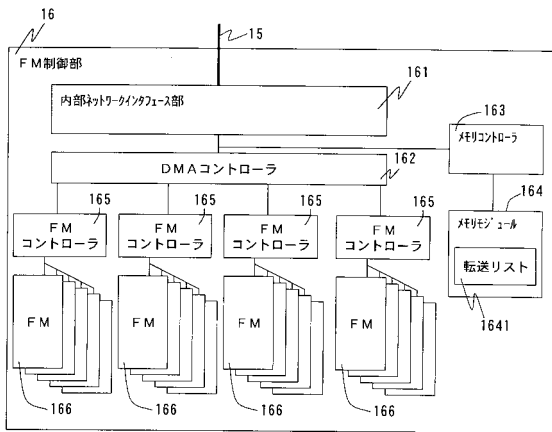
【図1】



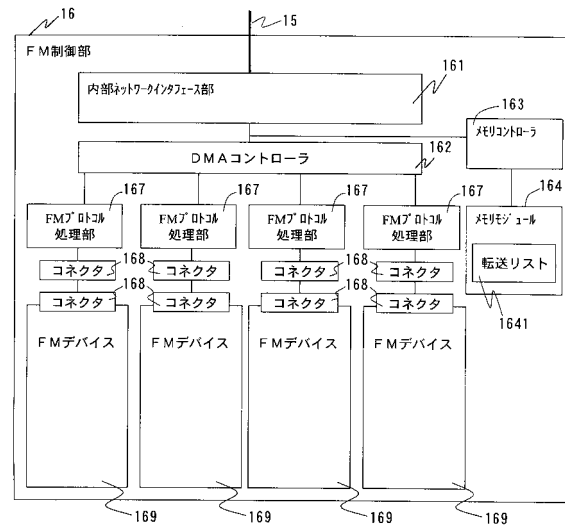
【図2】



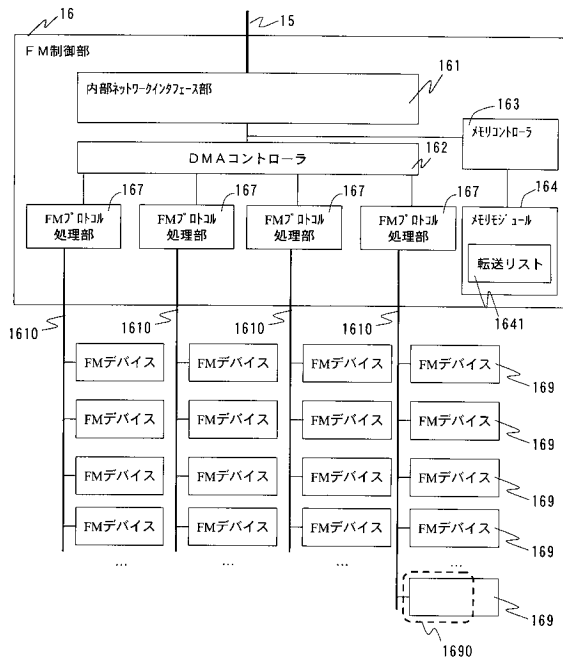
【図3】



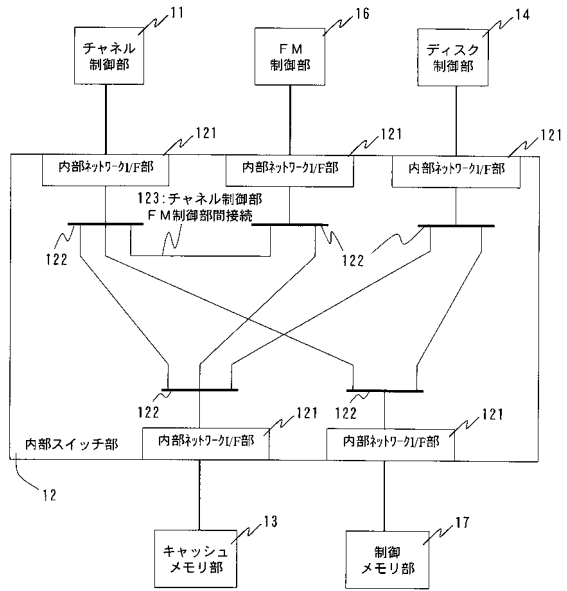
【図4】



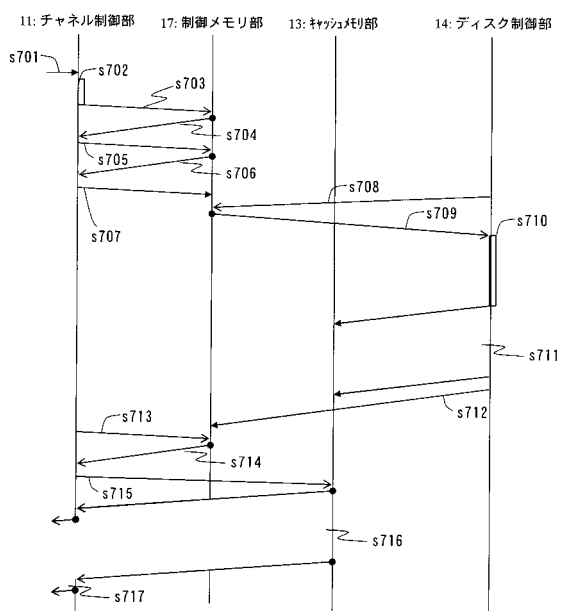
【図5】



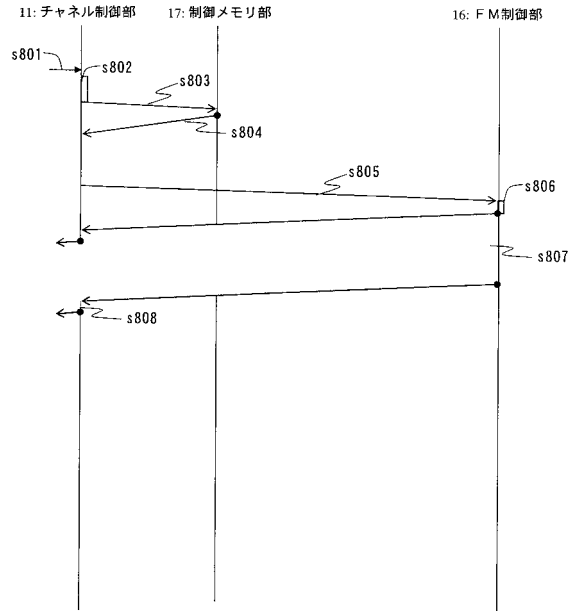
【図6】



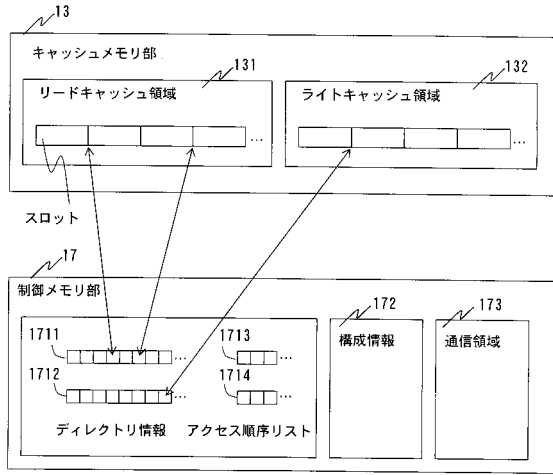
【図7】



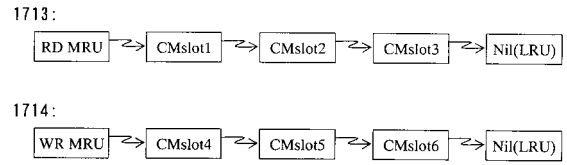
【図8】



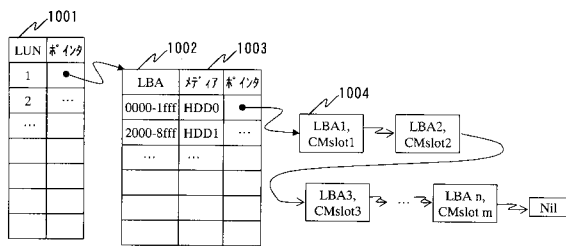
【図9】



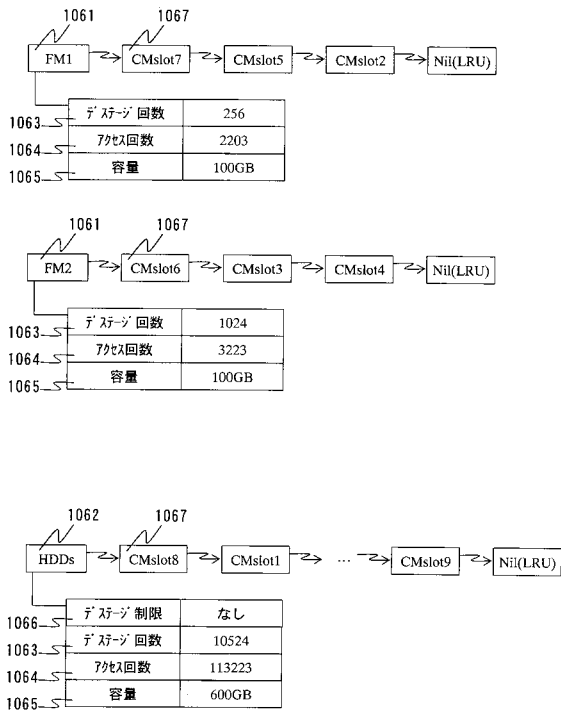
【図10B】



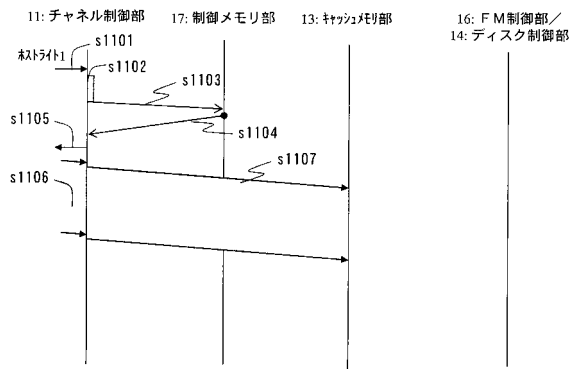
【図10A】



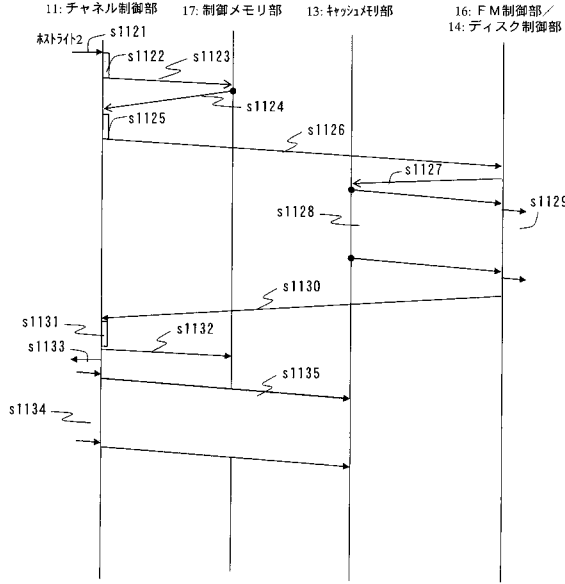
【図10C】



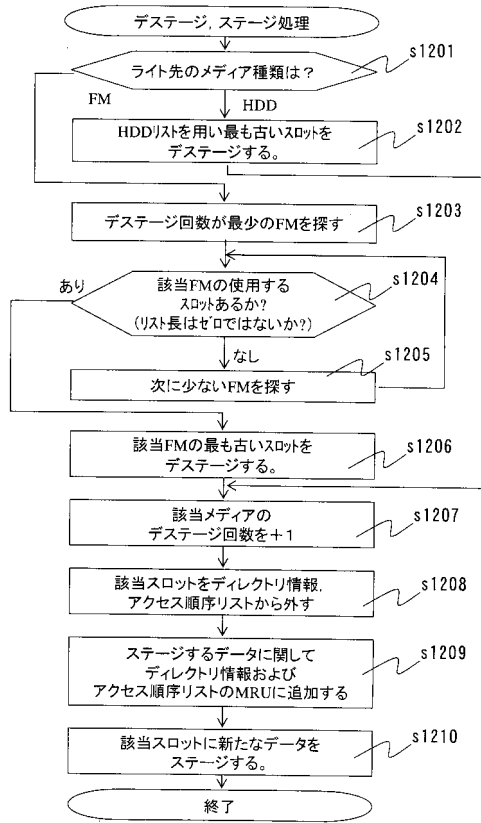
【図11A】



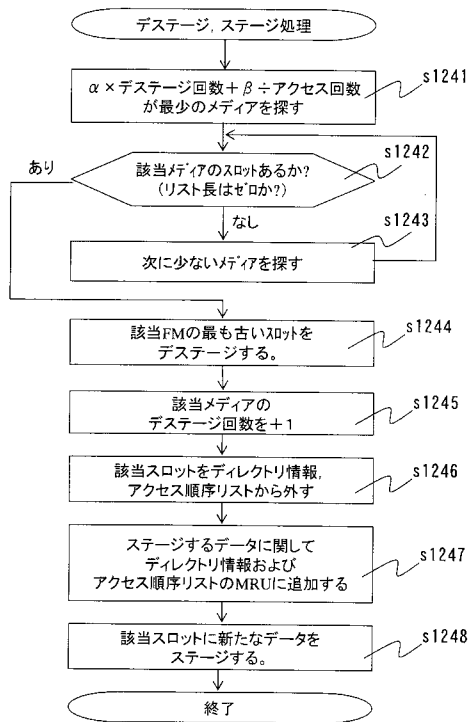
【図11B】



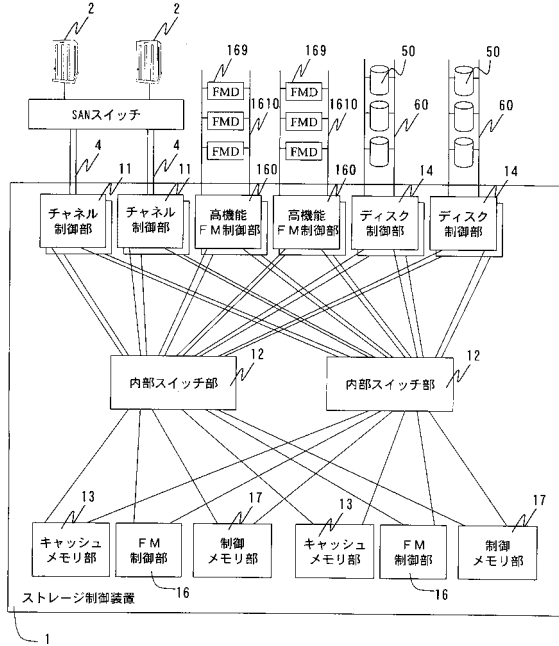
【図12A】



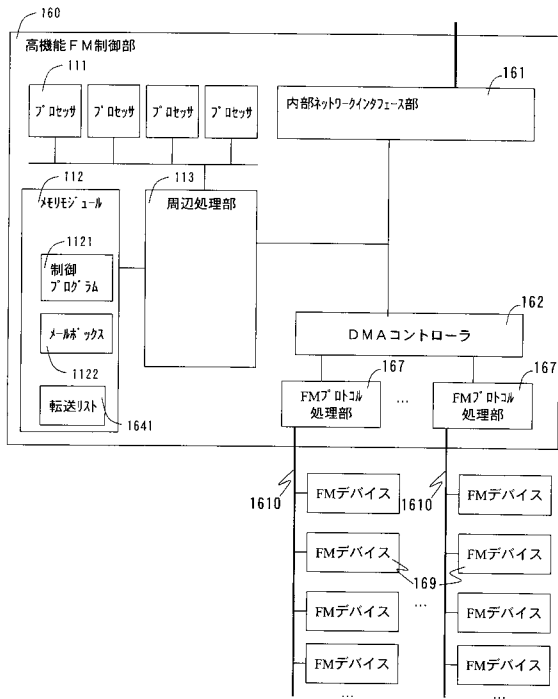
【図12B】



【図13】



【図14】



---

フロントページの続き

(72)発明者 藤林 昭

神奈川県川崎市麻生区王禅寺1099番地 株式会社 日立製作所 システム開発研究所内

審査官 菅原 浩二

(56)参考文献 特開2004-021811(JP,A)  
特開2000-276402(JP,A)  
特開2005-222550(JP,A)  
特開2003-228513(JP,A)  
特開2003-271444(JP,A)  
特開2004-078963(JP,A)  
特開平06-075836(JP,A)  
特開2004-102539(JP,A)  
特開2005-115603(JP,A)  
特開2005-115857(JP,A)  
特表2007-525753(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06