

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2005-293478

(P2005-293478A)

(43) 公開日 平成17年10月20日(2005. 10. 20)

(51) Int.Cl.⁷

F I

テーマコード (参考)

G06F 12/00

G06F 12/00 545A

5B005

G06F 3/06

G06F 12/00 514E

5B014

G06F 12/08

G06F 3/06 301S

5B065

G06F 13/10

G06F 3/06 302A

5B082

G06F 12/08 501E

審査請求 未請求 請求項の数 7 O L (全 21 頁) 最終頁に続く

(21) 出願番号 特願2004-111096 (P2004-111096)

(22) 出願日 平成16年4月5日(2004. 4. 5)

(71) 出願人 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 100095371

弁理士 上村 輝之

(74) 代理人 100089277

弁理士 宮川 長夫

(74) 代理人 100104891

弁理士 中村 猛

(72) 発明者 金井 宏樹

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 飯田 純一

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

最終頁に続く

(54) 【発明の名称】 記憶制御システム、記憶制御システムに備えられるチャネル制御装置、データ転送装置

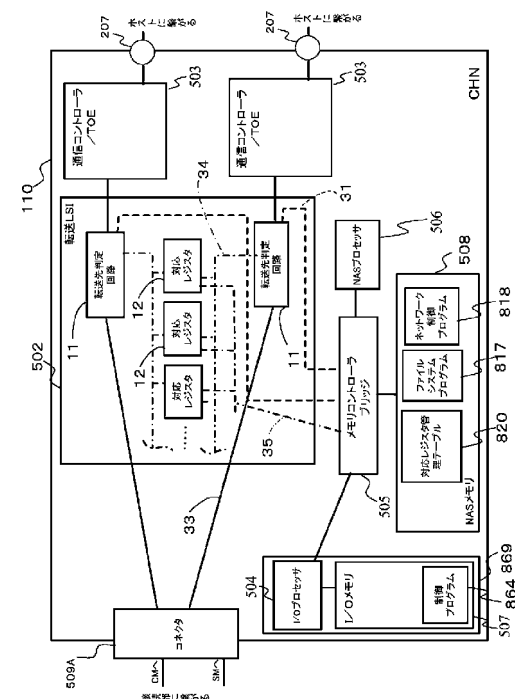
(57) 【要約】

【課題】 上位装置と下位装置との間のデータ転送を高速化する。

【解決手段】 CHN110は、ホスト端末200から受信したデータを記憶することができるNASメモリ508と、転送先判定回路11とを備える。転送先判定回路11は、ホスト端末200からアクセス要求を受信した場合、そのアクセス要求が、ホスト端末200とキャッシュ領域131との間で行われるデータ転送を伴うデータアクセス要求であり、データアクセス要求に含まれるファイルシステムアドレスに対応付けられたキャッシュアドレスが識別されたならば、そのキャッシュアドレスが表すキャッシュ領域131内のリード対象データを、NASメモリ508を経由することなくダイレクトにホスト端末200に転送するか、又は、アクセス要求に含まれるライト対象データを、NASメモリ508を経由することなくキャッシュ領域131に転送する。

【選択図】 図2

図2



【特許請求の範囲】

【請求項 1】

データを記憶する記憶デバイスと、
外部装置と前記記憶デバイスとの間でやり取りされるデータが格納されるキャッシュ領域を有するキャッシュメモリと、

外部装置からデータを受信して前記キャッシュ領域に格納したり、前記キャッシュ領域に格納されているデータを読み出して前記外部装置に転送したりするチャネル制御部と、

前記キャッシュ領域に格納されているデータを読み出して前記記憶デバイスに格納したり、前記記憶デバイスから読み出したデータを前記キャッシュ領域に格納したりする記憶デバイス制御部と

10

を備え、

前記チャネル制御部は、

前記外部装置から受信したデータを記憶することができるデータ記憶メモリと、

データファイル中のデータのファイルレベルのアドレスであるファイルシステムアドレスと、前記ファイルシステムアドレスに対応付けられたキャッシュアドレスとを記憶する対応メモリと、

前記外部装置から受信したアクセス要求が、ファイルシステムアドレスを有するファイルアクセス要求である場合、前記ファイルアクセス要求に含まれるファイルシステムアドレスに対応付けられたキャッシュアドレスが前記対応メモリから識別されたならば、前記識別されたキャッシュアドレスが表すキャッシュ領域であるターゲットキャッシュ領域内のアクセス対象データを、前記データ記憶メモリを経由することなく前記外部装置に転送する、又は、前記ファイルアクセス要求が有するアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲットキャッシュ領域に転送する転送先判定回路とを備える、

20

記憶制御システム。

【請求項 2】

前記チャネル制御部は、前記ファイルシステムアドレスを、前記記憶デバイスの記憶領域の管理単位であるブロックレベルのアドレスに変換する第 1 のプロセッサを更に備え、

前記チャネル制御部及び前記記憶デバイス制御部の少なくとも一方が、前記変換されたブロックレベルのアドレスに基づいて、前記キャッシュメモリ上に前記キャッシュ領域を確保する第 2 のプロセッサを備え、

30

前記第 1 のプロセッサが、前記ファイルアクセス要求が発行されることに先行して前記外部装置から先行コマンドを受信し、前記先行コマンドに応答して、アクセス対象となるファイルシステムアドレスを前記対応メモリに書込み、前記ファイルシステムアドレスを前記ブロックレベルのアドレスに変換し、

前記第 2 のプロセッサが、前記第 1 のプロセッサによって変換された前記ブロックレベルのアドレスに基づいて、前記キャッシュ領域を確保し、前記確保したキャッシュ領域を表すキャッシュアドレスを、前記ファイルシステムアドレスが書き込まれた前記対応メモリに書込む、

請求項 1 記載の記憶制御システム。

40

【請求項 3】

前記チャネル制御部は、前記転送先判定回路よりも前記外部装置側に、前記外部装置から受信するアクセス要求についてインターネットプロトコルの解釈を行う IP 回路と、前記外部装置から受信するアクセス要求についてトランスミッションコントロールプロトコルの解釈を行う TCP 回路とを備え、前記外部装置から前記 IP 回路及び前記 TCP 回路を経由して前記アクセス要求が前記転送先判定回路に入力されるようになっている、

請求項 1 記載の記憶制御システム。

【請求項 4】

前記対応メモリは、複数のレジスタであり、

前記複数のレジスタの各々には、1 つのファイルアクセス要求に含まれる 1 つのファイ

50

ルシステムアドレスと、そのファイルシステムアドレスに対応付けられた 1 又は複数のキャッシュアドレスとが記憶される、
請求項 1 記載の記憶制御システム。

【請求項 5】

データを記憶する記憶デバイスと、外部装置と前記記憶デバイスとの間でやり取りされるデータが格納されるキャッシュ領域を有するキャッシュメモリと、外部装置からデータを受信して前記キャッシュ領域に格納したり、前記キャッシュ領域に格納されているデータを読み出して前記外部装置に転送したりするチャンネル制御装置と、前記キャッシュ領域に格納されているデータを読み出して前記記憶デバイスに格納したり、前記記憶デバイスから読み出したデータを前記キャッシュ領域に格納したりする記憶デバイス制御部とを備える記憶制御システムに搭載される前記チャンネル制御装置であって、

10

前記外部装置から受信したデータを記憶することができるデータ記憶メモリと、

前記外部装置からアクセス要求を受信した場合、前記アクセス要求が、前記外部装置と前記キャッシュ領域との間で行われるデータ転送を伴うデータアクセス要求であり、前記データアクセス要求に含まれる情報要素に対応付けられたキャッシュアドレスが識別されたならば、前記識別されたキャッシュアドレスが表すキャッシュ領域であるターゲットキャッシュ領域内のアクセス対象データを、前記データ記憶メモリを経由することなく前記外部装置に転送する、又は、前記アクセス要求に含まれるアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲットキャッシュ領域に転送する転送先判定回路と

20

を備えるチャンネル制御装置。

【請求項 6】

データファイル中のデータのファイルレベルのアドレスであるファイルシステムアドレスと、前記ファイルシステムアドレスに対応付けられたキャッシュアドレスとを記憶する対応メモリを更に備え、

前記転送先判定回路は、前記外部装置から受信したアクセス要求が、ファイルシステムアドレスを有するファイルアクセス要求である場合、前記ファイルアクセス要求に含まれるファイルシステムアドレスに対応付けられたキャッシュアドレスが前記対応メモリから識別されたならば、前記識別されたキャッシュアドレスが表すキャッシュ領域であるターゲットキャッシュ領域内のアクセス対象データを、前記データ記憶メモリを経由することなく前記ホスト端末に転送する、又は、前記ファイルアクセス要求が有するアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲットキャッシュ領域に転送する、

30

請求項 5 記載のチャンネル制御装置。

【請求項 7】

通信装置に搭載可能なデータ転送装置において、

前記データ転送装置の上位に存在する外部装置から受信したアクセス要求が、前記外部装置と、前記データ転送装置の下位に存在する下位記憶装置の記憶領域との間で行われるデータ転送を伴うデータアクセス要求であり、前記データアクセス要求に含まれる情報要素に対応付けられた記憶アドレスが識別されたならば、前記識別された記憶アドレスが表す記憶領域であるターゲット記憶領域内のアクセス対象データを、前記通信装置内のデータ記憶メモリを経由することなく前記外部装置に転送する、又は、前記アクセス要求に含まれるアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲット記憶領域に転送する転送先判定回路を備えたデータ転送装置。

40

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、上位装置と下位装置との間でのデータ転送のための技術に関わり、例えば、RAIDシステムのような記憶制御システムや、そのような記憶制御システムに搭載可能なチャンネル制御装置に関する。

50

【背景技術】

【0002】

例えば、データセンタ等のような大規模なデータを取り扱うデータベースシステムでは、ホストコンピュータとは別に構成された記憶制御システムを用いてデータを管理する。この記憶制御システムは、例えば、多数の記憶デバイスをアレイ状に配設して構成された RAID (Redundant Array of Independent Inexpensive Disks) のようなディスクアレイシステムである。

【0003】

このような記憶制御システムには、例えば、特開 2003 - 316713 号公報に開示されているように、ファイル単位の I/O 要求を処理するネットワークチャネルアダプタ (以下、CHN) を備えることにより、NAS (Network Area Storage) となることができ 10
るものがある。また、この記憶制御システムは、ディスクと、ディスクアダプタ (DKA) と、CHN と DKA とに共有されるデータであって、ディスクに保存されるデータが格納されるキャッシュメモリとが備えられる。

【0004】

【特許文献 1】特開 2003 - 316713 号公報。

【発明の開示】

【発明が解決しようとする課題】

【0005】

上述した公報によれば、CHN にはメモリが搭載されており、そのメモリには、ファ 20
イルデータをキャッシュするため等に用いられるデータバッファが存在する。従って、NAS クライアントとキャッシュメモリとの間のデータ転送は、そのメモリ内のデータバッファを経由して行われることになる。このデータ転送を高速化することができれば、より利便性が高まると考えられる。

【0006】

従って、本発明の目的は、上位装置と下位装置との間のデータ転送を高速化することにある。

【0007】

本発明の更なる目的は、後の記載から明らかになるであろう。

【課題を解決するための手段】

【0008】

本発明の第 1 の観点に従う記憶制御システムは、データを記憶する記憶デバイスと、外部装置と前記記憶デバイスとの間でやり取りされるデータが格納されるキャッシュ領域を有するキャッシュメモリと、外部装置からデータを受信して前記キャッシュ領域に格納したり、前記キャッシュ領域に格納されているデータを読み出して前記外部装置に転送したりするチャネル制御部と、前記キャッシュ領域に格納されているデータを読み出して前記記憶デバイスに格納したり、前記記憶デバイスから読出したデータを前記キャッシュ領域に格納したりする記憶デバイス制御部とを備える。前記チャネル制御部は、前記外部装置から受信したデータを記憶することができるデータ記憶メモリと、データファイル中のデータのファイルレベルのアドレスであるファイルシステムアドレスと、前記ファイルシステム 40
アドレスに対応付けられたキャッシュアドレスとを記憶する対応メモリと、転送先判定回路を備える。転送先判定回路は、前記外部装置から受信したアクセス要求が、ファイルシステムアドレスを有するファイルアクセス要求である場合、前記ファイルアクセス要求に含まれるファイルシステムアドレスに対応付けられたキャッシュアドレスが前記対応メモリから識別されたならば、前記識別されたキャッシュアドレスが表すキャッシュ領域であるターゲットキャッシュ領域内のアクセス対象データを、前記データ記憶メモリを経由することなく前記外部装置に転送する、又は、前記ファイルアクセス要求が有するアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲットキャッシュ領域に転送する。

【0009】

10

20

30

40

50

本発明の第1の観点に従う記憶制御システムの第1の実施態様では、前記チャンネル制御部は、前記ファイルシステムアドレスを、前記記憶デバイスの記憶領域の管理単位であるブロックレベルのアドレスに変換する第1のプロセッサを更に備える。前記チャンネル制御部及び前記記憶デバイス制御部の少なくとも一方が、前記変換されたブロックレベルのアドレスに基づいて、前記キャッシュメモリ上に前記キャッシュ領域を確保する第2のプロセッサを備える。前記第1のプロセッサが、前記ファイルアクセス要求が発行されることに先行して前記外部装置から先行コマンドを受信し、前記先行コマンドに応答して、アクセス対象となるファイルシステムアドレスを前記対応メモリに書込み、前記ファイルシステムアドレスを前記ブロックレベルのアドレスに変換する。前記第2のプロセッサが、前記第1のプロセッサによって変換された前記ブロックレベルのアドレスに基づいて、前記

10

【0010】

本発明の第2の観点に従う記憶制御システムの第2の実施態様では、前記チャンネル制御部は、前記転送先判定回路よりも前記外部装置側に、前記外部装置から受信するアクセス要求についてインターネットプロトコルの解釈を行うIP回路と、前記外部装置から受信するアクセス要求についてトランスミッションコントロールプロトコルの解釈を行うTCP回路とを備える。それにより、前記外部装置から前記IP回路及び前記TCP回路を経由して前記アクセス要求が前記転送先判定回路に入力されるようになっている。

【0011】

20

本発明の第1の観点に従う記憶制御システムの第3の実施態様では、前記対応メモリは、複数のレジスタである。前記複数のレジスタの各々には、1つのファイルアクセス要求に含まれる1つのファイルシステムアドレスと、そのファイルシステムアドレスに対応付けられた1又は複数のキャッシュアドレスとが記憶される。

【0012】

本発明の第2の観点に従うチャンネル制御装置は、以下の記憶制御システム、すなわち、データを記憶する記憶デバイスと、外部装置と前記記憶デバイスとの間でやり取りされるデータが格納されるキャッシュ領域を有するキャッシュメモリと、外部装置からデータを受信して前記キャッシュ領域に格納したり、前記キャッシュ領域に格納されているデータを読み出して前記外部装置に転送したりするチャンネル制御装置と、前記キャッシュ領域に格納されているデータを読み出して前記記憶デバイスに格納したり、前記記憶デバイスから読み出したデータを前記キャッシュ領域に格納したりする記憶デバイス制御部とを備える記憶制御システムに搭載される前記チャンネル制御装置である。チャンネル制御装置は、前記外部装置から受信したデータを記憶することができるデータ記憶メモリと、転送先判定回路とを備える。転送先判定回路は、前記外部装置からアクセス要求を受信した場合、前記アクセス要求が、前記外部装置と前記キャッシュ領域との間で行われるデータ転送を伴うデータアクセス要求であり、前記データアクセス要求に含まれる情報要素に対応付けられたキャッシュアドレスが識別されたならば、前記識別されたキャッシュアドレスが表すキャッシュ領域であるターゲットキャッシュ領域内のアクセス対象データを、前記データ記憶メモリを経由することなく前記外部装置に転送する、又は、前記アクセス要求に含まれる

30

40

【0013】

本発明の第2の観点に従うチャンネル制御装置の第1の実施態様では、チャンネル制御装置は、データファイル中のデータのファイルレベルのアドレスであるファイルシステムアドレスと、前記ファイルシステムアドレスに対応付けられたキャッシュアドレスとを記憶する対応メモリを更に備える。その場合、前記転送先判定回路は、前記外部装置から受信したアクセス要求が、ファイルシステムアドレスを有するファイルアクセス要求である場合、前記ファイルアクセス要求に含まれるファイルシステムアドレスに対応付けられたキャッシュアドレスが前記対応メモリから識別されたならば、前記識別されたキャッシュアド

50

レスが表すキャッシュ領域であるターゲットキャッシュ領域内のアクセス対象データを、前記データ記憶メモリを経由することなく前記ホスト端末に転送する、又は、前記ファイルアクセス要求が有するアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲットキャッシュ領域に転送する。

【 0 0 1 4 】

本発明の第 3 の観点に従うデータ転送装置は、通信装置に搭載可能なデータ転送装置であり、転送先判定回路を備える。転送先判定回路は、前記データ転送装置の上位に存在する外部装置から受信したアクセス要求が、前記外部装置と、前記データ転送装置の下位に存在する下位記憶装置の記憶領域との間で行われるデータ転送を伴うデータアクセス要求であり、前記データアクセス要求に含まれる情報要素に対応付けられた記憶アドレスが識別されたならば、前記識別された記憶アドレスが表す記憶領域であるターゲット記憶領域内のアクセス対象データを、前記通信装置内のデータ記憶メモリを経由することなく前記外部装置に転送する、又は、前記アクセス要求に含まれるアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲット記憶領域に転送する。

10

【 0 0 1 5 】

本発明の第 4 の観点に従うデータ転送方法は、第 1 と第 2 のステップを有する。前記データ転送方法は、第 1 のステップでは、上位装置からアクセス要求を受信する。また、前記データ転送方法は、前記第 2 のステップでは、前記受信したアクセス要求が、前記上位装置と、下位装置の記憶領域との間で行われるデータ転送を伴うデータアクセス要求であり、前記データアクセス要求に含まれる情報要素に対応付けられた記憶アドレスが識別されたならば、前記識別された記憶アドレスが表す記憶領域であるターゲット記憶領域内のアクセス対象データを、通信装置内のデータ記憶メモリを経由することなく前記外部装置に転送するか、又は、前記アクセス要求に含まれるアクセス対象データを、前記データ記憶メモリを経由することなく前記ターゲット記憶領域に転送する。

20

【発明の効果】

【 0 0 1 6 】

本発明によれば、上位装置（例えば外部装置）と下位装置（例えばキャッシュメモリ）との間のデータ転送を高速化することができる。

【発明を実施するための最良の形態】

【 0 0 1 7 】

以下、図面を参照して、本発明の一実施形態について説明する。

30

【 0 0 1 8 】

図 1 は、本実施形態に係る記憶制御システムを備えるコンピュータシステムの構成を示すブロック図である。

【 0 0 1 9 】

このコンピュータシステム 1 では、通信ネットワーク 8 2 0 に、1 以上のホスト端末 2 0 0 と、記憶制御システム 6 0 0 とが接続されている。通信ネットワーク 8 2 0 は、ファイルレベルのデータがやり取りされる通信ネットワークであり、例えば、LAN、インターネット、専用回線、公衆回線等の場合に応じて適宜用いることができる。

【 0 0 2 0 】

40

1 以上のホスト端末 2 0 0 の各々は、例えば、CPU (Central Processing Unit) やメモリ等の情報処理資源を備えたコンピュータ装置であり、例えば、パーソナルコンピュータ、ワークステーション、メインフレーム等として構成される。ホスト端末 2 0 0 は、例えば、キーボードスイッチやポインティングデバイス、マイクロフォン等の情報入力装置（図示せず）と、例えば、モニタディスプレイやスピーカー等の情報出力装置（図示せず）とを備えている。さらに、ホスト端末 2 0 0 は、例えば、NAS を利用するための NAS 利用ソフトウェア 2 0 0 B と、Windows（登録商標）又は UNIX（登録商標）等の OS（オペレーティングシステム）2 0 0 C と、ネットワークドライバ 2 0 0 F とを備える。ホスト端末 2 0 0 は、例えば、ファイル名を指定してファイル単位でのデータ入出力を記憶制御システム 6 0 0 に要求する。NAS 用ソフトウェア 2 0 0 B としては、例えば

50

、OS 2 0 0 C が U N I X (登録商標) の場合、N F S (Network File System) であり、OS 2 0 0 C が W i n d o w s (登録商標) の場合、C I F S (Common Interface File System) である。また、ネットワークドライバ 2 0 0 F には、例えば、T C P (Transmission Control Protocol) に基づくデータ処理を行う T C P ドライバ 2 0 0 D と、I P (Internet Protocol) に基づくデータ処理を行う I P ドライバ 2 0 0 E とが含まれる。

【0 0 2 1】

記憶制御システム 6 0 0 は、例えば、アレイ状に配列された多数の物理記憶デバイス 3 0 0 を備える R A I D システムである。記憶制御システム 6 0 0 は、記憶制御装置 1 0 0 と記憶装置ユニット 1 0 1 とに大別することができる。記憶制御装置 1 0 0 は、例えば、複数のチャネルアダプタ N A S (以下、「CHN」と略記) 1 1 0 と、複数のディスクアダプタ (以下、D K A) 1 4 0 と、キャッシュメモリ 1 3 0 と、共有メモリ 1 2 0 と、接続部 1 5 0 とを備えている。

10

【0 0 2 2】

CHN 1 1 0 は、ホスト端末 2 0 0 との間のデータ通信を行うものである。CHN 1 1 0 は、例えば、C P U やメモリ等を備えたマイクロコンピュータシステムとして構成されており、ホスト端末 2 0 0 から受信した各種コマンドを解釈して実行する。CHN 1 1 0 には、自分を識別するためのネットワークアドレス (例えば、I P アドレスや W W N) が割り当てられている。CHN 1 1 0 は、ホスト端末 2 0 0 から通信ネットワーク 8 2 0 を介してファイル単位での I / O 要求 (例えば、ファイル名と、そのファイル名を持つファイルをリード又はライトする命令とを含んだコマンド、以下、「ファイル I / O 要求」と言う) を受けて、そのファイル I / O 要求を処理する N A S (Network Attached Storage) として振る舞うことができるようになっている。CHN 1 1 0 は、ホスト端末 2 0 0 から受信したデータをキャッシュメモリ 1 3 0 に格納したり、D K A 1 4 0 によってキャッシュメモリ 1 3 0 に格納されたデータを取得してホスト端末 2 0 0 に送信したりする。

20

【0 0 2 3】

D K A 1 4 0 は、記憶装置ユニット 1 0 1 内の論理的な記憶ユニット (以下、L U) 3 1 0 との間のデータ授受を行うものである。D K A 1 4 0 は、L U 3 1 0 を備える物理記憶デバイス 3 0 0 に接続するための図示しない通信ポートを備えている。また、D K A 1 4 0 は、C P U やメモリ等を備えたマイクロコンピュータシステムとして構成されている。D K A 1 4 0 は、CHN 1 1 0 からキャッシュメモリ 1 3 0 に書き込まれたデータを取得して L U 3 1 0 に書込んだり、また、L U 3 1 0 から読み出したデータをキャッシュメモリ 1 3 0 に格納したりする。D K A 1 4 0 は、L U 3 1 0 との間でデータ入出力を行う場合、論理的なアドレスを物理的なアドレスに変換する。

30

【0 0 2 4】

キャッシュメモリ (以下、「CM」と略記する場合有り) 1 3 0 は、例えば揮発性又は不揮発性のメモリであり、ホスト端末 2 0 0 から受信して L M 3 1 0 へ書き込まれるデータや、L U 3 1 0 から読出されてホスト端末 2 0 0 へ転送されるデータを一時的に記憶するものである。

【0 0 2 5】

共有メモリ (以下、「SM」と略記する場合有り) 1 2 0 は、例えば不揮発性のメモリであり、ホスト端末との間でやり取りされるデータに関する制御情報 (例えば、CM 1 3 0 上に確保されたどのキャッシュ領域にどのデータが格納されるべきかを示す情報) 等が格納される。また、共有メモリ 1 2 0 には、例えば、ワーク領域 (例えば、CHN 1 1 0 及び D K A 1 4 0 の C P U 間でやり取りされるメッセージを一時的に記憶する領域) が設定される。なお、図示の例では、CM 1 3 0 と SM 1 2 0 は、物理的に分離しているが、一つのメモリであっても良い。その場合、そのメモリ上のメモリ空間が、論理的に、CM 用の空間と SM 用の空間とに分けられても良い。

40

【0 0 2 6】

接続部 1 5 0 は、各 CHN 1 1 0、各 D K A 1 4 0、キャッシュメモリ 1 3 0 及び共有

50

メモリ 120 を相互に接続させる。接続部 150 は、例えば、高速スイッチング動作によってデータ伝送を行う超高速クロスバススイッチ等のような高速バスとして構成することができる。

【0027】

記憶装置ユニット 101 には、アレイ状に配列された複数の物理記憶デバイス 300 が含まれている。物理記憶デバイス 300 としては、例えば、ハードディスク、フレキシブルディスク、磁気テープ、半導体メモリ、光ディスク等のようなデバイスを用いることができる。物理記憶デバイス 300 の記憶領域上には、論理的な記憶デバイスである複数の論理ユニット（以下、「LU」と略記）310 が備えられている。

【0028】

以下、本実施形態についてより詳細に説明する。なお、以下の説明では、NAS 用ソフトウェア 200B は、NFS であるとする。

【0029】

図 2 は、CHN 110 の構成例を示すブロック図である。

【0030】

CHN 110 は、複数（又は 1 つ）の通信ポート 207 と、2 以上（又は 1 つ）の通信コントローラ / TOE 503 と、I/O プロセッサ 504 及び I/O メモリ 507 を備えた 1 つ（又は複数）の入出力制御部 869 と、メモリコントローラブリッジ 505 と、NAS プロセッサ 506 と、NAS メモリ 508 と、コネクタ 509 と、転送 LSI 502 とを備えている。

【0031】

通信コントローラ / TOE 503 は、通信コントローラと TOE（TCP/IP オフロードエンジン）とが一体になったハードウェア回路であり、通信ポート 207 及び転送 LSI 502 に接続される。通信コントローラは、例えば LAN コントローラであり、IP に基づくデータ処理を行う。TOE は、TCP に基づくデータ処理を行う。なお、通信コントローラと TOE は、物理的に分離していても良い。

【0032】

メモリコントローラブリッジ 505 は、NAS プロセッサ 506、NAS メモリ 508、I/O プロセッサ 504 及び転送 LSI 502 に接続されており、それらの間の通信を中継する LSI（Large-Scale Integrated circuit）である。メモリコントローラブリッジ 505 には、後述するメモリブリッジ転送バス 31 及びプロセッサレジスタバス 35 が接続される。

【0033】

NAS メモリ 508 は、NAS プロセッサ 506 の制御を司るプログラムを記憶することができる。NAS メモリ 508 は、例えば、ファイルシステムプログラム 817、ネットワーク制御プログラム 818、対応レジスタ管理テーブル 820 等を記憶することができる。ファイルシステムプログラム 817 は、例えば、ファイル I/O 要求に含まれているファイル名と、そのファイル名を有するファイルが格納されている場所のアドレス情報（例えば LUN 及び先頭論理ブロックアドレス）との対応づけを管理し、その対応付けに基づいて、ファイル I/O 要求をブロック I/O 要求に変換する。ネットワーク制御プログラム 818 は、例えば、NFS（Network File System）と Samba の 2 つのファイルシステムプロトコルを含んで構成される。NFS は、NFS が動作する UNIX（登録商標）オペレーティングシステムを搭載したホスト端末からのファイル I/O 要求を受け付ける。一方、Samba は、CIFS（Common Interface File System）が動作する Windows（登録商標）オペレーティングシステムを搭載したホスト端末からのファイル I/O 要求を受け付ける。対応レジスタ管理テーブル 820 には、複数の対応レジスタ 12 にそれぞれ対応した複数のレジスタ ID 及び複数の使用状態データが登録されている。

【0034】

NAS プロセッサ 506 は、例えば CPU 又はマイクロプロセッサであり、メモリコン

10

20

30

40

50

トローラブリッジ505に接続されている。NASプロセッサ506は、NASメモリ508に格納されているファイルシステムプログラム817及びネットワーク制御プログラム818等を読み出し、読み出したコンピュータプログラムに従う処理を実行することができる。例えば、NASプロセッサ506は、ネットワーク制御プログラム818により、ホスト端末200からのファイルI/O要求を受け付ける。また、NASプロセッサ506は、ファイルシステムプログラム817により、ホスト端末200から受信しNASメモリ508に格納されたファイルI/O要求をブロックI/O要求に変換してI/Oプロセッサ504に出力することができる（なお、ブロックとは、LU310上の記憶領域におけるデータの管理単位である）。また、NASプロセッサ506は、データ転送を伴うアクセス要求がホスト端末200から出力される場合、それに先行して、所定の先行コマンドをホスト端末200から受ける。その場合、NASプロセッサ506は、メモリコントローラブリッジ505及びプロセッサレジスタバス35を介して、複数の対応レジスタ12の中から選択された対応レジスタ12に、ファイルシステムアドレスを登録する。また、NASプロセッサ506は、上記選択された対応レジスタ12のレジスタIDを、I/Oプロセッサ504に通知する。 10

【0035】

I/Oプロセッサ504は、例えばCPU又はマイクロプロセッサであり、I/Oメモリ507から読み出した制御プログラム864により、接続部150との間のデータの授受や、NASプロセッサ506と接続部150との間のデータ通信の中継や、キャッシュメモリ130上のキャッシュアドレスの管理等を実行することができる。また、I/Oプロセッサ504は、NASプロセッサ506からのブロックI/O要求に応答して、キャッシュメモリ130上にキャッシュ領域を確保すると共に、確保されたキャッシュ領域を表すキャッシュアドレスを、メモリコントローラブリッジ505及びプロセッサレジスタバス35を介して、NASプロセッサ506から通知されたレジスタIDに対応した対応レジスタ12に登録することができる。 20

【0036】

I/Oメモリ507は、I/Oプロセッサ504の制御を司るコンピュータプログラム等を格納する。

【0037】

コネクタ509は、接続部150に接続される。具体的には、例えば、コネクタ509は、接続部150に含まれるCM転送経路（CM130に接続された転送経路）及びSM転送経路（SM120に接続されて転送経路）に接続される。なお、必ずしも、一つのコネクタ509に、CM転送経路とSM転送経路とが混在する必要はない。例えば、コネクタ509の代わりに、CM転送経路に接続される第1のコネクタと、SM転送経路に接続される第2のコネクタとが備えられても良い。 30

【0038】

転送LSI502は、例えば、通信コントローラ/TOE503、メモリコントローラブリッジ505、キャッシュメモリ130及び共有メモリ120の相互の通信を可能とするためのLSI（Large-Scale Integrated circuit）である。転送LSI502は、複数（又は1つ）の対応レジスタ12と、1以上（例えば通信コントローラ/TOE503と同数）の転送先判定回路11とを備えている。 40

【0039】

複数の対応先レジスタ12の各々には、その対応先レジスタ12とプロセッサ504又は506との間でやり取りされるデータが経由するプロセッサレジスタバス35と、後述する判定レジスタバス34とが接続される。複数の対応レジスタ12の各々には、後に詳述するように、ファイルシステムアドレスと、そのファイルシステムアドレスに対応付けられたキャッシュアドレスとが登録される。なお、転送先LSI502に備えられる対応先レジスタ12の数は、例えば、CHN110が同時に受けることができるホスト端末200からのアクセス数と同数である。

【0040】

転送先判定回路 11 は、例えば純粋なハードウェア回路である。転送先判定回路 11 には、複数の対応先レジスタ 12 に対するアクセス経路となる判定レジスタバス 34 と、コネクタ 590 を介してキャッシュメモリ 130 や共有メモリ 120 との間でやり取りされるデータが経由するコネクタ転送バス 33 と、メモリコントローラブリッジ 505 との間でやり取りされるデータが経由するメモリブリッジ転送バス 35 とが接続される。転送先判定回路 11 は、例えば、ホスト端末 200 から通信コントローラ / T O E 503 を介して受信したアクセス要求が、データ転送を伴うアクセス要求であるか否かを判別し、その判別の結果、データ転送を伴うアクセス要求であり、且つ、複数の対応レジスタ 12 の少なくとも 1 つからアクセス先キャッシュアドレスを認識できる場合には、認識されたアクセス先キャッシュアドレスとの間で、N A S メモリ 508 を介すること無くダイレクトにデータ転送を行う。具体的には、例えば、転送先判定回路 11 は、ホスト端末 200 から通信コントローラ / T O E 503 を介してリード要求（例えば、「NFSPROC_READ」）を受信した場合、そのリード要求に含まれているファイルシステムアドレスを記憶した対応レジスタ 12 が存在するならば、その対応レジスタ 12 に記憶されているキャッシュアドレスが示すキャッシュ領域からコネクタ転送バス 33 を介してリード対象データを読み出し、そのリード対象データを、N A S メモリ 508 を介することなくダイレクトに通信コントローラ / T O E 503 を介してホスト装置 200 に転送する。また、例えば、転送先判定回路 11 は、ホスト端末 200 から通信コントローラ / T O E 503 を介してライト要求（例えば、「NFSPROC_WRITE」）を受信した場合、そのライト要求に含まれているファイルシステムアドレスを記憶した対応レジスタ 12 が存在するならば、そのライト要求が有するライト対象データを、その対応レジスタ 12 に記憶されているキャッシュアドレスが示すキャッシュ領域に、N A S メモリ 508 を介することなくコネクタ転送バス 33 を介してダイレクトに転送する。

【0041】

以上が、CHN 110 の構成例である。なお、上述したプロセッサレジスタバス 35、判定レジスタバス 34 及びメモリブリッジ転送バス 31 は、論理的なバスを示す。転送 L S I 502 内の各構成要素とメモリコントローラブリッジ 505 とは、P C I バス等の物理的な単一バス（図示せず）で互いに接続される（この点は、後述の図 7 についても同様である）。

【0042】

図 3 は、対応先レジスタ 12 の構成例を示す。

【0043】

図 3 に示すように、対応先レジスタ 12 には、ホスト端末 200 の N F S 200 B から出力されるファイルシステムアドレスと、そのファイルシステムアドレスに対応付けられた 1 又は複数のキャッシュアドレスとが登録される。

【0044】

ファイルシステムアドレスとは、アクセス対象のデータファイル（例えば、ホスト端末 200 が L U 310 に対して入出力するユーザファイル）のアドレスに関する情報、例えば、アクセス対象のデータファイル中のどの場所からのどれぐらいの長さのデータのことを指すかを示す情報である。具体的には、例えば、ファイルシステムアドレスには、ファイルハンドルと、オフセットと、データ長とが含まれている。ファイルハンドルは、例えば、アクセス対象のデータファイルの先頭アドレスを表す。オフセットとは、ファイルハンドルからのオフセットのことであり、具体的には、アクセス対象のデータファイルにおける実際のアクセス対象先頭アドレスを表す。データ長とは、アクセス対象のデータファイルにおける上記オフセットからのデータ長のことであり、例えば、図示のように、データファイル 30 が M バイト単位でアクセス（書込み又は読込み）される場合には、データ長は M を表す。

【0045】

1 又は複数のキャッシュアドレスとは、キャッシュメモリ 130 上の記憶領域を示すアドレスのことであり、例えば、1 つのファイルシステムアドレスに基づいてキャッシュメ

10

20

30

40

50

メモリに書込まれる又は読み出されるデータサイズ分の記憶領域を表す。具体的には、例えば、キャッシュメモリ 130 がキャッシュブロック 130 S 単位で管理されている場合、1 つのキャッシュアドレスは 1 つのキャッシュブロック 130 S を表し、キャッシュアドレスは、1 つのファイルシステムアドレスに対して K ($K \geq 1$) 個登録される。具体的には、例えば、データファイル 30 が M バイト (例えば 4096 バイト) 単位でアクセスされ、キャッシュメモリ 130 上の各キャッシュブロック 130 S の領域サイズが N バイト (例えば 512 バイト) の場合、キャッシュアドレスの数 K は、 $M \div N$ の商の値 (例えば $4096 / 512 = 8$) となる。なお、余りが生じた場合は、商の値に 1 が加えられてもよい。

【0046】

図 4 は、対応レジスタ管理テーブル 820 の構成例を示す。

【0047】

対応レジスタ管理テーブル 820 には、複数の対応レジスタ 12 にそれぞれ対応した複数のレジスタ ID 及び複数の使用状態データが登録される。

【0048】

レジスタ ID とは、対応する対応レジスタの識別情報 (例えば番号) である。

【0049】

使用状態データとは、対応する対応レジスタ 12 の使用状態を示すデータであり、例えば、使用中か否かを示す情報 (例えば、使用中であれば「1」であり、未使用であれば「0」) である。

【0050】

以下、本実施形態で行われるデータ転送処理流れの例を説明する。

【0051】

図 5 は、ホスト端末 200 のユーザが NFS 200 B に対してファイルリードを要求した場合に行なわれる処理流れを示す。

【0052】

ホスト端末 200 の NFS 200 B は、ユーザからファイルリード要求を受けた場合 (ステップ S1)、要求されたファイル名を有するデータファイルを取得するためにどこにアクセスすれば良いかを要求する先行コマンドを送信する。具体的には、例えば、ホスト端末 200 の NFS 200 B は、ファイルハンドラの要求 (例えば「NFSPROC_LOOKUP」) を送信する (S2)。

【0053】

NAS プロセッサ 506 は、ホスト端末 200 の NFS 200 B から通信コントローラ / TOE 503 を介して受けたファイルハンドラ要求に従って、ファイルシステムプログラム 817 に基づくファイル処理を行う (S3)。そして、NAS プロセッサ 506 は、そのファイル処理によって取得されたファイルハンドラ (受信したファイルハンドラ要求に対応したファイルハンドラ) を、ホスト端末 200 の NFS 200 B に送信する (S4)。

【0054】

また、NAS プロセッサ 506 は、複数の対応レジスタ 12の中から未使用の対応レジスタ 12 を探す (S5)。具体的には、NAS プロセッサ 506 は、メモリコントローラブリッジ 505 を介して対応レジスタ管理テーブル 820 を参照し、対応レジスタ管理テーブル 820 に記録されている複数の使用状態データの中から未使用を示す使用状態データを探す。

【0055】

S5において、未使用の対応レジスタ 12 が探し出された場合 (S5でY)、NAS プロセッサ 506 は、探し出された未使用の対応レジスタ 12 を確保する (S6)。具体的には、NAS プロセッサ 506 は、対応レジスタ管理テーブル 820 において、探し出された対応レジスタ 12 に対応する使用状態データを「未使用」から「使用中」に変更する。

10

20

30

40

50

【 0 0 5 6 】

また、N A S プロセッサ 5 0 6 は、S 5 で Y の場合、上記取得されたファイルハンドラに基づくファイルシステムアドレスを、メモリコントローラブリッジ 5 0 5 及びプロセッサレジスタバス 3 5 を介して、S 6 で確保した対応レジスタ 1 2 に書込む (S 7) 。

【 0 0 5 7 】

また、N A S プロセッサ 5 0 6 は、S 5 で Y の場合、上記取得されたファイルハンドラに基づくファイルシステムアドレスを、L U 3 1 0 の記憶領域の管理単位であるブロックレベルのアドレスに変換し、そのブロックレベルのアドレスを有するリード要求 (以下、ブロックリード要求) と、S 6 で確保された対応レジスタ 1 2 のレジスタ I D とを、I / O プロセッサ 5 0 4 に送信する (S 8) 。

10

【 0 0 5 8 】

I / O プロセッサ 5 0 4 は、レジスタ I D 及びブロックリード要求を受信した場合、ブロックリード要求に基づいて、キャッシュメモリ 1 3 0 上に、そのブロックリード要求に従って読み出すデータ (以下、リード対象データ) を格納するためのキャッシュ領域 1 3 1 を確保する (S 9) 。そして、I / O プロセッサ 5 0 4 は、D K A 1 4 0 上の図示しない I / O プロセッサとプロセッサ間通信を行うことにより、受信したブロックリード要求に含まれるアドレスが表す場所 (L U 3 1 0 内の場所) からリード対象データを読み出し、そのリード対象データを、上記確保したキャッシュ領域 1 3 1 に格納する (S 1 0) 。

【 0 0 5 9 】

また、I / O プロセッサ 5 0 4 は、確保したキャッシュ領域 1 3 1 を表す K 個のキャッシュアドレスを、メモリコントローラブリッジ 5 0 5 及びプロセッサレジスタバス 3 5 を介して、上記受信したレジスタ I D に対応する対応レジスタ 1 2 に書込む (S 1 1) 。

20

【 0 0 6 0 】

さて、ホスト端末 2 0 0 の N F S 2 0 0 B は、S 4 の処理によってファイルハンドラを受信した場合、受信したファイルハンドラに基づくファイルシステムアドレスを含んだファイルリード要求を、C H N 1 1 0 に送信する (S 1 2) 。

【 0 0 6 1 】

C H N 1 1 0 の転送先判定回路 1 1 は、ファイルシステムアドレスを含んだファイルリード要求を通信コントローラ / T O E 5 0 3 を介して受信する。受信したファイルリード要求は、通信コントローラ / T O E 5 0 3 において、既に、I P 及び T C P に基づくデータ処理が行われている。

30

【 0 0 6 2 】

転送先判定回路 1 1 は、受信したアクセス要求はキャッシュメモリ 1 3 0 へのデータ転送を伴うアクセス要求であり、且つ、そのアクセス要求に対応したキャッシュアドレスが複数の対応レジスタ 1 2 のいずれかに記録されているか否かを判断する (S 1 3) 。具体的には、例えば、転送先判定回路 1 1 は、受信したアクセス要求が「NFSPROC_READ」又は「NFSPROC_WRITE」であり、且つ、そのアクセス要求に含まれるファイルシステムアドレスを記憶する対応レジスタ 1 2 を複数の対応レジスタ 1 2 の中に存在するか否かを判断する。

【 0 0 6 3 】

40

S 1 3 において、否定的な判断結果が得られた場合、すなわち、以下の (1) 又は (2) の場合、

(1) 受信したアクセス要求がキャッシュメモリ 1 3 0 へのデータ転送を伴うアクセス要求では無い (例えば「NFSPROC_READ」及び「NFSPROC_WRITE」のいずれでもない) と判断された場合、

(2) 受信したアクセス要求がキャッシュメモリ 1 3 0 へのデータ転送を伴うアクセス要求 (例えば「NFSPROC_READ」) であっても、受信したアクセス要求に含まれるファイルシステムアドレスを記憶する対応レジスタ 1 2 が複数の対応レジスタ 1 2 の中から見つからなかった場合、

転送先判定回路 1 1 は、N A S プロセッサ 5 0 6 にデータ読出し指令を出力する。それに

50

より、N A S プロセッサ 5 0 6 から I / O プロセッサ 5 0 4 に、ブロックリード要求が出力される。そして、そのブロックリード要求を受信した I / O プロセッサ 5 0 4 によって、キャッシュメモリ 1 3 0 にステージングされているリード対象データが読み出されて N A S メモリ 5 0 8 に格納され (S 1 5 及び S 1 6)、N A S プロセッサ 5 0 6 によって、N A S メモリ 5 0 8 に格納されたリード対象データが読み出されて、ホスト端末 2 0 0 に転送される (S 1 7 及び S 1 8)。

【 0 0 6 4 】

一方、S 1 3 において、肯定的な判断結果が得られた場合、すなわち、受信したアクセス要求がキャッシュメモリ 1 3 0 へのデータ転送を伴うアクセス要求 (例えば「NFSPROC_READ」) であり、且つ、受信したアクセス要求に含まれるファイルシステムアドレスを記憶する対応レジスタ 1 2 が複数の対応レジスタ 1 2 の中から見つかった場合、転送先判定回路 1 1 は、ダイレクトデータ転送処理を行う。具体的には、転送先判定回路 1 1 は、その見つかった対応レジスタ 1 2 に記憶されている K 個のキャッシュアドレスを判定レジスタバス 3 4 を介して取得し、取得された K 個のキャッシュアドレスが示すキャッシュ領域 1 3 1 (例えば K 個のキャッシュブロック 1 3 0 S) から、コネクタ転送バス 3 3 を経由してリード対象データを読み出し、読出したリード対象データを、N A S メモリ 5 0 8 を介することなく通信コントローラ / T O E 5 0 3 を介してホスト端末 2 0 0 に転送する (S 1 9)。

【 0 0 6 5 】

以上が、ホスト端末 2 0 0 のユーザが N F S 2 0 0 B に対してファイルリードを要求した場合に行なわれる処理流れである。なお、この処理流れにおいて、S 4 の処理は、S 1 1 の後に行われても良い。また、S 1 2 のファイルリード要求に含まれるファイルシステムアドレスは、例えば、転送先判定回路 1 1 によって、N A S メモリ 5 0 8 に格納されても良いし、I / O プロセッサ 5 0 4 によって、I / O メモリ 5 0 7 に格納されても良い。

【 0 0 6 6 】

図 6 は、ホスト端末 2 0 0 のユーザが N F S 2 0 0 B に対してファイルライトを要求した場合に行なわれる処理流れを示す。なお、以下の説明において、図 5 を参照して説明した部分と重複する部分については説明を省略或いは簡略する。

【 0 0 6 7 】

ホスト端末 2 0 0 の N F S 2 0 0 B は、ユーザからファイルライト要求を受けた場合 (ステップ S 2 1)、要求されたファイル名を有するデータファイルを書込むためにどこにアクセスすれば良いかを要求するためのファイルハンドラ要求を送信する (S 2 2)。

【 0 0 6 8 】

N A S プロセッサ 5 0 6 は、受信したファイルハンドラ要求に従ってファイル処理を行い (S 2 3)、そのファイル処理によって取得されたファイルハンドラを、ホスト端末 2 0 0 の N F S 2 0 0 B に送信する (S 2 4)。

【 0 0 6 9 】

また、N A S プロセッサ 5 0 6 は、複数の対応レジスタ 1 2 の中から未使用の対応レジスタ 1 2 を探す (S 2 5)。

【 0 0 7 0 】

S 2 5 において、未使用の対応レジスタ 1 2 が探し出された場合 (S 2 5 で Y)、N A S プロセッサ 5 0 6 は、探し出された未使用の対応レジスタ 1 2 を確保する (S 2 6)。

【 0 0 7 1 】

また、N A S プロセッサ 5 0 6 は、S 2 5 で Y の場合、上記取得されたファイルハンドラに基づくファイルシステムアドレスを、メモリコントローラブリッジ 5 0 5 及びプロセッサレジスタバス 3 5 を介して、S 2 6 で確保した対応レジスタ 1 2 に書込む (S 2 7)。

【 0 0 7 2 】

また、N A S プロセッサ 5 0 6 は、S 2 5 で Y の場合、上記取得されたファイルハンドラに基づくファイルシステムアドレスをブロックレベルのアドレスに変換し、そのブロッ

10

20

30

40

50

クレベルのアドレスを有するライト要求（以下、ブロックライト要求）と、S 2 6で確保された対応レジスタ12のレジスタIDとを、I/Oプロセッサ504に送信する（S 2 8）。

【0073】

I/Oプロセッサ504は、レジスタID及びブロックライト要求を受信した場合、ブロックライト要求に基づいて、キャッシュメモリ130上に、そのブロックライト要求に従って書込むデータ（以下、ライト対象データ）を格納するためのキャッシュ領域131を確保する（S 2 9）。そして、I/Oプロセッサ504は、DKA140上の図示しないI/Oプロセッサとプロセッサ間通信を行うことにより、受信したブロックライト要求に含まれるアドレスが表す場所（LU310内の場所）から空データ（例えば全て「0」のビットから成るデータ）を読み出し、その空データを、上記確保したキャッシュ領域131に格納する（S 3 0）。なお、このS 3 0の処理は行わなくても良い。

【0074】

また、I/Oプロセッサ504は、確保したキャッシュ領域131を表すK個のキャッシュアドレスを、メモリコントローラブリッジ505及びプロセッサレジスタバス35を介して、上記受信したレジスタIDに対応する対応レジスタ12に書込む（S 3 1）。

【0075】

さて、ホスト端末200のNFS200Bは、S 2 4の処理によってファイルハンドラを受信した場合、受信したファイルハンドラに基づくファイルシステムアドレスを含んだファイルライト要求を、CHN110に送信する（S 3 2）。

【0076】

CHN110の転送先判定回路11は、ファイルシステムアドレス及びライト対象データを含んだファイルライト要求を通信コントローラ/TOE503を介して受信する。転送先判定回路11は、受信したアクセス要求はキャッシュメモリ130へのデータ転送を伴うアクセス要求であり、且つ、そのアクセス要求に対応したキャッシュアドレスが複数の対応レジスタ12のいずれかに記録されているか否かを判断する（S 3 3）。

【0077】

S 3 3において、否定的な判断結果が得られた場合、転送先判定回路11は、NASプロセッサ506にデータ書込み指令を出力する（S 3 4）。それにより、NASプロセッサ506によってNASメモリ508にライト対象データが書込まれ、且つ、NASプロセッサ506からブロックライト要求がI/Oプロセッサ504に出力される。そして、ブロックライト要求を受けたI/Oプロセッサ504によって、NASメモリ508に格納されているライト対象データが読み出されてキャッシュ領域131に格納される（S 3 5及びS 3 6）。

【0078】

一方、S 3 3において、肯定的な判断結果が得られた場合、すなわち、受信したアクセス要求がキャッシュメモリ130へのデータ転送を伴うアクセス要求（例えば「NFSPROC_WRITE」）であり、且つ、受信したアクセス要求に含まれるファイルシステムアドレスを記憶する対応レジスタ12が複数の対応レジスタ12の中から見つかった場合、転送先判定回路11は、ダイレクトデータ転送処理を行う。具体的には、転送先判定回路11は、その見つかった対応レジスタ12に記憶されているK個のキャッシュアドレスを判定レジスタバス34を介して取得し、取得されたK個のキャッシュアドレスが示すキャッシュ領域131（例えばK個のキャッシュブロック130S）に、NASメモリ508を介することなく、コネクタ転送バス33を経由してライト対象データを転送する（S 3 7）。

【0079】

以上が、ホスト端末200のユーザがNFS200Bに対してファイルライトを要求した場合に行なわれる処理流れである。なお、この処理流れにおいて、S 2 4の処理は、S 3 1の後に行われても良い。また、S 3 2のファイルライト要求に含まれるファイルシステムアドレスは、例えば、転送先判定回路11によって、NASメモリ508に格納されても良いし、I/Oプロセッサ504によって、I/Oメモリ507に格納されても良い

10

20

30

40

50

。

【 0 0 8 0 】

以上、上述した実施形態によれば、データ転送を伴うアクセス要求がホスト端末 2 0 0 から発行されることに先行して、所定の先行コマンドがホスト端末 2 0 0 から発行される。CHN 1 1 0 において所定の先行コマンドが検出された場合、そのアクセス要求でやり取りされるデータのキャッシュ領域が確保されると共に、そのキャッシュ領域を表すキャッシュアドレスが対応レジスタ 1 2 上でファイルシステムに対応付けられる。また、ホスト端末 2 0 0 から出力されたアクセス要求は、ホスト端末 2 0 0 の T C P ドライバ 2 0 0 D 及び I P ドライバ 2 0 0 E によって T C P 及び I P に基づく処理が施されるが、そのアクセス要求の I P 層及び T C P 層の部分は、通信コントローラ / T O E 5 0 3 によって解釈されるので、転送先判定回路 1 1 が受信するアクセス要求は、N A S 用ソフトウェア（例えば N F S ） 2 0 0 B から出力されたアクセス要求と同じものである。これらのことにより、ホスト端末 2 0 0 から、実際に、データ転送を伴うアクセス要求が発行された場合は、そのアクセス要求に従うデータは、そのアクセス要求に含まれるファイルシステムアドレスとそれに対応付けられたキャッシュアドレスとに基づいて、CHN 1 1 0 上のメモリを介することなく、ホスト端末 2 0 0 とキャッシュメモリ 1 3 0 との間で直接やり取りされる。これにより、ホスト端末 2 0 0 とキャッシュメモリ 1 3 0 との間のデータのやり取りが高速化される。これは、ホスト端末 2 0 0 から出力されたライト対象データをキャッシュメモリ 1 3 0 に転送する場合には特に効果的であると考えられる。また、これは、ホスト端末 2 0 0 とキャッシュメモリ 1 3 0 との間でやり取りされるデータが長くなればなるほど効果的であると考えられる。

【 0 0 8 1 】

また、上述した実施形態によれば、N A S プロセッサ 5 0 6 が、ファイルレベルの I / O 要求をブロックレベルの I / O 要求に変換する仕様になっており、且つ、I / O プロセッサ 5 0 4 が、キャッシュメモリ 1 3 0 を管理する（例えばキャッシュメモリ 1 3 0 上にキャッシュ領域を確保する）仕様になっている場合、それらの仕様を変更することなく、上述した高速化を実現することができる。

【 0 0 8 2 】

ところで、上述した実施形態について、例えば幾つかの変形例が考えられる。以下、各変形例について述べる。

【 0 0 8 3 】

（ 1 ）第 1 の変形例。

【 0 0 8 4 】

図 7 は、本実施形態の第 1 の変形例における CHN 1 1 0 の構成例を示す。なお、以下の説明では、上述した実施形態の説明と重複する部分については説明を省略或いは簡略する。

【 0 0 8 5 】

第 1 の変形例では、転送先判定回路 1 1 は、通信コントローラ / T O E 3 に搭載されている。また、この場合、通信コントローラ / T O E 3 において、I P や T C P の解釈或いは処理を行う I P / T C P 部 4 は、転送先判定回路 1 1 の場所よりも上位の場所（つまりホスト端末側の場所）に搭載されている。通信コントローラ / T O E 3 がアクセス要求を受信した場合に、I P 及び T C P の解釈をした後のアクセス要求が転送先判定回路 1 1 に受信されるようにするためである。

【 0 0 8 6 】

この第 1 の変形例では、上述した実施形態と同様に、転送先判定回路 1 1 に、メモリブリッジ転送バス 3 1 や、判定レジスタバス 3 4 を接続することができる。また、コネクタ転送バス 3 3 を介して、転送 L S I 5 と転送先判定回路 1 1 とを接続することができる。この場合、例えば、転送先判定回路 1 1 からダイレクトデータ転送処理により出力されたライト対象データは、コネクタ転送バス 3 3 及びデータ転送 L S I 5 を介して、キャッシュメモリ 1 3 0 に転送される。

10

20

30

40

50

【 0 0 8 7 】

(2) 第 2 の 変 形 例

図 8 は、本実施形態の第 2 の変形例における、CHN 1 1 0 と DKA 1 4 0 との間のデータ通信例を示す。

【 0 0 8 8 】

第 2 の変形例では、CHN 1 1 0 には、I/O プロセッサは搭載されず、CHN 1 1 0 内の I/O プロセッサ 5 0 4 が行う処理を、DKA 1 4 0 の I/O プロセッサ 6 0 3 に行わせる。

【 0 0 8 9 】

例えば、図 8 (A) に示すように、CHN 1 1 0 と DKA 1 4 0 との間に、専用の割込み線 5 1 0 が備えられる。この場合、NAS プロセッサ 5 0 6 から出力されたブロックレベルの I/O 要求は、専用の割込み線 5 1 0 を介して、DKA 1 4 0 の I/O プロセッサ 6 0 3 に送信される。これにより、DKA 1 4 0 の I/O プロセッサ 6 0 3 によって、例えば、図 5 の S 9 ~ S 1 1 の処理や、図 6 の S 2 9 ~ S 3 1 の処理が行われる。

【 0 0 9 0 】

また、例えば、図 8 (B) に示すように、CHN 1 1 0 上の NAS メモリ 5 0 8 に、コマンドキュー 5 1 1 が設けられる。NAS プロセッサ 5 0 6 から出力されたブロックレベルの I/O 要求は、そのコマンドキュー 5 1 1 に格納される (S 5 1)。DKA 1 4 0 の I/O プロセッサ 6 0 3 は、コマンドキューをポーリングし (S 5 2)、I/O 要求が存在することを検出した場合には、その I/O 要求をコマンドキューから取得する (S 5 3)。

【 0 0 9 1 】

(3) 第 3 の 変 形 例。

【 0 0 9 2 】

図 9 は、本実施形態の第 3 の変形例に係る対応レジスタ管理テーブル 1 8 2 0 の構成例を示す。

【 0 0 9 3 】

対応レジスタ管理テーブル 8 2 0 には、各対応レジスタ毎に、更に、ファイルシステムアドレスのアドレス登録領域が存在する。各対応レジスタに対応したアドレス登録領域には、その対応レジスタに読み出されるデータのファイルシステムアドレスが登録される。

【 0 0 9 4 】

例えば、I/O プロセッサ 5 0 4 は、図 5 の S 1 1 の後に、キャッシュアドレスの書込みが済んだことを NAS プロセッサ 5 0 6 に通知する (S 1 1 A)。NAS プロセッサ 5 0 6 は、その通知を受けた場合、S 6 で確保された対応レジスタ 1 2 に対応する対応レジスタ管理テーブル 8 2 0 上のアドレス登録領域に、S 7 で書き込んだファイルシステムアドレスを書込む (1 1 B)。これにより、そのファイルシステムアドレスに対応するデータを記憶するキャッシュ領域 1 3 0 のキャッシュアドレスが登録済みの対応レジスタ 1 2 がどれであるかがわかるようになる。

【 0 0 9 5 】

このような処理を行っておくことで、以後、ファイルハンドラ要求に応答してファイルハンドラが取得された場合には、そのファイルハンドラに基づくファイルシステムアドレスに対応したキャッシュアドレスが既に登録されている対応レジスタ 1 2 が存在するか否かの判定 (以下、「登録判定処理」) を実行することができる。なお、S 1 1 B の処理は、例えば、S 7 の際に行なわれても良い。

【 0 0 9 6 】

図 1 0 は、この第 3 の変形例に係る登録判定処理が行われるタイミングの例を示す。

【 0 0 9 7 】

図 1 0 (A) は、リードの場合の例である。すなわち、例えば、図 5 の S 3 又は S 4 の後に、NAS プロセッサ 5 0 6 は、登録判定処理を実行する (S 4 A)。具体的には、例えば、NAS プロセッサ 5 0 6 は、取得されたファイルハンドラに基づくファイルシステ

10

20

30

40

50

ムアドレスと同一のファイルシステムアドレスが対応レジスタ管理テーブル 1 8 2 0 に存在するか否かを判断する。その結果、その同一のファイルシステムアドレスが対応レジスタ管理テーブル 1 8 2 0 から見つからない場合 (S 4 A で N)、上述した S 5 が行われる。一方、 S 4 A で、同一のファイルシステムアドレスが対応レジスタ管理テーブル 1 8 2 0 から見つかった場合 (S 4 A で Y)、 S 5 ~ 1 1 が行われることなく S 1 2 が実行される。

【 0 0 9 8 】

図 1 0 (B) は、ライトの場合の例である。すなわち、例えば、図 6 の S 2 3 又は S 2 4 の後に、 N A S プロセッサ 5 0 6 は、登録判定処理を実行する (S 2 4 A)。その結果、取得されたファイルハンドラに基づくファイルシステムアドレスと同一のファイルシステムアドレスが対応レジスタ管理テーブル 1 8 2 0 から見つからない場合 (S 2 4 A で N)、上述した S 2 5 が行われる。一方、 S 2 4 A で、上記同一のファイルシステムアドレスが対応レジスタ管理テーブル 1 8 2 0 から見つかった場合 (S 2 4 A で Y)、 S 2 5 ~ 3 1 が行われることなく S 3 2 が実行される。

10

【 0 0 9 9 】

ホスト端末 2 0 0 から、同一のファイルリード要求が連続して発行される場合がある。また、ホスト端末 2 0 0 から発行されるファイルライト要求に含まれるライト対象データは、過去にファイルリード要求を発行することによって読み出されたデータである場合がある。それらの場合、既に、図 5 に示した処理流れによって、ファイルリード要求或いはファイルライト要求に含まれるファイルシステムアドレスに対応したキャッシュアドレスが既に登録された対応レジスタ 1 2 が存在することになる。このような場合に、わざわざ、図 5 の S 6 ~ S 1 1 の処理や、図 6 の S 2 6 ~ S 3 1 の処理が行われるのは無駄である。上述した第 3 の変形例によれば、その無駄を防ぐことができる。

20

【 0 1 0 0 】

以上、本発明の実施形態及び変形例を説明したが、これらは本発明の説明のための例示であって、本発明の範囲をこの実施形態及び変形例にのみ限定する趣旨ではない。本発明は、他の種々の形態でも実施することが可能である。例えば、 C H N 1 1 0 と D K A 1 4 0 は一体に作られていても良い。また、例えば、複数の対応レジスタ 1 2 に代えて、同一のメモリ上で、複数のファイルシステムアドレスと、複数のファイルシステムアドレスの各々に対応する K 個のキャッシュアドレスとが記録されても良い。

30

【図面の簡単な説明】

【 0 1 0 1 】

【図 1】本発明の一実施形態に係る記憶制御システムを備えるコンピュータシステムの構成を示すブロック図である。

【図 2】 C H N 1 1 0 の構成例を示すブロック図である。

【図 3】対応レジスタ 1 2 の構成例を示す。

【図 4】対応レジスタ管理テーブル 8 2 0 の構成例を示す。

【図 5】ホスト端末 2 0 0 のユーザが N F S 2 0 0 B に対してファイルリードを要求した場合に行なわれる処理流れを示す。

【図 6】ホスト端末 2 0 0 のユーザが N F S 2 0 0 B に対してファイルライトを要求した場合に行なわれる処理流れを示す。

40

【図 7】本実施形態の第 1 の変形例における C H N 1 1 0 の構成例を示す。

【図 8】本実施形態の第 2 の変形例における、 C H N 1 1 0 と D K A 1 4 0 との間のデータ通信例を示す。

【図 9】本実施形態の第 3 の変形例に係る対応レジスタ管理テーブル 1 8 2 0 の構成例を示す。

【図 1 0】本実施形態の第 3 の変形例に係る登録判定処理が行われるタイミングの例を示す。

【符号の説明】

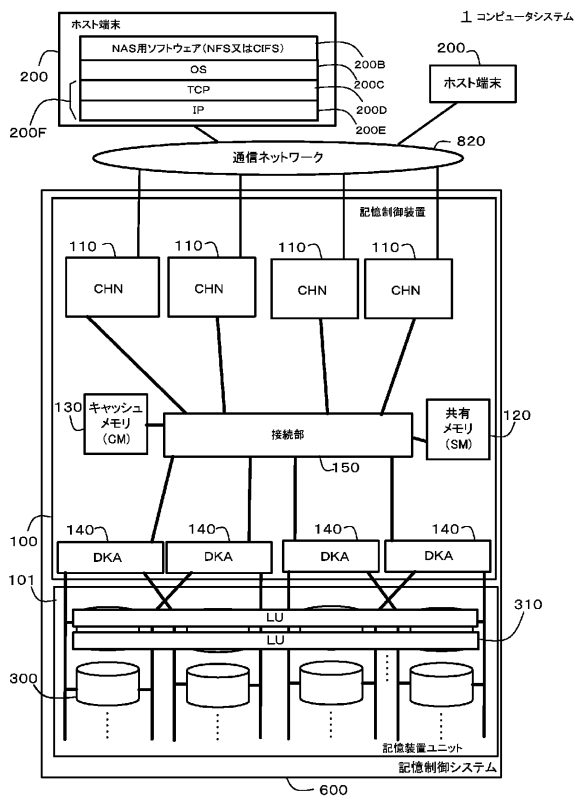
【 0 1 0 2 】

50

100 ... 記憶制御装置 101 ... 記憶装置ユニット 103 ... データベース 104 ... データファイル
 110 ... チャンネルアダプタNAS (CHN) 120 ... 共有メモリ 130 ... キャッシュメモリ
 140 ... ディスクアダプタ 150 ... 接続部 200 ... ホスト端末
 200B ... NAS用ソフトウェア 200C ... OS 200D ... TCPドライバ 200E ... IPドライバ
 200F ... ネットワークドライバ 310 ... 論理ユニット 600 ... 記憶制御システム

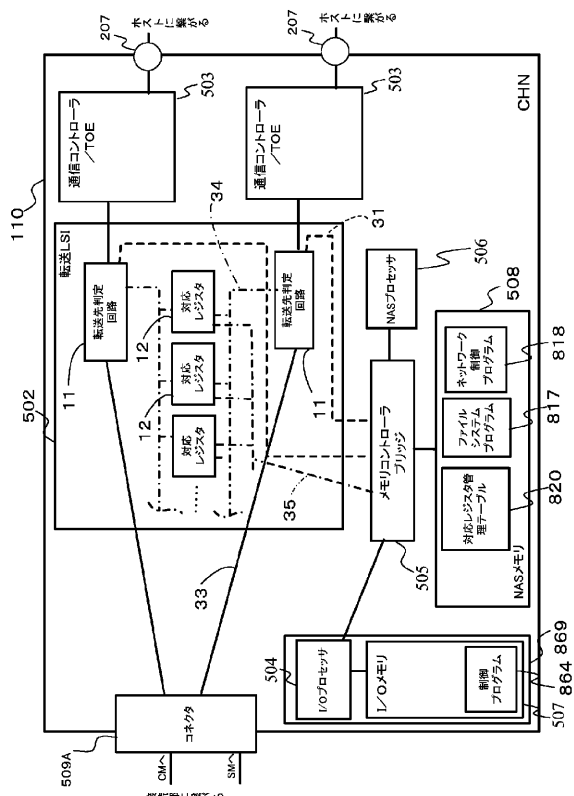
【図1】

図1

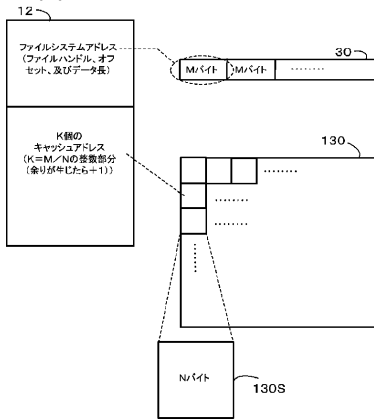


【図2】

図2



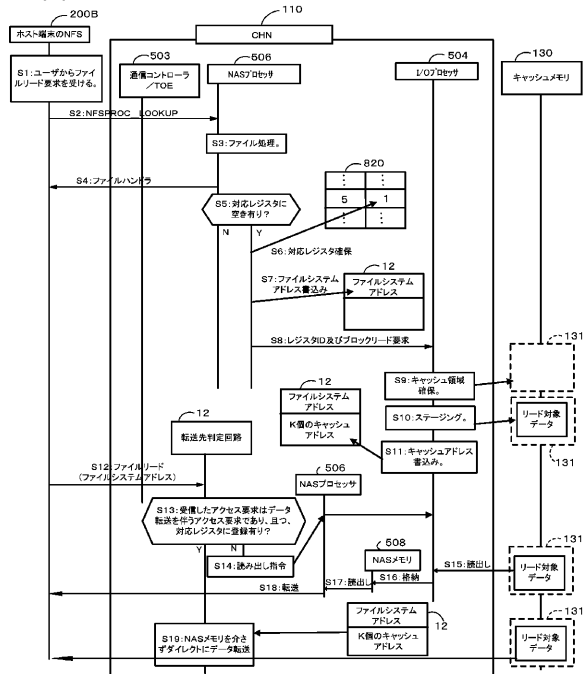
【 図 3 】



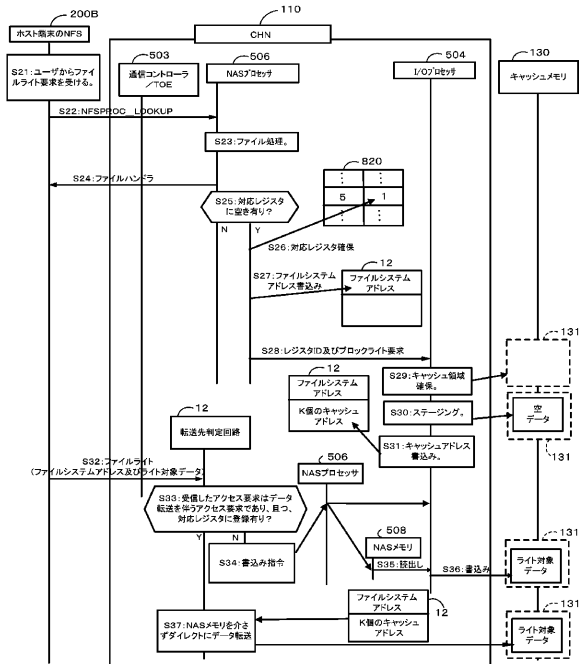
【 図 4 】

対応レジスタID	使用状態データ
1	1
2	0
3	1
⋮	⋮

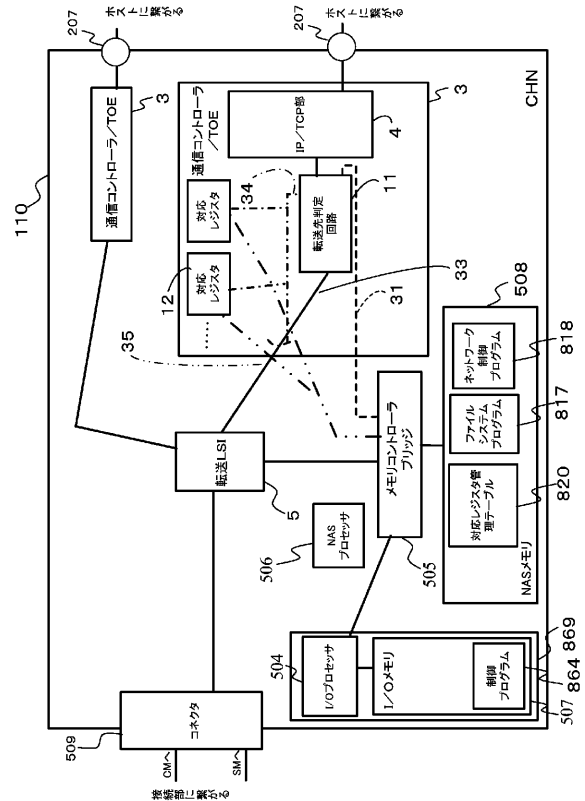
【 図 5 】



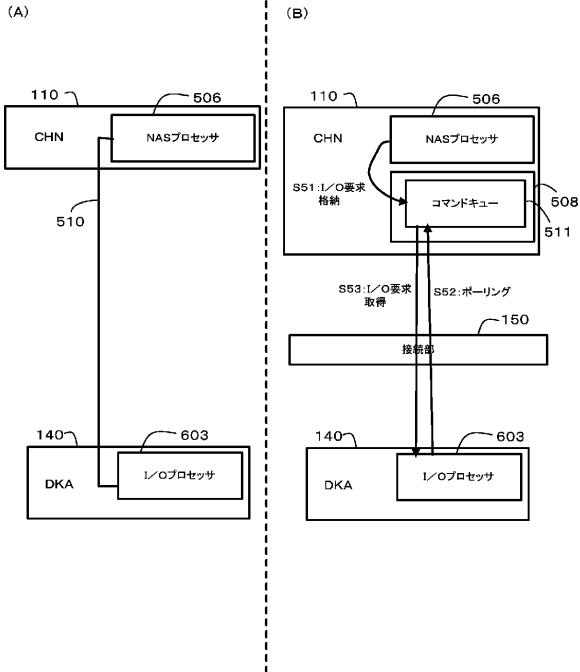
【 図 6 】



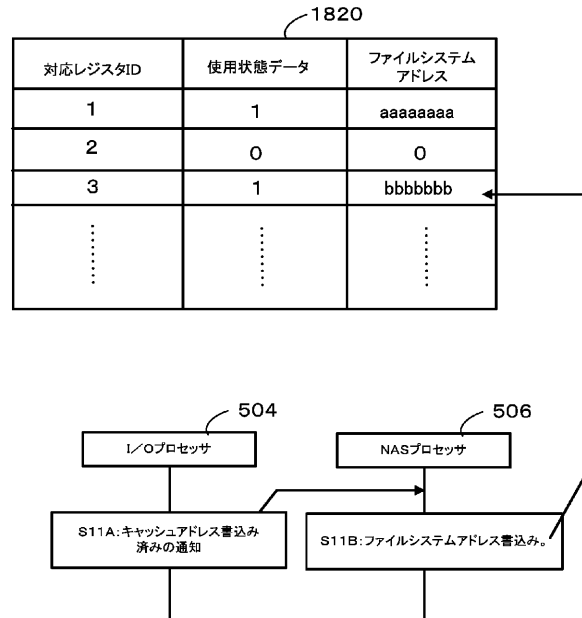
【圖 7】



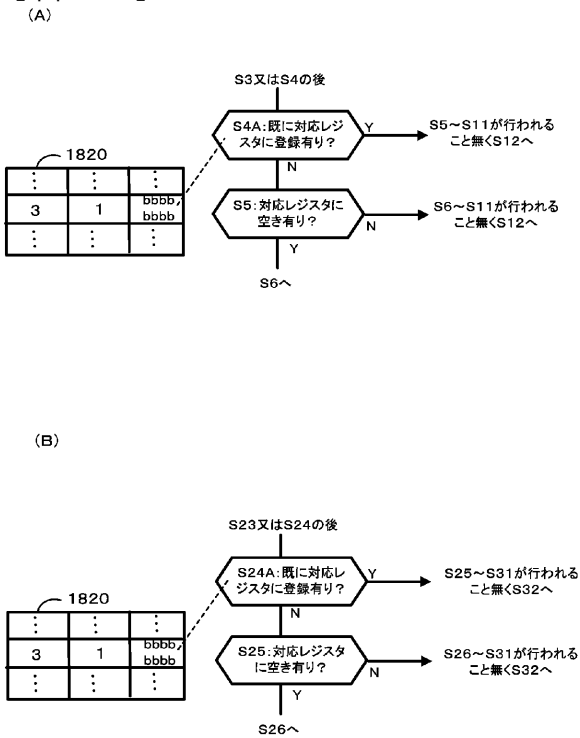
【図 8】



【図 9】



【図 10】



フロントページの続き(51) Int.Cl.⁷

F I

テーマコード(参考)

G 0 6 F	12/08	5 0 3 Z
G 0 6 F	12/08	5 1 7 B
G 0 6 F	12/08	5 5 1 Z
G 0 6 F	12/08	5 5 7
G 0 6 F	13/10	3 4 0 B

F ターム(参考) 5B005 JJ12 MM11 NN12

5B014 EB05 GC07

5B065 BA01 BA06 CA12 CA30 CC08 CH01 CH20

5B082 FA04 FA12