

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6294569号  
(P6294569)

(45) 発行日 平成30年3月14日 (2018. 3. 14)

(24) 登録日 平成30年2月23日 (2018. 2. 23)

(51) Int. Cl.	F I				
<b>G06F 13/10</b>	<b>(2006.01)</b>	G06F 13/10	340A		
<b>G06F 3/06</b>	<b>(2006.01)</b>	G06F 3/06	302A		

請求項の数 15 (全 18 頁)

(21) 出願番号	特願2017-524248 (P2017-524248)	(73) 特許権者	000005108
(86) (22) 出願日	平成27年6月19日 (2015. 6. 19)		株式会社日立製作所
(86) 国際出願番号	PCT/JP2015/067676		東京都千代田区丸の内一丁目6番6号
(87) 国際公開番号	W02016/203629	(74) 代理人	110000279
(87) 国際公開日	平成28年12月22日 (2016. 12. 22)		特許業務法人ウィルフォート国際特許事務所
審査請求日	平成29年8月10日 (2017. 8. 10)		所
		(72) 発明者	井澤 信介
			東京都千代田区丸の内一丁目6番6号 株
			株式会社日立製作所内
		(72) 発明者	杉本 定広
			東京都千代田区丸の内一丁目6番6号 株
			株式会社日立製作所内
		(72) 発明者	坂下 悠貴
			東京都千代田区丸の内一丁目6番6号 株
			株式会社日立製作所内

最終頁に続く

(54) 【発明の名称】 ストレージシステム及びキャッシュ制御方法

(57) 【特許請求の範囲】

【請求項1】

ホストシステムに接続されたストレージシステムであって、  
 複数のストレージ装置を有し、  
 前記複数のストレージ装置の各々が、前記ホストシステムに接続されており、  
 前記複数のストレージ装置の各々が、記憶デバイスと、その記憶デバイスに対して入出力されるデータがキャッシュされるキャッシュメモリ領域とを有し、  
 各ストレージ装置が、そのストレージ装置内のキャッシュメモリ領域だけでなく、そのストレージ装置に接続されているいずれの他のストレージ装置内のキャッシュメモリ領域も管理しており、いずれの他のストレージ装置のキャッシュメモリ領域にデータを入出力することができ、

論理ボリュームのアドレスを指定したI/O (Input/Output) 要求を前記複数のストレージ装置のうちのいずれかが受信した場合、前記複数のストレージ装置のうちの少なくとも1つが、対象I/Oパターンと接続形態とのうちの少なくとも1つに基づいて、前記受信したI/O要求に従うデータであるI/Oデータがキャッシュされるキャッシュ先ストレージ装置を決定し、

前記対象I/Oパターンは、複数のI/Oパターンのうち、前記受信したI/O要求に従うI/Oが属するI/Oパターンであり、

前記複数のI/Oパターンの各々は、前記論理ボリュームにおけるI/O先アドレスがランダムとなるかシーケンシャルとなるかに関するパターンであり、

前記接続形態は、I/O要求を受信するストレージ装置が確定しているか否かである、ストレージシステム。

【請求項2】

前記対象I/Oパターンが、ランダムライト及びランダムリードのうちのいずれかの場合、前記キャッシュ先ストレージ装置が、前記I/O要求を受信したストレージ装置である、請求項1記載のストレージシステム。

【請求項3】

前記対象I/Oパターンが、シーケンシャルリードの場合、前記キャッシュ先ストレージ装置が、前記I/Oデータを格納している記憶デバイスを有するストレージ装置である、請求項1記載のストレージシステム。

10

【請求項4】

前記接続形態が、I/O要求を受信するストレージ装置が確定していることを意味している場合、前記キャッシュ先ストレージ装置が、前記I/O要求を受信したストレージ装置であるホスト側ストレージ装置である、請求項1記載のストレージシステム。

【請求項5】

前記ホスト側ストレージ装置が、複数のストレージコントローラを有し、複数のストレージコントローラの各々が、前記ホストシステムに接続されており、複数のストレージコントローラの各々が、キャッシュメモリ領域を有し、前記接続形態が、I/O要求を受信するストレージコントローラが確定していることを意味している場合、前記I/Oデータがキャッシュされるキャッシュ先ストレージコントローラは、前記複数のストレージコントローラのうち、I/O要求を受信することが確定しているストレージコントローラである確定ストレージコントローラである、請求項4記載のストレージシステム。

20

【請求項6】

I/O要求を受信するストレージ装置が確定していないことを前記接続形態が意味している場合、前記複数のストレージ装置のうちの少なくとも1つが、前記対象I/Oパターンに基づいて前記キャッシュ先ストレージ装置を決定する、請求項4記載のストレージシステム。

30

【請求項7】

前記対象I/Oパターンが、ランダムライト及びランダムリードのうちのいずれかの場合、前記キャッシュ先ストレージ装置が、前記ホスト側ストレージ装置である、請求項6記載のストレージシステム。

【請求項8】

前記ホスト側ストレージ装置が、複数のストレージコントローラを有し、複数のストレージコントローラの各々が、前記ホストシステムに接続されており、複数のストレージコントローラの各々が、キャッシュメモリ領域を有し、前記接続形態が、I/O要求を受信するストレージコントローラが確定していることを意味している場合、前記I/Oデータがキャッシュされるキャッシュ先ストレージコントローラは、前記複数のストレージコントローラのうち、I/O要求を受信することが確定しているストレージコントローラである確定ストレージコントローラであり、

40

前記接続形態が、I/O要求を受信するストレージコントローラが確定していないことを意味しており、且つ、前記対象I/Oパターンが、ランダムライト及びランダムリードのうちのいずれかの場合、前記キャッシュ先ストレージコントローラは、前記複数のストレージコントローラからのいずれかである、請求項7記載のストレージシステム。

【請求項9】

前記対象I/Oパターンが、シーケンシャルリードの場合、前記キャッシュ先ストレ

50

ジ装置が、前記 I / O データを格納している記憶デバイスを有するストレージ装置であり、

前記キャッシュ先ストレージ装置のキャッシュメモリ領域に、シーケンシャルリードに従い読み出され得るデータが先読みされる、  
請求項 8 記載のストレージシステム。

【請求項 10】

前記対象 I / O パターンが、ランダムライト、ランダムリード及びシーケンシャルリードのうちのいずれでもない場合、前記キャッシュ先ストレージ装置が、前記複数のストレージ装置のうちのいずれかである、  
請求項 9 記載のストレージシステム。 10

【請求項 11】

前記複数のストレージ装置の各々が、A L U A (Asymmetric Logical Unit Access) を使用して前記ホストシステムと通信するようになっており、  
前記確定ストレージコントローラは、前記複数のストレージコントローラと前記ホストシステム間の複数のパスのうちの最も優先度が高いパスに接続されているストレージコントローラである、  
請求項 10 記載のストレージシステム。

【請求項 12】

前記確定ストレージコントローラは、前記複数のストレージコントローラと前記ホストシステム間の複数のパスのうち唯一の通信可能状態パスに接続されているストレージコントローラである、  
請求項 10 記載のストレージシステム。 20

【請求項 13】

前記確定ストレージコントローラは、前記複数のストレージコントローラと前記ホストシステム間の複数のパスのうち A L U A (Asymmetric Logical Unit Access) に従う優先度が最も高いパスに接続されているストレージコントローラである、  
請求項 11 記載のストレージシステム。

【請求項 14】

スケールアウト型のストレージシステムである請求項 10 記載のストレージシステム。

【請求項 15】 30

それぞれホストシステムに接続された複数のストレージ装置を有するストレージシステムのキャッシュ制御方法であって、

前記複数のストレージ装置の各々が、記憶デバイスと、その記憶デバイスに対して入出力されるデータがキャッシュされるキャッシュメモリ領域とを有し、

各ストレージ装置が、そのストレージ装置内のキャッシュメモリ領域だけでなく、そのストレージ装置に接続されているいずれの他のストレージ装置内のキャッシュメモリ領域も管理しており、いずれの他のストレージ装置のキャッシュメモリ領域にデータを入出力することができ、

前記キャッシュ制御方法が、

論理ボリュームのアドレスを指定した I / O (Input/Output) 要求を前記複数のストレージ装置のうちのいずれかにより受信し 40

対象 I / O パターンと接続形態とのうちの少なくとも 1 つに基づいて、前記受信した I / O 要求に従うデータである I / O データがキャッシュされるキャッシュ先ストレージ装置を決定し、

前記対象 I / O パターンは、複数の I / O パターンのうち、前記受信した I / O 要求に従う I / O が属する I / O パターンであり、

前記複数の I / O パターンの各々は、前記論理ボリュームにおける I / O 先アドレスがランダムとなるかシーケンシャルとなるかに関するパターンであり、

前記接続形態は、I / O 要求を受信するストレージ装置が確定しているか否かである、  
キャッシュ制御方法。 50

**【発明の詳細な説明】****【技術分野】****【0001】**

本発明は、概して、ストレージシステムのキャッシュ制御に関する。

**【背景技術】****【0002】**

近年、ソーシャルネットワーキングサービスや金融、医療及び交通といった社会インフラに関する膨大なデータを分析することにより、新たな価値を生み出すビッグデータ分析と呼ばれる技術が実用化されつつある。

10

**【0003】**

ビッグデータ分析においては、社会インフラから収集される入力データ及び分析結果である出力データの容量が共に非常に大きく、時間と共に増加し続ける。ストレージシステムは、このような膨大なビッグデータを安全に保管及び管理するプラットフォームとして、企業ITシステムにおいて重要な位置づけとされている。

**【0004】**

例えば、クラウドサービスを提供する企業は、導入コストを削減するために、サービスの初期稼働時は当面必要となる資源でストレージシステムを構築する。ストレージシステムとして、例えば、スケールアウト型ストレージシステムが採用される。すなわち、サービスが稼働し、資源の利用量が増加してきたら、ストレージ装置（ノード）の台数を増やすことで、システム全体の処理能力を向上させることが行われる。

20

**【0005】**

ストレージシステムの処理能力を向上させる1つの手段として、ストレージシステムのキャッシュ制御が考えられる。キャッシュ制御に関して、例えば特許文献1の技術が知られている。特許文献1には、キャッシュ制御として、割り当てられるキャッシュ容量の制御が行われる。

**【先行技術文献】****【特許文献】****【0006】**

【特許文献1】特開2010-286923号公報

30

**【発明の概要】****【発明が解決しようとする課題】****【0007】**

一般に、スケールアウト型ストレージシステムにおけるストレージ装置間接続は、「疎結合」である。本明細書において、「疎結合」とは、一方のストレージ装置から他方のストレージ装置のキャッシュメモリ領域（CM領域）にデータを入出力できないストレージ装置間接続を意味する。疎結合のストレージ装置では、一方のストレージ装置は、自分のCM領域を管理するが、他方のストレージ装置のCM領域を管理しない。このため、一方のストレージ装置は、自分のCM領域からデータの格納先キャッシュセグメント（CMサブ領域の一例）を確保できるが、他のストレージ装置のCM領域からキャッシュセグメントを確保できない。結果として、疎結合のストレージ装置間でI/Oデータ（I/O要求に従い入出力される対象のデータ）が転送された場合、そのI/Oデータは、転送元のストレージ装置のCM領域にも転送先のストレージ装置のCM領域にも格納されることになる。

40

**【0008】**

そこで、スケールアウト型ストレージシステムにおけるストレージ装置間接続を、「密結合」にすることを検討する。本明細書において、「密結合」とは、一方のストレージ装置から他方のストレージ装置のCM領域にI/Oデータを入出力できるストレージ装置間接続を意味する。密結合のストレージ装置では、一方のストレージ装置は、自分のCM領

50

域だけでなく、他方のストレージ装置のCM領域を管理する。このため、一方のストレージ装置は、自分のCM領域と他のストレージ装置のCM領域のいずれからもキャッシュセグメントを確保できる。結果として、密結合のストレージ装置間でI/Oデータが転送された場合、そのI/Oデータは、転送元のストレージ装置のCM領域と転送先のストレージ装置のCM領域とのいずれか一方にのみ格納されることになる。

【0009】

このように、密結合では、1つのI/Oデータについて、キャッシュセグメントが確保されるCM領域は1つでよい(キャッシュセグメントの二重化が行われることもあり得るが、その場合、メインのキャッシュセグメントが確保されるCM領域は1つでよい)。

【0010】

しかし、そのために、CM領域からホストへのI/Oデータがストレージシステムにおいて転送される回数が多くなり得る。以下、具体例を説明する。なお、課題の具体例の説明では、各ストレージ装置において、ストレージコントローラ(CTL)が二重化されており、各CTLが、CM領域を有する。I/O要求を受信したストレージ装置を、「ホスト側ストレージ装置」と言い、ホスト側ストレージ装置の要素XXを、「ホスト側XX」と言い、I/O要求を受信したホスト側CTLを「I/O受信CTL」と言うことができる。また、I/Oデータを記憶しているPDEV(物理記憶デバイス)を有するストレージ装置を、「PDEV側ストレージ装置」と言い、PDEV側ストレージ装置の要素XXを、「PDEV側XX」と言うことができる。

【0011】

(課題の具体例1(図8))

【0012】

ホスト側CTL5104Aaが、ホスト5102からのリード要求を受信したときに、PDEV側CTL5104Ba内のCM領域5205Baからキャッシュセグメントを確保したとする。この場合、PDEV5110からPDEV側CM領域5205Baにリードデータが格納され、その後2段転送が行われる。すなわち、リードデータが、PDEV側CM領域5205Baからホスト側CTL5104Aa内のバッファ領域5204Aaに転送され(矢印901)、そのホスト側バッファ領域5204Aaからホスト5102にリードデータが転送される(矢印902)。

【0013】

その後、ホスト側CTL5104Aaが同じリードデータのリード要求を受信しても、同様の2段転送が発生する。そのリードデータのCM領域はPDEV側CM領域5205Baであり、且つ、ホスト側バッファ領域5204Aaには、CM領域と異なり、転送されたリードデータが残らないためである。

【0014】

このように、リードデータのCM領域がPDEV側であると、同一データのリード要求を受信する度に、ストレージ装置間データ通信を含んだ2段転送が生じる。このため、性能低下が起こり得る。2段転送を避けるために、リードデータのCM領域をホスト側CM領域にすることが考えられるが、それは容易ではない。いずれのストレージ装置がホスト側ストレージ装置になるかは、リード要求を受信するまでわからないためである。言い換えれば、いずれのストレージ装置もホストからリード要求を受信し得るためである。

【0015】

(課題の具体例2(図9))

【0016】

同一ストレージ装置5101A内のCTL5104Aa及び5104Abのいずれも、ホスト5102からI/O要求を受信し得るとする。言い換えれば、ストレージ装置5101Aがホスト5102にActive/Activeで接続されているとする。また、ストレージ装置5101Aが、少なくとも1つのリード要求の処理としてシーケンシャルリードを実行することを決定した場合、CM領域5205Aa又は5205Abに、シーケンシャルリードによりやがて読み出されることになると予測されるデータを先読みする

10

20

30

40

50

(前もって読み出す) ようになっているとする。

【0017】

一方のCTL5104Aaが、ホスト5102からリード要求を受信して、CM領域5205Aaからキャッシュセグメントを確保し、シーケンシャルリードの実行を決定したとする。この場合、PDEV5110からPDEV側CTL5104Baのバッファ領域5204Baにリードデータが転送され、そのリードデータが、そのPDEV側バッファ領域5204Baからホスト側CM領域5205Aa(確保されたキャッシュセグメント)に転送され、ホスト側CM領域5205Aaからホスト5102へ転送される。その一連の処理に並行して(又はその一連の処理の後に)、シーケンシャルに読み出されると予測されるデータが、PDEV5110からPDEV側バッファ領域5204Baに転送され(矢印1001)、そのデータが、そのPDEV側バッファ領域5204Baからホスト側CM領域5205Aaに転送される(矢印1002)。つまり、ホスト側CM領域5205Aaにデータが先読みされる。

10

【0018】

しかし、その後、先読みされたデータを指定したリード要求を、他方のCTL5104Abが受信したとする。この場合、データが、CM領域5205AaからCTL5104Abのバッファ領域5204Abに転送され(矢印1003)、そのデータが、バッファ領域5204Abからホスト5102へ転送されることになる(矢印1004)。

【0019】

以上の通り、先のリード要求を受信した一方のCTL5104AaのCM領域5205Aaにデータが先読みされた後に、他方のCTL5104Abがその先読みされたデータのリード要求を受信すると、CTL5104Aa及び5104Ab間のデータ転送が発生する。この結果、CTL5104Aa及び5104Ab間のパス5109Aの帯域が消費されることになる。シーケンシャルリードが必要になる分析又はバックアップ等のユースケースでは、CTL間帯域の消費が大きくなりがちであり、CTL間のデータ転送を避けることが望ましい。

20

【0020】

以上の課題は、特許文献1のようなキャッシュ容量制御を利用しても解決することはできない。この種の課題は、背景技術で述べた用途のストレージシステム以外の用途のストレージシステムについてもあり得る。

30

【課題を解決するための手段】

【0021】

論理ボリュームのアドレスを指定したI/O(Input/Output)要求を複数のストレージ装置のうちの一つが受信した場合、対象I/Oパターンと接続形態とのうちの少なくとも一つに基づいて、受信したI/O要求に従うデータであるI/Oデータがキャッシュされるキャッシュ先ストレージ装置が決定される。キャッシュ先ストレージ装置のCM領域にI/Oデータがキャッシュされることになる。対象I/Oパターンは、複数のI/Oパターンのうち、受信したI/O要求に従うI/Oが属するI/Oパターンである。複数のI/Oパターンの各々は、論理ボリュームにおけるI/O先アドレスがランダムとなるかシーケンシャルとなるかに関するパターンである。接続形態は、受信したI/O要求で指定されているアドレスと同一アドレスが指定されたI/O要求を受信するストレージ装置が確定しているか否かである。

40

【発明の効果】

【0022】

対象I/Oパターン及び接続形態のうちの一つに基づいてキャッシュ先ストレージ装置が決定される。このため、CM領域からホストへのI/Oデータがストレージシステムにおいて転送される回数を減らすことができる。

【図面の簡単な説明】

【0023】

【図1】実施例に係るコンピュータシステムの構成例を示す図である。

50

【図2】メモリの論理構成例を示す図である。

【図3】プログラム領域に格納されるプログラムの一例を示す図である。

【図4】共有メモリ領域に格納される情報の一例を示す図である。

【図5】ローカルメモリ領域に格納される情報の一例を示す図である。

【図6】キャッシュセグメント属性管理テーブルの構成例を示す図である。

【図7】キャッシュセグメント管理テーブルの構成例を示す図である。

【図8】比較例の具体例1の概要図である。

【図9】比較例の具体例2の概要図である。

【図10】キャッシュ先CTL決定処理の流れを示すフローチャートである。

【図11】確定ポートを有するCTLがキャッシュ先CTLであることの利点の一例の説明図である。 10

【図12】PDEV側CTLがキャッシュ先CTLであることの利点の一例の説明図である。

【発明を実施するための形態】

【0024】

以下、図面を参照して一実施例を説明する。

【0025】

以下の説明では、「abcテーブル」の表現にて情報を説明することがあるが、情報は、テーブル以外のデータ構成で表現されていてもよい。データ構成に依存しないことを示すために「abcテーブル」のうちの少なくとも1つを「abc情報」と呼ぶことができる。また、以下の説明において、各テーブルの構成は一例であり、1つのテーブルは、2以上のテーブルに分割されてもよいし、2以上のテーブルの全部又は一部が1つのテーブルであってもよい。 20

【0026】

また、以下の説明では、要素のIDは、要素の識別情報（例えば識別子）のことであり、識別情報は、文字、数字、記号又はそれらの組合せ等（例えば名前）で表現可能である。

【0027】

また、以下の説明では、同種の要素を区別しないで説明する場合には、参照符号における共通符号（又は参照符号）を使用し、同種の要素を区別して説明する場合は、参照符号（又は要素のID）を使用することがある。 30

【0028】

また、以下の説明では、「記憶部」は、メモリを含んだ1以上の記憶デバイスでよい。例えば、記憶部は、主記憶デバイス（典型的には揮発性のメモリ）及び補助記憶デバイス（典型的には不揮発性の記憶デバイス）のうちの少なくとも主記憶デバイスでよい。また、記憶部は、CM領域（例えばキャッシュメモリ又はその一部領域）とバッファ領域（例えばバッファメモリ又はその一部領域）とのうちの少なくとも1つを含んでもよい。

【0029】

また、以下の説明では、「PDEV」は、物理的な記憶デバイスを意味し、典型的には、不揮発性の記憶デバイス（例えば補助記憶デバイス）でよい。PDEVは、例えば、HDD（Hard Disk Drive）又はSSD（Solid State Drive）でよい。 40

【0030】

また、以下の説明では、「RAID」は、Redundant Array of Independent (or Inexpensive) Disksの略である。RAIDグループは、複数のPDEVで構成され、そのRAIDグループに関連付けられたRAIDレベルに従いデータを記憶する。RAIDグループは、パリティグループと呼ばれてもよい。パリティグループは、例えば、パリティを格納するRAIDグループのことでよい。

【0031】

また、以下の説明では、「プログラム」を主語として処理を説明する場合があるが、プログラムは、プロセッサ（例えばCPU（Central Processing Unit））によって実行 50

されることで、定められた処理を、適宜に記憶部（例えばメモリ）及び／又はインターフェースデバイス（例えば通信ポート）等を用いながら行うため、処理の主語がプロセッサとされてもよい。プログラムを主語として説明された処理は、プロセッサあるいはそのプロセッサを有する装置又はシステムが行う処理としてもよい。また、プロセッサは、処理の一部または全部を行うハードウェア回路を含んでもよい。プログラムは、プログラムソースから計算機のような装置にインストールされてもよい。プログラムソースは、例えば、プログラム配布サーバまたは計算機が読み取り可能な記憶メディアであってもよい。プログラムソースがプログラム配布サーバの場合、プログラム配布サーバはプロセッサ（例えばCPU）と記憶部を含み、記憶部はさらに配布プログラムと配布対象であるプログラムとを記憶してよい。そして、プログラム配布サーバのプロセッサが配布プログラムを実行することで、プログラム配布サーバのプロセッサは配布対象のプログラムを他の計算機に配布してよい。また、以下の説明において、2以上のプログラムが1つのプログラムとして実現されてもよいし、1つのプログラムが2以上のプログラムとして実現されてもよい。

10

**【0032】**

また、以下の説明では、「ホストシステム」は、ストレージシステムにI/O要求を送信するシステムであり、インターフェースデバイスと、記憶部（例えばメモリ）と、それらに接続されたプロセッサとを有してよい。ホストシステムは、1以上のホスト計算機で構成されてよい。少なくとも1つのホスト計算機は、物理的な計算機でよく、ホストシステムは、物理的なホスト計算機に加えて仮想的なホスト計算機を含んでよい。

20

**【0033】**

以下、I/Oパターン及び接続形態の少なくとも1つに応じてキャッシュ先ストレージ装置を決定するストレージシステムの例を説明する。なお、以下に説明する実施例は、請求の範囲にかかる発明を限定するものではなく、また実施例の中で説明されている特徴の組合せの全てが発明の解決手段に必須であるとは限らない。

**【0034】**

図1は、実施例に係る計算機システムの構成例を示す図である。

**【0035】**

計算機システムは、スケールアウト型ストレージシステムと、ホスト102と、それらを接続する外部ネットワーク103から構成される。スケールアウト型ストレージシステムは、複数のストレージ装置101を有する。ストレージ装置101間接続が密結合である。ホスト102は複数あってもよい。

30

**【0036】**

複数のストレージ装置101の各々を、1つのストレージ装置101を例に取り説明する。

**【0037】**

ストレージ装置101は、複数のPDEV110と、コントローラ部50とを有する。コントローラ部50は、複数のストレージコントローラ（以下、CTL）104を有する。複数のCTL104は、例えば二重化されたCTLである。複数のCTL104の各々に、複数のPDEV110に接続されている。

40

**【0038】**

CTL104は、メモリ106と、FEI/F（フロントエンドインターフェース）107と、BEI/F（バックエンドインターフェース）108と、HCA（Host Channel Adapter）111とを、それらに接続されたプロセッサ（例えば、マイクロプロセッサ）105とを有する。CTL104が有する各種要素の数は1以上でよい。

**【0039】**

また、CTL104はCTL間パス109を介して他のCTL104に接続されている。各CTL104は、他のCTL104のメモリ106に、CTL間パス109を介して、プロセッサ105（又は、図示されていないDMA（Direct Memory Access）回路等）によってアクセス可能である。

50



## 【 0 0 4 0 】

これ以降、説明の簡単のため、1つのストレージ装置101にCTL104は二つであるとする。しかし、CTL104は、三つ以上でもよい。また、プロセッサ105がマルチコアプロセッサの場合、プロセッサ内のいくつかのコアをグループとして、論理的に複数のプロセッサとして管理してもよい。

## 【 0 0 4 1 】

FEI/F107は、例えば、SAN (Storage Area Network) などの外部ネットワーク103を通してホスト102に接続する。外部ネットワーク103は、データ通信ができるプロトコルであれば何でもよい。

## 【 0 0 4 2 】

CTL104は、BEI/F108を介し、PDEV110が接続される。PDEV110は、例えばHDD (Hard Disk Drive)、SSD (Solid State Drive)、又はテープであってよい。

## 【 0 0 4 3 】

CTL104は、論理ボリューム(以下、VOL)をホスト102に提供する。VOLは、1以上のPDEV110に基づく実体的なVOLであってもよいし、1以上のPDEV110に基づくプールから記憶領域が動的に割り当てられる仮想的なVOL(例えばThin Provisioningに従うVOL)であってもよい。複数のPDEV110は、冗長化のためにRAIDグループを構成していてもよく、実体的なVOLは、構成されたRAIDグループに基づくVOLでもよい。VOLは、ホスト102に対して、LU (Logical Unit) として提供可能であり、ホスト102が指定するアドレスに対してライト要求及びリード要求を受け付けることが可能である。なお、各VOLには、I/O処理を担当するプロセッサ105が一つ設定されていてよい。また、ホスト102は、ホストシステムの一例である。

## 【 0 0 4 4 】

ストレージ装置101は、HCA111と一以上のSW (Switch) 112を介し、一以上の他のストレージ装置101と接続される。

## 【 0 0 4 5 】

本実施例では、ストレージ装置101間接続は、密結合である。密結合のストレージ装置では、ストレージ装置101間の通信で使用される通信プロトコルと、ストレージ装置101内のデバイス(要素)間の通信で使用される通信プロトコルは、同じである。どちらの通信プロトコルも、例えば、PCIe (PCI-Express) である。一方、疎結合のストレージ装置では、ストレージ装置間の通信で使用される通信プロトコルと、ストレージ装置内のデバイス間の通信で使用される通信プロトコルは、異なる。前者の通信プロトコルは、例えば、FC (Fibre Channel) 又はIP (Internet Protocol) であり、後者の通信プロトコルは、例えば、PCIe (PCI-Express) である。なお、密結合及び疎結合の定義は、既に説明した通りである。

## 【 0 0 4 6 】

図2は、メモリ106の論理構成例を示す図である。

## 【 0 0 4 7 】

メモリ106には、プログラム領域201と、共有メモリ(SM)領域202と、ローカルメモリ(LM)領域203と、バッファ領域204と、キャッシュメモリ(CM)領域205とが確保されている。

## 【 0 0 4 8 】

プログラム領域201は、プロセッサ105が処理を実施するための各プログラムが格納されている領域である。

## 【 0 0 4 9 】

SM領域202は、全てのプロセッサからアクセスされる領域で、各種管理テーブルが格納される領域である。

## 【 0 0 5 0 】

LM領域203は、プロセッサ105毎に存在する領域であり、特定のプロセッサ105のみからアクセスされる領域である。

【0051】

バッファ領域204及びCM領域205は、FE I/F107、BE I/F108又はCTL間パス109等によるデータ転送の際に、一時的にデータが格納される領域である。なお、CM領域205は、複数のキャッシュセグメント（CMサブ領域の一例）で構成され、キャッシュセグメント単位で領域が確保される。また、バッファ領域204から読み出されてデータは、バッファ領域204には残らない。一方、CM領域205から読み出されたデータは、CM領域205に残る。

【0052】

図3は、プログラム領域201に格納されるプログラムの一例を示す図である。

【0053】

プログラム領域201には、例えば、コマンド振分けプログラム301と、I/Oプログラム302と、新規セグメント割当てプログラム304と、フリーセグメント確保プログラム305と、フリーセグメント解放プログラム306と、CTL間データ転送プログラム309と、ポート判定プログラム310と、I/Oパターン判定プログラム311とが格納される。

【0054】

ホスト102からのリード/ライト要求は、コマンド振分けプログラム301により、リード/ライト対象のVOLのI/O処理を担当するプロセッサ105に振り分けられる。次に、そのリード/ライト要求を受けたプロセッサ105（I/O処理を担当するプロセッサ105）が、I/Oプログラム302を実行することにより、そのリード/ライト要求に従い、対象領域のリード/ライト処理を実行する。その際、必要に応じて、ステージング（PDEV110からCM領域205へのデータの読出し）及びデステージ（CM領域205からPDEV110へのデータの書込み）のうちの少なくとも1つが、I/Oプログラム302により実行される。

【0055】

I/Oプログラム302が、ポート判定プログラム310の判定結果（図10のS1701の判定結果）と、I/Oパターン判定プログラム311の判定結果（図10のS1702の判定結果）とに基づいて、キャッシュ先CTL（新たにキャッシュセグメントが確保されるCM領域205を有するCTL）104を決定する。そして、CM領域205からの新規キャッシュセグメントの割当て処理を新規セグメント割当てプログラム304が実行する。その割当て処理において、フリーセグメント確保プログラム305が、共有フリーキャッシュセグメントを確保することでフリーキャッシュセグメントを確保する。フリーセグメント解放プログラム306は、フリーキャッシュセグメントを解放することで、共有フリーキャッシュセグメントを増やす。

【0056】

CTL104間のデータ転送は、CTL間データ転送プログラム309により実行される。

【0057】

図4は、SM領域202に格納される情報の一例を示す図である。

【0058】

SM領域202には、PDEVテーブル501と、VOLテーブル502と、キャッシュセグメント属性管理テーブル503と、キャッシュセグメント管理テーブル504とが格納される。PDEVテーブル501は、ストレージ装置101とPDEV110との対応関係を示す。VOLテーブル502は、PDEV110とVOLとの対応関係を示す。キャッシュセグメントは、キャッシュセグメント属性管理テーブル503及びキャッシュセグメント管理テーブル504を用いて管理される。SM領域202内の情報は、複数のCTL104のメモリ106にコピーされてもよい。

【0059】

10

20

30

40

50

図5は、LM領域203に格納される情報の一例を示す図である。

【0060】

LM領域203には、キャッシュセグメント属性管理テーブル(ローカル)601と、キャッシュセグメント管理テーブル(ローカル)602とが格納される。

【0061】

図6は、キャッシュセグメント属性管理テーブル503の構成例を示す図である。

【0062】

キャッシュセグメント属性管理テーブル503が有する各エントリ(レコード)は、ストレージ番号、CTL番号、プロセッサ番号、及び、キャッシュセグメント属性を記憶する。

10

【0063】

ストレージ番号は、ストレージ装置101の識別番号ある。CTL番号は、CTL104の識別番号である。プロセッサ番号は、プロセッサ105の識別番号である。キャッシュセグメント属性は、キャッシュセグメントの属性(クリーン、ダーティ、フリー及び共有フリーのいずれか)を示す。「クリーン」とは、PDEV101に格納済みのデータが記憶された確保済キャッシュセグメントを意味する(クリーンキャッシュセグメント内のデータを「クリーンデータ」と呼ぶことができる)。「ダーティ」は、PDEV110に未格納のデータを含んだ確保済キャッシュセグメントを意味する(ダーティキャッシュセグメント内のデータを「ダーティデータ」と呼ぶことができる)。「フリー」は、格納先キャッシュセグメントとして割り当てられる候補である確保済キャッシュセグメント(但しデータが論理的には存在しない)を意味する。「共有フリー」は、未確保のキャッシュセグメント(データが論理的に存在しない)ことを意味する。

20

【0064】

キャッシュセグメント属性管理テーブル(ローカル)601は、キャッシュセグメント属性管理テーブル503と同じ構成であるため説明は省略する。

【0065】

本実施例では、新規セグメント割当てプログラム304が、データをキャッシュ領域に格納する(キャッシュする)ために、I/Oを担当するプロセッサ105に対して、フリーキャッシュセグメントを、データ格納先として割り当てる。そのキャッシュセグメントの属性は、格納されたデータの状態に応じて、フリーからクリーン又はダーティに変化する。

30

【0066】

図7は、キャッシュセグメント管理テーブル504の構成例を示す図である。

【0067】

キャッシュセグメント管理テーブル504は、キャッシュセグメント毎にエントリを有する。各エントリは、キャッシュセグメント番号、VOL番号、VOL内セグメント番号、キャッシュセグメント属性、ダーティビットマップ、CTL番号及びストレージ番号を記憶する。

【0068】

キャッシュセグメント番号は、キャッシュセグメントの識別番号である。VOL番号は、キャッシュセグメントに格納されるデータの格納先VOLの識別番号である。VOLセグメント番号は、キャッシュセグメントに格納されるデータの格納先VOLセグメント(VOLを構成する領域の一例)の識別番号である。VOLは、複数のVOLセグメントで構成され、各VOLセグメントのサイズは、キャッシュセグメントのサイズと同じである。キャッシュセグメント属性は、キャッシュセグメントの属性(クリーン、ダーティ、フリー及び共有フリーのいずれか)を示す。ダーティビットマップは、キャッシュセグメント内のデータを構成する複数の各々がダーティか否かを示す。キャッシュセグメント内のデータが1つでもダーティビットを含んでいれば、そのデータはダーティデータである。CTL番号は、キャッシュセグメントが属するCTLの識別番号である。ストレージ番号は、キャッシュセグメントが属するストレージ装置の識別番号である。

40

50

## 【 0 0 6 9 】

キャッシュセグメント管理テーブル（ローカル）6 0 2 は、キャッシュセグメント管理テーブル5 0 4 と同じ構成であるため説明は省略する。

## 【 0 0 7 0 】

参照は、テーブル（ローカル）6 0 1 及び6 0 2 に対して実行され、更新は、テーブル5 0 3 及び5 0 4 と、テーブル（ローカル）6 0 1 及び6 0 2 の両方に対して実行される。但し、共有フリーキャッシュセグメントに関しては、LM領域2 0 3 へのコピーが存在しないため、SM領域2 0 2 へのアクセスが実行される。

## 【 0 0 7 1 】

さて、密結合のストレージ装置を含んだストレージシステムの比較例には、図8 及び図9 を参照して説明した具体例のような課題がある。

## 【 0 0 7 2 】

そこで、本実施例では、VOLのアドレスを指定したI/O要求（ライト/リード要求）を複数のストレージ装置1 0 1 のうちのいずれかが受信した場合、キャッシュ先CTL決定処理が行われる。キャッシュ先CTL決定処理では、対象I/Oパターンと接続形態とのうちの少なくとも1つに基づいて、受信したI/O要求に従うデータであるI/Oデータがキャッシュされるキャッシュ先CTL1 0 4 が決定される。キャッシュ先CTL1 0 4 のCM領域2 0 5 にI/Oデータ（ライト/リードデータ）がキャッシュされることになる。対象I/Oパターンは、複数のI/Oパターンのうち、受信したI/O要求に従うI/Oが属するI/Oパターンである。複数のI/Oパターンの各々は、VOLにおけるI/O先アドレスがランダムとなるかシーケンシャルとなるかに関するパターンである。接続形態は、受信したI/O要求で指定されているアドレスと同一アドレスが指定されたI/O要求を受信するストレージ装置が確定しているか否かである。

## 【 0 0 7 3 】

図1 0 は、キャッシュ先CTL決定処理の流れを示すフローチャートである。キャッシュ先CTL決定処理は、ストレージシステムにおける複数のCTL1 0 4 のうちの少なくとも1つにより実行される。本実施例では、I/O要求を受信したCTLによりキャッシュ先CTL決定処理が実行されるとする。

## 【 0 0 7 4 】

ポート判定プログラム3 1 0 が、ホスト1 0 2 からI/O要求を受信するポートが確定しているか否か（例えば、今回受信したI/O要求で指定されているアドレスと同一アドレスを指定したI/O要求を受信するポートが確定しているか否か）を判定する（S 1 7 0 1）。「ポート」は、CTLのFE I/F 1 0 7 が有するポートである。1つのCTL1 0 4 の1以上のFE I/F 1 0 7 が、1以上のポートを有する。I/O要求を受信することが確定していると判断されたポートを「確定ポート」と言うことができる。

確定ポートの具体例は、図1 1 に示す通りである。すなわち、ホスト側ストレージ装置（I/O要求を受信したストレージ装置）1 0 1 AのCTL1 0 4 A a及び1 0 4 A bとホスト1 0 2 間の複数のパスのうちALUA（Asymmetric Logical Unit Access）に従う優先度が最も高いパス1 3 0 1 に接続されているポートが確定ポートの一例である。或いは、CTL1 0 4 A a及び1 0 4 A bとホスト1 0 2 間の複数のパスのうち唯一の通信可能状態パス（例えばシングルパス）に接続されているポートが確定ポートの一例である。

## 【 0 0 7 5 】

S 1 7 0 1 の判定結果が肯定の場合（S 1 7 0 1 : Y E S）、新規セグメント割当てプログラム3 0 4 が、テーブル5 0 3 及び5 0 4 を参照することで、図1 1 に示すように、確定ポートを有するCTL1 0 4 A a内のCM領域2 0 5 A aからフリーキャッシュセグメントを選択し、選択したフリーキャッシュセグメントを、当該I/O要求に従うI/O先VOLセグメントに割り当てる。これにより、そのキャッシュセグメントの属性が、フリーからダーティ又はクリーンに更新され、且つ、割当て先VOLセグメントの識別番号が、そのキャッシュセグメントに関連付けられる（図7のキャッシュセグメントテーブルが更新される）。以下、キャッシュセグメントの割当てに関しては、同様の処理が行われ

10

20

30

40

50

るため、キャッシュセグメントの割当ての説明を簡略化する。S 1 7 0 1 : Y E Sの後のS 1 7 0 3によれば、確定ポートを有するC T L 1 0 4 A a (すなわち、必ずI / O要求を受信することになるC T L 1 0 4 A a)のC M領域2 0 5 A aがキャッシュ先となるので、図8を参照して説明したような無駄なデータ転送の発生を回避できる。例えば、今回受信したリード要求で指定されているアドレスと同一アドレスを指定したリード要求が受信された場合、そのリード要求を受信したC T L 1 0 4 A a内のC M領域2 0 5 A aからホスト1 0 2へデータを転送でき、図8を参照して説明した2段転送は不要である。なお、キャッシュセグメントの二重化が行われる場合は、キャッシュ先とされたC T Lと同じストレージ装置内の他のC T Lからキャッシュセグメントが割り当てられてもよいし、S W 1 1 2を介して接続された他のストレージ装置内のC T Lからキャッシュセグメントが

10

#### 【 0 0 7 6 】

S 1 7 0 1の判定結果が否定の場合(S 1 7 0 1 : N O)、I / Oパターン判定プログラム3 1 1が、対象I / Oパターンを判定する。I / Oパターン判定よりもポート判定が優先される。この結果、I / Oパターン判定に関わらず、確定ポートがあれば、確定ポートを有するC T L 1 0 4がキャッシュ先となる。これにより、上述した無駄なデータ転送の発生を回避できる。

#### 【 0 0 7 7 】

S 1 7 0 2の判定結果、対象I / Oパターン(複数のI / Oパターンのうち、受信したI / O要求に従うI / Oが属するI / Oパターン)がランダムリード及びランダムライトのいずれかの場合、新規セグメント割当てプログラム3 0 4が、図11に示すように、ホスト側ストレージ装置1 0 1 AのいずれかのC T L (例えば1 0 4 A a)のC M領域(2 0 5 A a)が有するフリーキャッシュセグメントを、当該I / O要求に従うI / O先V O Lセグメントに割り当てる(S 1 7 0 3)。つまり、ここでは、S 1 7 0 1 : Y E Sの場合と異なり、I / O要求を受信したストレージ装置1 0 1から任意のC T L 1 0 1が(例えばランダム又はラウンドロビンで)選択されてよい。ホスト側ストレージ装置1 0 1がキャッシュ先とされる理由は、そのストレージ装置1 0 1が提供したV O Lを指定したI / O要求はそのストレージ装置1 0 1が受信するためである。I / O要求を受信したストレージ装置1 0 1のうちの任意のC T L 1 0 1がキャッシュ先とされる理由は、いずれのC T LがI / O要求を受信するかわからないためである。

20

30

#### 【 0 0 7 8 】

S 1 7 0 2の判定結果、対象I / Oパターンがシーケンシャルリードの場合、新規セグメント割当てプログラム3 0 4が、図12に示すように、当該リード要求に従うリードデータが格納されているP D E V (指定されたV O Lの基になっているP D E V) 1 1 0を有するストレージ装置(P D E V側ストレージ装置) 1 0 1 BのうちのいずれかのC T L (例えば1 0 4 B a)内のC M領域(2 0 5 B a)が有するフリーキャッシュセグメントを、当該リード要求に従うリード元V O Lセグメントに割り当てる(S 1 7 0 4)。S 1 7 0 4では、そのリード元V O Lセグメントに連続する1以上のV O Lセグメントにも、それぞれ、P D E V側ストレージ装置1 0 1内の同じC T L 1 0 4 (又は別のC T L 1 0 4)のC M領域2 0 5から1以上のキャッシュセグメントを割り当てる。リードデータの他にシーケンシャルに読み出され得るデータが、P D E V側ストレージ装置1 0 1内のC M領域2 0 5 (割り当てられたキャッシュセグメント)に、例えばI / Oプログラム3 0 2により先読みされる。S 1 7 0 4によれば、シーケンシャルリードのようにデータが先読みされるI / Oパターンでは、P D E V側C T L 1 0 4 B aがキャッシュ先となるので、図9を参照して説明したような無駄なC T L間データ転送の発生を回避できる。すなわち、図12に示すように、先読みされたデータのうちのいずれかをリード対象としたリード要求をホスト側C T L 1 0 4 A a及び1 0 4 A bのいずれが受信しても、リード対象のデータを、P D E V側C M領域2 0 5 B aから、S W 1 1 2を介して、そのリード要求を受信したホスト側C T L (例えば1 0 4 A b)のバッファ領域(2 0 4 A b)に転送できるからである。このため、ホスト側ストレージ装置1 0 1 AでのC T L間データ転送が不

40

50

要である。

【0079】

S1702の判定結果、対象I/Oパターンが、ランダムリード、ランダムライト及びシーケンシャルリードのいずれにも該当しない場合、新規セグメント割当てプログラム304が、任意のCTL104のCM領域が有するフリーキャッシュセグメントを、当該I/O要求に従うI/O先VOLセグメントに割り当てる(S1705)。任意のCTL104は、ランダムに選択されてもよいし、ラウンドロビンで選択されてもよい。

【0080】

本実施例によれば、I/Oパターン及び接続形態の少なくとも1つに基づきキャッシュ先CTLが決定される。これにより、無駄なデータ転送を回避でき、以って、高性能化を図ることができる。

10

【0081】

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

【0082】

また、上記の各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。また、上記の各構成、機能等は、プロセッサがそれぞれの機能を実現するプログラムを解釈し、実行することによりソフトウェアで実現してもよい。各機能を実現するプログラム、テーブル、ファイル等の情報は、メモリや、ハードディスク、SSD等の記憶デバイス、または、ICカード、SDカード等の記憶デバイスに置くことができる。

20

【0083】

また、制御線や情報線は説明上必要と考えられるものを示しており、製品上必ずしも全ての制御線や情報線を示しているとは限らない。実際には殆ど全ての構成が相互に接続されていると考えてもよい。

【符号の説明】

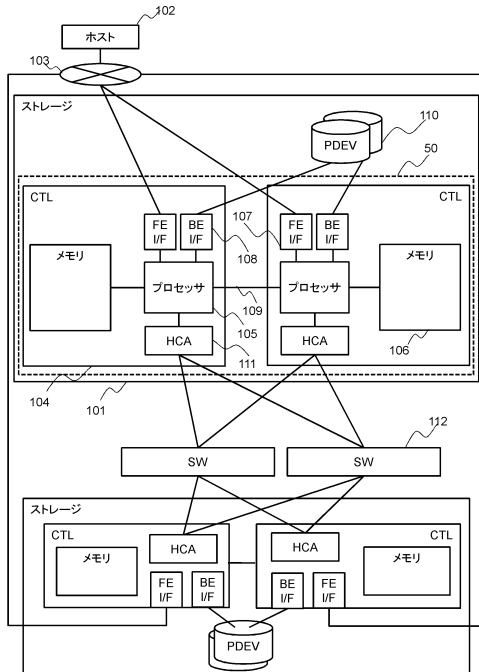
30

【0084】

101...ストレージ装置

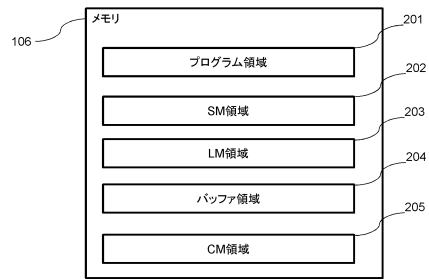
【 図 1 】

FIG. 1



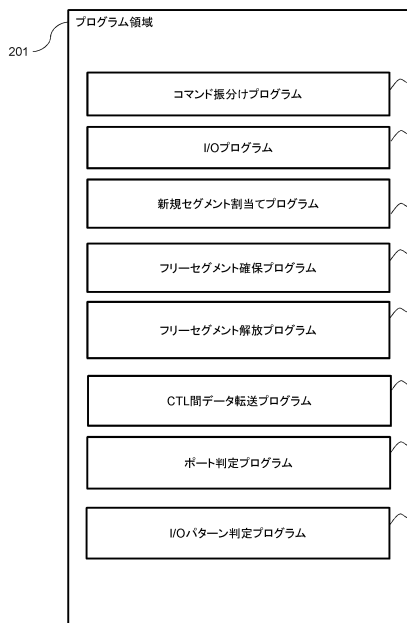
【 図 2 】

FIG. 2



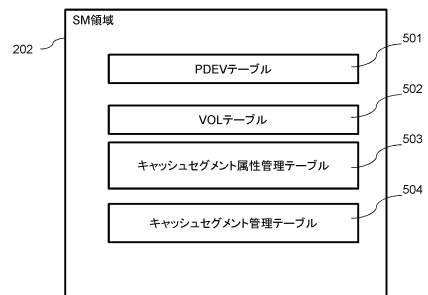
【 図 3 】

FIG. 3



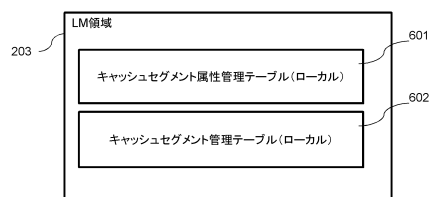
【 図 4 】

FIG. 4



【 図 5 】

FIG. 5



【 図 6 】

FIG. 6

キャッシュセグメント属性管理テーブル  
503

ストレージ番号	CTL番号	プロセス番号	キャッシュセグメント属性
1	1	1	クリーン
1	1	1	ダーティ
1	1	1	フリー
1	2	1	クリーン
1	2	1	ダーティ
2	2	NULL	共有フリー

【 図 7 】

FIG. 7

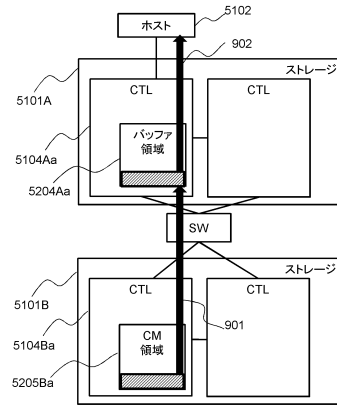
キャッシュセグメント管理テーブル  
504

キャッシュセグメント番号	VOL番号	VOL内セグメント番号	キャッシュセグメント属性	ダーティビットマップ	CTL番号	ストレージ番号
1	なし	なし	フリー	00000000	1	1
2	2	1	クリーン	00000000	1	2
3	1	1	クリーン	00000000	2	1
4	3	4	ダーティ	00100010	2	2
5	4	5	ダーティ	10010011	1	1

【 図 8 】

FIG. 8

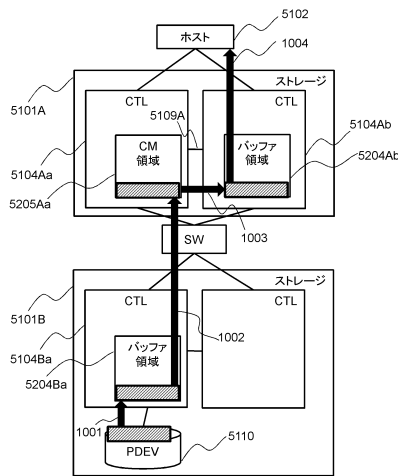
比較例



【 図 9 】

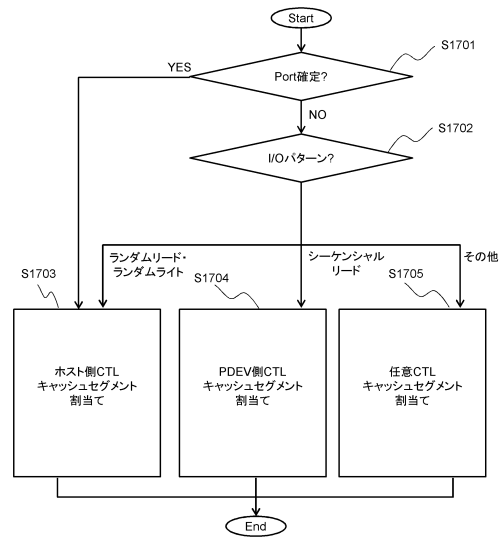
FIG. 9

比較例



【 図 10 】

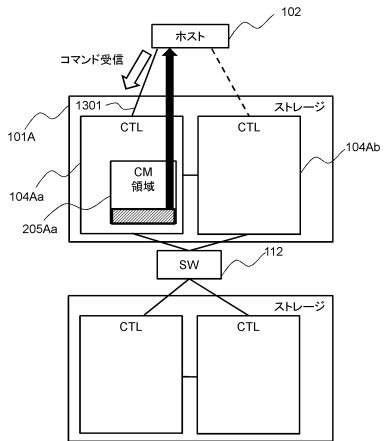
FIG. 10





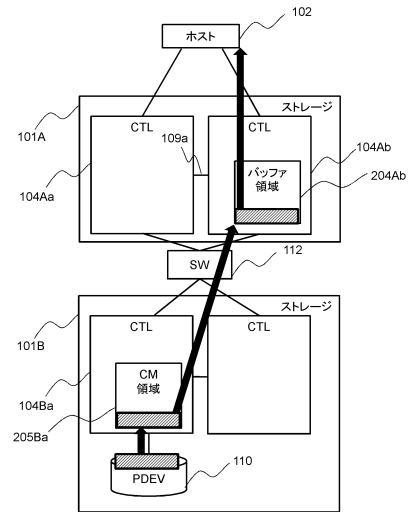
【 図 1 1 】

FIG. 11



【 図 1 2 】

FIG. 12



---

フロントページの続き

審査官 桜井 茂行

- (56)参考文献 国際公開第2015/056301(WO, A1)  
特開2005-275525(JP, A)  
米国特許出願公開第2005/0216692(US, A1)  
特開2006-252019(JP, A)  
米国特許出願公開第2006/0215682(US, A1)  
特表2007-537548(JP, A)  
国際公開第2005/114428(WO, A1)  
特開2000-148587(JP, A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 13/10  
G06F 3/06