

(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2020-140325
(P2020-140325A)

(43) 公開日 令和2年9月3日(2020.9.3)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 302A	5B205
G06F 13/10 (2006.01)	G06F 3/06 302Z	
G06F 12/0866 (2016.01)	G06F 13/10 340A	
G06F 12/0868 (2016.01)	G06F 12/0866 100	
	G06F 12/0868 105	

審査請求 有 請求項の数 11 O L (全 33 頁)

(21) 出願番号 特願2019-33920 (P2019-33920)
(22) 出願日 平成31年2月27日 (2019.2.27)

(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 110001678
特許業務法人藤央特許事務所
(72) 発明者 鶴谷 昌弘
東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(72) 発明者 吉原 朋宏
東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(72) 発明者 達見 良介
東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内

最終頁に続く

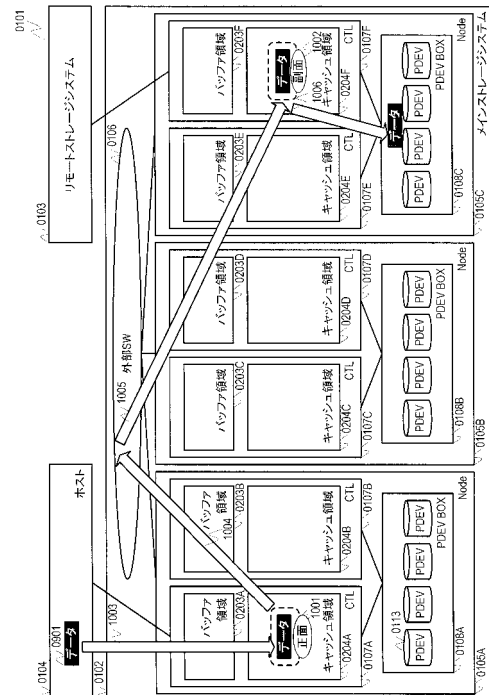
(54) 【発明の名称】 複数のストレージノードを含むストレージシステム

(57) 【要約】

【課題】ストレージシステムの性能を向上する。

【解決手段】ストレージシステムは、ネットワークを介して通信する複数のストレージノードを含む。複数のストレージノードのそれぞれは、1以上のコントローラを含む。コントローラにおける少なくとも1つのコントローラは、ホストからライトデータを受信するコントローラと、ライトデータを処理するコントローラとに基づいて、ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定し、指定したコントローラにおいて、キャッシュサブ領域を確保する。

【選択図】 図10



【特許請求の範囲】**【請求項 1】**

ネットワークを介して通信する複数のストレージノードを含むストレージシステムであって、

前記複数のストレージノードのそれぞれは、1以上のコントローラを含み、

前記コントローラにおける少なくとも1つのコントローラは、

ホストからライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとに基づいて、前記ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定し、

指定したコントローラにおいて、キャッシュサブ領域を確保する

ストレージシステム。

10

【請求項 2】

請求項 1 に記載のストレージシステムであって、

前記少なくとも1つのコントローラは、さらに、各コントローラが属するストレージノードの、前記ライトデータを格納する記憶デバイスへの接続の有無に基づいて、前記二つのコントローラを指定する

ストレージシステム。

【請求項 3】

請求項 1 に記載のストレージシステムであって、

前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとは、別のコントローラであり、

前記少なくとも1つのコントローラは、前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定する

ストレージシステム。

20

【請求項 4】

請求項 3 に記載のストレージシステムであって、

前記指定した二つのコントローラのうち、前記ライトデータを格納する記憶デバイスに接続されたコントローラは、その前記キャッシュサブ領域に記憶された前記ライトデータを前記記憶デバイスに格納する

ストレージシステム。

30

【請求項 5】

請求項 4 に記載のストレージシステムであって、

前記指定した二つのコントローラのいずれもが、前記記憶デバイスに接続されていない場合、二つのコントローラのうちいずれかが、前記記憶デバイスに接続された他のコントローラのバッファ領域に前記ライトデータを転送し、

前記他のコントローラが、前記バッファ領域に格納された前記ライトデータを、前記記憶デバイスに格納する

ストレージシステム。

40

【請求項 6】

請求項 2 に記載のストレージシステムであって、

前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとは、同じコントローラであり、

前記少なくとも1つのコントローラは、前記同じコントローラと、他の1つのコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定し、

前記指定する二つのコントローラのうち少なくとも1つは、前記記憶デバイスに接続されている

ストレージシステム。

【請求項 7】

請求項 2 に記載のストレージシステムにおいて、

50

前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを格納する記憶デバイスへの接続された他のコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定する

ストレージシステム。

【請求項 8】

請求項 2 に記載のストレージシステムであって、

前記少なくとも 1 つのコントローラは、前記ホストから前記ライトデータを受信するコントローラと、前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを格納する記憶デバイスへの接続された他のコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定する

ストレージシステム。

【請求項 9】

請求項 2 に記載のストレージシステムであって、

前記記憶デバイスは、前記ストレージノードの内に設けられた記憶媒体、またはリモートストレージシステムである

ストレージシステム。

【請求項 10】

請求項 1 に記載のストレージシステムであって、

前記少なくとも 1 つのコントローラは、

前記ホストからのライト要求のアクセスパターンがシーケンシャルで場合に、前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを処理する前記コントローラとに基づいて、前記コントローラの指定を行い、

前記ホストからのライト要求のアクセスパターンがランダムで場合に、キャッシュ利用率に基づいて、前記コントローラの指定を行う

ストレージシステム。

【請求項 11】

ネットワークを介して通信する複数のストレージノードを含むストレージシステムの制御方法であって、

前記複数のストレージノードのそれぞれは、1 以上のコントローラを含み、

前記制御方法は、

ホストからライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとに基づいて、前記ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定し、

指定したコントローラにおいて、キャッシュサブ領域を確保する

制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数のストレージノードを含むストレージシステムに関する。

【背景技術】

【0002】

近年、IoT (Internet of Things) や AI (Artificial Intelligence) に代表されるような膨大なデータを蓄積、及び分析することにより、新たな価値を生み出す技術の重要性が増している。これらの技術では、膨大なデータを蓄積可能な容量だけでなく、蓄積したデータを分析するために高性能な I/O (Input/Output) 性能を有するストレージシステムが必要となる。

【0003】

一方、ストレージシステムの導入にあたっては、導入コストを抑えるため、導入初期は小規模な構成のシステムを導入し、事業規模の拡大に合わせてシステムを拡張していくことが望ましい。これを実現する方式の一つとして、スケールアウト型ストレージシステム

10

20

30

40

50

がある。スケールアウト型ストレージシステムでは、事業規模が拡大し、必要となる I / O 性能が増加してきたら、ストレージ装置 (Node) の台数を増やすことで、ストレージシステムの I / O 性能を向上することが可能である。

【 0 0 0 4 】

スケールアウト型ストレージシステムの I / O 性能を向上させる 1 つの手段として、スケールアウト型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 では、キャッシュ制御として、キャッシュを割り当てる Node を制御することで、リード性能を向上させる方式が述べられている。

【 先行技術文献 】

【 特許文献 】

【 0 0 0 5 】

【 特許文献 1 】 特願 2 0 1 7 - 5 2 4 2 4 8 号 公 報

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

一般に、スケールアウト型ストレージシステムにおける Node 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 Node から他方の Node のキャッシュメモリ領域 (以下単に、キャッシュ領域) にデータを入出力できない Node 間接続を意味する。疎結合のストレージシステムでは、一方の Node は、自分のキャッシュ領域を管理するが、他方の Node のキャッシュ領域を管理しない。

【 0 0 0 7 】

このため、一方の Node は、自分のキャッシュ領域からデータの格納先キャッシュセグメント (キャッシュサブ領域) を確保できるが、他の Node のキャッシュ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 Node 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 Node のキャッシュ領域にも転送先の Node のキャッシュ領域にも格納されることになる。

【 0 0 0 8 】

そこで、スケールアウト型ストレージシステムにおける Node 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 Node から他方の Node のキャッシュ領域に I / O データを入出力できる Node 間接続を意味する。密結合のストレージシステムでは、一方の Node は、自分のキャッシュ領域だけでなく、他方の Node のキャッシュ領域を管理する。このため、一方の Node は、自分のキャッシュ領域と他の Node のキャッシュ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 Node 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 Node のキャッシュ領域と転送先の Node のキャッシュ領域とのいずれか一方にのみ格納されることになる。

【 0 0 0 9 】

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キャッシュ領域は 1 つでよい。なお、ライトにおいて、キャッシュセグメントの二重化 (冗長化) が行われることがある。その場合、メインのキャッシュセグメント (正面) が確保されるキャッシュ領域は 1 つでよく、サブのキャッシュセグメント (副面) が確保されるキャッシュ領域は 1 つ以上でよい。以降、明細書においては、「密結合」のスケールアウト型ストレージシステムを前提とする。

【 0 0 1 0 】

スケールアウト型ストレージシステムは、前述のとおり、性能を拡張性可能な点で、従来の非スケールアウト型ストレージに対する利点がある一方、Node 間で I / O データを複製する回数が多くなり、Node 間を接続する経路の帯域がボトルネックとなり、ストレージシステムの I / O 性能が下がる場合がある。そこで、特許文献 1 では、Node

10

20

30

40

50

間での I/O データ複製回数が小さくなるよう、I/O パターンと接続形態の一方又は両方を用いて、キャッシュを確保する Node (Node に複数のストレージコントローラ (CTL) が含まれる場合、CTL) を決定する方式が述べられており、これによりストレージシステムのリード性能を向上させることができる。

【0011】

一方、ライト性能については、特許文献 1 では改善方式が述べられていない。一般的にライトバック方式を採用するストレージシステムでは、キャッシュメモリに I/O データを格納した時点で I/O 要求元のホストにライト完了を通知する。このため、単一の Node 障害で I/O データを消失しないよう、I/O データを最終記憶媒体に格納するまでの間、2 つ以上の Node で I/O データを複製してキャッシュする方式が用いられる。このため、リードと比較し、ライトの場合は Node 間のデータ複製回数が多くなり、Node 間のデータバス帯域がボトルネックとなった場合、従来の非スケールアウト型ストレージよりも性能が低下する。

10

【0012】

一方、ストレージシステムでは、I/O のレスポンス性能を改善するため、高頻度でアクセスされるデータを可能な限りキャッシュに保持しておくことが望ましい。このため、単純に CTL 間のデータ複製が少なくなるようにキャッシュを確保しても、ホストからのアクセスパターンによっては、特定の Node でキャッシュヒット率の高いデータがキャッシュに保持できず、キャッシュヒット率が低下することで、レスポンス性能の悪化を招いてしまう。よって、レスポンス性能の低下を招くことなく、ライトスループット性能を向上させることが課題である。

20

【課題を解決するための手段】

【0013】

本発明の一態様は、ネットワークを介して通信する複数のストレージノードを含むストレージシステムであって、前記複数のストレージノードのそれぞれは、1 以上のコントローラを含み、前記コントローラにおける少なくとも 1 つのコントローラは、ホストからライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとに基づいて、前記ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定し、指定したコントローラにおいて、キャッシュサブ領域を確保する。

30

【発明の効果】

【0014】

本発明の一態様により、ストレージシステムの性能を向上させる。

【図面の簡単な説明】

【0015】

【図 1】メインストレージシステムを含む情報システムの構成例を示す図である。

【図 2】主記憶の構成例を示す図である。

【図 3】プログラム領域の構成例を示す図である。

【図 4】管理情報領域の構成例を示す図である。

【図 5】PDEV 管理テーブルの構成例を示す図である。

【図 6】LDEV 管理テーブルの構成例を示す図である。

40

【図 7】バッファセグメント管理テーブルの構成例を示す図である。

【図 8】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 9】メインストレージシステムにおいて、従来方式におけるシーケンシャルライト時のデータ転送の一例を示す図である。

【図 10】メインストレージシステムにおいて、本発明に基づく方式におけるシーケンシャルライト時のデータ転送の一例を示す図である。

【図 11】メインストレージシステムにおけるホストらの I/O 要求に基づく処理流れの一例を示すシーケンス図である。

【図 12】メインストレージシステムにおけるライト要求を処理する過程でセグメントを確保するストレージコントローラを決定する処理の一例を示すシーケンス図である。

50

【図13】メインストレージシステムにおける特定のストレージコントローラを指定してセグメントを確保するストレージコントローラを決定する処理の一例を示すシーケンス図である。

【図14】メインストレージシステムにおけるストレージコントローラ間のデータ複写回数を最小化するようにセグメントを確保するストレージコントローラを決定する処理の一例を示すシーケンス図である。

【図15】メインストレージシステムにおけるセグメントを確保する処理の一例を示すシーケンス図である。

【図16】メインストレージシステムにおけるライト処理の一例を示すシーケンス図である。

【図17】メインストレージシステムにおけるセグメントを解放する処理の一例を示すシーケンス図である。

【発明を実施するための形態】

【0016】

以下、図面を用いて実施例を説明する。なお、実施例は本発明を実現するための一例に過ぎず、本発明の技術的範囲を限定するものではない。以下の説明では、「xxxテーブル」の表現にて各種情報を説明することがあるが、各種情報は、テーブル以外のデータ構造で表現されていてもよい。データ構造に依存しないことを示すために「xxxテーブル」を「xxx情報」と呼ぶことができる。

【0017】

以下の説明では、要素の識別情報として番号が使用されるが、他種の識別情報（例えば名前、識別子）が使用されて良い。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。

【0018】

以下の説明では、「主記憶」は、メモリを含んだ1以上の記憶デバイスでよい。例えば、主記憶は、主記憶デバイス（典型的に揮発性の記憶デバイス）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、キャッシュ領域（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。

【0019】

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。以下の説明では、「RAID」は、Redundant Array of Independent（or Inexpensive）Disksの略である。

【0020】

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グループのことである。以下の説明では、「LDEV」は、論理的な記憶デバイスを意味し、RAIDグループ、またはPDEVから構成され、ホストは「LDEV」に対してI/O要求を実行する。

【0021】

以下、「プログラム」を主語として処理を説明することがある場合、プログラムは、ストレージコントローラに含まれるプロセサ（例えばCPU（Central Processing Unit））によって実行されることで、定められた処理を、適宜に記憶資源（例えば主記憶）及び/又は通信インタフェース装置（例えばHCA）を用いながら行うため、処理の主語がストレージコントローラ或いはプロセサとされてもよい。また、ストレージコントローラは、処理の一部又は全部を行うハードウェア回路を含んでもよい。

10

20

30

40

50

コンピュータプログラムは、プログラムソースから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、又は、コンピュータ読取可能な記憶メディアであってもよい。

【0022】

以下の説明では、「ホスト」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。

【0023】

以下、I/Oパターン及び接続形態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。一般的に、ストレージシステムにおいて、ランダムアクセス時のI/O処理は、プロセッサ処理がボトルネックとなり、シーケンシャルアクセス時のI/O処理は、データバス帯域（Node間を含む）がボトルネックとなる。また、シーケンシャルアクセスされるデータは、再度アクセスされる可能性が低い。

【0024】

そこで、以下に説明する例において、ストレージシステムは、Node間のデータバスがボトルネックとなるシーケンシャルライト時、Node間のデータ複写回数が少なくなるようにキャッシュ領域を確保するNodeを決定する。さらに、一例において、ストレージシステムは、最終記憶媒体にI/Oデータを格納した後、すぐにキャッシュ領域を解放する。以下に説明する例において、ストレージシステムは、ランダムライト時、Node間でキャッシュ利用効率が高くなるようにキャッシュ領域を確保するNodeを決定する。

【0025】

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。たとえば、ストレージシステムは、ライト要求のアクセスパターン（シーケンシャルライト/ラインダムライト）と独立に、以下に説明するようにNode間のデータ複写回数が少なくなるようにキャッシュ領域を確保するNodeを決定してもよい。

【0026】

図1は、情報システム0101の構成例を示す図である。情報システム0101は、1つ以上のメインストレージシステム0102、1つ以上のリモートストレージシステム0103、1つ以上のホスト0104から構成されている。リモートストレージシステム0103は、含んでいなくてもよい。また、図1では、メインストレージシステム0102とリモートストレージシステム0103、及びメインストレージシステム0102とホスト0104は直接接続されているが、SAN（SAN: Storage Area Network）、LAN（LAN: Local Area Network）、及びWAN（WAN: World Area Network）などのネットワークを介して接続されていてもよい。

【0027】

メインストレージシステム0102は、1つ以上のNode0105（ストレージノードとも呼ぶ）、及び1つ以上の外部SW（SW: Switch）0106から構成されている。図1では、Node0105は、外部SW0106を介して接続されているが、Node0105間で直接接続されていてもよい。

【0028】

複数のNode0105の各々を、1つのNode0105を例に取り説明する。Node0105は、1つ以上のストレージコントローラ（CTL）0107、及び1つ以上のPDEV BOX0108を含む。また、PDEV BOX0108は、1つ以上のPDEV0113を有し、CTL0107とPDEV0113の接続を仲介している。

【0029】

10

20

30

40

50

図1では、CTL0107は、直接PDEV BOX0108に接続されているが、Switchなどを介して接続されていてもよい。また、PDEV BOX0108を介さず、CTL0107とPDEV0113は直接接続されていてもよい。さらに、同一Node0105内のCTL0107間は、外部SW0106を介して接続されているが、直接接続されていてもよい。

【0030】

複数のCTL0107の各々を、1つのCTL0107を例に取り説明する。CTL0107は、プロセサ0106、主記憶0109、FEI/F(FrontEnd Interface)0110、BEI/F(BackEnd Interface)0111、及びHCA(Host Channel Adapter)0112を含む。CTL0107を構成する各種要素の数は、1以上でよい。

10

【0031】

プロセサ0106は、CTL0107全体を制御し、主記憶0109に格納されたマイクロプログラムに基づき動作する。FEI/F0110は、プロセサ0106により制御され、ホスト0104、及びリモートストレージシステム0103とI/O要求やI/Oデータの送受信などを実施する。BEI/F0111は、プロセサ0106により制御され、PDEV BOX0106を介し、PDEV0113とI/Oデータの送受信などを実施する。HCA0111は、プロセサ0106により制御され、外部SW0106を介して他のCTL0107と制御情報やI/Oデータの送受信などを実施する。

20

【0032】

本実施例では、Node0105間接続は、密結合である。密結合のストレージシステムでは、Node0105間の通信で使用される通信プロトコルと、CTL0107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。

【0033】

一方、疎結合のストレージ装置では、Node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。

30

【0034】

図2は、主記憶0109の構成例を示す図である。主記憶0109には、プログラム領域0201、管理情報領域0202、バッファ領域0203、及びキャッシュ領域0204が確保されている。プログラム領域0201は、同一CTL0107のプロセサ0106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。管理情報領域0202は、メインストレージシステム0102内の全てのプロセ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。

【0035】

バッファ領域0203、及びキャッシュ領域0204は、同一CTL0107のFEI/F0110、BEI/F0111、及びHCA0112等によるデータ転送の再、一時的にデータが格納される領域である。なお、バッファ領域0203とキャッシュ領域0204は、複数のセグメント(キャッシュ領域を区切った単位であり、キャッシュサブ領域とも呼ぶ)で構成され、セグメント単位で領域が確保される。また、バッファ領域0203から読みだされたデータは、バッファ領域0203には残らない。一方、キャッシュ領域0204から読み出されたデータは、キャッシュ領域0204に残る。

40

【0036】

図3は、プログラム領域0201に格納されるプログラムの一例を示す図である。プログラム領域0201には、例えば、I/O要求のCTL間振り分けプログラム0301、I/Oプログラム0302、I/Oパターン判定プログラム0303、セグメント確保プログラム0304、セグメント解放プログラム0305、フリーセグメント確保プログラム

50

0306、及びストレージコントローラ間データ転送プログラム0307が格納される。

【0037】

ホスト0104からのI/O要求を受けると、I/O要求のCTL間振り分けプログラム0301は、当該I/O要求を担当するCTL0107を決定し、振り分ける。振り分けは、あらかじめLDEVごとにI/O要求を処理するCTL0107を決めておいてもよいし、動的に決定してもよい。次に、I/O要求を割り振られたCTL0107のプロセッサ0106は、I/Oプログラム0302を実行することにより、I/O要求に従い、対応する処理を実行する。

【0038】

I/Oパターン判定プログラム0303は、例えばI/Oプログラム0302がI/O要求を処理する過程で呼び出され、I/O要求のアクセスパターン（アクセスパターン又はI/Oパターン）が、シーケンシャルであるか、ランダムであるかを判定する。セグメント確保プログラム0304は、例えばI/Oプログラム0302がI/O要求を処理する過程で呼び出され、バッファセグメント、及びキャッシュセグメントを確保する。セグメント解放プログラム0305は、例えばI/Oプログラム0302がI/O要求を処理する過程で呼び出され、バッファセグメント、及びキャッシュセグメントを解放する。

【0039】

フリーセグメント確保プログラム0306は、I/O要求とは非同期的に実行され、バッファセグメント、及びキャッシュセグメントのうち、確保可能な状態のセグメント（フリーセグメント）の量を一定以上に保つ。ストレージコントローラ間データ転送プログラム0307は、例えばI/Oプログラム0302がI/O要求を処理する過程などで呼び出され、CTL0107間でのデータ転送を実行する。

【0040】

図4は、管理情報領域0202に格納される情報の一例を示す図である。管理情報領域0202には、PDEV管理テーブル0401、LDEV管理テーブル0402、バッファセグメント管理テーブル0403、及びキャッシュセグメント管理テーブル0404が格納される。

【0041】

PDEV管理テーブル0401は、メインストレージシステム102内の全てのPDEV0114の状態や、CTL0107との対応関係を示す。LDEV管理テーブル0402は、メインストレージシステム102内の全てのPDEV0114とLDEVの対応関係を示す。バッファセグメント管理テーブル0403は、メインストレージシステム102内の全てのバッファ領域0203を管理するのに用いられる。キャッシュセグメント管理テーブル0404は、メインストレージシステム102内の全てのキャッシュ領域0204を管理するのに用いられる。

【0042】

図5は、PDEV管理テーブル0401の構成例を示す図である。PDEV管理テーブル0401は、PDEV#0501、容量0502、状態0503、及び接続CTL#0504のエントリを有する。PDEV#0501は、PDEV0114の識別子である。容量0502は、PDEV0114のデータを格納可能な容量を示す。状態0503は、PDEV0114が正常に動作中か否か（故障していないか）の状態を示す。接続CTL#0504は、PDEV0114に接続され、アクセスするCTL0107を示す。

【0043】

図6は、LDEV管理テーブル0402の構成例を示す図である。LDEV管理テーブル0402は、LDEV#0601、容量0602、状態0603、冗長構成0604、及び構成情報0605のエントリを有する。

【0044】

LDEV#0601は、LDEVの識別子である。容量0602は、LDEVにデータを格納可能な容量を示す。状態0503は、LDEVへ正常にI/O可能か否かを示す。冗長構成0604は、RAIDによる冗長化、またはリモートストレージシステム010

10

20

30

40

50

3とのストレージシステム冗長化の状態を示す。構成情報0605は、LDEVに属するPDEV0114、及びストレージシステム冗長化の対応リモートストレージシステム0103を示す。

【0045】

図7は、バッファセグメント管理テーブル0403の構成例を示す図である。バッファセグメント管理テーブル0403は、セグメント#0701、CTL#0702、状態0703を有する。セグメント#0701は、バッファセグメントの識別子である。CTL#0702は、CTL0107の識別子であり、当該バッファセグメントの利用権利を有したCTL0107を示している。状態0703は、バッファセグメントが確保されているか否か(ロック、フリー)を示す。

【0046】

図8は、キャッシュセグメント管理テーブル0404の構成例を示す図である。キャッシュセグメント管理テーブル(0404)は、セグメント#0801、LDEV#0802、LDEV内セグメント#0803、CTL#0804、状態0805、及び属性0806を有する。

【0047】

セグメント#0801は、キャッシュセグメントの識別子である。LDEV#は、LDEVの識別子である。LDEV内セグメント#0803は、当該キャッシュセグメントがLDEVの論理アドレス空間上のどこに割り当てられているかを一意に確定する識別子である。CTL#0804は、CTL0107の識別子であり、当該バッファセグメントの利用権利を有したCTL0107を示している。

【0048】

状態0805は、キャッシュセグメントの状態を表し、「フリー」は未使用、「ロック」は確保済みだが未使用、「クリーン」確保済みでPDEVにデータ格納済み、及び「ダーティー」は確保済みでPDEVにデータ未格納であることを示す。属性0806は、当該キャッシュセグメントが正面、または副面のどちらで確保されているかを示す。

【0049】

図9は、メインストレージシステム0102において、従来方式を用いた、シーケンシャルライトにおけるライトシーケンスの一例を示したものである。本図は、ライトシーケンスを説明するため、情報システム0101の構成要素を一部省略して記載している。また、メインストレージシステム0102を構成する要素のうち、複数存在し、かつ説明上識別する必要のあるものについて、4桁の通番に加え、アルファベット一文字を付加し、区別している。

【0050】

ライト処理で必要となる、キャッシュセグメントの正面0902、及び副面0903は、それぞれCTL0107Dのキャッシュ領域0204D、及びCTL0107Bのキャッシュ領域0204Bに確保されている。データ0901の流れを示す矢印0904、0905、0906、0907、0908、0909、0910、及び0911より、データ0901のCTL間データ複写の回数は、3回である。

【0051】

図10は、メインストレージシステム0102において、本実施例の方式を用いた、シーケンシャルライトにおけるライトシーケンスの一例を示したものである。本図は、ライトシーケンスを説明するため、情報システム0101の構成要素を一部省略して記載している。

【0052】

図9との差分は、正面1001と副面1002を確保するCTL0107である。データ0901の流れを示す矢印1003、1004、1005、及び1006に示すとおり、データ0901のCTL間複写回数は、図9の3回から1回に減少している。

【0053】

図11は、ストレージシステム(メイン)0102がホスト0104から受領したI /

10

20

30

40

50

要求を処理する流れの一例を示したフロー図である。CTL0107は、ホスト0104からのI/O要求を受領する(1101)。ホスト0104からI/O要求を受領したCTL0107は、当該I/Oを担当するCTL0107を決定し、その担当CTL0107が、処理を継続する(1102)。

【0054】

担当CTL0107は、I/O要求を解析し、アクセス先のLDEV、I/O要求種別(リード要求、ライト要求、他)、及びアクセスパターン(ランダムアクセスまたはシーケンシャルアクセス)を判定する。例えば、アクセス先の連続アドレスが所定値より大きい場合に、シーケンシャルと判定される。アクセスパターンは、LDEVごとに直近のI/O履歴を取得して判定することにより、ストレージシステム(メイン)が判定してもよいし、ホスト0104からアクセスパターンに関するヒント情報を取得して判定してもよい(1103)。

10

【0055】

担当CTL0107は、ホスト0104からのI/O要求の解析1103で明らかにしたI/O要求種別を用いて、I/O要求種別がライト要求であるか否かを判定する(1104)。I/O要求種別がライト要求の場合(1104:Yes)、担当CTL0107は、ステップ1105へ進む(A)。一方、I/O要求種別がライト要求でない場合(1104:No)、担当CTL0107は、ステップ1110へ進む(B)。

【0056】

(A)

20

I/O要求種別がライト要求の場合(1104:Yes)、担当CTL0107は、PDEV管理テーブル0401とLDEV管理テーブル0402を参照し、当該I/OのLDEVに関する情報、及び当該LDEVを構成するPDEVの情報を取得する(1105)。担当CTL0107は、セグメント確保CTL決定処理を呼び出し、ライト処理に必要なキャッシュセグメントを確保するCTL0105を決定する(1106)。なお、セグメント確保CTL決定処理1106の詳細は、後述する。

【0057】

担当CTL0107は、セグメント確保CTL決定処理1106での決定に基づき、キャッシュセグメント、及びバッファセグメントを確保する(1107)。なお、セグメント確保処理1107の詳細は、後述する。担当CTL0107は、ライト処理を呼び出し、セグメント確保処理1107で確保したセグメントを使用し、ライト処理を実行する(1108)。なお、ライト処理1108の詳細は、後述する。

30

【0058】

担当CTL0107は、セグメント解放処理を呼びだし、セグメント確保処理1107で確保したセグメントの一部、または全てを解放し、処理を終了する(1109)。なお、セグメント解放処理1109の詳細は、後述する。

【0059】

(B)

I/O要求種別がライト要求でない場合(1104:No)、担当CTL0107は、I/O要求に基づき処理を実行した(1110)後、処理を終了する。当該処理は、本実施例の影響を受けないため、詳細な説明を省略する。

40

【0060】

図12は、図11で説明したフローから呼び出される、セグメント確保CTL決定処理1106の流れの一例を示したフロー図である。担当CTL0107は、ホスト0104からのI/O要求の解析1103で明らかにしたアクセスパターンを参照し、アクセスパターンがシーケンシャルであるか否かを判定する(1201)。アクセスパターンがシーケンシャルである場合(1201:Yes)、担当CTL0107は、ステップ1102に進む(A)。一方、アクセスパターンがシーケンシャルでない場合(1201:No)、担当CTL0107は、ステップ1208に進む(G)。

【0061】

50

(A)

アクセスボタンがシーケンシャルである場合 (1 2 0 1 : Y e s)、担当 C T L 0 1 0 7 は、ステップ 1 1 0 5 で取得した L D E V 情報を用いて、ストレージ間冗長化の L D E V に対するライトであるか否かを判定する (1 2 0 2)。ストレージ間冗長化の L D E V に対するライトデータは、リモートストレージシステム 0 1 0 3 に接続されている C T L 0 1 0 7 によって、リモートストレージシステム 0 1 0 3 に転送される。ストレージ間冗長化の L D E V に対するライトの場合 (1 2 0 2 : Y e s)、担当 C T L 0 1 0 7 は、ステップ 1 2 0 3 に進む (B)。一方、ストレージ間冗長化の L D E V に対するライトでない場合 (1 2 0 2 : N o)、担当 C T L 0 1 0 7 は、ステップ 1 2 0 5 に進む (C)。

【 0 0 6 2 】

10

(B)

ストレージ間冗長化の L D E V に対するライトの場合 (1 2 0 2 : Y e s)、担当 C T L 0 1 0 7 は、当該 L D E V がストレージ間冗長化をしているストレージコントローラ (リモート) 0 1 0 3 が接続されている C T L の C T L # を取得し (1 2 0 3)、ステップ 1 2 0 4 に進む (E)。なお、取得方法は、あらかじめストレージシステム同士の接続状態を管理するテーブルを作成しておき、これを参照してもよいし、ストレージシステム (メイン) 0 1 0 2 を構成する全ての C T L 0 1 0 7 に問い合わせてもよい。

【 0 0 6 3 】

(C)

ストレージ間冗長化の L D E V に対するライトでない場合 (1 2 0 2 : N o)、担当 C T L 0 1 0 7 は、当該 L D E V へのライトについて、担当 C T L 0 1 0 7 でホスト 0 1 0 4 からライトされたデータを参照する処理 (例えば、スナップショットなどライトに起因する処理) が存在するか否かを判定する (1 2 0 5)。担当 C T L 0 1 0 7 でデータを参照する処理が存在する場合 (1 2 0 5 : Y e s)、担当 C T L 0 1 0 7 は、ステップ 1 2 0 6 へ進む (D)。一方、担当 C T L 0 1 0 7 でデータを参照する処理が存在しない場合 (1 2 0 5 : N o)、担当 C T L 0 1 0 7 は、ステップ 1 2 0 7 へ進む (F)。

20

【 0 0 6 4 】

(D)

担当 C T L 0 1 0 7 でデータを参照する処理が存在する場合 (1 2 0 5 : Y e s)、担当 C T L 0 1 0 7 は、担当 C T L 自身の C T L # を取得し (1 2 0 6)、ステップ 1 2 0 4 へ進む (E)。

30

【 0 0 6 5 】

(E)

担当 C T L 0 1 0 7 は、ステップ 1 2 0 3、またはステップ 1 2 0 6 で取得した C T L # を引数とし、C T L # 指定のセグメント確保 C T L 決定処理 1 1 0 6 を呼び出し、指定した C T L # でセグメントを確保したうえで、C T L 間のデータ複写回数が最小となるように、ライト処理に必要なセグメントを確保し (1 2 0 4)、処理を終了する。なお、C T L # 指定のセグメント確保処理 1 2 0 4 の詳細は、後述する。

【 0 0 6 6 】

(F)

担当 C T L 0 1 0 7 でデータを参照する処理が存在しない場合 (1 2 0 5 : N o)、担当 C T L 0 1 0 7 は、C T L 間データ複写回数最小セグメント確保 C T L 決定処理を呼び出し、C T L 間のデータ複写回数が最小となるように、ライト処理に必要なセグメントを確保し (1 2 0 7)、処理を終了する。なお、C T L 間データ複写回数最小セグメント確保 C T L 決定処理 1 2 0 7 の詳細は、後述する。

40

【 0 0 6 7 】

(G)

アクセスボタンがシーケンシャルでない場合 (1 2 0 1 : N o)、担当 C T L 0 1 0 7 は、キャッシュセグメントを「C T L の指定なしで確保」することに決定し、結果を返却する (1 2 0 8)。

50

【 0 0 6 8 】

図 1 3 は、図 1 2 で説明したフローから呼び出される、CTL # 指定のセグメント確保 CTL 決定処理 1 2 0 4 の流れの一例を示したフロー図である。担当 CTL 0 1 0 7 は、キャッシュセグメントの正面を「担当 CTL で確保」することに決定する (1 3 0 1)。

【 0 0 6 9 】

担当 CTL 0 1 0 7 は、引数として渡された指定の CTL # を用いて、指定の CTL と担当 CTL 0 1 0 7 が一致するか否かを判定する (1 3 0 2)。指定の CTL と担当 CTL 0 1 0 7 が一致する場合 (1 3 0 2 : Y e s)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 3 に進む (A)。一方、指定の CTL と担当 CTL 0 1 0 7 が一致しない場合 (1 3 0 2 : N o)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 9 に進む (G)。

10

【 0 0 7 0 】

(A)

指定の CTL と担当 CTL が一致する場合 (1 3 0 2 : Y e s)、担当 CTL 0 1 0 7 は、担当 CTL 0 1 0 7 にライト先の P D E V が接続されているか否かを判定する (1 3 0 3)。担当 CTL 0 1 0 7 にライト先の P D E V が接続されている場合 (1 3 0 3 : Y e s)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 4 に進む (B)。一方、担当 CTL にライト先の P D E V が接続されていない場合 (1 3 0 3 : N o)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 6 に進む (D)。

【 0 0 7 1 】

(B)

担当 CTL 0 1 0 7 にライト先の P D E V が接続されている場合 (1 3 0 3 : Y e s)、担当 CTL 0 1 0 7 は、担当 CTL 0 1 0 7 にホスト 0 1 0 4 が接続されているか否かを判定する (1 3 0 4)。担当 CTL にホスト 0 1 0 4 が接続されている場合 (1 3 0 4 : Y e s)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 5 に進む (C)。一方、担当 CTL 0 1 0 7 にホスト 0 1 0 4 が接続されていない場合 (1 3 0 4 : N o)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 7 に進む (E)。

20

【 0 0 7 2 】

(C)

担当 CTL にホスト 0 1 0 4 が接続されている場合 (1 3 0 4 : Y e s)、担当 CTL 0 1 0 7 は、キャッシュセグメントの副面を「担当 CTL 以外で確保」することに決定し (1 3 0 5)、結果を返却し、処理を終了する。

30

【 0 0 7 3 】

(D)

担当 CTL にライト先の P D E V が接続されていない場合 (1 3 0 3 : N o)、担当 CTL 0 1 0 7 は、担当 CTL 0 1 0 7 にホスト 0 1 0 4 が非接続であるか否かを判定する (1 3 0 6)。担当 CTL にホスト 0 1 0 4 が非接続である場合 (1 3 0 6 : Y e s)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 7 に進む (E)。一方、担当 CTL 0 1 0 7 にホスト 0 1 0 4 が非接続でない (接続されている) 場合 (1 3 0 6 : N o)、担当 CTL 0 1 0 7 は、ステップ 1 3 0 8 に進む (F)。

40

(E)

【 0 0 7 4 】

担当 CTL にホスト 0 1 0 4 が非接続である場合 (1 3 0 6 : Y e s)、担当 CTL 0 1 0 7 は、キャッシュセグメントの副面を「ホスト接続 CTL で確保」することに決定し (1 3 0 7)、結果を返却し、処理を終了する。

【 0 0 7 5 】

(F)

担当 CTL にホスト 0 1 0 4 が接続されている場合 (1 3 0 6 : N o)、担当 CTL 0 1 0 7 は、キャッシュセグメントの副面を「P D E V 接続 CTL で確保」することに決定し (1 3 0 8)、結果を返却し、処理を終了する。

50

【 0 0 7 6 】

(G)

指定のCTLと担当CTL0107が一致しない場合(1302:No)、担当CTL0107は、キャッシュセグメントの副面を「指定CTLで確保」することに決定し(1309)、結果を返却し、処理を終了する。

【0077】

上述のように、CTL#指定のセグメント確保CTL決定処理1204の開始時、リモートストレージシステムに接続されたCTLまたは担当CTLが、キャッシュセグメントを確保するCTLとして指定されている。リモートストレージシステムに接続されたCTLが指定されている場合(1302:No)、当該CTLのキャッシュセグメントが確保される(1309)。

【0078】

担当CTLがライトデータを参照する処理が存在する場合(1205:Yes)、担当CTLのキャッシュセグメントが確保される(1301)。上記処理のために担当CTLはライトデータの複製を必要とする。また、担当CTLの自キャッシュセグメントの確保の負荷は、他のCTLのキャッシュセグメントの確保よりも小さい。したがって、CTL間のデータ複製の回数を低減しつつ、レスポンスを改善できる。

【0079】

ストレージ間冗長化のライトのため、リモートストレージシステムに接続されているCTLはライトデータの複製を必要とする。ストレージ間冗長化のライトは長い時間を必要とし、さらに、当該ライトが完了した後にホストに対してライト完了の応答が送信される。したがって、リモートストレージシステムに接続されているCTLのキャッシュセグメントを確保することで、CTL間のデータ複製の回数を低減しつつ、レスポンスを改善できる。

【0080】

CTL#指定のセグメント確保CTL決定処理1204の開始時、担当CTLが指定され(1302:Yes)、担当CTLがPDEVに接続されておらず(1303:No)、担当CTLがホストに接続されている場合(1306:No)、PDEVに接続されているCTLのキャッシュセグメントが確保される(1308)。PDEVはライトデータの最終格納媒体であり、PDEVに接続されているCTLは、ホスト又は他のCTLからライトデータを受信する。そのため、CTL間でのライトデータの複製回数を低減できる。

【0081】

CTL#指定のセグメント確保CTL決定処理1204の開始時、担当CTLが指定され(1302:Yes)、担当CTLがホストに接続されていない場合(1304:Noまたは1306:Yes)、ホストに接続されているCTLのキャッシュセグメントが確保される(1307)。ホストに接続されているCTLがライトデータをホストから受信するため、CTL間でのライトデータの複製回数を低減できる。

【0082】

CTL#指定のセグメント確保CTL決定処理1204の開始時、担当CTLが指定され(1302:Yes)、担当CTLがPDEVに接続されており(1303:Yes)、担当CTLがホストに接続されている場合(1304:Yes)、担当CTL以外のCTLのキャッシュセグメントが確保される(1305)。担当CTL以外のCTLは、ライトデータの複製を必要としておらず、キャッシュデータの二重化のために、任意のCTLを選択可能である。

【0083】

なお、CTL#指定のセグメント確保CTL決定処理1204の開始時リモートストレージシステムに接続されたCTLが指定されている場合、ステップ1301は、担当CTLと異なるCTLを指定してもよい。上記フローにおいて、担当CTLのキャッシュセグメントが副面で、もう一方のCTLのキャッシュセグメントが正面であってもよい。

【0084】

10

20

30

40

50

図14は、図12で説明したフローから呼び出される、CTL間データ複写回数最小セグメント確保CTL決定処理1207の流れの一例を示したフロー図である。

【0085】

担当CTL0107は、担当CTL0107にライト先のPDEVが接続されているかを判定する(1401)。担当CTL0107にライト先のPDEVが接続されている場合(1401:Yes)、担当CTL0107は、ステップ1402に進む(A)。一方、担当CTL0107にライト先のPDEVが接続されていない場合(1401:No)、担当CTL0107は、ステップ1406に進む(D)。

【0086】

(A)

担当CTL0107にライト先のPDEVが接続されている場合(1401:Yes)、担当CTL0107は、ホスト接続CTLにライト先のPDEVが接続されているかを判定する(1402)。ホスト接続CTLにライト先のPDEVが接続されている場合(1402:Yes)、担当CTL0107は、ステップ1403に進む(B)。一方、ホスト接続CTLにライト先のPDEVが接続されていない場合(1402:No)、担当CTL0107は、ステップ1407に進む(E)。

【0087】

(B)

ホスト接続CTLにライト先のPDEVが接続されている場合(1402:Yes)、担当CTLは、担当CTL0107にホスト0104が接続されているかを判定する(1403)。担当CTL0107にホスト0104が接続されている場合(1403:Yes)、担当CTL0107は、ステップ1404に進む(C)。一方、担当CTL0107にホスト0104が接続されていない場合(1403:No)、担当CTL0107は、ステップ1407に進む(E)。

【0088】

(C)

担当CTL0107にホスト0104が接続されている場合(1403:Yes)、担当CTL0107は、キャッシュセグメントの正面を「担当CTLで確保」することに決定する(1404)。担当CTL0107は、キャッシュセグメントの副面を「担当CTL以外で確保」することに決定し(1405)、結果を返却し、処理を終了する。

【0089】

(D)

担当CTL0107にライト先のPDEVが接続されていない場合(1401:No)、担当CTL0107は、ホスト接続CTLにライト先のPDEVが接続されているかを判定する(1406)。ホスト接続CTLにライト先のPDEVが接続されている場合(1406:Yes)、担当CTL0107は、ステップ1407に進む(E)。一方、ホスト接続CTLにライト先のPDEVが接続されていない場合(1406:No)、担当CTL0107は、ステップ1409に進む(F)。

【0090】

(E)

ホスト接続CTLにライト先のPDEVが接続されている場合(1406:Yes)、担当CTL0107は、キャッシュセグメントの正面を「担当CTLで確保」することに決定する(1407)。担当CTL0107は、キャッシュセグメントの副面を「ホスト接続CTLで確保」することに決定し(1408)、結果を返却し、処理を終了する。

【0091】

(F)

ホスト接続CTLにライト先のPDEVが接続されていない場合(1406:No)、担当CTL0107は、キャッシュセグメントの正面を「ホスト接続CTLで確保」することに決定する(1409)。担当CTL0107は、キャッシュセグメントの副面を「PDEV接続CTLで確保」することに決定し(1410)、結果を返却し、処理を終了

10

20

30

40

50

する。

【 0 0 9 2 】

図 1 4 を参照した上記フローは、シーケンシャルアクセス時のライトに限定して、ホストが接続された C T L にキャッシュの正面及び副面の一方の面を配置し、P D E V が接続されている C T L に他方の面を配置する。一つの C T L がホスト及び P D E V に接続されている場合、当該 C T L 及び他の C T L が選択される。ホストに接続された C T L はホストからライトデータを受信する。P D E V に接続された C T L は、ライトデータを P D E V に格納する。したがって、キャッシュデータの二重化のための C T L 間の複製回数を低減することができる。

【 0 0 9 3 】

図 1 2、1 3 及び 1 4 を参照して説明したように、リモートストレージシステムに接続された C T L、ライトデータの処理を行う担当 C T L、ホストからライトデータを受信する C T L、及び P D E V にデータを格納する C T L は、それぞれ、ライトデータのキャッシュと独立して、異なる目的のために、ライトデータをホストまたは他の C T L から受信することが決まっている。したがって、これらがキャッシュセグメントを確保する C T L として指定可能であることで、ライトデータのキャッシュ多重化のためのデータ転送回数を低減できる。

【 0 0 9 4 】

図 1 5 は、図 1 1 で説明したフローから呼び出される、セグメント確保処理 1 1 0 7 の流れの一例を示したフロー図である。担当 C T L 0 1 0 7 は、キャッシュセグメント確保 C T L 決定処理 1 1 0 6 での決定に基づき、キャッシュセグメントが C T L 指定されているか否かを判定する (1 5 0 1)。キャッシュセグメントが C T L 指定されている場合 (1 5 0 1 : Y e s)、担当 C T L 0 1 0 7 は、ステップ 1 5 0 2 に進む (A)。一方、キャッシュセグメントが C T L 指定されていない場合 (1 5 0 1 : N o)、担当 C T L 0 1 0 7 は、ステップ 1 5 0 3 に進む (B)。

【 0 0 9 5 】

(A)

キャッシュセグメントが C T L 指定されている場合 (1 5 0 1 : Y e s)、担当 C T L 0 1 0 7 は、指定された C T L にてキャッシュセグメントの正面と副面を確保し、担当 C T L 0 1 0 7 は、ステップ 1 5 0 4 に進む (C)。

【 0 0 9 6 】

(B)

キャッシュセグメントが C T L 指定されていない場合 (1 5 0 1 : N o)、担当 C T L 0 1 0 7 は、C T L ごとのキャッシュセグメント使用量を算出し、C T L 間で使用量が偏らないように正面と副面を選択した異なる C T L にそれぞれ確保 (例えば、最も使用量の少ない C T L から順に正面と副面を確保) し、処理 1 5 0 4 に進む (C)。これにより、キャッシュヒット率を上げることができる。なお、キャッシュセグメントを確保する C T L は、必ずしも C T L ごとのキャッシュセグメント使用量に基づいて決定する必要はなく、例えば、前回に確保した C T L を記憶しておき、ラウンドロビンで正面と副面を確保する C T L を決定してもよい。

【 0 0 9 7 】

(C)

担当 C T L 0 1 0 7 は、ホスト 0 1 0 4 が接続された C T L において正面又は副面が確保されているか否かを判定する (1 5 0 4)。ホスト 0 1 0 4 が接続された C T L において正面又は副面が確保されている場合 (1 5 0 4 : Y e s)、担当 C T L 0 1 0 7 は、ステップ 1 5 0 6 に進む (E)。一方、ホスト 0 1 0 4 が接続された C T L において正面又は副面が確保されていない場合 (1 5 0 4 : N o)、担当 C T L 0 1 0 7 は、ステップ 1 5 0 5 に進む (D)。

【 0 0 9 8 】

(D)

10

20

30

40

50

ホスト0104が接続されたCTLにおいて正面又は副面が確保されていない場合(1504:No)、担当CTL0107は、ホスト0104が接続されたCTLにおいてバッファセグメントを確保し(1505)、担当CTL0107は、ステップ1506に進む(E)。

【0099】

(E)

ホスト0104が接続されたCTLにおいて正面又は副面が確保されている場合(1504:Yes)、担当CTL0107は、ライト先のPDEVが接続されたCTLにおいて正面又は副面が確保されているか否かを判定する(1506)。ライト先のPDEVが接続されたCTLにおいて正面又は副面が確保されている場合(1506:Yes)、処理を終了する。

10

【0100】

一方、ライト先のPDEVが接続されたCTLにおいて正面又は副面が確保されていない場合(1506:No)、担当CTL0107は、ステップ1507に進む。担当CTL0107は、ライト先のPDEVが接続されたCTLにおいてバッファセグメントを確保し(1507)、処理を終了する。

【0101】

図16は、図11で説明したシーケンスから呼び出される、ライト処理1108の流れの一例を示したシーケンス図である。以下において、担当CTL0107は、他のCTL0107に必要な指示を行うことで、各ステップを実行する。担当CTL0107は、ホスト0104に接続されているCTL0107を介して、ホスト0104にデータ転送を要求する(1601)。

20

【0102】

担当CTL0107は、ホスト0104に接続されているCTL0107に確保したセグメントがキャッシュであるか否かを判定する(1602)。ホスト0104に接続されているCTL0107に確保したセグメントがキャッシュである場合(1602:Yes)、担当CTL0107は、ステップ1603に進む(A)。一方、ホスト0104に接続されているCTLに確保したセグメントがキャッシュでない(バッファセグメントである)場合(1602:No)、担当CTL0107は、ステップ1605に進む(B)。

【0103】

30

(A)

ホスト0104に接続されているCTL0107に確保したセグメントがキャッシュである場合(1602:Yes)、担当CTL0107は、ホスト0104から転送されたI/Oデータを当該キャッシュセグメントに格納する(1603)。担当CTL0107は、ステップ1603でデータを格納したキャッシュセグメントでない方のキャッシュセグメントへ、当該I/Oデータを複写し(1604)、担当CTL0107は、ステップ1607に進む(C)。

【0104】

(B)

ホスト0104に接続されているCTLに確保したセグメントがバッファセグメントである場合(1602:No)、担当CTL0107は、ホスト0104から転送されたI/Oデータを当該バッファセグメントに格納する(1605)。担当CTL0107は、キャッシュセグメントの正面と副面の両方に、当該I/Oデータを複写し(1606)、担当CTL0107は、ステップ1607に進む(C)。

40

【0105】

(C)

担当CTL0107は、ストレージ冗長化のLDEVに対するライトであるか否かを判定する(1607)。ストレージ冗長化のLDEVに対するライトである場合(1607:Yes)、担当CTL0107は、ステップ1608に進む(D)。一方、ストレージ冗長化のLDEVに対するライトでない場合(1607:No)、担当CTL0107は、

50

ステップ1609に進む(E)。

【0106】

(D)

ストレージ冗長化のLDEVに対するライトである場合(1607:Yes)、担当CTL0107は、リモートストレージシステム0103に接続されているCTL0107より、当該I/Oデータを転送する(1608)。

【0107】

(E)

ストレージ冗長化のLDEVに対するライトでない場合(1607:No)、該当CTLは0107、ホスト0104にライトの完了を応答する(1609)。担当CTL0107は、担当CTL0107で固有の処理を実行する(1610)。例えば、スナップショットに関する処理であり、RAIDグループを構成している場合、担当CTL0107は、パリティの生成を実施する。

【0108】

担当CTL0107は、前述の処理でデータを複製したバッファセグメントのほかに、確保したバッファセグメントが存在するか否かを判定する(1611)。他に確保したバッファセグメントが存在する場合(1611:Yes)、担当CTL0107は、ステップ1612に進む(F)。一方、他に確保したバッファセグメントが存在しない場合(1611:No)、担当CTL0107は、ステップ1613に進む(G)。

【0109】

(F)

他に確保したバッファセグメントが存在する場合(1611:Yes)、担当CTL0107は、ホスト0104から転送されたI/Oデータを当該バッファセグメントに格納し(1612)し、担当CTL0107は、ステップ1613に進む(E)。

【0110】

(G)

担当CTL0107は、ライト先のPDEVにホスト0104から転送されたI/Oデータを転送し(1613)、処理を終了する。

【0111】

図17は、図11で説明したフローから呼び出される、セグメント解放処理1109の流れの一例を示したフロー図である。担当CTL0107は、当該I/Oのアクセスボタンがシークエンシャルであるか否かを判定する(1701)。当該I/Oのアクセスボタンがシークエンシャルであった場合(1701:Yes)、担当CTL0107は、確保した全てのキャッシュセグメント、及びバッファセグメントを解放し(1702)、処理を終了する。

【0112】

このように、アクセスボタンがシークエンシャルである場合、ライト処理1108の終了(最終記憶媒体へのライトデータの格納)に応答して、確保した全てのキャッシュセグメントが解放される。シークエンシャルアクセスされるデータは、再度アクセスされる可能性が低く、効率的にキャッシュ領域を利用できる。

【0113】

一方、当該I/Oのアクセスボタンがシークエンシャルでなかった場合(1702:No)、担当CTL0107は、正面を除く確保した副面のキャッシュセグメント、及び確保した全てのバッファセグメントを解放し(1703)、処理を終了する。これにより、キャッシュヒット率を高めることができる。

【0114】

上述のように、ストレージシステムは、ホストからライトデータを受信するコントローラと、ライトデータを処理するコントローラとに基づいて、ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定する。ストレージシステムは、さらに、各コントローラが属するストレージノードの、前記ライトデータを格納す

10

20

30

40

50

る記憶デバイスへの接続の有無に基づいて、前記二つのコントローラを指定する。

【0115】

具体的には、上記例において、キャッシュセグメントを確保するCTLとして指定されるCTLの候補は、リモートストレージシステムにライトデータを転送するCTL（存在する場合）、ライトデータを処理する担当CTL（存在する場合）、ライトデータをホストから受信するCTL及びPDEVにライトデータを格納するCTLで構成されている。

【0116】

これらのCTLの候補の一部、例えば、リモートストレージシステムにライトデータを転送するCTL（存在する場合）及びライトデータを処理する担当CTL（存在する場合）の双方または一方は、省略されてもよく、ライトデータをホストから受信するCTL及びPDEVにライトデータを格納するCTLの双方または一方が省略されてもよい。CTLの候補が一つであり、キャッシュセグメントの確保のために常に指定されてもよい。上記例は最大二つのCTLを指定するが、他の例は、最大3以上のCTLを指定してもよい。

10

【0117】

指定されるCTLの候補の優先度は、設計により決定してもよい。上記例は、リモートストレージシステムにライトデータを転送するCTL（存在する場合）及びライトデータを処理する担当CTL（存在する場合）を、ライトデータをホストから受信するCTL及びPDEVにライトデータを格納するCTLよりも優先して選択している。他の例は、これと異なる優先度が候補に与えられていてもよい。

20

【0118】

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、必ずしも説明したすべ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。

【0119】

また、上記の各構成・機能・処理部等は、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（Solid State Drive）等の記録装置、または、ICカード、SDカード等の記録媒体に置くことができる。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、製品上必ずしもすべての制御線や情報線を示しているとは限らない。実際には殆どすべての構成が相互に接続されていると考えてもよい。

30

【符号の説明】

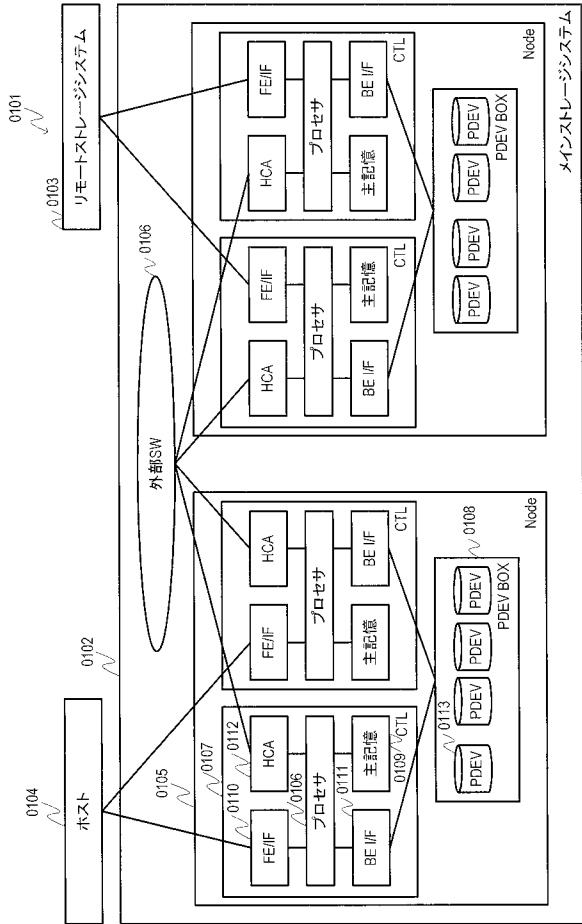
【0120】

0101 情報システム、0102 メインストレージシステム、0103 リモートストレージシステム、0104 ホスト、0105 Node、0106 外部スイッチ（SW）、0106 プロセッサ、0107 コントローラ（CTL）、0108 PDEV BOX、0109 主記憶、0110 FE I/F（Front End Interface）、0111 BE I/F（Back End Interface）、0112 HCA（Host Channel Adapter）、0113 PDEV、0201 プログラム領域、0202 管理情報領域、0203 バッファ領域、0204 キャッシュ領域、0301 CTL間振り分けプログラム、0302 I/Oプログラム、0303 I/Oパターン判定プログラム、0304 セグメント確保プログラム、0305 セグメント解放プログラム、0306 フリーセグメント確保プログラム、0307 ストレージコントローラ間データ転送プログラム

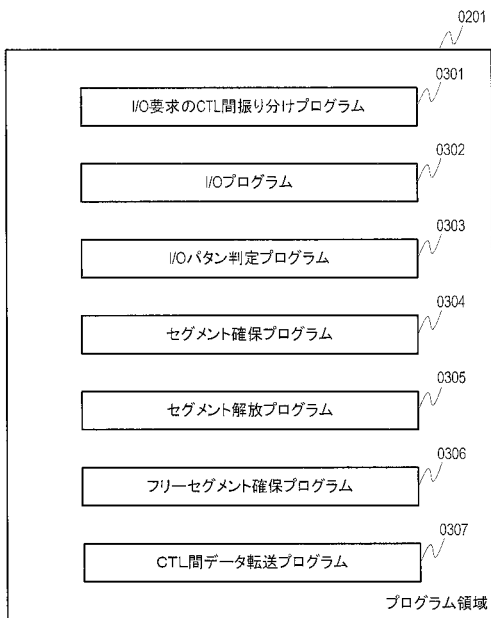
40

50

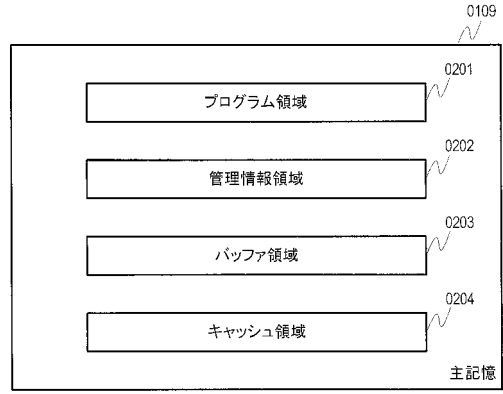
【 図 1 】



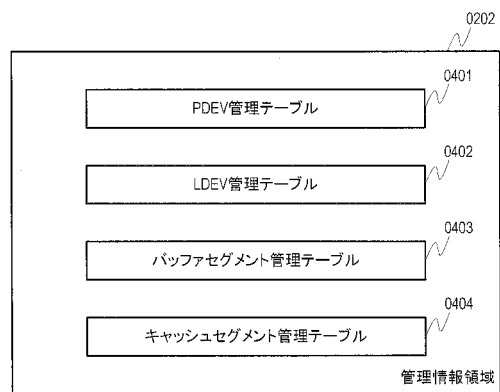
【 図 3 】



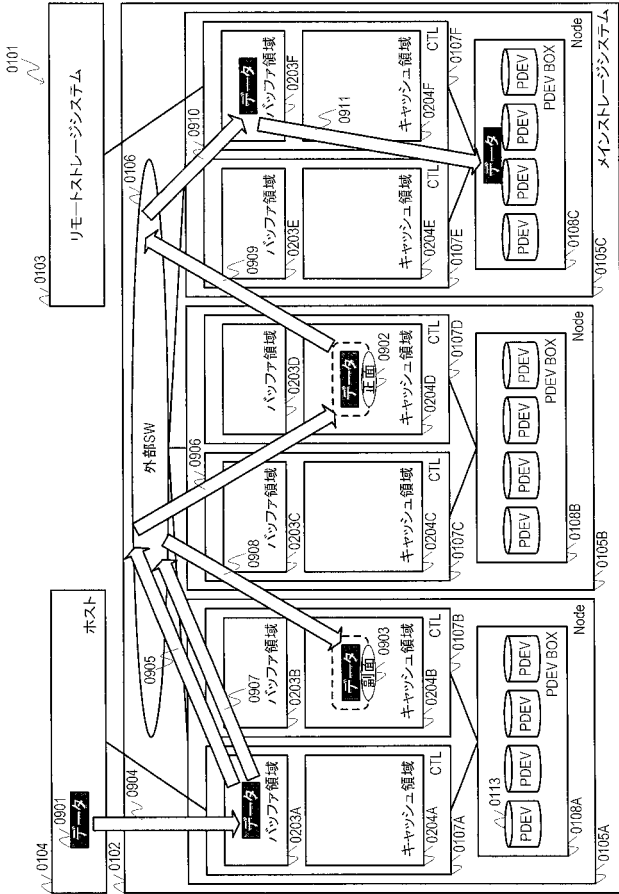
【 図 2 】



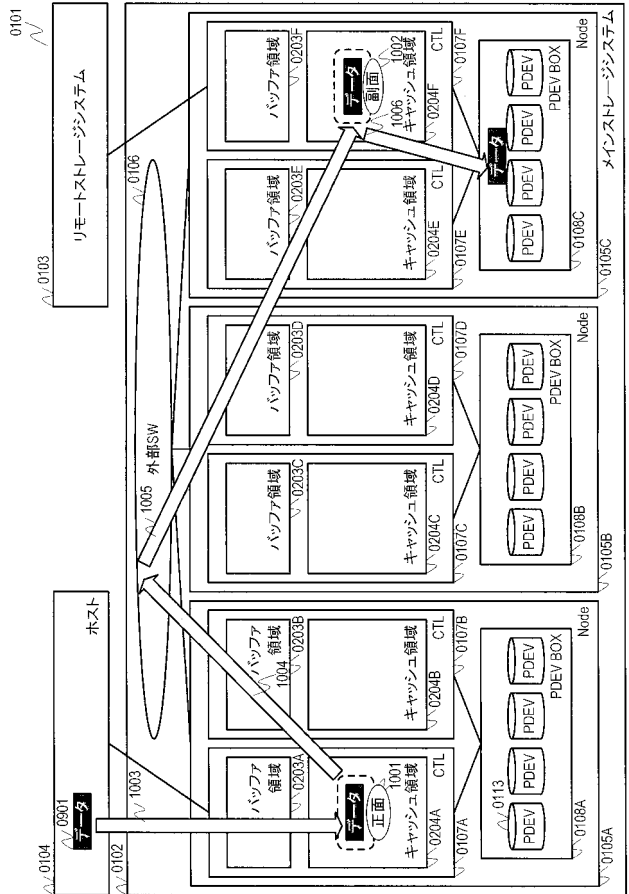
【 図 4 】



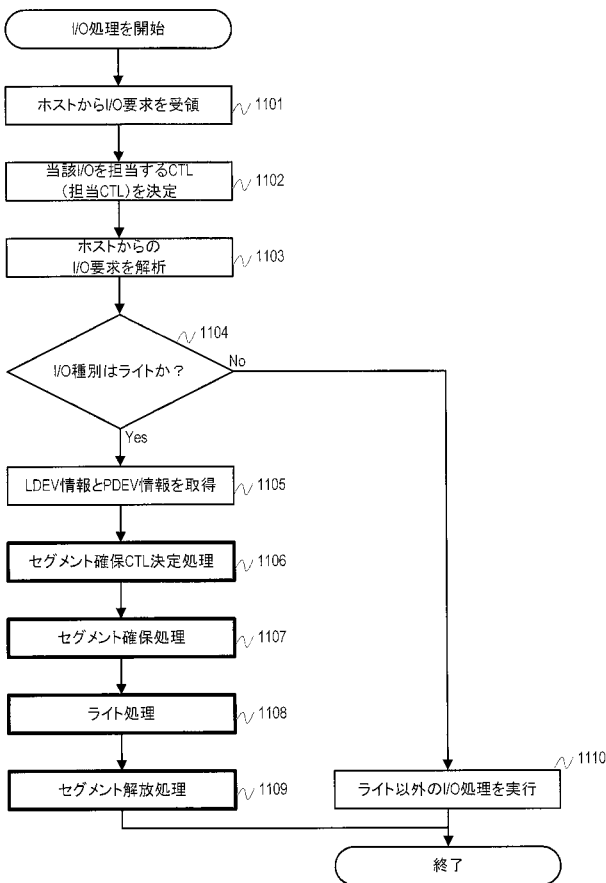
【図 9】



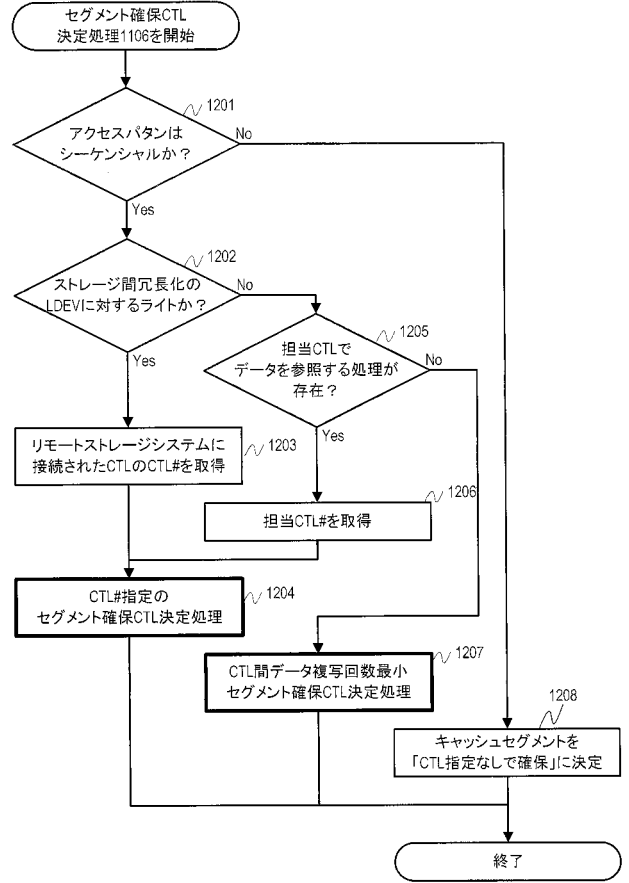
【図 10】



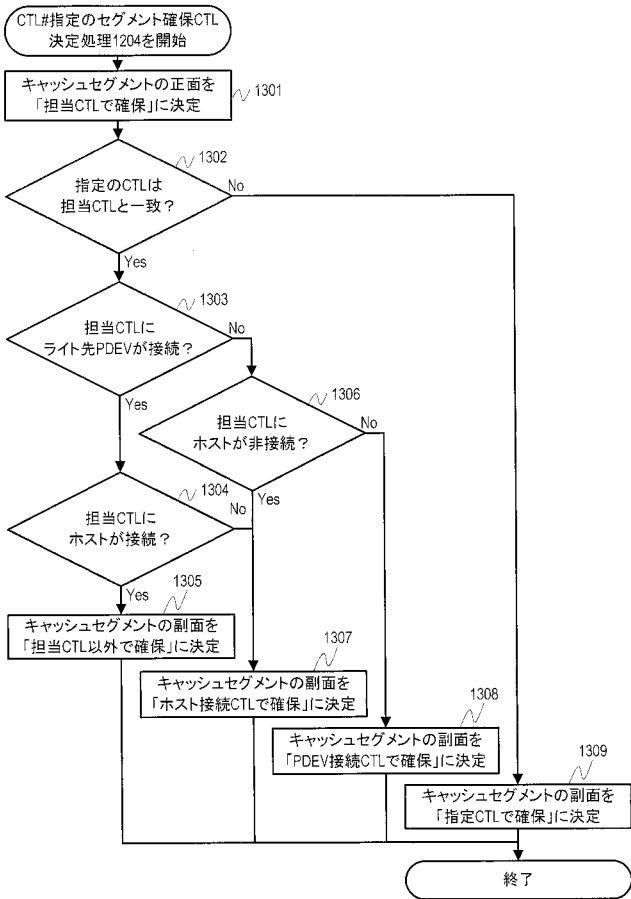
【図 11】



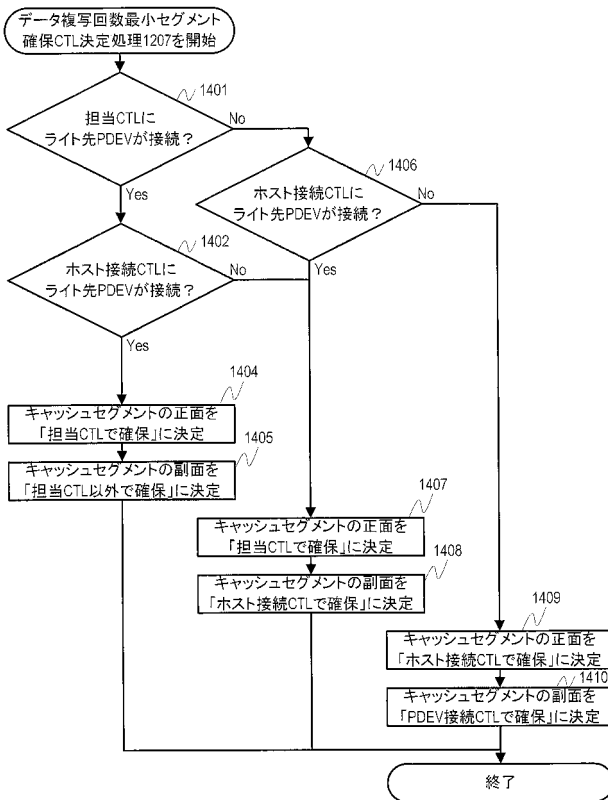
【図 12】



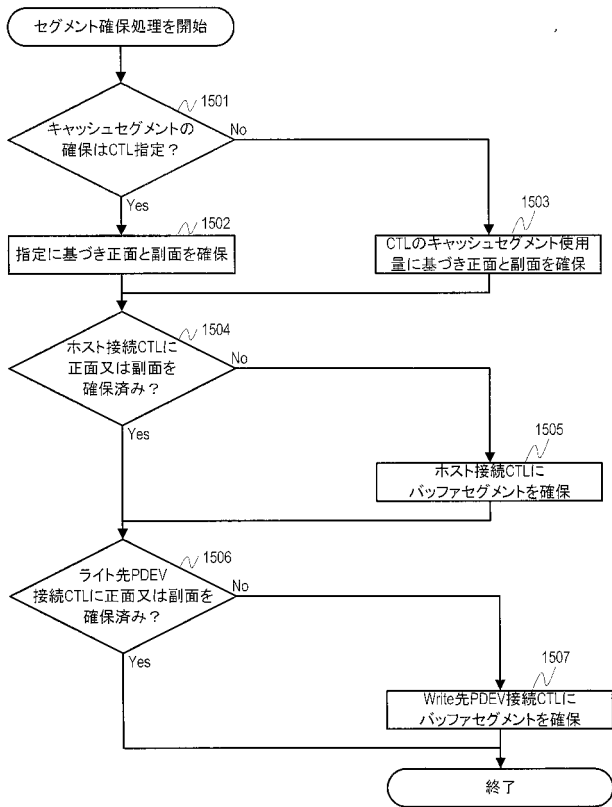
【 図 1 3 】



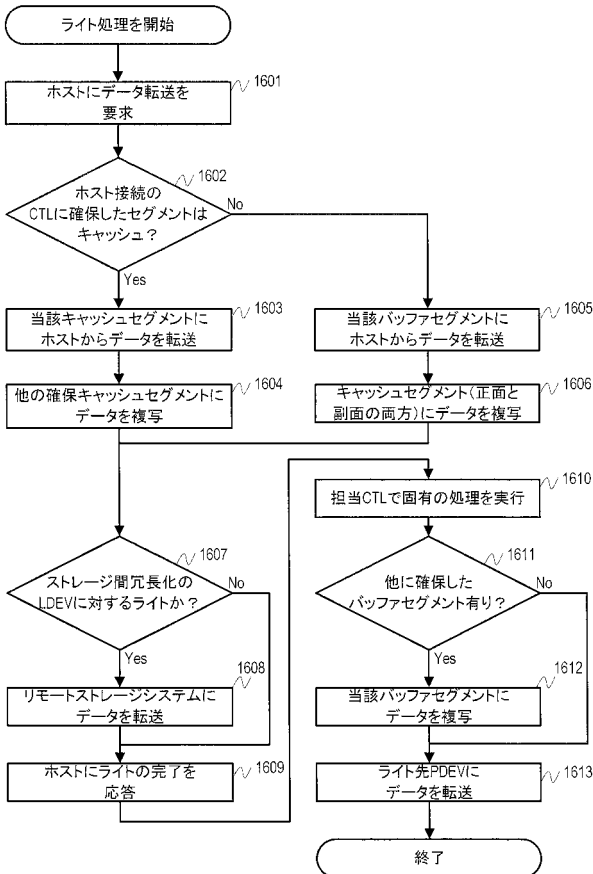
【 図 1 4 】



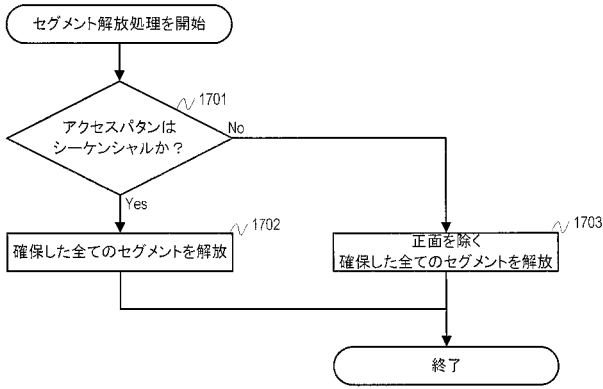
【 図 1 5 】



【 図 1 6 】



【 図 1 7 】



【 手続補正書 】

【 提出日 】 令和2年1月24日 (2020.1.24)

【 手続補正 1 】

【 補正対象書類名 】 明細書

【 補正対象項目名 】 0 0 1 9

【 補正方法 】 変更

【 補正の内容 】

【 0 0 1 9 】

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。以下の説明では、「RAID」は、Redundant Array of Independent（or Inexpensive）Disksの略である。

【 手続補正 2 】

【 補正対象書類名 】 明細書

【 補正対象項目名 】 0 0 2 7

【 補正方法 】 変更

【 補正の内容 】

【 0 0 2 7 】

メインストレージシステム0102は、1つ以上のNode0105（ストレージノードとも呼ぶ）、及び1つ以上の外部SW（SW：SWitch）0116から構成されている。図1では、Node0105は、外部SW0106を介して接続されているが、Node0105間で直接接続されていてもよい。

【 手続補正 3 】

【 補正対象書類名 】 明細書

【補正対象項目名】0029

【補正方法】変更

【補正の内容】

【0029】

図1では、CTL0107は、直接PDEV BOX0108に接続されているが、Switchなどを介して接続されていてもよい。また、PDEV BOX0108を介さず、CTL0107とPDEV0113は直接接続されていてもよい。さらに、同一Node0105内のCTL0107間は、外部SW0116を介して接続されているが、直接接続されていてもよい。

【手続補正4】

【補正対象書類名】明細書

【補正対象項目名】0031

【補正方法】変更

【補正の内容】

【0031】

プロセサ0106は、CTL0107全体を制御し、主記憶0109に格納されたマイクロプログラムに基づき動作する。FEI/F0110は、プロセサ0106により制御され、ホスト0104、及びリモートストレージシステム0103とI/O要求やI/Oデータの送受信などを実施する。BEI/F0111は、プロセサ0106により制御され、PDEV BOX0108を介し、PDEV0113とI/Oデータの送受信などを実施する。HCA0112は、プロセサ0106により制御され、外部SW0116を介して他のCTL0107と制御情報やI/Oデータの送受信などを実施する。

【手続補正5】

【補正対象書類名】明細書

【補正対象項目名】0041

【補正方法】変更

【補正の内容】

【0041】

PDEV管理テーブル0401は、メインストレージシステム102内の全てのPDEV0113の状態や、CTL0107との対応関係を示す。LDEV管理テーブル0402は、メインストレージシステム102内の全てのPDEV0113とLDEVの対応関係を示す。バッファセグメント管理テーブル0403は、メインストレージシステム102内の全てのバッファ領域0203を管理するのに用いられる。キャッシュセグメント管理テーブル0404は、メインストレージシステム102内の全てのキャッシュ領域0204を管理するのに用いられる。

【手続補正6】

【補正対象書類名】明細書

【補正対象項目名】0042

【補正方法】変更

【補正の内容】

【0042】

図5は、PDEV管理テーブル0401の構成例を示す図である。PDEV管理テーブル0401は、PDEV#0501、容量0502、状態0503、及び接続CTL#0504のエントリを有する。PDEV#0501は、PDEV0113の識別子である。容量0502は、PDEV0113のデータを格納可能な容量を示す。状態0503は、PDEV0113が正常に動作中か否か（故障していないか）の状態を示す。接続CTL#0504は、PDEV0113に接続され、アクセスするCTL0107を示す。

【手続補正7】

【補正対象書類名】明細書

【補正対象項目名】0044

【補正方法】変更

【補正の内容】

【0044】

LDEV#0601は、LDEVの識別子である。容量0602は、LDEVにデータを格納可能な容量を示す。状態0503は、LDEVへ正常にI/O可能か否かを示す。冗長構成0604は、RAIDによる冗長化、またはリモートストレージシステム0103とのストレージシステム冗長化の状態を示す。構成情報0605は、LDEVに属するPDEV0113、及びストレージシステム冗長化の対応リモートストレージシステム0103を示す。

【手続補正8】

【補正対象書類名】明細書

【補正対象項目名】0053

【補正方法】変更

【補正の内容】

【0053】

図11は、ストレージシステム(メイン)0102がホスト0104から受領したI/O要求を処理する流れの一例を示したフロー図である。CTL0107は、ホスト0104からのI/O要求を受領する(1101)。ホスト0104からI/O要求を受領したCTL0107は、当該I/Oを担当するCTL0107を決定し、その担当CTL0107が、処理を継続する(1102)。

【手続補正9】

【補正対象書類名】明細書

【補正対象項目名】0056

【補正方法】変更

【補正の内容】

【0056】

(A)

I/O要求種別がライト要求の場合(1104:Yes)、担当CTL0107は、PDEV管理テーブル0401とLDEV管理テーブル0402を参照し、当該I/OのLDEVに関する情報、及び当該LDEVを構成するPDEVの情報を取得する(1105)。担当CTL0107は、セグメント確保CTL決定処理を呼び出し、ライト処理に必要なキャッシュセグメントを確保するCTL0107を決定する(1106)。なお、セグメント確保CTL決定処理1106の詳細は、後述する。

【手続補正10】

【補正対象書類名】明細書

【補正対象項目名】0060

【補正方法】変更

【補正の内容】

【0060】

図12は、図11で説明したフローから呼び出される、セグメント確保CTL決定処理1106の流れの一例を示したフロー図である。担当CTL0107は、ホスト0104からのI/O要求の解析1103で明らかにしたアクセスパターンを参照し、アクセスパターンはシーケンシャルであるか否かを判定する(1201)。アクセスパターンがシーケンシャルである場合(1201:Yes)、担当CTL0107は、ステップ1202に進む(A)。一方、アクセスパターンがシーケンシャルでない場合(1201:No)、担当CTL0107は、ステップ1208に進む(G)。

【手続補正11】

【補正対象書類名】明細書

【補正対象項目名】0063

【補正方法】変更

【補正の内容】

【0063】

(C)

ストレージ間冗長化のLDEVに対するライトでない場合(1202:No)、担当CTL0107は、当該LDEVへのライトについて、担当CTL0107でHOST0104からライトされたデータを参照する処理(例えば、スナップショットなどライトに起因する処理)が存在するか否かを判定する(1205)。担当CTL0107でデータを参照する処理が存在する場合(1205:Yes)、担当CTL0107は、ステップ1206へ進む(D)。一方、担当CTL0107でデータを参照する処理が存在しない場合(1205:No)、担当CTL0107は、ステップ1207へ進む(F)。

【手続補正12】

【補正対象書類名】明細書

【補正対象項目名】0087

【補正方法】変更

【補正の内容】

【0087】

(B)

HOST接続CTLにライト先のPDEVが接続されている場合(1402:Yes)、担当CTLは、担当CTL0107にHOST0104が接続されているか否かを判定する(1403)。担当CTL0107にHOST0104が接続されている場合(1403:Yes)、担当CTL0107は、ステップ1404に進む(C)。一方、担当CTL0107にHOST0104が接続されていない場合(1403:No)、担当CTL0107は、ステップ1407に進む(E)。

【手続補正13】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

ネットワークを介して通信する複数のストレージノードを含むストレージシステムであって、

前記複数のストレージノードのそれぞれは、1以上のコントローラを含み、

前記コントローラにおける少なくとも1つのコントローラは、

HOSTからライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとに基づいて、前記ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定し、

指定したコントローラにおいて、キャッシュサブ領域を確保するストレージシステム。

【請求項2】

請求項1に記載のストレージシステムであって、

前記少なくとも1つのコントローラは、さらに、各コントローラが属するストレージノードの、前記ライトデータを格納する記憶デバイスへの接続の有無に基づいて、前記二つのコントローラを指定する

ストレージシステム。

【請求項3】

請求項1に記載のストレージシステムであって、

前記HOSTから前記ライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとは、別のコントローラであり、

前記少なくとも1つのコントローラは、前記HOSTから前記ライトデータを受信するコ

ントローラと、前記ライトデータを処理するコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定する

ストレージシステム。

【請求項 4】

請求項 3 に記載のストレージシステムであって、

前記指定した二つのコントローラのうち、前記ライトデータを格納する記憶デバイスに接続されたコントローラは、前記ライトデータを格納する記憶デバイスに接続されたコントローラの前記キャッシュサブ領域に記憶された前記ライトデータを、前記記憶デバイスに格納する

ストレージシステム。

【請求項 5】

請求項 4 に記載のストレージシステムであって、

前記指定した二つのコントローラのいずれもが、前記記憶デバイスに接続されていない場合、前記指定した二つのコントローラのうちいずれかが、前記記憶デバイスに接続された他のコントローラのバッファ領域に前記ライトデータを転送し、

前記他のコントローラが、前記バッファ領域に格納された前記ライトデータを、前記記憶デバイスに格納する

ストレージシステム。

【請求項 6】

請求項 2 に記載のストレージシステムであって、

前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとは、同じコントローラであり、

前記少なくとも一つのコントローラは、前記同じコントローラと、他の一つのコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定し、

前記指定した二つのコントローラのうち少なくとも一つは、前記記憶デバイスに接続されている

ストレージシステム。

【請求項 7】

請求項 2 に記載のストレージシステムであって、

前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを格納する記憶デバイスに接続された他のコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定する

ストレージシステム。

【請求項 8】

請求項 2 に記載のストレージシステムであって、

前記少なくとも一つのコントローラは、前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを格納する記憶デバイスに接続された他のコントローラとを、前記キャッシュサブ領域を確保するコントローラに指定する

ストレージシステム。

【請求項 9】

請求項 2 に記載のストレージシステムであって、

前記記憶デバイスは、前記ストレージノードの内に設けられた記憶媒体、またはリモートストレージシステムである

ストレージシステム。

【請求項 10】

請求項 1 に記載のストレージシステムであって、

前記少なくとも一つのコントローラは、

前記ホストからのライト要求のアクセスボタンがシーケンシャルである場合に、前記ホストから前記ライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとに基づいて、前記コントローラの指定を行い、

前記ホストからのライト要求のアクセスパターンがランダムである場合に、キャッシュ利用効率に基づいて、前記コントローラの指定を行う
ストレージシステム。

【請求項 11】

ネットワークを介して通信する複数のストレージノードを含むストレージシステムの制御方法であって、

前記複数のストレージノードのそれぞれは、1以上のコントローラを含み、

前記制御方法は、

ホストからライトデータを受信するコントローラと、前記ライトデータを処理するコントローラとに基づいて、前記ライトデータを格納するキャッシュサブ領域を確保する少なくとも二つのコントローラを指定し、

指定したコントローラにおいて、キャッシュサブ領域を確保する
制御方法。

【手続補正 14】

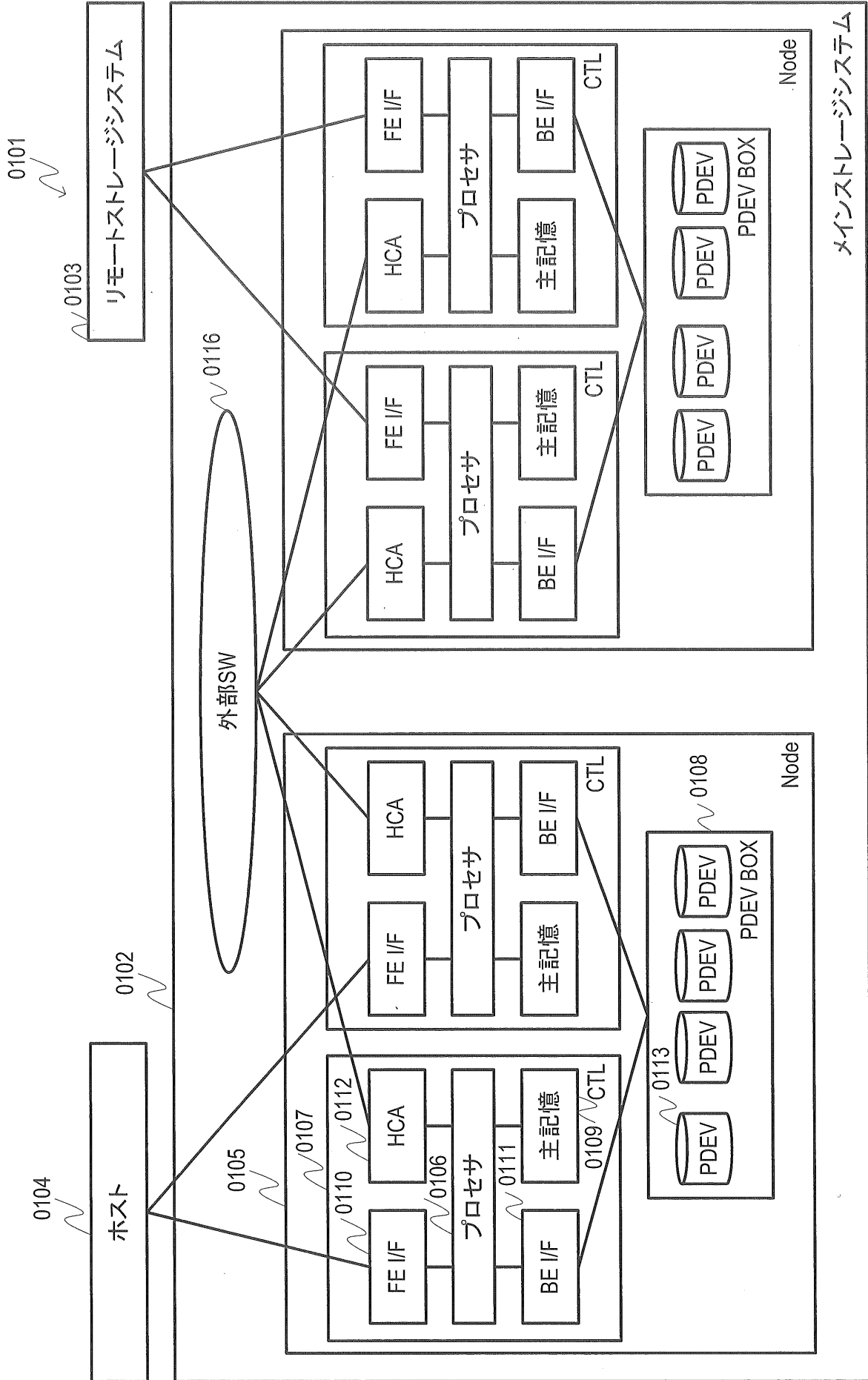
【補正対象書類名】図面

【補正対象項目名】図 1

【補正方法】変更

【補正の内容】

【図1】



【手続補正 1 5】

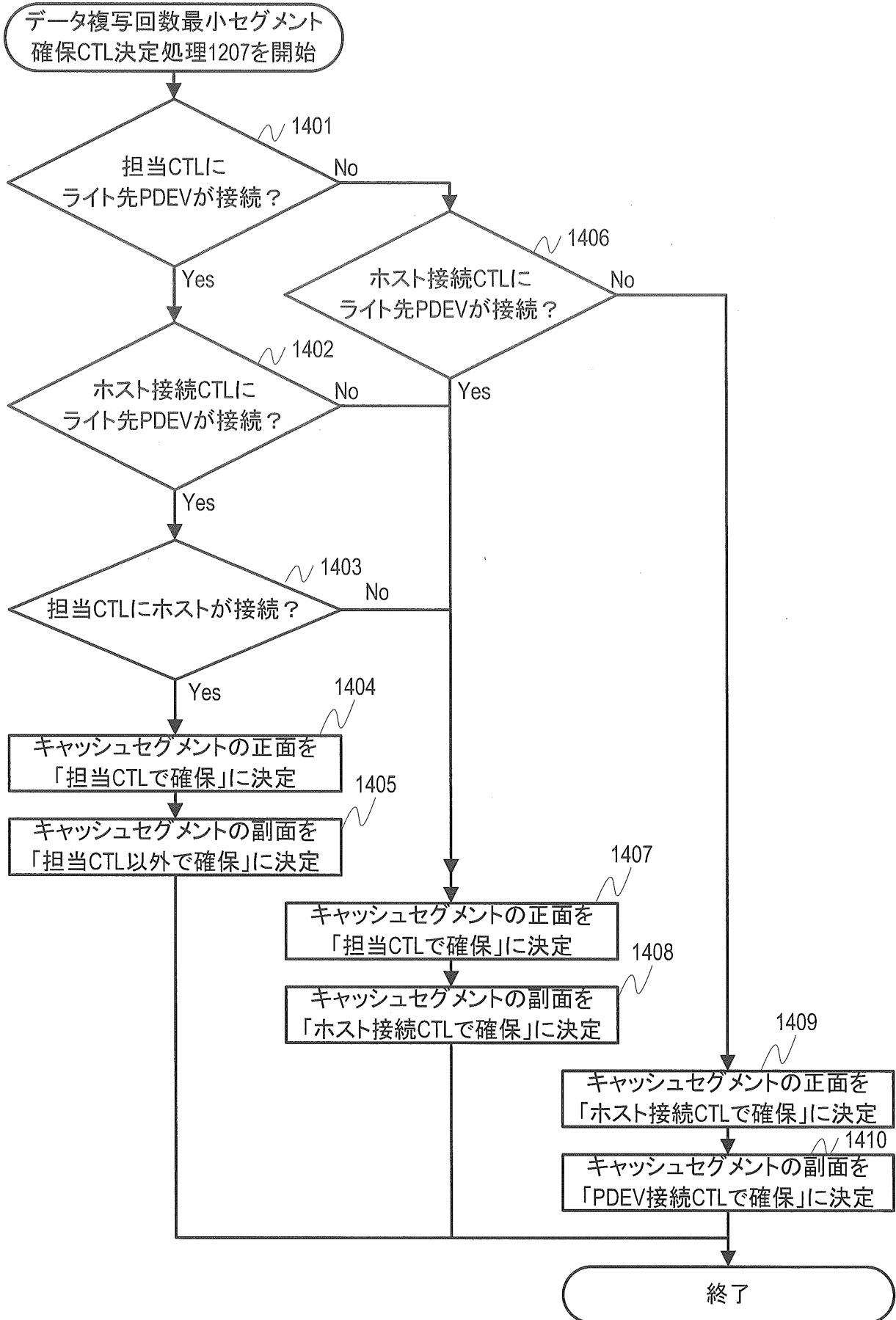
【補正対象書類名】図面

【補正対象項目名】図 1 4

【補正方法】変更

【補正の内容】

【図14】



フロントページの続き

(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 株式会社日立製作所内

Fターム(参考) 5B205 MM12