

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5806776号
(P5806776)

(45) 発行日 平成27年11月10日(2015.11.10)

(24) 登録日 平成27年9月11日(2015.9.11)

(51) Int. Cl. F I
G06F 3/06 (2006.01) G O 6 F 3/06 3 O 1 Z
G06F 13/10 (2006.01) G O 6 F 13/10 3 4 O A

請求項の数 17 (全 57 頁)

(21) 出願番号	特願2014-505435 (P2014-505435)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(86) (22) 出願日	平成23年10月7日(2011.10.7)	(74) 代理人	110001678 特許業務法人藤央特許事務所
(65) 公表番号	特表2014-527204 (P2014-527204A)	(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
(43) 公表日	平成26年10月9日(2014.10.9)	(72) 発明者	出口 彰 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
(86) 国際出願番号	PCT/JP2011/005649	(72) 発明者	坏 弘明 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
(87) 国際公開番号	W02013/051069		
(87) 国際公開日	平成25年4月11日(2013.4.11)		
審査請求日	平成26年2月3日(2014.2.3)		

最終頁に続く

(54) 【発明の名称】 ストレージシステム

(57) 【特許請求の範囲】

【請求項1】

それぞれが、複数の不揮発性半導体メモリを含む複数の記憶ドライブに基づいて構成される複数のボリュームと、

前記複数のボリュームのデータを一時的に格納するキャッシュメモリと、

複数のプロセッサパッケージと、

共有メモリと、を含み、

前記複数のプロセッサパッケージそれぞれは、前記複数のボリュームに含まれる1つのボリュームに対する入出力を担当するプロセッサと、前記1つのボリュームのデータキャッシング制御情報を格納するローカルメモリと、を含み、

前記共有メモリは、前記複数のプロセッサパッケージそれぞれのローカルメモリに格納されたデータキャッシング制御情報を格納し、前記プロセッサによってアクセス可能であり、

第1プロセッサパッケージに含まれ、第1ボリュームを担当する第1プロセッサは、計算機からの前記第1ボリュームへのリードコマンドを受信すると、前記第1プロセッサパッケージに含まれる第1ローカルメモリに格納されたデータキャッシング制御情報を更新し、

前記第1プロセッサは、前記第1ローカルメモリにおける前記データキャッシング制御情報の更新を前記共有メモリに反映するか否か、及び前記リードコマンドにおけるリードデータを前記キャッシュメモリに格納するか否かを決定し、

前記第1プロセッサは、前記第1ボリュームの記憶領域を提供する記憶ドライブに前記リードデータが格納されている場合に、前記リードデータを前記キャッシュメモリに格納することなく、かつ前記共有メモリに前記データキャッシング制御情報の更新を反映することなく、前記第1ボリュームの記憶領域を提供する記憶ドライブから読み出された前記リードデータを前記計算機に送信することを決定する、ストレージシステム。

【請求項2】

異なる制御情報を担当する複数のプロセッサと、
前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、

前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第1プロセッサが担当する制御情報を格納する共有メモリと、を含み、

前記第1プロセッサは、割り当てられている第1ローカルメモリにおいて制御情報を更新し、

前記第1プロセッサは、前記第1ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、

前記第1プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第1ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、

前記ストレージシステムは、複数のボリュームを提供するアクセス性能が異なる複数種類の不揮発性記憶領域と、キャッシュ領域とを含み、

前記第1ローカルメモリにおける前記制御情報及び前記共有メモリにおける前記制御情報は、それぞれ、前記複数のボリュームにおける前記第1プロセッサが担当する第1ボリュームのデータキャッシング制御情報を含み、

前記第1プロセッサは、前記第1ボリュームを提供する不揮発性記憶領域の種別に基づいて、前記第1ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定する、ストレージシステム。

【請求項3】

前記第1プロセッサは、当該第1プロセッサの負荷、前記キャッシュ領域の負荷及び前記第1ボリュームのキャッシュヒット率の少なくとも一つに基づいて、前記第1ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定する、請求項2に記載のストレージシステム。

【請求項4】

異なる制御情報を担当する複数のプロセッサと、
前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、

前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第1プロセッサが担当する制御情報を格納する共有メモリと、を含み、

前記第1プロセッサは、割り当てられている第1ローカルメモリにおいて制御情報を更新し、

前記第1プロセッサは、前記第1ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、

前記第1プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第1ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、

前記共有メモリにおける前記制御情報はカウントされる数値を含み、

前記第1ローカルメモリにおける前記制御情報は、前記数値の前の更新からの変化を示す差分値を含み、

前記第1プロセッサは、前記差分値が規定数に達すると、前記第1ローカルメモリにおける前記制御情報に基づいて前記共有メモリにおける前記制御情報に含まれる前記数値を更新する、ストレージシステム。

10

20

30

40

50

【請求項 5】

異なる制御情報を担当する複数のプロセッサと、
前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、
前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第 1 プロセッサが担当する制御情報を格納する共有メモリと、を含み、
前記第 1 プロセッサは、割り当てられている第 1 ローカルメモリにおいて制御情報を更新し、
前記第 1 プロセッサは、前記第 1 ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、
前記第 1 プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第 1 ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、
 前記ストレージシステムはボリュームを提供する 1 以上の不揮発性記憶領域を含み、
 前記第 1 ローカルメモリにおける前記制御情報及び前記共有メモリにおける前記制御情報は、それぞれ、前記ボリュームにおける記憶領域へのアクセス数の情報を含み、
 前記第 1 プロセッサは、前記ボリュームにおける前記記憶領域へのアクセスに応答して前記第 1 ローカルメモリにおける前記アクセス数の情報を更新し、
 前記第 1 プロセッサは、前記第 1 ローカルメモリにおけるアクセス数の情報の更新回数が規定値に達すると、前記第 1 ローカルメモリにおける前記アクセス数の情報の更新を、
 前記共有メモリにおける前記アクセス数の情報に反映する、ストレージシステム。

10

20

【請求項 6】

異なる制御情報を担当する複数のプロセッサと、
前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、
前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第 1 プロセッサが担当する制御情報を格納する共有メモリと、を含み、
前記第 1 プロセッサは、割り当てられている第 1 ローカルメモリにおいて制御情報を更新し、
前記第 1 プロセッサは、前記第 1 ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、
前記第 1 プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第 1 ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、
 前記ストレージシステムは、プライマリボリュームと、当該プライマリボリュームとコピーペアを構成するセカンダリボリュームと、前記プライマリボリュームの更新データを前記セカンダリボリュームにコピーする前に前記更新データを更新順序に従って格納するジャーナルボリュームと、前記ジャーナルボリュームにおける更新データの順序を示すシーケンス番号を含むジャーナル管理情報と、をさらに含み、
 前記共有メモリにおける前記制御情報は、前記ジャーナル管理情報における先頭シーケンス番号を示す値を含み、
 前記第 1 ローカルメモリにおける前記制御情報は、前記ジャーナル管理情報における先頭シーケンス番号を示す値と、前記第 1 ローカルメモリにおける前記値が示す先頭シーケンス番号と前記共有メモリにおける前記値が示す先頭シーケンス番号との差分を示す値と、を含み、
 前記第 1 プロセッサは、前記ジャーナルボリュームへの更新データの格納に応答して、前記第 1 ローカルメモリにおける前記先頭シーケンス番号を示す値と前記差分を示す値とを更新し、
 前記第 1 プロセッサは、前記差分を示す値が規定値に達すると、前記第 1 ローカルメモリにおける前記先頭シーケンス番号を示す値の更新を、前記共有メモリにおける前記先頭

30

40

50

シーケンス番号を示す値に反映する、ストレージシステム。

【請求項 7】

障害発生に起因して前記第 1 プロセッサの担当を引きついだ第 2 プロセッサは、前記共有メモリにおける前記先頭シーケンス番号を示す値を取得し、前記ジャーナル管理情報において、前記取得した値が示す先頭シーケンス番号より先のシーケンス番号領域を検索して、前記ジャーナル管理情報における先頭シーケンス番号を特定する、請求項 6 に記載のストレージシステム。

【請求項 8】

異なる制御情報を担当する複数のプロセッサと、
前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、

前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第 1 プロセッサが担当する制御情報を格納する共有メモリと、を含み、

前記第 1 プロセッサは、割り当てられている第 1 ローカルメモリにおいて制御情報を更新し、

前記第 1 プロセッサは、前記第 1 ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、

前記第 1 プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第 1 ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、

前記ストレージシステムは、プライマリボリュームと、前記プライマリボリュームとコピーペアを構成するセカンダリボリュームと、を含み、

前記共有メモリにおける前記制御情報は、それぞれが前記プライマリボリュームの複数の部分領域のそれぞれに対応し、前記プライマリボリュームと前記セカンダリボリュームとの間に差が存在するか否かを示す複数の差分フラグを含み、

前記第 1 ローカルメモリにおける制御情報は、それぞれが前記プライマリボリュームの前記複数の部分領域のそれぞれに対応し、前記プライマリボリュームと前記セカンダリボリュームとの間に差が存在するか否かを示す複数の差分フラグを含み、

前記第 1 ローカルメモリにおける制御情報は、それぞれが前記第 1 ローカルメモリにおける前記複数の差分フラグの一部の複数の差分フラグに対応し、当該一部の複数の差分フラグの更新を前記共有メモリにおける前記制御情報に反映するか否かを示す、複数の反映制御フラグを含み、

前記第 1 プロセッサは、前記反映制御フラグが反映を指示する差分フラグの更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステム。

【請求項 9】

前記複数の反映制御フラグのそれぞれは、対応する複数の差分フラグにおいて、差が存在することを示す差分フラグの比率が規定値に達している場合に、前記対応する複数の差分フラグの更新を前記共有メモリにおける前記制御情報に反映しないことを示す、請求項 8 に記載のストレージシステム。

【請求項 10】

障害発生に起因して前記第 1 プロセッサの担当を引きついだ第 2 プロセッサは、前記共有メモリにおける前記複数の反映制御フラグを参照し、前記参照した反映制御フラグにおいて、前記共有メモリにおける前記制御情報に更新を反映しないことを示す反映制御フラグを特定し、

前記特定した反映制御フラグに対応する前記プライマリボリュームの領域における全データを、前記セカンダリボリュームにコピーする、請求項 8 に記載のストレージシステム。

【請求項 11】

前記ストレージシステムは、スイッチを含むパスにより接続された第 1 ストレージモジュールと第 2 ストレージモジュールとを含み、

10

20

30

40

50

前記第 1 プロセッサ及び前記第 1 ローカルメモリは前記第 1 ストレージモジュール内に実装され、

前記共有メモリは前記第 2 ストレージモジュール内に実装され、

前記第 1 プロセッサは、前記バスにおける負荷に基づいて前記第 1 ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定する、請求項 3 に記載のストレージシステム。

【請求項 1 2】

異なる制御情報を担当する複数のプロセッサと、

前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、

前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第 1 プロセッサが担当する制御情報を格納する共有メモリと、を含み、

前記第 1 プロセッサは、割り当てられている第 1 ローカルメモリにおいて制御情報を更新し、

前記第 1 プロセッサは、前記第 1 ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、

前記第 1 プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第 1 ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、

前記第 1 プロセッサは、前記共有メモリの負荷に基づいて、前記共有メモリにおける制御情報の少なくとも一部の情報の格納領域を、前記共有メモリから異なる種別のデバイスの記憶領域に変更することを決定し、

前記第 1 プロセッサは、前記他のデバイスの前記記憶領域における前記少なくとも一部の情報を、前記第 1 ローカルメモリにおける更新に同期して更新する、ストレージシステム。

【請求項 1 3】

前記少なくとも一部の情報は、前記第 1 プロセッサが担当する第 1 ボリュームのデータキャッシング制御情報を含み、

前記第 1 ローカルメモリにおける前記第 1 ボリュームのデータキャッシング制御情報の更新を前記共有メモリに反映しないことを決定する条件は、前記第 1 プロセッサの負荷が第 1 閾値以上であることを含み、

前記他のデバイスの前記記憶領域は、前記ストレージシステムにおいてボリュームを提供する不揮発性記憶装置の記憶領域であり、

前記少なくとも一部の情報の格納場所を、前記第 2 共有メモリに変更することを決定する条件は、前記第 1 プロセッサの負荷が前記第 1 閾値よりも小さい第 2 閾値よりも小さいことを含む、請求項 1 2 に記載のストレージシステム。

【請求項 1 4】

異なる制御情報を担当する複数のプロセッサと、

前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、

前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第 1 プロセッサが担当する制御情報を格納する共有メモリと、を含み、

前記第 1 プロセッサは、割り当てられている第 1 ローカルメモリにおいて制御情報を更新し、

前記第 1 プロセッサは、前記第 1 ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定し、

前記第 1 プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第 1 ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する、ストレージシステムであって、

前記第 1 ローカルメモリにおける前記制御情報及び前記共有メモリにおける前記制御情

10

20

30

40

50

報は、それぞれ、前記複数のボリュームにおける前記第1プロセッサが担当する第1ボリュームのデータキャッシング制御情報を含み、

前記第1プロセッサは、前記第1ボリュームのデータキャッシングによるアクセス性能の向上に基づき、前記第1ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定する、ストレージシステム。

【請求項15】

前記複数の不揮発性半導体メモリそれぞれは、フラッシュメモリである、請求項1に記載のストレージシステム。

【請求項16】

それぞれが、複数の不揮発性半導体メモリを含む複数の記憶ドライブに基づいて構成される複数のボリュームと、

前記複数のボリュームのデータを一時的に格納するキャッシュメモリと、

複数のプロセッサパッケージと、

共有メモリと、を含み、

前記複数のプロセッサパッケージそれぞれは、前記複数のボリュームに含まれる1つのボリュームに対する入出力を担当するプロセッサと、前記1つのボリュームのデータキャッシング制御情報を格納するローカルメモリと、を含み、

前記共有メモリは、前記複数のプロセッサパッケージそれぞれのローカルメモリに格納されたデータキャッシング制御情報を格納し、前記プロセッサによってアクセス可能であり、

第1プロセッサパッケージに含まれ、第1ボリュームを担当する第1プロセッサは、計算機からの前記第1ボリュームへのリードコマンドを受信すると、前記第1プロセッサパッケージに含まれる第1ローカルメモリに格納されたデータキャッシング制御情報を更新し、

前記第1プロセッサは、前記第1ボリュームの記憶領域を提供する記憶ドライブに、前記リードコマンドにおけるリードデータが格納されているか否かを判定し、

前記第1プロセッサは、前記第1ボリュームの記憶領域を提供する記憶ドライブに前記リードデータが格納されている場合に、前記リードデータを前記キャッシュメモリに格納することなく、かつ前記共有メモリに前記データキャッシング制御情報の更新を反映することなく、前記第1ボリュームの記憶領域を提供する記憶ドライブから読み出された前記リードデータを前記計算機に送信することを決定する、ストレージシステム。

【請求項17】

前記複数の不揮発性半導体メモリそれぞれは、フラッシュメモリである、請求項16に記載のストレージシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はストレージシステムに関し、特に、ストレージシステムの制御に関する。

【背景技術】

【0002】

国際公開第2010/131373号パンフレット(特許文献1)は、各ボリュームのI/O担当プロセッサが、共有メモリ上のデータキャッシング制御情報をローカルメモリへキャッシング(制御キャッシング)することで、ストレージシステムを高性能化する技術を開示している。

【0003】

プロセッサは、ローカルメモリの制御情報を更新する時、共有メモリの制御情報も同期して更新する。これにより、障害が起きたプロセッサから担当を引き継いだ他のプロセッサは、共有メモリから最新のデータキャッシング制御情報を取得することができ、キャッシュヒット率低下によるストレージシステムの性能低下、を抑止することができる。

【0004】

10

20

30

40

50

この他、ストレージシステムでは、不揮発性メディアからユーザデータをキャッシュメモリにキャッシングすることでストレージシステムを高性能化する、データキャッシングが広く利用されている。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】国際公開第2010/131373号パンフレット

【発明の概要】

【発明が解決しようとする課題】

【0006】

10

しかし、性能向上が目的である共有メモリにおける制御情報の更新が、アクセス対象である共有メモリとアクセスを制御するプロセッサのオーバーヘッドを増加させている。性能向上が目的であるデータキャッシングが、アクセス対象であるキャッシュメモリとアクセスを制御するプロセッサのオーバーヘッドを増加させている。

【課題を解決するための手段】

【0007】

本発明の一態様のストレージシステムは、異なる制御データを担当する複数のプロセッサと、前記複数のプロセッサのそれぞれに割り当てられており、割り当てられたプロセッサが担当する制御情報を格納するローカルメモリと、前記複数のプロセッサがアクセス可能であり、前記複数のプロセッサにおける第1プロセッサが担当する制御情報を格納する共有メモリと、を含む。前記第1プロセッサは、割り当てられている第1ローカルメモリにおいて制御情報を更新する。前記第1プロセッサは、前記第1ローカルメモリにおける前記制御情報の更新を前記共有メモリにおける前記制御情報に反映するか否かを決定する。前記第1プロセッサは、前記共有メモリにおける前記制御情報に反映することを決定した前記第1ローカルメモリにおける前記制御情報の更新を、前記共有メモリにおける前記制御情報に反映する。

20

【発明の効果】

【0008】

本発明の一態様は、ストレージシステムにおけるオーバーヘッドを低減しストレージシステムの性能を向上する。

30

【図面の簡単な説明】

【0009】

【図1】第1実施形態において、計算機システムの全体構成を模式的に示すブロック図である。

【図2】第1実施形態において、ストレージシステムのローカルメモリに格納されている情報を示す図である。

【図3】第1実施形態において、ストレージシステムの共有メモリに格納されている情報を示す図である。

【図4】第1実施形態において、管理計算機の構成を模式的に示す図である。

【図5】第1実施形態において、性能ブースト機能有効化テーブルの一例を示す図である

40

【図6】第1実施形態において、ボリューム毎性能ブースト機能有効化テーブルの一例を示す図である。

【図7】第1実施形態において、メディア種別テーブルの一例を示す図である。

【図8】第1実施形態において、RAIDレベルテーブルの一例を示す図である。

【図9】第1実施形態において、ボリューム毎ヒット率テーブルの一例を示す図である。

【図10】第1実施形態において、ヒット率閾値テーブルの一例を示す図である。

【図11】第1実施形態において、MP稼働率テーブルの一例を示す図である。

【図12】第1実施形態において、MP稼働率閾値テーブルの一例を示す図である。

【図13】第1実施形態において、CM稼働率テーブルの一例を示す図である。

50

【図14】第1実施形態において、CM稼働率閾値テーブルの一例を示す図である。

【図15】第1実施形態におけるホストからのリードコマンドの処理のフローチャートである。

【図16】第1実施形態におけるデータキャッシングに関するSM制御情報更新判定処理のフローチャートである。

【図17】第1実施形態におけるホストデータキャッシング判定処理のフローチャートである。

【図18A】第1実施形態において、ホストからのライトコマンドの処理のフローチャートの一部である。

【図18B】第1実施形態において、ホストからのライトコマンドの処理のフローチャートの他の一部である。

【図19】第1実施形態における管理計算機からの設定処理のフローチャートである。

【図20】第1実施形態において、管理計算機における設定メニュー画面の一例を示す図である。

【図21】第1実施形態におけるメディア種別テーブルの更新処理のフローチャートである。

【図22】第1実施形態におけるCMPK稼働率更新処理のフローチャートである。

【図23】第1実施形態におけるヒット率更新処理のフローチャートである。

【図24】第1実施形態におけるMP稼働率更新処理のフローチャートである。

【図25】第1実施形態におけるオーナー移動時SM更新処理のフローチャートである。

【図26】第2実施形態において、ローカルメモリに格納されている情報を示す図である。

【図27】第2実施形態において、共有メモリに格納されている情報を示す図である。

【図28】第2実施形態において、ダイナミックマッピングテーブルの一例を示す図である。

【図29】第2実施形態において、ページ毎モニタテーブルの一例を示す図である。

【図30】第2実施形態において、ページ毎モニタ差分テーブルの一例を示す図である。

【図31】第2実施形態におけるストレージ階層仮想化機能モニタ更新処理のフローチャートである。

【図32】第3実施形態において、計算機システム全体構成を模式的に示す図である。

【図33】第3実施形態において、非同期リモートコピーを説明する図である。

【図34】第3実施形態において、ローカルメモリに格納されている情報を示す図である。

【図35】第3実施形態において、共有メモリに格納されている情報を示す図である。

【図36】第3実施形態において、LM非同期リモートコピーシーケンス番号管理テーブルの一例を示す図である。

【図37】第3実施形態において、SM非同期リモートコピーシーケンス番号管理テーブルの一例を示す図である。

【図38】第3実施形態における非同期リモートコピーシーケンス番号更新処理のフローチャートである。

【図39】第3実施形態におけるMPK障害時非同期リモートコピーシーケンス番号回復処理のフローチャートである。

【図40】第4実施形態において、ローカルメモリに格納されている情報を示す図である。

【図41】第4実施形態において、共有メモリに格納されている情報を示す図である。

【図42】第4実施形態において、LMローカルコピー差分管理テーブルの一例を示す図である。

【図43】第4実施形態において、SMローカルコピー差分管理テーブルの一例を示す図である。

【図44】第4実施形態において、LMローカルコピー差分領域間引き動作管理テーブル

10

20

30

40

50

の一例を示す図である。

【図45】第4実施形態において、SMローカルコピー差分領域間引き動作管理テーブルの一例を示す図である。

【図46】第4実施形態における非同期ローカルコピー差分管理情報更新処理のフローチャートである。

【図47】第4実施形態におけるMPPK障害時ローカルコピー差分コピー処理のフローチャートである。

【図48】第4実施形態において、管理計算機における設定メニュー画面の一例を示す図である。

【図49】第5実施形態において、計算機システムの全体構成を模式的に示す図である。 10

【図50】第5実施形態において、ローカルメモリに格納されている情報を示す図である。

【図51】第5実施形態において、Xパス稼働率テーブルの一例を示す図である。

【図52】第5実施形態において、Xパス稼働率閾値テーブルの一例を示す図である。

【図53】第5実施形態におけるXパスを考慮したデータキャッシングに関するSM制御情報更新判定処理のフローチャートである。

【図54】第5実施形態におけるXパス稼働率更新処理のフローチャートである。

【図55】第6実施形態において、計算機システムの全体構成を模式的に示す図である。

【図56】第6実施形態において、ローカルメモリに格納されている情報を示す図である。 20

【図57】第6実施形態において、MP稼働率テーブルの一例を示す図である。

【図58】第6実施形態において、MP稼働率閾値テーブルの一例を示す図である。

【図59】第6実施形態において、共有メモリ領域管理テーブルの一例を示す図である。

【図60A】第6実施形態におけるデータキャッシングに関するSM制御情報更新判定処理のフローチャートの一部である。

【図60B】第6実施形態におけるデータキャッシングに関するSM制御情報更新判定処理のフローチャートの他の一部である。

【図61】第6実施形態におけるMP稼働率更新処理のフローチャートである。

【図62】第7実施形態において、ローカルメモリに格納されている情報を示す図である。 30

【図63】第7実施形態において、レスポンステーブルの一例を示す図である。

【図64】第7実施形態において、CM利用閾値テーブルの一例を示す図である。

【図65】第7実施形態におけるヒット率更新処理のフローチャートである。

【発明を実施するための形態】

【0010】

本発明は、ストレージシステムの性能を向上するための技術に関する。以下、添付図面を参照して本発明の実施形態を説明する。説明の明確化のため、以下の記載及び図面の詳細は、適宜、省略及び簡略化がなされており、必要に応じて重複説明は省略されている。本実施形態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。 40

【0011】

第1実施形態

本実施形態のストレージシステムは、それぞれが異なるボリュームの入出力(I/O)を担当するプロセッサを含む。各プロセッサは、ローカルメモリが割り当てられている。本実施形態のストレージシステムは、異なるボリュームを担当する複数のプロセッサがアクセス可能な共有メモリを有する。ローカルメモリ及び共有メモリは、典型的には揮発性半導体メモリである。

【0012】

プロセッサが担当するボリュームのデータキャッシング制御情報は、当該プロセッサのローカルメモリに格納されている(制御データキャッシング)。さらに、共有メモリは、 50

当該ボリュームのデータキャッシング制御情報を格納する。

【0013】

プロセッサは、ローカルメモリ上のキャッシング制御情報を参照及び更新して、担当ボリュームのデータキャッシング制御を行う。これにより、データキャッシング制御の処理を高速化することができる。

【0014】

上述のように、共有メモリは、異なるボリュームを担当する複数のプロセッサがアクセスすることができる。いずれかのボリュームの担当プロセッサについて障害が発生した場合、他のプロセッサがその担当を引き継ぎ、共有メモリから対応するデータキャッシング制御情報を自身のローカルメモリにロードする。上記他のプロセッサは、共有メモリから取得したデータキャッシング制御情報を使用して、引き継いだボリュームのデータキャッシングを制御する。

10

【0015】

本実施形態において、プロセッサは、予め定められている条件に従って、ローカルメモリにおけるキャッシング制御情報の更新を、共有メモリにおける制御情報に反映するか否かを決定する。ローカルメモリにおける制御情報の更新において必要な更新のみを共有メモリにおける制御情報に反映することで、プロセッサと共有メモリの通信におけるオーバーヘッドを低減し、ストレージシステムの性能を向上することができる。

【0016】

さらに、本実施形態のストレージシステムは、リードデータ及びライトデータのキャッシングの有無を、予め定められている条件に従って決定する。リードデータ及びライトデータを選択的にキャッシングすることで、キャッシュ領域を効率的に利用し、さらに、キャッシュメモリ及びデータキャッシングを行うプロセッサのオーバーヘッドを低減することで、ストレージシステムの性能を向上する。

20

【0017】

以下において、図1から図25を参照して本実施形態を具体的に説明する。図1は、本実施形態のストレージシステム10、データ処理及び演算を行うホスト計算機180及びストレージシステム10を管理する管理計算機20を含む、計算機システムの一例を示す。計算機システムは、複数のホスト計算機180を含むことができる。

【0018】

ストレージシステム10とホスト計算機180とは、データネットワーク190を介して互いに接続される。データネットワーク190は、例えば、SAN (Storage Area Network) である。データネットワーク190は、IPネットワーク若しくはその他のいかなる種類のデータ通信用のネットワークであってもよい。

30

【0019】

ストレージシステム10、ホスト計算機180及び管理計算機20は、管理ネットワーク(不図示)を介して互いに接続される。管理ネットワークは、例えば、IPネットワークである。管理ネットワークは、SAN若しくはその他のいかなる種類のネットワークであってもよい。データネットワーク190と管理ネットワークとが同一のネットワークであってもよい。

40

【0020】

ストレージシステム10は、複数の記憶ドライブ170を収容している。記憶ドライブ170は、不揮発性の磁気ディスクを有するハードディスクドライブ(HDD)、不揮発半導体メモリ(例えばフラッシュメモリ)を搭載したSolid State Drive(SSD)を含む。記憶ドライブ170は、ホスト計算機180から送られたデータ(ユーザデータ)を格納する。複数の記憶ドライブ170がRAID演算によるデータの冗長化を行うことで、1つの記憶ドライブ170に障害が発生した場合のデータ消失を防ぐことができる。

【0021】

ストレージシステム10は、ホスト計算機180と接続するためのフロントエンドパッ

50

ケージ (FEPK) 100、記憶ドライブ 170 と接続するためのバックエンドパッケージ (BEPK) 140、キャッシュメモリを搭載するキャッシュメモリパッケージ (CMPK) 130、内部処理を行うマイクロプロセッサを搭載するマイクロプロセッサパッケージ (MPPK) 120、及びそれらを接続する内部ネットワーク 150 を有する。図 1 に示すように、本例のストレージシステム 10 は、複数の FEPK 100、複数の BEPK 140、複数の CMPK 130、そして複数の MPPK 120 を含む。

【0022】

各 FEPK 100 は、ホスト計算機との接続を行うためのインタフェース 101、ストレージシステム 10 内のデータ転送を行うための転送回路 112 を基板上に有する。インタフェース 101 は複数のポートを含むことができ、各ポートがホスト計算機と接続することができる。インタフェース 101 は、ホスト計算機 180 とストレージシステム 10 との間の通信に用いられるプロトコル、例えば Fibre Channel Over Ethernet (FCoE) を、内部ネットワーク 150 で用いられるプロトコル、例えば PCI-Express に変換する。

10

【0023】

各 BEPK 140 は、ドライブ 170 と接続するためにインタフェース 141、ストレージシステム 10 内のデータ転送を行うための転送回路 142 を基板上に有する。インタフェース 141 は複数ポートを含むことができ、各ポートがドライブ 170 と接続することができる。インタフェース 141 は、記憶ドライブ 170 との通信に用いられるプロトコル、例えば FC を、内部ネットワーク 150 で用いられるプロトコルに変換する。

20

【0024】

各 CMPK 130 は、ホスト計算機 180 から読み書きされるユーザデータを一時的に格納するキャッシュメモリ (CM) 131 及び 1 又は複数の MPPK 120 が扱う制御情報を格納する共有メモリ (SM) 132 を基板上に有する。異なるボリュームを担当する複数の MPPK 120 (のマイクロプロセッサ) が、共有メモリ 132 にアクセスすることができる。MPPK 120 が扱うデータやプログラムは、ストレージシステム 10 内の不揮発性メモリ (不図示) 又は記憶ドライブ 170 からロードされる。関連づけられるキャッシュメモリ 131 と共有メモリ 132 とは、別の基板上 (パッケージ内) に実装されていてもよい。

30

【0025】

各 MPPK 120 は、1 以上のマイクロプロセッサ 121、ローカルメモリ (LM) 122 及びそれらを接続するバス 123 を有する。本例は複数のマイクロプロセッサ 121 が実装されている。マイクロプロセッサ 121 の数は 1 つでもよい。複数のマイクロプロセッサ 121 を一つのプロセッサと見ることができる。ローカルメモリ 122 は、マイクロプロセッサ 121 が実行するプログラム及びマイクロプロセッサ 121 が使用する制御情報を格納する。

【0026】

上述のように、一つの共有メモリ 132 は、MPPK 120 が扱う制御情報を格納する。MPPK 120 は、共有メモリ 132 から、自身が必要とする制御情報を自身のローカルメモリ 122 に格納する (制御キャッシング)。これにより、マイクロプロセッサ 121 による制御情報への高速アクセスが実現され、ストレージシステム 10 の性能を向上することができる。

40

【0027】

マイクロプロセッサ 121 は、ローカルメモリ 122 の制御情報を更新すると、その更新を、必要により、共有メモリ 132 上の制御情報に反映する。本実施形態の特徴の一つは、この更新の制御である。マイクロプロセッサ 121 は、予め定められた条件が満たされている場合に、ローカルメモリ 122 における制御情報の更新を共有メモリ 132 における制御情報に反映する。

【0028】

本構成例において、マイクロプロセッサ 121 には、ストレージシステム 10 がホスト

50

計算機 180 に提供するボリュームの担当が割り当てられる。マイクロプロセッサ 121 に割り当てられているローカルメモリ 122 及び共有メモリ 132 が、上記マイクロプロセッサが I/O を担当するボリュームのデータキャッシング制御情報を格納する。

【0029】

なお、本発明を適用できる制御情報は、共有メモリ 132 における制御情報を更新しない場合でも MP 障害時にホストデータロスにつながらない制御情報全般である。本実施形態におけるデータキャッシング制御情報以外の制御情報の例は、他の実施形態で説明する。本実施形態はマイクロプロセッサがボリュームを担当する例を記載しているが、担当マイクロプロセッサが割り当てられる対象はボリュームに限定されず、担当マイクロプロセッサは制御情報毎に存在してもよい。

10

【0030】

図 2 は、ローカルメモリ 122 内に格納されている情報を示すブロック図である。ローカルメモリ 122 は、性能ブースト機能有効化テーブル 210、ボリューム毎性能ブースト機能有効化テーブル 220、メディア種別テーブル 230、RAID レベルテーブル 240、ボリューム毎ヒット率テーブル 250、ヒット率閾値テーブル 260、マイクロプロセッサ (MP) 稼働率テーブル 270 を格納する。

【0031】

ローカルメモリ 122 は、さらに、マイクロプロセッサ (MP) 稼働率閾値テーブル 280、キャッシュメモリ (CM) 稼働率テーブル 290、キャッシュメモリ (CM) 稼働率閾値テーブル 300 を含む。例えば、マイクロプロセッサ 121 は、記憶ドライブ 170 その他のストレージシステム 10 内の不揮発性記憶領域からこれらテーブルの少なくとも一部を取得して、ローカルメモリ 122 に格納し、いくつかのテーブルを新たに作成する。これらテーブルの詳細な説明は後述する。

20

【0032】

ローカルメモリ 122 は、さらに、キャッシュディレクトリ 310 を格納している。図 3 は、共有メモリ 132 内のキャッシュディレクトリ 510 を示すブロック図である。マイクロプロセッサ 121 は、共有メモリ 132 からキャッシュディレクトリ 510 を自身のローカルメモリ 122 にキャッシュし、ローカルメモリ 122 上のキャッシュディレクトリ 310 の更新を、必要により、共有メモリ 132 のキャッシュディレクトリ 510 に反映する。キャッシュディレクトリ 510 は、キャッシュディレクトリ 310 のバックアップデータである。

30

【0033】

マイクロプロセッサ 121 は、ホスト計算機 180 からリードコマンドを受信すると、そのローカルメモリ 122 のキャッシュディレクトリ 310 を参照して、対象データがキャッシュメモリ 131 にキャッシュされているか (キャッシュヒットか) を決定する。このように、キャッシュディレクトリ 310 は、キャッシュメモリ 131 に格納されているキャッシュデータを検索するための情報を与える。

【0034】

キャッシュディレクトリ 310 は、参照テーブル GRPP (GRouP Pointer)、GRPT (GRouP Table) 1、GRPT 2、管理テーブルとしてのスロットコントロールブロック (SLCB) から構成されている。参照テーブル GRPP、GRPT 1、GRPT 2 は、キャッシュセグメントを検索するときにマイクロプロセッサ 121 により参照されるテーブルであって、ディレクトリ構造を有する。参照テーブル GRPP が最上位に位置し、参照テーブル GRPT 2 が最下位に位置している。上位のテーブルは、次のテーブルのポインタを含む。GRPT 2 には、SLCB へのポインタが含まれている。

40

【0035】

SLCB は、キャッシュ管理の最小単位であるセグメントに関する制御情報を管理するテーブルであって、キャッシュメモリ 131 上にリードコマンドの指定データが存在するか否か、キャッシュされているデータのキャッシュメモリ 131 上のアドレス情報等、が格納されている。

50

【 0 0 3 6 】

1つのスロットには、1又は複数のセグメントを関連付けることができる。1つのセグメントには、例えば、64KBのデータを格納することができる。キャッシュ管理の最小単位はセグメントであるが、キャッシュをスロット単位で管理してもよい。典型的には、ダーティデータ（物理ディスクへの書込み前の状態）と、クリーンデータ（物理ディスクへの書込み後の状態）の各状態間の遷移は、スロット単位で行われる。キャッシュ領域のリザーブ及びリリースは、スロット単位又はセグメント単位で行われる。

【 0 0 3 7 】

ホスト計算機180からリードアクセスがあった場合は、マイクロプロセッサ121は、それに含まれる論理ブロックアドレス（LBA）に基づいて、各階層テーブルを順番に辿ることにより、要求されたデータがキャッシュメモリ131に存在するか、存在する場合にはそのアドレスを知ることができる。

10

【 0 0 3 8 】

要求されたデータがキャッシュメモリ131に存在する場合、マイクロプロセッサ121は、そのデータをホスト計算機180に送信する。要求されたデータがキャッシュメモリ131に存在しない場合、マイクロプロセッサ121は、ホスト計算機180が要求しているデータを記憶ドライブ170から読み出して、キャッシュ領域上の1つまたは複数のスロットに格納する。ライトデータも同様にキャッシュされる。なお、キャッシュディレクトリを使用したキャッシュデータの検索は広く知られた技術であり、ここでの詳細な説明を省略する。

20

【 0 0 3 9 】

図4は、管理計算機20の構成を模式的に示すブロック図である。管理計算機20は、入力インタフェース22、入力デバイス28、表示インタフェース23、表示デバイス29、CPU26、通信インタフェース21、メモリ24、HDD25を備える。入力デバイス28の典型的な例は、キーボード及びポインタデバイスであるが、これらと異なるデバイスでもよい。表示デバイス29は、典型的には、液晶表示装置である。

【 0 0 4 0 】

管理者（ユーザ）は、表示デバイス29によって処理結果を視認しながら、入力デバイス28によって必要なデータを入力する。管理者が入力する情報及び表示デバイス29による表示例は、後述する。図1の計算機システムにおいて、管理システムは一つの管理計算機20で構成されているが、管理システムは、管理計算機20に加え、管理コンソールを含むことができる。管理コンソールは、入力デバイス及び表示デバイスを含み、ネットワークを介して管理計算機20に接続する。

30

【 0 0 4 1 】

管理者は管理コンソールから管理計算機20にアクセスし、管理計算機20に処理を指示するとともに、管理コンソールに管理計算機20の処理結果を取得、表示させる。管理システムは、それぞれが管理計算機20の機能の一部又は全部を備える複数の計算機を含むこともできる。CPU26は、メモリ24に格納されたプログラムを実行するプロセッサである。通信I/F21は、管理ネットワークとのインタフェースであって、システム管理のためにホスト計算機180及びストレージシステム10と、データや制御命令の送受信を行う。

40

【 0 0 4 2 】

図5は、性能ブースト機能有効化テーブル210の構成例を示している。性能ブースト機能有効化テーブル210は、性能ブースト機能有効フラグのカラム211を有する。性能ブースト機能有効フラグは、ストレージシステム10全体の性能ブースト機能がアクティブであるか否かを示す。このフラグが1である場合、ストレージシステム10全体の性能ブースト機能がアクティブである。

【 0 0 4 3 】

本実施形態において、性能ブースト機能は、ローカルメモリ122に格納された制御情報更新の共有メモリ132への反映（バックアップ）の制御及びデータキャッシング制御

50

の機能である。この機能については後述する。性能ブースト機能有効化テーブル 2 1 0 のデータは、例えば、管理者が管理計算機 2 0 から設定する。

【 0 0 4 4 】

図 6 は、ボリューム毎性能ブースト機能有効化テーブル 2 2 0 の構成例を示している。ボリューム毎性能ブースト機能有効化テーブル 2 2 0 は、ボリューム毎の性能ブースト機能を管理する。ボリューム毎性能ブースト機能有効化テーブル 2 2 0 は、論理ボリューム番号のカラム 2 2 1 及び性能ブースト機能有効化フラグのカラム 2 2 2 を有する。論理ボリューム番号は、論理ボリュームの識別子である。

【 0 0 4 5 】

性能ブースト機能有効化フラグが 1 である場合、そのボリュームの性能ブースト機能がアクティブであることを示す。システム全体及びボリュームの性能ブースト機能有効化フラグの双方が ON (1) である場合、そのボリュームの性能ブースト機能が有効化される。このように、ボリューム毎に性能ブースト機能を管理、制御することで、ボリューム特性に応じた制御が実現される。ボリューム毎性能ブースト機能有効化テーブル 2 2 0 のデータは、例えば、管理者が管理計算機 2 0 から設定する。

10

【 0 0 4 6 】

図 7 は、メディア種別テーブル 2 3 0 の構成例を示している。メディア種別テーブル 2 3 0 は、RAID グループのメディア種別を管理する。本実施形態において、1 又は複数の記憶ドライブ 1 7 0 が提供する記憶領域及びそのインタフェースを含む構成をメディアと呼ぶ。メディア種別テーブル 2 3 0 は、RAID グループ番号のカラム 2 3 1 及びメディア種別のカラム 2 3 2 を含む。

20

【 0 0 4 7 】

RAID グループ番号は、RAID グループを一意に識別する識別子である。なお、本明細書において、対象を識別する識別情報のために、識別子、名、ID 等の表現を用いることができ、これらは置換可能である。メディア種別テーブル 2 3 0 のデータは、例えば、管理者が管理計算機 2 0 から設定する。

【 0 0 4 8 】

図 8 は、RAID レベルテーブル 2 4 0 の構成例を示している。RAID レベルテーブル 2 4 0 は、RAID グループの RAID レベルを管理する。RAID グループ番号のカラム 2 4 1 及び RAID レベルのカラム 2 4 2 を有する。RAID レベルテーブル 2 4 0 のデータは、例えば、管理者が管理計算機 2 0 から設定する。

30

【 0 0 4 9 】

図 9 は、ボリューム毎ヒット率テーブル 2 5 0 の構成例を示している。ボリューム毎ヒット率テーブル 2 5 0 は、各ボリュームのキャッシュヒット率を管理する。ボリューム毎ヒット率テーブル 2 5 0 は、論理ボリューム番号のカラム 2 5 1、ヒット率のカラム 2 5 2、I/O 数のカラム 2 5 3、ヒット数のカラム 2 5 4、低ヒット率フラグのカラム 2 5 5 を有する。

【 0 0 5 0 】

I/O 数は、論理ボリュームに対して発行されたリードコマンド数である。ヒット数は、キャッシュヒットしたリードコマンド数である。低ヒット率フラグが 1 である場合、そのエントリのヒット率が、規定閾値未満であることを示す。マイクロプロセッサ 1 2 1 は、ボリュームへのリードアクセス及びキャッシュヒット数をカウントし、ボリューム毎ヒット率テーブル 2 5 0 の各フィールドのデータを更新する。

40

【 0 0 5 1 】

なお、マイクロプロセッサ 1 2 1 がヒット率をモニタする単位は、論理ボリュームより小さい単位でもよい。例えば、仮想ボリューム機能や階層化機能で使用されるページを単位としてもよい。

【 0 0 5 2 】

ヒット率の算出は、リードキャッシュのヒット率の他にライトキャッシュのヒット率を含んでもよい。ライトキャッシュのヒット率 [%] は、 $100 \times (1 - (\text{記憶ドライブへのラ$

50

イト処理回数/ライトコマンドI/O数)で表すことができる。マイクロプロセッサ121は、リードキャッシュのヒット率とライトキャッシュのヒット率を個別に管理してもよい。例えば、マイクロプロセッサ121は、後述のリードキャッシング制御及びライトキャッシング制御において、それぞれのヒット率を参照する。

【0053】

図10は、ヒット率閾値テーブル260の構成例を示している。ヒット率閾値テーブル260は、ヒット率閾値のカラム261を有する。ヒット率がここに登録されている閾値以下である場合、ボリューム毎ヒット率テーブル250におけるそのエントリの低ヒット率フラグが1(ONフラグ)に設定される。ヒット率閾値は、例えば、管理者が管理計算機20から設定する。

10

【0054】

図11は、マイクロプロセッサ121の稼働率を管理するMP稼働率テーブル270の構成例を示している。MP稼働率は、単位時間内のマイクロプロセッサ121の処理時間であり、マイクロプロセッサの負荷を表す。MP稼働率テーブル270は、マイクロプロセッサ番号のカラム271、稼働率のカラム272、過負荷判定フラグのカラム273、稼働時間のカラム274を有する。マイクロプロセッサ番号は、ストレージシステム10内で一意にマイクロプロセッサを識別する。

【0055】

各マイクロプロセッサ121は、自身の稼働状況を監視し、稼働率及び稼働時間の値を、自身のエントリの稼働率のカラム272及び稼働時間のカラム274のフィールドに格納する。稼働時間は、単位時間(本例で1秒)当たりの稼働時間である。稼働率は、単位時間で稼働時間を割った値である。マイクロプロセッサ121は、自己の稼働率を規定の閾値と比較し、その閾値以上である場合に、自己エントリの過負荷判定フラグのフィールドの値を1(ONフラグ)に設定する。

20

【0056】

図12は、上記閾値を格納するカラム281を有する、MP稼働率閾値テーブル280の構成例を示している。本例において、MP稼働率閾値は、全てのマイクロプロセッサに共通であるが、異なる閾値を使用してもよい。

【0057】

図13は、キャッシュメモリの稼働率を管理する、CM稼働率テーブル290の構成例を示している。CM稼働率は、単位時間内のキャッシュメモリ131へのアクセス時間である。CM稼働率テーブル290は、CMPK番号のカラム291、稼働率のカラム292、過負荷判定フラグのカラム293を有する。CMPK番号は、ストレージシステム10内のCMPKの識別子である。

30

【0058】

マイクロプロセッサ121は、CMPK130上のコントローラから、その稼働率の値を取得し、稼働率のカラム292の該当フィールドにそれを格納する。マイクロプロセッサ121は、取得した稼働率の値を規定閾値と比較し、稼働率の値が閾値以上である場合に、そのエントリの過負荷判定フラグのフィールドに1(ONフラグ)を設定する。

【0059】

図14は、上記閾値を格納するCM稼働率閾値テーブル300の構成例を示している。本例において、CM稼働率閾値は、全てのCMPKに共通であるが、異なる閾値を使用してもよい。

40

【0060】

図15のフローチャートを参照して、ストレージシステム10がホスト計算機180から受けたリードコマンドに対して行う処理を説明する。ホスト計算機180からのリードコマンドを受けた(S101)マイクロプロセッサ121は、リードコマンドが示す論理ボリューム(LDEVとも呼ぶ)に、自身がアクセス権を有するか判定する(S102)。アクセス権を有していない場合(S102:NO)、そのマイクロプロセッサ121は、アクセス権を有するMPK120にリードコマンドを転送する(S103)。

50

【 0 0 6 1 】

マイクロプロセッサ 1 2 1 がアクセス権を有する場合 (S 1 0 2 : Y E S)、そのマイクロプロセッサ 1 2 1 は、同一 M P P K 1 2 0 上のローカルメモリ 1 2 2 内で、キャッシュディレクトリ 3 1 0 を検索する (S 1 0 4)。リードコマンドが指定するアドレス (データ) が見つかった場合 (S 1 0 5 : Y E S)、マイクロプロセッサ 1 2 1 は、キャッシュディレクトリ 3 1 0 の情報に従って、キャッシュメモリ 1 3 1 からリードデータを読み出し、ホスト計算機 1 8 0 に送信する (S 1 0 6)。

【 0 0 6 2 】

リードコマンドが指定するアドレス (データ) が見つからなかった (キャッシュミス) 場合 (S 1 0 5 : N O)、マイクロプロセッサ 1 2 1 は、ローカルメモリ 1 2 2 の未キャッシュフラグを確認する (S 1 0 7)。未キャッシュフラグは、共有メモリ 1 3 2 のキャッシュディレクトリ 5 1 0 の全てのデータが、ローカルメモリ 1 2 2 にキャッシュされているかを示すフラグであり、ローカルメモリ 1 2 2 内に格納されている。一部のデータが読み込まれていない場合、その値は O N である。例えば、障害フェイルオーバ直後で共有メモリ 1 3 2 からローカルメモリ 1 2 2 へ制御情報が読み込まれていない場合、未キャッシュフラグは O N である。

10

【 0 0 6 3 】

未キャッシュフラグが O N である場合 (S 1 0 7 : Y E S)、共有メモリ 1 3 2 のキャッシュディレクトリ 5 1 0 の一部データがキャッシュされていない。マイクロプロセッサ 1 2 1 は、C M P K 1 3 0 のコントローラを介して、共有メモリ 1 3 2 からローカルメモリ 1 2 2 へキャッシュディレクトリ (制御情報) を転送する (S 1 0 8)。

20

【 0 0 6 4 】

マイクロプロセッサ 1 2 1 は、ローカルメモリ 1 2 2 内で、キャッシュディレクトリ 3 1 0 を検索する (S 1 0 9)。リードコマンドが指定するデータが見つかった場合 (S 1 1 0 : Y E S)、マイクロプロセッサ 1 2 1 は、キャッシュディレクトリ 3 1 0 の情報に従って、キャッシュメモリ 1 3 1 からリードデータを読み出し、ホスト計算機 1 8 0 に送信する (S 1 1 1)。

【 0 0 6 5 】

キャッシュミスの場合 (S 1 1 0 : N O) 又は未キャッシュフラグが O F F の場合 (S 1 0 7 : N O)、マイクロプロセッサ 1 2 1 は、キャッシュメモリ 1 3 1 のセグメントをリードデータのためのスロットとして確保し、さらに、ローカルメモリ 1 2 2 のキャッシュディレクトリ 3 1 0 を更新する (S 1 1 2)。

30

【 0 0 6 6 】

マイクロプロセッサ 1 2 1 は、データキャッシングに関する制御情報であるキャッシュディレクトリ 3 1 0 の更新を、共有メモリ 1 3 2 のデータ 5 1 0 に反映するか否かを判定する (S 1 1 3)。この判定の具体的な方法については後に詳述する。共有メモリ 1 3 2 の制御情報の更新を行うと判定した場合 (S 1 1 4 : Y E S)、マイクロプロセッサ 1 2 1 は、共有メモリ 1 3 2 のキャッシュディレクトリ 5 1 0 を更新して (S 1 1 5)、次のステップ S 1 1 6 に進む。

【 0 0 6 7 】

共有メモリ 1 3 2 の制御情報の更新を行わないと判定した場合 (S 1 1 4 : N O)、マイクロプロセッサ 1 2 1 は、共有メモリ 1 3 2 の制御情報を更新することなく、ステップ S 1 1 6 に進む。ステップ S 1 1 6 において、マイクロプロセッサ 1 2 1 は、リードデータ (ホストデータ) をキャッシングするか否かを判定する。この判定方法については後述する。

40

【 0 0 6 8 】

リードデータをキャッシュメモリ 1 3 1 に格納してからホスト計算機 1 8 0 に送信すると判定した場合 (S 1 1 7 : Y E S)、マイクロプロセッサ 1 2 1 は、B E P K 1 4 0 及び C M P K 1 3 0 により、記憶ドライブ 1 7 0 (永続メディア) からリードデータを読み出し、キャッシュメモリ 1 3 1 上の確保したセグメントに格納する。その後、マイクロプロ

50

ロセッサ 121 は、そのキャッシュデータを、CMPK130 及び FEPK100 により、ホスト計算機 180 に送信する (S118)。

【0069】

リードデータをキャッシュすることなくホスト計算機 180 に送信すると判定した場合 (S117: NO)、マイクロプロセッサ 121 は、BEPK140 及び FEPK100 により、ドライブ 170 (永続メディア) から読みだしたリードデータを、CMPK130 を介することなくホスト計算機 180 に転送する (S119)。キャッシュメモリ 131 を経由しないデータのセグメントは、キャッシュメモリ 131 を経由するデータのセグメントと比較して、他のデータ用に再利用されやすくなるように管理するとより効率的である。例えば、LRU キュー管理をしている場合、LRU に接続することが考えられる。

10

【0070】

図 16 を参照して、図 15 のフローチャートにおける、共有メモリ 132 内のデータキャッシング制御情報の更新についての判定 (S113) を説明する。マイクロプロセッサ 121 は、このステップ S113 を開始し (S121)、リードコマンドの指定する論理ボリュームの性能ブースト機能が ON であるか否かを、性能ブースト機能有効化テーブル 210 及びボリューム毎性能ブースト機能有効化テーブル 220 を参照して判定する (S122)。一方のテーブルが、性能ブースト機能が OFF であることを示す場合、当該ボリュームの性能ブースト機能は OFF である。

【0071】

当該論理ボリュームの性能ブースト機能が ON ではない場合 (S122: NO)、マイクロプロセッサ 121 は、共有メモリ 132 の制御情報 (キャッシュディレクトリ) を更新することを決定する (S128)。当該論理ボリュームの性能ブースト機能が ON である場合 (S122: YES)、マイクロプロセッサ 121 は、次に、指定データが格納されている RAID グループのメディア種別が SSD であるか否かを、RAID グループ番号をキーとしてメディア種別テーブル 230 を参照し、判定する (S123)。

20

【0072】

マイクロプロセッサ 121 は、ローカルメモリ 122 内に、各ボリュームの構成管理情報を有しており、各ボリュームの各領域がいずれの RAID グループに属するかをその情報を参照して知ることができる。

【0073】

メディア種別が SSD である場合 (S123: YES)、マイクロプロセッサ 121 は、共有メモリ 132 の制御情報 (キャッシュディレクトリ) を更新しないことを決定する (S127)。そのメディア種別が SSD ではない場合 (S123: NO)、マイクロプロセッサ 121 は、次に、指定データが格納されている論理ボリュームの低ヒット率フラグが ON であるか否かを、論理ボリューム番号をキーとしてボリューム毎ヒット率テーブル 250 を参照し、判定する (S124)。

30

【0074】

その低ヒット率フラグが ON である場合 (S124: YES)、マイクロプロセッサ 121 は、共有メモリ 132 の制御情報 (キャッシュディレクトリ) を更新しないことを決定する (S127)。低ヒット率フラグが OFF である場合 (S124: NO)、マイクロプロセッサ 121 は、次に、自身の過負荷フラグが ON であるか否かを、マイクロプロセッサ番号をキーとして MP 稼働率テーブル 270 を参照し、判定する (S125)。

40

【0075】

過負荷フラグが ON である場合 (S125: YES)、マイクロプロセッサ 121 は、共有メモリ 132 の制御情報 (キャッシュディレクトリ) を更新しないことを決定する (S127)。過負荷フラグが OFF である場合 (S125: NO)、マイクロプロセッサ 121 は、次に、アクセス先の CMPK130 の過負荷フラグが ON であるか否かを、CMPK 番号をキーとして CM 稼働率テーブル 290 を参照し、判定する (S126)。

【0076】

過負荷フラグが ON である場合 (S126: YES)、マイクロプロセッサ 121 は、

50

共有メモリ132の制御情報(キャッシュディレクトリ)を更新しないことを決定する(S127)。過負荷フラグがOFFである場合(S126:NO)、マイクロプロセッサ121は、共有メモリ132の制御情報(キャッシュディレクトリ)を更新することを決定する(S128)。

【0077】

このように、規定条件を満たす場合、マイクロプロセッサ121は、ローカルメモリ122でのキャッシュディレクトリ310の更新を、共有メモリ132のキャッシュディレクトリ510に反映しないことを決定する。これにより、マイクロプロセッサ121及びCMPK130の負荷を低減し、システムのスループットを向上することができる。

【0078】

ローカルメモリの制御情報(本例ではキャッシュディレクトリ)の更新を共有メモリ132に反映していないことは、その制御情報の担当MPPK120に障害が発生した場合に問題となる。通常動作において、マイクロプロセッサ121は、自身のローカルメモリ122を参照するため、更新された最新の制御情報を参照することができる。一方、担当MPPK120に障害が発生した場合、他のMPPK120が担当を引き継ぐ(フェイルオーバー)。

【0079】

障害発生したMPPK120のローカルメモリ122上のデータは消失するため、引き継いだMPPK120(のマイクロプロセッサ121)は、共有メモリ132に格納されている古い制御情報しか得ることができない。そのため、共有メモリ132に格納されており、更新(共有メモリ132へのバックアップ)を省略することができるデータは、MPPK120の障害時にユーザデータロストにつながらない制御情報である。

【0080】

上記好ましい構成は、MPPK120で障害が発生した場合に影響が小さい共有メモリ132での更新を省略する。具体的には、キャッシュミスによりリードデータが読みだされる記憶ドライブ170がSSDである場合(S123:YES)、マイクロプロセッサ121は、共有メモリ132での更新を行わないことを決定する(S127)。

【0081】

MPPK120の障害により、SSDから読みだされた上記データがキャッシュされていることを示す情報が失われる。しかし、SSDは、他のメディア種別のドライブ170よりもアクセス性能が高く、失われた制御情報に起因するキャッシュミスの影響は小さく、MPPK120及びCMPK130のオーバーヘッド低減によるシステム性能向上効果の方が大きい。

【0082】

本構成においては、共有メモリ132での更新を省略するメディア種別はSSDであるが、この種別は、システム設計に依存する。システムに実装されるメディア(ドライブ)の種別は、SSD及びHDDに限らず、これらに加え又はこれらに代えて異なる種別のドライブを実装することができる。実装されている複数のメディア種別において、共有メモリ132での更新省略の条件を満たす種別は、設計に従って選択される。最もアクセス性能が高い種別を含む、1又は複数の他の種別よりもアクセス性能が高い種別が選択される。

【0083】

本構成において、リードコマンド指定データを格納する論理ボリュームのキャッシュヒット率が低い場合(S124:YES)、マイクロプロセッサ121は、共有メモリ132での更新を行わないことを決定する(S127)。ヒット率が低いボリュームのデータのキャッシュ制御情報が失われても、そのボリュームのアクセス性能への影響は小さく、MPPK120及びCMPK130のオーバーヘッド低減によるシステム性能向上効果の方が大きい。

【0084】

本構成は、さらに、MPPK120及びCMPK130の現状負荷に基づいて、共有メ

10

20

30

40

50

メモリ132での更新の有無を決定する(S125、S126)。MPPK120又はCMPK130の負荷が高い場合、共有メモリ132での更新を省略することによる性能向上の効果が大きい。

【0085】

このように、本構成は、対象ボリュームの性能ブースト機能がONであり、上記4つの条件のいずれかが満たされる場合、共有メモリ132でのキャッシュ制御情報の更新を省略する。マイクロプロセッサ121は、これらと異なる条件に基づき共有メモリ132での更新の有無を決定してもよい。マイクロプロセッサ121は、上記4条件のうち複数の条件が満たされることを、共有メモリ132での制御情報更新省略の条件としてもよい。

【0086】

図17は、図15のフローチャートにおける、ホストデータ(リードデータ)キャッシングについての判定(S116)のフローチャートを示している。本ステップのフローチャートは、図16に示すフローチャートと略同様である。従って、主にこれと異なる点について具体的に説明する。

【0087】

図17において、ステップS132からステップS136は、それぞれ、図15のフローチャートにおけるステップ122からステップS126と同様である。ステップ137において、マイクロプロセッサ121は、記憶ドライブ170から読みだしたホストデータ(リードデータ)をキャッシュメモリ131に格納することなく、ホスト計算機180に送信することを決定する。一方、ステップS138において、マイクロプロセッサ121は、記憶ドライブ170から読みだしたホストデータをキャッシュメモリ131に格納する(キャッシュする)ことを決定する。

【0088】

本例において、リードデータをキャッシュするか否かの判定条件は、キャッシュ制御情報の更新を共有メモリ132で行うか否かの判定条件と同一である。このように、リードデータキャッシングを制御することで、MPPK120及びCMPK130のオーバヘッド低減によりシステム性能を向上することができる。キャッシュ制御の判定条件と制御情報更新制御の判定条件とは、異なってもよい。

【0089】

次に、ホスト計算機180から受信したライトコマンドに対する処理を、図18A及び図18Bに示すフローチャートを参照して説明する。マイクロプロセッサ121は、ホスト計算機180からライトコマンドを受け(S141)、その指定アドレスのボリューム(LDEV)に、アクセス権を有するか否かを判定する(S142)。

【0090】

そのマイクロプロセッサ121がアクセス権を有しない場合(S142:NO)、マイクロプロセッサ121は、他の担当MPPK120にライトコマンドを転送する(S143)。そのマイクロプロセッサ121がアクセス権を有している場合(S142:YES)、マイクロプロセッサ121は、同一基板上のローカルメモリ122内でキャッシュディレクトリ310を検索する(S144)。

【0091】

ライトコマンドが指定するアドレスが見つかった場合(S145:YES)、マイクロプロセッサ121は、キャッシュディレクトリ310の情報に従って、キャッシュメモリ131にライトデータを書き込み、ホスト計算機180にコマンド完了を通知する(S146)。

【0092】

ライトコマンドが指定するアドレスが見つからなかった(キャッシュミス)場合(S145:NO)、マイクロプロセッサ121は、ローカルメモリ122への未キャッシュフラグを確認する(S147)。未キャッシュフラグがONである場合(S147:YES)、マイクロプロセッサ121は、CMPK130のコントローラを介して、共有メモリ132からローカルメモリ122へキャッシュディレクトリ(制御情報)を転送する(S

10

20

30

40

50

148)。

【0093】

マイクロプロセッサ121は、ローカルメモリ122内で、キャッシュディレクトリ310を検索する(S149)。ライトコマンドが指定するアドレスが見つかった場合(S150: YES)、マイクロプロセッサ121は、キャッシュディレクトリ310の情報に従って、キャッシュメモリ131にライトデータを書き込み、ホスト計算機180にコマンド完了を通知する(S151)。

【0094】

キャッシュミスの場合(S150: NO)又は未キャッシュフラグがOFFの場合(S147: NO)、マイクロプロセッサ121は、キャッシュメモリ131のセグメントをライトデータのためのスロットとして確保し、さらに、ローカルメモリ122のキャッシュディレクトリ310を更新する(S152)。

【0095】

マイクロプロセッサ121は、データキャッシングに関する制御情報であるキャッシュディレクトリ310の更新を、共有メモリ132のデータ510に反映するか否かを判定する(S153)。この判定の具体的な方法は、図15を参照して説明した方法と同様である。マイクロプロセッサ121は、さらに、ライトデータ(ホストデータ)をキャッシングするか否かを判定する(S154)。この判定方法は、図16を参照して説明した方法と同様である。

【0096】

マイクロプロセッサ121がライトデータをキャッシュすると判定した場合(S155: YES)、マイクロプロセッサ121は、キャッシュメモリ131に新たに確保した領域にライトデータを書き込み、ホスト計算機180にコマンド完了を通知する(S156)。マイクロプロセッサ121は、ステップS153での判定結果に関わらず、ローカルメモリ122におけるキャッシュディレクトリ310の更新に同期して、共有メモリ132におけるキャッシュディレクトリ510を更新する。

【0097】

マイクロプロセッサ121がライトデータをキャッシュしないと判定した場合(S155: NO)、マイクロプロセッサ121は、ステップS153における判定結果に基づいて、共有メモリ132における制御情報の更新を行う又は省略する。マイクロプロセッサ121が、共有メモリ132におけるキャッシュ制御情報(キャッシュディレクトリ510)を更新すると判定した場合(S157: YES)、マイクロプロセッサ121は、ローカルメモリ122のキャッシュディレクトリ310の更新を、共有メモリ132におけるキャッシュディレクトリ510に反映し(S158)、次のステップS159に進む。

【0098】

マイクロプロセッサ121が、共有メモリ132におけるキャッシュ制御情報を更新しないと判定した場合(S157: NO)、マイクロプロセッサ121は、書き込み先のRAIDレベルを、RAIDレベルテーブル240を参照して特定する(S159)。そのRAIDレベルが1である場合(S159: YES)、マイクロプロセッサ121は、キャッシュメモリ131にライトデータを格納することなく、BEPK140により記憶ドライブ170にデータを書き込み、ホスト計算機180にコマンド完了を通知する(S160)。

【0099】

そのRAIDレベルが1と異なる場合(S159: NO)、マイクロプロセッサ121は、パリティを生成し、キャッシュメモリ131にライトデータを格納することなく、BEPK140により記憶ドライブ170にパリティ及びライトデータを書き込む。さらに、マイクロプロセッサ121はホスト計算機180にコマンド完了を通知する(S161)。

【0100】

このように、本例において、ライトコマンドのハンドリングにおいては、共有メモリ1

10

20

30

40

50

32におけるキャッシュディレクトリ510の更新を省略するためには、キャッシュメモリ131へのライトデータの格納が省略されることが必要である。キャッシュされたライトデータのデステージ(ドライブ170への書き込み)前にそのキャッシュ制御情報が失われると、キャッシュメモリ131でそのライトデータを特定することができないからである。

【0101】

上述のように、本例において、ステップS154におけるライトデータをキャッシュするか否かの判定条件は、図15におけるステップS116の判定条件と同一である。また、ステップS153におけるキャッシュ制御情報の更新を共有メモリ132で行うか否かの判定条件は、図15におけるステップS113の判定条件と同一である。これらは異な

10

【0102】

このように、ライトデータのキャッシング及びキャッシュ制御情報の更新を制御することによって、MPPK120及びCMPK130のオーバヘッドを低減し、ストレージシステム10の性能を向上することができる。

【0103】

次に、図19のフローチャートを参照して、管理計算機20からの設定処理を説明する。管理計算機20は、その上で実行される管理プログラムに従って動作する。したがって、管理計算機20を主語とする記載は、管理プログラム又はCPU26を主語とすることができる。管理計算機20は設定処理を開始し(S171)、設定データ入力のためのメ

20

【0104】

全ての必要なデータが入力されると(S174: YES)、管理計算機20は、保存ボタンの選択に 응답して設定データを保存する(S175)。設定データは、ストレージシステム10からの要求に応じて、管理計算機20からストレージシステム10に送信される。管理者は、キャンセルボタンを選択することで、入力をやり直すことができる。

【0105】

図20は、メニュー画面の一例2000を示している。メニュー画面2000は、性能ブースト機能設定エリア2001及びボリューム毎性能ブースト機能設定エリア2004

30

【0106】

管理者は、性能ブースト機能設定エリア2001における“ENABLE”又は“DISABLE”の一方を入力デバイス28で選択することで、ストレージシステム10の性能ブースト機能(上記制御情報の更新制御及びユーザデータのキャッシング制御の機能)をイネーブル又はディセーブルすることができる。この設定が、性能ブースト機能有効化テーブル210に反映される。これがディセーブルされると、ストレージシステム10の全ての性能ブースト機能が使用されない。

【0107】

ボリューム毎性能ブースト機能設定エリア2004は、論理ボリューム番号のカラム2005及び性能ブースト機能設定カラム2006を含む。管理者は、ボリューム毎性能ブースト機能設定エリア2004において、各論理ボリュームの性能ブースト機能のイネーブル/ディセーブルを入力デバイス28で選択することができる。

40

【0108】

この設定が、ボリューム毎性能ブースト機能有効化テーブル220に反映される。システムの性能ブースト機能がイネーブルされており、かつ、ボリュームの性能ブースト機能がイネーブルされているボリュームに対して、本実施形態の性能ブースト機能が使用される。

【0109】

図20は、性能ブースト機能の設定画面を例示しているが、この他、管理計算機20は

50

、例えば、判定条件に含まれる閾値の設定画面を表示し、管理者によって入力された設定データをストレージシステム10に送信する。典型的には、ストレージシステム10は、管理者により設定可能な項目のデフォルト値を予め有しており、管理者により設定された項目のデータを、入力データにより更新する。

【0110】

次に、図21から図24を参照して、ストレージシステム10内のテーブル更新を説明する。図21は、メディア種別テーブル230の更新のフローチャートである。RAIDグループが増減されると(S201)、BEPK140が、その情報をいずれかのマイクロプロセッサ121に送信する。更新情報を受信したマイクロプロセッサ121は、ローカルメモリ122のメディア種別テーブル230及びRAIDレベルテーブル240を
10

【0111】

更新すると共に、不揮発性記憶領域のこれらテーブルを更新し(S202)、それを他のMPPK120に通知する。管理理計算機20がMPPK120に情報を渡してもよい。

【0112】

図22を参照して、CM稼働率テーブル290の更新を説明する。MPPK120の任意のマイクロプロセッサ121がこの処理を行う。典型的には、定期的(例えば1秒毎)にこの処理が行われる。マイクロプロセッサ121は、アクセス先のCMPK130から稼働率の情報を取得する(S212)。具体的には、マイクロプロセッサ121は、CMPK130内のコントローラ(不図示)に、CMPK130の稼働率(CM稼働率)を示す値を要求し、それをCMPK130内のコントローラから取得する。
20

【0113】

稼働率が閾値以上である場合(S214: YES)、マイクロプロセッサ121は、CM稼働率テーブル290における対応エントリの稼働率カラム292のフィールドを更新する(S213)。さらに、マイクロプロセッサ121は、更新した稼働率の値が、CM稼働率閾値テーブル300の閾値以上であるか判定する(S214)。
30

【0114】

稼働率が閾値以上である場合(S214: YES)、マイクロプロセッサ121は、CM稼働率テーブル290における、該当エントリの過負荷フラグを1(ON)に設定する(S215)。稼働率が閾値未満である場合(S214: NO)、マイクロプロセッサ121は、該当エントリの過負荷フラグを0(OFF)に設定する(S216)。マイクロプロセッサ121は、アクセスする全てのCMPK130について、ステップS212からステップS216を行う(S217)。
40

【0115】

ヒット率が閾値以下である場合(S224: YES)、マイクロプロセッサ121は、当該エントリの低ヒット率フラグを1(ON)に設定する(S225)。一方、ヒット率が閾値より大きい場合(S224: NO)、マイクロプロセッサ121は、当該エントリの低ヒット率フラグを0(OFF)に設定する(S226)。マイクロプロセッサ121
50

【0116】

ヒット率が閾値以下である場合(S224: YES)、マイクロプロセッサ121は、当該エントリの低ヒット率フラグを1(ON)に設定する(S225)。一方、ヒット率が閾値より大きい場合(S224: NO)、マイクロプロセッサ121は、当該エントリの低ヒット率フラグを0(OFF)に設定する(S226)。マイクロプロセッサ121

は、担当する全ての論理ボリュームについて、ステップ S 2 2 2 からステップ S 2 2 6 を行う (S 2 2 7)。

【 0 1 1 7 】

図 2 4 を参照して、MP 稼働率テーブル 2 7 0 の更新を説明する。各マイクロプロセッサ 1 2 1 がこの処理を行う。典型的には、定期的 (例えば 1 秒毎) にこの処理が行われる。マイクロプロセッサ 1 2 1 は、自身の単位時間 (本例で 1 秒) 当たりの稼働時間を監視し、その値をローカルメモリ 1 2 2 内に格納する。マイクロプロセッサ 1 2 1 は、ローカルメモリ 1 2 2 からその値を取得する (S 2 3 2)。

【 0 1 1 8 】

マイクロプロセッサ 1 2 1 は、取得した値を使用して、該当エントリの稼働率のフィールドを更新し、稼働時間のフィールドに 0 をセットする (S 2 3 3)。さらに、マイクロプロセッサ 1 2 1 は、更新された稼働率と MP 稼働率閾値テーブル 2 8 0 の閾値とを比較する (S 2 3 4)。稼働率が閾値以上である場合 (S 2 3 4 : Y E S)、マイクロプロセッサ 1 2 1 は、当該エントリの過負荷フラグを 1 (O N) に設定する (S 2 3 5)。稼働率が閾値未満である場合 (S 2 3 4 : N O)、マイクロプロセッサ 1 2 1 は、当該エントリの過負荷フラグを 0 (O F F) に設定する (S 2 3 6)。

【 0 1 1 9 】

図 2 5 を参照して、論理ボリュームのオーナー権の現在 M P P K 1 2 0 から他の M P P K 1 2 0 への移動を説明する。オーナー権が移動する前に、現在 M P P K 1 2 0 は、ローカルメモリ 1 2 2 に格納するキャッシュディレクトリ 3 1 0 における未反映部分を、共有メモリ 1 3 2 に反映する。これにより、次の M P P K 1 2 0 が、最新のキャッシュディレクトリを使用してキャッシュ制御することができ、キャッシュヒット率を高めることができる。

【 0 1 2 0 】

現在オーナー M P P K のマイクロプロセッサ 1 2 1 は、キャッシュディレクトリ 3 1 0 において検索する対象を、オーナー権を移動する論理ボリュームの論理アドレスの 0 番に設定する (S 2 4 2)。マイクロプロセッサ 1 2 1 は、そのアドレスを、キャッシュディレクトリ 3 1 0 で検索する (S 2 4 3)。

【 0 1 2 1 】

そのアドレスが、共有メモリ未反映フラグが O N に設定されているディレクトリに存在する場合 (S 2 4 4 : Y E S)、マイクロプロセッサ 1 2 1 は、その共有メモリ 1 3 2 における当該ディレクトリを更新し (S 2 4 5)、ステップ S 2 4 6 に進む。共有メモリ未反映フラグは、対象ディレクトリの更新が共有メモリ 1 3 2 に反映済みであるか否かを示すフラグであり、それが O N である場合、対象ディレクトリの更新が共有メモリ 1 3 2 に未反映であることを示す。

【 0 1 2 2 】

上記アドレスが、共有メモリ未反映フラグが O F F に設定されているディレクトリに存在する場合 (S 2 4 4 : N O)、マイクロプロセッサ 1 2 1 は、共有メモリ 1 3 2 上のそのディレクトリを更新することなく、ステップ S 2 4 6 に進む。

【 0 1 2 3 】

ステップ S 2 4 6 において、マイクロプロセッサ 1 2 1 は、当該ボリュームについてのキャッシュディレクトリ 3 1 0 の探索が終了したが否かを判定する。全てのアドレスの探索を終了している場合 (S 2 4 6 : Y E S)、マイクロプロセッサ 1 2 1 はこの処理を終了する。未探索のアドレスが残っている場合 (S 2 4 6 : N O)、マイクロプロセッサ 1 2 1 は対象アドレスを次の論理アドレスに変更し (S 2 4 7)、ステップ S 2 4 3 からステップ S 2 4 6 を繰り返す。

【 0 1 2 4 】

第 2 実施形態

本実施形態は、ストレージ階層仮想化機能を有するストレージシステム 1 0 を説明する。本実施形態のストレージシステム 1 0 は、複数のプールボリューム (実ボリューム) を

10

20

30

40

50

含むプールを構築する。プールは、ストレージシステム 10 内の性能の異なる複数のメディアを含み、アクセス性能によって複数の階層に階層化される。各階層は、1 又は複数のプールボリュームで構成されている。

【0125】

ストレージシステム 10 は、そのプールから構築した仮想ボリュームをホスト計算機 180 に提供する。ストレージシステム 10 は、プールを、特定容量のページ単位で管理する。各プールボリュームは複数ページに分割され、各ページにデータが格納される。ストレージシステム 10 は、仮想ボリュームに対するホスト計算機 180 からの書き込みに対して、プールから必要な容量の 1 又は複数ページを割り当てる。

【0126】

ストレージシステム 10 は、ホスト計算機 180 により認識される仮想ボリュームの容量を、仮想ボリュームに割り当てられている実容量よりも大きくすることができ、ホスト計算機 180 に割り当てられる容量を実現するために必要な実容量を、それよりも小さくすることができる（シンプロビジョニング）。

【0127】

ストレージシステム 10 は、仮想ボリュームに対するホスト計算機 180 からの I/O 負荷を分析し、I/O 負荷の高いページを、性能の高い高価なメディアで構成されたりソースから成る上位階層に、そうでないページを性能の低い安価なメディアで構成されたりソースから成る下位階層に自動配置する。これにより、仮想ボリュームへのアクセス性能を維持しつつ、システムのコストを低減することができる。

【0128】

以下において、第 1 実施形態との差異を主に説明する。図 26 は、本実施形態のローカルメモリ 122 が格納している情報を示している。ローカルメモリ 122 における制御情報は、第 1 実施形態で説明した情報に加え、ページ毎モニタ差分テーブル 320 を含む。図 27 は、本実施形態の共有メモリ 132 が格納するデータを示している。共有メモリ 132 の制御情報は、第 1 実施形態で説明した情報に加え、ダイナミックマッピングテーブル 520 及びページ毎モニタテーブル 530 を含む。

【0129】

図 28 は、ダイナミックマッピングテーブル 520 の一例を示す。ダイナミックマッピングテーブル 520 は、各仮想ボリュームにおいて、アクセス数をカウントするエントリ（記憶領域のエントリ）を管理するテーブルである。例えば、1 ページが、ダイナミックマッピングテーブル 520 の 1 エントリである。ここでは、この例を説明する。

【0130】

ダイナミックマッピングテーブル 520 は、プール番号のカラム 521、仮想ボリューム番号のカラム 522、論理アドレスのカラム 523、プールボリューム番号のカラム 524、論理アドレスのカラム 525、モニタ情報インデックス番号のカラム 526 を有する。プール番号及び仮想ボリューム番号は、それぞれ、ストレージシステム 10 内で、プールと仮想ボリュームを一意に識別する識別子である。モニタ情報インデックス番号は、ダイナミックマッピングテーブル 520 におけるエントリ識別子である。

【0131】

論理アドレスのカラム 523 は、各エントリの仮想ボリュームにおける開始論理アドレスを格納する。論理アドレスのカラム 525 は、各エントリのプールボリュームにおける開始論理アドレスを格納する。本例においてエントリの容量は一定であるが、一定でなくともよい。

【0132】

図 29 は、ページ毎モニタテーブル 530 の一例を示す。ページ毎モニタテーブル 530 は、各ページの I/O 数を管理する。マイクロプロセッサ 121 は、このテーブル 530 を参照して、当該ページのデータを格納する階層を決定する。

【0133】

ページ毎モニタテーブル 530 は、モニタ情報インデックス番号のカラム 531、I/O

10

20

30

40

50

0カウンタ（現在）のカラム532、I/Oカウンタ（前回）のカラム533を有する。マイクロプロセッサ121は、ページへのアクセスを監視し、所定の監視期間（例えば1秒）内のI/O数（アクセス数）をカウントして、ページ毎モニタテーブル530に格納する。監視期間は連続して続く。

【0134】

I/Oカウンタ（前回）のカラム533は、前回監視期間におけるI/O数を格納する。I/Oカウンタ（現在）のカラム532は、現在監視期間におけるI/O数を格納する。マイクロプロセッサ121は、現在監視期間内において、I/Oカウンタ（現在）のカラム532の値を繰り返し更新する。

【0135】

本構成において、マイクロプロセッサ121は、ローカルメモリ122におけるページ毎モニタ差分テーブル320を使用してI/O数をカウントし、その値の更新を共有メモリ132におけるページ毎モニタテーブル530に反映する。この点は後述する。現在監視期間が終了すると、マイクロプロセッサ121は、現在監視期間におけるI/O数を、前回監視期間におけるI/O数のフィールドに移す。

【0136】

図30は、ページ毎モニタ差分テーブル320の一例を示す。ページ毎モニタ差分テーブル320は、各ページへのアクセスをカウントするために使用される。ページ毎モニタ差分テーブル320は、モニタ情報インデックス番号のカラム321及びI/O差分カウンタのカラム322を有する。マイクロプロセッサ121は、各ページのアクセスを監視し、アクセスがあると、I/O差分カウンタのカラム322の該当フィールドの値をインクリメントする。

【0137】

I/O差分カウンタのカラム322のフィールドの値が規定値（本例で最大値）に達すると、マイクロプロセッサ121は、ページ毎モニタテーブル530の対応エントリのI/Oカウンタ（現在）のカラム532のフィールドの値にその値を加算して、当該フィールドを更新する。マイクロプロセッサ121は、最大値に達したI/O差分カウンタのカラム322のフィールドの値を初期値（0値）に戻す。I/O差分カウンタは、このようにページ毎モニタテーブル530の前回更新からの、I/O数の差分を示す。

【0138】

図30及び図29に示すように、ページ毎モニタ差分テーブル320のI/O差分カウンタのカラム322は8ビットデータを格納し、ページ毎モニタテーブル530のI/Oカウンタ（現在）のカラム532は、8ビットよりも大きい32ビットのデータを格納する。データのビット数は設計に依存し、8又は32ビットに限定されない。

【0139】

図31のフローチャートを参照して、上記ストレージ階層仮想化機能モニタ更新の具体的な方法を説明する。マイクロプロセッサ121は、ページへのアクセスを受けると、ページ毎モニタ差分テーブル320におけるそのページのI/O差分カウンタをインクリメントする（S302）。ページに対するアクセスは、ページにマッピングされている記憶デバイスへのI/O又はページに対するホストI/Oを示し、マイクロプロセッサ121は、いずれかのI/O数をカウントする。

【0140】

マイクロプロセッサ121は、当該論理ボリューム性能ブースト機能がONであるか判定する（S303）。このステップは、図16におけるステップS122と同様である。ボリューム性能ブースト機能がOFFである場合（S303：NO）、マイクロプロセッサ121は、ステップS307に進む。

【0141】

ボリューム性能ブースト機能がONである場合（S303：YES）、マイクロプロセッサ121は、自身の過負荷フラグがONであるか否かを判定する（S304）。このステップは、図16におけるステップS125と同様である。

10

20

30

40

50

【0142】

過負荷フラグがONである場合（S304：YES）、マイクロプロセッサ121は、ステップS306に進む。過負荷フラグがOFFである場合（S304：NO）、マイクロプロセッサ121は、アクセス先のCMPK130の過負荷フラグがONであるか否かを判定する（S305）。このステップは、図16におけるステップS126と同様である。

【0143】

CMPK130の過負荷フラグがOFFである場合（S305：NO）、マイクロプロセッサ121は、ステップS307に進む。CMPK130の過負荷フラグがONである場合（S305：YES）、マイクロプロセッサ121は、ステップS306に進む。ステップS306において、マイクロプロセッサ121は、ページ毎モニタ差分テーブル320の上記I/O差分カウンタの値が、最大値であるかを判定する。

10

【0144】

I/O差分カウンタの値が最大値未満である場合（S306：NO）、このフローは終了する。I/O差分カウンタの値が最大値である場合（S306：YES）、マイクロプロセッサ121は、ページ毎モニタテーブル530の対応エントリのI/Oカウンタ（現在）のカラム532のフィールドの値にその最大値を加算して、当該フィールドを更新する（S307）。マイクロプロセッサ121は、さらに、最大値に達したI/O差分カウンタのカラム322のフィールドの値を0値（初期値）に設定する（S308）。

【0145】

本例は、マイクロプロセッサ121及びCMPK130の負荷が小さい場合、ローカルメモリ122におけるI/O差分カウンタの更新に同期して、共有メモリ132のI/Oカウンタを更新する。これらの負荷が小さいためシステム性能の低下が問題とならず、障害発生時に正確なI/Oカウント数を得ることができる。これら二つのデバイスの負荷条件は省略してもよく、双方の成立をI/Oカウンタ値の非同期更新の条件としてもよい。これらと異なる条件を使用してもよい。

20

【0146】

上述のように、マイクロプロセッサ121は、ローカルメモリ122内のカウンタでページのI/O数をカウントし、その値が規定値に達すると、その規定値を共有メモリ132のカウンタに反映する。これにより、マイクロプロセッサ121とCMPK130との間の通信によるオーバーヘッドを低減する。

30

【0147】

ページ毎モニタ差分テーブル320のカウンタのビット数が、ページ毎モニタテーブル530のカウンタのビット数より小さい。このように、ローカルメモリ上で差分をカウントすることで、I/O数カウントのためにローカルメモリ122で必要とされる容量を削減することができる。MPPK120の障害時には、所定期間のI/Oカウント数の情報が失われるが、ページI/Oカウント数にI/Oカウント数の差分の反映がなされないだけであるので、ページのI/O解析に実質的な影響を与えることはない。

【0148】

なお、本実施形態の性能モニタ方法は、階層仮想化機能のモニタに限らず、そのほかの性能モニタにも適用可能である。例えば、HDDやSSDなどのドライブのモニタに適用できる。上記例は、カウンタ数が最大値に達するとカウンタを初期化するが、初期化においてI/Oをカウントしてもよい。マイクロプロセッサ121は、例えば、I/O差分カウンタの初期化と共に、その最大数に1を加えた値をページ毎モニタテーブル530のI/Oカウンタの値に加算する。これは、他の実施形態におけるカウント方法と同様である。

40

【0149】

第3実施形態

以下において、本発明を非同期リモートコピーに適用した例を説明する。以下においては、第1実施形態及び第2実施形態との差異を主に説明する。図32は、本実施形態の計

50

算機システムの構成を模式的に示すブロック図である。本実施形態のストレージシステムは、第1ストレージシステム10A及び第2ストレージシステム10Bを含む。典型的には、第1ストレージシステム10A及び第2ストレージシステム10Bは異なるサイトに設置されており、データネットワーク(例えばSAN)190A、データネットワーク(例えばSAN)190B及び広域ネットワークを介して通信可能に接続する。

【0150】

第1ストレージシステム10A及び第2ストレージシステム10Bは、図1を参照して説明したハードウェア構成と同様の構成を有する。具体的には、第1ストレージシステム10Aは、複数のFEPK110A、複数のMPPK120A、複数のCMPK130A、複数のBEPK140Aを含み、これらは内部ネットワーク150Aを介して接続する。第1管理計算機20Aは、第1ストレージシステム10Aを管理する。

10

【0151】

同様に、第2ストレージシステム10Bは、複数のFEPK110B、複数のMPPK120B、複数のCMPK130B、複数のBEPK140Bを含み、これらは内部ネットワーク150Bを介して接続する。第2管理計算機20Bは、第2ストレージシステム10Bを管理する。

【0152】

第1ストレージシステム10A及び第2ストレージシステム10Bは、非同期リモートコピー機能を有する。第1ストレージシステム10Aのプライマリボリューム(PVOL)171Pと、第2ストレージシステム10Bのセカンダリボリューム(SVOL)171Sが、コピーペアを構成する。ボリュームは、典型的には、1又は複数のRAIDグループにおける1又は複数の記憶領域からなる。

20

【0153】

プライマリボリューム171Pがコピー元ボリューム、セカンダリボリューム171Sがコピー先ボリュームであり、プライマリボリューム171Pのデータが、セカンダリボリューム171Sにコピーされる。プライマリボリューム171Pへデータ書き込み順序と、セカンダリボリューム171Sへのデータコピー順序は一致する(順序保障)。

【0154】

同期コピーは、ホスト計算機180がプライマリボリューム171Pに書き込みをおこなった場合、セカンダリボリューム171Sへのコピーの完了後(典型的にはキャッシュメモリへの書き込み後)に、ホスト計算機180にI/O成功を通知する。これに対して、非同期コピーは、プライマリボリューム171Pへの書き込み完了後、セカンダリボリューム171Sへのコピー完了前に、ホスト計算機180にI/O成功を通知する。

30

【0155】

本実施形態のストレージシステムは、プライマリボリューム171Pからセカンダリボリューム171Sへのコピー用のバッファとして、ジャーナルボリューム(JVOL)171JP、171JSを使用する。第1ストレージシステム10Aにおいて、プライマリボリューム171Pとジャーナルボリューム171JPとがグループ化されている。第2ストレージシステム10Bにおいて、セカンダリボリューム171Sとジャーナルボリューム171JSとがグループ化されている。

40

【0156】

プライマリボリューム171Pにおける更新データは、ジャーナルボリューム171JP、171JSを介して、セカンダリボリューム171Sに送信される。これにより、リモートコピーのデータ転送において、性能が不安定な広域ネットワークを使用することができる。

【0157】

図33を参照して、ホスト計算機180からのプライマリボリューム171Pへのデータ書き込み及びその更新データのセカンダリボリューム171Sへのコピーの流れを説明する。FEPK110Aは、ホスト計算機180からのライトコマンド及びライトデータを受信する。MPPK120A(のマイクロプロセッサ121)は、ライトコマンドを解

50

析し、FEPK110A及びBEPK140A(不図示)に、プライマリボリューム171P及びジャーナルボリューム171JPにライトデータを書き込むことを指示する。

【0158】

具体的には、MPPK120Aは、FEPK110A及びBEPK140Aにライトデータを指定した次の転送先に転送することを指示する。最終的な転送先はプライマリボリューム171P及びジャーナルボリューム171JPであり、ライトデータは、プライマリボリューム171P及びジャーナルボリューム171JPのそれぞれに書き込まれる。ジャーナルボリューム171JPへの書き込み順序をもって、プライマリボリューム171Pへの書き込み順序とする。

【0159】

本図において、ライトデータのキャッシュメモリ131への書き込みの説明は省略されている、又はライトデータはキャッシュメモリ131を介することなくボリュームに格納される。MPPK120Aは、ライトデータのキャッシュメモリ131への書き込み完了又はボリュームへの書き込み完了に回答して、ホスト計算機180に書き込み完了を通知する。

【0160】

MPPK120Aは、ジャーナルボリューム171JPの更新に従って、ジャーナルボリューム171JPの管理データを更新する。図33に示すように、ジャーナルボリューム171JPは、管理領域611とデータ領域612を有し、それぞれが、ジャーナルボリューム管理データ及び更新データを格納する。ジャーナルボリューム管理データはジャーナルボリューム外に格納されていてもよい。

【0161】

ジャーナルボリューム管理データは、シーケンス番号601及びポインタ602のペアを含む。これらの値のペアが、各ライトデータ(更新データ)に付与される。本図の例において、シーケンス番号601は、1からnの値のいずれかの値であり、データ領域に格納された順に、各ライトデータに昇順で付与される。シーケンス番号は循環的であり、nが付与されたライトデータの次のデータには1が付与される。ポインタ602は、データ領域612において対応するシーケンス番号が付与されているライトデータが格納されている位置(アドレス)を示す。

【0162】

管理領域611は、シーケンス番号601とポインタ602のペアが書き込まれている領域と、未使用領域604を含む。未使用領域604は初期値を格納しており、本例において初期値は0値である。マイクロプロセッサ121は、データ領域612に格納されている更新データを第2ストレージシステム10Bに転送すると、そのデータのシーケンス番号601とポインタ602を格納している領域の値を初期値(無効値)に更新する。更新データの転送順序は、シーケンス番号601の順に一致し、シーケンス番号601の順をもって、更新データのジャーナルボリューム171JPへの書き込み順とする。

【0163】

管理領域611において、シーケンス番号601とポインタ602の次の新たなペアを書き込む位置は決まっており、例えば、ペアは、管理領域611におけるアドレス昇順で書き込まれる。終点アドレスに書き込まれているペアの次のペアは開始アドレスに書き込まれる。

【0164】

シーケンス番号601とポインタ602とを格納する領域(ジャーナル領域とも呼ぶ)において、初期値を格納している領域の直前位置のシーケンス番号601、つまりジャーナル領域の先頭のシーケンス番号が最も新しい更新データを示す。一方、初期値を格納している領域の直後位置のシーケンス番号601、つまりジャーナル領域の最後尾のシーケンス番号が最も古い更新データを示す。

【0165】

第2ストレージシステム10BのMPPK120Bが第1ストレージシステム10Aへ更

10

20

30

40

50

新データのコピーの要求を送ったとき、第1ストレージシステム10AのMPPK120Aは、ジャーナルボリューム171JPに格納されている更新データおよびシーケンス番号を、更新順（書き込み順）で、第2ストレージシステム10Bに転送する。第2ストレージシステム10BのMPPK120Bは、そのFEPK110Bが受信した更新データを、順次、ジャーナルボリューム171JSに格納する。本図においてキャッシュメモリ131へのキャッシングが省略されている。MPPK120Bは、規定のタイミングで、ジャーナルボリューム171JSに格納されている更新データを、更新順序でセカンダリボリューム171Sに書き込む。

【0166】

第2ストレージシステム10Bのジャーナルボリューム171JSは、ジャーナルボリューム171JPと同様に、管理領域とデータ領域とを含み、それぞれが、ジャーナル管理データと更新データを格納する。

【0167】

MPPK120Bは、転送された更新データをジャーナルボリューム171JSに格納してから、転送されたシーケンス番号及びポインタを書き込み、管理データを更新する。管理データの構成はジャーナルボリューム171JPと同様である。ジャーナルボリューム171JS内の更新データがセカンダリボリューム171Sに書き込まれると、MPPK120Bは、対応するシーケンス番号とポインタの値を初期値（無効値）に変更する。

【0168】

図34は、第1ストレージシステム10A及び第2ストレージシステム10Bにおけるローカルメモリ122が格納している制御情報を示している。本実施形態において、LM非同期リモートコピーシーケンス番号管理テーブル330が、ローカルメモリ122内に格納されている。図35は、第1ストレージシステム10A及び第2ストレージシステム10Bにおける共有メモリ132が格納している制御情報を示している。本実施形態において、非同期リモートコピー管理テーブル540及びSM非同期リモートコピーシーケンス番号管理テーブル530が格納されている。

【0169】

非同期リモートコピー管理テーブル540は、ペア管理のための管理情報を格納している。具体的には、プライマリボリュームとセカンダリボリュームの各ペアを管理する管理情報、リモートコピーのパスの情報、そして、プライマリボリューム及びセカンダリボリュームのそれぞれとグループ化されるジャーナルボリュームの情報を含む。マイクロプロセッサ121は、この管理テーブル540を参照して、リモートコピーの実行を制御する。

【0170】

図36は、LM非同期リモートコピーシーケンス番号管理テーブル330の一例を示す。LM非同期リモートコピーシーケンス番号管理テーブル330は、ローカルメモリ122において、各ジャーナルボリュームの最新シーケンス番号を管理する。MPPK120Aのマイクロプロセッサ121は、LM非同期リモートコピーシーケンス番号管理テーブル330を参照して、新たにジャーナルボリューム171JSに書き込まれる更新データのシーケンス番号を決定することができる。

【0171】

LM非同期リモートコピーシーケンス番号管理テーブル330は、JVOL番号のカラム331、シーケンス番号のカラム332、そしてシーケンス番号差分のカラム333を有する。JVOL番号は、第1ストレージシステム10Aにおけるジャーナルボリュームの識別子である。シーケンス番号のカラム332は、JVOLにおける先頭シーケンス番号を示すデータを格納する。シーケンス番号差分については後述する。

【0172】

図37は、SM非同期リモートコピーシーケンス番号管理テーブル530の一例を示す。SM非同期リモートコピーシーケンス番号管理テーブル530は、共有メモリ132において、各ジャーナルボリュームのシーケンス番号を管理する。SM非同期リモートコピ

10

20

30

40

50

シーケンス番号管理テーブル530は、JVOL番号のカラム531及びシーケンス番号のカラム532を有する。

【0173】

シーケンス番号のカラム532は、JVOLにおける先頭シーケンス番号を示すデータを格納する。1エン트리におけるシーケンス番号のカラム532の値は、ローカルメモリ122において対応するシーケンス番号のカラム332の値と一致する又は異なる(図36及び図37の例では全てのエントリの値が異なる)。それらの更新は、同期又は非同期である。

【0174】

図36及び図37に示すように、各JVOLのエン트리において、シーケンス番号差分カラム333のフィールドの値は、LM非同期リモートコピーシーケンス番号管理テーブル330のシーケンス番号カラム332の対応フィールドの値と、SM非同期リモートコピーシーケンス番号管理テーブル530のシーケンス番号カラム532の対応フィールドの値との差分である。

10

【0175】

このように、シーケンス番号差分カラム333のフィールドの値は、シーケンス番号カラム532における対応フィールドの前回更新からのJVOLにおけるシーケンス番号の更新を示し、共有メモリ132に格納されている前回更新時の先頭シーケンス番号と最新の先頭シーケンス番号との差分を示す。

【0176】

20

MPPK120Aのマイクロプロセッサ121は、ジャーナルボリュームに更新データの書き込みがあるたびに、そのジャーナルボリュームのエン트리において、シーケンス番号カラム332及びシーケンス番号差分カラム333の値をインクリメントする。シーケンス番号カラム332の各フィールドは、対応するジャーナルボリュームの最新のシーケンス番号(最後に割り当てたシーケンス番号)を示している。シーケンス番号カラム332の各フィールドの値は、最大値からインクリメントされると最小値に戻る。

【0177】

シーケンス番号差分カラム333のビット数(最大値)は、シーケンス番号カラム332のビット数(最大値)よりも小さい。マイクロプロセッサ121は、シーケンス番号差分カラム333のフィールドの値が最大値に達すると、LM非同期リモートコピーシーケンス番号管理テーブル330における当該エントリの更新を、SM非同期リモートコピーシーケンス番号管理テーブル530の対応エントりに反映する。

30

【0178】

具体的には、SM非同期リモートコピーシーケンス番号管理テーブル530における対応エントリのシーケンス番号を、LM非同期リモートコピーシーケンス番号管理テーブル330の対応エントリのシーケンス番号に一致させる。SM非同期リモートコピーシーケンス番号管理テーブル530における更新値は、更新前の値にシーケンス番号差分カラム333における対応フィールドの値を加算した値である。

【0179】

このように、ローカルメモリ122においてシーケンス番号の最大数よりも小さい所定数までシーケンス番号の変化をカウントし、ローカルメモリ122におけるシーケンス番号の変化を共有メモリ132のシーケンス番号に反映することで、マイクロプロセッサ121によるCMPK130へのアクセス回数を低減し、それらの間の通信によるマイクロプロセッサ121及びCMPK130の負荷を低減することができる。

40

【0180】

図38のフローチャートを参照して、本実施形態の非同期リモートコピーシーケンス番号の更新を説明する。この処理は、ジャーナルボリューム171JPの担当MPPK120Aのマイクロプロセッサ121が実行する。

【0181】

マイクロプロセッサ121は、ジャーナルボリューム171JPへの更新データ書き込

50

みがあると、LM非同期リモートコピーシーケンス番号管理テーブル330を参照して、当該ジャーナルボリューム171JPの管理領域611に、新たなシーケンス番号及びポインタを追加する。さらに、マイクロプロセッサ121は、LM非同期リモートコピーシーケンス番号管理テーブル330において、当該ジャーナルボリューム171JPのエントリのシーケンス番号及びシーケンス番号差分の値を更新する（本例においてそれら値をインクリメントする）（S412）。

【0182】

マイクロプロセッサ121は、当該ボリュームの性能ブースト機能がONであるか判定する（S413）。性能ブースト機能がOFFである場合（S413:NO）、マイクロプロセッサ121は、ステップS417に進む。性能ブースト機能がONである場合（S413:YES）、マイクロプロセッサ121は、自身の過負荷フラグがONであるか判定する（S414）。

10

【0183】

過負荷フラグがONである場合（S414:YES）、マイクロプロセッサ121は、ステップS416に進む。過負荷フラグがOFFである場合（S414:NO）、マイクロプロセッサ121は、アクセス先のCMPK130Aの過負荷フラグがONであるか判定する（S415）。

【0184】

CMPK130Aの過負荷フラグがOFFである場合（S415:NO）、マイクロプロセッサ121は、ステップS417に進む。CMPK130Aの過負荷フラグがONである場合（S415:YES）、マイクロプロセッサ121は、ステップS416に進む。ステップS413からステップS415の詳細は、第2実施形態で既に説明した通りである。マイクロプロセッサ121及び/又はCMPK130Aの負荷に応じて制御情報の更新反映を制御することで、システム性能の低下を抑えつつ、共有メモリの更新をより適切に行うことができる。

20

【0185】

ステップS416において、マイクロプロセッサ121は、LM非同期リモートコピーシーケンス番号管理テーブル330において、当該ジャーナルボリューム171JPのシーケンス番号差分が、最大値であるか判定する。その値が最大値ではない場合（S416:NO）、マイクロプロセッサ121は、本処理を終了する。

30

【0186】

上記値が最大値である場合（S416:YES）、マイクロプロセッサ121は、SM非同期リモートコピーシーケンス番号管理テーブル530において、当該ジャーナルボリューム171JPのシーケンス番号を更新する（S417）。具体的には、マイクロプロセッサ121は、現在のシーケンス番号の値にシーケンス番号差分の値を加算した値に更新する。ステップS418において、マイクロプロセッサ121は、最大値に達しているシーケンス番号差分のフィールドの値を0値に更新（初期化）する。

【0187】

上記シーケンス番号差分を使用した共有メモリ132におけるシーケンス番号の更新（性能ブースト機能）を使用しない場合、LM非同期リモートコピーシーケンス番号管理テーブル330及びSM非同期リモートコピーシーケンス番号管理テーブル530の更新は同期する。

40

【0188】

MPPK120Aに障害が発生した場合、ローカルメモリ122上のLM非同期リモートコピーシーケンス番号管理テーブル330が失われる。上述のように、このテーブル330は、各ジャーナルボリュームの最新の先頭シーケンス番号を示す情報を有している。正常なりモートコピーを行うためには、ジャーナル管理データにおける最新の先頭シーケンス番号が必要である。

【0189】

本実施形態の第1ストレージシステム10Aは、障害発生したMPPK120Aと異なる

50

るM P P K 1 2 0 Aが、ジャーナルボリューム1 7 1 J Pの管理領域6 1 1を参照して、ジャーナル領域の先頭を示す最新の先頭シーケンス番号を確認する。図3 9のフローチャートを参照して、M P P K障害発生時の非同期リモートコピーシーケンス番号回復処理を説明する。

【0 1 9 0】

担当を引き継いだ正常なM P P K 1 2 0 Aのマイクロプロセッサ1 2 1は、共有メモリ1 2 3に格納されているS M非同期リモートコピーシーケンス番号管理テーブル5 3 0から、一つのジャーナルボリュームを選択し、そのシーケンス番号を読み出す(S 4 2 2)。マイクロプロセッサ1 2 1は、そのジャーナルボリュームから、上記シーケンス番号の領域の次のシーケンス番号領域からデータを読み出す(S 4 2 3)。

10

【0 1 9 1】

マイクロプロセッサ1 2 1は、ステップS 4 2 3で読み出したシーケンス番号が0値(無効値)であるか判定する(S 4 2 4)。そのシーケンス番号が0値ではない場合(S 4 2 4 : N O)、マイクロプロセッサ1 2 1は、その読み出したシーケンス番号をテンポラル領域(典型的にそのローカルメモリ1 2 2内の領域)に格納する(S 4 2 5)。

【0 1 9 2】

そのシーケンス番号が0値である場合(S 4 2 4 : Y E S)、その領域は未使用領域であり、マイクロプロセッサ1 2 1は、テンポラル領域に格納されているシーケンス番号で、S M非同期リモートコピーシーケンス番号管理テーブル5 3 0における対応ジャーナルボリュームのシーケンス番号を更新する。S M非同期リモートコピーシーケンス番号管理テーブル5 3 0のシーケンス番号が最新の先頭シーケンス番号である場合、更新は不要である。マイクロプロセッサ1 2 1は、障害発生したM P P K 1 2 0 Aが担当であった全てのジャーナルボリュームについて、上記更新を行う。

20

【0 1 9 3】

上記フローにより、S M非同期リモートコピーシーケンス番号管理テーブル5 3 0が最新情報を含むように更新され、他のM P P K 1 2 0 Aが、障害が起きたM P P K 1 2 0 Aの担当を引き継ぎ、正常な非同期リモートコピーを続けることができる。

【0 1 9 4】

上記シーケンス番号管理テーブル3 3 0、5 3 0が格納する値は一例であって、それらは、先頭シーケンス番号又はそれらテーブル3 3 0、5 3 0の先頭シーケンス番号間の差分を示すことができれば、どのような値を格納していてもよい。

30

【0 1 9 5】

第4実施形態

以下において、本発明を非同期ローカルコピーに適用した例を説明する。以下においては、上記他の実施形態と異なる点を主に説明する。図4 0は、本実施形態のローカルメモリ1 2 2に格納されている制御情報を示している。ローカルメモリ1 2 2には、L Mローカルコピー差分管理テーブル3 4 0及びL Mローカルコピー差分領域間引き動作管理テーブル3 5 0が格納されている。

【0 1 9 6】

図4 1は、本実施形態の共有メモリ1 3 2内の制御情報を示している。S Mローカルコピー差分管理テーブル5 6 0、S Mローカルコピー差分領域間引き動作管理テーブル5 7 0、ローカルコピー管理テーブル5 8 0が、共有メモリ1 3 2における制御情報に含まれている。複数のM P P K 1 2 0が、共有メモリ1 3 2内テーブル5 6 0、5 7 0、5 8 0を参照可能である。特に、S Mローカルコピー差分管理テーブル5 6 0及びS Mローカルコピー差分領域間引き動作管理テーブル5 7 0は、M P P K 1 2 0の障害時に、他のM P P K 1 2 0により参照される。

40

【0 1 9 7】

ローカルコピー管理テーブル5 8 0は、プライマリボリュームとセカンダリボリュームの各ペアを管理する管理情報を含む。例えば、ペアを構成するプライマリボリュームとセカンダリボリュームの識別情報、それらのアドレス情報及びコピーポリシの情報を含む。

50

マイクロプロセッサ 121 は、ローカルコピー管理テーブル 580 を参照して、ローカルコピーの実行を制御する。

【0198】

共有メモリ 132 内の SM ローカルコピー差分管理テーブル 560 及び SM ローカルコピー差分領域間引き動作管理テーブル 570 は、それぞれ、ローカルメモリ 122 内の LM ローカルコピー差分管理テーブル 340 及び LM ローカルコピー差分領域間引き動作管理テーブル 350 のバックアップである。マイクロプロセッサ 121 は、予め定められた規則に従って、ローカルメモリ 122 でのテーブル 340、350 の更新を、共有メモリ 132 のテーブル 560、570 に反映する。

【0199】

図 42 は、LM ローカルコピー差分管理テーブル 340 の一例を示す。LM ローカルコピー差分管理テーブル 340 は、ボリューム番号のカラム 341、論理アドレスのカラム 342、差分有ビット列のカラム 343 を有する。ボリューム番号は、ストレージシステム内でのプライマリボリュームの識別子である。各エントリは、ボリューム内の所定広さの記憶領域（アドレス範囲）を示している。論理アドレスは、各エントリの記憶領域の開始論理アドレスを示す。本例において、エントリの記憶領域の広さは共通である。

【0200】

差分有ビット列は、そのエントリの記憶領域において、プライマリボリュームとセカンダリボリュームとの間にデータの相違が存在するか否か、つまり、プライマリボリュームでの更新がセカンダリボリュームに反映されているか否かを示す。

【0201】

差分有ビット列の各ビット（差分有ビットとも呼ぶ）は、エントリの記憶領域における各部分領域のデータがプライマリボリュームとセカンダリボリュームとの間で異なるか否かを示す。本例では、各ビットに対応する領域の広さは共通である。本例において、差分有ビット列のビットが 1 である場合、その領域のデータは、プライマリボリュームとセカンダリボリュームとで異なることを示す。

【0202】

マイクロプロセッサ 121 は、所定のタイミングで、プライマリボリュームの更新データをセカンダリボリュームにコピーする（非同期ローカルコピー）。非同期ローカルコピーにおいて、マイクロプロセッサ 121 は、LM ローカルコピー差分管理テーブル 340 を参照し、プライマリボリュームにおける差分有ビットが 1 である領域のデータを、セカンダリボリュームにコピーする。

【0203】

この非同期ローカルコピーにตอบสนองして、マイクロプロセッサ 121 は、LM ローカルコピー差分管理テーブル 340 において、更新がセカンダリボリュームに反映された領域の差分有ビットを 0 値に更新する。

【0204】

図 43 は、SM ローカルコピー差分管理テーブル 560 の一例を示す。SM ローカルコピー差分管理テーブル 560 は、LM ローカルコピー差分管理テーブル 340 のバックアップテーブルであり、LM ローカルコピー差分管理テーブル 340 と同一の構成を有する。具体的には、ボリューム番号のカラム 561、論理アドレスのカラム 562、差分有ビット列のカラム 563 を有する。

【0205】

マイクロプロセッサ 121 は、所定規則に従って、LM ローカルコピー差分管理テーブル 340 における更新を、SM ローカルコピー差分管理テーブル 560 にコピーする。本例において、プライマリボリュームからセカンダリボリュームへの非同期ローカルコピーによる LM ローカルコピー差分管理テーブル 340 の更新と SM ローカルコピー差分管理テーブル 560 更新は同期する。プライマリボリュームへのデータライトによる更新に対する SM ローカルコピー差分管理テーブル 560 の更新については後述する。

【0206】

10

20

30

40

50

図 4 4 は、L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 の一例を示す。L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 は、ボリューム番号のカラム 3 5 1、論理アドレスのカラム 3 5 2、間引き中ビット列のカラム 3 5 3 を有する。各エントリは、ボリューム内の所定広さの記憶領域（アドレス範囲）を示している。

【 0 2 0 7 】

論理アドレスは、各エントリの記憶領域の開始論理アドレスを示す。本例において、エントリの記憶領域の広さは共通である。L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 におけるエントリの記憶領域は、L Mローカルコピー差分管理テーブル 3 4 0 のエントリの記憶領域よりも広い。

【 0 2 0 8 】

間引き中ビット列は、L Mローカルコピー差分管理テーブル 3 4 0 における差分有ビット列の更新を、S Mローカルコピー差分管理テーブル 5 6 0 の対応する差分有ビット列に反映するか否かを示す。上述のように、L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 において、間引き中ビット列は、論理ボリューム内の記憶領域に関連づけられている。

【 0 2 0 9 】

間引き中ビット列の各ビット（間引き中ビットとも呼ぶ）は、その間引き中ビット列に関連付けられている記憶領域の部分領域に関連付けられている。間引き中ビット列の各ビットは、それが関連づけられている部分領域を介して、1 又は複数の差分有ビットに関連づけられる。

【 0 2 1 0 】

好ましい本例において、間引き中ビットは複数の差分有ビットに関連づけられている。また、L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 におけるエントリの記憶領域（アドレス範囲）は、L Mローカルコピー差分管理テーブル 3 4 0 におけるエントリの記憶領域（アドレス範囲）よりも広い。間引き中ビット列のビット数は、差分有ビット列のビット数と同一又は異なる（図 4 3、図 4 4 の例において同一）。

【 0 2 1 1 】

上述のように、L Mローカルコピー差分管理テーブル 3 4 0 において、各差分有ビットは、記憶領域に関連づけられている。間引き中ビットに関連づけられている記憶領域の少なくとも一部が差分有ビットの記憶領域と一致する場合、その間引き中ビットはその差分有ビットに関連づけられている。

【 0 2 1 2 】

間引き中ビットが 1 である場合、ローカルメモリ 1 2 2 においてそれに関連づけられている差分有ビットの、プライマリボリュームの更新（データ書き込み）に回答した更新は、共有メモリ 1 3 2 における差分有ビットに反映されない。具体的には、プライマリボリュームへのライトコマンドの受信に回答して、マイクロプロセッサ 1 2 1 は、L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 において、ライトコマンドが指示する領域の間引き中ビットを参照する。

【 0 2 1 3 】

間引き中ビットが 1 である場合、マイクロプロセッサ 1 2 1 は、L Mローカルコピー差分管理テーブル 3 4 0 において対応する差分有ビットの更新を、S Mローカルコピー差分管理テーブル 5 6 0 に反映しない。これにより、M P P K 1 2 0 と C M P K 1 3 0 との間の通信による M P P K 1 2 0 と C M P K 1 3 0 の負荷を低減する。

【 0 2 1 4 】

図 4 5 は、S Mローカルコピー差分領域間引き動作管理テーブル 5 7 0 の一例を示す。S Mローカルコピー差分領域間引き動作管理テーブル 5 7 0 は、L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 のバックアップテーブルであり、それと同じ構成を有する。具体的には、ボリューム番号のカラム 5 7 1、論理アドレスのカラム 5 7 2、間引き中ビット列のカラム 5 7 3 を有する。マイクロプロセッサ 1 2 1 は、L Mローカルコピー差分領域間引き動作管理テーブル 3 5 0 の更新に同期して、S Mローカルコピー差分領

10

20

30

40

50

域間引き動作管理テーブル570を更新する。

【0215】

図46のフローチャートを参照して、非同期ローカルコピー差分管理情報の更新を説明する。プライマリボリュームにデータが書き込まれると、マイクロプロセッサ121は、LMローカルコピー差分管理テーブル340を更新する(S502)。具体的には、プライマリボリュームにおいて更新された領域に関連づけられている差分有ビットを更新する。

【0216】

マイクロプロセッサ121は、当該ボリュームの性能ブースト機能がONであるか判定する(S503)。性能ブースト機能がOFFである場合(S503:NO)、マイクロプロセッサ121は、ステップS509に進み、SMローカルコピー差分管理テーブル560を更新する(同期更新)。性能ブースト機能がONである場合(S503:YES)、マイクロプロセッサ121は、自身の過負荷フラグがONであるか判定する(S504)。

【0217】

過負荷フラグがONである場合(S504:YES)、マイクロプロセッサ121は、ステップS506に進む。過負荷フラグがOFFである場合(S504:NO)、マイクロプロセッサ121は、アクセス先のCMPK130の過負荷フラグがONであるか判定する(S505)。

【0218】

CMPK130の過負荷フラグがOFFである場合(S505:NO)、マイクロプロセッサ121は、ステップS509に進み、SMローカルコピー差分管理テーブル560を更新する。CMPK130の過負荷フラグがONである場合(S505:YES)、マイクロプロセッサ121は、ステップS506に進む。ステップS503からステップS505の詳細は、第2実施形態で既に説明した通りであり、システム性能の低下を抑えつつ、共有メモリ132の制御情報を適切に更新する。

【0219】

ステップS506において、マイクロプロセッサ121は、プライマリボリュームにおいて更新された領域が間引き中であるか判定する。具体的には、マイクロプロセッサ121は、LMローカルコピー差分領域間引き動作管理テーブル350を参照し、上記更新領域の各間引き中ビットを確認する。間引き中ビットが1である場合(S506:YES)、マイクロプロセッサ121は、SMローカルコピー差分管理テーブル560において、その間引き中ビットに対応する差分有ビットの更新を省略する。

【0220】

間引き中ビットが0である場合(S506:NO)、マイクロプロセッサ121は、その間引き中ビットが関連づけられている領域の差分が閾値以上であるか判定する(S507)。具体的には、マイクロプロセッサ121は、LMローカルコピー差分管理テーブル340を参照し、当該間引き中ビットに対応する差分有ビットにおいて、1のビット数が閾値以上であるか判定する。この判定基準については、図47を参照して後述するMPPK障害時の処理において説明する。

【0221】

差分が閾値未満である場合(S507:NO)、マイクロプロセッサ121は、SMローカルコピー差分管理テーブル560を更新する(S509)。差分が閾値以上である場合(S507:YES)、マイクロプロセッサ121は、LMローカルコピー差分領域間引き動作管理テーブル350及びSMローカルコピー差分領域間引き動作管理テーブル560を更新する(S508)。具体的には、マイクロプロセッサ121は、上記2つのテーブル350、560において、上記間引き中ビットを0から1に変更する。

【0222】

次に、図47のフローチャートを参照して、MPPK120障害時における、ローカルコピー差分のコピーを説明する。MPPK120で障害発生した場合、他のMPPK12

10

20

30

40

50

0 が、障害発生した M P P K 1 2 0 が担当していたコピーペアにおいて、プライマリボリュームからセカンダリボリュームへ、それらの差分をコピーする。これにより、コピーペアの同一性を確保し、その後の正常な非同期ローカルコピーを実現する。

【 0 2 2 3 】

上記他の M P P K 1 2 0 におけるマイクロプロセッサ 1 2 1 は、S M ローカルコピー差分領域間引き動作管理テーブル 5 7 0 を参照し (S 5 1 2)、間引き中領域が残っているか否かを判定する (S 5 1 3)。間引き中領域は、その間引き中ビットが 1 である領域である。間引き領域が残っていなければ (S 5 1 3 : N O)、このフローは終了する。間引き中領域が残っている場合 (S 5 1 3 : Y E S)、マイクロプロセッサ 1 2 1 は、プライマリボリュームにおけるその領域のデータを、セカンダリボリュームにコピーする (S 5 1 4)。

10

【 0 2 2 4 】

上述のように、共有メモリ 1 3 2 は、「 1 」の間引き中ビットに対応する最新の差分有ビット列を格納していない。そのため、M P P K 1 2 0 での障害発生時には、間引き中ビットが 1 (O N) である領域の全てのデータを、プライマリボリュームからセカンダリボリュームにコピーする。これにより、セカンダリボリュームのデータをプライマリボリュームのデータに正確に一致させることができる。

【 0 2 2 5 】

図 4 6 のフローチャートを参照して説明したように、本例は、間引き中ビットに対応する差分有ビットの内の「 1 」のビットが閾値以上である場合に、間引き中ビットを O N (1) に設定する。障害時には、対応する間引き中ビットが O N である全てのデータをプライマリボリュームからセカンダリボリュームにコピーするため、要更新データが多い領域の更新を間引くことで、更新による負荷を低減すると共に障害時の処理を効率化することができる。

20

【 0 2 2 6 】

本実施形態において、差分管理テーブル及び間引き動作管理テーブルの構成は一例であり、差分領域及び間引き中領域を示すことができれば、どのようなデータによりそれらを示してもよい。S M ローカルコピー差分領域間引き動作管理テーブル 5 7 0 を使用する代わりに、間引き領域の L M / S M の差分有ビットを全て 1 にする処理を加えることでも同様の障害時の回復を実現できる。

30

【 0 2 2 7 】

図 4 8 は、第 2 実施形態から第 4 実施形態で使用可能な、性能ブースト機能設定のためのメニュー画面の例 4 8 0 0 を示している。メニュー画面 4 8 0 0 は、性能ブースト機能設定エリア 4 8 0 1、ボリューム毎性能ブースト機能設定エリア 4 8 0 2 及び機能毎性能ブースト機能設定エリア 4 8 0 3 含む。

【 0 2 2 8 】

管理者は、性能ブースト機能設定エリア 4 8 0 1 における“ E N A B L E ”又は“ D I S A B L E ”の一方を入力デバイス 2 8 で選択することで、ストレージシステム 1 0 の性能ブースト機能をイネーブル又はディセーブルすることができる。この設定が、性能ブースト機能有効化テーブル 2 1 0 に反映される。

40

【 0 2 2 9 】

ボリューム毎性能ブースト機能設定エリア 4 8 0 2 は、各論理ボリュームの性能ブースト機能のイネーブル/ディセーブルを可能とする。管理者は、ボリューム毎性能ブースト機能設定エリア 4 8 0 2 において、各論理ボリュームの性能ブースト機能のイネーブル/ディセーブルを入力デバイス 2 8 で選択することができる。この設定が、ボリューム毎性能ブースト機能有効化テーブル 2 2 0 に反映される。

【 0 2 3 0 】

機能毎性能ブースト機能設定エリア 4 8 0 3 は、各性能ブースト機能のイネーブル/ディセーブルを可能とする。管理者は、機能毎性能ブースト機能設定エリア 4 8 0 3 において、各機能のイネーブル/ディセーブルを入力デバイス 2 8 で選択することができる。こ

50

の設定が、ストレージシステム 10 内の機能毎性能ブースト機能有効化テーブル（不図示）に反映される。システム、ボリュームそして機能のブースト機能の全てがイネーブルされている場合に、その性能ブースト機能はそのボリュームにおいて使用される。

【0231】

第5実施形態

本実施形態において、スイッチにより結合した複数のストレージモジュールを含むストレージシステムに本発明を適用した例を説明する。本実施形態は、主に上記他の実施形態と異なる点を説明する。図49は、本実施形態の計算機システムの構成を模式的に示す。ストレージモジュール10C及びストレージモジュール10Dは、モジュール間パス（スイッチ）195（Xパスとも呼ぶ）により通信可能に接続されている。

10

【0232】

図49におけるストレージモジュール10C、10Dの構成は、図1を参照して説明したストレージシステム10の構成と同様である。本例においては、2つの結合したモジュールが一つのストレージシステムを構成するが、3以上のモジュールが一つのストレージシステムを構成してもよい。

【0233】

ストレージモジュール10C及びストレージモジュール10Dを結合するXパス（スイッチ）195は、内部ネットワーク150のパスと同様のパスとして機能し、一方のモジュールの任意のパッケージは、他方のモジュールの任意のパッケージ及びメディアと、Xパス195により通信することができる。また、ホスト計算機180は、いずれのストレージモジュールにもアクセスすることができる。

20

【0234】

Xパスは、内部ネットワーク150よりも帯域が狭く、データ転送能力が低い。そのため、Xパスは、パッケージ間のデータ転送においてボトルネックとなりやすい。そのため、Xパスの負荷に基づいて性能ブースト機能のON/OFFを判定することで、ストレージシステムの性能の低下を小さくすることができる。

【0235】

本実施形態のマイクロプロセッサ121は、性能ブースト機能のイネーブル/ディセーブル制御において、Xパス195の稼働率を参照する。これにより、複数のモジュールからなるストレージシステムにおいて適切にシステム性能を向上することができる。

30

【0236】

図50は、本実施形態のローカルメモリ122が格納している制御情報を示している。図50において、Xパス稼働率テーブル360及びXパス稼働率閾値テーブル370がローカルメモリ122内に格納されている。図51は、Xパス稼働率テーブル360の一例を示す。図52は、Xパス稼働率閾値テーブル370の一例を示す。

【0237】

Xパス稼働率テーブル360は、Xパスの稼働率を管理する。本例において、Xパス稼働率テーブル360は、Xパス番号のカラム361、稼働率のカラム361、そして過負荷判定フラグのカラム363を有する。Xパス番号は、システム内でXパスを一意的に識別する識別子である。図51の例において、Xパス稼働率テーブル360は、複数のXパスを管理している。つまり、複数のXパスが2以上のストレージモジュールを結合している。複数のXパスは、同一又は異なるスイッチを通過する。

40

【0238】

稼働率は、単位時間当たりのデータ転送時間である。Xパスの稼働率は、そのXパスのコントローラが計算し、レジスタに格納する。マイクロプロセッサ121は、各Xパスの稼働率を、スイッチのレジスタから取得して、Xパス稼働率テーブル360に格納する。

【0239】

マイクロプロセッサ121は、Xパス稼働率テーブル360の各エントリ稼働率と、予め定められているXパス稼働率閾値とを比較して、過負荷判定フラグの値を決定する。Xパス稼働率が閾値以上である場合、マイクロプロセッサ121は過負荷判定フラグを1に

50

設定する。Xパス稼働率閾値は、Xパス稼働率閾値テーブル370のXパス稼働率閾値コラムに格納されている。例えば、Xパス稼働率閾値テーブル370は、ストレージシステム内の不揮発性記憶領域からロードされ、その値は管理者により設定される。

【0240】

図53のフローチャートを参照して、Xパスの稼働率を考慮したデータキャッシングに関する制御情報の共有メモリ132における更新についての判定を説明する。基本的な部分は、第1実施形態と同様である。図53のフローチャートにおいて、ステップS607以外のステップは、第1実施形態における図16に示すフローチャートと同様であり、その説明を省略する。

【0241】

ステップS607において、マイクロプロセッサ121は、Xパス稼働率テーブル360を参照し、共有メモリ132へのアクセスに使用するXパスの過負荷フラグが1(ON)であるか判定する。アクセスするCMPK130と使用するXパスとの関係を示す制御情報は、ローカルメモリ122内に格納されており、それにより、マイクロプロセッサ121は、使用するXパスを特定することができる。

【0242】

過負荷フラグがONである場合(S607: YES)、マイクロプロセッサ121は、共有メモリ132の制御情報を更新しないことを決定する(S609)。過負荷フラグがOFF(0)である場合(S607: NO)、マイクロプロセッサ121は、共有メモリ132の制御情報を更新することを決定する(S608)。本例はデータキャッシング制御情報の更新判定においてXパスの稼働率を参照するが、他の実施形態で説明した他の判定処理も、Xパスの稼働率を参照することができる。

【0243】

次に、図54のフローチャートを参照して、Xパス稼働率テーブル360におけるXパス稼働率の更新を説明する。典型的には、この処理は、定期的に、例えば1秒毎に実行される。マイクロプロセッサ121は、一つのXパス、一例としてXパス195を選択し、そのXパス195の稼働率を取得する(S612)。

【0244】

マイクロプロセッサ121は、取得した稼働率の値により、Xパス稼働率テーブル360の該当エントリの稼働率の値を更新する(S613)。マイクロプロセッサ121は、取得した稼働率の値が、Xパス稼働率閾値テーブル370におけるXパス稼働率閾値以上であるか判定する(S614)。稼働率が閾値以上である場合(S614: YES)、マイクロプロセッサ121は、Xパス稼働率テーブル360における当該エントリの過負荷フラグを1(ON)に設定する(S615)。

【0245】

一方、稼働率が閾値未満である場合(S614: NO)、マイクロプロセッサ121は、Xパス稼働率テーブル360における当該エントリの過負荷フラグを0(OFF)に設定する(S616)。マイクロプロセッサ121は、全てのXパスの稼働率を更新したか判定し(S617)、全てのXパスについて判定している場合(S617: YES)にこのフローを終了し、未判定のXパスが残っている場合(S617: NO)には、残りのXパスから一つのXパスを選択して、このフローを繰り返す。

【0246】

第6実施形態

本実施形態は、MPPK120が、複数の異なる種別のデバイスに分散している複数の共有メモリ領域にアクセス可能な構成を説明する。本実施形態において、上記他の実施形態と異なる点について主に説明する。

【0247】

図55は、本実施形態の計算機システムの構成を模式的に示している。ストレージシステム10において、複数の異なるデバイスに共有メモリ(記憶領域)が存在している。具体的には、CMPK130上の共有メモリ132の他、MPPK120上に共有メモリ1

10

20

30

40

50

24、そして記憶ドライブ170に共有メモリ178が存在している。MPPK120上の共有メモリ124の領域は、ローカルメモリ122内の記憶領域である。記憶ドライブ170上の共有メモリ178の領域は、記憶ドライブにおける不揮発性記憶媒体の記憶領域である。

【0248】

図56は、本実施形態のローカルメモリ122が格納している制御情報を示している。図56において、MP稼働率テーブル380、MP稼働率閾値テーブル390、SM領域管理テーブル400がローカルメモリ122内に格納されている。

【0249】

図57は、MP稼働率テーブル380の一例を示す。MP稼働率テーブル380は、MP番号のカラム381、稼働率のカラム382、過負荷判定フラグ1のカラム383、過負荷判定フラグ2のカラム384、稼働時間のカラム385を有する。過負荷判定フラグ2のカラム384以外のカラムは、図11に示すMP稼働率テーブル270と同様である。過負荷判定フラグ1のカラム383は、過負荷判定フラグのカラム273に相当する。

【0250】

図58は、MP稼働率閾値テーブル390の一例を示す。MP稼働率閾値テーブル390は、MP稼働率閾値1のカラム391及びMP稼働率閾値2のカラム392を有する。MP稼働率閾値1の値は、MP稼働率閾値2の値より高い。MP稼働率閾値1は、図12に示すMP稼働率閾値に相当する。

【0251】

図59は、SM領域管理テーブル400の一例を示す。SM領域管理テーブル400は、複数のデバイスに分散している共有メモリ領域を管理する。SM領域管理テーブル400は、種別のカラム401、番号のカラム402、先頭アドレスのカラム403、空き容量のカラム404を有する。「種別」は、共有メモリ領域が存在するデバイスの種別を示す。「番号」は、同一種別のデバイスにおける識別子である。「先頭アドレス」は、各デバイスにおける共有メモリ領域の先頭アドレスを示す。「空き容量」は、共有メモリ領域の空き容量である。

【0252】

種別のカラム401、番号のカラム402、先頭アドレスのカラム403には、予め値が設定されている。マイクロプロセッサ121は、各デバイスのコントローラ(MPPKにおいてはマイクロプロセッサ121)から、共有メモリ領域の空き容量の値を取得し、それを空き容量のカラム404に格納する。

【0253】

図60A及び60Bを参照して、データキャッシングに関する共有メモリ領域に格納された制御情報の更新についての判定を説明する。図60AのフローチャートにおけるステップS702からステップS707は、図16のフローチャートにおけるステップS122からステップS127までと同様である。ただし、ステップS706において、当該CMPK130の過負荷フラグがONである場合(S706: YES)、マイクロプロセッサは図60BにおけるステップS709に進む。

【0254】

ステップS706において当該CMPK130の過負荷フラグがOFFである場合(S706: NO)又はステップS702において当該論理ボリュームの性能ブースト機能がOFFである場合(S702: NO)、マイクロプロセッサ121は、当該CMPK130の共有メモリの制御情報を更新すると決定する。

【0255】

図60BにおけるステップS709において、マイクロプロセッサ121は、SM領域管理テーブル400を参照し、必要な空き共有メモリ領域を有するMPPK120が存在するか判定する。いずれかのMPPK120が必要な空き共有メモリ領域を有する場合(S709: YES)、マイクロプロセッサ121は、そのMPPK120の番号を特定し、キャッシング制御情報を当該MPPK120の共有メモリ124に格納し、その更新を

10

20

30

40

50

行うことを決定する (S 7 1 0)。この M P P K 1 2 0 は、マイクロプロセッサ 1 2 1 が実装された M P P K 1 2 0 と異なる M P P K である。

【 0 2 5 6 】

必要な空き共有メモリ領域を有する M P P K 1 2 0 が存在しない場合 (S 7 0 9 : N O)、マイクロプロセッサ 1 2 1 は、自身の過負荷フラグ 2 が 1 (O N) であるか判定する (S 7 1 1)。過負荷フラグ 2 が O N である場合 (S 7 1 1 : Y E S)、マイクロプロセッサ 1 2 1 は、共有メモリ領域における制御情報の更新を行わないことを決定する (S 7 1 6)。

【 0 2 5 7 】

過負荷フラグ 2 が O F F である場合 (S 7 1 1 : N O)、マイクロプロセッサ 1 2 1 は、 S M 領域管理テーブル 4 0 0 を参照し、必要な空き共有メモリ領域を有する S S D R A I D グループが存在するか判定する (S 7 1 2)。

【 0 2 5 8 】

いずれかの S S D R A I D グループが必要な空き共有メモリ領域を有する場合 (S 7 1 2 : Y E S)、マイクロプロセッサ 1 2 1 は、当該 S S D R A I D グループの番号を特定し、キャッシュ制御情報を当該 S S D R A I D グループの共有メモリ領域に格納し、その更新を行うことを決定する (S 7 1 3)。

【 0 2 5 9 】

必要な空き共有メモリ領域を有する S S D R A I D グループが存在しない場合 (S 7 1 2 : N O)、マイクロプロセッサ 1 2 1 は、 S M 領域管理テーブル 4 0 0 を参照し、必要な空き共有メモリ領域を有する H D D R A I D グループが存在するか判定する (S 7 1 4)。必要な空き共有メモリ領域を有する H D D R A I D グループが存在しない場合 (S 7 1 4 : N O)、マイクロプロセッサ 1 2 1 は、共有メモリ 1 3 2 における制御情報を更新しないことを決定する (S 7 1 6)。

【 0 2 6 0 】

必要な空き共有メモリ領域を有する H D D R A I D グループが存在する場合 (S 7 1 4 : Y E S)、マイクロプロセッサ 1 2 1 は、当該 H D D R A I D グループの番号を特定し、キャッシュ制御情報を当該 H D D R A I D グループの共有メモリ領域に格納し、その更新を行うことを決定する (S 7 1 5)。

【 0 2 6 1 】

マイクロプロセッサ 1 2 1 は、共有メモリ 1 3 2 以外のいずれかの共有メモリに制御情報を格納し、その制御情報を更新することを決定すると、ローカルメモリ 1 2 2 におけるデータキャッシング制御情報を、選択した共有メモリにコピーする。

【 0 2 6 2 】

このように、制御情報を現在の共有メモリ 1 3 2 の領域から他の共有メモリ領域に移動することで、共有メモリにおける制御情報の更新を、ローカルメモリにおける更新に同期させることができ、障害発生時のキャッシュヒット率を向上することができる。上記フローは、アクセス性能が高いデバイスから、空き共有メモリ領域の有無を判定する。これにより、よりアクセス性能が高い共有メモリに制御情報を格納することができ、システム性能の低下を抑えることができる。

【 0 2 6 3 】

本実施形態の共有メモリ領域管理は、データキャッシング制御情報の格納及び更新管理の他、上記他の実施形態で説明した他の制御情報の格納及び更新管理に適用することができる。 M P P K 障害時には、他の M P P K 1 2 0 は、共有メモリ領域管理テーブル 4 0 0 を参照し、分散している共有メモリ領域において対応する制御情報を検索することができる。

【 0 2 6 4 】

図 6 1 のフローチャートを参照して、 M P 稼働率の更新を説明する。このフローは、 1 秒などの周期で呼び出される。マイクロプロセッサ 1 2 1 は、自身の M P 稼働時間を取得し (S 7 2 2)、 M P 稼働率テーブル 3 8 0 の稼働率の値を更新する (S 7 2 3)。ステ

10

20

30

40

50

ップ S 7 2 2、S 7 2 3 は、図 2 4 におけるステップ S 2 3 2、S 2 3 3 と同様である。

【 0 2 6 5 】

次に、ステップ S 7 2 4 において、マイクロプロセッサ 1 2 1 は、更新した稼働率の値が、MP 稼働率閾値 1 の値以上であるか判定する。稼働率の値が MP 稼働率閾値 1 以上である場合 (S 7 2 4 : Y E S)、マイクロプロセッサ 1 2 1 は、MP 稼働率テーブル 3 8 0 の過負荷フラグ 1 を 1 (O N) に設定する (S 7 2 5)。稼働率の値が MP 稼働率閾値 1 未満である場合 (S 7 2 4 : N O)、マイクロプロセッサ 1 2 1 は、MP 稼働率テーブル 3 8 0 の過負荷フラグ 1 を 0 (O F F) に設定する (S 7 2 6)。

【 0 2 6 6 】

次に、ステップ S 7 2 7 において、マイクロプロセッサ 1 2 1 は、更新した稼働率の値が、MP 稼働率閾値 2 以上であるか判定する。稼働率の値が MP 稼働率閾値 2 以上である場合 (S 7 2 7 : Y E S)、マイクロプロセッサ 1 2 1 は、MP 稼働率テーブル 3 8 0 の過負荷フラグ 2 を 1 (O N) に設定する (S 7 2 8)。稼働率の値が MP 稼働率閾値 2 未満である場合 (S 7 2 7 : N O)、マイクロプロセッサ 1 2 1 は、MP 稼働率テーブル 3 8 0 の過負荷フラグ 2 を 0 (O F F) に設定する (S 7 2 9)。

【 0 2 6 7 】

第 7 実施形態

本実施形態のストレージシステムは、ホストデータのキャッシングによるアクセス性能の向上に基づき、低ヒット率フラグの O N / O F F を決定する。低ヒット率フラグは第 1 実施形態で説明した通りである。アクセス性能は、例えば、レスポンスタイムやスループットで表される。以下に説明する構成は、レスポンスタイムを使用する。

【 0 2 6 8 】

データキャッシングの使用によるレスポンスタイムの向上が大きい場合に低ヒット率フラグ (第 1 実施形態参照) は O F F に設定され、データキャッシングの使用によるレスポンスタイムの向上が小さい場合に低ヒット率フラグは O N に設定される。これにより、平均レスポンスタイムを向上することができる。

【 0 2 6 9 】

以下において、本実施形態を具体的に説明する。主に上記他の実施形態と異なる点を説明する。図 6 2 は、本実施形態のローカルメモリ 1 2 2 に格納されている制御情報を示している。レスポンステーブル 4 1 0 及び C M 利用閾値テーブル 4 2 0 がローカルメモリ 1 2 2 に格納されている。図 6 3 はレスポンステーブル 4 1 0 の一例を示し、図 6 4 は C M 利用閾値テーブル 4 2 0 の一例を示す。

【 0 2 7 0 】

レスポンステーブル 4 1 0 は、メディアのレスポンスタイムを管理するテーブルである。図 6 3 において、レスポンステーブル 4 1 0 は、メディア種別のカラム 4 1 1 及びレスポンスタイムのカラム 4 1 2 を有する。本例のレスポンステーブル 4 1 0 は、メディア種別によりレスポンスタイムを管理するが、R A I D グループや論理ボリュームによりレスポンスタイムを管理してもよい。

【 0 2 7 1 】

本例において、レスポンスタイムは、メディアからデータを読み出すために要する時間である。レスポンスタイムのカラム 4 1 2 には、予め値が格納されている、又は、マイクロプロセッサ 1 2 1 は、レスポンスタイムのカラム 4 1 2 の値を更新してもよい。マイクロプロセッサ 1 2 1 は、データ読み出しにおけるレスポンスタイムを測定し、例えば測定値の平均値をレスポンスタイムのカラム 4 1 2 に格納する。

【 0 2 7 2 】

図 6 4 において、C M 利用閾値テーブル 4 2 0 は、レスポンス向上のカラム 4 2 1 において、レスポンス向上を示す値の閾値を格納している。閾値は予め設定されている。例えば、管理者により設定された値が、ストレージシステム内の不揮発性記憶領域に格納されている。

【 0 2 7 3 】

10

20

30

40

50

後述するように、マイクロプロセッサ121は、メディアのレスポンスタイムとCMPK130（キャッシュメモリ131）のレスポンスタイムとの差を使用して、レスポンス向上を表す値を算出する。この値が上記閾値より大きい場合、レスポンス向上がデータキャッシングに見合うレベルにあることを示す。

【0274】

図65のフローチャートを参照して、本実施形態のレスポンス向上に基づく低ヒット率フラグ更新を含むヒット率更新処理を説明する。MPPK120は、定期的に、例えば、1秒毎にこの処理を実行する。図65のフローチャートにおけるステップS802、S803、S805～S807は、それぞれ、図23のフローチャートにおけるステップS222、S223、S225～S227と同様である。

10

【0275】

ステップS804において、マイクロプロセッサ121は、下記の式1に従って、レスポンス向上を表す値を算出する。

$$\text{ヒット率} \times (\text{当該メディアのレスポンスタイム} - \text{CMPKレスポンスタイム}) / 100$$

【0276】

マイクロプロセッサ121は、当該ボリュームのRAIDグループから、メディア種別テーブル230を参照して、当該メディアの種別を特定することができる。レスポンスタイムの値は、上述のように、レスポンステーブル410に格納されている。マイクロプロセッサ121は、算出した値とCM利用閾値テーブル420のCM利用閾値とを比較する。

20

【0277】

算出した値がCM利用閾値以下である場合（S804：YES）、マイクロプロセッサ121は、当該ボリュームの低ヒット率フラグを1（ON）に設定する（S805）。算出した値がCM利用閾値より大きい場合（S804：NO）、マイクロプロセッサ121は、当該ボリュームの低ヒット率フラグを0（OFF）に設定する（S806）。

【0278】

以上、本発明の実施形態を説明したが、本発明が上記の実施形態に限定されるものではない。当業者であれば、上記の実施形態の各要素を、本発明の範囲において容易に変更、追加、変換することが可能である。ある実施形態の構成の一部を他の実施形態の構成に置き換えることが可能であり、ある実施形態の構成に他の実施形態の構成を加えることも可能である。各実施形態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。

30

【0279】

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されたハード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、不揮発性半導体メモリ、ハードディスクドライブ、SSD等の記憶デバイス、または、ICカード、SDカード、DVD等の計算機読み取り可能な非一時的データ記憶媒体に格納することができる。

【0280】

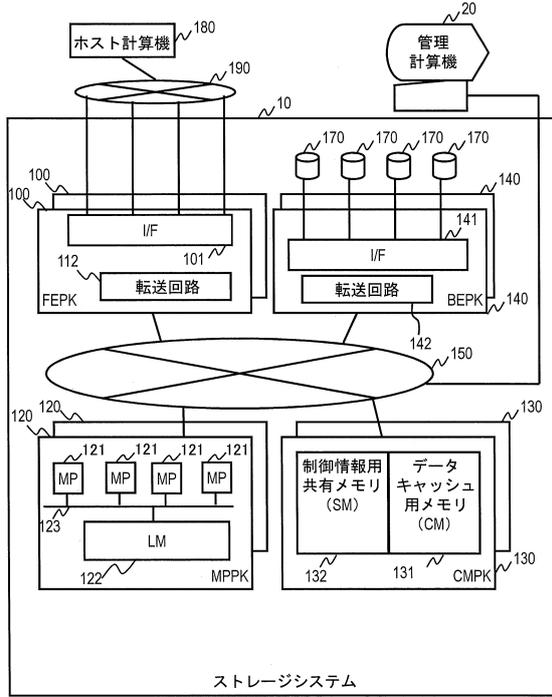
上記実施形態において、制御情報は複数のテーブルにより表されているが、本発明が使用する制御情報は、データ構造に依存しない。制御情報は、テーブルの他、例えば、データベース、リスト、キュー等のデータ構造で表現することができる。上記実施形態において、識別子、名、ID等の表現は、互いに置換が可能である。

40

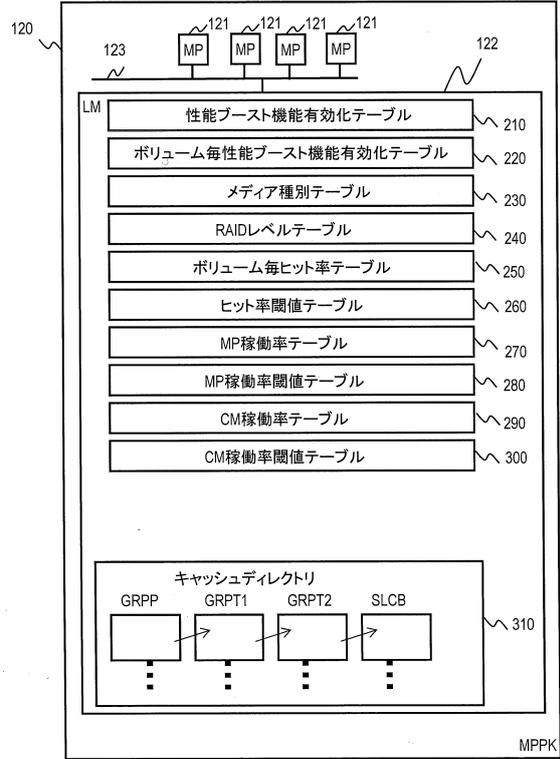
【0281】

プロセッサである、CPU、マイクロプロセッサ又は複数のマイクロプロセッサのグループは、プログラムに従って動作することで、定められた処理を実行する。従って、本実施形態においてプロセッサを主語とする説明は、プログラムを主語とした説明でもよく、プロセッサが実行する処理は、そのプロセッサが実装された装置及びシステムが行う処理である。

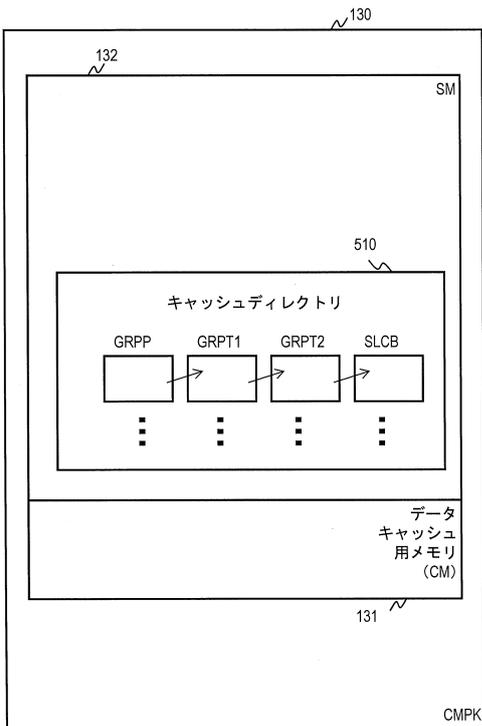
【図1】



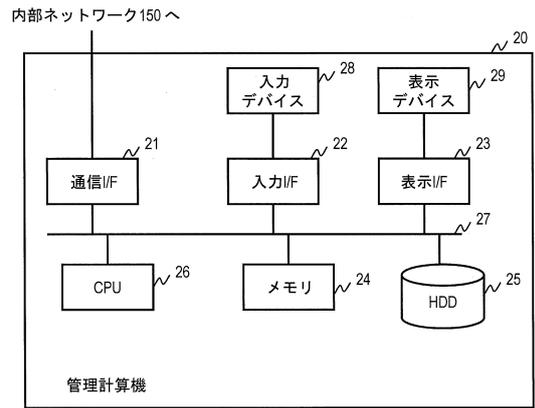
【図2】



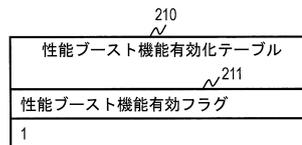
【図3】



【図4】



【図5】



【図 6】

220
ボリューム毎性能ブースト機能有効化テーブル

221 論理ボリューム 番号	222 性能ブースト機能 有効化フラグ
1	1
2	0
3	0
...	...

【図 8】

240
RAIDレベルテーブル

241 RAIDグループ番号	242 RAIDレベル
1	5
2	6
3	1
...	...

【図 7】

230
メディア種別テーブル

231 RAIDグループ番号	232 メディア種別
1	SSD
2	SATA
3	SAS
...	...

【図 9】

250
ボリューム毎ヒット率テーブル

251 論理ボリューム 番号	252 ヒット率 [%]	253 I/O数	254 ヒット数	255 低ヒット率 フラグ
1	1	135	2	1
2	20	355	40	0
3	80	1000	800	0
...

【図 10】

260
ヒット率閾値テーブル

261 ヒット率閾値[%]
5

【図 11】

270
MP稼働率テーブル

271 マイクロ プロセッサ 番号	272 稼働率[%]	273 過負荷判定フラグ	274 稼働時間[ms]
1	65	1	650
2	20	0	200
...

【図 14】

300
CM稼働率閾値テーブル

301 閾値CMPK稼働率[%]
80

【図 12】

280
MP稼働率閾値テーブル

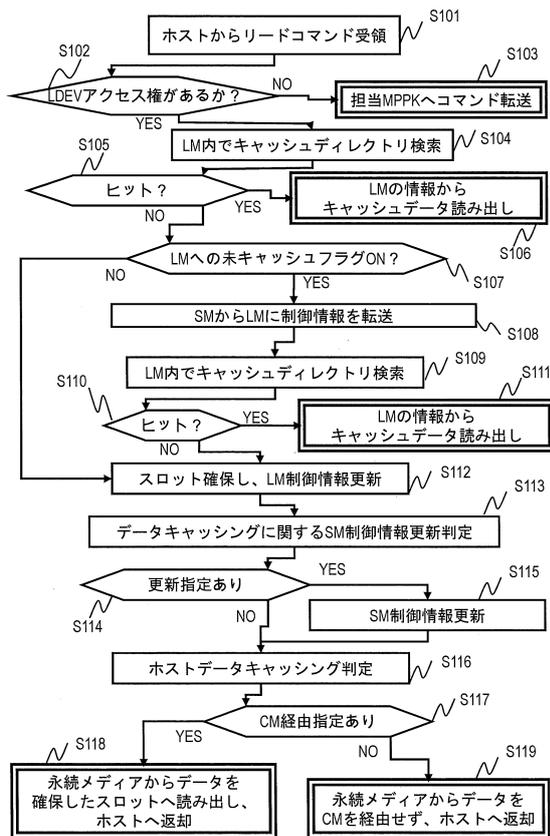
281 マイクロプロセッサ稼働率閾値[%]
80

【図 13】

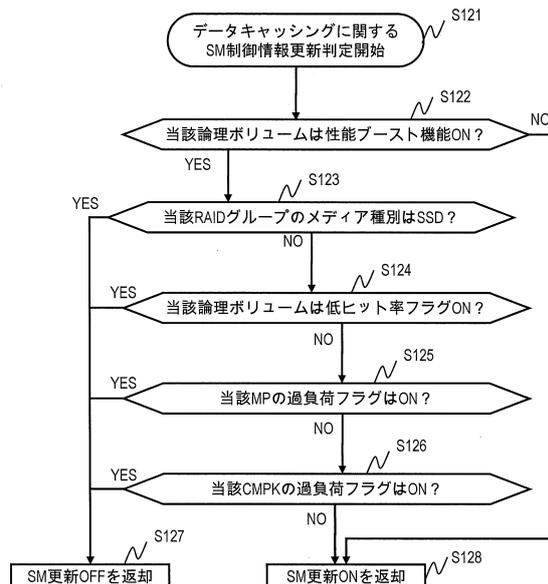
290
CM稼働率テーブル

291 CMPK番号	292 稼働率[%]	293 過負荷判定フラグ
1	44	0
2	84	1
...

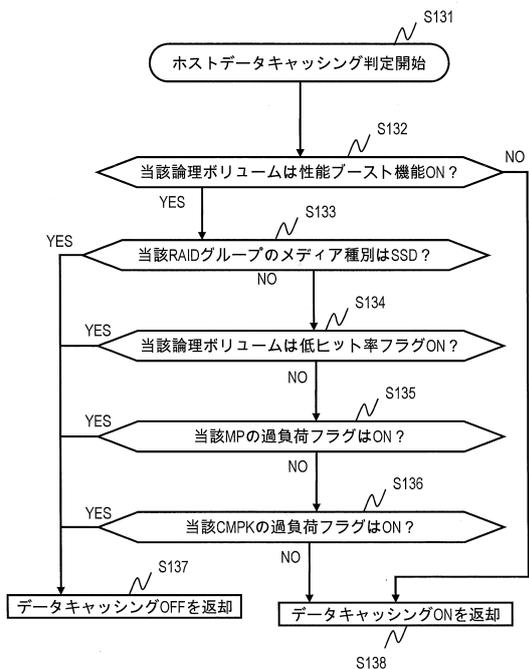
【図15】



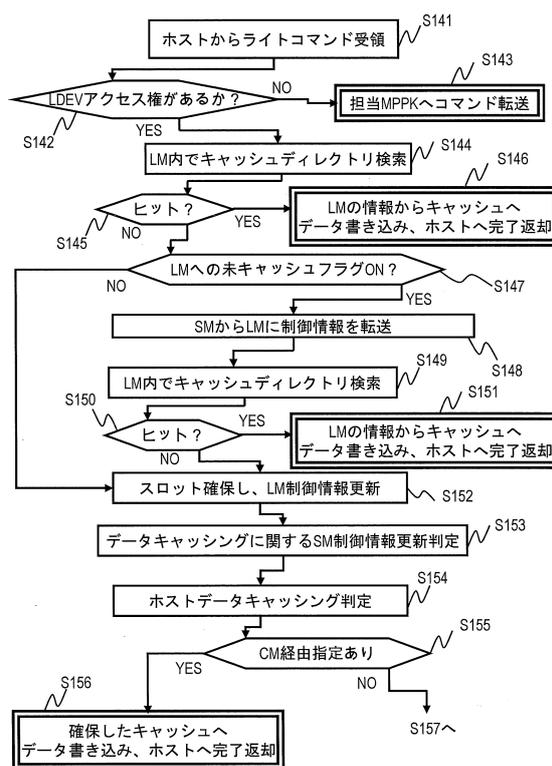
【図16】



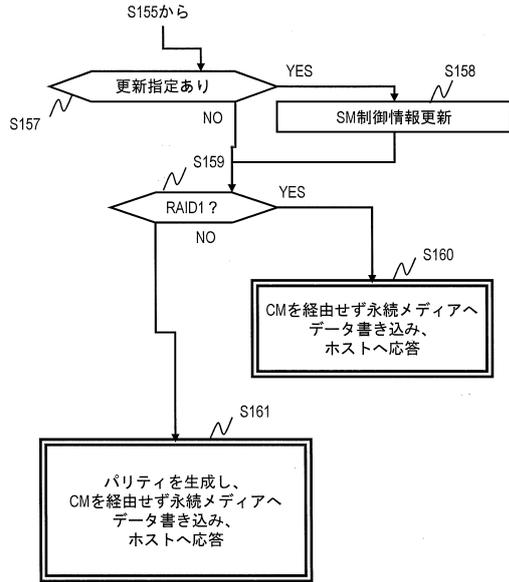
【図17】



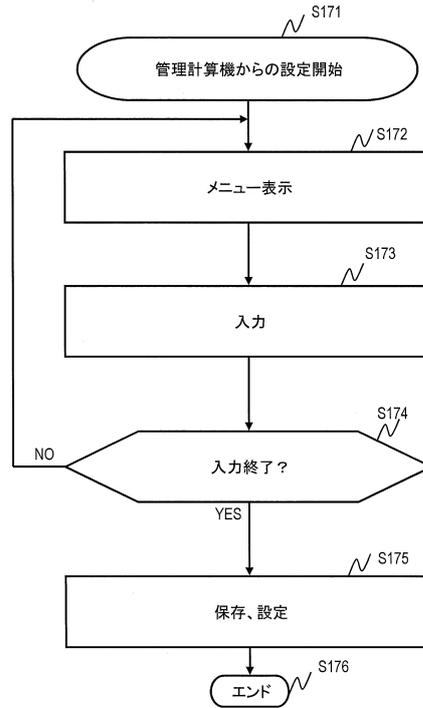
【図18A】



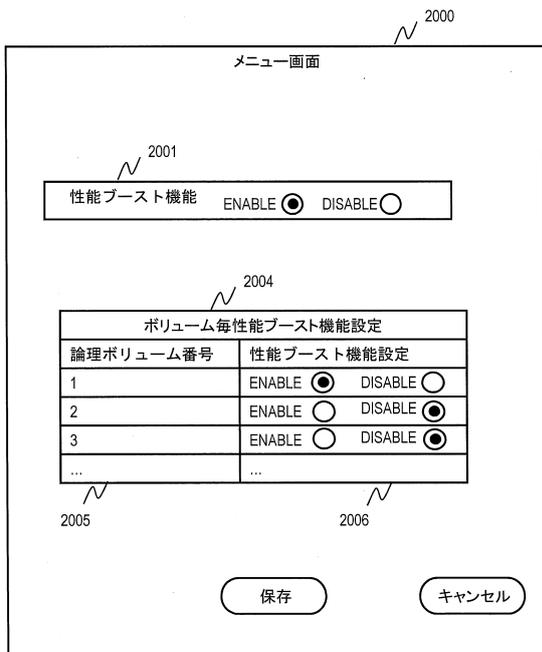
【図18B】



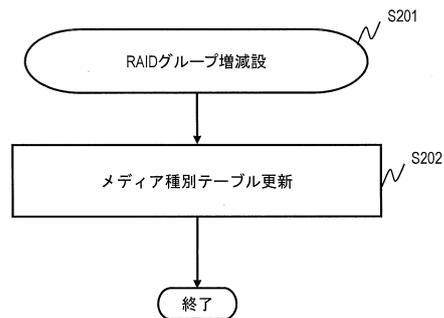
【図19】



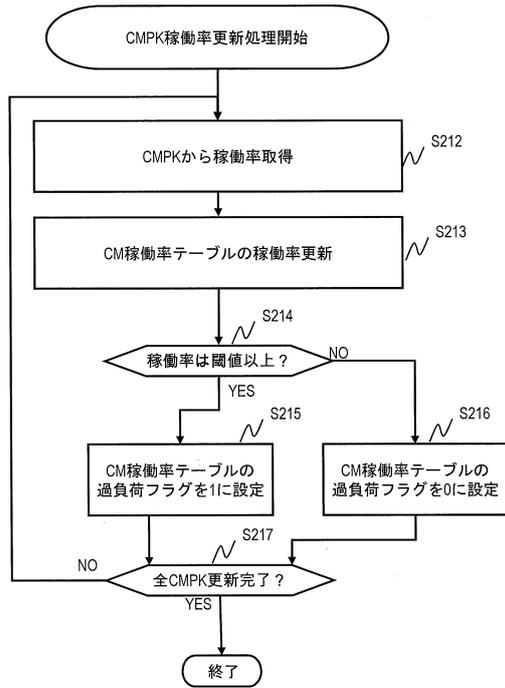
【図20】



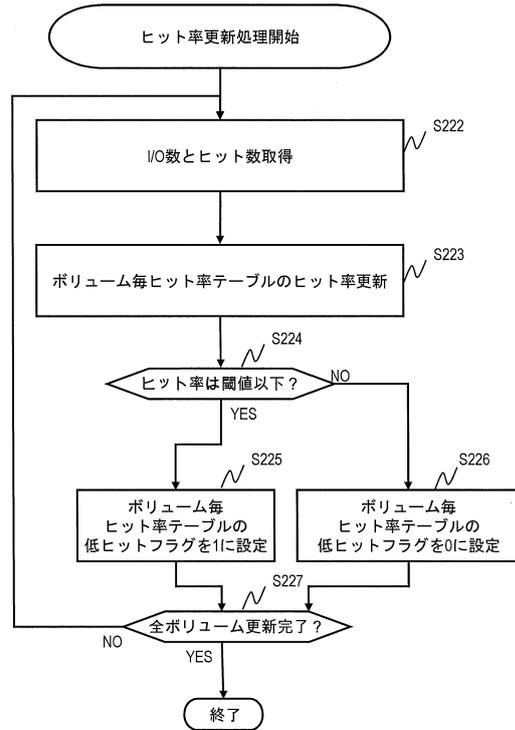
【図21】



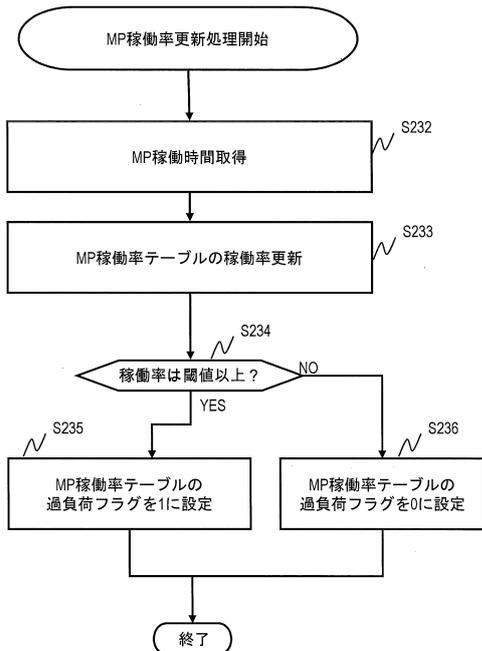
【図22】



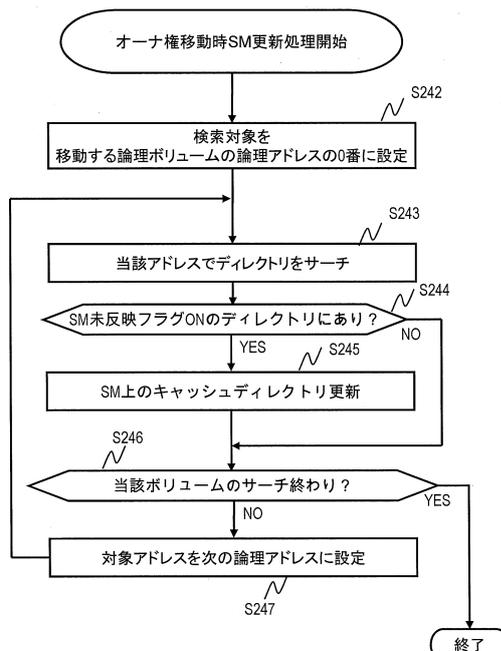
【図23】



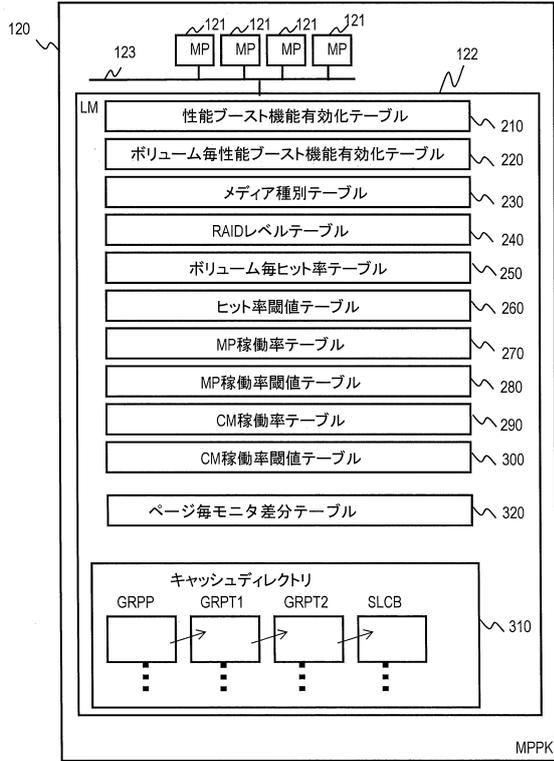
【図24】



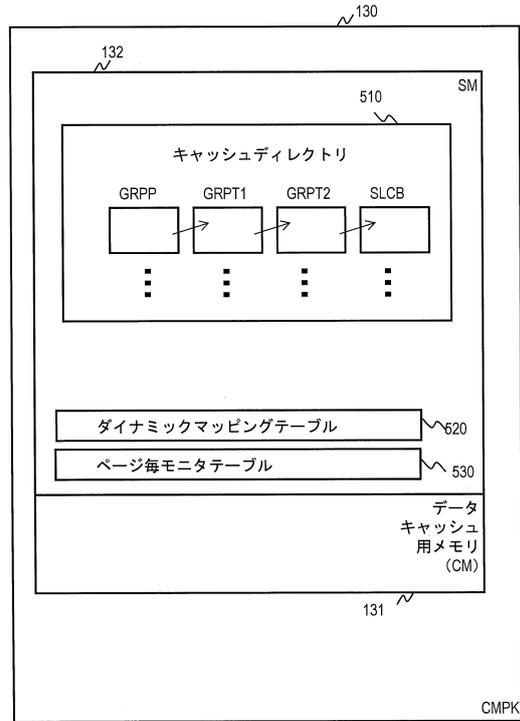
【図25】



【図 26】



【図 27】



【図 28】

プール番号	仮想ボリューム番号	論理アドレス	プールボリューム番号	論理アドレス	モニタ情報インデックス番号
1	1	0x0000	204	0x0040	1
1	1	0x0010	201	0x0050	2
1	2	0x0000	203	0x0020	3
1	2	0x0030	202	0x0040	4
1	2	0x0020	201	0x0010	5

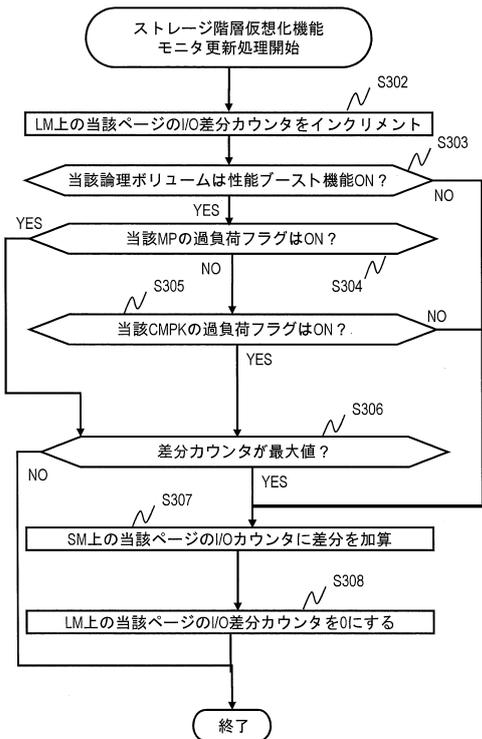
【図 30】

モニタ情報インデックス番号	I/O差分カウンタ (8bit)
1	56
2	61
3	23
4	35
5	2

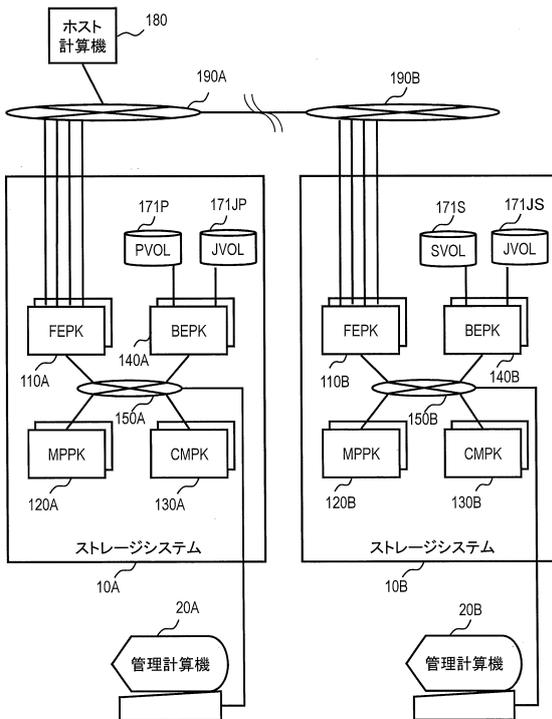
【図 29】

モニタ情報インデックス番号	I/Oカウンタ 現在 (32bit)	I/Oカウンタ 前回 (32bit)
1	1000	2000
2	613	12124
3	232	123
4	35	32
5	2	325

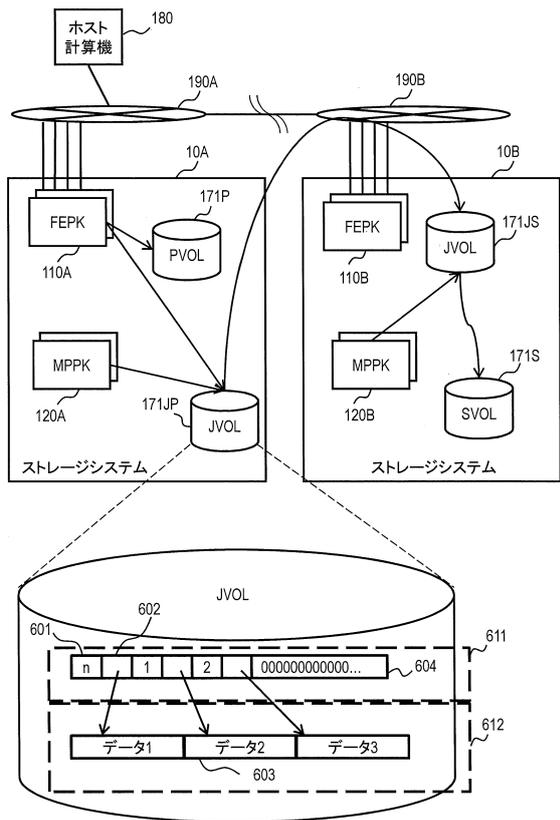
【図 3 1】



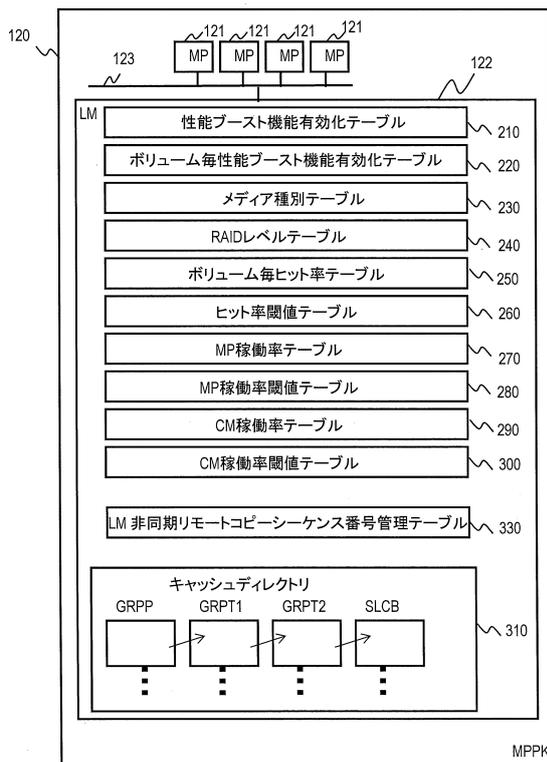
【図 3 2】



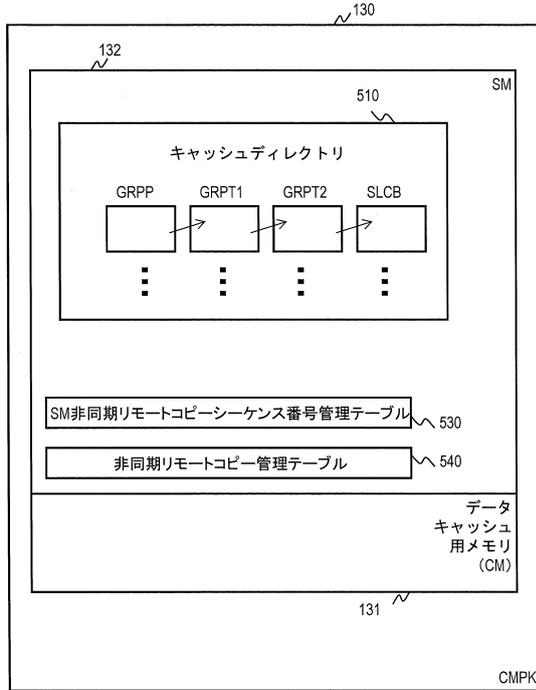
【図 3 3】



【図 3 4】



【図35】



【図36】

LM非同期リモートコピーシーケンス番号管理テーブル

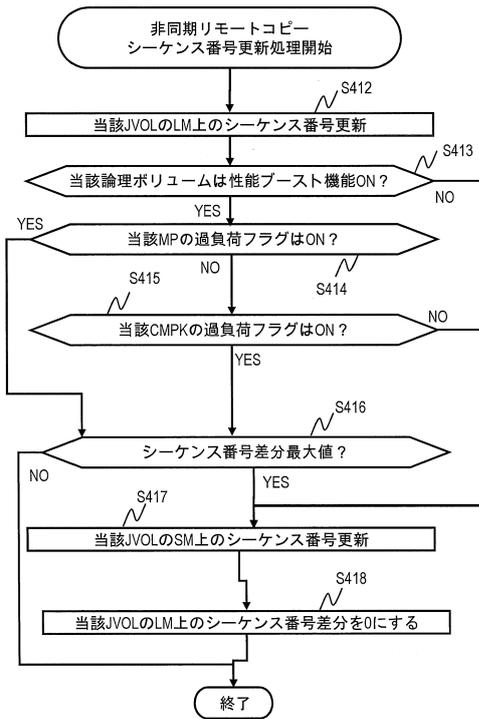
JVOL番号	シーケンス番号	シーケンス番号差分
1	312	56
2	422	61
3	546	23
4	470	35
...

【図37】

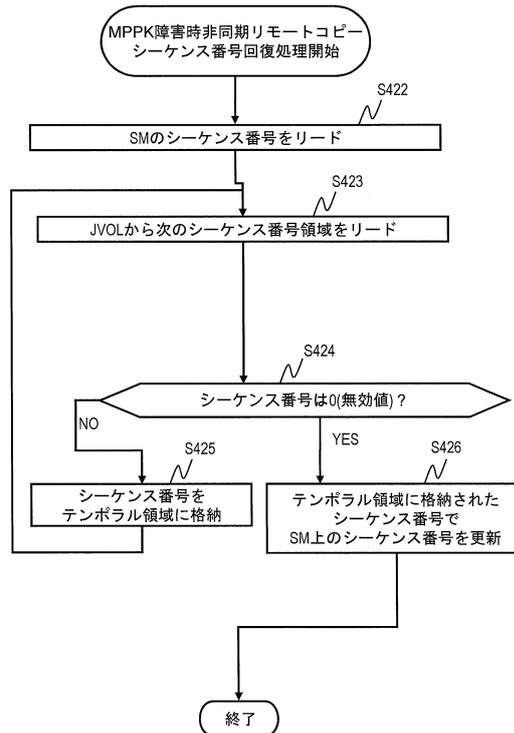
SM非同期リモートコピーシーケンス番号管理テーブル

JVOL番号	シーケンス番号
1	256
2	361
3	523
4	435
...	...

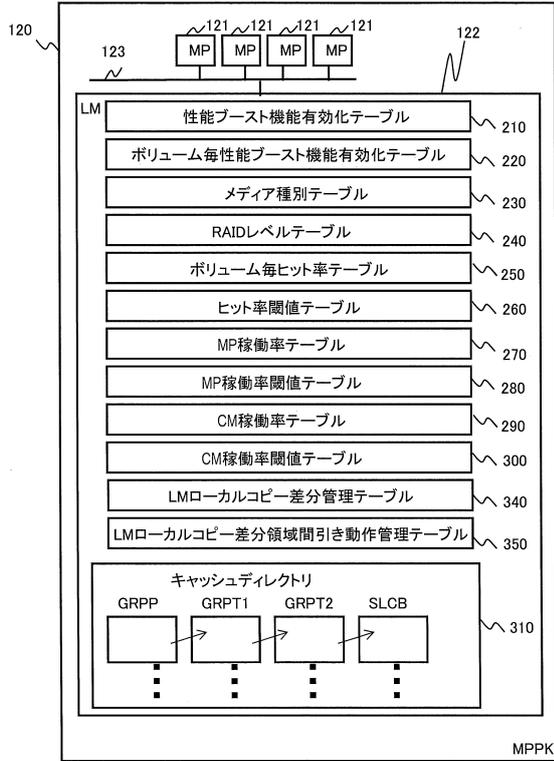
【図38】



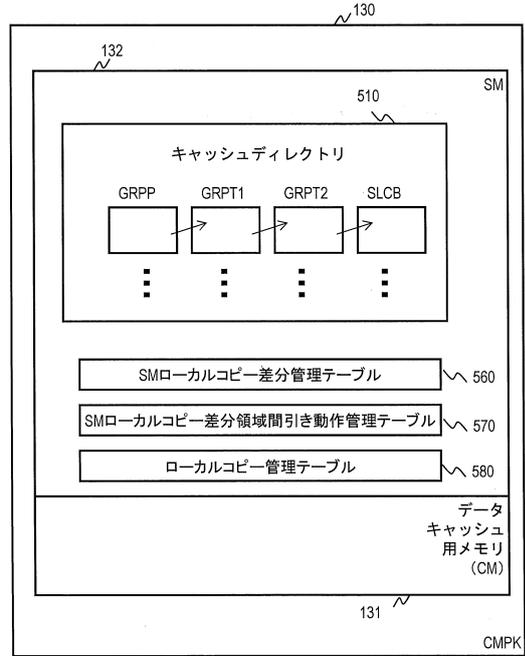
【図39】



【図40】



【図41】



【図42】

340

LMローカルコピー差分管理テーブル		
341	342	343
ボリューム番号	論理アドレス	差分有ビット列
1	0x0000	01101111
1	0x0008	00000000
1
1	0xffff	00000000
2	0x0000	11110110
...

【図44】

350

LMローカルコピー差分領域間引き動作管理テーブル		
351	352	353
ボリューム番号	論理アドレス	間引き中ビット列
1	0x0000	01101111
1	0x0040	00000000
1
1	0xffff	00000000
2	0x0000	11110110
...

【図43】

560

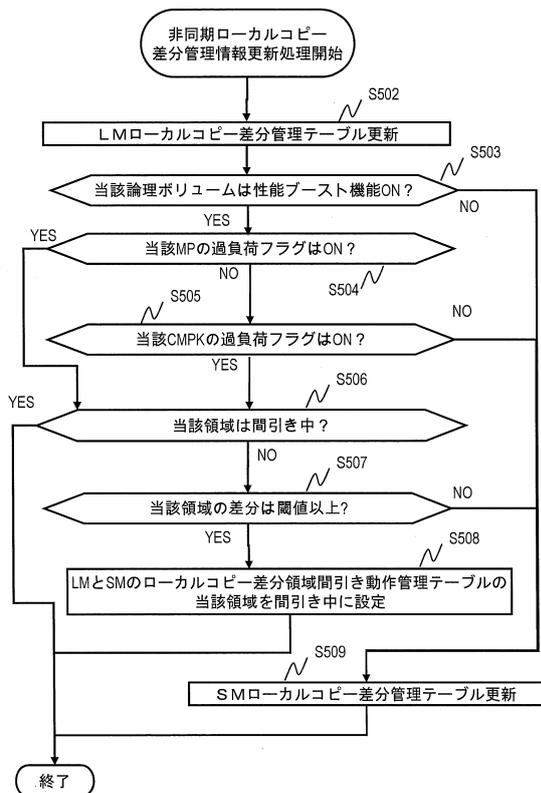
SMローカルコピー差分管理テーブル		
561	562	563
ボリューム番号	論理アドレス	差分有ビット列
1	0x0000	01101111
1	0x0008	00000000
1
1	0xffff	00000000
2	0x0000	11110110
...

【図45】

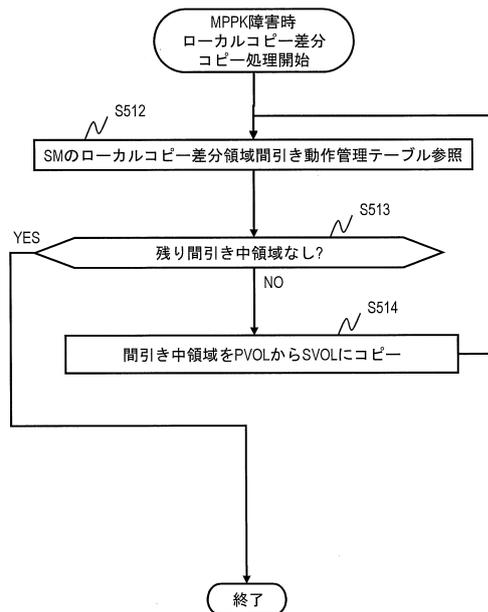
570

SMローカルコピー差分領域間引き動作管理テーブル		
571	572	573
ボリューム番号	論理アドレス	間引き中ビット列
1	0x0000	01101111
1	0x0040	00000000
1
1	0xffff	00000000
2	0x0000	11110110
...

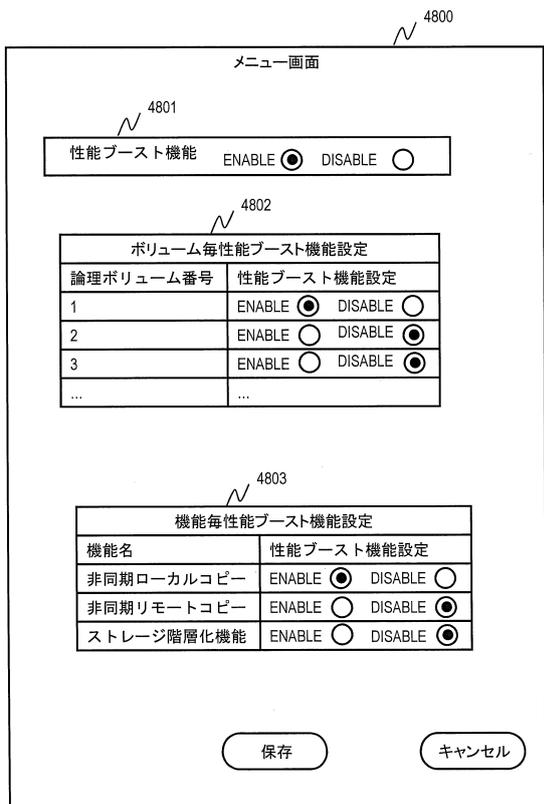
【図46】



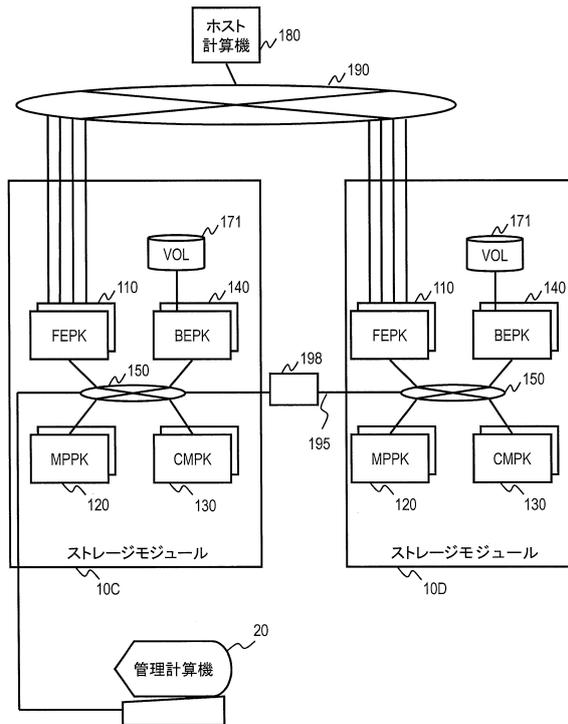
【図47】



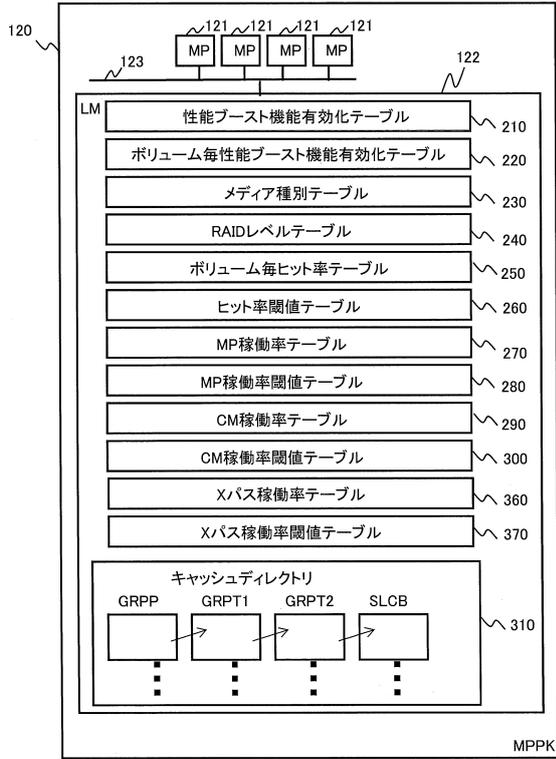
【図48】



【図49】



【図50】



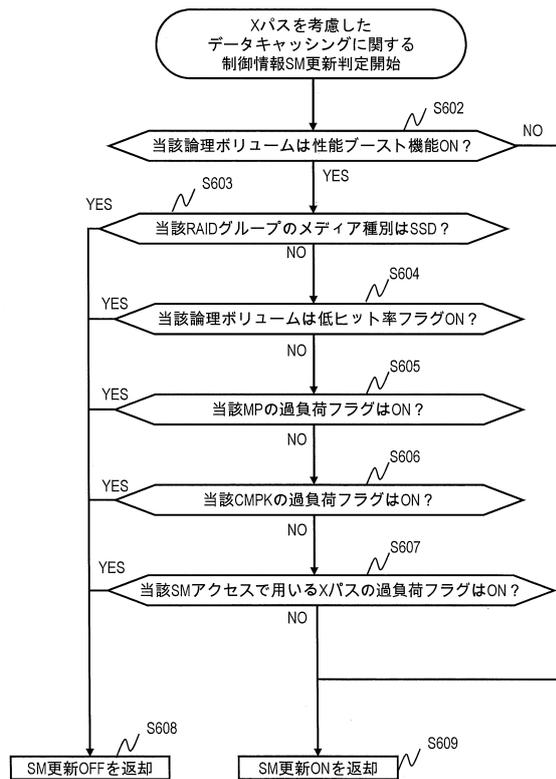
【図51】

Xバス稼働率テーブル		
Xバス番号	稼働率[%]	過負荷判定フラグ
1	44	0
2	84	1
...

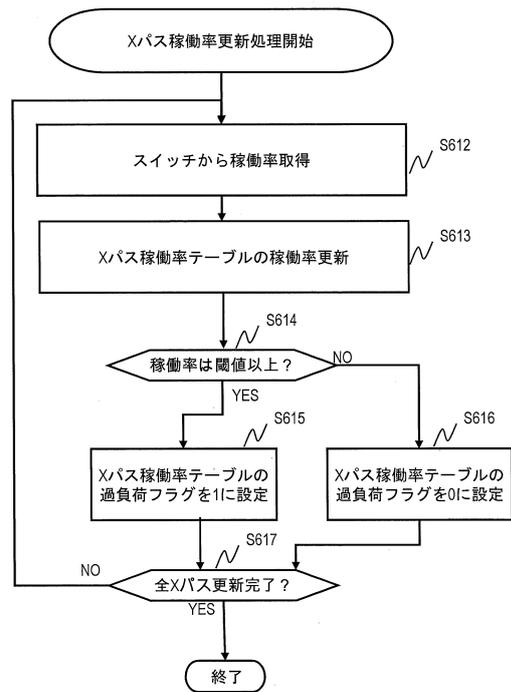
【図52】

Xバス稼働率閾値テーブル
Xバス稼働率閾値[%]
80

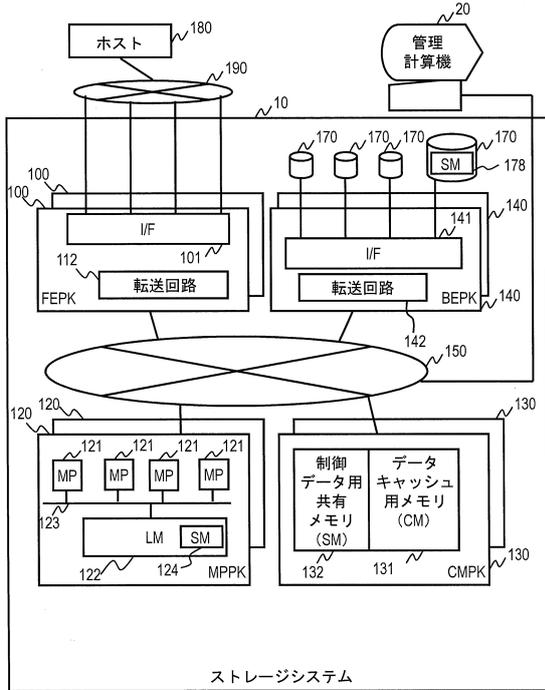
【図53】



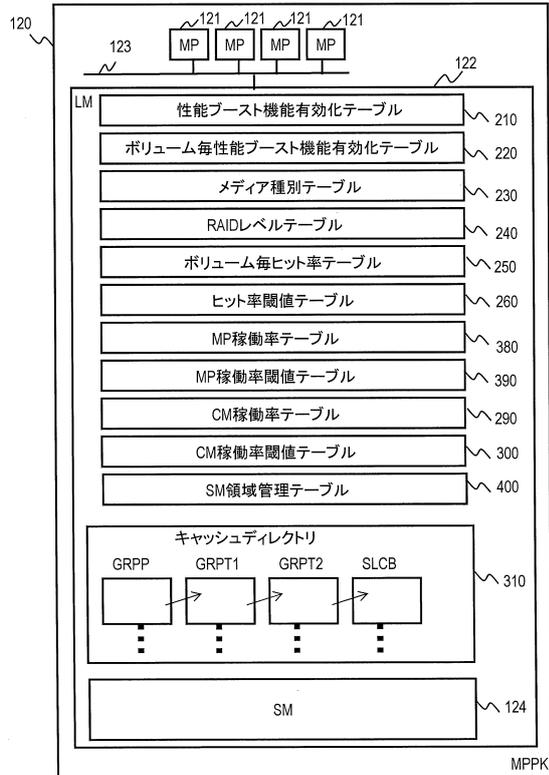
【図54】



【図55】



【図56】



【図57】

MP番号 (381)	稼働率 [%] (382)	過負荷判定フラグ1 (383)	過負荷判定フラグ2 (384)	稼働時間 [ms] (385)
1	85	1	1	650
2	40	0	0	200
3	68	0	1	340
...

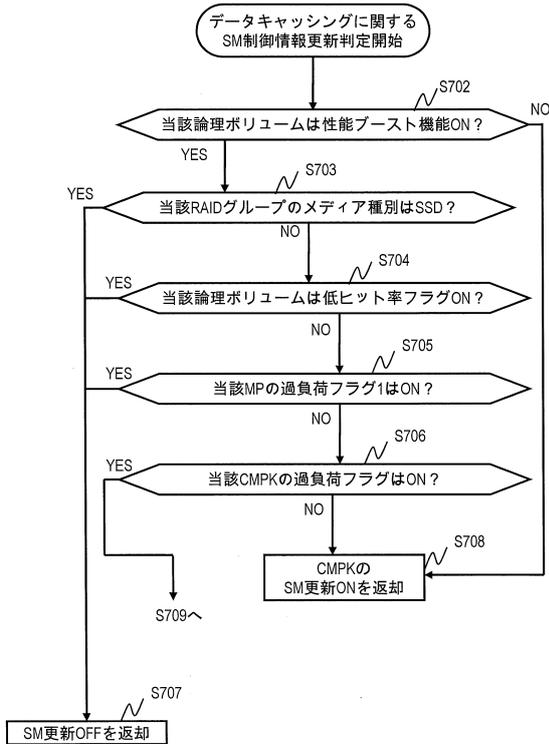
【図59】

種別 (401)	番号 (402)	先頭アドレス (403)	空き容量 (404)
MPPK	0	0x40000000	1
MPPK	1	0x40000000	0
...	0
MPPK	16	0x40000000	3
SSD	0	0x80000000	10
SSD	1	0x80000000	10
...
SSD	4	0x80000000	10
HDD	0	0x80000000	10
HDD	1	0x80000000	10
...
HDD	128	0x80000000	10

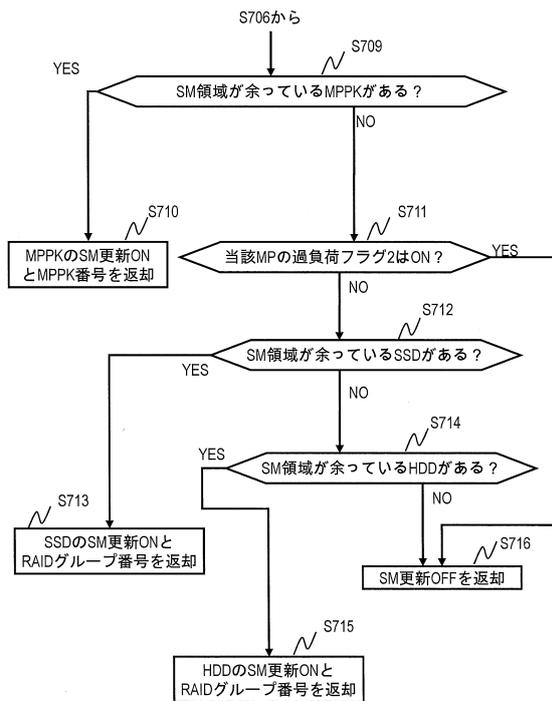
【図58】

MP稼働率閾値1[%] (391)	MP稼働率閾値2[%] (392)
80	60

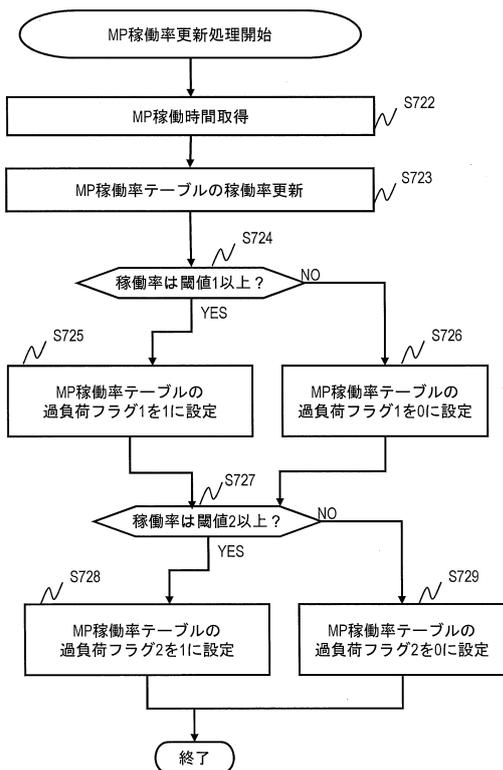
【図60A】



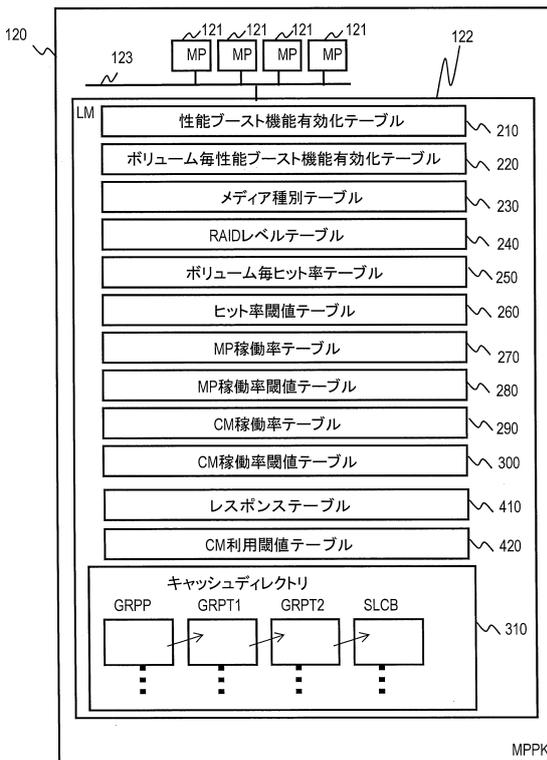
【図60B】



【図61】



【図62】



【図63】

410

レスポンステーブル	
種別	レスポンスタイム
CMPK	0.1ms
SSD	1ms
SATA 7.2krpm	10ms
SAS 15krpm	5ms
...	...

411 412

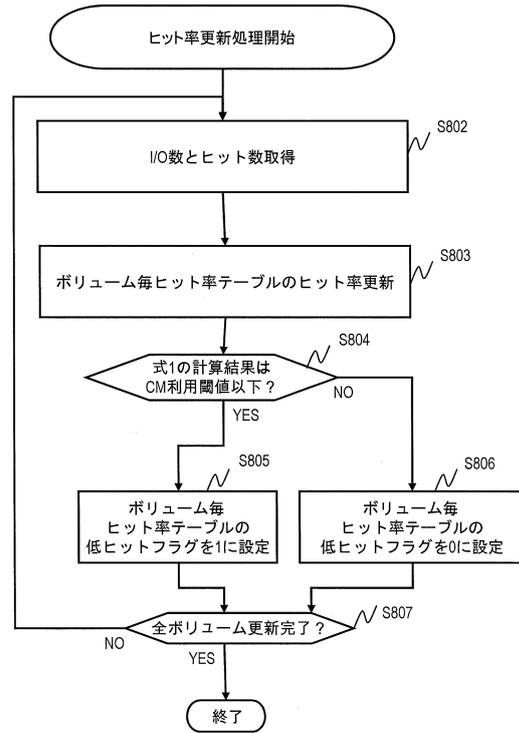
【図64】

420

CM利用閾値テーブル
レスポンス向上
0.8ms

421

【図65】



フロントページの続き

審査官 宮久保 博幸

- (56)参考文献 特開2003-228461(JP,A)
特開2005-196673(JP,A)
特開2007-079958(JP,A)
米国特許出願公開第2011/0153954(US,A1)
特開2010-086211(JP,A)
特表2012-504792(JP,A)

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