

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4124348号
(P4124348)

(45) 発行日 平成20年7月23日(2008.7.23)

(24) 登録日 平成20年5月16日(2008.5.16)

(51) Int.Cl.	F I		
G06F 12/00 (2006.01)	G06F 12/00	531J	
G06F 3/06 (2006.01)	G06F 12/00	531D	
G06F 12/08 (2006.01)	G06F 3/06	304F	
	G06F 3/06	304P	
	G06F 12/08	551Z	
請求項の数 21 (全 33 頁) 最終頁に続く			

(21) 出願番号 特願2003-183734 (P2003-183734)
 (22) 出願日 平成15年6月27日(2003.6.27)
 (65) 公開番号 特開2005-18506 (P2005-18506A)
 (43) 公開日 平成17年1月20日(2005.1.20)
 審査請求日 平成18年6月12日(2006.6.12)

早期審査対象出願

(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
 (74) 代理人 110000279
 特許業務法人ウィルフォート国際特許事務所
 (72) 発明者 平川 裕介
 神奈川県小田原市中里322番地2号 株式会社日立製作所 RAIDシステム事業部内
 (72) 発明者 荒川 敬史
 神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究所内

最終頁に続く

(54) 【発明の名称】 記憶システム

(57) 【特許請求の範囲】

【請求項1】

ホストコンピュータに接続される正ストレージシステムに接続可能であり、複数のディスクドライブを備える副ストレージシステムであって、

少なくとも一部の前記複数のディスクドライブに対応付けられており、前記ホストコンピュータから前記正ストレージシステムの複数の第1の正論理ボリュームへ送信される更新データと対応するデータを含み且つ前記正ストレージシステムから受信されるジャーナルデータを、複数記憶する少なくとも一つの第1のジャーナル論理ボリュームと、

少なくとも一部の前記複数のディスクドライブに対応付けられており、前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータから取得される複数のデータを記憶する複数の第1の副論理ボリュームと、

少なくとも一部の前記複数のディスクドライブに対応付けられており、前記ホストコンピュータから前記正ストレージシステムの複数の第2の正論理ボリュームへ送信される更新データと対応するデータを含み且つ前記正ストレージシステムから受信されるジャーナルデータを、複数記憶する少なくとも一つの第2のジャーナル論理ボリュームと、

少なくとも一部の前記複数のディスクドライブに対応付けられており、前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータから取得される複数のデータを記憶する複数の第2の副論理ボリュームと、

前記複数の第1の副論理ボリューム及び前記少なくとも一つの第1のジャーナル論理ボリュームが割当てられ、前記複数の第1の副論理ボリューム間でデータ整合性が維持され

10

20

る第1のグループと、前記複数の第2の副論理ボリューム及び前記少なくとも一つの第2のジャーナル論理ボリュームが割当てられ、前記複数の第2の副論理ボリューム間でデータ整合性が維持される第2のグループと、を管理する制御部と、
を有することを特徴とする副ストレージシステム。

【請求項2】

請求項1に記載の副ストレージシステムであって、
前記第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第1のグループにおける更新順序に関する情報を含み、
前記第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第2のグループにおける更新順序に関する情報を含む
ことを特徴とする副ストレージシステム。

10

【請求項3】

請求項1又は2に記載の副ストレージシステムであって、
前記少なくとも一つの第1のジャーナル論理ボリュームおよび前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、時刻情報を含む
ことを特徴とする副ストレージシステム。

【請求項4】

請求項1乃至3のいずれかに記載の副ストレージシステムであって、
前記第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータは、前記第1のグループの識別子を含み、
前記第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータは、前記第2のグループの識別子を含む
ことを特徴とする副ストレージシステム。

20

【請求項5】

請求項1乃至4のいずれかに記載の副ストレージシステムであって、
前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータが、前記複数の第1の副論理ボリュームのうちの少なくとも一つに書込まれ、
前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータが、前記複数の第2の副論理ボリュームのうちの少なくとも一つに書込まれる
ことを特徴とする副ストレージシステム。

30

【請求項6】

請求項5に記載の副ストレージシステムであって、
前記少なくとも一つの第1のジャーナル論理ボリュームの一部は、前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される最古のジャーナルデータから取得されるデータが前記複数の第1の副論理ボリュームのうちの少なくとも一つに書込まれた後に解放され、
前記少なくとも一つの第2のジャーナル論理ボリュームの一部は、前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される最古のジャーナルデータから取得されるデータが前記複数の第2の副論理ボリュームのうちの少なくとも一つに書込まれた後に解放される
ことを特徴とする副ストレージシステム。

40

【請求項7】

請求項1乃至6のいずれかに記載の副ストレージシステムであって、
(1) 前記副ストレージシステムから受信され前記制御部により前記少なくとも一つの第1のジャーナル論理ボリューム又は前記少なくとも一つの第2のジャーナル論理ボリュームに書き込まれるジャーナルデータ、及び、
(2) 前記少なくとも一つの第1のジャーナル論理ボリューム又は前記少なくとも一つの

50

第2のジャーナル論理ボリュームに記憶されているジャーナルデータ内のデータを前記複数の第1の副論理ボリュームのうち少なくとも一つ又は前記複数の第2の副論理ボリュームのうち少なくとも一つに書き込むために、前記制御部により、前記少なくとも一つの第1のジャーナル論理ボリューム又は前記少なくとも一つの第2のジャーナル論理ボリュームから読み出されたジャーナルデータ、を一時的に記憶するキャッシュメモリを更に備える、ことを特徴とする副ストレージシステム。

【請求項8】

ホストコンピュータに接続される正ストレージシステムに接続可能であり、複数のディスクドライブを備える副ストレージシステムの制御部であって、

10

前記ホストコンピュータから前記正ストレージシステムの複数の第1の正論理ボリュームへ送信される更新データと対応するデータを含み且つ前記正ストレージシステムから受信されるジャーナルデータの複数を、少なくとも一部の前記複数のディスクドライブに対応付けられる少なくとも一つの第1のジャーナル論理ボリュームに記憶させるように制御するものであり、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータから取得される複数のデータを、少なくとも一部の前記複数のディスクドライブに対応付けられる複数の第1の副論理ボリュームに記憶させるように制御するものであり、

前記ホストコンピュータから前記正ストレージシステムの複数の第2の正論理ボリュームへ送信される更新データと対応するデータを含み且つ前記正ストレージシステムから受信されるジャーナルデータの複数を、少なくとも一部の前記複数のディスクドライブに対応付けられる少なくとも一つの第2のジャーナル論理ボリュームに記憶させるように制御するものであり、

20

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータから取得される複数のデータを、少なくとも一部の前記複数のディスクドライブに対応付けられる複数の第2の副論理ボリュームに記憶させるように制御するものであり、

前記複数の第1の副論理ボリューム及び前記少なくとも一つの第1のジャーナル論理ボリュームが割当てられ、前記複数の第1の副論理ボリューム間でデータ整合性が維持される第1のグループと、前記複数の第2の副論理ボリューム及び前記少なくとも一つの第2のジャーナル論理ボリュームが割当てられ、前記複数の第2の副論理ボリューム間でデータ整合性が維持される第2のグループと、を管理するものである

30

ことを特徴とする副ストレージシステムの制御部。

【請求項9】

請求項8に記載の副ストレージシステムの制御部であって、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータは、前記第1のグループの識別子を含み、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータは、前記第2のグループの識別子を含む

ことを特徴とする副ストレージシステムの制御部。

【請求項10】

40

請求項8又は9に記載の副ストレージシステムの制御部であって、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第1のグループにおける更新順序に関する情報を含み、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第2のグループにおける更新順序に関する情報を含む

ことを特徴とする副ストレージシステムの制御部。

【請求項11】

請求項8乃至10のいずれかに記載の副ストレージシステムの制御部であって、

前記少なくとも一つの第1のジャーナル論理ボリュームおよび前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、時刻情報

50

を含む

ことを特徴とする副ストレージシステムの制御部。

【請求項 1 2】

請求項 8 乃至 1 1 のいずれかに記載の副ストレージシステムの制御部であって、

前記少なくとも一つの第 1 のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを、前記複数の第 1 の副論理ボリュームのうちの少なくとも一つに記憶させるように制御するものであり、

前記少なくとも一つの第 2 のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを、第 2 の副論理ボリュームのうちの少なくとも一つに記憶させるように制御するものである

10

ことを特徴とする副ストレージシステムの制御部。

【請求項 1 3】

請求項 1 2 に記載の副ストレージシステムの制御部であって、

前記少なくとも一つの第 1 のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを前記複数の第 1 の副論理ボリュームのうちの少なくとも一つに記憶させた後に、前記少なくとも一つの第 1 のジャーナル論理ボリュームのうちの前記最古のジャーナルデータが記憶されていた記憶領域を解放するように制御するものであり、

前記少なくとも一つの第 2 のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを前記複数の第 2 の副論理ボリュームのうちの少なくとも一つに記憶させた後に、前記少なくとも一つの第 2 のジャーナル論理ボリュームのうちの前記最古のジャーナルデータが記憶されていた記憶領域を解放するように制御するものである、

20

ことを特徴とする副ストレージシステムの制御部。

【請求項 1 4】

請求項 8 乃至 1 3 のいずれかに記載の副ストレージシステムの制御部であって、

(1) 前記正ストレージシステムから受信され前記少なくとも一つの第 1 のジャーナル論理ボリューム又は前記少なくとも一つの第 2 のジャーナル論理ボリュームに書き込むジャーナルデータ、及び、

(2) 前記少なくとも一つの第 1 のジャーナル論理ボリューム又は前記少なくとも一つの第 2 のジャーナル論理ボリュームに記憶されているジャーナルデータ内のデータを前記複数の第 1 の副論理ボリュームのうちの少なくとも一つ又は前記複数の第 2 の副論理ボリュームのうちの少なくとも一つに書き込むために、前記少なくとも一つの第 1 のジャーナル論理ボリューム又は前記少なくとも一つの第 2 のジャーナル論理ボリュームから読み出したジャーナルデータ、

30

を一時的に記憶するキャッシュメモリを更に備える、

ことを特徴とする副ストレージシステムの制御部。

【請求項 1 5】

ホストコンピュータに接続される正ストレージシステムに接続可能であり、複数のディスクドライブを備える副ストレージシステムの制御部のプログラムであって、

40

前記副ストレージシステムの制御部に、

前記ホストコンピュータから前記正ストレージシステムの複数の第 1 の正論理ボリュームへ送信される更新データと対応するデータを含み且つ前記正ストレージシステムから受信されるジャーナルデータの複数を、少なくとも一部の前記複数のディスクドライブに対応付けられる少なくとも一つの第 1 のジャーナル論理ボリュームに記憶させるように制御する処理と、

前記少なくとも一つの第 1 のジャーナル論理ボリュームに記憶される複数のジャーナルデータから取得される複数のデータを、少なくとも一部の前記複数のディスクドライブに対応付けられる複数の第 1 の副論理ボリュームに記憶させるように制御する処理と、

前記ホストコンピュータから前記正ストレージシステムの複数の第 2 の正論理ボリューム

50

ムへ送信される更新データと対応するデータを含み且つ前記正ストレージシステムから受信されるジャーナルデータの複数を、少なくとも一部の前記複数のディスクドライブに対応付けられる少なくとも一つの第2のジャーナル論理ボリュームに記憶させるように制御する処理と、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータから取得される複数のデータを、少なくとも一部の前記複数のディスクドライブに対応付けられる複数の第2の副論理ボリュームに記憶させるように制御する処理と、

前記複数の第1の副論理ボリューム及び前記少なくとも一つの第1のジャーナル論理ボリュームが割当てられ、前記複数の第1の副論理ボリューム間でデータ整合性が維持される第1のグループと、前記複数の第2の副論理ボリューム及び前記少なくとも一つの第2のジャーナル論理ボリュームが割当てられ、前記複数の第2の副論理ボリューム間でデータ整合性が維持される第2のグループとを管理する処理と、

を実行させることを特徴とする副ストレージシステムの制御部のプログラム。

【請求項16】

請求項15に記載の副ストレージシステムの制御部のプログラムであって、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第1のグループの識別子を含み、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第2のグループの識別子を含む

ことを特徴とする副ストレージシステムの制御部のプログラム。

【請求項17】

請求項15又は16に記載の副ストレージシステムの制御部のプログラムであって、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第1のグループにおける更新順序に関する情報を含み、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、前記第2のグループにおける更新順序に関する情報を含む

ことを特徴とする副ストレージシステムの制御部のプログラム。

【請求項18】

請求項15乃至17のいずれかに記載の副ストレージシステムの制御部のプログラムであって、

前記少なくとも一つの第1のジャーナル論理ボリュームおよび前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータの各々は、時刻情報を含む

ことを特徴とする副ストレージシステムの制御部のプログラム。

【請求項19】

請求項15乃至18のいずれかに記載の副ストレージシステムの制御部のプログラムであって、

前記副ストレージシステムの制御部に、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを前記複数の第1の副論理ボリュームのうちの少なくとも一つに記憶させるように制御する処理と、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを前記複数の第2の副論理ボリュームのうちの少なくとも一つに記憶させるように制御する処理とを、

実行させることを特徴とする副ストレージシステムの制御部のプログラム。

【請求項20】

請求項19に記載の副ストレージシステムの制御部のプログラムであって、

前記副ストレージシステムの制御部に、

前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータから取得されるデータを前記複数の第1の副論理

10

20

30

40

50

ボリュームのうちの少なくとも一つに記憶させた後に、前記少なくとも一つの第1のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータが記憶されていた記憶領域を解放する処理と、

前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの前記最古のジャーナルデータから取得されるデータを前記複数の第2の副論理ボリュームのうちの少なくとも一つに記憶させた後に、前記少なくとも一つの第2のジャーナル論理ボリュームに記憶される複数のジャーナルデータのうちの最古のジャーナルデータが記憶されていた記憶領域を解放する処理とを、

実行させることを特徴とする副ストレージシステムの制御部のプログラム。

【請求項21】

請求項15乃至20のいずれかに記載の副ストレージシステムの制御部のプログラムであって、

前記制御部が、キャッシュメモリを備えており、

前記制御部に、以下の(1)及び(2)のジャーナルデータをキャッシュメモリに一時的に記憶させる処理、

(1)前記正ストレージシステムから受信され前記少なくとも一つの第1のジャーナル論理ボリューム又は前記少なくとも一つの第2のジャーナル論理ボリュームに書き込むジャーナルデータ、

(2)前記少なくとも一つの第1のジャーナル論理ボリューム又は前記少なくとも一つの第2のジャーナル論理ボリュームに記憶されているジャーナルデータ内のデータを前記複数の第1の副論理ボリュームのうちの少なくとも一つ又は前記複数の第2の副論理ボリュームのうちの少なくとも一つに書き込むために、前記少なくとも一つの第1のジャーナル論理ボリューム又は前記少なくとも一つの第2のジャーナル論理ボリュームから読み出したジャーナルデータ、

を実行させることを特徴とする副ストレージシステムの制御部のプログラム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は記憶システムに関し、特に複数の記憶システム間でのデータの複製に関する。

【0002】

【従来の技術】

近年、常に顧客に対して継続したサービスを提供するために、第一の記憶システムに障害が発生した場合でもデータ処理システムがサービスを提供できるよう、記憶システム間でのデータの複製に関する技術が重要になっている。第一の記憶システムに格納された情報を第二および第三の記憶システムに複製する技術として、以下の特許文献に開示された技術が存在する。

【0003】

米国特許5170480号公報には、第一の記憶システムに接続された第一の計算機が、第一の記憶システムに格納されたデータを、第一の計算機と第二の計算機間の通信リンクを介し、第二の計算機に転送し、第二の計算機が第二の計算機に接続された第二の記憶システムに転送する技術が開示されている。

【0004】

米国特許6209002号公報には、第一の記憶システムが、第一の記憶システムに格納されたデータを、第二の記憶システムに転送し、さらに、第二の記憶システムが、第三の記憶システムに転送する技術が開示されている。計算機と第一の記憶システムとは通信リンクにより接続され、第一の記憶システムと第二の記憶システムとは通信リンクにより接続され、さらに、第二の記憶システムと第三の記憶システムとは通信リンクにより接続されている。第一の記憶システムは、複製対象の第一の論理ボリュームを保持する。第二の記憶システムは、第一の論理ボリュームの複製である第二の論理ボリューム、及び、第二の論理ボリュームの複製である第三の論理ボリュームを保持する。第三の記憶システムは

10

20

30

40

50

、第三の論理ボリュームの複製である第四の論理ボリュームを保持する。この特許文献において、第二の記憶システムは、第二の論理ボリュームから第三の論理ボリュームへのデータ複製処理と、第三の論理ボリュームから第四の論理ボリュームへのデータ複製処理とを排他的に実行する。

【0005】

【特許文献1】

米国特許5170480号公報

【特許文献2】

米国特許6209002号公報

【発明が解決しようとする課題】

米国特許5170480号公報に開示された技術は、データの複製のために、第一の計算機および第二の計算機を常に使用する。第一の計算機は、通常業務を行っており、第一の計算機にかかるデータ複製処理の負荷は無視できない。さらに、複製のためのデータは、第一の計算機と第一の記憶システム間の通信リンクを使用するため、通常業務のために必要なデータ転送と衝突し、通常業務に必要なデータ参照、データ更新時間が長くなるという課題がある。

10

【0006】

米国特許6209002号公報に開示された技術は、複製を行うデータ量の倍の記憶容量が第二の記憶システム及び第三の記憶システムに必要となる。また、複製対象のデータ量が多いことにより、データ複製処理に費やされる時間が長く、第三の記憶システムのデータは古いものとなってしまふ。その結果、第三の記憶システムのデータを用いて業務が再開される場合、第三の記憶システムのデータを最新とするまでの時間が長くなり、業務再開までの時間が延びるといふ課題がある。さらに、かかる文献において、第一の記憶システムは、第一の記憶システム内のデータ更新処理に加えて第二の記憶システムとの間でのデータ行進処理が終了した時点で、上位の計算機にデータ更新完了報告を行う。したがって、計算機からのデータ更新に費やされる時間が長く、第一の記憶システムと第二の記憶システムとの間の距離が遠くなればなるほど、データ更新に費やされる時間は長くなる。その結果、かかる文献に開示された技術によれば、各記憶システム間の距離をあまり遠くすることができないという課題もある。

20

【0007】

本発明の目的は、記憶システムの上位の計算機に影響を与えず、複数の記憶システム間でデータ転送又はデータの複製をするものである。さらに、記憶システムと計算機との間の通信にも影響を与えないものである。

30

【0008】

さらに、本発明の目的は、複数の記憶システム内に保持するデータ格納領域を少なくすることができるものである。さらに、複数の記憶システムの上位の計算機の業務に影響を与えることのないように、高速かつ効率的に複数の記憶システム間でデータ転送又はデータの複製をするものである。

【0009】

【課題を解決するための手段】

これらの課題を解決するために、本発明において、第一の記憶システムは、第一の記憶システムに格納されたデータの更新に関する情報をジャーナルとして格納する。ジャーナルは、具体的には、更新に用いられたデータのコピーと更新時のライト命令等の更新情報とによって構成される。さらに、第二の記憶システムは、第一の記憶システムと第二の記憶システム間の通信線を介して、前記ジャーナルを取得する。第二の記憶システムは、第一の記憶システムが保持するデータの複製を保持しており、前記ジャーナルを用いて、第一の記憶システムでのデータ更新順に、第一の記憶システムのデータと対応するデータを更新する。

40

【0010】

本発明においては、第二の記憶システムは、第一の記憶システムが格納する第一の記憶領

50

域に格納されたデータの複製を格納する第二の記憶領域を保持し、第二の記憶領域のデータの更新に関する情報をジャーナルとしてジャーナル専用の第三の記憶領域に格納する。第三の記憶領域の記憶容量は、第二の記憶領域より少ない記憶容量とすることが可能である。さらに、第三の記憶システムは、第二の記憶システムと第三の記憶システム間の通信線を介して、前記ジャーナルを取得し、ジャーナル専用の第四の記憶領域に格納する。第四の記憶領域の記憶容量は、第二の記憶領域より少ない記憶容量とすることが可能である。第三の記憶システムは、第二の記憶領域に格納されたデータの複製を格納する第五の記憶領域を保持し、前記ジャーナルを用いて、第二の記憶領域でのデータ更新順に、第二の記憶領域に対応する第五の記憶領域のデータを更新する。

【0011】

10

【発明の実施の形態】

本発明によるデータ処理システムの実施形態を図面により詳細に説明する。

【0012】

図1は、本発明の一実施形態の論理的な構成を示すブロック図である。

【0013】

本発明の一実施形態は、ホストコンピュータ180と記憶システム100Aを接続バス190により接続し、記憶システム100Aと記憶システム100Aに保存されたデータの複製を保持する記憶システム100Bを接続バス200により接続した構成である。以下の説明において、複製対象のデータを保持する記憶システム100と複製データを保持する記憶システム100との区別を容易とするために、複製対象のデータを保持する記憶システム100を正記憶システム100A、複製データを保持する記憶システム100を副記憶システム100Bとよぶこととする。なお、記憶システムの記憶領域は、分割して管理されており、分割された記憶領域を論理ボリュームと呼ぶこととする。

20

【0014】

記憶システム100の記憶領域は、分割して管理されており、分割した記憶領域を論理ボリューム230とよぶこととする。論理ボリューム230の容量および記憶システム100内の物理的な格納位置(物理アドレス)は、記憶システム100に接続したコンピュータ等の保守端末もしくはホストコンピュータ180を用いて指定できる。各論理ボリューム230の物理アドレスは、後述するボリューム情報400に保存する。物理アドレスは、例えば、記憶システム100内の記憶装置150を識別する番号(記憶装置番号)と記憶装置内の記憶領域を一意に示す数値、例えば、記憶装置の記憶領域の先頭からの位置である。以下の説明では、物理アドレスは、記憶装置番号と記憶装置の記憶領域の先頭からの位置の組とする。以下の説明では、論理ボリュームは、1つの記憶装置の記憶領域であるが、論理アドレスと物理アドレスの変換により、1つの論理ボリュームを複数の記憶装置の記憶領域に対応づけることも可能である。

30

【0015】

記憶システム100が保存しているデータの参照、更新は、論理ボリュームを識別する番号(論理ボリューム番号)と記憶領域を一意に示す数値、例えば、論理ボリュームの記憶領域の先頭からの位置により一意に指定することができ、以下、論理ボリューム番号と論理ボリュームの記憶領域の先頭からの位置(論理アドレス内位置)の組を論理アドレスとよぶ。

40

【0016】

以下の説明において、複製対象のデータと複製データとの区別を容易とするために、複製対象の論理ボリューム230を正論理ボリューム、複製データである論理ボリューム230を副論理ボリュームとよぶこととする。一对の正論理ボリュームと副論理ボリュームをペアとよぶ。正論理ボリュームと副論理ボリュームの関係および状態等は後述するペア情報500に保存する。

【0017】

論理ボリューム間のデータの更新順序を守るために、グループという管理単位を設ける。例えば、ホストコンピュータ180が、正論理ボリューム1のデータ1を更新し、その後

50

、データ1を読み出し、データ1の数値を用いて、正論理ボリューム2のデータ2を更新する処理を行うとする。正論理ボリューム1から副論理ボリューム1へのデータ複製処理と、正論理ボリューム2から副論理ボリューム2へのデータ複製処理とが独立に行われる場合、副論理ボリューム1へのデータ1の複製処理より前に、副論理ボリューム2へのデータ2の複製処理が行われる場合がある。副論理ボリューム2へのデータ2の複製処理と副論理ボリューム1へのデータ1の複製処理との間に、故障等により副論理ボリューム1へのデータ1の複製処理が停止した場合、副論理ボリューム1と副論理ボリューム2のデータの整合性がなくなる。このような場合にも副論理ボリューム1と副論理ボリューム2のデータの整合性を保つために、データの更新順序を守る必要のある論理ボリュームは、同じグループに登録し、データの更新毎に、後述するグループ情報600の更新番号を割り当て、更新番号順に、副論理ボリュームに複製処理を行う。例えば、図1では、正記憶システム100Aの論理ボリューム(DATA1)と論理ボリューム(DATA2)がグループ1を構成する。論理ボリューム(DATA1)の副製である論理ボリューム(COPY1)と論理ボリューム(DATA2)の複製である論理ボリューム(COPY2)は、副記憶システム内でグループ1を構成する。

10

【0018】

データ複製対象である正論理ボリュームのデータを更新する場合、副論理ボリュームのデータを更新するために、後述するジャーナルを作成し、正記憶システム100A内の論理ボリュームに保存する。本実施例の説明では、グループ毎にジャーナルのみを保存する論理ボリューム(以下、ジャーナル論理ボリュームとよぶ)を割り当てる。図1では、グループ1に論理ボリューム(JNL1)を割り当てている。

20

【0019】

副記憶システム100Bのグループにもジャーナル論理ボリュームを割り当てる。ジャーナル論理ボリュームは、正記憶システム100Aから副記憶システム100Bに転送したジャーナルを保存するために使用する。ジャーナルをジャーナル論理ボリュームに保存することにより、例えば、副記憶システム100Bの負荷が高い場合、ジャーナル受信時に副論理ボリュームのデータ更新を行わず、暫く後、副記憶システム100Bの負荷が低い時に、副論理ボリュームのデータを更新することもできる。さらに、接続線200が複数ある場合、正記憶システム100Aから副記憶システム100Bへのジャーナルの転送を多重に行い、接続線200の転送能力を有効に利用することができる。更新順番のため、副記憶システム100Bに多くのジャーナルが溜まる可能性があるが、副論理ボリュームのデータ更新に直ぐに使用できないジャーナルは、ジャーナル論理ボリュームに退避することにより、キャッシュメモリを開放することができる。図1では、副記憶システム内のグループ1に論理ボリューム(JNL2)を割り当てている。

30

【0020】

ジャーナルは、ライトデータと更新情報とから構成する。更新情報は、ライトデータを管理するための情報で、ライト命令を受信した時刻、グループ番号、後述するグループ情報600の更新番号、ライト命令の論理アドレス、ライトデータのデータサイズ、ライトデータを格納したジャーナル論理ボリュームの論理アドレス等からなる。更新情報は、ライト命令を受信した時刻と更新番号のどちらか一方のみを保持してもよい。ホストコンピュータ180からのライト命令内にライト命令の作成時刻が存在する場合は、ライト命令を受信した時刻の代わりに、当該ライト命令内の作成時刻を使用してもよい。図3と図21を用いて、ジャーナルの更新情報の例を説明する。更新情報310は、1999年3月17日の22時20分10秒に受信したライト命令を記憶する。当該ライト命令は、ライトデータを論理ボリューム番号1の記憶領域の先頭から700の位置に格納する命令であり、データサイズは300である。ジャーナルのライトデータは、論理ボリューム番号4(ジャーナル論理ボリューム)の記憶領域の先頭から1500の位置から格納されている。論理ボリューム番号1の論理ボリュームはグループ1に属し、グループ1のデータ複製開始から4番目のデータ更新であることがわかる。

40

【0021】

50

ジャーナル論理ボリュームは、例えば、図3に示すように、更新情報を格納する記憶領域（更新情報領域）とライトデータを格納する記憶領域（ライトデータ領域）に分割して使用する。更新情報領域は、更新情報領域の先頭から、更新番号の順に格納し、更新情報領域の終端に達すると、更新情報領域の先頭から格納する。ライトデータ領域は、ライトデータ領域の先頭からライトデータを格納し、ライトデータ領域の終端に達すると、ライトデータ領域の先頭から格納する。更新情報領域およびライトデータ領域の比は、固定値でもよいし、保守端末あるいはホストコンピュータ180により設定可能としてもよい。これらの情報は、後述するポインタ情報700内に保持する。以下の説明では、ジャーナル論理ボリュームを更新情報とライトデータの領域に分割して、ジャーナル論理ボリュームを使用するが、論理ボリュームの先頭から、ジャーナル、つまり、更新情報とライトデータを連続に格納する方式を採用してもよい。

10

【0022】

図1を用いて、記憶システム100Aの正論理ボリュームへのデータ更新を副記憶システム100Bの副論理ボリュームに反映する動作について概説する。

【0023】

(1) 記憶システム100Aは、ホストコンピュータ180から正論理ボリューム内のデータに対するライト命令を受信すると、後述する命令受信処理210およびリードライト処理220によって、正論理ボリューム(DATA1)内のデータ更新と、ジャーナル論理ボリューム(JNL1)にジャーナルの保存を行う(図1の270)。

【0024】

(2) 記憶システム100Bは、後述するジャーナルリード処理240によって、記憶システム100Aからジャーナルをリードし、リードライト処理220によって、ジャーナル論理ボリューム(JNL2)にジャーナルを保存する(図1の280)。

20

【0025】

(3) 記憶システム100Aは、記憶システム100Bからジャーナルをリードする命令を受信すると、後述する命令受信処理210およびリードライト処理220によって、ジャーナル論理ボリューム(JNL1)からジャーナルを読み出し、記憶システム100Bに送信する(図1の280)。

【0026】

(4) 記憶システム100Bは、後述するリストア処理250およびリードライト処理220によって、ポインタ情報700を用いて、更新番号の昇順に、ジャーナル論理ボリューム(JNL2)からジャーナルを読み出し、副論理ボリューム(COPY1)のデータを更新する(図1の290)。

30

【0027】

記憶システム100の内部構造を図2に示す。記憶システム100は、1つ以上のホストアダプタ110、1つ以上のディスクアダプタ120、1つ以上のキャッシュメモリ130、1つ以上の共有メモリ140、1つ以上の記憶装置150、1つ以上のコモンバス160、1つ以上の接続線170を備えて構成される。ホストアダプタ110、ディスクアダプタ120、キャッシュメモリ130、共有メモリ140はコモンバス160により相互間が接続されている。コモンバス160は、コモンバス160の障害時のために2重化されてもよい。ディスクアダプタ120と記憶装置150とは接続線170によって接続されている。また、図示していないが、記憶システム100の設定、監視、保守等を行うための保守端末が全てのホストアダプタ110とディスクアダプタ120とに専用線を用いて接続されている。

40

【0028】

ホストアダプタ110は、ホストコンピュータ180とキャッシュメモリ130間のデータ転送を制御する。ホストアダプタ110は、接続線190および接続線200によりホストコンピュータ180もしくは、他の記憶システム100と接続される。ディスクアダプタ120は、キャッシュメモリ130と記憶装置150との間のデータ転送を制御する。キャッシュメモリ130は、ホストコンピュータ180から受信したデータあるいは記

50

憶装置 150 から読み出したデータを一時的に保持するメモリである。共有メモリ 140 は、記憶システム 100 内の全てのホストアダプタ 110 とディスクアダプタ 120 とが共有するメモリである。

【0029】

ボリューム情報 400 は、論理ボリュームを管理する情報であり、ボリューム状態、フォーマット形式、容量、ペア番号、物理アドレスを保持する。図 4 にボリューム情報 400 の一例を示す。ボリューム情報 400 は、ホストアダプタ 110 およびディスクアダプタ 120 から参照可能なメモリ、例えば管理メモリ 140 に保存される。ボリューム状態は、“正常”、“正”、“副”、“異常”、“未使用”のいずれかを保持する。ボリューム状態が“正常”もしくは“正”である論理ボリューム 230 は、ホストコンピュータ 180 から正常にアクセス可能な論理ボリューム 230 であることを示す。ボリューム状態が“副”である論理ボリューム 230 は、ホストコンピュータ 180 からのアクセスを許可してもよい。ボリューム状態が“正”である論理ボリューム 230 は、データの複製が行われている論理ボリューム 230 であることを示す。ボリューム状態が“副”である論理ボリューム 230 は、複製に使用されている論理ボリューム 230 であることを示す。ボリューム状態が“異常”の論理ボリューム 230 は、障害により正常にアクセスできない論理ボリューム 230 であることを示す。障害とは、例えば、論理ボリューム 230 を保持する記憶装置 150 の故障である。ボリューム状態が“未使用”の論理ボリューム 230 は、使用していない論理ボリューム 230 であることを示す。ペア番号は、ボリューム状態が“正”もしくは“副”の場合に有効であり、後述するペア情報 500 を特定するためのペア番号を保持する。図 4 に示す例では、論理ボリューム 1 は、フォーマット形式が OPEN3、容量が 3GB、記憶装置番号 1 の記憶装置 150 の記憶領域の先頭からデータが格納されており、アクセス可能であり、データの複製対象であることを示す。

【0030】

ペア情報 500 は、ペアを管理する情報であり、ペア状態、正記憶システム番号、正論理ボリューム番号、副記憶システム番号、副論理ボリューム番号、グループ番号、コピー済みアドレスを保持する。図 5 にペア情報 500 の一例を示す。ペア情報 500 は、ホストアダプタ 110 およびディスクアダプタ 120 から参照可能なメモリ、例えば管理メモリ 140 に保存する。ペア状態は、“正常”、“異常”、“未使用”、“コピー未”、“コピー中”のいずれかを保持する。ペア状態が“正常”の場合は、正論理ボリューム 230 のデータ複製が正常に行われていることを示す。ペア状態が“異常”の場合は、障害により正論理ボリューム 230 の複製が行えないことを示す。障害とは、例えば、接続パス 200 の断線などである。ペア状態が“未使用”の場合は、当該ペア番号の情報は有効でないことを示す。ペア状態が“コピー中”の場合は、後述する初期コピー処理中であることを示す。ペア状態が“コピー未”の場合は、後述する初期コピー処理が未だ行われていないことを示す。正記憶システム番号は、正論理ボリューム 230 を保持する正記憶システム 100A を特定する番号を保持する。副記憶システム番号は、副論理ボリューム 230 を保持する副記憶システム 100B を特定する番号を保持する。グループ番号は、正記憶システムの場合は、正論理ボリュームが属するグループ番号を保持する。副記憶システムの場合は、副論理ボリュームが属するグループ番号を保持する。コピー済みアドレスは、後述する初期コピー処理にて説明する。図 5 のペア情報 1 は、データ複製対象が正記憶システム 1 の正論理ボリューム 1 であり、データ複製先が副記憶システム 2 の副論理ボリューム 1 であり、正常にデータ複製処理が行われていることを示す。

【0031】

グループ情報 600 は、グループ状態、ペア集合、ジャーナル論理ボリューム番号、更新番号を保持する。図 6 にグループ情報 600 の一例を示す。グループ情報 600 は、ホストアダプタ 110 およびディスクアダプタ 120 から参照可能なメモリ、例えば管理メモリ 140 に保存する。グループ状態は、“正常”、“異常”、“未使用”のいずれかを保持する。グループ状態が“正常”の場合は、ペア集合の少なくともひとつのペア状態が“正常”であることを示す。グループ状態が“異常”の場合は、ペア集合の全てのペア状態

10

20

30

40

50

が“異常”であることを示す。グループ状態が“未使用”の場合は、当該グループ番号の情報は有効でないことを示す。ペア集合は、正記憶システムの場合は、グループ番号が示すグループに属する全ての正論理ボリュームのペア番号を保持する。副記憶システムの場合は、グループ番号が示すグループに属する全ての副論理ボリュームのペア番号を保持する。ジャーナル論理ボリューム番号は、当該グループ番号のグループに属するジャーナル論理ボリューム番号を示す。更新番号は、初期値は1であり、グループ内の正論理ボリュームに対しデータの書き込みが行われると、1を加える。更新番号は、ジャーナルの更新情報に記憶し、副記憶システム100Bにて、データの更新順を守るために使用する。例えば、図6のグループ情報1は、ペア情報1, 2から、正論理ボリューム1, 2と、ジャーナル論理ボリューム4から構成されており、正常にデータの複製処理が行われていることを示す。

10

【0032】

ポインタ情報700は、グループ毎に保持し、当該グループのジャーナル論理ボリュームを管理する情報であり、更新情報領域先頭アドレス、ライトデータ領域先頭アドレス、更新情報最新アドレス、更新情報最古アドレス、ライトデータ最新アドレス、ライトデータ最古アドレス、リード開始アドレス、リトライ開始アドレスを保持する。図7および図8にポインタ情報700の一例を示す。更新情報領域先頭アドレスは、ジャーナル論理ボリュームの更新情報を格納する記憶領域（更新情報領域）の先頭の論理アドレスを保持する。ライトデータ領域先頭アドレスは、ジャーナル論理ボリュームのライトデータを格納する記憶領域（ライトデータ領域）の先頭の論理アドレスを保持する。更新情報最新アドレスは、次にジャーナルを格納する場合に、更新情報の保存に使用する先頭の論理アドレスを保持する。更新情報最古アドレスは、最古の（更新番号が小さい）ジャーナルの更新情報を保存する先頭の論理アドレスを保持する。ライトデータ最新アドレスは、次にジャーナルを格納する場合に、ライトデータの保存に使用する先頭の論理アドレスを保持する。ライトデータ最古アドレスは、最古の（更新番号が小さい）ジャーナルのライトデータを保存する先頭の論理アドレスを保持する。リード開始アドレスとリトライ開始アドレスは、正記憶システム100Aのみで使用し、後述するジャーナルリード受信処理にて使用する。図7および図8のポインタ情報700の例では、ジャーナルの管理情報を保存する領域（更新情報領域）は、論理ボリューム4の記憶領域の先頭から699の位置までであり、ジャーナルのライトデータを保存する領域（ライトデータ領域）は、論理ボリューム4の記憶領域の700の位置から2699の位置までである。ジャーナルの管理情報は、論理ボリューム4の記憶領域の200の位置から499の位置まで保存されており、次のジャーナルの管理情報は、論理ボリューム4の記憶領域の500の位置から保存する。ジャーナルのライトデータは論理ボリューム4の記憶領域の1300の位置から2199の位置まで保存されており、次のジャーナルのライトデータは、論理ボリューム4の記憶領域の2200の位置から保存する。

20

30

【0033】

下記の説明では、1つのグループに1つのジャーナル論理ボリュームを割り当てる形態で説明しているが、1つのグループに複数のジャーナル論理ボリュームを割り当ててもよい。例えば、1つのグループに2つのジャーナル論理ボリュームを割り当て、ジャーナル論理ボリューム毎にポインタ情報700を設け、交互にジャーナルを格納する。これにより、ジャーナルの記憶装置150への書き込みが分散でき、性能の向上が見込める。さらに、ジャーナルのリード性能も向上する。別の例としては、1つのグループに2つのジャーナル論理ボリュームを割り当て、通常は、1つのジャーナル論理ボリュームのみを使用する。もう一方のジャーナル論理ボリュームは、使用しているジャーナル論理ボリュームの性能が低下した場合に使用する。性能が低下する例は、ジャーナル論理ボリュームが、複数の記憶装置150から構成され、RAID5の方式でデータを保持しており、構成する記憶装置150の一台が故障中の場合である。

40

【0034】

なお、上述のボリューム情報400、ペア情報500、グループ情報600、及びポイン

50

タ情報 700 等は、共有メモリ 140 に格納されていることが好ましい。しかし、本実施例は、この場合に限られず、これらの情報を、キャッシュメモリ 130、ホストアダプタ 110、ディスクアダプタ 120、その他記憶装置 150 に集中して格納または分散して格納することもよい。

【0035】

次に、正記憶システム 100A から副記憶システム 100B に対して、データ複製を開始する手順を図 9、図 10 を用いて説明する。

【0036】

(1) グループ作成について説明する(ステップ 900)。ユーザは、保守端末あるいはホストコンピュータ 180 を使用して、正記憶システム 100A のグループ情報 600 を参照し、グループ状態が“未使用”のグループ番号 A を取得する。ユーザは、保守端末あるいはホストコンピュータ 180 を使用して、グループ番号 A を指定し、グループ作成指示を正記憶システム 100A に行う。

10

【0037】

グループ作成指示を受けて、正記憶システム 100A は、指定されたグループ番号 A のグループ状態を“正常”に変更する。

【0038】

同様に、ユーザは、副記憶システム 100B のグループ情報 600 を参照し、グループ状態が“未使用”のグループ番号 B を取得する。ユーザは、保守端末あるいはホストコンピュータ 180 を使用して、副記憶システム 100B とグループ番号 B を指定し、グループ作成指示を正記憶システム 100A に行う。正記憶システム 100A は、受信したグループ作成指示を副記憶システム 100B に転送する。副記憶システム 100B は、指定されたグループ番号 B のグループ状態を“正常”に変更する。

20

【0039】

ユーザは、副記憶システム 100B の保守端末あるいは副記憶システム 100B に接続したホストコンピュータ 180 を使用して、グループ番号 B を指定し、グループ作成指示を副記憶システム 100B に行ってもよい。

【0040】

(2) ペア登録について説明する(ステップ 910)。ユーザは、保守端末あるいはホストコンピュータ 180 を使用して、データ複製対象を示す情報とデータ複製先を示す情報を指定し、ペア登録指示を正記憶システム 100A に行う。データ複製対象を示す情報は、データ複製対象のグループ番号 A と正論理ボリューム番号である。データ複製先を示す情報は、複製データを保存する副記憶システム 100B とグループ番号 B、副論理ボリューム番号である。

30

【0041】

前記ペア登録指示を受けて、正記憶システム 100A は、ペア情報 500 からペア情報が“未使用”のペア番号を取得し、ペア状態を“コピー未”に、正記憶システム番号に、正記憶システム 100A を示す正記憶システム番号を、正論理ボリューム番号に、指示された正論理ボリューム番号を、副記憶システム番号に、指示された副記憶システム番号を、副論理ボリューム番号に、指示された副論理ボリューム番号を、グループ番号に、指示されたグループ番号 A を設定する。正記憶システム 100A は、指示されたグループ番号 A のグループ情報 600 のペア集合に前記取得したペア番号を追加し、正論理ボリューム番号のボリューム状態を“正”に変更する。

40

【0042】

正記憶システム 100A は、正記憶システム 100A を示す正記憶システム番号、ユーザから指定されたグループ番号 B、正論理ボリューム番号、および副論理ボリューム番号を副記憶システム 100B に指示する。副記憶システム 100B は、ペア情報 500 から未使用のペア番号を取得し、ペア状態を“コピー未”に、正記憶システム番号に、記憶システム 100A を示す正記憶システム番号を、正論理ボリューム番号に、指示された正論理ボリューム番号を、副記憶システム番号に、副記憶システム B を示す副記憶システム番号

50

を、副論理ボリューム番号に、指示された副論理ボリューム番号を、グループ番号に、指示されたグループ番号Bを設定する。

【0043】

副記憶システム100Bは、指示されたグループ番号Bのグループ情報600のペア集合に前記取得したペア番号を追加し、副論理ボリューム番号のボリューム状態を“副”に変更する。

【0044】

以上の動作を全てのデータ複製対象のペアに対して行う。

【0045】

前記の説明では、論理ボリュームのグループへの登録と、論理ボリュームのペアの設定を同時に行う処理を説明したが、それぞれ個別に行ってもよい。

10

【0046】

(3) ジャーナル論理ボリューム登録について説明する(ステップ920)。ユーザは、保守端末あるいはホストコンピュータ180を使用して、ジャーナルの保存に使用する論理ボリューム(ジャーナル論理ボリューム)をグループに登録する指示(ジャーナル論理ボリューム登録指示)を正記憶システム100Aに行う。ジャーナル論理ボリューム登録指示は、グループ番号と論理ボリューム番号からなる。

【0047】

正記憶システム100Aは、指示されたグループ番号のグループ情報600のジャーナル論理ボリューム番号に指示された論理ボリューム番号を登録する。当該論理ボリュームのボリューム情報400のボリューム状態を“正常”に設定する。

20

【0048】

同様に、ユーザは、保守端末あるいはホストコンピュータ180を使用して、副記憶システム100Bのボリューム情報400を参照し、副記憶システム100B、グループ番号B、ジャーナル論理ボリュームとして使用する論理ボリューム番号を指定し、ジャーナル論理ボリューム登録を正記憶システム100Aに行う。正記憶システム100Aは、ジャーナル論理ボリューム登録指示を副記憶システム100Bに転送する。副記憶システム100Bは、指示されたグループ番号Bのグループ情報600のジャーナル論理ボリューム番号に指示された論理ボリューム番号を登録する。当該論理ボリュームのボリューム情報400のボリューム状態を“正常”に設定する。

30

【0049】

ユーザは、副記憶システム100Bの保守端末あるいは副記憶システム100Bに接続したホストコンピュータ180を使用して、グループ番号、ジャーナル論理ボリュームとして使用する論理ボリューム番号を指定し、ジャーナル論理ボリューム登録指示を副記憶システム100Bに行ってもよい。

【0050】

以上の動作を全てのジャーナル論理ボリュームとして使用する論理ボリュームに対して行う。ステップ910とステップ920の順は不同ではない。

【0051】

(4) データ複製処理開始について説明する(ステップ930)。ユーザは、保守端末あるいはホストコンピュータ180を使用して、データ複製処理を開始するグループ番号を指定し、データ複製処理の開始を正記憶システム100Aに指示する。正記憶システム100Aは、指示されたグループに属する全てのペア情報400のコピー済みアドレスを0に設定する。

40

【0052】

正記憶システム100Aは、副記憶システム100Bに、後述するジャーナルリード処理およびリストア処理の開始を指示する。

【0053】

正記憶システム100Aは、後述する初期コピー処理を開始する。

【0054】

50

(5) 初期コピー終了について説明する (ステップ 9 4 0) 。

【 0 0 5 5 】

初期コピーが終了すると、正記憶システム 1 0 0 A は、初期コピー処理の終了を副記憶システム 1 0 0 B に通知する。副記憶システム 1 0 0 B は、指示されたグループに属する全ての副論理ボリュームのペア状態を “ 正常 ” に変更する。

【 0 0 5 6 】

図 1 0 は、初期コピー処理のフローチャートである。初期コピー処理は、データ複製対象の正論理ボリュームの全記憶領域に対し、ペア情報 5 0 0 のコピー済みアドレスを用い、記憶領域の先頭から順に、単位サイズ毎にジャーナルを作成する。コピー済みアドレスは、初期値は 0 であり、ジャーナルの作成毎に、作成したデータ量を加算する。論理ボリュームの記憶領域の先頭から、コピー済みアドレスの一つ前までは、初期コピー処理にてジャーナルは作成済みである。初期コピー処理を行うことにより、正論理ボリュームの更新されていないデータを副論理ボリュームに転送することが可能となる。以下の説明では、正記憶システム 1 0 0 A 内のホストアダプタ A が処理を行うように記載しているが、ディスクアダプタ 1 2 0 が行ってもよい。

【 0 0 5 7 】

(1) 正記憶システム 1 0 0 A 内のホストアダプタ A は、処理対象のグループに属するペアでペア状態が “ コピー未 ” である正論理ボリューム A を得、ペアの状態を “ コピー中 ” に変更し、以下の処理を繰り返す (ステップ 1 0 1 0 、 1 0 2 0) 。正論理ボリューム A が存在しない場合は、処理を終了する (ステップ 1 0 3 0) 。

【 0 0 5 8 】

(2) ステップ 1 0 2 0 にて、論理ボリューム A が存在した場合、ホストアダプタ A は、単位サイズ (例えば、 1 M B) のデータを対象にジャーナルを作成する。ジャーナル作成処理は後述する (ステップ 1 0 4 0) 。

【 0 0 5 9 】

(3) ホストアダプタ A は、コピー済みアドレスに作成したジャーナルのデータサイズを加算する (ステップ 1 0 5 0) 。

【 0 0 6 0 】

(4) コピー済みアドレスが、正論理ボリューム A の容量に達するまで、上記処理を繰り返す (ステップ 1 0 6 0) 。コピー済みアドレスが、正論理ボリューム A の容量と等しくなった場合、正論理ボリューム A の全記憶領域に対してジャーナルを作成したため、ペア状態を “ 正常 ” に更新し、他の正論理ボリュームの処理を開始する (ステップ 1 0 7 0) 。

【 0 0 6 1 】

前記のフローチャートでは、論理ボリュームを 1 つずつ対象とするように説明したが、複数の論理ボリュームを同時に処理してもよい。

【 0 0 6 2 】

図 1 1 は命令受信処理 2 1 0 の処理を説明する図、図 1 2 は命令受信処理 2 1 0 のフローチャート、図 1 3 はジャーナル作成処理のフローチャートである。以下、これらを用いて、正記憶システム 1 0 0 A が、ホストコンピュータ 1 8 0 からデータ複製対象の論理ボリューム 2 3 0 にライト命令を受信した場合の動作について説明する。

【 0 0 6 3 】

(1) 記憶システム 1 0 0 A 内のホストアダプタ A は、ホストコンピュータからアクセス命令を受信する。アクセス命令は、リード、ライト、後述するジャーナルリード等の命令、命令対象の論理アドレス、データ量等を含んでいる。以下、アクセス命令内の論理アドレスを論理アドレス A、論理ボリューム番号を論理ボリューム A、論理ボリューム内位置を論理ボリューム内位置 A、データ量をデータ量 A とする (ステップ 1 2 0 0) 。

【 0 0 6 4 】

(2) ホストアダプタ A は、アクセス命令を調べる (ステップ 1 2 1 0 、 1 2 1 5) 。ステップ 1 2 1 5 の調べで、アクセス命令がジャーナルリード命令の場合は、後述するジャ

10

20

30

40

50

ーナルリード受信処理を行う（ステップ1220）。アクセス命令がジャーナルリード命令およびライト命令以外、例えば、リード命令の場合は、従来技術と同じようにリード処理を行う（ステップ1230）。

【0065】

（3）ステップ1210の調べで、アクセス命令がライト命令の場合は、論理ボリュームAのボリューム情報400を参照し、ボリューム状態を調べる（ステップ1240）。ステップ1240の調べで、論理ボリュームAのボリューム状態が、“正常”もしくは“正”以外の場合は、論理ボリュームAへのアクセスは不可能なため、ホストコンピュータ180に異常終了を報告する（ステップ1245）。

【0066】

（4）ステップ1240の調べで、論理ボリュームAのボリューム状態が、“正常”、“正”のいずれかの場合は、ホストアダプタAは、キャッシュメモリ130を確保し、ホストコンピュータ180にデータ受信の準備ができたことを通知する。ホストコンピュータ180は、その通知を受け、ライトデータを正記憶システム100Aに送信する。ホストアダプタAは、ライトデータを受信し、当該キャッシュメモリ130に保存する（ステップ1250、図11の1100）。

【0067】

（5）ホストアダプタAは、論理ボリュームAのボリューム状態を参照し、論理ボリュームAがデータ複製対象かどうかを調べる（ステップ1260）。ステップ1260の調べで、ボリューム状態が、“正”である場合は、論理ボリュームAがデータ複製対象であるため、後述するジャーナル作成処理を行う（ステップ1265）。

【0068】

（6）ステップ1260の調べで、ボリューム状態が、“正常”である場合、もしくはステップ1265のジャーナル作成処理の終了後、ホストアダプタAは、ディスクアダプタ120にライトデータを記憶装置150に書き込むことを命令し（図11の1140）、ホストコンピュータ180に終了報告する（ステップ1270、1280）。その後、当該ディスクアダプタ120は、リードライト処理により、記憶装置150にライトデータを保存する（図11の1110）。

【0069】

次に、ジャーナル作成処理について説明する。

【0070】

（1）ホストアダプタAは、ジャーナル論理ボリュームのボリューム状態を調べる（ステップ1310）。ステップ1310の調べで、ジャーナル論理ボリュームのボリューム状態が、“異常”の場合は、ジャーナル論理ボリュームにジャーナルの格納が不可能なため、グループ状態を“異常”に変更し、処理を終了する（ステップ1315）。この場合、ジャーナル論理ボリュームを正常な論理ボリュームに変更する等を行う。

【0071】

（2）ステップ1310の調べで、ジャーナル論理ボリュームが正常である場合、ジャーナル作成処理を継続する。ジャーナル作成処理は、初期コピー処理内の処理であるか、命令受信処理内の処理であるかによって処理が異なる（ステップ1320）。ジャーナル作成処理が命令受信処理内の処理の場合は、ステップ1330からの処理を行う。ジャーナル作成処理が初期コピー処理内の場合は、ステップ1370からの処理を行う。

【0072】

（3）ジャーナル作成処理が命令受信処理内の処理の場合、ホストアダプタAは、ライト対象の論理アドレスAが、初期コピー処理の処理対象となったかを調べる（ステップ1330）。論理ボリュームAのペア状態が“コピー未”の場合は、後に初期コピー処理にてジャーナル作成処理が行われるため、ジャーナルを作成せずに処理を終了する（ステップ1335）。論理ボリュームAのペア状態が“コピー中”の場合は、コピー済みアドレスが論理アドレス内位置Aと等しいもしくは、小さい場合は、後に初期コピー処理にてジャーナル作成処理が行われるため、ジャーナルを作成せずに処理を終了する（ステップ13

10

20

30

40

50

35)。上記以外、つまり、論理ボリュームAのペア状態が“コピー中”かつコピー済みアドレスが論理アドレス内位置A以上の場合もしくは、論理ボリュームAのペア状態が“正常”の場合は、既に初期コピー処理が完了しているため、ジャーナル作成処理を継続する。

【0073】

(4)次に、ホストアダプタAは、ジャーナルがジャーナル論理ボリュームに格納可能であるを調べる。ポインタ情報700を用い、更新情報領域の未使用領域の有無を調べる(ステップ1340)。ポインタ情報700の更新情報最新アドレスと更新情報最古アドレスが等しい場合は、更新情報領域に未使用領域が存在しないため、ジャーナル作成失敗として処理を終了する(ステップ1390)。

10

【0074】

ステップ1340の調べで、更新情報領域に未使用領域が存在する場合は、ポインタ情報700を用い、ライトデータ領域にライトデータが格納できるかを調べる(ステップ1345)。ライトデータ最新アドレスとデータ量Aの和が、ライトデータ最古アドレスと等しいもしくは大きい場合は、ライトデータ領域に格納できないため、ジャーナル作成失敗として処理を終了する(ステップ1390)。

【0075】

(5)ジャーナルが格納可能である場合、ホストアダプタAは、更新番号と更新情報を格納する論理アドレスとライトデータを格納する論理アドレスを取得し、更新情報をキャッシュメモリ130内に作成する。更新番号は、対象グループのグループ情報600から取得し、1を足した数値をグループ情報600の更新番号に設定する。更新情報を格納する論理アドレスは、ポインタ情報700の更新情報最新アドレスであり、更新情報のサイズを足した数値をポインタ情報700の更新情報最新アドレスに設定する。ライトデータを格納する論理アドレスは、ポインタ情報700のライトデータ最新アドレスであり、ライトデータ最新アドレスにデータ量Aを足した数値をポインタ情報700のライトデータ最新アドレスに設定する。

20

【0076】

ホストアダプタAは、上記取得した数値とグループ番号、ライト命令を受信した時刻、ライト命令内の論理アドレスA、データ量Aを更新情報に設定する(ステップ1350、図11の1120)。例えば、図6に示すグループ情報600、図7に示すポインタ情報700の状態、グループ1に属する正論理ボリューム1の記憶領域の先頭から800の位置にデータサイズ100のライト命令を受信した場合、図22に示す更新情報を作成する。グループ情報の更新番号は5、ポインタ情報の更新情報最新アドレスは600(更新情報のサイズは100とする)、ライトデータ最新アドレスは2300となる。

30

【0077】

(6)ホストアダプタAは、ディスクアダプタ120に、ジャーナルの更新情報とライトデータを記憶装置150に書き込むことを命令し、正常終了する(ステップ1360、図11の1130、1140、1150)。

【0078】

(7)ジャーナル作成処理が、初期コピー処理内の処理の場合は、ステップ1370からの処理を行う。ホストアダプタAは、ジャーナルが作成可能であるを調べる。ポインタ情報700を用い、更新情報領域の未使用領域の有無を調べる(ステップ1370)。ポインタ情報700の更新情報最新アドレスと更新情報最古アドレスが等しい場合は、更新情報領域に未使用領域が存在しないため、ジャーナル作成失敗として処理を終了する(ステップ1390)。本実施例で示した初期コピー処理の場合、ジャーナルのライトデータは、正論理ボリュームからリードし、ライトデータ領域は使用しないため、ライトデータ領域の未使用領域の確認は不要である。

40

【0079】

(8)ステップ1370の調べで、ジャーナルが作成可能である場合、ホストアダプタAは、更新情報に設定する数値を取得し、更新情報をキャッシュメモリ130内に作成する

50

。更新番号は、対象グループのグループ情報600から取得し、1を足した数値をグループ情報600の更新番号に設定する。更新情報を格納する論理アドレスは、ポインタ情報700の更新情報最新アドレスの位置であり、更新情報のサイズを足した数値をポインタ情報700の更新情報最新アドレスに設定する。

【0080】

ホストアダプタAは、上記取得した数値とグループ番号、本処理の開始時刻、初期コピー処理対象の論理アドレス、初期コピーの1回の処理量、ライトデータを格納したジャーナル論理ボリュームの論理アドレスに初期コピー処理対象の論理アドレスを設定する（ステップ1380、図11の1120）。

【0081】

（9）ホストアダプタAは、ディスクアダプタ120に、更新情報を記憶装置150に書き込むことを命令し、正常終了する（ステップ1385、図11の1140、1160）。

【0082】

上記説明では、更新情報をキャッシュメモリ130内に存在するように記載しているが、共有メモリ140内等に格納してもよい。

【0083】

ライトデータの記憶装置150への書き込みは、非同期、つまり、ステップ1360およびステップ1385の直後でなくともよい。ただし、ホストコンピュータ180が、論理アドレスAにライト命令を再び行った場合、ジャーナルのライトデータが上書きされるため、ホストコンピュータ180からライトデータを受信する前に、ジャーナルのライトデータは、更新情報のジャーナル論理ボリュームの論理アドレスに対応する記憶装置150に書き込む必要がある。もしくは、別のキャッシュメモリに退避し、後に更新情報のジャーナル論理ボリュームの論理アドレスに対応する記憶装置150に書き込みを行ってもよい。

【0084】

前述したジャーナル作成処理では、ジャーナルを記憶装置150に保存するとしていたが、ジャーナル用に予め一定量のキャッシュメモリ130を用意しておき、当該キャッシュメモリを全て使用してから、記憶装置150にジャーナルを保存してもよい。ジャーナル用のキャッシュメモリ量は、例えば、保守端末から指定する。

【0085】

リードライト処理220は、ディスクアダプタ120が、ホストアダプタ110もしくはディスクアダプタ120から命令を受け、実施する処理である。実施する処理は、指定されたキャッシュメモリ130のデータを指定された論理アドレスに対応する記憶装置150内の記憶領域に書き込む処理、指定された論理アドレスに対応する記憶装置150内の記憶領域から指定されたキャッシュメモリ130にデータを読み込む処理等である。

【0086】

図14はジャーナルリード命令を受信した正記憶システム100AのホストアダプタAの動作（ジャーナルリード受信処理）を説明する図、図15はフローチャートである。以下、これらを用いて、正記憶システム100Aが、副記憶システム100Bからジャーナルリード命令を受信した場合の動作について説明する。

【0087】

（1）正記憶システム100A内のホストアダプタAは、副記憶システム100Bからアクセス命令を受信する。アクセス命令は、ジャーナルリード命令であることを示す識別子、命令対象のグループ番号、リトライ指示の有無を含んでいる。以下、アクセス命令内のグループ番号をグループ番号Aとする（ステップ1220、図14の1410）。

【0088】

（2）ホストアダプタAは、グループ番号Aのグループ状態が“正常”であるかを調べる（ステップ1510）。ステップ1510の調べで、グループ状態が“正常”以外、例えば、“障害”の場合は、副記憶システム100Bにグループ状態を通知し、処理を終了す

10

20

30

40

50

る。副記憶システム 100B は、受信したグループ状態に応じて処理を行う。例えば、グループ状態が“障害”の場合は、ジャーナルリード処理を終了する(ステップ 1515)。

【0089】

(3) ステップ 1510 の調べで、グループ番号 A のグループ状態が“正常”の場合、ホストアダプタ A は、ジャーナル論理ボリュームの状態を調べる(ステップ 1520)。ステップ 1520 の調べで、ジャーナル論理ボリュームのボリューム状態が“正常”でない場合は、例えば、“障害”の場合は、グループ状態を“障害”に変更し、副記憶システム 100B にグループ状態を通知し、処理を終了する。副記憶システム 100B は、受信したグループ状態に応じて処理を行う。例えば、グループ状態が“障害”の場合は、ジャーナルリード処理を終了する(ステップ 1525)。

10

【0090】

(4) ステップ 1520 の調べで、ジャーナル論理ボリュームのボリューム状態が“正常”の場合は、ジャーナルリード命令がリトライ指示かを調べる(ステップ 1530)。

【0091】

(5) ステップ 1530 の調べで、ジャーナルリード命令がリトライ指示の場合、ホストアダプタ A は、前回送信したジャーナルを再度、副記憶システム 100B に送信する。ホストアダプタ A は、キャッシュメモリ 130 を確保し、ディスクアダプタに、ポインタ情報 700 のリトライ開始アドレスから、更新情報のサイズの情報をキャッシュメモリに読み込むことを命令する(図 14 の 1420)。

20

【0092】

ディスクアダプタのリードライト処理は、記憶装置 150 から更新情報を読み込み、キャッシュメモリ 130 に保存し、ホストアダプタ A に通知する(図 14 の 1430)。

【0093】

ホストアダプタ A は、更新情報のリード終了の通知を受け、更新情報から、ライトデータの論理アドレスおよびライトデータのサイズを取得し、キャッシュメモリ 130 を確保し、ディスクアダプタにライトデータをキャッシュメモリに読み込むことを命令する(ステップ 1540、図 14 の 1440)。

【0094】

ディスクアダプタのリードライト処理は、記憶装置 150 からライトデータを読み込み、キャッシュメモリ 130 に保存し、ホストアダプタ A に通知する(図 14 の 1450)。

30

【0095】

ホストアダプタ A は、ライトデータのリード終了の通知を受け、更新情報とライトデータを副記憶システム 100B に送信し、ジャーナルを保持しているキャッシュメモリ 130 を開放し、処理を終了する(ステップ 1545、図 14 の 1460)。

【0096】

(6) ステップ 1530 の調べで、リトライ指示でない場合、ホストアダプタ A は、送信していないジャーナルが存在するかを調べ、存在すれば、ジャーナルを副記憶システム 100B に送信する。ホストアダプタ A は、ポインタ情報 700 のリード開始アドレスと更新情報最新アドレスを比較する(ステップ 1550)。

40

【0097】

リード開始アドレスが更新情報最新アドレスと等しい場合は、全てのジャーナルを副記憶システム 100B に送信済みであるため、副記憶システム 100B に“ジャーナル無”を送信し(ステップ 1560)、前回のジャーナルリード命令の時に、副記憶システム 100B に送信したジャーナルの記憶領域を開放する(ステップ 1590)。

【0098】

ジャーナルの記憶領域の開放処理は、ポインタ情報 700 の更新情報最古アドレスに、リトライ開始アドレスを設定する。更新情報最古アドレスがライトデータ領域先頭アドレスとなった場合は、更新情報最古アドレスは 0 とする。ポインタ情報 700 のライトデータ最古アドレスは、前回のリードジャーナル命令に応じて送信したライトデータのサイズを

50

足した数値に変更する。ライトデータ最古アドレスが、ジャーナル論理ボリュームの容量以上の論理アドレスとなった場合は、ライトデータ領域先頭アドレスを減じ、補正する。

【0099】

(7) ステップ1550の調べで、未送信のジャーナルが存在する場合、ホストアダプタAは、キャッシュメモリ130を確保し、ディスクアダプタにポインタ情報700のリード開始アドレスから、更新情報のサイズの情報をキャッシュメモリに読み込むことを命令する(図14の1420)。

【0100】

ディスクアダプタAのリードライト処理は、記憶装置150から更新情報を読み込み、キャッシュメモリ130に保存し、ホストアダプタに通知する(図14の1430)。

10

【0101】

ホストアダプタAは、更新情報のリード終了の通知を受け、更新情報から、ライトデータの論理アドレスおよびライトデータのサイズを取得し、キャッシュメモリ130を確保し、ディスクアダプタAにライトデータをキャッシュメモリに読み込むことを命令する(ステップ1570、図14の1440)。

【0102】

ディスクアダプタAのリードライト処理は、記憶装置150からライトデータを読み込み、キャッシュメモリ130に保存し、ホストアダプタに通知する(図14の1450)。

【0103】

ホストアダプタAは、ライトデータのリード終了の通知を受け、更新情報とライトデータを副記憶システム100Bに送信(ステップ1580)し、ジャーナルを保持しているキャッシュメモリ130を開放する(図14の1460)。そして、ポインタ情報700のリトライ開始アドレスにリード開始アドレスを設定し、リード開始アドレスに送信したジャーナルの更新情報サイズを足した数値を設定する。

20

【0104】

(8) ホストアダプタAは、前回のジャーナルリード命令の処理時に、副記憶システム100Bに送信したジャーナルの記憶領域を開放する(ステップ1590)。

【0105】

前述したジャーナルリード受信処理では、正記憶システム100Aは、ジャーナルを一つずつ副記憶システム100Bに送信していたが、複数同時に副記憶システム100Bに送信してもよい。1つのジャーナルリード命令で、送信するジャーナル数は、副記憶システム100Bがジャーナルリード命令内に指定してもよいし、グループ登録の際等に、ユーザが正記憶システム100Aもしくは、副記憶システム100Bに指定してもよい。さらに、正記憶システム100Aと副記憶システム100Bの接続パス200の転送能力もしくは、負荷等により、動的に1つのジャーナルリード命令で送信するジャーナル数を変更してもよい。また、ジャーナル数でなく、ジャーナルのライトデータのサイズを考慮し、ジャーナルの転送量を指定してもよい。

30

【0106】

前述したジャーナルリード受信処理では、ジャーナルを記憶装置150からキャッシュメモリ130に読み込んでいたが、キャッシュメモリ130に存在する場合は、当該処理は不要である。

40

【0107】

前述したジャーナルリード受信処理内のジャーナルの記憶領域の開放処理は、次のジャーナルリード命令の処理時としたが、副記憶システム100Bにジャーナルを送信した直後に開放してもよい。また、副記憶システム100Bが、ジャーナルリード命令内に開放してよい更新番号を設定し、正記憶システム100Aは、その指示に従って、ジャーナルの記憶領域を開放してもよい。

【0108】

図16はジャーナルリード命令処理240を説明する図、図17はフローチャート、図18はジャーナル格納処理のフローチャートである。以下、これらを用いて、副記憶システ

50

ム 1 0 0 B のホストアダプタ B が、正記憶システム 1 0 0 A からジャーナルを読み出し、ジャーナル論理ボリュームに格納する動作について説明する。

【 0 1 0 9 】

(1) 副記憶システム 1 0 0 B 内のホストアダプタ B は、ジャーナルを格納するキャッシュメモリ 1 3 0 を確保し、ジャーナルリード命令であることを示す識別子、命令対象の正記憶システム 1 0 0 A のグループ番号、リトライ指示の有無を含むアクセス命令を正記憶システム 1 0 0 A に送信する。以下、アクセス命令内のグループ番号をグループ番号 A とする (ステップ 1 7 0 0 、 図 1 6 の 1 6 1 0) 。

【 0 1 1 0 】

(2) ホストアダプタ B は、正記憶システム 1 0 0 A の応答およびジャーナルを受信する (図 1 6 の 1 6 2 0) 。ホストアダプタ B は応答を調べ、正記憶システム 1 0 0 A からの応答が、“ジャーナル無”の場合は、正記憶システム 1 0 0 A には、指定したグループのジャーナルが存在しないため、一定時間後、正記憶システム 1 0 0 A にリードジャーナル命令を送信する (ステップ 1 7 2 0 、 1 7 2 5) 。

【 0 1 1 1 】

(4) 正記憶システム 1 0 0 A の応答が、“グループ状態は障害”もしくは“グループ状態は未使用”の場合は、副記憶システム 1 0 0 B のグループ状態を受信した状態に変更し、ジャーナルリード処理を終了する (ステップ 1 7 3 0 、 1 7 3 5) 。

【 0 1 1 2 】

(5) 正記憶システム 1 0 0 A の応答が、上記以外、つまり、正常終了の場合は、ジャーナル論理ボリュームのボリューム状態を調べる (ステップ 1 7 4 0) 。ジャーナル論理ボリュームのボリューム状態が“異常”の場合は、ジャーナル論理ボリュームにジャーナルの格納が不可能なため、グループ状態を“異常”に変更し、処理を終了する (ステップ 1 7 4 5) 。この場合、ジャーナル論理ボリュームを正常な論理ボリュームに変更する等を行い、グループの状態を正常に戻す。

【 0 1 1 3 】

(6) ステップ 1 7 4 0 の調べで、ジャーナル論理ボリュームのボリューム状態が“正常”の場合は、後述するジャーナル格納処理 1 8 0 0 を行う。ジャーナル格納処理 1 8 0 0 が正常に終了した場合は、次のジャーナルリード命令を送信する。もしくは、一定時間経過後、次のジャーナルリード命令を送信する (ステップ 1 7 6 0) 。次のジャーナル命令を送信するタイミングは、一定の時間間隔で定期的にもよく、受信したジャーナルの個数もしくは、接続線 2 0 0 の通信量、副記憶システム 1 0 0 B が保持しているジャーナルの記憶容量、副記憶システム 1 0 0 B の負荷等によって決めてもよい。さらに、正記憶システム 1 0 0 A が保持しているジャーナルの記憶容量、もしくは正記憶システム 1 0 0 A のポインタ情報を副記憶システム 1 0 0 B から読み出し、その数値に基づいて決めてもよい。上記情報の転送は、専用のコマンドで行ってもよいし、ジャーナルリード命令の応答に含んでもよい。その後の処理は、ステップ 1 7 1 0 以降と同じである。

【 0 1 1 4 】

(7) ステップ 1 8 0 0 のジャーナル格納処理が正常に終了しない場合は、ジャーナル論理ボリュームの未使用領域が足りないため、受信したジャーナルを破棄し、一定時間後にリトライ指示のジャーナルリード命令を送信する (ステップ 1 7 5 5) 。もしくは、ジャーナルをキャッシュメモリに保持しておき、一定時間後に、再度ジャーナル格納処理を行う。これは、後述するリストア処理 2 5 0 が行われることにより、一定時間後には、ジャーナル論理ボリュームに未使用領域が増える可能性があるためである。この方式の場合は、ジャーナルリード命令にリトライ指示の有無は不要である。

【 0 1 1 5 】

次に、図 1 8 に示すジャーナル格納処理 1 8 0 0 について説明する。

【 0 1 1 6 】

(1) ホストアダプタ B は、ジャーナルがジャーナル論理ボリュームに格納可能であることを調べる。ポインタ情報 7 0 0 を用い、更新情報領域に未使用領域の有無を調べる (ステ

10

20

30

40

50

ップ1810)。ポインタ情報700の更新情報最新アドレスと更新情報最古アドレスが等しい場合は、更新情報領域に未使用領域が存在しないため、ジャーナル作成失敗として処理を終了する(ステップ1820)。

【0117】

(2)ステップ1810の調べで、更新情報領域に未使用領域が存在する場合は、ポインタ情報700を用い、ライトデータ領域にライトデータが格納できるかを調べる(ステップ1830)。ライトデータ最新アドレスと受信したジャーナルのライトデータのデータ量の和が、ライトデータ最古アドレスと等しいもしくは大きい場合は、ライトデータ領域にライトデータを格納できないため、ジャーナル作成失敗として処理を終了する(ステップ1820)。

10

【0118】

(3)ジャーナルが格納可能である場合、ホストアダプタBは、受信した更新情報のグループ番号とジャーナル論理ボリュームの論理アドレスを変更する。グループ番号は、副記憶システム100Bのグループ番号に変更し、ジャーナル論理ボリュームの論理アドレスはポインタ情報700のライトデータ最新アドレスに変更する。ホストアダプタBは、ポインタ情報700の更新情報最新アドレスを更新情報最新アドレスに更新情報のサイズを足した数値に変更する。ホストアダプタBは、ポインタ情報700のライトデータ最新アドレスを、ライトデータ最新アドレスにライトデータのサイズを足した数値に変更する(ステップ1840)。

【0119】

(4)ホストアダプタBは、ディスクアダプタ120に、更新情報とライトデータを記憶装置150に書き込むことを命令し、ジャーナル作成成功として処理を終了する(ステップ1850、図16の1630)。その後、ディスクアダプタ120は、リードライト処理により、記憶装置150に更新情報とライトデータを書き込み、キャッシュメモリ130を開放する(図16の1640)。

20

【0120】

前述したジャーナル格納処理では、ジャーナルを記憶装置150に保存するとしていたが、ジャーナル用に予め一定量のキャッシュメモリ130を用意しておき、当該キャッシュメモリを全て使用してから、記憶装置150にジャーナルを保存してもよい。ジャーナル用のキャッシュメモリ量は、例えば、保守端末から指定する。

30

【0121】

図19はリストア処理250を説明する図、図20はフローチャートである。以下、これらを用いて、副記憶システム100BのホストアダプタBが、ジャーナルを利用し、データの更新を行う動作について説明する。リストア処理250は副記憶システム100Bのディスクアダプタ120が処理を行ってもよい。

【0122】

(1)ホストアダプタBは、グループ番号Bのグループ状態が“正常”であるかを調べる(ステップ2010)。ステップ2010の調べで、グループ状態が“正常”以外、例えば、“障害”の場合は、リストア処理を終了する(ステップ2015)。

【0123】

(2)ステップ2010の調べで、グループ状態が“正常”の場合は、ジャーナル論理ボリュームのボリューム状態を調べる(ステップ2020)。ステップ2020の調べで、ジャーナル論理ボリュームのボリューム状態が、“異常”の場合は、アクセス不可能なため、グループ状態を“異常”に変更し、処理を終了する(ステップ2025)。

40

【0124】

(3)ステップ2020の調べで、ジャーナル論理ボリュームのボリューム状態が、“正常”の場合は、リストア対象のジャーナルが存在するかを調べる。ホストアダプタBは、ポインタ情報700の更新情報最古アドレスと更新情報最新アドレスを取得する。更新情報最古アドレスと更新情報最新アドレスが等しい場合、ジャーナルは存在しないため、リストア処理は一旦終了し、一定時間後、リストア処理を再開する(ステップ2030)。

50

【 0 1 2 5 】

(4) ステップ 2 0 3 0 の調べで、リストア対象のジャーナルが存在する場合、最古 (最小) の更新番号を持つジャーナルに対して次の処理を行う。最古 (最小) の更新番号を持つジャーナルの更新情報は、ポインタ情報 7 0 0 の更新情報最古アドレスから保存されている。ホストアダプタ B は、キャッシュメモリ 1 3 0 を確保し、ディスクアダプタに更新情報最古アドレスから、更新情報のサイズの情報をキャッシュメモリ 1 3 0 に読み込むことを命令する (図 1 9 の 1 9 1 0) 。

【 0 1 2 6 】

ディスクアダプタのリードライト処理は、記憶装置 1 5 0 から更新情報を読み込み、キャッシュメモリ 1 3 0 に保存し、ホストアダプタ B に通知する (図 1 9 の 1 9 2 0) 。

10

【 0 1 2 7 】

ホストアダプタ B は、更新情報のリード終了の通知を受け、更新情報から、ライトデータの論理アドレスおよびライトデータのサイズを取得し、キャッシュメモリ 1 3 0 を確保し、ディスクアダプタにライトデータをキャッシュメモリに読み込むことを命令する (図 1 9 の 1 9 3 0) 。

【 0 1 2 8 】

ディスクアダプタのリードライト処理は、記憶装置 1 5 0 からライトデータを読み込み、キャッシュメモリ 1 3 0 に保存し、ホストアダプタに通知する (ステップ 2 0 4 0 、 図 1 9 の 1 9 4 0) 。

【 0 1 2 9 】

(5) ホストアダプタ B は、更新情報から更新する副論理ボリュームの論理アドレスを求め、ディスクアダプタに副論理ボリュームにライトデータを書き込むことを命令する (ステップ 2 0 5 0 、 図 1 9 の 1 9 5 0) 。ディスクアダプタのリードライト処理は、副論理ボリュームの論理アドレスに対応する記憶装置 1 5 0 にデータを書き込み、キャッシュメモリ 1 3 0 を開放し、ホストアダプタに通知する (図 1 9 の 1 9 6 0) 。

20

【 0 1 3 0 】

(6) ホストアダプタ B は、ディスクアダプタのライト処理完了の通知を受け、ジャーナルの記憶領域を開放する。ジャーナルの記憶領域の開放処理は、ポインタ情報 7 0 0 の更新情報最古アドレスを更新情報のサイズを足した数値に変更する。更新情報最古アドレスが、ライトデータ領域先頭アドレスとなった場合は、ライトデータ領域先頭アドレスは 0 とする。ポインタ情報 7 0 0 のライトデータ最古アドレスは、ライトデータのサイズを足した数値に変更する。ライトデータ最古アドレスが、ジャーナル論理ボリュームの容量以上の論理アドレスとなった場合は、ライトデータ領域先頭アドレスを減じ、補正する。その後、ホストアダプタ B は、次のリストア処理を開始する (ステップ 2 0 6 0) 。

30

【 0 1 3 1 】

前述したリストア処理 2 5 0 では、記憶装置 1 5 0 からキャッシュメモリ 1 3 0 にジャーナルを読み込んでいたが、キャッシュメモリ 1 3 0 に存在する場合は、当該処理は不要である。

【 0 1 3 2 】

前述したジャーナルリード受信処理とジャーナルリード命令処理 2 4 0 では、正記憶システム 1 0 0 A が送信するジャーナルをポインタ情報 7 0 0 により決めていたが、副記憶システム 1 0 0 B が送信するジャーナルを決めてもよい。例えば、ジャーナルリード命令に更新番号を追加する。この場合、ジャーナルリード受信処理にて、副記憶システム 1 0 0 B が指定した更新番号の更新情報の論理アドレスを求めるために、正記憶システム 1 0 0 A の共有メモリ 1 4 0 内に、更新番号から更新情報を格納した論理アドレスを求めるテーブルもしくは検索方法を設ける。

40

【 0 1 3 3 】

前述したジャーナルリード受信処理とジャーナルリード命令処理 2 4 0 では、ジャーナルリード命令を用いていたが、通常のリード命令を用いてもよい。例えば、正記憶システム 1 0 0 A のグループ情報 6 0 0 とポインタ情報 7 0 0 を予め副記憶システム 1 0 0 B に転

50

送しておき、副記憶システム100Bは、正記憶システム100Aのジャーナル論理ボリュームのデータ(つまり、ジャーナル)をリードする。

【0134】

前述したジャーナルリード受信処理では、更新番号の順に、正記憶システム100Aから副記憶システム100Bにジャーナルを送信すると説明したが、更新番号の順に送信せずともよい。また、正記憶システム100Aから副記憶システム100Bに複数のジャーナルリード命令を送信してもよい。この場合、リストア処理にて更新番号順にジャーナルを処理するために、副記憶システム100Bに、更新番号から更新情報を格納した論理アドレスを求めるテーブルもしくは検索方法を設ける。

【0135】

前述した本発明のデータ処理システムでは、正記憶システムがジャーナルを取得し、副記憶システムがデータの複製を行う。これにより、正記憶システムに接続したホストコンピュータは、データの複製に関する負荷を負わない。さらに、正記憶システムと副記憶システム間でジャーナルを転送することにより、正記憶システムと正記憶システムに接続したホストコンピュータの通信線を使用しない。

【0136】

図23は、第二の実施形態の論理的な構成を示す図である。

【0137】

ホストコンピュータ180と記憶システム100Cを接続バス190により接続し、記憶システム100Cと正記憶システム100Aを接続バス200により接続し、正記憶システム100Aと副記憶システム100Bを接続バス200により接続した構成である。記憶システム100Cは、記憶システム100Cの論理ボリューム(ORG1)へのデータ更新時、論理ボリューム(ORG1)のデータ更新と正記憶システム100A内の論理ボリューム(DATA1)のデータ更新を行う。

【0138】

正記憶システム100Aは、第1の実施例で説明した通り、正論理ボリューム(DATA1)へのデータ更新時、前述した命令受信処理210およびリードライト処理220によって、ジャーナル論理ボリューム(JNL1)にジャーナルの保存を行う(2310)。

【0139】

副記憶システム100Bは、前述したジャーナルリード処理240によって、正記憶システム100Aからジャーナルをリードし、リードライト処理220によって、ジャーナル論理ボリューム(JNL2)にジャーナルを保存する(2320)。

【0140】

正記憶システム100Aは、副記憶システム100Bからジャーナルをリードする命令を受信すると、命令受信処理210およびリードライト処理220によって、ジャーナル論理ボリューム(JNL1)からジャーナルを読み出し、副記憶システム100Bに送信する(2320)。

【0141】

副記憶システム100Bは、前述したリストア処理250およびリードライト処理220によって、更新番号に従い、ジャーナル論理ボリューム(JNL2)からジャーナルを読み出し、正論理ボリューム(DATA1)の複製である副論理ボリューム(COPY1)のデータを更新する(2330)。このように、更新番号の順にデータを更新することにより、論理ボリューム間のデータの整合性を保つことが可能となる。

【0142】

前述した本発明のデータ処理システムでは、正記憶システムがジャーナルを取得し、ジャーナル専用の記憶領域に格納する。さらに、副記憶システムは、正記憶システムから受信したジャーナルをジャーナル専用の記憶領域に格納する。ジャーナル専用の記憶領域はデータ複製対象の記憶領域より少なくすることが可能であり、より少ない記憶容量で、副記憶システムに正記憶システムのデータの複製が可能となる。

【0143】

10

20

30

40

50

図 2 4 は、第三の実施形態の論理的な構成を示す図である。

【 0 1 4 4 】

ホストコンピュータ 1 8 0 と記憶システム 1 0 0 C を接続バス 1 9 0 により接続し、記憶システム 1 0 0 C と正記憶システム 1 0 0 A を接続バス 2 0 0 により接続し、正記憶システム 1 0 0 A と副記憶システム 1 0 0 B を接続バス 2 0 0 により接続した構成である。記憶システム 1 0 0 C は、従来技術で説明した通り、記憶システム 1 0 0 C の論理ボリューム (O R G 1) へのデータ更新時、論理ボリューム (O R G 1) のデータ更新と正記憶システム 1 0 0 A 内の論理ボリューム (D A T A 1) のデータ更新を行う。

【 0 1 4 5 】

正記憶システム 1 0 0 A は、記憶システム 1 0 0 C に対し、正論理ボリューム (D A T A 1) があるように見せるが、実際の記憶領域、つまり記憶装置 1 5 0 は割り当てない。例えば、ボリューム情報 4 0 0 の物理アドレスに記憶装置 1 5 0 を割り当てていないことを示す数値を設定する。正記憶システム 1 0 0 A は、記憶システム 1 0 0 C から正論理ボリューム (D A T A 1) のデータへのライト命令受信時、前述した命令受信処理 2 1 0 内のステップ 1 2 7 0 の処理を行わず、ジャーナル論理ボリューム (J N L 1) にジャーナルの保存のみを行う (2 4 1 0)。

10

【 0 1 4 6 】

副記憶システム 1 0 0 B は、前述したジャーナルリード処理 2 4 0 によって、正記憶システム 1 0 0 A からジャーナルをリードし、リードライト処理 2 2 0 によって、ジャーナル論理ボリューム (J N L 2) にジャーナルを保存する (2 4 2 0)。

20

【 0 1 4 7 】

正記憶システム 1 0 0 A は、副記憶システム 1 0 0 B からジャーナルをリードする命令を受信すると、命令受信処理 2 1 0 およびリードライト処理 2 2 0 によって、ジャーナル論理ボリューム (J N L 1) からジャーナルを読み出し、記憶システム 1 0 0 B に送信する (2 4 2 0)。

【 0 1 4 8 】

副記憶システム 1 0 0 B は、前述したリストア処理 2 5 0 およびリードライト処理 2 2 0 によって、更新番号に従い、ジャーナル論理ボリューム (J N L 2) からジャーナルを読み出し、論理ボリューム (O R G 1) の複製である副論理ボリューム (C O P Y 1) のデータを更新する (2 4 3 0)。このように、更新番号の順にデータを更新することにより、論理ボリューム間のデータの整合性を保つことが可能となる。

30

【 0 1 4 9 】

前述した本発明のデータ処理システムでは、記憶システム 1 0 0 C もしくは、記憶システム 1 0 0 C に接続したホストコンピュータ 1 8 0 に障害が生じた場合、副記憶システム 1 0 0 B の論理ボリューム (C O P Y 1) に対し、正記憶システム 1 0 0 A 内のジャーナル (J N L 1) を反映することにより、記憶システム 1 0 0 B に接続したホストコンピュータにより最新データの参照、更新が可能となる。さらに、正記憶システム 1 0 0 A にデータの複製を保持せず、ジャーナルのみを格納することで、データ複製に必要な記憶容量が少なくすることが可能となる。

【 0 1 5 0 】

以上、本発明者によってなされた発明を実施例の形態に基づき具体的に説明したが、本発明は前記実施の形態に限定されるものでなく、その要旨を逸脱しない範囲で種々変更可能であることはいうまでもない。

40

【 0 1 5 1 】

【発明の効果】

本発明によれば、記憶システムの上位の計算機に影響を与えず、複数の記憶システム間でデータ転送又はデータの複製をすることが可能な記憶システムを提供できる。さらに、記憶システムと計算機との間の通信にも影響を与えない記憶システムを提供できる。

【 0 1 5 2 】

さらに、本発明によれば、複数の記憶システム内に保持するデータ格納領域を少なくする

50

ことができる。さらに、複数の記憶システムの上位の計算機の業務に影響を与えることのないように、高速かつ効率的に複数の記憶システム間でデータ転送又はデータの複製をすることができる。

【図面の簡単な説明】

【図 1】本発明の一実施形態の論理的な構成を示すブロック図である。

【図 2】本発明の一実施形態の記憶システムのブロック図である。

【図 3】本発明の一実施形態の更新情報とライトデータの関係を示す図である。

【図 4】本発明の一実施形態のボリューム情報の例を示す図である。

【図 5】本発明の一実施形態のペア情報の例を示す図である。

【図 6】本発明の一実施形態のグループ情報の例を示す図である。

【図 7】本発明の一実施形態のポイント情報の例を示す図である。

【図 8】本発明の一実施形態のジャーナル論理ボリュームの構造を示す図である。

【図 9】本発明の一実施形態のデータの複製を開始する手順を示すフローチャートである。

【図 10】本発明の一実施形態の初期コピー処理を示すフローチャートである。

【図 11】本発明の一実施形態の命令受信処理を示す図である。

【図 12】本発明の一実施形態の命令受信処理のフローチャートである。

【図 13】本発明の一実施形態のジャーナル作成処理のフローチャートである。

【図 14】本発明の一実施形態のジャーナルリード受信処理を示す図である。

【図 15】本発明の一実施形態のジャーナルリード受信処理のフローチャートである。

【図 16】本発明の一実施形態のジャーナルリード命令処理を示す図である。

【図 17】本発明の一実施形態のジャーナルリード命令処理のフローチャートである。

【図 18】本発明の一実施形態のジャーナル格納処理のフローチャートである。

【図 19】本発明の一実施形態のリストア処理を示す図である。

【図 20】本発明の一実施形態のリストア処理のフローチャートである。

【図 21】本発明の一実施形態の更新情報の例を示す図である。

【図 22】本発明の一実施形態のジャーナル作成処理時の更新情報の例を示す図である。

【図 23】本発明の第 2 の実施形態を示す図である。

【図 24】本発明の第 3 の実施形態を示す図である。

【符号の説明】

1 0 0 記憶システム

1 1 0 ホストアダプタ

1 2 0 ディスクアダプタ

1 3 0 キャッシュメモリ

1 4 0 管理メモリ

1 5 0 記憶装置

1 6 0 コモンバス

1 7 0 ディスクアダプタと記憶装置間の接続線

1 8 0 ホストコンピュータ

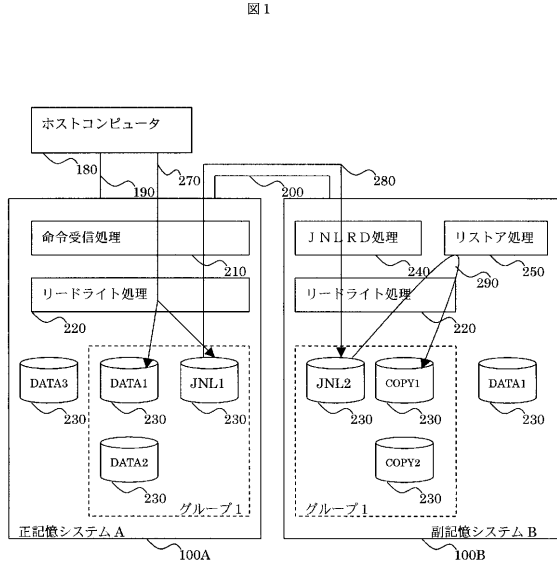
10

20

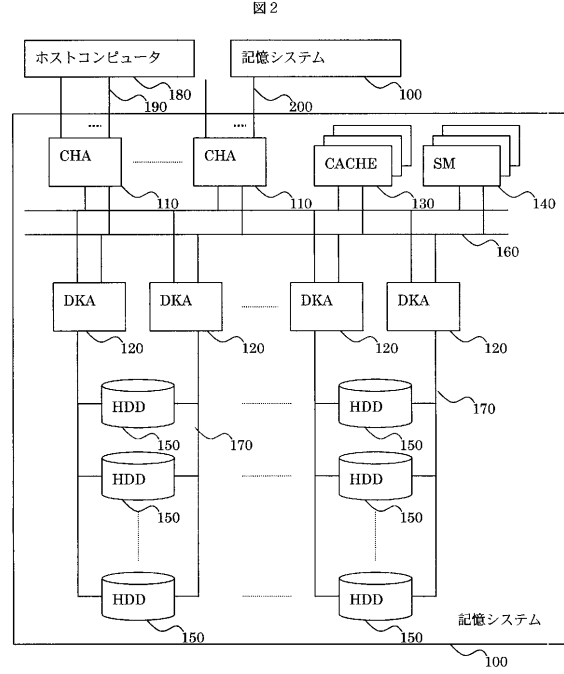
30

40

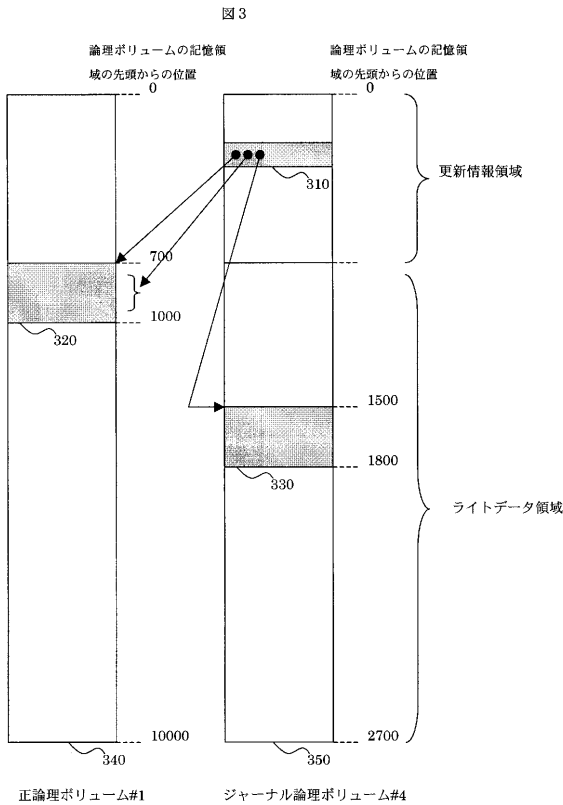
【図1】



【図2】



【図3】



【図4】

図4

論理ボリューム番号	ボリューム状態	フォーマット形式	容量	ベア番号	物理アドレス	
					記憶番号	先頭から位置
1	正	OPEN3	3	1	1	0
2	正	OPEN6	6	2	1	3
3	未使用	OPEN6	6	0	1	9
4	正常	OPEN9	9	0	2	0
5	正常	OPEN3	3	0	2	9
6	未使用	OPEN6	6	0	2	12

400 ボリューム情報

【図5】

図5

ベア番号	ベア状態	正記憶システム番号	正論理ボリューム番号	副記憶システム番号	副論理ボリューム番号	グループ番号	コピー済みアドレス
1	正常	1	1	2	1	1	0
2	正常	1	2	2	3	1	0
3	未使用	0	0	0	0	0	0
4	未使用	0	0	0	0	0	0
5	未使用	0	0	0	0	0	0

500 ベア情報

【図6】

図6

グループ番号	グループ状態	ベア集合	ジャーナル論理ボリューム番号	更新番号
1	正常	1,2	4	1
2	未使用	0	0	0

600 グループ情報

【図7】

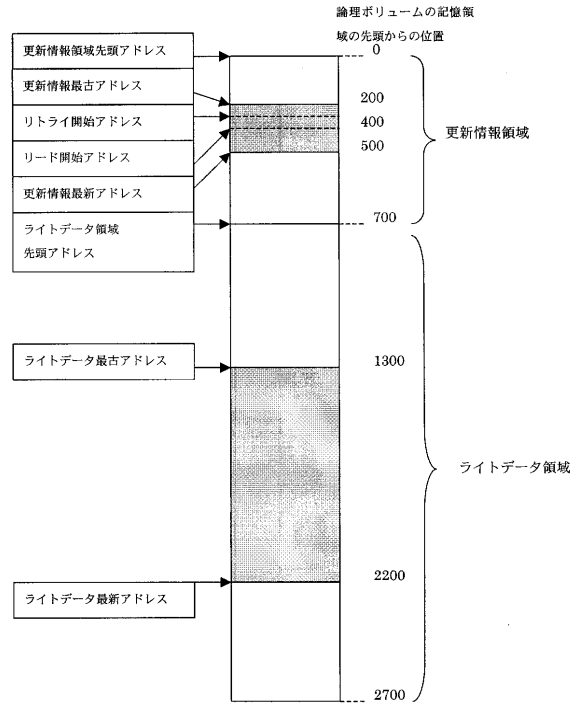
図7

	論理アドレス	論理ボリューム番号	論理ボリュームの記憶領域の先頭からの位置
更新情報領域先頭アドレス	4		0
ライトデータ領域先頭アドレス	4		700
更新情報最新アドレス	4		500
更新情報最古アドレス	4		200
ライトデータ最新アドレス	4		2200
ライトデータ最古アドレス	4		1300
リード開始アドレス	4		400
リトライ開始アドレス	4		300

700 ポインタ情報

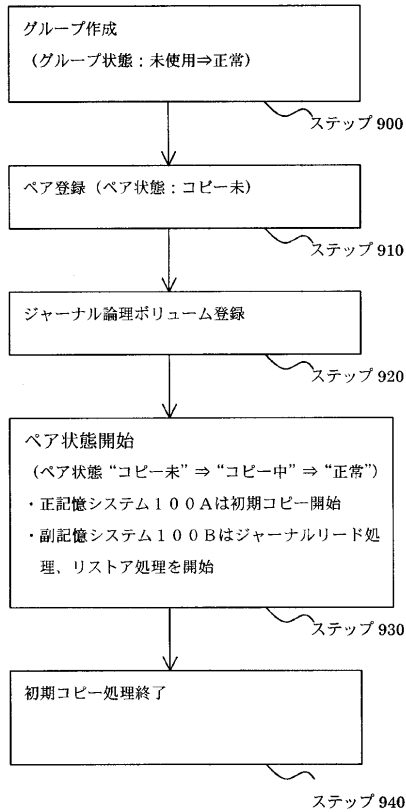
【図8】

図8



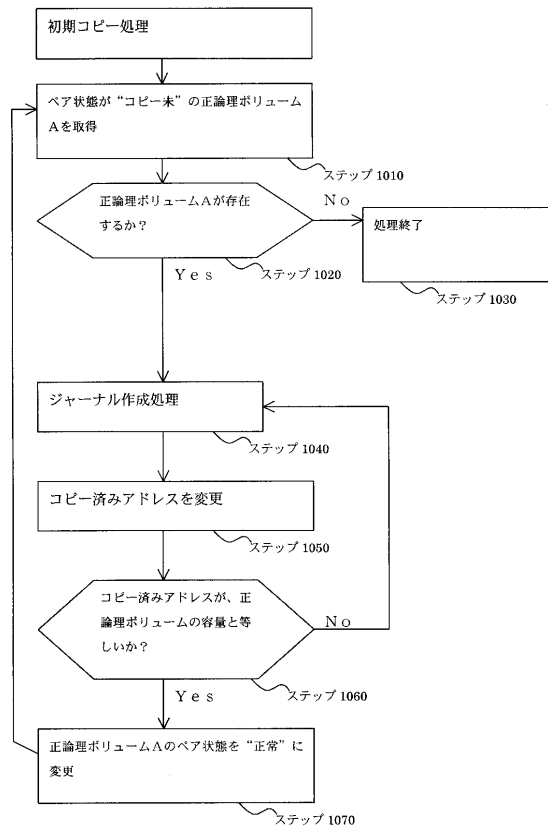
【図9】

図9

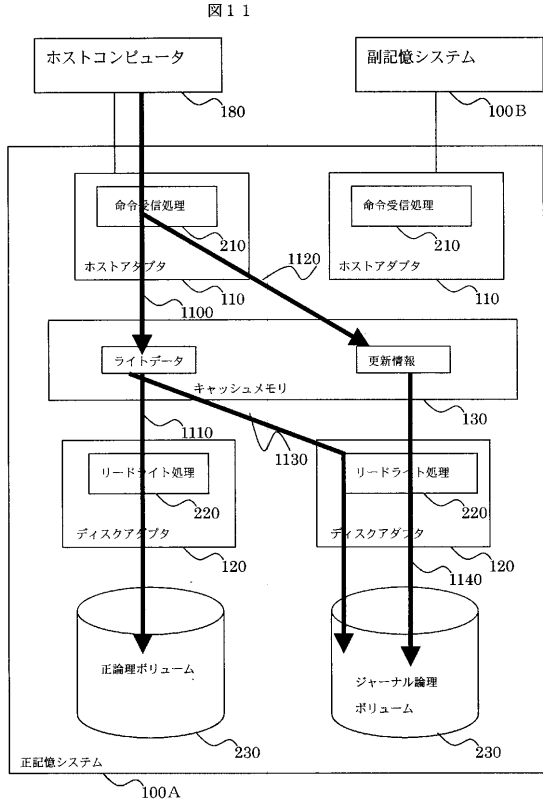


【図10】

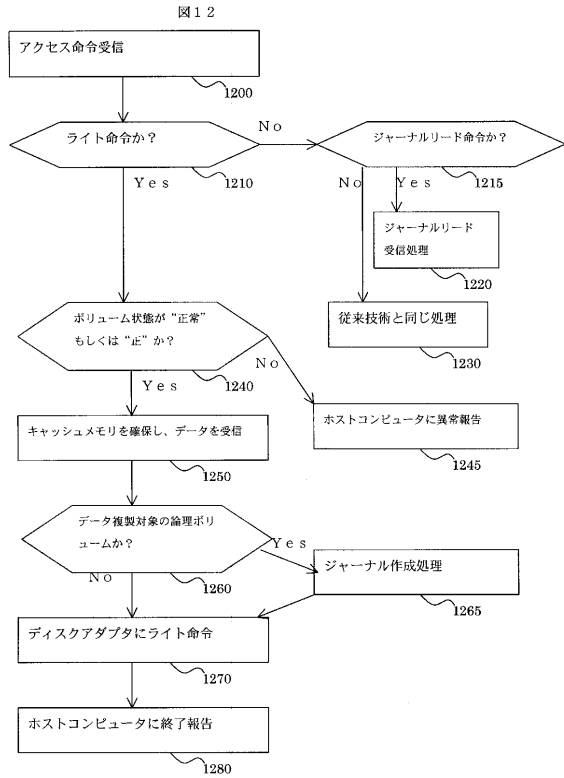
図10



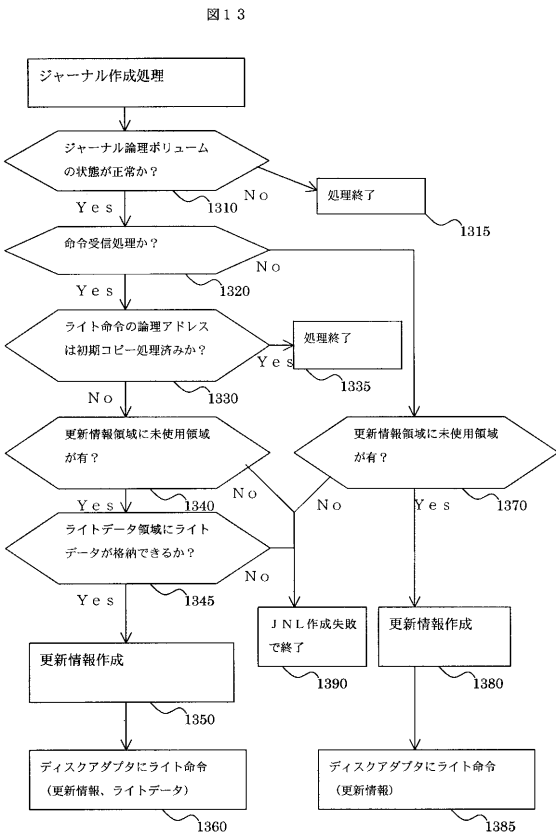
【図11】



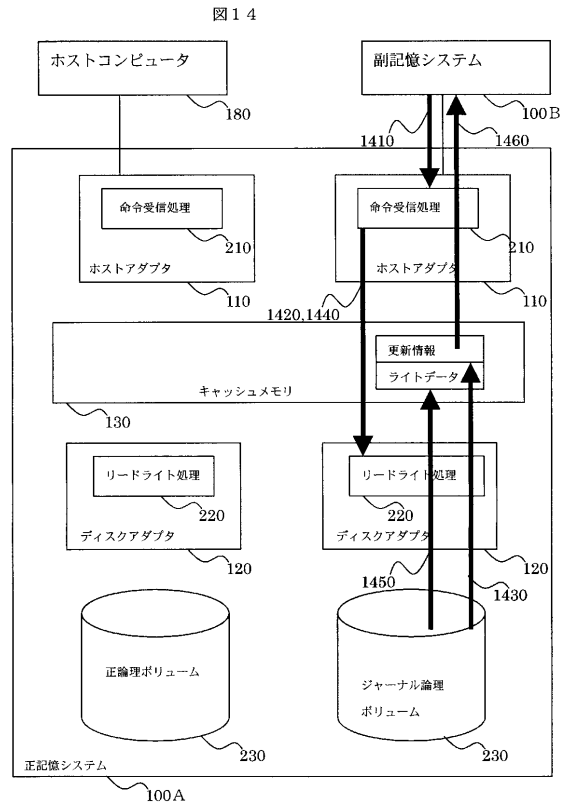
【図12】



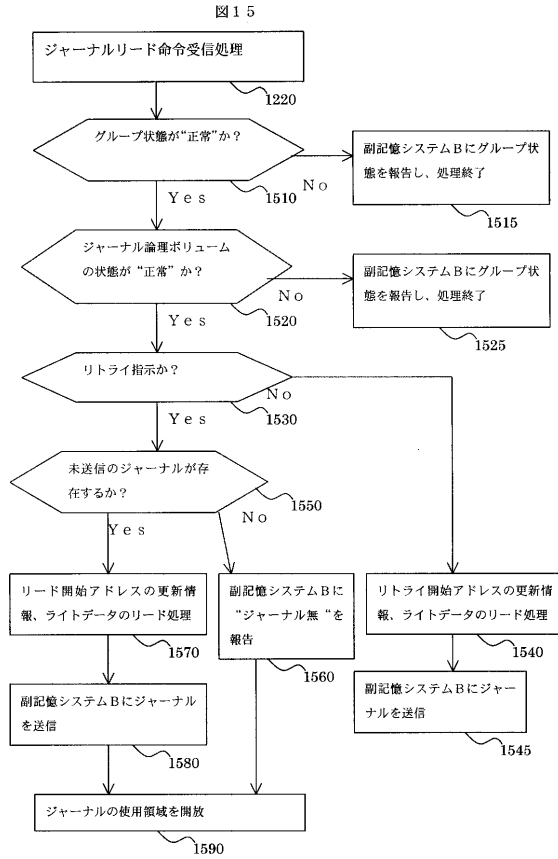
【図13】



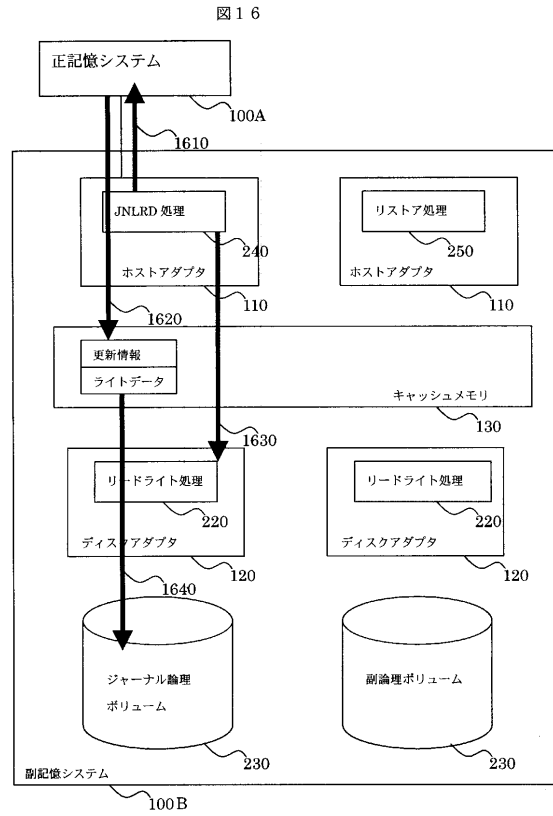
【図14】



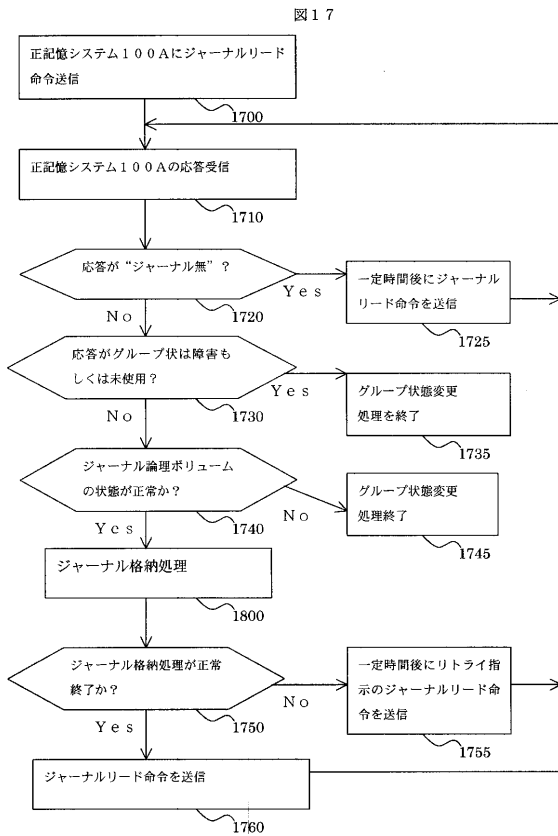
【図15】



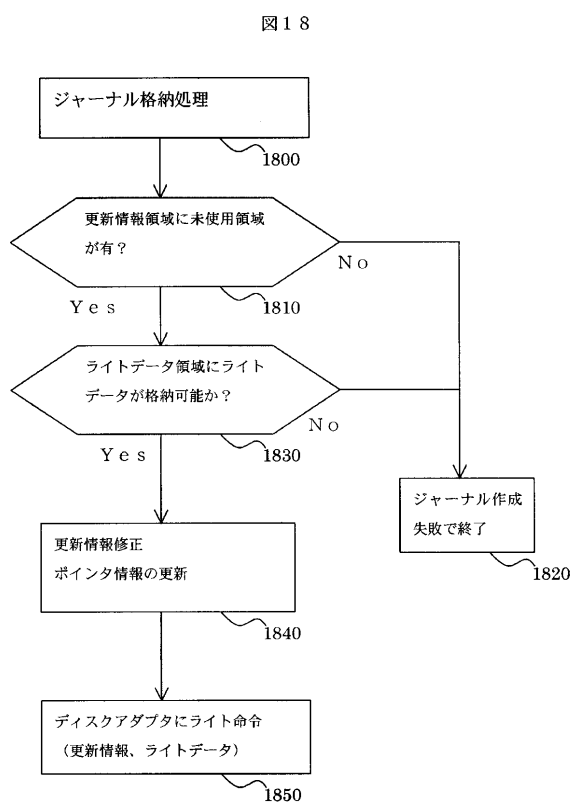
【図16】



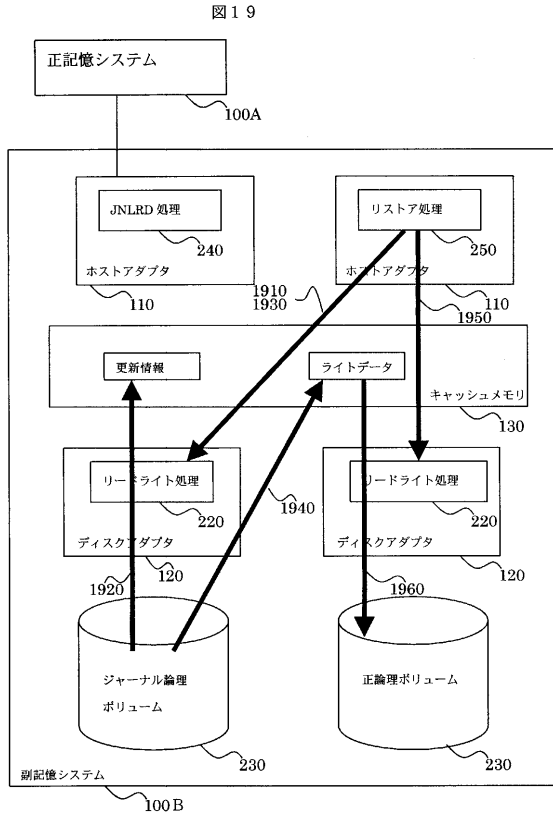
【図17】



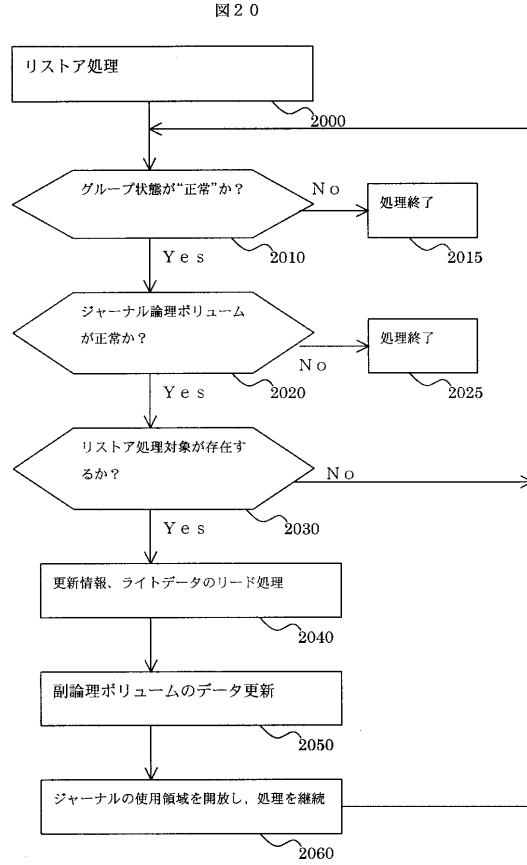
【図18】



【図19】



【図20】



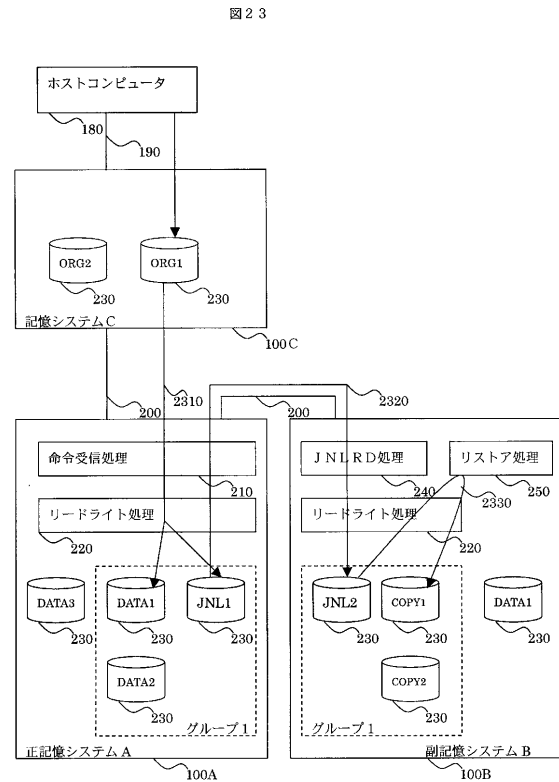
【図21】

図21

設定項目	設定値例
ライト命令を受信した時刻	1999/3/17 22:20:10
グループ番号	1
更新番号	4
ライト命令の論理アドレス	論理ボリューム番号: 1 論理ボリュームの記憶領域の先頭からの位置: 700
ライトデータのデータサイズ	300
ライトデータを格納したジャーナル論理ボリュームの論理アドレス	論理ボリューム番号: 4 論理ボリュームの記憶領域の先頭からの位置: 1500

310 更新情報

【図23】



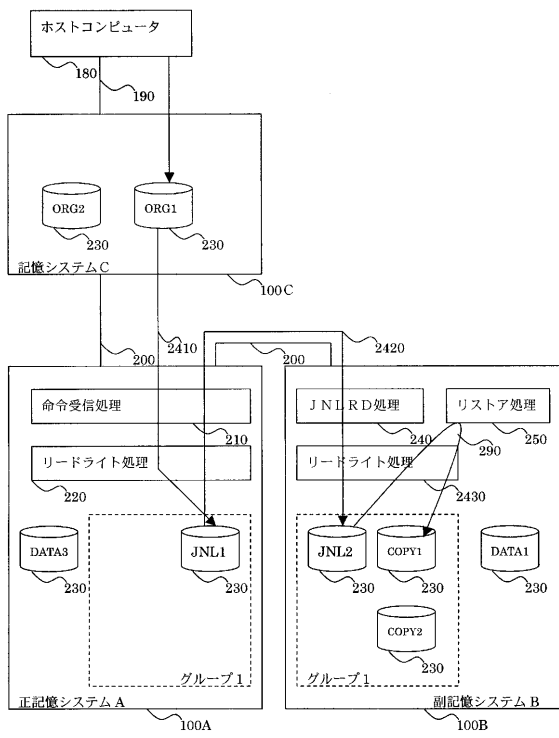
【図22】

図22

設定項目	設定値例
ライト命令を受信した時刻	1999/3/17 22:22:10
グループ番号	1
更新番号	4
ライト命令の論理アドレス	論理ボリューム番号: 1 論理ボリュームの記憶領域の先頭からの位置: 800
ライトデータのデータサイズ	100
ライトデータを格納したジャーナル論理ボリュームの論理アドレス	論理ボリューム番号: 4 論理ボリュームの記憶領域の先頭からの位置: 2200

【図24】

図24



フロントページの続き

(51)Int.Cl. F I
G 0 6 F 12/08 5 5 7

(72)発明者 武田 貴彦
神奈川県小田原市中里322番地2号 株式会社日立製作所 RAIDシステム事業部内

(72)発明者 佐藤 孝夫
神奈川県小田原市中里322番地2号 株式会社日立製作所 RAIDシステム事業部内

審査官 高瀬 勤

(56)参考文献 特開2000-181634(JP,A)
国際公開第02/031696(WO,A1)
特表2004-511854(JP,A)
特開2001-282628(JP,A)
特開2002-189570(JP,A)
特開平11-306058(JP,A)
特開平06-195250(JP,A)
特開平07-244597(JP,A)
特開2003-006016(JP,A)
特開2003-099306(JP,A)
特開昭58-175066(JP,A)
鈴木,中嶋,ストレージテクノロジー, FUJITSU, 日本, 富士通株式会社, 1995年
7月10日, 第46巻, 第4号, p.389-397

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

G06F 3/06

G06F 12/08

JSTPlus(JDreamII)