

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4477370号  
(P4477370)

(45) 発行日 平成22年6月9日 (2010.6.9)

(24) 登録日 平成22年3月19日 (2010.3.19)

(51) Int. Cl.

F I

G O 6 F 12/00 (2006.01)

G O 6 F 12/00 5 3 1 J

G O 6 F 3/06 (2006.01)

G O 6 F 12/00 5 4 5 A

G O 6 F 3/06 3 0 4 F

請求項の数 7 (全 27 頁)

(21) 出願番号 特願2004-23090 (P2004-23090)  
 (22) 出願日 平成16年1月30日 (2004.1.30)  
 (65) 公開番号 特開2005-216067 (P2005-216067A)  
 (43) 公開日 平成17年8月11日 (2005.8.11)  
 審査請求日 平成18年10月10日 (2006.10.10)

(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 110000350  
 ポレール特許業務法人  
 (74) 代理人 100068504  
 弁理士 小川 勝男  
 (74) 代理人 100086656  
 弁理士 田中 恭助  
 (74) 代理人 100094352  
 弁理士 佐々木 孝  
 (72) 発明者 川村 俊二  
 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】 データ処理システム

(57) 【特許請求の範囲】

【請求項 1】

第一の計算機と前記第一の計算機に接続された第一の記憶装置システムとを備えるプライマリサイトと、第二の計算機と前記第二の計算機に接続された第二の記憶装置システムとを備えるセカンダリサイトとを有するデータ処理システムにおいて、

前記第一の記憶装置システムと前記第二の記憶装置システムとは、通信線により接続され、

前記第一の記憶装置システムは、データの更新履歴をジャーナルとして記憶装置の複数の論理ボリュームに格納し、前記ジャーナルを第二の記憶装置システムに前記通信線を介して転送し、

閾値条件として、未転送ジャーナル量が設定閾値以上であるとき、または、未転送ジャーナル時間差が設定閾値以上であるとき、前記ジャーナルをある論理ボリュームに格納中に、他の論理ボリュームに格納用の論理ボリュームを切り替え、

切り替えのタイミングが、前記第二の記憶装置システムから、前記ジャーナルの送付を要求するコマンドを受信したタイミングであり、

前記第二の記憶装置システムの転送されたきた前記ジャーナルは、記憶装置の複数の論理ボリュームに格納され、

前記第二の記憶装置システムは、閾値条件として、未復元ジャーナル量が設定閾値以上であるとき、または、未復元ジャーナル時間差が設定閾値以上であるとき、前記ジャーナルをある論理ボリュームに転送中に、他の論理ボリュームに転送先の論理ボリュームを切

り替え、

切り替えのタイミングが、前記第一の記憶装置システムから、前記ジャーナルの転送が開始されたタイミングであることを特徴とするデータ処理システム。

【請求項 2】

前記第二の記憶装置システムは、格納された前記ジャーナルに基づいて、データ復元をおこなうことを特徴とする請求項 1 記載のデータ処理システム。

【請求項 3】

前記第二の記憶装置システムは、前記第一の記憶装置システムにおいて記録されたジャーナルに関する情報を取得し、

前記第二の記憶装置システムは、前記第一の記憶装置システムに対して、前記ジャーナルの送付を要求するコマンドを発行することを特徴とする請求項 1 記載のデータ処理システム。

10

【請求項 4】

第一の計算機と前記第一の計算機に接続された記憶装置システムとを備えるプライマリサイトと、第二の計算機と前記第二の計算機に接続された記憶装置システムとを備えるセカンダリサイトとを有するデータ処理システムにおいて、

前記第一の計算機と前記第二の計算機とは、第一の通信線により接続され、

前記第二の記憶装置システムと前記第二の記憶装置システムとは、第二の通信線により接続され、

前記第一の記憶装置システムは、データの更新履歴をジャーナルとして記憶装置の複数の論理ボリュームに格納し、

20

前記第一の計算機は、前記第一の記憶装置システムから前記ジャーナルに関する情報を取得して、前記第一の通信線を介して前記第二の記憶装置システムに送信し、

前記第一の記憶装置システムは、前記ジャーナルを第二の記憶装置システムに前記第二の通信線を介して転送し、

閾値条件として、未転送ジャーナル量が設定閾値以上であるとき、または、未転送ジャーナル時間差が設定閾値以上であるとき、前記ジャーナルをある論理ボリュームに格納中に、他の論理ボリュームに格納用の論理ボリュームを切り替え、

切り替えのタイミングが、前記第二の記憶装置システムから、前記ジャーナルの送付を要求するコマンドを受信したタイミングであり、

30

前記第二の記憶装置システムの転送されたきた前記ジャーナルは、記憶装置の複数の論理ボリュームに格納され、

前記第二の記憶装置システムは、閾値条件として、未復元ジャーナル量が設定閾値以上であるとき、または、未復元ジャーナル時間差が設定閾値以上であるとき、前記ジャーナルをある論理ボリュームに転送中に、他の論理ボリュームに転送先の論理ボリュームを切り替え、

切り替えのタイミングが、前記第一の記憶装置システムから、前記ジャーナルの転送が開始されたタイミングであることを特徴とするデータ処理システム。

【請求項 5】

前記第二の記憶装置システムは、前記第一の記憶装置システムに対して、前記ジャーナルの送付を要求するコマンドを発行することを特徴とする請求項 4 記載のデータ処理システム。

40

【請求項 6】

前記第二の記憶装置システムにおけるデータの復元は、前記転送されたジャーナルに基づき、前記第二の計算機で実行される復元プログラムによりおこなわれることを特徴とする請求項 4 記載のデータ処理システム。

【請求項 7】

第一の計算機と前記第一の計算機に接続された前記第一の記憶装置システムとを備えるプライマリサイトと、第二の計算機と前記第二の計算機に接続された前記第二の記憶装置システムとを備えるセカンダリサイトとを有するデータ処理システムにおいて、

50

前記第一の記憶装置システムと前記第二の記憶装置システムとは、通信線により接続され、

前記第一の記憶装置システムは、第一の記憶制御装置と第一の記憶装置とを有し、

前記第一の記憶制御装置は、データの更新履歴をジャーナルとして前記第一の記憶装置に記録するジャーナル取得プログラムと、前記ジャーナルを第二の記憶装置システムに前記通信線を介して転送するジャーナル転送プログラムとを実行し、

前記第二の記憶装置システムは、第二の記憶制御装置と第二の記憶装置とを有し、

前記第二の記憶装置システムは、データをジャーナルに基づいて復元するジャーナル復元プログラムと、前記ジャーナルを第一の記憶装置システムから転送された受け取るジャーナル転送プログラムとを実行し、

前記第一の記憶装置システムから前記第二の記憶装置システムに前記ジャーナル転送時に、

前記第一の記憶制御装置は、閾値条件として、未転送ジャーナル量が設定閾値以上であるとき、または、未転送ジャーナル時間差が設定閾値以上であるとき、前記第一の記憶装置のある論理ボリュームに格納中に、前記第一の記憶装置の他の論理ボリュームに格納用の論理ボリュームを切り替え、

切り替えのタイミングが、前記第二の記憶装置システムから、前記ジャーナルの送付を要求するコマンドを受信したタイミングであり、

前記第二の記憶制御装置は、閾値条件として、未復元ジャーナル量が設定閾値以上であるとき、または、未復元ジャーナル時間差が設定閾値以上であるとき、前記第二の記憶装置のある論理ボリュームに転送中に、前記第二の記憶装置の他の論理ボリュームに転送先の論理ボリュームを切り替え、

切り替えのタイミングが、前記第一の記憶装置システムから、前記ジャーナルの転送が開始されたタイミングであることを特徴とするデータ処理システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ処理システムに係り、特に、複数のサイトにジャーナルを分散して保持する技術に用いて好適なデータ処理システムに関する。

【背景技術】

【0002】

データ処理システムにおいて災害等による記憶装置システム内のデータ損失を避けることが要請されている。そのために、遠隔地の記憶装置システムにデータを複製する技術が重要になっている。ここで、記憶装置システムは、記憶制御装置とディスク装置などの記憶装置を含むものとする。

【0003】

記憶装置システムに格納されたデータを、別の記憶装置システムに複製する技術としては、特許文献1に開示された技術がある。特許文献1には、第一システムに含まれる計算機（以下、「プライマリホスト」という）のOSのデバイスドライバがライトシステムコールを受けると、データをローカルデータデバイスに書き込むとともに更新ログをライトログデバイスに格納し、プライマリホストのプログラムが第二のシステムに含まれる計算機（以下、「セカンダリホスト」という）のプログラムに更新ログ情報を転送し、セカンダリホストのプログラムが受け取った更新ログ情報を元に第二のシステムのデータデバイスのデータを更新する技術が開示されている。

【0004】

【特許文献1】米国特許6324654号公報

【発明の開示】

【発明が解決しようとする課題】

【0005】

上記従来技術では、プライマリホストとセカンダリホストの間で双方の記憶装置システ

10

20

30

40

50

ムに格納されたデータの転送がおこなわれる。このときに、ホストがデータ転送の経路として使用される。そして、各記憶装置システムに格納されたデータがホスト間の通信リンクを介して転送されるため、各ホストのCPU負荷、チャネル負荷、および、ホスト間を接続する回線のトラフィックが増加するという問題点があった。また、ログを採取することによるデバイスの負荷分散について考慮されていないため、ログ書き込み処理とログ読み出し処理がライトログデバイスへ集中するという問題点があった。

【0006】

本発明は、上記従来技術の問題点を解決するためになされたもので、その目的は、データの更新記録を採取してデータの復元をおこなうことのできるデータ処理システムにおいて、ホストとネットワークの負荷をかけることなく、しかも、データの更新復元に伴う特定の記憶装置への負荷集中を避けて、システムの処理能力が低下せずに、複数のサイトでデータの整合性を保証できるデータ処理システムを提供することにある。

10

【課題を解決するための手段】

【0007】

本発明のデータ処理システムは、プライマリサイトとセカンダリサイトで構成されたシステムであって、各サイトには、ホストと記憶装置システムを備えている。

【0008】

そして、プライマリサイトの記憶装置システム（以下、「プライマリ記憶装置システム」という）に格納されたデータの更新の情報をジャーナル（更新履歴）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メタデータの記録である。

20

【0009】

プライマリ記憶装置システムは、セカンダリサイトの記憶装置システム（以下、「セカンダリ記憶装置システム」という）と接続された通信線を介して、このジャーナルをセカンダリ記憶装置システムに転送する。セカンダリ記憶装置システムは、プライマリ記憶装置システムから受け取ったジャーナルを用いてセカンダリ記憶装置システムに格納されたデータを復元する（以下では、ジャーナルによりデータ復元をおこなうことを、「ジャーナル復元」ともいうことにする）。

【0010】

本発明では、単にデータではなく、ジャーナルをセカンダリサイトに転送して復元をおこなうので、障害が起こっても任意の時点のデータを速やかに復旧でき、データの整合性を保証することができる。

30

【0011】

また、本発明においては、プライマリ記憶装置システムは、ジャーナルを格納している論理ボリュームを複数持ち、現時点においてジャーナルの記録に用いている格納用論理ボリュームを切り替えることによつて、ジャーナルの転送に用いるジャーナル転送元の論理ボリュームとアクセスの集中を避け、負荷分散を図ることができる。

【0012】

同様にセカンダリ記憶装置システムではジャーナル転送に用いるジャーナル転送先の論理ボリュームを切り替え、ジャーナル復元に用いている論理ボリュームと別ボリュームとすることでアクセスの集中を避け、負荷分散を図ることができる。

40

【0013】

なお、ジャーナルの転送は、セカンダリ記憶装置システムがプライマリ記憶装置システムにジャーナル転送要求を発行することで実現する構成としても良い。

【0014】

また、プライマリホストおよびセカンダリホストが、各ホスト上で動作するプログラムに基づいて、各ホストに接続された記憶装置システムの状態を監視しておき、必要に応じて、それぞれのサイトのホストが、記憶装置システム間のデータ転送を、そのサイトの記憶装置システムに指示する構成としても良い。

【発明の効果】

50

## 【 0 0 1 5 】

本発明によれば、データの更新記録を採取してデータの復元をおこなうことのできるデータ処理システムにおいて、ホストとネットワークの負荷をかけることなく、しかも、データの更新・復元に伴う特定の記憶装置への負荷集中を避けて、システムの処理能力の低下が少なく、複数のサイトでデータの整合性を保証できるデータ処理システムを提供することができる。

## 【発明を実施するための最良の形態】

## 【 0 0 1 6 】

以下、本発明に係る各実施形態を、図 1 ないし図 1 7 を用いて説明する。

## 【 0 0 1 7 】

## 〔実施形態 1〕

以下、本発明に係る第一の実施形態を、図 1 ないし図 1 3 を用いて説明する。

## (1) データ処理システムの構成

先ず、図 1 を用いて本発明の第一の実施形態に係るデータ処理システムの構成について説明する。

図 1 は、本発明の第一の実施形態に係るデータ処理システムのハードウェア構成を示す図である。

図 2 は、本発明の第一の実施形態に係るデータ処理システムの機能構成を示す図である。

## 【 0 0 1 8 】

本実施形態のデータ処理システムは、各サイトが連携して処理をおこない、各サイトには、ホストと記憶装置システムを備えている。

## 【 0 0 1 9 】

ここで、第一のサイトを「プライマリサイト」といい、第二のサイトを「セカンダリサイト」といい、プライマリサイトからセカンダリサイトにジャーナルを転送する例を説明する。また、プライマリサイトに属するホストを「プライマリホスト」、プライマリサイトに属する記憶装置システムを「プライマリ記憶装置システム」、プライマリ記憶装置システムの記憶装置に格納されるジャーナルを、「プライマリジャーナル」、プライマリジャーナルを格納するためのボリュームを「プライマリジャーナルボリューム」ということにする。セカンダリサイトについても同様である。

## 【 0 0 2 0 】

さて、図 1 に示されるように、プライマリサイト 1 0 0 A は、プライマリホスト 1 1 0 A、プライマリ記憶装置システム 1 2 0 A を有していて、セカンダリサイト 1 0 0 B は、セカンダリホスト 1 1 0 B、セカンダリ記憶装置システム 1 2 0 B を有していて、それらに管理端末 1 3 0 が接続されている。

## 【 0 0 2 1 】

各ホスト 1 1 0 (プライマリホスト 1 1 0 A、セカンダリホスト 1 1 0 B) は、プロセッサ 1 1 1、主記憶装置 1 1 2、および、入出力処理装置 1 1 3 を有する計算機である。具体的には、ワークステーション、パーソナルコンピュータ、メインフレーム等である。

## 【 0 0 2 2 】

各記憶装置システム 1 2 0 は、記憶制御装置 1 4 0、一つ以上の記憶装置 1 2 1 および、保守端末 1 2 2 を有する。記憶装置 1 2 1 は、磁気ディスク記憶装置や光ディスク記憶装置などの補助記憶装置である。また、保守端末 1 2 2 は備えていないこともある。

## 【 0 0 2 3 】

記憶制御装置 1 4 0 は、ホスト入出力処理装置 1 4 1、キャッシュメモリ 1 4 2、ディスク入出力処理装置 1 4 3、プロセッサ 1 4 4、および、制御メモリ 1 4 5 を有する。

## 【 0 0 2 4 】

ホストは、記憶装置システムと LAN (Local Area Network) や SAN (Storage Area Network) といったホスト - 記憶装置システム間ネットワーク 1 5 0 によって互いに接続される。各ホストのプロセッサ 1 1 1 と主記憶装置 1 1 2 は、入出力処理装置 1 1 3 とホ

10

20

30

40

50

スト - 記憶装置システム間ネットワーク 150 を介して記憶装置システム 120 のホスト入出力処理装置 141 に接続される。

【0025】

記憶装置システム同士は、記憶装置システム間ネットワーク 160 を介して接続される。記憶装置システム間ネットワーク 160 は、一般に公衆回線などのグローバルネットワークである場合が多く、通信サービス提供者から有料で借用することが多い。システムの安全性のために（両サイトが同時に障害にあわないように）距離をおくためにグローバルネットワークを使用することが多いが、記憶装置システム同士が同じ部屋の中、同じビルの中、近隣ビルの中にある場合などはローカルネットワークを使用してもよい。ただし、このようなネットワークの形態によって本発明が限定されるものではない。

10

【0026】

管理端末 130 も、プロセッサや主記憶装置を有する計算機である。管理端末 130、プライマリホスト 110A、セカンダリホスト 110B、プライマリ記憶装置システム 120A、および、セカンダリ記憶装置システム 120B は、LAN や WAN 等のネットワーク 170 を介して相互に接続される。

【0027】

このようなデータ処理システムの機能的な構成を示すと図 2 のようになる。

【0028】

各記憶制御装置 140 では、記憶装置システム 120 間のデータ転送を制御するプログラムであるジャーナル処理管理プログラム 221、ジャーナル転送プログラム 222、ジャーナル取得・復元プログラム 223 がプロセッサ 144 上で実行される。これらのプログラムは制御メモリ 145 に格納される。

20

【0029】

ジャーナル処理管理プログラム 221 は各記憶装置システムが実行するジャーナル処理（ジャーナル取得処理、ジャーナル転送処理、ジャーナル復元処理）の管理をおこなう。詳細は後述する。また、ジャーナル処理中は各記憶制御装置 140 のジャーナル処理管理プログラム 221 間で随時通信がおこなわれ、ジャーナル処理に必要な管理情報を交換する。

【0030】

ジャーナル取得・復元プログラム 223 は、プロセッサ 144 にジャーナルの取得とジャーナルの復元を実行させるプログラムであり、ジャーナル取得プログラムとジャーナル復元プログラムから構成される。

30

【0031】

また、記憶制御装置 140 は、このようなジャーナル処理管理プログラム 221、ジャーナル転送プログラム 222、ジャーナル取得・復元プログラム 223 に関する処理の他にも、各ホストからの指示に基づいて、記憶装置 121 への入出力処理を実行する。

【0032】

記憶装置 121 には、一つ以上の論理記憶領域（論理ボリューム）が作成される。論理ボリュームは、記憶装置 121 が有する物理的な記憶領域と関連付けられる。これらの論理ボリュームはユーザの指定によりデータボリューム領域 225 とジャーナルボリューム領域 227 として使用される。なお、各ホスト 110 では、ユーザが使用するアプリケーションプログラム 211 や、記憶装置システムとのインタフェース制御をおこなう記憶装置制御プログラム 212 もホスト 110 が有するプロセッサ 111 で実行される。また、ジャーナル処理管理プログラムと記憶装置制御プログラム 212 は相互に情報をやり取りする。

40

【0033】

なお、データ複製のコピー元となるプライマリ記憶装置システム 120A が有するデータボリュームを、以下「PVOL」といい、PVOL に格納されるデータの複製先となるセカンダリ記憶装置システム 120B が有するデータボリュームを、以下「SVOL」ということにする。

50

## 【 0 0 3 4 】

管理端末 1 3 0 で実行される管理プログラム 2 3 1 は、本実施形態のデータ処理システムの各構成要素、具体的には各ホスト 1 1 0 や各記憶装置システム 1 2 0 の管理をおこなう。管理プログラム 2 3 1 は管理端末 1 3 0 の主記憶装置に格納される。

## 【 0 0 3 5 】

なお、ここまで説明した各プログラムは、コンパクトディスクや光磁気ディスクといった可搬媒体を用いて、あるいは、ネットワーク 1 7 0 を介して、各装置が有する記憶媒体にインストールされる。

## (11) データ処理システムに用いられるデータ構造

次に、図 3 ないし図 9 を用いて本実施形態のデータ処理システムに用いられるデータ構造について説明する。

図 3 は、ジャーナルグループ管理情報を示す図である。

図 4 は、データボリューム管理情報を示す図である。

図 5 は、ジャーナルメタ情報を示す図である。

図 6 は、データボリュームとジャーナルボリュームの対応を示す図である。

図 7 は、プライマリジャーナルボリューム内のジャーナルデータ領域 8 2 0 の内部構造を示す図である。

図 8 は、セカンダリジャーナルボリューム内のジャーナルデータ領域 8 2 0 の内部構造を示す図である。

図 9 は、ジャーナルボリュームの切り替え設定情報を示す図である。

## 【 0 0 3 6 】

ジャーナルグループ管理情報は、ジャーナルグループを管理するための情報であり、各記憶装置システム 1 2 0 の制御メモリ 1 4 5 に格納される。

## 【 0 0 3 7 】

ここで、「ジャーナルグループ」とは、データボリュームとそのデータのジャーナルを格納するジャーナルボリュームを関連付けたペアをいう。

## 【 0 0 3 8 】

ジャーナルグループ管理情報には、図 3 に示されるようにジャーナルグループ ID 4 1 0、最新ジャーナルシーケンス番号 4 2 0、データボリューム情報 4 3 0、ジャーナルボリューム情報 4 4 0、および、転送グループ情報 4 5 0 が含まれる。

## 【 0 0 3 9 】

ジャーナルグループ ID 4 1 0 は、ジャーナルグループを一意に決定する識別子である。最新ジャーナルシーケンス番号 4 2 0 は、ジャーナルグループのジャーナルに付与される連続する番号のうち、最も新しく付与された番号である。

## 【 0 0 4 0 】

データボリューム情報 4 3 0 には、ジャーナルグループに含まれる各データボリュームのデータボリューム管理情報 4 3 1 が含まれる。データボリューム管理情報の詳細については後述する。

## 【 0 0 4 1 】

ジャーナルボリューム情報 4 4 0 には、ジャーナルグループに含まれる各ジャーナルボリュームの情報と格納用ジャーナルボリューム ID 4 4 4 が含まれる。各ジャーナルボリューム毎の情報には、ボリュームを識別するボリューム ID 4 4 1、ジャーナルボリュームに格納されているジャーナルのうち、最も古いジャーナルのシーケンス番号を示す先頭ジャーナルシーケンス番号 4 4 2、および、最も新しいジャーナルのシーケンス番号を示す末尾ジャーナルシーケンス番号 4 4 3 が含まれる。

## 【 0 0 4 2 】

格納用ジャーナルボリューム ID 4 4 4 は、複数のジャーナルボリュームがジャーナルグループに含まれる場合に、次のジャーナルが格納されるジャーナルボリュームを示すためのものである。プライマリ記憶装置システム 1 2 0 A においてはジャーナル取得処理でジャーナル格納に用いられるジャーナルボリュームであり、セカンダリ記憶装置システム

10

20

30

40

50

120Bにおいてはジャーナル転送処理でジャーナルの転送先に用いられるジャーナルボリュームを指す。なお、ジャーナルボリュームが1つの場合には格納用ジャーナルボリュームID444には格納用かつ転送用に使用されるジャーナルボリュームのIDが設定される。

【0043】

転送グループ情報450には、ペアとなるジャーナルグループの識別子であるジャーナルグループID451、ジャーナル転送処理222によりどのジャーナルまでセカンダリ記憶装置システムに転送が完了したかを示す転送済みジャーナルシーケンス番号452、セカンダリ記憶装置システム120Bにおいてどのジャーナルまでジャーナル復元処理が完了したかを示す復元済みジャーナルシーケンス番号453、および、転送グループの状態を示すペア状態情報454が含まれる。なお、転送済みジャーナルシーケンス番号452の次のシーケンス番号を有するジャーナルが一番古いジャーナルとなる。セカンダリ記憶装置システム（第二および第四の実施形態ではプライマリ記憶装置システム）は、この一番古いジャーナルが含まれるジャーナルボリュームを転送用のジャーナルボリュームとして特定し、使用する。転送用のジャーナルボリュームと格納用のジャーナルボリュームが一致する場合は、この一番古いジャーナルを含むジャーナルボリュームと格納用ジャーナルボリュームID444で示されるジャーナルボリュームとは一致する。一方、切り替え処理の結果またはあらかじめ転送用のジャーナルボリュームと格納用のジャーナルボリュームとを異ならせた場合、一番古いジャーナルを含むジャーナルボリュームと格納用ジャーナルボリュームID444で示されるジャーナルボリュームとは一致しない。これはセカンダリ記憶装置システム120Bでも同様である。

【0044】

ここで、「転送グループ」とは、転送元のジャーナルグループと転送先のジャーナルグループを関連付けてペアとしたものである。

【0045】

復元済みジャーナルシーケンス番号453は、セカンダリ記憶装置システム120Bのジャーナル処理管理プログラム221からプライマリ記憶装置システム120Aのジャーナル処理管理プログラム221へ通知される。

【0046】

ペア状態454には、ジャーナルグループ内の全データボリュームが複製状態にある「PAIR」、ジャーナルグループ内の一つ以上のデータボリュームが差分コピー（差分コピーについては後述）をおこなっている「COPY」がある。また、ジャーナルグループ内の全データボリュームの複製がおこなわれず、サスペンド状態にある状態で、ジャーナルグループ内でのデータ一貫性がある状態「SUSPEND」と、ジャーナルグループ内でのデータ一貫性がない状態「SUSPEND-E」がある。

【0047】

データボリューム管理情報は、データボリュームを管理するための情報であり、図4に示されるようにデータボリュームの記憶装置システム内で識別するためのボリュームID510、および、ペアであるデータボリュームの情報であるペアボリューム情報520が含まれる。

【0048】

ペアボリューム情報520には、記憶装置システムID521、ペアであるデータボリュームを識別するデータボリュームID522、複製状態を示すペア状態523、差分ビットマップ有効フラグ524、および、差分ビットマップ525が含まれる。

【0049】

記憶装置システムID521はペアであるデータボリュームが存在する記憶装置システムの識別子であり、データボリュームID522は、その記憶装置システムの中でのボリューム識別子である。記憶装置システムID521とデータボリュームID522を組み合わせることにより、データボリュームが一意に決まる。

【0050】

10

20

30

40

50



ペア状態 5 2 3 には、データボリュームが複製状態にある（ボリューム内でデータ一貫性がある状態）「PAIR」、差分コピーをおこなっている「COPY」、ボリュームやパスの閉塞等によりコピー処理が停止されサスペンド状態にある「SUSPEND」のうちいずれかの状態を指す。

【0051】

差分ビットマップ有効フラグ 5 2 4 は、差分ビットマップの値が有効であるかどうかを示す。差分ビットマップ 5 2 5 は、PVOLとSVOLのデータに差異がある領域を示す情報である。データボリューム領域を複数の領域に分割し、SUSPEND中にデータボリュームに更新があった場合には、更新のあった領域を示すビットをONにする。SUSPEND後に、PVOLとSVOLの各々の差分ビットマップ 5 2 5 のORを取ったビットマップに基づいてビットがONの領域のみをコピーする（差分コピー）ことで、ペアを複製状態に戻すことができる。差分コピーによりコピー転送量を減らすことができる。差分コピーではコピーが完了した領域はビットをOFFにし、全ビットがOFFになると差分コピーが完了したこととなる。また、ペア生成時には差分ビットマップの全てONにして差分コピーをおこなう（形成コピー）ことで、PVOLの全領域をSVOLにコピーすることができる。

【0052】

ジャーナルメタ情報は、データとジャーナルを関連付けるための管理情報であり、図 5 に示されるようにデータボリューム情報 7 1 0 とジャーナル情報 7 2 0 を含む。

【0053】

データボリューム情報 7 1 0 には、データが更新された時刻を示す更新時刻 7 1 1、データが更新されるデータボリュームがジャーナルグループ内のどのデータボリュームであるかを示す格納データボリュームのジャーナルグループ内オフセット 7 1 2、および、データボリューム上のデータを格納する先頭アドレスを示すデータ格納アドレス 7 1 3 を含む。

【0054】

ジャーナル情報 7 2 0 には、ジャーナルボリューム上のジャーナルデータを格納する先頭アドレスを示すジャーナルデータ格納アドレス 7 2 1、ジャーナルデータのデータ長 7 2 2、および、ジャーナル取得の際に割り当てられたジャーナルグループ内でのジャーナルの通番であるジャーナルシーケンス番号 7 2 3 を含む。

【0055】

このジャーナルメタ情報により関連付けられるデータボリュームとジャーナルボリュームは、図 6 に示される如くである。

【0056】

一般に、PVOL、SVOLおよびジャーナルボリュームは各々予め定められた論理ブロック単位で管理される（例えば 5 1 2 K B）。論理ブロックの各々には、論理ブロックアドレス（以下「LBA」と記す）が付与されている。

【0057】

プライマリジャーナルボリューム 2 6 7 A は、メタデータ領域 8 1 0、および、ジャーナルデータ領域 8 2 0 を有する。ジャーナルデータ領域 8 2 0 には、先に説明したジャーナルデータ 8 2 1 A、すなわち、ライトコマンドによってPVOLに書き込まれたデータ 8 3 1 A のコピーが格納される。メタデータ領域 8 1 0 には、先に説明したメタデータ 8 1 1 A が格納される。メタデータには、更新データのデータ格納アドレス 8 1 2 A、および、ジャーナルデータの格納アドレス 8 1 3 A が含まれる。

【0058】

セカンダリジャーナルボリューム 2 6 7 B も、プライマリジャーナルボリューム 2 6 7 A と同様にメタデータ領域 8 1 0、および、ジャーナルデータ領域 8 2 0 を有する。メタデータ領域 8 1 0 には、プライマリジャーナルボリュームのメタデータ領域から転送されたメタデータ 8 1 1 B が格納される。ジャーナルデータ領域 8 2 0 には、プライマリジャーナルボリューム 2 6 7 A のジャーナルデータ領域から転送されたジャーナルデータ 8 2

10

20

30

40

50

1 B (メタデータに対応する) が格納される。

【0059】

メタデータ 8 1 1 B は P V O L でおこなわれたデータ更新の情報を持ち、そのアドレス情報 8 1 3 B は対応するジャーナルデータ 8 2 1 B のアドレスを示す。さらに、ジャーナルデータ 8 2 1 B を、セカンダリジャーナルボリューム 2 6 7 B のジャーナルデータ領域 8 2 0 から、アドレス 8 1 2 B に対応する S V O L 2 6 6 のアドレスへコピーすることによって、P V O L 2 6 5 での更新を S V O L 2 6 6 へ反映することができる。

【0060】

各アドレスは L B A により、データ長は論理ブロック数により、各々表わすことができる。また、データが格納されている場所は、データが格納された領域 (ジャーナルデータ領域またはメタデータ領域) のベースアドレス (先頭 L B A ) との差分 (オフセット) で表されても良い。本実施形態において、メタデータのデータ長は一定 (例えば 6 4 バイト) であるが、ジャーナルデータのデータ長は、ライトコマンドで更新されるデータに依存するので一定ではない。

【0061】

ジャーナルグループ定義時に、各記憶装置システム 1 2 0 は、設定されるジャーナルボリューム 2 6 7 に対して、メタデータ領域 8 1 0 およびジャーナルデータ領域 8 2 0 の設定をおこなう。具体的には、各領域の先頭 L B A およびブロック数が設定される。

【0062】

次に、プライマリジャーナルボリューム内のジャーナルデータ領域 8 2 0 の内部構造について説明する。

【0063】

プライマリジャーナルボリュームが有するジャーナルデータ領域 8 2 0 は、図 7 に示されるようにジャーナルデータが格納されているジャーナル格納済み領域 9 1 0 と、ジャーナルデータが格納されていない、あるいは、パージ可能なジャーナルデータが格納されているパージ済み領域 9 2 0 とに区別される。パージ済み領域 9 1 0 は、ジャーナルデータをセカンダリサイトに転送したために、そこに格納されているジャーナルデータを開放しても良いことになった領域であり、P V O L 6 2 5 の新たなジャーナルデータの格納に使用することができる。

【0064】

図 7 では、ジャーナルデータ領域 9 0 0 A とジャーナルデータ領域 9 0 0 B は、別個の論理ボリュームに格納されていることを示している。

【0065】

ジャーナルデータ領域 9 0 0 A には、ボリューム内先頭ジャーナルシーケンス番号 9 3 1 A からボリューム内末尾ジャーナルシーケンス番号 9 3 2 A までのジャーナルが、ジャーナルデータ領域 9 0 0 B にはボリューム内先頭ジャーナルシーケンス番号 9 3 1 B からボリューム内末尾ジャーナルシーケンス番号 9 3 2 B までのジャーナルが格納されている。ジャーナルを格納するときには、ボリューム内先頭に書き込まれていく。

【0066】

ジャーナルボリュームはサイクリックバッファと同じように繰り返し使用される。つまりジャーナルの末尾の論理ブロックまで使用すると、先頭の論理ブロックが再度使用される。ただし、ジャーナルグループに複数のジャーナルボリュームが含まれる場合は、ジャーナルの末尾の論理ブロックまで使用すると、次のジャーナルボリュームの先頭論理ブロックが使用される。最後のジャーナルボリュームの末尾論理ブロックまで使用されると最初のジャーナルボリュームの先頭論理ブロックに戻る。なお、ジャーナルボリュームの末尾論理ブロックまで使用する前に、途中の論理ブロックから次のジャーナルボリュームの先頭論理ブロックに移っても良い。ジャーナル格納先を次のジャーナルボリュームに移すことを「ジャーナルボリューム切り替え処理」という。

【0067】

図 7 では、ボリューム内先頭ジャーナルシーケンス番号 9 3 1 A までジャーナル領域 9

10

20

30

40

50

00Aを使用した後、ジャーナルボリュームが切り替わり、ジャーナルデータ領域900Bの先頭に移ったことになる。このためボリューム内先頭ジャーナルシーケンス番号931Aとボリューム内末尾ジャーナルシーケンス番号932Bは連続していることになる。なお、このジャーナルボリューム切り替え処理については、後に詳述する。

#### 【0068】

格納済みジャーナルシーケンス番号911は、最も新しいジャーナルを示す。次に取得されるジャーナルは、格納済みジャーナルシーケンス番号911に1を足した値のジャーナルシーケンス番号が付与され、ジャーナルデータ領域900Bのページ領域に格納される。ここでは格納用ジャーナルボリュームはジャーナルデータ領域900Bを有するジャーナルボリュームをいう。転送済みジャーナルシーケンス番号921に1を足したシーケンス番号のジャーナルが最も古いジャーナルである。転送済みジャーナルシーケンス番号921と格納済みジャーナルシーケンス番号911が等しい場合には、ジャーナルが空であることを意味する。

10

#### 【0069】

次に、セカンダリジャーナルボリューム内のジャーナルデータ領域820の内部構造について説明する。

#### 【0070】

セカンダリジャーナルボリュームが有するジャーナルデータ領域820は、図8に示されるように既にSVOL626へのジャーナル復元に使用されたジャーナルデータが格納されている（あるいは、ジャーナルデータが格納されていない）ページ済み領域1030、SVOL626へのジャーナル復元の対象として指定されたジャーナルデータが格納されている復元中領域1020、ジャーナル復元の対象となっておらず、プライマリジャーナルボリュームからジャーナルの転送が終わったジャーナルデータが格納されている転送済み領域1011、および、プライマリジャーナルボリュームから転送中のジャーナルデータが格納される転送中領域1010とに区別される。

20

#### 【0071】

ジャーナルデータ領域1000Aには、ボリューム内先頭ジャーナルシーケンス番号931Cからボリューム内末尾ジャーナルシーケンス番号932Cまでのジャーナルが、ジャーナルデータ領域1000Bにはボリューム内先頭ジャーナルシーケンス番号931Dからボリューム内末尾ジャーナルシーケンス番号932Dまでのジャーナルが格納されている。ここではボリューム内先頭ジャーナルシーケンス番号931Cを格納した後にジャーナルボリューム切り替えが起こったことになり、ボリューム内先頭ジャーナルシーケンス番号931Cとボリューム内末尾ジャーナルシーケンス番号932Dは連続する。

30

#### 【0072】

転送予定ジャーナルシーケンス番号1012は、プライマリジャーナルボリュームから転送中のジャーナルの先頭ジャーナルシーケンス番号を示す。ジャーナルは転送予定ジャーナルシーケンス番号1012以降のジャーナルが次に転送され、ジャーナルデータ領域1000Bのページ領域に格納される。ここでは、図3の格納用ジャーナルボリューム444に格納されるボリュームIDは、転送先のボリュームであるジャーナルデータ領域1000Bを有するジャーナルボリュームを指す。転送済みジャーナルシーケンス番号1013は、最後に転送処理が完了したジャーナルのシーケンス番号を指す。

40

#### 【0073】

復元予定ジャーナルシーケンス番号1021は、ジャーナル復元処理の対象とされたジャーナルの先頭シーケンス番号を指す。また、復元済みジャーナルシーケンス番号1022より最後に復元処理が完了したジャーナルのシーケンス番号を指す。

#### 【0074】

次に、ジャーナルボリューム切り替え設定情報について説明する。

#### 【0075】

ジャーナルボリューム切り替え設定情報は、ジャーナルボリュームを切り替えるための管理情報である。

50

## 【 0 0 7 6 】

本発明の処理では、プライマリサイトとセカンダリサイトの両方でジャーナルボリュームを切り替えることが可能である。

## 【 0 0 7 7 】

プライマリサイト 1 0 0 A のジャーナルボリュームを切り替えは、転送用に用いている（転送元）論理ボリュームと現時点でのジャーナルの格納に用いている論理ボリュームの負荷を分散するために、現時点でのジャーナルの格納に用いている論理ボリュームを切り替える。

## 【 0 0 7 8 】

また、セカンダリサイト 1 0 0 B のジャーナルボリュームを切り替えは、転送用に用いている（転送先）論理ボリュームと現時点でのジャーナルによるデータ復元に用いている論理ボリュームの負荷を分散するために、転送用に用いている論理ボリュームを切り替える。

10

## 【 0 0 7 9 】

ジャーナルボリューム切り替え設定情報は、図 9 に示されるように切り替え判定有効フラグ 1 1 1 0 と、判定情報 1 1 2 0 と、条件外時動作 1 1 3 0 よりなる。

## 【 0 0 8 0 】

切り替え判定有効フラグ 1 1 1 0 は、切り替え判定処理をおこなうか否かのフラグであり、ジャーナルボリューム切り替え判定をおこなわない場合には、判定有効フラグ 1 1 1 0 は OFF になる。

20

## 【 0 0 8 1 】

プライマリサイト 1 0 0 A での切り替え判定処理は、任意のタイミングでおこなうことができるが、あるボリュームのジャーナルの転送が終わって、次に転送すべきジャーナルの格納されているジャーナルボリュームが、現時点でのジャーナルを格納しているボリュームに使われているときに、切り替え判定処理がおこなうことが有効的である。また、セカンダリサイト 1 0 0 B での切り替え判定処理も、任意のタイミングでおこなうことができるが、あるボリュームのジャーナルによるデータ復元が終わって、次に復元に用いられるジャーナルの格納されているジャーナルボリュームが、転送先のボリュームに使われているときに、切り替え判定処理がおこなうことが有効的である。ここで、ジャーナルデータは、図 7 および図 8 に用いられるように、シーケンシャルな構造を持つことに注意する。なお、このような切り替え判定処理は、後に詳説する。

30

## 【 0 0 8 2 】

判定情報 1 1 2 0 には、未転送ジャーナル量に関する設定値と、未転送ジャーナルのうち最古のジャーナルが更新された時刻と判定時の時間差に関する設定値が含まれる。各々の設定値には、その判定基準を有効にするか否かの有効フラグ 1 1 2 1、1 1 2 3、および、格納用ボリューム切り替えするための閾値条件 1 1 2 2、1 1 2 4 が含まれる。

## 【 0 0 8 3 】

例えば、プライマリジャーナルボリュームの場合には、未転送ジャーナル量が少ない場合には、現時点でのジャーナルの格納に用いている論理ボリュームと転送に用いている論理ボリュームが同一の論理ボリュームとなってもジャーナル処理の負荷はあまり問題にならないと考えられる。また、セカンダリジャーナルボリュームの場合には、復元予定のジャーナル量が少ない場合には、現時点でのジャーナルによるデータ復元に用いている論理ボリュームと転送に用いている論理ボリュームが同一の論理ボリュームとなってもジャーナル処理の負荷はあまり問題にならないと考えられる。

40

## 【 0 0 8 4 】

したがって、プライマリジャーナルボリュームの場合には、閾値条件 1 1 2 2 として、未転送ジャーナル量（格納済みジャーナルシーケンス番号 9 1 1 と転送済みジャーナルシーケンス番号 9 2 1 の差）が設定閾値以上であること、閾値条件 1 1 2 4 として、未転送ジャーナル時間差（転送済みジャーナルシーケンス番号 9 2 1 に 1 加えたシーケンス番号のジャーナルの更新時刻と判定時の時刻の差）が設定閾値以上であることなどが考えられ

50

る。また、セカンダリジャーナルボリュームの場合には、閾値条件 1 1 2 2 として、未復元ジャーナル量（転送済みジャーナルシーケンス番号 1 0 1 3 と復元済みジャーナルシーケンス番号 1 0 2 2 の差）が設定閾値以上であることや、閾値条件 1 1 2 4 として、未復元ジャーナル時間差（復元済みジャーナルシーケンス番号 1 0 2 2 に 1 加えたシーケンス番号のジャーナルの更新時刻と判定時の時刻の差）が設定閾値以上であることがある。

【 0 0 8 5 】

このような場合には、プライマリジャーナルボリュームの場合では、格納用ジャーナルボリューム切り替えをおこない、セカンダリジャーナルボリュームでは、転送用ボリュームの切り替えをおこなう。

【 0 0 8 6 】

条件外動作 1 1 3 0 には、条件にあわないときのシステム動作を記述する。条件に合わない場合の動作としては、プライマリジャーナルボリュームの場合では、（ 1 ）ジャーナル転送を中断する、あるいは、（ 2 ）ジャーナル格納に用いているジャーナルボリュームを用いてジャーナル転送をおこない、また、セカンダリジャーナルボリュームの場合では、（ 1 ）ジャーナル復元を中断する、（ 2 ）転送先に用いているジャーナルボリュームを用いてジャーナル復元をおこなう。

【 0 0 8 7 】

プライマリサイト 1 0 0 A では、ジャーナル切り替え設定情報は、ジャーナルグループ設定時等にホスト 1 1 0、管理端末 1 3 0、あるいは、保守端末 1 2 2 の G U I を介してユーザによって設定され、プライマリ記憶装置システム 1 2 0 A の制御メモリ 1 4 5 に格納される。セカンダリサイト 1 0 0 B でも、ジャーナル切り替え設定情報は、同様に設定され、セカンダリ記憶装置システム 1 2 0 B の制御メモリ 1 4 5 に格納される。

（ III ）データ処理システムの処理概要

（ III - 1 ）データ処理システムの処理概要

まず、図 1 0 を用いて本発明の第一の実施形態に係るデータ処理システムの処理概要について説明する。

【 0 0 8 8 】

図 1 0 は、第一の実施形態に係るデータ処理システムの概略処理を示すフローチャートである。

【 0 0 8 9 】

まず、ユーザは、ホスト 1 1 0、管理端末 1 3 0、あるいは、保守端末 1 2 2 が持つ G U I （ Graphical User Interface ）等を用いて記憶装置システムにペア生成コマンドを入力する（ステップ 3 0 1 ）。

【 0 0 9 0 】

ペア生成コマンドは、データ複製のコピー元となるプライマリ記憶装置システム 1 2 0 A が有する P V O L 6 2 5 と、 P V O L に格納されるデータの複製先となるセカンダリ記憶装置システム 1 2 0 B が有する S V O L 6 2 6 とを、ペアとして関連付けるコマンドである。

【 0 0 9 1 】

次に、各々のサイトのジャーナル処理管理プログラム 2 2 1 を用いて、プライマリサイト 1 0 0 A において、記憶装置システム 1 2 0 A で指定された P V O L に対応するジャーナルを格納するボリュームを割り当てるように制御し、また、セカンダリサイト 1 0 0 B において、セカンダリ記憶装置システム 1 2 0 B で指定された S V O L 6 2 6 に対応するジャーナルを格納するジャーナルボリュームを割り当てるように制御する（ステップ 3 0 2、ステップ 3 0 3 ）。

【 0 0 9 2 】

プライマリサイト 1 0 0 A では、 P V O L 6 2 5 とそれに割り当てられたジャーナルボリューム 6 2 7 A により、ジャーナルグループが形成され、セカンダリサイト 1 0 0 B では、 S V O L 6 2 6 とそれに割り当てられたジャーナルボリューム 6 2 7 B により、ジャーナルグループが形成される。ジャーナルグループを形成する際に、ジャーナルボリュー

10

20

30

40

50

ムに複数のボリュームの集合を割り当てることができる。

【0093】

ペア生成コマンドは、P V O LのジャーナルグループとS V O L 6 2 6のジャーナルグループも関連付け、転送グループとする(ステップ304)。ジャーナルグループの設定時に、ジャーナルボリューム切り替え情報を設定しても良い。ジャーナルボリューム切り替え設定の詳細は後述する。

【0094】

なお、ジャーナルグループ形成の際には、単一のデータボリュームに限らず複数のデータボリュームの集合を割り当てることができる。このP V O L集合におけるデータの更新と同様にS V O L集合でもデータが更新されるため、このデータボリューム集合内でデータの一貫性が保持される。

10

【0095】

次に、転送グループが作成された後、ジャーナル処理が実行される(ステップ305)。ジャーナル処理は、ジャーナル取得処理、ジャーナル転送処理、ジャーナル復元処理の各処理をいう。ジャーナル取得処理を開始は、ユーザからのジャーナルの取得を指示するコマンド(以下「ジャーナル取得開始コマンド」)をプライマリ記憶装置システムが受信したことを契機に、プライマリ記憶装置システム120でおこなわれる。ジャーナル処理の詳細については後述する。

【0096】

一方、ジャーナル取得処理が開始される前にP V O L 6 2 5に格納されていたデータは、ジャーナル転送処理が開始されてもセカンダリ記憶装置システムには転送されない。別途P V O L 6 2 5からS V O L 6 2 6へこれらのデータ(以下「イニシャルデータ」)をコピーする必要がある。イニシャルデータをコピーする処理を「形成コピー」という。本実施形態においては、イニシャルデータをP V O L 6 2 5からS V O L 6 2 6に転送する形成コピー処理が実行される(ステップ306)。イニシャルデータは、P V O L 6 2 5のボリューム先頭領域から末尾まで転送される。

20

【0097】

(III-2)ジャーナル処理の詳細

次に、図11を用いてジャーナル処理の詳細について説明する。

図11は、本発明の第一の実施形態に係るジャーナル処理の動作を示す図である。

30

【0098】

記憶装置システム120Aおよび120Bにより、それぞれジャーナル処理管理プログラム221を実行してジャーナル処理を制御する。

【0099】

プライマリ記憶装置システム120Aは、ジャーナル取得・復元プログラム223の内、ジャーナル取得プログラム623を実行する。ジャーナル取得プログラム623を実行することによって、プライマリ記憶装置システム120Aは、P V O L 6 2 5に書き込まれるデータのコピーをジャーナルデータとして、ジャーナルボリューム627Aに格納する。また、プライマリ記憶装置システム120Aはジャーナルの一部としてメタデータもジャーナルボリューム627Aに格納する。これらの処理がジャーナル取得処理である。

40

【0100】

一方、セカンダリ記憶装置システム120Bは、ジャーナル取得・復元プログラム223の内、ジャーナル復元プログラム624を実行してジャーナル復元処理をおこなう。ジャーナル復元プログラム624は、ジャーナルボリューム627Bに格納されたジャーナルに基づきデータを復元して、P V O L 6 2 5で更新されたデータをS V O L 6 2 6に反映する。

【0101】

これらのような処理を、図2に示したシステムで実行する場合には、図6に示されるようになる。

【0102】

50

プライマリ記憶装置システム120AでPVOL625のジャーナル取得処理が開始されると、プライマリ記憶装置システム120Aは、プライマリホスト110AからPVOL625への書き込み処理（矢印601）に応じて、ジャーナルを作成し、作成したジャーナルをジャーナルボリューム627Aに格納する（矢印602）。ここでジャーナル取得プログラムは、プライマリ記憶装置システム120Aの制御メモリ145上に格納されているジャーナルグループ管理情報から、最新ジャーナルシーケンス番号や格納用ジャーナルボリュームID等の情報を取得し、ジャーナル格納先を決定し、かつ、メタデータ情報を作成する。

#### 【0103】

セカンダリ記憶装置120Bはジャーナル処理管理プログラム222を実行してプライマリ記憶装置120Aから、ジャーナル作成状況に関する情報（例えば、ジャーナルボリューム内のジャーナルの容量、ジャーナルの一番古い時刻等）を取得する（矢印603）。

10

#### 【0104】

セカンダリ記憶装置システム120Bは、ジャーナル処理管理プログラムを実行し、GUIを介したユーザからの指示の入力あるいは予め定められていたスケジュール（例えば、プライマリ記憶装置システム120Aで一定量のジャーナルがジャーナルボリュームに格納されたとき、または、一定期間ごと等）に従い、セカンダリ記憶装置システム120Bのジャーナル転送プログラム222に対し、ジャーナル転送要求を出す（矢印604）。

20

#### 【0105】

ジャーナル転送要求には、コピーすべきジャーナル（複数でも良い）、そのジャーナルが格納されているジャーナルボリューム、そのジャーナルボリュームを有する記憶装置システム120（ここではプライマリ記憶装置システム120A）を指定する情報、および、コピーしたジャーナルが格納されるジャーナルボリュームを指定する情報が含まれている。これらの情報はジャーナル処理管理プログラム221が制御メモリ145上のジャーナルグループ管理情報から取得した情報から作成される。

#### 【0106】

ジャーナル転送要求を受けたジャーナル転送プログラム222は、リードコマンドをプライマリ記憶装置システム120Aに対して発行する（矢印605）。このリードコマンドを受けたプライマリ記憶装置システム120Aは、リードコマンドで指定されたジャーナルをセカンダリ記憶装置システム120Bに送信する（矢印606）。

30

#### 【0107】

セカンダリ記憶装置システム120Bに送信されたジャーナルが格納されていた、プライマリ記憶装置システム120Aのジャーナルボリュームの領域はページ（開放）され、新たなジャーナルの格納に利用することができる。なお、ページは送信後すぐにおこなう必要はない。ページを定期的におこなっても良いし、ユーザからの指示に従ってページしても良い。

#### 【0108】

ジャーナルを受信したセカンダリ記憶装置システム120Bは、ジャーナル転送要求で指定されたジャーナルボリューム627Bに、受信したジャーナルを格納する。

40

#### 【0109】

その後、セカンダリ記憶装置システム120Bのジャーナル処理管理プログラム221は、セカンダリ記憶装置システム120Bのジャーナル復元プログラム624に対して、ジャーナル復元要求を発行する（矢印607）。ジャーナル復元要求を受けたジャーナル復元プログラム624は、ジャーナルボリューム627Bのジャーナルに基づき、SVOL626にデータの復元をおこなう（矢印608）。なお、復元に用いられたジャーナルが格納されていた領域はページされ、新たなジャーナルの格納に利用することができる。

#### 【0110】

（III-3）ジャーナルボリューム切り替え処理

50

次に、図 1 2 ないし図 1 4 を用いてジャーナルボリューム切り替え処理について説明する。

図 1 2 は、ジャーナルボリュームの切り替えを説明する概念図である。

図 1 3 は、プライマリジャーナルボリュームの切り替え処理を示すフローチャートである。

図 1 4 は、セカンダリジャーナルボリュームの切り替え処理を示すフローチャートである。

【 0 1 1 1 】

既に述べてきているように、本発明では、プライマリサイト 1 0 0 A からセカンダリサイト 1 0 0 B にジャーナルを転送するために、ジャーナルの格納、ジャーナル復元と転送に用いている論理ボリュームが同一となることを回避してシステムの負荷を軽減するために、ジャーナルボリュームの切り替え処理をおこなう。

10

【 0 1 1 2 】

今、図 1 2 に示されるように、プライマリサイト 1 0 0 A には、論理ボリューム P A、論理ボリューム P B、論理ボリューム P C があり、論理ボリューム P B がジャーナル格納用として使用され、論理ボリューム P A が転送元の論理ボリュームとして転送用に使用されているとする。

【 0 1 1 3 】

ジャーナルの順序は、ボリュームの上の方が古く、論理ボリューム P A のジャーナルの最後が、論理ボリューム P B のジャーナルの最初につながっているものとする。

20

【 0 1 1 4 】

さて、論理ボリューム P A 内のジャーナルが転送され終わると、論理ボリューム P B のジャーナルを転送することになるが、このときに、ジャーナル格納用の論理ボリュームを、論理ボリューム P B から論理ボリューム P C に切り替える。

【 0 1 1 5 】

また、セカンダリサイト 1 1 0 B には、論理ボリューム S A、論理ボリューム S B、論理ボリューム S C があり、論理ボリューム S C がジャーナル復元用として使用され、論理ボリューム S A が転送先の論理ボリュームとして転送用に使用されているとする。

【 0 1 1 6 】

ジャーナルの順序は、ボリュームの上の方が古く、論理ボリューム S C のジャーナルの最後が、論理ボリューム S A のジャーナルの最初につながっているものとする。

30

【 0 1 1 7 】

このときには、論理ボリューム S C 内のジャーナルが復元され終わると、論理ボリューム S A のジャーナルを復元することになるが、このときに、ジャーナル転送先の論理ボリュームを、論理ボリューム S A から論理ボリューム S B に切り替える。

【 0 1 1 8 】

上記では、ジャーナル転送用の論理ボリュームとジャーナル復元用の論理ボリュームにあるジャーナルが全て転送されたり、復元されたりしたタイミングで切り替える例を示したが、ジャーナルボリューム切り替えは、ユーザによって入力された切り替えコマンドに従っておこなわれても良い。また、一定期間毎や、設定時刻におこなわれても良い。また、予めホスト 1 1 0、管理端末 1 3 0、あるいは、保守端末 1 2 2 のインタフェースを介して設定された条件を満たした場合におこなわれても良い。

40

【 0 1 1 9 】

また、プライマリサイト 1 0 0 A において、図 1 1 に示した矢印 6 0 5 のジャーナル転送を要求するコマンドを受け付けたタイミングで切り替えることにしてもよい。セカンダリサイト 1 0 0 B において、ジャーナル転送が開始されたタイミングで切り替えることにしてもよい。

【 0 1 2 0 】

このようなジャーナルボリュームの切り替えは、具体的には、プライマリサイト 1 0 0 A では、複数のプライマリジャーナルボリュームがある場合に、図 7 に示した格納済みジ

50



ジャーナルシーケンス番号 9 1 1 と転送済みジャーナルシーケンス番号 9 2 1 の差を制御することによって、ジャーナル取得処理に用いるジャーナルボリュームとジャーナル転送処理に用いるジャーナルボリュームを異なるボリュームにすることが可能である。

【 0 1 2 1 】

また、セカンダリサイト 1 0 0 B では、複数のセカンダリジャーナルボリュームがある場合に、図 8 に示した復元済みジャーナルシーケンス番号 1 0 2 2 と転送予定ジャーナルシーケンス番号 1 0 1 2 の差を制御することによって、ジャーナル復元に用いるジャーナルボリュームとジャーナル転送処理に用いるジャーナルボリュームを異なるボリュームにすることが可能である。

【 0 1 2 2 】

なお、プライマリサイト 1 0 0 A で、ジャーナル転送処理対象となるジャーナル（以下「未転送ジャーナル」）がジャーナル取得処理に用いられているジャーナルボリューム（格納用ジャーナルボリューム）以外のボリュームにない場合（すなわち、格納用ジャーナルボリュームのボリューム内末尾ジャーナルシーケンス番号 9 3 2 B が、転送済みジャーナルシーケンス番号 9 2 1 に 1 加えた番号である場合）、格納用ジャーナルボリュームの切り替えをおこなわなければ、ジャーナル取得とジャーナル転送が同一ジャーナルボリュームに対しておこなわれることになる。

【 0 1 2 3 】

次に、図 1 3 のフローチャートを追いながら、プライマリジャーナルボリューム切り替え処理について説明する。

【 0 1 2 4 】

ジャーナル転送処理をおこなう際に、ジャーナルボリューム切り替え判定フラグをチェックする（ステップ 1 2 0 1 ）。

【 0 1 2 5 】

切り替え判定が有効でなければジャーナル転送処理をおこなう（ステップ 1 2 0 8 ）。切り替え判定が有効であれば、格納用ジャーナルボリューム以外のジャーナルボリュームに未転送ジャーナルがあるか否かをチェックし（ステップ 1 2 0 2 ）、未転送ジャーナルがあれば、そのジャーナルボリュームのジャーナル転送処理をおこなう（ステップ 1 2 0 8 ）。未転送ジャーナルがなければ、未転送ジャーナル量による判定が有効であり、かつ閾値条件を満たしているか否かをチェックし（ステップ 1 2 0 3 ）、有効かつ条件合致であれば、格納用ジャーナルボリュームを次のジャーナルボリュームに切り替えて（ステップ 1 2 0 7 ）、ジャーナル転送をおこなう（ステップ 1 2 0 8 ）。ジャーナル量判定が無効または条件に合わなければ、未転送ジャーナルの最古更新時刻と現在時刻との時間差による判定が有効であり、かつ、閾値条件を満たすか否かをチェックし（ステップ 1 2 0 4 ）、有効かつ条件合致であれば格納用ジャーナルボリュームを次のジャーナルボリュームに切り替えて（ステップ 1 2 0 7 ）、ジャーナル転送をおこなう（ステップ 1 2 0 8 ）。時間差判定が無効または条件に合わなければ、設定情報に合わせて（ステップ 1 2 0 5 ）、格納用ジャーナルボリュームを用いてジャーナル転送をおこなうか（ステップ 1 2 0 8 ）、次のジャーナル転送要求があるまでジャーナル転送を中断する（ステップ 1 2 0 6 ）。ステップ 1 2 0 7 で格納用ジャーナルボリュームを切り替える際、プライマリ記憶装置システム 1 2 0 A の記憶制御装置 1 4 0 は、制御メモリ 1 4 5 に格納されているジャーナルグループ管理情報の格納用ジャーナルボリューム ID を、切り替え先の新たなジャーナルボリュームを示す情報に書き換える。

【 0 1 2 6 】

次に、図 1 4 のフローチャートを追いながら、セカンダリジャーナルボリューム切り替え処理について説明する。

【 0 1 2 7 】

ジャーナル復元処理をおこなう際に、ジャーナルボリューム切り替え判定フラグをチェックする（ステップ 1 3 0 1 ）。切り替え判定が有効でなければ、ジャーナル復元処理をおこなう（ステップ 1 3 0 8 ）。切り替え判定が有効であれば、格納用ジャーナルボリ

10

20

30

40

50

ーム以外のジャーナルボリュームに未復元ジャーナルがあるか否かをチェックし（ステップ1302）、未復元ジャーナルがあれば、ジャーナル復元処理をおこなう（ステップ1308）。未復元ジャーナルがなければ、未復元ジャーナル量による判定が有効であり、かつ、閾値条件を満たしているか否かをチェックし（ステップ1303）、有効かつ条件合致であれば、転送先ジャーナルボリュームを次のジャーナルボリュームに切り替えて（ステップ1307）、ジャーナル復元をおこなう（ステップ1308）。ステップ1307で転送先ジャーナルボリュームを切り替える際、セカンダリ記憶装置システム120Bの記憶制御装置140は、制御メモリ145に格納されているジャーナルグループ管理情報の格納先ジャーナルボリュームIDを、切り替え先の新たなジャーナルボリュームを示す情報に書き換える。ジャーナル量判定が無効または条件に合わなければ未復元ジャーナルの最古更新時刻と現在時刻との時間差による判定が有効であり、かつ閾値条件を満たすかチェックし（ステップ1304）、有効かつ条件合致であれば転送先ジャーナルボリュームを次のジャーナルボリュームに切り替えて（ステップ1307）、ジャーナル復元をおこなう（ステップ1308）。時間差判定が無効または条件に合わなければ、設定情報に合わせて（ステップ1305）、格納用ジャーナルボリュームを用いてジャーナル復元をおこなう（ステップ1308）、次のジャーナル復元要求があるまでジャーナル復元を中断する（1306）。

【0128】

〔実施形態2〕

次に、本発明に係る第二の実施形態を、図15を用いて説明する。

図15は、本発明の第二の実施形態に係るジャーナル処理の動作を示す図である。

【0129】

第一の実施形態のジャーナル処理では、図11に示されるようにジャーナル転送処理においてセカンダリ記憶装置システム120Bがプライマリ記憶装置システム120Aにジャーナル転送を要求するリードコマンドを発行して（矢印605）、転送をおこなった。本実施形態では、プライマリ記憶装置システム120Aがセカンダリ記憶装置システム120Bからのリードコマンドを待つのではなく、プライマリ記憶装置システム120Aからセカンダリ記憶装置120Bに対してデータを書き込むライトコマンドを発行することにより、ジャーナル転送処理をおこなう。

【0130】

まず、プライマリサイト100Aにおいて、PVOLデータの更新（矢印601）に対するジャーナル取得処理は（矢印602）、第一の実施形態と同様である。

【0131】

プライマリ記憶装置システム120Aのジャーナル処理管理プログラム221Aは、ジャーナル転送プログラム222にジャーナル転送要求を発行する（矢印1404）。ジャーナル転送要求には、記憶装置システム120Bへ送信すべきジャーナルが格納されているジャーナルボリューム、記憶装置システム120Bを指定する情報、および、そのジャーナルを記憶装置システム120Bで格納するべきジャーナルボリュームを指定する情報等が含まれる。これらの情報は制御メモリ145上に格納されたジャーナルグループ管理情報から取得する。

【0132】

ジャーナル転送要求を受け取ったジャーナル転送プログラム222は、ライトコマンドをセカンダリ記憶装置システム120Bに発行することで指定されたジャーナルをセカンダリ記憶装置システムに送信する（矢印1406）。セカンダリ記憶装置システム120Bは、プライマリ記憶装置システムからライトコマンドとして受信したジャーナルを、そのコマンドで指定されたセカンダリジャーナルボリュームの領域に格納する。

【0133】

その後、セカンダリサイト100Bにおけるジャーナル復元処理（ステップ608）は、第一の実施形態と同様である。また、本実施形態におけるジャーナルボリューム切り替え処理も、第一の実施形態におけるジャーナルボリューム切り替え処理と同様である。

## 【 0 1 3 4 】

## 〔 実施形態 3 〕

次に、本発明に係る第三の実施形態を、図 1 6 を用いて説明する。

図 1 6 は、本発明の第三の実施形態に係るジャーナル処理の動作を示す図である。

## 【 0 1 3 5 】

本実施形態のデータ処理システムは、図 1 6 に示されるように、ジャーナル処理管理プログラム 1 5 2 1 が、記憶装置システム 1 2 0 に含まれるのではなく、ホスト 1 1 0 内に含まれる点で第一の実施形態と異なる。そして、ジャーナル処理管理プログラム 1 5 2 1 がプライマリホスト 1 1 0 A とセカンダリサイト 1 1 0 B とを接続する通信線により通信をおこなう。

10

## 【 0 1 3 6 】

まず、プライマリサイト 1 0 0 A において、P V O L データの更新（矢印 6 0 1 ）に対するジャーナル取得処理は（矢印 6 0 2 ）、第一の実施形態と同様である。

## 【 0 1 3 7 】

プライマリホスト 1 1 0 A は、ジャーナル処理管理プログラム 1 5 2 1 を実行して特定のコマンド（以下「ジャーナル作成状況取得コマンド」）を発行することで、プライマリ記憶装置システム 1 2 0 A の制御メモリ 1 4 5 上に格納されているジャーナルボリューム管理情報からジャーナル作成状況に関する情報（例えば、ジャーナルの容量）を取得する（矢印 1 5 0 9 ）。

## 【 0 1 3 8 】

プライマリホスト 1 1 0 A が取得したジャーナル作成状況に関する情報は、セカンダリホスト 1 1 0 B へ通知される（矢印 1 5 0 3 ）。

20

## 【 0 1 3 9 】

セカンダリホスト 1 1 0 B は、ジャーナル処理管理プログラム 1 5 2 1 を実行して、G U I を介したユーザからの指示の入力あるいは予め定められたスケジュール（例えば、プライマリ記憶装置システム 1 2 0 A で一定量以上のジャーナルが格納されたとき、または一定期間毎）に従い、セカンダリ記憶装置システム 1 2 0 B に対して、ジャーナル転送要求を発行する（矢印 1 5 0 4 ）。

## 【 0 1 4 0 】

ジャーナル転送要求には、コピーすべきジャーナル、そのジャーナルが格納されているジャーナルボリューム、ジャーナルボリュームを有する記憶装置システム 1 2 0 を指定する情報、および、コピーしたジャーナルが格納されるジャーナルボリュームを指定する情報が含まれる。

30

## 【 0 1 4 1 】

ジャーナル転送要求を受信したセカンダリ記憶装置システム 1 2 0 B は、ジャーナル転送プログラム 2 2 2 を実行することで、リードコマンドをプライマリ記憶装置システム 1 2 0 A に発行する。リードコマンドを受け取ったプライマリ記憶装置 1 2 0 A は、リードコマンドで指定されたジャーナルをセカンダリ記憶装置システム 1 2 0 B に送信する（矢印 6 0 6 ）。セカンダリ記憶装置システム 1 2 0 B に送信されたジャーナルが格納されていた領域は、ページ（開放）され、新たなジャーナルの格納に利用することができる。

40

## 【 0 1 4 2 】

ジャーナルを受信したセカンダリ記憶装置システム 1 2 0 B は、ジャーナル転送要求で指定されたジャーナルボリューム 6 2 7 B に、受信したジャーナルを格納する。

## 【 0 1 4 3 】

その後、セカンダリホスト 1 1 0 B は、セカンダリ記憶装置システム 1 2 0 B に対して、ジャーナル復元要求を発行する（矢印 1 5 0 7 ）。

## 【 0 1 4 4 】

ジャーナル復元要求を受信したセカンダリ記憶装置システム 1 2 0 B は、ジャーナル復元プログラム 6 2 4 を実行して、ジャーナルボリューム 6 2 7 B から S V O L 6 2 6 にデータの復元をおこなう（矢印 6 0 8 ）。復元が終わったジャーナルが格納されていた領域

50

はページされ、新たなジャーナルの格納に利用することができる。

【0145】

本実施形態におけるジャーナルボリューム切り替え処理は、第一の実施形態におけるジャーナルボリューム切り替え処理と同様である。

【0146】

〔実施形態4〕

次に、本発明に係る第四の実施形態を、図17を用いて説明する。

図17は、本発明の第四の実施形態に係るジャーナル処理の動作を示す図である。

【0147】

本実施形態のデータ処理システムは、図17に示されるようにジャーナル転送処理においてプライマリ記憶装置システム120Aがセカンダリ記憶像値システム120Bからのリードコマンドを待つのではなく、プライマリ記憶装置システム120Aからセカンダリ記憶装置120Bに対してデータを書き込むライトコマンドを発行する点で第一の実施形態と異なる。また、ジャーナル処理管理プログラム1621が、記憶装置システム120上ではなく、ホスト110上で実行される点でも第一の実施形態と異なる。また、セカンダリ記憶装置システム120Bがジャーナル復元処理を実行するのではなく、セカンダリホスト110Bがセカンダリジャーナルボリューム627Bから復元に使用するジャーナルを読み出し、SVOL626のデータを復元する点でも第一の実施形態と異なる。本実施形態において、ジャーナル復元プログラムはセカンダリホスト110B上で実行されるプログラムである。

【0148】

本実施形態において、ジャーナル転送処理の主体がプライマリ記憶装置システム120Aであること、ジャーナル復元処理をおこなうのがセカンダリホスト110Bであることから、セカンダリ記憶装置システム120Bには、特殊な機能を有しない一般的な記憶装置を用いることができる。

【0149】

プライマリサイト100Aにおいて、PVOLデータの更新(矢印601)に対するジャーナル取得処理(矢印602)は第一の実施形態と同様である。

【0150】

プライマリホスト110Aは、ジャーナル処理管理プログラム1621を実行してジャーナル作成状況取得コマンドを発行することで、プライマリ記憶装置システム120Aの制御メモリ145上に格納されているジャーナルグループ管理情報からジャーナル作成状況に関する情報(例えば、ジャーナルの容量)を取得する(矢印1609)。

【0151】

プライマリホスト110Aが取得したジャーナル作成状況に関する情報は、セカンダリホスト110Bへ通知される(矢印1603)。

【0152】

プライマリホスト110Aは、ジャーナル処理管理プログラム1521を実行して、GUIを介したユーザからの指示の入力あるいは予め定められたスケジュール(例えば、プライマリ記憶装置システム120Aで一定量以上のジャーナルが格納されたとき、または、一定期間毎)に従い、プライマリ記憶装置システム120Aに対して、ジャーナル転送要求を発行する(矢印1604)。

【0153】

ジャーナル転送要求には、セカンダリ記憶装置システム120Bへ送信すべきジャーナルが格納されているジャーナルボリューム、セカンダリ記憶装置システム120Bを指定する情報、および、そのジャーナルを指定する情報等が含まれている。

【0154】

ジャーナル転送要求を受け取ったプライマリ記憶装置システム120Aは、ライトコマンドをセカンダリ記憶装置システム120Bに発行することで、指定されたジャーナルをセカンダリ記憶装置システム120Bに送信する(矢印1606)。

## 【 0 1 5 5 】

セカンダリ記憶装置システム 1 2 0 B は、プライマリ記憶装置システム 1 2 0 A からのライトコマンドとして受信したジャーナルを、ライトコマンドで指定されたセカンダリジャーナルボリュームの領域に格納する。

## 【 0 1 5 6 】

セカンダリホスト 1 1 0 B は、ジャーナル復元プログラム 1 6 2 4 を実行することで、セカンダリジャーナルボリューム 6 2 7 B からジャーナルを読み出して S V O L 6 2 6 にデータの復元をおこなう（矢印 1 6 0 8 ）。

## 【 0 1 5 7 】

セカンダリジャーナルグループの管理は、セカンダリホスト 1 1 0 B がおこない、ジャーナル転送要求の作成必要な情報（格納用ジャーナルボリュームの情報など）をプライマリホスト 1 1 0 A に通知する。復元が終わったジャーナルが格納されていた領域はパーズされ、新たなジャーナルの格納に利用することができる

本実施形態におけるジャーナルボリューム切り替え処理は、第一の実施形態におけるジャーナルボリューム切り替え処理と同様である。

## 【 0 1 5 8 】

〔上記の実施形態による本発明のデータ処理システムの特徴〕

上述した本発明のデータ処理システムでは、記憶装置システムがジャーナル取得・復元および転送処理をおこない、ジャーナル管理、コピー状態管理は、ホストあるいは記憶装置システムがおこなう構成とした。これによりプライマリサイトとセカンダリサイト間でのデータ複製の実際のデータ転送は、記憶装置システム間のファイバケーブル等で実施される。このことによりホスト間の一般回線のトラフィックを最小限に抑え、またデータの転送は高速な回線となりコピー処理性能の向上が可能である。

## 【 0 1 5 9 】

さらに、記憶装置システムがジャーナルをライトコマンドにより他の記憶装置システムに書き出す機能を有し、そのジャーナルをホストが読み出し、ジャーナル復元することで、セカンダリサイトの記憶装置システムに特別な機能を持たせなくてもデータ複製を実現することが可能となる。

## 【 0 1 6 0 】

さらに、ジャーナルグループは複数のジャーナルボリュームを有し、格納用ジャーナルボリュームを切り替え、プライマリサイトにおいてジャーナル取得処理とジャーナル転送処理を、セカンダリサイトにおいてジャーナル転送処理とジャーナル復元処理を異なるジャーナルボリュームに対しておこなう。これによりジャーナルボリュームへのアクセス負荷を分散することができる。このことによりボリュームへの負荷集中によるボリュームリード処理やライト処理の遅延が軽減され、システムトータル性能の向上が可能である。

## 【図面の簡単な説明】

## 【 0 1 6 1 】

【図 1】本発明の第一の実施形態に係るデータ処理システムのハードウェア構成を示す図である。

【図 2】本発明の第一の実施形態に係るデータ処理システムの機能構成を示す図である。

【図 3】ジャーナルグループ管理情報を示す図である。

【図 4】データボリューム管理情報を示す図である。

【図 5】ジャーナルメタ情報を示す図である。

【図 6】データボリュームとジャーナルボリュームの対応を示す図である。

【図 7】プライマリジャーナルボリューム内のジャーナルデータ領域 8 2 0 の内部構造を示す図である。

【図 8】セカンダリジャーナルボリューム内のジャーナルデータ領域 8 2 0 の内部構造を示す図である。

【図 9】ジャーナルボリュームの切り替え設定情報を示す図である。

【図 10】第一の実施形態に係るデータ処理システムの概略処理を示すフローチャートで

10

20

30

40

50

ある。

【図 1 1】本発明の第一の実施形態に係るジャーナル処理の動作を示す図である。

【図 1 2】ジャーナルボリュームの切り替えを説明する概念図である。

【図 1 3】プライマリジャーナルボリュームの切り替え処理を示すフローチャートである。

。

【図 1 4】セカンダリジャーナルボリュームの切り替え処理を示すフローチャートである。

。

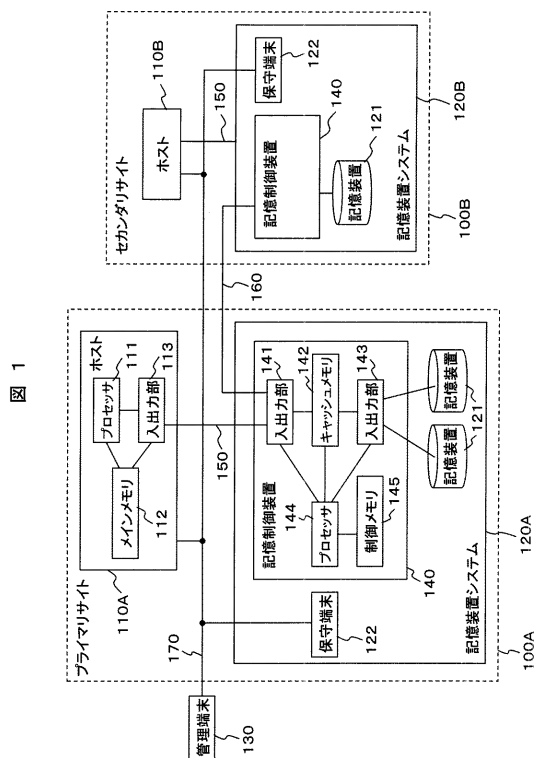
【図 1 5】本発明の第二の実施形態に係るジャーナル処理の動作を示す図である。

【図 1 6】本発明の第三の実施形態に係るジャーナル処理の動作を示す図である。

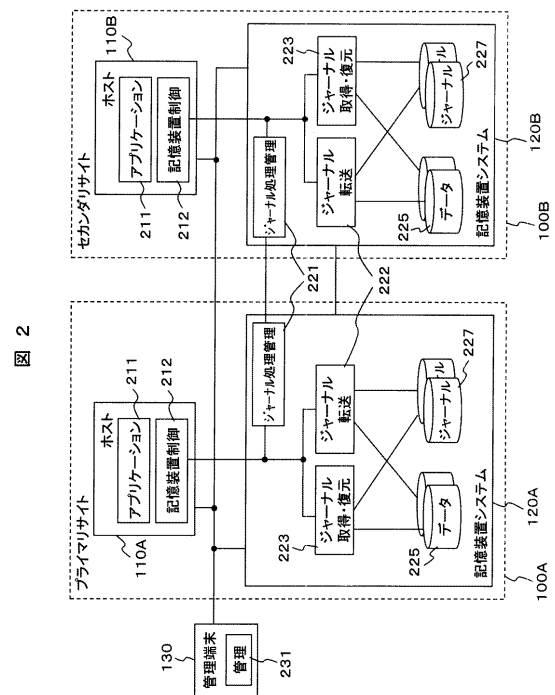
【図 1 7】本発明の第四の実施形態に係るジャーナル処理の動作を示す図である。

10

【図 1】

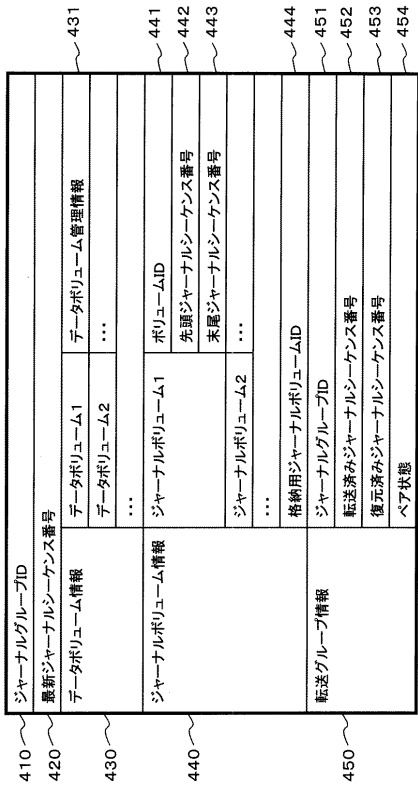


【図 2】



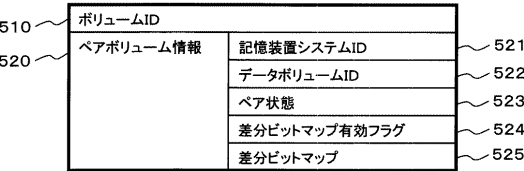
【図 3】

図 3



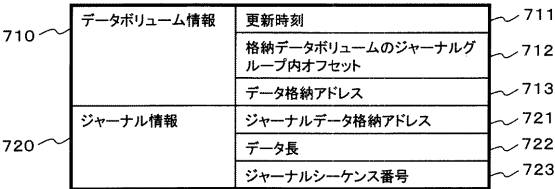
【図 4】

図 4



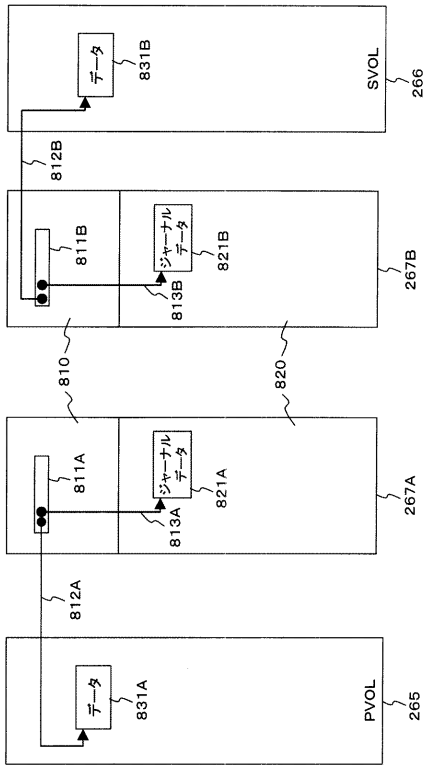
【図 5】

図 5



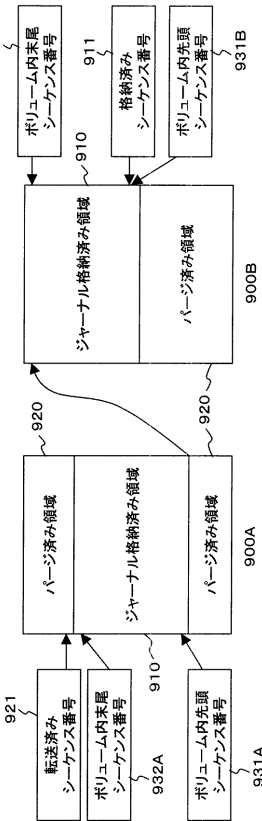
【図 6】

図 6

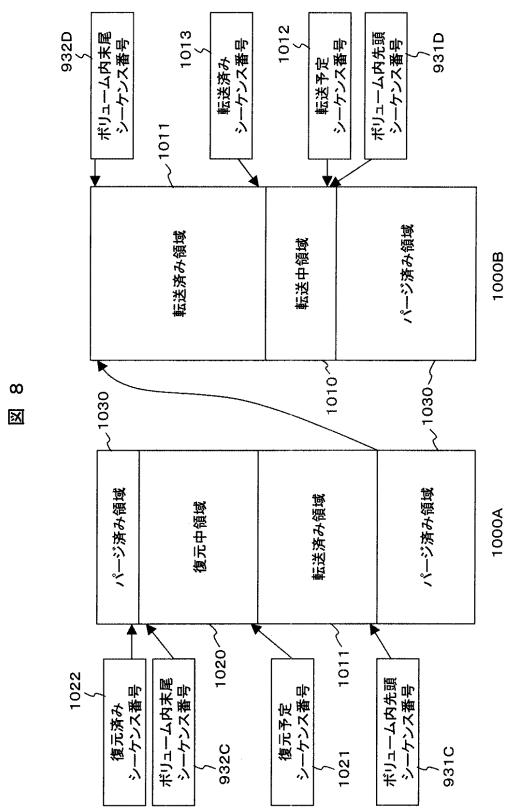


【図 7】

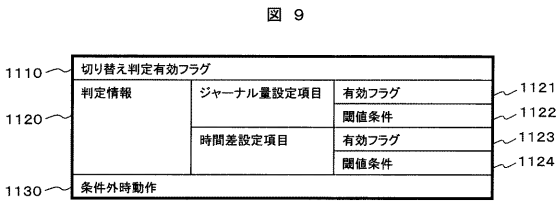
図 7



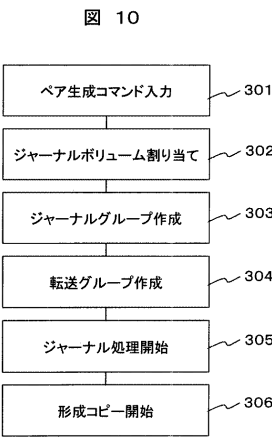
【図 8】



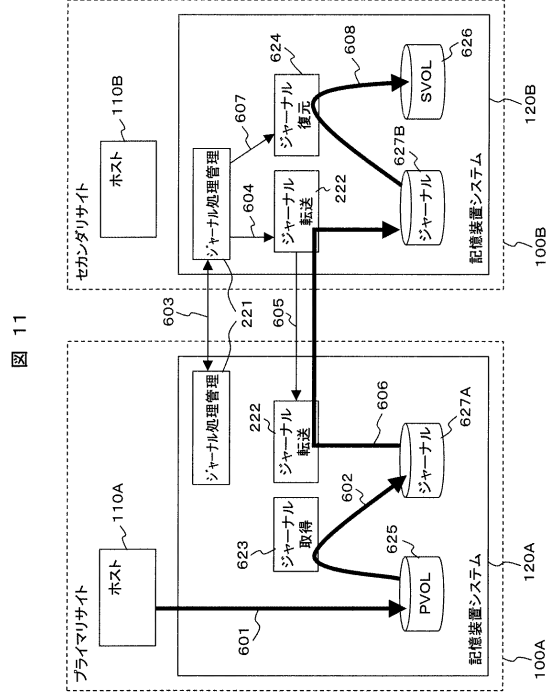
【図 9】



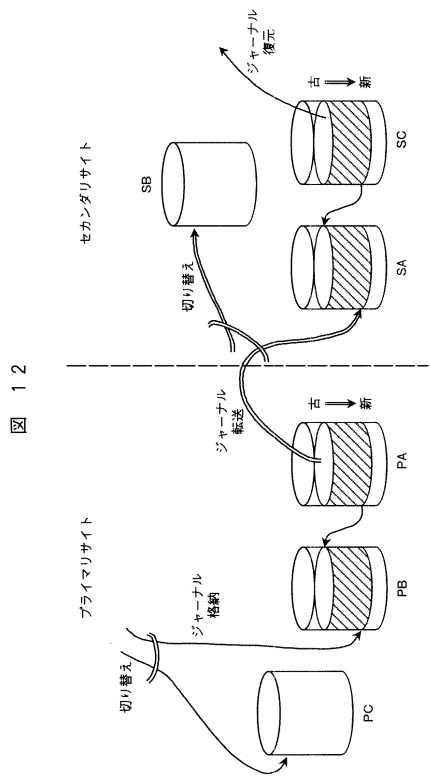
【図 10】



【図 11】

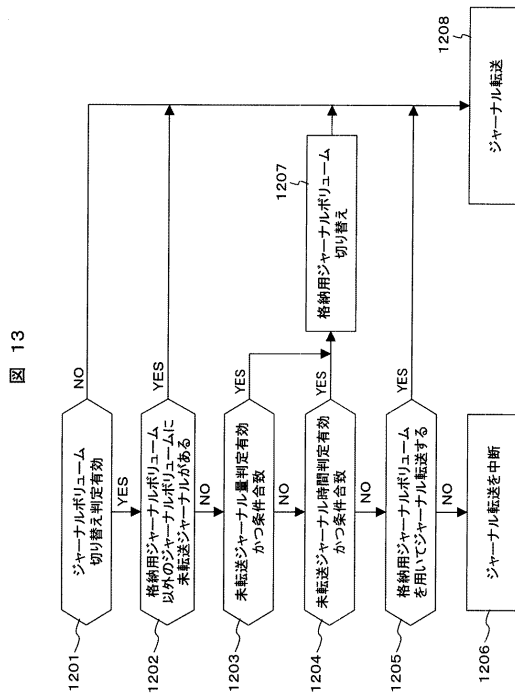


【図 12】

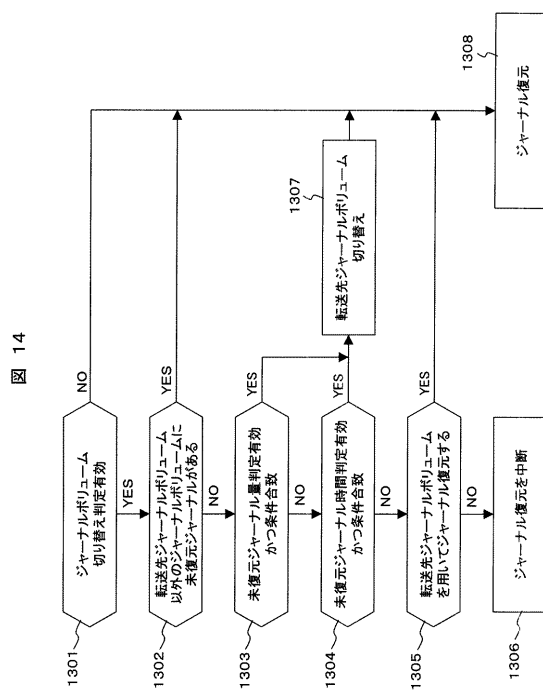




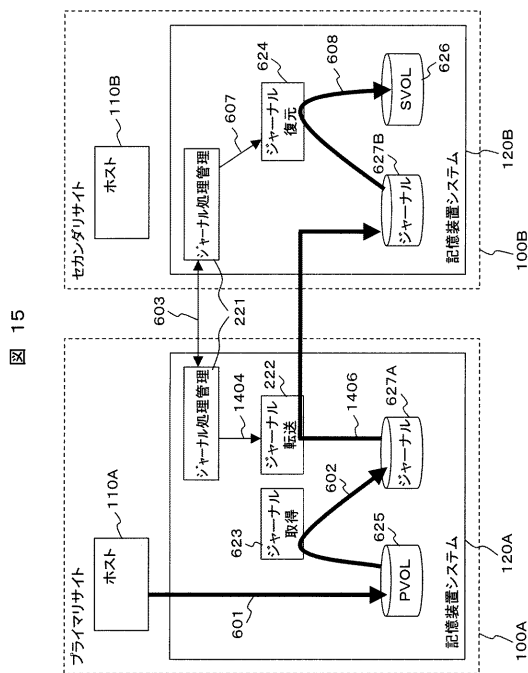
【 図 1 3 】



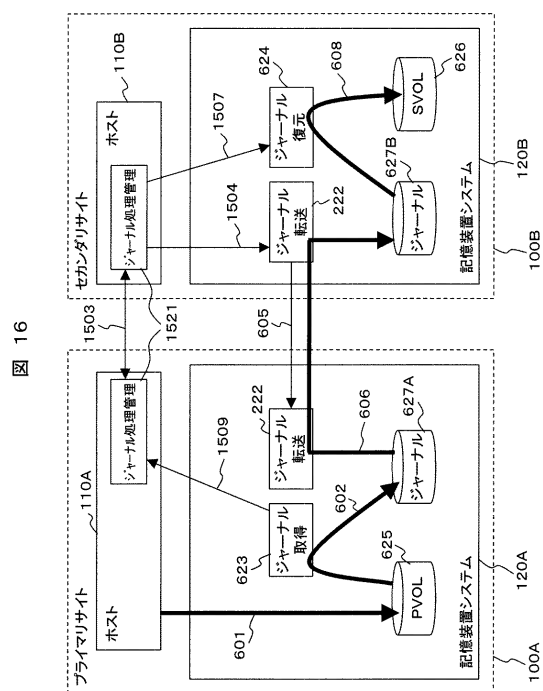
【 図 1 4 】



【 図 1 5 】

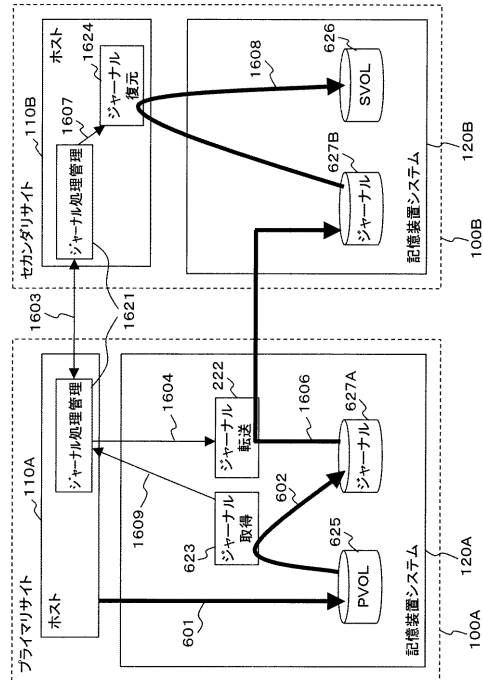


【 図 1 6 】



【図 17】

図 17



---

フロントページの続き

(72)発明者 江口 賢哲

神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所システム開発研究所内

審査官 上嶋 裕樹

(56)参考文献 特開 2 0 0 0 - 1 8 1 6 3 4 ( J P , A )

特開平 0 5 - 0 0 2 5 1 7 ( J P , A )

特開平 0 3 - 2 3 8 5 3 7 ( J P , A )

特開平 1 1 - 0 3 9 2 7 3 ( J P , A )

特開平 1 0 - 0 4 9 4 1 8 ( J P , A )

特開 2 0 0 0 - 2 5 9 5 0 5 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 F 1 2 / 0 0

G 0 6 F 3 / 0 6