

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4294692号
(P4294692)

(45) 発行日 平成21年7月15日(2009.7.15)

(24) 登録日 平成21年4月17日(2009.4.17)

(51) Int.Cl.

F I

G O 6 F 12/00 (2006.01)

G O 6 F 3/06 (2006.01)

G O 6 F 13/10 (2006.01)

G O 6 F 12/00 5 3 1 R

G O 6 F 12/00 5 4 5 A

G O 6 F 12/00 5 1 4 E

G O 6 F 3/06 3 0 4 Z

G O 6 F 3/06 5 4 0

請求項の数 1 (全 31 頁) 最終頁に続く

(21) 出願番号 特願2007-2432 (P2007-2432)
 (22) 出願日 平成19年1月10日(2007.1.10)
 (62) 分割の表示 特願2003-41986 (P2003-41986)
 の分割
 原出願日 平成15年2月20日(2003.2.20)
 (65) 公開番号 特開2007-179552 (P2007-179552A)
 (43) 公開日 平成19年7月12日(2007.7.12)
 審査請求日 平成19年2月6日(2007.2.6)

(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 110000062
 特許業務法人第一国際特許事務所
 (72) 発明者 江口 賢哲
 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
 (72) 発明者 茂木 和彦
 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】 情報処理システム

(57) 【特許請求の範囲】

【請求項 1】

計算機に接続し、制御部とキャッシュメモリと複数の記憶装置とを有する記憶装置システムにおけるジャーナルデータ生成方法であって、

前記複数の記憶装置から第一の論理記憶装置と第二の論理記憶装置とを構成する構成ステップと、

前記制御部にて、計算機から指示を取得し、ジャーナルモードとなるジャーナルモード開始ステップと、

前記制御部にて、前記計算機から複数のライト要求と前記複数のライト要求に対応する複数のライトデータを受信する、ライト要求受信ステップと、

前記複数のライト要求の対象がジャーナルモードの前記第一の論理記憶装置である場合の方法として、

前記制御部にて、前記複数のライト要求が前記第一の論理記憶装置の同じアドレスを指定した場合は前記複数のライトデータの各々を前記キャッシュメモリの異なる領域へ格納する、第一のライト要求格納ステップと、

前記制御部にて、前記複数のライト要求の各々に対応する複数のジャーナルデータを生成するためのジャーナルデータ生成領域を確保する、ジャーナルデータ生成確保ステップと、

前記制御部にて、前記キャッシュメモリの異なる領域に格納された前記複数のライトデータの各々を前記ジャーナルデータ生成領域に格納し、前記複数のライト要求が更新対象

とする従前データを前記キャッシュメモリ又は前記第一の論理記憶装置に対応する前記複数の記憶装置の一部から前記ジャーナルデータ生成領域へ格納することで、各々がライトデータ及び従前データを含む前記複数のジャーナルデータを前記キャッシュメモリに生成する、ジャーナルデータ生成ステップと、

前記ジャーナルデータ生成ステップとは非同期に、前記複数のジャーナルデータを前記キャッシュメモリから前記第二の論理記憶装置に対応する前記複数の記憶装置の一部へ書き込む、ジャーナルデータ格納ステップと、

を有し、

前記複数のライト要求の対象がジャーナルモードの前記第一の論理記憶装置でない場合、前記制御部にて、前記ライトデータの各々を前記キャッシュメモリに格納し、前記キャッシュメモリから前記ライトデータの各々を前記複数の記憶装置に格納する、第二のライト要求格納ステップと、

を有することを特徴とするジャーナルデータ生成方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、計算機や記憶装置システムを含む情報処理システムに関し、特に、障害などによって破壊された記憶装置システムに格納されたデータを復旧する情報処理システム及びそのデータ復旧方法に関する。

【背景技術】

【0002】

情報処理システムで行われるオンライン処理やバッチ処理では、プログラムのバグや記憶装置システムの障害などによってこれらの処理が異常終了し、情報処理システムが有する記憶装置システムに格納されたデータが矛盾した状態になってしまいうことがある。また、人為的ミスによって記憶装置システムに格納されたデータが消去されてしまうことも多い。

【0003】

このような状態になった情報処理システムのデータを回復させる目的で、データの矛盾を解消して途中で止まった処理を再開させたり、あるいは、途中で止まった処理をもう一度やり直したりするための技術の一つとして、データのバックアップとリストアによるデータ回復技術がある。

【0004】

バックアップおよびリストアに関する従来技術の一つが、特許文献1に開示されている。本文献には、ユーザが指定した時点における記憶装置システムに格納されたデータを、記憶装置システムに接続された計算機（以下「ホスト」）からのデータの入出力（以下「I/O」）を止めることなく磁気テープに複製し（以下データの複製を「データのバックアップ」と称する）、その複製されたデータ（以下、「バックアップデータ」）を用いてデータの回復（以下「リストア」）する技術が開示されている。

【0005】

一方、特許文献2には、データのリストアにかかる時間を短縮するために、データのバックアップが実行された後、データが更新された個所についての情報を差分情報として保持し、記憶装置システムに格納されたデータをバックアップデータでリストアする際に、バックアップデータのうち、差分情報で示されるデータの部分のみをデータのリストアに用いる技術が記載されている。

【特許文献1】米国特許番号5,263,154号公報

【特許文献2】特開2001-216185号公報

【発明の開示】

【発明が解決しようとする課題】

【0006】

特許文献1に記載されたリストア処理では、磁気テープからバックアップデータを読み

10

20

30

40

50

出す際、バックアップデータを取得した時点から更新されていない部分（記憶装置システムのデータと磁気テープのデータの内容が一致している部分）も磁気テープから読み出され、記憶装置システムに書き込まれる。このようなデータの転送は、無駄が多く、リストアに要する時間を長びかせる。

【 0 0 0 7 】

一方、特許文献 2 に開示されている技術では、特許文献 1 に比べ、重複したデータの読み出しが発生しない分、リストアに係る時間は少なくなる。しかし、双方の技術をもってしても、データのバックアップの後から記憶装置システムが故障するまでの間に更新されたデータについては、データのリストアを行うことができない。データのバックアップ後に更新されたデータまでリストアしようとする、そのデータの更新の内容等をホスト側がログ等で管理する必要がある、ホストへの負荷が大きく、かつ処理に長い時間がかかる。

10

【 0 0 0 8 】

本発明の目的は、障害発生前までの任意の時点におけるデータのリストア処理を高速に行う記憶装置システム並びに情報処理システムを提供することである。

【課題を解決するための手段】

【 0 0 0 9 】

上記目的を達成するために、本発明は以下の構成を有する。すなわち、計算機及び計算機に接続された記憶装置システムを有する情報処理システムであり、記憶装置システムは制御部及び複数の記憶装置を有する。そして、記憶装置システムは、所定の指示にしたがって、一つの記憶装置に格納されたデータを他の記憶装置に複製する。その後、記憶装置システムは、複製元となった記憶装置へのデータ更新を更新履歴として他の記憶装置に格納する。一方、計算機は、複製が作成された後の任意の時間において、ある識別情報を作成し、記憶装置システムへ送信する。識別情報を受信した記憶装置システムは、その識別情報を更新履歴と関連させて記憶装置へ格納する。

20

【 0 0 1 0 】

データを復元させたい場合、計算機は、記憶装置システムへ識別情報を送信する。識別情報を受信した記憶装置システムは、記録した識別情報から受信した識別情報と一致する識別情報を検索する。一致する識別情報を発見したら、記憶装置システムは、複製先の記憶装置に格納されたデータと、一致した識別情報と関連付けられる更新履歴より前に記録された更新履歴の内容を用いて、複製元の記憶装置にデータを復元する。

30

【 0 0 1 1 】

尚、本発明では、記憶装置システムへデータの更新を要求する計算機は、識別情報を作成する計算機と異なる構成も考えられる。

【 0 0 1 2 】

また、本発明では、識別情報を作成する計算機は、その識別情報を自計算機に格納する構成も考えられる。

【 0 0 1 3 】

更に、本発明では、計算機に格納された識別情報に関する情報をユーザに提示し、ユーザの指定した識別情報を記憶装置システムへ送信する構成も考えられる。

40

【 0 0 1 4 】

更に、本発明の構成として、以下が考えられる。すなわち、中央処理装置を備えた計算機と、記憶装置を備えた記憶装置システムとを有する構成とする。計算機は、記憶装置システムに対して記憶装置に格納されているデータの複製の作成保存を要求する手段、計算機の処理によるデータの更新部分の記録を要求する手段、及びシステムのある時点の状態を識別する識別情報を記憶装置システムに送信する手段とを保持する。記憶装置システムは、計算機の要求に応答して、記憶装置のデータの複製を作成保存する手段、記憶装置の内容が更新されたときに更新前後のデータ及び更新場所をジャーナルデータとして保存する手段、計算機より送信される識別情報を保持識別する手段、並びにジャーナルデータと識別情報を関連付ける手段を有する。更に、計算機は、記憶装置の内容をある時点の状態

50

に復旧する必要がある場合、状態識別情報を指定してデータの復旧要求を記憶装置システムに送信する手段を有し、記憶装置システムは送信された状態識別情報を識別し、前記データの複製とジャーナルデータを用いてデータをリストアする手段を有する。

【 0 0 1 5 】

更に本発明は、以下の構成を有する。即ち、計算機及び記憶装置システムで一つの識別情報を共有し、記憶装置システムではその識別情報と更新履歴を関連付けて管理し、計算機の指示に応じて、特定の識別情報で示される更新履歴まで、記憶装置に格納されたデータを復元するデータの復元方法である。

【発明の効果】

【 0 0 1 6 】

10

本発明によれば、記憶装置システムに格納されたデータを復旧する場合に、ホストに負担をかけず、短時間でデータを所定の状態までリストアすることができる。また、ユーザは、任意のシステム状態までデータをリストアすることができる。

【発明を実施するための最良の形態】

【 0 0 1 7 】

以下、図面を用いて、本発明の第一の実施形態について説明する。尚、これにより本発明が限定されるものではない。以下、「記憶装置システム」には、ディスク装置等の記憶装置、ディスクアレイ等のように複数の記憶装置を有するシステムが含まれるものとする。

【 0 0 1 8 】

20

図 1 は、本発明を適用した情報処理システムの第一の実施形態を示す図である。情報処理システムは、ホスト 1、記憶装置システム 2、管理端末 3、ホスト 1 と記憶装置システム 2 とを接続するネットワーク 4、並びにホスト 1、記憶装置システム 2 及び管理端末 3 とを接続するネットワーク 5 を有する。

【 0 0 1 9 】

ホスト 1 は、パーソナルコンピュータ、ワークステーション、メインフレーム等の計算機である。ホスト 1 では、その計算機の種類に応じたオペレーティングシステム（以下「OS」）と様々な業務、用途に対応したアプリケーションプログラム（AP）、たとえばデータベース（DB）プログラム等、が動作する。本実施形態では、簡単のため、ホスト 1 を 2 つ記載しているが、ネットワーク 4 及び 5 に接続されるホスト 1 は幾つあってもよい。

30

【 0 0 2 0 】

管理端末 3 は、記憶装置システム 2 の障害、保守、構成、性能情報等の管理を行うために使用される計算機である。例えば、情報処理システムの管理者が、記憶装置システム 2 に論理的な記憶装置を設定する場合、データをバックアップするための記憶領域を設定する場合、又はデータを複製する際の記憶領域の対を設定する場合に、管理端末 3 が使用される。情報処理システムの管理者は、記憶装置システム 2 の保守・管理、記憶装置システム 2 が有する物理記憶装置 10 の設定、及び記憶装置システム 2 と接続されるホスト 1 の設定等を行う場合に、管理端末 3 に設定したい内容を入力する。管理端末 3 は、ネットワーク 5 を介して記憶装置システム 2 及びホスト 1 に管理者が入力した内容を送信する。

40

【 0 0 2 1 】

ネットワーク 4 は、ホスト 1 が記憶装置システム 2 へ I / O の処理要求等を伝送するために使用される。ネットワーク 4 には、光ケーブルや銅線等が用いられる。又、ネットワーク 4 で使用される通信プロトコルには、イーサネット（登録商標）、FDDI、ファイバチャネル、SCSI、Infiniband、TCP / IP、iSCSI などがある。

【 0 0 2 2 】

ネットワーク 5 は、記憶装置システム 2 が、自身の障害、保守、構成、性能等の管理情報を管理端末 3 やホスト 1 に送信したり、管理端末 3 やホスト 1 が、記憶装置システム 2 から管理情報を取得する際に使用される。ネットワーク 5 に使用されるケーブル及び通信プロトコルはネットワーク 4 と同一でも異なってもよい。

50

【0023】

図2は、本実施形態における記憶装置システム2の構成を示す図である。記憶装置システム2は、ホスト1が使用するデータやプログラムを格納し、ホスト1のI/O処理要求を受信し、I/O処理要求に対応した処理を行い、その結果を所定のホスト1に送信する。

【0024】

記憶装置システム2は、記憶装置制御装置11、物理記憶装置10、キャッシュメモリ14、共有メモリ19及びLocal Network18とを有する。

【0025】

物理記憶装置10には、ユーザが使用するデータが格納される。物理記憶装置10は、電氣的に不揮発な記憶媒体である磁気ディスクや不揮発性半導体メモリで構成される、シリコンディスク、光ディスク、光磁気ディスク又はハードディスク等である。尚、物理記憶装置10は、物理記憶装置10が有する記憶領域に障害がおきてもデータが損失しないように、冗長性を持つRAID(Redundancy Array Independent Disk)構成になっていてもよい。

【0026】

記憶装置制御装置11は、ホスト1からのI/O要求の処理及び物理記憶装置10の制御を行う装置である。記憶装置制御装置11は、物理記憶装置10と接続される物理記憶装置アダプタ13、所定のプログラムを実行するプロセッサ12、プロセッサ12で実行されるプログラム、プログラムが動作する上で必要な情報、記憶装置システム2の設定情報及び構成情報等が格納される不揮発性メモリ15、記憶装置システム2とネットワーク5とを接続するためのネットワークアダプタ17、記憶装置システム2とネットワーク4とを接続するためのI/Oネットワークアダプタ16とを有する。

【0027】

尚、記憶装置制御装置11は記憶装置システム2に複数存在しても良い。また記憶装置システム2の冗長性を確保するために、システム内の各装置、例えば、記憶装置制御装置11内の各構成要素への電源供給のための回路、キャッシュメモリ14、不揮発性メモリ15、Local Network18、物理記憶装置アダプタ13等は、夫々2重化された冗長構成になっていても良い。

【0028】

キャッシュメモリ14は、記憶装置システム2にホスト1から入力されるデータ又は記憶装置システム2からホスト1へ転送されるデータが一時的に格納される記憶媒体である。

【0029】

共有メモリ19は、複数の記憶装置制御装置11、複数のプロセッサ12間で共有される情報を格納するための不揮発性メモリである。例えばI/O処理のためにキャッシュメモリ14のある領域へアクセスを行うための排他処理用ビットや物理記憶装置10とキャッシュメモリ14との対応関係を示す情報等が格納される。Local Network18は、記憶装置制御装置11、キャッシュメモリ14、及び物理記憶装置10を相互に接続する。Local Network18は、共有バス型の構成でもよいし、スター型等のネットワーク構成となっても良い。

【0030】

図3は、ホスト1の構成を示す図である。ホスト1は、所定のプログラムを実行するプロセッサ20、プロセッサ20が実行するOSやAP及びAPが使用するデータを格納するために使用されるメモリ21、OSやAP、APが使用するデータが格納されるローカルディスク装置22、ネットワーク4とホスト1とを接続するホストバスアダプタ23、ネットワーク5とホスト1とを接続するためのネットワークアダプタ24、フロッピー(登録商標)ディスク等の可搬記憶メディアからのデータの読み出し等を制御するリムーバブル記憶ドライブ装置26、及びこれらの構成部品間を接続し、OSやAPのデータや制御データの転送に用いられるLocal I/O Network25とを有する。

【 0 0 3 1 】

リムーバブル記憶ドライブ装置 2 6 で使用される可搬記憶媒体としては、C D - R O M、C D - R、C D - R W、D V D や M O 等の光ディスク、光磁気ディスクや、ハードディスクやフロッピー（登録商標）ディスク等の磁気ディスク等がある。尚、以下に説明される各プログラムは、可搬記憶媒体からリムーバブル記憶ドライブ装置 2 6 を介して読み出されることで、あるいはネットワーク 4 又は 5 を経由することで、ホスト 1 のローカルディスク装置 2 2 にインストールされる。

【 0 0 3 2 】

ホスト 1 は、冗長性確保のために、プロセッサ 2 0 等の構成部品を複数有していても良い。

10

【 0 0 3 3 】

図 4 は、記憶装置システム 2 が有するプログラムの構成及びシステムの論理的構成を示す図である。記憶装置システム 2 では、単数又は複数の物理記憶装置 1 0（図で点線で表示）が組み合わせられ、冗長性を有するパリティグループ 4 0 7 が構成される。パリティグループ 4 0 7 は、データを格納する物理記憶装置 1 0 及び格納されたデータから作成される冗長データが格納される物理記憶装置 1 0 の組である。また、記憶装置システム 2 は、ホスト 1 に対して、パリティグループ 4 0 7 を構成する複数の物理記憶装置 1 0 が作る記憶領域空間から、論理的な記憶領域を論理記憶装置 4 0 8 として提供する。したがって、ホスト 1 は、記憶装置システム 2 には、図 4 に示すような、記憶装置制御装置 1 1 に接続された記憶装置（論理記憶装置 4 0 8）が存在すると認識する。

20

【 0 0 3 4 】

記憶装置制御装置 1 1 は、記憶装置システム 2 内の処理を制御するために、I / O 処理プログラム 4 0 3、レプリケーション制御処理プログラム 4 0 4、ストレージサブシステム構成管理プログラム 4 0 2、リストア制御処理プログラム 4 0 6 及びジャーナル制御部 4 0 5 の各プログラムを不揮発性メモリ 1 5 に有する。記憶装置制御装置 1 1 は、これらのプログラムをプロセッサ 1 2 で実行することで、以下に説明する処理を制御する。

【 0 0 3 5 】

I / O 処理プログラム 4 0 3 は、更に、コマンド処理プログラム 4 1 5 及びリードライト処理プログラム 4 1 6 からなる。記憶装置制御装置 1 1 は、ホスト 1 からの I / O 処理要求をネットワークインターフェース 1 7 で受信すると、コマンド処理プログラム 4 1 5 を実行して、受信した I / O 処理要求の内容を解析する。解析の結果、I / O 処理要求の内容がデータの読み出し I / O（以下「リード I / O」）要求やデータの書き込み I / O（以下「ライト I / O」）処理要求であれば、記憶装置制御装置 1 1 は、リードライト処理プログラム 4 1 6 を実行する。

30

【 0 0 3 6 】

ライト I / O 処理要求の場合、記憶装置制御装置 1 1 は、ホスト 1 からのライト I / O 処理要求に対する応答処理（実際にホスト 1 から転送されるデータを受領できる状態にあるかどうかの応答）を行い、更に転送されてくる更新用のデータ（以下「ライトデータ」）をキャッシュメモリ 1 4 又は物理記憶装置 1 0 の所定の箇所への書き込み、またはキャッシュメモリ 1 4 に格納されたライトデータを物理記憶装置 1 0 に書き込む制御等を行う。リード I / O 処理要求の場合、記憶装置制御装置 1 1 は、リード I / O 処理要求に対応するデータ（以下、「リードデータ」）を、キャッシュメモリ 1 4 もしくは物理記憶装置 1 0 の所定の箇所から読み出してホスト 1 に転送したり、物理記憶装置 1 0 からリードデータを読み出してキャッシュメモリ 1 4 に格納する処理を制御する。

40

【 0 0 3 7 】

その他の処理の場合、たとえば S C S I の I n q u i r y コマンド（デバイスサーチを指示するコマンド）等の場合、記憶装置制御装置 1 1 は、コマンド処理プログラム 4 1 5 を実行することによって、処理内容に対応した動作の制御を行う。

【 0 0 3 8 】

ストレージサブシステム構成管理プログラム 4 0 2 は、デバイス管理情報 4 1 0 及びデ

50

バイス管理プログラム 409 から構成される。デバイス管理情報 410 は、論理記憶装置 408 のアドレスと物理記憶装置 10 のアドレスとの対応関係を示すマッピング情報、パリティグループ 407 を構成する物理記憶装置 10 に関する情報、スナップショットペア 450 に関する情報、及びジャーナルデータ格納対象情報等とを保持するテーブルである。

【0039】

デバイス管理プログラム 409 は、記憶装置制御装置 11 がデバイス管理情報 410 を管理する際に実行されるプログラムである。記憶装置制御装置 11 は、デバイス管理プログラム 409 を実行することによって、管理端末 3 等から入力される論理記憶装置 408 の定義やスナップショットが格納される対象となる論理記憶装置 408 の設定、ジャーナルデータ格納対象情報の登録等を行う。

10

【0040】

記憶装置制御装置 11 がデータのリードライト I/O 処理を実行する際は、デバイス管理プログラム 409 を実行することによって、リードライト I/O 処理要求が指定するリード又はライトデータが読み出され又は格納されるべき個所の論理記憶装置 408 のアドレスがどの物理記憶装置 10 のアドレスに対応するかを計算し、その結果に基づいて、物理記憶装置 10 へのアクセスを行う。

【0041】

ジャーナル制御プログラム 405 は、記憶装置制御装置 11 がジャーナルデータを作成する際に実行するジャーナル作成プログラム 419、記憶装置制御装置 11 が作成したジャーナルデータを読み出す際に実行するジャーナル読み出しプログラム 420、ジャーナル取得の対象となる論理記憶装置 408 についての情報が登録されたジャーナル管理情報 418、及び記憶装置制御装置 11 がジャーナル管理情報 418 の設定等を行う際に実行するジャーナル管理プログラム 417 から構成される。

20

【0042】

記憶装置制御装置 11 は、ジャーナルデータ取得を行うとき（以下、「ジャーナルモード時」）にホスト 1 からライト I/O 処理要求を受信した場合、ジャーナル作成プログラム 419 を実行することで、ライトデータをキャッシュメモリ 14 に書き込むとともに、ライトデータの格納される個所に存在している従前のデータ（以下「ライト対象データ」）及びライトデータを、キャッシュメモリ 14 に確保されたジャーナルデータ作成用の所定の領域に書き込む。

30

【0043】

尚、キャッシュメモリ 14 に格納されたライト対象データ及びライトデータは、更新履歴であるジャーナルデータとして、ジャーナルデータを格納するための論理記憶装置 408（以下「ジャーナル論理記憶装置」）に格納される。又、記憶装置制御装置 11 は、リストアマネージャ 406 及びジャーナル読み込みプログラム 420 を実行することで、ホスト 1 からの指示に基づき、ジャーナル論理記憶装置に格納されたジャーナルデータを順次読み出し、読み出したジャーナルデータが有するアドレスで示される、複製先となる論理記憶装置 408 又は複製元である論理記憶装置 408 の記憶領域にデータを上書きする。

40

【0044】

スナップショット制御プログラム 404 は、コピー処理プログラム 413、差分情報 414、ペア制御管理プログラム 411 及びペア管理情報 412 から構成される。記憶装置制御装置 11 は、ペア制御管理プログラム 411 を実行することで、ホスト 1 からの指示に従って、ある論理記憶装置 408（以下、「正論理記憶装置」）及び正論理記憶装置に格納されたデータの複製を格納する論理記憶装置 408（以下、「副論理記憶装置」）について、ペア形成（Pair Create）、ペア分離（Pair Split）、ペア再結合（Pair Resync）、ペア削除（Pair Delete）の処理を行う。ここで、「ペア」とは、正論理記憶装置と、正論理記憶装置に対応する副論理記憶装置の組（以下「スナップショットペア 450」）を指す。

50

【0045】

尚、1つの正論理記憶装置に対して、複数の副論理記憶装置を設定・作成することもできる。また、副論理記憶装置を新たな正論理記憶装置とし、新たな正論理記憶装置とペアになる副論理記憶装置を設定・作成することもできる。

【0046】

ペア管理情報412には、ある論理記憶装置のスナップショットペア450がペア結合状態(Pair Duplex)のペア同期状態(Pair Synchronus)、ペア結合状態(Pair Duplex)のペア非同期状態(Pair Asynchronus)、ペア形成状態(Pair Create)、ペア分離状態(Pair Symplex)にあるかどうかを示す情報が登録される。Pair Synchronus状態とは、ホスト1のライトI/Oによる正論理記憶装置の更新と副論理記憶装置の更新が同期して行われる状態を示す。Pair Asynchronus状態とは、ホスト1のライトI/Oによる正論理記憶装置の更新と副論理記憶装置の更新が非同期に行われる状態を示す。尚、Pair Asynchronus状態の場合は、副論理記憶装置に正論理記憶装置の更新が反映されるまで、ライトデータは、差分情報414で管理される。

10

【0047】

差分情報414には、あるペアがペア非同期状態(Pair Asynchronus)又は分離状態(Pair Symplex)の場合に、正論理記憶装置にデータの書き込みが発生することによって生ずる正論理記憶装置と副論理記憶装置との間の差異が有る部分を示すアドレス情報等が保持される。

20

【0048】

記憶装置制御装置11は、コピー処理プログラム413を実行することによって、ペア作成(Pair Create)時に正論理記憶装置の先頭アドレスから順次副論理記憶装置にデータを複写することで、正論理記憶装置に格納されたデータを副論理記憶装置にバックアップする。さらに記憶装置制御装置11は、差分情報414を参照して、差異が有る部分のデータを正論理記憶装置から副論理記憶装置にコピーしたり、逆に、差分情報414を参照して、差異があるデータを副論理記憶装置から正論理記憶装置へコピーする。

【0049】

バックアップ/リストア制御プログラム406は、リストアプログラム421とバックアッププログラム422から構成される。記憶装置制御装置11は、リストアプログラム421を実行することで、ホスト1からのリストア要求に基づいて、指定された論理記憶装置408のデータをリストアする。尚、リストア処理の詳細は後述する。

30

バックアッププログラム422は、記憶装置制御装置11が、ホスト1の指示等に従って、論理記憶装置408の複製を作成したり、記憶装置システム2のデータを他の記憶装置、例えばテープに転送したりする際に実行される。

【0050】

図5は、ホスト1で動作するプログラム及び使用されるデータの例を示す図である。これらのプログラムは、ホスト1のローカルディスク装置22又はメモリ21に格納され、プロセッサ20で実行される。ホスト1は、OS500の下で動作するAPとして、データベースマネジメントソフトウェア(以下「DBMS」)501を有する。DBMS501は、OS500、ファイルシステム(FS)530、ボリュームマネージャ(VM)540等を介して記憶装置システム2にアクセスする。また、DBMS501は、ユーザが使用する他のAP520との間で、トランザクション処理等のI/O処理の遣り取りを行う。

40

【0051】

DBMS501は、DBファイル505、LOGファイル506、INDEXファイル507、DBバッファ509、LOGバッファ510、デバイス情報ファイル511、状態ファイル508、DB定義ファイル512、トランザクションマネージャ502、ログマネージャ503、バッファマネージャ513、及びリソースマネージャ504から構成

50

されている。

【 0 0 5 2 】

D B バッファ 5 0 9 は、D B M S 5 0 1 の処理性能を上げる目的で、ホスト 1 のメモリ 2 1 に確保される D B M S 5 0 1 専用の領域である。このバッファ 5 0 9 には、D B M S 5 0 1 によって良くアクセスされるデータが一時的に保持される。ログバッファ 5 1 0 も D B バッファ 5 0 9 と同様にメモリ 2 1 上に確保された領域で、D B M S 5 0 1 の処理記録（以下「ログ」）が一時的に格納される。

【 0 0 5 3 】

D B ファイル 5 0 5 は、D B のテーブル等 D B のデータそのものであり、実際には記憶装置システム 2 の物理記憶装置 1 0 内に格納されている。そして、良く使用されるテーブル等のデータが D B バッファ 5 0 9 に一時格納され、D B M S 5 0 1 は、そのデータでトランザクション処理を行う。D B バッファ 5 0 9 に要求されるデータが無い場合、D B M S 5 0 1 は、データを記憶装置システム 2 から読み上げる。

【 0 0 5 4 】

ログファイル 5 0 6 も実際には記憶装置システム 2 の物理記憶装置 1 0 に格納されている。ログファイル 5 0 6 には、トランザクション処理等の D B M S 5 0 1 が D B に対して行った処理のログ（処理を行った A P の識別子、処理順序識別子、処理を行った時間や処理を行ったデータ及び処理対象前データ等を含む）が順次記録される。記録の際には、ログバッファ 5 1 0 を用いて順次追記される。ログファイル 5 0 6 には、A P 5 2 0 が一連の処理を行い整合性が取れた状態でコミットした際及び D B M S 5 0 1 が一定時間間隔やトランザクション数等毎に物理記憶装置 1 0 にバッファに格納されたダーティデータを格納するシンク処理を行った際にも、それを示す情報が記録される。

【 0 0 5 5 】

ホスト 1 は、トランザクションマネージャ 5 0 2 を実行することで、D B に対するトランザクション処理や、ログファイル 5 0 6 に格納されたデータを読み出してデータのリカバリを実行したり、チェックポイントの制御を行ったりする。又、ホスト 1 は、ログマネージャ 5 0 3 を実行することで、D B に対するデータの入出力を制御する。

【 0 0 5 6 】

以下、本実施形態の動作概要について説明する。本実施形態の情報処理システムでは、まず、記憶装置システム 2 において、正論理記憶装置と副論理記憶装置のある時点に有するデータのバックアップデータ（以下「スナップショットデータ」）を有する副論理記憶装置を作成し保持する。スナップショットが作成された時点以降にホスト 1 からのライト I / O 処理要求がある度に、記憶装置システム 2 は、ライト I / O 処理前後のデータ（ライトデータ及びライト対象データ）をジャーナルデータ（「更新履歴」）として記録する。

【 0 0 5 7 】

さらに、ホスト 1 は、自身が作成する任意の識別情報であるチェックポイント情報（以下「C P 情報」）を記憶装置システム 2 に対して通知する。具体的には、ホスト 1 は、任意の時点、例えば記憶装置システム 2 との間でのデータを一致させる処理（シンク処理）時に、C P 情報を記憶装置システム 2 のジャーナルデータに書込む。これにより、記憶装置システム 2 は、ホスト 1 で作成されたものと同じの C P 情報を保持する。つまり、従来ホスト 1 でのみ管理されていた C P 情報をホスト 1 と記憶装置システム 2 の双方で管理する。これによって、ホスト 1 が指示する C P 情報及び記憶装置システム 2 内のジャーナルデータに格納された C P 情報を利用して、記憶装置システム 2 は、ホスト 1 が意図した時（C P 情報作成時）の記憶装置システム 2 が有していたデータの状態に高速にリストアを行う。

【 0 0 5 8 】

このような処理を実行するために、ホスト 1 は、あらかじめ、ジャーナルデータを取得する準備指示（ジャーナル取得開始準備指示）、及びジャーナル取得開始指示を記憶装置システム 2 に送信する。これにより、記憶装置システム 2 は、ジャーナルデータの取得を

開始し、ジャーナルモードとなる。その後、情報処理システムは、上述したC P情報の遣り取りを行う。

【0059】

以下、ホスト1がジャーナル取得開始準備指示を記憶装置システム2に発行した際に、記憶装置システム2で行われる処理について説明する。

【0060】

ジャーナル取得開始準備指示指示には、ジャーナル論理記憶装置を指定する情報や、正論理記憶装置及び副論理記憶装置の作成指示等が含まれる。ジャーナル取得開始準備指示を受領した記憶装置システム2は、指示に従い、データ格納領域の割当等を実行する。正副論理記憶装置は、ジャーナル開始準備指示を受領する前からスナップショットペア450 10
0になっても良いが、本実施形態では、記憶装置システム2が、ジャーナル取得開始準備指示に基づいて新たに論理記憶装置408をスナップショットペア450に設定する。

【0061】

記憶装置システム2は、次に、正論理記憶装置のスナップショットデータを指定された副論理記憶装置に作成する。具体的には、記憶装置システム2がジャーナル取得開始準備指示を受取った時点で正論理記憶装置に格納されているデータを副論理記憶装置に複製し、正論理記憶装置と副論理記憶装置の状態を同期させる。尚、ジャーナル取得開始準備指示以前から正論理記憶装置とスナップショットペア450 20
0になっている副論理記憶装置が指定された場合は、記憶装置システム2は、副論理記憶装置と正論理記憶装置とを同期させた状態にするだけで良い。

【0062】

更に、記憶装置システム2は、ホスト1の指示に基づいて、正論理記憶装置に対応するジャーナル論理記憶装置の設定も行う。

【0063】

次に、ホスト1は、記憶装置システム2に、ジャーナル取得開始指示を出す。ジャーナル取得開始指示には、ジャーナルデータ取得開始を示す最初のC P情報であるチェックポイント識別子(以下「C P I D」)が含まれている。記憶装置システム2は、受信した最初のC P I Dを記録し、その後、ジャーナルデータの取得を開始する。尚、その後にホスト1から送信されるチェックポイントコマンドにも最初のC P I Dとは別のC P I Dが含まれている。C P I Dは、記憶装置システム2でジャーナルデータとして記録される。 30

【0064】

図6は、ホスト1からジャーナル取得開始準備指示及びジャーナル取得開始指示を受領した記憶装置システム2における処理の詳細手順を示す図である。

【0065】

ホスト1は、D B M S 501を実行することで、記憶装置システム2に対して、ジャーナル取得開始準備指示を送信する。尚、本実施形態では、D B M S 501が使用するD Bのテーブルが格納された論理記憶装置408が正論理記憶装置として指定される。ジャーナル取得開始準備指示には、正論理記憶装置を示す識別子、ジャーナル取得開始準備指示を記憶装置システム2が受領した瞬間のある正論理記憶装置に格納されたデータのスナップショットデータを格納するための副論理記憶装置を示す識別子、ジャーナル論理記憶装置を示す識別子が含まれる(ステップ601)。 40

【0066】

ジャーナルデータは、スナップショットデータが作成された後のライトI / O処理要求に基づくライト対象データ、ライトデータ及びこれらのデータの正論理記憶装置内における格納位置を示すアドレス情報等から構成される。構成の具体例は後述する。

【0067】

尚、スナップショットデータが格納される副論理記憶装置やジャーナル論理記憶装置の設定は、ジャーナル取得開始準備指示とは別の指示に基づいて、予め行われていても良い。この場合、ジャーナル取得開始準備指示には、これらの論理記憶装置408を示す識別 50

子は含まれなくても良い。

【0068】

ホスト1からジャーナル取得開始準備指示を受領した記憶装置制御装置11は、指示に含まれていている副論理記憶装置を示す識別子を用いてデバイス管理情報410を参照し、無効なデバイスの指定の有無、例えば、指定された副論理記憶装置の存在の有無や障害発生の有無、論理記憶装置の状態の確認、例えば指定された副論理記憶装置が、既に他の処理に使用されている等、の確認を行う。確認の結果、指定された副論理記憶装置が使用可能である場合、記憶装置制御装置11は、指定された副論理記憶装置がジャーナル作成中であることを示す情報をデバイス管理情報410に設定するとともに、指定された副論理記憶装置に関するジャーナル管理情報をジャーナル管理情報418に設定し、かつPa

10

【0069】

同様に、記憶装置制御装置11は、ジャーナル論理記憶装置を示す識別子を用いてデバイス管理情報410を参照し、指定されたジャーナル論理記憶装置の無効なデバイスの指定の有無及び状態の確認を行う。指定されたジャーナル論理記憶装置が使用できる場合、指定されたジャーナル論理記憶装置がジャーナル作成中とする情報をデバイス管理情報410に登録する(ステップ603)。

【0070】

次に、記憶装置制御装置11は、副論理記憶装置に正論理記憶装置のスナップショットデータを作成する処理(以下「スナップショット作成処理」)を行う。スナップショット作成処理においては、ジャーナル取得開始準備処理指示のコマンド受領時に正論理記憶装置に格納されていたデータが、副論理記憶装置に順次転送される。尚、ジャーナル取得開始準備処理指示に副論理記憶装置の指示が含まれず、予めDuplex状態のPairである副論理記憶装置が管理端末3で指定されていた場合や、副論理記憶装置の指示が含まれていても、指定された副論理記憶装置が既に正論理記憶装置とDuplex状態にある場合は、スナップショット作成処理は行わなくても良い。

20

【0071】

尚、記憶装置システム2がスナップショット作成処理を実行している最中に、ホスト1から正論理記憶装置に格納されたデータに対するライトI/O処理要求があった場合、記憶装置制御装置11は、要求時点でライト対象データが未だ副論理記憶装置にコピーされていなかったら正論理記憶装置にライトデータを書込み、要求時点で既にライト対象データが副論理記憶装置にコピーされていたら、ライトデータを正論理記憶装置に書き込むとともに、副論理記憶装置にも書きこむ(ステップ604)。

30

【0072】

スナップショット作成処理が終了したら、記憶装置制御装置11は、ペア管理情報をDuplex状態にし(ステップ605)、ジャーナル取得準備処理の完了を、ジャーナル取得開始準備指示を発行したホスト1に報告する。尚、Duplex状態にあるスナップショットペア450では、正論理記憶装置に書き込まれたデータは、副論理記憶装置にも反映される(ステップ606)。

【0073】

ジャーナル取得準備処理の完了報告を受領したホスト1は、任意のタイミング、例えば情報処理システムの状態が整合性が取れている時、指定時間又はあるランザクション処理の前や後で、ジャーナル取得開始指示を記憶装置システム2に送信する(ステップ607)。

40

【0074】

ジャーナル取得開始指示を受領した記憶装置制御装置11は、先に準備したジャーナル論理記憶装置、正副論理記憶装置に障害が発生していないかを確認して、ジャーナル取得開始指示に対してReady応答を返す(ステップ608)。

【0075】

その後、記憶装置制御装置11は、正副論理記憶装置をPair Split状態にす

50

る。具体的には、ホスト 1 からライト I / O 処理要求を受取っても、正論理記憶装置の更新が副論理記憶装置には一切反映されない状態にする（ステップ 6 0 9）。

【 0 0 7 6 】

一方、R e a d y 応答を受領したホスト 1 は、チェックポイントコマンドを用いて、C P I D を含む C P 情報を送信する（ステップ 6 1 0）。

【 0 0 7 7 】

C P 情報を受領した記憶装置システム 2 は、ジャーナル論理記憶装置に、受信した C P 情報、具体的には、C P I D、記憶装置システム 2 内の処理シーケンス番号及び処理時間をジャーナルデータとして格納する。もしくは、記憶装置制御装置 1 1 にある不揮発性メモリ 1 5 又は共有メモリ 1 9 に C P 情報を格納する（ステップ 6 1 1）。

10

【 0 0 7 8 】

チェックポイントコマンドを送信したホスト 1 は、ホスト 1 のメモリ 2 1 に格納されているライトデータを記憶装置システム 2 に送信する（ステップ 6 1 2）。

【 0 0 7 9 】

ライトデータを受領した記憶装置制御装置 1 1 は、ライトデータを正論理記憶装置に書き込むと共に、ライト対象データ及びライトデータをジャーナル論理記憶装置に書きこむ（ステップ 6 1 3）。

【 0 0 8 0 】

チェックポイントコマンド受領以降、記憶装置システム 2 はジャーナルデータの取得を継続するジャーナルモードとなる。また、これ以降、一定時間毎や一定トランザクション数毎等、D B 管理者が設定した間隔で、ホスト 1 は、その時点に D B バッファ 5 0 9 上のデータ全てを記憶装置システム 2 に送信する。更に、記憶装置システム 2 とホスト 1 とで C P 情報を共有するタイミングである場合には、ホスト 1 は、C P 情報を共有するタイミングであることを示すチェックポイントコマンドを送信する。

20

【 0 0 8 1 】

ジャーナルモード中にチェックポイントコマンドを受領した記憶装置制御装置 1 1 は、C P 情報をジャーナルデータとして、ジャーナル論理記憶装置、不揮発性メモリ 1 5 又は共有メモリ 1 9 に格納する。

【 0 0 8 2 】

図 7 は、ジャーナルモード中の記憶装置システム 2 が、ホスト 1 よりリードライト I / O 処理要求を受信した場合の処理手順を示す図である。

30

【 0 0 8 3 】

ホスト 1 よりリードまたはライト I / O 処理要求を受領した記憶装置システム 2 の記憶装置制御装置 1 1 は（ステップ 7 0 1）、受信した処理要求がライト I / O 処理要求であるかどうかを判断する（ステップ 7 0 2）。ライト I / O 処理要求でない場合、記憶装置制御装置 1 1 は、デバイス管理情報 4 1 0 を用いて、リード I / O 処理要求の対象となっているリードデータを、対応する物理記憶装置 1 0 又はキャッシュメモリ 1 4 から読み出して I / O インタフェース 1 6 を介してホスト 1 に転送する（ステップ 7 0 9）。

【 0 0 8 4 】

ステップ 7 0 2 でライト I / O 処理要求と判断した場合は、記憶装置制御装置 1 1 は、デバイス管理情報 4 1 0 を参照し、ライト I / O 処理要求で指定される論理記憶装置 4 0 8 が、ジャーナルモードである正論理記憶装置であるかを判断する（ステップ 7 0 3）。ジャーナルモードの正論理記憶装置でなければ、記憶装置制御装置 1 1 は、キャッシュメモリ 1 4 にライト I / O 処理要求に伴うライトデータを格納する領域を確保する（ステップ 7 0 7）。その後、記憶装置制御装置 1 1 は、ライトデータをキャッシュメモリ 1 4 の確保された領域に格納して、ライト I / O 処理が終了したことをホスト 1 に通知する（ステップ 7 0 8）。

40

【 0 0 8 5 】

尚、記憶装置制御装置 1 1 は、キャッシュメモリ 1 4 から物理記憶装置 1 0 にデータを格納した後にライト I / O 処理の終了をホスト 1 に報告してもよく、又ライトデータをキ

50

キャッシュメモリ 14 を介さず直接物理記憶装置 10 に格納してもよい。

【0086】

一方、ステップ 703 でライト I/O 処理の対象となる論理記憶装置 408 がジャーナルモードの正論理記憶装置であった場合、記憶装置制御装置 11 は、ライトデータを格納するための領域をキャッシュメモリ 14 に確保し、ホスト 1 から送信されるライトデータを当該領域に格納する。

【0087】

尚、通常の論理記憶装置 408 へのライトデータの書き込みとは違い、記憶装置制御装置 11 は、同じアドレスが指定される複数のライトデータの連続した書き込みの際は、各々のライトデータをキャッシュメモリ 14 の異なる領域に格納しなければならない。これは、ライト I/O 処理要求の対象となるライト対象データがキャッシュメモリ 14 に存在するが物理記憶装置 10 にそのライトデータが反映されていない場合、通常書き込み処理の様にキャッシュメモリ 14 に存在するライト対象データを更新してしまうと、更新前のライト対象データが失われ、ライト対象データをジャーナル論理記憶装置に格納することができなくなるからである（ステップ 705）。その後、記憶装置制御装置 11 は、ジャーナルデータの作成処理を行い、処理を終了する（ステップ 706）。

【0088】

図 8 は、図 7 のステップ 706 のジャーナルデータ作成処理の手順を示す図である。ライトデータをキャッシュメモリ 14 に格納した記憶装置制御装置 11 は、ジャーナルデータを一時的に格納するための領域をキャッシュメモリ 14 に確保する（ステップ 901）。

【0089】

その後、記憶装置制御装置 11 は、キャッシュメモリ 14 に格納されているライトデータを、CP 情報、処理シーケンス番号、処理時間とともに、キャッシュメモリ 14 に確保されたジャーナルデータ格納用の領域にコピーする（ステップ 902、903）。ただし、CP 情報の CPID 1007 エントリには、ホスト 1 からのチェックポイントコマンド受領時にのみ CPID が格納されるので、それ以外の場合は、CPID 1007 エントリには無効データが格納される。処理シーケンス番号は、プロセッサ 12 が処理を行うごとに付ける処理通番号である。

【0090】

同時に、記憶装置制御装置 11 は、キャッシュメモリ 14 に格納されたライトデータによって更新されるライト対象データを格納するための領域をキャッシュメモリ 14 に確保し、そのライト対象データを物理記憶装置 10 あるいはキャッシュメモリ 14 から読みだして、キャッシュメモリ 14 の確保された記憶領域に格納する（ステップ 904、905）。これにより、ライトデータ、ライト対象データ、CP 情報、処理シーケンス番号及び処理時間を含むジャーナルデータが作成される。

【0091】

全ての処理が終了した後、記憶装置制御装置 11 は、図 7 の処理に戻る。尚、キャッシュメモリ 14 で作成されたジャーナルデータは、キャッシュメモリ 14 にジャーナルデータが作成されるのとは非同期に、キャッシュメモリ 14 から物理記憶装置 10 に書き込まれる（ステップ 906）。

【0092】

図 9 は、ジャーナルデータのデータ形式を示す図である。

【0093】

ジャーナルデータは、図 6 で説明したように、ジャーナル取得開始指示受信後、記憶装置システム 2 が正論理記憶装置に対するライト I/O 処理要求を処理する毎にキャッシュメモリ 14 上に作成され、その後物理記憶装置 10 に格納される。ジャーナルデータは、ホスト 1 と記憶装置システム 2 でシステムの状態を一意に識別する CP 情報を格納するエントリ 1001、データが更新される箇所を示すブロックアドレスが格納されるエントリ 1002、更新に用いられるライトデータの長さが格納されるエントリ 1003、データ

が更新される個所に格納されていたライト対象データが格納されるエントリ 1 0 0 4、及びライトデータが格納されるエントリ 1 0 0 5 とから構成される。C P 情報エントリ 1 0 0 1 には更に、チェックポイントフラグエントリ 1 0 0 6、C P I D が格納されるエントリ 1 0 0 7、処理順序番号エントリ 1 0 0 8、及び時刻エントリ 1 0 0 9 が含まれている。

【 0 0 9 4 】

記憶装置システム 2 がホスト 1 よりチェックポイントコマンドを受領して C P 情報を受信した場合、記憶装置制御装置 1 1 は、受信した際に作成されるジャーナルデータの C P 情報エントリ 1 0 0 1 に含まれるチェックポイントフラグエントリ 1 0 0 6 に「ON」を示す情報を登録し、C P I D エントリ 1 0 0 7 に、送信されてきた C P I D を格納する。C P I D エントリ 1 0 0 7 に格納される C P I D は、ホスト 1 が管理するログファイルに記録されている C P 情報に含まれる特定の C P I D と対応する一意の値を持っている。したがって、ホスト 1 がある C P I D を指定すると、指定された C P I D に対応する、ジャーナルデータに格納された C P I D を指定することができる。

10

【 0 0 9 5 】

図 1 0 は、ホスト 1 が C P 情報を記憶装置システム 2 に送信する処理手順を示す図である。ホスト 1 は、チェックポイントコマンドを発行し記憶装置システム 2 に C P 情報を送信することによって、D B が有するデータの状態を確定しログファイルにチェックポイントを記録した (C P I D 等の情報が記録される) ことを記憶装置システム 2 に通知することが出来る。

20

【 0 0 9 6 】

先ず、ホスト 1 は、D B バッファ 5 0 9 及びログバッファ 5 1 0 等メモリ 2 1 にあるバッファに格納されたデータを、記憶装置システム 2 へ強制的に書き込むためのライト I / O 処理要求を記憶装置システム 2 に送信する。本処理によって、ホスト 1 は、これらのバッファにのみ格納されていて記憶装置システム 2 には格納されていないデータ (以下「ダーティデータ」) を記憶装置システム 2 に反映して、D B のデータを確定することができる (ステップ 1 1 0 1) 。

【 0 0 9 7 】

ライト I / O 処理要求を受信した記憶装置制御装置 1 1 は、ホスト 1 から送信されるデータをキャッシュメモリ 1 4 に書き込む (ステップ 1 1 0 2) 。転送されたデータを全てキャッシュメモリ 1 4 に書き込んだら、記憶装置制御装置 1 1 は、ライト I / O 処理の終了をホスト 1 に通知する。この際、記憶装置制御装置 1 1 は、これらのデータに対応するジャーナルデータの作成も行う (ステップ 1 1 0 3) 。

30

【 0 0 9 8 】

尚、ライト I / O 処理の終了の通知を受信したホスト 1 は、以下のステップで実行される C P I D 書き込み処理の完了が記憶装置システム 2 から報告されるまでは、記憶装置システム 2 へのデータの書き込みを行わないが、データの読み出しは実行してもよい。

【 0 0 9 9 】

ライト I / O 処理の終了が通知されたホスト 1 は、トランザクションマネージャ 5 0 2 を実行して、C P 情報及び C P 処理に用いられるログを作成する。具体的には、ログファイル 5 0 6 に C P I D 等の C P 情報をログとして格納する。尚、ログの C P 情報には、C P I D、リソースマネージャの数、リソースマネージャの状態、動作中のトランザクションの数及び各々のトランザクション記述なども含まれる。尚、リソースマネージャに関しては、詳細を割愛する (ステップ 1 1 0 4 ~ 1 1 0 5) 。同時に、ホスト 1 は、チェックポイントコマンドを記憶装置システム 2 に対して発行する。チェックポイントコマンドには C P I D が含まれている (ステップ 1 1 0 5) 。

40

【 0 1 0 0 】

ホスト 1 からのチェックポイントコマンドを受信した記憶装置システム 2 は (ステップ 1 1 0 6) 、受信した C P I D をジャーナルデータとしてジャーナル論理記憶装置に記録する。この場合、ジャーナルデータのエントリ 1 0 0 4 及び 1 0 0 5 に対応するライト対

50

象データ及びライトデータは存在しないので、これらのエントリには、データが格納されないか、無効データ（例えば - 1）が格納される（ステップ 1107）。記録が完了したら、記憶装置制御装置 11 は、記録の完了をホスト 1 に通知する（ステップ 1108）。

【0101】

ホスト 1 は、記憶装置システム 2 から C P I D 記録完了の報告を受領すると、C P 情報に関する処理を終了する（ステップ 1109）。

【0102】

図 11 は、管理端末 3 やホスト 1 からリストア指示を受領した記憶装置システム 2 における処理手順を示す図である。尚、以下の処理は、記憶装置制御装置 11 が、リストアプログラム 421 を実行することで行われる。

10

【0103】

本実施形態では、D B を使用する A P 540 のバグやユーザのオペレーションミス等により論理記憶装置 408 にホスト 1 にとって論理的不整合等の障害が起き、かつ障害が発生した論理記憶装置 408 がジャーナルモードの正論理記憶装置であった場合を考える。この場合、管理端末 3 又はホスト 1 からは、障害が発生した正論理記憶装置に対応する副論理記憶装置及びジャーナル論理記憶装置に格納されたデータを使用して記憶装置システム 2 内で正論理記憶装置に格納されたデータをリストアする指示が送信される。

【0104】

ホスト 1 は、A P 540 のログ情報等を参照し、オペミスや誤ったデータを送信した A P 等の誤った操作を起こした時点を解析し、その時点の直前のチェックポイントコマンド送信時を検索し、記憶装置システム 2 でリストアする際に使用される C P I D を決定する。尚、ホスト 1 のユーザは、障害発生直前の C P I D ではなく、ホスト 1 から C P 情報を記憶装置システム 2 に送信する際にホスト 1 に記録される C P I D のリストから、任意の C P I D を選択することができる。これにより、本システムのユーザは、任意の C P I D を選択することで、選択された C P I D が作成された時点に記憶装置システム 2 の正論理記憶装置が格納していたデータの状態まで、正論理記憶装置に格納されたデータをリストアすることができる（ステップ 1201）。

20

【0105】

次に、ホスト 1 は、ステップ 1201 で選択した C P I D までのデータのリストア処理要求を記憶装置システム 2 に発行する。リストア処理要求には、リストア処理の対象となる正論理記憶装置の識別子（例えば W W N と L U N 等）、正論理記憶装置に対応する副論理記憶装置を指定する識別子、ジャーナル論理記憶装置を指定する識別子、及び選択された C P I D の情報等が含まれる。尚、正論理記憶装置に対応する副論理記憶装置が複数有る場合は、その内のいずれかを指定する情報もリストア処理要求に含まれる（ステップ 1202）。

30

【0106】

ホスト 1 より発行されたリストア処理要求を受領した記憶装置制御装置 11 は、リストアプログラム 421 を実行して、リストア処理要求に含まれる副論理記憶装置を示す識別子とペア管理情報 412 を比較参照し、指定された副論理記憶装置が正論理記憶装置に対する正しい副論理記憶装置であるかを確認する。また同様に、リストア処理要求に含まれるジャーナル論理記憶装置を示す識別子とジャーナル管理情報とを比較参照し、指定されたジャーナル論理記憶装置が正論理記憶装置に対応する正しいジャーナル論理記憶装置であるかを確認する（ステップ 1203）。

40

【0107】

更に、記憶装置制御装置 11 は、リストア処理要求の内容から、正論理記憶装置にリストア処理を行うのか、副論理記憶装置にリストア処理を行うのか、もしくは異なった未使用の論理記憶装置 408 にリストア処理を行うのかを確認する。尚、リストア処理対象に正論理記憶装置が指定されていても、正論理記憶装置が使用不可能であれば、論理記憶装置 408 の障害により処理続行が出来ない旨をホスト 1 に通知し、処理を中止する。また同様に副論理記憶装置やその他の論理記憶装置 408 にデータをリストアする指示であって

50

も、指定された論理記憶装置 408 に何らかの障害がある場合は障害により処理続行が出来ない旨をホストに通知し、処理を中止する（ステップ 1204）。

【0108】

正論理記憶装置もしくはその他の空き論理記憶装置 408 にリストア処理を行う場合、記憶装置制御装置 11 は、副論理記憶装置に格納されていたスナップショットデータを先頭から順次読み出して正論理記憶装置へコピーし、正論理記憶装置が有するディスクイメージを副論理記憶装置と同一にする。尚、副論理記憶装置にデータをリストアする場合は、本コピー処理は不要である（ステップ 1206）。

【0109】

副論理記憶装置からのコピー処理が終了したら、あるいは副論理記憶装置へデータをリストアする場合、記憶装置制御装置 11 は、キャッシュメモリ 14 にデータ格納領域を確保する。その後、記憶装置制御装置 11 は、正論理記憶装置に対応するジャーナル論理記憶装置の先頭から、具体的には、処理シーケンス番号順に、順次ジャーナルデータをキャッシュメモリ 14 に確保された領域に読み出す。尚、ジャーナル論理記憶装置からのジャーナルデータの読み出しの先頭は、ホスト 1 から指定されても、記憶装置システム 2 が処理シーケンス番号で特定しても良い（ステップ 1207）。

10

【0110】

その際、読み出されたジャーナルデータに CP 情報が含まれるかどうかを確認する。具体的には、ジャーナルデータのチェックポイントフラグ 1006 が ON になっているかどうかを確認する（ステップ 1208）。

20

【0111】

読み出されたジャーナルデータが、CP 情報を含むジャーナルデータである場合、記憶装置制御装置 11 は更に、読み出されたジャーナルデータの CPID1007 に含まれる CPID がホスト 1 から指定された CPID かどうかを確認する（ステップ 1209）。

【0112】

CPID1007 に含まれる CPID がホスト 1 から指定された CPID でない場合又は CPID1007 に CPID が格納されていない場合（チェックポイントフラグが ON になっていない場合）、記憶装置制御装置 11 は、読み出したジャーナルデータのアドレス 1002 に格納された情報から、読み出されたジャーナルデータが、指定されたリストア対象である正論理記憶装置に関するジャーナルデータであるかどうかを確認する（ステップ 1210）。

30

【0113】

読み出されたジャーナルデータがリストア対象の正論理記憶装置に関するジャーナルデータであれば、記憶装置制御装置 11 は、読み出されたジャーナルデータに含まれるライトデータを、正論理記憶装置又は副論理記憶装置の対応するアドレスに書き込む。ただし、CPID に対応するジャーナルデータである場合には、ライトデータが存在しないので、データの書き込みは行われない（ステップ 1211）。

【0114】

その後、記憶装置制御装置 11 は、ステップ 1207 に戻り次のジャーナルデータの読み出し処理を行う。また、ステップ 1210 で読み出されたジャーナルデータが指定された正論理記憶装置に対応するジャーナルデータでない場合、記憶装置制御装置 11 は、ジャーナルデータをリストア先である論理記憶装置 408 に書き込まずに、ステップ 1207 の処理に戻る。以下、記憶装置制御装置 11 は、ステップ 1207 ~ 1211 の処理を繰り返すことで、指示された CPID までのジャーナルデータをリストアする。

40

【0115】

ステップ 1209 で、CPID1007 の CPID が指定された CPID と一致した場合、記憶装置制御装置 11 は、リストアすべきデータをすべて正論理記憶装置、副論理記憶装置や他の論理記憶装置 408 に書き込んだと判断して、リストア処理の終了をホスト 1 に通知する。尚、正論理記憶装置以外にリストア処理を行う場合は、ホスト 1 への通知前に、論理物理マッピング情報を書き換えて、正論理記憶装置と副論理記憶装置またはそ

50

他のリストア先となる論理記憶装置 4 0 8 とを交換し、ホスト 1 からアクセスする論理記憶装置 4 0 8 の識別子（たとえば F C の W W N と L U 番号の組合せ）は変わらないようにする（ステップ 1 2 1 2）。

【 0 1 1 6 】

尚、正論理記憶装置毎にジャーナル論理記憶装置が割り当てられている場合には、前記ステップ 1 2 1 0 の処理、すなわち、読み出されたジャーナルデータと正論理記憶装置との対応関係の確認は不要である。

【 0 1 1 7 】

ホスト 1 または管理端末 3 は、記憶装置システム 2 から終了報告を受領したら、ホスト 1 が指定した C P I D 時点までのデータが回復されたと判断して、他の処理を継続する（1 2 1 3）。

10

【 0 1 1 8 】

図 1 2 は、デバイス管理情報 4 1 0 の一例を示した図である。

【 0 1 1 9 】

デバイス管理情報 4 1 0 は、論理記憶装置 4 0 8 のアドレス情報を登録するエン트리 1 3 0 1 及び物理記憶装置 1 0 のアドレス情報を登録するエン트리 1 3 0 4 とを有するテーブル 1 3 0 0、ホスト 1 に提供される論理記憶装置番号を登録するエン트리 1 3 3 1、記憶装置システム 2 で論理記憶装置 4 0 8 を統一的に識別する記憶装置内論理記憶装置番号を登録するエン트리 1 3 3 2、記憶装置システム 2 内で管理する P a r i t y G r o u p の通し番号を登録するエン트리 1 3 3 3、論理記憶装置 4 0 8 のペア情報を登録するエン트리 1 3 3 4 及びジャーナル情報を登録するエン트리 1 3 3 5 を有するテーブル 1 3 3 0 並びに、記憶装置システム 2 内の論理記憶装置番号が登録されるエン트리 1 3 5 1、空き / リザーブ情報が登録されるエン트리 1 3 5 2、P a t h 定義情報が登録されるエン트리 1 3 5 3、E m u l a t i o n T y p e / サイズが登録されるエン트리 1 3 5 4 及び障害情報が登録されるエン트리 1 3 5 5 とを有するテーブル 1 3 5 0 とを保持する。

20

【 0 1 2 0 】

テーブル 1 3 0 0 のエン트리 1 3 0 1 は、更に、ホスト 1 に提供される論理記憶装置 4 0 8 の番号が登録されるエン트리 1 3 1 1、その論理記憶装置 4 0 8 に対応する内部アドレスが登録されるエン트리 1 3 1 2、記憶装置システム 2 内部で論理記憶装置 4 0 8 を統一的に識別する論理記憶装置番号が登録されるエン트리 1 3 1 3 及びその内部論理記憶装置アドレスが登録されるエン트리 1 3 1 4 を有する。また、テーブル 1 3 0 0 のエン트리 1 3 0 4 は、更に、エン트리 1 3 0 1 に登録された論理記憶装置 4 0 8 に対応する物理記憶装置 1 0 の P a r i t y G r o u p 4 0 7 の番号を登録するエン트리 1 3 2 1、物理記憶装置 1 0 の番号を登録するエン트리 1 3 2 2 及びその物理記憶装置 1 0 のアドレス情報を登録するエン트리 1 3 2 3 を有する。

30

【 0 1 2 1 】

テーブル 1 3 3 0 のペア情報エン트리 1 3 3 4 には、論理記憶装置 4 0 8 がスナップショットペア状態にあるかどうかを示す情報が登録される。ジャーナル対象モードエン트리 1 3 3 5 には、論理記憶装置 4 0 8 がジャーナル取得の対象、すなわちジャーナルモードの対象であるかどうかを示す情報が登録される。

40

【 0 1 2 2 】

テーブル 1 3 5 0 の空き / リザーブ情報エン트리 1 3 5 2 には、論理記憶装置 4 0 8 が、副論理記憶装置やジャーナル論理記憶装置に用いるために予約されている状態にあるかどうかを示す情報が登録される。リザーブ情報が登録されている論理記憶装置 4 0 8 は、その他の用途、例えば新たに業務用論理記憶装置として割り当てるなどが出来ない。P a t h 定義情報エン트리 1 3 5 3 には、論理記憶装置 4 0 8 がホスト 1 に提供されるために外部に公開されているかどうかを示す情報が登録される。例えば I / O N e t w o r k が F C だったら、論理記憶装置 4 0 8 と F C の P o r t との関連付けに関する情報が登録される。

【 0 1 2 3 】

E m u l a t i o n T y p e エン트리 1 3 5 4 には、論理記憶装置 4 0 8 が O S が認

50

識できる記憶装置のいずれに擬似化されている（エミュレートされる）かを示す情報及びその記憶容量が登録される。例えば、具体的には、オープン系システムのOSが認識できる記憶装置であることを示す「OPEN」や、メインフレーム系のOSが認識できる記憶装置であることを示す「3990」等の情報が登録される。

【0124】

障害情報エントリ1355には、論理記憶装置408が何らかの障害になったかどうかを示す情報が登録される。ここで、障害とは、主に論理記憶装置408が存在する物理記憶装置10の物理的障害や管理者が意識的に記憶装置システム2を閉塞状態にした場合等の論理的障害がある。

【0125】

図13は、ペア管理情報情報412のテーブルの一例を示した図である。

【0126】

ペア管理情報412は、ホスト1に提供される論理記憶装置番号を登録するエントリ1401、記憶装置システム2内での論理記憶装置番号を登録するエントリ1402、Emulation Type/サイズを登録するエントリ1403、ペア状態を登録するエントリ1404、世代情報を登録するエントリ1405及びペア管理情報を登録するエントリ1406とを有する。

【0127】

ペア状態エントリ1404には、先に記したペア結合状態等のペアの状態を示す情報が登録される。ペア管理情報エントリ1406には、論理記憶装置408が正論理記憶装置か副論理記憶装置かを示す情報が登録される。論理記憶装置408が正論理記憶装置に指定されていれば、正側エントリ1411には0が登録され、対応する副側エントリ1412にはペアとなる副論理記憶装置の番号を示す値が登録される。一方、論理記憶装置408が副論理記憶装置に指定されていれば、副側エントリ1411には0の値が登録され、対応する正側エントリ1412にはペアとなる正論理記憶装置の番号を示す情報が登録される。

【0128】

また、論理記憶装置408が正副論理記憶装置として指定されていない場合には、正側エントリ1411及び副側エントリ1412の双方に無意味な値を示す「-1」が登録される。また、論理記憶装置408がスナップショットペア450のカスケード構成の真ん中、すなわち、一つのペアの副論理記憶装置でもあり、同時に他のペアの正論理記憶装置である場合は、正側エントリ1411、副側エントリ1412双方にペアを形成する他方の論理記憶装置408の番号を示す情報が登録される。また、正側エントリ1411、副側エントリ1412に複数の論理記憶装置番号が登録される場合もある。

【0129】

図14は、ジャーナル管理情報418の一例を示した図である。

【0130】

ジャーナル管理情報418は、テーブル1500及びCP情報を管理するためのジャーナル管理テーブル1520を有する。テーブル1500は、CPIDを格納するエントリ1501、エントリ1501に格納されたCPIDが記録されたジャーナルデータが格納された位置を示すアドレスが登録される1502及びエントリ1501に格納されたCPIDがジャーナル論理記憶装置に記録された時間を示す時間情報1503とを有する。また、ジャーナル管理テーブル1520は、デバイス番号を登録するエントリ1521ごとに、CPIDを登録するエントリ1522及びチェックポイント管理テーブルの格納アドレスを登録するエントリ1523を有する。

【0131】

次に、第二の実施形態として、ホスト1ではなく、管理端末3と記憶装置システム2との間でCP情報を共有し、記憶装置システム2に障害が起きた場合のデータのリカバリを行う場合について述べる。

【0132】

本実施形態では、ホスト1が記憶装置システム2との間のログやチェックポイントを管理するプログラム、例えばDBMS501を有しない場合に、ホスト1にエージェントというプログラムを導入する。以下エージェントが導入されたホストをホスト1'と称する。

【0133】

図22は、ホスト1'が有するプログラムの構成を例示した図である。ホスト1と異なる点は、DBMS501が存在せず、代わりにAgentプログラム2200が含まれている点である。Agentプログラム2200は、モード情報2210、FS Agent 820、I/O制御プログラム2230、チェックポイントAgent 2250、VM Agent 2240、及び構成管理Agent 2260から構成されている。

10

【0134】

モード情報2210には、ホスト1'が管理端末3から受信した、スナップショットを取る時期やジャーナルデータを取る期間の状態が、モード情報として保持されている。FS Agent 2220は、FS530に対してファイルの排他制御やファイルを閉じる処理を指示し、かつFS530が管理するダーティデータをメモリ21のアドレスとして管理する際に実行される。

【0135】

VM Agent 2240は、VM540に対して、VM540で設定される論理記憶領域への読み出し/書き込みの可否を制御し、かつVM540が管理するダーティデータをメモリ21のアドレスとして管理するために実行される。

20

【0136】

I/O制御プログラム2230は、ホスト1'が、記憶装置システム2に強制的にダーティデータを転送する処理を行う際に実行される。構成管理Agent 2260は、記憶装置システム2がホスト1'に提供する論理記憶装置408とVM540が構成する論理記憶領域との対応関係、及びVM540が構成する論理記憶領域とFSが構成する論理記憶領域との関係を管理する際に実行される。

【0137】

チェックポイントAgent 2250は、管理端末3からチェックポイントについて指示された際に、ホスト1'が、モード情報2210の設定、FS Agent 2220、VM Agent 2240、及びI/O制御プログラム2230等に所定の動作を指示する際に実行される。

30

【0138】

ホスト1'は、管理端末3からの指示により、エージェントプログラム2200を実行して、ホスト1'のメモリ21に存在するダーティデータを記憶装置システム2に送信する。一方、ホスト1'からのダーティデータの送信に合わせて、管理端末3は、チェックポイントコマンドを記憶装置システム2に送る。記憶装置システム2は、ホスト1'から送信されたダーティデータを処理する。記憶装置システム2は、また、管理端末3から送信されたCP情報を、第一の実施形態で説明したホスト1から送信されたCP情報と同様に扱って、自システム2内で管理する。このようにすることで、正論理記憶装置に論理的な障害が発生した際に、ホスト1にチェックポイント作成等の機能が無い場合でも、管理

40

【0139】

図15は、管理端末3の詳細な構成を示した図である。尚、本構成は、他の実施形態で使用されてもよい。

【0140】

管理端末3は、プロセッサ1601、電氣的に不揮発なメモリ1602、ネットワークI/F1605、入力部1604及び表示部1603とを有する。また、各々の構成部品は、データや制御命令等を伝送する伝送路1612で接続されている。

【0141】

50

プロセッサ 1601 は、管理端末 3 が有するプログラムを実行する。メモリ 1602 には、プロセッサ 1601 が実行するプログラムおよびそのプログラムが使用する情報等が格納される。例えば、表示部制御プログラム 1610、入力部制御プログラム 1611、記憶装置システム 2 の構成を管理する記憶装置制御情報 1606、記憶装置制御情報 1606 に登録された情報を使用して記憶装置システム 2 を制御・管理するための記憶装置管理プログラム 1607、記憶装置システム 2 に送信した CP 情報が含まれるシステム確定情報 1608、及びシステム確定情報 1608 に登録された情報を用いて記憶装置システム 2 の状態を所定の時点に復旧する等の制御処理等を行うためのシステム状態管理プログラム 1609 等がメモリ 1602 に登録される。

【0142】

10

ネットワーク I/F 1605 はネットワーク 5 に接続されている。管理端末 3 は、ネットワーク 5 を介して記憶装置システム 2 のシステム構成、例えばデバイス管理情報 410、ペア管理情報 412 及びジャーナル管理情報 418 を取得する。又、管理端末 3 は、ネットワーク 5 を介して、構成定義処理（例えば Parity Group 407 に論理記憶装置 408 を定義し、記憶装置システム 2 内部の論理記憶装置番号を割り振ることや論理記憶装置 408 をホスト 1' に使用可能にするためにパスを定義してホスト 1' が使用する論理記憶装置番号を割り振ること）をしたり、記憶装置システム 2 のリストア処理の実行を制御したりする。

【0143】

また、記憶装置システム 2 のユーザ又は管理者は、入力部 1604 及び表示部 1603 を使用して、記憶装置システム 2 の保守/管理やリストア処理の指示等を行う。

20

【0144】

図 16 は、メモリ 1602 に格納されるシステム確定情報 1608 の一例を示す図である。管理端末 3 は、ホスト 1' の状態が確定する時点を経験装置システム 2 に指示する際に、管理端末 3 自身で記憶装置システム 2 に指示した内容をシステム確定情報 1608 としてメモリ 1602 に記録する。システム確定情報 1608 は、システムの状態が確定する時点の CPID を登録するエントリ 1701、論理記憶装置を示す番号が登録されるエントリ 1702 及びシステムの状態が確定する時点の時間を登録するエントリ 1703 を有する。

【0145】

30

図 17 は、表示部 1603 における表示の一例を示す図である。本図では、表示部 1603 に、図 16 に示したシステム確定情報 1608 の内容が GUI を用いて表示されたものを例示している。このように、表示部 1603 は、システム状態が確定された時間を複数表示し、表示された複数の時間からユーザがある時間を選択したことを表示することができる。これにより、ユーザの利便性が向上する。

【0146】

具体的には、表示部 1603 は、管理情報を表示する領域 1802 を有する。その領域 1802 には、論理記憶装置番号を表示する領域 1803 及び領域 1803 に表示された論理記憶装置 408 の状態を確定した時間が表示される領域 1804 が含まれる。ユーザは、マウス等で操作可能なポインタ 1805 で、表示された論理記憶装置 408 について、チェックポイントコマンドによって状態が確定された時間を指定することができる。

40

【0147】

又、ユーザは、ある論理記憶装置 408 に障害が起きた場合、記憶装置システム 2 に対して、管理端末 3 の GUI 1603 を介してリストア処理の指示を行う。例えば、本図では、領域 1803 に表示された論理記憶装置 408 の内容を、領域 1804 で示された時刻中、2002 年 5 月 5 日 14:00 の時点にリストアするための指示例を示している。ユーザは、ポインタ 1805 を用いて 2002 年 5 月 5 日 14:00 を示す領域 1804 を選択し、それを領域 1803 へ Drag & Drop 等を行うことで、論理記憶装置 408 のリストア時刻を指示する。

【0148】

50

管理端末3は、ユーザによって指定された論理記憶装置408及びリストア時間に基づいて、図16に示されたシステム確定情報1608を検索し、リストアに使用するチェックポイントを特定する。その後、管理端末3は、記憶装置システム2に、検索の結果得られたCP情報をリストアコマンドを用いて送信する。

【0149】

図18は、ユーザが、管理端末3を介してジャーナルデータ開始準備処理を情報処理システムに指示する処理の流れを示した図である。

【0150】

まず、ユーザは、管理端末3の表示部1603及び入力部1604を用いて、ジャーナルデータを取得すべき対象となる正論理記憶装置や副論理記憶装置を指定する。管理端末3は、ユーザの指定に基づいて、ジャーナル取得準備指示コマンドを記憶装置システム2にネットワーク5を介し送信する。ジャーナル取得開始準備指示には、ユーザが指定した正論理記憶装置を示す識別子、当該正論理記憶装置と対になる副論理記憶装置を示す識別子、ジャーナル論理記憶装置を示す識別子が含まれる(ステップ1901)。

【0151】

ジャーナル準備処理指示を受領した記憶装置システム2は(ステップ1961)、ジャーナル準備処理を実行する。本処理は、図6のステップ602～ステップ606で説明された処理と同様の処理である(ステップ1962)。ジャーナル準備処理を終了した記憶装置システム2は、ネットワーク5を介して、管理端末3に終了報告を送信する(ステップ1963)。

【0152】

完了報告を受信した管理端末3は(ステップ1902)、ホスト1'にジャーナル開始モード指示のコマンドをネットワーク5を介して送信する(1903)。

【0153】

ジャーナル開始モード指示のコマンドを受領したホスト1'は、Agent800を実行することで、ジャーナルデータ取得の対象となる正論理記憶装置に対応するモード情報810をジャーナル開始モードに設定する。更に、ホスト1'は、ジャーナル開始モードに設定された正論理記憶装置に格納されるべきダーティデータを確定するために、ファイルの使用を終了する。尚、ジャーナル開始モード中は、ジャーナル開始モードを設定された正論理記憶装置に関連する記憶領域は書き込み禁止となる(ステップ1921)。

【0154】

次にホスト1'は、FSが管理するメモリ21に格納されたダーティデータをすべて記憶装置システム2に送信するため、記憶装置システム2にライトI/O処理要求を出す(ステップ1922)。

【0155】

ホスト1'からライトI/O処理要求を受け付けた記憶装置システム2は、ユーザが指定した正論理記憶装置への書き込み処理であれば、ジャーナル作成処理を行う。処理が終了すると、記憶装置システム2は、ホスト1'へ完了を報告する(ステップ1965)。

【0156】

完了の報告を受取ったホスト1'は、FSが管理する全ダーティデータを記憶装置システム2に書き込んだかどうかを判断する(ステップ1923)。全ダーティデータの書き込みが完了していない場合、ホスト1'は、ステップ1922から処理を繰り返す。全ダーティデータの書き込みが終了した場合、ホスト1'は管理端末3に、完了報告をネットワーク5を介し送信する(ステップ1925)。

【0157】

ダーティデータのライト完了報告を受領した管理端末3は、記憶装置システム2に対しチェックポイントコマンドを発行するとともに、システム確定情報1608の更新を行う。具体的には、管理端末3は、ジャーナルデータを取得する論理記憶装置408を指定するデバイス番号に対応するエントリに、送信したCPIDと送信した時間を記録する(1905)。

10

20

30

40

50

【 0 1 5 8 】

チェックポイントコマンドを受領した記憶装置システム 2 は (ステップ 1 9 6 6)、受領したチェックポイントコマンド中の C P I D をジャーナルデータとしてジャーナル論理記憶装置に記録する (ステップ 1 9 6 7)。記録が完了したら、記憶装置システム 2 は、完了報告を管理端末 3 にネットワーク 5 を介し送信する (ステップ 1 9 6 8)。

【 0 1 5 9 】

完了報告を受領した管理端末 3 は (1 9 0 6)、ホスト 1' に対してジャーナル開始モード解除指示をネットワーク 5 を介し送信する (ステップ 1 9 0 7)。ジャーナル開始モード解除指示を受領したホスト 1' は、ステップ 1 9 2 1 で設定された、正論理記憶装置に対応するモード情報 8 1 0 のジャーナル開始モードを解除する。その後、ホスト 1' は、正論理記憶装置に対応する記憶領域への書き込み禁止も解除する (ステップ 1 9 2 7)。

10

【 0 1 6 0 】

その後、ユーザは、管理端末 3 を用いて、所定のタイミングでジャーナルモード開始指示をホスト 1' 及び記憶装置システム 2 に送信する。ジャーナルモード開始指示を受信したホスト 1' は、指示で指定される正論理記憶装置に対応するモード情報 8 1 0 にジャーナルモードを設定する。一方、ジャーナルモード開始指示を受信した記憶装置システムは、先に指定されたジャーナル論理記憶装置にジャーナルデータの記録を開始する。

【 0 1 6 1 】

図 1 9 は、ジャーナルデータを取得している正論理記憶装置の内容を後にリストアできるように、ユーザの指示等に基づいて、ホスト 1' の代わりに管理端末 3 がチェックポイントコマンドを記憶装置システム 2 に送信し、記憶装置システム 2 と管理端末 3 双方で一意の C P I D を格納する際の処理手順を示した図である。

20

【 0 1 6 2 】

管理端末 3 は、ユーザの指示もしくは管理端末 3 自身のプログラムの実行に基づいて、チェックポイントモード指示をホスト 1' にネットワーク 5 を介し送信する。チェックポイントモード指示には、チェックポイント取得の対象となる論理記憶装置 4 0 8 を示す番号が含まれている (ステップ 2 0 0 1)。

【 0 1 6 3 】

チェックポイントモード指示を受け取ったホスト 1' は、A g e n t プログラム 8 0 0 を実行して、指示に含まれる論理記憶装置の番号及びモード情報 8 1 0 に登録された情報とを参照し、指示された論理記憶装置 4 0 8 がジャーナルモードであることを確認する。指示された論理記憶装置 4 0 8 がジャーナルモードである場合、ホスト 1' は、メモリ 2 1 にあるダーティデータを記憶装置システム 2 へ強制的に転送する。

30

【 0 1 6 4 】

具体的には、ホスト 1' は、構成定義 A g e n t F S A g e n t 2 2 2 0 を実行して、指定された論理記憶装置 4 0 8 を使用しているファイルが使用されているかを確認する。その後、ホスト 1' は、F S A g e n t 8 2 0 を実行して、使用しているファイルを終了する又は使用しているファイルへの書き込み要求が実行されないようにする。その後、ホスト 1' は、メモリ 2 1 に格納されたダーティデータの転送を要求するライト I / O 処理要求を記憶装置システム 2 に送信する。尚、ホスト 1' が V M を使用している場合は、ホスト 1' は、上述と同様の処理を、V M A g e n t 2 2 4 0 を実行して行う (ステップ 2 0 2 2)。

40

【 0 1 6 5 】

ライト I / O 処理要求を受けた記憶装置システム 2 は、図 7 で説明したフローに従って、ジャーナルデータをジャーナル論理記憶装置に格納する処理を行う (ステップ 2 0 6 1、ステップ 2 0 6 2)。

【 0 1 6 6 】

ジャーナル作成の完了を受信したホスト 1' は、全てのダーティデータが記憶装置システム 2 に格納されたかどうかを確認する。全てのダーティデータが記憶装置システム 2 に

50

格納されていない場合、ホスト 1' は、ステップ 2022 からの処理を繰り返す（ステップ 2023）。

【0167】

全てのダーティデータが記憶装置システム 2 に格納されたと確認した場合、ホスト 1' は、管理端末 3 に、チェックポイントモード指示に対する応答メッセージとしてダーティデータのライト完了報告を送信する（ステップ 2025）。

【0168】

ライト完了報告を受領した管理端末 3（ステップ 2002）は、記憶装置システム 2 に対し、チェックポイントコマンドを発行するとともに、システム確定情報 1608 の更新を行い、処理対象である論理記憶装置を示すデバイス番号に対応するエントリに、送信した C P I D と送信した時間を記録する（ステップ 2003）。

10

【0169】

チェックポイントコマンドを受領した記憶装置システム 2 は（ステップ 2063）、受領したチェックポイントコマンドに含まれる C P I D をジャーナルデータとして記録する（ステップ 2064）。その後、記憶装置システム 2 は、完了報告を管理端末 3 にネットワーク 5 を介し送信する（ステップ 2065）。

【0170】

完了報告を受領した管理端末 3 は、ホスト 1' に対して、チェックポイントモード解除指示をネットワーク 5 を介し送信する（ステップ 2004）。

【0171】

20

チェックポイントモード解除指示を受領したホスト 1' は、ステップ 2021 でチェックポイントモードを設定された論理記憶装置 408 に対応するモード情報 810 に登録されたチェックポイントモードを解除する（ステップ 2026）。その後、ホスト 1' は、ファイルへの書き込みを再開するか、ファイルを使用可能状態とする（ステップ 2027）。

【0172】

図 21 は、ユーザが、管理端末 3 を介して、記憶装置システム 2 へリストア指示を出す際の処理手順を示す図である。本実施形態では、スナップショットペア 450 が既に形成され、副論理記憶装置に正論理記憶装置のスナップショットが取得されていて、ジャーナル論理記憶装置には、副論理記憶装置にスナップショットをとる時点より後もしくは前後のジャーナルデータが格納されているものとする。

30

【0173】

この場合において、正論理記憶装置を使用していたホスト 1' の A P が、使用しているファイルに誤った編集をした等の理由で、バックアップデータに基づくリストアが必要になった場合を考える。

【0174】

まず、ユーザは、図 17 で説明したように、管理端末 3 の入力部 1604 及び表示部 1603 を用いて、リストアの対象となる論理記憶装置 408 及びどの時点までリストアを行うかを指示する（ステップ 2101）。

【0175】

40

ユーザの指示を受けた管理端末 3 は、ユーザが画面上で指定した入力情報が、システム確定情報 1608 のどのエントリに登録された情報と一致するかを判断し、一致したエントリに登録されている C P I D を決定する（ステップ 2102）。その後、管理端末 3 は、リストアを行う論理記憶装置 408 を示す識別子（番号）及び C P I D を含むリストアコマンドを、記憶装置システム 2 に送信する（ステップ 2103）。

【0176】

管理端末 3 からリストアコマンドを受信した記憶装置システム 2 は、指定された論理記憶装置 408 について、図 11 で説明したリストア処理を実行する。その後、リストア処理完了報告を管理端末 3 に送信する（ステップ 2104）。完了報告を受領した管理端末 3 は、記憶装置制御情報 1606 を更新する（ステップ 2105）。

50

【 0 1 7 7 】

本実施形態によれば、第一の実施形態と比較し、ホスト 1' が C P 情報を管理することが無いので、その分ホスト 1' の負荷を低減することが出来る。また、ホスト 1 等がチェックポイント作成機能を有さない場合でも、C P 情報を用いたリストア処理を行うことが出来る。

【 0 1 7 8 】

尚、本実施形態では、管理端末 3 から記憶装置システム 2 に対してチェックポイントコマンドを発行する前に、管理端末 3 は、ホスト 1' に格納されているダーティデータを記憶装置システム 2 に反映させるために、ホスト 1' に対してダーティデータをフラッシュさせる指示（ジャーナルモード開始指示、チェックポイントモード指示）を送信した。しかし、この場合、上述したように、ホスト 1' にエージェントというプログラムを用意しなければならない。したがって、全てのホスト 1' にエージェントを用意するのが困難な場合、本実施形態は採用しづらい。そこで、ホスト 1' に存在するダーティデータを無視し、管理端末 3 と記憶装置システム 2 のみでジャーナルモードの設定、C P 情報の遣り取り及びリストア処理の実行を行う第三の実施形態を考える。

【 0 1 7 9 】

本実施形態は、第二の実施形態と以下の点で異なる。すなわち、図 1 8 において、ジャーナル作成準備処理の完了を報告された管理端末 3（ステップ 1 9 0 2）は、ステップ 1 9 0 3 の処理を行わずに、直接記憶装置システム 2 に対してチェックポイントコマンドを送付するステップ 1 9 0 4 の処理を行う。記憶装置システム 2 では、受信したチェックポイントコマンドに従って、ステップ 1 9 6 6 以降の処理を行う。

【 0 1 8 0 】

また、図 1 9 においては、管理端末 3 は、ステップ 2 0 0 1 のチェックポイントモード指定をホスト 1' に送信せず、直接チェックポイントコマンドを記憶装置システム 2 に送信する（ステップ 2 0 0 3）。チェックポイントコマンドを受信した記憶装置システム 2 は、ステップ 2 0 6 3 以降の処理を行う。

【 0 1 8 1 】

尚、本実施形態で使用されるホストは、ホスト 1 のように D B のログを有する計算機でも、ホスト 1' のようにエージェントを有する計算機でも、あるいは、何ら特別なプログラムを有さない通常の計算機でも良い。他の構成及び処理、例えばリストア処理等は、第二の実施形態と同様である。

【 0 1 8 2 】

本実施形態によれば、ホストの種別に関わらず、管理端末 3 及び記憶装置システム 2 との遣り取りだけで、記憶装置システム 2 の記憶装置を、任意のシステム状態までにリストアすることができる。

【図面の簡単な説明】

【 0 1 8 3 】

【図 1】本発明を適用した情報処理システムの構成を示す図である。

【図 2】記憶装置システム 2 の構成を示す図である。

【図 3】ホスト 1 の構成を示す図である。

【図 4】記憶装置システム 2 が有するプログラム等の構成を示す図である。

【図 5】ホスト 1 が有するプログラム等の構成を示す図である。

【図 6】ジャーナル取得準備の処理手順を示す図である。

【図 7】ジャーナルモード中の I / O 処理要求の手順を示す図である。

【図 8】ジャーナルデータ作成処理の手順を示す図である。

【図 9】ジャーナルデータの形式を示す図である。

【図 1 0】C P 情報の送信処理の手順を示す図である。

【図 1 1】リストア処理の手順を示す図である。

【図 1 2】デバイス管理情報の構成例を示す図である。

【図 1 3】ペア管理情報の構成例を示す図である。

10

20

30

40

50

【図 1 4】ジャーナル管理情報の構成例を示す図である。

【図 1 5】管理端末 3 の構成を示す図である。

【図 1 6】システム確定情報 1 6 0 8 の構成例を示す図である。

【図 1 7】管理端末 3 の表示部の構成例を示す図である。

【図 1 8】第二の実施形態におけるジャーナルデータ取得指示の処理手順を示す図である。

【図 1 9】第二の実施形態における C P I D の送受信の処理手順を示す図である。

【図 2 0】第二の実施形態におけるリストア指示処理の手順を示す図である。

【図 2 1】第二の実施形態におけるホスト 1 ' の論理的構成を示す図である。

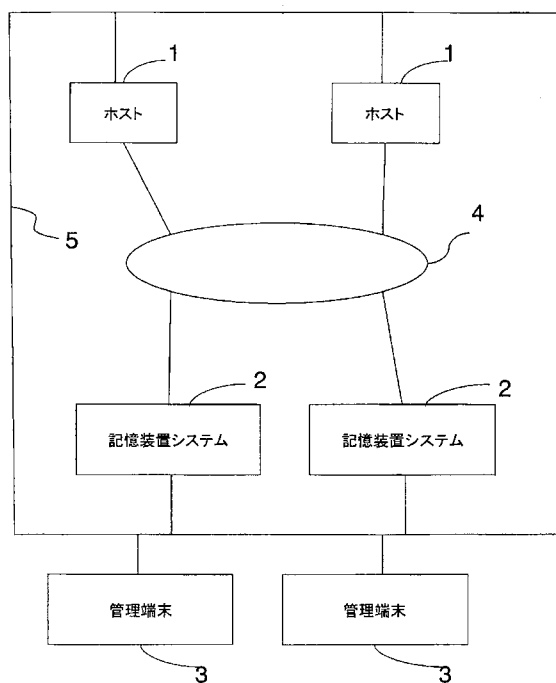
【符号の説明】

10

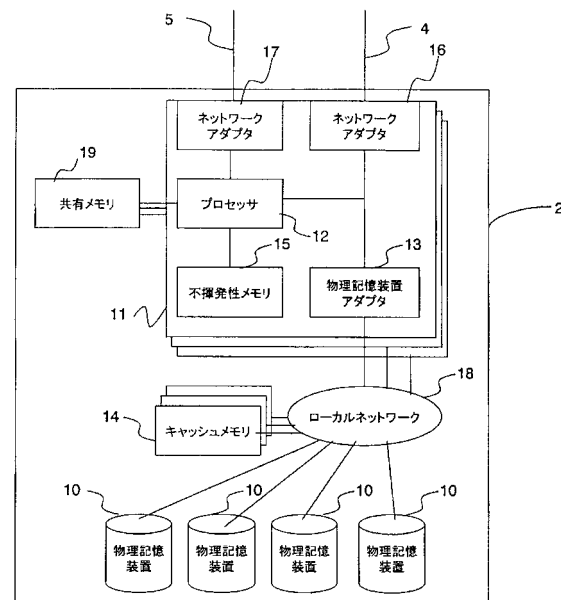
【 0 1 8 4 】

- 1 ホスト
- 2 記憶装置システム
- 3 管理端末
- 4 ネットワーク
- 5 ネットワーク
- 1 0 物理記憶装置
- 1 1 記憶装置制御装置。

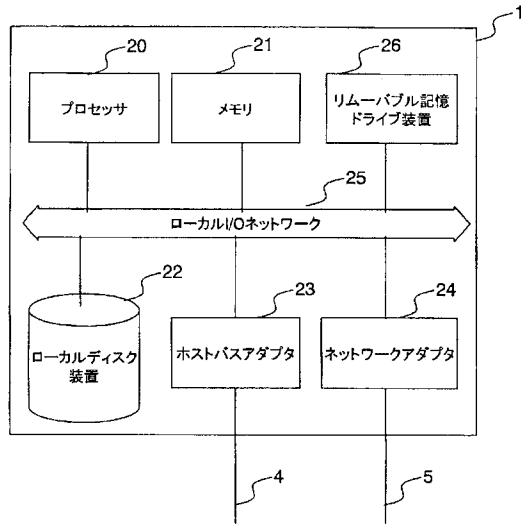
【図 1】



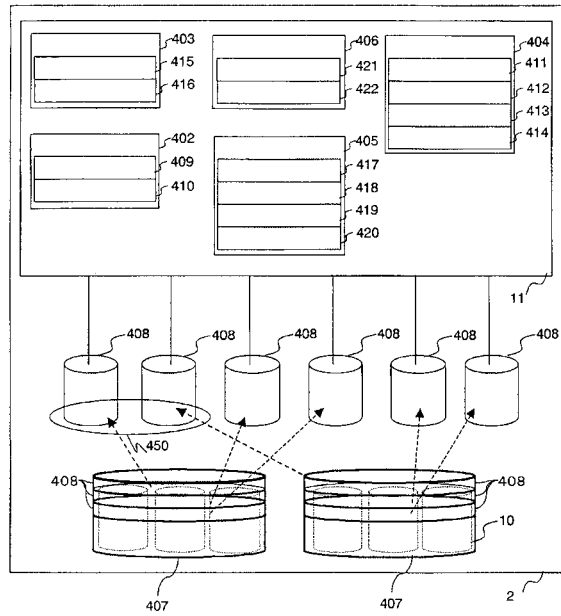
【図 2】



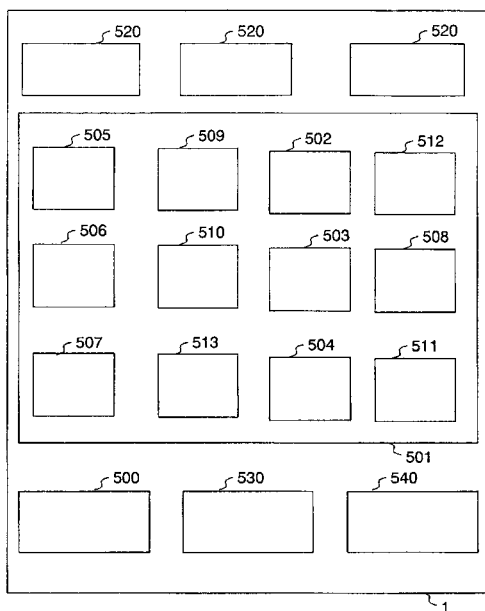
【図 3】



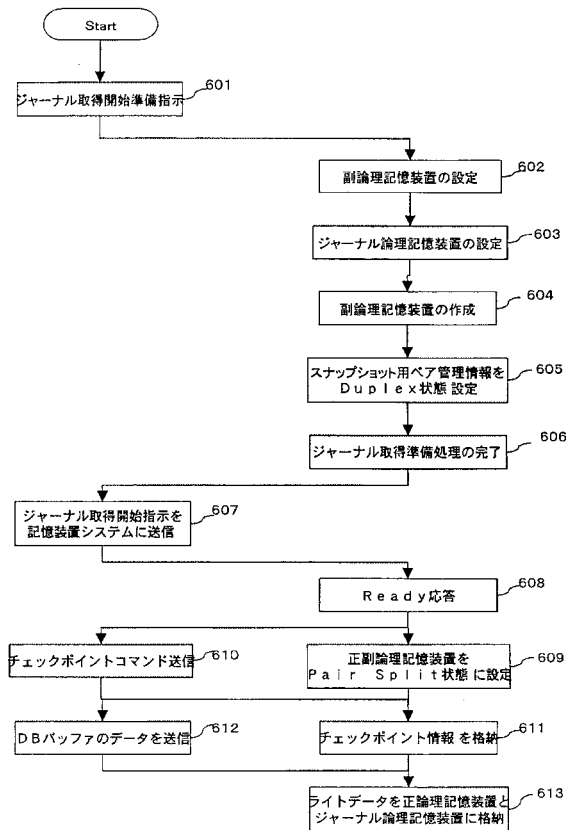
【図 4】



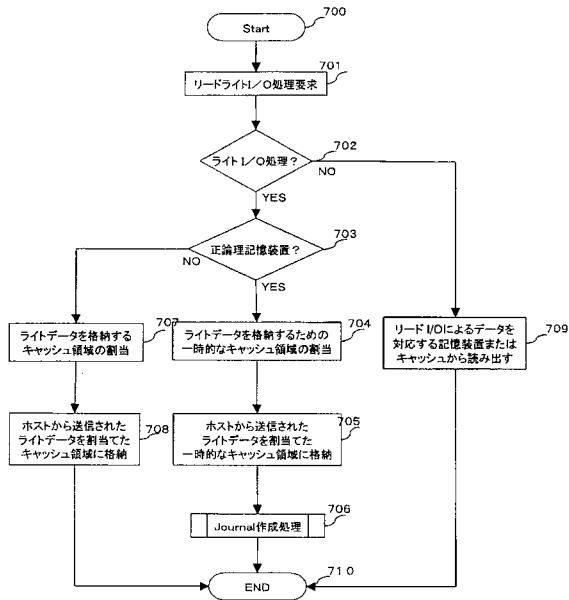
【図 5】



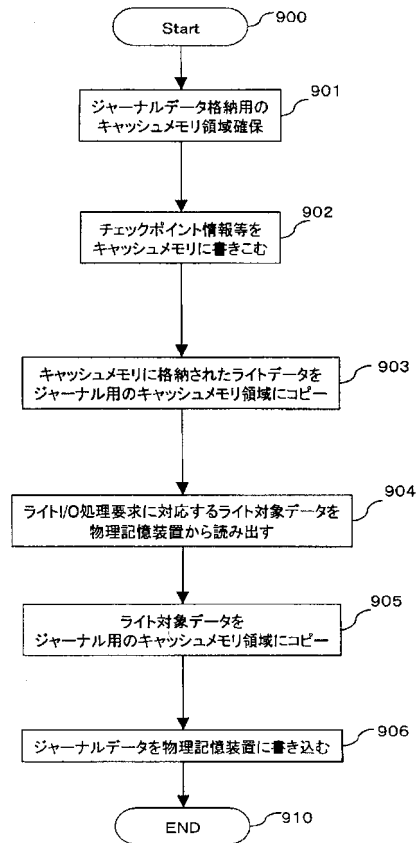
【図 6】



【図 7】



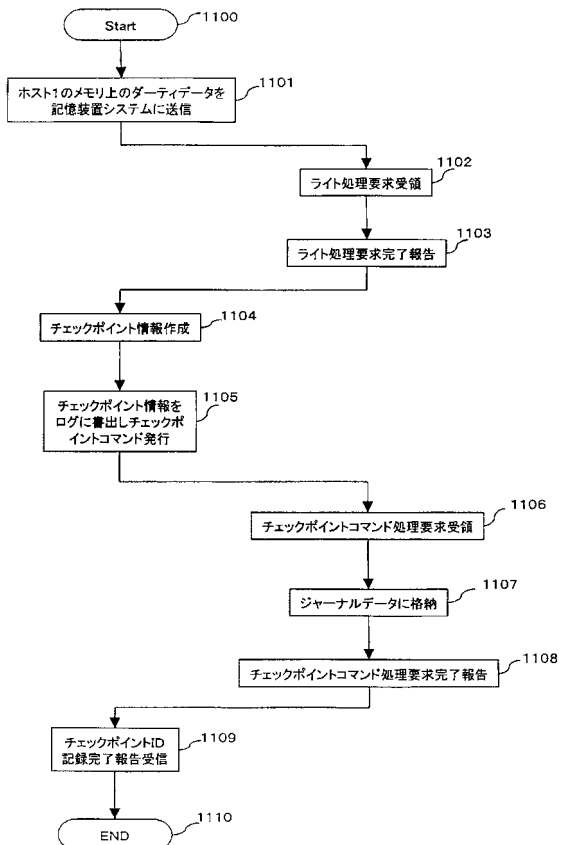
【図 8】



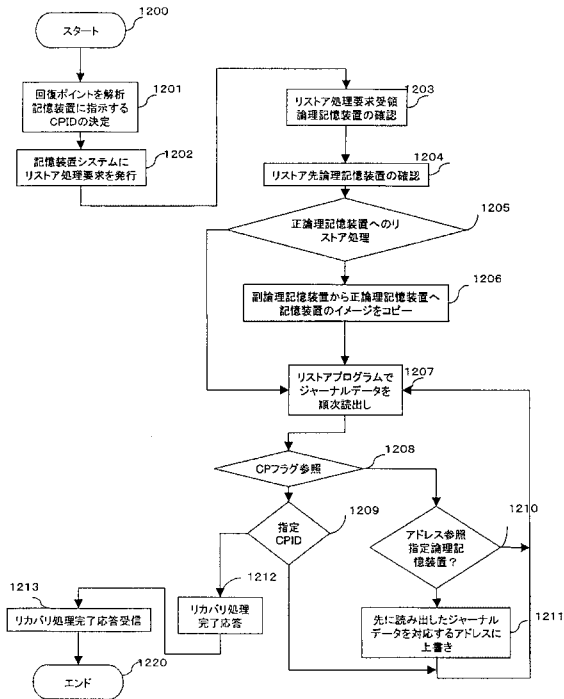
【図 9】

1001	1002	1003	1004	1005
チェックポイント 情報	アドレス (論理記憶装置番号、論理記憶装 置内アドレス)	データ長	ライト対象データ (Before)	ライトデータ (After)
チェックポイントフ ラグ	CPIID	処理順序番 号	処理時間	
1006	1007	1008	1009	

【図 10】



【図 1 1】



【図 1 2】

1301	1302	1303	1304	1305	1306	1307
1311	1312	1313	1314	1321	1322	1323
ホスト提供論理記憶装置番号	ホスト提供論理記憶装置内アドレス	記憶装置システム内論理記憶装置番号	論理記憶装置内アドレス	RAID Group番号	物理記憶装置番号	物理記憶装置内アドレス
Port#, LU#	0~512	Cu#, LDEV#	0~512	00	00	0~512
...

1331	1332	1333	1334	1335
ホスト提供論理記憶装置番号	記憶装置システム内論理記憶装置番号	RAID Group番号	Pair情報	ジャーナル対象モード
01234567.01	00.01	00.01	Pair Duplex	モードOFF
01234567.02	00.02	01.01	No Pair	モードOFF
01234567.03	00.03	02.01	Pair Symplex	モードON
...

1351	1352	1353	1354	1355
記憶装置システム内論理記憶装置番号	書き/リザープ情報	Path定義情報	Emulation Type, サイズ	障害情報
00.01	無し	パス定義済 接続Port #	OPEN 172G	無し
00.02	無し	パス定義済 接続Port #	OPEN 172G	無し
...
0A.01	スタンバイ用リザープ	無し	OPEN 36G	無し
0A.02	空き	無し	OPEN 36G	無し
...

【図 1 3】

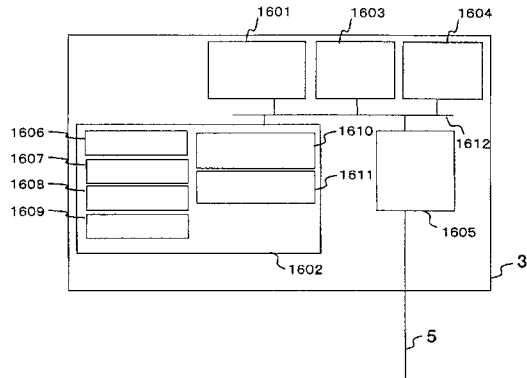
1401	1402	1403	1404	1405
ホスト提供論理記憶装置番号	記憶装置システム内論理記憶装置番号	Emulation Type サイズ	ペア状態	ペア管理情報
				1411 1412
				正 副
11234567.01	01.01	OPEN 9G	Pair Duplex	0 02.01
11234567.02	01.02	OPEN 9G	Pair Symplex	02.02 0
11234567.03	01.03	OPEN 16G	No Pair	-1 -1
11234567.04	01.04	OPEN 16G	Pair Resync	02.04 0
...
11234568.0A	0D.0A	3390 3G	Pair Create	0A.0A 0F.0A
...

【図 1 4】

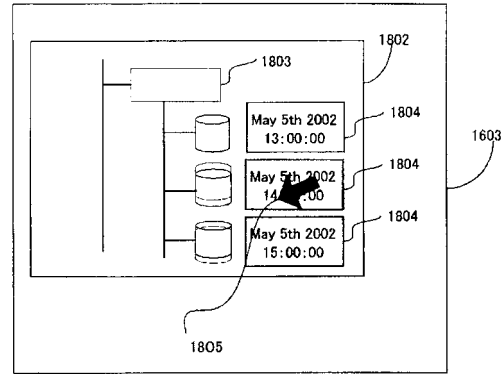
1501	1502	1503
チェックポイントID	アドレス	時間
ABCDEF12	12	2002.06.02 09:00:00
ABCDEF19	1A2B3C56	2002.06.02 21:00:00
...
0345FCA9	3ADE68B0	2002.06.12 09:00:00

1521	1522	1523
ジャーナルデータ格納論理記憶装置番号	チェックポイントID	アドレス
01.0F	ABCDEF12	0
01.0F	ABCDEF19	FEDCBA12
...

【図15】



【図17】

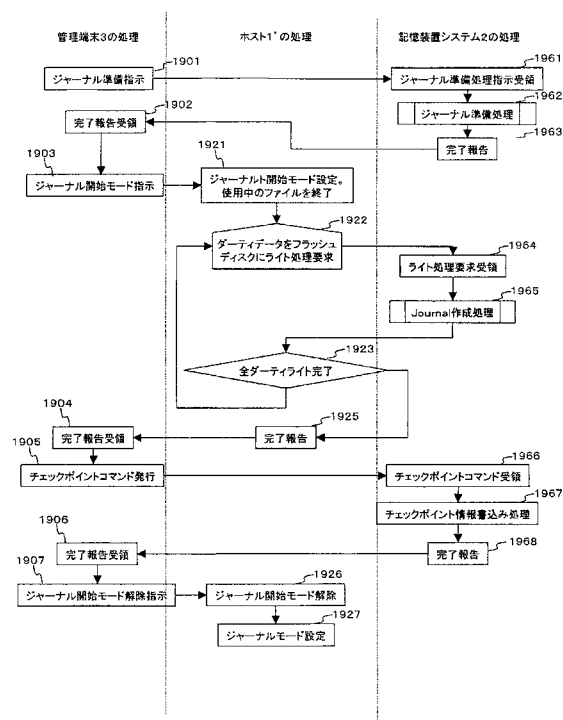


【図16】

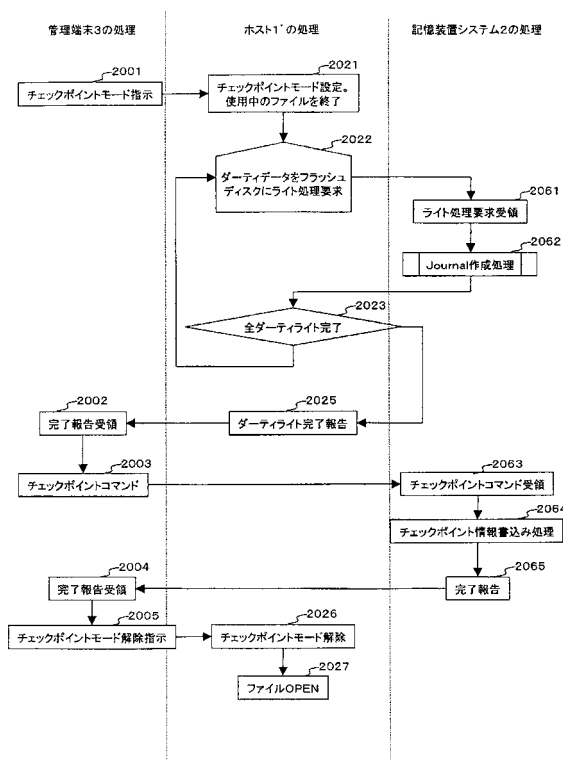
	1702	1701	1703
デバイス番号: 010F	チェックポイント識別子		時間
	I		2002.05.05 14:00:00
	J		2002.05.05 15:00:00

デバイス番号: 0B03	チェックポイント識別子		時間
	I		2002.05.05 14:00:00
	J		2002.05.05 15:00:00

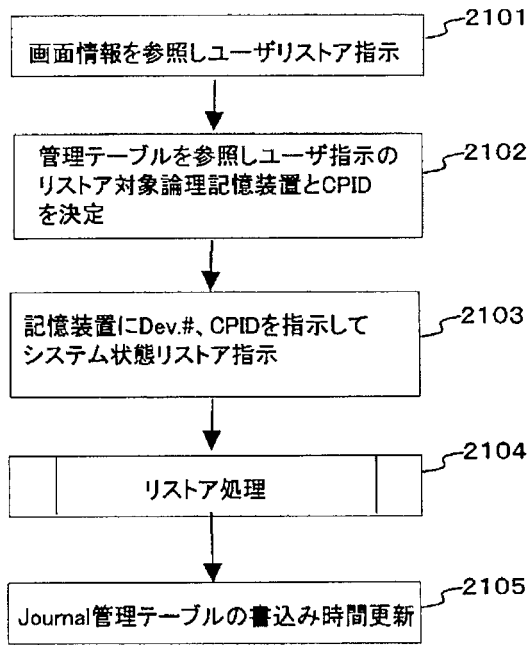
【図18】



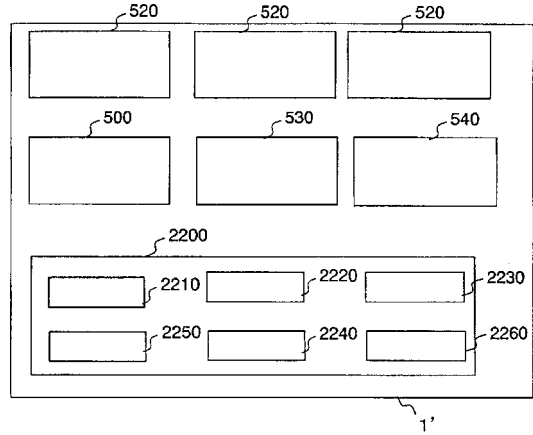
【図19】



【図 20】



【図 21】



フロントページの続き

(51)Int.Cl. F I
G 0 6 F 13/10 3 4 0 A

- (72)発明者 山本 康友
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
- (72)発明者 大枝 高
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
- (72)発明者 荒井 弘治
神奈川県小田原市中里 3 3 2 番地 2 号 株式会社 日立製作所 R A I D システム事業部内

審査官 高瀬 勤

- (56)参考文献 特開昭 5 8 - 1 0 1 3 5 6 (J P , A)
特開平 0 2 - 1 7 6 9 4 9 (J P , A)
特開 2 0 0 0 - 3 3 0 7 2 9 (J P , A)
特開 2 0 0 1 - 2 1 6 1 8 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
G 0 6 F 1 3 / 1 0
J S T P l u s (J D r e a m I I)