

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2004-13367

(P2004-13367A)

(43) 公開日 平成16年1月15日(2004.1.15)

(51) Int. Cl.<sup>7</sup>

G06F 3/06

G06F 12/00

F I

G06F 3/06

3 O 4 F

G06F 12/00

5 2 O E

G06F 12/00

5 4 5 M

テーマコード (参考)

5 B O 6 5

5 B O 8 2

審査請求 未請求 請求項の数 8 O L (全 16 頁)

(21) 出願番号 特願2002-163705 (P2002-163705)

(22) 出願日 平成14年6月5日(2002.6.5)

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(74) 代理人 100075096

弁理士 作田 康夫

(72) 発明者 大野 洋

神奈川県小田原市中里322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

(72) 発明者 占部 喜一郎

神奈川県小田原市中里322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

最終頁に続く

(54) 【発明の名称】 データ記憶サブシステム

(57) 【要約】

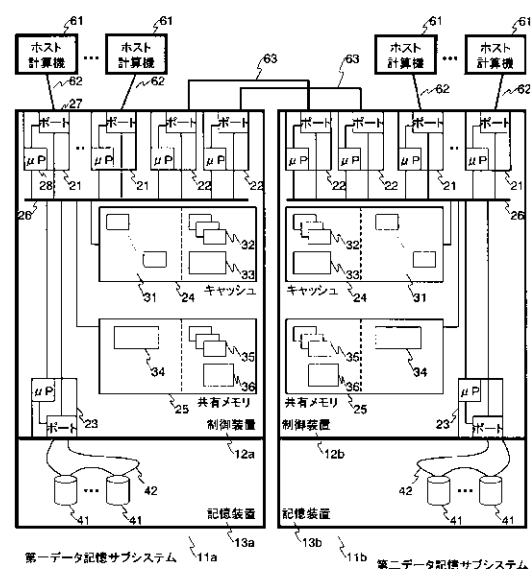
【課題】データ記憶サブシステム間のリモートコピーにおいて、コピー先で論理ボリュームに対する属性情報を使用するためにはホスト計算機間で別途属性情報の交換の通信を行わなければならなかった。また、属性情報の更新と通常のデータの更新の間の順序関係を維持することは困難であった。

【解決手段】データ記憶サブシステム間で、データのリモートコピー機能と同一のデータ転送パスを使用して、属性情報をコピーする。また、更新データおよび更新属性情報を合わせた通しの更新順序番号を付与してコピーを行い、受信側のデータ記憶サブシステムにおいてこの更新順序番号に従ってデータならびに属性情報を更新する。

ホスト計算機側では全く意識することなしに、属性情報がバックアップ側のデータ記憶サブシステムにコピーされ、またデータおよび属性情報の更新の順序を維持したままコピー先のデータ記憶サブシステム上でこれらの更新が反映される。

【選択図】 図1

【図1】 全体構成図



**【特許請求の範囲】****【請求項 1】**

2 以上のデータ記憶サブシステムが相互に接続された場合において、  
前記データ記憶サブシステムが有する、データのリモートコピー機能を用いて、  
各データ記憶サブシステムが、自己の属性情報を、他のデータ記憶サブシステムへ、コピーすることを特徴とするデータ記憶サブシステム。

**【請求項 2】**

磁気ディスク装置その他の記憶装置と、前記記憶装置に対して、データの入出力を行う制御装置とからなるデータ記憶サブシステムであって、  
10 ホスト計算機に接続された一のデータ記憶サブシステムを含む、複数のデータ記憶サブシステムが前記制御装置を介して接続された構成において、  
前記ホスト計算機から前記一のデータ記憶サブシステムに対し書込まれたデータを、別のデータ記憶サブシステムに対してコピーする際に、  
前記ホスト計算機から、前記一のデータ記憶サブシステムに対して発行されたデータ書込みコマンド以外の制御指示コマンドの内容に従い、当該データ記憶サブシステム内の制御装置が処理した結果発生する属性情報を、前記別のデータ記憶サブシステムに対してコピーすることを特徴とするデータ記憶サブシステム。

**【請求項 3】**

請求項 2 記載のデータ記憶サブシステムにおいて、  
前記属性情報のコピーは、前記データのコピーと同じ転送パスを用いて行なわれるデータ 20 記憶サブシステム。

**【請求項 4】**

請求項 2 記載のデータ記憶サブシステムにおいて、  
前記記憶装置は、一の論理ボリューム又は 2 以上の区切られた論理ボリュームを有し、  
前記ホスト計算機は、前記論理ボリュームを指定して、データの読み書きを行うことが可能であり、  
前記属性情報は、次のうち少なくともいづれか 1 の情報を含むデータ記憶サブシステム：  
1) 論理ボリュームを識別するために付与するタグ情報、  
2) 論理ボリュームへアクセス中のホスト計算機のホスト ID 情報、  
3) 論理ボリュームに対して特定のアクセス権を持つホスト計算機の情報、 30  
4) 論理ボリュームに対して特定のアクセス権を持つホスト計算機のポートの情報、  
5) 論理ボリュームのリザーブ制御情報、又は、  
6) 論理ボリュームへの更新を許可するキー情報。

**【請求項 5】**

請求項 2 記載のデータ記憶サブシステムにおいて、  
前記記憶装置は、一の論理ボリューム又は 2 以上の区切られた論理ボリュームを有し、  
前記ホスト計算機は、前記論理ボリュームを指定して、データの読み書きを行うことが可能であり、  
前記一のデータ記憶サブシステムは、自己に接続された前記ホスト計算機から、当該データ記憶サブシステムの論理ボリュームに対して発行された、データ書込みコマンド及び前 40  
記制御指示コマンドを送出し、  
前記別のデータ記憶サブシステムは、前記データ書込みコマンドに対応するデータ及び前記制御指示コマンドに対応する属性情報を、これらの順序性を保って、当該別のデータ記憶サブシステムの論理ボリュームに対して格納するデータ記憶サブシステム。

**【請求項 6】**

請求項 2 記載のデータ記憶サブシステムにおいて、  
前記記憶装置は、一の論理ボリューム又は 2 以上の区切られた論理ボリュームを有し、  
前記ホスト計算機は、前記論理ボリュームを指定して、データの読み書きを行うことが可能であり、  
前記一のデータ記憶サブシステムは、自己が接続されている前記ホスト計算機上で動作す 50

るプログラムが、前記論理ボリュームに対して、データ以外の任意の情報を付与する機能と、前記プログラムからの指示で付与された前記データ以外の任意の情報を格納する機能とを有するデータ記憶サブシステム。

【請求項 7】

請求項 2 記載のデータ記憶サブシステムにおいて、前記制御指示コマンドに代わり、前記データ書込みコマンドが一部に有する制御指示情報が用いられるデータ記憶サブシステム。

【請求項 8】

請求項 2 記載のデータ記憶サブシステムにおいて、前記記憶装置は、一の論理ボリューム又は 2 以上の区切られた論理ボリュームを有し、前記一のデータ記憶サブシステムは、前記制御指示コマンドを、特定の論理ボリュームに対する書込みデータとして扱う機能を有し、前記ホスト計算機からの制御指示コマンドに代わり、前記特定の論理ボリュームに書き込まれたデータを用いるデータ記憶サブシステム。

10

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、ホスト計算機のデータを格納する外部記憶装置群であって、相互に附近地又は遠隔地に設置された複数の外部記憶装置の間で、情報の複製を行なう技術に係り、特に、各外部記憶装置の属性情報の取扱いに関する。

20

【0002】

【従来の技術】

近年の計算機システムでは、使用されるデータの容量が指数関数的に増大し、データの更新頻度も高くなっている。このデータをどのようにバックアップし、また障害発生時にどれだけ迅速に正常稼働状態に復旧できるかが、データストレージ技術における重要課題となっている。この課題に対する一つの解として、磁気ディスクアレイを搭載したデータ記憶サブシステム（外部記憶装置）を遠隔地に複数台設置して、それらの間を通信パスで接続し、一方のデータ記憶サブシステム上で更新されたデータを、ホスト計算機を経由せずに、自動的に他のデータ記憶サブシステムに常時コピーする、リモートコピー技術がある。

30

【0003】

このリモートコピー技術では、コピー先データ記憶サブシステムにもホスト計算機を接続することにより、第一のデータ記憶サブシステムが使用不能となる故障や災害が発生した場合に、コピー先のデータ記憶サブシステム上にコピーされているデータを使用して業務を引き継ぐディザスタリカバリが実現できる。この引き継ぎの際、データの整合性が重要になるが、コピーの中断がどの時点で発生するかは予測不能である。このため、当該データを使用するアプリケーションが整合性を回復できるようにするため、コピー先のデータ記憶サブシステム上のデータの更新順序を、コピー元のデータ記憶サブシステムに対して行われた更新順序と同一に維持することが必要になる。

【0004】

このようなデータ更新順序維持を実現する従来技術として、特開平 6 - 290125 号公報記載の技術があげられる。また、このようなデータ更新順序維持を、ホスト計算機の介在無く、かつ各データのコピー完了確認を待たない非同期方式で実現する従来技術として、特開平 11 - 85408 号公報記載の技術が挙げられる。

40

【0005】

一方、データ記憶サブシステムに接続されるホスト計算機上で動作するプログラムは、データの取扱い方法が高機能化しており、データ記憶サブシステムに対しても、単なるデータの格納以外に付加的な制御を要求するようになってきている。このような制御は、ホスト計算機がデータ記憶サブシステムにアクセスする際のデータ格納領域の単位となる、論理ボリュームを対象として行われるのが一般的である。

50

## 【 0 0 0 6 】

付加的な制御の例として、複数のホスト計算機からアクセスされる可能性のある論理ボリュームにおけるアクセス制御があげられる。これは当該論理ボリュームに対する参照や更新といった操作の許可・不許可を、ホスト計算機単位あるいはプログラムのインスタンス単位に行い、なおかつ動的に設定変更が可能であることが望ましい。アクセス制御をホスト計算機単位で行う場合、データ記憶サブシステム側では、ホスト計算機のIDやホスト計算機に接続されたポートのIDによりホスト計算機を識別し、アクセス制御を行う。また、最も簡単なホスト計算機単位のアクセス制御としては、単一のホスト計算機にのみアクセスを許可するリザーブ機能がある。一方、プログラムのインスタンス単位で行う場合、アクセスを許可するためのキー値を登録する機能の提供と、その後のデータ読み書き等のアクセスコマンドに当該キー値が付加されている要求のみアクセスを許可する制御方法で実現できる。

10

## 【 0 0 0 7 】

別な制御の例として、論理ボリュームの状態等を報告する機能がある。例えば、ある時点で当該論理ボリュームを使用中のホスト計算機のID情報が登録されていて、別なホスト計算機からの要求に応じてその情報を返すことで、後からアクセスしたホスト計算機が当該論理ボリュームの使用状況を確認でき、その後の動作の決定に利用できる。また、論理ボリュームを認識するのに、データ記憶サブシステム固有のLUN ( Logical Unit Number ) の代わりに、接続されたホスト計算機側で都合の良い適当なタグ情報を付与して利用することがある。

20

## 【 0 0 0 8 】

一方で、データ記憶サブシステムが提供する制御機能ではなく、プログラム側での独自の処理を行う場合がある。たとえば、複数のホスト計算機上に分散した、プログラムの複数インスタンスが連携して動作する場合に、共有するデータ記憶サブシステム経由で、連携に必要な情報を交換することがある。この場合、データ記憶サブシステムの論理ボリュームへデータを書込む代わりに、データ記憶サブシステム上の作業メモリ領域上に必要な情報を読み書きする方法があり、これにより物理的なディスクへの書込み処理が無い高速なデータ交換を実現することができる。この場合、当該データ記憶サブシステムは、任意の作業用情報の読み書きを行うという制御コマンドを提供することになる。

## 【 0 0 0 9 】

これらのデータ記憶サブシステムの制御処理で使用するあるいは処理の結果発生する情報は、データ記憶サブシステム内の処理プロセッサ間で共有されるメモリ領域に保持される。以下、本明細書中では、これらの情報のことを「属性情報」と呼ぶ。属性情報は、前記の各例の場合、アクセス設定のリスト、リザーブ制御情報、アクセス許可のキー値、使用中ホストのホストID、論理ボリュームのタグ情報、プログラムが読み書きする任意の作業用情報そのもの、などとなる。

30

## 【 0 0 1 0 】

これらの付加的な制御は各データ記憶サブシステムの中で閉じたものであり、この結果発生する属性情報はリモートコピーの対象外であった。しかし、リモートコピーのコピー先でディザスタリカバリ等の引継ぎを実行するときに、より正確にデータを復元するためには、引継ぎ直前時点でのデータ記憶サブシステム上の属性情報も利用できた方が好ましい。

40

## 【 0 0 1 1 】

従来このような属性情報の利用を行う際には、各データ記憶サブシステムに接続されたホスト計算機間で通信を行うことにより、これらの情報を交換していた。例えば、特開平11-305947では、ホスト計算機から発するアテンション報告指示情報を、磁気ディスク制御装置が受け、これをリモートの磁気ディスク制御装置へ転送し、これを受け取ったリモートの磁気ディスク制御装置が、リモートのホスト計算機へ通知する技術が開示されている。

## 【 0 0 1 2 】

50

**【発明が解決しようとする課題】**

従来技術では、属性情報を使用するためには、ホスト計算機間で、別途、属性情報の交換の通信を行わなければならない、またその実装にはリモートコピーの状態を考慮することが必要となり、処理が複雑になる。本発明の実施例では、ホスト計算機側では全く意識することなしに、属性情報をバックアップ側のデータ記憶サブシステムにもコピーし、バックアップ側で自動的に利用可能としている。

**【0013】**

また従来技術では、属性情報の更新はホスト計算機間で、データの更新はデータ記憶サブシステム間で各々個別に行われるため、属性情報の更新とデータの更新の間の順序関係を維持することは、基本的に不可能であった。ホスト計算機上のプログラムでデータ記憶サブシステムと密に連携して管理することで実現は可能だが、このような連携は実装が難しいだけでなく、オーバーヘッドが大きい性能も出ない。本発明の実施例では、ホスト計算機側では全く意識することなしに、データおよび属性情報の更新の順序を維持したまま、しかも高速に、コピー先のデータ記憶サブシステム上にこれらの更新が反映される機能を実現している。

**【0014】**

上述した、属性情報とリモートコピーの状態とを完全に連携させるためには、リモートコピー中断後の再開時やコピー方向反転時での属性情報の再同期の方法も解決されなければならない。本発明の実施例では、ホスト計算機側では全く意識することなしに、リモートコピー中断後の再開時やコピー方向反転時での属性情報の差分コピーを実行できるようにし、短時間で属性情報の再同期を終了している。

**【0015】**

さらに上述した機能を、複数のデータ記憶サブシステムにまたがるリモートコピーに適用して、耐障害性をより高くしたいというニーズがある。本発明の実施例では、複数のデータ記憶サブシステムにまたがる属性情報のコピーを実現している。

**【0016】****【課題を解決するための手段】**

データ記憶サブシステム間で、データのリモートコピー機能と同一のデータ転送パスを使用して、属性情報をコピーする機能を設ける。

**【0017】**

更新データおよび更新属性情報を合わせた通しの更新順序番号を付与してコピーを行い、受信側のデータ記憶サブシステムにおいてこの更新順序番号に従って、データならびに属性情報を更新する機能を設ける。

**【0018】**

リモートコピーの中断中に更新された属性情報を、コピー元・コピー先両方のデータ記憶サブシステムにおいて記憶し、リモートコピーの再開後に更新部分のみをコピーして、再同期する機能を設ける。

**【0019】**

コピー先となるデータ記憶サブシステムの数だけコピー元で属性情報更新情報を複製し、各コピー先に別々に送信する機能を設ける。また1次コピー先で受け取った属性更新情報を、2次コピー先にリモートコピー機能により送信する機能を設ける。

**【0020】****【発明の実施の形態】****<はじめに>**

図9により、本発明の前提となる同期転送方式のリモートコピー、および非同期転送方式のリモートコピーについて説明する。

**【0021】**

同期転送方式によるリモートコピーでは、第一データ記憶サブシステム11aは、ホスト計算機61から更新データブロックを受領し、そのデータブロックを第二データ記憶サブシステム11bに転送する。11bは、当該データブロックの受領を11aに通知する。

11aは、当該通知受領後ホスト計算機に対し更新データブロックの書込み完了通知を行う。以上の処理のいずれかを失敗した場合には、11aは、ホスト計算機に書込み障害として報告する。

【0022】

非同期転送方式によるリモートコピーでは、第一データ記憶サブシステム11aは、ホスト計算機61から更新データブロックを受領し、その時点で当該ホスト計算機に更新データブロックの書込み完了通知を行う。11aは、自身のスケジュールで、ホスト計算機側の処理と非同期に、第二データ記憶サブシステム11bへ当該データブロックを転送する。遠隔地へのデータ経路の複雑化、中途経路のボトルネック化により、データ転送中の当該データの順序性は保証されない。

10

【0023】

図9では、例えば、データ転送経路上の順序は#1、#4、#3、#2である。更新順序の保存は、11bにおいて、11bのソート機能により、第一データ記憶サブシステムがホストからデータを受領した順序を保って、更新することにより達成される。つまり、11bにおいては、#1、#2、#3、#4の順序であり、この更新処理の途中に、第一データ記憶サブシステムやデータ経路に不慮の災害が発生しても、第二データ記憶サブシステムにおけるデータ更新順序は守られている。このため、11bに接続されるホスト計算機61では、データベースやジャーナルファイルシステムなどは、矛盾なく回復処理を行なうことができる。非同期転送方式はホスト計算機の高性能処理、データ記憶サブシステム間の距離拡大、という特長を持ち、かつ順序性を保証したりリモートサイトへの書込み制御により任意時点のデータベースやジャーナルファイルシステムの整合性を確保できる特長を有す。コピー元からの送信順序が任意であっても、リモート側で一旦、受信し、格納した上で、対応する情報が揃ったところからリモート側の更新処理を行なうためである。

20

【0024】

<実施例1>

図1は第1の実施例の全体構成を示している。二つのデータ記憶サブシステム(以下、単に適宜「サブシステム」と略記する)11a、11bが、データ転送パス63を介して接続され、互いに附近地又は遠隔地に設置されている。

【0025】

各サブシステムは、ホストアクセスパス62を介して、複数のホスト計算機61と接続されている。サブシステムは、データの読み書き要求の処理や、更新データの格納を制御する制御装置12(以下、適宜、12a、12bの総称として記す。以下同様)と、実際のデータを最終的に格納する記録媒体を含む記憶装置13の二つの部分から構成される。

30

【0026】

制御装置12は、ホスト計算機と接続されるチャンネルアダプタ21、他のサブシステムと接続されるチャンネルアダプタ22、記憶装置と接続されるディスクアダプタ23を持つ。各アダプタは、実際のデータ送受を行なうポート27と、データ送受を制御するマイクロプロセッサ28から構成される。

【0027】

図1では、一つのアダプタに一つのポートと一つのマイクロプロセッサが搭載されているが、いずれも個数は任意であり、個数が増えるほど処理能力が高まることになる。

40

【0028】

また制御装置12には、記憶装置13から取り出した、あるいは記憶装置13にこれから格納すべきデータを保持するキャッシュメモリ24(以下、適宜「キャッシュ」と略記する)と、各アダプタのマイクロプロセッサが処理を実行するためのワーキングエリア等として使用される共有メモリ25がある。各アダプタは、キャッシュおよび共有メモリとバス26で接続されていて、これらのメモリ領域とデータをやり取りすることができる。

【0029】

記憶装置13は、複数の磁気ディスクドライブ41を持つ。これらの磁気ディスクドライブは、ディスク入出力バス42により制御装置12側のディスクアダプタ23に接続され

50

る。

#### 【0030】

サブシステムは、ホスト計算機に対して論理ボリュームと呼ばれるデータ記憶領域が存在するように動作し、ホスト計算機側では対象となる論理ボリュームを指定して、データの読み書きを行なう。この論理ボリュームの識別子のことをLUN (Logical Unit Number) と呼ぶ。制御装置12では、記憶装置41内の磁気ディスクドライブの記憶領域を分割して、各論理ボリュームに対応づけ、最終的にこれらの対応領域にデータを格納する。

#### 【0031】

< 一般的な処理 >

制御装置12の処理について簡単に説明する。ここでは、第一サブシステム側から第二サブシステム側に対して、リモートコピーを行なっている場合を説明する。ただし、リモートコピーの向きは論理ボリューム毎に設定・変更でき、ある時点で両方向のリモートコピーが共存することもある。

#### 【0032】

ホスト計算機61からデータ書込み要求を受け取ったチャンネルアダプタ21は、キャッシュ24上の対応領域にキャッシュエントリ31としてデータを格納する。これらのキャッシュエントリは、ディスクアダプタ23上のマイクロプロセッサが、他のマイクロプロセッサの処理とは独立したスケジューリングにより、対応する磁気ディスクドライブの領域にデータを書込み、データを最終的に保存する。

#### 【0033】

同期転送方式によるリモートコピーを行なっている場合は、書込み処理を行なったチャンネルアダプタ21から、第二サブシステムに接続されるチャンネルアダプタ22に転送要求が出され、前記で説明した同期転送が実行される。この際、第二サブシステム側のチャンネルアダプタ22は、転送されたデータ内容により、キャッシュ上の該当キャッシュエントリ31を更新していく。

#### 【0034】

非同期転送方式によりリモートコピーを行なっている場合は、第一サブシステム側のチャンネルアダプタ21上のマイクロプロセッサは、更新データ内容をキャッシュエントリ31以外にデータ更新情報32としてキャッシュ24上の別領域に保存し、この時点で書込み要求をしたホスト計算機61には処理完了を報告する。第二サブシステムに接続されるチャンネルアダプタ22上のマイクロプロセッサは、他のマイクロプロセッサと独立したスケジューリングにより、データ更新情報32を第二サブシステムに、ホスト計算機からの書込み要求順序の発行順序番号を付加して、転送する。

#### 【0035】

第二サブシステム側のチャンネルアダプタ22は、転送されたデータをキャッシュ31上の別領域に更新データ32として一旦保存し、さらに別なスケジューリングの処理により、発行順序番号が揃った部分について、当該発行順序番号の順で、キャッシュエントリ31の更新を行なっていく。

#### 【0036】

いずれの方式のリモートコピーでも、第二サブシステム側でのキャッシュエントリ31は、ディスクアダプタ23上のマイクロプロセッサが、他のマイクロプロセッサとは独立したスケジューリングにより、対応する磁気ディスクドライブの領域にデータを書込み、データを最終的に保存する。

#### 【0037】

以上の説明で明らかなように、いずれのサブシステムでも、磁気ディスクドライブへの書込みは、対ホスト計算機の処理やリモートコピー処理とは独立して動作する。このため、以下、本明細書中では、ディスクアダプタ23ならびに記憶装置13の構成要素に関する議論は行なわない。

#### 【0038】

10

20

30

40

50

一方、ホスト計算機 61 は、サブシステムに対して、磁気ディスク装置に記憶するためのデータの読み書き以外の処理要求を行う場合がある。これらの要求はチャネルアダプタ上のマイクロプロセッサが、必要に応じて他のマイクロプロセッサと協調しながら、処理を行なう。処理の結果は、一般に論理ボリューム毎に管理される属性情報として、共有メモリ上の属性情報テーブル 34 上に格納される。このようなデータ以外の処理要求としては、従来の技術の項で述べたように、論理ボリュームのリザーブ設定などがあげられる。

#### 【0039】

< 本発明における処理 >

本発明に係る制御装置固有の処理について説明する。本発明による制御装置では、前記のホスト計算機からの書込みデータの他に、前記属性情報をもリモートコピー機能の一部として転送し、異なるサブシステム間で属性情報の常時一致化を行なうものである。ホスト計算機は、属性情報を発生させるような要求について、通常のデータの書込み要求を含めて、途中で打ち切らずに終わりまで、要求発行順にサブシステム側で処理されることを期待する。このため、リモートコピーにおいても、属性情報を発生させるような要求と、通常のデータの書込み要求とを通した、要求発行順に対応データ又は属性情報が更新されるように構成する。

#### 【0040】

以後の説明では、ホスト計算機から行われる通常のデータの書込み要求のことを「データ更新」、属性情報を発生させるような、データ書込み以外の要求のことを「制御指示」と呼ぶこととする。また、リモートコピーのコピー元（以下の説明の大部分では第一サブシステム側）を「正側」、コピー先（以下の説明の大部分では第二サブシステム側）を「副側」と呼ぶ場合がある。

#### 【0041】

図 2 に属性情報テーブル 34 の構成を示す。テーブルは、論理ボリューム毎の行に、当該論理ボリュームの属性情報が格納される構成となっている。一番左の列に論理ボリュームの ID である LUN が格納され、続いて、リザーブ状態を表すフラグ、アクセスを許可するキー、使用中のホスト ID、その他の属性情報を格納する列がある。

#### 【0042】

表の一部の領域は、ユーザ用領域の列として予約され、サブシステムに接続されたホスト計算機上のプログラムが、専用のアプリケーションインタフェース（以下、「API」と略記する。）を使って、当該論理ボリュームに関して設定・記憶させたい任意の属性情報を格納することができる。これにより、サブシステムの制御装置ではサポートしていない、プログラム独自の競合管理等を実現することが容易になる。

#### 【0043】

図 3 に、同期転送によるリモートコピー構成での、ホスト計算機 61 からの要求に対する処理フローを示す。各要求は、ホスト計算機 61 から第一サブシステム 11a のチャネルアダプタ 21 に対して、SCSI プロトコルのコマンドとして発行される。チャネルアダプタのマイクロプロセッサ 28（図 1）は、コマンドの種別を判断する。データ更新コマンドであれば正側キャッシュ 24 上のキャッシュエントリ 31 として書込み（S11）、制御指示コマンドであれば、制御指示の内容を実行し（S16）、その結果発生する論理ボリュームに対する属性情報は、共有メモリ 25 上の属性情報テーブル 34 に格納される（S17）。続いて、チャネルアダプタ 21 からの要求により、チャネルアダプタ 22 上のマイクロプロセッサ 28 は、データ転送パス 63 を経由して、副側のチャネルアダプタ 22 に更新情報を転送する（S12）。副側のチャネルアダプタ 22 上のマイクロプロセッサ 28 は、更新情報にあるフラグから当該更新内容がデータか属性情報か判断できるので、データならば対応する副側キャッシュ 24 上のキャッシュエントリ 31 として書込み（S13）、属性情報ならば副側共有メモリ 25 上の属性情報テーブル 34 に格納される（S18）。これらの処理が完了すると、副側のチャネルアダプタから正側のチャネルアダプタに転送完了通知を行い（S14）、正側のチャネルアダプタは、要求を発行したホスト計算機に完了通知を行なう。ここで説明した一連の処理は、ある論理ボリュームに対

10

20

30

40

50



し同時には一つしか実行されないように管理され、後着の処理は保留ないしエラーとなる。これにより、副側でも、データならびに属性情報の更新が、両者を通じた要求発行順に更新される。

#### 【 0 0 4 4 】

##### < 更新情報の転送 >

更新情報の転送について、さらに説明する。本実施例では、データ転送パス 6 3 上でも、S C S I プロトコルを使用する。正側チャネルアダプタ 2 1 上のポート 2 7 がイニシエータとなり、副側のチャネルアダプタ 2 1 上のポート 2 7 に対して、更新情報を送信する。

#### 【 0 0 4 5 】

図 5 には、データの更新情報を送信する場合の S C S I コマンドの構造を示す。オペレーションコードにはリモートコピーライトがセットされ、その他に、宛先の論理ボリュームを示すターゲット L U N 番号、更新データの論理ボリューム内データ格納先頭アドレス、転送長（データサイズ）が格納される。同期転送の場合、順序番号は使用しない。本実施例では、更新情報としてデータ以外に属性情報も送られることになるので、これを区別するためのフラグが追加されており、ここに「データ」を示すフラグを設定する。このコマンドが副側で受け付けられると、続くデータ転送フェーズでデータ本体が送られることになる。

#### 【 0 0 4 6 】

図 6 には、属性情報の更新情報を送信する場合の S C S I コマンドの構造を示す。基本的にはデータの更新情報と同様の構造を持ち、オペレーションコードにはリモートコピーライトがセットされ、その他に、宛先の論理ボリュームを示すターゲット L U N 番号、更新情報の属性情報テーブル内アドレス、転送長（属性情報サイズ）が格納される。同期転送の場合、順序番号は使用しない。前記のフラグに「属性情報」を示すフラグを設定する。このコマンドが副側で受け付けられると、続くデータ転送フェーズで属性情報本体が送られることになる。なお、属性情報は一般に容量が小さいので、S C S I コマンドの特定の領域（ベンダ固有領域やアドレスフィールドの一部など）に更新情報本体も格納することも可能である。また、テーブル内アドレスと転送長の組み合わせの代わりに属性情報の種別を示す I D を使うことも可能である。

#### 【 0 0 4 7 】

図 4 に、非同期転送によるリモートコピー構成での、ホスト計算機 6 1 からの要求に対する処理フローを示す。同期転送の場合と同様に、第一サブシステム 1 1 a のチャネルアダプタ 2 1（図 1）のマイクロプロセッサ 2 8 は、コマンドの種別を判断し、データ更新コマンドであれば正側キャッシュ 2 4 上のキャッシュエントリ 3 1 として書込み（S 2 1）、制御指示コマンドであれば制御指示の内容を実行して（S 2 4）発生する属性情報を正側共有メモリ 2 5 上に格納する（S 2 5）。さらにデータの更新の場合は正側キャッシュ 2 4 上の別領域にデータ更新情報 3 2 を書込み（S 2 2）、制御指示の場合には発生した属性情報の更新内容に関する情報を正側共有メモリ 2 5 上に属性情報更新情報 3 5 として書込む（S 2 6）。S 2 2、S 2 6 で付けられる順序番号は、データの更新情報、属性情報の更新情報の双方を合わせた通し番号とする。論理ボリューム毎、または設定された論理ボリュームグループ毎に、通し番号を設定する。

#### 【 0 0 4 8 】

ここで各更新情報の格納内容は、図 5、図 6 で示した更新情報の転送コマンドに含まれる内容と更新内容本体を合せたものである。各更新情報を格納する際、順序番号として、当該論理ボリュームに対するデータ更新および属性更新の両者を併せた通しでの処理順序を示す連続番号を書込む。また、論理ボリューム毎に順序番号を付与する代わりに、あらかじめ正側・副側で同一に設定された論理ボリュームを複数含む「論理ボリュームグループ」の中で、当該グループに対する更新処理の番号を付与することもできる。これにより複数論理ボリュームを併用するプログラムに対しても、リモートコピーの更新の順序性を保証する。

#### 【 0 0 4 9 】

10

20

30

40

50

順序性を保証するとは、送信されたデータや属性情報の全てが、順序正しく受信側でソート処理等により復元されることを意味する。言い換えれば、ホストから要求された順序を維持しつつ、コピー先のデータ及び属性情報が更新されていくことをいう。つまり、何らかの原因でデータや属性情報のコピーが中断した場合であっても、コピー元で、ある特定時点までに更新されたデータや属性情報の全てが欠落なくコピー先に反映されていて、かつ、コピー元で当該時点以降に更新されたデータ又は属性情報はコピー先では一切反映されていないことを意味している。伝送の途中で、データ又は属性情報が喪失した場合には、その伝送（コピー動作）をエラーと見なして処理を中断するか、何らかの手段で再送して復旧させる。

#### 【 0 0 5 0 】

続いて、チャンネルアダプタ 2 2 上のマイクロプロセッサ 2 8 は、独自のスケジューリングにより動作し、キャッシュ上のデータ更新情報 3 2 および共有メモリ上の属性情報更新情報 3 5 を、任意の順序で副側のチャンネルアダプタ 2 2 に転送する（S 2 7）。副側のチャンネルアダプタ 2 2 上のマイクロプロセッサ 2 8 は、受信した更新情報を、データならば副側キャッシュ 2 4 上のデータ更新情報 3 2 として、属性情報ならば副側共有メモリ 2 5 上の属性情報更新情報 3 5 として、一旦、書込む（S 2 8）。

#### 【 0 0 5 1 】

副側チャンネルアダプタ上のマイクロプロセッサは、さらに別な独自のスケジューリングで、前記の各更新情報を論理ボリュームないし論理ボリュームグループ内での順序番号によりソートし（S 2 9）、順序番号の不連続が無くなった部分から、同期処理での処理（S 1 3 ないし S 1 8）と同様の更新処理を行う（S 3 0）。この処理により、副側でも、論理ボリュームないしは論理ボリュームグループ内で、データならびに属性情報の更新が、両者を通した要求発行順で更新されることとなる。

#### 【 0 0 5 2 】

なお、一つの制御指示の結果、複数の属性情報を生成する場合もあり得る。この場合、同期転送であれば全ての属性情報の転送処理（図 3：S 1 7、S 1 2、S 1 8 の繰り返し）を完了させ、非同期転送であれば全ての属性情報の格納処理（図 4：S 2 5、S 2 6 の繰り返し）を完了させてから、ホスト計算機に完了通知を行う。

#### 【 0 0 5 3 】

ここまでの説明では、ホスト計算機からのデータ更新ならびに制御指示は、各々、単一のコマンドに単一の要求が含まれることを前提とした。その他に、単一のコマンドに複数の制御指示が含まれていたり、データ更新コマンドに制御指示が含まれていたりという場合もあり得る。この場合、同期転送であれば、全ての制御指示に関する処理（S 1 6、S 1 7、S 1 2、S 1 8、S 1 4 の繰り返し）を行なった後、データ更新の処理（S 1 1、S 1 2、S 1 3、S 1 4）を行い、ホスト計算機に完了通知を行う（S 1 5）。非同期転送であれば、全ての制御指示に関する処理（S 2 4、S 2 5、S 2 6 の繰り返し）を行なった後、データ更新の処理（S 2 1、S 2 2）を行い、ホスト計算機に完了通知を行う（S 2 3）。

#### 【 0 0 5 4 】

##### < 中断と逆転 >

図 7 を使用して、リモートコピーの中断と逆転に関する処理について説明する。リモートコピーのコピー状態は、論理ボリューム毎に管理されている。「正常」状態とは、正副で同期が取れている、すなわち非更新データは全て一致していて更新データのみコピーが実行されている状態を指す（図 7（a））。この状態では、データの一貫性を保つため、副側のホスト計算機は、当該論理ボリュームの内容を更新できない。

#### 【 0 0 5 5 】

正常状態にあるとき、プログラム等からの指示により、リモートコピーを中断させることができる（図 7（b））。この「中断」状態では、副側のホスト計算機からも当該論理ボリュームの内容を更新可能であり、正側でプログラムは動作したまま、副側でバックアップを取得したり、開発中のプログラムを実行テストしたり、といったことが可能である。

10

20

30

40

50

中断状態からは、プログラム等からの指示により、同じコピーの向きでリモートコピーを再開させることと（図7（c））、逆向きでリモートコピーを再開させるテイクオーバー（図7（d））が可能である。

【0056】

いずれの場合も、両サブシステム上で「中断」状態の間に更新されたデータないし属性情報の位置をそれぞれ記録しており（図7（b）、33、36）、いずれかのサブシステムにおいて更新されたデータないし属性情報についてのみ正側の現在値を副側にコピーすることにより、当該論理ボリュームのコピー状態を「正常」状態に復旧させる。この動作を再同期と呼び、コピー内容を更新情報に限定することで、「正常」状態復旧までの時間が短縮される。再同期が完了するまでは、副側の論理ボリュームは整合性はないため、更新順序に意味は無く、順序は無関係にコピーを行う。

10

【0057】

「中断」状態の間の更新位置の記録について、さらに説明する。キャッシュ24上にデータの更新位置を記録するビットマップ33が存在し（図1）、対応する論理ボリュームの管理単位（トラック・シリンダ・複数のシリンダを含むグループ）ごとに更新フラグビットを持つ。「中断」状態の間にデータが更新された場合、当該データを含む管理単位に対応する更新フラグビットをセットする。再同期時には、いずれかのサブシステムで更新フラグビットのセットされている管理単位について、当該管理単位全体のデータを正側から副側へコピーする。

【0058】

20

一方、本実施例ではこれに属性情報についての更新管理・再同期が追加される。属性情報の更新管理については、共有メモリ25上に属性情報の更新位置を記録する属性情報更新テーブル36が存在する。これは図2で示した属性情報テーブル34と同様の行と列の構成を持っていて、各列の内容が属性情報そのものの代わりに更新されたことを示す更新フラグビットになっている。

【0059】

各サブシステムでは、リモートコピーが「中断」状態に遷移するタイミングで、当該論理ボリュームに対応する行の全ての更新フラグをクリアする。「中断」状態の間は、属性情報の更新を行うチャネルアダプタ上のプロセッサが、属性情報本体を更新すると同時に、対応する属性情報更新テーブル上の更新フラグをセットする。再同期時には、いずれかのサブシステムで属性情報更新テーブル上の更新フラグがセットされている全てについて、当該フラグが指す属性情報の正側の現在値を副側へコピーして、内容の一致化を行なう。

30

【0060】

図7（c）及び（d）いずれの場合も、いずれかのデータ記憶サブシステムで更新された、データ及び属性情報を、全て正側から副側へコピーして再同期を行う。再同期完了後、正常状態に遷移する。

【0061】

<実施例2>

次に本発明の第二の実施例を図8に示す。

40

図8では、第一から第四までのデータ記憶サブシステム11a～dが配置されている。第一の実施例と全く同一の構成要素、すなわち制御装置12・記憶装置13・ディスクアダプタ23・各アダプタ上のポート27およびマイクロプロセッサ28・磁気ディスクドライブ41・ディスク入出力バス42については、説明を省略する。また、以後の説明では、第一サブシステムをコピー元とするリモートコピーの動作のみを説明するため、第一サブシステム以外へのホスト計算機の接続も図示しない。

【0062】

本実施例では、第一サブシステム11aと第二サブシステム11b・第三サブシステム11cとの間、および第三サブシステム11cと第四サブシステム11dとの間が、データ転送バス63にて接続されている。このような構成で、第一データ記憶サブシステム11

50

aに接続されたホスト計算機が使用する論理ボリュームを、第二から第四までの三台のサブシステムにリモートコピーすることが可能である。具体的には、第一サブシステムから第二・第三サブシステムに並行してリモートコピーを実施し、また第三サブシステムは、リモートコピー機能で第一サブシステムから受信した更新情報を、さらに第四サブシステムに転送する。

【0063】

第一サブシステムの動作を説明する。

同期転送の場合、図3に示した処理フローのうち、データ転送パスにて更新情報を送信する処理で、リモートコピーのコピー先となる複数の副側サブシステムへ全て送信し(S12)、全てのサブシステムから転送完了通知が来た時点でホスト計算機に完了通知を行なう(S15)。

【0064】

一方、非同期転送の場合、図4に示した処理フローのうち、正側キャッシュ24へのデータ更新情報32あるいは正側共有メモリ25への属性情報更新情報35への書込み処理で、それぞれリモートコピーのコピー先となるサブシステムの台数分だけ領域を確保して同一内容を書込み(S22、S26)、以後の更新情報を送信する処理はコピー先毎に独立して行なう。

【0065】

いずれの場合でも、リモートコピーの中断については、コピー先の各サブシステム毎に独立に行われるため、中断中の更新情報を管理する更新ビットマップ33ならびに更新属性情報テーブル36は、コピー先の副側サブシステムの台数分だけ用意される。

【0066】

以上の処理により、複数台のサブシステムに対して、データならびに属性情報のリモートコピーを実現する。

【0067】

続いて第三サブシステムの動作を説明する。

第一サブシステムとのリモートコピーの関係が同期転送ならば、第一サブシステムから受信した更新情報を、ホスト計算機からのデータ更新ないしは制御指示とみなし第四サブシステムに対するリモートコピーを実施する。第四サブシステムに対するリモートコピーの処理は、第一の実施例で説明した同期ないしは非同期転送の場合の処理と同一である。

【0068】

一方、第一サブシステムとのリモートコピーの関係が非同期転送ならば、既に要求元のホスト計算機との同期関係は崩れているので、第四サブシステムに対しては非同期転送によるリモートコピーが一般的に使用される。この場合は、第一サブシステムから受信した更新情報を、自サブシステムへの更新に使うために一時記憶する他に、二次のリモートコピーを実施するコピー先の副側サブシステム台数分だけ別に記憶する。この別に記憶した更新情報を使用して、第一の実施例の図4の処理フローのうち、更新情報をデータ転送パスで転送する処理(S27)以降の処理を実施する。キャッシュ24上には自サブシステムでの更新用と第四サブシステムへの送信用の二種類の受信更新情報32が格納され、共有メモリ25上には自サブシステムでの更新用と第四サブシステムへの送信用の二種類の受信更新情報35が格納される。

【0069】

いずれの場合でも、リモートコピーの中断については、受信側(第一サブシステムとの間)・送信側(第四サブシステムとの間)で独立に行われるため、中断中の更新情報を管理する更新ビットマップ33ならびに更新属性情報テーブル36は、リモートコピーの相手サブシステム台数分だけ用意される。

【0070】

以上の処理により、リモートコピーの副側サブシステムにおいて二次的に別のサブシステムに対して、データならびに属性情報のリモートコピーを実現することができる。

【0071】

本実施例の仕組みを複数組み合わせることにより、任意の数のサブシステム間で、データならびに属性情報のリモートコピーを実現することができる。

#### 【 0 0 7 2 】

##### 【 発明の効果 】

本発明によれば、ホスト計算機側では全く意識することなしに、属性情報がバックアップ側のデータ記憶サブシステムにコピーされ、バックアップ側で自動的に利用可能となる。これにより、ホスト計算機間での複雑な属性情報の交換が不要となる。

#### 【 0 0 7 3 】

本発明によれば、ホスト計算機側では全く意識することなしに、データおよび属性情報の更新の順序を維持したまま、しかも高速に、コピー先のデータ記憶サブシステム上にこれらの更新が反映される機能を実現できる。これにより、属性情報の更新とデータの更新の間の順序関係を維持することができ、障害時の復旧処理等が容易になる。

10

#### 【 0 0 7 4 】

本発明によれば、ホスト計算機側では全く意識することなしに、リモートコピー中断後の再開時やコピー方向反転時での属性情報の差分コピーを実現できる。これにより、短時間で属性情報の再同期を終了できる。

#### 【 0 0 7 5 】

本発明によれば、複数のデータ記憶サブシステムにまたがる属性情報のコピーを実現できる。これにより、複数のデータ記憶サブシステムにまたがるリモートコピーの適用が容易となり、耐障害性をより高くしつつ、システムの構築が容易となる。

20

#### 【 図面の簡単な説明 】

【 図 1 】 本発明の複合データ記憶サブシステムの全体構成図である。

【 図 2 】 属性情報テーブルの構成である。

【 図 3 】 同期転送リモートコピーの場合の処理フローである。

【 図 4 】 非同期転送リモートコピーの場合の処理フローである。

【 図 5 】 データ転送パス上のデータ更新情報のライトコマンドである。

【 図 6 】 データ転送パス上の属性情報更新情報のライトコマンドである。

【 図 7 】 リモートコピーの中断に関する処理を示す図である。

【 図 8 】 本発明の別な実施例の全体構成図である。

【 図 9 】 同期転送方式と非同期転送方式を説明する図である。

30

#### 【 符号の説明 】

1 1 a、1 1 b、1 1 c、1 1 d ... データ記憶サブシステム、

1 2 a、1 2 b ... 制御装置、

1 3 a、1 3 b ... 記憶装置、

2 3 ... ディスクアダプタ、

2 5 ... 共有メモリ、

2 7 ... ポート、

3 1 ... キャッシュエントリ、

3 3 ... 更新情報、

3 5 ... 属性情報更新テーブル、

4 1 ... 磁気ディスクドライブ、

6 1 ... ホスト計算機、

6 3 ... データ転送パス。

2 1、2 2 ... チャンネルアダプタ、

2 4 ... キャッシュ、

2 6 ... バス、

2 8 ... マイクロプロセッサ、

3 2 ... データ更新情報、

3 4 ... 属性情報テーブル、

3 6 ... 更新属性情報テーブル、

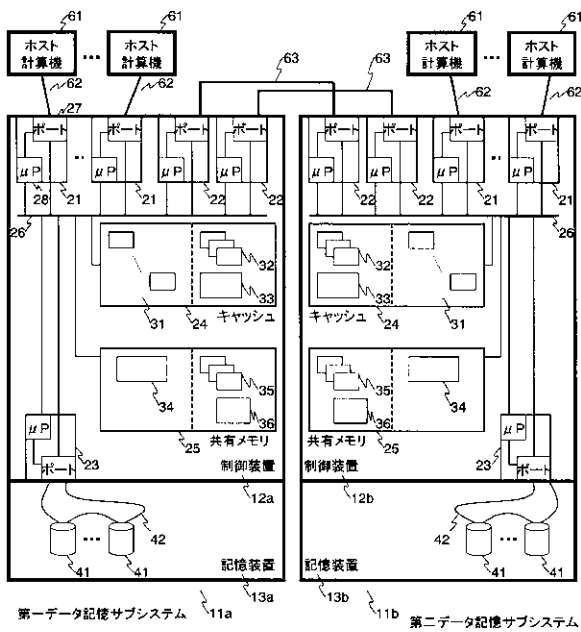
4 2 ... ディスク入出力パス、

6 2 ... ホストアクセスパス、

40

【図 1】

【図1】 全体構成図



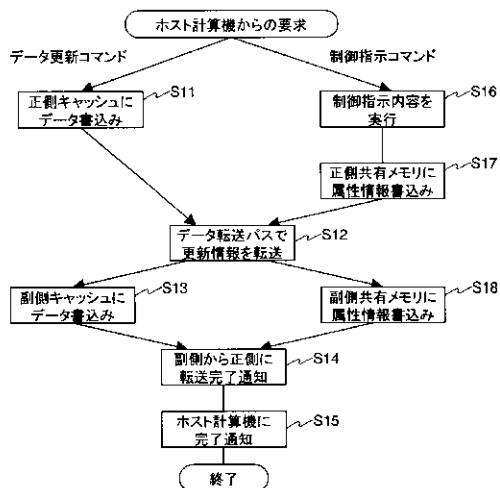
【図 2】

【図2】 属性情報テーブルの構成

LUN	リザーブ 制御情報	更新 アクセス 許可キー	使用中 ホストID	.....	ユーザ領域
00:00	0	321A	—	.....	
00:01	1, TID=7	—	—	.....	3FC00129 .....
⋮	⋮	⋮	⋮	⋮	⋮
FF:FF	0	—	00-00-3C	.....	

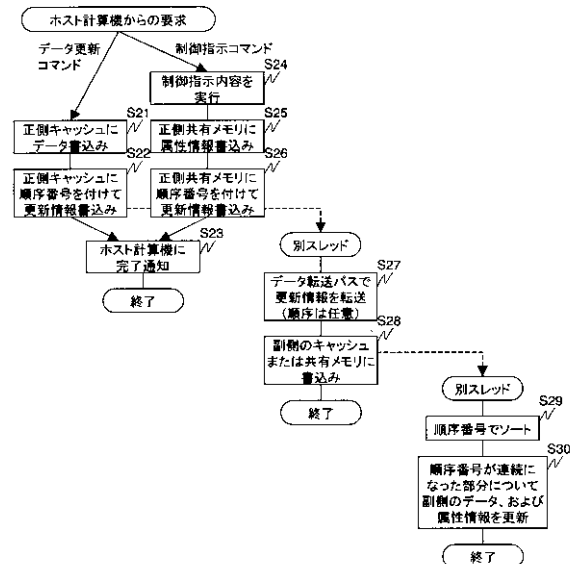
【図 3】

【図3】 同期転送の場合の処理フロー



【図 4】

【図4】 非同期転送の場合の処理フロー



【図5】

【図5】 データ転送バス上のコマンド(データ)

オペレーションコード=リモートコピーライト
ターゲットLUN番号
データ格納 先頭アドレス
転送長(データサイズ)
順序番号
フラグ=データ

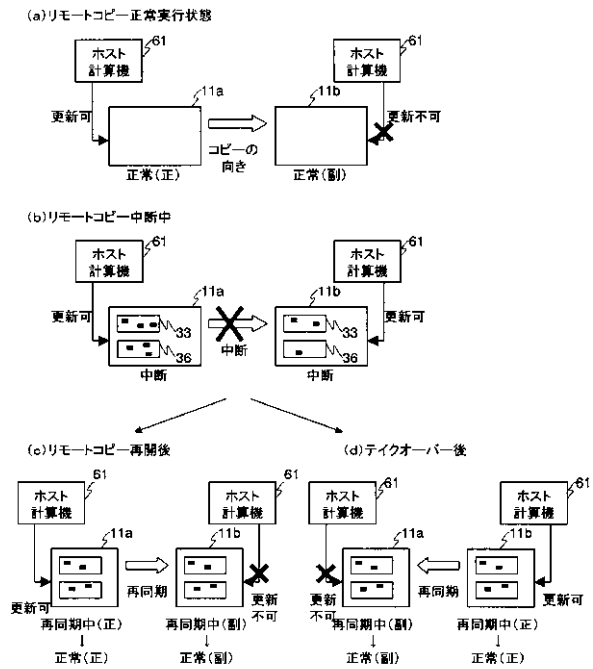
【図6】

【図6】 データ転送バス上のコマンド(属性情報)

オペレーションコード=リモートコピーライト
ターゲットLUN番号
属性情報のテーブル内アドレス
転送長(属性情報サイズ)
順序番号
フラグ=属性情報

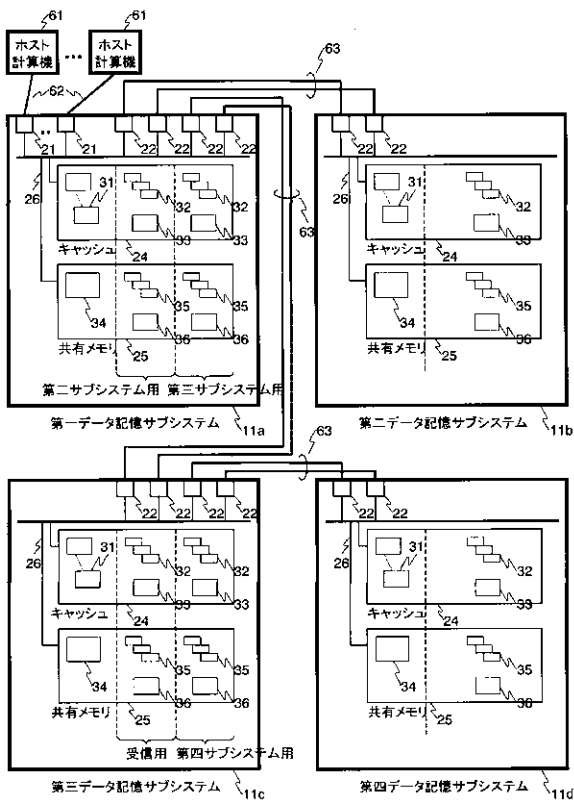
【図7】

【図7】 リモートコピーの中断



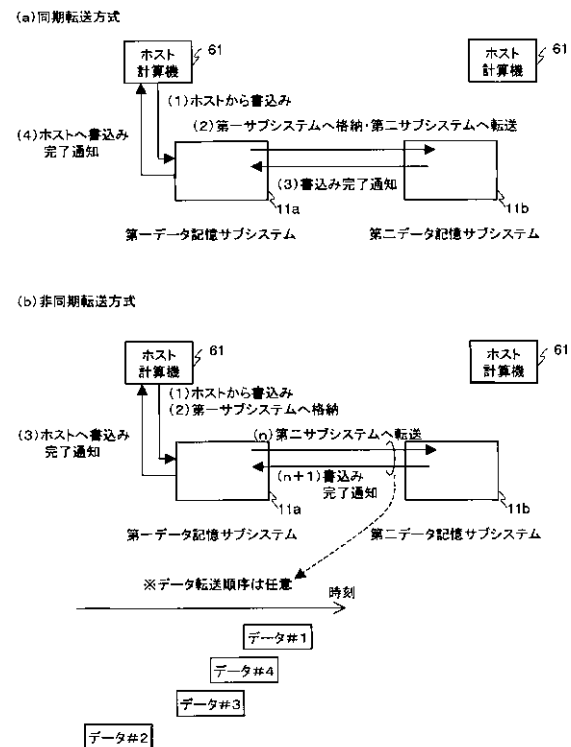
【図8】

【図8】 別な実施例の全体構成図



【図9】

【図9】 同期転送方式と非同期転送方式の説明



---

フロントページの続き

(72)発明者 中野 俊夫

神奈川県小田原市中里 3 2 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

(72)発明者 田淵 英夫

神奈川県小田原市中里 3 2 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

F ターム(参考) 5B065 BA01 CA12 CC08 CE22 EA34 EA35

5B082 EA11 HA03