

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第4892185号
(P4892185)

(45) 発行日 平成24年3月7日 (2012.3.7)

(24) 登録日 平成23年12月22日 (2011.12.22)

(51) Int.Cl.

G 0 6 F 3 / 0 6 (2 0 0 6 . 0 1)

F 1

G 0 6 F 3 / 0 6 3 0 4 F

請求項の数 11 外国語出願 (全 25 頁)

(21) 出願番号	特願2004-338341 (P2004-338341)	(73) 特許権者	000005108
(22) 出願日	平成16年11月24日 (2004.11.24)		株式会社日立製作所
(65) 公開番号	特開2005-276162 (P2005-276162A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成17年10月6日 (2005.10.6)	(74) 代理人	100093861
審査請求日	平成19年9月27日 (2007.9.27)		弁理士 大賀 真司
(31) 優先権主張番号	10/758971	(72) 発明者	山神 憲司
(32) 優先日	平成16年1月15日 (2004.1.15)		アメリカ合衆国カリフォルニア州ロスガト
(33) 優先権主張国	米国 (US)		ス カルニベル 1 0 8

審査官 菅原 浩二

最終頁に続く

(54) 【発明の名称】 分散リモートコピーシステム

(57) 【特許請求の範囲】

【請求項 1】

第一のプライマリボリュームを含む第一のプライマリストレージサブシステム、及び第二のプライマリボリュームを含む第二のプライマリストレージサブシステムであって、前記第一及び第二のプライマリボリュームは、複数の書込みデータを当該複数の書き込みデータの書込み順に保存する、第一及び第二のプライマリストレージサブシステムと、

前記第一及び第二のプライマリストレージサブシステムに接続されるとともに、前記第一及び第二のプライマリストレージサブシステムから前記書込みデータを同期的に受信するように構成された、中間ストレージサブシステムであって、前記第一及び第二のプライマリストレージサブシステムから受信した前記書込みデータの保存の順番を反映する書込み順情報を生成するカウンタを含み、

データ完全性が保証される整合性グループとして定義された第一及び第二の中間ボリュームを含み、前記第一の中間ボリュームは、前記第一のプライマリボリュームから前記書込みデータを受信するように構成され、前記第二の中間ボリュームは前記第二のプライマリボリュームから前記書込みデータを受信するように構成される、中間ストレージサブシステムと、前記中間ストレージサブシステムの有効カウンタであって、前記有効カウンタは、前記中間ストレージサブシステムに接続されるとともに、前記中間ストレージサブシステムから前記書込みデータを非同期的に受信するように構成された、第一及び第二のセカンダリストレージサブシステムでのコピーのために有効化される準備ができていて前記書込みデータの最大順番番号を保持するように構成される有効カウンタと、

10

20

前記第一及び第二のセカンダリストレージサブシステムであって、前記第一のセカンダリストレージサブシステムは、前記中間ストレージサブシステムとは別のストレージサブシステムである前記第一のプライマリストレージサブシステムの前記第一のプライマリボリュームをミラーするように構成された第一のセカンダリボリュームを含み、前記第二のセカンダリストレージサブシステムは、前記中間ストレージサブシステムとは別のストレージサブシステムである前記第二のプライマリストレージサブシステムの前記第二のプライマリボリュームをミラーするように構成された第二のセカンダリボリュームを含み、前記第一のプライマリボリュームを有する前記第一のプライマリストレージサブシステム及び前記第二のプライマリボリュームを有する前記第二のプライマリストレージサブシステムとは別のストレージサブシステムである前記中間ストレージサブシステムの前記書込み順情報プロバイダによって生成されたとおりの、前記書込みデータと対応付けられた前記書込み順情報に従って、前記書込みデータは、前記第一及び第二のセカンダリストレージサブシステムに保存され、前記中間ストレージサブシステムから前記準備要求受信した場合、前記準備要求受信前に受信した書き込みデータを回復するために必要な情報である回復用情報を保持し、前記中間ストレージサブシステムから前記有効化を受信した場合、前記回復用情報を破棄する第一及び第二のセカンダリストレージサブシステムと、

を備え、

前記中間ストレージサブシステムは、前記書込みデータを準備するようにという準備要求を送信し、かつ当該準備要求に従って準備された前記書込みデータを前記有効カウンタの値に基づいて有効化するようにという有効化要求を送信する、

リモートコピーシステム。

【請求項 2】

各々前記第一及び第二のプライマリサブシステムに設けられる第一及び第二のプライマリビットマップと、

前記中間サブシステムに設けられる第一及び第二の中間ビットマップと、
各々前記第一及び第二のセカンダリサブシステムに設けられる第一及び第二のセカンダリビットマップと

をさらに備え、

前記第一及び第二のプライマリビットマップは前記第一の中間ビットマップに対応付けられ、前記第一及び第二のセカンダリビットマップは前記第二の中間ビットマップに対応付けられ、

前記ビットマップは、再同期プロセスで、ペアボリュームのミラー処理が中断した後に変更のあったデータのみを決定し、かつ当該データをコピーするために使用される、請求項 1 に記載のリモートコピーシステム。

【請求項 3】

前記中間サブシステムは、前記第一及び第二のプライマリサブシステムから前記書込みデータを受信するジャーナルボリュームを含む、請求項 1 に記載のリモートコピーシステム。

【請求項 4】

前記第一及び第二のプライマリサブシステムはディスクアレイユニットである、請求項 1 に記載のリモートコピーシステム。

【請求項 5】

前記プライマリサブシステムはプライマリサイトに設けられ、前記セカンダリサブシステムはセカンダリサイトに設けられ、前記プライマリサイトはプライマリホストを含み、前記セカンダリサイトはセカンダリホストを含み、

前記プライマリサブシステム又は前記プライマリホストのいずれかが障害を被るか又はオフラインにされた場合に、前記セカンダリサブシステムは、第一の保存領域として、前記プライマリサブシステムに取って代わるように構成されている、請求項 1 に記載のリモートコピーシステム。

【請求項 6】

第一のプライマリストレージサブシステムとは別のストレージサブシステムである中間ストレージサブシステムで、第一のプライマリストレージサブシステムから同期的に送信される第一の書込みデータを受信することと、

データ完全性が保証される整合性グループとして定義される複数のボリュームからなる中間ボリュームを含む前記中間ストレージサブシステムで、第一の書込み順情報を前記第一の書込みデータに対応付けることと、

前記第二のプライマリストレージサブシステムとは別のストレージサブシステムである前記中間サブシステムで、第二のプライマリストレージサブシステムから同期的に送信される第二の書込みデータを受信することと、

前記中間ストレージサブシステムで、第二の書込み順情報を前記第二の書込みデータに対応付けることと、

前記中間ストレージサブシステムとは別のストレージサブシステムである第一のセカンダリストレージサブシステムに、前記第一の書込みデータ及び前記第一の書込み順情報を非同期的に送信することと、

前記中間ストレージサブシステムとは別のストレージサブシステムである第二のセカンダリストレージサブシステムに、前記第二の書込みデータ及び前記第二の書込み順情報を非同期的に送信することと、

前記第一及び第二のセカンダリサブシステムに保存するために前記第一、第二及び第三の書込みデータを準備するようという要求であって、参照用の順番番号を含む準備要求を、前記中間サブシステムから前記第一及び第二のセカンダリサブシステムへ送信することと、

前記準備要求に従って準備された前記書込みデータを有効化するようという要求を送信することであって、当該有効化要求は、前記中間サブシステムから前記第一及び第二のセカンダリサブシステムへ送信され、準備すべき前記書込みデータを識別し、前記参照用の順番番号よりも小さい又は等しい前記順番番号で前記書込みデータを有効化する、前記有効化要求を送信することと、

を含む、リモートコピーシステムを操作する方法であって、

前記第一及び第二の書込みデータは、各々前記第一及び第二の書込み順情報に従って前記第一及び第二のセカンダリサブシステムに保存され、

前記中間ストレージサブシステムは、ホストユニットに直接的には接続されていない、
リモートコピーシステムを操作する方法。

【請求項 7】

前記中間ストレージサブシステムで、前記第一のプライマリサブシステムから第三の書込みデータを、前記第一の書込みデータの後に、同期的に受信することと、

前記第三の書込みデータを第三の書込み順情報に対応付けることと、

前記第三の書込みデータ及び前記第三の書込み順情報を前記第一のセカンダリサブシステムに非同期的に送信することと、

をさらに含み、前記第一及び第三の書込みデータは、前記第一のセカンダリサブシステムでの保存領域を識別する保存先アドレスと同一の保存先アドレスを有し、

前記第一及び第三の書込み順情報は、前記第三の書込みデータを前記識別された保存領域に保存する前に、前記第一の書込みデータを前記識別された保存領域に保存するために使用される、請求項 6 に記載の方法。

【請求項 8】

リモートコピーシステムに設けられるとともに、複数のプライマリストレージサブシステムに接続され、かつ複数のセカンダリサブシステムに接続される、中間ストレージサブシステムであって、

前記中間ストレージサブシステムとは別のストレージサブシステムである少なくとも一つのプライマリサブシステムから書込みデータを受信するように構成された第一の中間保存領域であって、前記書込みデータは、前記少なくとも一つのプライマリサブシステムから同期的に受信され、前記第一の中間保存領域は、データ完全性が保証される整合性グル

10

20

30

40

50

ープとして定義される複数のボリュームの第一の中間保存領域と、

前記少なくとも一つのプライマリサブシステムから受信された前記書込みデータのために書込み順情報を生成するように構成された、カウンタであって、前記書込み順情報は前記書込みデータに対応付けられる、カウンタと、

前記中間サブシステムに設けられた有効カウンタであって、前記セカンダリサブシステムでのコピーのために有効化される準備ができている前記書込みデータの最大順番番号を保持するように構成された、有効カウンタと、

第二の中間保存領域と、

を備える中間ストレージサブシステムであり、

前記中間ストレージサブシステムとは別のストレージサブシステムである前記セカンダリサブシステムのうちの少なくとも一つに前記書込みデータを保存するために、前記書込み順情報を使用することにより、前記少なくとも一つのセカンダリサブシステムは、前記少なくとも一つのプライマリサブシステムをミラーし、

前記中間ストレージサブシステムは、前記書込みデータを準備するようにという要求と、当該準備要求に従って準備された前記書込みデータを、前記有効カウンタの値に基づいて有効化するようにという要求とを送信し、

前記複数のプライマリサブシステムは、第一のプライマリサブシステムに設けられた第一のプライマリボリュームと、第二のプライマリサブシステムに設けられた第二のプライマリボリュームとを含み、

前記第一及び第二の中間保存領域は第一及び第二の中間ボリュームであって、前記第一の中間ボリュームは、前記第一のプライマリボリュームから書込みデータを受信するように構成され、前記第二の中間ボリュームは、前記第二のプライマリボリュームから書込みデータを受信するように構成され、

前記第一の中間ボリュームは、前記第一のプライマリボリュームから受信した前記書込みデータを、第一のセカンダリサブシステムに設けられた第一のセカンダリボリュームに送信するように構成され、前記第二の中間ボリュームは、前記第二のプライマリボリュームから受信した前記書込みデータを、第二のセカンダリサブシステムに設けられた第二のセカンダリボリュームに送信するように構成され、

前記中間ストレージサブシステムは、ホストユニットに直接的には接続されてはいない、

中間ストレージサブシステム。

【請求項 9】

前記第一の保存領域は、前記少なくとも一つのプライマリサブシステムから所定の順番で第一及び第二の書込みデータを受信するように構成され、前記第一及び第二の書込みデータには各々、

前記中間サブシステムによって第一及び第二の書込み順情報が提供され、前記第一及び第二の書込みデータは、前記第一及び第二の書込み順情報を用いて、前記所定の順番に従って前記少なくとも一つのセカンダリサブシステムに保存される、請求項 8 に記載のストレージサブシステム。

【請求項 10】

前記書込み順情報は、順番番号を生成するように構成されたカウンタによって生成され、当該生成された順番番号は、前記書込みデータが前記少なくとも一つのプライマリサブシステムから受信された順番に従って、前記書込みデータに対応付けられる、請求項 8 に記載のストレージサブシステム。

【請求項 11】

前記第一の保存領域は、前記複数のプライマリサブシステムから書込みデータを受信するように構成されたジャーナルボリュームである、請求項 8 に記載のストレージサブシステム。

【発明の詳細な説明】

【技術分野】

【0001】

0001 本発明はストレージシステムに関わり、より具体的には、リモートコピー機能を果たす分散ストレージシステムに関連する。

【背景技術】

【0002】

0002 データは全ての計算プロセスの基礎となる資源である。最近のインターネットやe-ビジネスの爆発的成長により、データストレージシステムに対する需要は凄まじく増大している。一般的にネットワーク結合型ストレージはNAS(Network Attached Storage)又はSAN(Storage Area Network)の二つの形態をとっている。NASでは、ファイルフォーマットのデータをストレージサーバとクライアント間で転送するのにイーサネット(イーサネットは登録商標です。)上のIPを使用する。NASに於いては、ディスクレイヤやテープデバイスの如き集合ストレージシステムは、TCP/IPのようなメッセージ通信プロトコルを用いて、イーサネット等のLAN(Local Area Network)を通してメッセージ交換ネットワークに直結される。本ストレージシステムはクライアント-サーバシステムでのサーバとしての働きをする。

【0003】

0003 一般的にSANは種々のサーバとストレージ資源の間でデータ移動を行う専用の高速ネットワークである。NASとは異なって、クライアントとサーバ間での慣用的メッセージ交換での通信衝突を避ける為に、分離した専用ネットワークが用いられる。SANでは、ストレージ資源とプロセッサ又はサーバとの間の直接結合が可能になる。SANは単一サーバ専用でも又サーバ間で共用することもできる。単一のローカルサイトに集中することも又地理的に離れたサイトに拡張することもできる。SANインタフェースでは、FC(Fibre Channel)、ESCON(Enterprise System Connection)、SCSI(Small Computer System Interface)、SSA(Serial Storage Architecture)、HIPPI(High Performance Parallel Interface)その他将来登場するものも含めて多様で異なったプロトコルが使用可能である。例えば、IETF(Internet Engineering Task Force)はTCP/IP上でブロックレベル転送を可能にするiSCSIの新プロトコルを開発中である。又ある会社は、iSCSIをSANでの主要スタンダードにすべく、iSCSI TCP/IPプロトコルスタックをホストプロセッサから分離しようとしている。

【0004】

0004 ストレージシステムのタイプに関わらず、データストレージシステムのユーザは、貴重なデータがストレージユニット(又はストレージシステム)の障害で損失することを防ぐ為に、バックアップデータを維持する事に強い関心がある。従って、データストレージシステムは、プライマリユニットが障害になった時の緊急回復データを保存するバックアップユニットを持つことが普通である。然しながら、障害は、ユニット自身の問題によるばかりではなく、ストレージユニットの設置場所での地震や嵐等の自然災害で発生する事もある。バックアップユニットがプライマリユニットの近くに設置されている場合には、両ユニットとも自然災害で損傷を受ける事があり得る。従って、多くのストレージシステムのユーザは、プライマリユニットとバックアップユニットを例えば100マイル以上も離れて設置している。実際には、プライマリとバックアップユニットを異なる大陸に設置しているユーザもある。

【0005】

0005 今日では、同期モードと非同期モードの二つの動作モードが、ストレージシステムがデータをバックアップ又はセカンダリサイトにコピーするのに用いられる。同期モードでは、ホストからのプライマリストレージシステムへの書込み要求は、書込みデータがセカンダリストレージシステムにコピーされ、そこからの応答が返って初めて完了する。従って、本モードでは、セカンダリストレージシステムからの応答が返るまでは、プライマリシステムのキャッシュにホストからのデータが保存されている為に、セカンダリストレージシステムでのデータ損失は発生しないことが保証されている。加えて、プライマリ

ストレージシステムのプライマリボリューム(PVOL)とセカンダリストレージシステムのセカンダリボリューム(SVOL)は同一に保たれている為に、PVOLが損失した場合にはSVOLが直ちに代わりを務める事ができる。然しながら、このモードでは、プライマリとセカンダリは、例えば、100マイルも離れて設置する事はできない。距離を離すと、ストレージシステムは、ホストからの書き込み要求を効率的に実行することはできない。

【0006】

0006 非同期モードでは、ホストからのプライマリストレージシステムへの書き込み要求は、データを指定プライマリシステムに書き込むのみで終了する。書き込みデータはその後、セカンダリストレージシステムにコピーされる。即ち、プライマリストレージシステムへのデータ書き込みは、セカンダリストレージシステムへのデータコピーとは独立したプロセスである。この結果として、プライマリとセカンダリは、例えば、100マイル以上も離れて設置することができる。然しながら、PVOLとSVOLは常に同じデータを維持している訳ではない為、プライマリシステムがダウンしたらデータが失われる可能性がある。従って、同期モードと非同期モードの長所、即ち、データを失う事無くプライマリとセカンダリシステムが離れて設置できるデータストレージシステムが望まれるところである。非同期モードの好適なリモートコピー方法については、参照記載のYamamotoその他による、米国特許NO. 6,408,370に公開されている。

【0007】

【特許文献1】米国特許第6408370号

【発明の開示】

【発明が解決しようとする課題】

【0008】

同期モードでは、ホストからのプライマリストレージシステムへの書き込み要求は、書き込みデータがセカンダリストレージシステムにコピーされ、そこからの応答が返って初めて完了する。従って、本モードでは、セカンダリストレージシステムからの応答が返るまでは、プライマリシステムのキャッシュにホストからのデータが保存されている為に、セカンダリストレージシステムでのデータ損失は発生しないことが保証されている。加えて、プライマリストレージシステムのプライマリボリューム(PVOL)とセカンダリストレージシステムのセカンダリボリューム(SVOL)は同一に保たれている為に、PVOLが損失した場合にはSVOLが直ちに代わりを務める事ができる。然しながら、このモードでは、プライマリとセカンダリは、例えば、100マイルも離れて設置する事はできない。距離を離すと、ストレージシステムは、ホストからの書き込み要求を効率的に実行することはできない。

【0009】

非同期モードでは、ホストからのプライマリストレージシステムへの書き込み要求は、データを指定プライマリシステムに書き込むのみで終了する。書き込みデータはその後、セカンダリストレージシステムにコピーされる。即ち、プライマリストレージシステムへのデータ書き込みは、セカンダリストレージシステムへのデータコピーとは独立したプロセスである。この結果として、プライマリとセカンダリは、例えば、100マイル以上も離れて設置することができる。然しながら、PVOLとSVOLは常に同じデータを維持している訳ではない為、プライマリシステムがダウンしたらデータが失われる可能性がある。従って、同期モードと非同期モードの長所、即ち、データを失う事無くプライマリとセカンダリシステムが離れて設置できるデータストレージシステムが望まれるところである。

【課題を解決するための手段】

【0010】

0007 リモートコピーシステム(リモートコピー機能を備えるストレージシステム)は一つ又は複数のプライマリストレージサブシステムを持つことができる。分散リモートコピーシステムは、複数のプライマリストレージサブシステムを持つシステムを言う。リモートコピーの目的の為に、分散リモートコピーシステムでの複数のプライマリストレージサ

10

20

30

40

50

ブシステムのデータは、複数のセカンダリストレージサブシステムにコピーされる。例えば、第一のプライマリサブシステムは第一のセカンダリサブシステムとペアを組みミラーされ、又第二のプライマリサブシステムは第二のセカンダリサブシステムとペアを組みミラーされる。

【 0 0 1 1 】

0008 このようなリモートコピーシステムで、データインテグリティ(データの完全性)を保つ事は、現在使用可能なストレージデバイスと技術レベルのもとでは、困難な課題である。プライマリストレージサブシステムを使用するデータベースでは、書込み順又はデータインテグリティを保証する為に、一般的に、書込み順を記録する書込みログを書き込む。第一のログが第一のプライマリサブシステムの為に、第二のログが第二のプライマリサブシステムの為に書き込まれる。第一と第二のログが共に、プライマリサブシステムでのデータベースのデータインテグリティの為に使用される。同様に、第一と第二のログは、それぞれ、セカンダリサブシステムでの書込み順を保証する為に、セカンダリサブシステムにコピーする必要がある。プライマリサブシステムでの実際の書込み順をミラーする為に、二つのログ間の関係情報を維持しながら、第一と第二のログをセカンダリサブシステムにコピーする必要がある。然しながら、現行技術によるリモートコピーでは、第一と第二のプライマリサブシステムの間には、殆ど協調関係が存在しない為、このことは容易には実現できないことである。

10

【 0 0 1 2 】

0009 ここで用いる“ストレージシステム”は、データ保存の為に構成されたコンピュータシステムであり、一つ以上のストレージユニット又はディスクアレイの如きストレージサブシステムを備える。従って、ストレージシステムは、一つ以上のホストと一つ以上のストレージサブシステム、又は一つだけのストレージサブシステム又はユニット、又は通信回線で複数のホストに結合された複数のストレージサブシステム又はユニットを備えるコンピュータシステムである。

20

【 0 0 1 3 】

0010 ここで用いる“ストレージサブシステム”は、データ保存の為に構成され、ストレージ領域及び一つ以上のホストからの要求を処理するためのストレージコントローラを備えるコンピュータシステムである。ストレージサブシステムは、ストレージデバイス、ストレージユニット、ストレージ装置などと呼ばれる事がある。ストレージサブシステムの一例はディスクアレイユニットである。

30

【 0 0 1 4 】

0011 ここで用いる“ホスト”は、一つ以上のストレージシステム又はストレージサブシステムに結合し、本ストレージシステム又はストレージサブシステムに要求を送信するように構成されたコンピュータシステムである。ホストはサーバ又はクライアントとしての役割を果たす。

【 0 0 1 5 】

0012 ここで用いる“リモートコピーシステム”は、リモートコピー機能を果たすように構成されたコンピュータシステムである。リモートコピーシステムは、ネットワーク又は通信回線で結合した、単一のストレージシステム、ストレージサブシステム又はユニット、あるいは複数のストレージユニット、ストレージシステム、又はストレージサブシステムを指すこともある。従って、リモートコピーシステムは、プライマリストレージシステム、セカンダリストレージシステム、中間ストレージシステム、又はこれらの組合せを示すことがある。更にリモートコピーシステムは、一つ以上のホストを含むこともできる。

40

【 0 0 1 6 】

0013 一実施例では、プライマリストレージサブシステム(以降プライマリサブシステムで表す)は中間ストレージサブシステム(以降中間サブシステムで表す)にデータを同期的に送信する。書込みデータの受信を契機に、中間サブシステムは書込み順を指定する情報を含む制御データを生成する。一実施例では、制御データの一部例えば順序番号が生成

50

され、書込みデータの制御データに付加される。中間サブシステムは、この書込みデータを制御データと共に、セカンダリストレージサブシステム(以降セカンダリサブシステムで表す)に非同期的に送信する。セカンダリサブシステムは、制御データに基づいてこの書込みデータをセカンダリボリュームに保存する。書込み順は、中間サブシステムにより生成された順序番号により維持される。

【 0 0 1 7 】

0014 一実施例では、リモートコピーシステムは、複数の書込みデータを与えられた順で保存する各々第一と第二のプライマリボリュームを備える第一と第二のプライマリサブシステムと；第一と第二のプライマリサブシステムからの書込みデータを受信し、本書込みデータに対して、本第一と第二のプライマリサブシステムで与えられた書込み順を反映し、各書込みデータに付与される書込み順情報、を生成する書込み順情報生成部を備える、中間サブシステムと；各々第一と第二のプライマリボリュームをミラーする第一と第二のセカンダリボリュームを備え、本中間サブシステムからの書込みデータを受信する第一と第二のセカンダリサブシステムと；から成り、各データに付与された書込み順情報に基づいて、本書込みデータを第一と第二のセカンダリサブシステムに保存する。

10

【 0 0 1 8 】

0015 他の実施例では、リモートコピーシステムで、複数のプライマリサブシステムと複数のセカンダリサブシステムに結合される中間サブシステムは、少なくとも一式のプライマリサブシステムからの書込みデータを同期的に受け取る第一のストレージ領域と；少なくとも一式のプライマリサブシステムから受けとった各書込みデータに付与される書込みデータの書込み順情報、を生成する書込み順情報生成部とを持ち、本書込み順情報は、セカンダリサブシステムの少なくとも一式でデータ書込みの為に使用され、少なくとも一式のセカンダリサブシステムが少なくとも一式のプライマリサブシステムのミラーをする。

20

【 0 0 1 9 】

0016 更に他の実施例では、リモートコピーシステムを操作する方法が提供され、中間サブシステムで第一と第二のプライマリサブシステムから各々第一と第二の書込みデータを同期的に受信し；本第一と第二の書込みデータに各々第一と第二の書込み順情報を付与し；本第一と第二の書込みデータと第一と第二の書込み順情報を各々第一と第二のセカンダリサブシステムに非同期的に送信する；ステップから成り、第一と第二の書込みデータは、各々第一と第二の書込み順情報に従って、第一と第二のセカンダリサブシステムに保存される。

30

【 0 0 2 0 】

0017 更に他の実施例は、リモートコピーシステムを操作するコンピュータプログラムを格納するコンピュータ読み取り可能な媒体を提供し、本媒体は、中間サブシステムで、第一と第二のプライマリサブシステムでの各々第一と第二のプライマリボリュームから第一と第二の書込みデータを同期的に受信するコードと；本第一と第二の書込みデータに各々第一と第二の書込み順情報を付与するコードと；本第一と第二の書込みデータと本第一と第二の書込み順情報を、各々第一と第二のセカンダリボリュームより成る第一と第二のセカンダリサブシステムに、非同期的に転送するコードとを、格納し、第一と第二の書込みデータは各々第一と第二の書込み順情報に従って、第一と第二のセカンダリサブシステムに保存され、第一と第二のセカンダリボリュームは各々第一と第二のプライマリボリュームのミラーをする。

40

【 0 0 2 1 】

0018 更にもう一つの実施例では、分散リモートコピーシステムでの中間サブシステムは、第一と第二のプライマリサブシステムで定義される各々第一と第二のプライマリボリュームから書込みデータを同期的に受信する手段と；プライマリサブシステムから受信した書込みデータに書込み順情報を与える書込み順情報生成手段と、から成り、本書込み順情報は、各々第一と第二のセカンダリサブシステムで定義される第一と第二のセカンダリボリュームへの書込みデータの書込みに使用され、本第一と第二のセカンダリボリューム

50

が各々第一と第二のプライマリボリュームをミラーする。

【発明の効果】

【0022】

本リモートコピーシステムは好適な分散リモートコピーシステムで、セカンダリストレージサブシステム(以下セカンダリサブシステムと呼ぶ)に亘って、データのインテグリティ(データの完全性)を維持すべく構成されている。

【発明を実施するための最良の形態】

【0023】

0033 図1A(図1(a))は、本発明の一実施例による、複数のプライマリストレージサブシステム(以下プライマリサブシステムと呼ぶ)100aを持つリモートコピーシステム50を示す。本リモートコピーシステムは好適な分散リモートコピーシステムで、セカンダリストレージサブシステム(以下セカンダリサブシステムと呼ぶ)に亘って、データのインテグリティ(データの完全性)を維持すべく構成されている。

【0024】

0034 システム50は、複数のプライマリサブシステム100a、複数のセカンダリサブシステム100b、及び一つの間接サブシステム(以下中間サブシステムと呼ぶ)100cを含む。プライマリサブシステムは第一と第二のプライマリサブシステム100a-1と100a-2を含む。セカンダリサブシステムは第一と第二のセカンダリサブシステム100b-1と100b-2を含む。一実施例では、単一の中間サブシステムが用いられる。

【0025】

0035 本実施例では、プライマリサブシステムはデータを中間サブシステムに同期的にコピーし、中間サブシステムは本データをセカンダリサブシステムに非同期的にコピーする。中間サブシステム100cは、プライマリサブシステム100aの近くに設置され、セカンダリサブシステム100bとは十分離れて設置される。例えば、プライマリと中間サブシステムは10マイル以内に設置され、中間とセカンダリサブシステムは互いに100マイル以上離れて設置される。一実施例では、プライマリと中間サブシステムは、同じビル又は複合建築体内に設置される。

【0026】

0036 図1B(図1(b))は、読み書き要求を処理するストレージコントローラ62と書込み要求に対応してデータ保存の為に記録媒体を備えるストレージユニット63を含む好適なストレージサブシステム60を示す。ストレージコントローラ62は、ホストコンピュータに結合する為のホストチャネルアダプタ64と他のサブシステムに結合する為のサブシステムチャネルアダプタ66とストレージサブシステム60内のストレージユニット63に結合する為のディスクアダプタ68を有する。本実施例では、これらの各アダプタはデータを送受信する為のポート(図示していない)を備え、本ポートを通してデータ転送をコントロールする為のマイクロプロセッサ(図示していない)を備える。

【0027】

0037 ストレージコントローラ62は、更に、ストレージユニット63との間で読み書きするデータを一時的に保存する為のキャッシュメモリ70を備える。一実施例ではストレージユニットは複数の磁気ディスクドライブ(図示していない)である。

【0028】

0038 本サブシステムはホストコンピュータに対する記憶領域として、複数の論理ボリュームを提供する。ホストコンピュータはストレージサブシステムとの間でデータを読み書きする為に、これら論理ボリュームの識別子を使用する。これら論理ボリュームの識別子はLUN(Logical Unit Number)と呼ばれる。これら論理ボリュームは単一の物理ストレージデバイスに含まれてもよく又複数の物理ストレージデバイスに跨って定義されても良い。同様に複数の論理ボリュームが単一の物理ストレージデバイスに収容されても良い。ストレージサブシステムに関するより詳細な記述は、2002/6/5に出願受理された日本出願No. 2002-163705に対する優先出願、標題"Data Storage Subsys

10

20

30

40

50

tem", 2003/3/21に出願受理され本代理人の担当する米国特許NO. 10/394,631に見られ、ここに参考として含まれる。

【0029】

0039 図1Aに戻って、リモートコピーシステム50は更にプライマリホスト110aとセカンダリホスト110bを含む。ホストは、各ストレージサブシステムに対して、データの読み書き要求を発行する。プライマリホストは、プライマリサブシステムのボリュームにアクセスするアプリケーションプログラム102aを持つ。セカンダリホストは、セカンダリサブシステムのボリュームにアクセスするアプリケーションプログラム102bを持つ。プライマリホスト110a又はプライマリサブシステム100aの何れか又は双方が使用不能になれば、アプリケーションプログラム102aはセカンダリホスト110bにフェイルオーバーされ、例えば、エンタプライズビジネスにデータベースを継続使用できる様にする。

10

【0030】

0040 各ストレージサブシステムはデータ保存の為にボリュームを有する。プライマリサブシステムはプライマリボリューム即ちPVOL101aを有する。プライマリボリュームはプライマリホストがアクセスする生産用データを保存する。プライマリボリュームは第一と第二のPVOL101a-1と101a-2で成る。セカンダリサブシステムはセカンダリボリューム即ちSVOL101bを有する。セカンダリボリュームはPVOLの生産用データのバックアップデータを保存する。セカンダリボリュームは第一と第二のSVOL101b-1と101b-2で成る。中間サブシステムは中間ボリュームIVOL101cを有する。中間ボリュームはセカンダリボリュームに保存されるべき生産用データのコピーを一時的に保存する。従って、中間ボリュームは、プライマリボリューム101aに対して“セカンダリボリューム”で、セカンダリボリューム101bに対して“プライマリボリューム”である。中間ボリュームは、第一と第二のIVOL101c-1と101c-2で成る。

20

【0031】

0041 整合性グループ120が中間サブシステム100cにて定義される。整合性グループは、データインテグリティが保証されるボリュームの集合である。ユーザは、典型的に、データインテグリティが望まれる所定のアプリケーション、例えばデータベース等、に関連するデータを格納するボリュームに対して整合性グループを定義する。

30

【0032】

0042 タイマ130が中間サブシステム100cに備えられ、データ要求に対して時刻情報を提供する。例えば、タイマはプライマリサブシステム100aより受信した書込み要求の制御データにタイムスタンプを与える。書込み要求は、通常、所定のストレージロケーションのアドレス指定と共にここに書き込まれるデータを伴う。一般的に書込み要求は、書き込まれるデータ(書込みデータ)とこれに関連するヘッダ又は管理情報を含む。制御データは、所与の要求(読み出し、または書込み)に関連する管理情報を提供し、当該要求のヘッダの一部とも考えられる。

【0033】

0043 カウンタ121は、中間サブシステム100cに配備され、プライマリサブシステム100aから書込み要求を受信した時に制御データに順序番号を付与する。順序番号は書込み順を保持する為に使用される。一般的に、各整合性グループ120は自分のカウンタ121を持ち、整合性グループ又はデータベースに対する書込み順を維持する。通常、中間サブシステムは書込みデータと共にヘッダ又は制御データを受け取り、順序番号はこのヘッダ又は制御データに付加される。

40

【0034】

0044 タイムスタンプ又は順序番号又は双方は、セカンダリサブシステムが書込み順をプライマリサブシステム100aと同じにする為に使用することができる。本発明では、本順序番号が書込み順を保証する為に使用される。

【0035】

50

0045 有効カウンタ 1 2 2 が、セカンダリサブシステム 1 0 0 b で有効化可能な最新順序番号を保存する為に中間サブシステム 1 0 0 c に用意される。

【 0 0 3 6 】

0046 図 5 を参照するに、各ストレージサブシステムは、ボリューム間のペア関係の情報を保存する為に、整合性グループテーブルを備える。プライマリサブシステム 1 0 1 a はプライマリ整合性グループテーブル 5 0 0 a、即ち各々第一と第二のプライマリサブシステムの為の第一と第二のテーブル 5 0 0 a - 1 と 5 0 0 a - 2 を備える。中間サブシステム 1 0 1 c は第一と第二の整合性グループテーブル 5 0 0 c - 1 と 5 0 0 c - 2 を備える。セカンダリサブシステム 1 0 1 b は、セカンダリ整合性グループテーブル 5 0 0 b、即ち各々第一と第二のセカンダリサブシステムの為の第一と第二のテーブル 5 0 0 b - 1 と 5 0 0 b - 2 を備える。

10

【 0 0 3 7 】

0047 各整合性グループテーブルは、グループ ID 欄 G I D 5 1 0、グループ属性欄 G R A T T R 5 2 0、グループ状態欄 G R S T S 5 3 0、中間ポインタ欄 I P T R 5 4 0、及びボリューム情報欄 V O L I N F O # 5 5 0 を有する。

【 0 0 3 8 】

0048 G I D 5 1 0 は整合性グループの識別情報を格納する。識別情報は一般に当該ストレージサブシステムでの固有番号である。G I D 5 1 0 はペアを構成する整合性グループを識別する為に使用される。例えば、二つの G I D 欄が同じ値を持っていれば、二つの整合性グループはペアになっているとされる。

20

【 0 0 3 9 】

0049 G R A T T R 5 2 0 は、当該グループが、プライマリボリュームか又はセカンダリボリュームのどちらとして機能しているかを表示する。例えば、プライマリ又はセカンダリサブシステムの整合性グループは各々 P R I M A R Y 又は S E C O N D A R Y を示す。P R I M A R Y は、当該整合性グループはデータの送信者で、S E C O N D A R Y は、当該整合性グループはデータの受信者であることを示す。従って、中間サブシステムの整合性グループテーブル 5 0 0 c は、プライマリグループに関連する第一グループテーブル 5 0 0 c - 1 は S E C O N D A R Y を、セカンダリグループに関連する第二グループテーブル 5 0 0 c - 2 は P R I M A R Y を示す。

【 0 0 4 0 】

30

0050 G R S T S 5 3 0 は整合性グループの状態を示す。一般的に、整合性グループは C O P Y (コピー中)、P A I R (ペア中)、又は S U S P (サスペンド中)の何れかの状態にある。C O P Y 状態は、現在、データが、一つのボリュームからもう一つのボリュームにコピー中、例えば P V O L から S V O L にコピー中、であることを示す。P A I R 状態は、例えば、P V O L と S V O L はペア間のミラーが完成し、同一の情報を持つことを示す。S U S P 状態は、例えば、P V O L から S V O L へのコピーがサスペンド(中断)している事を示す。

【 0 0 4 1 】

0051 I P T R 5 4 0 は、中間サブシステムでペアと成っている整合性グループテーブルへのポインタを含む。中間サブシステムでの本テーブルは、第一テーブル 5 0 0 c - 1 と第二テーブル 5 0 0 c - 2 を関連つける相応しい値が格納される。即ち、第一テーブル 5 0 0 c - 1 の I P T R 5 4 0 は第二テーブル 5 0 0 c - 2 をポイントし、第二テーブル 5 0 0 c - 2 の I P T R 5 4 0 は第一テーブル 5 0 0 c - 1 をポイントする。プライマリ又はセカンダリサブシステムは、通常唯一つの整合性グループテーブルしか持たない為、本サブシステムでは本テーブルの本欄の値は N U L L である。

40

【 0 0 4 2 】

0052 V O L I N F O # 5 5 0 はボリューム情報を格納する。一ボリュームに一エリアが割り当てられる。ボリューム識別欄 V O L I D 5 5 1 は所定のストレージサブシステム 1 0 0 でのボリューム識別子を表示する。各ボリュームはそのストレージサブシステム 1 0 0 での固有の番号を持つ。ボリューム状態欄 V O L S T S 5 5 2 は、ボリュームの状態

50

即ち、COPY, PAIR, 又はSUSPを示す。シリアル番号欄PAIRDKC553は、ペアを組むストレージサブシステム100の識別子又はシリアル番号を格納する。各ストレージサブシステムは各々固有の番号を持つ。ペアボリューム欄PAIRVOLID554は、PAIRDKC553で指定されるペアを組むストレージサブシステム100のボリューム識別子を格納する。

【0043】

0053 図5は、一実施例による、リモートコピーシステム50での下記のみラー構成を説明する。各プライマリサブシステムは単一のPVOLを持つとする。各プライマリサブシステム100aは、整合性グループテーブルを持ち、ここでは、

- ・0054 GID = 0、
- ・0055 GRATTR = PRIMARYは、当ボリュームがPVOLである事を示し、

・0056 CTGはVOLINF#0で示される如く唯一つのPVOLを持ち、
 ・0057 IPTTR = NULLは、本CTG項目はIVOL101cにて定義されていない事を示す。

【0044】

0058 各セカンダリサブシステム100bは単一のSVOLを持つとする。各セカンダリサブシステム100bは整合性グループテーブルを持ち、ここでは、

- ・0059 GID = 0, は中間サブシステム100cの整合性グループテーブルのGIDと同じであり、
- ・0060 GRATTR = SECONDARYは、当該ボリュームがSVOLであることを示し、

・0061 本CTGはVOLINFO#0に示される如く唯一つのSVOLを含み、
 ・0062 IPTTR = NULLは、本CTG項目はIVOL101cにて定義されていない事を示す。

【0045】

0063 中間サブシステム100cは、プライマリサブシステム100aとセカンダリサブシステム100bに結合する。本中間サブシステム100cは二つの整合性グループテーブルを持つ。ここでは、

- ・0064 GID = 0は、プライマリサブシステム100aの整合性グループテーブルに含まれるのと同じGIDであることを示し、

・0065 GRATTR = SECONDARYは、第一テーブル500c-1に対して表示され、GRATTR = PRIMARYは第二テーブル500c-2に対して表示され、

- ・0066 各テーブルは、VOLINFO#0とVOLINFO#1の二つのボリューム情報欄を含み、

・0067 IPTTRはIVOL101cでの対応する整合性グループテーブルをポイントする。

【0046】

0068 図2を参照するに、本発明の一実施例による同期型リモートコピー方法に関するプロセス200は、プライマリサブシステムと中間サブシステムの間で実行される。本プロセス200は、書込み要求がプライマリホスト110aからプライマリサブシステム100aに到着する事により起動される。

【0047】

0069 プライマリサブシステムが書込み要求を受信すると、本書込みデータを不揮発(又は安定な)記憶領域に格納する(ステップ205)。本不揮発記憶領域は、書込みデータを一時的に記憶する、キャッシュメモリ又は磁気ディスク上に割り当てられた記憶領域で良い。本プライマリサブシステムは次いで、書込みデータを中間サブシステム100cに送信し、中間サブシステム100cからの応答を待つ(ステップ210)。中間サブシステム100cがプライマリサブシステム100aから書込みデータを受信すると、本サブシ

10

20

30

40

50

システムは本書込みデータの制御データの少なくとも一部を生成する(ステップ215)。本実施例では、本生成ステップで、書込み要求のヘッダに対して順序番号を生成しこれを付与する。タイムスタンプが更にヘッダに対して付与されてもよい。

【0048】

0070 一実施例では、制御データには、シリアル番号291(PVOL-ID291)、書込みアドレス292(ADDR292)、長さ293(LEN293)、時刻294、及び順序番号(SEQ295)が含まれる。PVOL-ID291は、プライマリサブシステムと本サブシステム内のPVOLを特定する。プライマリサブシステムと本PVOLはシリアル番号を用いて表示される。

【0049】

0071 ADDR292は、PVOLでの書込みの開始アドレスを示し、データが本位置より書き込めるようにする。LEN293は書込みデータの長さを示す。従って、ADDR292とLEN293で書込みデータの正確な領域が指定される。

【0050】

0072 時刻294は、プライマリサブシステム100aから中間サブシステム100cに書込み要求が到着した時に制御データに付与されるタイムスタンプである。本タイムスタンプは中間サブシステム100cのタイマ130により与えられる。

【0051】

0073 SEQ295は書込み要求に付与された順序番号である。この値は、当該中間サブシステム100cの整合性グループに関連するカウンタ121により提供される。本ステップでは、カウンタ121からの値を+1(つまりカウンタを加算)してカウンタ121に戻し、得られた値を書込み要求の制御データに付加する。

【0052】

0074 プロセス200に戻って、中間サブシステム100cは、プライマリサブシステム100aからの書込みデータとこの制御データを不揮発又は安定な記憶領域に保存する(ステップ220)。第一の応答がプライマリサブシステム100aに送信される(ステップ225)。本第一の応答は、書込みデータが中間サブシステムにコピーされたことを示す。不揮発記憶領域に保存されたデータは、続いてIVOL101cに、より長期的に保存される。プライマリサブシステム100aは、第二の応答をプライマリホスト110aに送信する(ステップ230)。図示されている様に、プロセス200では、第二の応答は、第一の応答がプライマリサブシステムで受信されて始めて返送される為、同期型コピー動作である。

【0053】

0075 図3は、本発明の一実施例による、中間サブシステム100cとセカンダリサブシステム100bとの間の非同期型コピー動作を行うプロセス300を示す。プロセス300の特徴は、プライマリホスト110aからプライマリサブシステム100aに発行された書込み順序に従って、セカンダリサブシステム100bが受信した書込みデータの書込み順序を維持する事である。このデータインテグリティの保証は、セカンダリサブシステムが受信した書込みデータを先ずINVALIDとして表示し、次いで有効化することを含む。本有効化プロセスは図4に関連して説明する。

【0054】

0076 プロセス300は一定間隔、例えば10秒毎に起動される。中間サブシステム100cの不揮発記憶領域に保存された書込みデータとその制御データは、セカンダリサブシステム100bに送信される(ステップ301)。本発明では、送信する書込みデータは対応する制御データのSEQ295の昇順に従って選択される。

【0055】

0077 セカンダリサブシステム100bは、書込みデータとその制御データを保存し、当該データをINVALIDとしてマークする(ステップ305)。セカンダリサブシステムは、中間サブシステム100cに応答を返す(ステップ310)。

【0056】

10

20

30

40

50

0078 本応答を受け取ると、中間サブシステム 100c は、有効カウンタ 122 の現在値が送信済みの制御データの SEQ 295 より小さいかを判定する(ステップ 315)。小さければこの値は更新される(ステップ 320)。プロセス 300 は次いで、再度起動される迄、一定時間休眠する(ステップ 325)。有効カウンタ 122 は、かくして、セカンダリサブシステム 100b で有効化可能な書込みデータの最新の順序番号を保持する事になる。別の言い方をすれば、有効カウンタ 122 は、セカンダリサブシステムに保存された INVALID ID データ(有効化待ちのデータ)の最大順序番号数を示す、と言える。

【0057】

0079 図 4 は、本発明の一実施例による、セカンダリサブシステム 100b でのデータ有効化の為にプロセス 400 を示す。本プロセス 400 は PREPARE 及び VALIDATE 状態即ちフェーズを使用する。PREPARE フェーズでは、全セカンダリサブシステムは、ロールバック動作に必要な回復用情報を保持したまま、中間サブシステム 100c より受信したデータを有効化する為の全ての必要な動作を実行する。セカンダリサブシステム 100b がプロセス 400 の処理を実行できない条件が発生すれば、本回復情報を用いてロールバックする。VALIDATE フェーズでは、セカンダリサブシステム 100b は、データを有効化し回復情報を破棄する。これら二つのフェーズと回復機構は、所定のデータに対する有効化プロセス 400 がセカンダリサブシステム 100b の何れかで障害になるかもしれない為、有益なものである。

【0058】

0080 本発明の一実施例では、プロセス 400 は中間サブシステムにて一定時間間隔で起動される(ステップ 401)。有効カウンタ 122 は、その値が前回のプロセス以降変化があったか否かを決定するために、チェックされる(ステップ 402)。変化がなかったら、プロセスはステップ 401 に戻り、所定の時間休眠する。有効カウンタ 122 が変化しなかったと言うことは、前回のセッション以後、セカンダリサブシステムには書込みデータが送られていないことを示す。値が変化していれば、プロセスは先に進む。

【0059】

0081 中間サブシステム 100c は、有効カウンタ 122 の現在値と共に、全てのセカンダリサブシステム 100b に対して PREPARE 要求を送信する(ステップ 403)。この要求により、セカンダリサブシステム 100b は、送られてきた有効カウンタ 122 以下の順序番号を持つデータについて有効化作業を開始する。中間サブシステム 100c は、次いで、セカンダリサブシステム 100b からの応答を待つ(ステップ 406)。

【0060】

0082 セカンダリサブシステムでは、PREPARE 要求の受信を契機に、INVALID ID としてマークされ、有効カウンタの現在値以下の順序番号を割り当てられたデータの有効化作業を開始する(ステップ 405)。現在値は基準値として参照される。データは PREPARED としてマークされる(ステップ 410)。上記のステップには、データを不揮発ストレージの一時的領域から長期的領域にコピーし、本一時的領域を開放し、制御データを開放し、INVALID ID 及び PREPARED 状態に関する情報を更新することを含む。

【0061】

0083 データを回復可能に維持する為に、セカンダリサブシステムは旧データに上書きしないで新データが有効化される迄、旧データを保持する。一つの方法は、新データを永久ストレージ(例えばディスク)の最終ターゲットに書き込まないことである。セカンダリサブシステム 100b は、上記が成功裏に実行できたらデータを PREPARED としてマークする。

【0062】

0084 セカンダリサブシステムは、全ての指定された書込みデータが PREPARED としてマークされたか否かを判定する(ステップ 415)。判定が成立すれば、セカンダリサブシステム 100b は、中間サブシステム 100c に全ての指定されたデータは PREPARED としてマークされたことを示す応答を返す(ステップ 420)。セカンダリサブ

10

20

30

40

50

システム 1 0 0 b は、タイマをセットして、V A L I D A T E 要求が中間サブシステム 1 0 0 c から受信されるのを待つ(ステップ 4 2 1)。

【 0 0 6 3 】

0085 中間サブシステム 1 0 0 c が、全てのセカンダリサブシステム 1 0 0 b から、全ての指定された書込みデータが P R E P A R E D としてマークされたことを示す応答を受信したら、中間サブシステムはセカンダリサブシステム 1 0 0 b に V A L I D A T E 要求を送信する(ステップ 4 2 5)。中間サブシステムは、次いで、セカンダリサブシステムから V A L I D A T E 要求に対する応答を待つ(ステップ 4 2 6)。

【 0 0 6 4 】

0086 セカンダリサブシステムに於いては、各サブシステムは、P R E P A R E D としてマークされたデータを V A L I D に更新して、当該データの回復に関する情報を破棄する(ステップ 4 3 0)。旧データは通常この時点で破棄される。各セカンダリサブシステムは中間サブシステム 1 0 0 c に応答を返す(ステップ 4 3 5)。

【 0 0 6 5 】

0087 ステップ 4 1 5 に戻って、幾つかの書込みデータを P R E P A R E D としてマークすることに失敗した場合には、F A I L U R E 通知が中間サブシステムに送信される(ステップ 4 1 5)。中間サブシステムは、以下により詳細に記述される、エラー処理を実行する(ステップ 4 2 4)。

【 0 0 6 6 】

0088 プロセス 4 0 0 は、参照順序番号、即ち有効カウンタ 1 2 2 の値、以下の順序番号を持つ書込みデータを有効化するのに用いられる。即ち、セカンダリサブシステム 1 0 0 b の S V O L 1 0 1 b の集合は、S V O L 1 0 1 b を通して整合しており、参照順序番号以下の順序番号で送信された全てのデータが受信され有効化されたことをセカンダリサブシステムが表示した以降には、書込みデータの損失は存在しない。

【 0 0 6 7 】

0089 時によっては、プロセス 4 0 0 は失敗し回復処理が必要になる。障害の例は、ステップ 4 2 4 で、書込みデータの P R E P A R E 操作のときに発生する。本実施例では、回復処理は中間サブシステムで実施される。

【 0 0 6 8 】

0090 一般的に、中間サブシステムは、ステップ 4 2 4 でのエラーか、セカンダリサブシステムに要求を出してからタイムアウト(ステップ 4 0 6 か 4 2 6)で障害発生を検出する。

【 0 0 6 9 】

0091 障害が P R E P A R E フェーズで発生した場合には、中間サブシステム 1 0 0 c は、セカンダリサブシステム 1 0 0 b に A B O R T 要求を送信し、セカンダリサブシステムはロールバックを行い、P R E P A R E D としてマークされたデータの回復は、本データとこれの制御データを I N V A L I D に戻すことにより開始される。

【 0 0 7 0 】

0092 V A L I D A T E フェーズの場合には、障害はセカンダリサブシステムからの応答が既定時間内に受信できない場合に発生する。このような事態は起きそうもないが、もし実際に起きたなら、中間サブシステムは V A L I D A T E 要求を再発行する。

【 0 0 7 1 】

0093 中間サブシステム 1 0 0 c が、P R E P A R E か V A L I D A T E フェーズ中にダウンした場合には、セカンダリサブシステム 1 0 0 b は、リポートされた中間サブシステム 1 0 0 c かホスト 1 0 0 の一つからのコマンドを待つ。一般に、リポート又は回復した中間サブシステム 1 0 0 c は、プロセス 4 0 0 を最初からではなく、最新の成功ステップから再開する。

【 0 0 7 2 】

0094 図 8 は、本発明の一実施例に従い、プロセス 4 0 0 が障害なく成功裏に実行された場合のプロセスフローを示す。第一の状態 8 0 0 で示すように、中間サブシステム 1 0

10

20

30

40

50

0 c は、セカンダリサブシステム 1 0 0 b - 1 , 1 0 0 b - 2、及び 1 0 0 b - 3 に P R E P A R E 要求を発行する(ステップ 8 0 0)。本 P R E P A R E 要求は、1 0 4 2 までの順序番号の書込みデータを P R E P A R E することをセカンダリサブシステムに要求する。

【 0 0 7 3 】

0095 第一の状態 8 0 0 で P R E P A R E 要求を受信する前には、セカンダリサブシステム 1 0 0 b の第一、第二、及び第三の整合性グループ 1 2 0 - 1 , 1 2 0 - 2、及び 1 2 0 - 3 は、1 0 2 4 までの順序番号の全ての書込みデータは P R E P A R E D になっていることを示す。1 0 1 0 までの順序番号の全ての書込みデータは V A L I D になっている。1 0 2 4 より高い順序番号の書込みデータは I N V A L I D が表示されている。例えば、第一の整合性グループ 1 2 0 - 1 は、サブシステム 1 0 0 b - 1 は I N V A L I D の順序番号 1 0 2 9 の書込みデータを持っている、ことを示す。

10

【 0 0 7 4 】

0096 第二の状態 8 1 0 は、P R E P A R E 要求が完了した後の整合性グループを示す。第一の整合性グループ 1 2 0 - 1 は順序番号 1 0 4 2 までの書込みデータは P R E P A R E されていることを示す。従って、順序番号 1 0 2 9 までの書込みデータは最早 I N V A L I D ではなく、整合性グループ 1 2 0 - 1 は、サブシステム 1 0 0 b - 1 は I N V A L I D データを持っていないことを示す。整合性グループは、順序番号 1 0 1 0 迄の書込みデータは V A L I D のままであることを示す。

【 0 0 7 5 】

0097 中間サブシステム 1 0 0 c は V A L I D A T E 要求をセカンダリサブシステムに送信し、P R E P A R E された書込みデータは本要求により V A L I D になる。第三の状態 8 2 0 は、V A L I D A T E 要求がセカンダリサブシステム 1 0 0 b にて実行された後の整合性グループを示す。本整合性グループは、全ての P R E P A R E された書込みデータは V A L I D になった事を示す。即ち、1 0 4 2 迄の順序番号の全ての書込みデータは V A L I D になっている。

20

【 0 0 7 6 】

0098 図 9 は、プロセス 4 0 0 で障害が発生した場合を示す。第一の状態 9 0 0 は、図 8 の第一の状態 8 0 0 のミラーイメージである。中間サブシステム 1 0 0 c はセカンダリサブシステム 1 0 0 b に P R E P A R E 要求を発行する。

【 0 0 7 7 】

30

0099 第二の状態 9 1 0 は、P R E P A R E 処理がセカンダリサブシステムの 1 つ、即ちセカンダリサブシステム 1 0 0 b - 2 で成功裏には実行できなかったことを示す。セカンダリサブシステム 1 0 0 b - 2 は、要求されたすべての、即ち順序番号 1 0 4 2 迄の、書込みデータを成功裏には P R E P A R E できなかった。セカンダリサブシステム 1 0 0 b - 2 の整合性グループ 1 2 0 - 2 は、順序番号 1 0 3 8 迄しか書込みデータは P R E P A R E されていないことを示す。本 P R E P A R E 処理はセカンダリサブシステム 1 0 0 b - 1 と 1 0 0 b - 3 では成功裏に終了している。これらの整合性グループは、1 0 4 2 迄は P R E P A R E されている。

【 0 0 7 8 】

0100 セカンダリサブシステム 1 0 0 b - 2 の障害の結果、中間サブシステムはセカンダリサブシステムに参照順序番号 1 0 3 8 (1 0 4 2 でなく) で V A L I D A T E 要求を発行する。本 V A L I D A T E 要求は、セカンダリサブシステム 1 0 0 b - 2 での書込みデータの P R E P A R E 処理での問題を反映して、修正されている。

40

【 0 0 7 9 】

0101 第三の状態 9 2 0 は、修正 V A L I D A T E 要求がセカンダリサブシステム 1 0 0 b で成功裏に実行できたことを示す。整合性グループ 1 2 0 - 1 と 1 2 0 - 3 は、依然として、P R E P A R E D だが V A L I D でない書込みデータを持っていることを示す。然しながら、1 2 0 - 2 の整合性グループは、全ての P R E P A R E D データは V A L I D になっていることを示す。V A L I D になった書込みデータは、順序番号 1 0 3 8 迄のものであることを示す。

50

【 0 0 8 0 】

0102 図 1 0 A (図 1 0 (a)) と 1 0 B (図 1 0 (b)) は、本発明の一実施例に従う、プライマリサイトからセカンダリサイトへのフェイルオーバー処理を示す。フェイルオーバーはプライマリサイトがユーザに対して使用不能になったときに実行される。これは、プライマリサブシステム 1 0 0 a かプライマリホスト 1 1 0 a の何れか又は何れも使用不能になったときに発生する。プライマリサイトが障害になる、又は保守の為に一時的にオフラインにされた場合には、セカンダリサイトがアプリケーションを実行可能になる。

【 0 0 8 1 】

0103 一般的に、セカンダリサブシステムのボリューム、即ち S V O L 1 0 1 b は、プライマリサブシステム 1 0 0 a にペア又はミラーしている場合には、書込み不能である。従って、本発明の一実施例に拠れば、フェイルオーバーではプライマリサブシステムへのミラー動作はサスペンドする必要がある。

10

【 0 0 8 2 】

0104 図 1 0 (a) は、プライマリサイトのみが使用不能で、中間サブシステム 1 0 0 c が使用可能な場合でのフェイルオーバー処理を示す。この状況では、I V O L 1 0 1 c がプライマリサブシステムから受信した全てのデータを保持している為に、データは損失していない。これらのデータは、中間サブシステム 1 0 0 c が S V O L 1 0 1 b に転送する。この後に、セカンダリサブシステムは使用可能になる。

【 0 0 8 3 】

0105 最初にユーザが、セカンダリホスト 1 0 0 b を使用してフェイルオーバーを起動する(ステップ 1 0 0 0)。セカンダリサブシステム 1 0 0 b は F A I L O V E R 要求を中間サブシステム 1 0 0 c に送信する。中間サブシステム 1 0 0 c は全ての未転送データをセカンダリサブシステム 1 0 0 b に送信する(ステップ 1 0 0 5)。本発明の一実施例では、中間サブシステムからセカンダリサブシステムにデータを送信するのに、プロセス 3 0 0 が使用される。中間サブシステムは、全ての未転送データがセカンダリサブシステムに転送されたら、ペア状態を P S U S に変更する(ステップ 1 0 0 6)。P S U S はペア動作がサスペンドしていることを示す。ひとたびセカンダリサブシステムが F A I L O V E R 要求の完了を知らされると、セカンダリサブシステム 1 0 0 b は、ペア状態を P S U S に変更する(ステップ 1 0 1 0)。

20

【 0 0 8 4 】

0106 図 1 0 (b) は中間サブシステム 1 0 0 c が使用不能の場合のフェイルオーバー動作を示す。プライマリサイトは使用不能でも可能でもよい。この状況では、中間サブシステムが、セカンダリサブシステム 1 0 0 b への未転送データを保持している可能性がある為、データは損失する可能性がある。ユーザは、セカンダリホスト 1 1 0 b を使用して、フェイルオーバー動作を起動して、セカンダリサブシステムに F A I L O V E R 要求を送信する(ステップ 1 0 2 0)。中間サブシステム 1 0 0 c が使用不能の為、F A I L O V E R 要求は失敗して、セカンダリサブシステム 1 0 0 b はセカンダリホストにエラーメッセージを返す(ステップ 1 0 2 5)。

30

【 0 0 8 5 】

0107 セカンダリホスト 1 1 0 b は、全てのセカンダリサブシステムの整合性グループからある種の情報を取得する(ステップ 1 0 3 0)。これらの情報には、I N V A L I D データ、P R E R A R E D データ、及び V A L I D データの各順序番号が含まれる。全ての P R E R A R E D データは、プロセス 4 0 0 により V A L I D データに更新する(ステップ 1 0 3 5)。

40

【 0 0 8 6 】

0108 セカンダリホスト 1 1 0 b は、有効化プロセスが成功したことを確認後、セカンダリサブシステム 1 0 0 b のペア状態を P S U S に変更開始する(ステップ 1 0 4 0)。セカンダリサブシステムの状態は P S U S に変更される(ステップ 1 0 4 5)。

【 0 0 8 7 】

0109 以上で説明したように、リモートコピーシステム 5 0 は三つのデータセンタを持

50

っている。プライマリと中間サブシステムの双方が同時に障害にならない限り、データは損失しない。プライマリとセカンダリのデータセンタが各々複数のサブシステムを持っていてもデータの整合性は保証される。中間サブシステムは、プライマリサブシステムからの全ての書込みデータを受信しセカンダリサブシステムに転送する。中間サブシステムは、複数のプライマリサブシステムからの書込みデータの書込み順を保持する為の管理情報（例えば、順序番号）を生成する。

【0088】

0110 図6は、本発明のもう一つの実施例に従う、リモートコピーシステム50'を示す。システム50'は、図1のシステム50に比べて、中間サブシステム100c'が減少したボリューム数で構成されている。アプリケーションが中間サブシステムにフェイルオーバーしないなら、中間サブシステムはプライマリサイトのプライマリボリュームより少数のボリュームで構成することができる。ボリューム数を減らす事でコストを低減できる。中間サブシステム100c'は、プライマリサブシステム100aから受信したデータを保存するために、一つのボリュームを備えている。中間サブシステム100c'がボリュームを持たないで、一時的にデータを保存する為にキャッシュメモリを使用する、実施例も存在する。

【0089】

0111 中間サブシステム100c'はジャーナルボリューム(JNLVOL)700を持つ。JNLVOL700はIVOL101cに相当する。システム50のIVOL100cとは異なり、JNLVOL700はPVOL101aのミラーではない。JNLVOL700は、"Remote Copy System"と標題を有し、参照としてここに含まれる米国特許NO.10/602,223(2003/6/23出願受理)にてより詳しく説明されているように、プライマリサブシステム100aから受信した書込みデータとその制御データの為のバッファ又は一時的記憶領域として働く。

【0090】

0112 図7は、本発明の一実施例に従う、リモートコピーシステム50'に対する整合性グループテーブル702を示す。本整合性グループテーブル702は、ボリューム情報(VOLINFO#)750、ジャーナルボリューム情報(JNL__VOLID#)760、ジャーナル入力情報(JNLIN)770、及びジャーナル出力情報(JNLOUT)780を含む。VOLINFO#750は、中間サブシステム100cで定義されるIVOL101cは存在しない事を示す為に、NULLを保存するボリュームID(VOLID)751を含む。同様に、プライマリサブシステム100aとセカンダリサブシステム100bで定義される、PAIRVOLID(図示していない)は、関連する中間サブシステムで定義されるIVOL101cが存在しない為、NULLを表示する。

【0091】

0113 JNL__VOLID760は、プライマリサブシステム100aより送信された書込みデータと制御データを保存する為に使用されるボリュームを識別する情報を含む。JNLIN770は、JNL__VOLID760エントリの識別子とジャーナルボリューム700でのアドレスを含む。JNLIN770は、次の書込みデータと制御データのセットが保存されるアドレスを表示する。JNLOUT780は、JNL__VOLID760エントリの識別子とジャーナルボリューム700でのアドレスを含む。JNLOUT780は、次の書込みデータと制御データのセットがセカンダリサブシステム100bに送信されるソースアドレスを表示する。有効なJNL__VOLID760エントリが存在しない(例えば、全てのJNL__VOLID760がNULL)場合には、整合性グループにはジャーナルボリューム700は割り当てられていない。

【0092】

0114 リモートコピーシステム50'では、図2のプロセス200と実質的に同じプロセスを使用して同期コピー動作を実行する。一つの相違点はステップ220関するもので整合性グループテーブルを参照して不揮発記憶媒体を選択することである。本操作は、(1) IVOL101c(即ち、VOLID751!=NULLか)とJNLVOL700(

10

20

30

40

50

即ち、JNL__VOLID760!=NULLか)が存在するかをチェックする;(2)(1)の結果に従い、下記の手順を実行する。

【0093】

・0115 プライマリサブシステム100aから受信した書込みデータと制御データを中間サブシステム100cのキャッシュメモリに格納する、

・0116 JNLVOL700が存在すれば、書込みデータと制御データをJNLIN770で指定されるJNLVOL700のアドレスに格納する、

・0117 IVOL101cが存在すれば、書込みデータをIVOL101cに格納する。

【0094】

0118 一般的に、中間サブシステム100cは、同期コピー動作を促進する為に、書込みデータをキャッシュメモリに格納し次第、プライマリサブシステム100aに応答を返す。

【0095】

0119 リモートコピーシステム50'では、図3のプロセス300と実質的に同じプロセスを使用して非同期コピー動作を実行する。相違点の一つは、図3のステップ301に関し、データをセカンダリサブシステム100bに送信すべく選択する場合に、

0120 (1)IVOL101c(即ち、VOLID751!=NULLか)とJNLVOL700(即ち、JNL__VOLID760!=NULLか)が存在するかをチェックし、

0121 (2)(1)の結果に従い、下記の手順を実行する。

【0096】

・0122 キャッシュメモリから最小の順序番号を持つデータを選択する。

【0097】

・0123 キャッシュメモリにデータが存在せず、更にJNLVOL700が存在すれば、JNLVOL700からデータを抽出する。更に、JNLVOL700が存在しなければ、セカンダリサブシステムはミラー処理を維持できず、状態をPSUSにする。

【0098】

0124 リモートコピーシステム50'は、(1)プライマリサブシステム100aからの書込みデータとその制御データを不揮発のランダムメモリ(NVRAM)のみに格納し、(2)書込みデータとその制御データをJNLVOL700に格納する、ことによって、中間サブシステムでのボリューム数を低減できる。NVRAMのみを使用するなら中間サブシステムには、ボリュームは不要である。このような構成では、NVRAMの容量が限られ(例えば数GB)、又NVRAMの容量を増加する事は高価につく為、プライマリボリュームのミラーを維持する事は困難である。NVRAMと共にJNLVOL700を使用することは、書込みデータとその制御データをJNLVOL700に保存できる為に、より受け入れやすい構成である。この後者の構成でも、必要なボリュームはIVOL101cに比べて低減できる。

【0099】

0125 図11は、本発明の一実施例に従う、ミラー動作のサスペンド終了後の再同期時間を最小にするリモートコピーシステム50"を示す。システム50"は、プライマリサブシステムに配備される複数のビットマップ1100a、セカンダリサブシステムに配備される複数のビットマップ1110b、及び中間サブシステムに配備されるビットマップ1100cと1110cを備える。

【0100】

0126 一般的に、ミラー動作は、プライマリと中間サブシステムの間、又はセカンダリと中間サブシステムの間で結合が障害になるか、人的に中断されると、サスペンドされる。障害が回復するかユーザがサスペンドしたミラーボリュームの再同期の指示をしたときに、ミラー動作は回復するか再同期される。再同期に於いては、再同期時間を最小にする為に、サスペンド後に変更があったデータのみが、プライマリからセカンダリサブシステムにコピーされる。本実施例では、この目的の為にビットマップを使用する。同様に、サ

10

20

30

40

50

スPEND後に変更があったデータのみが中間からセカンダリサブシステムにコピーされる。

【 0 1 0 1 】

0127 ミラーボリュームの各ペアに二つのビットマップが用意される。ビットマップ 1 1 0 0 はプライマリボリュームにアサインされ、ビットマップ 1 1 1 0 はセカンダリボリュームにアサインされる。従って、本発明の一実施例では、ビットマップ 1 1 0 0 a はプライマリサブシステムの各プライマリボリュームにアサインされる。対応するビットマップ 1 1 1 0 c は中間サブシステムでのセカンダリボリューム(又はセカンダリボリュームとして稼動している中間ボリューム 1 0 1 c)にアサインされる。

【 0 1 0 2 】

0128 中間サブシステムでは、中間ボリュームは、セカンダリサブシステムに対してプライマリボリュームである為、ビットマップ 1 1 0 0 c を有する。セカンダリサブシステムは、中間サブシステムのビットマップ 1 1 0 0 c に対応するビットマップ 1 1 1 0 b を有する。

【 0 1 0 3 】

0129 ビットマップ 1 1 0 0 と 1 1 1 0 は、ミラーがサスペンドしてからのボリュームになされた変更を追跡する。これらビットマップは、対応するミラーがサスペンドして初めて活性化する。例えば、ミラー(PVOL 1 0 1 a、IVOL 1 0 1 c)のみがサスペンドしている場合には、PVOL 1 0 1 aのビットマップ 1 1 0 0 aとIVOL 1 0 1 cのビットマップ 1 1 1 0 cのみが活性化され、IVOL 1 0 1 cのビットマップ 1 1 0 0 cとSVOL 1 0 1 bのビットマップ 1 1 1 0 bは非活性のままである。

【 0 1 0 4 】

0130 再同期化に於いては、全ての活性化しているビットマップを、再同期化するソースストレージサブシステムに転送し、次いでマージ(OR)される。ソースストレージサブシステムは、そのときのアプリケーションに従って、プライマリサブシステムか中間サブシステムである。再同期でのソースストレージサブシステムは、マージされたビットマップに基づいてデータをコピーする。例えば、図 1 2 に於いては、PVOL 1 0 1 aとIVOL 1 0 1 c間のミラーは継続され、IVOL 1 0 1 cとSVOL 1 0 1 bがサスペンドしている。再同期化では、中間サブシステム 1 0 1 c がセカンダリサブシステムのビットマップ 1 1 1 0 b を取得して、これをIVOL 1 0 1 cのビットマップ 1 1 0 0 cとマージ(OR)する。マージされたビットマップに基づく差分データがIVOL 1 0 1 cからSVOL 1 0 1 bにコピーされる。リモートコピーシステム 5 0 "では、2セットの異なるビットマップが用意されているので、第一のミラー(PVOL 1 0 1 aとIVOL 1 0 1 c)は、第二のミラー(IVOL 1 0 1 cとSVOL 1 0 1 b)とは独立に再同期化される。

【 0 1 0 5 】

0131 これまでの詳細な記述は、本発明の具体的実施例を説明する為のもので、制約を加えることを意図しているものではない。本発明の範囲内で多数の改造と変形が可能である。従って、本発明は、添付の請求範囲によって定義されるものである。

【図面の簡単な説明】

【 0 1 0 6 】

【図 1 (a)】0019 図 1 (a)は、本発明の一実施例による、複数のプライマリサブシステムを持つリモートコピーシステムを示す。

【図 1 (b)】0020 図 1 (b)は、データの読み書き要求を処理するストレージコントローラと、本書込み要求に従って、データを保存する記録媒体を持つストレージユニットを含むストレージサブシステムの好適な一実施例を示す。

【図 2】0021 図 2 によれば、本発明の一実施例による、同期的リモートコピー方法に関するプロセスは、プライマリサブシステムと中間サブシステムとの間で実行される。

【図 3】0022 図 3 は、本発明の一実施例による、中間サブシステムとセカンダリサブシステムとの間の非同期リモートコピー動作の為のプロセスを説明する。

【図 4】0023 図 4 は、本発明の一実施例による、セカンダリサブシステムでのデータ有効化の為のプロセスを説明する。

【図 5】0024 図 5 は、本発明の一実施例による、リモートコピーシステムの為のミラー構成を説明する。

【図 6】0025 図 6 は、本発明のもう一つの実施例による、リモートコピーシステムを説明する。

【図 7】0026 図 7 は、本発明の一実施例による、リモートコピーシステムの為の整合性グループテーブルを説明する。

【図 8】0027 図 8 は、本発明の一実施例による、障害なしに成功裏に実行されたプロセスのプロセスフローを示す。

10

【図 9】0028 図 9 は、プロセス中に障害が発生した場合のプロセスフローを示す。

【図 10 (a)】0029 図 10 (a) は、プライマリサイトのみが使用不能で、中間サブシステムが使用可能な場合のフェイルオーバー動作を示す。

【図 10 (b)】0030 図 10 (b) は、中間サブシステムが使用不能になった場合のフェイルオーバー動作を示す。

【図 11】0031 図 11 は、本発明の一実施例による、ミラー動作がサスペンドした後の再同期時間を最小限にするリモートコピーシステムを示す。

【図 12】0032 図 12 は、P V O L 1 0 1 a と I V O L 1 0 1 c がミラー継続し、I V O L 1 0 1 c と S V O L 1 0 1 b 間のミラーがサスペンドした状態を説明する。

【符号の説明】

20

【0107】

1 0 0 a-1、1 0 0 a-2、1 0 0 b-1、1 0 0 b-2、1 0 0 c... D K C

1 0 1 a-1、1 0 1 a-2... P V O L

1 0 1 b-1、1 0 1 b-2... S V O L

1 0 1 c-1、1 0 1 c-2... I V O L

1 0 2 a、1 0 2 b... アプリケーション

1 1 0 a、1 1 0 b... ホスト

1 2 1 ... カウンタ

1 2 2 ... 有効カウンタ

1 3 0 ... タイマ

30

5 0 ... リモートコピーシステム

6 0 ... ストレージサブシステム

6 2 ... ストレージコントローラ

6 3 ... ストレージユニット

6 4 ... ホストチャネルアダプタ

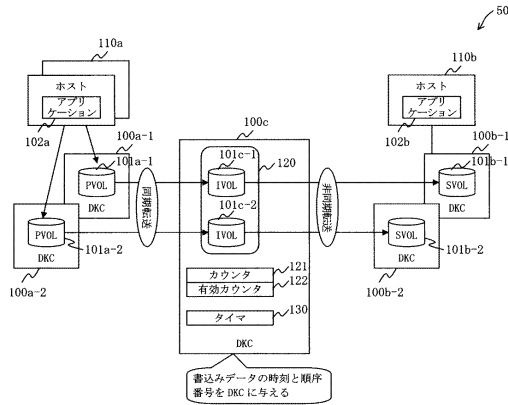
6 6 ... ストレージチャネルアダプタ

6 8 ... ディスクアダプタ

7 0 ... キャッシュメモリ

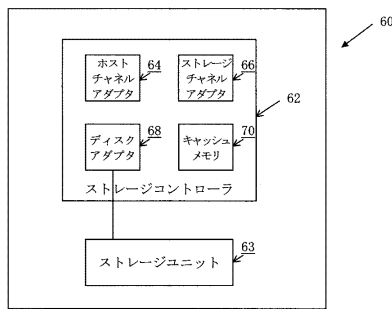
【図 1 (a)】

【図 1 (a)】



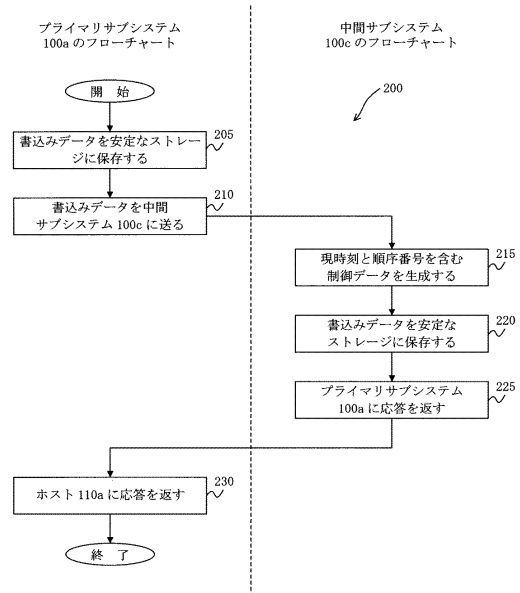
【図 1 (b)】

【図 1 (b)】



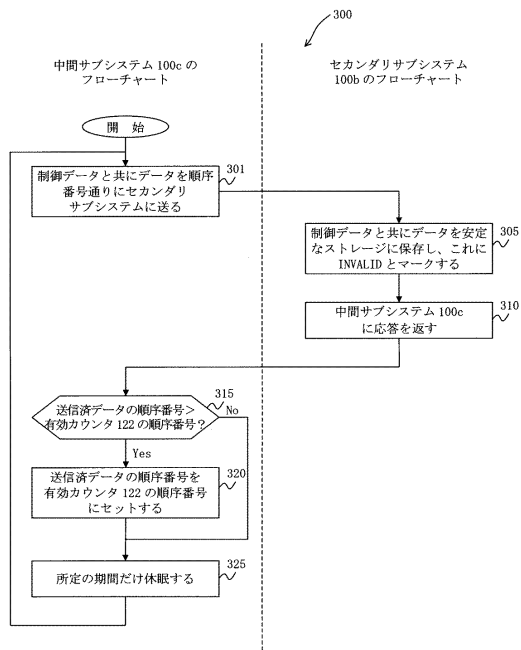
【図 2】

【図 2】



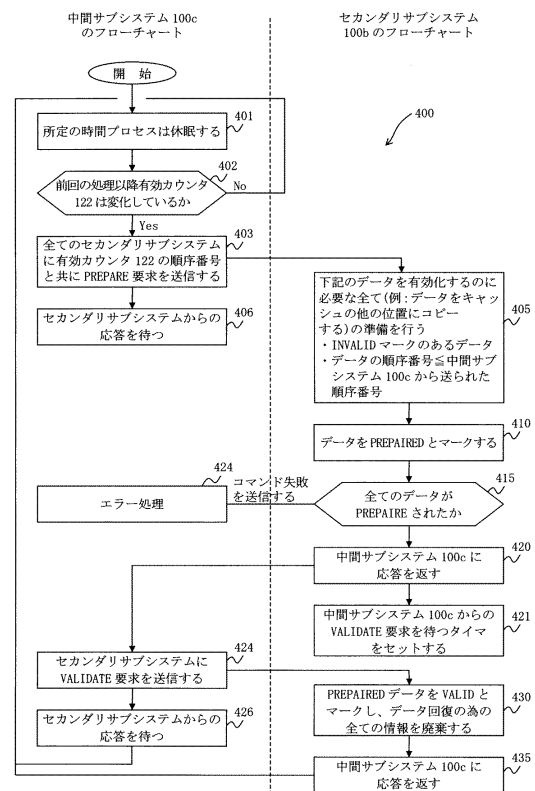
【図 3】

【図 3】

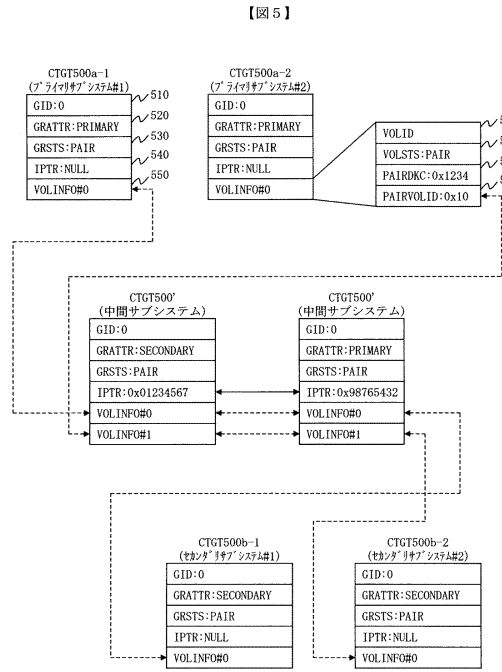


【図 4】

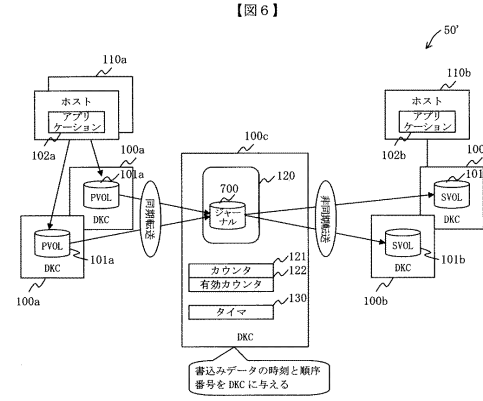
【図 4】



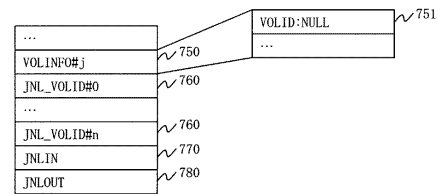
【図 5】



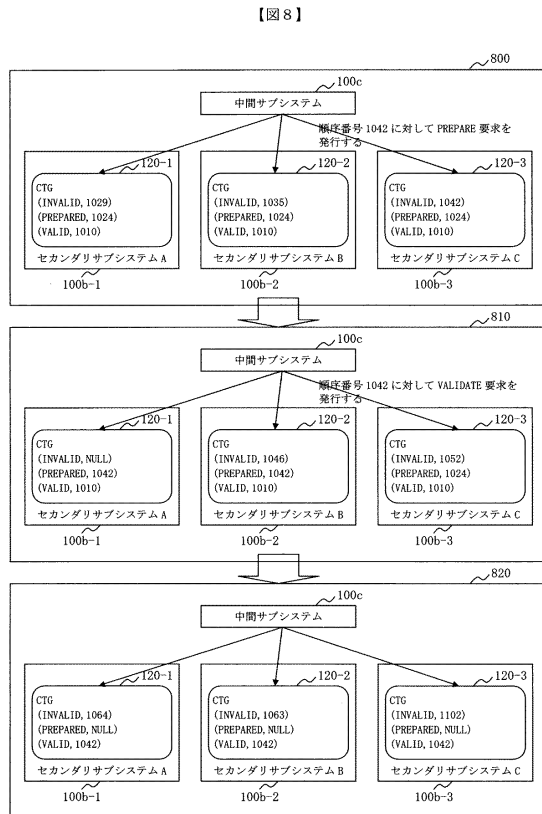
【図 6】



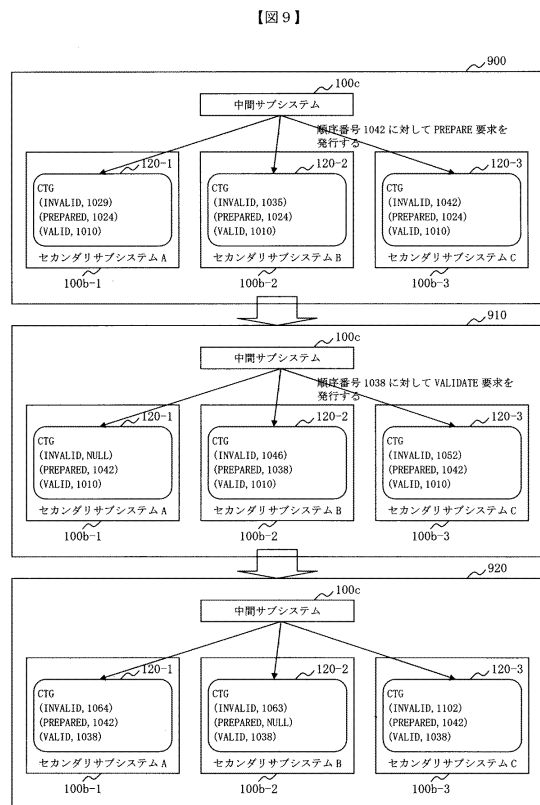
【図 7】



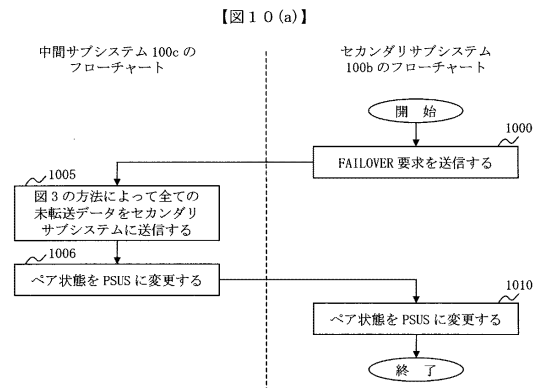
【図 8】



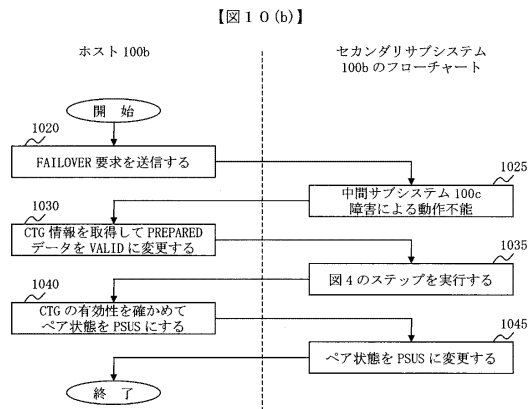
【図 9】



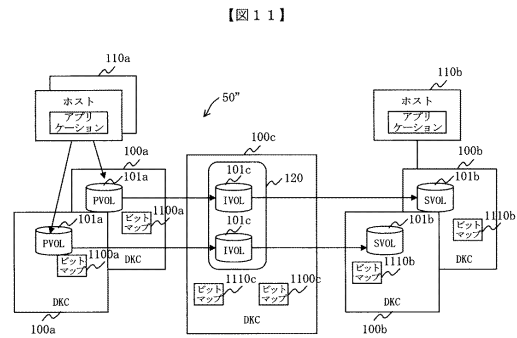
【図 10 (a)】



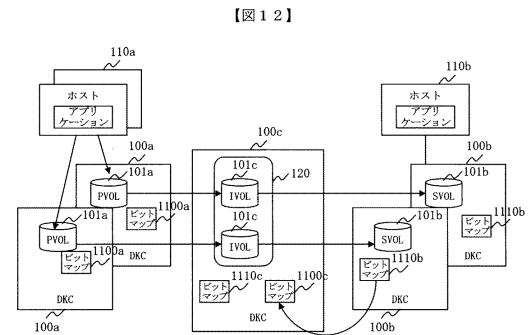
【図 10 (b)】



【図 11】



【図 12】



フロントページの続き

(56)参考文献 特開2000-305856(JP,A)
特開2003-122509(JP,A)
特開2001-282628(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 3/06