

(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2006-285991  
(P2006-285991A)

(43) 公開日 平成18年10月19日(2006.10.19)

(51) Int. Cl.	F I	テーマコード (参考)
<b>G06F 12/00 (2006.01)</b>	G06F 12/00 531R	5B065
<b>G06F 3/06 (2006.01)</b>	G06F 3/06 304F	5B082
	G06F 12/00 533J	

審査請求 未請求 請求項の数 30 O L (全 17 頁)

(21) 出願番号	特願2006-89117 (P2006-89117)	(71) 出願人	390009531 インターナショナル・ビジネス・マシー ズ・コーポレーション INTERNATIONAL BUSIN ESS MASHINES CORPO RATION アメリカ合衆国10504 ニューヨーク 州 アーモンク ニュー オーチャード ロード
(22) 出願日	平成18年3月28日 (2006.3.28)	(74) 代理人	100086243 弁理士 坂口 博
(31) 優先権主張番号	11/093521	(74) 代理人	100091568 弁理士 市位 嘉宏
(32) 優先日	平成17年3月30日 (2005.3.30)	(74) 代理人	100108501 弁理士 上野 剛史
(33) 優先権主張国	米国 (US)		

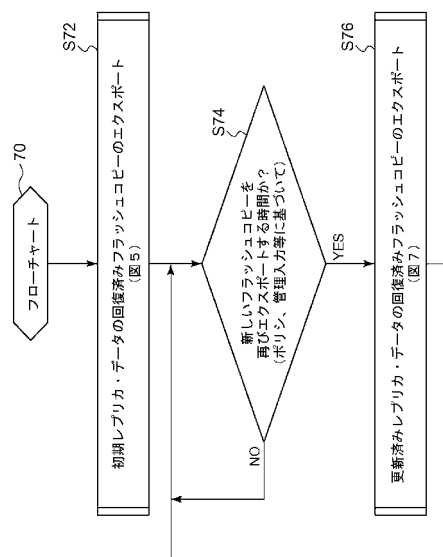
最終頁に続く

(54) 【発明の名称】 ブロック複製によってファイルシステムの可用性を高めるための方法、サーバおよびプログラム

(57) 【要約】

【課題】 ブロック・レベルの複製方式の利点を活用してファイル・システム複製方式を実施するための技法を提供する。

【解決手段】 ソース・サイトは、ソース・サーバおよびソース・ストレージ・システムを用い、ターゲット・サイトは、ターゲット・サーバおよびターゲット・ストレージ・システムを用いる。ソース・サーバは、ソース・ストレージ・システムのソース論理ユニット内に記憶されているデータを操作する。ターゲット・サーバは、レプリカ・ボリュームのフラッシュコピーを作成し、ストレージ・ボリュームの複製である第1のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、レプリカ・ボリュームのフラッシュコピーの回復を実行する。ターゲット・サーバは、更に、レプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートし、これによって、レプリカ・ボリュームの回復したフラッシュコピーをクライアントに利用可能とする。



## 【特許請求の範囲】

## 【請求項 1】

ターゲット・サーバのプロセッサに、

第 1 のレプリカ・ボリュームのフラッシュコピーを作成し、第 1 のストレージ・ボリュームの複製である前記第 1 のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、前記第 1 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行するステップと、

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートして、前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーをクライアントに利用可能とするステップと、  
を実行させるための、プログラム。

10

## 【請求項 2】

前記第 1 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された初期ソース・データを含み、

前記第 1 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された初期レプリカ・データを含む、請求項 1 に記載のプログラム。

## 【請求項 3】

前記第 1 のストレージ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信された後に、第 2 のレプリカ・ボリュームのフラッシュコピーを作成し、第 2 のストレージ・ボリュームの複製である前記第 2 のレプリカ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信されたことに応答して、前記第 2 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行するステップと、

20

前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートし、これによって、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーを前記クライアントに利用可能とするステップと、  
を更に実行させる、請求項 1 に記載のプログラム。

## 【請求項 4】

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーを消去して、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータの管理を容易にするステップを更に実行させる、請求項 3 に記載のプログラム。

30

## 【請求項 5】

前記第 2 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された更新済みソース・データを含み、

前記第 2 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された更新済みレプリカ・データを含む、請求項 3 に記載のプログラム。

## 【請求項 6】

ターゲット・サーバであって、  
プロセッサと、

前記プロセッサと共に動作可能な命令を記憶するメモリであって、前記命令が、

40

第 1 のレプリカ・ボリュームのフラッシュコピーを作成し、第 1 のストレージ・ボリュームの複製である前記第 1 のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、前記第 1 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行し、

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートして、前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーをクライアントに利用可能とする、  
ために実行される、メモリと、  
を含む、ターゲット・サーバ。

## 【請求項 7】

50

前記第 1 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された初期ソース・データを含み、

前記第 1 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された初期レプリカ・データを含む、請求項 6 に記載のターゲット・サーバ。

【請求項 8】

前記命令が、更に、

前記第 1 のストレージ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信された後に、第 2 のレプリカ・ボリュームのフラッシュコピーを作成し、第 2 のストレージ・ボリュームの複製である前記第 2 のレプリカ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信されたことに応答して、前記第 2 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行し、

10

前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートして、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーを前記クライアントに利用可能とする、

ために実行される、請求項 6 に記載のターゲット・サーバ。

【請求項 9】

前記命令が、更に、

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーを消去して、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータの管理を容易にするために実行される、請求項 8 に記載のターゲット・サーバ。

20

【請求項 10】

前記第 2 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された更新済みソース・データを含み、

前記第 2 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された更新済みレプリカ・データを含む、請求項 8 に記載のターゲット・サーバ。

【請求項 11】

ターゲット・サーバであって、

第 1 のレプリカ・ボリュームのフラッシュコピーを作成し、第 1 のストレージ・ボリュームの複製である前記第 1 のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、前記第 1 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行するための手段と、

30

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートし、これによって、前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーをクライアントに利用可能とするための手段と、

を含む、ターゲット・サーバ。

【請求項 12】

前記第 1 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された初期ソース・データを含み、

前記第 1 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された初期レプリカ・データを含む、請求項 11 に記載のターゲット・サーバ。

40

【請求項 13】

前記第 1 のストレージ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信された後に、第 2 のレプリカ・ボリュームのフラッシュコピーを作成し、第 2 のストレージ・ボリュームの複製である前記第 2 のレプリカ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信されたことに応答して、前記第 2 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行するための手段と、

前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートして、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーを前記クラ

50

クライアントに利用可能とするための手段と、  
を更に含む、請求項 1 1 に記載のターゲット・サーバ。

【請求項 1 4】

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーを消去して、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータの管理を容易にするための手段を更に含む、請求項 1 3 に記載のターゲット・サーバ。

【請求項 1 5】

前記第 2 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された更新済みソース・データを含み、

前記第 2 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された更新済みレプリカ・データを含む、請求項 1 3 に記載のターゲット・サーバ。 10

【請求項 1 6】

第 1 のレプリカ・ボリュームのフラッシュコピーを作成し、第 1 のストレージ・ボリュームの複製である前記第 1 のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、前記第 1 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行するステップと、

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートし、これによって、前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーをクライアントに利用可能とするステップと、

を含む、方法。 20

【請求項 1 7】

前記第 1 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された初期ソース・データを含み、

前記第 1 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された初期レプリカ・データを含む、請求項 1 6 に記載の方法。

【請求項 1 8】

前記第 1 のストレージ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信された後に、第 2 のレプリカ・ボリュームのフラッシュコピーを作成し、第 2 のストレージ・ボリュームの複製である前記第 2 のレプリカ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信されたことに応答して、前記第 2 のレプリカ・ボリュームの前記フラッシュコピーの回復を実行するステップと、

前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータをエクスポートして、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーを前記クライアントに利用可能とする、ステップと、

を更に含む、請求項 1 6 に記載の方法。 30

【請求項 1 9】

前記第 1 のレプリカ・ボリュームの前記回復したフラッシュコピーを消去して、前記第 2 のレプリカ・ボリュームの前記回復したフラッシュコピーのデータの管理を容易にするステップを更に含む、請求項 1 8 に記載の方法。 40

【請求項 2 0】

前記第 2 のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された更新済みソース・データを含み、

前記第 2 のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された更新済みレプリカ・データを含む、請求項 1 8 に記載の方法。

【請求項 2 1】

ターゲット・サーバのプロセッサに、

第 1 のストレージ・ボリュームの複製である第 1 のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、前記第 1 のレプリカ・ボリュームの回復したフラッシュコピーのデータをエクス 50

ポートするステップと、

前記第1のストレージ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信された後に、第2のストレージ・ボリュームの複製である第2のレプリカ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信されたことに応答して、前記第2のレプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートするステップと、  
 を実行させ、前記第1のレプリカ・ボリュームの前記回復したフラッシュコピーを消去して、前記第2のレプリカ・ボリュームの前記回復したフラッシュコピーのデータの管理を容易にする、プログラム。

【請求項22】

前記第1のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された初期ソース・データを含む、請求項21に記載のプログラム。

【請求項23】

前記第1のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された初期レプリカ・データを含む、請求項21に記載のプログラム。

【請求項24】

前記第2のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された更新済みソース・データを含む、請求項21又は請求項22に記載のプログラム。

【請求項25】

前記第2のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された更新済みレプリカ・データを含む、請求項21又は請求項23に記載のプログラム。

【請求項26】

ターゲット・サーバであって、  
 プロセッサと、

前記プロセッサと共に動作可能な命令を記憶するメモリであって、前記命令が、

第1のストレージ・ボリュームの複製である第1のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、前記第1のレプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートし、

前記第1のストレージ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信された後に、第2のストレージ・ボリュームの複製である第2のレプリカ・ボリュームが前記ソース・ストレージ・システムから前記ターゲット・ストレージ・システムによって受信されたことに応答して、前記第2のレプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートする、  
 ために実行される、メモリと、  
 を含み、前記第1のレプリカ・ボリュームの前記回復したフラッシュコピーを消去して、前記第2のレプリカ・ボリュームの前記回復したフラッシュコピーのデータの管理を容易にする、ターゲット・サーバ。

【請求項27】

前記第1のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された初期ソース・データを含む、請求項26に記載のターゲット・サーバ。

【請求項28】

前記第1のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された初期レプリカ・データを含む、請求項26に記載のターゲット・サーバ。

【請求項29】

前記第2のストレージ・ボリュームが、ストレージ論理ユニット内に記憶された更新済みソース・データを含む、請求項26又は請求項27に記載のターゲット・サーバ。

【請求項30】

前記第2のレプリカ・ボリュームが、レプリカ論理ユニット内に記憶された更新済みレプリカ・データを含む、請求項26又は請求項28に記載のターゲット・サーバ。

10

20

30

40

50

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、一般に、ストレージ・ボリューム (storage volume) と呼ぶ S C S I 論理ユニット (「LU: logical unit」) のブロック・レベルの複製 (replication) に関する。本発明は、具体的には、ブロック・レベルの複製を利用してファイル・システム・レベルの複製システムを構築し、ファイル・システム・データに対する高可用性のリード・オンリー・アクセス (read-only access) を提供することに関する。

## 【背景技術】

## 【0002】

一般に、複数のストレージ・システム間での既存の (同期および非同期の) ブロック複製方式は、障害回復 (disaster recovery) のために有用である。具体的には、あるソース・サイトが故障した場合、ターゲット・サイトが、そのソース・サイトの機能を引き継ぐことができる。障害によってソース・サーバが故障した場合、キャッシュしたデータを含めて、ソースにおける動的 / 揮発性サーバ状態は、必ずしもターゲット・サイトに提示されるわけではないので、ターゲット・サイト上のサーバは、一般に、このターゲット・サイトにおけるサービスを開始する前に、複製したデータのシステム・レベルまたはアプリケーション・レベルの回復を実行しなければならない。例えば、ロギング・ファイル・システムにおいて、遠隔サイトにおけるストレージ・ボリュームのレプリカが、作動中のファイル・システムとして使用可能になる前に、ファイル・システムのログを再生 (replay) しなければならない。特に、ストレージ・ボリューム上のデータ構造は、更新されているので、常に相互に一貫しているわけではない。概して、いかなる更新も2つ以上の記憶されたデータ構造に影響を与える可能性があり、記憶された構造を全て同時に更新することはできない。ストレージ・アプリケーションのための異なるアーキテクチャは、この一時的な不一致を管理するために異なる方法を用いる。一部のものは、べき等の演算 (idempotent operations) のログを用いて対象の更新を記録し、再起動の後にそれらの演算を再生する。あるものは、その代わりに、再起動の後に、全てのデータ構造をスイープ (sweep) して、修復する不一致を探す。いずれの場合も、その目的は、システム再起動の後にデータ構造の一貫性を修復することである。

## 【0003】

データ複製方式によって理想的に解決されると思われる別の問題は、ソース・データのリード・オンリー・バージョンに対するアクセスを与える多数のレプリカを有することによって、同時に多数のサイトでデータを利用可能とすることである。この問題を解決するためには、ブロック複製のみでは不十分である。なぜなら、サーバによって基礎のデータに対するアクセスを与えることを可能とするために、回復ステップが必要であるからである。サーバ・ソフトウェアは、それが用いるストレージ・ボリュームにおけるデータに対する変更を追跡することができず、そのデータのそれ自身のキャッシュを更新する。多数のサイトでデータを複製することは有益である。なぜなら、これによって、ネットワーク故障に直面した場合、例えば、複製したデータを有するサイトの全てではないが一部に接続できない場合に、データの可用性を高めることができるからである。

## 【0004】

これに対して、rsync等のファイル・システム複製方式は、同時データ更新および高可用性のリード・オンリー・アクセスを可能とする。なぜなら、データ更新はファイル・システム自体を流れ、サーバ・キャッシュおよびディスク・データを更新するからである。「rsync」では、ソース・サイトにおけるファイル・システム・アプリケーションとして実行しているプロセスは、ネットワークを介して、宛先サイトにおいてファイル・システム・アプリケーションとして実行しているプロセスと通信する。ソース・サイトにおけるプロセスは、1組のファイルの現在の状態を読み取り、その現在の状態を宛先サイトにおけるプロセスに送信し、宛先サイトにおけるプロセスは、宛先サイトにおける対応するファイルの状態を調整して、ソース・サイトにおけるそれらをミラーリングする。

10

20

30

40

50

しかしながら、かかる方式は、通常、特別なハードウェアおよびファームウェアを利用するために、ブロック複製方式に比べて性能が制限される。

【発明の開示】

【発明が解決しようとする課題】

【0005】

クラッシュの後の「回復」ステップに必要な場合、利用可能なデータ・コピーを用いて、障害回復のために、ブロック複製またはファイル複製のいずれかを利用可能である。しかしながら、ブロック複製は、単独では可用性も負荷バランシング (load-balancing) も増大させない。これは、ファイル・システム複製では可能であるものの、ブロック複製によって可能であるよりも性能が低い。

10

【0006】

従って、ブロック・レベルの複製方式の利点を活用してファイル・システム複製方式を実施するための技法を開発することは、コンピュータ業界にとって課題である。

【課題を解決するための手段】

【0007】

本発明の第1の実施形態は、ターゲット・サーバのプロセッサによって実行可能な機械読み取り可能命令のプログラムを明確に具現化して動作を実行する信号搬送媒体であって、その動作が、レプリカ・ボリュームのフラッシュコピーを作成し、第1のストレージ・ボリュームの複製であるレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、レプリカ・ボリュームのフラッシュコピーの回復を実行するステップと、レプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートし、これによって、レプリカ・ボリュームの回復したフラッシュコピーをクライアントに利用可能とするステップと、を含む。

20

【0008】

本発明の第2の実施形態は、ターゲット・サーバであって、プロセッサと、このプロセッサと共に動作可能な命令を記憶するメモリであって、この命令が、レプリカ・ボリュームのフラッシュコピーを作成し、第1のストレージ・ボリュームの複製である第1のレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、レプリカ・ボリュームのフラッシュコピーの回復を実行し、レプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートし、これによって、レプリカ・ボリュームの回復したフラッシュコピーをクライアントに利用可能とするために実行される、メモリと、を含む。

30

【0009】

本発明の第3の実施形態は、ターゲット・サーバであって、レプリカ・ボリュームのフラッシュコピーを作成し、第1のストレージ・ボリュームの複製であるレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、レプリカ・ボリュームのフラッシュコピーの回復を実行するための手段と、レプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートし、これによって、レプリカ・ボリュームの回復したフラッシュコピーをクライアントに利用可能とするための手段と、を含む。

40

【0010】

本発明の第4の実施形態は、レプリカ・ボリュームのフラッシュコピーを作成し、第1のストレージ・ボリュームの複製であるレプリカ・ボリュームがソース・ストレージ・システムからターゲット・ストレージ・システムによって受信されたことに応答して、レプリカ・ボリュームのフラッシュコピーの回復を実行するステップと、クライアントによる要求に応じてレプリカ・ボリュームの回復したフラッシュコピーのデータをエクスポートし、これによって、レプリカ・ボリュームの回復したフラッシュコピーをクライアントに利用可能とするステップと、を含む。

【0011】

本発明の前述の実施形態および他の実施形態、目的、および態様、ならびに特徴および

50

利点は、本発明の様々な実施形態の以下の詳細な説明から、更に明らかになる。詳細な説明および図面は、単に本発明の例示にすぎず、その特許請求の範囲および均等物によって規定される本発明の範囲を限定するものではない。

【発明を実施するための最良の形態】

【0012】

本発明は、新しい独自のファイル・システム複製機構を提供し、ブロック・レベルの複製によってデータに対する高可用性のリード・オンリー・アクセスを提供する。本発明は、高性能の同期および非同期のブロック・レベル複製方式を利用した、1つ以上のサイトにおけるリード・オンリー・ファイル・システム複製をサポートするための方式を提案し、基礎にあるターゲット・ブロック・デバイスのフラッシュコピーを定期的に生成し、それらのフラッシュコピーの内容を調整してターゲット・サイトにおける安定したファイル・システム状態を反映し、次いで、ネットワーク・ファイル・システム（「NFS」）等の分散型ファイル・システム・プロトコルを用いてその内容を再びエクスポート（reexport）する。

10

【0013】

図1は、本発明のブロック・レベル複製システムを例示的に示す。これは、ソース・サイト10、ターゲット・サイト20、ターゲット・サイト30、およびクライアント・サイト40を用いる。ソース・サイト10は、ソース・サーバ11およびソース・ストレージ・システム13を含む。ターゲット・サイト20は、ターゲット・サーバ21およびターゲット・ストレージ・システム23を含む。ターゲット・サイト30は、ターゲット・サーバ31およびターゲット・ストレージ・システム33を含む。クライアント・サイト40は、アプリケーション（「APP」）42を実行するクライアント・システム41を含み、これは、ソース・サーバ11、ターゲット・サーバ21、あるいはターゲット・サーバ31またはそれら全てが応じることができるファイルI/O要求を行う。

20

【0014】

ソース・サーバ11の動作は、データ変更器（「DM」）12が、ソース・サーバ11上のローカル・アプリケーション（図示せず）からの要求あるいはネットワーク60上のアプリケーション（例えばアプリケーション42）が駆動する分散型ファイル・システムのクライアントからの要求またはその双方に回答することを含み、これによって、データ変更器12はデータを作成し、これは、ソース・ストレージ・システム13に記憶された1つまたはそれ以上のX個のソース論理ユニット15に記憶される。ここで、Xは1以上である。ソース・ストレージ・システム13のソース・ブロック複製モジュール（「SBR」）14は、ネットワーク50（例えばストレージ・エリア・ネットワーク）を介して、ソース・データ・ブロックを、ターゲット・ストレージ・システム23のターゲット・ブロック複製器（「TBR」）24に伝達し、これによって、ターゲット・ブロック複製器24は、ソース・データ・ブロックを処理して、レプリカ・データを作成し、1つまたはそれ以上のY個のレプリカ論理ユニット25内に記憶する。ここで、Yは1以上である。同様に、ソース・ブロック複製モジュール14は、ネットワーク50を介して、ソース・データ・ブロックを、ターゲット・ストレージ・システム33のターゲット・ブロック複製器（「TBR」）34に伝達して、これによって、ターゲット・ブロック複製器34は、ソース・データ・ブロックを処理して、レプリカ・データを作成し、1つまたはそれ以上のZ個のレプリカ論理ユニット35内に記憶する。ここで、Zは1以上である。サイト30および40間では生じないがサイト20および40間で生じる等の部分的なネットワーク故障の場合にも、サイト40上の最終的なクライアントに対してデータのいくらかのコピーの可用性を保証するために、多数のレプリカが望ましい。また、40のような多くのクライアントからの負荷のバランスを取るためにも、多数のレプリカが有用である。

30

40

【0015】

例えば、図2に示すように、最初にソース論理ユニット15（1）内に記憶されているソース・データSD(IN)のソース・データ・ブロックSDB(IN)は、ソース・ブロック複製器14によって、ターゲット・ブロック複製器24および34に伝達すること

50

ができ、これらが、初期ソース・データ・ブロックSDB(IN)を処理して、レプリカ・データRD(IN)を作成し、レプリカ論理ユニット25(1)およびレプリカ論理ユニット35(1)内にそれぞれ記憶する。その後、図3に示すように、ソース論理ユニット15(1)内に記憶されている更新済みのソース・データSD(UP)を表す更新済みのソース・データ・ブロックSDB(UP)は、ソース・ブロック複製器14によって、ターゲット・ブロック複製器24および34に伝達することができ、これらが、更新済みのソース・データ・ブロックDSB(UP)を処理して、更新済みのレプリカ・データRD(UP)を作成し、各レプリカ論理ユニット25(1)およびレプリカ論理ユニット35(1)内に記憶する。レプリカ論理ユニット25(1)および35(1)内でのレプリカ・データRD(UP)の更新は、無限に繰り返すことができる。

10

**【0016】**

図1から図3を参照すると、レプリカ論理ユニット25(1)およびレプリカ論理ユニット35(1)内に記憶されているレプリカ・データRDのネットワーク60を介したクライアント41に対する可用性(例えばレプリカ・データRDのエクスポート)は、従来、実行不可能である。これは、レプリカ・データRDはネットワーク60を介してエクスポートされているが、本発明以前の各ターゲット・サーバ21および31が、各レプリカ論理ユニット25(1)および35(1)内に記憶されているレプリカ・データRDに対する更新を処理することができないことによる。特に、レプリカ論理ユニット25(1)およびレプリカ論理ユニット35(1)内に記憶されたレプリカ・データRDに対して更新を行うたびに、ネットワーク60を介したレプリカ・データRDのエクスポートにおいて大きな不一致が生じる。

20

**【0017】**

この欠点を克服するため、本発明は、ターゲット・サーバの各々において、フラッシュコピー・エクスポート(FCE: flashcopy exporter)を設ける。ここで、flashcopyはIBM Corporationの商標である。図1の例では、フラッシュコピー・エクスポート22および32は、ターゲット・サーバ21および31上にそれぞれインストールされている。動作時に、ターゲット20および30の各々におけるレプリカ・データRDが更新のプロセス中である場合にもそうでない場合にも、フラッシュコピー・エクスポート22およびフラッシュコピー・エクスポート32は、レプリカ論理ユニット25(1)およびレプリカ論理ユニット35(1)内にそれぞれ記憶されているようなレプリカ・データRDに対するクライアント41上のアプリケーション42による安定した高可用性のリード・オンリー・アクセスを提供する。このために、フラッシュコピー・エクスポート22およびフラッシュコピー・エクスポート32は、本発明のレプリカ・データ・エクスポート方法を表す図4に示すようなフローチャート70を実施する。

30

**【0018】**

図4を参照すると、フローチャート70の段階72は、フラッシュコピー・エクスポート22およびフラッシュコピー・エクスポート32が、最初に各レプリカ論理ユニット25(1)および35(1)内に記憶されていたようなレプリカ・データの回復したフラッシュコピーのエクスポートを実施することを含む。フローチャートの段階S76は、フラッシュコピー・エクスポート22およびフラッシュコピー・エクスポート32が、各レプリカ論理ユニット25(1)および35(1)内で更新されているような記憶レプリカ・データの回復した各フラッシュコピーのエクスポートを実行することを含む。使用およびエクスポートの準備のために、レプリカ・データのフラッシュコピーの回復を行って、データの様々な部分が相互に一貫していることを確実にすることは、当業者には認められよう。

40

**【0019】**

各レプリカ論理ユニット25(1)および35(1)内に記憶されているレプリカ・データの更新を無限に反復可能であるという事実を考慮して、段階S76は、フラッシュコピー・エクスポート22およびフラッシュコピー・エクスポート32によって無限に繰り返すことができる。フローチャート70の決定段階S74は、いつ段階S76を実行する

50

かを制御する。段階 S 7 4 の意思決定プロセスの背後にある論理の一実施形態については、以下で説明する。

【 0 0 2 0 】

図 5 に示すように、本発明の初期のレプリカ・データ・エクスポート方法を表すフローチャート 8 0 は、段階 S 7 2 の一実施形態である。図 7 に示すように、本発明の更新済みレプリカ・データ・エクスポート方法を表すフローチャート 9 0 は、段階 S 7 6 の一実施形態である。実際には、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が段階 S 7 2 および S 7 6 を実行する方法は、限定されない。このため、以下のフローチャート 8 0 および 9 0 の説明は、フローチャート 7 0 の範囲を限定するものではない。

10

【 0 0 2 1 】

図 1、図 2、および図 5 を参照すると、フローチャート 8 0 の段階 S 8 2 は、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、各レプリカ論理ユニット 2 5 ( 1 ) およびレプリカ論理ユニット 3 5 ( 1 ) 内に記憶されている初期レプリカ・データ R D ( I N ) のフラッシュコピーを生成することを含む。フローチャート 8 0 の段階 S 8 4 は、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、各レプリカ論理ユニット 2 5 ( 1 ) およびレプリカ論理ユニット 3 5 ( 1 ) 内に記憶されている初期レプリカ・データ R D ( I N ) のフラッシュコピーの回復を実行することを含む。レプリカ論理ユニット 2 5 ( 1 ) および 3 5 ( 1 ) 内のデータは同一となり、これによって、レプリカ論理ユニット 2 5 ( 1 ) および 3 5 ( 1 ) が異なる地理的位置にあるという事実を考慮して、データの多数の同時コピーを提供して、アプリケーション・サイト 4 0 に対するデータのいくらかのコピーの可用性を保証する。フローチャート 8 0 の段階 S 8 6 は、クライアント 4 1 上にインストールされたアプリケーション 4 2 のコピーによって要求があると、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、各レプリカ論理ユニット 2 5 ( 1 ) およびレプリカ論理ユニット 3 5 ( 1 ) 内に記憶されている初期レプリカ・データ R D ( I N ) の回復したフラッシュコピーを、ネットワーク 6 0 を介してクライアント 4 1 にエクスポートすることを含む。

20

【 0 0 2 2 】

例えば、図 2 および図 6 に示すように、各レプリカ論理ユニット 2 5 ( 1 ) およびレプリカ論理ユニット 3 5 ( 1 ) 内に記憶されている初期レプリカ・データ R D ( I N ) を、段階 S 8 2 および S 8 4 ( 図 5 ) に従って処理し、これによって、回復済みフラッシュコピー論理ユニット 2 5 ( 2 ) および回復済みフラッシュコピー論理ユニット 3 5 ( 2 ) を与える。これらは双方とも、各レプリカ論理ユニット 2 5 ( 1 ) および 3 5 ( 1 ) 内に記憶された初期レプリカ・データ R D ( I N ) の回復済みフラッシュコピー R F D ( I N ) を記憶する。回復済みフラッシュコピー R F D ( I N ) のデータは、クライアント 4 1 ( 図 1 ) 上のアプリケーション 4 2 による要求に応じて、各フラッシュコピー・エクスポート 2 2 および 3 2 によってエクスポートされ、これによって、R F D ( I N ) 内のファイル・システム・データのコピーに対する高可用性のリード・オンリー・アクセスを得る。

30

【 0 0 2 3 】

図 1、図 5、および図 6 を参照すると、フローチャート 7 0 の段階 S 7 4 は、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、初期レプリカ・データ R D ( I N ) の回復済みフラッシュコピー R F D ( I N ) の更新が認可されているか否かを判定することを含む。一実施形態では、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 は、各ターゲット・ブロック複製器 2 4 および 3 4 に問い合わせ、各レプリカ論理ユニット 2 5 ( 1 ) および 3 5 ( 1 ) 内に記憶されているレプリカ・データに対していずれかの新しい更新が行われたか否かを判定する。フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 は、( 1 ) 各レプリカ論理ユニット 2 5 ( 1 ) および 3 5 ( 1 ) 内に記憶されているレプリカ・データに行われた更新の量、および、( 2 ) 初期レプリカ・データ R D ( I N ) の回復

40

50

済みフラッシュコピー RFD (IN) を更新するために適切な場合を指定する更新ポリシーに基づいて、段階 S 7 4 における判定を行う。更新ポリシーが、レプリカ・データに行われた更新の量以外のファクタを含むことは、当業者には認められよう。また、段階 S 7 4 の間に、それらの他のファクタを、フラッシュコピー・エクスポート 2 2 および 3 2 によって考慮することができる。

**【 0 0 2 4 】**

段階 S 7 4 の間に、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、初期レプリカ・データ RD (IN) の回復済みフラッシュコピー RFD (IN) に対する更新が認可されていないと判定すると、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 は、段階 S 7 4 に戻って、後の時点で、レプリカ・データ RD (IN) の回復済みフラッシュコピー RFD (IN) に対する更新が認可されているか否かの判定を繰り返す。

10

**【 0 0 2 5 】**

段階 S 7 4 の間に、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、初期レプリカ・データ RD (IN) の回復済みフラッシュコピー RFD (IN) の更新が認可されていると判定すると、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 は、上述のように、段階 S 7 6 においてフローチャート 9 0 を実施する。

**【 0 0 2 6 】**

図 1、図 3、および図 7 を参照すると、フローチャート 9 0 の段階 S 9 2 は、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、各レプリカ論理ユニット 2 5 (1) およびレプリカ論理ユニット 3 5 (1) 内に記憶されている更新済みレプリカ・データ RD (UP) のフラッシュコピーを生成することを含む。この場合も、多数の位置においてデータを冗長的に提供して、可用性に対応し、単一のサイトが提供可能なよりも大きな負荷のサービスを容易にするという目的のために、レプリカ論理ユニット 2 5 (1) および 3 5 (1) に記憶されているデータは同一である。フローチャート 9 0 の段階 S 9 4 は、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、各レプリカ論理ユニット 2 5 (1) およびレプリカ論理ユニット 3 5 (1) 内に記憶されている更新済みレプリカ・データ RD (UP) のフラッシュコピーの回復を実行することを含む。フローチャート 9 0 の段階 S 9 6 は、クライアント 4 1 上にインストールされたアプリケーション 4 2 のコピーによる要求に応じて、フラッシュコピー・エクスポート 2 2 およびフラッシュコピー・エクスポート 3 2 が、各レプリカ論理ユニット 2 5 (1) およびレプリカ論理ユニット 3 5 (1) 内に記憶されている更新済みレプリカ・データ RD (UP) の回復済みフラッシュコピーを、ネットワーク 6 0 を介してクライアント 4 1 にエクスポートすることを含む。段階 S 9 6 は、更に、各レプリカ論理ユニット 2 5 (1) およびレプリカ論理ユニット 3 5 (1) 内に記憶された初期のまたは更新済みのレプリカ・データ RD の以前の回復済みフラッシュコピーを消去することを含む。

20

30

**【 0 0 2 7 】**

例えば、図 8 に示すように、レプリカ論理ユニット 2 5 (1) およびレプリカ論理ユニット 3 5 (1) 内に記憶されている第 1 の組の更新済みレプリカ・データ RD (UP 1) を、段階 S 9 2 および S 9 4 (図 7) に従って処理し、これによって、回復済みフラッシュコピー論理ユニット 2 5 (3) および回復済みフラッシュコピー論理ユニット 3 5 (3) を与える。これらは双方とも、各レプリカ論理ユニット 2 5 (1) および 3 5 (1) 内に記憶された第 1 の組の更新済みレプリカ・データ RD (UP 1) の回復済みフラッシュコピー RFD (UP 1) を記憶する。回復済みフラッシュコピー RFD (UP 1) のデータは、クライアント 4 1 (図 1) 上のアプリケーション 4 2 による要求に応じて、各フラッシュコピー・エクスポート 2 2 および 3 2 によってエクスポートされ、これによって、回復済みフラッシュコピー RFD (UP 1) 内のファイル・システム・データのコピーに対する高可用性のリード・オンリー・アクセスを得る。クライアント 4 0 によるこのアク

40

50

セスは、レプリカ論理ユニット25(1)および35(1)内に記憶された第1の組の更新済みレプリカ・データRD(UP1)がそれ自体更新されるまで続く。

【0028】

例えば、図9に示すように、レプリカ論理ユニット25(1)およびレプリカ論理ユニット35(1)内に記憶されている第2の組の更新済みレプリカ・データRD(UP2)を、段階S92およびS94(図7)に従って処理し、これによって、回復済みフラッシュコピー論理ユニット25(4)および回復済みフラッシュコピー論理ユニット35(4)を与える。これらは双方とも、各レプリカ論理ユニット25(1)および35(1)内に記憶された第2の組の更新済みレプリカ・データRD(UP2)の回復済みフラッシュコピーRFD(UP2)を記憶する。回復済みフラッシュコピーRFD(UP2)のデータは、クライアント41(図1)上のアプリケーション42による要求に応じて、各フラッシュコピー・エクスポート22および32によってエクスポートされ、これによって、回復済みフラッシュコピーRFD(UP2)内のファイル・システム・データのコピーに対する高可用性のリード・オンリー・アクセスを得る。クライアント40によるこのアクセスは、レプリカ論理ユニット25(1)および35(1)内に記憶された第2の組の更新済みレプリカ・データRD(UP2)がそれ自体更新されてフローチャート90が完了するまで続く。更に、回復済みフラッシュコピー論理ユニット25(3)および回復済みフラッシュコピー論理ユニット35(3)を消去する。

10

【0029】

再び図7を参照すると、更新済みレプリカ・データの前の回復済みフラッシュコピーから、更新済みレプリカ・データの新しい回復済みフラッシュコピーへの切り替えは、動作上の問題を生じる恐れがあるが、これを防ぐためには、更新済みレプリカ・データの前の回復済みフラッシュコピーのエクスポートを順次休止させ、切り替えを実行し、更新済みレプリカ・データの新しい回復済みフラッシュコピーのエクスポートを可能とすれば良いことは、当業者には認められよう。アプリケーションに基づいたシステムの実施形態では、故障指示をアプリケーションに与えて、切り換えの後に開いたファイルが無効である場合、切り換えの前にファイルを開く。ファイルシステムの実施形態では、クライアントをプログラミングして、原則的に、切り換えを行う時を知らせ、これによっていかなる動作上の問題も防ぐ。これは、従来、NFSプロトコルを用いたアプリケーションによって達成される。

20

30

【0030】

本発明の理解を容易にする目的のため、本明細書では、図1に示したブロック・レベルの複製システムの文脈において、本発明の様々な方法を説明している。しかしながら、1つ以上のターゲット・サイトを有し、ネットワーク・アプリケーション(図1に示すアプリケーション42)に1つ以上のデータ・ソースを提供する他のブロック・レベルの複製システムの文脈において、本発明の様々な方法がどのように実施されるかは、当業者には認められよう。

【0031】

図1から図9の前述の説明から、当業者には、本発明の多数の利点が認められよう。かかる利点のうち最も重要なものは、ブロック・レベルの複製による、データに対する高可用性のリード・オンリー・アクセスである。

40

【0032】

図1を参照すると、1つの実際の実施形態において、モジュール14、24、および34は、従来のソフトウェア・アプリケーションであるが、フラッシュコピー・エクスポート22および32は、各サーバ21および31のメモリ内にインストールされた新しいソフトウェア・モジュールにおいてハードウェア資源との協働関係によって具現化され、これによって、各サーバ21および31のプロセッサ(複数のプロセッサ)は、各フラッシュコピー・エクスポート22および32を実行して、図4、図5、および図7に例示的に示すようなフローチャート70、80、および90を実施することができる。フラッシュコピー・エクスポート22および32は、ソフトウェア・モジュールとして具現化される

50

場合、図 1 から図 9 の本明細書における説明を評価する当業者によって、いずれかのプログラミング言語で、プロセッサ等を実行させるプログラムとして記述することができる。

【 0 0 3 3 】

本明細書において開示した本発明の実施形態は、現在好適な実施形態と考えられるが、本発明の精神および範囲から逸脱することなく、様々な変更および変形が可能である。本発明の範囲は、特許請求の範囲に指示され、均等物の意味および範囲内に該当する全ての変更は、それに包含することが意図される。

【 図面の簡単な説明 】

【 0 0 3 4 】

【 図 1 】 本発明に従ったブロック複製システムの例示的な実施形態を示す。 10

【 図 2 】 従来技術において既知であるようなソース・データの例示的なブロック複製を示す。

【 図 3 】 従来技術において既知であるようなソース・データの例示的なブロック複製を示す。

【 図 4 】 本発明に従ったレプリカ・データ・エクスポート方法を表すフローチャートを示す。

【 図 5 】 本発明に従った初期レプリカ・データ・エクスポート方法を表すフローチャートを示す。

【 図 6 】 図 5 に示すフローチャートの例示的な動作を示す。

【 図 7 】 本発明に従った更新済みレプリカ・データ・エクスポート方法を表すフローチャートを示す。 20

【 図 8 】 図 7 に示すフローチャートの例示的な以降の動作を示す。

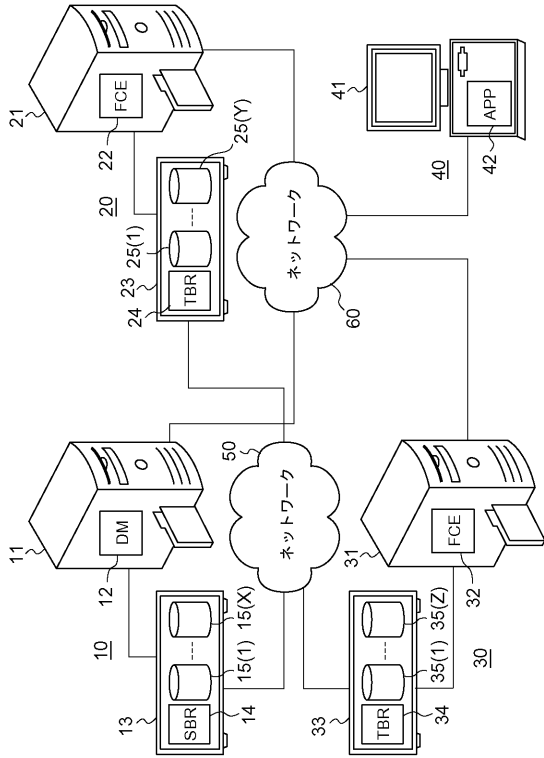
【 図 9 】 図 7 に示すフローチャートの例示的な以降の動作を示す。

【 符号の説明 】

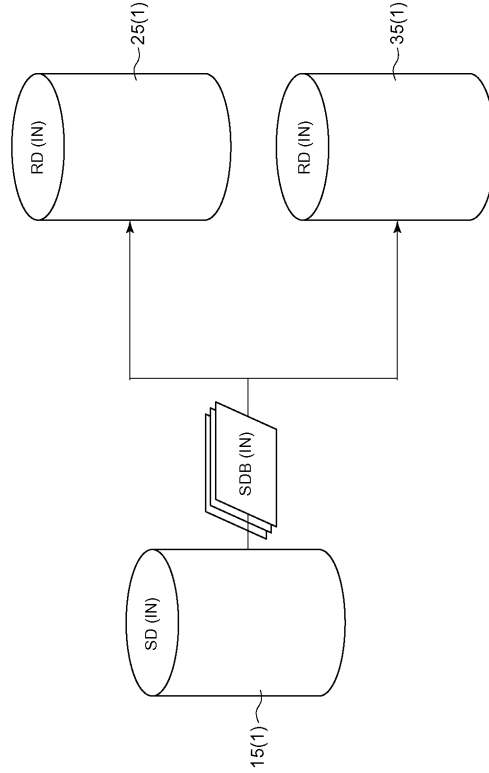
【 0 0 3 5 】

- 1 0 ソース・サイト
- 1 1 ソース・サーバ
- 1 2 データ変更器
- 1 3 ソース・ストレージ・システム
- 1 4 ソース・ブロック複製モジュール 30
- 1 5 ソース論理ユニット
- 2 0、3 0 ターゲット・サイト
- 2 1、3 1 ターゲット・サーバ
- 2 2、3 2 フラッシュコピー・エクスポータ
- 2 3、3 3 ターゲット・ストレージ・システム
- 2 4、3 4 ターゲット・ブロック複製器
- 2 5、3 5 レプリカ論理ユニット
- 4 0 クライアント・サイト
- 4 1 クライアント・システム
- 4 2 アプリケーション 40
- 5 0、6 0 ネットワーク

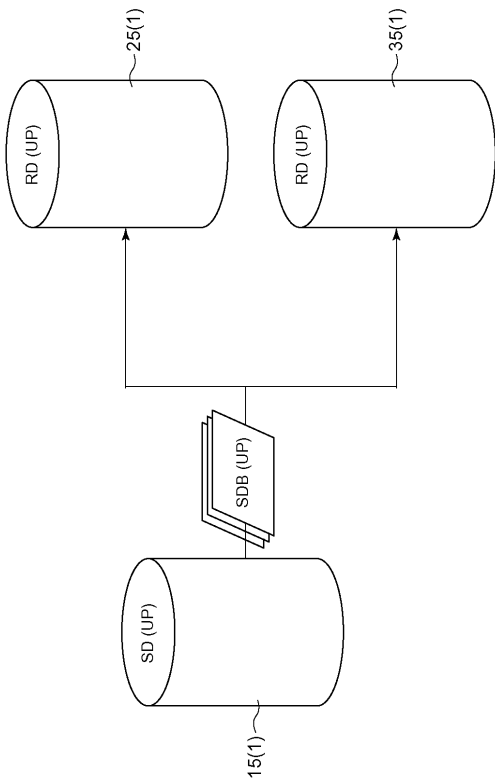
【 図 1 】



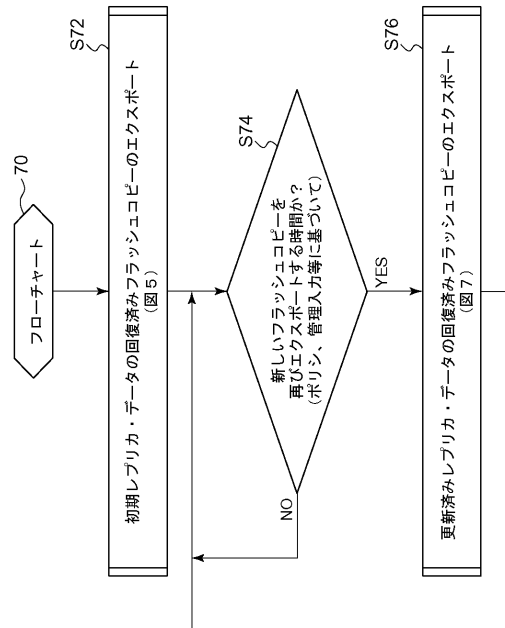
【 図 2 】



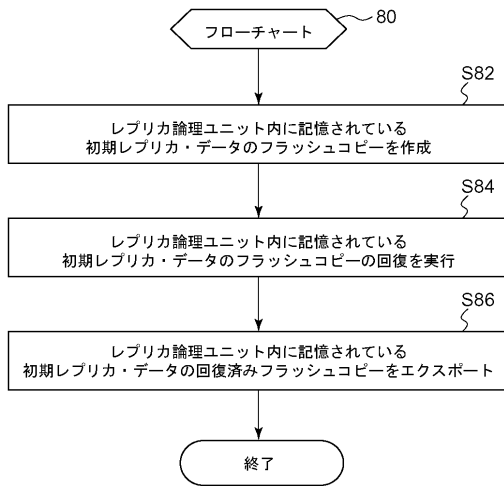
【 図 3 】



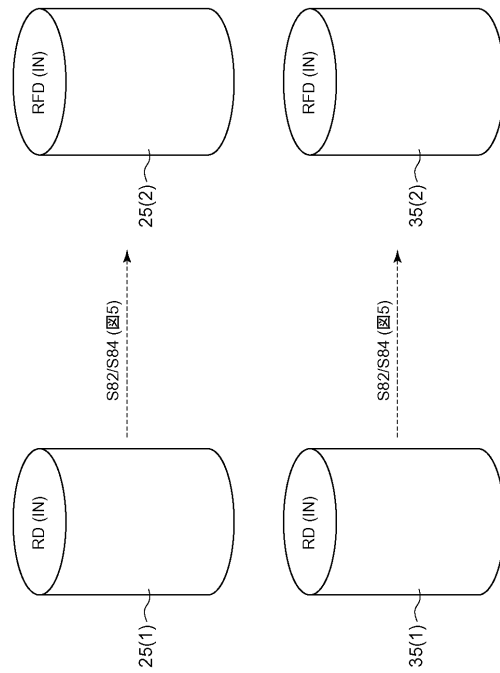
【 図 4 】



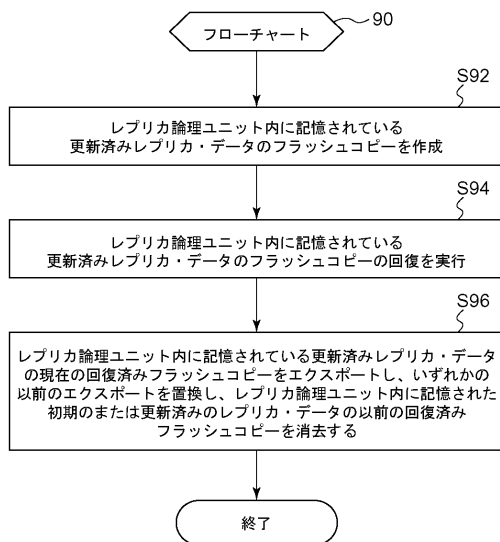
【 図 5 】



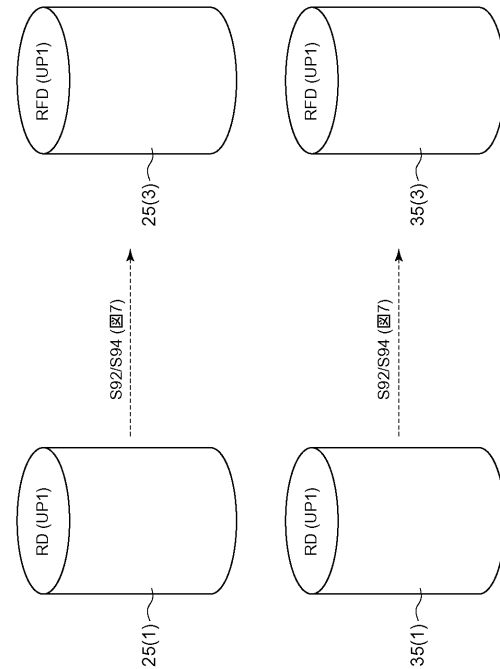
【 図 6 】



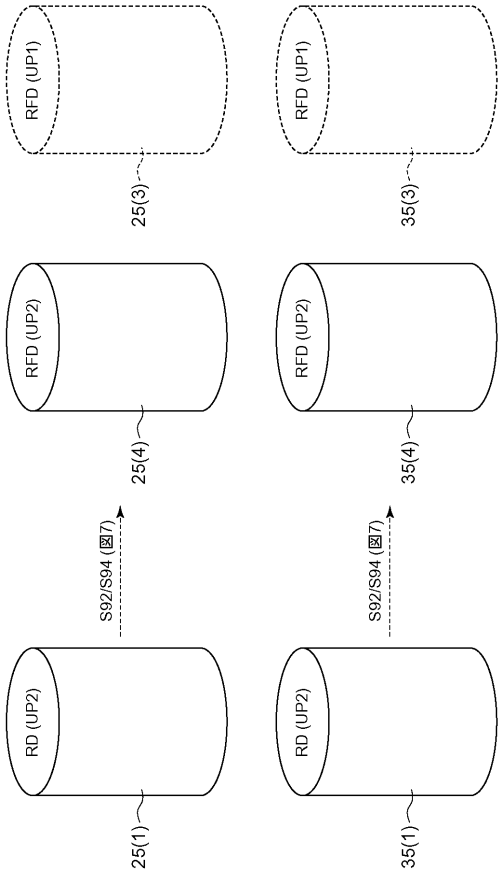
【 図 7 】



【 図 8 】



【 図 9 】



---

フロントページの続き

(74)代理人 100112690

弁理士 太佐 種一

(72)発明者 クレイグ・フルマー・エヴァーハート

アメリカ合衆国 2 7 5 1 4 ノースカロライナ州チャペル・ヒル ハンティントン・ドライブ 2 2  
5

(72)発明者 ショーミットロ・サーカー

アメリカ合衆国 2 7 5 1 3 ノースカロライナ州キャリー パーセル・ドライブ 3 1 4

Fターム(参考) 5B065 EA33

5B082 DE06

【要約の続き】

【選択図】 図 4