

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4108074号  
(P4108074)

(45) 発行日 平成20年6月25日 (2008. 6. 25)

(24) 登録日 平成20年4月11日 (2008. 4. 11)

(51) Int. Cl.	F I
<b>G 0 6 F 3/06 (2006. 01)</b>	G O 6 F 3/06 3 O 4 F
<b>G 0 6 F 13/00 (2006. 01)</b>	G O 6 F 13/00 5 2 O C

請求項の数 1 外国語出願 (全 13 頁)

(21) 出願番号	特願2004-307582 (P2004-307582)	(73) 特許権者	000005108
(22) 出願日	平成16年10月22日 (2004. 10. 22)		株式会社日立製作所
(65) 公開番号	特開2005-301976 (P2005-301976A)		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27日 (2005. 10. 27)	(74) 代理人	100079108
審査請求日	平成18年10月24日 (2006. 10. 24)		弁理士 稲葉 良幸
(31) 優先権主張番号	10/792550	(74) 代理人	100093861
(32) 優先日	平成16年3月2日 (2004. 3. 2)		弁理士 大賀 真司
(33) 優先権主張国	米国 (US)	(72) 発明者	大崎 伸之
早期審査対象出願			アメリカ合衆国カリフォルニア州 キャン ベル エラムアベニュー 1 2 8 1
		審査官	梅景 篤

最終頁に続く

(54) 【発明の名称】 多重リモートストレージでのデータ同期

(57) 【特許請求の範囲】

【請求項 1】

プライマリストレージサブシステムと、少なくとも一のセカンダリストレージサブシステムを有するストレージシステムにおいて、

前記各ストレージサブシステムはそれぞれ、受信した一つ以上のデータブロックと各データブロックに付与された順序番号を格納するキューを備え、前記各セカンダリストレージサブシステムは、前記プライマリストレージサブシステムに直接、又は、他のセカンダリストレージサブシステムを介して間接的に接続されており、

前記プライマリストレージサブシステムは、ホストシステムで生成された一つ以上のデータブロックを受信すると、各データブロックに対して順序番号を付与し、該プライマリストレージサブシステムから直接データブロックを受信するストレージサブシステムである後続ストレージサブシステムに対して、前記各データブロックと各データブロックに付与された順序番号を送信し、

前記各セカンダリストレージサブシステムは、前記プライマリストレージサブシステムから送信された各データブロックと各データブロックに付与された順序番号を、前記プライマリストレージサブシステム、又は、他のセカンダリストレージサブシステムのいずれか一つから直接受信し、該セカンダリストレージサブシステムから直接データブロックを受信するストレージサブシステムである後続ストレージサブシステムに対して、前記各データブロックと各データブロックに付与された順序番号を送信し、

該セカンダリストレージサブシステムが受信した最新のデータブロックに付与された順

10

20

序番号と該セカンダリストレージサブシステムの後続ストレージサブシステムから受信した順序番号の中で最も古いデータブロックに付与された順序番号を、該セカンダリストレージサブシステムに直接データブロックを送信するストレージサブシステムである先行ストレージサブシステムに対して送信し、

前記プライマリストレージサブシステムは、該プライマリストレージサブシステムの後続ストレージサブシステムから受信した順序番号の中から最も古いデータブロックに付与された順序番号を検索して、該順序番号を該後続ストレージサブシステムに報告し、

前記各セカンダリストレージサブシステムは、該セカンダリストレージサブシステムの先行ストレージサブシステムから、前記プライマリストレージサブシステムにより検索された順序番号を受信し、当該順序番号を該セカンダリストレージサブシステムの後続ストレージサブシステムに報告することを特徴とするストレージシステム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般的に、プライマリ(又はローカル)ストレージ施設と本プライマリストレージ施設に保存されているデータの少なくとも一部をミラーする複数のセカンダリ(又はリモート)ストレージ施設より成る、データ処理ストレージシステムに関連する。より具体的には、本発明は、一つのストレージロケーションから他のストレージロケーションにデータをコピー中に中断が発生した場合に、生存ストレージ施設間でデータを同期化する方法及び本方法を備えた装置に関連する。

20

【背景技術】

【0002】

企業、政府及び他の組織体でのデータ処理業務の拡大により、夥しい量のデータが保存されるに至っており、その多くは組織体の日々の運営にとって極めて重要なものである。例えば、金融上の大量の取引は今や全面的に電子的に処理されている。航空会社でのビジネスでは、予約情報が万一喪失した場合には、大混乱になるリスクを背負っている。この結果として、データの信頼性を確保する為に、ローカルデータはデータが喪失した場合に備えて、一式以上のデータのコピーを通常しばしばリモートロケーションにバックアップしている。データが重要である程、バックアップ方法も精巧になる。例えば、影響の大きい重要なデータを守る一方法は、データをローカルストレージ施設とは地理的に離れたサイトにバックアップコピーを保存することである。各リモートストレージ施設は、ローカルストレージ施設のミラーイメージを維持して、ローカルストレージ施設のデータイメージが変更される度に、自施設のデータを変更してミラーを維持する。ローカルストレージシステムのデータをミラーするリモートストレージシステムの一例は、"Method and Apparatus for Mirroring Data in a Remote Data Storage System"のタイトルの米国特許5,933,653号に記されている。

30

【0003】

【特許文献1】米国特許5,933,653号

【特許文献2】米国特許6,092,066号

【発明の開示】

40

【発明が解決しようとする課題】

【0004】

リモートストレージ施設に転送される更新データは、しばしばキューされ、リモートコピー処理のオーバーヘッドを減らす為にグループ化して、インターネット等のネットワーク転送媒体を通して、送信される。この為、リモートサイトでミラーされるデータイメージは、ローカルサイトのものと常に同一であるわけではない。ローカルデータをミラーするのに、二つ以上のリモートストレージが使用された場合には、少なくとも全てが更新される迄は、リモートストレージのデータイメージは互いに異なっている。このような異なったデータイメージの存在は、ローカル施設が障害になった場合には、問題になる。ローカルストレージ施設が障害になると、あるリモートストレージ施設は、完全に正確

50

では無くともローカルストレージ施設の障害前のデータをより最新まで反映し、他の施設はより古い中間データイメージのままで、最新の更新処理は全く追従できない状態であることがあり得る。かくして、ローカルストレージ施設で障害が発生した場合には、システムを再起動する前に、全てのリモートストレージ施設が同一の最新データイメージを持つように、再同期することが必要になる。

【 0 0 0 5 】

もう一つの検討しなければならない問題は、リモートコピー動作中に“ 停止 (Suspension) ”が発生した場合での回復である。例えば、キャッシュオーバーフロー、コピー中のストレージシステム障害、リモートコピー動作中のネットワーク停止又はその他の介入障害等の予期せぬ事態での停止が発生した場合も、再同期が必要になる。リモートコピーでの再同期の一方法は、“Method and Apparatus for Independent Operation of a Remote Data Facility”のタイトルの米国特許 6,092,066 号に記されている。しかしながら、この特許に記されている技術では、限定された状況での再同期が可能になるのみである。例えば、リンク障害、キャッシュオーバーフロー、更にドライブ障害等の二つの障害が重なった、より以上に複雑なシステムサスペンションが発生した場合には、システム全体の再初期化を回避する為にすぐ使用できる再同期方法は存在しない。このような事態の場合、これらの技術は少なくとも二セットのコピーが残るように保障していない為、再同期の為には通常、ボリューム全体のコピーが必要になってしまう。

【 0 0 0 6 】

プライマリサイトが災害等の問題で障害になり、プライマリストレージサブシステムのデータが使用不能になったら、コンピュータシステムは、セカンダリストレージサブシステムのデータを使用してジョブを開始する。何れかのセカンダリストレージサブシステムのデータを使用してジョブを開始するのに先立って、他のセカンダリストレージサブシステムのデータ全てが同じデータを持つように同期を行う必要がある。ストレージシステムが同期されていないと、本ストレージシステム内のデータは不整合状態になってしまう。幾つかのストレージシステムが使用されている場合には、各ストレージシステムは他のストレージシステムについてコピーの進行状態（即ち、どのデータのコピーが終了しているか）について知る手段がない。ストレージシステム間のデータの相違については、人手によってつかむことは実質上不可能である。この結果として、生産用に使用する一つのストレージシステムの全データを他の全てのストレージシステムにコピーすることが必要になり得る。結果として、不必要に大量のデータが転送され、完了するのに長時間が必要になってしまう。

【課題を解決するための手段】

【 0 0 0 7 】

本発明の一様態では、プライマリストレージサブシステムと、少なくとも一のセカンダリストレージサブシステムを有するストレージシステムは、前記各ストレージサブシステムはそれぞれ、受信した一つ以上のデータブロックと各データブロックに付与された順序番号を格納するキューを備え、前記各セカンダリストレージサブシステムは、前記プライマリストレージサブシステムに直接、又は、他のセカンダリストレージサブシステムを介して間接的に接続されており、前記プライマリストレージサブシステムは、ホストシステムで生成された一つ以上のデータブロックを受信すると、各データブロックに対して順序番号を付与し、該プライマリストレージサブシステムから直接データブロックを受信するストレージサブシステムである後続ストレージサブシステムに対して、前記各データブロックと各データブロックに付与された順序番号を送信し、前記各セカンダリストレージサブシステムは、前記プライマリストレージサブシステムから送信された各データブロックと各データブロックに付与された順序番号を、前記プライマリストレージサブシステム、又は、他のセカンダリストレージサブシステムのいずれか一つから直接受信し、該セカンダリストレージサブシステムから直接データブロックを受信するストレージサブシステムである後続ストレージサブシステムに対して、前記各データブロックと各データブロックに付与された順序番号を送信し、該セカンダリストレージサブシステムが受信した最新の

データブロックに付与された順序番号と該セカンダリストレージサブシステムの後続ストレージサブシステムから受信した順序番号の中で最も古いデータブロックに付与された順序番号を、該セカンダリストレージサブシステムに直接データブロックを送信するストレージサブシステムである先行ストレージサブシステムに対して送信し、前記プライマリストレージサブシステムは、該プライマリストレージサブシステムの後続ストレージサブシステムから受信した順序番号の中から最も古いデータブロックに付与された順序番号を検索して、該順序番号を該後続ストレージサブシステムに報告し、前記各セカンダリストレージサブシステムは、該セカンダリストレージサブシステムの先行ストレージサブシステムから、前記プライマリストレージサブシステムにより検索された順序番号を受信し、当該順序番号を該セカンダリストレージサブシステムの後続ストレージサブシステムに報告するものである。

10

**【 0 0 0 8 】**

他の実施例では、複数のストレージサブシステムでは、第二のストレージサブシステムに接続され、第二のサブシステムから一つ以上のデータブロックを受信し保存する少なくとも一式の第三のストレージサブシステムを含む。第二のストレージサブシステムは第三のストレージサブシステムの先行ストレージサブシステムである。第三のストレージサブシステムは第二のストレージサブシステムの後続ストレージサブシステムである。

**【 0 0 0 9 】**

本発明の具体的な実施例では、ストレージサブシステムは、テーブル内の最小最終順序番号を検索するプロセッサを備える。プロセッサは、テーブル内の最小最終順序番号と各後続ストレージサブシステムで特定されている最小最終順序番号を比較して、本ストレージサブシステムの最小最終順序番号を決定する。プロセッサは、自サブシステムの最小最終順序番号を、先行ストレージサブシステムが存在すれば先行サブシステムに報告し、存在しなければ各後続ストレージサブシステムに報告する。ストレージサブシステムに先行ストレージサブシステムが存在すれば、ストレージサブシステムのプロセッサは、先行ストレージサブシステムからの最小最終順序番号の受信を契機に、受信した最小最終順序番号以下の順序番号を持つデータブロックをキューから削除する。ストレージサブシステムに先行ストレージサブシステムが存在しなければ、ストレージサブシステムのプロセッサは、ストレージシステムの最小最終順序番号を決定次第、決定された最小最終順序番号以下の順序番号を持つデータブロックをキューから削除する。いずれの場合でも、受信又は決定した最小最終順序番号を後続(存在すれば)するストレージサブシステムに伝える。

20

30

**【 0 0 1 0 】**

本発明のある実施例では、各ストレージサブシステムはプロセッサを持ち、一つ以上のストレージサブシステムの障害に於いては、自システムの最終順序番号レコードに保存されている最終順序番号を他の生存ストレージサブシステム内の最終順序番号レコードに保存されている最終順序番号と比較し、何れかの生存ストレージサブシステムに保存されている最終順序番号が、自システムの最終順序番号より大きい場合には、当該システムの自システムより大きい最終順序番号を持つデータブロックを、自ストレージサブシステムのキューにコピーする。

**【 0 0 1 1 】**

40

本発明の他の様態は、複数のストレージサブシステムを持つストレージシステムに保存されたデータを管理する方法に向けられる。本方法は、一つのストレージサブシステムから一つ以上のデータブロックをサブシステムに接続されている後続ストレージサブシステムにコピーし、各コピーされたデータブロックとこれに付与された順序番号を後続するストレージサブシステムのメモリ内のキューに保存し；後続するストレージサブシステムのメモリ内の最終順序番号レコードに最終順序番号を保存し；各ストレージサブシステムに各後続するストレージサブシステムで最小最終順序番号として特定された各最終順序番号をメモリ内のテーブルの各対応する欄に保存する。

**【 0 0 1 2 】**

本発明の他の様態は、複数のストレージサブシステムを持つストレージシステムに保存

50

されたデータを管理する方法に向けられる。本方法は、一式のストレージサブシステムからこれに接続された後続ストレージサブシステムの一つ以上のデータブロックをコピーし、データブロックとこれに付与された順序番号を後続ストレージサブシステムのメモリ内のキューに保存し；ストレージサブシステムのデータブロックの最終順序番号から最小最終順序番号を決定し；決定された最小最終順序番号以下の順序番号を持つデータブロックを各ストレージサブシステム内のキューから削除する。

【発明の効果】

【0013】

本発明を実施する事により、より進歩したデータ処理ストレージシステムが可能になり、本ストレージシステムでは、データは一つのストレージサブシステムから他の後続するストレージサブシステムに下り方向にコピーされ、コピー状態に関する情報は一つのストレージサブシステムから他の先行するストレージサブシステムに上り方向に伝えられる。本ストレージサブシステムは、典型的には、一つのプライマリストレージサブシステムと複数のセカンダリストレージサブシステムで構成される。データがプライマリストレージサブシステムから一つ以上のセカンダリストレージサブシステムに下り方向にコピーされる時には、各セカンダリストレージサブシステムは、コピー状態の更新情報を時刻順に最終順序番号を用いてキューの中に保存し、更に受信データの最終順序番号に関する情報を先行するストレージサブシステムに上り方向に報告する。ストレージサブシステム又はサイトが障害になると、本最終データを受信したストレージサブシステムのキューが、最終順序番号として決定され特定される。他のストレージサブシステムの各々に於いては、受信されていないデータの順序番号が検出され、これらのデータが当該ストレージサブシステムに送信される。各ストレージサブシステム内の最終順序番号は共有される為、各ストレージシステムは、他のストレージサブシステムとの同期に不必要な消去可能データをキュー内で、認識することが出来る。手短に言うと、各ストレージサブシステムのキューに保存された受信データの最終順序番号は、他のストレージサブシステムと交換又は共有され、障害時のデータ同期の為にストレージサブシステム間でデータをコピーするのに使用される。障害時のデータ同期の為に各ストレージサブシステムのキュー内に保存されるデータ量は最少に維持され、障害時のデータ同期の為にストレージサブシステム間でコピーしなければならないデータ量も減少する。

【発明を実施するための最良の形態】

【0014】

図1は、ホストシステム101と複数のストレージサブシステムを持つストレージシステム100を示す。ストレージシステムには一式のプライマリストレージサブシステム102と複数のセカンダリストレージサブシステム111、112が存在する。ホストシステム101は、入出力が転送される入出力インタフェース103を通してプライマリストレージサブシステム102に結合している。プライマリストレージサブシステム102に入力されたデータ又は情報は、メモリ105を使用してCPU106により処理され、ディスクやRAID又はJBOD(Just a Bunch of Disks)と呼ばれるストレージ媒体に保存される。メモリ

105は、以下に詳細に述べるように、データブロックとこれに付与された順序番号を保存する為のキュー108、最終順序番号を保存する為のレコード109、及び各後続ストレージシステムの最終順序番号を保存する為のテーブル113を収容する。内部バス110は内部データの転送に用いられる。ネットワークインタフェース104は、ネットワークに接続され、セカンダリストレージサブシステム等の他のシステムがデータをコピーし、情報を送受信する為の通信を可能にする。データは、プライマリストレージサブシステム102からセカンダリストレージサブシステム111、112への下り方向にコピーされる。ディスク107からのデータは、プライマリストレージサブシステム102から、ネットワークに結合するネットワークインタフェース104を通して、一つ以上の対象セカンダリストレージサブシステムに流れる。他のストレージシステム(セカンダリストレージサブシステム)は、メモリ105、プロセッサ106、及び記憶媒体107を含めて

同様な構成を持つことが望ましい。

【 0 0 1 5 】

図 2 は、ホストシステム 2 0 8、プライマリストレージサブシステム 2 0 1、及び複数のセカンダリストレージサブシステム 2 0 2 ~ 2 0 5 を含むストレージシステム 2 0 0 の論理システム構成を示す概略図である。図 1 で示すように、プライマリストレージサブシステム 2 0 1 は、ホストシステム 2 0 8 からの入出力が実行されるシステムである。このストレージシステム 2 0 0 では、プライマリストレージサブシステム 2 0 1 のデータは、セカンダリストレージサブシステム 2 0 2 ~ 2 0 5 にコピーされる事になっている。ストレージシステム 2 0 0 は、枝分かれ形式のストレージシステム構造になっており、セカンダリストレージサブシステム 2 0 2、2 0 3 は第一レベルのセカンダリストレージサブシステムで、セカンダリストレージサブシステム 2 0 4、2 0 5 は第二レベルのセカンダリストレージサブシステムである。データは、プライマリストレージサブシステム 2 0 1 から、プライマリストレージサブシステム 2 0 1 に直接後続するシステムとして、第一レベルのセカンダリストレージサブシステム 2 0 2 と 2 0 3 にコピーされる。データは次いで、第一レベルのセカンダリストレージサブシステム 2 0 2 に直接後続するシステムとして、第二レベルのセカンダリストレージサブシステム 2 0 4 と 2 0 5 に、第一レベルのセカンダリストレージサブシステム 2 0 2 からコピーされる。別の言い方をすると、プライマリストレージサブシステム 2 0 1 は、第二ストレージシステム 2 0 2 と 2 0 3 に先行する第一ストレージシステムで、一方、第三ストレージシステム 2 0 4 と 2 0 5 は、第二ストレージシステム 2 0 2 に後続するストレージシステムである。各ストレージシステムは、結合を自動検出するか人為入力によって、結合しているストレージサブシステムを認識している。

【 0 0 1 6 】

コピープロセスに於いては、データブロックは、一つのストレージシステムから他のストレージシステムに、ホストシステム 2 0 8 がプライマリストレージサブシステム 2 0 1 に書き込んだ順序通りに、転送される。各データブロックには、プライマリストレージサブシステム 2 0 1 にて、典型的には # 1、# 2・・・の昇順に、順序番号が付与される。

【 0 0 1 7 】

図 3 は、セカンダリストレージサブシステムの一つなどのストレージサブシステムでのコピープロセスを示す。サブシステムは、データブロックとその順序番号をサブシステムに直接接続する先行サブシステムから受信する（ステップ 3 0 1）。ここで使用されているように、ストレージサブシステム間の関係を記述するのに、便宜上、“先行する”なる用語は“直接に先行する”意味で、“後続する”なる用語は“直接に後続する”意味とする。ステップ 3 0 2 に於いて、データブロックとその順序番号は、書き込み順序を保持する為に、最初に、サブシステムのメモリ中のキュー（例えば、図 1 のメモリ 1 0 5 内のキュー 1 0 8 を参照）に保存し、データは保存媒体（例えば、図 1 のディスク 1 0 7 を参照）に保存する。キュー 1 0 8 に収めたデータは、この直後か又は一定時間経過後にディスク 1 0 7 に収める。本データは又、キュー 1 0 8 に収める前にディスク 1 0 7 に格納しても良い。入力データの最終ブロックに伴う最終順序番号は、最終順序番号レコード（例えば、図 1 のレコード 1 0 9 を参照）に収める（ステップ 3 0 3）。ストレージサブシステムが、一式以上の後続ストレージサブシステムを持っておれば、入力データとこの順序番号は、この各後続ストレージサブシステムに転送する（ステップ 3 0 5）。例えば、図 2 のストレージサブシステム 2 0 2 は、後続ストレージサブシステム 2 0 4 と 2 0 5 を持ち、ストレージサブシステム 2 0 1 はストレージサブシステム 2 0 2 の先行ストレージサブシステムである。ストレージサブシステムが後続ストレージサブシステムを持っていないければ、受信した最終順序番号は、上りの流れとして先行ストレージサブシステムに戻される（ステップ 3 0 6）。

【 0 0 1 8 】

図 4 は図 3 のコピープロセスに関連するコピーステータス更新プロセスを説明するフロー図である。コピーステータス更新プロセスは、ストレージサブシステムが後続ストレ

10

20

30

40

50

ジサブシステムから最終順序番号を受信したときに開始する（ステップ401）。ストレージサブシステムは、テーブル内の、各後続ストレージサブシステムに対応する最終順序番号欄（図1のテーブル113を参照）を受信した順序番号で更新する（ステップ402）。このテーブル113は、各後続ストレージサブシステムの最終順序番号を保有する。図6はテーブルの一例を示す。ストレージサブシステムは、このテーブル内の最小最終順序番号を検索する（ステップ403）。最小最終順序番号は、ストレージシステムの配下に於ける最古のデータブロックに伴う順序番号である。ステップ403でのテーブル検索は、ストレージサブシステムが後続ストレージサブシステムの一つから最終順序番号を受信し次第、又は、例えば、プライマリストレージサブシステム201又は図2のホストシステム208で事前設定されたある一定時間たってから開始しても良い。ストレージシステムが構成されてコピーが開始した直後の期間に於いては、テーブル113の各最終順序番号欄は、対応する後続ストレージサブシステムの最終順序番号で満たされていない、事態が発生し得る。この場合には、本検索は、各欄が対応する後続ストレージサブシステムの最終順序番号で満たされてから、開始される。もとより、ストレージサブシステムがただ一つの後続ストレージサブシステムしか持っていない場合には、本検索及び比較処理は省略することができる。次いで、ストレージサブシステムがプライマリストレージサブシステムの場合には、決定された最小最終順序番号を全ての後続ストレージサブシステムに報告し（ステップ404）、先行ストレージサブシステムが存在すれば、先行ストレージサブシステムに本決定された最小最終順序番号を連絡する（ステップ405）。最小最終順序番号には、これを識別する為にヘッダーを加えても良い。図4に示すステップ404と405はこの順番に実行する必要はなく、この順序は逆転可能なことに注意すること。

10

20

#### 【0019】

図5は、コピーステータス更新プロセスでのコピーステータス更新情報に基づいて、データを削除するプロセスを示すフロー図である。データ削除プロセスは、ストレージサブシステムがこれの先行ストレージサブシステムから最小最終順序番号を受信することを契機に開始する（ステップ501）。ストレージサブシステムは、キュー（例えば、図1のキュー108を参照）内のデータブロックで、受信した最小最終順序番号以下の順序番号を持つ全てのデータブロックを削除しても良い（ステップ502）。このデータ削除は、ストレージサブシステムが自らの先行ストレージサブシステムから最小最終順序番号を受信した直後に実行しても、一定時間経過してから実行しても良い。先行ストレージサブシステムを持たないプライマリストレージサブシステム201に於いては、最小最終順序番号は、テーブルを検索して、各後続ストレージサブシステムから受信した最小最終順序番号の最小値を見つけることにより、図4のステップ403で決定される。プライマリストレージサブシステムは、上記で決まった最小最終順序番号以下の順序番号を持つキュー108内のデータブロックを削除できる。何れの場合でも、受信した又は決定した最小最終順序番号は各後続（存在すれば）するセカンダリストレージサブシステムに伝えられ（ステップ503）、本最小最終順序番号がステップ501で受信される。

30

#### 【0020】

本実施例では、コピーステータス更新プロセスとデータ削除プロセスの採用によって、サブシステム障害に於いて、不必要なデータ転送を行うことなく複数ストレージサブシステム間でデータ同期を行う方法が提供される。好都合なことに、プライマリストレージサブシステムのデータバックアップに、何台のストレージサブシステムが配置されているか、又これらのお互いの接続構成に関する全体のストレージシステム構造について、各ストレージサブシステムは認識する必要はない。各ストレージサブシステムは、サブシステム障害時に必要な情報を交換し同期に必要な冗長情報（即ち、最終順序番号を保存しサブシステム障害時のデータ同期に必要なデータをコピーする）を保持するのに、（コピーデータを受け取る）直接先行するストレージサブシステムと、（コピーデータを送信する）直接後続するストレージサブシステムの情報のみを知ればよい。

40

#### 【0021】

50

図7はコピープロセス中の各ストレージサブシステムのある瞬間の構成を抽出した構成図である。ホストシステム701はプライマリストレージサブシステム702に接続され、本システム702にデータを書き込む。プライマリストレージサブシステム702のデータは、これに接続されている第一レベルの後続セカンダリストレージサブシステム703と704にコピーされる。第一レベルのセカンダリストレージサブシステム704のデータは、これに接続されている第二レベルの後続セカンダリストレージサブシステム705と706にコピーされる。

#### 【0022】

図7はコピー動作が進行中の各ストレージサブシステムのある瞬間の構成を示す。データブロック70101～70105は発生順に並んでおり、ホストシステム701で生成されプライマリストレージサブシステム702のキュー108に書き込まれる。図7に示す通り、データブロック70201～70205は、キュー108内の書き込み済ブロックである。この時点では、データブロック70201～70204は、セカンダリストレージサブシステム703にコピー済で、データブロック70301～70304に対応するキューに格納済である。データブロック70201～70203は、セカンダリストレージサブシステム704にコピー済で、データブロック70401～70403に対応するキューに格納済である。セカンダリストレージサブシステム704内のデータブロック70401～70402は、セカンダリストレージサブシステム706にコピー済で、データブロック70601～70602に対応するキューに格納済である。データブロック70401は、セカンダリストレージサブシステム705にコピー済で、データブロック70501に対応するキューに格納済である。

#### 【0023】

最終順序番号レコード109と各後続装置の最終順序番号のテーブル113は、それぞれ70404と70406として、ストレージサブシステム704のみにについて示されている。ストレージサブシステム704が後続ストレージサブシステム705、706から最終順序番号（各々、レコード70407内の#1、レコード70408内の#2）を受信すると、本ストレージサブシステム704は、テーブル70406内の最小最終順序番号（レコード70407内の#1）を先行するストレージサブシステム702に報告する。

#### 【0024】

ストレージサブシステム702に於いては、最小最終順序番号#1はストレージサブシステム704から受け取ったもので、最終順序番号#4（これはストレージサブシステム703での最終順序番号でもある）はストレージサブシステム703から受け取ったものである。ストレージサブシステム702での最小最終順序番号はかくして#1で、本番号#1は後続ストレージサブシステム703、704に報告される。従って、ストレージサブシステム702、703、704は報告された最小最終順序番号（#1）以下の順序番号を持つデータブロックをいつでも削除可能である。次いで、ストレージサブシステム704はこの最小最終順序番号を第二レベルの後続ストレージサブシステム705と706に報告し、ストレージサブシステム705と706は報告された最小最終順序番号（#1）以下の順序番号を持つデータブロックをいつでも削除可能である。

#### 【0025】

図8は図7で示すデータ状態にデータ削除プロセスを適用した後の各ストレージサブシステムのデータ状態を抽出した構成図である。より具体的には、図8は、各データブロック#(80201、80301、80401、80501、80601)が各対応するストレージサブシステム(802、803、804、805、806)の各キューより削除された後の状態を示す。ストレージサブシステム障害の場合での、ストレージサブシステム間でのデータ同期のメカニズムが示されている。

#### 【0026】

一例として、プライマリストレージサブシステム802が障害又は使用不能になったとする。ストレージシステム構成で存在している最終データブロックは、ストレージサブシ

10

20

30

40

50



ステム 803 のデータブロック 80304 で、順序番号 #4 (即ち、最大の最終順序番号) を持っている。本順序番号は他の全てのストレージサブシステムに配布される。この配布は、データコピートランザクションに使用されるネットワークを通して、各ストレージサブシステムに配布論理を導入することによって行うか、あるいは各ストレージサブシステムと結合するメカニズムを有するサーバを追加して行うことができる。他のストレージサブシステム (804、805、806) は順序番号 #4 迄のデータブロックが必要である。ストレージサブシステム 804 は自らの最終順序番号 (#3) をレコード 109 内に持っており、データブロック 80304 のみをコピーすればよいことを知っている。このことをネットワーク又はサーバを通して連絡して、データブロック 80304 のコピー要求が発行され、本データブロックがストレージシステム 804 にコピーされる。ストレージサブシステム 805 はデータブロック 80304、80303、80302 を、ストレージサブシステム 806 はデータブロック 80304、80303 を、各々コピーする必要があることを知っている。コピーは、ストレージサブシステム 804 の場合と同じメカニズムを使用して実行される。

10

#### 【0027】

他の例として、ストレージサブシステム 804 が障害が使用不能だとする。更に、ストレージサブシステム 804、806 が使用不能だと仮定する。ストレージシステム構成に残っている最終データブロックは、ストレージシステム 802 の順序番号 #5 のデータブロック 80205 である。ストレージサブシステム 803 はデータブロック 80205 の、ストレージサブシステム 805 はデータブロック 80202 ~ 80205 のコピーをそれぞれ取得する必要がある。

20

#### 【0028】

各ストレージサブシステムは、最新データを持っているストレージサブシステムによって同期化される。データブロックはキュー内に冗長に保存され、どのストレージサブシステムが障害になっても、実質的にデータが失われることはない。キューの大きさを減少させる為に、不要になったデータブロックは削除される。

#### 【0029】

以上に述べた装置及び方法は、本発明の原理の応用を説明したもので、本発明の精神と請求項で定義された範囲を逸脱する事無く、多数の異なった実装や改変が可能である。従って、本発明の範囲は、これまでの記述を参照して決めるのではなく、添付する請求項の全範囲に従って判断する必要がある。

30

#### 【図面の簡単な説明】

#### 【0030】

【図1】図1は、本発明の一実施例による複数のストレージサブシステムを持つストレージシステムを説明する概略図である。

【図2】図2は、本発明の一実施例によるストレージシステムの論理的システム構成を説明する概略図である。

【図3】図3は、本発明の一実施例によるコピープロセスを示すフロー図である。

【図4】図4は、本発明の一実施例によるコピープロセスに関連するコピーステータス更新プロセスを示すフロー図である。

40

【図5】図5は、本発明の一実施例によるコピーステータス更新プロセスのコピーステータス更新情報に基づくデータ削除プロセスを示すフロー図である。

【図6】図6は、本発明の一実施例によるコピーステータス更新情報を含むテーブルを示す。

【図7】図7は、本発明の一実施例によるコピープロセス中でのストレージシステム内の各ストレージサブシステムのある瞬間の構成を抽出した図である。

【図8】図8は、本発明の一実施例による、コピープロセスと関連コピーステータス更新プロセス及びデータ削除プロセスにおけるストレージシステム内での各ストレージサブシステムのデータ状態を抽出した図である。

#### 【符号の説明】

50

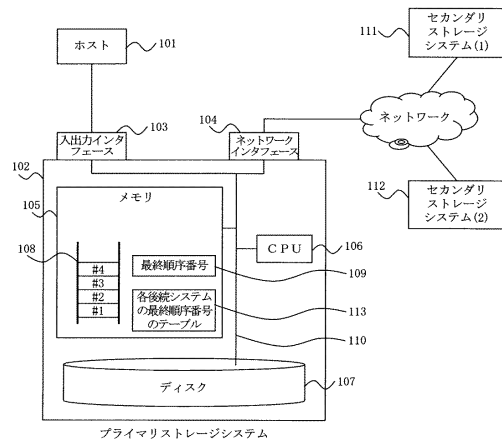
## 【 0 0 3 1 】

1 0 1	ホスト
1 0 2	プライマリストレージサブシステム
1 0 3	入出力インタフェース
1 0 4	ネットワークインタフェース
1 0 5	メモリ
1 0 6	C P U
1 0 7	ディスク
1 0 8	キュー
1 0 9	最終順序番号
1 1 0	内部バス
1 1 1	セカンダリストレージサブシステム( 1 )
1 1 2	セカンダリストレージサブシステム( 2 )
1 1 3	各後続システムの最終順序番号のテーブル

10

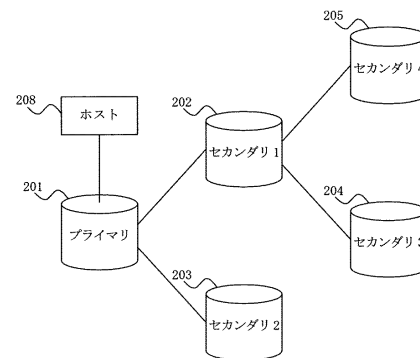
【 図 1 】

図 1



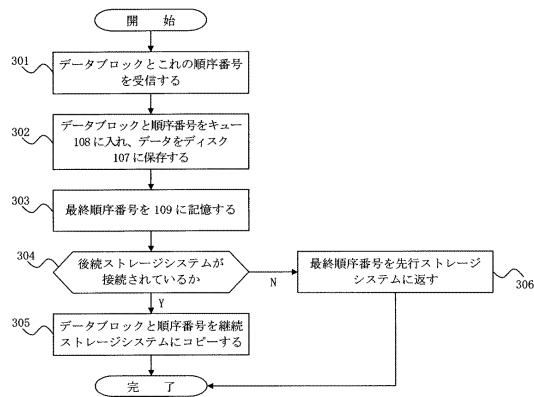
【 図 2 】

図 2



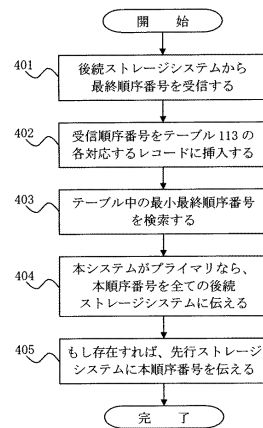
【図 3】

図 3



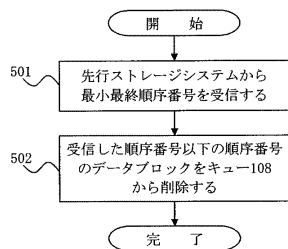
【図 4】

図 4



【図 5】

図 5



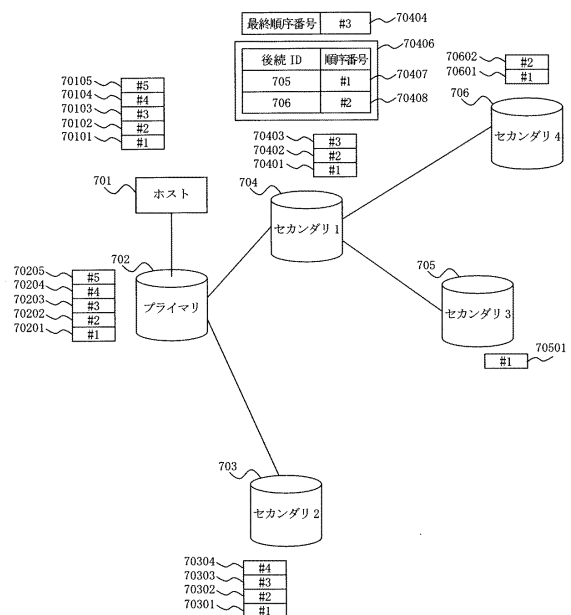
【図 6】

図 6

後続システム番号	最終順序番号
204	#3
205	#2

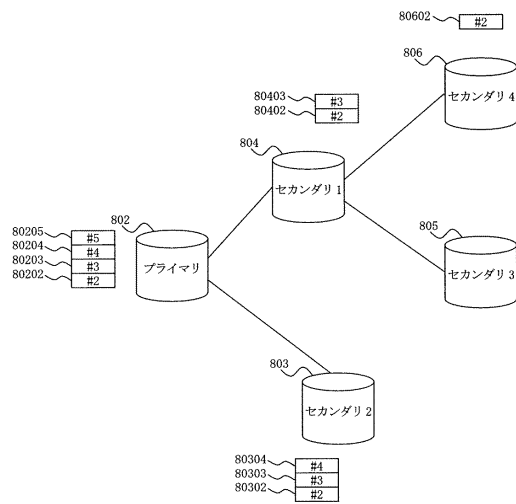
【図 7】

図 7



## 【図 8】

図 8



---

フロントページの続き

(56)参考文献 特開2003-263280(JP,A)  
特開2003-345524(JP,A)  
特開2003-122659(JP,A)  
特開平06-053973(JP,A)  
特開2001-168907(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06 - 3/08  
12/00  
13/00  
13/10 - 13/14