

(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2005-18738

(P2005-18738A)

(43) 公開日 平成17年1月20日(2005.1.20)

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

G06F 12/00

F I

G06F 12/00

5 3 1 J

G06F 12/00

5 0 1 B

テーマコード (参考)

5 B 0 8 2

審査請求 未請求 請求項の数 31 O L 外国語出願 (全 45 頁)

(21) 出願番号 特願2004-85961 (P2004-85961)  
(22) 出願日 平成16年3月24日 (2004.3.24)  
(31) 優先権主張番号 10/608391  
(32) 優先日 平成15年6月26日 (2003.6.26)  
(33) 優先権主張国 米国 (US)

(71) 出願人 000005108  
株式会社日立製作所  
東京都千代田区丸の内一丁目6番6号  
(74) 代理人 100075096  
弁理士 作田 康夫  
(74) 代理人 100100310  
弁理士 井上 学  
(72) 発明者 山神 憲司  
アメリカ合衆国カリフォルニア州ロスガト  
ス カイルニベル 108  
Fターム(参考) 5B082 CA14 DD08

(54) 【発明の名称】 ストレージベースのジャーナリングを用いてバックアップ及びリカバリを行う方法と装置

(57) 【要約】

【課題】

然しながら、データベースシステムを除くと、データを任意時点で回復する手段は用意されていない。データベースシステムにしても、ジャーナルを適用するには下記の手順の為に時間が必要である。

・ストレージ (例えばディスク) からジャーナルデータを読み込む。

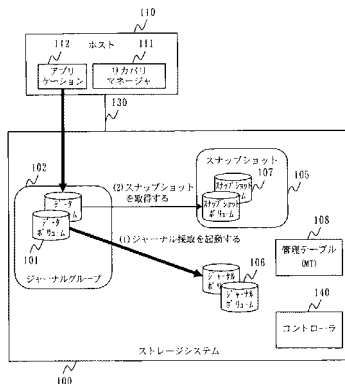
・ジャーナル内で目的のデータを特定するのに解析が必要である。

・データになされた活動を再生するには、ジャーナルデータを適切なデータイメージに適用しなければならない。このことは、イメージへのアクセスが通常必要でジャーナルが適用される度にデータの書き出しが必要である。

【解決手段】

ストレージシステムは、複数のジャーナルエントリと一つ以上のデータボリュームの少なくとも一つのスナップショットのジャーナルを維持する。ジャーナルとスナップショットの各々に発生順に固有な順序番号を割り当

図1 バックアップの概要



**【特許請求の範囲】****【請求項 1】**

アプリケーションデータの保存において、ホストからの書き込み要求の受信によって保存されるアプリケーションデータに対する少なくとも最初のスナップショットを生成するステップと、前記ホスト要求による各書き込み動作に対応してジャーナルエントリを生成するステップと、各ジャーナルエントリをジャーナルデータ領域に保存し、該ジャーナルエントリのリストを集積するステップと、前記ジャーナルデータ領域内の未使用容量をモニタするステップと、前記未使用容量が第一の限界値より減少したら、該未使用容量が第二の限界値より増大するまで、前記ジャーナルデータ領域から 1 ないしはそれ以上のジャーナルエントリを除去することで、未使用容量を増大させるステップと、から成ることを特徴とするデータ処理方法。 10

**【請求項 2】**

1 ないしはそれ以上のジャーナルエントリを除去する為に、前記最初のスナップショットにその時の最古のジャーナルエントリより始めて順にジャーナルエントリを適用することにより該最初のスナップショットを更新し、前記最初のスナップショットに適用されたジャーナルエントリが、前記ジャーナルエントリリストから除去されることで前記ジャーナルデータ領域の未使用容量を増大させることを特徴とする請求項 1 の方法。

**【請求項 3】**

1 ないしはそれ以上のジャーナルエントリを除去する為に、その時の最古のジャーナルエントリより古くてかつ最も近いスナップショットを特定して、該スナップショットを該最古のジャーナルエントリより始めて順に 1 ないしはそれ以上のジャーナルエントリを適用して更新することを特徴とする請求項 1 の方法。 20

**【請求項 4】**

1 ないしはそれ以上のジャーナルエントリを除去する為に、最新のスナップショットを見つけて、該最新のスナップショットより古いジャーナルエントリを前記ジャーナルデータ領域より除去することを特徴とする請求項 1 の方法。

**【請求項 5】**

前記第一の限界値と第二の限界値は異なった値であることを特徴とする請求項 1 の方法。

**【請求項 6】**

各ジャーナルエントリは固定長のヘッダ部分と可変長のデータ部分から構成され、前記ジャーナルデータ領域は複数のヘッダ部分が定義される第一のストレージ領域と複数のデータ部分を保存する第二のストレージ領域で構成され、ジャーナルエントリを生成する為にジャーナルヘッダのひとつを割り当て書き込み動作に対応するデータを収容するためにデータ部に十分なスペースを割り当てることを特徴とする請求項 1 の方法。 30

**【請求項 7】**

各ジャーナルエントリを保存する為に、ジャーナルヘッダが開始位置から始まって常に順番に割り付けられ、最終位置にジャーナルヘッダが割り当てられたら、該開始位置に戻って再度順番に割り当てられることを特徴とする請求項 6 の方法。

**【請求項 8】**

ホストからの書き込み要求の受信によって保存されるアプリケーションデータの少なくとも一部分に対する少なくとも最初のスナップショットを生成するステップと、ホストからの書き込み要求動作に対応してジャーナルデータ領域に保存され、該ジャーナルデータ領域の未使用容量を消費する複数のジャーナルエントリを記録するステップと、少なくとも前記最初のスナップショットを少なくとも一つのジャーナルエントリで更新して、該少なくとも一つのジャーナルエントリが消費している記憶スペースを未使用化して前記ジャーナルデータ領域の未使用容量を増大させるステップと、から成ることを特徴とするデータの処理方法。 40

**【請求項 9】**

前記更新ステップは、前記ジャーナルデータ領域の未使用容量を定期的にモニタして、 50

該未使用容量が第一の限界値以下になったら、少なくとも前記最初のスナップショットを更新するステップを含むことを特徴とする請求項 8 の方法。

【請求項 10】

少なくとも前記最初のスナップショットを更新するステップを前記未使用容量が第二の限界値を超えるまでジャーナルエントリに対して繰り返すことを特徴とする請求項 9 の方法。

【請求項 11】

更に、前記最初のスナップショットとジャーナルエントリに各々順序番号を割り当てるステップと、目標時刻を受信するステップと、前記ジャーナルエントリと前記最初のスナップショットの各々に割り当てられた順序番号に基づいて、開始ジャーナルエントリを決定するステップと、前記最初のスナップショットに前記開始ジャーナルエントリを適用することにより前記最初のスナップショットを更新するステップと、前記最初のスナップショットに、前記開始ジャーナルエントリに引き続く前記目標時刻以前のジャーナルエントリを適用することにより、前記最初のスナップショットの追加更新を実行するステップと、を含むことを特徴とする請求項 8 の方法。

10

【請求項 12】

更に、追加のスナップショットを生成して、複数のスナップショットを集積するステップと、前記スナップショットとジャーナルエントリに各々順序番号を割り当てるステップと、目標時刻を受信するステップと、前記目標時刻に基づいて一つのスナップショットを選択するステップと、前記ジャーナルエントリと前記選択されたスナップショットの各々の順序番号に基づいて、一つの開始ジャーナルエントリを決定するステップと、前記選択されたスナップショットに前記開始ジャーナルエントリを適用することにより前記選択されたスナップショットを更新するステップと、前記開始ジャーナルエントリに引き続き、前記目標時刻以前のジャーナルエントリを前記選択されたスナップショットに適用することにより、前記選択されたスナップショットの追加更新を実行するステップと、から成ることを特徴とする請求項 8 の方法。

20

【請求項 13】

前記選択されたスナップショットは、最初のスナップショットであることを特徴とする請求項 12 の方法。

【請求項 14】

前記選択されたスナップショットは前記目標時刻以前でかつ最も近いスナップショットであることを特徴とする請求項 12 の方法。

30

【請求項 15】

ホストからの書き込み要求の受信によって保存されるアプリケーションデータの 1 ないしはそれ以上の部分で構成される複数のスナップショットを生成するステップと、ホストから要求される各書き込み動作に対応して生成され、ジャーナルデータ領域に保存されリスト状に集積され、該ジャーナルデータ領域の未使用容量を消費するジャーナルエントリを記録するステップと、前記ジャーナルデータ領域内の残未使用容量をモニタするステップと、前記残未使用容量が第一の限界値より少なくなったら、一つのスナップショットを選択し該選択されたスナップショットより古いジャーナルエントリを除去することにより 1 ないしはそれ以上のジャーナルエントリを除去するステップと、から成ることを特徴とするデータ処理方法。

40

【請求項 16】

前記選択されたスナップショットは、その時の最古のジャーナルエントリより少なくとも決められた時間だけ新しいスナップショットであることを特徴とする請求項 15 の方法。

【請求項 17】

更に、前記スナップショットとジャーナルエントリに各々順序番号を割り当てるステップを含み、前記選択されたスナップショットは、その時の最古のジャーナルエントリの順序番号より決められた量だけ大きい順序番号を持つスナップショットであることを特徴と

50

する請求項 15 方法。

【請求項 18】

ホストからの書き込み要求の受信によって保存されるアプリケーションデータの複数のスナップショットを生成するステップと、ホストから要求される各書き込み動作に対応してジャーナルエントリを生成するステップと、各ジャーナルエントリをジャーナルデータ領域に保存し、ジャーナルエントリのリストを集積するステップと、前記ジャーナルデータ領域内の未使用容量をモニタするステップと、前記未使用容量が第一の限界値より少ないときには、その時の最古のジャーナルエントリより古くてかつ最も近いスナップショットを選択するステップと、前記最古のジャーナルエントリより始まる 1 ないしはそれ以上のジャーナルエントリにより前記選択したスナップショットを更新し、前記ジャーナルエントリのリストを減少させるステップと、前記未使用容量が第二の限界値を超えるまで追加のジャーナルエントリによる更新を繰り返すステップとから成ることを特徴とするデータ処理方法。

10

【請求項 19】

更に、前記スナップショットとジャーナルエントリに各々順序番号を割り当てるステップを含み、前記選択されたスナップショットは、該スナップショットの順序番号とその時の最古のジャーナルエントリの順序番号とに基づいて決定されることを特徴とする請求項 18 の方法。

【請求項 20】

各ジャーナルエントリは、ジャーナルエントリのシーケンシャルリストで表現されるように、順番に連続してジャーナルデータ領域に保存され、ジャーナルエントリのリストが終端に達したら、最初のリスト部分に折り重ねて保存し該領域を再使用することを特徴とする請求項 18 の方法。

20

【請求項 21】

ホストからの書き込み要求の受信によって保存されるアプリケーションデータの複数のスナップショットを生成するステップと、ホストから要求される各書き込み動作に対応して生成され、ジャーナルデータ領域に保存され、該ジャーナルデータ領域内でリスト上に集積されるジャーナルエントリを記録するステップと、複数の順序番号を生成するステップと、前記スナップショットとジャーナルエントリの各々に、前記順序番号を順に割り当てるステップと、前記ジャーナルデータ領域内の未使用容量をモニタするステップと、前記未使用容量が第一の限界値より少なくなったら、各スナップショットの順序番号をその時点の最古のジャーナルエントリの順序番号と比較して、一つのスナップショットを選択するステップと、前記選択したスナップショットを前記最古のジャーナルエントリに始まる 1 ないしはそれ以上のジャーナルエントリで順に更新するステップと、前記未使用容量が第二の限界値を超えるまで、追加のジャーナルエントリによる更新を繰り返すステップと、から成ることを特徴とするデータ処理方法。

30

【請求項 22】

ジャーナルエントリを記録するステップはスナップショット生成に先立って起動され、最初のスナップショットが生成されている期間に発生した各書き込み動作はジャーナルエントリに記録されることを特徴とする請求項 21 の方法。

40

【請求項 23】

前記選択されたスナップショットは、前記最古のジャーナルエントリより以前のスナップショットであることを特徴とする請求項 21 の方法。

【請求項 24】

前記選択されたスナップショットは、前記最古のジャーナルエントリより以前でかつ他のスナップショットに比べて前記最古のジャーナルエントリに最も近いスナップショットであることを特徴とする請求項 21 の方法。

【請求項 25】

前記第一の限界値は前記第二の限界値と異なっていることを特徴とする請求項 21 の方法。

50

## 【請求項 26】

データストレージシステムに送られた各書き込み動作に対応してジャーナルエントリを生成し、複数のジャーナルエントリをバックアップストレージシステムに保存する、ジャーナル処理を起動するステップと、前記ジャーナル処理起動に引き続いて、前記データストレージシステムに保存されたデータの一部分に対して一ないしはそれ以上のスナップショットを取得するステップと、各ジャーナルエントリとスナップショットが生成される度に、同じシーケンス系列で構成される順序番号を各々順に割り当てるステップと、前記バックアップストレージシステム内の使用可能ストレージ容量に基づいて、該バックアップストレージシステム内で前記ジャーナルエントリによって消費されているスペース量を減少させる為に、一つのスナップショットに幾つかのジャーナルエントリを適用すべきかどうか決定するステップと、から成ることを特徴とするデータストレージシステムにおけるデータ処理方法。 10

## 【請求項 27】

一つのスナップショットに幾つかのジャーナルエントリを適用すると決定されたら、該ジャーナルエントリとスナップショットの各々に割り当てられた順序番号に基づいて、一つのスナップショットを選択するステップと、前記選択されたスナップショットより新しくかつ最古のジャーナルエントリを特定するステップと、前記最古のジャーナルエントリから始まる一つ以上のジャーナルエントリで前記選択されたスナップショットを更新するステップと、から成ることを特徴とする請求項 26 の方法。 20

## 【請求項 28】

ホストからの書き込み動作を受信するように構成された生産用データ領域と、前記生産用データ領域の少なくとも一部分のスナップショットを一つ以上保存するように構成されたスナップショットデータ領域と、ホストからの各書き込み動作に対応してジャーナルエントリを生成し、一つ以上のジャーナルエントリを保存するように構成されたジャーナルデータ領域と、前記生産用データ領域とスナップショットデータ領域にアクセスするステップと、前記ジャーナルデータ領域にアクセスするステップと、全体の中の幾つかのジャーナルエントリにより一つのスナップショットを更新し、幾つかのジャーナルエントリをデータ回復の為に温存するステップとを実行するコントローラと、から構成されることを特徴とするデータ処理の為にストレージシステム。 30

## 【請求項 29】

前記コントローラは更に、前記スナップショットとジャーナルエントリに各々順序番号を割り当てるステップと、目標時刻を受信するステップと、前記目標時刻に基づき一つのスナップショットを選択するステップと、前記ジャーナルエントリと選択されたスナップショットの各順序番号に基づき一つの開始ジャーナルエントリを決定するステップと、前記開始ジャーナルエントリを前記選択されたスナップショットに適用することにより前記選択されたスナップショットを更新するステップと、前記開始ジャーナルエントリより新しく、前記目標時刻より以前のジャーナルエントリを前記選択されたスナップショットに適用することにより、前記選択されたスナップショットに追加の更新を実行するステップと、を実行するように構成されていることを特徴とする請求項 28 のストレージシステム。 40

## 【請求項 30】

一つ以上の書き込み動作を含む交信をホストから受信する手段と、生産用データを保存する為の第一の保存手段と、スナップショットを保存する為の第二の保存手段と、ジャーナルエントリを保存する為の第三の保存手段と、前記第二の保存手段に、生産用データの少なくとも一部の一つ以上のスナップショットを保存し、各書き込み動作に対応してジャーナルエントリを前記第三の保存手段に記録し、複数のジャーナルエントリを生成し、該第三の保存手段の未使用容量を減少させ、一つ以上のジャーナルエントリに従ってスナップショットの一つを更新して、前記第三の保存手段の未使用容量を増大させる、コントロール手段と、から構成されることを特徴とするデータ処理の為にストレージシステム。 50

**【請求項 3 1】**

前記更新が、周期的又は未使用容量に基づいた間隔で実施されることを特徴とする請求項 3 0 のストレージシステム。

**【発明の詳細な説明】****【技術分野】****【0001】**

本発明はコンピュータストレージに関わり、特に具体的にはデータのバックアップ及びリカバリに関連する。

**【背景技術】****【0002】**

データ損失を防止する為に幾つかの方法が慣用的に実施されている。典型的には、データはシステム管理者によって定期的に（例えば毎日一回）バックアップされる。ペリタス社のNetBackup、レガート社のNetworker等の多数のシステムがデータのバックアップとリカバリの為に出回っている。もう一つの技術はボリュームシャドウイングである。本技術では、データがプライマリストレージシステムに書き込まれているときにセカンダリストレージシステムにも書き込まれ、データのミラーイメージが生成される。

**【0003】**

ジャーナリングは、データベースシステムで一般に使用されているバックアップとリストアの技術である。バックアップされるデータのイメージが取得される。次いで、データが更新される度に更新ジャーナルが維持される。データのリカバリは、ジャーナルを回復対象データのイメージに適用する事によって行うことができ、任意時点のデータに回復される。オラクル等の典型的なデータベースシステムはジャーナリングを使用することができる。

**【発明の開示】****【発明が解決しようとする課題】****【0004】**

然しながら、データベースシステムを除くと、データを任意時点に回復する手段は用意されていない。データベースシステムにしても、ジャーナルを適用するには下記の手順の為に時間が必要である。

- ・ストレージ（例えばディスク）からジャーナルデータを読み込む。
- ・ジャーナル内で目的のデータを特定するのに解析が必要である。
- ・データになされた活動を再生するには、ジャーナルデータを適切なデータイメージに適用しなければならない。このことは、イメージへのアクセスが通常必要でジャーナルが適用される度にデータの書き出しが必要である。

**【0005】**

任意時点でのデータ回復は、次のような管理上の要求に応えなければならない。例えば、典型的な要求は、“誤って昨日10:00AM頃にファイルを消去してしまった。そのファイルを消去直前に回復しなければならない”といったことである。

**【0006】**

データがデータベースシステムに存在しなければ、この種類の要求は都合よく解決されているとは言えない。従って、損失データの回復を容易にするデータ処理方法のニーズが存在する。更にまた、データベースアプリケーション以外のユーザ環境でも、損失データの回復を容易にするデータ処理を可能とするニーズが存在する。

**【課題を解決するための手段】****【0007】**

本ストレージシステムが、スナップショット動作とジャーナリングを実行して、損失データを回復させる追加のデータ処理を実行する。スナップショットとジャーナルエントリは、ユーザの生産用データボリュームとは別に保存される。新しいジャーナルエントリを作る為に、より古いジャーナルエントリは除去される。本除去は、一つ以上の古いジャーナルエントリを対応するスナップショットに適用して、本スナップショットを更新するこ

10

20

30

40

50

とにより達成される。その後の損失データの回復は、所定のスナップショットにアクセスして、本スナップショットにジャーナルエントリを適用して希望するデータ状態を再生することにより行われる。

【発明の効果】

【0008】

ストレージシステムが、ユーザがアプリケーションを実行する為のデータストレージサービスを可能にする。

【発明を実施するための最良の形態】

【0009】

図1は本発明に従うバックアップとリカバリシステムの一実施例を説明する上位階層の一般的ブロックダイアグラムである。システムが立ち上がると、スナップショットが生産用データボリューム(DVOL)101に対して取得される。ここで、“スナップショット”とは、本明細書では慣例に従い、ある指定時点に於けるデータボリュームのデータイメージのことを言う。システムに対する要求や実装等に対応して、スナップショットは、データボリュームの全体、又はその或る部分、又は複数の部分に対してなされる。本発明では、システムの正常運用の間にジャーナルエントリがホストから本データボリュームに発せられた全ての書き込み動作に対して作成される。以下に説明するように、所定のスナップショットに対して一連のジャーナルエントリを適用する事により、データは如何なる時点に対しても回復され得る。

【0010】

図1のバックアップとリカバリシステムには、少なくとも一式のストレージシステム100が存在する。図示されてはいないものの、本ストレージシステムには適切なプロセッサ、メモリ、及びホスト110とそのストレージ媒体(例えばディスク)との間の入出力動作を実行するコントロール回路が存在する事は当業者には明らかである。バックアップとリカバリのシステムには、更に少なくとも一式のホスト110が必要である。適切な通信バス130がホストとストレージシステムの間に存在する。ホスト110には、通常一つ以上のユーザアプリケーション(APP)112が稼動している。これらのアプリケーションは本ストレージシステム100のデータボリューム101に存在するストレージ媒体に対してデータの読み書きを実行する。

【0011】

かくして、アプリケーション112とデータボリューム101が保護されるべき対象リソースとなる。ユーザのアプリケーションで使用されるデータは、一つあるいはそれ以上のデータボリュームに保存され得るのは明らかである。本発明では、ジャーナルグループ(JNLG)102を定義する。データボリューム101は本ジャーナルグループの中に編成される。本発明では、ジャーナルグループは、ホスト110からデータボリュームへの書き込み動作のジャーナリングが保証されるデータボリュームの最小ユニットである。本関連するジャーナルは、ホストからデータボリュームへの書き込み動作を適切な順序で記録する。ジャーナリング処理によって生成されたジャーナルデータは、一つ以上のジャーナルボリューム(JVOL)106に保存される。ホスト110には更にリカバリマネージャ(RM)111が存在する。本コンポーネントは、バックアップとリカバリ動作の上位階層での協調を可能にする。本リカバリマネージャに関しては更に後に記す。

【0012】

ストレージシステム100はジャーナルグループを構成するデータボリュームのスナップショット(SS)105を生成する。例えば、スナップショット105は、本スナップショットが取得された時点のジャーナルグループ102内のデータボリューム101を反映している。スナップショットイメージを取得する為の慣用的方法は公知である。スナップショットデータを含むストレージシステムの中には、一つ以上のスナップショットボリューム(SVOL)107が存在する。一つのスナップショットを一つ以上のスナップショットボリュームに保存することも出来る。本発明の実施例では、ジャーナルデータとスナップショットデータを別々のストレージコンポーネントに保存しているが、これらデー

10

20

30

40

50

タを単一のストレージコンポーネントに保存しても良い事は勿論である。

【0013】

管理テーブル(MT)108はジャーナルグループ102、スナップショット105、及びジャーナルボリューム106に関する情報を保存する。管理テーブルの更なる詳細は図3とこれに関する以下の記述で説明する。コントローラコンポーネント140は、書き込み動作のジャーナリングとデータボリュームのスナップショット、及びこれに対応する各異なったストレージコンポーネント101、106、107間のデータ移動を調整する。ここで、コントローラコンポーネントは、ストレージシステム100内に分散する一つ以上のサブコンポーネントを物理的に含んだものの論理的表現である事に注意願いたい。

【0014】

図2はジャーナル実装で使用するデータを示す。書き込み要求がホスト110からストレージシステム100に到達すると、これに対応してジャーナルが生成される。ジャーナルはジャーナルヘッダ219とジャーナルデータ225で構成される。ジャーナルヘッダ219は、対応するジャーナルデータ225に関する情報を保持する。ジャーナルデータ225は書き込み動作で用いられるデータ(書き込みデータ)を含む。この種類のジャーナルは“事後ジャーナル”とも呼ばれる。ジャーナルヘッダ219にはオフセット番号(JH\_\_OFS)211が存在する。本オフセット番号はジャーナルグループ102内の本ジャーナルに対応するデータボリューム101を指定する。今回の実装では、本データボリュームは、データボリューム0、データボリューム1、データボリューム2・・・と番号付けられる。これに対応して本オフセット番号は0、1、2・・・等となる。

【0015】

ジャーナルヘッダ219内のアドレス(JH\_\_ADR)212フィールドは、(オフセット番号211で指定され、)データが書き込まれるデータボリューム内の開始アドレスを保存する。例えば、本アドレスはブロック番号LBA(Logical Block Address)で表現される。ジャーナルヘッダ219内のデータ長(JH\_\_LEN)213フィールドは、書き込みデータのデータ長を保存する。典型的にはこれはブロック数で表す。ジャーナルヘッダ219内の書き込み時刻(JH\_\_TIME)214フィールドは、書き込み要求がストレージシステム100に到着した時刻を保存する。本時刻はカレンダー日、時、分、秒、更にはミリ秒を含む。本時刻はディスクコントローラ140又はホスト110により提供される。例えば、メインフレーム環境では複数のメインフレームホストが、タイマを共有し書き込みコマンドを発行する時刻を提供する。

【0016】

各書き込み要求には順序番号が割り当てられ、ジャーナルヘッダ219のフィールド(JH\_\_SEQ)215に保存される。ジャーナルグループ102内では各順序番号は固有である。本順序番号は、ジャーナルエントリが生成された時に本ジャーナルエントリに対して割り当てられる。ジャーナルヘッダ219内のジャーナルボリューム識別子(JH\_\_VOL)216は、ジャーナルデータ225を保存するジャーナルボリューム106を指定する。即ち、本識別子はジャーナルデータを保持しているジャーナルボリュームを示す。ジャーナルデータは、ジャーナルヘッダを保存するジャーナルボリュームとは異なったジャーナルボリュームに保存され得る事に注意願いたい。ジャーナルヘッダ219内のジャーナルデータアドレス(JH\_\_JADR)217フィールドは、ジャーナルデータを保持している本ジャーナルボリューム106内でのジャーナルデータ225の開始アドレスを保存する。

【0017】

図2は、ジャーナルボリューム106は、ジャーナルヘッダ領域210とジャーナルデータ領域220の二つのデータ領域で構成される事を示す。ジャーナルヘッダ領域210はジャーナルヘッダ219のみを、ジャーナルデータ領域220はジャーナルデータ225のみを保存する。各ジャーナルヘッダは固定長のデータ構造である。ジャーナルヘッダはジャーナルヘッダ領域の開始点より順番に割り当てられる。本順番はジャーナルエントリの発生順に対応して割り当てられる。後に説明するが、各時点の“最古の”ジャーナル

10

20

30

40

50



エントリであるリスト中の最初のジャーナルエントリをポイントするデータが用意される。指定された順序番号（順序番号フィールド 215 内に存在）又は指定された書き込み時刻（時刻フィールド 214 内に存在）に対するジャーナルヘッダ 219 を検出する事が、一般的に必要なになる。

#### 【0018】

ジャーナルヘッダ 219 とジャーナルデータ 225 はジャーナルボリューム 106 の対応領域に発生順に収容されている。かくして、ジャーナルヘッダとジャーナルデータがジャーナルボリュームに保存される順番は、割り当てられた順序番号と同じ順番である。後に説明する様に、本発明の一態様として、本ジャーナル情報 219、225 は各対応領域 210、220 内で折り返される。

10

#### 【0019】

図 3 は管理テーブル 108（図 1）の詳細を示す。ジャーナルヘッダ領域 210 とジャーナルデータ領域 220 を管理するには、各領域に対するポインタが必要である。既に述べたように、本管理テーブルは、ジャーナルグループ 102 についての攻勢情報及び、ジャーナルグループと本ジャーナルグループに所属するジャーナルボリューム 106 及びスナップショットイメージ 105 間の関係についての構成情報を維持する。図 3 に示す管理テーブル 300 は管理テーブルの一例とその内容を示す。本管理テーブルは、ストレージシステム 100 内の特定の一つのジャーナルグループ 102 を指定するジャーナルグループ ID（GRID）310 を保持する。ジャーナルグループ名（GRNAME）311 はジャーナルグループを指定する、人が認識できる識別子である。

20

#### 【0020】

ジャーナルグループ属性（GRATTR）312 がジャーナルグループ 102 に対して定義される。本実装例ではマスタ及びリストアの二つの属性が定義される。本マスタ属性はジャーナルグループがジャーナルされていること示す。本リストア属性はジャーナルグループがジャーナルによりリストアされていることを示す。ジャーナル状態（GRSTS）315 がジャーナルグループ 102 に対して定義される。これには活性及び非活性の二つの状態がある。管理テーブルには順序カウンタ（SEQ）313 を保持するフィールドが存在する。本カウンタはジャーナルヘッダ 219 内で使用される順序番号の源となる。新しいジャーナルが生成される時に、本順序番号 313 が読み込まれ本新ジャーナルに割り当てられる。次いで本順序番号は加算され本管理テーブルに戻される。管理テーブル内の数（NUM\_DVOL）314 は、ジャーナルグループ 102 内に存在するデータボリューム 101 の数を示す。データボリュームリスト（DVOL\_LIST）320 はジャーナルグループ内のデータボリュームをリストする。

30

#### 【0021】

一実装例では図 3 に見られる通り、DVOL\_LIST はデータボリューム情報を保持するデータ構造への最初のエントリのポインタである。各データボリューム情報にはオフセット番号（DVOL\_OFFSET）321 が存在する。例えば、ジャーナルグループ 102 が三式のデータボリュームを持つなら、本オフセット値は 0、1、2 となる。データボリューム識別子（DVOL\_ID）322 は全ストレージシステム 100 内でデータボリュームを固有に識別する。ポインタ（DVOL\_NEXT）324 は、ジャーナルグループ内で次のデータボリューム情報を保持しているデータ構造をポイントする。該当する情報が無ければナル値が与えられる。本管理テーブルは、ジャーナルグループ 102 に属するデータ（ジャーナルヘッダとジャーナルデータ）を保持する為に用いられているジャーナルボリューム数を保存する（NUM\_JVOL）330 フィールドを含んでいる。

40

#### 【0022】

図 2 で述べたように、ジャーナルヘッダ領域 210 は、各ジャーナルに対するジャーナルヘッダ 219 を保持する。ジャーナルデータ領域 220 についても同様である。既に述べたように、本発明の一態様として本データ領域 210、220 は折り返される。これにより、各データ領域のスペースは有限でもジャーナリングを続ける事が出来る。本管理テーブルは、データ領域 210、220 の異なった部分に対するポインタを保存する為のフ

50

フィールドを持ち、折り返しを可能にする。各フィールドは次のジャーナルエントリの格納場所を指示する。フィールド( J I \_ H E A D \_ V O L ) 3 3 1 は、次の新ジャーナルヘッダ 2 1 9 を保存するジャーナルヘッダ領域 2 1 0 を有するジャーナルボリューム 1 0 6 を指定する。フィールド( J I \_ H E A D \_ A D R ) 3 3 2 は、次のジャーナルヘッダを保存するジャーナルヘッダ領域の位置についてジャーナルボリューム 1 0 6 上のアドレスを指定する。

【 0 0 2 3 】

フィールド( J I \_ D A T A \_ V O L ) 3 3 5 は、次のジャーナルデータを保存するジャーナルデータ領域 2 2 0 を有するジャーナルボリュームを指定する。フィールド( J I \_ D A T A \_ A D R ) 3 3 6 は、次のジャーナルデータを保存するジャーナルボリューム 10 上のアドレスを指定する。かくして、次に書き込まれるジャーナルエントリは、“ J I \_ ” フィールド 3 3 1、3 3 2、3 3 5、3 3 6 の情報で“ ポイント ” される。管理テーブルは更に、その時の“ 最古の ” ジャーナルエントリを識別するフィールドを有する。

【 0 0 2 4 】

本情報の使用法を以下に記す。フィールド( J O \_ H E A D \_ V O L ) 3 3 3 は、最古のジャーナルヘッダ 2 1 9 を保持するジャーナルヘッダ領域 2 1 0 を保存するジャーナルボリュームを指定する。フィールド( J O \_ H E A D \_ A D R ) 3 3 4 は、ジャーナルヘッダ領域内の最古のジャーナルのジャーナルヘッダの位置のアドレスを指定する。フィールド( J O \_ D A T A \_ V O L ) 3 3 7 は、最古のジャーナルデータを保持するジャーナルデータ領域 2 2 0 を保存するジャーナルボリュームを指定する。本データを保存する 20 ジャーナルボリューム内のジャーナルデータ領域のデータアドレスは、フィールド( J O \_ D A T A \_ A D R ) 3 3 8 に保存される。管理テーブルには、特定のジャーナルグループ 1 0 2 に属するジャーナルボリュームのリスト( J V O L \_ L I S T ) 3 4 0 が存在する。本実装例では、J V O L \_ L I S T はジャーナルボリューム情報のデータ構造へのポインタである。

【 0 0 2 5 】

図 3 に示すように、各データ構造は、所定のジャーナルグループ 1 0 2 に属する特定のジャーナルボリューム 1 0 6 を指定するオフセット番号( J V O L \_ O F F S ) 3 4 1 を持つ。例えば、ジャーナルグループが二式のジャーナルボリューム 1 0 6 を持っていたら、各ジャーナルボリュームは 0 又は 1 で指定される。ジャーナルボリューム識別子( J V O L \_ I D ) 3 4 2 は、ストレージシステム 1 0 0 内でのジャーナルボリュームを固有に 30 指定する。最後にポインタ( J V O L \_ N E X T ) 3 4 4 は、本ジャーナルグループ内の次のジャーナルボリュームに関係する次のデータ構造のエントリをポイントする。存在しなければナル( N U L L ) 値が指定される。管理テーブルには、ジャーナルグループ 1 0 2 に属するスナップショットイメージ 1 0 5 のリスト( S S \_ L I S T ) 3 5 0 が存在する。

【 0 0 2 6 】

本実装例では図 3 に示すように、S S \_ L I S T はスナップショット情報のデータ構造へのポインタである。各スナップショット情報のデータ構造には、スナップショットが取得された時の順序番号( S S \_ S E Q ) 3 5 1 が与えられる。既に述べたように本番号は 40 順序カウンタ 3 1 3 よりもたらされる。時刻値( S S \_ T I M E ) 3 5 2 は本スナップショットが取得された時の時刻を示す。状態( S S \_ S T S ) 3 5 8 は各スナップショットの状態( 有効及び無効 ) の値を示す。ポインタ( S S \_ N E X T ) 3 5 3 は次のスナップショット情報のデータ構造をポイントする。存在しなければナル( N U L L ) 値が示される。各スナップショット情報のデータ構造は、更にスナップショットイメージ 1 0 5 を保存するスナップショットボリューム 1 0 7 ( 図 1 ) のリストを有する。

【 0 0 2 7 】

図 3 で分かるように、スナップショットボリューム情報のデータ構造へのポインタ( S V O L \_ L I S T ) 3 5 4 が、各スナップショット情報のデータ構造に存在する。各スナップショットボリューム情報のデータ構造には、スナップショットイメージの少なくとも 50

一部を含む各スナップショットボリュームを指定するオフセット番号 (SVOL\_\_OFFS) 355 が存在する。一式のスナップショットイメージはセグメント化又は分割されて、一つ以上のスナップショットボリュームに保存することが出来る。本実装例では、オフセット値はスナップショットイメージの一部分 (セグメント又は分割) を含む i 番目のスナップショットボリュームを指定する。i 番目のスナップショットイメージのセグメントは、i 番目のスナップショットボリュームに保存しても良い。各スナップショットボリューム情報のデータ構造には更に、ストレージシステム 100 内でのスナップショットボリュームを固有に指定するスナップショットボリューム識別子 (SVOL\_\_ID) 356 が存在する。ポインタ (SVOL\_\_NEXT) 357 は、指定されたスナップショットイメージに対する次のスナップショットボリューム情報のデータ構造をポイントする。

10

#### 【0028】

図 4 は本発明の実装例に従う、リカバリマネージャ 111 とストレージシステム 100 がバックアッププロセスを開始するプロセスを説明するフローチャートである。スナップショットが取得されていても、ジャーナルエントリが記録されていない書き込み動作は、データ回復処理中にデータは損失するか破損する。従って、本発明では最初のスナップショット取得前にジャーナリングプロセスを起動する。これを行うことにより、スナップショット取得後に発生する全ての書き込み動作がジャーナルされることが保障される。スナップショット完成前に記録されたジャーナルエントリは全て無視してもよい事を注記する。更に本発明に従えば、単一系列の順序番号 (SEQ) 313 が各スナップショットとジャーナルエントリに、生成される毎に割り当てられる。各スナップショットとジャーナルエントリに同一系列の順序番号を割り当てる目的は以下に述べる。

20

#### 【0029】

図 4 に戻って、リカバリマネージャ 111 は、ジャーナルグループ (JNLG) 102 が定義済でなければ、これをステップ 410 で定義する。

#### 【0030】

図 1 で示した通り、本定義では、ジャーナルの対象となる一台以上のデータボリューム (DVOL) 101 を指定し、ジャーナル関連情報を保存する一台以上のジャーナルボリューム (JVOL) 106 を指定する。リカバリマネージャは、目的を達成する為に、ストレージシステム 100 との間で適切な一連のやり取りを実行する。

#### 【0031】

ステップ 415 にて、ストレージシステムは図 3 のテーブル 300 の詳細図に示される多様な情報を収容する管理テーブル 108 (図 1) を生成する。とりわけ、本プロセスでは、ジャーナルグループ 102 を構成するジャーナルボリュームをリストする JVOL\_\_LIST 340 の初期化を行う。同様に、データボリュームのリスト DVOL\_\_LIST 320 を生成する。次のジャーナルエントリ (又はテーブルが最初に生成される場合には最初のジャーナルエントリ) を指定するフィールドを初期化する。かくして、JI\_\_HEAD\_\_VOL 331 はジャーナルボリュームリスト内の最初を指定し、JI\_\_HEAD\_\_ADR 332 は、最初のジャーナルボリュームに存在するジャーナルヘッダ領域 210 の最初のエントリをポイントする。

30

#### 【0032】

同様に、JI\_\_DATA\_\_VOL 335 は、ジャーナルボリュームリスト内の最初を指定し、JI\_\_DATA\_\_ADR 336 は、最初のジャーナルボリュームに存在するジャーナルデータ領域 220 の開始点をポイントする。ヘッダ領域 210 とデータ領域 220 は異なるジャーナルボリュームに収容でき、JI\_\_DATA\_\_VOL は最初のジャーナルヘッダボリュームとは異なったジャーナルボリュームを指定することがあり得ることに注意していただきたい。ステップ 420 にて、リカバリマネージャ 111 はジャーナリングプロセスを起動する。ジャーナリングを実行する為に、ストレージシステム 100 に対して適切な通信がなされる。ステップ 425 にて、ストレージシステムはホスト 110 からの今後の各書き込み動作に対して、ジャーナルエントリ (これは、“事後ジャーナル”とも呼ばれる) を作成する。

40

50

## 【 0 0 3 3 】

図 3 を参照するに、ジャーナルエントリを作成するには、とりわけ、次のジャーナルエントリを保存するロケーションを指定することが必要になる。フィールド J I \_ H E A D \_ V O L 3 3 1 と J I \_ H E A D \_ A D R 3 3 2 は各々、次のジャーナルヘッダ 2 1 9 のジャーナルボリューム 1 0 6 とジャーナルヘッダ領域 2 1 0 内での位置を指定する。管理テーブルからの順序カウンタ ( S E Q ) 3 1 3 は、次のヘッダの J H \_ S E Q 2 1 5 フィールドにコピーされ割り当てられる。本順序カウンタは次いで加算され、管理テーブルに戻される。勿論、本順序カウンタは最初に加算し、ついで J H \_ S E Q にコピーし、最後に管理テーブルに戻してもよい。

## 【 0 0 3 4 】

管理テーブルの J I \_ D A T A \_ V O L 3 3 5 と J I \_ D A T A \_ A D R 3 3 6 フィールドは、各々書き込み動作に付随するデータを保存する、ジャーナルボリュームとジャーナルデータ領域 2 2 0 の開始点を指定する。 J I \_ D A T A \_ V O L と J I \_ D A T A \_ A D R は各々、ジャーナルヘッダの J H \_ J V O L 2 1 6 と J H \_ A D R 2 1 2 にコピーされ、この結果ジャーナルヘッダがジャーナルデータに対応するポインタを得る事になる。書き込み動作のデータは保存される。

## 【 0 0 3 5 】

J I \_ H E A D \_ V O L 3 3 1 と J I \_ H E A D \_ A D R 3 3 2 フィールドは、次のジャーナルエントリに対する次のジャーナルヘッダ 2 1 9 をポイントするべく更新される。このことはジャーナルヘッダ領域 2 1 0 内の次に続くジャーナルヘッダエントリを用いる結果となる。同様に、 J I \_ D A T A \_ V O L フィールドと J I \_ D A T A \_ A D R フィールドは、各々次のジャーナルエントリに対するジャーナルデータボリュームとジャーナルデータ領域での開始点を反映するべく更新される。このことは、ジャーナルデータ領域における次の保存可能位置に進めることを意味する。従ってこれらのフィールドは、ジャーナルエントリのリストをポイントすると見なすことが出来る。リストのジャーナルエントリは、ジャーナルヘッダ領域 2 1 0 内のジャーナルヘッダ 2 1 9 のシーケンシャル構成によって、互いにリンクすることが出来る。

## 【 0 0 3 6 】

ジャーナルヘッダ領域 2 1 0 の最終点に達したら、次のジャーナルエントリに対するジャーナルヘッダ 2 1 9 は、ジャーナルヘッダ領域の開始点に戻される。ジャーナルデータ 2 2 5 に対しても同様である。以前のジャーナルエントリに上書きする事を防止する為に、本発明はジャーナルボリューム 1 0 6 のエントリを開放する処置を実施する。本件については、以下に述べる。ジャーナル起動後の最初のジャーナルエントリに対しては、 J O \_ H E A D \_ V O L 3 3 3 、 J O \_ H E A D \_ A D R 3 3 4 、 J O \_ D A T A \_ V O L 3 3 7 、 及び J O \_ D A T A \_ A D R 3 3 8 の各フィールドは、各対応する “ J I \_ ” フィールドの初期値がセットされる。後に説明するように、“ J O \_ ” フィールドはその時の最古のジャーナルエントリをポイントする。かくして、新しいジャーナルエントリが作成されると、“ J I \_ ” フィールドは前進するが “ J O \_ ” フィールドは前進しない。“ J O \_ ” フィールドの更新については以下で説明する。

## 【 0 0 3 7 】

図 4 のフローチャートに戻って、ジャーナリングプロセスが起動している場合には、ホストからのすべての書き込み動作はジャーナルされる。そこでステップ 4 3 0 に於いて、リカバリマネージャ 1 1 1 はデータボリューム 1 0 1 のスナップショットの取得を起動する。ストレージシステム 1 0 0 は、リカバリマネージャからスナップショットを取得する指示を受け取る。ステップ 4 3 5 にて、ストレージシステムはデータボリュームのスナップショットを取得するプロセスを開始する。この結果特に、管理テーブル ( 図 3 ) 上の S S \_ L I S T 3 5 0 にアクセスすることになる。次のスナップショットを扱う為に必要な量のメモリがフィールド 3 5 1 - 3 5 4 に割り当てられる。順序カウンタ ( S E Q ) 3 1 3 は、上記の J H \_ S E Q 2 1 5 で説明したように、 S S \_ S E Q 3 5 1 フィールドにコピーされ、加算される。

10

20

30

40

50

## 【 0 0 3 8 】

かくて、順序番号は S E Q 3 1 3 より生成され、順序カウンタ中の各番号はジャーナルエントリ又はスナップショットエントリのいずれかに処理中を通して割り当てられる。スナップショットは一台（以上）のスナップショットボリューム（S V O L）1 0 7 に保存される。フィールド 3 5 5 - 3 5 7 に必要な量のメモリが割り当てられる。スナップショットを保存する為の S V O L に関連する情報がフィールド 3 5 5 - 3 5 7 に保存される。スナップショット保存に追加のボリュームが必要なら、フィールド 3 5 5 - 3 5 7 に追加のメモリが割り当てられる。

## 【 0 0 3 9 】

図 5 は、ジャーナルエントリとスナップショットの関係を示す。スナップショット 5 2 0 は、ジャーナルグループ 1 0 2 に属するデータボリューム 1 0 1 の最初のスナップショットイメージを表す。順序番号 S E Q 0 と S E Q 1 のジャーナルエントリ（5 1 0）が作られ、これらは二つの書き込み動作に対するジャーナルエントリを表していることに注意していただきたい。これらのエントリは、ジャーナリングはスナップショットが取得（ステップ 4 2 0）されるに先立って起動されている事を示す。かくして、順序番号 S E Q 2 に対応する時点で、リカバリマネージャ 1 1 1 はスナップショット取得を開始し、更にジャーナル採取は既に起動されているため、スナップショット取得開始後の全ての書き込み動作は、ジャーナルされている。かくして、順序番号 S E Q 3 以上の書き込み動作は、全てジャーナルされている。見て分かる通り、順序番号 S E Q 0 と S E Q 1 のジャーナルエントリは破棄するか無視しなければならない。

## 【 0 0 4 0 】

データのリカバリでは、指定時点のデータボリューム 1 0 1 の少なくとも一部分のデータ状態に回復する事が通常要求される。このことは必要量のジャーナルエントリを本ジャーナルエントリより以前に取得されたスナップショットに適用する事で一般的に実行される。本実施例で開示する方法では、順序番号 S E Q 3 1 3 は、ジャーナルエントリ又はスナップショットに割り当てられる度に、加算される。従って、所定のスナップショットに対して適用可能なジャーナルエントリを特定することは単純である。即ち、（J H \_ S E Q）2 1 5 の順序番号が本スナップショットの順序番号（S S \_ S E Q）3 5 1 より大きなジャーナルエントリを適用すればよい。

## 【 0 0 4 1 】

管理者は或る時点、例えば、データボリュームのデータが失われたか又は破損した時（“目標時刻”）以前の時刻を、データ回復の時点に指定することが出来る。目標時刻以前の時刻が見つかるまで、各スナップショットの S S \_ T I M E 3 5 2 の時刻フィールドが検索される。次に、その時の“最古の”ジャーナルヘッダから始めて、ジャーナルヘッダ領域 2 1 0 内のジャーナルヘッダ 2 1 9 が検索される。本最古のジャーナルヘッダは、管理テーブルの“J O \_”フィールド 3 3 3、3 3 4、3 3 7 及び 3 3 8 によって指定されている。ジャーナルヘッダは、順序番号 J H \_ S E Q 2 1 5 が対象スナップショットの順序番号 S S \_ S E Q 3 5 1 より大きい最初のヘッダを目指して領域 2 1 0 内を順番に検索される。選択されたスナップショットには、一時には一ずつ、各ジャーナルエントリを順次適用する事により、順次加算更新され一連の書き込み動作が再現される。本処理は、ジャーナルエントリの J H \_ T I M E 2 1 4 フィールドの値が目標時刻より以前である限り継続される。本更新作業は時刻フィールド 2 1 4 の値が目標時刻を過ぎると停止する。

## 【 0 0 4 2 】

本発明の一実施例では、単一のスナップショットが取得される。本スナップショットに引き続く全てのジャーナルエントリは、所定時刻におけるデータ状態を再構築する為の適用対象となる。本発明の他の実施例として、複数スナップショットが取得される。本件は図 5 A に示され、複数スナップショット 5 2 0 ' が取得される。本発明では、各スナップショットとジャーナルエントリは各々が記録される順に順序番号が割り当てられる。各スナップショット 5 2 0 ' の間には通常、多数のジャーナルエントリ 5 1 0 が存在する事に注意願いたい。複数スナップショットを使用する事により、データ回復の時間を短縮する

10

20

30

40

50

事が出来る。目標回復時刻に対して最寄りのスナップショットが選択される。本スナップショットに続くジャーナルエントリが所望のデータ回復の為に適用可能である。

#### 【0043】

図6は本発明のもう一つの態様を説明する。本発明によれば、ジャーナルエントリはホストからの書き込み動作の度に作成され、これはかなり大量のジャーナルエントリ数になる。ジャーナル採取が繰り返されると時間の経過と共に、ジャーナルグループ102に対してリカバリマネージャ111が定義した各ジャーナルボリューム106は、何時かは満杯になってしまう。このままでは、この時点でこれ以上のジャーナルエントリは作成できなくなり、結果として、後続する書き込み動作はジャーナルされず、ジャーナルボリュームが満杯になった以降のデータ状態の回復は不可能になってしまう。

10

#### 【0044】

図6は、ストレージシステム100が、“オーバーフロー”条件の検出を契機に、相応しいスナップショットへジャーナルエントリを適用する状況を示している。ここで、“オーバーフロー”は、ジャーナルボリュームでの未使用容量が設定された限界値を下回った場合に発生したと見做される。オーバーフロー限界値の設定には多様な基準があることに注意願いたい。一つの分かりやすい限界値は、ジャーナルグループに割り当てられたジャーナルボリュームの合計ストレージ容量を基準にするものである。未使用容量が全ストレージ容量のある割合、例えば10%、に達したら、オーバーフローとする。他の限界値としては各ジャーナルボリューム対応に決めても良い。本発明の一態様では、ジャーナルボリューム内の未使用容量は定期的にモニタされる。反対に未使用容量は非定期的にモニタしても良い。例えば、モニタ間隔はランダムであっても良い。更に他の方法として、例えば、モニタ間隔は未使用容量の関数として変動させても良い。或いは又、ジャーナル採取を契機にモニタしても良い。

20

#### 【0045】

図7は、ストレージシステム100内でオーバーフロー条件の検出の為にを行うプロセスを説明する。かくして、ステップ710にて、ストレージシステムはジャーナルボリューム106の全未使用容量を定期的に、例えば10秒毎にチェックする。未使用容量は、管理テーブル300内のポインタ(JI\_\_HEAD\_\_VOL331、JI\_\_HEAD\_\_ADR332等)がジャーナルボリュームで消費しているストレージの現在の状態を維持しているので、容易に計算可能である。未使用容量が限界値以上であれば、次の周期を待って未使用容量のチェックを繰り返すのみである。未使用容量が前もって決定した限界値以下になったら、ステップ720にて、幾つかの、特に最古のジャーナルエントリがスナップショットに適用され、スナップショットを更新する。

30

#### 【0046】

図3を参照するに、本最古のジャーナルエントリのジャーナルヘッダ219は、JO\_\_HEAD\_\_VOLフィールド333とJO\_\_HEAD\_\_ADRフィールド334で指定される。これらのフィールドは、本最古のジャーナルエントリを保存するジャーナルボリュームとジャーナルヘッダ領域210に対する本ジャーナルボリューム内の位置を指定する。同様に、本最古のジャーナルエントリのジャーナルデータは、JO\_\_DATA\_\_VOLフィールド337とJO\_\_DATA\_\_ADRフィールド338で指定される。これらのフィールドで指定されるジャーナルエントリはスナップショットに適用される。ここで選択されるスナップショットは、ジャーナルエントリより以前でかつ本ジャーナルエントリに順序番号が最も近いスナップショットである。

40

#### 【0047】

かくして、本実装例では、順序番号が加算される度に、ジャーナルエントリに対して最も近くてより小さい(時刻的にはより早い)順序番号を持つスナップショットが選択される。ジャーナルエントリがスナップショットに適用されてスナップショットが更新されると、適用済のジャーナルエントリは開放される。このことは単に、JO\_\_HEAD\_\_VOLフィールド333、JO\_\_HEAD\_\_ADRフィールド334、JO\_\_DATA\_\_VOLフィールド337、及びJO\_\_DATA\_\_ADRフィールド338を次のジャーナルエ

50

ントリをポイントすべく更新するのみである。順序番号は最終的には折り返され 0 から再度スタートすることは当業者には明らかである。順序番号を比較する時に、本再スタートを把握するための相応しいメカニズムも当業者の技術範囲である。

#### 【0048】

図 7 に戻って、スナップショットにジャーナルエントリを適用してスナップショットを更新した後は、この結果として適用済ジャーナルエントリが開放され（ステップ 730）て、ジャーナルボリューム内の未使用容量が増大した事をチェックする。本未使用容量は、更にステップ 730 にて限界値と比較される。この場合第二の異なった限界値が使用できる。例えば、本プロセスを終了するには開始する時よりより高いレベルの未使用容量を必要としても良い。

10

#### 【0049】

これにより本プロセスが過剰に使用される事を防止する事が可能になるが、一旦開始されたら、第二の高いレベルの限界値により合理的な範囲で出来るだけ多くの未使用容量を確保して回復を促進する。これらの限界値を決めるには管理者の時間をかけた経験が要求される事に注意願いたい。かくして、ステップ 730 により未使用容量が第二の限界値を超え停止条件に達したら、本プロセスは終了する。そうでなければ、ステップ 720 に戻り次に古いジャーナルエントリに対して繰り返される。ステップ 730 と 720 が、未使用容量がステップ 730 にて本限界値に達するまで繰り返される。

#### 【0050】

図 7 A は、図 7 のステップ 720 に対する代替の実施例のサブステップを示す。ステップ 720 では、ジャーナルエントリを、本ジャーナルエントリより以前の最新のスナップショットに適用する事によって、開放する。然しながら、複数のスナップショットが使用可能なら、スナップショットを更新する為にジャーナルエントリを適用するプロセスで消費する時間を省略する事が出来る。

20

#### 【0051】

図 7 A は、図 7 のステップ 720 に対する代替ステップ 720' の詳細を示す。ステップ 721 にて、本最古のジャーナルエントリより新しいスナップショットが存在するかどうかを判定する。本判定は、順序番号が本最古のジャーナルエントリより大きい最初のスナップショットを検索する事により、可能である。代わって本判定は、その時の最古のジャーナルエントリより所定時間（例えば少なくとも 1 時間）だけ新しいスナップショットを見つけてこれを基準にしても良い。更に又他の方法として、スナップショットとジャーナルエントリの時間差に代わって、順序番号を用いても良い。例えば、最古のジャーナルエントリの順序番号より、N 分だけ進んだ順序番号を持つスナップショットを選択しても良い。

30

#### 【0052】

該当するスナップショットがステップ 721 で見つかったら、これより古いジャーナルエントリはスナップショットに適用しないまま除去しても良い。かくして、ステップ 722 にて、“JO\_\_”フィールド（JO\_\_HEAD\_\_VOL333, JO\_\_HEAD\_\_ADR334, JO\_\_DATA\_\_VOL337, 及び JO\_\_DATA\_\_ADR338）は、選択されたスナップショットより新しい最初のジャーナルエントリをポイントする値に単純に書き換えられる。該当するスナップショットが見つからなかったら、ステップ 720 で議論したように、ステップ 723 にて、残った最古のジャーナルエントリがこれより古いスナップショットに適用される。

40

#### 【0053】

ステップ 721 に対する更なる代替案として、単純に最も新しいスナップショットを選択する案がある。本スナップショットより小さな順序番号を持つジャーナルエントリは開放可能である。このことはさらに単純に、最新のスナップショットより大きな順序番号を持つ最初のジャーナルエントリをポイントするように、“JO\_\_”フィールドを更新するのみである。ここで、本発明の一態様はデータを希望する時点の状態に回復する事が出来る事を思い出してほしい。このことは、生成されたジャーナルエントリを必要範囲は保存

50

して、本ジャーナルエントリを所定のスナップショットに適用して、書き込み動作を再現する事により達成される。この最後の実施例は大量のジャーナルエントリを除去する事が出来るが、これによってデータ状態が回復できる時間帯を制限する結果になってしまう。然しながら、構成によっては、大量のジャーナルエントリを除去する事は与えられた動作環境にとっては望ましい事もあり得る。

#### 【 0 0 5 4 】

これまで述べたステップは全体として、コントローラ 1 4 0 (例えばディスクコントローラ) で実現される事に注意願いたい。このことは、特定の実装次第で純粋のソフトウェア、カスタムロジック、ソフトウェアとハードウェアの相応しい組合せで可能である。もっと一般的に言えば、前記の実現は典型的に、ソフトウェアとハードウェアの組合せでな  
10  
されると言ってよい。当業者には、最終的に採用する技術はシステムコスト、システム性能、慣用的ソフトウェアとハードウェア、動作環境その他を含めた要素等によって決まる事、しかしこれらに限定、制約されるものでない事を容易に理解してもらえらるであろう。これまで述べられ意図された実施例は、当業者により過大な実験を行う事無く具体的実装に直ちに適用可能である。

#### 【図面の簡単な説明】

#### 【 0 0 5 5 】

本発明の態様、効果と新規な機能が、本発明の以下の付随図面と共になされる説明によって明らかになる。

【図 1】は、本発明の一実施例を説明する上位階層の一般的ブロックダイアグラムである  
20

【図 2】は、本発明に従うジャーナルエントリを保存するデータ構造の一実施例を一般的に説明する。

【図 3】は、本発明に従うスナップショットボリュームとジャーナルエントリボリュームを管理するためのデータ構造の一実施例を一般的に説明する。

【図 4】は、リカバリマネージャとストレージシステム内のコントローラとの間での処理を説明する上位階層のフローチャートである。

【図 5】は、一つのスナップショットと複数のジャーナルエントリ間の関係を説明する。

【図 5 A】は、複数のスナップショットと複数のジャーナルエントリ間の関係を説明する  
30

【図 6】は、オーバーフロー状態が発生した場合のデータフローを説明する上位階層の図である。

【図 7】は、オーバーフロー状態を処理する為のストレージシステム内のコントローラ動作を説明する上位階層のフローチャートである。

【図 7 A】は、図 7 で示す処理ステップの別案を示す。

#### 【符号の説明】

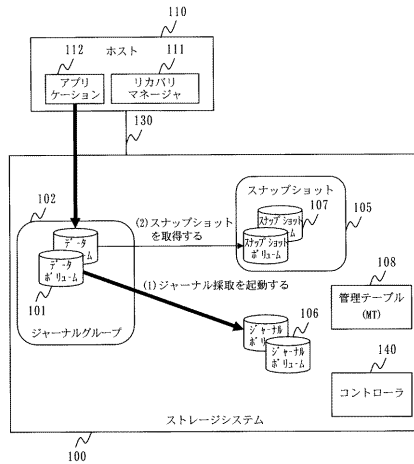
#### 【 0 0 5 6 】

1 0 0 . . . ストレージシステム  
1 0 1 . . . データボリューム ( D V O L )  
1 0 5 . . . スナップショット ( S S )  
1 0 6 . . . ジャーナルボリューム ( J V O L )  
1 0 7 . . . スナップショットボリューム ( S V O L )  
1 0 8 . . . 管理テーブル ( M T )  
1 0 2 . . . ジャーナルグループ ( J N L G )  
1 1 0 . . . ホスト  
1 1 1 . . . リカバリマネージャ ( R M )  
1 1 2 . . . アプリケーション ( A P P )  
1 4 0 . . . コントローラ  
40



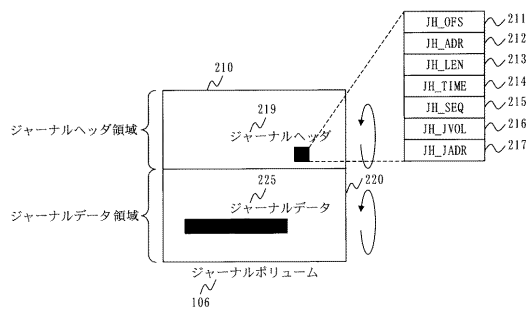
【図 1】

図1 バックアップの概要



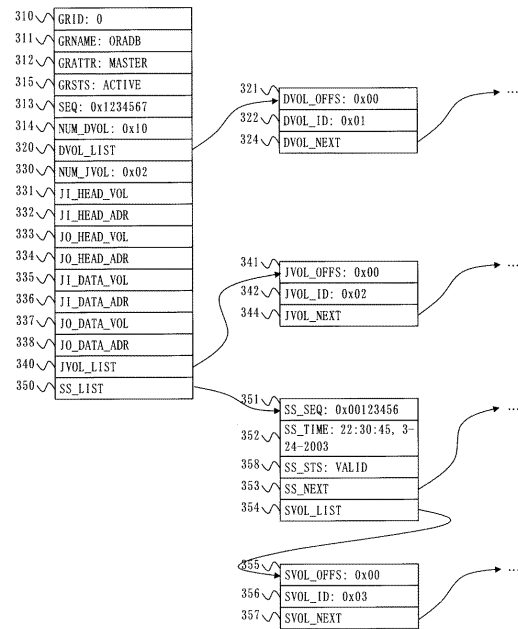
【図 2】

図2 ジャーナル用のコントロールデータ



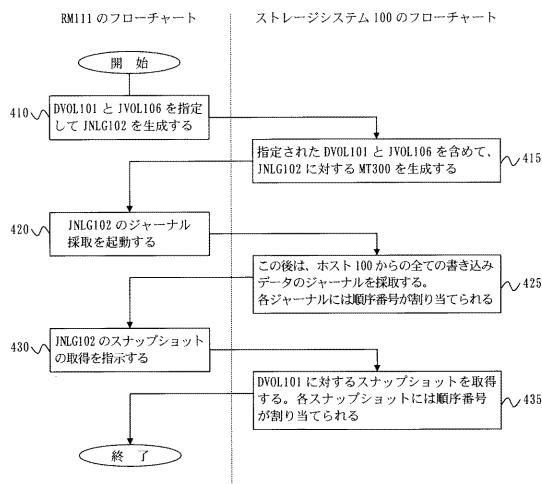
【図 3】

図3 MT（管理テーブル）300



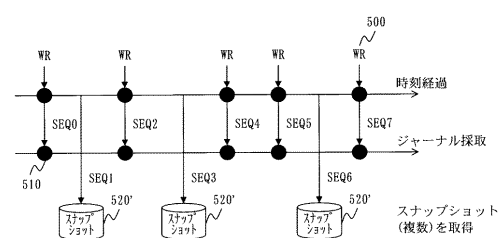
【図 4】

図4 ジャーナルの起動



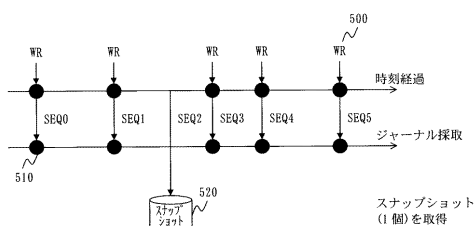
【図 5 A】

図5A 複数スナップショットとジャーナルの関係



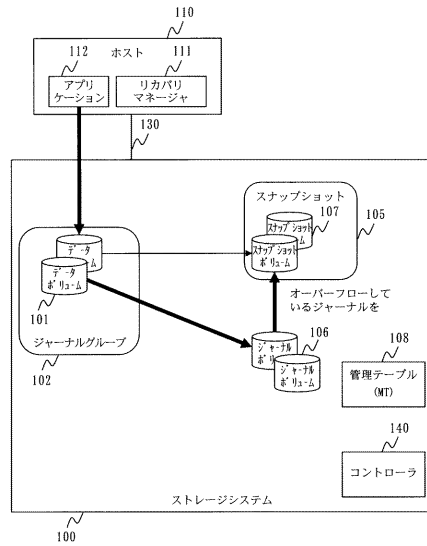
【図 5】

図5 スナップショットとジャーナルの関係



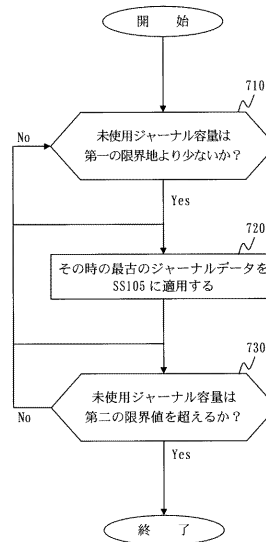
【図 6】

図 6 ジャーナルのオーバーフロー



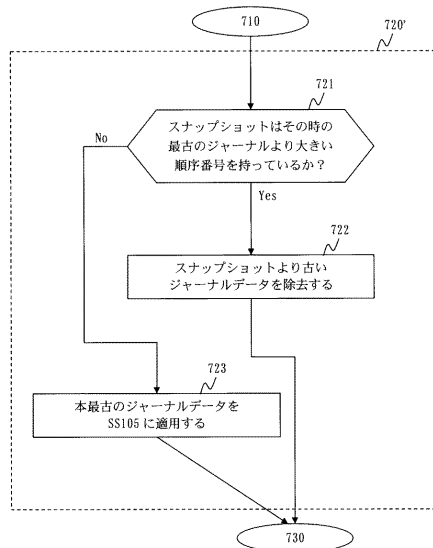
【図 7】

図 7 オーバフロージャーナル



【図 7 A】

図 7 A



---

フロントページの続き

【要約の続き】

ることにより、スナップショットに適用すべきジャーナルを見つける事が容易になる。ジャーナルスペースが枯渇するオーバーフロー条件を検出し、ジャーナルスペースを回復する技術が述べられる。

【選択図】図1

## 【 外国語明細書 】

PATENT

Attorney Docket No.: 16869B-082700US

Client ref. No.: HAL 279

**METHOD AND APPARATUS FOR BACKUP AND RECOVERY USING  
STORAGE BASED JOURNALING**

## CROSS-REFERENCES TO RELATED APPLICATIONS

5 [01] This application is related to the following commonly owned and co-pending U.S.  
applications:

“Method and Apparatus for Data Recovery Using Storage Based Journaling,”

Attorney Docket Number 16869B-082800US, and

10 “Method and Apparatus for Synchronizing Applications for Data Recovery Using  
Storage Based Journaling,” Attorney Docket Number 16869B-082900US,  
both of which are herein incorporated by reference for all purposes.

## BACKGROUND OF THE INVENTION

15 [02] The present invention is related to computer storage and in particular to backup and  
recovery of data.

[03] Several methods are conventionally used to prevent the loss of data. Typically, data is  
backed up in a periodic manner (e.g., once a day) by a system administrator. Many systems  
are commercially available which provide backup and recovery of data; e.g., Veritas  
NetBackup, Legato/Networker, and so on. Another technique is known as volume  
20 shadowing. This technique produces a mirror image of data onto a secondary storage system  
as it is being written to the primary storage system.

[04] Journaling is a backup and restore technique commonly used in database systems. An  
image of the data to be backed up is taken. Then, as changes are made to the data, a journal  
of the changes is maintained. Recovery of data is accomplished by applying the journal to an  
25 appropriate image to recover data at any point in time. Typical database systems, such as  
Oracle, can perform journaling.

[05] Except for database systems, however, there are no ways to recover data at any point  
in time. Even for database systems, applying a journal takes time since the procedure  
includes:

- 30
- reading the journal data from storage (e.g., disk)
  - the journal must be analyzed to determine at where in the journal the desired data  
can be found

- apply the journal data to a suitable image of the data to reproduce the activities performed on the data - this usually involves accessing the image, and writing out data as the journal is applied

5 [06] Recovering data at any point in time addresses the following types of administrative requirements. For example, a typical request might be, "I deleted a file by mistake at around 10:00 am yesterday. I have to recover the file just before it was deleted."

[07] If the data is not in a database system, this kind of request cannot be conveniently, if at all, serviced. A need therefore exists for processing data in a manner that facilitates

10 recovery of lost data. A need exists for being able to provide data processing that facilitates data recovery in user environments other than in a database application.

#### SUMMARY OF THE INVENTION

[08] A storage system provides data storage services for users and their applications. The

15 storage system performs additional data processing to provide for recovery of lost data, including performing snapshot operations and journaling. Snapshots and journal entries are stored separately from the production data volumes provided for the users. Older journal entries are cleared in order to make for new journal entries. This involves updating a snapshot by applying one or more of the older journal entries to an appropriate snapshot.

20 Subsequent recovery of lost data can be provided by accessing an appropriate snapshot and applying journal entries to the snapshot to reproduce the desired data state.

#### BRIEF DESCRIPTION OF THE DRAWINGS

[09] Aspects, advantages and novel features of the present invention will become apparent

25 from the following description of the invention presented in conjunction with the accompanying drawings:

Fig. 1 is a high level generalized block diagram of an illustrative embodiment of the present invention;

Fig. 2 is a generalized illustration of a illustrative embodiment of a data

30 structure for storing journal entries in accordance with the present invention;

Fig. 3 is a generalized illustration of an illustrative embodiment of a data structure for managing the snapshot volumes and the journal entry volumes in accordance with the present invention;

Fig. 4 is a high level flow diagram highlighting the processing between the

35 recovery manager and the controller in the storage system;

Fig. 5 illustrates the relationship between a snapshot and a plurality of journal entries;

Fig. 5A illustrates the relationship among a plurality of snapshots and a plurality of journal entries;

5 Fig. 6 is a high level illustration of the data flow when an overflow condition arises;

Fig. 7 is a high level flow chart highlighting an aspect of the controller in the storage system to handle an overflow condition; and

Fig. 7A illustrates an alternative to a processing step shown in Fig. 7.

10

#### DESCRIPTION OF THE SPECIFIC EMBODIMENTS

[10] Fig. 1 is a high level generalized block diagram of an illustrative embodiment of a backup and recovery system according to the present invention. When the system is activated, a snapshot is taken for production data volumes (DVOL) 101. The term  
15 "snapshot" in this context conventionally refers to a data image of at the data volume at a given point in time. Depending on system requirements, implementation, and so on, the snapshot can be of the entire data volume, or some portion or portions of the data volume(s). During the normal course of operation of the system in accordance with the invention, a journal entry is made for every write operation issued from the host to the data volumes. As  
20 will be discussed below, by applying a series of journal entries to an appropriate snapshot, data can be recovered at any point in time.

[11] The backup and recovery system shown in Fig. 1 includes at least one storage system 100. Though not shown, one of ordinary skill can appreciate that the storage system includes suitable processor(s), memory, and control circuitry to perform IO between a host 110 and its  
25 storage media (e.g., disks). The backup and recovery system also requires at least one host 110. A suitable communication path 130 is provided between the host and the storage system.

[12] The host 110 typically will have one or more user applications (APP) 112 executing on it. These applications will read and/or write data to storage media contained in the data  
30 volumes 101 of storage system 100. Thus, applications 112 and the data volumes 101 represent the target resources to be protected. It can be appreciated that data used by the user applications can be stored in one or more data volumes.

[13] In accordance with the invention, a journal group (JNLG) 102 is defined. The data volumes 101 are organized into the journal group. In accordance with the present invention,

a journal group is the smallest unit of data volumes where journaling of the write operations from the host 110 to the data volumes is guaranteed. The associated journal records the order of write operations from the host to the data volumes in proper sequence. The journal data produced by the journaling activity can be stored in one or more journal volumes(JVOL) 106.

5 [14] The host 110 also includes a recovery manager (RM) 111. This component provides a high level coordination of the backup and recovery operations. Additional discussion about the recovery manager will be discussed below.

[15] The storage system 100 provides a snapshot (SS) 105 of the data volumes comprising a journal group. For example, the snapshot 105 is representative of the data volumes 101 in  
10 the journal group 106 at the point in time that the snapshot was taken. Conventional methods are known for producing the snapshot image. One or more snapshot volumes (SVOL) 107 are provided in the storage system which contain the snapshot data. A snapshot can be contained in one or more snapshot volumes. Though the disclosed embodiment illustrates separate storage components for the journal data and the snapshot data, it can be appreciated  
15 that other implementations can provide a single storage component for storing the journal data and the snapshot data.

[16] A management table (MT) 108 is provided to store the information relating to the journal group 102, the snapshot 105, and the journal volume(s) 106. Fig. 3 and the accompanying discussion below reveal additional detail about the management table.

20 [17] A controller component 140 is also provided which coordinates the journaling of write operations and snapshots of the data volumes, and the corresponding movement of data among the different storage components 101, 106, 107. It can be appreciated that the controller component is a logical representation of a physical implementation which may comprise one or more sub-components distributed within the storage system 100.

25 [18] Fig. 2 shows the data used in an implementation of the journal. When a write request from the host 110 arrives at the storage system 100, a journal is generated in response. The journal comprises a Journal Header 219 and Journal Data 225. The Journal Header 219 contains information about its corresponding Journal Data 225. The Journal Data 225 comprises the data (write data) that is the subject of the write operation. This kind of journal  
30 is also referred to as an "AFTER journal."

[19] The Journal Header 219 comprises an offset number (JH\_OFS) 211. The offset number identifies a particular data volume 101 in the journal group 102. In this particular implementation, the data volumes are ordered as the 0<sup>th</sup> data volume, the 1<sup>st</sup> data volume, the 2<sup>nd</sup> data volume and so on. The offset numbers might be 0, 1, 2, etc.

- [20] A starting address in the data volume (identified by the offset number 211) to which the write data is to be written is stored to a field in the Journal Header 219 to contain an address (JH\_ADR) 212. For example, the address can be represented as a block number (LBA, Logical Block Address).
- 5 [21] A field in the Journal Header 219 stores a data length (JH\_LEN) 213, which represents the data length of the write data. Typically it is represented as a number of blocks.
- [22] A field in the Journal Header 219 stores the write time (JH\_TIME) 214, which represents the time when the write request arrives at the storage system 100. The write time can include the calendar date, hours, minutes, seconds and even milliseconds. This time can  
10 be provided by the disk controller 140 or by the host 110. For example, in a mainframe computing environment, two or more mainframe hosts share a timer and can provide the time when a write command is issued.
- [23] A sequence number (JH\_SEQ) 215 is assigned to each write request. The sequence number is stored in a field in the Journal Header 219. Every sequence number within a given  
15 journal group 102 is unique. The sequence number is assigned to a journal entry when it is created.
- [24] A journal volume identifier (JH\_JVOL) 216 is also stored in the Journal Header 219. The volume identifier identifies the journal volume 106 associated with the Journal Data 225. The identifier is indicative of the journal volume containing the Journal Data. It is noted that  
20 the Journal Data can be stored in a journal volume that is different from the journal volume which contains the Journal Header.
- [25] A journal data address (JH\_JADR) 217 stored in the Journal Header 219 contains the beginning address of the Journal Data 225 in the associated journal volume 106 that contains the Journal Data.
- 25 [26] Fig. 2 shows that the journal volume 106 comprises two data areas: a Journal Header Area 210 and a Journal Data Area 220. The Journal Header Area 210 contains only Journal Headers 219, and Journal Data Area 220 contains only Journal Data 225. The Journal Header is a fixed size data structure. A Journal Header is allocated sequentially from the beginning of the Journal Header Area. This sequential organization corresponds to the chronological  
30 order of the journal entries. As will be discussed, data is provided that points to the first journal entry in the list, which represents the "oldest" journal entry. It is typically necessary to find the Journal Header 219 for a given sequence number (as stored in the sequence number field 215) or for a given write time (as stored in the time field 214).



- [27] Journal Header 219 and Journal Data 225 are contained in chronological order in their respective areas in the journal volume 106. Thus, the order in which the Journal Header and the Journal Data are stored in the journal volume is the same order as the assigned sequence number. As will be discussed below, an aspect of the present invention is that the journal information 219, 225 wrap within their respective areas 210, 220.
- [28] Fig. 3 shows detail about the management table 108 (Fig. 1). In order to manage the Journal Header Area 210 and Journal Data Area 220, pointers for each area are needed. As mentioned above, the management table maintains configuration information about a journal group 102 and the relationship between the journal group and its associated journal volume(s) 106 and snapshot image 105.
- [29] The management table 300 shown in Fig. 3 illustrates an example management table and its contents. The management table stores a journal group ID (GRID) 310 which identifies a particular journal group 102 in a storage system 100. A journal group name (GRNAME) 311 can also be provided to identify the journal group with a human recognizable identifier.
- [30] A journal attribute (GRATTR) 312 is associated with the journal group 102. In accordance with this particular implementation, two attributes are defined: MASTER and RESTORE. The MASTER attribute indicates the journal group is being journaled. The RESTORE attribute indicates that the journal group is being restored from a journal.
- [31] A journal status (GRSTS) 315 is associated with the journal group 102. There are two statuses: ACTIVE and INACTIVE.
- [32] The management table includes a field to hold a sequence counter (SEQ) 313. This counter serves as the source of sequence numbers used in the Journal Header 219. When creating a new journal, the sequence number 313 is read and assigned to the new journal. Then, the sequence number is incremented and written back into the management table.
- [33] The number (NUM\_DVOL) 314 of data volumes 101 contained in a give journal group 102 is stored in the management table.
- [34] A data volume list (DVOL\_LIST) 320 lists the data volumes in a journal group. In a particular implementation, DVOL\_LIST is a pointer to the first entry of a data structure which holds the data volume information. This can be seen in Fig. 3. Each data volume information comprises an offset number (DVOL\_OFFS) 321. For example, if the journal group 102 comprises three data volumes, the offset values could be 0, 1 and 2. A data volume identifier (DVOL\_ID) 322 uniquely identifies a data volume within the entire storage

system 100. A pointer (DVOL\_NEXT) 324 points to the data structure holding information for the next data volume in the journal group; it is a NULL value otherwise.

[35] The management table includes a field to store the number of journal volumes (NUM\_JVOL) 330 that are being used to contain the data (journal header and journal data)

5 associated with a journal group 102.

[36] As described in Fig. 2, the Journal Header Area 210 contains the Journal Headers 219 for each journal; likewise for the Journal Data components 225. As mentioned above, an aspect of the invention is that the data areas 210, 220 wrap. This allows for journaling to continue despite the fact that there is limited space in each data area.

10 [37] The management table includes fields to store pointers to different parts of the data areas 210, 220 to facilitate wrapping. Fields are provided to identify where the next journal entry is to be stored. A field (JI\_HEAD\_VOL) 331 identifies the journal volume 106 that contains the Journal Header Area 210 which will store the next new Journal Header 219. A field (JI\_HEAD\_ADR) 332 identifies an address on the journal volume of the location in the  
15 Journal Header Area where the next Journal Header will be stored. The journal volume that contains the Journal Data Area 220 into which the journal data will be stored is identified by information in a field (JI\_DATA\_VOL) 335. A field (JI\_DATA\_ADR) 336 identifies the specific address in the Journal Data Area where the data will be stored. Thus, the next journal entry to be written is "pointed" to by the information contained in the "JI\_" fields  
20 331, 332, 335, 336.

[38] The management table also includes fields which identify the "oldest" journal entry. The use of this information will be described below. A field (JO\_HEAD\_VOL) 333 identifies the journal volume which stores the Journal Header Area 210 that contains the oldest Journal Header 219. A field (JO\_HEAD\_ADR) 334 identifies the address within the  
25 Journal Header Area of the location of the journal header of the oldest journal. A field (JO\_DATA\_VOL) 337 identifies the journal volume which stores the Journal Data Area 220 that contains the data of the oldest journal. The location of the data in the Journal Data Area is stored in a field (JO\_DATA\_ADR) 338.

[39] The management table includes a list of journal volumes (JVOL\_LIST) 340  
30 associated with a particular journal group 102. In a particular implementation, JVOL\_LIST is a pointer to a data structure of information for journal volumes. As can be seen in Fig. 3, each data structure comprises an offset number (JVOL\_OFS) 341 which identifies a particular journal volume 106 associated with a given journal group 102. For example, if a journal group is associated with two journal volumes 106, then each journal volume might be

identified by a 0 or a 1. A journal volume identifier (JVOL\_ID) 342 uniquely identifies the journal volume within the storage system 100. Finally, a pointer (JVOL\_NEXT) 344 points to the next data structure entry pertaining to the next journal volume associated with the journal group; it is a NULL value otherwise.

- 5 [40] The management table includes a list (SS\_LIST) 350 of snapshot images 105 associated with a given journal group 102. In this particular implementation, SS\_LIST is a pointer to snapshot information data structures, as indicated in Fig. 3. Each snapshot information data structure includes a sequence number (SS\_SEQ) 351 that is assigned when the snapshot is taken. As discussed above, the number comes from the sequence counter 313.
- 10 A time value (SS\_TIME) 352 indicates the time when the snapshot was taken. A status (SS\_STS) 358 is associated with each snapshot; valid values include VALID and INVALID. A pointer (SS\_NEXT) 353 points to the next snapshot information data structure; it is a NULL value otherwise.
- [41] Each snapshot information data structure also includes a list of snapshot volumes 107 (Fig. 1) used to store the snapshot images 105. As can be seen in Fig. 3, a pointer (SVOL\_LIST) 354 to a snapshot volume information data structure is stored in each snapshot information data structure. Each snapshot volume information data structure includes an offset number (SVOL\_OFFS) 355 which identifies a snapshot volume that contains at least a portion of the snapshot image. It is possible that a snapshot image will be segmented or
- 20 otherwise partitioned and stored in more than one snapshot volume. In this particular implementation, the offset identifies the  $i^{\text{th}}$  snapshot volume which contains a portion (segment, partition, etc) of the snapshot image. In one implementation, the  $i^{\text{th}}$  segment of the snapshot image might be stored in the  $i^{\text{th}}$  snapshot volume. Each snapshot volume information data structure further includes a snapshot volume identifier (SVOL\_ID) 356 that
- 25 uniquely identifies the snapshot volume in the storage system 100. A pointer (SVOL\_NEXT) 357 points to the next snapshot volume information data structure for a given snapshot image.
- [42] Fig 4 shows a flowchart highlighting the processing performed by the recovery manager 111 and Storage System 100 to initiate backup processing in accordance with the illustrative embodiment of the invention as shown in the figures. If journal entries are not
- 30 recorded during the taking of a snapshot, the write operations corresponding to those journal entries would be lost and data corruption could occur during a data restoration operation. Thus, in accordance with an aspect of the invention, the journaling process is started prior to taking the first snapshot. Doing this ensures that any write operations which occur during the

taking of a snapshot are journaled. As a note, any journal entries recorded prior to the completion of the snapshot can be ignored.

[43] Further in accordance with the invention, a single sequence of numbers (SEQ) 313 are associated with each of one or more snapshots and journal entries, as they are created. The purpose of associating the same sequence of numbers to both the snapshots and the journal entries will be discussed below.

[44] Continuing with Fig. 4, the recovery manager 111 might define, in a step 410, a journal group (JNLG) 102 if one has not already been defined. As indicated in Fig. 1, this may include identifying one or data volumes (DVOL) 101 for which journaling is performed, and identifying one or journal volumes (JVOL) 106 which are used to store the journal-related information. The recovery manager performs a suitable sequence of interactions with the storage system 100 to accomplish this. In a step 415, the storage system may create a management table 108 (Fig. 1), incorporating the various information shown in the table detail 300 illustrated in Fig. 3. Among other things, the process includes initializing the JVOL\_LIST 340 to list the journal volumes which comprise the journal group 102. Likewise, the list of data volumes DVOL\_LIST 320 is created. The fields which identify the next journal entry (or in this case where the table is first created, the first journal entry) are initialized. Thus, JI\_HEAD\_VOL 331 might identify the first in the list of journal volumes and JI\_HEAD\_ADR 332 might point to the first entry in the Journal Header Area 210 located in the first journal volume. Likewise, JI\_DATA\_VOL 335 might identify the first in the list of journal volumes and JI\_DATA\_ADR 336 might point to the beginning of the Journal Data Area 220 in the first journal volume. Note, that the header and the data areas 210, 220 may reside on different journal volumes, so JI\_DATA\_VOL might identify a journal volume different from the first journal volume.

[45] In a step 420, the recovery manager 111 will initiate the journaling process. Suitable communication(s) are made to the storage system 100 to perform journaling. In a step 425, the storage system will make a journal entry (also referred to as an "AFTER journal") for each write operation that issues from the host 110.

[46] With reference to Fig. 3, making a journal entry includes, among other things, identifying the location for the next journal entry. The fields JI\_HEAD\_VOL 331 and JI\_HEAD\_ADR 332 identify the journal volume 106 and the location in the Journal Header Area 210 of the next Journal Header 219. The sequence counter (SEQ) 313 from the management table is copied to (associated with) the JH\_SEQ 215 field of the next header. The sequence counter is then incremented and stored back to the management table. Of

course, the sequence counter can be incremented first, copied to JH\_SEQ, and then stored back to the management table.

[47] The fields JI\_DATA\_VOL 335 and in the management table identify the journal volume and the beginning of the Journal Data Area 220 for storing the data associated with the write operation. The JI\_DATA\_VOL and JI\_DATA\_ADR fields are copied to JH\_JVOL 216 and to JH\_ADR 212, respectively, of the Journal Header, thus providing the Journal Header with a pointer to its corresponding Journal Data. The data of the write operation is stored.

[48] The JI\_HEAD\_VOL 331 and JI\_HEAD\_ADR 332 fields are updated to point to the next Journal Header 219 for the next journal entry. This involves taking the next contiguous Journal Header entry in the Journal Header Area 210. Likewise, the JI\_DATA\_ADR field (and perhaps JI\_DATA\_VOL field) is updated to reflect the beginning of the Journal Data Area for the next journal entry. This involves advancing to the next available location in the Journal Data Area. These fields therefore can be viewed as pointing to a list of journal entries. Journal entries in the list are linked together by virtue of the sequential organization of the Journal Headers 219 in the Journal Header Area 210.

[49] When the end of the Journal Header Area 210 is reached, the Journal Header 219 for the next journal entry wraps to the beginning of the Journal Header Area. Similarly for the Journal Data 225. To prevent overwriting earlier journal entries, the present invention provides for a procedure to free up entries in the journal volume 106. This aspect of the invention is discussed below.

[50] For the very first journal entry, the JO\_HEAD\_VOL field 333, JO\_HEAD\_ADR field 334, JO\_DATA\_VOL field 337, and the JO\_DATA\_ADR field 338 are set to contain their contents of their corresponding "JI\_" fields. As will be explained the "JO\_" fields point to the oldest journal entry. Thus, as new journal entries are made, the "JO\_" fields do not advance while the "JI\_" fields do advance. Update of the "JO\_" fields is discussed below.

[51] Continuing with the flowchart of Fig. 4, when the journaling process has been initiated, all write operations issuing from the host are journaled. Then in a step 430, the recovery manager 111 will initiate taking a snapshot of the data volumes 101. The storage system 100 receives an indication from the recovery manager to take a snapshot. In a step 435, the storage system performs the process of taking a snapshot of the data volumes. Among other things, this includes accessing SS\_LIST 350 from the management table (Fig. 3). A suitable amount of memory is allocated for fields 351 - 354 to represent the next snapshot. The sequence counter (SEQ) 313 is copied to the field SS\_SEQ 351 and

incremented, in the manner discussed above for JH\_SEQ 215. Thus, over time, a sequence of numbers is produced from SEQ 313, each number in the sequence being assigned either to a journal entry or a snapshot entry.

[52] The snapshot is stored in one (or more) snapshot volumes (SVOL) 107. A suitable amount of memory is allocated for fields 355 - 357. The information relating to the SVOLs for storing the snapshot are then stored into the fields 355 - 357. If additional volumes are required to store the snapshot, then additional memory is allocated for fields 355 - 357.

[53] Fig. 5 illustrates the relationship between journal entries and snapshots. The snapshot 520 represents the first snapshot image of the data volumes 101 belonging to a journal group 102. Note that journal entries (510) having sequence numbers SEQ0 and SEQ1 have been made, and represent journal entries for two write operations. These entries show that journaling has been initiated at a time prior to the snapshot being taken (step 420). Thus, at a time corresponding to the sequence number SEQ2, the recovery manager 111 initiates the taking of a snapshot, and since journaling has been initiated, any write operations occurring during the taking of the snapshot are journaled. Thus, the write operations 500 associated with the sequence numbers SEQ3 and higher show that those operations are being journaled. As an observation, the journal entries identified by sequence numbers SEQ0 and SEQ1 can be discarded or otherwise ignored.

[54] Recovering data typically requires recover the data state of at least a portion of the data volumes 101 at a specific time. Generally, this is accomplished by applying one or more journal entries to a snapshot that was taken earlier in time relative to the journal entries. In the disclosed illustrative embodiment, the sequence number SEQ 313 is incremented each time it is assigned to a journal entry or to a snapshot. Therefore, it is a simple matter to identify which journal entries can be applied to a selected snapshot; i.e., those journal entries whose associated sequence numbers (JH\_SEQ, 215) are greater than the sequence number (SS\_SEQ, 351) associated with the selected snapshot.

[55] For example, the administrator may specify some point in time, presumably a time that is earlier than the time (the "target time") at which the data in the data volume was lost or otherwise corrupted. The time field SS\_TIME 352 for each snapshot is searched until a time earlier than the target time is found. Next, the Journal Headers 219 in the Journal Header Area 210 is searched, beginning from the "oldest" Journal Header. The oldest Journal Header can be identified by the "JO\_" fields 333, 334, 337, and 338 in the management table. The Journal Headers are searched sequentially in the area 210 for the first header whose sequence number JH\_SEQ 215 is greater than the sequence number SS\_SEQ 351 associated

with the selected snapshot. The selected snapshot is incrementally updated by applying each journal entry, one at a time, to the snapshot in sequential order, thus reproducing the sequence of write operations. This continues as long as the time field JH\_TIME 214 of the journal entry is prior to the target time. The update ceases with the first journal entry whose time field 214 is past the target time.

5 [56] In accordance with one aspect of the invention, a single snapshot is taken. All journal entries subsequent to that snapshot can then be applied to reconstruct the data state at a given time. In accordance with another aspect of the present invention, multiple snapshots can be taken. This is shown in Fig. 5A where multiple snapshots 520' are taken. In accordance with  
10 the invention, each snapshot and journal entry is assigned a sequence number in the order in which the object (snapshot or journal entry) is recorded. It can be appreciated that there typically will be many journal entries 510 recorded between each snapshot 520'. Having multiple snapshots allows for quicker recovery time for restoring data. The snapshot closest in time to the target recovery time would be selected. The journal entries made subsequent to  
15 the snapshot could then be applied to restore the desired data state.

[57] Fig. 6 illustrates another aspect of the present invention. In accordance with the invention, a journal entry is made for every write operation issued from the host; this can result in a rather large number of journal entries. As time passes and journal entries accumulate, the one or more journal volumes 106 defined by the recovery manager 111 for a  
20 journal group 102 will eventually fill up. At that time no more journal entries can be made. As a consequence, subsequent write operations would not be journaled and recovery of the data state subsequent to the time the journal volumes become filled would not be possible.

[58] Fig. 6 shows that the storage system 100 will apply journal entries to a suitable snapshot in response to detection of an "overflow" condition. An "overflow" is deemed to  
25 exist when the available space in the journal volume(s) falls below some predetermined threshold. It can be appreciated that many criteria can be used to determine if an overflow condition exists. A straightforward threshold is based on the total storage capacity of the journal volume(s) assigned for a journal group. When the free space becomes some percentage (say, 10%) of the total storage capacity, then an overflow condition exists.  
30 Another threshold might be used for each journal volume. In an aspect of the invention, the free space capacity in the journal volume(s) is periodically monitored. Alternatively, the free space can be monitored in an aperiodic manner. For example, the intervals between monitoring can be randomly spaced. As another example, the monitoring intervals can be

spaced apart depending on the level of free space; i.e., the monitoring interval can vary as a function of the free space level.

[59] Fig. 7 highlights the processing which takes place in the storage system 100 to detect an overflow condition. Thus, in a step, 710, the storage system periodically checks the total  
5 free space of the journal volume(s) 106; e.g., every ten seconds. The free space can easily be calculated since the pointers (e.g., JI\_CTL\_VOL 331, JI\_CTL\_ADDR 332) in the management table 300 maintain the current state of the storage consumed by the journal volumes. If the free space is above the threshold, then the monitoring process simply waits for a period of time to pass and then repeats its check of the journal volume free space.

10 [60] If the free space falls below a predetermined threshold, then in a step 720 some of the journal entries are applied to a snapshot to update the snapshot. In particular, the oldest journal entry(ies) are applied to the snapshot.

[61] Referring to Fig. 3, the Journal Header 219 of the "oldest" journal entry is identified by the JO\_HEAD\_VOL field 333 and the JO\_HEAD\_ADR field 334. These fields identify  
15 the journal volume and the location in the journal volume of the Journal Header Area 210 of the oldest journal entry. Likewise, the Journal Data of the oldest journal entry is identified by the JO\_DATA\_VOL field 337 and the JO\_DATA\_ADR field 338. The journal entry identified by these fields is applied to a snapshot. The snapshot that is selected is the snapshot having an associated sequence number closest to the sequence number of the journal  
20 entry and earlier in time than the journal entry. Thus, in this particular implementation where the sequence number is incremented each time, the snapshot having the sequence number closest to but less than the sequence number of the journal entry is selected (i.e., "earlier in time). When the snapshot is updated by applying the journal entry to it, the applied journal entry is freed. This can simply involve updating the JO\_HEAD\_VOL field 333,  
25 JO\_HEAD\_ADR field 334, JO\_DATA\_VOL field 337, and the JO\_DATA\_ADR field 338 to the next journal entry.

[62] As an observation, it can be appreciated by those of ordinary skill, that the sequence numbers will eventually wrap, and start counting from zero again. It is well within the level of ordinary skill to provide a suitable mechanism for keeping track of this when comparing  
30 sequence numbers.

[63] Continuing with Fig. 7, after applying the journal entry to the snapshot to update the snapshot, a check is made of the increase in the journal volume free space as a result of the applied journal entry being freed up (step 730). The free space can be compared against the threshold criterion used in step 710. Alternatively, a different threshold can be used. For



example, here a higher amount of free space may be required to terminate this process than was used to initiate the process. This avoids invoking the process too frequently, but once invoked the second higher threshold encourages recovering as much free space as is reasonable. It can be appreciated that these thresholds can be determined empirically over time by an administrator.

[64] Thus, in step 730, if the threshold for stopping the process is met (i.e., free space exceeds threshold), then the process stops. Otherwise, step 720 is repeated for the next oldest journal entry. Steps 730 and 720 are repeated until the free space level meets the threshold criterion used in step 730.

10 [65] Fig. 7A highlights sub-steps for an alternative embodiment to step 720 shown in Fig. 7. Step 720 frees up a journal entry by applying it to the latest snapshot that is not later in time than the journal entry. However, where multiple snapshots are available, it may be possible to avoid the time consuming process of applying the journal entry to a snapshot in order to update the snapshot.

15 [66] Fig. 7A shows details for a step 720' that is an alternate to step 720 of Fig. 7. At a step 721, a determination is made whether a snapshot exists that is later in time than the oldest journal entry. This determination can be made by searching for the first snapshot whose associated sequence number is greater than that of the oldest journal entry. Alternatively, this determination can be made by looking for a snapshot that is a  
20 predetermined amount of time later than the oldest journal entry can be selected; for example, the criterion may be that the snapshot must be at least one hour later in time than the oldest journal entry. Still another alternate is to use the sequence numbers associated with the snapshots and the journal entries, rather than time. For example, the criterion might be to select a snapshot whose sequence number is N increments away from the sequence number of  
25 the oldest journal entry.

[67] If such a snapshot can be found in step 721, then the earlier journal entries can be removed without having to apply them to a snapshot. Thus, in a step 722, the "JO\_" fields (JO\_HEAD\_VOL 333, JO\_HEAD\_ADR 334, JO\_DATA\_VOL 337, and JO\_DATA\_ADR 338) are simply moved to a point in the list of journal entries that is later in time than the  
30 selected snapshot. If no such snapshot can be found, then in a step 723 the oldest journal entry is applied to a snapshot that is earlier in time than the oldest journal entry, as discussed for step 720.

[68] Still another alternative for step 721 is simply to select the most recent snapshot. All the journal entries whose sequence numbers are less than that of the most recent snapshot can

- be freed. Again, this simply involves updating the "JO\_" fields so they point to the first journal entry whose sequence number is greater than that of the most recent snapshot. Recall that an aspect of the invention is being able to recover the data state for any desired point in time. This can be accomplished by storing as many journal entries as possible and then
- 5 applying the journal entries to a snapshot to reproduce the write operations. This last embodiment has the potential effect of removing large numbers of journal entries, thus reducing the range of time within which the data state can be recovered. Nevertheless, for a particular configuration it may be desirable to remove large numbers of journal entries for a given operating environment.
- 10 [69] It can be appreciated that the foregoing described steps can be embodied entirely in the controller 140 (e.g., a disk controller). This can take on the form of pure software, custom logic, or some suitable combination of software and hardware, depending on the particular implementation. More generally, the foregoing disclosed embodiments typically can be provided using a combination of hardware and software implementations. One of
- 15 ordinary skill can readily appreciate that the underlying technical solution will be determined based on factors including but not limited or restricted to system cost, system performance, legacy software and legacy hardware, operating environment, and so on. The described current and contemplated embodiments can be readily reduced to specific implementations without undue experimentation by those of ordinary skill in the relevant art.

WHAT IS CLAIMED IS:

- 1                   1.     A method for processing data in an application data store comprising:  
2                   producing at least a first snapshot of an application data store, the application  
3     data store configured to receive data by way of write operations issued from a host device;  
4                   producing a journal entry for each write operation issued from the host device;  
5                   storing each journal entry in a journal data store, thereby accumulating a list of  
6     journal entries;  
7                   monitoring an amount of free space on the journal data store; and  
8                   when the free space falls below a first threshold value, then removing one or  
9     more journal entries from the journal data store, thereby increasing the free space, wherein  
10    enough of the journal entries are removed so that the free space rises above a second  
11    threshold value.
- 1                   2.     The method of claim 1 wherein removing one or more journal entries  
2     includes updating the first snapshot by applying one or more journal entries to the first  
3     snapshot, beginning with an oldest journal entry, wherein journal entries applied to the first  
4     snapshot are removed from the list of journal entries thereby increasing the free space of the  
5     journal data store.
- 1                   3.     The method of claim 1 wherein removing one or more journal entries  
2     includes identifying a selected snapshot that is earlier in time than an oldest journal entry and  
3     is closest in time to the oldest journal entry than other snapshots and updating the selected  
4     snapshot with one or more journal entries beginning with the oldest journal entry.
- 1                   4.     The method of claim 1 wherein removing one or more journal entries  
2     includes looking for a most recent snapshot, removing journal entries in the journal data store  
3     that are earlier in time than the most recent snapshot.
- 1                   5.     The method of claim 1 wherein the first threshold value and the second  
2     threshold value are different.
- 1                   6.     The method of claim 1 wherein each journal entry comprises a fixed-  
2     size header portion and a variable-size data portion, wherein the journal data store comprises  
3     a first storage area within which a plurality of header portions are defined and a second  
4     storage area for storing a plurality of data portions, wherein producing a journal entry

5 includes allocating one of the journal headers and allocating a sufficient amount of space in  
6 the data portion to contain data associated with the write operation.

1                   7.     The method of claim 6 wherein storing each journal entry includes  
2 allocating one of the journal headers from the first portion such that the journal headers are  
3 always allocated in sequential manner and when the last journal header is allocated then  
4 allocating journal headers from the beginning of the first portion in sequential manner.

1                   8.     A data processing method comprising:  
2                   producing at least a first snapshot of at least a portion of an application data  
3 store, the application data store being configured to receive write operations issued from a  
4 host device;  
5                   recording a plurality of journal entries, each journal entry being recorded for a  
6 write operation issued from the host device, each journal entry being stored in a journal data  
7 store thereby consuming an amount of free space of the journal data store; and  
8                   updating the at least first snapshot with at least one journal entry, including  
9 associating space consumed by the at least one journal entry with the free space thereby  
10 increasing the free space of the journal data store.

1                   9.     The method of claim 8 wherein the step of updating includes  
2 periodically monitoring the free space of the journal data store and if the free space falls  
3 below a first threshold then updating the at least first snapshot.

1                   10.    The method of claim 9 wherein the step of updating the at least first  
2 snapshot is repeated for a number of journal entries so that the free space rises above a  
3 second threshold value.

1                   11.    The method of claim 8 further comprising:  
2                   associating sequence numbers to the first snapshot and to the journal entries;  
3                   receiving a target time;  
4                   determining a beginning journal entry based on the sequence numbers  
5 associated with the journal entries and the sequence number associated with the first  
6 snapshot;  
7                   updating the first snapshot by applying the beginning journal entry to it; and  
8                   performing additional updates to the first snapshot by applying journal entries  
9 that are subsequent in time to the beginning journal entry and prior in time to the target time.

1           12.    The method of claim 8 further comprising:  
2           producing additional snapshots thereby accumulating a plurality of snapshots;  
3           associating sequence numbers to the snapshots and to the journal entries;  
4           receiving a target time;  
5           determining a selected snapshot based on the target time;  
6           determining a beginning journal entry based on the sequence numbers  
7   associated with the journal entries and the sequence number associated with the selected  
8   snapshot;  
9           updating the selected snapshot by applying the beginning journal entry to it;  
10   and  
11           performing additional updates to the selected snapshot by applying journal  
12   entries that are subsequent in time to the beginning journal entry and prior in time to the  
13   target time.

1           13.    The method of claim 12 wherein the selected snapshot is a first  
2   snapshot.

1           14.    The method of claim 12 wherein the selected snapshot is closest in  
2   time to the target time and prior in time to the target time.

1           15.    A method for processing data comprising:  
2           producing a plurality of snapshots of one or more portions of an application  
3   data store, the data volume configured to receive data by way of write operations issued from  
4   a host device;  
5           recording a journal entry for each write operation issued from the host device,  
6   each journal entry being stored in a journal data store thereby accumulating a list of journal  
7   entries and consuming free space on the journal data store;  
8           monitoring a remaining free space on the journal data store; and  
9           when the remaining free space falls below a first threshold value, then  
10   removing one or more journal entries by identifying a selected snapshot and removing journal  
11   entries earlier in time than the selected snapshot.

1           16.    The method of claim 15 wherein the selected snapshot is at least a  
2   predetermined length of time later than an oldest journal entry.

1           17.    The method of claim 15 further comprising associating sequence  
2 numbers to the snapshots and to the journal entries, wherein the selected snapshot is  
3 associated with a sequence number that is greater than a sequence number associated with an  
4 oldest journal entry by a predetermined amount.

1           18.    A method for processing data comprising:  
2           producing a plurality of snapshots of an application data store, the data volume  
3 configured to receive data by way of write operations issued from a host device;  
4           producing a journal entry for each write operation issued from the host device;  
5           storing each journal entry in a journal data store, thereby accumulating a list of  
6 journal entries;  
7           monitoring an amount of free space on the journal data store; and  
8           when the free space falls below a first threshold value, then:  
9           identifying a selected snapshot that is earlier in time than an oldest  
10 journal entry and is closest in time to the oldest journal entry than other snapshots;  
11           updating the selected snapshot with one or more journal entries  
12 beginning with the oldest journal entry, thereby reducing the list of journal entries;  
13 and  
14           repeating the updating step for additional journal entries until the free  
15 space rises above a second threshold value.

1           19.    The method of claim 18 further comprising associating sequence  
2 numbers to the snapshots and to the journal entries, wherein the selected snapshot is  
3 determined based on the sequence numbers of the snapshots and the sequence number of the  
4 oldest journal entry.

1           20.    The method of claim 18 wherein journal entries are sequentially  
2 allocated and contiguously arranged in the journal data store in a manner representative of a  
3 sequential list of journal entries, wherein the step of storing includes wrapping to a beginning  
4 of the list of journal entries to reuse earlier journal entries when an end of the list of journal  
5 entries is reached.

1           21.    A method for processing data comprising:  
2           producing a plurality of snapshots of an application data store, the application  
3 data store configured to receive data by way of write operations issued from a host device;

4                    recording a journal entry for each write operation issued from the host device,  
5                    the journal entry being stored in a journal data store, thereby accumulating a list of journal  
6                    entries in the journal data store;  
7                    producing a plurality of sequential numbers;  
8                    for each snapshot and for each journal entry, associating one of the numbers  
9                    therewith;  
10                   monitoring an amount of free space on the journal data store; and  
11                   when the free space falls below a first threshold value, then:  
12                                     identifying a selected snapshot by comparing sequence numbers  
13                   associated with the snapshots against a sequence number associated with an oldest  
14                   journal entry;  
15                                     updating the selected snapshot with one or more journal entries  
16                   beginning with the oldest journal entry; and  
17                                     repeating the updating step with additional journal entries until the free  
18                   space rises above a second threshold value.

1                   22.        The method of claim 21 wherein the step of recording is initiated prior  
2                   to producing a snapshot, so each write operation that issues during a period of time that a first  
3                   snapshot is being produced will be recorded in a journal entry.

1                   23.        The method of claim 21 wherein the selected snapshot is earlier in time  
2                   than the oldest journal entry.

1                   24.        The method of claim 21 wherein the selected snapshot is earlier in time  
2                   than the oldest journal entry and closest in time to the oldest journal entry than other  
3                   snapshots.

1                   25.        The method of claim 21 wherein the first threshold is different from  
2                   the second threshold.

1                   26.     A method for processing data in a data storage system comprising:  
2                   initiating journal processing wherein a journal entry is created for each write  
3                   operation that is sent to the data storage system, thus producing a plurality of journal entries,  
4                   the journal entries being stored in a backup storage system;  
5                   subsequent to the step of initiating journal processing, taking one or more  
6                   snapshots of a portion of the data contained in the data storage system;  
7                   associating a sequence number to each journal entry when it is created and to  
8                   each snapshot when it is taken, wherein sequence numbers associated with the journal entries  
9                   and with the snapshots are part of the same sequence; and  
10                  determining whether to apply some of the journal entries to one of the  
11                  snapshots in order to reduce the amount of space in the backup storage system that is  
12                  consumed by the journal entries, based on the amount of available storage capacity in the  
13                  backup storage system.

1                   27.     The method of claim 26 wherein if a determination is made to apply  
2                   some of the journal entries to one of the snapshots, then:  
3                   identifying an oldest journal entry;  
4                   identifying a selected snapshot based on a sequence number associated with  
5                   the oldest journal entry and on sequence numbers associated with the snapshots; and  
6                   updating the selected snapshot with one or more journal entries, beginning  
7                   with the oldest journal entry.

1                   28.     A storage system for processing data comprising:  
2                   a production data store configured to receive write operations from a host  
3                   device;  
4                   a snapshot data store configured to store one or more snapshots of at least a  
5                   portion of the production data store;  
6                   a journal data store configured to store one or more journal entries; and  
7                   a controller configured to:  
8                   access the production data store and the snapshot data store to store  
9                   one or more snapshots of at least a portion of the production data store on the  
10                  snapshot data store;  
11                  access the journal data store to record a journal entry for each write  
12                  operation from the host device; and



13                   update one of the snapshots with some, but not all, of the journal  
14           entries thereby leaving some journal entries for a data recovery operation.

1                   29.     The storage system of claim 28 wherein the controller is further  
2   configured to:  
3                   associate sequence numbers to the snapshots and to the journal entries;  
4                   receive a target time;  
5                   select a snapshot based on the target time;  
6                   determine a beginning journal entry based on the sequence numbers associated  
7   with the journal entries and the sequence number of the selected snapshot;  
8                   update the selected snapshot by applying the beginning journal entry to it; and  
9                   perform additional updates to the selected snapshot by applying journal entries  
10   that are later in time than the beginning journal entry and earlier in time than the target time.

1                   30.     A storage system for processing data comprising:  
2                   means for receiving communications from a host device, the communications  
3   including one or more write operations;  
4                   first storage means for storing production data;  
5                   second storage means for storing snapshots;  
6                   third storage means for storing journal entries;  
7                   control means for:  
8                   storing onto the second storage means one or more snapshots of at least  
9   a portion of the production data;  
10                  recording onto the third storage means a journal entry for each write  
11   operation thereby producing a plurality of journal entries and reducing an amount of  
12   free space on the third storage means; and  
13                  updating one of the snapshots according to one or more journal entries  
14   thereby increasing the amount of free space on the third storage means.

1                   31.     The storage system of claim 30 wherein the updating can be performed  
2   periodically or at intervals based on the amount of free space.

PATENT

Attorney Docket No.: 16869B-082700US

Client ref. No.: HAL 279

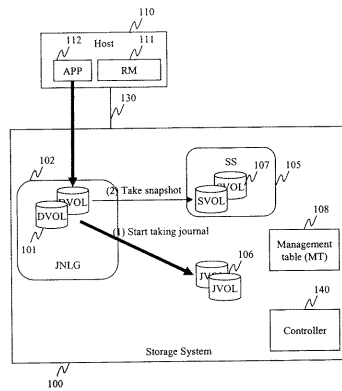
**METHOD AND APPARATUS FOR BACKUP AND RECOVERY USING  
STORAGE BASED JOURNALING****ABSTRACT OF THE DISCLOSURE**

A storage system maintains a journal of journal entries and at least one snapshot of one or more data volumes. By assigning a unique sequence number to journal and snapshot, it is easy to find a journal which can be applied to the snapshot. A technique is described for detecting an overflow condition of running out of journal space and recovering the journal space.

PA 3309031 v1

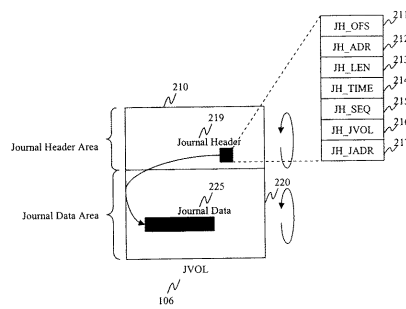
【 図 1 】

Fig. 1 Overview of backup



【 図 2 】

Fig. 2 Control Data for Journal



【 図 3 】

Fig. 3 MT(Management Table) 300

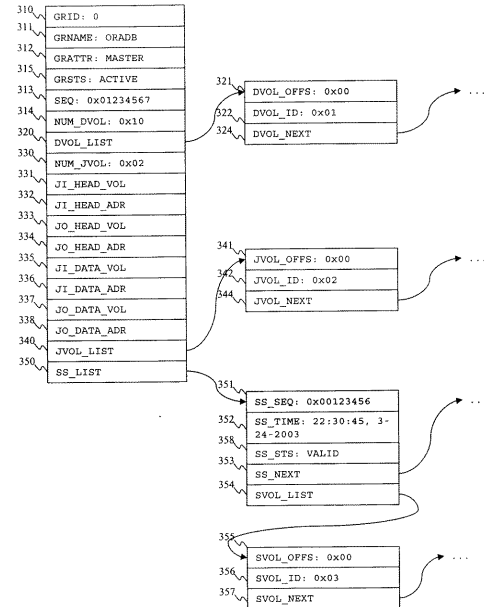
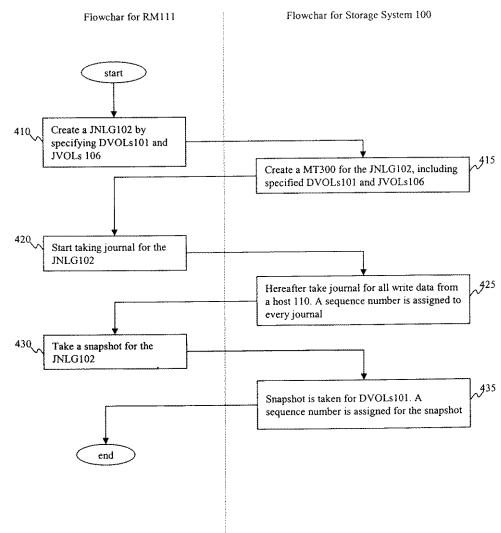


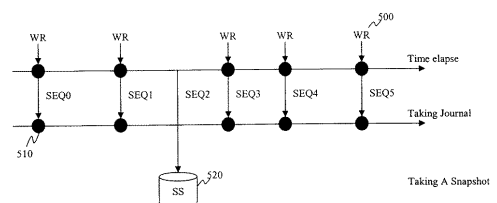
Fig. 4 Starting journal



【 図 4 】

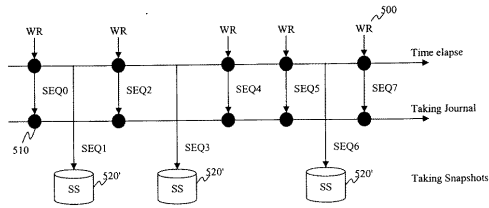
【 図 5 】

Fig. 5 Relationship between snapshot and journal



## 【 図 5 A 】

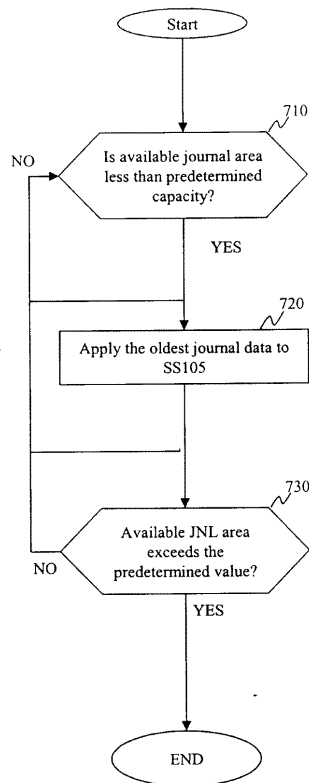
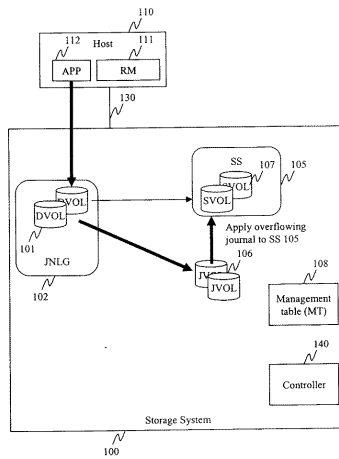
Fig. 5A Relationship between multiple snapshots and journal



## 【 図 7 】

## 【 図 6 】

Fig. 6 Overflowing JNL



## 【 図 7 A 】

Fig. 7 Overflowing JNL

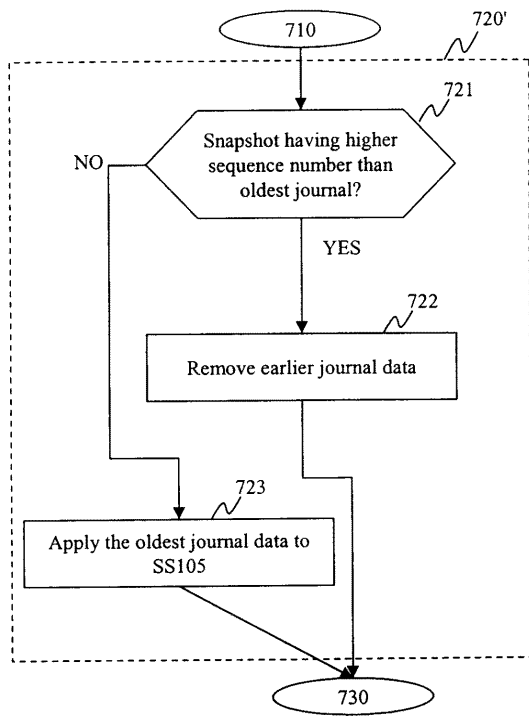


Fig. 7A