

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4363676号
(P4363676)

(45) 発行日 平成21年11月11日(2009.11.11)

(24) 登録日 平成21年8月28日(2009.8.28)

(51) Int.Cl.		F I			
G06F 3/06	(2006.01)		G06F 3/06	304F	
G06F 11/14	(2006.01)		G06F 11/14	310Z	
G06F 12/00	(2006.01)		G06F 12/00	531J	

請求項の数 7 (全 15 頁)

(21) 出願番号	特願平9-300916	(73) 特許権者	000003078
(22) 出願日	平成9年10月31日(1997.10.31)		株式会社東芝
(65) 公開番号	特開平11-134117		東京都港区芝浦一丁目1番1号
(43) 公開日	平成11年5月21日(1999.5.21)	(74) 代理人	100058479
審査請求日	平成16年10月13日(2004.10.13)		弁理士 鈴江 武彦
		(74) 代理人	100084618
			弁理士 村松 貞男
		(74) 代理人	100092196
			弁理士 橋本 良郎
		(74) 代理人	100091351
			弁理士 河野 哲
		(74) 代理人	100088683
			弁理士 中村 誠
		(74) 代理人	100070437
			弁理士 河井 将次

最終頁に続く

(54) 【発明の名称】 コンピュータシステム

(57) 【特許請求の範囲】

【請求項1】

不揮発性記憶装置への書き込みデータと、その論理アドレスとを一時的に保持するバッファ手段と、

前記バッファ手段によって一時的に保持される書き込みデータの数が $n - 1$ 個に達した場合に、当該 $n - 1$ 個の書き込みデータと、当該 $n - 1$ 個の書き込みデータの論理アドレスとタイムスタンプとを含む1個の管理データとから n 個のデータ分のサイズを持つストライプを作成し、この作成したストライプを前記不揮発性記憶装置に書き込むと共に、スナップショット採取の指示を受けた場合に、その時点で前記バッファ手段によって一時的に保持されている $n - 1$ 個以下の書き込みデータと、当該 $n - 1$ 個以下の書き込みデータの論理アドレスとタイムスタンプとを含む1個の管理データとから前記ストライプを作成し、この作成したストライプを前記不揮発性記憶装置に書き込むスナップショット採取手段と、

システム起動時に、前記不揮発性記憶装置に書き込まれたすべてのストライプの管理データを用いて、各論理アドレスについて最新のデータが格納された前記不揮発性記憶装置上の物理アドレスを取得するための変換マップを作成する変換マップ作成手段と、

スナップショット参照の指示を受けた場合に、当該スナップショットのタイムスタンプ値よりも小さいタイムスタンプ値を持つストライプの管理データを用いて、各論理アドレスについて当該指示されたスナップショット採取時において最新のデータが格納された前記不揮発性記憶装置上の物理アドレスを取得するための変換マップを作成し、この作成し

た変換マップに基づき、前記不揮発性記憶装置上のデータを参照するスナップショット参照手段と、

を具備したことを特徴とするコンピュータシステム。

【請求項 2】

今回のシステム起動が前回のシステム再起動失敗による再起動である旨を検出する検出手段と、

前回のシステム起動時に適用したスナップショットの識別情報を管理する識別情報管理手段と、

システム起動が失敗したときに、前記識別情報管理手段により管理された識別情報に基づき、スナップショットの世代を遡ったシステム再起動を自動的に実行する自動再起動手段と、

をさらに具備したことを特徴とする請求項 1 記載のコンピュータシステム。

【請求項 3】

前記検出手段は、前回のシステム起動時からの経過時間によって前記前回のシステム起動が失敗した旨を検出する請求項 2 記載のコンピュータシステム。

【請求項 4】

システム再起動を予め定められた条件で自動的に実行する自動モードおよびオペレータへの問い合わせを介在させながらシステム再起動を実行する手動モードのいずれかを設定する再起動モード設定手段をさらに具備してなることを特徴とする請求項 1、2 または 3 記載のコンピュータシステム。

【請求項 5】

前記自動モードは、最新のスナップショットまたは特定のスナップショットのいずれを参照するかを条件として定める請求項 4 記載のコンピュータシステム。

【請求項 6】

最新のスナップショットを参照するとした場合、システムの再起動後、予め定められた時間内に再度システムがダウンしたときに、各世代ごとに何度再起動を試行するか、および何世代まで遡るかの少なくとも一方をさらに定める請求項 5 記載のコンピュータシステム。

【請求項 7】

システム再起動に関する履歴情報を管理する履歴情報管理手段と、前記履歴情報管理手段により管理された履歴情報を表示または通知する履歴情報呈示手段をさらに具備してなることを特徴とする請求項 4、5 または 6 記載のコンピュータシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

この発明は、たとえばディスク装置などに格納されたファイル群の所定の時点における内容が保持されるスナップショットを適宜に採取するコンピュータシステムに係り、特に任意のスナップショットのディスクイメージでのシステム再起動を容易に実行することのできるコンピュータシステムに関する。

【0002】

【従来の技術】

従来より、ディスク単体の障害に対しては、冗長化ディスク構成を採用することによって対処が可能であるが、この冗長化技術は、ソフトウェアバグ、誤操作およびウイルス感染などによるプログラムまたはデータの喪失・改変、もしくはディスク記憶装置そのものを失なうような大規模な障害には役立たない。これらの問題に対処するためには、スナップショットやバックアップなどの採取が必須となる。

【0003】

スナップショットとは、ある時点のファイルやディスクなどのコピーイメージのことであり、定期的にファイルやディスク全体を同じディスク記憶装置や別のディスク記憶装置にコピーして、スナップショットを作成する。そして、プログラムまたはデータが喪失・改

10

20

30

40

50

変された場合には、直前に作成されたスナップショット（コピーイメージ）にプログラムまたはデータ、もしくはディスク全体を戻すことによりこれらの障害に対処する。

【0004】

また、バックアップとは、このスナップショットを別の大容量記憶媒体（磁気テープなど）に保存して保管したものである。これにより、ディスク記憶装置そのものを失った場合（中のデータも全部喪失）であっても、新しいディスク記憶装置にこのバックアップをロードすることによって、スナップショット採取時のプログラムまたはデータを復元することができる。

【0005】

通常、スナップショットやバックアップの採取時には、それらのファイルを変更する可能性があるアプリケーションをすべて停止する必要がある。そうしないと、コピー作成中にファイルやデータが変更されてしまい、正しいスナップショットやバックアップが採取できないからである。すなわち、不正なスナップショットやバックアップをディスク記憶装置に戻しても、アプリケーションプログラム側でエラーや誤動作を発生させる危険性があるため、先の障害対策として意味をなさない。

【0006】

なお、一般に大容量記憶媒体よりもディスク記憶装置への書き込みの方が高速であるので、ファイルやディスクイメージを直接バックアップするよりも、スナップショットを作成し、それを大容量記憶媒体にバックアップする方がアプリケーションの停止期間を短く押さえられる。また、プログラムまたはデータの喪失・改変からの復元もスナップショットの方が短期に行なえる。バックアップが磁気テープなどのシーケンシャルアクセスの媒体の場合、特定のファイルを見つけて読み出すには多くの時間を要してしまい非常に効率が悪い。

【0007】

また、このスナップショットやバックアップを採取するためにアプリケーションを長時間停止しなければならないといった問題を解決する、スナップショット機能を持ったファイルシステムも開示されている（文献1：「The Episode File System」, Proceedings of the Winter 1992 USENIX Conference, pp. 43 - 60, San Francisco, CA、文献2：「File System Design for an NFS File Server Appliance」, Proceedings of the Winter 1994 USENIX Conference, pp. 235 - 244, San Francisco, CA）。

【0008】

ところで、これらのファイルシステム（たとえば文献2）では、ファイルがルート・ノードを起点とするブロック木構造であることを前提にしている。そして、スナップショット作成時には、このルート・ノードのコピーを作成する。このコピーしたルート・ノード（スナップショット・ルートノード）は、スナップショットとしてのブロック木構造を表現することになる。また、コピー元のルート・ノードは、読み出し/書き込みが行なわれるアクティブなファイルシステムのブロック木構造を表現する。スナップショット・ルートノードは、生成されたときにはルート・ノードとまったく同じブロックを指しているので、新しいスナップショットのためにコピーしたルート・ノード以外はまったくディスクスペースを消費しない。

【0009】

ここで、ユーザがいずれかのデータブロックを変更したとすると、新しいデータブロックをディスク上に書き込み、ルート・ノードで表現するアクティブなファイルシステムが、この新しいブロックを指すように変更される。ディスク上の元のデータブロックは変更されずに残っており、スナップショット・ルートノードは、元のデータブロックをまだ指している。したがって、スナップショット・ルートノードを指定することにより、スナップショット作成時のデータブロックをそのまま参照することができる。

10

20

30

40

50

【 0 0 1 0 】

このように、スナップショット機能を持ったファイルシステムを使えば、ルート・ノードをコピーするだけで簡単にスナップショットを作成でき、アプリケーションプログラムを停止することもない。また、スナップショットをアプリケーションプログラムの実行と並行して大容量記憶媒体に保存できるので、バックアップの採取もアプリケーションを停止することなく行なえる。したがって、スナップショットやバックアップのデータ退避のためにアプリケーションプログラムを停止する必要がなくなる。

【 0 0 1 1 】

【 発明が解決しようとする課題 】

しかしながら、前述した手法においては、専用のファイルシステムを新しく開発する必要があり、既存のコンピュータシステムにそのまま適用できるものではなかった。また、ファイルシステム自身もブロック木構造で構成されていることを前提にしており、たとえばマイクロソフト社のNTFSなどのエクステンベースのファイルシステムには適用できる技術ではなかった。また、一般にファイルシステムは大きな木構造をしており、実際に前述のデータブロックを更新するには、その経路に位置する中間ノードすべてをコピーする必要があり、スナップショット作成後の更新性能が大きく低下するといった問題もあった。さらに、ファイルシステムという複雑なソフトウェアモジュールにスナップショット機能を付加したため、スナップショットは読み出しだけで非常に硬直なものになってしまっていた。

【 0 0 1 2 】

この発明はこのような実情に鑑みてなされたものであり、既存のコンピュータシステムやファイルシステムにそのまま適用することが可能であって、スナップショットを効率的に採取するとともに、任意のスナップショットのディスクイメージでのシステム再起動を容易に実行することのできるコンピュータシステムを提供することを目的とする。

【 0 0 1 3 】

【 課題を解決するための手段 】

前述した目的を達成するために、この発明のコンピュータシステムは、不揮発性記憶装置への書き込みデータと、その論理アドレスとを一時的に保持するバッファ手段と、前記バッファ手段によって一時的に保持される書き込みデータの数が $n - 1$ 個に達した場合に、当該 $n - 1$ 個の書き込みデータと、当該 $n - 1$ 個の書き込みデータの論理アドレスとタイムスタンプとを含む 1 個の管理データとから n 個のデータ分のサイズを持つストライブを作成し、この作成したストライブを前記不揮発性記憶装置に書き込むと共に、スナップショット採取の指示を受けた場合に、その時点で前記バッファ手段によって一時的に保持されている $n - 1$ 個以下の書き込みデータと、当該 $n - 1$ 個以下の書き込みデータの論理アドレスとタイムスタンプとを含む 1 個の管理データとから前記ストライブを作成し、この作成したストライブを前記不揮発性記憶装置に書き込むスナップショット採取手段と、システム起動時に、前記不揮発性記憶装置に書き込まれたすべてのストライブの管理データを用いて、各論理アドレスについて最新のデータが格納された前記不揮発性記憶装置上の物理アドレスを取得するための変換マップを作成する変換マップ作成手段と、スナップショット参照の指示を受けた場合に、当該スナップショットのタイムスタンプ値よりも小さいタイムスタンプ値を持つストライブの管理データを用いて、各論理アドレスについて当該指示されたスナップショット採取時において最新のデータが格納された前記不揮発性記憶装置上の物理アドレスを取得するための変換マップを作成し、この作成した変換マップに基づき、前記不揮発性記憶装置上のデータを参照するスナップショット参照手段と、を具備したことを特徴とする。

【 0 0 1 4 】

この発明によれば、既存のオペレーティングシステムやファイルシステム、あるいはアプリケーションプログラムになんらの改造を伴うことなく、スナップショットを効率的に採取することが可能となり、また、システムの稼働状況などに応じて選択される任意のスナップショットのディスクイメージでシステムを再起動させることが容易に行なえるため

10

20

30

40

50

、システムの可用性向上とシステム管理の省力化とを図ることが可能となる。

【0016】

【発明の実施の形態】

以下、図面を参照してこの発明の実施の形態を説明する。

まず、この実施形態におけるディスクスナップショットの管理手法について説明する。図1は、この実施形態のコンピュータシステムでディスクスナップショットの管理を司るディスクスナップショット部の概念構成を示す図である。図1に示すように、この実施形態のディスクスナップショット部2は、ディスクスナップショット制御部3および揮発性メモリ5から構成される。そして、この揮発性メモリ5内には、書き込みの時間的順序を維持するためのタイムスタンプ6、ディスク装置4に書き込むデータをログ構造化して保持する書き込みバッファ7、および書き込みバッファ7内の空き領域ならびに保持されている書き込みデータの論理アドレスの情報を保持するバッファ管理テーブル8が格納される。さらに、ディスク装置4内には、通常のデータのほかにスナップショットを管理するためのスナップショット情報(SS情報)10が格納されている。

10

【0017】

ディスクスナップショット制御部3は、これらタイムスタンプ6、書き込みバッファ7、バッファ管理テーブル8およびスナップショット情報10を管理し、ディスク装置4への読み出し/書き込みおよびスナップショットの参照を制御する。

【0018】

ここで、データの書き込み動作について説明する。ディスク装置4は、それぞれデータブロックサイズの整数倍(k)であるストライプと呼ばれる単位で書き込みを行なう。ディスクスナップショット制御部3は、ファイルシステム1から書き込むデータとその論理アドレスとを受け取り、揮発性メモリ5上の書き込みバッファ7の空き領域にブロック単位に分割して詰める。また、受け取った論理アドレスをブロックごとのアドレスに変換して、バッファ管理テーブル8の対応するエントリに格納する。図2には、この書き込みバッファ7およびバッファ管理テーブル8の内容が例示されている。

20

【0019】

そして、ディスクスナップショット制御部3は、ファイルシステム1からの書き込みデータが、1ストライプ分に1ブロック少ない数だけ書き込みバッファ7に溜まった段階で、図3に示すように、最後の1ブロックにストライプ11上のブロックの論理アドレスとタイムスタンプとを格納して論理アドレスタグブロック12とし、そのストライプ11をディスク装置4に書き込む。また、タイムスタンプ6の値は、この書き込みが完了した段階でインクリメントされる。

30

【0020】

このように、データの更新に際して、ディスク装置4上の該当する物理領域を書き換えるのではなく、別の空き領域にタイムスタンプとともに書き込んでいくので、ディスク装置4上には同じ論理アドレスをもつデータブロックが複数存在する。通常の実データ参照要求に対しては、最新のデータブロックの物理位置を返す必要があるため、システム起動時に、全ストライプの論理アドレスタグを調べて揮発性メモリ5上に論理アドレスから物理アドレスへの変換マップ9を作成する。変換マップ9は、図4に示すように、各論理アドレスのブロックが格納されているストライプ番号ST#と、そのストライプ内のブロック番号BLK#と、そのタイムスタンプTS#とをテーブル方式で保持している。これにより、各論理アドレスブロックの物理位置が求められる。

40

【0021】

次に、スナップショット採取処理について説明する。スナップショット採取の指示を受けると、ディスクスナップショット制御部3は、その時点で書き込みバッファ7に溜まっているデータをディスク装置4に書き込む。このとき、書き込みバッファ7の空き領域に該当するバッファ管理テーブル8のエントリには、論理アドレスとしてヌル(空値: Null Value)を設定する。その後、通常の実データと同様に、最後の1ブロックとして論理アドレスタグブロック12を作成し、1ストライプ分を書き込む。さらに、このと

50

きの論理アドレスタグブロック 12 に付加したタイムスタンプ 6 の値をスナップショット採取時の情報として、ディスク装置 4 内のディスクスナップショット情報 10 に登録し、完了したらタイムスタンプ 6 の値をインクリメントする。

【 0 0 2 2 】

次に、スナップショットの参照処理について説明する。あるスナップショットのイメージを参照するためには、参照するスナップショットのタイムスタンプ値より小さいタイムスタンプ値をもつストライプに含まれるすべての論理アドレスブロックについて、システム起動時と同様の変換マップを作成し、そのスナップショット参照用変換マップを使ってアクセスすればよい。

【 0 0 2 3 】

このように、この実施形態のディスクスナップショットの管理手法によれば、既存のオペレーティングシステムやファイルシステム、あるいはアプリケーションプログラムになんらの改造を伴うことなく、スナップショットを効率的に採取することが可能となる。以下では、このように採取したスナップショットの中の任意のスナップショットのディスクイメージでシステムを起動させるための手法について説明する。

【 0 0 2 4 】

図 5 は、最新のスナップショットをシステム起動時のディスクイメージとするシステムの構成図である。

システム起動時に、ディスクスナップショット制御部 3 は、スナップショット情報 10 を読み込んで、最新のスナップショットのタイムスタンプ (TS) 6 を得る。また、ディスク装置 4 中の有効なストライプ 11 の位置が記録されているストライプ管理情報 (ST) 15 を揮発性メモリ 5 上に読み込む。そして、ディスクスナップショット制御部 3 は、このストライプ管理情報 15 をもとにディスク装置 4 から有効なストライプを読み込みバッファ 14 に読み込み、ストライプ 11 の論理アドレスタグブロック 12 のタイムスタンプ 13 と最新のスナップショットのタイムスタンプ 6 とを比較して、タイムスタンプ値 13 がタイムスタンプ 6 より小さい場合のみ、ストライプ 11 に含まれるデータブロックの変換マップ 9 への登録を行なう。このシステムの場合、最新のスナップショットよりも後に書き込まれたストライプは、以後利用されることはないので、タイムスタンプ 13 がタイムスタンプ 6 より大きい場合は、ストライプ管理テーブル 15 の該当するエントリを削除し、そのストライプを無効化する。変換マップ 9 への登録処理は、論理アドレスタグブロック 12 内の各ブロックの論理アドレスを変換マップ 9 に登録されている論理アドレスと比較し、登録されていないければ登録し、すでに登録されている場合には、タイムスタンプ値 13 が登録されているタイムスタンプ値よりも大きい場合のみ、エントリを更新して登録する。ディスク装置 4 内の有効な全ストライプに対してこの処理を行なって作成した変換マップ 9 を使用することにより、通常と同じ手順で最新のスナップショット時点のディスクイメージにアクセスできる。そして、この実施形態のディスクスナップショット部 2 では、これを擬似ディスク装置 4 で表す。すなわち、続いて行なわれるファイルシステム 1 の作成処理において、ディスクスナップショット制御部 3 がディスク装置 4 の代わりに擬似ディスク装置 4 を見せることにより、擬似ディスク装置 4 をディスクイメージとしてアクセスするファイルシステム 1 が作成され、アプリケーションプログラムは擬似ディスク装置であることを意識することなく実行することができる。

【 0 0 2 5 】

また、図 5 に示したシステムの構成に、システムダウンしたときにディスクイメージを戻したい日時が予め指定できる機能と、スナップショット情報 10 からキーとなるタイムスタンプ 6 を得る処理においてその指定された日時以前で最も新しいスナップショットのタイムスタンプ 13 を選択する機能とを付け加えることによって、指定した日時のディスクイメージまで戻してシステム再起動することが可能になる。

【 0 0 2 6 】

さらに、図 5 に示したシステムの構成に、スナップショット情報 10 を揮発性メモリ 5 内に読み込んだ後、有効なスナップショット 11 の日時を画面表示する機能と、その中から

10

20

30

40

50

オペレータが選択したスナップショット11のタイムスタンプ13を変換マップ6の作成時のキーとなるタイムスタンプ6として設定する機能とを付け加えることによって、オペレータの指定したスナップショット時点のディスクイメージでシステム起動することが可能になる。なお、このとき必ず選択画面を表示するというオプション情報を設けてチェックするようにすれば、システムダウンした場合だけでなく、通常の立ち上げ時にも任意のスナップショットを選択してブートさせることが可能となる。

【0027】

図6は、システム起動が失敗した場合にスナップショットの世代を遡った自動再起動が可能なシステムのブート時の処理の流れを示すフローチャートである。図5に示したシステム構成に付加する必要がある要素は、前回の起動が失敗したことを検出する機能と、前回の起動で使用したスナップショットの識別情報を保存する機能とである。

10

【0028】

この場合、ディスクスナップショット制御部3は、まず、前回システムが正常終了したかどうかを判定する(ステップA1)。この処理は、予め決まった場所にある情報を参照することによって行なう。次に、ディスクスナップショット制御部3は、この再起動が初めてのものであるか、または再起動失敗に伴う再起動かを判定する(ステップA2)。ここでは、保存されている前回のリスタート開始時刻と現在の時刻とを比較して、前回システム再起動した時刻から予め指定された時間または所定の時間を経過しているかどうかを調べると同時に、リスタート開始時刻を現在の時刻で更新する。

【0029】

再起動失敗の後の再起動の場合であった場合(ステップA2のN)、ディスクスナップショット制御部3は、参照するスナップショットを前回使用したスナップショットより一世代古いものにし、使用するスナップショットの識別情報を更新する(ステップA3)。具体的には、前回の再起動時に参照したスナップショットのタイムスタンプが保存されているので、スナップショット情報10を検索してその一世代前のスナップショット11を選択するとともに、このスナップショット11のタイムスタンプ13の値でスナップショット情報10のタイムスタンプ6を更新する。次に、ディスクスナップショット制御部3は、選択されたスナップショット11のタイムスタンプ13をキーとして、このスナップショット11の参照用変換マップ9を作成する(ステップA3)。この後は、ファイルシステム作成から通常の立ち上げ処理に入る(ステップA5,ステップA6)。

20

30

【0030】

このように、スナップショットの世代を遡った再起動を自動的に実行することが可能となる。

次に、システム再起動を予め定められた条件で自動的に実行する自動モードおよびオペレータへの問い合わせを介在させながらシステムの再起動を実行する手動モードのいずれかによるリスタートを実現する手法について説明する。図7は、このリスタートを実現するコンピュータシステムの構成図である。

【0031】

図7に示すように、このコンピュータシステム100は、端末装置101を備えており、この端末装置101は、オペレータがコンピュータシステム100に対して指示を行なうための入力装置と処理結果やその他の情報を表示するための出力装置とからなっている。

40

【0032】

このコンピュータシステム100は、前述のシステムが実現する機能に加え、さらにきめ細かな再起動指定を可能にするために、リスタート情報設定モジュール(RST)21を用意し、リスタート自動/手動、自動のときのディスクイメージの日時指定、再起動失敗を判断する際の前回再起動からの経過時間T、スナップショット一世代ごとの再起動リスタート回数N、およびスナップショットを遡るときの世代数の制限値Mなどのリスタート設定を、予めオペレータに設定させ、リスタート設定情報(RST)16としてディスク装置4に保存する。

【0033】

50

一方、これらの指定にしたがってリスタートを行なうためには、現在何回目のリスタートであるか、再起動失敗後なのか、などのリスタート履歴を保存し、参照する必要があるので、さらにリスタート履歴情報（RLOG）17をディスク装置4に保存する。そして、図6のステップA1で示した、システムが前回正常終了したかどうかを判定するリスタート判定モジュール31を設けて、リスタート回数K（正常終了からの）を設定するために用いる。また、カレントの世代のスナップショットでのリトライ回数n、現在何世代さかのぼったかを示すカレント世代数m、リスタート開始時刻、および起動に使用したスナップショットのタイムスタンプなどもリスタート履歴情報17として保存する。

【0034】

ブート時に使用するスナップショットを選択するシナップショット選択モジュール34は、リスタート判定モジュール31、リスタート設定情報参照モジュール32、およびリスタート履歴情報参照・更新モジュール33を利用して指定に応じたスナップショットを選択し、ディスクスナップショット制御部3に指示してそのスナップショットを参照するための変換マップ9を作成させ、擬似ディスク装置4をディスクイメージとしてファイルシステム1に見せる。

10

【0035】

オペレータが起動に使うスナップショットを選択する指定では、ファイルシステム1が立ち上がる前のブート処理の段階で、その選択画面の表示を行なう必要があるので、スナップショット検索・一覧モジュール35を用意してスナップショット情報10を取得し、オペレータ問い合わせモジュール36を介してスナップショットの一覧を端末装置101に出力し、オペレータの選択したスナップショットのタイムスタンプをスナップショット選択モジュール34に通知する。

20

【0036】

なお、これらのリスタート管理のために用意した制御モジュールをリスタート管理モジュール30と呼び、一方、リスタート管理のためのアプリケーション群をリスタート管理ユーティリティ20と呼ぶことにする。

【0037】

図8は、ブート時にリスタート管理モジュール30が行なう処理の流れを示すフローチャートである。

リスタート管理モジュール30は、まず、システムが前回正常に終了したかどうかを判定する（ステップB1）。この処理は、ある決まった領域にハードウェアまたは基本システムモジュールが設定する情報を参照することによって行なわれる。システムダウンだった場合（ステップB1のN）、リスタート管理モジュール30は、リスタート履歴情報17の中のリスタート回数Kをインクリメントする（ステップB2）。このリスタート回数Kは、システムのシャットダウン（正常終了処理）の際にリスタート履歴参照・更新モジュール33によってクリアされる。

30

【0038】

次に、リスタート管理モジュール30は、システムダウンの場合の再起動の指定が自動か手動かを判断する（ステップB3）。この処理は、リスタート設定情報18の中の該当する情報をリスタート設定情報参照モジュール32が参照することによって行なわれる。

40

【0039】

手動の場合は（ステップB3のN）、前述したようにオペレータにスナップショットを選択させた後（ステップB10）、後述するステップB8へと進む。一方、自動の場合は（ステップB3のY）、まず、起動ディスクイメージの日時指定があるかどうかを判定する（ステップB4）。この処理もステップB3と同様に、リスタート設定情報18を参照することによって行なわれ、指定があった場合には（ステップB4のY）、スナップショット選択モジュール34でその日時以前の中で最新のスナップショットを選択した後（ステップB12）、後述するステップB8へと進む。また、日時指定がない場合（ステップB4のN）、さらに前回の起動から一定時間が経過しているかどうかを調べる（ステップB5）。これは、正常実行後初めての再起動なのか、再起動が失敗した直後の再起動なのか

50

を判定するためであり、図 6 のステップ A 2 と同様に、リスタート履歴情報参照・更新モジュール 3 3 でリスタート履歴情報 1 7 の中のリスタート開始時刻の参照・更新を行なう。この判断に使用する前回リスタートからの経過時間の基準は、既定値またはリスタート設定情報 1 8 から取得する。

【 0 0 4 0 】

ここで、前回のリスタートから指定時間が経過しており、一回目の再起動と判断した場合は (ステップ B 5 の Y)、カレントリトライ回数 n を 1 に、カレント世代数 m を 0 にセットした後 (ステップ B 6)、最新のスナップショットを選択する (ステップ B 7)。

【 0 0 4 1 】

また、前回のリスタートから指定時間内に再びシステムダウンした場合は (ステップ B 5 の N)、スナップショット選択処理 (ステップ B 9) に進み、ここで条件にあったスナップショットを選択する。

【 0 0 4 2 】

そして、これまでの処理で選択されたスナップショットから起動用のディスクイメージ (擬似ディスク装置 4') を作成する (ステップ B 8)。具体的には、その選択されたスナップショットのタイムスタンプをキーとして、変換マップ 9 へ論理アドレスの登録を行ない、擬似ディスク装置 4' へのアクセス用の変換マップ 9 を作成する処理となる。

【 0 0 4 3 】

図 9 は、スナップショット選択処理 (図 8 のステップ B 9) の流れを示すフローチャートである。

このスナップショット選択処理では、まず、カレントリトライ回数 n とリスタート設定情報 1 6 にある最大リトライ回数 N (一世代のスナップショットに関して連続して何回リスタートを試みるか) とを比較し (ステップ C 1)、 $n < N$ ならば (ステップ C 1 の Y)、前回と同じスナップショットを選択するのでカレントリトライ回数をインクリメントして (ステップ C 4) 終了する。そうでなければ (ステップ C 1 の N)、一世代前のスナップショットを選択するために、カレント世代数 m とリスタート設定情報 1 6 の中の最大世代数 M (何世代さかのぼって自動リスタートを行なうか) とを比較し (ステップ C 2)、指定されている世代数を越えていたら (ステップ C 2 の N)、立ち上げを中断する (ステップ C 5)。この場合は、その旨のメッセージを表示または通知し、リスタート履歴情報 1 7 にログを残す。一方、 $m < M$ ならば (ステップ C 2 の Y)、カレント世代数 m をインクリメントし、カレントリトライ回数 n を 1 に設定して、前回使用したスナップショットよりも一世代古いスナップショットを選択する (ステップ C 3)。この処理は、スナップショット情報 1 0 に有効なスナップショット 1 1 とタイムスタンプ 6 とが格納されているので、リスタート履歴情報 1 7 に保存している、前回使用したスナップショット 1 1 のタイムスタンプ 1 3 よりも小さい中で最大のものを選び、この選んだスナップショット 1 1 のタイムスタンプ 1 3 でタイムスタンプ 6 を更新する。

【 0 0 4 4 】

以上の構成要素により、システム再起動を予め定められた条件で自動的に実行する自動モードおよびオペレータへの問い合わせを介在させながらシステム再起動を実行する手動モードのいずれかによるリスタートを実現できるが、最新のスナップショット以外のディスクイメージで再起動した場合は、起動に使用したスナップショットよりも後に採取されたスナップショットと再起動後のディスクイメージとの関係を考えなければならない。たとえば、図 1 0 に示すように、スナップショット 3 (SS 3) の後にシステムがダウンし、オペレータがスナップショット 1 (SS 1) を使用して再起動したとすると、再起動後にスナップショット 4 (SS 4) を採取するのであれば、その前にスナップショット 2 (SS 2) およびスナップショット 3 (SS 3) は無効化しなければならない。スナップショット 4 (SS 4) を採取した後に、スナップショット 2 (SS 2) およびスナップショット 3 (SS 3) が残っていると、そこでシステムダウンした場合にディスクイメージの整合性がとれなくなってしまうからである。そこで、リスタート履歴情報 1 7 に、図 8 のステップ B 9, B 1 1, B 1 2 で最新のスナップショット以外を選択した場合にオンとする

10

20

30

40

50

フラグを設け、ディスクスナップショット制御部 3 によるスナップショット採取処理でそのフラグをチェックし、オンの場合には、該当するスナップショット（図 10 の場合には、スナップショット 2（SS2）およびスナップショット 3（SS3））を削除する。スナップショットの削除は、スナップショット情報 10 に登録されているスナップショットの中でリスタートに使用したスナップショットのタイムスタンプよりも大きいタイムスタンプをもっているスナップショットのエントリを削除するとともに、該当するストライプも無効化する。この処理は、ストライプの物理的な位置とタイムスタンプとが登録されているストライプ管理情報 15 を検索し、リスタートに使用しているスナップショットのタイムスタンプより大きく、リスタート開始時刻よりも小さいタイムスタンプをもつストライプを無効化すればよい。

10

【0045】

なお、この例では、スナップショットの無効化は再起動後の最初のスナップショット採取直前に行なわれるので、スナップショット採取を行なわないように設定する手段があれば、スナップショット 2（SS2）およびスナップショット 3（SS3）は再起動後も残すことができる。これにより、たとえば実験などのために、およびスナップショット 1（SS1）やスナップショット 0（SS0）のイメージで起動して実行し、次の起動では、スナップショット 2（SS2）またはスナップショット 3（SS3）を使って起動することも可能になる。図 11 には、このスナップショット採取のタイミングを変更できるシステムの概念図が示されている。

【0046】

ここでは、スナップショット採取のタイミングの制御は、タイミング制御モジュール 18 が行なう。通常、決められた一定時間間隔でディスクスナップショット部 2 に対して実行を指示するが、次のような機能を用意することにより、タイミングを変更可能にする。すなわち、システム管理者が、スナップショット採取タイミングの指定を行なうためのスナップショット採取モジュール 23 を用意することにより、タイミング制御モジュール 18 にスナップショット採取実行や採取間隔の指示を行なう。このタイミング制御モジュール 18 は、指定された時間ごとにディスクスナップショット部 2 に対してスナップショット採取実行を指示するほかに、スナップショット採取モジュール 23 からの指示で即時的に実行を指示する。また、採取間隔が 0 に設定されたときは採取を行なわない。

20

【0047】

さらに、採取を行なう指定で実行していたシステムがシステムダウンした場合に、採取を行なわない指定で再起動したいとこのために、図 7 で示したオペレータ問い合わせモジュール 36 に相当するモジュールでスナップショット採取に関する指定を可能にし、システム実行が開始される前にその情報をタイミング制御モジュール 18 に通知する。

30

【0048】

また、システムの再起動を予め定められた条件で自動的に実行する自動モードおよびオペレータへの問い合わせを介在させてシステムの再起動を実行する手動モードのいずれかによるリスタートを実現するような、きめ細かなリスタート制御を可能とする、リスタート履歴情報 17 を有するシステムでは、リスタート管理ユーティリティ 20 の一部として、リスタート履歴情報 17 を表示したり通知したりするモジュールを用意することにより、リスタートの履歴をシステムの診断に利用することが可能となる。

40

【0049】**【発明の効果】**

以上詳述したように、この発明によれば、たとえばディスク装置などの不揮発性記憶装置とファイルシステムとの間に、ディスクレベルのスナップショットを採取するスナップショット管理手段を介在させる構造を採用するため、既存のオペレーティングシステムやファイルシステム、あるいはアプリケーションプログラムになんらの改造を伴うことなく、スナップショットを効率的に採取することが可能となり、また、システムの稼働状況などに応じて選択される任意のスナップショットのディスクイメージでシステムを再起動させることが容易に行なえるため、システムの可用性向上とシステム管理の省力化とを図る

50

ことが可能となる。

【図面の簡単な説明】

【図 1】この発明の実施形態に係るコンピュータシステムでディスクスナップショットの管理を司るディスクスナップショット部の概念構成を示す図。

【図 2】同実施形態の書き込みバッファおよびバッファ管理テーブルの内容を例示する図。

【図 3】同第実施形態の最後の 1 ブロックにストライプ上のブロックの論理アドレスとタイムスタンプとを格納して論理アドレスタグブロックとし、そのストライプをディスク装置に書き込む様子を示す図。

【図 4】同実施形態の変換マップの記憶形式を示す図。

【図 5】同実施形態の最新のスナップショットをシステム起動時のディスクイメージとするシステムの構成図。

【図 6】同実施形態のシステム再起動が失敗した場合にスナップショットの世代を遡った自動再起動が可能なシステムのブート時の処理の流れを示すフローチャート。

【図 7】同実施形態のシステム再起動を予め定められた条件で自動的に実行する自動モードおよびオペレータへの問い合わせを介在させながらシステム再起動を実行する手動モードのいずれかによるリスタートを実現するコンピュータシステムの構成図。

【図 8】同実施形態のブート時にリスタート管理モジュールが行なう処理の流れを示すフローチャート。

【図 9】同実施形態のスナップショット選択処理の流れを示すフローチャート。

【図 10】同実施形態の最新のスナップショット以外のディスクイメージで再起動した場合の考慮点を説明するための図。

【図 11】同実施形態のスナップショット採取のタイミングを変更できるシステムの概念図。

【符号の説明】

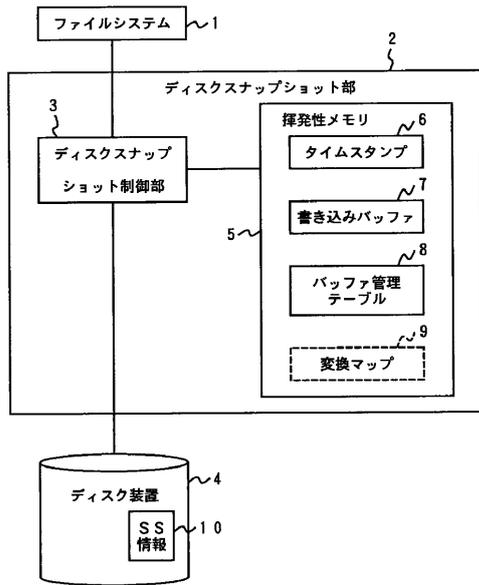
1 ... ファイルシステム、2 ... ディスクスナップショット部、3 ... ディスクスナップショット制御部、4 ... ディスク装置、5 ... 揮発性メモリ、6 ... タイムスタンプ、7 ... 書き込みバッファ、8 ... バッファ管理テーブル、9 ... 変換マップ、10 ... スナップショット情報、11 ... ストライプ、12 ... 論理アドレスタグブロック、13 ... タイムスタンプ、14 ... 読み込みバッファ、15 ... ストライプ管理情報、20 ... リスタート管理ユーティリティ、21 ... リスタート情報設定モジュール (RST)、30 ... リスタート管理モジュール、31 ... リスタート設定情報参照モジュール、33 ... リスタート履歴情報参照・更新モジュール、34 ... スナップショット選択モジュール、35 ... スナップショット検索・一覧モジュール、36 ... オペレータ問い合わせモジュール、100 ... コンピュータシステム。

10

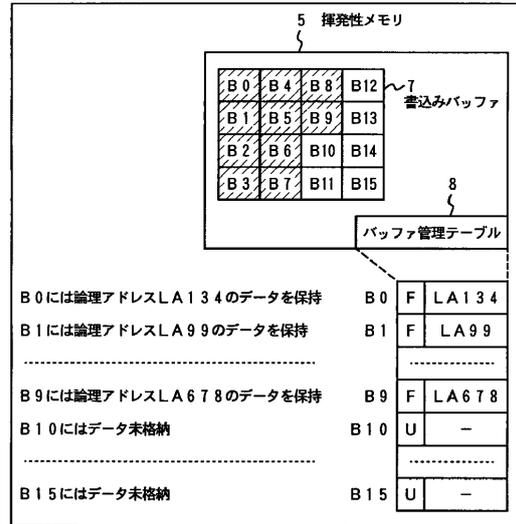
20

30

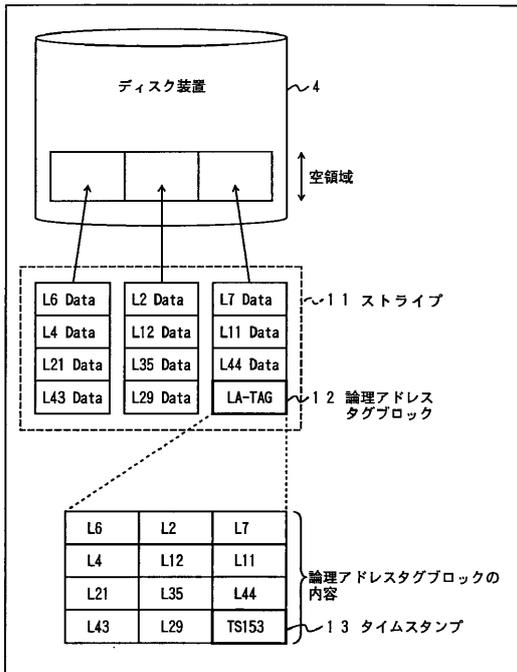
【図 1】



【図 2】



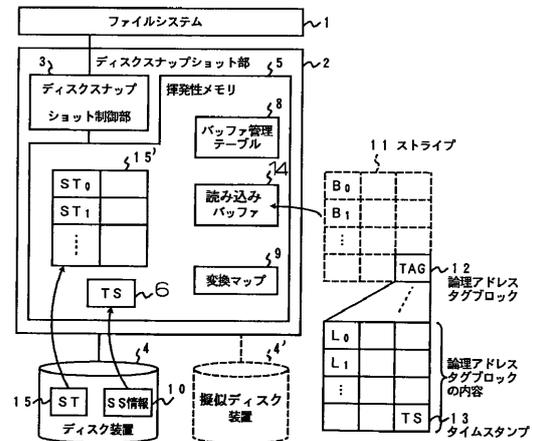
【図 3】



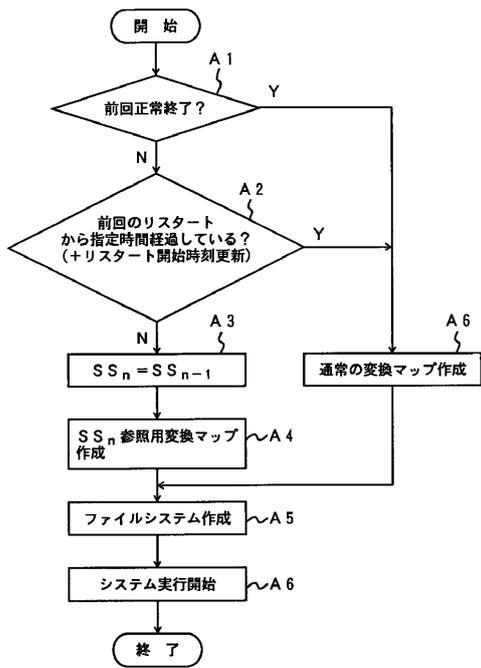
【図 4】

論理アドレス	ST#	BLK#	TS#
L ₀			
L ₁			
L ₂			
⋮			

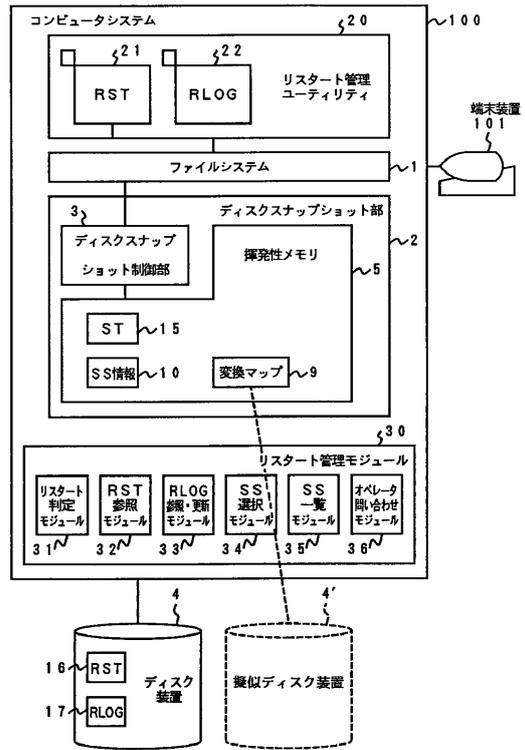
【図 5】



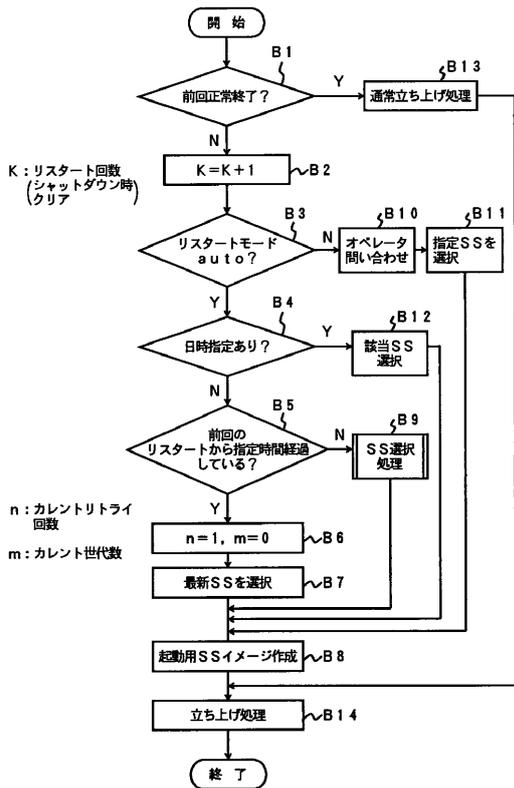
【図6】



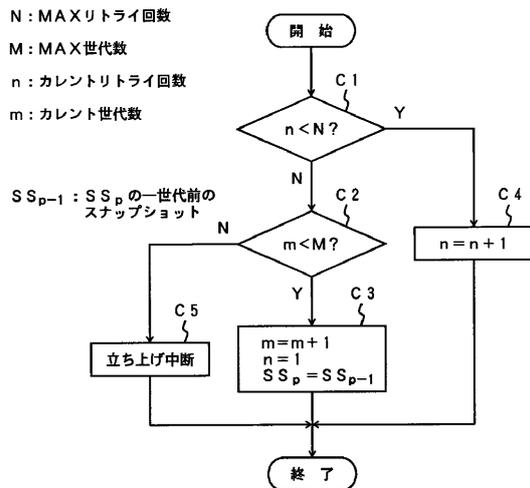
【図7】



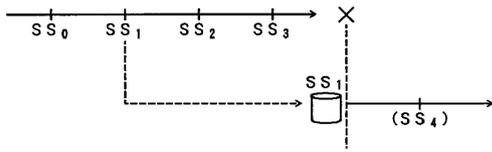
【図8】



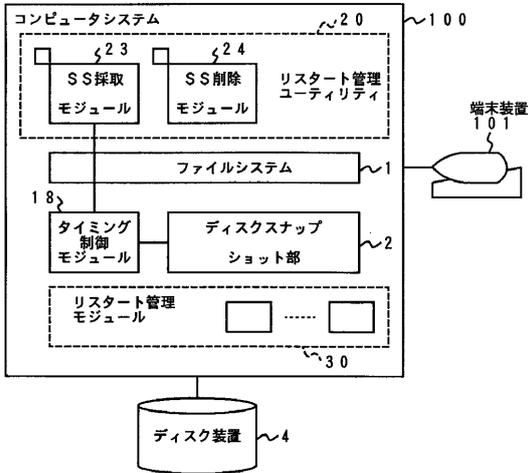
【図9】



【図10】



【図11】



フロントページの続き

(72)発明者 鹿目 奈美子
東京都青梅市末広町2丁目9番地 株式会社東芝青梅工場内

審査官 横山 佳弘

(56)参考文献 特開平07-261989(JP,A)
特開平05-289854(JP,A)
特開平07-281934(JP,A)
特開平06-110618(JP,A)
特開平06-324817(JP,A)
特開昭58-097724(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 3/06
G06F 11/14
G06F 12/00
G06F 1/00
G06F 9/06