

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6115331号  
(P6115331)

(45) 発行日 平成29年4月19日(2017.4.19)

(24) 登録日 平成29年3月31日(2017.3.31)

(51) Int. Cl. F 1  
**G 0 6 F** 11/14 (2006.01) G 0 6 F 11/14 6 0 7  
**G 0 6 F** 9/46 (2006.01) G 0 6 F 9/46 3 5 0

請求項の数 8 (全 25 頁)

(21) 出願番号	特願2013-120250 (P2013-120250)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成25年6月6日(2013.6.6)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2014-238677 (P2014-238677A)	(74) 代理人	100074099 弁理士 大菅 義之
(43) 公開日	平成26年12月18日(2014.12.18)	(74) 代理人	100133570 弁理士 ▲徳▼永 民雄
審査請求日	平成28年3月10日(2016.3.10)	(72) 発明者	大嶽 智裕 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	小高 敏裕 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 トランザクション再開プログラム、情報処理装置及びトランザクション再開方法

(57) 【特許請求の範囲】

【請求項1】

情報処理装置に、  
 仮想マシンに接続された仮想スイッチが、  
 前記仮想マシン宛ての packets を受信すると、該 packets を該仮想マシンに転送し、  
 受信した前記 packets であって、未完了のトランザクションに関連する該 packets を、  
 該未完了のトランザクションに対応付けて第1記憶部に記憶し、  
 前記仮想マシンを復元させる指示に応じて、前記仮想マシンへの前記 packets の転送を  
 停止し、  
 該仮想マシンの復元の完了に応じて、前記転送の再開指示を受信した場合に、前記未完  
 了のトランザクションに対応付けて前記第1記憶部に記憶された1または複数の packets  
 を前記仮想マシンに送信する、  
 ように処理を実行させることを特徴とするトランザクション再開プログラム。

【請求項2】

前記未完了のトランザクションに対応付けて該 packets を前記第1記憶部に記憶する場  
 合、該 packets に、該 packets の受信時刻または転送時刻を付与する  
 ことを特徴とする請求項1に記載のトランザクション再開プログラム。

【請求項3】

前記仮想スイッチが、さらに、前記 packets により形成される前記仮想マシンへの要求  
 に対する応答に関する packets を受信すると、該応答に関する packets を転送し、

前記応答に関するパケットの転送量を前記第 1 記憶部に記憶し、  
 該仮想マシンの復元完了に応じて、前記転送の再開指示を受信した場合に、さらに、前記応答に関するパケットを前記仮想マシンから受信すると、前記第 1 記憶部に記憶された前記転送量に対応するパケットを破棄し、前記転送量を超えるパケットを転送するように処理を実行させることを特徴とする請求項 1 または 2 に記載のトランザクション再開プログラム。

【請求項 4】

前記情報処理装置に、さらに、  
 仮想マシンが、  
 前記応答に関するパケットの送信を停止し、  
 前記仮想マシンが用いる記憶領域に対応する第 2 記憶部に書き込む前に該記憶領域に一時的に保持されている書込情報を、前記第 2 記憶部へ書き込み、  
 前記仮想スイッチから前記仮想マシンへの前記要求に関するパケットの転送の停止、該仮想マシンからの前記応答に関するパケットの送信の停止、及び前記書込情報の書き込みの後に、前記仮想マシンのスナップショットを作成する  
 ように処理を実行させることを特徴とする請求項 3 に記載のトランザクション再開プログラム。

10

【請求項 5】

前記情報処理装置に、さらに、  
 前記スナップショットに対応する前記仮想マシンが、前記スナップショット作成後に、  
 停止させていた前記要求に関するパケット及び前記応答に関するパケットの送信が再開している場合、実行中の処理を停止し、  
 前記仮想マシンが、前記記憶領域に一時的に保持されている書込情報を、前記第 2 記憶部へ書き込み、  
 前記仮想マシンが実行中の処理の停止、前記仮想スイッチから前記仮想マシンへの前記要求に関するパケットの転送の停止、及び前記書込情報の書き込みの後に、前記スナップショットを用いて、該スナップショットに対応する仮想マシンを復元する  
 ように処理を実行させることを特徴とする請求項 4 に記載のトランザクション再開プログラム。

20

【請求項 6】

前記情報処理装置に、さらに、  
 復元された前記仮想マシンが、  
 前記第 2 記憶部から読み出されて一時的に保持されている読出情報を破棄し、  
 前記スナップショット作成時に停止させていた前記応答に関するパケットの送信を再開し、  
 前記スナップショット作成時に行っていた処理を停止する  
 ように処理を実行させることを特徴とする請求項 5 に記載のトランザクション再開プログラム。

30

【請求項 7】

仮想マシンに接続された仮想スイッチが前記仮想マシン宛てのパケットを受信すると、  
 該仮想スイッチに、該パケットを転送させる第 1 転送制御部と、  
 受信された前記パケットであって、未完了のトランザクションに関連する該パケットを、該仮想スイッチに、該未完了のトランザクションに対応付けて第 1 記憶部に記憶させる記憶制御部と、  
 前記仮想マシンを復元させる指示に応じて、前記仮想スイッチに、前記仮想マシンへの前記パケットの転送を停止させる停止制御部と、  
 該仮想マシンの復元の完了に応じて、前記転送の再開指示を受信した場合に、前記仮想スイッチに、前記未完了のトランザクションに対応付けて前記第 1 記憶部に記憶された 1 または複数のパケットを前記仮想マシンへ送信させる送信制御部と、  
 を備えることを特徴とする情報処理装置。

40

50

## 【請求項 8】

情報処理装置は、

仮想マシンに接続された仮想スイッチに、

前記仮想マシン宛ての packets を受信すると、該 packets を該仮想マシンに転送し、

受信した前記 packets であって、未完了のトランザクションに関連する該 packets を、  
該未完了のトランザクションに対応付けて第 1 記憶部に記憶し、

前記仮想マシンを復元させる指示に応じて、前記仮想マシンへの前記 packets の転送を  
停止し、

該仮想マシンの復元の完了に応じて、前記転送の再開指示を受信した場合に、前記未完了  
のトランザクションに対応付けて前記第 1 記憶部に記憶された 1 または複数の packets  
を前記仮想マシンに送信する、

10

処理を実行させることを特徴とするトランザクション再開方法。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、トランザクションを再開させる技術に関する。

## 【背景技術】

## 【0002】

情報システムの運用者がシステムの設定変更やアプリケーションの更新をする時に設定  
変更ミスや新バージョンのバグなどによってサービスに不具合が出てしまうことがある。  
そのような状況が起きた時に即座に以前の状態に戻すために仮想マシンのスナップショット  
を用いることが行われている。

20

## 【0003】

仮想マシンのスナップショットを用いたリカバリーとは、仮想マシンの状態を保存して  
おき、後になってその状態にマシンを復元できる機能のことを指す。

## 【0004】

第 1 技術としては、通信を相互に行っている複数の仮想計算機のスナップショットを同  
時に作成する際に、送信中のデータを失うことなく整合性が取れた状態で行う技術がある  
(例えば、特許文献 1)。第 1 技術では、ゲスト OS (Operating System) A, B につい  
てのスナップショット保存を行う際に、保存実行部 A, B が、管理 OS 側仮想ネットワー  
クドライバ A, B に指示して、ゲスト OS A, B によるネットワーク送信キュー A, B へ  
のデータ格納を停止させ、ネットワーク送信キュー A, B 及びネットワーク受信キュー A  
, B が空になるのを待ち、これらキューが空になった際に、VM モニター A, B にスナ  
ップショット指示を出力して、ゲスト OS A, B の仮想ハードウェアの状態のスナップシ  
ョット保存を行わせる。

30

## 【0005】

第 2 技術としては、1 つまたは複数の仮想マシンをホストしているホストサーバは、ホ  
ストボリュームおよびそれにインストールされている 1 つまたは複数の仮想マシンを、ア  
プリケーション整合方法でバックアップする技術がある (例えば、特許文献 2)。第 2 技  
術では、ホストレベルのリクエストは、どの仮想マシンがアプリケーション整合バックア  
ップにふさわしいかを識別するようホストレベルのライタに指示する。次いで、ホストレ  
ベルのリクエストは、適切に構成された各仮想マシンにおけるゲストレベルのリクエスト  
を介して、仮想マシンバックアップを開始するようホストレベルのライタに指示する。仮  
想マシンは、仮想マシンボリューム内にアプリケーション整合バックアップを作成する。  
ホストレベルのリクエストは、次いで、ホストレベル上のサーバボリュームのスナップシ  
ョットを開始する。したがって、仮想マシンレベルのスナップショットを、サーバボリュ  
ームのホストレベルのスナップショット内から取り出すことができる。

40

## 【0006】

第 3 技術としては、システム外部への出力を待機させることなく、外部との整合性を保  
てる状態の如何を問わずに、仮想計算機の状態を障害発生時点まで復元できるようにする

50

技術がある（例えば、特許文献3）。第3技術では、通信記録装置は、外部システム及び仮想計算機との間で入出力される通信データを通信のログとして時系列順に記録する。スナップショット管理機構は、障害発生のために仮想計算機の復元が必要となった場合、スナップショット格納部に格納されている当該仮想計算機の最新のスナップショットに基づき、当該仮想計算機を当該スナップショットの取得時点（第1の時点）に復元する。ログ再生機構は、この仮想計算機に、通信記録装置によって記録されている通信のログに含まれている第1の時点から障害発生時である第2の時点までの入力データを時系列順に投入することにより、当該仮想計算機を第2の時点まで復元する。

【0007】

第4技術としては、連携中の仮想マシンを一時停止した場合でも、再開時に連携中の仮想マシン間の整合性を確保する技術がある（例えば、特許文献4）。第4技術では、管理サーバは、物理マシンを管理すると共に、物理マシン上で展開する複数の仮想マシンを管理し、複数の仮想マシンで相互に連携して制御処理を実行可能にする。管理サーバは、依存関係一覧作成部、対象管理テーブル作成部、対象提示部を有する。依存関係一覧作成部は、仮想マシンの依存関係を管理する依存関係一覧テーブルを作成する。対象管理テーブル作成部は、一時停止対象の仮想マシンの選択指示を検出すると、依存関係一覧テーブルに基づき、同仮想マシンと依存関係にある全仮想マシンを管理する一時停止順序管理テーブルを作成する。対象提示部は、一時停止順序管理テーブルに基づき、依存関係が大きい順に一時停止対象の全仮想マシンをクライアントに視覚的に提示する。

【0008】

第5技術としては、仮想マシンが動作している物理計算機に障害が発生した場合、別の物理計算機上で再生成または再起動される仮想マシンによりサービスを継続させる技術がある（例えば、特許文献5）。第5技術では、仮想マシンが動作しているサーバ計算機に障害が発生した場合、サーバ計算機の仮想マシンモニタは、障害発生時刻に最も近い時点でディスク装置に採取されたスナップショットに基づき、仮想マシンを仮想マシンとしてサーバ計算機上に再生成する。通信記録ユニットの状態再現部は、仮想マシンに対応付けられた通信履歴に基づき、スナップショットの採取時期から上記障害発生時刻までの期間における仮想マシンの状態を仮想マシンに再現させる。再起動部は、例えば仮想マシンの状態の再現に失敗した場合、仮想マシンをサーバ計算機上で再起動する。

【先行技術文献】

【特許文献】

【0009】

【特許文献1】特開2011-253350号公報

【特許文献2】特表2009-533777号公報

【特許文献3】特開2009-80705号公報

【特許文献4】特開2009-245317号公報

【特許文献5】特開2009-80692号公報

【発明の概要】

【発明が解決しようとする課題】

【0010】

しかしながら、第1技術では、スナップショットへの復元時に、スナップショット作成時点でのコネクション情報に復元されてしまうため、復元時に確立していたコネクションが維持されず、クライアントから見たらコネクションが切れてしまったかのように見えてしまう。

【0011】

本発明では、一側面として、未完了であったトランザクションが存在する場合であっても、仮想マシンの復元完了後にトランザクションを再開できる技術を提供する。

【課題を解決するための手段】

【0012】

トランザクション再開プログラムは、情報処理装置に、仮想マシンに接続された仮想ス

10

20

30

40

50

イッチが以下を行う処理を実行させる。仮想スイッチは、仮想マシン宛てのパケットを受信すると、パケットを仮想マシンに転送する。仮想スイッチは、受信したパケットであって、未完了のトランザクションに関連するパケットを、未完了のトランザクションに対応付けて第1記憶部に記憶する。仮想スイッチは、仮想マシンを復元させる指示に応じて、仮想マシンへのパケットの転送を停止する。仮想スイッチは、仮想マシンの復元の完了に応じて、転送の再開指示を受信した場合に、未完了のトランザクションに対応付けて第1記憶部に記憶された1または複数のパケットを仮想マシンに送信する。

【発明の効果】

【0013】

本発明の一側面によれば、未完了であったトランザクションが存在する場合であっても、仮想マシンの復元完了後にトランザクションを再開できる。

10

【図面の簡単な説明】

【0014】

【図1】本実施形態における情報処理装置のブロック図である。

【図2】本実施形態における物理マシンのハードウェア構成図である。

【図3】本実施形態における複数の仮想マシンを稼働させる物理マシンの構成の一例を示す。

【図4】本実施形態における仮想マシンの論理構成図である。

【図5】本実施形態における仮想システムの論理構成図である。

【図6】本実施形態におけるリクエスト再送信テーブルの一例を示す。

20

【図7】本実施形態におけるリクエスト再送信テーブルの1つの行の状態の遷移の一例を示す。

【図8】本実施形態におけるスナップショット管理テーブルの一例を示す。

【図9】本実施形態における仮想スイッチのリクエスト/レスポンス受信時の動作の一例を表したアクティビティ図である。

【図10】本実施形態におけるスナップショットを作成する場合の仮想システム内の動作の一例を表したアクティビティ図である。

【図11A】本実施形態におけるスナップショットを復元する場合の仮想システム内の動作の一例を表したアクティビティ図である。

【図11B】本実施形態におけるスナップショットを復元する場合の仮想システム内の動作の一例を表したアクティビティ図である。

30

【図12】本実施形態におけるスナップショット復元時の仮想スイッチのリクエスト再送信(S211~S213)に係るアクティビティ図である。

【図13】本実施形態の変形例におけるリクエスト再送信テーブルの一例を示す。

【発明を実施するための形態】

【0015】

個別の仮想マシンに関してはスナップショットを取ることができるが、複数の仮想マシンにまたがったシステムについてシステムレベルでスナップショットを取ることは容易ではない。その理由として、システム内のマシン間で同じタイミングでスナップショットが取れる保証がないことと、システム内のネットワークの状態が保存できていないことが挙げられる。

40

【0016】

ところが、現在の情報システムは高度化しており、単体のマシンだけではシステムは完結していないことが多い。そのため仮想マシンのスナップショットではシステム全体のスナップショットを取る必要があるとなっている。

【0017】

第1技術では、複数の仮想マシンにおいてネットワークドライバを制御することでネットワーク通信を止めてからスナップショットを作成することで、複数マシンにまたがって整合性が取れたスナップショットを作成する。しかしながら、第1技術では、スナップショットの対象のマシンとそれ以外との間の整合性について考慮されておらず、実運用を行

50

っているシステムで用いることが困難である。具体的には外部のネットワークにあるクライアントが考慮されておらず、閉じたシステム内の身を想定してしまっていること、および、データ用ディスクが考慮されていないことにより、スナップショットへの復元時にデータ変更が消失することである。また、スナップショットへの復元時に、スナップショット作成時点でのコネクション情報に復元されてしまうため、復元時に張られていたコネクションが維持されず、クライアントから見たらコネクションが切れてしまったかのように見えてしまう。

【 0 0 1 8 】

第2技術では、ホストサーバが仮想マシンに指示を出すことで、仮想マシンを停止せずに仮想マシンの整合性のあるスナップショットを作成できるようにする。NTFSのシャドーコピーを使うことでアプリケーションレベルでの整合性のあるスナップショットを作成する。しかしながら、第2技術は、個別のマシンに対するスナップショットについてであり、複数マシンの間での整合性までは保てない。

10

【 0 0 1 9 】

第3技術では、スナップショット作成から復元までの全ての通信を保存しておき、スナップショットに戻した後に保存しておいた通信を再投入することでスナップショット作成後のデータ変更に応じられるようにしている。第3技術ではスナップショットを作成時点からの通信を再現するため、スナップショットを作成したときからの時間に応じて復元に時間がかかってしまう。

【 0 0 2 0 】

20

第5技術では、仮想マシンの依存関係を用いて仮想マシンと物理マシンのシャットダウンを行うが、システムを止めることになるためサービスを継続したままスナップショットを取ることができない。

【 0 0 2 1 】

そこで、本実施形態では、未完了であったトランザクションが存在する場合であっても、仮想マシンの復元完了後にトランザクションを再開できる技術を提供する。さらに、本実施形態では、複数の仮想マシンを含めた仮想システム全体でのスナップショットの復元をサービス無停止でできるようにする技術を提供する。さらに、本実施形態では、スナップショットの作成から復元までの間にされたデータ更新を復元後も維持でき、復元に要する時間を短縮させる技術を提供する。

30

【 0 0 2 2 】

図1は、本実施形態における情報処理装置のブロック図である。情報処理装置1は、第1転送制御部2、記憶制御部3、停止制御部4、送信制御部5を含む。

【 0 0 2 3 】

第1転送制御部2は、仮想マシンに接続された仮想スイッチが仮想マシン宛ての packets を受信すると、仮想スイッチに、パケットを転送させる。第1転送制御部2の一例としては、ハイパバイザ32が挙げられる。

【 0 0 2 4 】

記憶制御部3は、受信されたパケットであって、未完了のトランザクションに関連するパケットを、仮想スイッチに、未完了のトランザクションに対応付けて第1記憶部に記憶させる。記憶制御部3の一例としては、ハイパバイザ32が挙げられる。第1記憶部の一例としては、リクエスト再送信テーブルが格納された外部記憶装置16が挙げられる。

40

【 0 0 2 5 】

停止制御部4は、仮想マシンを復元させる指示に応じて、仮想スイッチに、仮想マシンへのパケットの転送を停止させる。停止制御部4の一例としては、ハイパバイザ32が挙げられる。

【 0 0 2 6 】

送信制御部5は、仮想マシンの復元の完了に応じて、転送の再開指示を受信した場合に仮想スイッチに、未完了のトランザクションに対応付けて第1記憶部に記憶された1または複数のパケットを仮想マシンへ送信させる。

50

## 【 0 0 2 7 】

このように構成することにより、未完了であったトランザクションが存在する場合であっても、仮想マシンの復元完了後にトランザクションを再開できる。

## 【 0 0 2 8 】

記憶制御部 3 は、未完了のトランザクションに対応付けてパケットを第 1 記憶部に記憶する場合、パケットに、パケットの受信時刻または転送時刻を付与する。

## 【 0 0 2 9 】

このように構成することにより、再送信するパケットの順序を定めることができる。

情報処理装置 1 は、さらに、第 2 転送制御部 6、転送量記憶制御部 7 を含む。第 2 転送制御部 6 は、仮想スイッチがパケットにより形成される仮想マシンへの要求に対する応答に関するパケットを受信すると、仮想スイッチに、応答に関するパケットを転送させる。第 2 転送制御部 6 の一例としては、ハイパバイザ 3 2 が挙げられる。

10

## 【 0 0 3 0 】

転送量記憶制御部 7 は、仮想スイッチに、応答に関するパケットの転送量を第 1 記憶部に記憶させる。転送量記憶制御部 7 の一例としては、ハイパバイザ 3 2 が挙げられる。

## 【 0 0 3 1 】

第 2 転送制御部 6 は、仮想マシンの復元完了に応じて、転送の再開指示を受信した場合に、さらに、応答に関するパケットを仮想マシンから受信すると、仮想スイッチに、次の処理を行わせる。すなわち、第 2 転送制御部 6 は、仮想スイッチに、第 1 記憶部に記憶された転送量に対応するパケットを破棄させ、転送量を超えるパケットを転送させる。

20

## 【 0 0 3 2 】

このように構成することにより、既にクライアントに送ったパケットは送らず、未送信のパケットのみを送ることができる。

## 【 0 0 3 3 】

情報処理装置 1 は、さらに、仮想マシン制御部 8、作成部 9 を含む。仮想マシン制御部 8 は、仮想マシンに、応答に関するパケットの送信を停止させ、仮想マシンが用いる記憶領域に対応する第 2 記憶部に書き込む前に記憶領域に一時的に保持されている書込情報を、第 2 記憶部へ書き込ませる。仮想マシン制御部 8 の一例としては、ハイパバイザ 3 2 が挙げられる。

## 【 0 0 3 4 】

作成部 9 は、仮想スイッチから前記仮想マシンへの要求に関するパケットの転送の停止、仮想マシンからの応答に関するパケットの送信の停止、及び書込情報の書き込みの後に、仮想マシンのスナップショットを作成する。作成部 9 の一例としては、ハイパバイザ 3 2 が挙げられる。

30

## 【 0 0 3 5 】

このように構成することにより、仮想システムを停止させずに、スナップショットを作成できるので、仮想システムによるサービスを継続したままスナップショットを作成することができる。また、記憶領域に一時的に保持されている書込情報を、第 2 記憶部へ書き込ませ、第 2 記憶部についてはスナップショットの対象にならないから、スナップショット作成時間及びそのスナップショットを用いた復元時の時間を短縮することができる。

40

## 【 0 0 3 6 】

情報処理装置は、さらに、復元制御部 10 を含む。復元制御部 10 は、スナップショットに対応する仮想マシンを復元する。

このとき、仮想マシン制御部 8 は、スナップショット作成後に、停止させていた要求に関するパケット及び応答に関するパケットの送信が再開している場合、次の処理を行う。すなわち、仮想マシン制御部 8 は、スナップショットに対応する仮想マシンに、実行中の処理を停止させ、仮想マシンに、記憶領域に一時的に保持されている書込情報を、第 2 記憶部へ書き込ませる。

## 【 0 0 3 7 】

復元制御部 10 は、仮想マシンが実行中の処理の停止、仮想スイッチから仮想マシンへ

50

の前記要求に関するパケットの転送の停止、及び書込情報の書き込みの後に、スナップショットを用いて、該スナップショットに対応する仮想マシンを復元する。

【 0 0 3 8 】

仮想マシン制御部 8 は、さらに、

復元された前記仮想マシンに、第 2 記憶部から読み出されて一時的に保持されている読出情報を破棄させる。仮想マシン制御部 8 は、スナップショット作成時に停止させていた応答に関するパケットの送信を再開させる。仮想マシン制御部 8 は、スナップショット作成時に行っていた処理を停止させる。

【 0 0 3 9 】

このように構成することにより、スナップショットの作成から復元までの間に行われたデータの更新を復元後でも維持することができる。

10

【 0 0 4 0 】

図 2 は、本実施形態における物理マシンのハードウェア構成図である。物理マシン 1 1 は、中央演算装置 (Central Processing Unit : CPU) 1 2、メモリ 1 3、入力装置 1 4、出力装置 1 5、外部記憶装置 1 6、可搬記録媒体駆動装置 1 7、ネットワーク接続装置 1 8、バス 2 0 を含む。バス 2 0 には、CPU 1 2、メモリ 1 3、入力装置 1 4、出力装置 1 5、外部記憶装置 1 6、可搬記録媒体駆動装置 1 7、ネットワーク接続装置 1 8 が接続されている。

【 0 0 4 1 】

メモリ 1 3 は、ROM (Read Only Memory)、RAM (Random Access Memory) 等の記憶装置である。外部記憶装置 1 6 は、メモリ 1 3 に比べて大容量の記憶装置である。外部記憶装置 1 6 としては、ハードディスクドライブ、フラッシュメモリ装置、磁気ディスク装置など様々な形式の記憶装置を使用することができる。

20

【 0 0 4 2 】

CPU (物理 CPU) 1 2 は、ROM または外部記憶装置 1 6 等に格納した後述する処理を実現するプログラムを読み出し、当該プログラムを実行する。

【 0 0 4 3 】

可搬記録媒体駆動装置 1 7 は、可搬記録媒体 1 9 を読み出す装置である。可搬型記憶媒体 1 9 としては CD-ROM、フレキシブルディスク、光ディスク、光磁気ディスク、IC カード、DVD、USB メモリ装置など様々な形式の記憶媒体を使用することができる。

30

【 0 0 4 4 】

ネットワーク接続装置 1 8 は、ネットワーク 2 1 に接続された他の情報処理装置と通信するための通信インターフェースカードである。ネットワーク 2 1 は、インターネット、LAN、WAN、専用線、有線、無線等の通信網であってよい。

【 0 0 4 5 】

入力装置 1 4 には、キーボード、マウス、電子カメラ、ウェブカメラ、マイク、スキャナ、センサ、タブレット、タッチパネルなどを用いることが可能である。また、出力装置 1 5 には、ディスプレイ、プリンタ、スピーカなどを用いることが可能である。

【 0 0 4 6 】

40

後述する実施形態で説明する処理を実現するプログラムは、プログラム提供者側から通信ネットワーク 2 1、およびネットワーク接続装置 1 8 を介して、例えば外部記憶装置 1 6 に格納してもよい。また、後述する実施形態で説明する処理を実現するプログラムは、市販され、流通している可搬型記憶媒体に格納されていてもよい。この場合、この可搬型記憶媒体は可搬記録媒体駆動装置 1 7 に設定されて、CPU 1 2 によってそのプログラムが読み出されて、実行されてもよい。

【 0 0 4 7 】

図 3 は、本実施形態における複数の仮想マシンを稼働させる物理マシンの構成の一例を示す。物理マシン 1 1 にはハイパバイザ 3 2 がインストールされており、そのハイパバイザ 3 2 上で複数の仮想スイッチ 4 2 - 1, 4 2 - y と複数の仮想マシン 4 1 - 1 1 ~ 4 1

50



- 1 x , . . . . 4 1 - y 1 ~ 4 1 y z が稼働している。仮想スイッチと仮想マシンの論理的な集合をシステムと呼ぶことにする。また、複数の仮想マシン 4 1 - 1 1 ~ 4 1 - 1 x , . . . . 4 1 - y 1 ~ 4 1 y z を総称して、仮想マシン 4 1 という。複数の仮想スイッチ 4 2 - 1 , 4 2 - y を総称して、仮想スイッチ 4 2 という。

【 0 0 4 8 】

物理マシン 1 1 は、仮想化プログラムとしてのハイパバイザ 3 2 を実行することにより、複数の仮想マシン ( V M : Virtual Machine ) として機能する。物理マシン 1 1 は、図 2 で説明したように、 C P U 1 2、メモリ 1 3 等を含む物理的なデバイス群である。

【 0 0 4 9 】

物理マシン 1 1 上では、ハイパバイザ 3 2、及び複数の仮想システム 1 ~ y が稼働している。各仮想システムは、仮想スイッチ 4 2、複数の仮想マシン 4 1 を含む。例えば、仮想システム 1 は、仮想スイッチ S W \_ 1、仮想マシン V M \_ 1 1、. . .、仮想マシン V M \_ 1 x を含む。仮想システム y は、仮想スイッチ S W \_ y、仮想マシン V M \_ y 1 , . . . , V M \_ y z を含む。

10

【 0 0 5 0 】

ハイパバイザ 3 2 は、物理マシン 1 1 上で、仮想スイッチ 4 2、複数の仮想マシン 4 1 を稼働させて制御するために、仮想スイッチ 4 2、仮想マシン 4 1 に対して仮想的なハードウェア環境を提供するプログラムである。具体的は、ハイパバイザ 3 2 は、各仮想マシン 4 1 のオペレーティングシステム ( O S : Operating System ) のディスパッチ ( 実 C P U の制御権割り当て )、それらの O S が実行する特権命令のエミュレーション、物理 C P U 1 2 等のハードウェアの制御等を行う。

20

【 0 0 5 1 】

各仮想マシン 4 1 は、ハイパバイザ 3 2 上で、他の仮想マシン 4 1 と独立して稼働する仮想的な計算機である。各仮想マシン 4 1 は、それぞれの O S がハイパバイザ 3 2 を介して物理 C P U 1 2 の制御権を獲得して当該 C P U 1 2 上で実行されることにより実現される。

【 0 0 5 2 】

仮想スイッチ 4 2 は、物理マシン 1 1 が接続されるネットワークを介して、送信された通信パケット ( 以下、「パケット」という ) を受信すると、そのパケットの宛先を参照して、その宛先が示す仮想マシンへそのパケットを転送する。また、仮想スイッチ 4 2 は、物理マシン内のいずれかの仮想マシンからパケットを受信すると、その通信パケットの宛先を参照して、その宛先が示す当該物理マシン内の他の仮想マシン、他の物理マシン、または他のマシン内の仮想マシンへそのパケットを転送する。ここで、リクエストメッセージ ( 以下、「リクエスト」という )、レスポンスメッセージ ( 以下、「レスポンス」という ) 等のメッセージは、所定のサイズに分解されて、分解されたデータに所定の通信プロトコルに対応するヘッダが付与されたものをパケットという。

30

【 0 0 5 3 】

図 4 は、本実施形態における仮想マシンの論理構成図である。仮想マシン 4 1 は、 C P U 4 1 2、メモリ 4 1 3、入力装置 4 1 4、出力装置 4 1 5、システム用ディスク 4 1 6、データ用ディスク 4 1 7、ネットワーク接続装置 4 1 8、ハイパバイザ連携装置 4 1 9、バス 4 2 0 等の論理的 ( 仮想的 ) なデバイスを含む。

40

【 0 0 5 4 】

C P U 4 1 2 は、システム用ディスク 4 1 6 等に格納した後述する処理を実現するプログラムを読み出し、当該プログラムを実行する論理的な C P U である。メモリ 4 1 3 は、C P U 4 1 2 により使用される一時的にプログラムを保持したり、データを保持したりする論理的な記憶装置である。

【 0 0 5 5 】

入力装置 4 1 4 は、仮想マシン 4 1 に対してデータを入力する論理的な入力装置である。また、出力装置 4 1 5 は、仮想マシン 4 1 が用いる論理的な出力装置である。

【 0 0 5 6 】

50

システム用ディスク 4 1 6 には、OS やアプリケーションや設定ファイルなどが格納されている。データ用ディスク 4 1 7 には、アプリケーションソフトウェアが保存するデータが格納されている。システム用ディスク 4 1 6 もデータ用ディスク 4 1 7 も実体としては物理マシン 1 1 の外部記憶装置 1 6 上のファイルとして実装されていることが多い。

【 0 0 5 7 】

ネットワーク接続装置 4 1 8 は、ネットワークを介して、他の物理マシンまたは仮想マシンと通信するための論理的なネットワークインターフェースカード (NIC) である。ネットワークは、インターネット、LAN、WAN、専用線、有線、無線等の通信網であってよい。

【 0 0 5 8 】

ハイバイザ連携装置 (以下、HV 連携装置という) 4 1 9 は、ハイバイザ 3 2 と仮想マシン 1 1 の OS が連携を行うための仮想的なデバイスである。HV 連携装置 4 1 9 を使ってハイバイザ 3 2 と仮想マシン 4 1 の OS が連携することができる。また、HV 連携装置 4 1 9 は、仮想マシン 4 1 の電源管理やクリップボードの共有などを可能としている。

【 0 0 5 9 】

図 5 は、本実施形態における仮想システムの論理構成図である。仮想システム 5 1 は、ハイバイザ 3 2 上で駆動しているいずれかの仮想システムである。仮想システム 5 1 では、複数の仮想マシン 4 1 (VM\_1 ~ VM\_n) が内部ネットワーク 5 2 によって結合されている。また、仮想システム 5 1 は、内部ネットワーク 5 2 と外部ネットワーク 5 5 を

つなぐ仮想スイッチ 4 2 を含んでいる。

【 0 0 6 0 】

たとえば仮想サーバ VM\_1 (4 1 - 1) がアプリケーションサーバであり、仮想サーバ VM\_n (4 1 - n) がデータベースサーバという 2 階層システムの場合、仮想システム 5 1 は以下の動作を行う。クライアント装置 (以下、「クライアント」と称する) が外部ネットワーク 5 5 から発したリクエストは、仮想スイッチ 4 2 を通り、内部ネットワーク 5 2 を経由して仮想マシン VM\_1 に送信される。仮想マシン VM\_1 は内部ネットワーク 5 2 を通して仮想マシン VM\_n のデータベースに問合せを行い、仮想スイッチ 4 2 を通じてレスポンスをクライアントに返す。

【 0 0 6 1 】

本実施形態では、HV 連携装置 4 1 9 を通じて、仮想マシン 4 1 の受信バッファとデータ用ディスク 4 1 7 のライトバッファを制御する機構を追加することで、複数の仮想マシン 4 1 にまたがった一貫性のあるスナップショットを作成できるようにする。また、スナップショットを用いた復元においては、HV 連携装置 4 1 9 を通じて仮想マシン 4 1 に対して実行中のリクエストを中断させ、仮想スイッチ 4 2 から実行中のリクエストを再度送り直す。これにより、スナップショット復元時でもクライアントから見てリクエストが継続しているかのように見せることができる。

【 0 0 6 2 】

仮想スイッチ 4 2 は、外部ネットワーク仮想スイッチ 4 2 は、外部ネットワークから送信された、実行中のリクエストを記憶しておくバッファ 5 3 を有する。仮想スイッチ 4 2 は、外部から来るリクエストの内容を、リクエストに対するレスポンスが終わるまでバッファ 5 3 に記憶し続ける。バッファ 5 3 は、スナップショットの復元時に用いる。また、仮想スイッチ 4 2 の論理的なメモリ領域には、リクエスト再送信テーブル 5 4 が格納されている。

【 0 0 6 3 】

図 6 は、本実施形態におけるリクエスト再送信テーブルの一例を示す。図 7 は、本実施形態におけるリクエスト再送信テーブルの 1 つの行の状態の遷移の一例を示す。リクエスト再送信テーブル 5 4 は、「リクエスト ID」5 4 - 1、「リクエスト」5 4 - 2、「送信済みレスポンス長」5 4 - 3、「送信元」5 4 - 4、「送信先」5 4 - 5 のデータ項目を含む。

10

20

30

40

50

## 【 0 0 6 4 】

「リクエストID」54-1は、受信したリクエストを識別するための識別情報であり、図6の例では、追加された行順に識別番号が昇順に付与されている。「リクエスト」54-2には、受信したリクエストのボディのバイト列が記録される。ここで、リクエストボディとは、リクエストメッセージのボディ（ヘッダ以外の部分）を示す。リクエストメッセージ（リクエストボディ）は、パケットに分解されて送信され、仮想スイッチ42はその分解された各データを順次受信して、バッファに格納するが、そのバッファに格納されるデータの単位を、リクエストのボディのバイト列という。「送信済みレスポンス長」54-3は、そのリクエストに対応するレスポンスについて、送信済のレスポンス長（バイト）が格納される。「送信元」54-4には、リクエストを送信した送信元のアドレスが格納される。「送信先」54-5には、そのリクエストの送信先のアドレスが格納される。図7の詳細は、図9にて説明する。

10

## 【 0 0 6 5 】

図8は、本実施形態におけるスナップショット管理テーブルの一例を示す。スナップショット管理テーブル61は、ハイパバイザにより管理される記憶装置に格納されている。スナップショット管理テーブル61は、「システムスナップショットID」61-1、「システムID」61-2、「スナップショット作成日時」61-3、「個別マシンスナップショット」61-4のデータ項目を含む。

## 【 0 0 6 6 】

「システムスナップショットID」61-1には、各仮想システムのスナップショットを識別するための識別情報が格納される。「システムID」61-2には、スナップショットが作成された仮想システムを識別するための識別情報が格納される。「スナップショット作成日時」61-3には、スナップショットが作成された日時が格納される。「個別マシンスナップショット」61-4には、システムIDで特定される仮想システム内に含まれる仮想マシンのうち、スナップショットを作成した仮想マシンの名称とスナップショット名が格納される。図8では、システムID=1の仮想システムの場合、仮想マシン名「vm1」、「vm2」、「vm3」で特定される仮想マシンのスナップショットが作成されている。具体的には、仮想マシン「vm1」のスナップショットは、「snapshot-123」である。仮想マシン「vm2」のスナップショットは、「snapshot-234」である。仮想マシン「vm3」のスナップショットは、「snapshot-456」である。

20

30

## 【 0 0 6 7 】

図9は、本実施形態における仮想スイッチのリクエスト/レスポンス受信時の動作の一例を表したアクティビティ図である。アクティビティ図では、個別のリクエスト/レスポンスごとに状態を持たせた図にしてある。図7を参照しながら図9について説明する。

## 【 0 0 6 8 】

仮想スイッチ42は、外部ネットワーク55を介して、クライアントから新規リクエストに関するパケットが届くと、リクエスト再送信テーブル54に、新たな行を追加する。図7(A)において、追加した行の例が、リクエストID=10で示す行である。仮想スイッチ42は、新規リクエストのパケットのヘッダ情報を参照して、新たに追加した行の「送信元」54-4、「送信先」54-5にそれぞれ、送信元アドレス、送信先アドレスを登録する。

40

## 【 0 0 6 9 】

なお、新規リクエストの受信開始時点ではリクエストボディがまだ未達であるため、リクエスト再送信テーブル54の「リクエスト」54-2は空である。また、この時点では、レスポンスに関するパケットもまだ返されていないため、「送信済みレスポンス長」54-3も「0」である。「リクエスト」54-2の白い四角は、まだ受信していないリクエストボディを表し、黒い四角は受信済みのリクエストボディを表すものとする。なお、図6及び図7の例では、あらかじめ受信するリクエストボディのサイズが判明している例を記載しているが、これは説明の便宜のためであり、これに限定されず、あらかじめ受信するリクエストボディのサイズは、不明であってもよい。この場合、仮想スイッチ42は

50

、リクエストボディの最後のデータを有するパケットを解析する。

【0070】

仮想スイッチ42は、リクエストボディ待ちになっている状態で当該リクエストのリクエストボディを含むパケットを受信した場合、リクエストの宛先に応じて、そのリクエストボディを含むパケットを仮想マシンへ転送すると共に、次の処理を行う。すなわち、仮想スイッチ42は、リクエスト再送信テーブル54における、その受信したリクエストに対応する行の「リクエスト」54-2の値にリクエストボディを追記する(S001)。この時の行の例が、図7(B)において、リクエストID=10で示す行である。「リクエスト」54-2には、黒い四角で示すように、受信済みのリクエストボディが記録されている。なお、この時点では、「送信済みレスポンス長」54-3の値は「0」である。

10

【0071】

クライアントによるリクエストの送信が終了すると、当該リクエスト/レスポンスに対する仮想スイッチ42の状態は、リクエストボディ待ちからレスポンスボディ待ちになる(S002)。この時の行の例が、図7(C)において、リクエストID=10で示す行である。

【0072】

仮想スイッチ42は、内部ネットワーク52内の仮想マシン41から、上記で登録したリクエストボディに対するレスポンスボディを含むパケットを受信した場合、レスポンスの宛先に応じて、そのレスポンスボディを含むパケットを転送する。それと共に、仮想スイッチ42は、そのレスポンスボディを含むパケットを外部ネットワーク55のクライアントに送信し、送信したレスポンスボディのバイト数を「送信済みレスポンス長」54-3の値としてリクエスト再送信テーブル54に記録する。この時の行の例が、図7(D)において、リクエストID=10で示す行である。ここで、レスポンスボディとは、レスポンスメッセージのボディ(ヘッダ以外の部分)を示す。

20

【0073】

なお、受信したレスポンスが、S001で登録したリクエストボディに対するレスポンスであるかの判断の位置としては、次のものが考えられる。例えば、レスポンスのヘッダの送信元が、リクエスト再送信テーブル54の「送信先」54-5と一致し、レスポンスのヘッダの送信先が、「送信元」54-4と一致する場合、リクエストボディに対するレスポンスと判断することが考えられる。また、S001において、仮想マシンへ転送する際に、リクエストメッセージ(またはパケット)のヘッダに、リクエストIDを付与してカプセル化してもよい。この場合、仮想マシン41は、そのカプセル化されたリクエストを受信すると、そのリクエストに対する処理を行い、その処理結果(レスポンス)のヘッダ(またはパケット)に、そのリクエストIDを付与してカプセル化し、仮想スイッチ42へ送信する。仮想スイッチ42は、その受信したレスポンスのパケットを受信すると、リクエスト再送信テーブル54から、その受信したレスポンス(またはそのパケット)に付与されたリクエストIDと一致する行を検索する。付与されたリクエストIDと一致する行が検索された場合、仮想スイッチ42は、その行の「送信済みレスポンス長」54-3の値を更新する。

30

【0074】

仮想スイッチ42は、リクエストに対するレスポンスに関するパケットの送信が終了すると、リクエスト再送信テーブル54から、その送信したレスポンスに対応する行を削除して、当該リクエスト/レスポンスに対応する状態遷移を終える。

40

【0075】

図10は、本実施形態におけるスナップショットを作成する場合の仮想システム内の動作の一例を表したアクティビティ図である。アクティビティ図では、個別のリクエスト/レスポンスごとに状態を持たせた図にしてある。

【0076】

運用者がシステムの設定変更やアプリケーションの更新を行いたい場合に、仮想システムのスナップショットが作成されることについて説明する。ここで、アクティビティ図で

50

矢印を並行に複数本書いている個所は、複数の仮想マシンVMにまたがって指示を並行して出していることを示す。

【0077】

ハイバイザ32は、HV連携装置419を通じて、仮想システム51内の仮想マシン41、および仮想スイッチ42にスナップショットの作成準備の指示を出す。

【0078】

仮想マシンVM<sub>i</sub>は、ハイバイザ32からの指示を受けると、ネットワーク接続装置418の送信バッファから新たにパケットを送信しないように制御する(S101)。一方、仮想マシンVM<sub>i</sub>は、パケットの受信については今まで通り行う。

【0079】

仮想スイッチ42は、ハイバイザ32からの指示を受けると、外部ネットワーク55から内部ネットワーク52へ入るパケット(インバウンドパケット)をバッファ53に溜めておき、内部ネットワーク52には送信しないようにする(S102)。仮想スイッチ42は、内部ネットワーク52内だけで送受信されるパケット、及び内部ネットワーク52から外部ネットワーク55へのパケット(アウトバウンドパケット)については今まで通り通過させる。

【0080】

仮想スイッチ42は、内部ネットワーク52にパケットが流れなくなるまで待つ(S103)。ここでは、例えば、仮想スイッチ42は、内部ネットワーク52にパケットが流れなくなるまで所定時間待つ。

【0081】

ハイバイザ32は、HV連携装置419を通じて、仮想システム51内の仮想マシンVM<sub>i</sub>に対して、データ用ディスク417に書き込むためにバッファに保持している情報(ライトバッファ)の書き出しを指示する。この指示に基づいて、仮想マシンVM<sub>i</sub>はライトバッファをデータ用ディスク417に書き出す(S104)。S104の処理は、後にスナップショットに復元したときにライトバッファの中身をデータ用ディスク417に書き出させないようにするためである。

【0082】

その仮想マシンVM<sub>i</sub>が、HV連携装置419を介してバッファの書き出し終了を通知すると、ハイバイザ32は、その仮想マシンVM<sub>i</sub>のスナップショットを作成する(S105)。ハイバイザ32は、その作成したスナップショットを、ハイバイザ32が管理するディスク(例えば、外部記憶装置16)に格納する。このとき、ハイバイザ32は、外部記憶装置16に格納されているスナップショット管理テーブル61を更新する。ただし、データ用ディスク417はスナップショットの対象外である。すべての仮想マシンのスナップショット作成が終わった場合、ハイバイザ32は仮想システム51内の仮想マシンVM<sub>i</sub>と仮想スイッチ42に作成終了の指示を出す。

【0083】

仮想マシンVM<sub>i</sub>は、ハイバイザ32からのスナップショット作成終了指示を受信すると、ネットワーク接続装置418を用いて、送信処理を再開する(S106)。

【0084】

仮想スイッチ42は、ハイバイザ32からのスナップショット作成終了指示を受信すると、バッファ53に溜めていたインバウンドパケットを送り出し、続くインバウンドパケットも送り出すようにする(S107)。

【0085】

上記では、1台のハイバイザ32での動作例を示したが、複数のハイバイザにまたがってシステムが構築されている場合でもハイバイザ間で連携を行うことで、システム全体でのスナップショットの作成が可能である。例えば、ハイバイザ間で連携させるハイバイザ管理プログラムを用いて、各ハイバイザの動作タイミングを制御し、ハイバイザ間で連携を行うことができる。

【0086】

10

20

30

40

50

次に、スナップショットの破棄について説明する。スナップショットを作成した後に運用者が設定の変更やアプリケーションの更新を行った後、何も問題がなければ作成したスナップショットを破棄してもよい。スナップショットの破棄は各仮想マシン 4 1 のスナップショットをそのまま破棄するだけで実現できる。

【 0 0 8 7 】

ここで、スナップショットを作成した後に運用者が設定の変更やアプリケーションの更新を行った後、何らかの不具合が出てスナップショットに戻りたいということがある。しかしながら、単純にスナップショットに戻す場合には以下の問題がある。

【 0 0 8 8 】

まずは、クライアントとサーバ（仮想マシン）間で接続が確立している状態でスナップショットを用いて復元すると、クライアントとサーバも互いに相手を認知することができず、すなわち、処理中のリクエストが失われるという問題がある。

【 0 0 8 9 】

また、サーバ（仮想マシン）はクライアントからのリクエストを処理するためにリクエストのロックを外し、その旨をファイルシステムの管理情報に書き込む。この場合に、復元すると、復元された管理情報ではリクエストのロックがはずれていることが記録されていないので、リクエストが解放されていない状態（デッドロック状態）になるので、そのリクエストについて処理することができないという問題がある。

【 0 0 9 0 】

また、スナップショット作成時の状態に戻すことで、スナップショット作成時に実行していた処理が再実行されうるという問題がある。

【 0 0 9 1 】

これらの問題に対処するために、スナップショットを復元する場合に、図 1 1 の処理を行う。

【 0 0 9 2 】

図 1 1 は、本実施形態におけるスナップショットを復元する場合の仮想システム内の動作の一例を表したアクティビティ図である。アクティビティ図では、個別のリクエスト / レスポンスごとに状態を持たせた図にしてある。ここで、アクティビティ図で矢印を並行に複数本書いている個所は、複数の仮想マシン VM にまたがって指示を並行して出していることを示す。

【 0 0 9 3 】

ハイパバイザ 3 2 は、運用者から指示のあったシステムスナップショット ID に基づいて、スナップショット管理テーブル 6 1 から、そのシステムスナップショット ID を有する行を取得する。ハイパバイザ 3 2 は、その取得した行から、復元する仮想システム 5 1 の各仮想マシン 4 1 を特定する。ハイパバイザ 3 2 は、HV 連携装置 4 1 9 を通じて、その特定した仮想システム 5 1 内の仮想マシン 4 1、および仮想スイッチ 4 2 にスナップショットの復元準備の指示を出す。

【 0 0 9 4 】

仮想マシン VM<sub>i</sub> は、ハイパバイザ 3 2 からのスナップショット復元指示を受けると、リクエストに対する処理中の実行を中断する（S 2 0 1）。中断の手段としては、たとえばワーカースレッドに割り込みを入れるなどがある。例えば、Java（登録商標）アプリケーションであればリクエストを処理しているスレッドで例外を発生させ、アプリケーションの finally ブロックでリソースを解放させることができる。

【 0 0 9 5 】

仮想マシン VM<sub>i</sub> は、S 2 0 1 の中断処理が終わるまで待つ（S 2 0 2）。仮想スイッチ 4 2 からインバウンドリクエストはもう来ないため、中断処理はやがては終わる。

【 0 0 9 6 】

仮想スイッチ 4 2 は、ハイパバイザ 3 2 からのスナップショット復元指示を受けると、外部ネットワーク 5 5 から内部ネットワーク 5 2 に入っているインバウンドパケットをバッファ 5 3 に溜めておき、内部ネットワーク 5 2 には流さないようにする（S 2 0 3）。

10

20

30

40

50

仮想スイッチ 4 2 は、内部ネットワーク 5 2 内だけのパケットおよび内部ネットワーク 5 2 から外部ネットワーク 5 5 へのアウトバウンドパケットについて、今まで通り通す。

【 0 0 9 7 】

仮想スイッチ 4 2 は、内部ネットワーク 5 2 にパケットが流れなくなるまで待つ ( S 2 0 4 )。ここでは、例えば、仮想スイッチ 4 2 は、内部ネットワーク 5 2 にパケットが流れなくなるまで所定時間待つ。

【 0 0 9 8 】

ハイパバイザ 3 2 は、HV 連携装置 4 1 9 を通じて、仮想システム 5 1 内の仮想マシン VM\_i に対して、データ用ディスク 4 1 7 に書き込むためにバッファに保持している情報 (ライトバッファ) の書き出しを指示する。仮想マシン VM\_i は、その指示に基づいて、ライトバッファをデータ用ディスク 4 1 7 に書き出す ( S 2 0 5 )。

10

【 0 0 9 9 】

その仮想マシン VM\_i が、HV 連携装置 4 1 9 を介してバッファの書き出し終了を通知すると、ハイパバイザ 3 2 は、その仮想マシン VM\_i に対して、スナップショットへの復元を行う ( S 2 0 6 )。ただしデータ用ディスク 4 1 7 は、スナップショットの対象外であるため復元前の状態のままとなる。

【 0 1 0 0 】

データ用ディスク 4 1 7 は、スナップショット作成時とスナップショット復元時で内容が異なり得るため、ハイパバイザ 3 2 は、HV 連携装置 4 1 9 を通じて、仮想マシン VM\_i にデータ用ディスク 4 1 7 をマウントし直させる ( S 2 0 7 )。これにより、データ用ディスク 4 1 7 から読み出してバッファに保持されている情報 (リードバッファ) を破棄し、リードバッファとデータ用ディスク 4 1 7 との整合性を維持することができる。

20

【 0 1 0 1 】

スナップショットの状態に戻された仮想マシン VM\_i は、ネットワーク接続装置 4 1 8 の送信バッファからパケットが出ない状態になっている。この場合、ハイパバイザ 3 2 は、HV 連携装置 4 1 9 を通じて、仮想マシン VM\_i のネットワーク接続装置 4 1 8 の送信バッファからの送信を再開させる ( S 2 0 8 )。

【 0 1 0 2 】

スナップショットの状態に戻された仮想マシン VM\_i は、スナップショットを作成した時点の S 1 0 5 で実行中だったリクエストを処理しているままとなっている。そこで、ハイパバイザ 3 2 は、それらの処理を S 2 0 1 と同様の方法で中断させる ( S 2 0 9 )。

30

【 0 1 0 3 】

仮想マシン VM\_i は、S 2 0 9 の中断処理が終わるまで待つ ( S 2 1 0 )。仮想スイッチ 4 2 からインバウンドリクエストはもう来ないため、中断処理はやがては終わる。

【 0 1 0 4 】

仮想スイッチ 4 2 は、リクエスト再送信テーブル 5 4 におけるリクエストボディに、パケットヘッダを付与したパケットを仮想マシン VM\_i に再度送信する ( S 2 1 1 )。

【 0 1 0 5 】

仮想スイッチ 4 2 は、バッファ 5 3 に溜めていたインバウンドパケットを内部ネットワーク 5 2 へ送り出し、続くインバウンドパケットも内部ネットワーク 5 2 へ送り出すようにする ( S 2 1 2 )。

40

【 0 1 0 6 】

仮想スイッチ 4 2 はリクエスト再送信テーブル 5 3 における送信済みレスポンス長が正である時には、そのレスポンス長分のレスポンスに関するパケットを破棄し、それ以降のレスポンスに関するパケットのみをクライアントに返す ( S 2 1 3 )。

【 0 1 0 7 】

上記では、1 台のハイパバイザ 3 2 での動作例を示したが、複数のハイパバイザにまたがってシステムが構築されている場合でもハイパバイザ間で連携を行うことで、システム全体でのスナップショットへの復元が可能である。例えば、ハイパバイザ間で連携させるハイパバイザ管理プログラムを用いて、各ハイパバイザの動作タイミングを制御し、ハイ

50

パバイザ間で連携を行うことができる。

【0108】

図12は、本実施形態におけるスナップショット復元時の仮想スイッチのリクエスト再送信(S211~S213)に係るアクティビティ図である。図12は、図9と比較して「リクエストの再送信」と「レスポンスボディスキップ中」の動作状態が追加され、それに合わせて分岐も追加されたアクティビティ図となっている。図12では、図11のS211~S213の動作を具体的に説明する。ここでは、図6のリクエスト再送信テーブル54でのリクエストIDごとに状態遷移の例を挙げつつ説明する。

【0109】

まずは、仮想スイッチ42は、リクエストに関するパケットの再送信を行う(S301)。ここでは、図6において、リクエストIDが「1」の例では、リクエスト再送信テーブル54の「リクエスト」54-2にリクエストがないため、何もせずにS302の処理に進む。

10

【0110】

また、図6において、リクエストIDが「2」、「3」、「4」の例では、「リクエスト」54-2に保存されているリクエストボディ、すなわちバッファ53内に、再送信していないリクエストボディがある。この場合、仮想スイッチ42は、「リクエスト」54-2に保存されているリクエストをサーバ(仮想マシンVM<sub>i</sub>)に再送信する。具体的には、仮想スイッチ42は、リクエスト再送信テーブル54において、送信対象のリクエストボディの行に含まれる「送信元」54-4、「送信先」54-5を用いて、再送信用のリクエストに関するパケットヘッダを作成する。仮想スイッチ42は、その作成したヘッダを、各リクエストボディに付与してパケットを作成し、そのパケットを、送信先に再送信する。

20

【0111】

S302の条件分岐では、リクエスト/レスポンスのうちのリクエストが終了しているか否かで分岐する。S302において、図6にてリクエストIDが「3」、「4」の例では、リクエストの受信が完了であるため、処理がS305へ進む。

【0112】

S302において、リクエストIDが「1」、「2」の場合、リクエストの受信が未完なので、処理がS303へ進み、仮想スイッチ42は、クライアントからのリクエスト待ち状態になる(S303)。ここでは、クライアントからリクエストに関するパケットが送られて来たら、仮想スイッチ42は、そのパケットに含まれるリクエストボディをリクエスト再送信テーブル54に保存すると共に、サーバ(仮想マシンVM<sub>i</sub>)にそのパケットを送信する。リクエストの受信が完了した場合、処理はS304へ進む。

30

【0113】

S304の条件分岐では、リクエスト再送信テーブル54の送信済みレスポンス長によって分岐する。S304において、図6のリクエストIDが「1」、「2」、「3」の例では、送信済みレスポンス長が0であるため、処理がS306へ進む。S304において、リクエストIDが「4」の例では、送信済みレスポンス長が0よりも大きいため、処理はS305へ進む。

40

【0114】

S305において、仮想マシンVM<sub>i</sub>にはリクエストが既に送られているので、それに対するレスポンスに関するパケットが返ってくる。しかし、スナップショット作成時の状態への復元よりも前に既にクライアントにレスポンスに関するパケットの一部を送っているため、仮想スイッチ42は、送信済みレスポンス長までのレスポンスに関するパケットをクライアントに送らずに破棄する。

【0115】

仮想スイッチ42は、仮想サーバ41から、送信済みレスポンス長を超えるレスポンスに関するパケットを受信した場合、すなわち、クライアントにまだ送っていないレスポンスに関するパケットを受信した場合、仮想スイッチ42は、次の処理を行う。すなわち、

50



仮想スイッチ 4 2 は、クライアントにまだ送っていないレスポンスに関するパケットをクライアントに転送する。仮想スイッチ 4 2 は、そのレスポンスに関するパケットを全て転送した場合、そのレスポンスに関する行をリクエスト再送信テーブル 5 4 から削除する (S 3 0 6)。

【 0 1 1 6 】

次に、本実施形態の変形例について説明する。

図 1 3 は、本実施形態の変形例におけるリクエスト再送信テーブルの一例を示す。図 6 のリクエスト再送信テーブル 5 4 では、再送信するパケットの順序が定まらず、図 6 の 4 番目の行のレスポンスが再送信によって変わってしまう可能性がある。それを防ぐために、図 1 3 に示すように、リクエスト再送信テーブル 5 4 内の「リクエスト」5 4 - 2 a の各パケットに仮想マシン (サーバ) に送った時刻をタイムスタンプとして付け、再送信時にはそのタイムスタンプの順で再送信を行ってもよい。または、「リクエスト」5 4 - 2 a の各パケットに仮想スイッチが受信した時刻をタイムスタンプとして付け、再送信時にはそのタイムスタンプの順で再送信を行ってもよい。

【 0 1 1 7 】

具体的には、図 9 の S 0 0 1 において、仮想スイッチ 4 2 は、「リクエスト」5 4 - 2 の値にリクエストボディを追記する際に、仮想マシン (サーバ) 4 1 に送った時刻をタイムスタンプとして付与する。そして、リクエスト再送信時には、図 1 1 の S 2 1 1 において、仮想スイッチは、リクエスト再送信テーブル 5 4 の「リクエスト」5 4 - 2 を参照して、そのタイムスタンプの順でリクエストを仮想マシン 4 1 に再度送信する。

【 0 1 1 8 】

本実施形態によれば、ハイパバイザ 3 2 から仮想マシンのネットワーク送信バッファの送信を制御する機構が、HV 連携装置 4 1 9 に追加されている。また、ハイパバイザ 3 2 からの指示に従い、仮想マシン 4 1 の OS は、送信バッファからパケットの送信を行うか判断することができる。同様に、仮想スイッチ 4 1 のバッファ制御がハイパバイザ 3 2 から行うことができる。

【 0 1 1 9 】

ハイパバイザ 3 2 からの指示に基づきデータ用ディスク 4 1 7 へのライトバッファ書き出しを行うことにより、データの保全を図ることができる。スナップショットからの復元時にデータ用ディスク 4 1 7 を再マウントさせることにより、リードバッファ内のデータを破棄することができる。ハイパバイザ 3 2 からの指示に基づき、リクエストに対する実行中の処理を中断させる。また、スナップショットへの復元後にリクエストを再度仮想スイッチから送ることで、復元時点において貼られていたコネクションを維持することができる。

【 0 1 2 0 】

本実施形態によれば、クライアントから受信したリクエストからそのレスポンスまでの一連の処理が完結していない、すなわち未完了のトランザクションが存在する場合であっても、仮想マシンの復元完了後にトランザクションを再開することができる。また、複数マシンを含めたシステム全体でのスナップショットの復元がサービスを停止することなくできるようになる。また、スナップショットの作成から復元までに間に行われたデータ更新を復元後も維持することができる。

【 0 1 2 1 】

なお、本実施の形態は、以上に述べた実施の形態に限定されるものではなく、本実施の形態の要旨を逸脱しない範囲内で種々の構成または実施形態を取ることができる。

【 0 1 2 2 】

上記実施形態に関し、さらに、以下の付記を開示する。

(付記 1)

情報処理装置に、

仮想マシンに接続された仮想スイッチが、

前記仮想マシン宛てのパケットを受信すると、該パケットを該仮想マシンに転送し、

10

20

30

40

50

受信した前記パケットであって、未完了のトランザクションに関連する該パケットを、該未完了のトランザクションに対応付けて第 1 記憶部に記憶し、

前記仮想マシンを復元させる指示に応じて、前記仮想マシンへの前記パケットの転送を停止し、

該仮想マシンの復元の完了に応じて、前記転送の再開指示を受信した場合に、前記未完了のトランザクションに対応付けて前記第 1 記憶部に記憶された 1 または複数のパケットを前記仮想マシンに送信する、

処理を実行させることを特徴とするトランザクション再開プログラム。

(付記 2)

前記未完了のトランザクションに対応付けて該パケットを前記第 1 記憶部に記憶する場合、該パケットに、該パケットの受信時刻または転送時刻を付与する

ことを特徴とする付記 1 に記載のトランザクション再開プログラム。

(付記 3)

前記仮想スイッチが、さらに、前記パケットにより形成される前記仮想マシンへの要求に対する応答に関するパケットを受信すると、該応答に関するパケットを転送し、

前記応答に関するパケットの転送量を前記第 1 記憶部に記憶し、

該仮想マシンの復元完了に応じて、前記転送の再開指示を受信した場合に、さらに、前記応答に関するパケットを前記仮想マシンから受信すると、前記第 1 記憶部に記憶された前記転送量に対応するパケットを破棄し、前記転送量をを超えるパケットを転送する

ように処理を実行させることを特徴とする付記 1 または 2 に記載のトランザクション再開プログラム。

(付記 4)

前記情報処理装置に、さらに、

仮想マシンが、

前記応答に関するパケットの送信を停止し、

前記仮想マシンが用いる記憶領域に対応する第 2 記憶部に書き込む前に該記憶領域に一時的に保持されている書込情報を、前記第 2 記憶部へ書き込み、

前記仮想スイッチから前記仮想マシンへの前記要求に関するパケットの転送の停止、該仮想マシンからの前記応答に関するパケットの送信の停止、及び前記書込情報の書き込みの後に、前記仮想マシンのスナップショットを作成する

ように処理を実行させることを特徴とする付記 1 ~ 3 のうちいずれか 1 項に記載のトランザクション再開プログラム。

(付記 5)

前記情報処理装置に、さらに、

前記スナップショットに対応する前記仮想マシンが、前記スナップショット作成後に、停止させていた前記要求に関するパケット及び前記応答に関するパケットの送信が再開している場合、実行中の処理を停止し、

前記仮想マシンが、前記記憶領域に一時的に保持されている書込情報を、前記第 2 記憶部へ書き込み、

前記仮想マシンが実行中の処理の停止、前記仮想スイッチから前記仮想マシンへの前記要求に関するパケットの転送の停止、及び前記書込情報の書き込みの後に、前記スナップショットを用いて、該スナップショットに対応する仮想マシンを復元する

ように処理を実行させることを特徴とする付記 4 に記載のトランザクション再開プログラム。

(付記 6)

前記情報処理装置に、さらに、

復元された前記仮想マシンが、

前記第 2 記憶部から読み出されて一時的に保持されている読出情報を破棄し、

前記スナップショット作成時に停止させていた前記応答に関するパケットの送信を再開し、

10

20

30

40

50

前記スナップショット作成時に行っていた処理を停止する  
 ように処理を実行させることを特徴とする付記 5 に記載のトランザクション再開プログラム。

(付記 7)

仮想マシンに接続された仮想スイッチが前記仮想マシン宛ての packets を受信すると、  
 該仮想スイッチに、該 packets を転送させる第 1 転送制御部と、  
 受信された前記 packets であって、未完了のトランザクションに関連する該 packets を、  
 該仮想スイッチに、該未完了のトランザクションに対応付けて第 1 記憶部に記憶させる  
 記憶制御部と、

前記仮想マシンを復元させる指示に応じて、前記仮想スイッチに、前記仮想マシンへの  
 前記 packets の転送を停止させる停止制御部と、

該仮想マシンの復元の完了に応じて、前記転送の再開指示を受信した場合に、前記仮想  
 スイッチに、前記未完了のトランザクションに対応付けて前記第 1 記憶部に記憶された 1  
 または複数の packets を前記仮想マシンへ送信させる送信制御部と、

を備えることを特徴とする情報処理装置。

(付記 8)

前記記憶制御部は、前記未完了のトランザクションに対応付けて該 packets を前記第 1  
 記憶部に記憶する場合、該 packets に、該 packets の受信時刻または転送時刻を付与する  
 ことを特徴とする付記 7 に記載の情報処理装置。

(付記 9)

前記情報処理装置は、さらに、  
 前記仮想スイッチが前記 packets により形成される前記仮想マシンへの要求に対する応  
 答に関する packets を受信すると、前記仮想スイッチに、該応答に関する packets を転送  
 させる第 2 転送制御部と、

前記仮想スイッチに、前記応答に関する packets の転送量を前記第 1 記憶部に記憶させ  
 る転送量記憶制御部と、

を備え、

前記第 2 転送制御部は、該仮想マシンの復元完了に応じて、前記転送の再開指示を受信  
 した場合に、さらに、前記応答に関する packets を前記仮想マシンから受信すると、前記  
 仮想スイッチに、前記第 1 記憶部に記憶された前記転送量に対応する packets を破棄させ  
 、前記転送量を超える packets を転送させる

ことを特徴とする付記 7 または 8 に記載の情報処理装置。

(付記 10)

前記情報処理装置は、さらに、  
 仮想マシンに、前記応答に関する packets の送信を停止させ、前記仮想マシンが用いる  
 記憶領域に対応する第 2 記憶部に書き込む前に該記憶領域に一時的に保持されている書込  
 情報を、前記第 2 記憶部へ書き込ませる仮想マシン制御部と

前記仮想スイッチから前記仮想マシンへの前記要求に関する packets の転送の停止、該  
 仮想マシンからの前記応答に関する packets の送信の停止、及び前記書込情報の書き込み  
 の後に、前記仮想マシンのスナップショットを作成する作成部と、

を備えることを特徴とする付記 7 ~ 9 のうちいずれか 1 項に記載の情報処理装置。

(付記 11)

前記情報処理装置は、さらに、  
 該スナップショットに対応する仮想マシンを復元する復元制御部  
 を備え、

前記仮想マシン制御部は、前記スナップショット作成後に、停止させていた前記要求に  
 関する packets 及び前記応答に関する packets の送信が再開している場合、前記スナッ  
 プショットに対応する前記仮想マシンに実行中の処理を停止させ、前記仮想マシンに、前記  
 記憶領域に一時的に保持されている書込情報を、前記第 2 記憶部へ書き込ませ、

前記復元制御部は、前記仮想マシンが実行中の処理の停止、前記仮想スイッチから前記

10

20

30

40

50

仮想マシンへの前記要求に関するパケットの転送の停止、及び前記書込情報の書き込みの後に、前記スナップショットを用いて、該スナップショットに対応する仮想マシンを復元する

ことを特徴とする付記 10 に記載の情報処理装置。

(付記 12)

前記仮想マシン制御部は、さらに、

復元された前記仮想マシンに、前記第 2 記憶部から読み出されて一時的に保持されている読出情報を破棄させ、前記スナップショット作成時に停止させていた前記応答に関するパケットの送信を再開させ、前記スナップショット作成時に行っていた処理を停止させることを特徴とする付記 11 に記載の情報処理装置。

10

(付記 13)

情報処理装置は、

仮想マシンに接続された仮想スイッチに、

前記仮想マシン宛てのパケットを受信すると、該パケットを該仮想マシンに転送し、受信した前記パケットであって、未完了のトランザクションに関連する該パケットを、該未完了のトランザクションに対応付けて第 1 記憶部に記憶し、

前記仮想マシンを復元させる指示に応じて、前記仮想マシンへの前記パケットの転送を停止し、

該仮想マシンの復元の完了に応じて、前記転送の再開指示を受信した場合に、前記未完了のトランザクションに対応付けて前記第 1 記憶部に記憶された 1 または複数のパケット

20

を前記仮想マシンに送信する、

処理を実行させることを特徴とするトランザクション再開方法。

(付記 14)

前記情報処理装置は、

前記未完了のトランザクションに対応付けて該パケットを前記第 1 記憶部に記憶する場合、該パケットに、該パケットの受信時刻または転送時刻を付与する

ことを特徴とする付記 13 に記載のトランザクション再開方法。

(付記 15)

前記情報処理装置は、

前記仮想スイッチに、さらに、前記パケットにより形成される前記仮想マシンへの要求

30

に対する応答に関するパケットを受信すると、該応答に関するパケットを転送し、

前記応答に関するパケットの転送量を前記第 1 記憶部に記憶し、

該仮想マシンの復元完了に応じて、前記転送の再開指示を受信した場合に、さらに、前記応答に関するパケットを前記仮想マシンから受信すると、前記第 1 記憶部に記憶された前記転送量に対応するパケットを破棄し、前記転送量を超えるパケットを転送する

ように処理を実行させることを特徴とする付記 13 または 14 に記載のトランザクシ

ン再開方法。

【符号の説明】

【 0 1 2 3 】

- 1 情報処理装置
- 2 第 1 転送制御部
- 3 記憶制御部
- 4 停止制御部
- 5 送信制御部
- 6 第 2 転送制御部
- 7 転送量記憶制御部
- 8 仮想マシン制御部
- 9 作成部
- 10 復元制御部
- 11 物理マシン

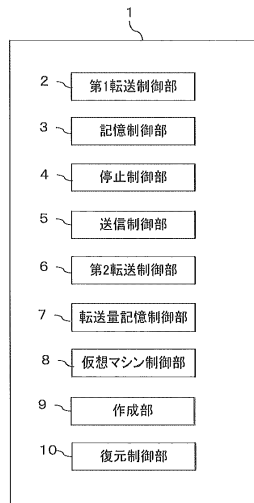
40

50

- 3 2 ハイパバイザ
- 4 1 仮想マシン
- 4 2 仮想スイッチ

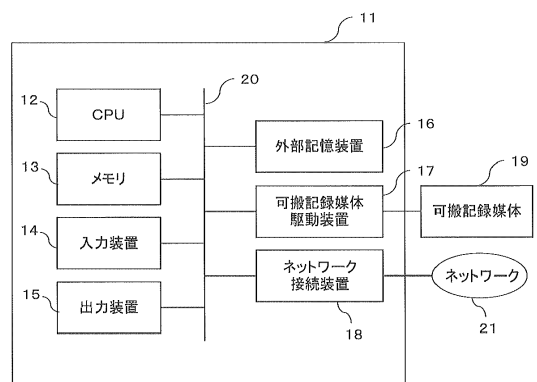
【図 1】

本実施形態における情報処理装置のブロック図



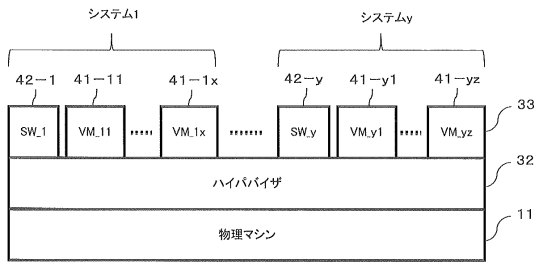
【図 2】

本実施形態における物理マシンのハードウェア構成図



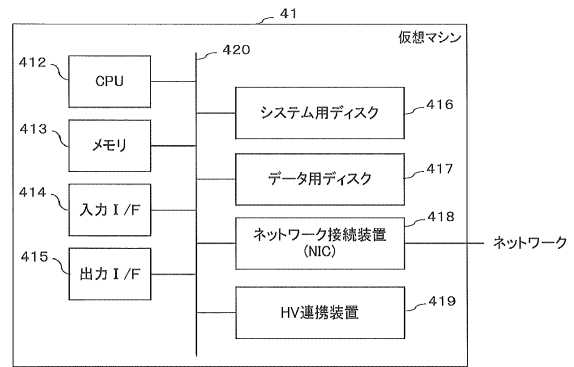
【図3】

本実施形態における複数の仮想マシンを稼働させる物理マシンの構成の一例



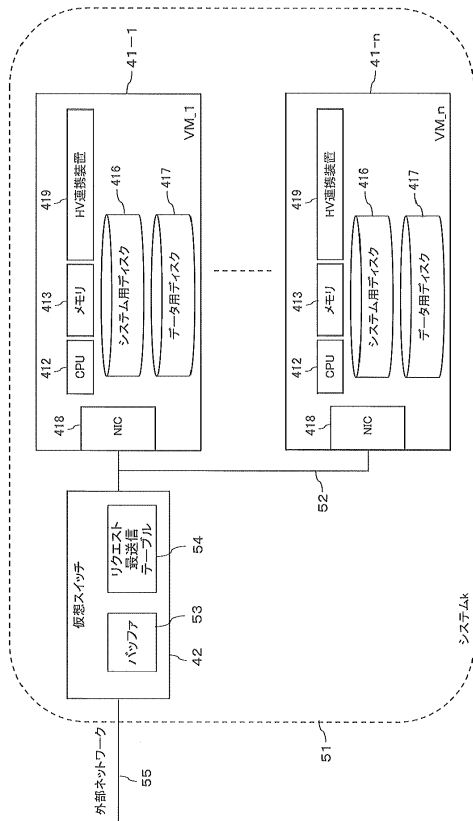
【図4】

本実施形態における仮想マシンの論理構成図



【図5】

本実施形態における仮想システムの論理構成図



【図6】

本実施形態におけるリクエスト再送信テーブルの一例

リクエストID	リクエスト	送信済みレスポンス長	送信元	送信先
1	□□□□	0	10.1.1.2	192.168.1.2
2	■□□□	0	10.1.3.4	192.168.1.2
3	■□□□	0	10.1.5.6	192.168.1.2
4	■□□□	123	10.1.7.8	192.168.1.2

【 図 7 】

本実施形態におけるリクエスト再送信テーブルの  
1つの行の状態の遷移の一例

	54-1	54-2	54-3	54-4	54-5
(A)	リクエストID	リクエスト	送信済みレスポンス長	送信元	送信先
	10	□□□□□□□□	0	10.1.1.2	192.168.1.2
	...	...	...	...	...
(B)	リクエストID	リクエスト	送信済みレスポンス長	送信元	送信先
	10	■□□□□□□□	0	10.1.1.2	192.168.1.2
	...	...	...	...	...
(C)	リクエストID	リクエスト	送信済みレスポンス長	送信元	送信先
	10	■□□□□□□□	0	10.1.1.2	192.168.1.2
	...	...	...	...	...
(D)	リクエストID	リクエスト	送信済みレスポンス長	送信元	送信先
	10	■□□□□□□□	123	10.1.1.2	192.168.1.2
	...	...	...	...	...

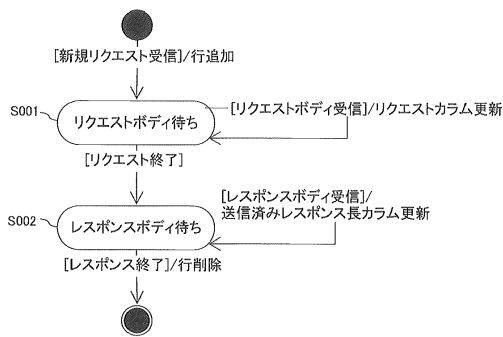
【 図 8 】

本実施形態におけるスナップショット管理テーブルの一例

	61-2	61-1	61-3	61-5
	システム スナップショットID	システムID	スナップショット 作成日時	個別マシンスナップショット
	4	system1	2013/2/4 1:23	["vm1": "snapshot-123", "vm2": "snapshot-234", "vm3": "snapshot-456"]
	7	system2	2013/3/6 12:34	["vm4": "snapshot-567", "vm5": "snapshot-789"]

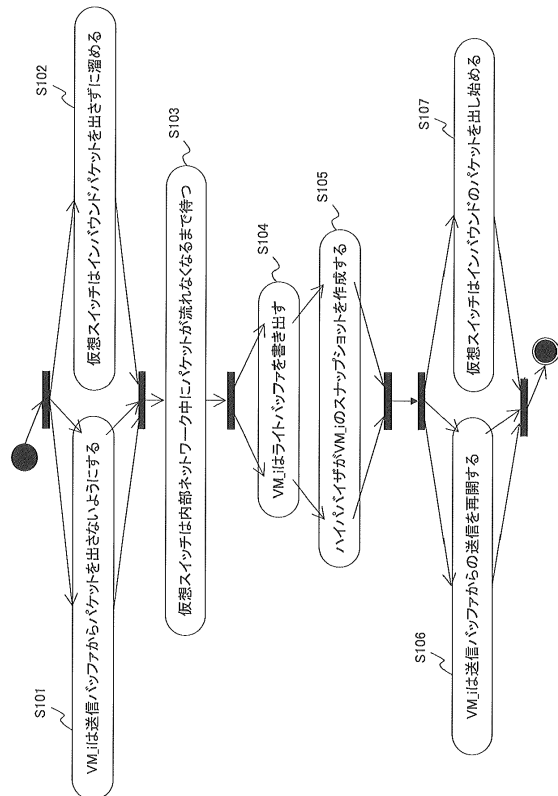
【 図 9 】

本実施形態における仮想スイッチのリクエスト/  
レスポンス受信時の動作の一例を表したアクティビティ図



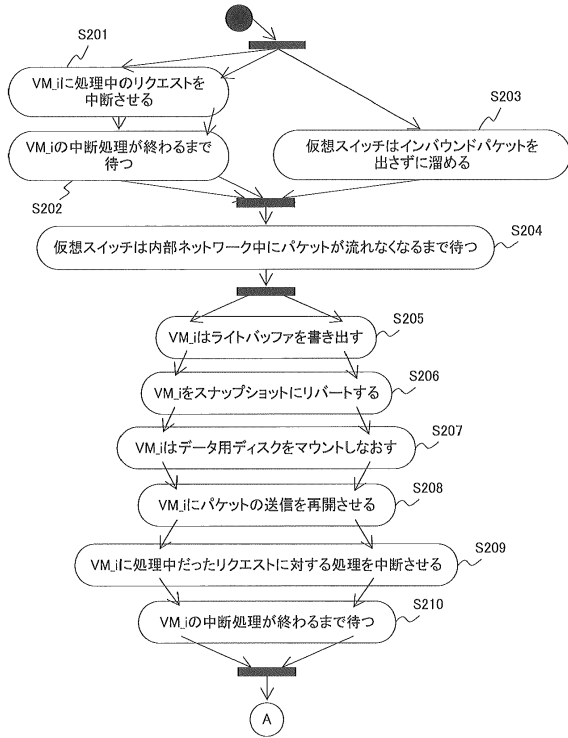
【 図 10 】

本実施形態におけるスナップショットを作成する場合の  
仮想システム内の動作の一例を表したアクティビティ図



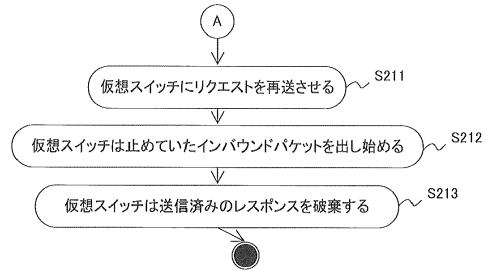
【図11A】

本実施形態におけるスナップショットを復元する場合の仮想システム内の動作の一例を表したアクティビティ図



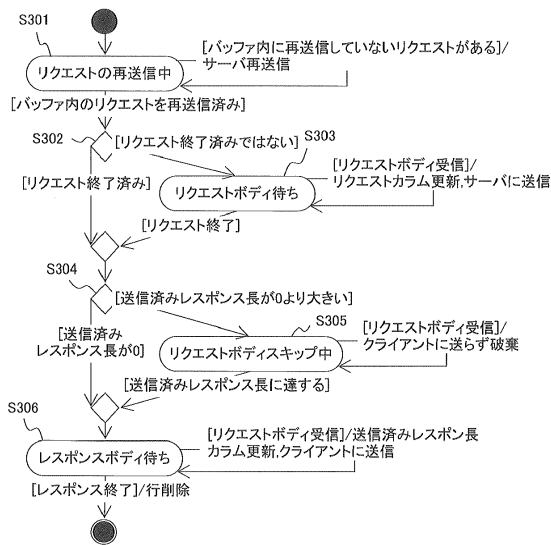
【図11B】

本実施形態におけるスナップショットを復元する場合の仮想システム内の動作の一例を表したアクティビティ図



【図12】

本実施形態におけるスナップショット復元時の仮想スイッチのリクエスト再送信(S211~S213)に係るアクティビティ図



【図13】

本実施形態の変形例におけるリクエスト再送信テーブルの一例

リクエストID	54-2a		54-3		54-4		54-5	
	送信済みレスポンス長	送信元	送信先	送信元	送信先	送信元	送信先	
1	0	10.1.1.2	192.168.1.2	0	10.1.1.2	192.168.1.2	0	10.1.1.2
2	0	10.1.3.4	192.168.1.2	0	10.1.3.4	192.168.1.2	0	10.1.3.4
3	0	10.1.5.6	192.168.1.2	0	10.1.5.6	192.168.1.2	0	10.1.5.6
4	123	10.1.7.8	192.168.1.2	123	10.1.7.8	192.168.1.2	123	10.1.7.8



---

フロントページの続き

- (72)発明者 関口 敦二  
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
- (72)発明者 堀田 勇次  
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
- (72)発明者 清水 智弘  
神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

審査官 三坂 敏夫

- (56)参考文献 特開平04-175830(JP,A)  
特開平11-265307(JP,A)  
特開2013-186745(JP,A)  
特開2012-169733(JP,A)  
米国特許出願公開第2011/0099267(US,A1)  
泉谷 洋三 他, 「運用・利用の容易さと無停止運転機能を取り込んだ超大型汎用コンピュータ」  
 , 日経コンピュータNIKKEI COMPUTER, 日本, 日経マグローヒル社, 1986年 4月28日,  
第120号, 165頁~178頁

- (58)調査した分野(Int.Cl., DB名)  
G06F 11/14  
G06F 9/46