

(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2009-80692
(P2009-80692A)

(43) 公開日 平成21年4月16日(2009.4.16)

(51) Int.Cl. F I テーマコード (参考)
G06F 11/20 (2006.01) G06F 11/20 310C 5B034
G06F 9/46 (2006.01) G06F 9/46 350

審査請求 有 請求項の数 6 O L (全 19 頁)

(21) 出願番号 特願2007-250062 (P2007-250062)
 (22) 出願日 平成19年9月26日 (2007.9.26)

(71) 出願人 000003078
 株式会社東芝
 東京都港区芝浦一丁目1番1号
 (71) 出願人 301063496
 東芝ソリューション株式会社
 東京都港区芝浦一丁目1番1号
 (74) 代理人 100058479
 弁理士 鈴江 武彦
 (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 100108855
 弁理士 蔵田 昌俊

最終頁に続く

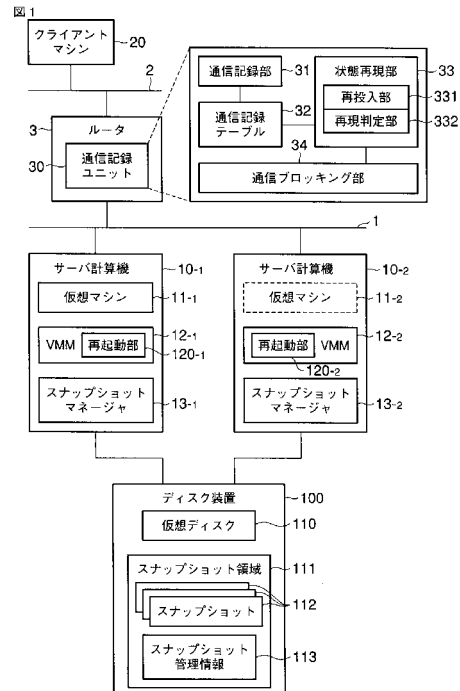
(54) 【発明の名称】 仮想計算機システム及び同システムにおけるサービス引き継ぎ制御方法

(57) 【要約】

【課題】 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。

【解決手段】 仮想マシン 11-1が動作しているサーバ計算機 10-1に障害が発生した場合、サーバ計算機 10-2の仮想マシンモタ 12-2は、障害発生時刻に最も近い時点でディスク装置 100に採取されたスナップショットに基づき、仮想マシン 11-1を仮想マシン 11-2としてサーバ計算機 10-2上に再生成する。通信記録ユニット 30の状態再現部 33は、仮想マシン 11-1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 11-1の状態を仮想マシン 11-2に再現させる。再起動部 120-2は、例えば仮想マシン 11-1の状態の再現に失敗した場合、仮想マシン 11-1をサーバ計算機 10-1上で再起動する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

仮想マシンがそれぞれ動作可能な、第 1 及び第 2 の物理計算機を含む複数の物理計算機を備えた仮想計算機システムにおいて、

前記複数の物理計算機によって共有されるディスク装置であって、前記複数の物理計算機のうちの任意の物理計算機で動作する仮想マシンが仮想ディスクとして使用可能なデータ領域を提供するディスク装置と、

前記任意の物理計算機で動作する仮想マシンによって提供されるサービスを利用するクライアントマシンと、

前記複数の物理計算機でそれぞれ動作する仮想マシンに対応付けられた通信記録テーブルに、対応する仮想マシンと当該仮想マシンによって提供されるサービスを利用するクライアントマシンとの間の通信の履歴を時系列順に記録する通信記録ユニットと

を具備し、

前記第 1 の物理計算機は、当該第 1 の物理計算機で仮想マシンが第 1 の仮想マシンとして動作する場合、当該第 1 の仮想マシンの動作状態及び当該第 1 の仮想マシンの使用前記仮想ディスクの状態を当該第 1 の仮想マシンに対応付けて定期的にスナップショットとして前記ディスク装置に採取するスナップショット管理手段を含み、

前記第 2 の物理計算機は、

当該第 2 の物理計算機で動作可能な仮想マシンを管理する仮想マシンモニタであって、当該第 2 の物理計算機とは別の前記第 1 の物理計算機上で前記第 1 の仮想マシンが動作している状態で当該第 1 の物理計算機に障害が発生した場合、当該第 1 の物理計算機の障害発生時刻に最も近い時点で当該第 1 の仮想マシンに対応付けて前記ディスク装置に採取された前記スナップショットに基づき、当該第 1 の仮想マシンを第 2 の仮想マシンとして当該第 2 の物理計算機上に再生成する仮想マシンモニタと、

当該第 2 の物理計算機上で、前記第 1 の仮想マシンを当該第 1 の仮想マシンによって使用されていた仮想ディスクに基づいて再起動する再起動手段とを含み、

前記通信記録ユニットは、前記スナップショットの採取時期から前記第 1 の物理計算機の障害発生時刻までの期間における前記第 1 の仮想マシンの状態を前記第 2 の仮想マシンに再現させるために、前記第 1 の仮想マシンに対応付けられた通信記録テーブルに記録された通信履歴のうち、前記スナップショットの採取時期から前記第 1 の物理計算機の障害発生時刻までの期間に前記第 1 の仮想マシンに送信された通信データを、前記第 2 の仮想マシンに時系列順に送信し、送信された通信データに対する前記第 2 の仮想マシンからの応答を、前記第 1 の仮想マシンに対応付けられた通信記録テーブルに記録された、当該通信データに対する前記第 1 の仮想マシンからの応答と比較することにより、前記第 1 の仮想マシンの状態の再現に成功したかを判定する状態再現手段を含み、

前記第 2 の物理計算機の前記仮想マシンモニタは、前記第 1 の仮想マシンの状態を再現できたと判定された場合、前記第 2 の仮想マシンによりサービスを継続させ、前記第 1 の仮想マシンの状態を再現できなかったと判定された場合、前記第 2 の物理計算機の前記再起動手段によって再起動される前記第 1 の仮想マシンによりサービスを継続させる

ことを特徴とする仮想計算機システム。

【請求項 2】

前記第 2 の物理計算機の前記仮想マシンモニタは、前記第 1 の仮想マシンを前記第 2 の仮想マシンとして再生成する際に、前記第 2 の物理計算機の前記再起動手段によって前記第 1 の仮想マシンを再起動させることにより、当該再起動手段の動作と前記第 2 の仮想マシンに前記第 1 の仮想マシンの状態を再現させるための前記状態再現手段の動作とを並行して実行させることを特徴とする請求項 1 記載の仮想計算機システム。

【請求項 3】

前記通信記録ユニットは、前記状態再現手段の動作期間中、前記第 2 の仮想マシンへの通信を当該状態再現手段による通信を除いてブロックする通信ブロッキング手段を更に含むことを特徴とする請求項 1 に記載の仮想計算機システム。

10

20

30

40

50

【請求項 4】

前記複数の物理計算機を接続するための第 1 のネットワークと、
 前記クライアントマシンを接続するための第 2 のネットワークと、
 前記第 1 及び第 2 のネットワークを接続するためのルータであって、前記通信記録ユニットを内蔵するルータと

を更に具備することを特徴とする請求項 3 記載の仮想計算機システム。

【請求項 5】

前記複数の物理計算機及び前記クライアントマシンを接続するためのネットワークと、
 前記通信記録ユニットを内蔵し、且つ前記ネットワークに接続されるプロキシサーバであって、前記クライアントマシンから前記複数の物理計算機の各々で動作する仮想マシンへのアクセスを代理するプロキシサーバと

を更に具備することを特徴とする請求項 3 記載の仮想計算機システム。

【請求項 6】

クライアントマシンにサービスを提供する第 1 の仮想マシンが配置される第 1 の物理計算機、及び前記第 1 の物理計算機に障害が発生した場合に、前記第 1 の仮想マシンを第 2 の仮想マシンとして再生成することが可能な第 2 の物理計算機を含む複数の物理計算機であって、当該物理計算機で仮想マシンが動作する場合、当該仮想マシンの動作状態及び当該仮想マシンの使用する仮想ディスクの状態を当該仮想マシンに対応付けて定期的にスナップショットとしてディスク装置に採取するスナップショット管理手段を含む複数の物理計算機と、前記複数の物理計算機でそれぞれ動作する仮想マシンに対応付けられた通信記録テーブルに、対応する仮想マシンと当該仮想マシンによって提供されるサービスを利用するクライアントマシンとの間の通信の履歴を時系列順に記録する通信記録ユニットとから構成される仮想計算機システムにおいて、前記第 1 の物理計算機の障害発生時に、前記第 1 の仮想マシンが提供していたサービスの引き継ぎを制御するためのサービス引き継ぎ制御方法であって、

前記第 1 の物理計算機上で前記第 1 の仮想マシンが動作している状態で当該第 1 の物理計算機に障害が発生した場合、前記第 2 の物理計算機が、前記第 1 の物理計算機の障害発生時刻に最も近い時点で当該第 1 の仮想マシンに対応付けて前記ディスク装置に採取されたスナップショットに基づき、当該第 1 の仮想マシンを前記第 2 の仮想マシンとして前記第 2 の物理計算機上に再生成するステップと、

前記スナップショットの採取時期から前記第 1 の物理計算機の障害発生時刻までの期間における前記第 1 の仮想マシンの状態を前記第 2 の仮想マシンに再現させるための状態再現処理であって、前記第 1 の仮想マシンに対応付けられた通信記録テーブルに記録された通信履歴のうち、前記スナップショットの採取時期から前記第 1 の物理計算機の障害発生時刻までの期間に前記第 1 の仮想マシンに送信された通信データを、前記第 2 の仮想マシンに時系列順に送信するステップ、当該送信された通信データに対する前記第 2 の仮想マシンからの応答を、前記第 1 の仮想マシンに対応付けられた通信記録テーブルに記録された、当該通信データに対する前記第 1 の仮想マシンからの応答と比較するステップ、及び前記比較の結果に基づいて前記第 1 の仮想マシンの状態の再現に成功したかを判定するステップとを含む状態再現処理を実行するステップと、

前記第 2 の物理計算機が、当該第 2 の物理計算機上で、前記第 1 の仮想マシンを当該第 1 の仮想マシンによって使用されていた仮想ディスクに基づいて再起動するステップと、前記状態再現処理に成功した場合、前記第 2 の仮想マシンによりサービスを継続させるステップと、

前記状態再現処理に失敗した場合、前記再起動される前記第 1 の仮想マシンによりサービスを継続させるステップと

を具備することを特徴とするサービス引き継ぎ制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

本発明は、仮想マシンが動作可能な複数の物理計算機を備えた仮想計算機システムに係り、特に、仮想マシンが動作する物理計算機の障害時のサービスの引き継ぎに好適な仮想計算機システム及び同システムにおけるサービス引き継ぎ制御方法に関する。

【背景技術】

【0002】

一般に計算機システムでは、計算機（またはプロセッサ）に障害が発生した場合に、当該計算機（またはプロセッサ）で実行されていた処理が継続可能なように、当該計算機（またはプロセッサ）の動作状態をスナップショットとして採取することが行われている（例えば、特許文献1参照）。

【0003】

また、このようなスナップショットの採取は、例えば特許文献2に記載されたようなクラスタ構成の計算機システム（クラスタシステム）においても行われている。例えば、第1及び第2のサーバから構成されるクラスタシステムにおいて、第1のサーバがクライアントに対してサービスを提供するものとする。このシステムでは、第1のサーバの動作状態（メモリの内容、CPUの状態、ディスクの内容）が定期的にスナップショットとして採取される。スナップショットには、スナップショット採取時点における第1のサーバの動作状態全てが保存される。

【0004】

したがって、第1のサーバに障害が発生した場合、その障害発生時に最も近い時点で採取された最新のスナップショットを用いることで、当該第1のサーバの最新のスナップショット採取時の状態を第2のサーバに復元することができる。つまり、障害が発生したサーバの動作を、当該サーバに関する最新のスナップショットに基づき、別のサーバで当該スナップショットの採取時の状態から再生することができる。

【特許文献1】特開2006-139621号公報

【特許文献2】特開2005-250626号公報

【発明の開示】

【発明が解決しようとする課題】

【0005】

一方、近年は、仮想マシンが動作可能な複数の物理計算機によってクラスタシステムが構成される仮想計算機システムが出現している。そこで、このような仮想計算機システムにおいても、仮想マシンの動作状態（メモリの内容、CPUの状態、ディスクの内容）をスナップショットという1つのファイルに定期的に保存することが考えられる。このようにすると、ある物理計算機に何らかの障害（例えばハードウェア障害）が発生した場合、その物理計算機上で動作していた仮想マシンの動作を、当該仮想マシンに関する最新のスナップショットを用いることで、別の物理計算機上の仮想マシンで当該スナップショットの採取時の状態から再生することが可能となる。

【0006】

しかしながら、仮想マシンがスナップショットの採取時（障害発生時に最も近いステップの採取時）から再生される場合、つまり仮想マシンの動作状態が障害発生時よりも前の状態に戻される場合、仮想マシンに接続（論理的に接続）して当該仮想マシンからのサービスの提供を受けていたクライアントマシンとの間に動作状態の不整合が発生してしまう。

【0007】

一方、複数の物理計算機（サーバ計算機）から構成される通常のクラスタシステムでは、クライアントに対してサービスを提供しているサーバ計算機の障害時には、当該サービスを引き継ぐ別のサーバ計算機上で、障害が発生したサーバ計算機で実行されていたOS（オペレーティングシステム）やアプリケーションの再起動が行われる。そこで、このような通常のクラスタシステムで適用されているサーバ計算機の障害時の処理を、仮想マシンによってクラスタシステムが構成される仮想計算機システムに適用することが考えられる。しかし、OSやアプリケーションの起動には、スナップショットからの仮想マシンの

10

20

30

40

50

再生に比べて長時間を要し、アプリケーションのロールバックが必要となる場合には更に時間を要する。

【 0 0 0 8 】

本発明は上記事情を考慮してなされたものでその目的は、クライアントマシンに対してサービスを提供する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に、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高速に且つ確実にサービスを継続させることができる仮想計算機システム及び同システムにおけるサービス引き継ぎ制御方法を提供することにある。

【 課題を解決するための手段 】

【 0 0 0 9 】

本発明の1つの観点によれば、仮想マシンがそれぞれ動作可能な、第1及び第2の物理計算機を含む複数の物理計算機を備えた仮想計算機システムが提供される。このシステムは、前記複数の物理計算機によって共有されるディスク装置であって、前記複数の物理計算機のうちの任意の物理計算機で動作する仮想マシンが仮想ディスクとして使用可能なデータ領域を提供するディスク装置と、前記任意の物理計算機で動作する仮想マシンによって提供されるサービスを利用するクライアントマシンと、前記複数の物理計算機でそれぞれ動作する仮想マシンに対応付けられた通信記録テーブルに、対応する仮想マシンと当該仮想マシンによって提供されるサービスを利用するクライアントマシンとの間の通信の履歴を時系列順に記録する通信記録ユニットとを具備する。前記第1の物理計算機は、当該第1の物理計算機で仮想マシンが第1の仮想マシンとして動作する場合、当該第1の仮想マシンの動作状態及び当該第1の仮想マシンの使用する前記仮想ディスクの状態を当該第1の仮想マシンに対応付けて定期的にスナップショットとして前記ディスク装置に採取するスナップショット管理手段を含む。前記第2の物理計算機は、当該第2の物理計算機で動作可能な仮想マシンを管理する仮想マシンモニタであって、当該第2の物理計算機とは別の前記第1の物理計算機上で前記第1の仮想マシンが動作している状態で当該第1の物理計算機に障害が発生した場合、当該第1の物理計算機の障害発生時刻に最も近い時点で当該第1の仮想マシンに対応付けて前記ディスク装置に採取された前記スナップショットに基づき、当該第1の仮想マシンを第2の仮想マシンとして当該第2の物理計算機上に再生成する仮想マシンモニタと、当該第2の物理計算機上で、前記第1の仮想マシンを当該第1の仮想マシンによって使用されていた仮想ディスクに基づいて再起動する再起動手段とを含む。前記通信記録ユニットは、前記スナップショットの採取時期から前記第1の物理計算機の障害発生時刻までの期間における前記第1の仮想マシンの状態を前記第2の仮想マシンに再現させるために、前記第1の仮想マシンに対応付けられた通信記録テーブルに記録された通信履歴のうち、前記スナップショットの採取時期から前記第1の物理計算機の障害発生時刻までの期間に前記第1の仮想マシンに送信された通信データを、前記第2の仮想マシンに時系列順に送信し、送信された通信データに対する前記第2の仮想マシンからの応答を、前記第1の仮想マシンに対応付けられた通信記録テーブルに記録された、当該通信データに対する前記第1の仮想マシンからの応答と比較することにより、前記第1の仮想マシンの状態の再現に成功したかを判定する状態再現手段を含む。前記仮想マシンモニタは、前記第1の仮想マシンの状態を再現できたと判定された場合、前記第2の仮想マシンによりサービスを継続させ、前記第1の仮想マシンの状態を再現できなかったと判定された場合、前記再起動手段によって再起動される前記第1の仮想マシンによりサービスを継続させる。

【 発明の効果 】

【 0 0 1 0 】

本発明によれば、クライアントマシンに対してサービスを提供する仮想マシン（第1の仮想マシン）が動作している物理計算機（第1の物理計算機）に障害が発生した場合に、別の物理計算機（第2の物理計算機）上で再生成される仮想マシン（第2の仮想マシン）または再起動される仮想マシン（第1の仮想マシン）により高速に且つ確実にサービスを継続させることができる。特に本発明においては、スナップショット及び通信履歴に基づ

10

20

30

40

50

き第1の仮想マシンの状態を再現できた場合には、障害発生直前の状態から極めて速やかにサービスを継続させることができる。

【発明を実施するための最良の形態】

【0011】

以下、本発明の実施の形態につき図面を参照して説明する。

図1は本発明の一実施形態に係る仮想計算機システムの構成を示すブロック図である。図1において、ネットワーク(第1のネットワーク)1には、複数のサーバ計算機(物理サーバ計算機、物理計算機)、例えば2台のサーバ計算機10-1及び10-2が接続されている。

【0012】

サーバ計算機10-1及び10-2は、CPU、I/O装置及びメモリのような周知のハードウェア資源(図示せず)を備えている。サーバ計算機10-1及び10-2は、当該計算機10-1及び10-2によって共有されるディスク装置100と接続されている。つまりディスク装置100は、サーバ計算機10-1及び10-2が共通に有するハードウェア資源である。

【0013】

サーバ計算機10-1及び10-2が有するハードウェア資源は、仮想化されることにより、仮想マシン(Virtual Machine: VM)が動作する環境(仮想マシン実行環境)を提供する。図1では、サーバ計算機10-1の仮想マシン実行環境で仮想マシン11-1が動作している状態が示されている。この仮想マシン実行環境は、当該実行環境がディスク装置100のうちの仮想マシン11-1に割り当てられる(仮想マシン11-1が利用可能な)仮想化されたディスク領域である仮想ディスク110を含む。仮想ディスク110の内容は、後述する仮想マシンモニタ12-1及び12-2からは、1つのファイルとして認識される。

【0014】

仮想マシン11-1が動作するサーバ計算機10-1に障害が発生した場合、当該仮想マシン11-1が提供するサービスを、別のサーバ計算機、例えばサーバ計算機10-2側に引き継がせるために、当該サーバ計算機10-2に仮想マシン11-1に相当する仮想マシン11-2が生成(再生成)される。図1では、仮想マシン11-2が破線のブロックで示されている。このことは、図1の状態では、未だ仮想マシン11-2がサーバ計算機10-2上に生成されていないことを示す。

【0015】

ネットワーク(第2のネットワーク)2には、クライアントマシン20が接続されている。ネットワーク1及びネットワーク2はルータ3によって接続されている。クライアントマシン20は、サーバ計算機10-1及び10-2上で仮想マシンが動作する場合に、当該仮想マシンの提供するサービスを利用するために、ネットワーク2、ルータ3及びネットワーク1を介して当該仮想マシンと通信を行う。図1の例では、クライアントマシン20は、サーバ計算機10-1上で動作する仮想マシン11-1と通信を行う。

【0016】

ルータ3は通信記録ユニット30を有する。通信記録ユニット30は通信記録部31、通信記録テーブル32、状態再現部33及び通信ブロッキング部34を含む。通信記録テーブル32は仮想マシン11-1に対応して用意される。

【0017】

通信記録部31は、クライアントマシン20とサーバ計算機10-1上で動作する仮想マシン11-1との間でネットワーク1及び2を介して行われる通信の履歴を通信記録テーブル32に時系列順に記録する。本実施形態において、通信記録ユニット30には、仮想マシン毎に通信記録テーブルが用意される。クライアントマシンと仮想マシンとの間の通信の履歴は、その仮想マシンに対応する通信記録テーブルに記録される。

【0018】

状態再現部33は、例えば仮想マシン11-1が動作するサーバ計算機10-1で障害が発生した場合に、障害発生直前のスナップショット採取時から障害発生時までの期間に通信

10

20

30

40

50

記録テーブル 3 2 に記録された通信の履歴に基づき、当該仮想マシン 1 1 -1 の障害発生時の状態を再現する。状態再現部 3 3 は、再投入部 3 3 1 及び再現判定部 3 3 2 を含む。

【 0 0 1 9 】

再投入部 3 3 1 は、通信記録テーブル 3 2 に記録された通信の履歴のうち、上述の障害発生直前のスナップショット採取時から障害発生時までの期間にクライアントマシン 2 0 から仮想マシン 1 1 -1 に送信された通信データを、状態再現の対象となる仮想マシンに時系列順に送信（投入）する。再現判定部 3 3 2 は、再投入部 3 3 1 によって送信（投入）された通信データに対する仮想マシンからの応答と、通信記録テーブル 3 2 に記録されている当該通信データに対する応答とを比較することにより、障害発生時の仮想マシン 1 1 -1 の状態が再現されたかを判定する。

10

【 0 0 2 0 】

通信ブロッキング部 3 4 は、状態再現部 3 3 が通信記録テーブル 3 2 に基いて仮想マシンの障害発生時の状態を再現する処理を行っている期間、当該通信記録テーブル 3 2 に記録されている通信データを送信していたクライアントマシン（ここではクライアントマシン 2 0 ）から当該仮想マシンへのアクセスをブロックする。

【 0 0 2 1 】

サーバ計算機 1 0 -1 及び 1 0 -2 上では、ハイパバイザである仮想マシンモニタ（Virtual Machine Monitor：VMM）1 2 -1 及び 1 2 -2 がそれぞれ動作する。仮想マシンモニタ 1 2 -1 及び 1 2 -2 は、仮想マシンマネージャとも呼ばれる。仮想マシンモニタ 1 2 -1 及び 1 2 -2 は、それぞれ、サーバ計算機 1 0 -1 及び 1 0 -2 が有する上述のハードウェア資源の利用を管理することで、サーバ計算機 1 0 -1 及び 1 0 -2 上で動作する仮想マシンを管理する。例えば仮想マシンモニタ 1 2 -1 及び 1 2 -2 は、サーバ計算機 1 0 -1 及び 1 0 -2 が有するハードウェア資源を仮想化することにより仮想マシンが動作する仮想マシン実行環境を提供する。つまり仮想マシンモニタ 1 2 -1 及び 1 2 -2 は、仮想化されたハードウェア資源を有する仮想マシンを構築する。

20

【 0 0 2 2 】

仮想マシンモニタ 1 2 -1 及び 1 2 -2 は、それぞれ再起動部 1 2 0 -1 及び 1 2 0 -2 を含む。再起動部 1 2 0 -i（ $i = 1, 2$ ）は、サーバ計算機 1 0 -j（ $j = 1, 2$ 、但し $j \neq i$ ）に障害が発生した場合に、当該サーバ計算機 1 0 -j で動作していた仮想マシンをサーバ計算機 1 0 -i 上で起動する。本実施形態では、サーバ計算機 1 0 -j で動作していた仮想マシンの状態再現（通信記録ユニット 3 0 内の状態再現部 3 3 による状態再現）に失敗したことをもって、再起動部 1 2 0 -i による仮想マシン起動処理が開始される。

30

【 0 0 2 3 】

サーバ計算機 1 0 -1 及び 1 0 -2 上ではまた、スナップショットマネージャ 1 3 -1 及び 1 3 -2 がそれぞれ動作する。スナップショットマネージャ 1 3 -1 及び 1 3 -2 は、サーバ計算機 1 0 -1 及び 1 0 -2 上で仮想マシンが動作する場合に、定期的に当該仮想マシンの動作状態及び当該仮想マシンが利用する仮想ディスクの内容をスナップショットとしてディスク装置 1 0 0 に採取（格納）する。仮想マシンの動作状態は、当該仮想マシンに割り当てられている CPU の状態（プログラムカウンタ及びレジスタの状態）及びメモリの状態を含む。

40

【 0 0 2 4 】

図 1 の例では、ディスク装置 1 0 0 には、サーバ計算機 1 0 -1 上で動作する仮想マシン 1 1 -1 に対応するスナップショット領域 1 1 1 が確保されている。このスナップショット領域 1 1 1 は、仮想マシン 1 1 -1 の動作状態及び仮想ディスク 1 1 0 の内容をスナップショット 1 1 2 として定期的に格納するのに用いられる。スナップショット領域 1 1 1 には、当該領域 1 1 1 に格納されたスナップショット 1 1 2 の列を管理するスナップショット管理情報 1 1 3 も格納される。

【 0 0 2 5 】

次に、図 1 の仮想計算機システムにおける動作を説明する。

今、クライアントマシン 2 0 が、サーバ計算機 1 0 -1 上で動作する仮想マシン 1 1 -1 の

50

提供するサービスを利用するために、ネットワーク 2、ルータ 3 及びネットワーク 1 を介して当該仮想マシン 1 1-1 との間で通信を行っているものとする。この場合、ルータ 3 に含まれている通信記録ユニット 3 0 内の通信記録部 3 1 は、クライアントマシン 2 0 と仮想マシン 1 1-1 との間の通信シーケンスで発生した全ての通信の履歴を通信記録テーブル 3 2 に時系列順に記録する。

【 0 0 2 6 】

図 2 は、通信記録テーブル 3 2 に記録された通信の履歴の例を示す。ここでは、クライアントマシン 2 0 と仮想マシン 1 1-1 との間の 1 回の通信毎に、通信記録部 3 1 によってシーケンシャルに割り当てられる通信番号、通信が行われた時刻（通信時刻）、通信の方向（通信データの流れる方向）及び通信データの組が、通信記録テーブル 3 2 に記録される。図の例では、通信の方向を、クライアントマシン 2 0 仮想マシン 1 1-1 を「 I N 」

10

、その逆を「 O U T 」で表記している
一方、仮想マシン 1 1-1 が動作するサーバ計算機 1 0-1 では、スナップショットマネージャ 1 3-1 が、当該仮想マシン 1 1-1 の動作状態と当該仮想マシン 1 1-1 が利用する仮想ディスク 1 1 0 の内容を、ディスク装置 1 0 0 に確保されている仮想マシン 1 1-1 用のスナップショット領域 1 1 1 にスナップショット 1 1 2 として定期的に（例えば時間 T 毎に）採取（格納）している。スナップショットマネージャ 1 3-1 は、スナップショット 1 1 2 を採取する都度、当該採取されたスナップショット 1 1 2 の世代管理のためにスナップショット管理情報 1 1 3 を更新する。

【 0 0 2 7 】

20

図 3 は、時刻 t_0 、 t_1 及び t_2 のそれぞれで、スナップショット領域 1 1 1 にスナップショット 1 1 2（# a）、1 1 2（# b）及び 1 1 2（# c）が採取された様子を示す。時刻 t_0 、 t_1 及び t_2 のそれぞれにおけるスナップショット 1 1 2（# a）、1 1 2（# b）及び 1 1 2（# c）は、スナップショット管理情報 1 1 3 によって世代管理される。

【 0 0 2 8 】

スナップショット 1 1 2（# a）は、時刻 t_0 における仮想マシン 1 1-1 の状態（動作状態）# 1 及び仮想ディスク 1 1 0 の内容 # A を含む。スナップショット 1 1 2（# b）は、時刻 t_1 における仮想マシン 1 1-1 の状態 # 2 及び仮想ディスク 1 1 0 の内容 # B を含む。スナップショット 1 1 2（# c）は、時刻 t_2 における仮想マシン 1 1-1 の動作状態 # 3 及び仮想ディスク 1 1 0 の内容 # C を含む。

30

【 0 0 2 9 】

このような状態で、時刻 t_2 と次にスナップショット 1 1 2 が採取されるべき時刻 t_3 との間の時刻 $t_2 3$ において、サーバ計算機 1 0-1 に障害（例えばハードウェア障害）が発生したものとする。

【 0 0 3 0 】

ここで、サーバ計算機 1 0-1 に障害が発生した場合の本実施形態における動作について説明する前に、従来技術の動作について説明する。ここでは便宜的に、図 1 に示す仮想計算機システムにおいて、サーバ計算機 1 0-1 上で動作する仮想マシン 1 1-1 で実行されていたサービスを、他のサーバ計算機 1 0-2 上に生成された仮想マシンに従来技術によって引き継がせることで、障害回復を図るものとする。

40

【 0 0 3 1 】

まず、サーバ計算機 1 0-2 上で動作する仮想マシンモニタ 1 2-2 は、サーバ計算機 1 0-1 の障害を検出すると、ディスク装置 1 0 0 のスナップショット領域 1 1 1 に保持されたスナップショット 1 1 2 の列のうち、障害発生時刻 $t_2 3$ に最も近い時刻で採取されたスナップショット、即ち時刻 t_2 で採取されたスナップショットを用いて、サーバ計算機 1 0-2 上に仮想マシンを生成する。つまり仮想マシンモニタ 1 2-2 は、障害発生時刻 $t_2 3$ に最も近いスナップショット採取時刻 t_2 における仮想マシン 1 1-1 と全く同じ状態の仮想マシン 1 1-2 をサーバ計算機 1 0-2 上に生成する。このときクライアントマシン 2 0 は、サーバ計算機 1 0-1 上の仮想マシン 1 1-1 の接続から、サーバ計算機 1 0-2 上に生成さ

50

れた仮想マシン10-2との接続に切り替えられる。一般に、この接続切り替えはクライアントマシン20から認識できず、当該クライアントマシン20は同一の仮想マシンに接続されているとして動作する。

【0032】

サーバ計算機10-2上に生成された仮想マシン11-2は、仮想マシン11-1によって実行されていたアプリケーションを起動して、時刻t2の状態から当該アプリケーションに従う動作を再開する。ところが、仮想マシン11-1とクライアントマシン20との通信は、時刻t2より先の時刻t23まで進んでいる。この場合、生成された仮想マシン11-2が時刻t2の状態から時刻t23の状態まで、クライアントマシン20との間で以前と同一の通信を再現できるとは限らない。そこで従来技術では、時刻t23で未完了のトランザクションがある場合、生成された仮想マシン11-2は、そのトランザクションから処理を再開する。このため、処理が著しく遅延する。

10

【0033】

これに対して本実施形態では、仮想マシン11-1が動作するサーバ計算機10-1に障害（ハードウェア障害）が発生した場合、障害発生時刻に最も近い時刻で採取されたスナップショットだけでなく、通信記録ユニット30内の通信記録部31によって通信記録テーブル32に記録された通信の履歴（通信記録）も用いて、サービスの引き継ぎが行われる。

【0034】

以下、サーバ計算機10-1に障害が発生した場合の障害回復のための動作について、図4乃至図7及び並びに先に挙げた図3を参照して説明する。図4は時系列に沿った仮想マシン、仮想ディスク、スナップショットの状態及び仮想マシン再生成/再起動を説明するための図、図5はサーバ計算機障害発生時のサービス引き継ぎのための手順を示すフローチャート、図6は通信記録ユニット30内の状態再現部33による状態再現処理の手順を示すフローチャート、図7は通信記録ユニット30内の通信ブロッキング部34による通信ブロッキング処理の手順を示すフローチャートである。

20

【0035】

まず、仮想マシン11-1が動作するサーバ計算機10-1では、スナップショットマネージャ11-1が、当該仮想マシン11の動作状態と当該仮想マシン11が利用する仮想ディスク110の内容を、前述のように仮想マシン11用のスナップショット領域111にスナップショット112として時間T毎に採取している。

30

【0036】

これにより、図3に示されているように、時刻t0、t1及びt2のそれぞれで、スナップショット領域111にスナップショット112（#a）、112（#b）及び112（#c）が採取されたものとする。

【0037】

そして時刻t2と次にスナップショット112（#D）が採取されるべき時刻t3との間の時刻t23において、サーバ計算機10-1に障害（ハードウェア障害）が発生したものとす。図4には、時刻t23における仮想マシン11-1の状態及び仮想ディスク110の内容が、それぞれ#3'及び#c'であることが表されている。また、サーバ計算機10-1に障害が発生したことが、サーバ計算機10-2上で動作する仮想マシンモニタ12-2によって検出されたものとする。

40

【0038】

サーバ計算機（10-1）の障害（ハードウェア障害）検出は、クラスタソフトウェアの持つハートビートによるサーバ死活チェックや、運用管理ソフトなどの機能をもって実現されることが、従来から知られている。そこで、仮想マシンモニタ12-2が、このような外部の管理手段から障害の通知を受ける構成としても、当該仮想マシンモニタ12-2自身に、当該外部の管理手段が有するのと同様のサーバ計算機障害検出機能を備える構成としても構わない。

【0039】

50

さて、サーバ計算機 10-2上で動作する仮想マシンモニタ 12-2は、サーバ計算機 10-1の障害を検出すると、ディスク装置 100内の(サーバ計算機 10-1上で動作する仮想マシン 11-1に対応する)スナップショット領域 111に保持されたスナップショットに基づき、サーバ計算機 10-2上に仮想マシン 11-2を生成する(ステップ S1)。更に具体的に述べるならば、仮想マシンモニタ 12-2は、スナップショット領域 111に保持されたスナップショット 112の列のうち、障害発生時刻 t_{23} に最も近い時刻 t_2 で採取されたスナップショット 112(#c)を用いて、サーバ計算機 10-2上に当該時刻(スナップショット採取時刻) t_2 における仮想マシン 11-1と全く同じ状態の仮想マシンを仮想マシン 11-2として生成(再生成)する。これにより仮想マシン 11-2は、時刻 t_2 における仮想マシン 11-1と同一の状態、当該仮想マシン 11-1が実行していたのと同様のアプリケーションプログラムを実行できる。

10

【0040】

次に仮想マシンモニタ 12-2は、通信記録ユニット 30に対して、サーバ計算機 10-1の障害発生時における仮想マシン 11-1の状態を仮想マシン 11-2に再現(復元)させるための状態再現(復元)処理を要求する(ステップ S2)。すると通信記録ユニット 30内の状態再現部 33は、要求された状態再現処理を、通信記録テーブル 32に記録された通信の履歴に基づき実行する(ステップ S3)。

【0041】

以下、状態再現部 33によって実行される状態再現処理の手順について説明する。まず状態再現部 33内の再投入部 331は、通信記録テーブル 32に記録された通信の履歴のうち、スナップショット採取時刻 t_2 から障害発生時刻 t_{23} までの期間にクライアントマシン 20から仮想マシン 11-1に送信された通信データを、当該クライアントマシン 20に代わって、時系列順にネットワーク 1を介して仮想マシン 11-2に順次送信する動作を開始する。この通信記録テーブル 32に基づく通信記録ユニット 30の送信動作を、ネットワーク通信再投入(またはネットワーク I/O再現)動作と呼ぶ。

20

【0042】

再投入部 331はネットワーク通信再投入動作の最初に、スナップショット採取時刻 t_2 から障害発生時刻 t_{23} までの期間にクライアントマシン 20から仮想マシン 11-1に送信された通信データのうち、1番目に送信された通信データを仮想マシン 11-2に送信する(ステップ S11)。

30

【0043】

仮想マシン 11-2は(通信記録ユニット 30内の)再投入部 331から送信される通信データを受け取ると、当該通信データに対応する処理を行い、当該通信データに対する応答(レスポンス)を通信記録ユニット 30に返す。通信記録ユニット 30内の状態再現部 33に含まれている再現判定部 332は、この仮想マシン 11-2からの応答を受け取ると(ステップ S12)、その応答(の通信データ)を、通信記録テーブル 32に保持されている、先に送信(投入)した通信データ(送信データ)に対する仮想マシン 11-1からの応答(の通信データ)と比較する(ステップ S13)。

【0044】

再現判定部 332は、上記両応答の比較結果から、当該両応答が一致しているかを判定する(ステップ S14)。もし、上記両応答が一致しているならば、再現判定部 332は、スナップショット採取時刻 t_2 から今回の応答までの期間における仮想マシン 11-1の動作状態が仮想マシン 11-2で正しく再現されたと判断する。この場合、再投入部 331は通信記録テーブル 32を参照して、スナップショット採取時刻 t_2 から障害発生時刻 t_{23} までの期間に仮想マシン 11-1に送信された全ての通信データについてネットワーク通信再投入動作が完了したかを判定する(ステップ S15)。

40

【0045】

もし、ネットワーク通信再投入動作(再投入)が未完了であるならば(ステップ S15)、再投入部 331はステップ S11に戻り、通信記録テーブル 32に保持されている通信データのうち、前回投入された通信データの次に仮想マシン 11-1に送信された通信デ

50

ータを仮想マシン 1 1-2に送信する。そして再現判定部 3 3 2 は、この通信データの送信に対する仮想マシン 1 1-2からの応答に関しても、通信記録テーブル 3 2 に保持されている、当該通信データ（送信データ）に対する仮想マシン 1 1-1からの応答と比較することで、当該両応答が一致しているかを判定する（ステップ S 1 2 ~ S 1 4）。

【 0 0 4 6 】

状態再現部 3 3 は、以上の動作を、ステップ S 1 4 で一致が判定されている限り、つまり仮想マシン 1 1-1の動作状態が仮想マシン 1 1-2に再現されていると判定されている限り繰り返す。やがて、再投入完了が判定されたものとする（ステップ S 1 5）。このことは、スナップショット採取時刻 t 2 から障害発生時刻 t 2 3までの仮想マシン 1 1-1の状態が時系列順に仮想マシン 1 1-2で再現されたことを表す。

10

【 0 0 4 7 】

明らかのように、再投入完了判定時点の仮想マシン 1 1-2の状態は、サーバ計算機 1 0-1の障害発生時における仮想マシン 1 1-1の状態 # 3 ' に一致している。つまり、再投入完了判定時点の仮想マシン 1 1-2は、サーバ計算機 1 0-1の障害発生時における仮想マシン 1 1-1の状態 # 3 ' に復元されている。

【 0 0 4 8 】

そこで再現判定部 3 3 2 は、再投入部 3 3 1 によって再投入完了が判定されると（ステップ S 1 5）、仮想マシン 1 1-2の再生成に成功したものと、その旨を、仮想マシン 1 1-2を管理する、サーバ計算機 1 0-2上の仮想マシンモニタ 1 2-2に通知する（ステップ S 1 6）。これにより状態再現部 3 3 における状態再現処理は終了する。このとき、仮想マシン 1 1-2は、サーバ計算機 1 0-1の障害発生時に仮想マシン 1 1-1と接続されていたクライアントマシン 2 0 の内部状態と整合性が取れた状態となっている。このため仮想マシン 1 1-2は、サーバ計算機 1 0-1の障害発生時 t 2 3における仮想マシン 1 1-1と同一の状態 # 3 ' でクライアントマシン 2 0 に対するサービスを継続することができる。

20

【 0 0 4 9 】

ここで、通信記録テーブル 3 2 に基づく上述の状態再現処理（ステップ S 3）の期間、クライアントマシン 2 0 から仮想マシン 1 1-2に対してアクセスがあったものとする。もし、このアクセスに対して仮想マシン 1 1-2が何らかの処理を行うならば、当該仮想マシン 1 1-2にサーバ計算機 1 0-1の障害発生時の状態を再現させることは困難となる。そこで通信記録ユニット 3 0 内の通信ブロッキング部 3 4 は、状態再現部 3 3 による状態再現処理（ステップ S 3）の期間、クライアントマシン 2 0 からの仮想マシン 1 1-2へのアクセスをブロックするための通信ブロッキング処理を実行する。

30

【 0 0 5 0 】

以下、通信ブロッキング部 3 4 によって実行される通信ブロッキング処理の手順について説明する。まず通信ブロッキング部 3 4 は、状態再現処理が開始されると（ステップ S 2 1）、クライアントマシン 2 0 から仮想マシン 1 1-2へのアクセス（通信）を監視して、そのアクセス（通信）を全てブロックする（ステップ S 2 2）。

【 0 0 5 1 】

やがて状態再現処理が終了すると（ステップ S 2 1）、通信ブロッキング部 3 4 は、通信ブロッキング状態を解除して、仮想マシン 1 1-2へのアクセス（通信）を通過させる（ステップ S 2 3）。これにより仮想マシン 1 1-1を利用するクライアントマシン 2 0 は、当該仮想マシン 1 1-1が別のサーバ計算機（ここではサーバ計算機 1 0-2）で仮想マシン 1 1-2として再現されたことを認識することなく、サービスを継続して利用できる。

40

【 0 0 5 2 】

次に、仮想マシン 1 1-2に対して投入した通信データに対する当該仮想マシン 1 1-2からの応答が、通信記録テーブル 3 2 に保持されている、当該投入した通信データに対する仮想マシン 1 1-1からの応答と異なっている場合（ステップ S 1 4）について説明する。

【 0 0 5 3 】

状態再現部 3 3（内の再現判定部 3 3 2）は、上述の両応答が一致していない場合（ステップ S 1 4）、サーバ計算機 1 0-1の障害発生時の仮想マシン 1 1-1の状態が仮想マシ

50

ン 1 1 -2 に再現されておらず、したがって当該仮想マシン 1 1 -2 によるサービスの継続はできないと判断する。そこで状態再現部 3 3 は、当該仮想マシン 1 1 -2 の再生成の失敗を、仮想マシン 1 1 -2 を管理する、サーバ計算機 1 0 -2 上の仮想マシンモニタ 1 2 -2 に通知して (ステップ S 1 7)、状態再現処理を終了する。

【 0 0 5 4 】

仮想マシンモニタ 1 2 -2 は、状態再現部 3 3 からの通知によって仮想マシン 1 1 -2 の再生成の失敗を判定すると (ステップ S 4)、当該仮想マシン 1 1 -2 を破棄する (ステップ S 5)。そして仮想マシンモニタ 1 2 -2 は、スナップショット 1 1 2 に基づき仮想マシンを再生成するのではなく、障害が発生したサーバ計算機 1 0 -1 上で動作していた仮想マシン 1 1 -1 を、サーバ計算機 1 0 -2 上で再起動する。(ステップ S 6)。このステップ S 6 における再起動処理は、仮想マシンモニタ 1 2 -2 が当該仮想マシンモニタ 1 2 -2 内の再起動部 1 2 0 -2 に指示することにより、当該再起動部 1 2 0 -2 によって行われる。再起動部 1 2 0 -2 は、仮想マシン 1 1 が使用していた仮想ディスク 1 1 0 の内容に基づき、当該仮想マシン 1 1 -1 で動作していた OS (ゲスト OS) をブートすることにより、当該仮想マシン 1 1 -1 を仮想マシン 1 1 -2 上で再起動する。

10

【 0 0 5 5 】

仮想ディスク 1 1 0 の内容は、時刻 t 2 におけるスナップショット 1 1 2 (# c) の採取後も障害発生時刻 t 2 3 まで、仮想マシン 1 1 -1 の動作により更新されている。明らかなように、仮想マシン 1 1 -1 の再起動に用いられる仮想ディスク 1 1 0 の内容は、障害発生時刻 t 2 3 における内容 # c' (図 4 参照) である。つまり図 4 の例では、仮想マシンモニタ 1 2 -2 から 1 つのファイルとして認識される仮想ディスク 1 1 0 の内容 # c' に基づき、仮想マシン 1 1 -1 がサーバ計算機 1 0 -2 上で再起動される。

20

【 0 0 5 6 】

この再起動された仮想マシン 1 1 -1 の状態は、障害発生時における当該仮想マシン 1 1 -1 の状態 # 3' とは必ずしも一致せず、図 4 に示すように例えば状態 # 3'' である。

【 0 0 5 7 】

さて、仮想マシン 1 1 -1 が使用していた仮想ディスク 1 1 0 に基づき当該仮想マシン 1 1 -1 が再起動されると、当該仮想マシン 1 1 -1 で実行されていた、サービスを提供するためのアプリケーション類も、ブートされた OS によって再起動される。このとき仮想ディスク 1 1 0 には、障害発生時刻 t 2 3 における内容 # c' が残されている。このため、再起動された仮想マシン 1 1 -1 でアプリケーションが再起動されると、当該アプリケーションの持つ障害回復機能が働く (ステップ S 7)。この障害回復機能としては、例えばデータベースのロールバック機能が知られている。

30

【 0 0 5 8 】

このように、再起動された仮想マシン 1 1 -1 でアプリケーションの障害回復機能 (による障害回復処理) が実行されると、サーバ計算機 1 0 -1 の障害発生時刻 t 2 3 まで仮想マシン 1 1 -1 からサービスの提供を受けていたクライアントマシン 2 0 も、障害発生を検知して障害回復処理 (リカバリ処理) を行う。これにより、再起動された仮想マシン 1 1 -1 で再起動されたアプリケーション (サービスを提供するためのアプリケーション) とクライアントマシン 2 0 との間の整合が取られ、再起動に成功する。この結果、クライアントマシン 2 0 に対するサービスを、再起動された仮想マシン 1 1 -1 によって継続することができる。但し、スナップショット 1 1 2 に基づいて仮想マシンを再生成して障害発生時刻 t 2 3 における仮想マシン 1 1 -1 の状態 # 3' を再現するのに比べ、OS ブート処理及びアプリケーションの回復処理のために、図 4 に示すように、サービスの再開までに要する時間が長くなる可能性が高い。

40

【 0 0 5 9 】

上述の本実施形態の動作を以下に整理する。本実施形態においてはまず、スナップショット 1 1 2 (# c) 及び通信記録テーブル 3 2 に基づき障害発生時刻 t 2 3 における仮想マシン 1 1 -1 の状態 # 3' の再現がサーバ計算機 1 0 -2 上で試みられる (ステップ S 1, S 2)。もし、仮想マシン 1 1 -1 の状態 # 3' の再現に失敗した場合 (ステップ S 4)、

50

仮想ディスク 110 の内容 # c ' に基づき当該仮想マシン 11-1 がサーバ計算機 10-2 上で再起動される (ステップ S5, S6)。この仮想マシン 11-1 の再起動に成功すると当該仮想マシン 11-1 は状態 3 " となり、この状態 3 " からサービスが再開される。

【0060】

上述の説明では、簡略化のために、1つのクライアントマシン 20 が仮想マシン 11-1 (によって提供されるサービス) を利用するものとしている。しかし、ネットワーク 2 にクライアントマシン 20 を含む複数のクライアントマシンが接続されていて、当該複数のクライアントマシンが当該仮想マシン 11-1 を利用する構成であっても構わない。

【0061】

このような構成では、通信記録ユニット 30 内の通信記録部 31 は、仮想マシン 11-1 に対応付けられた通信記録テーブル 32 に、当該仮想マシン 11-1 と当該仮想マシン 11-1 を利用する全てのクライアントマシンとの間の通信の履歴を記録すれば良い。通信記録ユニット 30 内の状態再現部 33 は、上記実施形態と同様に、通信記録テーブル 32 に記録された通信の履歴に基づきネットワーク通信再投入動作を行う。これにより状態再現部 33 は、複数のクライアントマシンに対してサービスを提供していた期間のうちのスナップショット時刻 t2 から障害発生時刻 t23 までの仮想マシン 11-1 の状態を、再生成された仮想マシン 11-2 に時系列順に再現することが可能となる。つまり、複数のクライアントマシンが仮想マシンを利用する構成においても、上記実施形態と同様の手順で当該複数のクライアントマシンに対するサービスを継続することができる。

10

【0062】

ここで、通信記録テーブル 32 には、通信時刻、通信の方向及び通信データに加えて、当該通信データを送受信するクライアントマシンの識別情報を記録すると良い。このようにすると、通信ブロッキング部 34 は、状態再現部 33 による状態再現処理 (ステップ S3) の期間、通信記録テーブル 32 に記録されている識別情報の示す複数のクライアントマシンから再生成された仮想マシン 11-2 へのアクセスを全てブロックすることができる。

20

【0063】

[第1の変形例]

次に上記実施形態の第1の変形例について説明する。

図8は、上記実施形態の第1の変形例に係る仮想計算機システムの構成を示すブロック図である。図8において、図1と同様の要素には同一参照番号を付してある。

30

【0064】

図8のシステムが図1のそれと相違するのは、サーバ計算機 10-1 及び 10-2 とクライアントマシン 20 (を含む複数のクライアントマシン) とが同一のネットワーク、例えばネットワーク 1 に接続されている点と、通信記録ユニット 30 が (ルータ 3 ではなくて)、プロキシサーバ (プロキシサーバ計算機) 300 に設けられている点とにある。

【0065】

図8のシステムにおいて、クライアントマシン 20 (を含む複数のクライアントマシン) とサーバ計算機 10-1 及び 10-2 上でそれぞれ動作する仮想マシンとの間の通信は必ずプロキシサーバ 300 を介して行われる。つまり、クライアントマシン 20 (を含む複数のクライアントマシン) からサーバ計算機 10-i (i = 1, 2) 上の仮想マシンへのアクセスは、当該クライアントマシン 20 (を含む複数のクライアントマシン) がプロキシサーバ 300 に接続することにより、当該プロキシサーバ 300 によって代理で行われる。したがって、通信記録ユニット 30 がプロキシサーバ 300 に設けられる第1の変形例では、当該通信記録ユニット 30 内の通信ブロッキング部 34 は、通信記録テーブル 32 に基づく上述の状態再現処理 (ステップ S3) の期間、プロキシサーバ 300 に接続されるクライアントマシン 20 (を含む複数のクライアントマシン) から再生成された仮想マシンへのアクセスを全てブロックすることができる。

40

【0066】

[変形例2]

50

次に、上記実施形態の第2の変形例について説明する。この第2の変形例の特徴は、サーバ計算機の障害発生に伴って実行されるサービス引き継ぎ処理（障害回復処理）において、仮想ディスクからの仮想マシン再起動を、再生成された仮想マシン上での状態再現処理の失敗を見越して、例えば当該状態再現処理と並行して投機的に行う点にある。

【0067】

以下、第2の変形例の動作について、当該動作が図1の仮想計算機システムで行われるものとして、図9のフローチャートを参照して説明する。図9はサーバ計算機障害発生時に仮想マシンモニタによって実行されるサービス引き継ぎ処理の手順を示すフローチャートである。

【0068】

今、仮想マシン11-1が動作するサーバ計算機10-1で障害が発生したことが、サーバ計算機10-2上で動作する仮想マシンモニタ12-2によって検出されたものとする。すると仮想マシンモニタ12-2は、上記実施形態のステップS1と同様に、障害発生時刻に最も近い時刻で採取されたスナップショット112を用いて、サーバ計算機10-2上に当該仮想マシンモニタ12-1が採取された時刻における仮想マシン11-1と全く同じ状態の仮想マシンを仮想マシン11-2として再生成する（ステップS31）。

【0069】

次に仮想マシンモニタ12-2は、上記実施形態のステップS2と同様に、通信記録ユニット30に対して状態再現処理を要求する（ステップS32）。同時に仮想マシンモニタ12-2は、再起動部120-2に対して再起動処理を指示する（ステップS33）。

【0070】

すると通信記録ユニット30内の状態再現部33は、上記実施形態と同様の手順（図6のフローチャート参照）でサーバ計算機10-1の障害発生時における仮想マシン11-1の状態を仮想マシン11-2に再現させるための状態再現処理を実行する。一方、再起動部120-2は、上記実施形態と同様に、障害が発生したサーバ計算機10-1上で動作していた仮想マシン11-1によって使用されていた仮想ディスク110の内容に基づいて、当該仮想マシン11-1をサーバ計算機10-2上で再起動するための再起動処理を実行する。

【0071】

その後、仮想マシンモニタ12-2は再起動部120-2によってサーバ計算機10-2上に再起動される仮想マシン11-1の状態を監視することにより、当該仮想マシン11-1によるサービスの再開が可能となったかを判定する（ステップS34）。ここでは、再起動された仮想マシン11-1でアプリケーションの障害回復処理が実行されると共に、クライアントマシン20で障害回復処理が実行された結果、当該仮想マシン11-1で再起動されたアプリケーションとクライアントマシン20との間の整合が取られた時点で、サービスの再開が可能となる。つまり、仮想マシン11-1の再起動に成功する。

【0072】

もし、再起動された仮想マシン11-1によるサービスの再開が可能な状態には未だなっていないならば、仮想マシンモニタ12-2は、通信記録ユニット30内の状態再現部33（に含まれている再現判定部332）から状態再現成功または失敗の通知（状態再現成功/失敗通知）が送られているかを判定する（ステップS35）。もし、状態再現成功/失敗通知が未だ送られていないならば、仮想マシンモニタ12-2は、再びステップS34の判定を行う。

【0073】

やがて、再起動された仮想マシン11-1によるサービスの再開が可能な状態となる前に（ステップS34）、状態再現部33から状態再現成功/失敗通知が送られたものとする（ステップS35）。この場合、仮想マシンモニタ12-2は、この通知が成功通知であるかを判定する（ステップS36）。もし、成功通知ならば、仮想マシンモニタ12-2は、再生成された仮想マシン11-2でサービスを継続させる（ステップS37）。このとき仮想マシンモニタ12-2は、再起動された仮想マシン11-1を停止または破棄する。

【0074】

10

20

30

40

50

これに対し、失敗通知ならば（ステップS36）、仮想マシンモニタ12-2は再生成された仮想マシン11-2を破棄する（ステップS38）。そして仮想マシンモニタ12-2は、再起動された仮想マシン11-1によるサービスの再開が可能な状態になるのを待つ（ステップS39）。やがて、再起動された仮想マシン11-1によるサービスの再開が可能な状態になると（ステップS39）、仮想マシンモニタ12-2は当該仮想マシン11-1でサービスを継続させる（ステップS40）。

【0075】

第2の変形例において、仮想マシン11-1をサーバ計算機10-2上で再起動するための再起動処理は、上記実施形態と異なって、サーバ計算機10-1の障害発生時における仮想マシン11-1の状態を仮想マシン11-2に再現させるための状態再現処理と並行して行われる。したがって第2の変形例においては、状態再現処理に失敗しても、再起動された仮想マシン11-1で速やかにサービスを継続することができるため、上記実施形態と比較してサービスの再開までの時間が短縮できる。

10

【0076】

一方、状態再現部33から状態再現成功/失敗通知が送られる前に（ステップS35）、再起動された仮想マシン11-1によるサービスの再開が可能な状態になったならば（ステップS34）、仮想マシンモニタ12-2は再生成された仮想マシン11-2を破棄する（ステップS41）。そして仮想マシンモニタ12-2は、再起動された仮想マシン11-1でサービスを継続させる（ステップS40）。このように第2の変形例においては、万が一、状態再現部33から状態再現成功/失敗が通知されるよりも先に、再起動された仮想マシン11-1によるサービスの再開が可能な状態になったならば、当該仮想マシン11-1でサービスが継続される。これにより上記実施形態と比較してサービスの再開までの時間が短縮できる。

20

【0077】

なお、本発明は、上記実施形態またはその変形例そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記実施形態またはその変形例に開示されている複数の構成要素の適宜な組み合わせにより種々の発明を形成できる。例えば、実施形態またはその変形例に示される全構成要素から幾つかの構成要素を削除してもよい。

30

【図面の簡単な説明】

【0078】

【図1】本発明の一実施形態に係る仮想計算機システムの構成を示すブロック図。

【図2】通信記録テーブルに記録された通信の履歴の例を示す図。

【図3】サーバ計算機障害発生時の仮想マシン及び仮想ディスクの状態と、当該障害発生時まで採取されたスナップショットの状態とを示す図。

【図4】時系列に沿った仮想マシン、仮想ディスク、スナップショットの状態及び仮想マシン再生成/再起動を説明するための図、

【図5】同実施形態で適用される、サーバ計算機障害発生時のサービス引き継ぎのための手順を示すフローチャート。

【図6】同実施形態で適用される、通信記録ユニット内の状態再現部による状態再現処理の手順を示すフローチャート。

40

【図7】同実施形態で適用される、通信記録ユニット内の通信ブロッキング部による通信ブロッキング処理の手順を示すフローチャート。

【図8】同実施形態の第1の変形例に係る仮想計算機システムの構成を示すブロック図。

【図9】同実施形態の第2の変形例において、仮想マシンモニタによって実行されるサービス引き継ぎ処理の手順を示すフローチャート。

【符号の説明】

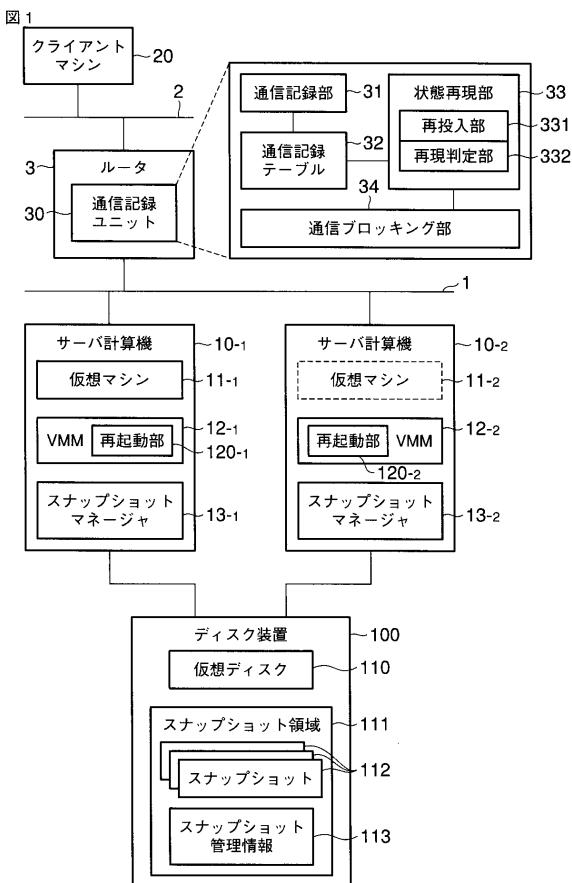
【0079】

1, 2...ネットワーク、3...ルータ、10-1, 10-2...サーバ計算機、11-1, 11-2...仮想マシン、12-1, 12-2...仮想マシンモニタ(VMM)、13-1, 13-2...スナッ

50

プッシュマネージャ（スナップショット管理手段）、30...通信記録ユニット、31...通信記録部、32...通信記録テーブル、33...状態再現部、34...通信ブロッキング部、100...ディスク装置、110...仮想ディスク、111...スナップショット領域、112...スナップショット、113...スナップショット管理情報、331...再投入部、332...再現判定部、300...プロキシサーバ。

【図1】

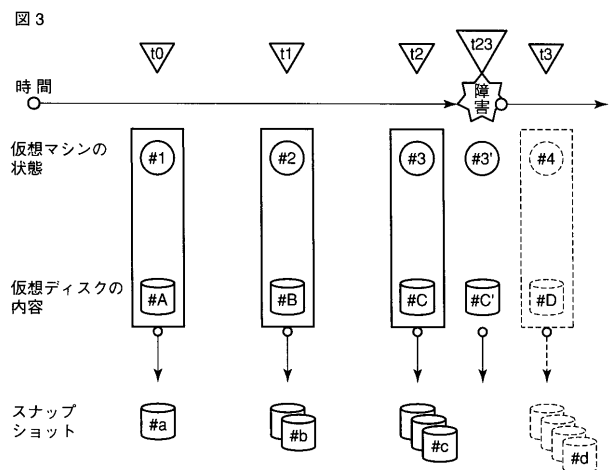


【図2】

通信記録テーブル 32

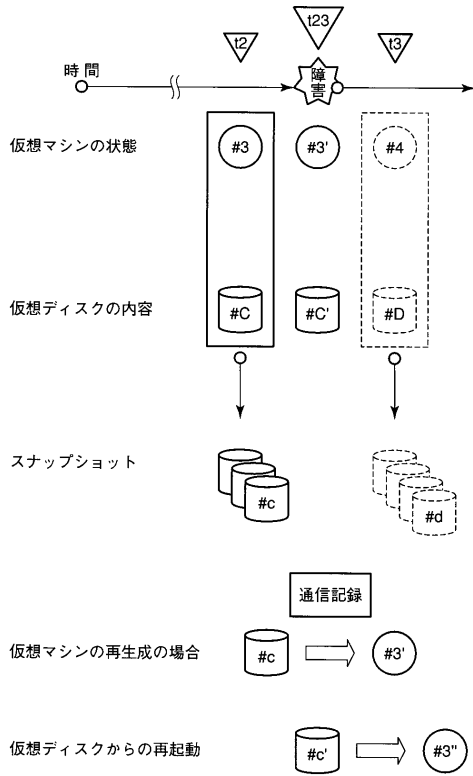
#	時刻	方向	データ
1	00:00:00	IN	AAAAA
2	00:00:02	IN	BBBBB
3	00:00:05	OUT	aaabbb
4	00:01:00	IN	CCCccc
5	00:05:00	OUT	DDDDDD
...

【図3】



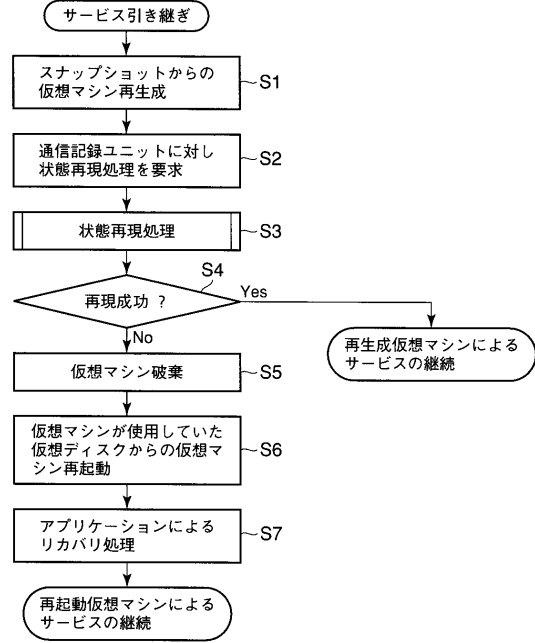
【 図 4 】

図 4



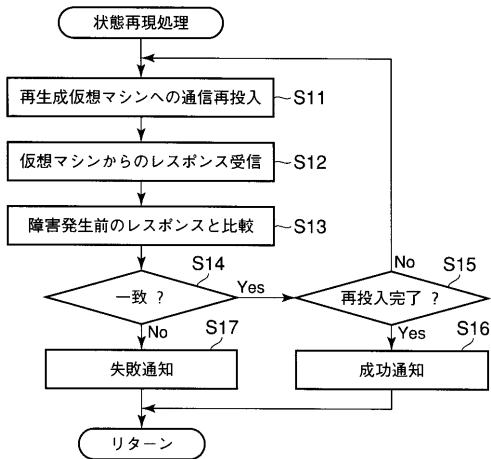
【 図 5 】

図 5



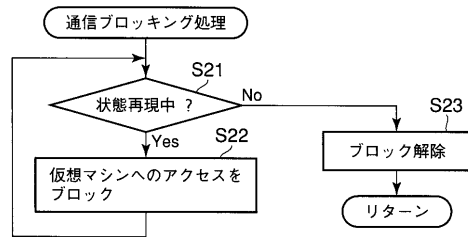
【 図 6 】

図 6

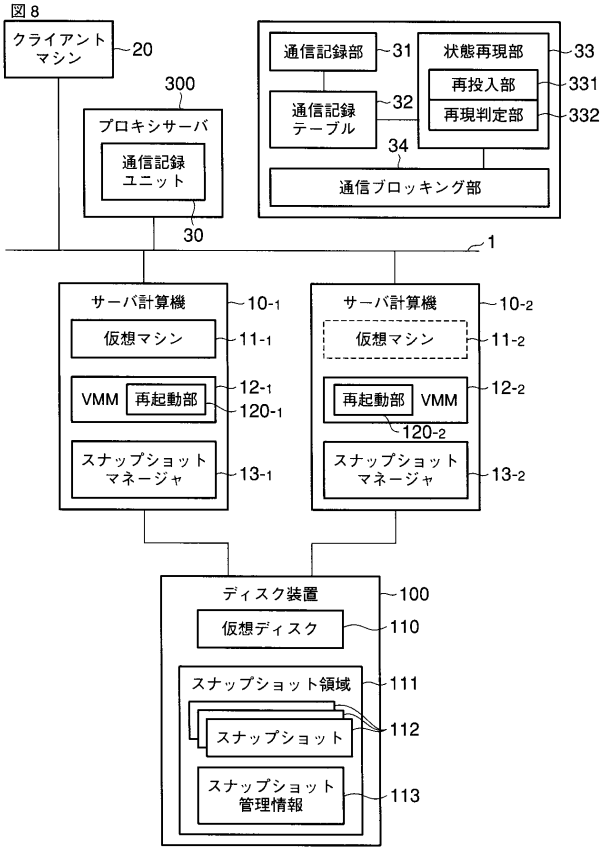


【 図 7 】

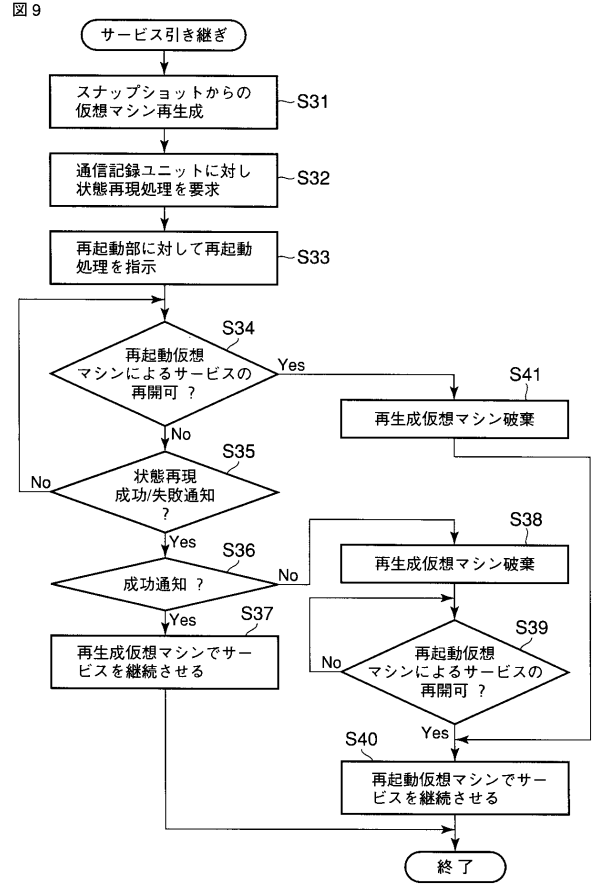
図 7



【 図 8 】



【 図 9 】



フロントページの続き

(74)代理人 100075672

弁理士 峰 隆司

(74)代理人 100109830

弁理士 福原 淑弘

(74)代理人 100084618

弁理士 村松 貞男

(74)代理人 100092196

弁理士 橋本 良郎

(72)発明者 飯沼 哲也

東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

Fターム(参考) 5B034 BB01 CC01 CC02 CC05 DD01 DD05