

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6421470号
(P6421470)

(45) 発行日 平成30年11月14日(2018.11.14)

(24) 登録日 平成30年10月26日(2018.10.26)

(51) Int.Cl.		F I			
G06F	9/455	(2006.01)	G06F	9/455	150
G06F	9/50	(2006.01)	G06F	9/50	120Z
G06F	9/44	(2018.01)	G06F	9/44	

請求項の数 7 (全 31 頁)

(21) 出願番号	特願2014-124545 (P2014-124545)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成26年6月17日(2014.6.17)		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2016-4432 (P2016-4432A)	(74) 代理人	100092152 弁理士 服部 毅巖
(43) 公開日	平成28年1月12日(2016.1.12)	(72) 発明者	長谷川 裕毅 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
審査請求日	平成29年3月9日(2017.3.9)	(72) 発明者	治部 将之 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	元藤 雄路 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】 仮想マシンマイグレーションプログラム、仮想マシンマイグレーションシステムおよび仮想マシンマイグレーション方法

(57) 【特許請求の範囲】

【請求項1】

コンピュータに、

第1のオペレーティングシステムを含む第1の仮想マシンが配置された他のコンピュータから、前記コンピュータ上に起動された第2の仮想マシンに含まれる第2のオペレーティングシステムを実行することで前記第2のオペレーティングシステムにより前記第1のオペレーティングシステムのデータを取得し、

前記第2のオペレーティングシステムを実行することで前記第2のオペレーティングシステムにより、前記第2の仮想マシンを再起動した際に実行されるオペレーティングシステムが、前記第2のオペレーティングシステムから、前記取得した第1のオペレーティングシステムのデータに基づいて実行される前記第1のオペレーティングシステムに変わるように、前記第2の仮想マシンの起動方法を示す起動情報を書き換え、

前記第2の仮想マシンを再起動する、

処理を実行させる仮想マシンマイグレーションプログラム。

【請求項2】

前記第2の仮想マシンには、前記コンピュータが備える記憶装置の中から、第1の記憶領域と前記第2のオペレーティングシステムのデータを記憶する第2の記憶領域と前記起動情報を記憶する第3の記憶領域とが割り当てられ、

前記第1のオペレーティングシステムのデータは前記第1の記憶領域に書き込まれ、

前記起動情報は、前記第2の仮想マシンの再起動の際に前記第1の記憶領域に記憶され

たデータが読み込まれるように書き換えられる、

請求項 1 記載の仮想マシンマイグレーションプログラム。

【請求項 3】

前記第 2 の仮想マシンへの前記第 2 の記憶領域の割り当ては、前記第 2 の仮想マシンで前記第 1 のオペレーティングシステムが実行された後に解除される、

請求項 2 記載の仮想マシンマイグレーションプログラム。

【請求項 4】

第 1 のオペレーティングシステムを含む第 1 の仮想マシンが配置され、前記第 1 のオペレーティングシステムのデータを送信する第 1 の情報処理装置と、

第 2 のオペレーティングシステムを含む第 2 の仮想マシンを起動し、

前記第 2 のオペレーティングシステムを実行することで前記第 2 のオペレーティングシステムにより前記第 1 の情報処理装置から前記第 1 のオペレーティングシステムのデータを取得し、

前記第 2 のオペレーティングシステムを実行することで前記第 2 のオペレーティングシステムにより、前記第 2 の仮想マシンを再起動した際に実行されるオペレーティングシステムが、前記第 2 のオペレーティングシステムから、前記取得した第 1 のオペレーティングシステムのデータに基づいて実行される前記第 1 のオペレーティングシステムに変わるように、前記第 2 の仮想マシンの起動方法を示す起動情報を書き換え、

前記第 2 の仮想マシンを再起動する、第 2 の情報処理装置と、

を有する仮想マシンマイグレーションシステム。

【請求項 5】

前記第 1 の情報処理装置は、前記第 1 のオペレーティングシステムを一時的にサスペンド状態に遷移させることで、メモリに記憶されたデータおよびプロセッサ内のレジスタに記憶されたデータの少なくとも一方を抽出する、

請求項 4 記載の仮想マシンマイグレーションシステム。

【請求項 6】

前記第 1 の情報処理装置は、前記第 1 の情報処理装置における単位時間当たりのデータ更新量が閾値を超えている場合、前記第 1 の情報処理装置で実行されている 1 または 2 以上のプロセスに対してプロセッサのリソースの割り当てを制限する、

請求項 4 または 5 記載の仮想マシンマイグレーションシステム。

【請求項 7】

第 1 の情報処理装置と第 2 の情報処理装置とを有するシステムが実行する仮想マシンマイグレーション方法であって、

前記第 1 の情報処理装置に配置された第 1 のオペレーティングシステムを含む第 1 の仮想マシンと異なる、第 2 のオペレーティングシステムを含む第 2 の仮想マシンを、前記第 2 の情報処理装置上に起動し、

前記第 1 のオペレーティングシステムのデータを、前記第 1 の情報処理装置から、前記第 2 のオペレーティングシステムを実行することで前記第 2 のオペレーティングシステムにより前記第 2 の情報処理装置にコピーし、

前記第 2 のオペレーティングシステムを実行することで前記第 2 のオペレーティングシステムにより、前記第 2 の仮想マシンを再起動した際に実行されるオペレーティングシステムが、前記第 2 のオペレーティングシステムから、前記取得した第 1 のオペレーティングシステムのデータに基づいて実行される前記第 1 のオペレーティングシステムに変わるように、前記第 2 の仮想マシンの起動方法を示す起動情報を書き換え、

前記第 2 の仮想マシンを再起動する、

仮想マシンマイグレーション方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は仮想マシンマイグレーションプログラム、仮想マシンマイグレーションシステ

10

20

30

40

50

ムおよび仮想マシンマイグレーション方法に関する。

【背景技術】

【0002】

現在、プロセッサやメモリやハードディスクなどのコンピュータのリソースを保持し、ネットワーク経由でユーザにリソースを使用させる情報処理サービスが存在する。このような情報処理サービスは、クラウドサービスと呼ばれることがある。

【0003】

データセンタなどクラウドサービスを提供する施設では、仮想化技術を用いて、1つの物理的なコンピュータ（物理マシン）上に複数の仮想的なコンピュータ（仮想マシン）を形成することがある。ハイパーバイザなど仮想マシンを制御する制御ソフトウェアは、複数の仮想マシンそれぞれにリソースを割り当てる。クラウドサービスの施設は、例えば、ユーザからの要求に応じてリソースを確保し、当該ユーザ用の仮想マシンを生成する。ユーザは、生成された仮想マシン上で、オペレーティングシステム（OS：Operating System）やアプリケーションソフトウェアなどの所望のソフトウェアを実行させることができる。すなわち、ユーザは、当該ユーザ用の仮想マシンのリソースを使用することができる。仮想マシンなど情報基盤（インフラストラクチャ）そのものをユーザに提供するクラウドサービスは、IaaS（Infrastructure as a Service）と呼ばれることがある。

【0004】

仮想マシンは、ハイパーバイザなどの制御ソフトウェアによる制御のもと、ある物理マシンから他の物理マシンに移行すること（マイグレーション）ができる場合がある。特に、移行元の物理マシンにおいて仮想マシンのOSをシャットダウンせずに、メモリデータやプロセッサの状態データなどを移行先の物理マシンにコピーし、移行先の物理マシンで仮想マシンの処理を引き継ぐことができる場合がある。このようなマイグレーション方法は、ライブマイグレーションと呼ばれることがある。仮想マシンのマイグレーションは、一部の物理マシンでリソース不足が生じたときや一部の物理マシンを停止して保守作業を行うときなどに、同一のクラウドサービスの施設内で行われることがある。

【0005】

なお、ある物理マシンから他の物理マシンに仮想マシンを移行するとき、仮想マシン上で実行されるプログラムのうち実行頻度の高いプログラムのメモリデータを優先的に物理マシン間で送信する仮想マシンの移行システムが提案されている。また、物理マシン間で処理を移行するとき、移行先の物理マシンに応じてプロセスマイグレーションと仮想マシンマイグレーションとを選択的に実行する分散処理システムが提案されている。メモリデータやレジスタデータの移行は、プロセスマイグレーションではプロセス単位で行われ、仮想マシンマイグレーションでは仮想マシン単位で行われる。

【0006】

また、物理マシン間でゲストOSを移行する情報処理システムが提案されている。この情報処理システムでは、移行元の物理マシンが参照するストレージと移行先の物理マシンが参照するストレージに、同じゲストOSのイメージデータを格納しておく。また、移行元の物理マシンは、ゲストOSのイメージデータへの書き込みを禁止し、書き込みデータをメモリ上にキャッシュしておく。ゲストOSの移行時には、移行元の物理マシンがメモリデータを移行先の物理マシンにコピーし、移行先の物理マシンは移行先のストレージに格納されたイメージデータと受信したメモリデータとを用いてゲストOSを起動する。

【0007】

また、ライブマイグレーションの際に、移行元の物理マシンにおける単位時間当たりのメモリへの書き込み量を、物理マシン間のデータ送信帯域よりも小さくなるように低減する仮想化システムが提案されている。この仮想化システムでは、メモリに対してデータを書き込んだ旨をプログラムに報告するときに発行される割り込みを一時的に保留して遅延させることで、メモリに対する連続的なデータ書き込みについて書き込み速度を低減する。

【先行技術文献】

10

20

30

40

50

【特許文献】

【0008】

【特許文献1】国際公開第2009/069573号

【特許文献2】国際公開第2010/035480号

【特許文献3】特開2011-210151号公報

【特許文献4】特開2013-250950号公報

【発明の概要】

【発明が解決しようとする課題】

【0009】

ところで、ユーザは、互いに離れた複数の情報処理施設（例えば、異なるサービス事業者がもつ複数のデータセンタ）を使い分けることがある。このとき、ある情報処理施設に配置された仮想マシンを他の情報処理施設に移行したいことが考えられる。例えば、ある仮想マシンの負荷が高くなったとき、リソースを豊富にもつ情報処理施設に当該仮想マシンを移行したいことが考えられる。また、例えば、ある仮想マシンを、課金の少ない情報処理施設やセキュリティの高い情報処理施設に移行したいことなども考えられる。

10

【0010】

一方で、情報処理施設には、ユーザの仮想マシン間のデータ通信またはユーザの仮想マシンと外部ネットワークとの間のデータ通信に用いられるユーザネットワークとは別に、仮想マシンの制御に用いられる管理ネットワークが設けられることがある。この場合、同一の情報処理施設内でのマイグレーションには、管理ネットワークが使用される。

20

【0011】

しかし、セキュリティなどの観点から、複数の情報処理施設の間で管理ネットワークを相互接続する（互いに相手の管理ネットワークにアクセス可能にする）ことは容易でない。特に、異なるサービス事業者がもつデータセンタの間で管理ネットワークを相互接続することは難しい。このため、同一の情報処理施設内でのマイグレーション方法を、異なる情報処理施設の間でのマイグレーションに適用することは難しいという問題がある。その結果、異なる情報処理施設の間では仮想マシンを円滑に移行することが難しかった。

【0012】

1つの側面では、本発明は、異なる情報処理施設の間で仮想マシンを円滑に移行できるようにする仮想マシンマイグレーションプログラム、仮想マシンマイグレーションシステムおよび仮想マシンマイグレーション方法を提供することを目的とする。

30

【課題を解決するための手段】

【0013】

1つの態様では、コンピュータに以下の処理を実行させる仮想マシンマイグレーションプログラムが提供される。第1のオペレーティングシステムを含む第1の仮想マシンが配置された他のコンピュータから、コンピュータ上に起動された第2の仮想マシンに含まれる第2のオペレーティングシステムを用いて第1のオペレーティングシステムのデータを取得する。第2の仮想マシンの再起動の際に取得した第1のオペレーティングシステムのデータに基づいて第1のオペレーティングシステムが実行されるように、第2の仮想マシンの起動方法を示す起動情報を書き換える。第2の仮想マシンを再起動する。

40

【0014】

また、1つの態様では、第1の情報処理装置と第2の情報処理装置とを有する仮想マシンマイグレーションシステムが提供される。

また、1つの態様では、第1の情報処理装置と第2の情報処理装置とを有するシステムが実行する仮想マシンマイグレーション方法が提供される。

【発明の効果】

【0015】

1つの側面では、異なる情報処理施設の間で仮想マシンを円滑に移行できる。

【図面の簡単な説明】

【0016】

50

【図 1】第 1 の実施の形態の仮想マシンマイグレーションシステムを示す図である。

【図 2】第 2 の実施の形態の情報処理システムを示す図である。

【図 3】ホストサーバのハードウェア例を示すブロック図である。

【図 4】仮想マシンとハイパーバイザの配置例を示すブロック図である。

【図 5】仮想マシンのマイグレーション例を示す第 1 の図である。

【図 6】仮想マシンのマイグレーション例を示す第 2 の図である。

【図 7】仮想マシンのマイグレーション例を示す第 3 の図である。

【図 8】仮想マシンのマイグレーション例を示す第 4 の図である。

【図 9】仮想マシンの機能例を示すブロック図である。

【図 10】状態テーブルの例を示す図である。

10

【図 11】ディスク更新情報とメモリ更新情報の例を示す図である。

【図 12】プロセステーブルの例を示す図である。

【図 13】リソース制限テーブルと重要プロセスリストの例を示す図である。

【図 14】仮想マシン生成要求の例を示す図である。

【図 15】マイグレーション開始の手順例を示すフローチャートである。

【図 16】ディスクデータ送信の手順例を示すフローチャートである。

【図 17】ディスク監視の手順例を示すフローチャートである。

【図 18】メモリデータ送信の手順例を示すフローチャートである。

【図 19】メモリ監視の手順例を示すフローチャートである。

【図 20】プロセス監視の手順例を示すフローチャートである。

20

【図 21】プロセス監視の手順例を示すフローチャート（続き）である。

【図 22】マイグレーション完了の手順例を示すフローチャートである。

【図 23】マイグレーション完了の手順例を示すフローチャート（続き）である。

【発明を実施するための形態】

【0017】

以下、本実施の形態を図面を参照して説明する。

[第 1 の実施の形態]

図 1 は、第 1 の実施の形態の仮想マシンマイグレーションシステムを示す図である。

【0018】

第 1 の実施の形態の仮想マシンマイグレーションシステムは、情報処理装置 1 , 2 を有する。情報処理装置 1 , 2 は、1 または 2 以上の仮想マシンを配置することができるコンピュータであり、例えば、ユーザに対して IaaS などのクラウドサービスを提供するサーバコンピュータである。情報処理装置 1 , 2 は、異なる情報処理施設（例えば、異なるデータセンタ）に設置される。情報処理装置 1 が属する情報処理施設と情報処理装置 2 が属する情報処理施設とは、異なるサービス事業者によって管理されてもよい。情報処理装置 1 , 2 は、インターネットなどの広域ネットワークを介して互いに通信できる。

30

【0019】

情報処理装置 1 には、仮想マシン 1a が配置される。仮想マシン 1a には、情報処理装置 1 が備えるプロセッサやメモリなどの情報処理のリソースが割り当てられる。仮想マシン 1a は、オペレーティングシステム 1b を含む。オペレーティングシステム 1b は、仮想マシン 1a に割り当てられたリソースを用いて実行される。

40

【0020】

上記のプロセッサは、CPU (Central Processing Unit) や DSP (Digital Signal Processor) であってもよい。また、プロセッサは、ASIC (Application Specific Integrated Circuit) や FPGA (Field Programmable Gate Array) などの特定用途の電子回路を含んでもよい。プロセッサは、例えば、メモリに記憶されたプログラムを実行する。複数のプロセッサの集合（マルチプロセッサ）を「プロセッサ」と呼ぶこともある。メモリは、例えば、RAM (Random Access Memory) などの半導体メモリである。

【0021】

また、情報処理装置 1 は、記憶部 1c を有する。記憶部 1c は、例えば、HDD (Hard

50

Disk Drive) やフラッシュメモリなど不揮発性の記憶装置である。記憶部 1 c は、記憶領域 1 d を含む。記憶領域 1 d は、仮想マシン 1 a に割り当てられたリソースであり、オペレーティングシステム 1 b のデータを記憶する。オペレーティングシステム 1 b のデータには、オペレーティングシステム 1 b の処理を定義したプログラムや、設定ファイルなどのプログラム以外のデータが含まれてもよい。情報処理装置 2 は、記憶部 2 c を有する。記憶部 2 c は、例えば、HDD やフラッシュメモリなど不揮発性の記憶装置である。

【0022】

ここで、情報処理装置 1 から情報処理装置 2 への仮想マシン 1 a のマイグレーションを考える。同一の情報処理施設の中でのマイグレーションでは、ユーザネットワークとは別個に設けられた管理ネットワークを介して、オペレーティングシステムのデータが転送されることがある。ユーザネットワークは、ユーザの仮想マシン間またはユーザの仮想マシンと情報処理施設の外との間のデータ通信に用いられるネットワークである。管理ネットワークは、仮想マシンの制御に用いられるネットワークである。

10

【0023】

しかし、異なる情報処理施設の間では、セキュリティなどの観点から、管理ネットワークへの直接アクセス(相互接続)が制限されることが想定される。このため、情報処理装置 1, 2 の間のマイグレーションは、管理ネットワークを用いて行うことが難しい場合がある。そこで、第 1 の実施の形態では、管理ネットワークを用いずに(例えば、ユーザネットワークを用いて)仮想マシン 1 a のマイグレーションを可能にする。

【0024】

仮想マシン 1 a のマイグレーションでは、情報処理装置 2 は、仮想マシン 2 a を起動する。仮想マシン 2 a は、オペレーティングシステム 2 b を含む。オペレーティングシステム 2 b は、マイグレーション中に一時的に用いられるソフトウェアであり、仮想マシン 2 a に対して割り当てられた情報処理装置 2 のプロセッサやメモリを用いて実行される。

20

【0025】

また、仮想マシン 2 a には、記憶部 2 c に含まれる記憶領域 2 d, 2 e, 2 f が割り当てられる。記憶領域 2 d は、空の記憶領域である。記憶領域 2 e は、オペレーティングシステム 2 b のデータを記憶する。記憶領域 2 f は、起動情報を記憶する。起動情報は、仮想マシン 2 a の起動方法を示す。例えば、仮想マシン 2 a が最初に起動されるときに起動情報は、起動時に記憶領域 2 e に記憶されたデータを読み込むことを示している。これにより、仮想マシン 2 a の起動時にオペレーティングシステム 2 b が実行される(S1)。

30

【0026】

仮想マシン 2 a が起動されると、情報処理装置 1 から情報処理装置 2 にオペレーティングシステム 1 b のデータがコピーされる。データ送信では、例えば、広域ネットワークおよび各情報処理施設内のユーザネットワークが用いられ、各情報処理施設内の管理ネットワークは迂回される。情報処理装置 1 は、記憶領域 1 d に記憶されたデータを情報処理装置 2 に送信する。情報処理装置 2 は、仮想マシン 2 a のオペレーティングシステム 2 b を用いて、情報処理装置 1 からオペレーティングシステム 1 b のデータを取得する。取得されたオペレーティングシステム 1 b のデータは、記憶領域 2 d に書き込まれる(S2)。

【0027】

オペレーティングシステム 1 b のデータのコピーが完了すると、情報処理装置 2 は、オペレーティングシステム 2 b を用いて、記憶領域 2 f に記憶された起動情報を書き換える(S3)。書き換えられた起動情報は、仮想マシン 2 a の再起動の際に、記憶領域 2 e に記憶されたデータに代えて記憶領域 2 d に記憶されたデータを読み込むことを示している。これにより、仮想マシン 2 a の再起動時に、オペレーティングシステム 2 b に代えてオペレーティングシステム 1 b が実行されることになる。

40

【0028】

起動情報が書き換えられると、情報処理装置 2 は、仮想マシン 2 a を再起動する。例えば、情報処理装置 2 は、仮想マシン 2 a のオペレーティングシステム 2 b をシャットダウンさせる。情報処理装置 2 は、記憶領域 2 f に記憶された起動情報を読み込み、読み込ん

50

だ起動情報に従って記憶領域 2 d に記憶されたデータを読み込む。これにより、仮想マシン 2 a ではオペレーティングシステム 1 b が実行される (S 4)。このとき、仮想マシン 1 a を停止させ、情報処理装置 1 から仮想マシン 1 a を削除してもよい。また、記憶領域 2 e に記憶されたオペレーティングシステム 2 b のデータを削除してもよい。移行元の仮想マシン 1 a の処理は、移行先の仮想マシン 2 a に引き継がれる。よって、仮想マシン 1 a は、実質的に情報処理装置 1 から情報処理装置 2 に移行したことになる。

【 0 0 2 9 】

第 1 の実施の形態の仮想マシンマイグレーションシステムによれば、移行先の情報処理装置 2 で仮想マシン 2 a が起動し、仮想マシン 2 a 上でオペレーティングシステム 2 b が一時的に実行される。オペレーティングシステム 1 b のデータのコピーは、仮想マシン 1 a と仮想マシン 2 a との間の通信として行うことができる。よって、情報処理施設外部からの管理ネットワークへの直接アクセスが制限されていても、管理ネットワークを迂回しユーザネットワークを用いてオペレーティングシステム 1 b のデータのコピーできる。

10

【 0 0 3 0 】

そして、起動情報を書き換えることで、仮想マシン 2 a が再起動すると、オペレーティングシステム 2 b に代えてオペレーティングシステム 1 b が仮想マシン 2 a 上で実行される。これにより、異なる情報処理施設の間 (例えば、異なるサービス事業者のデータセンタの間) でも、仮想マシン 1 a のマイグレーションを円滑に行うことができる。なお、上記のマイグレーションの制御は、例えば、オペレーティングシステム 1 b , 2 b に組み込まれるプログラムモジュールとして実装することが可能である。

20

【 0 0 3 1 】

[第 2 の実施の形態]

図 2 は、第 2 の実施の形態の情報処理システムを示す図である。

第 2 の実施の形態の情報処理システムは、データセンタ 1 0 , 2 0、インターネット 3 0 およびクライアント 3 1 を含む。データセンタ 1 0 , 2 0 は、クラウドサービスを提供する情報処理施設である。データセンタ 1 0 とデータセンタ 2 0 とは、異なるサービス事業者によって管理されている。クライアント 3 1 のユーザは、契約により、データセンタ 1 0 , 2 0 両方のクラウドサービスを利用できる。クライアント 3 1 は、広域ネットワークであるインターネット 3 0 を介して、データセンタ 1 0 , 2 0 にアクセスする。

【 0 0 3 2 】

30

データセンタ 1 0 は、ユーザネットワーク 1 1、管理ネットワーク 1 2、管理サーバ 1 3 およびホストサーバ 1 0 0 , 1 0 0 a を有する。データセンタ 2 0 も、データセンタ 1 0 と対応するように、ユーザネットワーク 2 1、管理ネットワーク 2 2、管理サーバ 2 3 およびホストサーバ 2 0 0 , 2 0 0 a を有する。以下では、代表してデータセンタ 1 0 について説明し、データセンタ 2 0 については説明を省略する。

【 0 0 3 3 】

ユーザネットワーク 1 1 は、ホストサーバ 1 0 0 , 1 0 0 a に配置されたユーザ用の仮想マシンがデータ通信に使用するローカルネットワークである。ユーザネットワーク 1 1 は、例えば、スイッチなどの通信装置を含む有線通信ネットワークである。ユーザネットワーク 1 1 は、データセンタ 1 0 の外部にあるインターネット 3 0 と接続されている。データセンタ 2 0 に属するコンピュータやクライアント 3 1 は、インターネット 3 0 およびユーザネットワーク 1 1 を介して、ユーザ用の仮想マシンと通信できる。

40

【 0 0 3 4 】

管理ネットワーク 1 2 は、ホストサーバ 1 0 0 , 1 0 0 a に配置された管理用の仮想マシンや管理サーバ 1 3 が、ユーザ用の仮想マシンの制御に使用するローカルネットワークである。管理ネットワーク 1 2 は、例えば、スイッチなどの通信装置を含む有線通信ネットワークである。仮想マシンの制御には、ホストサーバ 1 0 0 , 1 0 0 a 上に新たな仮想マシンを生成することや、データセンタ 1 0 内のホストサーバ 1 0 0 , 1 0 0 a 間で仮想マシンを移行することなどが含まれる。例えば、管理サーバ 1 3 が、管理ネットワーク 1 2 を介してホストサーバ 1 0 0 に、新たな仮想マシンの生成を指示する。データセンタ 2

50

0に属するコンピュータやクライアント31は、管理ネットワーク12には直接アクセスしない。

【0035】

管理サーバ13は、データセンタ10の外部に対してクラウドサービスの管理インタフェースを提供するサーバコンピュータである。管理サーバ13は、インターネット30と接続されている。データセンタ20内のコンピュータやクライアント31は、インターネット30を介して管理サーバ13にアクセスすることができる。管理サーバ13は、仮想マシンの生成・起動・停止・削除などのコマンドを受信することがある。

【0036】

例えば、仮想マシンの生成が要求されると、管理サーバ13は、データセンタ10内の
10
ホストサーバのうち要求量以上の空きリソースをもつホストサーバを選択し、選択したホストサーバのリソースを新たな仮想マシンに対して割り当てる。仮想マシンの起動が要求されると、管理サーバ13は、指定された仮想マシンを起動し、当該仮想マシン上でオペレーティングシステムを実行させる。仮想マシンの停止が要求されると、管理サーバ13は、指定された仮想マシン上のオペレーティングシステムをシャットダウンし、当該仮想マシンを停止させる。仮想マシンの削除が要求されると、管理サーバ13は、指定された仮想マシンに割り当てていたリソースを解放する。

【0037】

ホストサーバ100, 100aは、複数の仮想マシンを動作させる仮想化環境をもつサーバコンピュータである。ホストサーバ100, 100aそれぞれには、1つの管理用の
20
仮想マシンと1または2以上のユーザ用の仮想マシンとが配置される。管理用の仮想マシンはホストOSを実行し、ユーザ用の仮想マシンはゲストOSを実行する。

【0038】

また、ホストサーバ100, 100aそれぞれは、当該ホストサーバが備えるプロセッサやメモリなどのリソースを複数の仮想マシンに割り振るハイパーバイザを実行している。ハイパーバイザは、管理サーバ13からの指示に応じて、仮想マシンの生成・削除やデータセンタ10内での仮想マシンのマイグレーションを行う。また、ハイパーバイザは、仮想マシンの通信について、ユーザネットワーク11と管理ネットワーク12とを使い分ける。管理用の仮想マシンの通信については、主に管理ネットワーク12が使用される。ユーザ用の仮想マシンの通信については、主にユーザネットワーク11が使用される。
30

【0039】

クライアント31は、ユーザが操作する端末装置としてのクライアントコンピュータである。クライアント31は、データセンタ10のクラウドサービスを利用するにあたり、インターネット30を介して管理サーバ13にアクセスする。例えば、データセンタ10に新たな仮想マシンを配置する場合、クライアント31は、管理サーバ13に仮想マシンの生成を要求する。生成された仮想マシンには、インターネット30およびユーザネットワーク11を介してアクセスすることができる。同様に、クライアント31は、データセンタ20のクラウドサービスを利用するにあたり、管理サーバ23にアクセスする。

【0040】

第2の実施の形態では、データセンタ10とデータセンタ20との間で仮想マシンのマイグレーションを行うことができる。異なるデータセンタの間のマイグレーションにおけるデータ転送は、管理ネットワーク12, 22を用いず、ユーザネットワーク11, 21を用いて仮想マシン間のデータ通信として行われる。データセンタ10からデータセンタ20に仮想マシンを移行する場合、クライアント31は、ユーザネットワーク11を介して移行元の仮想マシンにマイグレーションを指示する。また、データセンタ20からデータセンタ10に仮想マシンを移行する場合、クライアント31は、ユーザネットワーク21を介して移行元の仮想マシンにマイグレーションを指示する。
40

【0041】

図3は、ホストサーバのハードウェア例を示すブロック図である。

ホストサーバ100は、CPU101、RAM102、HDD103、画像信号処理部
50

104、入力信号処理部105、媒体リーダー106および通信インタフェース107、108を有する。上記のユニットは、それぞれバス109に接続されている。なお、HDD103は、第1の実施の形態の記憶部1cの一例である。

【0042】

CPU101は、プログラムの命令を実行する演算回路を含むプロセッサである。CPU101は、HDD103に記憶されたプログラムやデータの少なくとも一部をRAM102にロードし、プログラムを実行する。なお、CPU101は複数のプロセッサコアを備えてもよく、ホストサーバ100は複数のプロセッサを備えてもよく、以下で説明する処理を複数のプロセッサまたはプロセッサコアを用いて並列に実行してもよい。また、複数のプロセッサの集合(マルチプロセッサ)を「プロセッサ」と呼んでもよい。

10

【0043】

RAM102は、CPU101が実行するプログラムやCPU101が演算に用いるデータを一時的に記憶する揮発性の半導体メモリである。なお、ホストサーバ100は、RAM以外の種類のメモリを備えてもよく、複数個のメモリを備えてもよい。

【0044】

HDD103は、OSやミドルウェアやアプリケーションソフトウェアなどのソフトウェアのプログラム、および、データを記憶する不揮発性の記憶装置である。プログラムには、仮想マシンのマイグレーションを制御するプログラムが含まれる。なお、ホストサーバ100は、フラッシュメモリやSSD(Solid State Drive)などの他の種類の記憶装置を備えてもよく、複数の不揮発性の記憶装置を備えてもよい。

20

【0045】

画像信号処理部104は、CPU101からの命令に従って、ホストサーバ100に接続されたディスプレイ111に画像を出力する。ディスプレイ111としては、CRT(Cathode Ray Tube)ディスプレイ、液晶ディスプレイ(LCD:Liquid Crystal Display)、プラズマディスプレイ(PDP:Plasma Display Panel)、有機EL(OEL:Organic Electro-Luminescence)ディスプレイなどを用いることができる。

【0046】

入力信号処理部105は、ホストサーバ100に接続された入力デバイス112から入力信号を取得し、CPU101に出力する。入力デバイス112としては、マウスやタッチパネルやタッチパッドやトラックボールなどのポインティングデバイス、キーボード、リモートコントローラ、ボタンスイッチなどを用いることができる。また、ホストサーバ100に、複数の種類の入力デバイスが接続されていてもよい。

30

【0047】

媒体リーダー106は、記録媒体113に記録されたプログラムやデータを読み取る読み取り装置である。記録媒体113として、例えば、フレキシブルディスク(FD:Flexible Disk)やHDDなどの磁気ディスク、CD(Compact Disc)やDVD(Digital Versatile Disc)などの光ディスク、光磁気ディスク(MO:Magneto-Optical disk)、半導体メモリなどを使用できる。媒体リーダー106は、例えば、記録媒体113から読み取ったプログラムやデータをRAM102またはHDD103に格納する。

【0048】

40

通信インタフェース107は、ユーザネットワーク11に接続され、ユーザネットワーク11を介してデータセンタ10内外のコンピュータと通信を行う。通信インタフェース108は、管理ネットワーク12に接続され、管理ネットワーク12を介してデータセンタ10内のコンピュータと通信を行う。通信インタフェース107は、スイッチなどの通信装置とケーブルで接続される有線通信インタフェースでもよいし、基地局またはアクセスポイントと無線リンクで接続される無線通信インタフェースでもよい。

【0049】

なお、ホストサーバ100は、媒体リーダー106を備えていなくてもよく、ユーザが操作する端末装置から制御可能である場合には画像信号処理部104や入力信号処理部105を備えていなくてもよい。また、ディスプレイ111や入力デバイス112が、ホスト

50

サーバ100の筐体と一体に形成されていてもよい。管理サーバ13, 23、クライアント31およびホストサーバ100a, 200, 200aも、ホストサーバ100と同様のハードウェア構成によって実現することができる。

【0050】

図4は、仮想マシンとハイパーバイザの配置例を示すブロック図である。

ホストサーバ100は、仮想マシン121, 123, 125を有する。仮想マシン121, 123は、クライアント31のユーザからの要求に応じて生成されたユーザ用の仮想マシンである。仮想マシン121は、ゲストOS122を実行する。仮想マシン123は、ゲストOS124を実行する。ゲストOS122, 124上では、ユーザから指示されたアプリケーションソフトウェアが実行される。仮想マシン125は、ホストサーバ100の管理に用いられる管理用の仮想マシンである。仮想マシン125は、ホストOS126を実行する。ホストOS126上では、管理ソフトウェアが実行される。

10

【0051】

また、ホストサーバ100は、ハイパーバイザ127を実行する。ハイパーバイザ127は、CPU101の処理時間やRAM102の記憶領域などのリソースを仮想マシン121, 123, 125に割り振り、仮想マシン121, 123, 125の動作を制御する。ハイパーバイザ127は、仮想スイッチ128, 129を有する。

【0052】

仮想スイッチ128は、通信インタフェース108の通信帯域を管理し、仮想マシン125と管理ネットワーク12との間の通信を実現する。例えば、仮想スイッチ128は、通信インタフェース108から届いたパケットをホストOS126に届け、ホストOS126が出力したパケットを通信インタフェース108から送信する。

20

【0053】

仮想スイッチ129は、通信インタフェース107の通信帯域を管理し、仮想マシン121, 123とユーザネットワーク11との間の通信や、仮想マシン121, 123の間の通信を実現する。例えば、仮想スイッチ129は、通信インタフェース107から届いたパケットを、宛先アドレスに応じてゲストOS122, 124に振り分ける。また、仮想スイッチ129は、ゲストOS122, 124が出力したパケットを、宛先アドレスに応じて通信インタフェース107から送信するかまたは他のゲストOSに届ける。

【0054】

次に、データセンタ10, 20の間の仮想マシンのマイグレーションを説明する。

30

図5は、仮想マシンのマイグレーション例を示す第1の図である。

以降の第2の実施の形態の説明では、データセンタ10に属するホストサーバ100からデータセンタ20に属するホストサーバ200へのマイグレーションを考える。

【0055】

前述の通り、ホストサーバ100は、CPU101、RAM102およびHDD103などのリソースを有する。これらのリソースを用いて、ホストサーバ100ではゲストOS122を含む仮想マシン121が動作する。CPU101が有するレジスタ101aには、仮想マシン121に関するデータが格納され得る。RAM102の中の仮想マシン121に対して割り当てられた記憶領域には、仮想マシン121に関するデータが一時的に格納される。この割り当てられた記憶領域は、仮想メモリと言うこともできる。

40

【0056】

HDD103の中の記憶領域のうち、パーティション131, 134が仮想マシン121に対して割り当てられる。パーティション131は、ゲストOS122に関するシステムデータを記憶する。システムデータには、ゲストOS122のプログラムや設定ファイルなどが含まれる。パーティション134は、仮想マシン121が使用するデータであってシステムデータ以外のものを記憶する。パーティション134は、ユーザが任意のデータを保存するために使用することが可能である。パーティション131, 134は、仮想マシン121にとって仮想ディスクと言うこともできる。

【0057】

50

パーティション 1 3 1 は、起動情報領域 1 3 2 および退避領域 1 3 3 を含む。起動情報領域 1 3 2 は、起動情報を記憶する。起動情報は、仮想マシン 1 2 1 の起動時に実行されるオペレーティングシステムのデータが記憶されたパーティションを示す。ここでは、起動情報はパーティション 1 3 1 を指し示している。よって、仮想マシン 1 2 1 の起動時にはパーティション 1 3 1 から RAM 1 0 2 にシステムデータがロードされることになる。

【 0 0 5 8 】

退避領域 1 3 3 は、レジスタ 1 0 1 a や RAM 1 0 2 に記憶された仮想マシン 1 2 1 に関するデータを退避する記憶領域である。例えば、仮想マシン 1 2 1 がサスペンド状態に遷移するなど処理を中断するときに、中断時点のレジスタ 1 0 1 a や RAM 1 0 2 やその他の周辺デバイスのデータを含む中断ファイルが退避領域 1 3 3 に書き込まれる。仮想マシン 1 2 1 が中断した処理を再開するとき、中断ファイルを用いて状態が再現される。

【 0 0 5 9 】

また、ホストサーバ 2 0 0 は、CPU 2 0 1、RAM 2 0 2 および HDD 2 0 3 などのリソースを有する。これらのリソースを用いて、ホストサーバ 2 0 0 ではゲスト OS 2 2 2 を含む仮想マシン 2 2 1 が動作する。仮想マシン 2 2 1 は、仮想マシン 1 2 1 の処理を引き継ぐ移行先の仮想マシンである。ゲスト OS 2 2 2 は、仮想マシン 1 2 1 のマイグレーションのために一時的に実行されるオペレーティングシステムである。

【 0 0 6 0 】

CPU 2 0 1 が有するレジスタ 2 0 1 a には、仮想マシン 2 2 1 に関するデータが格納され得る。RAM 2 0 2 の中の仮想マシン 2 2 1 に対して割り当てられた記憶領域には、仮想マシン 2 2 1 に関するデータが一時的に格納される。

【 0 0 6 1 】

HDD 2 0 3 の中の記憶領域のうち、パーティション 2 3 1、2 3 4、2 3 5 が仮想マシン 2 2 1 に対して割り当てられる。パーティション 2 3 1 は、ホストサーバ 1 0 0 から移行されるゲスト OS 1 2 2 に関するシステムデータを格納するために用意した記憶領域である。パーティション 2 3 4 は、ホストサーバ 1 0 0 から移行されるシステムデータ以外のデータを格納するために用意した記憶領域である。パーティション 2 3 1 はパーティション 1 3 1 に対応し、パーティション 2 3 4 はパーティション 1 3 4 に対応する。パーティション 2 3 5 は、ゲスト OS 2 2 2 に関するシステムデータを記憶する。

【 0 0 6 2 】

パーティション 2 3 1 は、起動情報領域 2 3 2 および退避領域 2 3 3 を含む。起動情報領域 2 3 2 は、起動情報領域 1 3 2 と同様に起動情報を記憶する。ただし、仮想マシン 2 2 1 が生成された時点における起動情報領域 2 3 2 の起動情報は、パーティション 2 3 5 を指し示している。よって、仮想マシン 2 2 1 が最初に起動するときは、パーティション 2 3 5 からゲスト OS 2 2 2 のシステムデータがロードされることになる。

【 0 0 6 3 】

マイグレーションを開始する場合、ホストサーバ 2 0 0 上に仮想マシン 2 2 1 が生成される。仮想マシン 2 2 1 に対しては、HDD 2 0 3 のパーティション 2 3 1、2 3 4、2 3 5 が割り当てられる。パーティション 2 3 1 の大きさはパーティション 1 3 1 と同じでよく、パーティション 2 3 4 の大きさはパーティション 1 3 4 と同じでよい。パーティション 2 3 1 の中の起動情報領域 2 3 2 以外の領域とパーティション 2 3 4 は空でよい。

【 0 0 6 4 】

ホストサーバ 2 0 0 は、ゲスト OS 2 2 2 のシステムデータをパーティション 2 3 5 に格納する。ゲスト OS 2 2 2 のシステムデータは、例えば、データセンタ 1 0 が備える所定の記憶装置に用意されている。ホストサーバ 2 0 0 は、所定の記憶装置からパーティション 2 3 5 にシステムデータをコピーする。所定の記憶装置は、何れかのホストサーバが備える不揮発性の記憶装置でもよいし、管理ネットワーク 1 2 に接続されたストレージ装置でもよい。また、ホストサーバ 2 0 0 は、起動パーティションとしてパーティション 2 3 5 を指定した起動情報を起動情報領域 2 3 2 に書き込む。

【 0 0 6 5 】

10

20

30

40

50

そして、ホストサーバ200は、仮想マシン221を起動する。このとき、ホストサーバ200は、起動情報領域232に記憶された起動情報をRAM202にロードし、起動情報に従って、パーティション235からRAM202にシステムデータをロードする。これにより、仮想マシン221上でゲストOS222が実行される。

【0066】

図6は、仮想マシンのマイグレーション例を示す第2の図である。

ホストサーバ200に仮想マシン221が起動すると、仮想マシン121のゲストOS122は、パーティション131に記憶されたシステムデータやパーティション134に記憶されたデータをホストサーバ200に送信する。仮想マシン221のゲストOS222は、受信したパーティション131のシステムデータをパーティション231に格納し、受信したパーティション134のデータをパーティション234に格納する。

10

【0067】

また、仮想マシン121のゲストOS122は、RAM102に記憶された仮想マシン121のデータをホストサーバ200に送信する。仮想マシン221のゲストOS222は、受信したRAM102のデータをHDD203の退避領域233に格納する。パーティション131、134およびRAM102のデータの送信は、仮想マシン121と仮想マシン221との間の通信として行われる。よって、上記のデータ送信は、管理ネットワーク12、22を使用せず、ユーザネットワーク11、21を使用して行われる。

【0068】

その後、仮想マシン121のゲストOS122は、ホストサーバ200に送信済のデータがホストサーバ100において更新されたか監視する。更新された場合、仮想マシン121のゲストOS122は、前回送信してからの差分データをホストサーバ200に送信する。仮想マシン221のゲストOS222は、受信した差分データをパーティション231、234の中の該当する記憶領域に上書き保存する。差分データの送信は、ホストサーバ100における前回差分データ送信時からのデータ更新量が閾値以下に低下するまで、あるいはデータ更新量の低下が限界に達する状態になるまで、継続する。

20

【0069】

図7は、仮想マシンのマイグレーション例を示す第3の図である。

ホストサーバ100における単位時間当たりのデータ更新量が低下すると、仮想マシン121のゲストOS122は、仮想マシン121を停止する前の最終データをホストサーバ200に送信する。最終データには、パーティション131、134の更新データやRAM102の更新データの他に、レジスタ101aに記憶されたCPU101の状態データや周辺デバイスの状態を示すデバイスデータなどが含まれる。レジスタ101aのデータは、ゲストOS122が一時的にサスペンド状態に遷移することで、退避領域133に格納される中断ファイルから抽出することができる。

30

【0070】

仮想マシン221のゲストOS222は、受信した最終データに含まれるパーティション131、134の更新データを、パーティション231、234に上書き保存する。また、仮想マシン221のゲストOS222は、受信した最終データに含まれるRAM102の更新データを退避領域233に上書き保存し、受信した最終データに含まれるCPU101の状態データなどを退避領域233に格納する。

40

【0071】

そして、仮想マシン221のゲストOS222は、最新のRAM102のデータとCPU101の状態データから、仮想マシン121の状態を再現するための中断ファイルを生成し退避領域233に格納する。また、仮想マシン221のゲストOS222は、起動情報領域232に記憶された起動情報を、起動時にパーティション235のシステムデータに代えてパーティション231のシステムデータがロードされるように書き換える。このとき、仮想マシン121のゲストOS122はシャットダウンしてよい。

【0072】

図8は、仮想マシンのマイグレーション例を示す第4の図である。

50

ホストサーバ200で中断ファイルが生成されて起動情報が書き換えられると、仮想マシン221のゲストOS222は、仮想マシン221を再起動するコマンドを発行する。すると、ゲストOS222がシャットダウンし、仮想マシン221は起動情報領域232に記憶された起動情報をロードする。起動情報はパーティション231を指し示しているため、仮想マシン221はパーティション231からRAM202にシステムデータをロードする。ロードされるシステムデータは、ゲストOS122のシステムデータである。

【0073】

仮想マシン221上でゲストOS122が起動すると、仮想マシン221のゲストOS122は、退避領域233に記憶された中断ファイルを用いて、仮想マシン121が停止したときの仮想マシン121の状態を再現する。仮想マシン221のゲストOS122は、中断ファイルに含まれるCPU101の状態データをCPU201のレジスタ201aにロードする。また、仮想マシン221のゲストOS122は、中断ファイルに含まれるRAM102のデータをRAM202にロードする。これにより、仮想マシン121が停止したときの仮想マシン121の処理を、仮想マシン221が引き継ぐことができる。

【0074】

次に、仮想マシン121, 221が有するマイグレーション機能について説明する。

図9は、仮想マシンの機能例を示すブロック図である。

仮想マシン121で実行されるゲストOS122は、マイグレーション制御部141、更新監視部142およびカーネル143を有する。マイグレーション制御部141と更新監視部142は、例えば、ゲストOS122に対して組み込まれたプログラムモジュールとして実装できる。上記のプログラムモジュールは、データセンタ10によって予め組み込まれていてもよいし、仮想マシン121のユーザが組み込んでよい。

【0075】

例えば、データセンタ10において仮想マシン121を生成したときに、マイグレーション制御部141と更新監視部142を有するゲストOS122が自動的に使用されるようにしてもよい。また、仮想マシン121のユーザは、データセンタ10によって用意されたオペレーティングシステムに対して上記のプログラムモジュールを追加してもよい。また、仮想マシン121のユーザは、仮想マシン121の生成後に、仮想マシン121が実行するオペレーティングシステムをゲストOS122に変更してもよい。

【0076】

マイグレーション制御部141は、仮想マシン121のマイグレーションを制御する。マイグレーション制御部141は、クライアント31からユーザネットワーク11を介してマイグレーション開始の指示を受信する。すると、マイグレーション制御部141は、ユーザネットワーク11を介してデータセンタ20の管理サーバ23にアクセスし、データセンタ20に移行先の仮想マシン221を生成させる。そして、マイグレーション制御部141は、仮想マシン121に関するHDD103のデータやRAM102のデータなどを、ユーザネットワーク11, 21を介して仮想マシン221に送信する。

【0077】

また、マイグレーション制御部141は、仮想マシン221へのデータ送信を継続している間、仮想マシン121で実行中のプロセスを監視する。仮想マシン121のプロセスによるデータ更新頻度が高い場合、マイグレーション制御部141は、データ更新頻度が下がるように各プロセスのCPU使用率を制限する。データ送信が完了して移行先の仮想マシン221が再起動すると、ゲストOS122は仮想マシン221上で実行されることになる。この場合、マイグレーション制御部141は、マイグレーションの後処理（例えば、移行元の仮想マシン121の削除など）を仮想マシン221において実行する。

【0078】

更新監視部142は、マイグレーション制御部141からの指示に応じて、HDD103のデータやRAM102のデータの更新を監視する。HDD103やRAM102へのデータの書き込みは、例えば、カーネル143を利用して検出できる。更新監視部142は、更新されたデータを示す情報をマイグレーション制御部141に提供する。

10

20

30

40

50

【 0 0 7 9 】

カーネル 1 4 3 は、ゲスト OS 1 2 2 の中核となる機能を実装したプログラムモジュールである。カーネル 1 4 3 は、CPU 1 0 1 や RAM 1 0 2 などのハードウェアへのアクセスを制御し、ゲスト OS 1 2 2 内のカーネル 1 4 3 以外のプログラムモジュールに対してハードウェアアクセスに関するインタフェースを提供する。例えば、カーネル 1 4 3 は、プロセスの情報や、RAM 1 0 2 や HDD 1 0 3 へのアクセスの情報を収集する。また、カーネル 1 4 3 は、各プロセスに対して CPU 1 0 1 の処理時間を割り当てる。

【 0 0 8 0 】

仮想マシン 1 2 1 は、制御情報記憶部 1 4 4 を有する。制御情報記憶部 1 4 4 は、例えば、仮想マシン 1 2 1 に割り当てられた RAM 1 0 2 の記憶領域または HDD 1 0 3 の記憶領域として実現できる。制御情報記憶部 1 4 4 は、マイグレーション制御部 1 4 1 および更新監視部 1 4 2 がマイグレーション処理に使用する制御情報を記憶する。

10

【 0 0 8 1 】

仮想マシン 2 2 1 で実行されるゲスト OS 2 2 2 は、マイグレーション制御部 2 4 1 を有する。マイグレーション制御部 2 4 1 は、例えば、ゲスト OS 2 2 2 に対して組み込まれたプログラムモジュールとして実装できる。このプログラムモジュールは、データセンタ 2 0 によって予め組み込まれていてもよい。例えば、マイグレーションの移行先として仮想マシン 2 2 1 が生成されたときに、マイグレーション制御部 2 4 1 を有するゲスト OS 2 2 2 がデータセンタ 2 0 によって自動的に選択されるようにしてもよい。

【 0 0 8 2 】

マイグレーション制御部 2 4 1 は、移行元のマイグレーション制御部 1 4 1 と連携して仮想マシン 1 2 1 のマイグレーションを制御する。マイグレーション制御部 2 4 1 は、マイグレーション制御部 1 4 1 から仮想マシン 1 2 1 に関するデータを受信し、受信したデータを HDD 2 0 3 に保存する。全てのデータを受信し終わると、マイグレーション制御部 2 4 1 は、仮想マシン 2 2 1 を再起動したときに、受信したデータに基づいてゲスト OS 1 2 2 が実行されるように、仮想マシン 2 2 1 の起動情報を書き換える。

20

【 0 0 8 3 】

また、マイグレーション制御部 2 4 1 は、ユーザネットワーク 2 1 を介してデータセンタ 1 0 の管理サーバ 1 3 にアクセスし、仮想マシン 1 2 1 を停止させる。そして、マイグレーション制御部 2 4 1 は、仮想マシン 2 2 1 を再起動する。再起動後の仮想マシン 2 2 1 では、ゲスト OS 2 2 2 に代えてゲスト OS 1 2 2 が実行されることになる。よって、移行先の仮想マシン 2 2 1 におけるマイグレーション処理は、マイグレーション制御部 1 4 1 がマイグレーション制御部 2 4 1 から引き継ぐことになる。

30

【 0 0 8 4 】

仮想マシン 2 2 1 は、制御情報記憶部 2 4 2 を有する。制御情報記憶部 2 4 2 は、例えば、仮想マシン 2 2 1 に割り当てられた RAM 2 0 2 の記憶領域または HDD 2 0 3 の記憶領域として実現できる。制御情報記憶部 2 4 2 は、マイグレーション制御部 2 4 1 がマイグレーション処理に使用する制御情報を記憶する。なお、ゲスト OS 1 2 2 のシステムデータが仮想マシン 1 2 1 から仮想マシン 2 2 1 に移行すると、制御情報記憶部 1 4 4 に記憶された制御情報も仮想マシン 1 2 1 から仮想マシン 2 2 1 に移行することになる。

40

【 0 0 8 5 】

図 1 0 は、状態テーブルの例を示す図である。

制御情報記憶部 1 4 4 は、状態テーブル 1 5 1 を記憶する。状態テーブル 1 5 1 は、パラメータ名およびフラグの項目を有する。パラメータ名の項目には、パラメータの名称として「移行済」、「ディスク更新低下」および「メモリ更新低下」が登録される。フラグの項目には、「1」(True) または「0」(False) が登録される。

【 0 0 8 6 】

移行済パラメータ = 1 は、仮想マシン 1 2 1 のデータを仮想マシン 2 2 1 にコピーし終えたことを示し、移行済パラメータ = 0 は、データをコピーし終わっていないことを示す。ディスク更新低下パラメータ = 1 は、未送信である HDD 1 0 3 の差分データの量が閾値

50

以下であることを示し、ディスク更新低下パラメータ = 0 は、閾値より大きいことを示す。メモリ更新低下パラメータ = 1 は、未送信である R A M 1 0 2 の差分データの量が閾値以下であることを示し、メモリ更新低下パラメータ = 0 は、閾値より大きいことを示す。全てのパラメータについてフラグの初期値は「0」である。

【0087】

図11は、ディスク更新情報とメモリ更新情報の例を示す図である。

制御情報記憶部144は、ディスク更新情報152およびメモリ更新情報153を記憶する。ディスク更新情報152は、更新監視部142がHDD103のデータの更新を監視するときに使用するものである。メモリ更新情報153は、更新監視部142がRAM102のデータの更新を監視するときに使用するものである。

10

【0088】

ディスク更新情報152は、仮想マシン121に割り当てられたHDD103の記憶領域への書き込みの有無を、セクタ単位で表現したビットマップである。セクタは、例えば、HDD103の記憶領域を細分化した固定長の領域である。あるセクタへデータが書き込まれ当該データが仮想マシン221に未送信である場合、当該セクタに対応するビットが「1」に設定される。それ以外のビットは「0」に設定される。

【0089】

メモリ更新情報153は、仮想マシン121に割り当てられたRAM102の記憶領域への書き込みの有無を、ページ単位で表現したビットマップである。ページは、例えば、RAM102の記憶領域を細分化した固定長の領域である。あるページへデータが書き込まれ当該データが仮想マシン221に未送信である場合、当該ページに対応するビットが「1」に設定される。それ以外のビットは「0」に設定される。

20

【0090】

図12は、プロセステーブルの例を示す図である。

制御情報記憶部144は、プロセステーブル154を記憶する。プロセステーブル154は、マイグレーション制御部141が各プロセスのCPU使用率を制限するときに使用するものである。プロセステーブル154は、プロセスID、プロセス名、開始状態、状態、CPU使用率、ディスク更新量およびメモリ更新量の項目を有する。

【0091】

プロセスIDの項目には、カーネル143によって付与されたプロセスの識別番号が登録される。プロセス名の項目には、プロセスの名称（例えば、当該プロセスの処理を定義したプログラムの名称など）が登録される。開始状態の項目には、仮想マシン121のマイグレーションの開始時点における各プロセスの状態が登録される。開始状態としては、「実行中」または「一時停止」が挙げられる。状態の項目には、現在の各プロセスの状態が登録される。状態としては、「実行中」、「一時停止」または「終了」が挙げられる。

30

【0092】

CPU使用率の項目には、カーネル143が各プロセスに対して現在割り当てているCPU時間（リソース量）の指標値が登録される。CPU時間は、例えば、1つのプロセッサまたはプロセッサコアの演算能力を100としたときの百分率で表現できる。ディスク更新量の項目には、各プロセスによる単位時間当たりのHDD103への書き込み量が登録される。メモリ更新量の項目には、各プロセスによる単位時間当たりのRAM102への書き込み量が登録される。プロセステーブル154に登録される情報は、例えば、カーネル143のインタフェースを用いて取得することができる。

40

【0093】

図13は、リソース制限テーブルと重要プロセスリストの例を示す図である。

制御情報記憶部144は、リソース制限テーブル155および重要プロセスリスト156を記憶する。リソース制限テーブル155および重要プロセスリスト156は、マイグレーション制御部141が各プロセスのCPU使用率を制限するときに使用するものである。リソース制限テーブル155および重要プロセスリスト156は、例えば、仮想マシン121のユーザによって作成されて制御情報記憶部144に格納される。

50

【 0 0 9 4 】

リソース制限テーブル 1 5 5 は、プロセス名およびリソース上限の項目を有する。プロセス名の項目には、プロセステーブル 1 5 4 と同様のプロセスの名称が登録される。リソース上限の項目には、各プロセスについて CPU 使用率を最も厳しく制限した場合の CPU 使用率の上限が登録される。例えば、図 1 3 の例ではプロセス B のリソース上限が 2 0 % に設定されている一方、図 1 2 の例ではプロセス B の現在の CPU 使用率は 2 5 % である。よって、CPU 使用率を制限する場合、プロセス B に対する CPU 時間の割り当てを減らすことにより、プロセス B の CPU 使用率が 2 0 % 以下に低下し得る。

【 0 0 9 5 】

重要プロセスリスト 1 5 6 には、1 または 2 以上のプロセス名が登録される。重要プロセスリスト 1 5 6 が示すプロセスは、仮想マシン 1 2 1 のユーザにとって重要なプロセスであり、CPU 時間の割り当てを減らさない (CPU 使用率を制限する対象から除外する) プロセスを示す。例えば、重要プロセスリスト 1 5 6 には、プロセス A , B , C , D , E のうちプロセス A , C が登録される。この場合、プロセス B , D , E の CPU 使用率は制限される可能性がある一方、プロセス A , C の CPU 使用率は制限されない。

10

【 0 0 9 6 】

図 1 4 は、仮想マシン生成要求の例を示す図である。

仮想マシン生成要求 1 5 7 は、仮想マシン 1 2 1 のマイグレーション制御部 1 4 1 からデータセンタ 2 0 の管理サーバ 2 3 に対して送信される。仮想マシン生成要求 1 5 7 は、移行先の仮想マシン 2 2 1 の生成の要求を示す。仮想マシン生成要求 1 5 7 は、パラメータ名と設定値の組を複数含む。パラメータには、CPU 数、メモリ容量、システムディスク容量、他データディスク容量、ネットワーク設定および初期 OS が含まれる。

20

【 0 0 9 7 】

「 CPU 数 」 は、仮想マシン 2 2 1 に割り当てる CPU の数を示す。通常、CPU 数は仮想マシン 1 2 1 と同じでよい。「メモリ容量」は、仮想マシン 2 2 1 に割り当てる RAM 2 0 2 の記憶領域の大きさを示す。通常、メモリ容量は仮想マシン 1 2 1 と同じでよい。「システムディスク容量」は、仮想マシン 2 2 1 に割り当てる HDD 2 0 3 のパーティション 2 3 1 , 2 3 5 の大きさを示す。通常、システムディスク容量は、仮想マシン 1 2 1 のパーティション 1 3 1 の大きさに、ゲスト OS 2 2 2 のシステムデータを記憶するパーティション 2 3 5 の大きさを加えたものである。なお、システムディスクのフォーマットは、仮想マシン 1 2 1 のパーティション 1 3 1 と同じでよい。

30

【 0 0 9 8 】

「他データディスク容量」は、仮想マシン 2 2 1 に割り当てる HDD 2 0 3 のパーティション 2 3 4 の大きさを示す。通常、他データディスク容量は、仮想マシン 1 2 1 のパーティション 1 3 4 と同じでよい。なお、他データディスクのフォーマットは、仮想マシン 1 2 1 のパーティション 1 3 4 と同じでよい。「ネットワーク設定」は、仮想マシン 2 2 1 が使用可能なネットワークを示す。通常、ネットワーク設定は仮想マシン 1 2 1 と同じでよい。ただし、仮想マシン 2 2 1 は仮想マシン 1 2 1 と通信するため、仮想マシン 1 2 1 が使用するネットワークにインターネット 3 0 が含まれていない場合、インターネット 3 0 を追加登録する (この場合、仮想マシン 1 2 1 の設定も変更する) 。

40

【 0 0 9 9 】

「初期 OS」は、仮想マシン 2 2 1 を最初に起動したときに実行されるオペレーティングシステムの種類を示す。初期 OS としてマイグレーション用 OS を指定した場合、マイグレーション制御部 2 4 1 を有するゲスト OS 2 2 2 のシステムデータがパーティション 2 3 5 に格納される。仮想マシン 2 2 1 の起動時には、ゲスト OS 2 2 2 が実行される。

【 0 1 0 0 】

次に、仮想マシン 1 2 1 , 2 2 1 が行うマイグレーション処理について説明する。

図 1 5 は、マイグレーション開始の手順例を示すフローチャートである。

クライアント 3 1 は、インターネット 3 0 およびユーザネットワーク 1 1 を介して、仮想マシン 1 2 1 に対してマイグレーション開始指示を送信する (S 1 1 0) 。

50

【 0 1 0 1 】

マイグレーション開始指示を受信すると、仮想マシン 1 2 1 上のマイグレーション制御部 1 4 1 は、ユーザネットワーク 1 1 およびインターネット 3 0 を介して、データセンタ 2 0 の管理サーバ 2 3 に仮想マシン生成要求 1 5 7 を送信する (S 1 1 1)。このとき、マイグレーション制御部 1 4 1 は、仮想マシン 1 2 1 のリソース割り当て状況から、仮想マシン生成要求 1 5 7 の各パラメータの設定値を算出する。原則として、仮想マシン 2 2 1 の構成は、パーティション 2 3 5 を除いて仮想マシン 1 2 1 と同じでよい。

【 0 1 0 2 】

仮想マシン生成要求 1 5 7 を受信すると、管理サーバ 2 3 は、仮想マシン生成要求 1 5 7 で指定された量のリソースを仮想マシン 2 2 1 に割り当てることで、ホストサーバ 2 0 0 上に仮想マシン 2 2 1 を生成する (S 1 1 2)。管理サーバ 2 3 は、仮想マシン 2 2 1 に対して付与されたアドレス (例えば、IP (Internet Protocol) アドレス) を、インターネット 3 0 およびユーザネットワーク 1 1 を介して仮想マシン 1 2 1 に通知する (S 1 1 3)。マイグレーション制御部 1 4 1 は、制御情報記憶部 1 4 4 に状態テーブル 1 5 1 を生成し、全てのフラグを「0」に初期化する (S 1 1 4)。

10

【 0 1 0 3 】

マイグレーション制御部 1 4 1 は、ユーザネットワーク 1 1 およびインターネット 3 0 を介して管理サーバ 2 3 に、ステップ S 1 1 2 で生成された仮想マシンを指定した仮想マシン起動要求を送信する (S 1 1 5)。仮想マシン起動要求を受信すると、管理サーバ 2 3 は、指定された仮想マシン 2 2 1 が配置されたホストサーバ 2 0 0 を特定し、管理ネットワーク 1 2 を介して、ホストサーバ 2 0 0 のハイパーバイザに対して仮想マシン起動指示を送信する (S 1 1 6)。ホストサーバ 2 0 0 のハイパーバイザは、仮想マシン 2 2 1 を起動する。このとき、起動情報領域 2 3 2 に記憶された起動情報に従って、パーティション 2 3 5 に記憶されたシステムデータがロードされてゲスト OS 2 2 2 が実行される。そして、マイグレーション制御部 2 4 1 が起動し、接続待ち状態になる (S 1 1 7)。

20

【 0 1 0 4 】

マイグレーション制御部 1 4 1 は、ステップ S 1 1 3 で通知されたアドレスを用いて、ユーザネットワーク 1 1 およびインターネット 3 0 を介して、仮想マシン 2 2 1 に対して接続要求を送信する (S 1 1 8)。マイグレーション制御部 2 4 1 は、接続処理を行い、仮想マシン 1 2 1 に対して接続完了通知を返信する (S 1 1 9)。マイグレーション制御部 1 4 1 は、以下で説明するディスクデータ送信、メモリデータ送信およびプロセス監視の 3 つのプロセスを起動する (S 1 2 0)。この 3 つのプロセスは並行に実行される。

30

【 0 1 0 5 】

図 1 6 は、ディスクデータ送信の手順例を示すフローチャートである。

このディスクデータ送信のプロセスは、上記のステップ S 1 2 0 で起動される。

マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にディスク更新情報 1 5 2 の初期化を指示する (S 2 1 0)。更新監視部 1 4 2 は、制御情報記憶部 1 4 4 に記憶されたディスク更新情報 1 5 2 の全てのビットを「0」に初期化する (S 2 1 1)。更新監視部 1 4 2 は、ビットマップ初期化完了をマイグレーション制御部 1 4 1 に通知する。マイグレーション制御部 1 4 1 は、ディスク書き込みの監視指示を更新監視部 1 4 2 に通知する (S 2 1 2)。更新監視部 1 4 2 は、ディスク監視のプロセスを起動することで、ディスク書き込みの監視を開始する (S 2 1 3)。ディスク監視のプロセスは、図 1 6 のプロセスなど他のプロセスと並行に実行される (S 2 1 4)。ディスク監視の処理は後述する。

40

【 0 1 0 6 】

更新監視部 1 4 2 は、ディスク監視開始をマイグレーション制御部 1 4 1 に通知する。マイグレーション制御部 1 4 1 は、ユーザネットワーク 1 1、インターネット 3 0 およびユーザネットワーク 2 1 を介して、HDD 1 0 3 のパーティション 1 3 1、1 3 4 に記憶されたデータを、セクタ単位で仮想マシン 2 2 1 に送信する (S 2 1 5)。このとき、マイグレーション制御部 1 4 1 は、ディスク更新情報 1 5 2 のビットが「1」になっているセクタのデータを送信しなくてもよい。データの送信時には、セクタアドレスが付加され

50

る。マイグレーション制御部 2 4 1 は、受信したデータを、パーティション 2 3 1 , 2 3 4 の中のセクタアドレスが示す記憶領域に保存する (S 2 1 6)。

【 0 1 0 7 】

マイグレーション制御部 2 4 1 は、データ受信完了を仮想マシン 1 2 1 に返信する。パーティション 1 3 1 , 1 3 4 全体について初回データ送信が終わると、マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にディスク更新を問い合わせる (S 2 1 7)。更新監視部 1 4 2 は、ディスク更新情報 1 5 2 のビットが「 1 」になっているセクタ、すなわち、データ書き込みのあったセクタを特定し、特定したセクタのセクタアドレスをマイグレーション制御部 1 4 1 に通知する (S 2 1 8)。

【 0 1 0 8 】

マイグレーション制御部 1 4 1 は、データ書き込みのあったセクタ数またはセクタ数に応じたデータ量 (ディスク更新量) が閾値以下であるか判断する (S 2 1 9)。ディスク更新量が閾値以下である場合、状態テーブル 1 5 1 のディスク更新低下フラグを「 1 」に設定し、ステップ S 2 2 6 に処理が進む。ディスク更新量が閾値より大きい場合、ディスク更新低下フラグを「 0 」に設定し、ステップ S 2 2 0 に処理が進む。

【 0 1 0 9 】

ディスク更新量が閾値より大きい場合、マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にディスク更新情報 1 5 2 の初期化を指示する (S 2 2 0)。更新監視部 1 4 2 は、ディスク更新情報 1 5 2 の全てのビットを「 0 」に初期化する (S 2 2 1)。更新監視部 1 4 2 は、ビットマップ初期化完了をマイグレーション制御部 1 4 1 に通知する。

【 0 1 1 0 】

マイグレーション制御部 1 4 1 は、ステップ S 2 1 8 または後述するステップ S 2 2 5 で通知されたセクタアドレスが示すデータ (更新されたデータ) を、ディスク差分データとしてパーティション 1 3 1 , 1 3 4 から読み出す。マイグレーション制御部 1 4 1 は、ディスク差分データにセクタアドレスを付加し、ユーザネットワーク 1 1、インターネット 3 0 およびユーザネットワーク 2 1 を介して仮想マシン 2 2 1 に送信する (S 2 2 2)。マイグレーション制御部 2 4 1 は、受信したデータを、パーティション 2 3 1 , 2 3 4 の中のセクタアドレスが示す記憶領域に上書き保存する (S 2 2 3)。

【 0 1 1 1 】

マイグレーション制御部 2 4 1 は、データ受信完了を仮想マシン 1 2 1 に返信する。マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にディスク更新を問い合わせる (S 2 2 4)。更新監視部 1 4 2 は、ディスク更新情報 1 5 2 のビットが「 1 」になっているセクタ、すなわち、前回の確認以後にデータ書き込みのあったセクタを特定し、特定したセクタのセクタアドレスをマイグレーション制御部 1 4 1 に通知する (S 2 2 5)。

【 0 1 1 2 】

マイグレーション制御部 1 4 1 は、ディスク更新量が閾値以下である場合、状態テーブル 1 5 1 のディスク更新低下フラグを「 1 」に設定し、ディスク更新量が閾値より大きい場合、ディスク更新低下フラグを「 0 」に設定する。そして、マイグレーション制御部 1 4 1 は、状態テーブル 1 5 1 のディスク更新低下フラグとメモリ更新低下フラグの両方が「 1 」であるか判断する (S 2 2 6)。2 つのフラグが「 1 」である場合、ディスクデータ送信のプロセスが終了し、後述するマイグレーション完了のプロセスが起動する。少なくとも一方のフラグが「 0 」である場合、ステップ S 2 2 0 に処理が進む。

【 0 1 1 3 】

図 1 7 は、ディスク監視の手順例を示すフローチャートである。

このディスク監視のプロセスは、上記のステップ S 2 1 4 で起動される。

更新監視部 1 4 2 は、HDD 1 0 3 に対する書き込みが発生したときに当該書き込みが検出されるよう、カーネル 1 4 3 のインタフェースを利用してフック設定を行う (S 2 3 0)。更新監視部 1 4 2 は、HDD 1 0 3 への書き込みが検出されると、書き込み先のセクタを特定する (S 2 3 1)。更新監視部 1 4 2 は、ディスク更新情報 1 5 2 の中のステップ S 2 3 1 で特定したセクタに対応するビットを「 1 」に設定する (S 2 3 2)。

10

20

30

40

50

【 0 1 1 4 】

更新監視部 1 4 2 は、状態テーブル 1 5 1 のディスク更新低下フラグとメモリ更新低下フラグの両方が「 1 」であるか判断する (S 2 3 3)。2 つのフラグが「 1 」である場合、ディスク監視のプロセスが終了する。少なくとも一方のフラグが「 0 」である場合、ステップ S 2 3 1 に処理が進み、次に HDD 1 0 3 への書き込みが検出されるのを待つ。

【 0 1 1 5 】

図 1 8 は、メモリデータ送信の手順例を示すフローチャートである。

このメモリデータ送信のプロセスは、上記のステップ S 1 2 0 で起動される。

マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にメモリ更新情報 1 5 3 の初期化を指示する (S 2 4 0)。更新監視部 1 4 2 は、制御情報記憶部 1 4 4 に記憶されたメモリ更新情報 1 5 3 の全てのビットを「 0 」に初期化する (S 2 4 1)。更新監視部 1 4 2 は、ビットマップ初期化完了をマイグレーション制御部 1 4 1 に通知する。マイグレーション制御部 1 4 1 は、メモリ書き込みの監視指示を更新監視部 1 4 2 に通知する (S 2 4 2)。更新監視部 1 4 2 は、メモリ監視のプロセスを起動することで、メモリ書き込みの監視を開始する (S 2 4 3)。メモリ監視のプロセスは、図 1 8 のプロセスなど他のプロセスと並行に実行される (S 2 4 4)。メモリ監視の処理は後述する。

10

【 0 1 1 6 】

更新監視部 1 4 2 は、メモリ監視開始をマイグレーション制御部 1 4 1 に通知する。マイグレーション制御部 1 4 1 は、ユーザネットワーク 1 1、インターネット 3 0 およびユーザネットワーク 2 1 を介して、仮想マシン 1 2 1 に割り当てられた RAM 1 0 2 の記憶領域に記憶されたデータを、ページ単位で仮想マシン 2 2 1 に送信する (S 2 4 5)。このとき、マイグレーション制御部 1 4 1 は、メモリ更新情報 1 5 3 のビットが「 1 」になっているページのデータを送信しなくてもよい。データの送信時には、ページアドレスが付加される。マイグレーション制御部 2 4 1 は、受信したデータを、RAM 2 0 2 の中のページアドレスが示す記憶領域に保存する (S 2 4 6)。

20

【 0 1 1 7 】

マイグレーション制御部 2 4 1 は、データ受信完了を仮想マシン 1 2 1 に返信する。仮想マシン 1 2 1 に割り当てられた RAM 1 0 2 の記憶領域全体について初回データ送信が終わると、マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にメモリ更新を問い合わせる (S 2 4 7)。更新監視部 1 4 2 は、メモリ更新情報 1 5 3 のビットが「 1 」になっているページ、すなわち、データ書き込みのあったページを特定し、特定したページのページアドレスをマイグレーション制御部 1 4 1 に通知する (S 2 4 8)。

30

【 0 1 1 8 】

マイグレーション制御部 1 4 1 は、データ書き込みのあったページ数またはページ数に応じたデータ量 (メモリ更新量) が閾値以下であるか判断する (S 2 4 9)。メモリ更新量が閾値以下である場合、状態テーブル 1 5 1 のメモリ更新低下フラグを「 1 」に設定し、ステップ S 2 5 6 に処理が進む。メモリ更新量が閾値より大きい場合、メモリ更新低下フラグを「 0 」に設定し、ステップ S 2 5 0 に処理が進む。

【 0 1 1 9 】

メモリ更新量が閾値より大きい場合、マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にメモリ更新情報 1 5 3 の初期化を指示する (S 2 5 0)。更新監視部 1 4 2 は、メモリ更新情報 1 5 3 の全てのビットを「 0 」に初期化する (S 2 5 1)。更新監視部 1 4 2 は、ビットマップ初期化完了をマイグレーション制御部 1 4 1 に通知する。

40

【 0 1 2 0 】

マイグレーション制御部 1 4 1 は、ステップ S 2 4 8 または後述するステップ S 2 5 5 で通知されたメモリアドレスが示すデータ (更新されたデータ) を、メモリ差分データとして RAM 1 0 2 から読み出す。マイグレーション制御部 1 4 1 は、メモリ差分データにページアドレスを付加し、ユーザネットワーク 1 1、インターネット 3 0 およびユーザネットワーク 2 1 を介して仮想マシン 2 2 1 に送信する (S 2 5 2)。マイグレーション制御部 2 4 1 は、受信したデータを、RAM 2 0 2 の中のページアドレスが示す記憶領域に

50

上書き保存する (S 2 5 3)。

【 0 1 2 1 】

マイグレーション制御部 2 4 1 は、データ受信完了を仮想マシン 1 2 1 に返信する。マイグレーション制御部 1 4 1 は、更新監視部 1 4 2 にメモリ更新を問い合わせる (S 2 5 4)。更新監視部 1 4 2 は、メモリ更新情報 1 5 3 のビットが「 1 」になっているページ、すなわち、前回の確認以後にデータ書き込みのあったページを特定し、特定したページのページアドレスをマイグレーション制御部 1 4 1 に通知する (S 2 5 5)。

【 0 1 2 2 】

マイグレーション制御部 1 4 1 は、メモリ更新量が閾値以下である場合、状態テーブル 1 5 1 のメモリ更新低下フラグを「 1 」に設定し、メモリ更新量が閾値より大きい場合、メモリ更新低下フラグを「 0 」に設定する。そして、マイグレーション制御部 1 4 1 は、状態テーブル 1 5 1 のディスク更新低下フラグとメモリ更新低下フラグの両方が「 1 」であるか判断する (S 2 5 6)。2 つのフラグが「 1 」である場合、メモリデータ送信のプロセスが終了し、後述するマイグレーション完了のプロセスが起動する。少なくとも一方のフラグが「 0 」である場合、ステップ S 2 5 0 に処理が進む。

【 0 1 2 3 】

図 1 9 は、メモリ監視の手順例を示すフローチャートである。

このメモリ監視のプロセスは、上記のステップ S 2 4 4 で起動される。

更新監視部 1 4 2 は、RAM 1 0 2 に対する書き込みが発生したときに当該書き込みが検出されるよう、カーネル 1 4 3 のインタフェースを利用してフック設定を行う (S 2 6 0)。更新監視部 1 4 2 は、RAM 1 0 2 への書き込みが検出されると、書き込み先のページを特定する (S 2 6 1)。更新監視部 1 4 2 は、メモリ更新情報 1 5 3 中のステップ S 2 6 1 で特定したページに対応するビットを「 1 」に設定する (S 2 6 2)。

【 0 1 2 4 】

更新監視部 1 4 2 は、状態テーブル 1 5 1 のディスク更新低下フラグとメモリ更新低下フラグの両方が「 1 」であるか判断する (S 2 6 3)。2 つのフラグが「 1 」である場合、メモリ監視のプロセスが終了する。少なくとも一方のフラグが「 0 」である場合、ステップ S 2 6 1 に処理が進み、次に RAM 1 0 2 への書き込みが検出されるのを待つ。

【 0 1 2 5 】

図 2 0 は、プロセス監視の手順例を示すフローチャートである。

このプロセス監視のプロセスは、上記のステップ S 1 2 0 で起動される。

マイグレーション制御部 1 4 1 は、ゲスト OS 1 2 2 上に存在するプロセスのリストを、カーネル 1 4 3 のインタフェースを利用して取得する (S 2 7 0)。プロセスのリストには、プロセス ID、プロセス名および状態 (実行中または一時停止中) が含まれる。マイグレーション制御部 1 4 1 は、取得したプロセス ID とプロセス名をプロセステーブル 1 5 4 に登録し、取得した状態を開始状態としてプロセステーブル 1 5 4 に登録する。

【 0 1 2 6 】

マイグレーション制御部 1 4 1 は、各プロセスについて、現在の状態、CPU 使用率、ディスク更新量およびメモリ更新量を、カーネル 1 4 3 のインタフェースを利用して取得する (S 2 7 1)。マイグレーション制御部 1 4 1 は、取得した状態、CPU 使用率、ディスク更新量およびメモリ更新量をプロセステーブル 1 5 4 に登録する。

【 0 1 2 7 】

マイグレーション制御部 1 4 1 は、プロセステーブル 1 5 4 に登録されたプロセスの中から、ディスク更新に関する以下の条件を満たすプロセスを高々 1 つ選択する (S 2 7 2)。条件を満たすプロセスは、重要プロセスリスト 1 5 6 に登録されておらず、CPU 使用率がリソース制限テーブル 1 5 5 に登録されたリソース上限を超えているプロセスの中で、単位時間当たりのディスク更新量が最大のものである。

【 0 1 2 8 】

マイグレーション制御部 1 4 1 は、プロセステーブル 1 5 4 に登録されたプロセスの中から、メモリ更新に関する以下の条件を満たすプロセスを高々 1 つ選択する (S 2 7 3)

10

20

30

40

50

。条件を満たすプロセスは、重要プロセスリスト156に登録されておらず、CPU使用率がリソース制限テーブル155に登録されたリソース上限を超えているプロセスの中で、単位時間当たりのメモリ更新量が最大のものである。

【0129】

マイグレーション制御部141は、ステップS272，S273で選択したプロセスがある場合、カーネル143のインタフェースを利用して、当該プロセスのCPU使用率の上限をリソース制限テーブル155に登録されたリソース上限に設定する(S274)。これにより、選択したプロセスに対するCPU時間の割り当てが減少し、単位時間当たりのディスク更新量またはメモリ更新量が小さくなると期待できる。

【0130】

マイグレーション制御部141は、各プロセスについて、現在の状態、CPU使用率、ディスク更新量およびメモリ更新量を、カーネル143のインタフェースを利用して取得する(S275)。マイグレーション制御部141は、取得した状態、CPU使用率、ディスク更新量およびメモリ更新量をプロセステーブル154に上書き保存する。マイグレーション制御部141は、ゲストOS122全体の単位時間当たりのディスク更新量とメモリ更新量を、カーネル143のインタフェースを利用して取得する(S276)。ただし、マイグレーション制御部141は、ステップS275で取得したプロセス毎のディスク更新量とメモリ更新量から、ゲストOS122全体の値を算出してもよい。

【0131】

図21は、プロセス監視の手順例を示すフローチャート(続き)である。

マイグレーション制御部141は、ステップS276で取得した全体ディスク更新量が閾値以下であるか判断する(S277)。全体ディスク更新量が閾値以下の場合にはステップS278に処理が進み、閾値より大きい場合はステップS279に処理が進む。全体ディスク更新量が閾値以下の場合、マイグレーション制御部141は、状態テーブル151のディスク更新低下フラグを「1」に設定する(S278)。そして、ステップS280に処理が進む。全体ディスク更新量が閾値より大きい場合、マイグレーション制御部141は、状態テーブル151のディスク更新低下フラグを「0」に設定する(S279)。

【0132】

マイグレーション制御部141は、ステップS276で取得した全体メモリ更新量が閾値以下であるか判断する(S280)。全体メモリ更新量が閾値以下の場合にはステップS281に処理が進み、閾値より大きい場合はステップS282に処理が進む。全体メモリ更新量が閾値以下の場合、マイグレーション制御部141は、状態テーブル151のメモリ更新低下フラグを「1」に設定する(S281)。そして、ステップS283に処理が進む。全体メモリ更新量が閾値より大きい場合、マイグレーション制御部141は、状態テーブル151のメモリ更新低下フラグを「0」に設定する(S282)。

【0133】

なお、ステップS277の閾値は、例えば、ユーザまたはデータセンタ10の管理者によって予め定義される。ステップS277の閾値は、前述のステップS219，S226で使用する閾値と同じでもよい。また、ステップS280の閾値は、例えば、ユーザまたはデータセンタ10の管理者によって予め定義される。ステップS280の閾値は、前述

【0134】

マイグレーション制御部141は、更にCPU使用率を制限可能なプロセスが存在するか判断する(S283)。更にCPU使用率を制限可能なプロセスは、重要プロセスリスト156に登録されておらず、かつ、現在のCPU使用率がリソース制限テーブル155に登録されたリソース上限を超えているプロセスである。該当するプロセスが存在する場合はステップS285に処理が進み、存在しない場合はステップS284に処理が進む。

【0135】

マイグレーション制御部141は、状態テーブル151のディスク更新低下フラグとメモリ更新低下フラグの両方を「1」に設定する(S284)。これは、単位時間当たりの

10

20

30

40

50

ディスク更新量またはメモリ更新量が依然として大きいものの、差分データの送信を打ち切って後述するマイグレーション完了のプロセスを開始することを意味する。

【0136】

マイグレーション制御部141は、ディスク更新低下フラグとメモリ更新低下フラグの両方が「1」であるか判断する(S285)。2つのフラグが「1」である場合は、プロセス監視のプロセスが終了し、少なくとも一方のフラグが「0」である場合はステップS272に処理が進む。なお、2つのフラグが「1」になると、前述のステップS226、S256の判断が「YES」になり、ディスクデータ送信とメモリデータ送信のプロセスも終了する。また、後述するマイグレーション完了のプロセスが開始される。

【0137】

図22は、マイグレーション完了の手順例を示すフローチャートである。

マイグレーション制御部141は、カーネル143に対してサスペンド状態に遷移するよう指示する(S121)。サスペンド状態では、ゲストOS122がシャットダウンせずに、CPU101による仮想マシン121のプログラムの実行が停止する。

【0138】

サスペンド指示を受けて、カーネル143は、現在の仮想マシン121の状態を示す中断ファイルを作成し、HDD103の退避領域133に保存する(S122)。中断ファイルには、CPU101のレジスタ101aに記憶されている状態データ、RAM102に記憶されているメモリデータ、その他の周辺デバイスの状態を示すデバイスデータなどが含まれる。また、ゲストOS122は、ディスク更新情報152を参照して、仮想マシン221に未送信の差分データを含む記憶領域のスナップショットを作成しておく。また、カーネル143は、すぐにサスペンド状態を解除するよう設定しておく。

【0139】

カーネル143は、サスペンド状態に遷移した後、ステップS122の設定に従ってサスペンド状態を解除する(S123)。このとき、ステップS122で作成した中断ファイルを用いて、サスペンド直前のCPU101の状態などが復元される。これにより、ゲストOS122内のマイグレーション制御部141の処理が再開する。

【0140】

マイグレーション制御部141は、3つの監視処理を停止させる(S124)。3つの監視処理には、更新監視部142によるディスク監視、更新監視部142によるメモリ監視、マイグレーション制御部141によるプロセス監視が含まれる。

【0141】

マイグレーション制御部141は、更新監視部142にメモリ更新を問い合わせ、更新されたページのページアドレスを取得する。そして、マイグレーション制御部141は、退避領域133に記憶された中断ファイルから、更新されたページのメモリデータを抽出する。また、マイグレーション制御部141は、中断ファイルから、CPU101の状態データや周辺デバイスのデバイスデータなどを抽出する。また、マイグレーション制御部141は、ステップS122で作成されたスナップショットから、更新されたセクタのディスクデータを抽出する。マイグレーション制御部141は、抽出した上記のデータを含む最終データを、ユーザネットワーク11、インターネット30およびユーザネットワーク21を介して仮想マシン221に送信する(S125)。

【0142】

マイグレーション制御部141は、状態テーブル151の移行済フラグを「1」に設定する(S126)。マイグレーション制御部241は、仮想マシン121から受信した最終データに含まれる各種データを、HDD203の適切な記憶領域に保存する(S127)。最終データに含まれるディスクデータは、セクタアドレスが示すHDD203内の記憶領域に上書き保存される。最終データに含まれるメモリデータ、CPU101の状態データ、周辺デバイスのデバイスデータなどは、退避領域233に保存される。

【0143】

マイグレーション制御部241は、退避領域233に記憶されたデータに基づいて、仮

10

20

30

40

50

想マシン 1 2 1 の状態を再現するための中断ファイルを生成する (S 1 2 8)。マイグレーション制御部 2 4 1 は、制御情報記憶部 2 4 2 に状態テーブル (状態テーブル 1 5 1 と同様のもの) を生成し、移行済フラグを「 1 」に設定する (S 1 2 9)。

【 0 1 4 4 】

マイグレーション制御部 2 4 1 は、ユーザネットワーク 2 1 およびインターネット 3 0 を介して、データセンタ 1 0 の管理サーバ 1 3 に対して仮想マシン 1 2 1 の停止要求を送信する (S 1 3 0)。停止要求を受信すると、管理サーバ 1 3 は、管理ネットワーク 1 2 を介して仮想マシン 1 2 1 に停止指示を送信する (S 1 3 1)。仮想マシン 1 2 1 は、ゲスト OS 1 2 2 をシャットダウンさせて仮想マシン 1 2 1 を停止させる (S 1 3 2)。

【 0 1 4 5 】

マイグレーション制御部 2 4 1 は、起動情報領域 2 3 2 に記憶された起動情報を編集する (S 1 3 3)。編集された起動情報は、仮想マシン 2 2 1 の起動時にパーティション 2 3 1 からシステムデータがロードされることを示す。また、マイグレーション制御部 2 4 1 は、退避領域 2 3 3 に記憶された中断ファイルを用いて途中状態から処理が開始されるように設定しておく。マイグレーション制御部 2 4 1 は、ゲスト OS 2 2 2 のカーネルに再起動を指示する (S 1 3 4)。ゲスト OS 2 2 2 がシャットダウンして仮想マシン 2 2 1 が再起動すると、起動情報に従ってパーティション 2 3 1 からシステムデータがロードされる。これにより、仮想マシン 2 2 1 上ではゲスト OS 1 2 2 が実行される。また、退避領域 2 3 3 に記憶された中断ファイルに基づいて、移行元の CPU 1 0 1 や RAM 1 0 2 の状態が、CPU 2 0 1 や RAM 2 0 2 上に再現される。

【 0 1 4 6 】

図 2 3 は、マイグレーション完了の手順例を示すフローチャート (続き) である。

仮想マシン 2 2 1 上で起動したマイグレーション制御部 1 4 1 は、制御情報記憶部 2 4 2 に記憶された状態テーブルを参照し、移行済フラグが「 1 」であるか判断する (S 1 3 5)。移行済フラグが「 1 」の場合はステップ S 1 3 6 に処理が進み、移行済フラグが「 0 」の場合はマイグレーション制御部 1 4 1 の処理が終了する。

【 0 1 4 7 】

移行済フラグが「 1 」の場合、マイグレーション制御部 1 4 1 は、カーネル 1 4 3 のインタフェースを利用して、前述のステップ S 2 7 4 で設定した CPU 使用率の上限を解除する (S 1 3 6)。マイグレーション制御部 1 4 1 は、制御情報記憶部 1 4 4 に記憶された状態テーブル 1 5 1 などの制御情報を削除する (S 1 3 7)。マイグレーション制御部 1 4 1 は、ゲスト OS 2 2 2 のシステムデータが記憶されたパーティション 2 3 5 を、ゲスト OS 1 2 2 から認識されないように解放する (S 1 3 8)。

【 0 1 4 8 】

マイグレーション制御部 1 4 1 は、ユーザネットワーク 2 1 およびインターネット 3 0 を介して、データセンタ 1 0 の管理サーバ 1 3 に対して仮想マシン 1 2 1 の削除要求を送信する (S 1 3 9)。削除要求を受信すると、管理サーバ 1 3 は、仮想マシン 1 2 1 に対するホストサーバ 1 0 0 のリソースの割り当てを解除することで、データセンタ 1 0 から仮想マシン 1 2 1 を削除する (S 1 4 0)。これにより、パーティション 1 3 1 , 1 3 4 が解放されて他の仮想マシンが使用可能になる。

【 0 1 4 9 】

第 2 の実施の形態の情報処理システムによれば、移行先のホストサーバ 2 0 0 で仮想マシン 2 2 1 が起動し、仮想マシン 2 2 1 上でゲスト OS 2 2 2 が一時的に実行される。仮想マシン 1 2 1 に関するデータのコピーは、仮想マシン 1 2 1 と仮想マシン 2 2 1 との間の通信として行うことができる。よって、データセンタ 1 0 の外部からの管理ネットワーク 1 2 への直接アクセスやデータセンタ 2 0 の外部からの管理ネットワーク 2 2 への直接アクセスが制限されていても、管理ネットワーク 1 2 , 2 2 を迂回しユーザネットワーク 1 1 , 2 1 を用いて仮想マシン 1 2 1 に関するデータをコピーできる。

【 0 1 5 0 】

そして、仮想マシン 2 2 1 の起動情報を書き換えることで、仮想マシン 2 2 1 が再起動

10

20

30

40

50

すると、ゲストOS 2 2 2 に代えてゲストOS 1 2 2 が仮想マシン 2 2 1 上で実行される。これにより、異なるサービス事業者が管理するデータセンタ 1 0 , 2 0 の間でも、仮想マシン 1 2 1 のマイグレーションを円滑に行うことができる。

【 0 1 5 1 】

また、移行元の仮想マシン 1 2 1 においてメモリデータやディスクデータの更新頻度が高い場合、少なくとも一部のプロセスのCPU使用率を制限することで、メモリデータやディスクデータの更新頻度が低下する。これにより、送信する差分データの量を低減することができ、インターネット 3 0 の通信遅延が大きい場合であっても、仮想マシン 1 2 1 のマイグレーションに要する時間を短縮することができる。また、CPU使用率を制限しない重要プロセスを定義しておくことで、業務への影響を軽減できる。

10

【 0 1 5 2 】

なお、前述のように、第 1 の実施の形態の情報処理は、情報処理装置 1 , 2 にプログラムを実行させることで実現することができる。また、第 2 の実施の形態の情報処理は、管理サーバ 1 3 , 2 3、クライアント 3 1、ホストサーバ 1 0 0 , 1 0 0 a , 2 0 0 , 2 0 0 a などにプログラムを実行させることで実現することができる。

【 0 1 5 3 】

プログラムは、コンピュータ読み取り可能な記録媒体（例えば、記録媒体 1 1 3）に記録しておくことができる。記録媒体としては、例えば、磁気ディスク、光ディスク、光磁気ディスク、半導体メモリなどを使用できる。磁気ディスクには、FDおよびHDDが含まれる。光ディスクには、CD、CD-R（Recordable）/RW（Rewritable）、DVDおよびDVD-R/RWが含まれる。プログラムは、可搬型の記録媒体に記録されて配布されることがある。その場合、可搬型の記録媒体からHDDなどの他の記録媒体（例えば、HDD 1 0 3）にプログラムをコピーして実行してもよい。

20

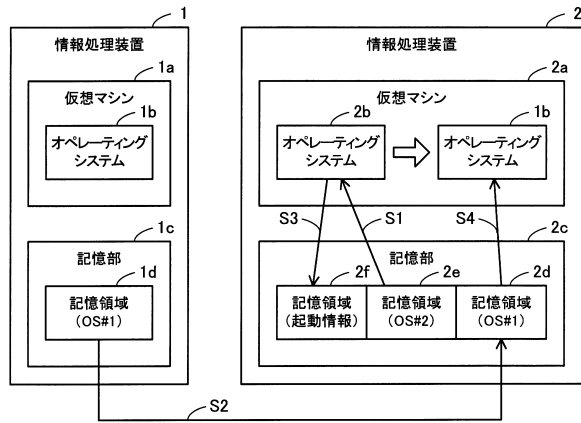
【符号の説明】

【 0 1 5 4 】

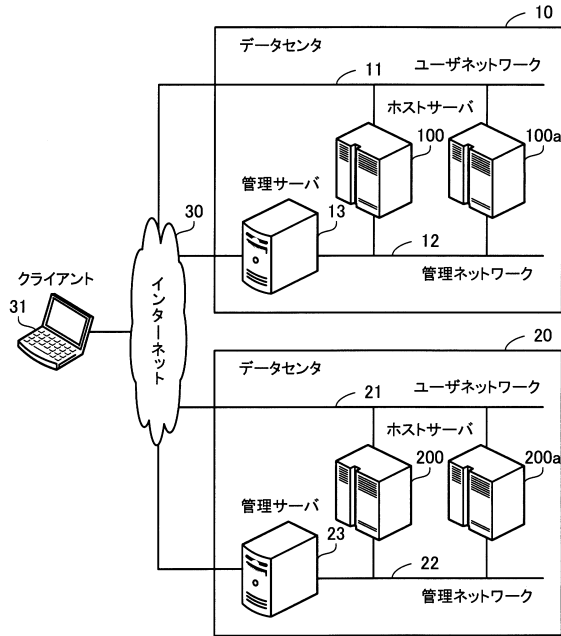
- 1 , 2 情報処理装置
- 1 a , 2 a 仮想マシン
- 1 b , 2 b オペレーティングシステム
- 1 c , 2 c 記憶部
- 1 d , 2 d , 2 e , 2 f 記憶領域

30

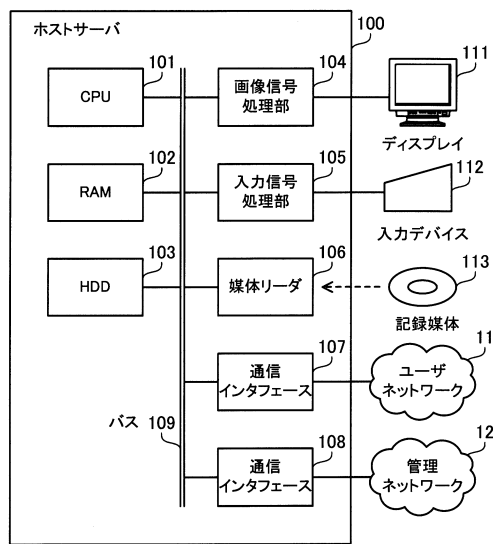
【図1】



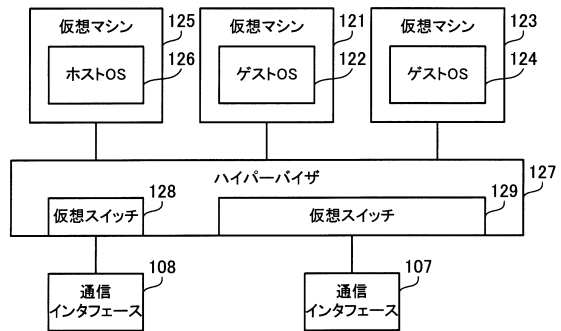
【図2】



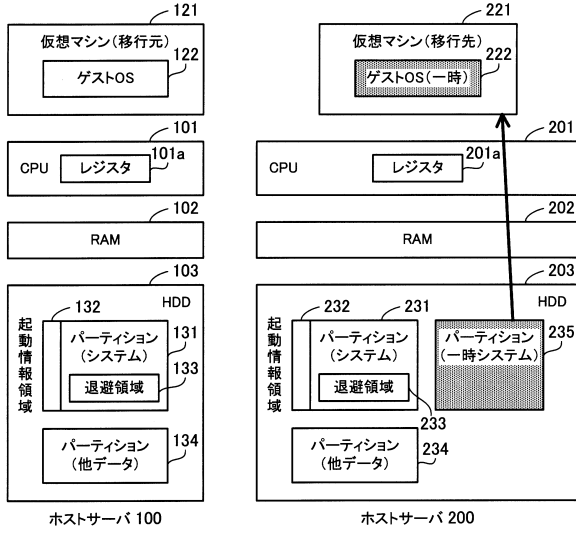
【図3】



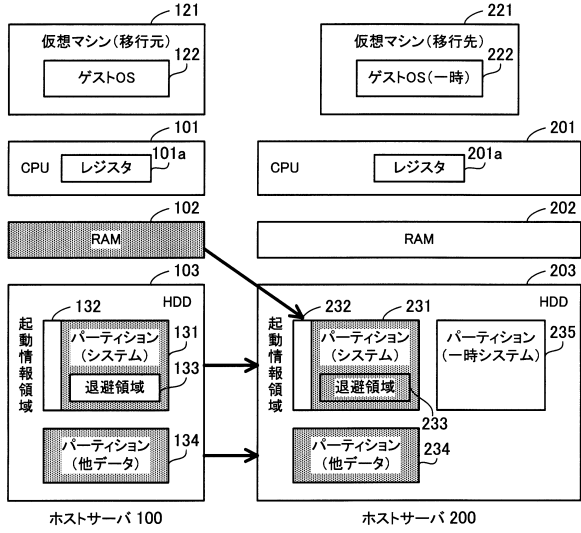
【図4】



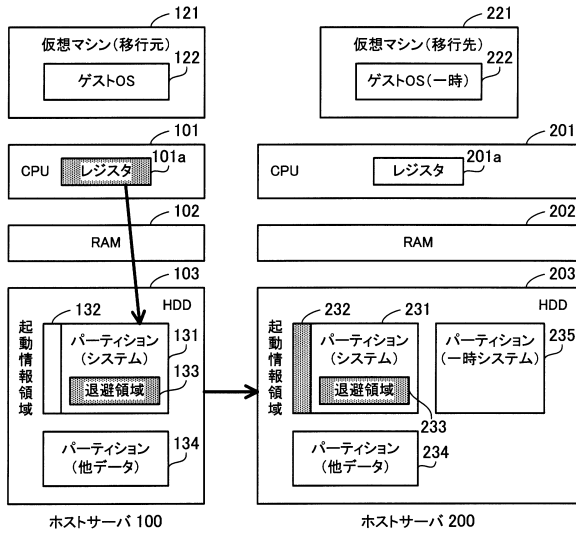
【図5】



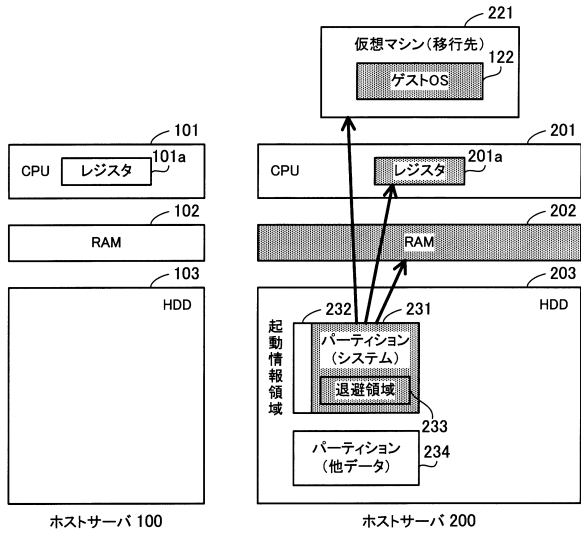
【図6】



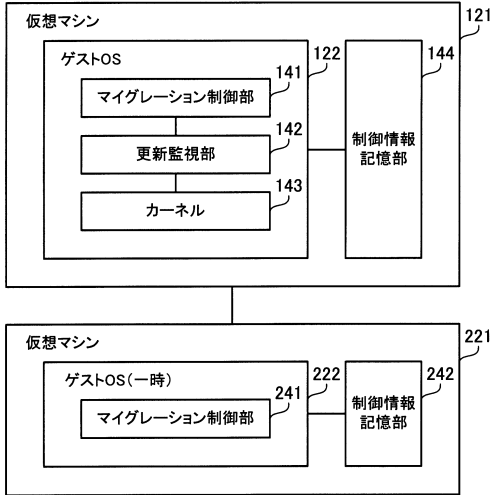
【図7】



【図8】



【図9】



【図11】

制御情報記憶部

ディスク更新情報

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
1000	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1010	0	0	1	1	1	1	1	1	0	0	0	0	0	0	0	0
1020	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0
⋮																

メモリ更新情報

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
00	0	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
20	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
⋮																

【図10】

制御情報記憶部

状態テーブル

パラメータ名	フラグ
移行済	0
ディスク更新低下	0
メモリ更新低下	0

【図12】

制御情報記憶部

プロセステーブル

プロセスID	プロセス名	開始状態	状態	CPU使用率	ディスク更新量	メモリ更新量
1	プロセスA	実行中	実行中	40 %	51216 Byte	30323 Byte
100	プロセスB	実行中	実行中	25 %	19531 Byte	1334 Byte
200	プロセスC	実行中	終了	10 %	9045 Byte	1214 Byte
300	プロセスD	一時停止	一時停止	5 %	15642 Byte	123 Byte
400	プロセスE	実行中	実行中	5 %	9 Byte	44 Byte

【図13】

制御情報記憶部

リソース制限テーブル

プロセス名	リソース上限
プロセスA	100 %
プロセスB	20 %
プロセスC	100 %
プロセスD	10 %
プロセスE	10 %

重要プロセスリスト

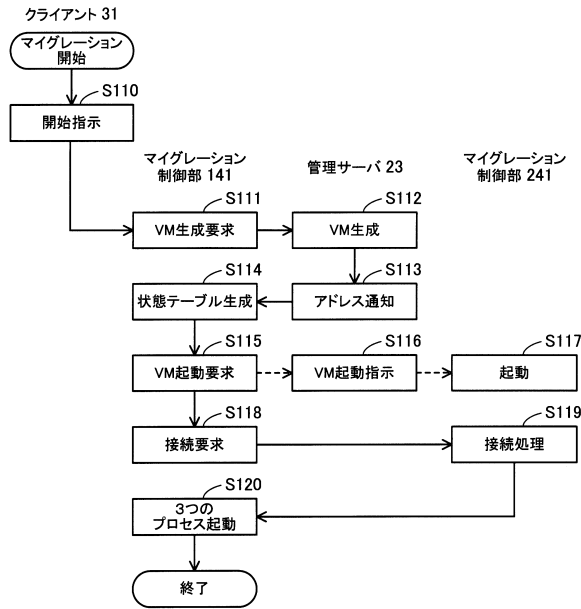
プロセスA
プロセスC

【図14】

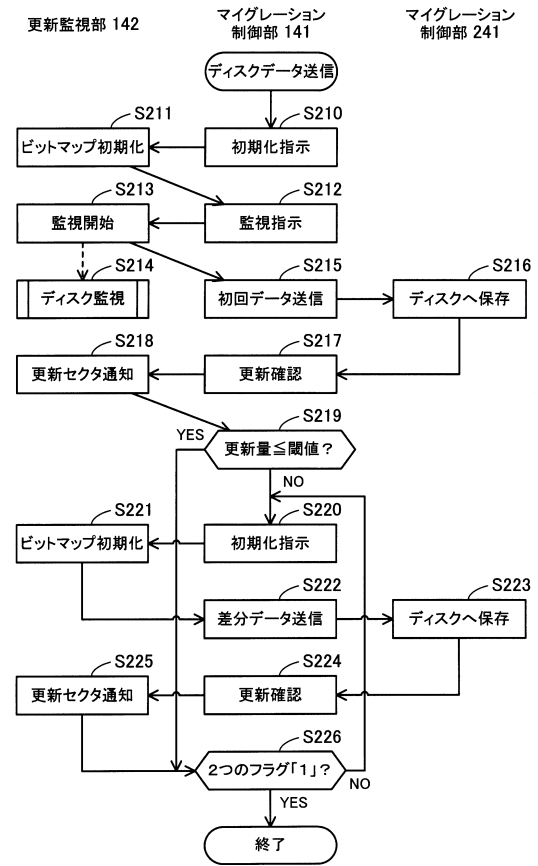
仮想マシン生成要求

パラメータ名	設定値
CPU数	2
メモリ容量	2 GByte
システムディスク容量	16 GByte
他データディスク容量	16 GByte
ネットワーク設定	インターネット
初期OS	マイグレーション用

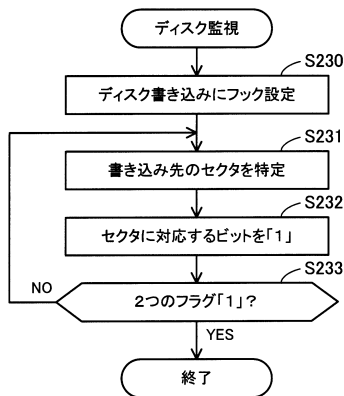
【図15】



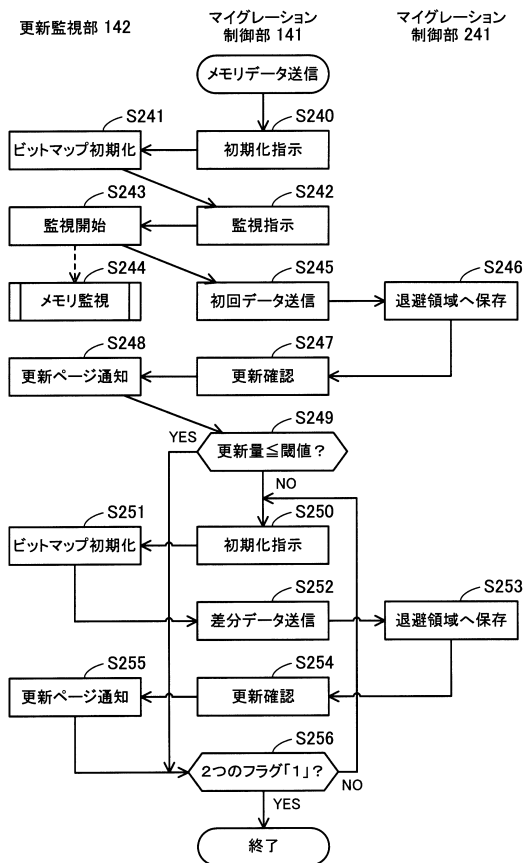
【図16】



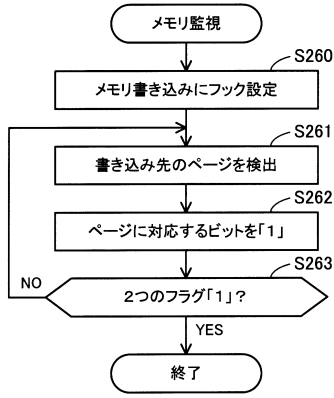
【図17】



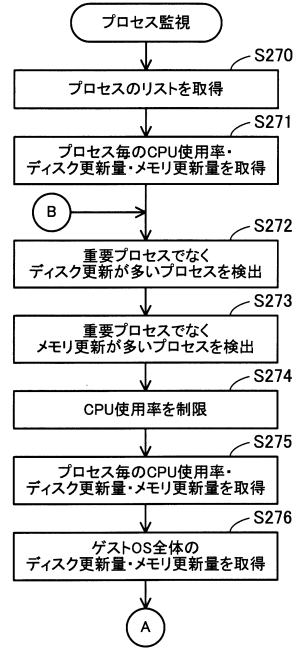
【図18】



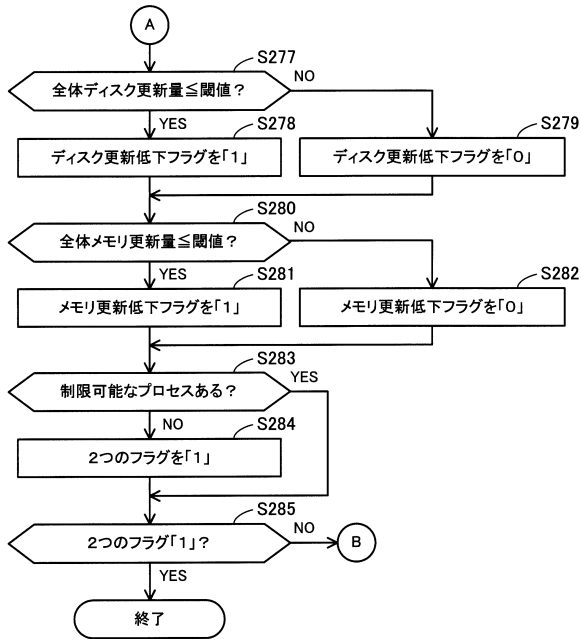
【図19】



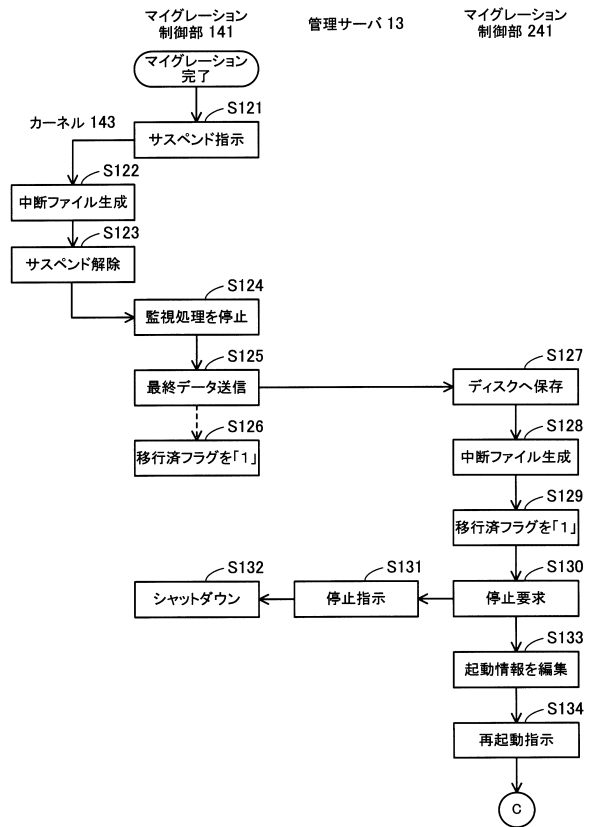
【図20】



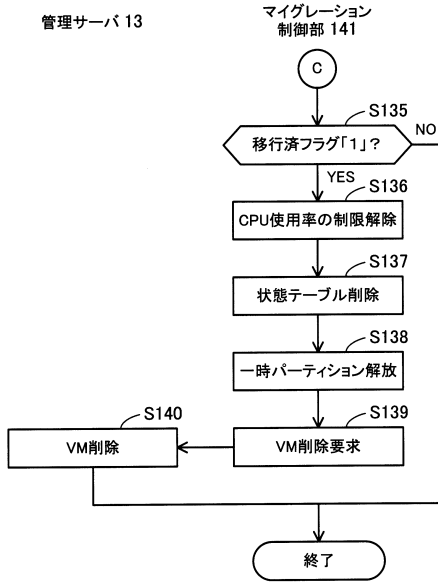
【図21】



【図22】



【図 23】



フロントページの続き

- (72)発明者 寺嶋 直矢
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 今枝 一英
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 原 忠

- (56)参考文献 特開2012-221321(JP,A)
国際公開第2009/034760(WO,A1)
特開2013-250950(JP,A)
特開2014-102724(JP,A)

- (58)調査した分野(Int.Cl., DB名)
- | | |
|---------|-------------------|
| G 0 6 F | 9 / 4 4 |
| G 0 6 F | 9 / 4 5 5 |
| G 0 6 F | 9 / 4 6 - 9 / 5 4 |