

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6277069号
(P6277069)

(45) 発行日 平成30年2月7日(2018.2.7)

(24) 登録日 平成30年1月19日(2018.1.19)

(51) Int. Cl. F 1
G 0 6 F 9/46 (2006.01) G 0 6 F 9/46 3 5 0
G 0 6 F 11/20 (2006.01) G 0 6 F 11/20

請求項の数 8 (全 23 頁)

(21) 出願番号	特願2014-124600 (P2014-124600)	(73) 特許権者	000004226
(22) 出願日	平成26年6月17日 (2014.6.17)		日本電信電話株式会社
(65) 公開番号	特開2016-4433 (P2016-4433A)		東京都千代田区大手町一丁目5番1号
(43) 公開日	平成28年1月12日 (2016.1.12)	(74) 代理人	100127535
審査請求日	平成26年12月4日 (2014.12.4)		弁理士 豊田 義元
審判番号	不服2016-19308 (P2016-19308/J1)	(74) 代理人	100159190
審判請求日	平成28年12月26日 (2016.12.26)		弁理士 渡部 比呂志
		(72) 発明者	山登 庸次
			東京都千代田区大手町一丁目5番1号 日 本電信電話株式会社内
		(72) 発明者	西澤 幸久
			東京都千代田区大手町一丁目5番1号 日 本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 仮想機器管理装置、仮想機器管理方法及び仮想機器管理プログラム

(57) 【特許請求の範囲】

【請求項 1】

障害の生じた物理機器を検出する検出部と、

前記障害の生じた物理機器以外の稼働中の物理機器のうち物理資源の空き容量のある物理機器を複数特定し、前記稼働中の物理機器の空き容量が無い場合は、予備の物理機器を特定し、特定した複数の前記稼働中又は予備の物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する選択部と、

選択された前記物理機器それぞれに前記仮想機器のそれぞれの再配置を、並行して行うように、物理機器への仮想機器の作成を制御するクラウドコントローラに依頼する依頼部と

を備えたことを特徴とする仮想機器管理装置。

【請求項 2】

前記物理資源には、物理機器に障害の生じた場合にのみ使用される確保領域を含み、

前記選択部は、前記障害の生じた物理機器以外の物理機器のうち確保領域を含んだ前記物理資源の空き容量のある物理機器を複数特定し、特定した複数の前記物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する

ことを特徴とする請求項 1 に記載の仮想機器管理装置。

【請求項 3】

前記選択部は、特定した複数の物理機器に順序付けを行い、前記障害の生じた物理機器に配置された仮想機器それぞれの再配置先として、前記順序に基づいて選択した物理機器

の物理資源を選択する処理を繰り返すことを特徴とする請求項 1 又は 2 に記載の仮想機器管理装置。

【請求項 4】

前記選択部は、特定した複数の物理機器の物理資源の空き容量が多い順に、前記特定した複数の物理機器に順序付けを行うことを特徴とする請求項 3 に記載の仮想機器管理装置。

【請求項 5】

仮想機器を新規に作成する依頼を受付ける受付部と、

新規に作成する依頼を受付けた前記仮想機器を配置する物理機器を、物理機器に障害の生じた場合のみ使用される確保領域を含まない前記物理資源の空き容量に基づいて選択する作成部とを更に備え、

前記依頼部は、前記作成部によって選択された物理機器に前記仮想機器の作成を依頼することを特徴とする請求項 2 ～ 4 のいずれか一つに記載の仮想機器管理装置。

【請求項 6】

前記作成部は、物理機器のうち前記確保領域を含まない前記物理資源の空き容量の最も多い物理機器を特定し、特定した前記物理機器の前記確保領域を含まない物理資源の一部を前記新規に作成する依頼を受付けた仮想機器のうち何れかの仮想機器の配置先として選択する処理を、前記新規に作成する依頼を受付けた全ての仮想機器の配置先を選択するまで繰り返すことを特徴とする請求項 5 に記載の仮想機器管理装置。

【請求項 7】

仮想機器管理装置で実行する仮想機器管理方法であって、

前記仮想機器管理装置が、

障害の生じた物理機器を検出する検出工程と、

前記障害の生じた物理機器以外の稼働中の物理機器のうち物理資源の空き容量のある物理機器を複数特定し、前記稼働中の物理機器の空き容量が無い場合は、予備の物理機器を特定し、特定した複数の前記稼働中又は予備の物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する選択工程と、

選択された前記物理機器それぞれに前記仮想機器のそれぞれの再配置を、並行して行うように、物理機器への仮想機器の作成を制御するクラウドコントローラに依頼する依頼工程と

を含んだことを特徴とする仮想機器管理方法。

【請求項 8】

障害の生じた物理機器を検出する検出手順と、

前記障害の生じた物理機器以外の稼働中の物理機器のうち物理資源の空き容量のある物理機器を複数特定し、前記稼働中の物理機器の空き容量が無い場合は、予備の物理機器を特定し、特定した複数の前記稼働中又は予備の物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する選択手順と、

選択された前記物理機器それぞれに前記仮想機器のそれぞれの再配置を、並行して行うように、物理機器への仮想機器の作成を制御するクラウドコントローラに依頼する依頼手順と

をコンピュータに実行させることを特徴とする仮想機器管理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想機器管理装置、仮想機器管理方法及び仮想機器管理プログラムに関する。

【背景技術】

【0002】

IaaS (Infrastructure as a Service) 型クラウドサービスの実施例として、A

10

20

30

40

50

m a z o n E l a s t i c C o m p u t e C l o u d (web site, <http://aws.amazon.com/ec2>)、R a c k s p a c e C l o u d S e r v e r (web site, <http://www.rackspacecloud.com/cloud-hosting-products/servers/>)がある。

【 0 0 0 3 】

I a a S 型クラウドサービスの基盤として、A m a z o n はプロプライエタリなプラットフォームを用いているが、R a c k S p a c e は O p e n S o u r c e の O p e n S t a c k (<http://www.openstack.org/>)を用いている。

【 0 0 0 4 】

しかし、O p e n S t a c k 等の I a a S 基盤は、仮想リソースの管理を行うプリミティブな A P I (Application Programming Interface) 提供がターゲットの中心であり、物理機器の管理はスコープ外であるため、事業者がクラウドサービスを提供する際は考慮が必要である。

10

【 0 0 0 5 】

具体的には、仮想リソースが動作する物理機器が故障した際の復旧は、O p e n S t a c k は特にサポートしていなく、サービス事業者にて対策が必要である。市中で採用されている方法として、H i g h A v a i l a b i l i t y クラスソフトウェアの P a c e m a k e r 等を用いて H A 構成を構築し、物理機器故障時はフェールオーバーする方法がある。

【 先行技術文献 】

【 非特許文献 】

20

【 0 0 0 6 】

【 非特許文献 1 】 Pacemaker web site、[平成26年5月30日検索]、インターネット (URL : <http://www.linux-ha.org/wiki/Pacemaker/>)

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 7 】

しかしながら、上記の従来技術では、仮想機器を復旧するまでの時間が長くなる場合があるという問題がある。

【 0 0 0 8 】

仮想機器が動作する物理機器が故障した際に、復旧する手段として、P a c e m a k e r 等の H A (High Availability) クラスソフトウェアを用いる方法がある。この方法では、N - A c t (Active)、M - S b y (standby) の冗長化構成を取り、A c t i v e の物理機器が故障した際に、S t a n d b y の物理機器に自動切り替え (フェールオーバー) を行う。なお、N 及び M は任意の整数である。例えば、O p e n S t a c k 上で仮想ルータが動作している物理機器であれば、新しく A c t 機になった物理機器の N e u t r o n エージェントは、O p e n S t a c k D B 上の構成情報を用いて、当該論理ホスト上に存在すべき仮想ルータを再構築する。再構築時間は、収容していた仮想ルータ数に比例した時間を要する。これは、仮想マシン等の場合も同様で、N o v a の機能を用いて再構築される。このため、収容している仮想機器が多い場合は、全仮想機器が復旧するまでの時間が長くなり、サービス断時間が1時間超に及ぶなど、サービス断時間が大きい問題がある。

30

40

【 0 0 0 9 】

また、H A クラスソフトウェアの P a c e m a k e r では、パケットサイズの制限から H A クラスの最大構成が8台程度までしか設定できない。H A クラスは N - A c t、M - S b y 構成であるため、クラスあたり最低1台の S t a n d b y 機が必要である。そのため、予備機比率が1/8以上となり、物理機器の利用効率が悪い。

【 0 0 1 0 】

このように、H A クラスソフトウェアを用いた方法では、仮想リソースを S t a n d b y 機に復旧するまでの時間が長時間かかる問題や、S t a n d b y 機を準備する必要があり物理機器数が増える問題がある。

50

【 0 0 1 1 】

開示の技術は、上述に鑑みてなされたものであって、仮想機器を復旧するまでの時間を短縮することを目的とする。

【課題を解決するための手段】

【 0 0 1 2 】

本願の開示する仮想機器管理装置は、検出部と、選択部と、依頼部とを有する。検出部は、障害の生じた物理機器を検出する。選択部は、前記障害の生じた物理機器以外の物理機器のうち物理資源の空き容量のある物理機器を複数特定し、特定した複数の前記物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する。依頼部は、選択された前記物理機器それぞれに前記仮想機器のそれぞれの再配置を依頼する。

10

【 0 0 1 3 】

また、本願の開示する仮想機器管理方法は、検出工程と、選択工程と、依頼工程とを含む。検出工程は、障害の生じた物理機器を検出する。選択工程は、前記障害の生じた物理機器以外の物理機器のうち物理資源の空き容量のある物理機器を複数特定し、特定した複数の前記物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する。依頼工程は、選択された前記物理機器それぞれに前記仮想機器のそれぞれの再配置を依頼する。

【 0 0 1 4 】

また、本願の開示する仮想機器管理プログラムは、検出手順と、選択手順と、依頼手順とを有する。検出手順は、障害の生じた物理機器を検出する。選択手順は、前記障害の生じた物理機器以外の物理機器のうち物理資源の空き容量のある物理機器を複数特定し、特定した複数の前記物理機器の物理資源を、前記障害の生じた物理機器に配置された仮想機器の再配置先として選択する。依頼手順は、選択された前記物理機器それぞれに前記仮想機器のそれぞれの再配置を依頼する。

20

【発明の効果】

【 0 0 1 5 】

開示する仮想機器管理装置の一つの態様によれば、仮想機器を復旧するまでの時間を短縮することができるという効果を奏する。

【図面の簡単な説明】

30

【 0 0 1 6 】

【図 1】図 1 は、第 1 の実施形態に係る仮想機器管理システムの構成の一例を示す図である。

【図 2】図 2 は、仮想機器配置スケジューラ機能部による仮想機器の作成処理を説明するための図である。

【図 3】図 3 は、仮想機器配置スケジューラ機能部による仮想機器の再配置処理を説明するための図である。

【図 4】図 4 は、仮想機器管理システムにおける仮想機器を作成する処理動作を説明するための図である。

【図 5】図 5 は、仮想機器管理システムにおける仮想機器を再配置する処理動作を説明するための図である。

40

【図 6】図 6 は、仮想機器管理装置が実現する仮想機器配置スケジューラ DB 及び仮想機器配置スケジューラ機能部を説明するための図である。

【図 7】図 7 は、仮想機器配置情報テーブルのデータ構造の一例を示す図である。

【図 8】図 8 は、物理資源情報テーブルのデータ構造の一例を示す図である。

【図 9】図 9 は、配置先選択部による処理動作を説明するための図である。

【図 10】図 10 は、障害検出部による処理動作を説明するための図である。

【図 11】図 11 は、再配置先選択部による処理動作を説明するための図である。

【図 12】図 12 は、仮想機器の作成を要求された場合の仮想機器配置スケジューラ機能部による処理手順を示すフローチャートである。

50

【図13】図13は、物理機器に障害が発生した場合の仮想機器配置スケジューラ機能部による処理手順を示すフローチャートである。

【図14】図14は、仮想機器管理プログラムを実行するコンピュータを示す図である。

【発明を実施するための形態】

【0017】

以下に、開示する仮想機器管理装置、仮想機器管理方法及び仮想機器管理プログラムの実施形態について、図面に基づいて詳細に説明する。なお、本実施形態により開示する発明が限定されるものではない。

【0018】

(第1の実施形態)

図1は、第1の実施形態に係る仮想機器管理システムの構成の一例を示す図である。図1に示すように、仮想機器管理システムは、ユーザ端末101、物理機器103a、物理機器103b、物理機器103c、クラウドコントローラ108、及び仮想機器管理装置109を有する。ここで言う「物理機器」とは、仮想機器を生成可能な物理サーバ、ストレージ装置、及びネットワーク機器等である。なお、物理機器103a、物理機器103b及び物理機器103cを区別しない場合には、物理機器103と記載する。また、仮想機器管理システムが有する物理機器103の数は図1に示す数に限定されるものではなく、任意に変更可能である。

【0019】

ユーザ端末101は、ユーザが利用する端末であり、ユーザの指示に応じて仮想機器の作成を仮想機器管理装置109に要求する。物理機器103は、クラウドコントローラ108から仮想機器の作成や削除依頼を受け、実際の仮想機器を作成したり削除したりする。例えば、物理機器103は、仮想機器を作成する指示をクラウドコントローラ108から受け、仮想機器を作成する。

【0020】

例えば、物理機器103aは、図示しない仮想ボリューム制御部を有し、仮想ボリューム104aと、仮想ボリューム105aとを作成する。また、物理機器103bは、図示しない仮想ネットワーク制御部を有し、仮想L2ネットワーク104bと、仮想ルータ105bと、仮想ロードバランサ106bとを作成する。なお、仮想ネットワーク制御部は、例えば「Neutron」によって実現される。また、物理機器103cは、図示しない仮想マシン制御部を有し、仮想マシン104cと、仮想マシン105cとを作成する。なお、仮想マシン制御部は、例えば「Nova」によって実現される。

【0021】

また、物理機器103の稼働状態には、「稼働中」、「予備」及び「故障中/メンテ中」3つの状態がある。「稼働中」は、物理機器が稼働中であることを示す。「予備」は、物理機器が予備系として設けられ稼働中ではないことを示す。「故障中/メンテ中」は、物理機器が故障中やメンテナンス中であることを示す。なお、仮想機器管理システムにおいて、「予備」の物理機器が設けられなくてもよい。

【0022】

また、物理機器103には、物理資源の容量に応じて、仮想機器を配置するために利用可能な物理資源の容量が定義される。ここで、物理資源には、例えば、物理メモリ、CPU (Central Processing Unit)、ネットワークポートなどが含まれる。なお、仮想マシンは、フレーバー (仮想マシンのスペック指定) に応じてメモリサイズが異なるため、作成する仮想マシンに応じて利用される物理資源の容量は異なる。しかしながら説明の便宜上、以下では、全ての仮想機器1つにつき、使用される物理資源の容量が同じであるものと仮定する。そして、1つの仮想機器を配置するために使用される物理資源の容量を1単位とし、「1スペース」と呼ぶ。言い換えると、1スペースには、1つの仮想機器を配置可能であり、1つの仮想機器を作成する場合には、いずれかの物理機器のスペースが1つ消費される。

【0023】

10

20

30

40

50

また、物理機器内のスペースの状態は、「空き」、「使用中」、及び「障害用バッファ」の3種類で管理されるものとする。ここで、「空き」は、仮想機器が配置されていないスペースであることを示す。「使用中」は、仮想機器が配置されているスペースであることを示す。「障害用バッファ」は、障害復旧用に確保されたスペースであることを示す。

【0024】

また、物理機器103aは、高可用ソフトウェア107aを備えている。同様に、物理機器103bは、高可用ソフトウェア107bを備えており、物理機器103cは、高可用ソフトウェア107cを備えている。なお、高可用ソフトウェア107a~107cを区別しない場合には高可用ソフトウェア107と記載する。この高可用ソフトウェア107には、例えば「Pacemaker」等が利用できる。高可用ソフトウェア107は、物理機器103の障害を検知し、仮想機器管理装置109に物理機器の障害を通知する。かかる場合、物理機器103は、仮想機器を再配置させる指示をクラウドコントローラ108から受け、障害の生じた物理機器103に配置された仮想機器を再配置する。

10

【0025】

なお、「Pacemaker」は、信頼性の高い故障検出メカニズムを備えており、スプリットブレイン対策が確立している。「Pacemaker」は、スプリットブレイン状態（孤立状態）を、Quorumモジュール等による多数決原理で検出する。

【0026】

クラウドコントローラ108は、物理機器103と仮想機器管理装置109とに接続されている。このクラウドコントローラ108は、CPU（Central Processing Unit）、メモリ、データ保持領域、及びネットワーク通信機能を有する装置である。クラウドコントローラ108は、仮想機器管理装置109からAPI（Application Programming Interface）経由で仮想機器の作成依頼を受け、受け取った作成依頼に基づいて、仮想機器の作成を物理機器103に指示する。例えば、クラウドコントローラ108は、OpenStack等である。

20

【0027】

仮想機器管理装置109は、ユーザ端末101と物理機器103とクラウドコントローラ108とに接続されている。仮想機器管理装置109は、CPU、メモリ、データ保持領域、及びネットワーク通信機能を有する装置であり、例えば、図1に示すように、仮想機器管理装置109は、仮想機器配置スケジューラDB（Data Base）110及び仮想機器配置スケジューラ機能部111を有する。

30

【0028】

仮想機器配置スケジューラDB110は、例えば、RAM（Random Access Memory）、フラッシュメモリ（Flash Memory）等の半導体メモリ素子、又は、ハードディスク、光ディスク等の記憶装置などである。仮想機器配置スケジューラDB110は、仮想機器配置情報及び物理資源情報を記憶する。仮想機器配置情報は、各仮想機器がどの物理機器上に配置されているかを示す情報である。物理資源情報は、各物理機器の稼働状態と物理機器が有する物理資源の空き容量とを示す情報である。なお、仮想機器配置情報の詳細については、図7を用いて後述し、物理資源情報の詳細については、図8を用いて後述する。

40

【0029】

仮想機器配置スケジューラ機能部111は、物理機器の稼働状態と物理機器が有する物理資源の使用状態とを参照して、ビジネス要件に応じた仮想機器を配置する。例えば、仮想機器配置スケジューラ機能部111は、ユーザ端末101から、仮想機器の作成を要求された場合、仮想機器の作成要求と仮想機器配置スケジューラDB110の情報とを用いて、仮想機器の作成を仲介する。ここで、仮想機器配置スケジューラ機能部111は、仮想マシンや仮想ルータ等の仮想機器を新規に作成する通常のオペレーション時に、仮想機器を配置する物理機器を決め、クラウドコントローラ108に物理機器を指定して仮想機器の作成を依頼する。

【0030】

50

図2は、仮想機器配置スケジューラ機能部111による仮想機器の作成処理を説明するための図である。図2では、3台の物理機器103a~103cに、仮想ルータであるLR11~LR16、LR21~LR26、及びLR31~LR36を仮想機器として作成する場合を示す。なお、3台の物理機器103a~103cはいずれも稼働中であるものとする。また、図2の例では、仮想機器が仮想ルータである場合を示すが、仮想機器は、仮想マシン等のその他の仮想機器であってもよい。

【0031】

図2に示すように、仮想機器配置スケジューラ機能部111は、仮想機器LR11~LR16の配置先として物理機器103aを選択し、仮想機器LR21~LR26の配置先として物理機器103bを選択し、仮想機器LR31~LR36の配置先として物理機器103cを選択する。そして、仮想機器配置スケジューラ機能部111は、クラウドコントローラ108に配置を依頼する。すなわち、仮想機器配置スケジューラ機能部111は、物理機器103aに仮想機器LR11~LR16を作成するようにクラウドコントローラ108に依頼する。また、仮想機器配置スケジューラ機能部111は、物理機器103bに仮想機器LR21~LR26を作成するようにクラウドコントローラ108に依頼し、物理機器103cに仮想機器LR31~LR36を作成するようにクラウドコントローラ108に依頼する。

【0032】

また、仮想機器配置スケジューラ機能部111は、例えば、いずれかの物理機器103に障害が生じた場合に、仮想機器配置スケジューラDB110の情報を用いて、仮想機器の再配置を仲介する。ここで、仮想機器配置スケジューラ機能部111は、高可用ソフトウェア107及びクラウドコントローラ108と連携することで障害復旧時に仮想機器を再配置する。図3は、仮想機器配置スケジューラ機能部111による仮想機器の再配置処理を説明するための図である。図3では、3台の物理機器103a~103cのうち、仮想機器としてLR21~LR26を配置する物理機器103bに障害が生じた場合を示す。なお、3台の物理機器103a~103cはいずれも稼働中であるものとする。

【0033】

図3に示すように、仮想機器配置スケジューラ機能部111は、物理機器103bに障害が生じたことを検出する。そして、仮想機器配置スケジューラ機能部111は、仮想機器配置スケジューラDB110の情報を用いて、仮想機器LR21~LR26の再配置先を決定する。図3に示す例では、仮想機器配置スケジューラ機能部111は、LR21、LR23、及びLR25の再配置先として物理機器103aを選択し、LR22、LR24、及びLR26の再配置先として物理機器103cを選択する。

【0034】

そして、仮想機器配置スケジューラ機能部111は、クラウドコントローラ108に再配置を依頼する。すなわち、仮想機器配置スケジューラ機能部111は、物理機器103aに仮想機器LR21、LR23、及びLR25を作成するようにクラウドコントローラ108に依頼する。また、仮想機器配置スケジューラ機能部111は、物理機器103cに仮想機器LR22、LR24、及びLR26を作成するようにクラウドコントローラ108に依頼する。

【0035】

この結果、物理機器103aは、仮想機器LR21、LR23、及びLR25を再構築し、仮想機器LR11~LR16に加えて、仮想機器LR21、LR23、及びLR25を配置する。また、物理機器103cは、仮想機器LR22、LR24、及びLR26を再構築し、仮想機器LR31~LR36に加えて、仮想機器LR22、LR24、及びLR26を配置する。

【0036】

このように、仮想機器配置スケジューラ機能部111は、物理機器103bに障害が生じた場合、物理機器103bに配置された仮想機器LR21~LR26を、物理機器103aと物理機器103cとに再配置する。すなわち、仮想機器配置スケジューラ機能部1

10

20

30

40

50

11は、複数台の物理機器を仮想機器の復旧先として利用するので、物理機器故障時の仮想機器復旧時間を短縮できる。

【0037】

続いて、このような仮想機器管理システムにおける処理動作について、図4及び図5を用いて説明する。図4は、仮想機器管理システムにおける仮想機器を作成する処理動作を説明するための図である。

【0038】

図4に示すように、ユーザ端末101は、仮想機器作成依頼を、仮想機器配置スケジューラ機能部111に送信する(ステップS1)。仮想機器配置スケジューラ機能部111は、仮想機器配置スケジューラDB110を参照し(ステップS2)、物理資源情報を確認する(ステップS3)。これにより仮想機器配置スケジューラ機能部111は、仮想機器を作成する物理機器103を決定し、APIパラメータを準備する(ステップS4)。

【0039】

次に、仮想機器配置スケジューラ機能部111は、決定した物理機器103に仮想機器を作成させるようにクラウドコントローラ108に依頼する(ステップS5)。続いて、クラウドコントローラ108は、物理機器103に仮想機器の作成を依頼する(ステップS6)。

【0040】

そして、物理機器103は、仮想機器を作成し(ステップS7)、仮想機器の作成が完了したことをクラウドコントローラ108に通知する(ステップS8)。続いて、クラウドコントローラ108は、仮想機器の作成が完了したことを仮想機器配置スケジューラ機能部111に通知する(ステップS9)。そして、仮想機器配置スケジューラ機能部111は、仮想機器の作成が完了したことをユーザ端末101に通知する(ステップS10)。

【0041】

図5は、仮想機器管理システムにおける仮想機器を再配置する処理動作を説明するための図である。図5では、いずれかの物理機器103に障害が生じた場合に、仮想機器配置スケジューラ機能部111が仮想機器の再配置を仲介する動作を説明する。図5に示すように、仮想機器管理システムでは、物理機器103a、物理機器103b及び物理機器103cが相互に機器状態を監視している(ステップS21、ステップS22)。以下では、物理機器103aに障害が生じた場合について説明する。ここで、物理機器103aで障害が起きた際は、物理機器103a上の高可用ソフトウェアは物理機器103a上のプロセスを停止し、障害を仮想機器配置スケジューラ機能部111に通知する。物理機器103b及び物理機器103cも同様に物理機器103aの障害を仮想機器配置スケジューラ機能部111に通知する。ここで、物理機器103aが完全に故障している場合は、物理機器103aから仮想機器配置スケジューラ機能部111に障害の発生を通知はできないが、物理機器103b及び物理機器103cは、物理機器103aの故障を仮想機器配置スケジューラ機能部111に通知できる。このため、仮想機器配置スケジューラ機能部111は、物理機器103aの故障を知ることができる。なお図5に示す例では、物理機器103aが完全に故障し、物理機器103aから障害の発生を仮想機器配置スケジューラ機能部111に通知できない場合を示す。

【0042】

かかる場合、物理機器103bは、物理機器103aに障害が生じたことを仮想機器配置スケジューラ機能部111に通知する(ステップS23)。そして、仮想機器配置スケジューラ機能部111は、物理機器103bにACK(ACKnowledgement)を応答する(ステップS24)。同様に、物理機器103cは、物理機器103aに障害が生じたことを仮想機器配置スケジューラ機能部111に通知する(ステップS25)。そして、仮想機器配置スケジューラ機能部111は、物理機器103cにACKを応答する(ステップS26)。ここで、仮想機器配置スケジューラ機能部111は、最初に受信した通知に従って仮想機器の再配置を始めるが、2番目以降に受信した通知に対してもACKを応答す

10

20

30

40

50

る。

【 0 0 4 3 】

仮想機器配置スケジューラ機能部 1 1 1 は、仮想機器配置スケジューラ DB 1 1 0 を参照し（ステップ S 2 7）、物理資源情報を確認する（ステップ S 2 8）。これにより仮想機器配置スケジューラ機能部 1 1 1 は、物理機器 1 0 3 b 及び物理機器 1 0 3 c の物理資源の空き容量を取得して、仮想機器を再配置する物理機器 1 0 3 を決定し、API パラメータを準備する（ステップ S 2 9）。ここで、仮想機器配置スケジューラ機能部 1 1 1 は、複数台の物理機器 1 0 3 を仮想機器の復旧先として選択することで、高速の復旧を可能とする。

【 0 0 4 4 】

次に、仮想機器配置スケジューラ機能部 1 1 1 は、再配置する仮想機器を物理機器 1 0 3 b に作成させるようにクラウドコントローラ 1 0 8 に依頼する（ステップ S 3 0）。続いて、クラウドコントローラ 1 0 8 は、物理機器 1 0 3 b に仮想機器の作成を依頼する（ステップ S 3 1）。同様に、仮想機器配置スケジューラ機能部 1 1 1 は、再配置する仮想機器を物理機器 1 0 3 c に作成させるようにクラウドコントローラ 1 0 8 に依頼する（ステップ S 3 2）。続いて、クラウドコントローラ 1 0 8 は、物理機器 1 0 3 c に仮想機器の作成を依頼する（ステップ S 3 3）。ここで、仮想機器配置スケジューラ機能部 1 1 1 は、選択した配置先を指定してクラウドコントローラ 1 0 8 の API を呼び出す。これにより、クラウドコントローラ 1 0 8 は、指定された物理機器 1 0 3 に対して仮想機器作成を依頼する。

【 0 0 4 5 】

そして、物理機器 1 0 3 b は、仮想機器を作成し（ステップ S 3 4）、仮想機器の作成が完了したことをクラウドコントローラ 1 0 8 に通知する（ステップ S 3 5）。続いて、クラウドコントローラ 1 0 8 は、仮想機器の作成が完了したことを仮想機器配置スケジューラ機能部 1 1 1 に通知する（ステップ S 3 6）。同様に、物理機器 1 0 3 c は、仮想機器を作成し（ステップ S 3 7）、仮想機器の作成が完了したことをクラウドコントローラ 1 0 8 に通知する（ステップ S 3 8）。続いて、クラウドコントローラ 1 0 8 は、仮想機器の作成が完了したことを仮想機器配置スケジューラ機能部 1 1 1 に通知する（ステップ S 3 9）。

【 0 0 4 6 】

続いて、図 6 を用いて、仮想機器管理装置 1 0 9 が実現する仮想機器配置スケジューラ DB 1 1 0 及び仮想機器配置スケジューラ機能部 1 1 1 について説明する。図 6 は、仮想機器管理装置 1 0 9 が実現する仮想機器配置スケジューラ DB 1 1 0 及び仮想機器配置スケジューラ機能部 1 1 1 を説明するための図である。

【 0 0 4 7 】

図 6 に示すように、仮想機器配置スケジューラ DB 1 1 0 は、仮想機器配置情報テーブル 1 1 0 a 及び物理資源情報テーブル 1 1 0 b を記憶する。仮想機器配置情報テーブル 1 1 0 a は、各仮想機器がどの物理機器上に配置されているかを示す仮想機器配置情報を記憶する。

【 0 0 4 8 】

図 7 は、仮想機器配置情報テーブル 1 1 0 a のデータ構造の一例を示す図である。図 7 に示すように、仮想機器配置情報テーブル 1 1 0 a は、「仮想機器 ID」と、「物理機器 ID」とを対応付けた仮想機器配置情報を記憶する。ここで、仮想機器配置情報テーブル 1 1 0 a が記憶する「仮想機器 ID」は、物理機器 1 0 3 に作成された仮想機器を一意に識別する識別子を示す。例えば、「仮想機器 ID」には、「仮想ボリューム # 1」、「仮想ボリューム # 2」等のデータ値が格納される。仮想機器配置情報テーブル 1 1 0 a が記憶する「物理機器 ID」は、物理機器 1 0 3 を一意に識別する識別子を示す。例えば、「物理機器 ID」には、「物理機器 # 1」、「物理機器 # 2」等のデータ値が格納される。

【 0 0 4 9 】

一例をあげると、図 7 に示す仮想機器配置情報テーブル 1 1 0 a は、識別子が「物理機

10

20

30

40

50

器 # 1」である物理機器 103 には、仮想機器「仮想ボリューム # 1」及び「仮想ボリューム # 2」が配置されていることを示す。また、図 7 に示す仮想機器配置情報テーブル 110a は、識別子が「物理機器 # 2」である物理機器 103 には、仮想機器「仮想 L2 ネットワーク # 1」、「仮想ルータ # 1」及び「仮想ロードバランサ # 1」が配置されていることを示す。また、図 7 に示す仮想機器配置情報テーブル 110a は、識別子が「物理機器 # 3」である物理機器 103 には、仮想機器「仮想マシン # 1」及び「仮想マシン # 2」が配置されていることを示す。

【0050】

図 6 に戻る。物理資源情報テーブル 110b は、各物理機器の稼働状態と物理機器が有する物理資源の空き容量とを示す物理資源情報を記憶する。図 8 は、物理資源情報テーブル 110b のデータ構造の一例を示す図である。図 8 に示すように、物理資源情報テーブル 110b は、「物理機器 ID」と「稼働状態」と「空き」と「使用中」と「障害用」とを対応付けた物理資源情報を記憶する。

10

【0051】

ここで、物理資源情報テーブル 110b が記憶する「物理機器 ID」は、物理機器 103 を一意に識別する識別子を示す。例えば、「物理機器 ID」には、「物理機器 # 1」、「物理機器 # 2」等のデータ値が格納される。

【0052】

また、物理資源情報テーブル 110b が記憶する「稼働状態」は、物理機器が稼働中であるか否かを示す。例えば、物理機器が稼働中である場合、「稼働状態」には「稼働中」が格納される。なお、図 8 では図示していないが、物理機器が予備系として設けられ稼働中ではない場合、「稼働状態」には「予備」が格納される。また、物理機器が故障中やメンテナンス中である場合、「稼働状態」には「故障中 / メンテ中」が格納される。

20

【0053】

また、「空き」は、物理機器が有する物理資源の容量のうち空き容量を示す。例えば、「空き」には、「3」、「5」、「4」等の値が格納される。また、「使用中」は、物理機器が有する物理資源の容量のうち使用中の容量を示す。例えば、「使用中」には、「1」、「3」、「2」等の値が格納される。また、「障害用」は、物理機器が有する物理資源の容量のうち復旧用に確保された容量を示す。例えば、「障害用」には、「2」等の値が格納される。

30

【0054】

一例をあげると、図 8 に示す物理資源情報テーブル 110b は、物理機器 # 1 は、稼働中であり、物理資源の空き容量が「3」であり、使用中の容量が「1」であり、復旧用に確保された容量が「2」であることを示す。また、図 8 に示す物理資源情報テーブル 110b は、物理機器 # 2 は、稼働中であり、物理資源の空き容量が「5」であり、使用中の容量が「3」であり、復旧用に確保された容量が「2」であることを示す。同様に、図 8 に示す物理資源情報テーブル 110b は、物理機器 # 3 は、稼働中であり、物理資源の空き容量が「4」であり、使用中の容量が「2」であり、復旧用に確保された容量が「2」であることを示す。

【0055】

図 6 に戻る。仮想機器配置スケジューラ機能部 111 は、作成依頼受付部 111a と、配置先選択部 111b と、障害検出部 111c と、再配置先選択部 111d と、作成要求部 111e とを有する。

40

【0056】

作成依頼受付部 111a は、仮想機器の作成要求をユーザ端末 101 から受け付ける。作成依頼受付部 111a は、受け付けた仮想機器の作成要求を配置先選択部 111b に受け渡す。

【0057】

配置先選択部 111b は、仮想機器を新規に作成する際に、仮想機器を配置する物理機器を選択する。ここで、配置先選択部 111b は、仮想機器を出来るだけ分散して配置す

50

るように物理機器 103 を選択する。言い換えると、配置先選択部 111b は、「稼働中」の「空き」スペースの数が平準化するように仮想機器を配置する。

【0058】

図9は、配置先選択部 111b による処理動作を説明するための図である。図9では、物理機器 #1 ~ 物理機器 #6 の6台の物理機器を有する仮想機器管理システムにおいて、仮想機器を新規に作成する場合について説明する。ここで、物理機器 #1 ~ 物理機器 #5 の稼働状態は「稼働中」であり、物理機器 #6 の稼働状態は「予備」である。また、物理機器 #1 のスペースの状態は、「空き」2、「使用中」1、「障害用バッファ」2であり、物理機器 #2 のスペースの状態は、「空き」0、「使用中」3、「障害用バッファ」2であり、物理機器 #3 のスペースの状態は、「空き」4、「使用中」0、「障害用バッファ」2である。また、物理機器 #4 のスペースの状態は、「空き」0、「使用中」5、「障害用バッファ」2であり、物理機器 #5 のスペースの状態は、「空き」2、「使用中」0、「障害用バッファ」2であり、物理機器 #6 のスペースの状態は、「空き」3、「使用中」0、「障害用バッファ」2である。

10

【0059】

例えば、配置先選択部 111b は、配置先選択時に、稼働状態が「稼働中」である物理機器の空きスペースの量をチェックし、最も空きスペースが多い稼働中の物理機器を特定する。より具体的には、配置先選択部 111b は、作成する仮想機器のうち1つの仮想機器（例えば、仮想機器 #1）を選択する。そして、図9に示すスペースの状態である場合には、「空き」が4である物理機器 #3 を、最も空きスペースが多い稼働中の物理機器に

20

【0060】

続いて、配置先選択部 111b は、配置先として選択する処理を、作成を依頼された全ての仮想機器の配置先を選択するまで繰り返す。一例をあげると、配置先選択部 111b は、図9に示す数字順に仮想機器を配置するように物理機器を選択する。このように、配置先選択部 111b は、最も空きスペースが多い稼働中の物理機器のスペースの一部を選択することで「空き」スペースの数を平準化する。

【0061】

また、配置先選択部 111b は、稼働状態が「稼働中」である物理機器のスペースが全て埋まった場合に、稼働状態が「予備」である物理機器に仮想機器を配置する。このため、配置先選択部 111b は、図9に示す8番のスペースまで仮想機器を配置したら、予備の物理機器に仮想機器を配置する。すなわち、配置先選択部 111b は、図9に示す例において、仮想機器を9台以上作成する場合には、稼働状態が「予備」である物理機器に仮想機器を配置する。

30

【0062】

図6に戻る。障害検出部 111c は、障害の生じた物理機器 103 を検出する。ここで、障害検出部 111c は、各物理機器 103 が有する高可用ソフトウェア 107 と連携することで、障害の生じた物理機器 103 を検出する。図10は、障害検出部 111c による処理動作を説明するための図である。

40

【0063】

図10では、物理機器 103 a ~ 103 c を図示しており、物理機器 103 a に障害が発生した場合について説明する。また、図10では、物理機器 103 が有する機能のうち、物理機器 103 a には、自装置の障害発生時に機能する構成部を示し、物理機器 103 b 及び物理機器 103 c には、他装置の障害を検知した場合に機能する構成部を示す。

【0064】

図10に示すように、物理機器 103 の障害の検知には、高可用ソフトウェア 107 a ~ 107 c が用いられる。全ての物理機器 103 は、CIB (Cluster Information Base) に、クラスタ内の全物理機器 103 の状態を保持する。高可用ソフトウェア 107 は

50

、R A (Resource Agent) を用いて自物理機器の状態を確認する。なお、R A とは、例えば、仮想ボリューム制御部や仮想マシン制御部に相当する。

【 0 0 6 5 】

また、高可用ソフトウェア 1 0 7 は、H e a r t b e a t により、クラスタ内のどの物理機器も他の物理機器の状態を知り得る。このため、高可用ソフトウェア 1 0 7 は、h e a r t b e a t パケットを使ってクラスタ内に状態を通知する。この仕組みにより、各物理機器は他の物理機器の状態を知る。高可用ソフトウェア 1 0 7 は、ある物理機器からの h e a r t b e a t パケットが継続的にロストすると、他の物理機器は当該物理機器がダウンしたとみなす。

【 0 0 6 6 】

物理機器 1 0 3 は、仮想機器配置スケジューラ機能部 1 1 1 に物理機器に生じた障害を通知するため、通知 R A と通知プロセスとを備える。例えば、P a c e m a k e r が自物理機器の故障を検出した場合、通知 R A を使用して仮想機器配置スケジューラ機能部 1 1 1 に自物理機器の故障を通知する。一方、通知プロセスは、常駐プロセスとして設定され、C I B の状態を定期的に確認する。そして、通知プロセスは、他物理機器の故障を検出すると、他物理機器に障害が生じたことを仮想機器配置スケジューラ機能部 1 1 1 に通知する。通知 R A による通知及び通知プロセスによる通知は、A C K が仮想機器配置スケジューラ機能部 1 1 1 から返るまで一定回数繰り返される。

【 0 0 6 7 】

続いて、仮想機器配置スケジューラ機能部 1 1 1 において、障害検出部 1 1 1 c は、通知を受信したら A C K を応答し、再配置先選択部 1 1 1 d に仮想機器の再配置処理を実行させる。また、障害検出部 1 1 1 c は、2 通目以降の通知を無視して A C K を応答する。これにより、複数の物理機器から通知を受けることで冗長化対策をとることができることも、復旧処理を繰り返さないようにする。

【 0 0 6 8 】

なお、高可用ソフトウェア 1 0 7 は、自物理機器の停止に失敗する場合がある。仮想マシンの場合、復旧により、複数の仮想マシンが同時に存在してしまい、データ領域への同時アクセスによりデータ破壊の可能性が出てしまう。そこで、高可用ソフトウェア 1 0 7 が「P a c e m a k e r」である場合、S T O N I T H モジュールを用いて、確実に故障物理機器を落とす。S T O N I T H は、I P M I (Intelligent Platform Management

Interface) 経由で、故障物理機器を停止することで、故障物理機器が動作し続けないことを保証する。Q u o r u m で過半数を形成した多数派の物理機器が、S T O N I T H を起動することで、誤発動を防止する。なお、Q u o r u m は過半数で判断するため、クラスタの物理機器数が少ない場合に、ある物理機器が故障したら、正常な物理機器が過半数を確保できなくなる。このため、クラスタから故障物理機器を切り離す減設作業が必要である。また、図 1 0 では、高可用ソフトウェア 1 0 7 が、P a c e m a k e r である場合を示しているが、他の高可用ソフトウェアでも同様のメカニズムで障害の発生を検知したり、障害の発生を通知したりすることが可能である。

【 0 0 6 9 】

図 6 に戻る。再配置先選択部 1 1 1 d は、物理機器に生じた障害を通知された場合に、仮想機器を再作成する物理機器を選択する。ここで、再配置先選択部 1 1 1 d は、障害の生じた物理機器以外の物理機器にできるだけ順番に割り振られるように物理機器を選択する。例えば、再配置先選択部 1 1 1 d は、障害の生じた物理機器以外の物理機器のうち物理資源の空き容量のある物理機器を複数特定する。そして、再配置先選択部 1 1 1 d は、特定した複数の物理機器の物理資源を、障害の生じた物理機器に配置された仮想機器の再配置先として選択する。

【 0 0 7 0 】

図 1 1 は、再配置先選択部 1 1 1 d による処理動作を説明するための図である。図 1 1 では、物理機器 # 1 ~ 物理機器 # 6 の 6 台の物理機器を有する仮想機器管理システムにおいて、物理機器 # 4 が故障した際の復旧について説明する。ここで、物理機器 # 1 ~ 物理

10

20

30

40

50

機器 # 5 の稼働状態は「稼働中」であり、物理機器 # 6 の稼働状態は「予備」である。また、物理機器 # 1 のスペースの状態は、「空き」2、「使用中」1、「障害用バッファ」2であり、物理機器 # 2 のスペースの状態は、「空き」0、「使用中」3、「障害用バッファ」2であり、物理機器 # 3 のスペースの状態は、「空き」4、「使用中」0、「障害用バッファ」2である。また、物理機器 # 4 のスペースの状態は、「空き」0、「使用中」10、「障害用バッファ」2であり、物理機器 # 5 のスペースの状態は、「空き」2、「使用中」0、「障害用バッファ」2であり、物理機器 # 6 のスペースの状態は、「空き」3、「使用中」0、「障害用バッファ」2である。

【 0 0 7 1 】

例えば、再配置先選択部 1 1 1 d は、障害の生じた物理機器以外の物理機器のうち物理資源の空き容量のある物理機器を複数特定する。ここで、再配置先選択部 1 1 1 d は、物理機器の障害発生時には、「空き」のスペースに加えて、「障害用バッファ」のスペースも使用する。これにより、再配置先選択部 1 1 1 d は、「空き」が 0 の物理機器も含めて、より多くの物理機器が復旧処理を分担できるようにする。図 1 1 に示す例では、再配置先選択部 1 1 1 d が、稼働中である物理機器 # 1 ~ 物理機器 # 3 及び物理機器 # 5 を特定した場合を示す。

10

【 0 0 7 2 】

そして、再配置先選択部 1 1 1 d は、特定した複数の物理機器の物理資源を、障害の生じた物理機器 # 4 に配置された仮想機器の再配置先として選択する。例えば、再配置先選択部 1 1 1 d は、特定した複数の物理機器に順序付けを行う。ここで、再配置先選択部 1 1 1 d は、特定した複数の物理機器の物理資源の空き容量が多い順に、特定した複数の物理機器に順序付けを行う。例えば、再配置先選択部 1 1 1 d は、「空き」のスペースと「障害用バッファ」のスペースとの合計スペースを物理資源の空き容量とし、合計スペースが多い順に物理機器に順序付けを行う。図 1 1 の例では、再配置先選択部 1 1 1 d が、合計スペースが 6 である物理機器 # 3、合計スペースが 5 である物理機器 # 1、合計スペースが 4 である物理機器 # 5、そして、合計スペースが 2 である物理機器 # 2 の順で順序付けした場合を示す。

20

【 0 0 7 3 】

続いて、再配置先選択部 1 1 1 d は、障害の生じた物理機器に配置された仮想機器それぞれの再配置先として、順序に基づいて選択した物理機器の物理資源を選択する処理を繰り返す。一例をあげると、再配置先選択部 1 1 1 d は、図 1 1 に示す数字順に仮想機器を再配置するように物理機器を選択する。より具体的には、再配置先選択部 1 1 1 d は、物理機器 # 3、物理機器 # 1、物理機器 # 5、そして、物理機器 # 2 の順で選択した物理機器の物理資源を仮想機器の再配置先として選択する処理を繰り返す。ここで、再配置先選択部 1 1 1 d は、「空き」のスペースや「障害用バッファ」のスペースが無くなるまでは、各物理機器に仮想機器を順番に配置する。また、再配置先選択部 1 1 1 d は、スペースが無くなった物理機器は飛ばすようにする。

30

【 0 0 7 4 】

なお、再配置先選択部 1 1 1 d は、稼働状態が「稼働中」である全ての物理機器の「空き」のスペース及び「障害用バッファ」のスペースが満たされるまで、稼働状態が「予備」である物理機器を選択しない。このように、仮想機器配置スケジューラ機能部 1 1 1 は、「空き」のスペースに加えて、仮想機器の作成時には利用されない「障害用バッファ」のスペースを予め準備しておき、障害時に多くの物理機器に仮想機器を再配置することで、高速の復旧を可能とする。また、再配置先選択部 1 1 1 d は、障害が発生した物理機器に配置された仮想機器の全てを再配置可能ではない場合には、特定した物理機器の「空き」のスペースと「障害用バッファ」のスペースとに再配置可能な範囲で、仮想機器ごとに再配置先を選択する。

40

【 0 0 7 5 】

また、P a c e m a k e r のクラスタ構成は、最大 8 台程度で組み、障害の検知を行う。また、仮想機器配置スケジューラ機能部 1 1 1 は、クラスタを跨いで別物理機器に仮想

50

機器を作成してもよいため、再作成が依頼される物理機器はクラスタのサイズ以上でも良い。また、全てが埋まった際に利用される予備機は存在してもしなくてもよい。クラスタ構成上はN - A c t、0 - S b yで、S t a n d b y機を準備する必要はないため、物理機器の利用効率を高めることも出来る。

【 0 0 7 6 】

図6に戻る。作成要求部111eは、配置先選択部111bにより選択された物理機器103に、仮想機器を作成するようにクラウドコントローラ108に依頼する。また、作成要求部111eは、再配置先選択部111dにより選択された物理機器103に、仮想機器を作成するようにクラウドコントローラ108に依頼する。

【 0 0 7 7 】

図12は、仮想機器の作成を要求された場合の仮想機器配置スケジューラ機能部111による処理手順を示すフローチャートである。図12に示すように、作成依頼受付部111aは、仮想機器の作成をユーザ端末101から依頼されたか否かを判定する(ステップS101)。ここで、作成依頼受付部111aは、仮想機器の作成をユーザ端末101から依頼されたと判定した場合(ステップS101、Y e s)、作成を依頼された仮想機器を特定する(ステップS102)。なお、作成依頼受付部111aは、仮想機器の作成をユーザ端末101から依頼されたと判定しなかった場合(ステップS101、N o)、繰り返し物理機器の作成をユーザ端末101から依頼されたか否かを判定する。

【 0 0 7 8 】

続いて、配置先選択部111bは、作成を依頼された仮想機器を1つ選択する(ステップS103)。そして、配置先選択部111bは、空き容量の最も多い物理機器を特定する(ステップS104)。また、配置先選択部111bは、空き容量の最も多い物理機器を特定できたか否かを判定する(ステップS105)。ここで、配置先選択部111bは、空き容量の最も多い物理機器を特定できたと判定した場合(ステップS105、Y e s)、ステップS109に移行する。一方、配置先選択部111bは、空き容量の最も多い物理機器を特定できたと判定しなかった場合(ステップS105、N o)、予備系の物理機器が存在するか否かを判定する(ステップS106)。ここで、配置先選択部111bは、予備系の物理機器が存在すると判定しなかった場合(ステップS106、N o)、処理を終了する。

【 0 0 7 9 】

一方、配置先選択部111bは、予備系の物理機器が存在すると判定した場合(ステップS106、Y e s)、空き容量の最も多い予備系の物理機器を特定する(ステップS107)。また、配置先選択部111bは、空き容量の最も多い予備系の物理機器を特定できたか否かを判定する(ステップS108)。ここで、配置先選択部111bは、空き容量の最も多い予備系の物理機器を特定できたと判定しなかった場合(ステップS108、N o)、処理を終了する。一方、配置先選択部111bは、空き容量の最も多い予備系の物理機器を特定できたと判定した場合(ステップS108、Y e s)、ステップS109に移行する。

【 0 0 8 0 】

ステップS109において、配置先選択部111bは、特定した物理機器を選択した仮想機器の配置先に選択する(ステップS109)。そして、配置先選択部111bは、作成を依頼された全ての仮想機器の配置先を選択したか否かを判定する(ステップS110)。ここで、配置先選択部111bは、作成を依頼された全ての仮想機器の配置先を選択したと判定しなかった場合(ステップS110、N o)、作成を依頼された全ての仮想機器の配置先を選択するまでステップS103からステップS110までの処理を繰り返し実行する。

【 0 0 8 1 】

一方、配置先選択部111bは、作成を依頼された全ての仮想機器の配置先を選択したと判定した場合(ステップS110、Y e s)、クラウドコントローラ108に配置を依頼し(ステップS111)、処理を終了する。

10

20

30

40

50

【0082】

図13は、物理機器103に障害が発生した場合の仮想機器配置スケジューラ機能部111による処理手順を示すフローチャートである。図13に示すように、障害検出部111cは、物理機器の障害を通知されたか否かを判定する(ステップS201)。ここで、障害検出部111cは、物理機器の障害を通知されたと判定した場合(ステップS201、Yes)、仮想機器配置情報テーブル110aを参照して、障害が発生した物理機器に配置された仮想機器を特定する(ステップS202)。なお、障害検出部111cは、物理機器の障害を通知されたと判定しなかった場合(ステップS201、No)、繰り返し物理機器の障害を通知されたか否かを判定する。

【0083】

続いて、再配置先選択部111dは、空き容量及び障害用容量の少なくともいずれかがある物理機器を複数特定する(ステップS203)。そして、再配置先選択部111dは、複数の物理機器を特定できたか否かを判定する(ステップS204)。ここで、再配置先選択部111dは、複数の物理機器を特定できなかった場合(ステップS204、No)、予備系の物理機器が存在するか否かを判定する(ステップS205)。ここで、再配置先選択部111dは、予備系の物理機器が存在すると判定しなかった場合(ステップS205、No)、ステップS207に移行する。

【0084】

一方、再配置先選択部111dは、予備系の物理機器が存在すると判定した場合(ステップS205、Yes)、空き容量及び障害用容量の少なくともいずれかがある予備系の物理機器を特定する(ステップS206)。そして、再配置先選択部111dは、稼働中及び予備系を合わせて1以上の物理機器を特定できたか否かを判定する(ステップS207)。ここで、再配置先選択部111dは、稼働中及び予備系を合わせて1以上の物理機器を特定できたと判定した場合(ステップS207、Yes)、ステップS211に移行する。一方、再配置先選択部111dは、稼働中及び予備系を合わせて1以上の物理機器を特定できたと判定しなかった場合(ステップS207、No)、処理を終了する。

【0085】

再配置先選択部111dは、ステップS204において、複数の物理機器を特定できたと判定した場合(ステップS204、Yes)、特定した複数の物理機器に、障害が発生した物理機器に配置された仮想機器を全て再配置可能であるか否かを判定する(ステップS208)。ここで、再配置先選択部111dは、障害が発生した物理機器に配置された仮想機器を全て再配置可能であると判定した場合(ステップS208、Yes)、ステップS211に移行する。一方、再配置先選択部111dは、障害が発生した物理機器に配置された仮想機器を全て再配置可能であると判定しなかった場合(ステップS208、No)、予備系の物理機器が存在するか否かを判定する(ステップS209)。ここで、再配置先選択部111dは、予備系の物理機器が存在すると判定した場合(ステップS209、Yes)、空き容量及び障害用容量の少なくともいずれかがある予備系の物理機器を特定する(ステップS210)。

【0086】

ステップS211において、再配置先選択部111dは、特定した物理機器に順序付けを行う(ステップS211)。例えば、再配置先選択部111dは、空き容量及び障害用容量が多い順に、特定した物理機器に順序付けを行う。なお、再配置先選択部111dは、特定した物理機器が1つである場合には、ステップS211の処理を省略してもよい。

【0087】

続いて、再配置先選択部111dは、仮想機器の再配置先を選択する(ステップS212)。例えば、再配置先選択部111dは、障害の生じた物理機器に配置された仮想機器それぞれの再配置先として、順序に基づいて選択した物理機器の物理資源を選択する処理を繰り返す。なお、再配置先選択部111dは、障害が発生した物理機器に配置された仮想機器の全てを再配置可能ではない場合には、特定した物理機器の空き容量及び障害用容量に再配置可能な範囲で、仮想機器ごとに再配置先を選択する。

10

20

30

40

50

【 0 0 8 8 】

そして、再配置先選択部 1 1 1 d は、クラウドコントローラ 1 0 8 に再配置を依頼し (ステップ S 2 1 3)、処理を終了する。なお、再配置先選択部 1 1 1 d は、仮想機器管理システムにおいて予備系の物理機器がない場合には、ステップ S 2 0 5、ステップ S 2 0 6、ステップ S 2 0 9 及びステップ S 2 1 0 の処理を省略してもよい。

【 0 0 8 9 】

上述したように、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、障害の生じた物理機器以外の物理機器のうち物理資源の空き容量のある物理機器を複数特定する。そして、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、特定した複数の物理機器の物理資源を、障害の生じた物理機器に配置された仮想機器の再配置先として選択する。すなわち、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、複数台の物理機器を仮想機器の復旧先として利用する。これにより、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、仮想機器を復旧するまでの時間を短縮することができる。

10

【 0 0 9 0 】

より具体的には、従来方式では、N - A c t、M - S b y でクラスタを組み物理機器に障害が起きた際に、P a c e m a k e r 等の高可用ソフトウェアの機能により S t a n d b y 機にフェールオーバーし、O p e n S t a c k 等のクラウドコントローラの D B を元に仮想機器を再構築していた。ここで、従来方式では、H A クラスタソフトウェアを用いたフェールオーバーは、1 台の S t a n d b y 機に仮想機器を新たに再構築するため、全仮想機器の復旧に時間がかかるという問題がある。

20

【 0 0 9 1 】

一方、第 1 の実施形態に係る仮想機器管理装置 1 0 9 では、N - A c t、0 - S b y でクラスタを組み、物理機器に障害が起きた際は、高可用ソフトウェアの機能により障害を検知するが、フェールオーバーせずに物理機器の障害を仮想機器管理装置 1 0 9 に通知する。仮想機器管理装置 1 0 9 は、各仮想機器に対して、再配置する複数の物理機器を決定し、配置する物理機器を指定してクラウドコントローラ 1 0 8 に再作成依頼を行う。そして、クラウドコントローラ 1 0 8 は、指定された物理機器に仮想機器を作成する。

【 0 0 9 2 】

このように、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、故障した物理機器上で動作していた仮想機器を、複数台の物理機器に再作成することで高速に復旧する。言い換えると、仮想機器管理装置 1 0 9 は、複数台の物理機器を復旧先として利用するため、物理機器故障時の仮想機器復旧時間が短縮される。例えば、移行先物理機器が 3 台の場合は、復旧処理時間が 1 / 3 に短縮できる。

30

【 0 0 9 3 】

また、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、P a c e m a k e r 等の高可用ソフトウェアで N - A c t、0 - S b y でクラスタを組み障害検知を行う。ここで、第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、クラスタの枠を超えて故障復旧を行うことが出来るため、移行先物理機器の台数をクラスタサイズ以上にとることもできる。これにより、復旧時間をより短縮できる。

【 0 0 9 4 】

更に、仮想機器管理システムでは、障害検知のためのクラスタは N - A c t、0 - S b y であるため、S t a n d b y 用の物理機器を準備する必要がなく、物理機器数の増大を抑えることができる。

40

【 0 0 9 5 】

第 1 の実施形態に係る仮想機器管理装置 1 0 9 は、物理機器に障害が起きた際の、仮想機器復旧を高速に行う事を狙っているが、実施形態はこれに限定されるものではない。例えば、物理機器が完全に故障した場合以外の、以下のユースケースにも、拡張して対応することが出来る。例えば、物理機器が完全には故障していなくても、複数あるファンの一つが故障した場合等は、サーバを停止してファン交換したい場合がある。かかる場合、仮想機器管理装置 1 0 9 は、物理機器上の仮想機器一括移動の A P I / G U I をオペレータ

50

に提供する。そして、仮想機器管理装置 109 は、API / GUI 経由でのオペレータの依頼を受けて、故障復旧時と同様に仮想機器を一括で別物理機器に移動する。これにより、例えば管理者は、ファンの一つが故障した物理機器をメンテナンスすることが出来る。

【0096】

なお、上述した実施形態では、配置先選択部 111b は、仮想機器を新規に作成する通常のオペレーション時に、仮想機器を出来るだけ分散して配置するように物理機器 103 を選択するものとして説明したが実施形態はこれに限定されるものではない。例えば、配置先選択部 111b は、仮想機器を新規に作成する通常のオペレーション時には、仮想機器を分散させることなく配置するように物理機器 103 を選択するようにしてもよい。

【0097】

また、図 11 に示す例では、再配置先選択部 111d が、特定した複数の物理機器の物理資源の空き容量が多い順に、特定した複数の物理機器に順序付けを行う場合について説明したが、実施形態はこれに限定されるものではない。例えば、再配置先選択部 111d は、物理資源の空き容量とは関係なく、特定した複数の物理機器に任意に順序付けを行うようにしてもよい。

【0098】

(第2の実施形態)

さて、これまで本発明の実施形態について説明したが、本発明は上述した実施形態以外にも、その他の実施形態にて実施されてもよい。そこで、以下では、その他の実施形態を示す。

【0099】

(システム構成)

また、本実施形態において説明した各処理のうち、自動的に行われるものとして説明した処理の全部又は一部を手動的に行うこともでき、あるいは、手動的に行われるものとして説明した処理の全部又は一部を公知の方法で自動的に行うこともできる。この他、上述の文書中や図面中で示した処理手順、制御手順、具体的名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。

【0100】

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部又は一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的又は物理的に分散・統合して構成することができる。

【0101】

(プログラム)

また、上記第1の実施形態に係る仮想機器管理装置 109 が実行する処理をコンピュータが実行可能な言語で記述した仮想機器管理プログラムを生成することもできる。この場合、コンピュータが仮想機器管理プログラムを実行することにより、上記実施形態と同様の効果を得ることができる。さらに、かかる仮想機器管理プログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録された仮想機器管理プログラムをコンピュータに読み込ませて実行することにより上記実施形態と同様の処理を実現してもよい。以下に、図 1 等に示した仮想機器管理装置 109 と同様の機能を実現する仮想機器管理プログラムを実行するコンピュータの一例を説明する。

【0102】

図 14 は、仮想機器管理プログラムを実行するコンピュータ 1000 を示す図である。図 14 に示すように、コンピュータ 1000 は、例えば、メモリ 1010 と、CPU 1020 と、ハードディスクドライブインタフェース 1030 と、ディスクドライブインタフェース 1040 と、シリアルポートインタフェース 1050 と、ビデオアダプタ 1060 と、ネットワークインタフェース 1070 とを有する。これらの各部は、バス 1080 によって接続される。

【0103】

10

20

30

40

50

メモリ 1010 は、ROM (Read Only Memory) 1011 および RAM (Random Access Memory) 1012 を含む。ROM 1011 は、例えば、BIOS (Basic Input Output System) 等のブートプログラムを記憶する。ハードディスクドライブインタフェース 1030 は、ハードディスクドライブ 1031 に接続される。ディスクドライブインタフェース 1040 は、ディスクドライブ 1041 に接続される。ディスクドライブ 1041 には、例えば、磁気ディスクや光ディスク等の着脱可能な記憶媒体が挿入される。シリアルポートインタフェース 1050 には、例えば、マウス 1051 およびキーボード 1052 が接続される。ビデオアダプタ 1060 には、例えば、ディスプレイ 1061 が接続される。

【0104】

ここで、図 14 に示すように、ハードディスクドライブ 1031 は、例えば、OS 1091、アプリケーションプログラム 1092、プログラムモジュール 1093 およびプログラムデータ 1094 を記憶する。上記実施形態で説明した仮想機器管理プログラムは、例えばハードディスクドライブ 1031 やメモリ 1010 に記憶される。

【0105】

また、仮想機器管理プログラムは、例えば、コンピュータ 1000 によって実行される指令が記述されたプログラムモジュールとして、例えばハードディスクドライブ 1031 に記憶される。具体的には、上記実施形態で説明した障害検出部 111c と同様の情報処理を実行する検出手順と、再配置先選択部 111d と同様の情報処理を実行する選択手順と、作成要求部 111e と同様の情報処理を実行する依頼手順とが記述されたプログラムモジュール 1093 が、ハードディスクドライブ 1031 に記憶される。

【0106】

また、仮想機器管理プログラムによる情報処理に用いられるデータは、プログラムデータ 1094 として、例えば、ハードディスクドライブ 1031 に記憶される。そして、CPU 1020 が、ハードディスクドライブ 1031 に記憶されたプログラムモジュール 1093 やプログラムデータ 1094 を必要に応じて RAM 1012 に読み出して、上述した各手順を実行する。

【0107】

なお、仮想機器管理プログラムに係るプログラムモジュール 1093 やプログラムデータ 1094 は、ハードディスクドライブ 1031 に記憶される場合に限られず、例えば、着脱可能な記憶媒体に記憶されて、ディスクドライブ 1041 等を介して CPU 1020 によって読み出されてもよい。あるいは、仮想機器管理プログラムに係るプログラムモジュール 1093 やプログラムデータ 1094 は、LAN (Local Area Network) や WAN (Wide Area Network) 等のネットワークを介して接続された他のコンピュータに記憶され、ネットワークインタフェース 1070 を介して CPU 1020 によって読み出されてもよい。

【0108】

(その他)

なお、本実施形態で説明した特定プログラムは、インターネットなどのネットワークを介して配布することができる。また、特定プログラムは、ハードディスク、フレキシブルディスク (FD)、CD-ROM、MO、DVD などのコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行することもできる。

【符号の説明】

【0109】

- 109 仮想機器管理装置
- 110 仮想機器配置スケジューラ DB
- 110 a 仮想機器配置情報テーブル
- 110 b 物理資源情報テーブル
- 111 仮想機器配置スケジューラ機能部

10

20

30

40

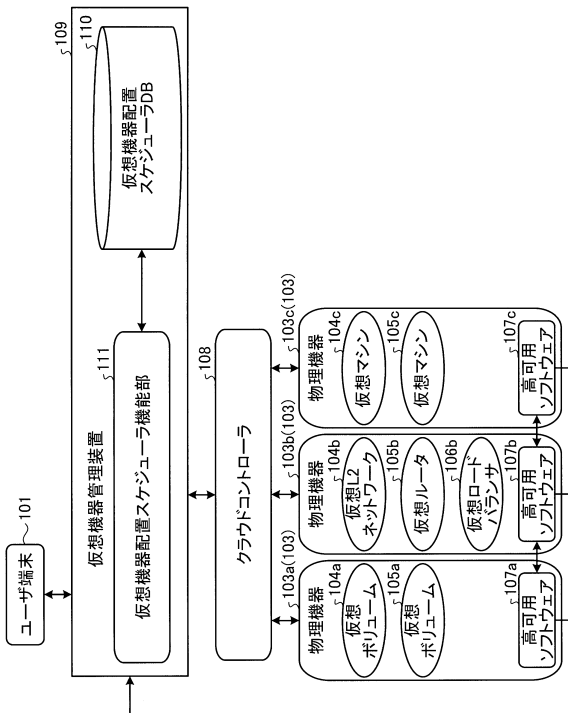
50

- 1 1 1 a 作成依頼受付部
- 1 1 1 b 配置先選択部
- 1 1 1 c 障害検出部
- 1 1 1 d 再配置先選択部
- 1 1 1 e 作成要求部
- 1 0 0 0 コンピュータ
- 1 0 1 0 メモリ
- 1 0 1 1 ROM
- 1 0 1 2 RAM
- 1 0 2 0 CPU
- 1 0 3 0 ハードディスクドライブインタフェース
- 1 0 3 1 ハードディスクドライブ
- 1 0 4 0 ディスクドライブインタフェース
- 1 0 4 1 ディスクドライブ
- 1 0 5 0 シリアルポートインタフェース
- 1 0 5 1 マウス
- 1 0 5 2 キーボード
- 1 0 6 0 ビデオアダプタ
- 1 0 6 1 ディスプレイ
- 1 0 7 0 ネットワークインタフェース
- 1 0 8 0 バス
- 1 0 9 1 OS
- 1 0 9 2 アプリケーションプログラム
- 1 0 9 3 プログラムモジュール
- 1 0 9 4 プログラムデータ

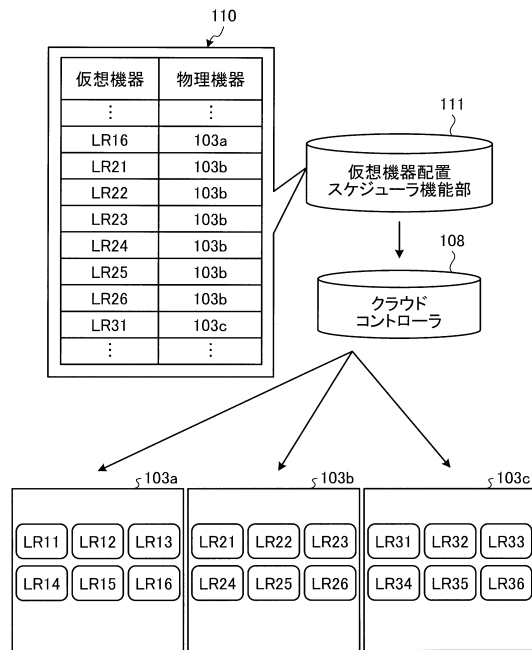
10

20

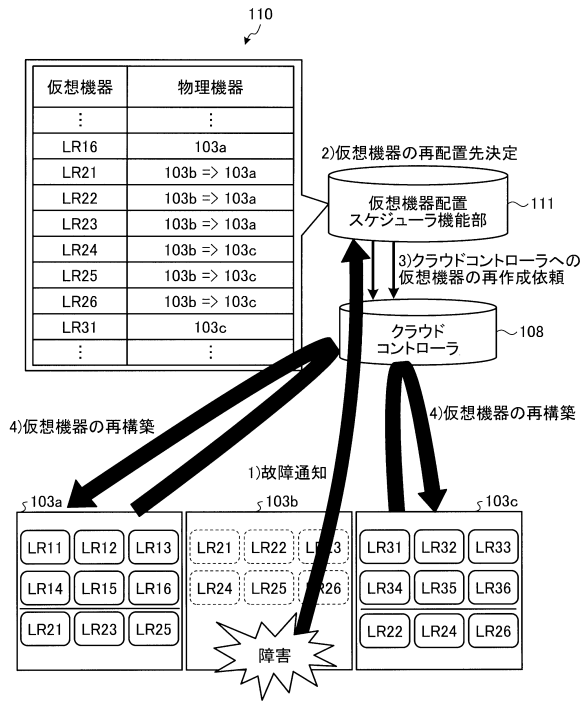
【図1】



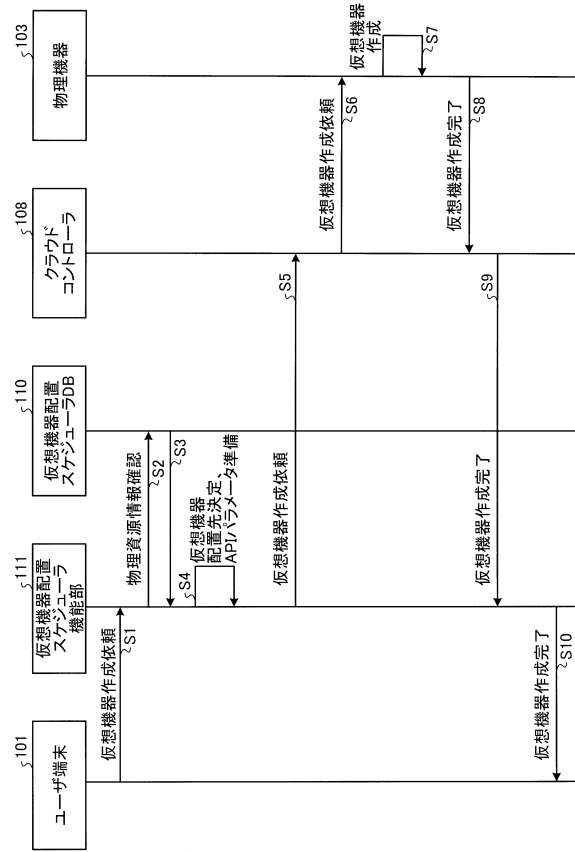
【図2】



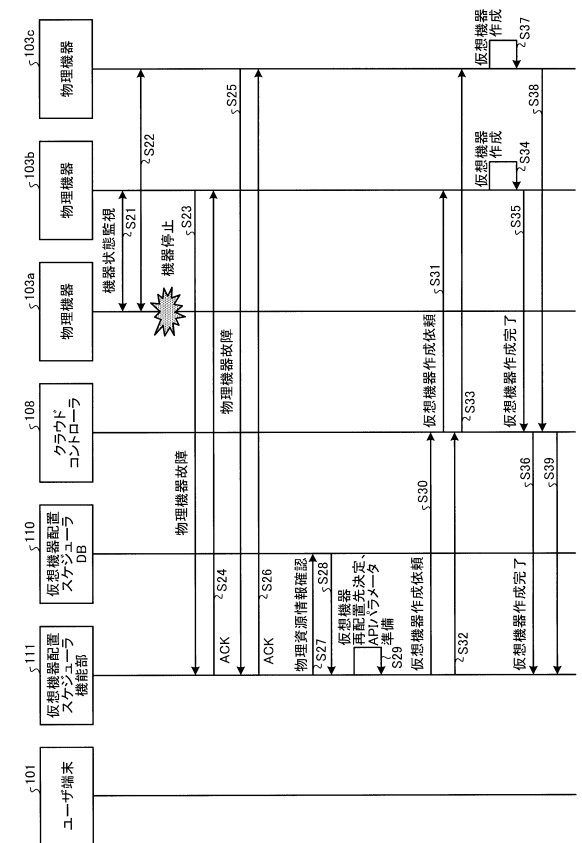
【図3】



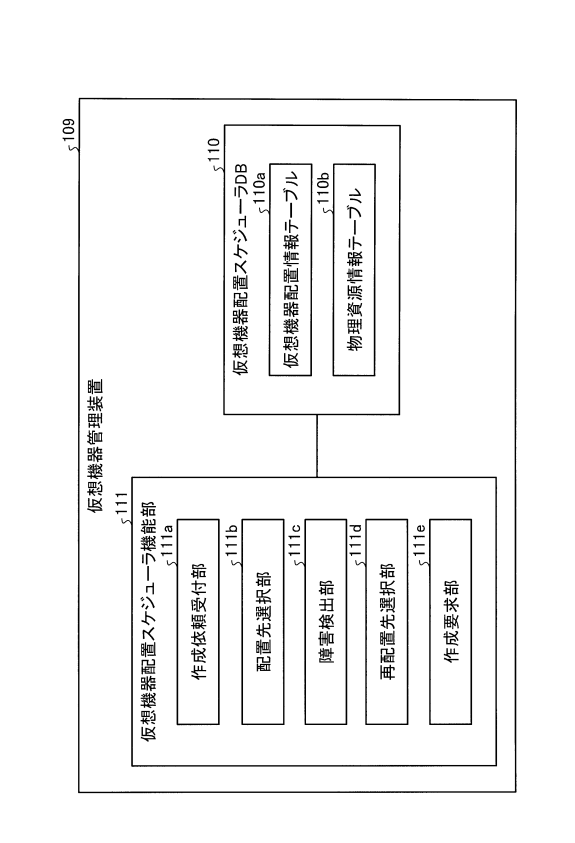
【図4】



【図5】



【図6】



【図7】

110a

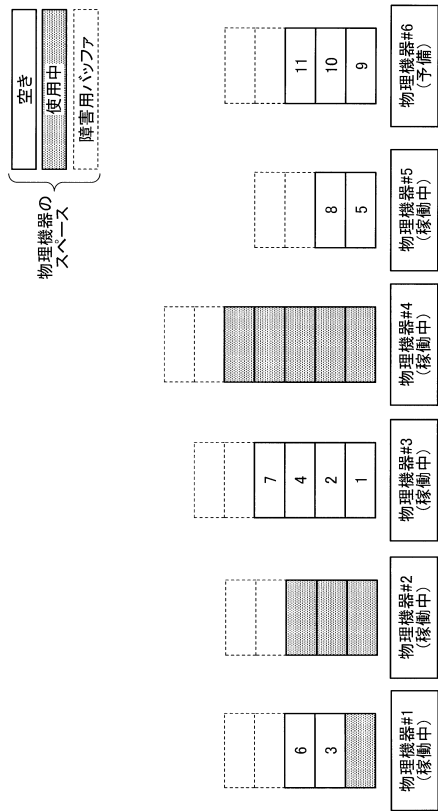
仮想機器ID	物理機器ID
仮想ボリューム#1	物理機器#1
仮想ボリューム#2	物理機器#1
仮想L2ネットワーク#1	物理機器#2
仮想ルータ#1	物理機器#2
仮想ロードバランサ#1	物理機器#2
仮想マシン#1	物理機器#3
仮想マシン#2	物理機器#3

【図8】

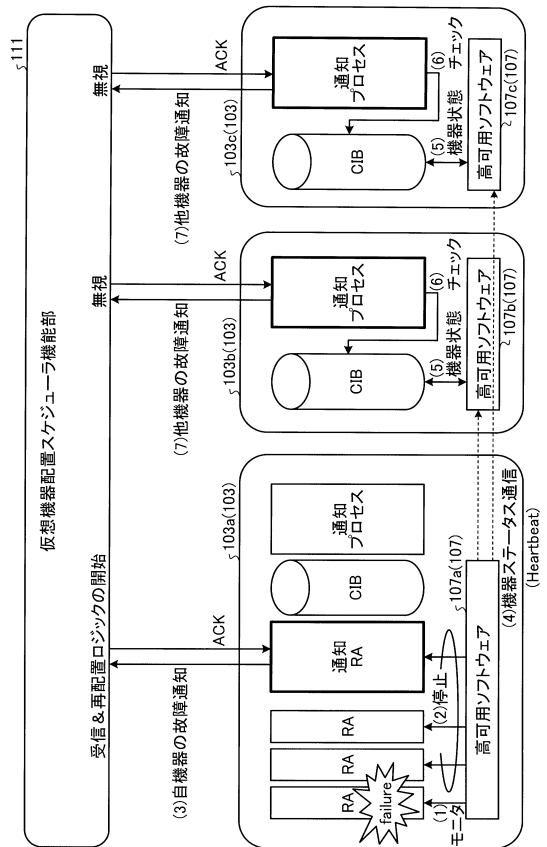
110b

物理機器ID	稼働状態	空き	使用中	障害用
物理機器#1	稼働中	3	1	2
物理機器#2	稼働中	5	3	2
物理機器#3	稼働中	4	2	2

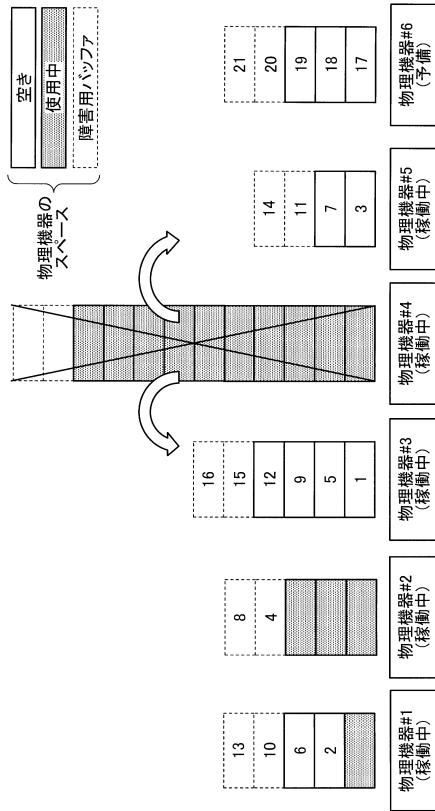
【図9】



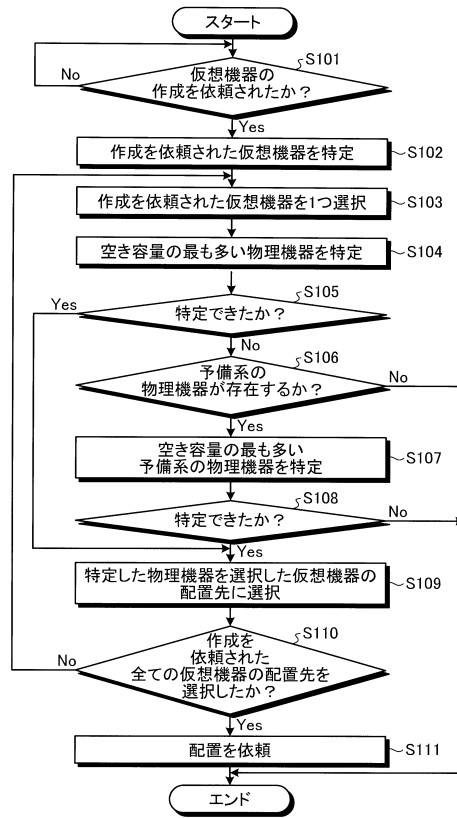
【図10】



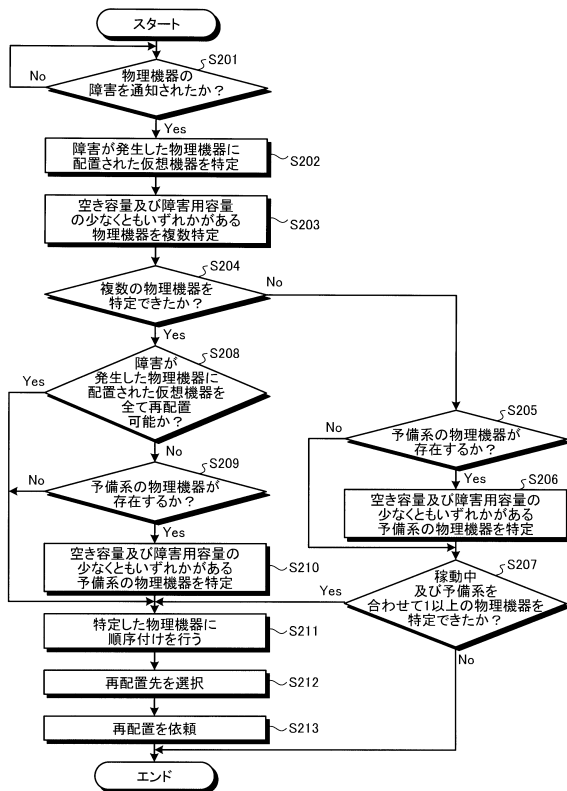
【図11】



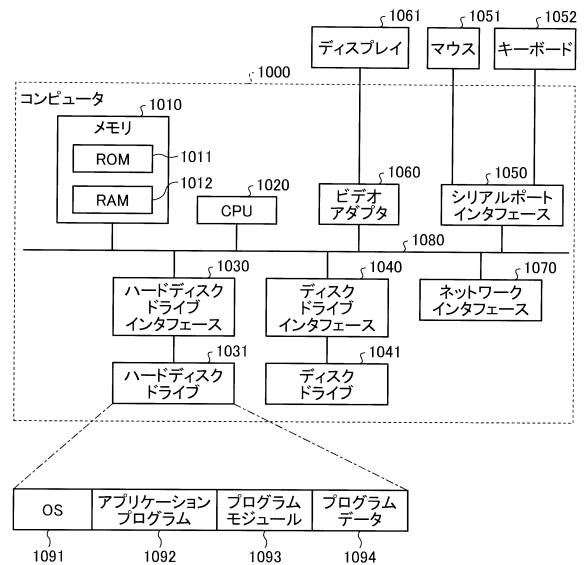
【図12】



【図13】



【図14】



フロントページの続き

- (72)発明者 長尾 伸二
東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
- (72)発明者 渡邊 拓磨
東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内

合議体

- 審判長 仲間 晃
審判官 辻本 泰隆
審判官 佐久 聖子

- (56)参考文献 特開2009-282714(JP,A)
米国特許出願公開第2013/0232504(US,A1)
特開2011-232916(JP,A)
特開2011-233146(JP,A)
特開2012-208598(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F 9/46- 9/54
G06F11/20