

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5868466号  
(P5868466)

(45) 発行日 平成28年2月24日(2016.2.24)

(24) 登録日 平成28年1月15日(2016.1.15)

(51) Int.Cl. F I  
G O 6 F 11/20 (2006.01) G O 6 F 11/20 3 1 0 A

請求項の数 12 外国語出願 (全 37 頁)

(21) 出願番号	特願2014-168666 (P2014-168666)	(73) 特許権者	392026693
(22) 出願日	平成26年8月21日 (2014. 8. 21)		株式会社NTTドコモ
(62) 分割の表示	特願2012-523039 (P2012-523039) の分割		東京都千代田区永田町二丁目11番1号
原出願日	平成22年7月29日 (2010. 7. 29)	(74) 代理人	100088155
(65) 公開番号	特開2014-238885 (P2014-238885A)		弁理士 長谷川 芳樹
(43) 公開日	平成26年12月18日 (2014. 12. 18)	(74) 代理人	100113435
審査請求日	平成26年8月22日 (2014. 8. 22)		弁理士 黒木 義樹
(31) 優先権主張番号	12/831, 119	(74) 代理人	100121980
(32) 優先日	平成22年7月6日 (2010. 7. 6)		弁理士 沖山 隆
(33) 優先権主張国	米国 (US)	(74) 代理人	100128107
(31) 優先権主張番号	61/230, 226		弁理士 深石 賢治
(32) 優先日	平成21年7月31日 (2009. 7. 31)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 信頼性保証のある仮想化インフラストラクチャのためのリソース割振りプロトコル

(57) 【特許請求の範囲】

【請求項1】

物理ノードと物理リンクからなる物理リソースを、プライマリ仮想ノード及び冗長仮想ノードを有する仮想インフラストラクチャに割り振るための方法であって、

1組の物理ノード、前記物理ノードを接続している物理リンク、及び前記仮想インフラストラクチャに対して要求された信頼性要件についての第1の要求を受信するステップと

、  
前記要求された信頼性を提供するための冗長仮想ノードの数を計算するステップであって、nがプライマリ仮想ノードの数、kが冗長仮想ノードの数である場合、n個のプライマリ仮想ノードについてk個の冗長仮想ノードを供給して、n:kの複製を達成するステップと、

1つの仮想インフラストラクチャのn+k個の仮想ノードをそれぞれ異なる物理ノードに割り振るステップと

を含む物理リソース割振り方法。

【請求項2】

前記物理リソースのそれぞれが物理ノードを備える請求項1に記載の方法。

【請求項3】

前記物理リソースが計算容量及び帯域幅を備える請求項1または2に記載の方法。

【請求項4】

前記n個のプライマリ仮想ノードのうちの任意のものについて、前記k個の冗長仮想ノ

10

20

ードがバックアップされる請求項 1 から 3 のいずれか一項に記載の方法。

【請求項 5】

仮想ノード間に挿入するための仮想冗長リンクを決定するステップをさらに含む請求項 1 から 4 のいずれか一項に記載の方法。

【請求項 6】

既存の割り振られた冗長仮想ノードが前記第 1 の要求により要求された冗長仮想ノードと共有されるとき、前記第 1 の要求により要求された信頼性、及び前記既存の割り振られた冗長仮想ノードの信頼性のそれぞれを計算するステップと、

前記計算された信頼性及び使用可能な物理リソースに基づいて、前記第 1 の要求により要求された冗長仮想ノードを前記既存の割り振られた冗長仮想ノードと結合するかどうかを決定するステップと、

前記第 1 の要求により要求された冗長仮想ノードと前記既存の割り振られた冗長仮想ノードとを結合するかどうかを決定する結果に基づいて、前記第 1 の要求により要求された冗長仮想ノードと前記既存の割り振られた冗長仮想ノードとを結合するステップと

をさらに含む請求項 1 から 5 のいずれか一項に記載の方法。

【請求項 7】

前記物理リソースがサーバリソース及びこれらのサーバリソースを接続する物理リンクを含む請求項 1 から 6 のいずれか一項に記載の方法。

【請求項 8】

前記信頼性要件が少なくとも 1 つの用途に関連する請求項 1 から 7 のいずれか一項に記載の方法。

【請求項 9】

仮想ノード間の帯域幅予約をフローとしてモデル化するステップをさらに含む請求項 1 から 8 のいずれか一項に記載の方法。

【請求項 10】

物理ノードと冗長仮想ノード又は冗長物理ノードとの間の双方向のマッピングを二値変数によってモデル化するステップをさらに含む請求項 1 から 9 のいずれか一項に記載の方法。

【請求項 11】

各仮想インフラストラクチャの 1 つの仮想ノードが単一の物理ノードにのみマップされ、1 つの仮想インフラストラクチャの  $n + k$  個の仮想ノードについて、1 つの物理ノードに 1 つ以下の仮想ノードしかマップされない請求項 1 から 10 のいずれか一項に記載の方法。

【請求項 12】

コンピュータによって実行されると、前記コンピュータが物理ノードと物理リンクからなる物理リソースをプライマリ仮想ノード及び冗長仮想ノードを有する仮想インフラストラクチャに割り振る方法を実行するプログラムを格納するコンピュータ読み取り可能な記録媒体であって、前記方法が、

1 組の物理ノード、前記物理ノードを接続している物理リンク、及び前記仮想インフラストラクチャに対して要求された信頼性要件についての第 1 の要求を受信するステップと

前記要求された信頼性を提供するための冗長仮想ノードの数を計算するステップであって、 $n$  がプライマリ仮想ノードの数、 $k$  が冗長仮想ノードの数である場合、 $n$  個のプライマリ仮想ノードについて  $k$  個の冗長仮想ノードを供給して、 $n : k$  の複製を達成するステップと、

1 つの仮想インフラストラクチャの  $n + k$  個の仮想ノードをそれぞれ異なる物理ノードに割り振るステップと

を含むコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【優先権】

10

20

30

40

50

## 【 0 0 0 1 】

[0001]本特許出願は、2009年7月31日に出願された「信頼性保証のある仮想化インフラストラクチャのためのリソース割振りプロトコル(A Resource Allocation Protocol for Virtualized Infrastructure with Reliability Guarantees)」という名称の対応する特許仮出願第61/230,226号の優先権を主張し、参照によって組み込む。

## 【 発 明 の 分 野 】

## 【 0 0 0 2 】

[0002]本発明は、信頼性、仮想化インフラストラクチャ、及びリソース割振りの分野に関し、より詳細には、本発明は、信頼性保証のある仮想インフラストラクチャにおけるリソースの割振りに関する。

10

## 【 発 明 の 背 景 】

## 【 0 0 0 3 】

[0003]通信ネットワークは、物理から仮想へとシフトしつつある。従来、通信ネットワークは、所与のネットワークをサポートするために物理的なインフラストラクチャを使用して構築されている。インフラストラクチャは、ますます仮想になっている。すなわち、専用の物理ネットワークを構築する代わりに、又は、特定の意図によって設計されていないネットワークを他者と共有する代わりに、共有の物理的な基体の上に構築された、特定のカスタマイズされたプロトコルを備える専用ネットワークの外観をそのユーザに提供するカスタマイズされたネットワークの仮想ネットワークが使用されている。仮想ネットワークは、その(仮想の)オペレータにとってのプライベートネットワークであり、アンダーレイは、様々なオペレータの間で共有される。

20

## 【 0 0 0 4 】

[0004]仮想化は、物理リソースが現在使用されている方法を急速に変えている。元々、サーバを切り離し、物理サーバにわたってリソースを共有するために設計されており、仮想化は、サーバをソフトウェアによって完全に定義できるようにすることによって、高速で素早い配置及びマイグレーションを提供する。これは、演算処理を弾性のあるリソースに変え、他の商業的なエンティティの間でも急速に人気が出つつある。仮想化のパラダイムは、ネットワークに及んでいる。例えば、仮想化のパラダイムによって、複数の研究グループは、惑星規模のネットワークの異なる仮想スライスにわたって複数のオーバーレイの叩き台を稼働させることができる。事実、仮想化インフラストラクチャがサービスとして提供される場合、仮想化によってもたらされる機敏さ及び柔軟性によって、インフラストラクチャプロバイダとサービスプロバイダとの間の効率の良い分離によって、次世代のインターネットを融通が利くものにすることができることを研究者は確信している。

30

## 【 0 0 0 5 】

[0005]こうした仮想化されたアーキテクチャの1つの鍵となる側面は、アンダーレイのリソースを上部の仮想ネットワークに適切に割り当てることである。使用するリソースは、仮想化されるので、物理的なアンダーレイにおける異なるスポットに配置することができ、ネットワークの最高のパフォーマンスのために、物理に対する仮想リソースの慎重な割振りが重要である。適切になされると、各仮想ネットワークはよりよく機能し、物理的なアンダーレイの使用率が増加する(及びしたがって、経費を低減する)。

40

## 【 0 0 0 6 】

[0006]インフラストラクチャが急速に仮想化され、共有され、動的に変化している状態で、強い信頼性を物理的なインフラストラクチャに提供することが重要である。その理由は、単一の物理サーバ又はリンクの故障がいくつかの共有される仮想化されたエンティティに影響を及ぼすからである。信頼性は、冗長性を使用することによって提供される。現在では、信頼性は、リソースを複製することによって提供される。この理由は、信頼性が物理層に設けられているからである。したがって、物理コンポーネントの故障は、別の物理的な要素を持ち出すことによって処理される。仮想化インフラストラクチャにおいて、

50

物理コンポーネントや物理的な要素はバックアップすることを必要とする仮要素であり、物理コンポーネントの故障は、何らかの仮コンポーネントの消失を意味し、これらの仮コンポーネントを他の物理コンポーネント上に配置し直さなければならない。

【 0 0 0 7 】

[0007]信頼性を提供することは、多くの場合、計算、ネットワーク、及びストレージの能力を過度に供給し、付加的な堅牢性のために負荷バランシングを使用することにつながる。こうした高可用性システムは、例えば、リンク又はノードの故障をルート変更する間のネットワークフローの再起動、又は、部分的なジョブがノード故障時に再開するなど、大きい不連続性が許容され得る用途に適している。高レベルのフォールトトレランスは、何らかの故障がシステムの現在の状態にかなりの影響を及ぼす用途で必要とされる。例えば、アドミッションコントロール、スケジューリング、負荷バランシング、帯域幅ブローキング ( b a n d w i d t h b r o k i n g )、AAA、又はネットワーク状態のスナップショットを維持する他のNOC動作を実行するサーバを有する仮想ネットワークは、全体的な故障を許容することができない。例えばMapReduce、PVMなど、マスター-スレーブ/ワーカーアーキテクチャ ( m a s t e r - s l a v e / w o r k e r a r c h i t e c t u r e ) では、マスターノードでの故障は、スレーブ/ワーカーにおいてリソースを浪費する。

10

【 0 0 0 8 】

[0008]ネットワークの仮想化は、ネットワークの運転費及び管理の複雑さを低減する有望な技術であり、研究対象の増加を受けている。インフラストラクチャのプロバイダがより単純な、より安価なコモディティハードウェア ( c o m m o d i t y h a r d w a r e ) 上に各自のネットワークを仮想化する方向に進むにつれて、信頼性はますます重要な問題にならざるを得ない。

20

【 0 0 0 9 】

[0009]ネットワークの信頼性を研究するために、「shadow VNet」、すなわち、パラレルの仮想化スライスの使用を考慮に入れているものもある。しかし、こうしたスライスは、バックアップとして使用されるのではなく、監視ツールとして使用され、故障時にはネットワークをデバッグする方法として使用される。

【 0 0 1 0 】

[0010]一方、サーバ仮想化レベルでのノードのフォールトトレランスに目標が定められる研究がいくつかある。少なくとも1つは、ハイパーバイザにおけるフォールトトレランスを導入した。同じ物理ノードにある2枚の仮想スライスがハイパーバイザを介して同期して動作するようにできる。しかし、これは、せいぜいソフトウェアの故障に対する信頼性を提供するだけである。その理由は、スライスが同じノードにあるからである。

30

【 0 0 1 1 】

[0011]仮想スライスがネットワークにわたって複製され、遷移されるようにするために発達したものもある。異なるタイプの用途 ( ウェブサーバ、ゲームサーバ、及びベンチマークの用途 ) のために、様々な複製技術及びマイグレーションプロトコルが提案された。別のシステムは、ある期間にわたる2つの仮想ノードの間の状態の同期を可能にする。したがって、信頼性のためにネットワークにわたって分散される冗長な仮想ノードを有することは実際に可能である。しかし、これらのソリューションは、ネットワークのどこかに冗長ノードが存在するのに、( 計算能力における ) リソース割振り問題に対処していない。

40

【 0 0 1 2 】

[0012]基本的なレベルで、ノード及びリンクの信頼性に対処する冗長ノードのトポロジを構成するための方法がある。一部の入力グラフに基づいて、追加のリンク ( 又は帯域幅予約 ) は、最小数だけで良いように、最適に導かれる。しかし、これは、大部分はステートレスであるマルチプロセッサシステムのフォールトトレランスを設計することに基づく。この場合、ノードの故障は、最初のトポロジを保護するために、残りのノードの中のマイグレーション又は回転を伴う。これは、マイグレーションが故障によって影響を受けな

50

いネットワークの一部に中断を引き起こす場合がある仮想化ネットワークのシナリオにおいて適切でない場合がある。

【 0 0 1 3 】

[0013]フォールトトレランスは、データセンターにおいても設けられる。冗長性は、過剰な大量のノード及びリンクに関する。故障の回復のために、いくつかのプロトコルが定義されるが、信頼性保証はほとんどない。

【発明の概要】

【 0 0 1 4 】

[0014]本明細書において、リソース割振りプロトコルのための方法及び装置が開示される。一実施形態において、この装置は、物理リソースをプライマリ及び冗長仮想インフラストラクチャに割り振るためのリソース割振りエンジンを備え、リソース割振りエンジンが仮想インフラストラクチャを割り振るとき、冗長な仮想インフラストラクチャの物理リソースが複数のプライマリ仮想インフラストラクチャにわたって共有される。

10

【 0 0 1 5 】

[0015]本発明は、以下の詳細な説明から、及び本発明の様々な実施形態の添付の図面から、より完全に理解されるが、こうした図面及び実施形態は、本発明を特定の実施形態に制限するものと見なされないものとし、説明及び理解のためのものである。

【図面の簡単な説明】

【 0 0 1 6 】

【図 1】フォールトトレラントアーキテクチャの比較を示す図である。

20

【図 2】99.999%のノードの信頼性のために必要とされる冗長ノードの数を示す図である。

【図 3】 $n : k$ の複製がサポートできるノードの数を示す図である。

【図 4】1つのバックアップノード及びそれぞれ割り振られたフェールオーバー帯域幅を備えるVIの一例を示す図である。

【図 5】冗長性をプールし、仮想ノードを拡散させる一例を、仮想化データセンターにある4つのVIを示すことによって示す図である。

【図 6】2つのVIのバックアップノードをプールするときのこのトレードオフを示す図である。

【図 7】管理アーキテクチャの一実施形態を示す図である。

30

【図 8】入ってくる各要求にサービスを提供するプロセスの一実施形態のフロー図である。

【図 9】重複する帯域幅予約の一例を示す図である。

【図 10】冗長ノードの位置が固定された拡張グラフの一例を示す図である。

【図 11】コンピュータシステムを示すブロック図である。

【発明を実施するための形態】

【 0 0 1 7 】

[0016]以下、 $n : k$ の冗長アーキテクチャが開示され、この場合、 $k$ 個の冗長リソースは、 $n$ 個のプライマリリソースのうちの任意のものについてのバックアップであり、複数の仮想インフラストラクチャ(VI)にわたってバックアップを共有することができる。例えば、 $n_1$ 及び $n_2$ のコンピューティングノードを有する2つのVIは、 $k_1$ 及び $k_2$ の冗長性がそれぞれ $r_1$ 及び $r_2$ の保証された信頼性となることを要求する。バックアップを共有することは、同じレベルの信頼性を有する $k_0 < k_1 + k_2$ の冗長性を達成し、フォールトトレランスのために供給されるリソースを低減する。さらに、冗長ノードが、保証された接続性、帯域幅で、ほとんど中断なく、故障したノードを引き継ぐことができるように、ジョイントノード及びリンク冗長性がある。リンクの故障は、同じ機構を介して回復することができる。

40

【 0 0 1 8 】

[0017]物理リソース(例えば、計算容量、ストレージ、及び帯域幅)を同時にプライマリ及び冗長VIに静的に割り振る方法も本明細書に開示される。この方法は、既存の冗長

50

ノードを使用して、冗長仮想リンクの帯域幅をできるだけ多く重複させることによって冗長性に割り振られるリソースを低減することを試みる。

【0019】

[0018]さらに、物理リソースの使用を最低限に抑える、又は大幅に低減し、対応することができる仮想リソースの数を最大にすることを試みる方法で信頼性保証を提供する、物理的な基体上に仮想インフラストラクチャのリソースを割り振る機構が開示される。

【0020】

[0019]信頼性を物理リソースの割振りに組み込み、冗長ノードをいくつかの仮想ネットワーク中で共有することは、信頼性のためのリソースの量を大幅に低減する。

【0021】

[0020]以下の説明では、本発明のより完全な説明を提供するために、多数の詳細が記載される。しかし、これらの具体的な詳細なしに本発明を実践できることは、当業者であれば明らかである。他の例において、周知の構造及びデバイスは、本発明を不明瞭にすることを回避するために、詳細にはではなく、ブロック図の形式で示される。

【0022】

[0021]以下の詳細な説明の一部は、コンピュータメモリ内のデータビット上の動作のアルゴリズム及び記号表現の形で示される。これらのアルゴリズムの説明及び表現は、他の当業者に最も効果的にその作業の要旨を伝えるために、データ処理技術に熟練した人々によって使用される手段である。アルゴリズムは、ここでは、また一般的に、所望の結果をもたらす自己矛盾のない一連のステップであると考えられる。こうしたステップは、物理量の物理的操作を必要とするものである。必須ではないが、通常、これらの量は、格納され、移動され、結合され、比較され、そうでなければ操作され得る電気又は磁気の信号の形をとる。主に一般的な使用の理由で、これらの信号をビット、値、要素、シンボル、文字、項、数字等と呼ぶことが時として便利であることがわかっている。

【0023】

[0022]しかし、これらの及び類似の用語はすべて、適切な物理量に関連付けられたものであり、こうした物理量に適用される便宜上のラベルに過ぎないことを理解されたい。特に明記しない限り、以下の説明からわかるように、説明の全体にわたって、「処理する」、「計算する」、「算出する」、「決定する」、又は「表示する」等の用語を使用している議論は、コンピュータシステムのレジスタ及びメモリ内の物理（電子）量として表されるデータを操作し、コンピュータシステムメモリ若しくはレジスタ、又は他のこうした情報ストレージ、送信又は表示デバイス内で物理量として同じように表される他のデータに変換する、コンピュータシステム又は類似の電子コンピューティング装置の動作及びプロセスを指すことを理解されたい。

【0024】

[0023]本発明は、本明細書において動作を実行するための装置にも関する。この装置は、要求される目的のために特別に構成することができ、又は、コンピュータに格納されたコンピュータプログラムによって選択的に稼働される、又は再構成される汎用コンピュータを備えることができる。こうしたコンピュータプログラムは、例えば、それだけには限定されないが、コンピュータシステムバスにそれぞれ結合される、フロッピーディスク、光ディスク、CD ROM、及び光磁気ディスクを含む任意のタイプのディスク、読み取り専用メモリ（ROM）、ランダムアクセスメモリ（RAM）、EPROM、EEPROM、磁気又は光学カード、又は電子命令の格納に適した任意のタイプの媒体など、コンピュータ可読記憶媒体に格納することができる。

【0025】

[0024]本明細書において提示されるアルゴリズム及び表示は、任意の特定のコンピュータ又は他の装置に本質的に関連がない。様々な汎用システムは、本明細書における教示によるプログラムによって使用することができ、又は、必要な方法ステップを実行するためにより専門の装置を構成することは便利であることがわかり得る。様々なこれらのシステムに必要な構造が下記の説明から明らかになる。さらに、本発明は、任意の特定のプログ

10

20

30

40

50

ラミング言語を参照して記載されていない。本明細書において記載されている本発明の教示を実施するために、様々なプログラミング言語を使用することができることを理解されたい。

【0026】

[0025]機械可読媒体は、マシン（例えば、コンピュータ）によって可読の形式で情報を格納又は送信するための任意の機構を含む。例えば、機械可読媒体は、読み取り専用メモリ（「ROM」）、ランダムアクセスメモリ（「RAM」）、磁気ディスク記憶媒体、光記憶媒体、フラッシュメモリ記憶デバイスなどを含む。

概要

【0027】

[0026]一実施形態において、リソースは、仮想インフラストラクチャ要求の信頼性要件を考慮する割振り方法を使用して割り振られる。物理リソースを仮想リソース要求に割り振るための方法が存在するのに対して、本明細書に記載される割振り方法は、まず、明確な信頼性保証を提供する。

【0028】

[0027]一実施形態において、割振り機構は、1組のリソース（例えば、サーバリソース）又はその一部、これらのリソースを接続しているリンクについての要求、及び信頼性要件、例えば99.999%の動作可能時間を受信する。一実施形態において、要求は、 $G = (V, E, r)$ として表され、この場合、 $V$ はノード、 $E$ はノードを接続しているリンク、 $r$ は信頼性である。次いで割振り機構は、要求された信頼性を提供するために要求に追加する冗長ノードの数を計算する。この要求を別のものと結合することができ、そうすることからの利点がある場合、割振り機構は2つの要求を結合する。一実施形態において、割振り機構は、割振り要求を結合することが有益かどうかを決定する。したがって、一実施形態において、割振り機構は、要求を集約し、冗長性のために取っておかれる物理リソースの量を低減し、場合によっては最低限に抑えることを試みる方法で物理的な冗長リソースを割り振る。

【0029】

[0028]一実施形態において、要求に追加する冗長ノードの量、及びエッジの間に挿入するリンクを決定した後、割振り機構は、新しい要求 $G' = (V', E')$ を計算し、従来のマルチフローコモディティ問題（*traditional multi-flow commodity problem*）を使用して、この要求を割り振る。マルチフローコモディティ問題は、公知技術である。

[ネットワークモデル]

【0030】

[0029]本明細書のために、コンピューティング及びネットワークのリソースを仮想化し、分離し、いくつかのエンティティにわたって共有することができる物理的なネットワークインフラストラクチャが使用される。物理的なインフラストラクチャからのリソースについての要求が、計算ノードの容量及び好ましい場所、ノード間の帯域幅、及び要求されたノード（及びそれらのリンク）のサブセットにおけるあるレベルの信頼性に関して定義される。入ってくる各リソース要求は、冗長なインフラストラクチャと共に静的に割り振られる。

【0031】

[0030]一実施形態において、物理ネットワークは、無向グラフ $G^P = (N^P, E^P)$ としてモデル化され、この場合、 $N^P$ は物理ノードの組であり、 $E^P$ はリンクの組である。各ノード $u \in N^P$ は、 $M_u$ の使用可能な計算容量を有する。各無向リンク $(u, v) \in E^P$ 、 $u, v \in N^P$ は、 $H_{u,v}$ の使用可能な帯域幅容量を有する。マルチフローコモディティ問題を単純化するために、物理ノードにおける故障は、確率 $p$ で、独立しており一様であると見なされる。

【0032】

[0031]リソース要求は、追加の特性を有する無向グラフ $G^V = (N^V, E^V)$ としてモ

10

20

30

40

50

デル化される。 $N^V$  は、1組の計算ノードであり、 $E^V$  は1組のエッジである。 $\mu_x$  は、ノードごとの計算能力要件、 $x \in N^V$  であり、ノードの間の帯域幅要件は  $\mu_{x,y}$ 、 $(x, y) \in E^V$  及び  $x, y \in N^V$  である。さらに、 $[x] \in N^P$  は、仮想ノード  $x$  をマップすることができる追加の制約である。すなわち、仮想ノードを物理ノード上に何らかの特定のマッピングを課すために、 $x$  を物理ノードのサブセットにしかマップできないように、制約  $[x]$  として指定される。このことは物理的位置の選好（本文で述べるように）又は物理ノードタイプ（CPUノード、ストレージノード、ルータノード）に起因し得ることに留意されたい。これは、任意の物理的位置の選好、例えば、入口及び出口の仮想ルータ、他のノードとの近さなどを表す。後述するように、この組は、すでに設けられている別のVIから冗長ノードを再利用する/共有するためにも利用される。各要求も1組のクリティカル仮想ノード  $C^V \subseteq N^V$  及びその関連のリンク  $\{(c, x) \mid (c, x) \in E^V, c \in C^V, x \in N^V\}$  から成り、信頼性  $r$  で保護される。本明細書のために、冗長ノードの組は、 $N^K$  と示される。

10

## 【0033】

[0032] 整合性のために、 $i, j$  を使用して任意のタイプのノードを表し、 $x, y, z \in N^V$  を使用して仮想ノードを表し、 $u, v, w \in N^P$  を使用して物理ノードを表し、 $c, d \in C^V$  を使用してクリティカルノードを表し、 $a, b \in N^K$  を使用して冗長ノードを表す。

[冗長性のための仮想アーキテクチャ]

## 【0034】

20

[0033] 一実施形態において、冗長性のためのアーキテクチャは、以下の特性を有する。

## 【0035】

$n : k$  の冗長性。最高で

## 【数1】

$$\frac{1}{k+1},$$

の使用率をもたらす  $1 : k$  レベルの冗長性を有することに対して、 $n$  個のプライマリソースの任意のものについて、 $k$  個の冗長ノードをバックアップすることができること、より良い粒度及び使用率を達成することができる。

30

## 【0036】

ジョイントノード及びリンク冗長性。ノードが故障すると、冗長ノードが保証された接続性、帯域幅で、ほとんど中断なく引き継ぐように、冗長ノード及びリンクが一緒に供給される。

## 【0037】

交わらない位置。同じ物理ノードにおいてホストすることができる仮想又は冗長ノードは1つ以下である。

## 【0038】

[0034] のように、 $|N^V| + k$  個から、物理ノードの故障の数が  $k$  個を上回る確率が  $1 - r$  を下回らないように、 $k$  個の冗長仮想ノードが供給される。言い換えれば、信頼性は、

40

## 【数2】

$$r \leq \sum_{i=0}^k \binom{n+k}{i} p^i \bar{p}^{n+k-i} \quad (1)$$

$$= I_{\bar{p}}(n, k+1)$$

として得られ、この場合、

## 【数3】

$$n = |N^V|, \bar{p} = 1 - p.$$

50



である。右辺の和は、公知技術である、規則化された不完全なベータ関数  $I_x(\cdot, \cdot)$  に相当する。

【0039】

【0035】一実施形態において、故障時に  $k$  個の冗長ノードすべてに対して、十分な量の計算能力及び帯域幅が使用可能である。したがって、リンク及びノードの故障について、回復手順は、 $k$  個の冗長ノードのうちの1つ又は複数を持ち出し、確保された冗長リソースを使用するように動作する。一実施形態において、これによってさらに中断が生じる場合があるため、仮想ノードのマイグレーション又はスワッピングは回復を援助することができない。さらに、冗長ノードも故障する場合があるため、冗長ノードは、式(1)において述べられる信頼性を達成するために、任意のノード  $c \in C^V$  の代わりとなることができなければならない。上述したように、帯域幅予約においてパス分割が使用され、リンクについての別の層の保護及びグレースフルデグラデーションが提供される。

10

【0040】

【0036】冗長性のための帯域幅予約は、 $N^K$  のノードから発せられる1組の重み付けされた無向仮想リンク  $L$  としてモデル化される。

【0041】

$$L = N^K \times (N^V \times N^K) = (N^K \times N^V) \times (N^K \times N^K) \quad (2)$$

すなわち、 $L$  は、 $N^V$  の頂点を有する、それ自体の間の冗長ノードからのリンクを含む、2つの2部グラフの結合である。これらのリンクは、仮想ネットワークの埋め込みのために  $G^V$  に追加される。より形式的に、 $L$  は、下の2つの定理によって定義される。

20

【0042】

定理1。  $a \in N^K$  及び  $x \in N^V$  が与えられる。その場合、

【数4】

$$(a, x) \notin L \text{ iff } \neg \{(c, x) \in E^V, c \in C^V\}.$$

である。これは、クリティカルリンク  $(c, x) \in E^V$  を  $L$  におけるリンク  $(a, x)$  によってバックアップする必要があり、したがって、故障のために  $c$  が  $a$  に遷移された場合、 $x$  は、リソースの新しい位置にまだ接続されていることを示す。

【0043】

証明：

30

【数5】

$$(a, x) \notin L$$

及び仮想リンク  $(c, x) \in E^V$  が存在し、 $c \in C^V$  であると仮定する。

その時、アーキテクチャは、 $n : k$  の冗長性を有していない。その理由は、 $c$  が故障した場合、 $a$  は、 $x$  に供給される帯域幅を有していないからである。同様に、

【数6】

$$\{(c, x) \in E^V, c \in C^V\}$$

であり、 $(a, x) \in L$  である場合、 $c$  が故障した場合、 $(a, x)$  のために供給される帯域幅は決して使用されない。

40

【0044】

推論1。

【数7】

$$(a, c) \notin L \text{ iff } (c, d) \notin E^V$$

であり、この場合、 $a \in N^K$  及び  $c, d \in C^V$  である。

【0045】

証明：これは、 $x$  の領域を  $C^V$  に限定することによる定理1からの直接的な結果である。

50

【 0 0 4 6 】

[0037]上記は、L が 2 部グラフ  $L^1$  から成ることを明示する。

【 0 0 4 7 】

$$L^1 = \{ (a, x) \mid a \in N^k, c \in C^v, (c, x) \in E^v, x \in N^v \} \quad (3)$$

【 0 0 4 8 】

定理 2。a, b  $\in N^k$  が与えられた場合、

【数 8】

$$(a, b) \notin L \text{ iff } \exists (c, d) \in E^v$$

c, d  $\in C^v$  である。これは、クリティカルノードの間にリンクがある場合、各冗長ノード間のリンクが存在しなければならないことを示す。

【 0 0 4 9 】

証明：L において a 及び b が接続されておらず、しかし、リンク  $(c, d) \in E^v$  が存在すると仮定する。その時、c 及び d が故障し、a 及び b に遷移される場合、帯域幅保証はない。逆に、 $(a, b) \in L$  であり、

【数 9】

$$\exists (c, d) \in E^v$$

である場合、 $(a, b)$  のために供給される帯域幅は決して使用されない。

【 0 0 5 0 】

[0038]これによって、任意の 2 つのクリティカルノードの間にリンクがある限り、 $N^k$  の冗長ノード間に完全グラフを含む L が得られる。冗長ノード間の完全グラフを  $L^2$  によって示す。

【 0 0 5 1 】

$$L^2 = \{ (a, b) \mid a \neq b, a, b \in N^k \} \quad (4)$$

【 0 0 5 2 】

[0039]  $L = (N^k \times N^v) \cup (N^k \times N^k)$  であるため、冗長リンクの最小の組は、

【数 10】

$$L = \begin{cases} L^1 & , (x, y) \in E^v, \forall x, y \in C^v \\ L^1 \cup L^2 & , \text{それ以外の場合} \end{cases} \quad (5)$$

によって与えられる。この結果は、他の提案されたアーキテクチャより多くのリンクを必要とする。しかし、後者の結果は、故障後の回復されたグラフが  $G^v$  を含むという仮定に基づく。このことは、故障によって影響を受けないノードを  $G^v$  の最初のトポロジの回復のために遷移する必要がないことを確実にしない。この追加の制約は、L を構成する際、考慮に入れられる。それにもかかわらず、この制約が必要ない場合、L を他のソリューションと置き換えることができる。

【 0 0 5 3 】

[0040]マルチコモディティフロー (MCF) によって帯域幅が供給されている下記の場合、冗長フローを可能な限り重複させることによって、帯域幅が低減され、又は最低限に抑えられる。これらの重複は、MCF モデルへの制約とし捕捉される。

[冗長性の共有]

【 0 0 5 4 】

[0041]以下、n : k のフォールトトレラントアーキテクチャの利点、及び冗長ノードの共有が使用率をどれだけ増加させ得るかの表示が開示される。説明を簡単にするために、本明細書では  $C^v = N^v$  と仮定する。

【 0 0 5 5 】

[0042]図 1 A の小さい 3 ノードの仮想ネットワークについて考慮する。図 1 B に示されるように、1 ノードの故障を許容する簡単で端的な方法は、1 : k の許容値を使用するこ

10

20

30

40

50

とであり、すなわちあらゆる仮想マシンをいったん複製し、複製及びプライマリノードへの論理リンクを作成することである。k<sub>s</sub>ノードの故障のフォールトトレランスは、k<sub>s</sub>層の複製によって達成することができる。各物理ノードの故障の確率がpであると仮定すると、n個のノード及びe個のリンクのネットワークについて、ノードの信頼性rを達成するために必要な冗長性の層の数は、

【数11】

$$k_s = \left\lceil \frac{\log \left( 1 - r^{\frac{1}{n}} \right)}{\log p - 1} \right\rceil \quad (6)$$

10

となる。残念なことに、簡単な複製は、システムにあまりに多くの冗長ノード及び論理リンクを追加し、それぞれk<sub>s</sub>n及びk<sub>s</sub>n + 3k<sub>s</sub>eとなる。

【0056】

[0043]この方法と、冗長ノードが、すなわちn:kである、図1Cにおけるの3つのノードのいずれかについてのバックアップである他の手法とを比較する。ノード信頼性rは、より細かい粒度、及びより少ない数の冗長ノード及びリンクによって提供することができ、(1)によって得られる。冗長ノードの数及びリンク(最大)は、k及び

20

【数12】

$$kn + \frac{k}{2}(k-1).$$

である。図2は、99.999%のノードの信頼性に必要とされる冗長ノードの数について、2つの手法の間の比較を示す。

【0057】

[0044]予想されるように、冗長ノードの数は、同じレベルの信頼性について、n:kの複製より1:kの複製でかなり速く増える。実際に、図3でわかるように、n:kの手法では、十分に拡張される。1つの興味深い傾向は、kの値が小さい場合、nが超線形であるということである。例えば、故障の確率p=0.01の曲線について、95ノードのVIは、99.999%の信頼性のためにk=7を必要とし、190ノードのVIは、k=10を必要とする。一見したところ、2つの95ノードVIを割り振るときに、7つずつの冗長ノードを供給するよりも、k=10の冗長ノードを共有することに価値がある。

30

【0058】

[0045]大きいkでは、nが直線的に増えることに留意されたい。共有はもはや、冗長ノードの数を減らさない。n対kの線形の挙動は、結合も有害でないことを意味することにも留意されたい。冗長リンクの数が少なくともnkであると想定すれば、冗長ノードを共有するときに、より多くの帯域幅が確保される。一方で、小さいkでは、より多くの冗長リンクについて、冗長ノードの数の減少がトレードオフされる。

40

【0059】

[0046]冗長ノードを共有するのに価値がある方法が2つある。

【0060】

1) kにおける離散型のジャンプ(discrete jump)の使用。例えば、12ノードのVIは、99.999%の信頼性のために4つの冗長ノードを使用する。同じ4つのノードは、同じレベルの信頼性のために、別の13個のプライマリノードをサポートすることができる。

【0061】

2) 異なるレベルの信頼性を必要とするVI間での非対称の共有。例えば、21ノード

50

の V I は、99.999% の信頼性のために、4 つの冗長ノードを使用する。4 つの冗長ノードのうちの一つを、99.9% の信頼性のために、別の 5 ノードの V I と共有することができる。( m 個の他の V I と共有するとき ) より大きい V I の信頼性は、下記のように計算することができる。

【数 13】

$$r_{00} = 1 - \sum_{x=0}^k Pr( k \text{ 個のうち } x \text{ 個のバックアップがダウンするまたは } V_{Inf-1}, \dots, V_{Inf-m} \text{ によって使用されている } ) \times$$

$Pr( k_0-k \text{ 個のバックアップを有する } V_{Inf-0} \text{ から } k_0-x \text{ 個超のノードが故障する } )$

10

【0062】

[0047] 第 1 の方法と比較すると、共有後、k は不変のままであるため、共有のこれら 2 つの方法はより良好である。このことは、V I が順次割り振られる場合、稼働している V I が再構成を必要としないことを確実にする。

[ 仮想データセンターに対するリソース割振りの管理及び適用 ]

【0063】

[0048] 一実施形態では、最初の管理アーキテクチャは、仮想化データセンターにおいて、仮想エンティティ (例えば、ホストされたサービス) の信頼性保証及びリソースを自律的に管理する。このアーキテクチャにおいて、追加の仮想バックアップノード及びその関連のリンクは、任意のレベルの信頼性保証のために適切に調整される。アイドルの冗長ノードを有するにもかかわらず、より多くの物理リソースが入ってくる新しいサービスが利用できるように、データセンター全体への冗長性のプールは集合的に管理される。さらに、一実施形態において、いくつかのコンポーネントの故障がデータセンター全体を低下させないように、アーキテクチャは、故障に対して障害許容量があるように設計される。

20

【0064】

[0049] 一実施形態において、冗長機構は、仮想化データセンターにおいて、顧客当たりのレベルで、フォールトトレランスをサポートする。以下は、主要な要求に使用されるリソースの管理の概要、及び仮想化データセンターに適用するための追加の冗長性を提供する。しかし、これらの技術を他の仮想化された環境に適用することができること、及びこ

30

仮想化データセンターのためのリソース要求モデルの一例

【0065】

[0050] リソース要求モデルは、例えばアマゾン EC2 クラウドサービス及び他のクラウドサービスプロバイダなど、その物理リソースをリースする仮想化データセンターのものなど、リソースを要求するためのものである。独立サーバインスタンスをリースするより、一実施形態において、リソース要求のモデルは、

【0066】

1) 最低限の CPU 容量の要件を有するワーカー及びマスターノード、並びに

【0067】

2) これらのノード間の帯域幅保証

を含む全仮想インフラストラクチャ ( V I ) に対応する。

40

【0068】

[0051] 一実施形態において、ワーカーノードは基本的にデータプロセッサ/ナンバークランチャであり、マスターノードはワーカーノードの機能を調整するサーバである。複数のサーバを有する V I は、複数のマスターノードを有する。さらに、マスターノードが故障のクリティカルポイントであるため、各 V I 要求がマスターノードに関する信頼性保証を要求する。これは、帯域幅保証を加重されたノード間のエッジとする加重グラフとしてモデル化することができ、マスターノードはサブグラフを形成する。このモデルは、様々なニーズを表すのに十分包括的である。

50

【 0 0 6 9 】

[0052]データセンターのオペレータは、すべての現在のリース及び入ってくる新しいV I 要求を管理することを必要とする。

信頼性のための仮想バックアップノード

【 0 0 7 0 】

[0053]一実施形態において、クリティカルなマスターノードに関する信頼性を保証するために、管理アーキテクチャは、空きのあるCPU及び帯域幅の容量を有する追加のバックアップノードを確保する。すべてのクリティカルなマスターノードの状態を複製し、周知の最適化された同期技術を使用してすべてのバックアップノードに同期することができる。ノードの故障の場合、どのバックアップノードも、故障したノードを置き換えるために「ホットスワップ」の用意ができています。

10

【 0 0 7 1 】

[0054]図4は、1つのバックアップノード及びそれぞれ割り振られたフェールオーバ帯域幅を備える4ノードのV I の一例を示す。図4を参照すると、クリティカルノードA及びBに関して、バックアップノード(黒)及び信頼性のために確保されている帯域幅(点線)が示されている。リンク上の数字は、確保されている帯域幅を表す。このようにして、バックアップノードeは、クリティカルノードが故障した場合、クリティカルノードに取って代わることができる。

【 0 0 7 2 】

[0055]これは、任意のk個のバックアップノードがn個のクリティカルノードをカバーするように、容易に拡張することができる。例えば、pによって物理ノードの故障の確率を定義する。pが物理ノードごとに独立同分布であると推定する。この場合、n個のクリティカルノードに関する信頼性rは、以下のように計算される。

20

【数14】

$$r \leq \sum_{i=0}^k \binom{n+k}{i} p^i \bar{p}^{n+k-i} \quad (1)$$

n + k 個のノードのうち1つ以下が同じ物理ノード上にホストされると仮定される。下記の表Iは、2%の物理ノードの故障率の場合に様々な信頼性保証の下でバックアップノードの数によってサポートすることができるクリティカルノードの最大数を示す。

30

【表1】

TABLE I  
NUMBER OF BACKUP NODES REQUIRED,  $p = 2\%$

Max. no. of critical nodes (n)				No. of backups (k)
99.9%	99.99%	99.999%	99.9999%	
1	0	0	0	1
8	3	1	0	2
19	10	4	2	3
34	20	11	6	4

40

バックアップノードの数は、バックアップノードの増加が準直線であるため、表Iに示される範囲について十分拡張する。したがって、冗長性を低減することができるように、いくつかのV I にわたってバックアップノードを一緒にプールのことは有益であり、このことは、より良いリソースの使用率につながる。

分散された割り振り、冗長マッピング、及び同期

50

## 【 0 0 7 3 】

[0056]一実施形態において、効果的に信頼性を管理するために、同じV Iのすべての仮想ノード、及びそれぞれのバックアップノードは、同じ物理ノード上にホストされない、すなわち、可能な限りデータセンターにわたって分散される。図5は、冗長性をプールの、仮想ノードを分散させる一例を、仮想化データセンターにある4つのV Iを示すことによって示す。仮想ノードは、同じ物理ノードにスタックされない。クリティカルノードは、その状態をバックアップノードに同期させる。一実施形態において、仮想ノードでのほぼシームレスな動作のために、同期機構がハイパーバイザ層に設けられる。信頼度のレベルに応じて、あるクリティカルノードは、図5のV I 4の場合と同様に、1を超える仮想バックアップノードを必要とする、又はV I 3のようにまったく必要としない場合がある。時として、逆もあり得る。仮想バックアップノード1は、複数のV Iからの複数のノードをカバーし、CPUリソースを保護する。このことは、本明細書において、冗長性のプールと呼び、後述する。

10

仮想データセンターの冗長性のプール

## 【 0 0 7 4 】

[0057]冗長ノードを一緒にプールし、それらをいくつかのV Iにわたって共有することによって、バックアップノードの総量及びしたがってアイドルのCPU容量を低減するために、表Iのnとkの間の準線形の関係が利用される。例示のために、99.999%の信頼性保証のある5つのクリティカルノードを有するV Iは、4つのバックアップノードが確保されていることを必要とする。同じレベルの信頼性のために、同じ4つのノードが最高11個のクリティカルノードをサポートすることができるため、最高6つのクリティカルノードを有する別のV Iは、追加のバックアップを確保することなく、同じ4つのノードを使用することができる。

20

## 【 0 0 7 5 】

[0058]しかし、冗長性のプールは、必ずしも「無料」ではない。無計画なプールは、バックアップノードに関係するフェールオーバー帯域幅を確保する際、かなりのコストをもたらす。n個のクリティカルノード及びk個のバックアップノードを含む新しいV I要求に追加される追加のリンクの数は、少なくとも

## 【数15】

$$kn + \frac{k}{2}(k-1),$$

30

であり、第1の項は、バックアップノードとクリティカルノードとの間に確保されるすべての帯域幅を表し、後者は、バックアップノードを相互接続している帯域幅を表す。したがって、冗長性をプールしながらバックアップノードの数を増加させることは、フェールオーバー帯域幅も増加させるため、逆効果である。

## 【 0 0 7 6 】

[0059]図6は、2つのV Iのバックアップノードをプールするときのこのトレードオフを示す。これらの領域の境界は、表Iのすべてのnについて、式 $x + y = n$ 及び $x, y > n$ である。見やすいように、これらの線は、境界点がどの線上にもないように、定数によってシフトされる。領域601は、V I - 1及びV I - 2のバックアップノードをプールすることが完全に価値があるケースを示す。例えば、(5, 6)では、5ノードのV I及び6ノードのV Iがいずれもそれぞれ4つのバックアップノードを必要とする。両方のV Iに同じ4つのバックアップノードを割り当てることによって、99.999%の同じ信頼性を保証する。これは、領域602における場合では、両方のV Iに対するバックアップノードの数を増やすことなしに行うことができず、例えば、(3, 2)では、各V Iは、個々に3つのバックアップノードのみを必要とする。しかし、2つを結合すると、5つのノードが追加のバックアップノードを必要とすることになるため、より多くのフェールオーバー帯域幅を必要とする。領域603は、1つのV Iだけが必要より多くのバックアップノードを必要とするケースを表す。

40

50

## 【 0 0 7 7 】

[0060]異なる信頼性保証をサポートしているバックアップノードは、一緒にプールのこともでき、類似のトレードオフ領域を有する。それはすべて、表 I に示されるバックアップノードの残りの「サポート容量」に依存する。

仮想化データセンターのための障害許容力のあるアーキテクチャ

## 【 0 0 7 8 】

[0061]図 7 は、仮想化データセンターの物理リソースの上にある管理アーキテクチャの一実施形態を示しており、リソース要求を管理する集中化したコントローラとして作用する。各コンポーネントは、個々に機能しているエンティティとなるように設計されており、故障に対して障害許容力を確実にするための方策を有する。

10

## 【 0 0 7 9 】

[0062]物理リソースアカウントングコンポーネント ( p h y s i c a l r e s o u r c e a c c o u n t i n g c o m p o n e n t ) 7 0 6 は、入ってくる新しい要求に対するリソースの割振り時に必要な、仮想化データセンターにおいて割り振られていない残りのリソースを追跡する。価格付けポリシー ( p r i c i n g p o l i c y ) 7 0 5 は、動的な価格付けを容易にするために、物理リソースアカウントングコンポーネント 7 0 6 からのその入力を引き出す。リソース割振りエンジン 7 0 0 及びリソース解放モジュール 7 0 9 のみが、物理リソースアカウントングコンポーネント 7 0 6 を更新することができる。一実施形態において、更新は、要求及びリーベイベントに回答して起きる。

20

## 【 0 0 8 0 】

[0063]一実施形態において、物理リソースアカウントングコンポーネント 7 0 6 のための障害許容力を確実にする方法が 2 つある。( i ) 周知の障害許容力のあるデータベースが使用される、又は ( i i ) データの複数のコピーは、データに対する書き込み及び読み取りをそれぞれマルチキャスト及びエニーキャストとして、別々に格納される。一実施形態において、データは、( P h y N o d e , r C P U ) 及び ( P h y L i n k , r B W ) の形の鍵 - 値の組として格納され、この場合、P h y N o d e 及び P h y L i n k は、物理ノード及びリンクをそれぞれ一意に識別し、r C P U 及び r B W は、使用可能な C P U 及び帯域幅リソースの量をそれぞれ提供する。

30

## 【 0 0 8 1 】

[0064]上述したように、時として、V I にわたってバックアップノードをプールするとき、C P U 及び帯域幅を保護する際のトレードオフが存在する。一実施形態において、信頼性ポリシー 7 0 5 は、入ってくる新しい V I のバックアップノードをデータセンターにおける別の既存の V I と共にプールするべきかどうかを指定する決定ルールを備える。図では、これらのルールは、トレードオフ領域の境界を表す(一例として図 6 を参照)。故障に対する障害許容力を確実にするための戦略は、物理リソースアカウントングコンポーネント 7 0 6 と同様である。

## 【 0 0 8 2 】

[0065]リソース割振りエンジン 7 0 0 は、入ってくる要求に対するリソースをマップし、確保する役割を果たす。図 8 は、入ってくる各要求にサービスを提供するプロセスの一実施形態のフロー図である。プロセスは、ハードウェア(例えば、論理回路、回路など)、ソフトウェア(汎用コンピュータシステム又は専用マシン上で稼働されるものなど)、又は両方の組み合わせを備えることができるリソース割振りエンジンにおいて処理論理回路によって実行される。図 8 を参照すると、プロセスは、処理論理回路が入ってくる要求を受信することによって開始する(処理ブロック 8 0 1)。入ってくる要求を受信することによって、処理論理回路は、入ってくる要求に対応するために必要なバックアップノードの数を計算する(処理ブロック 8 0 2)。その後、処理論理回路は、バックアップノードを既存の V I と共にプールするかどうかを検査する(処理ブロック 8 0 3)。処理論理回路がバックアップノードを既存の V I と共にプールすると決定する場合、処理論理回路は現在のバックアップに新しい制約及び帯域幅を追加し(処理ブロック 8 0 4)、プロセ

40

50

スは処理ブロック 807 に移行する。そうでない場合、処理論理回路は、新しいバックアップノード及び帯域幅を作り（処理ブロック 805）、これらの新しいリソースを、リソース割振りのための最初のリソース要求に追加し、プロセスは、処理ブロック 807 に移行する。処理ブロック 807 で、処理論理回路は、外部ツールを介して線形最適化を解決する。こうしたツールの例には、それだけには限定されないが、COIN CBC、ilog CPLEX、及びMOSEKなどがある。

【0083】

[0066]次に、処理論理回路は、外部ツールの出力を介した使用可能なリソースに基づいてソリューションが実行可能かどうかを決定する（処理ブロック 808）。より詳細には、ソルバー（solver）が稼働され、ソルバーが所与の制約を有するソリューションを探し出すことができない場合、ソリューションは実行可能ではない。ソルバーがソリューションを戻す場合、ソリューションは実行可能である。そうである場合、処理論理回路は、他のコンポーネントを更新し（処理ブロック 809）、プロセスが終了する。そうでない場合、処理論理回路は、リソース要求を拒否し（処理ブロック 810）、プロセスが終了する。

10

【0084】

[0067]一実施形態において、仮想ノードから物理ノード、及び仮想リンクから物理パスへのマッピングの問題は、仮想ノード間の帯域幅予約が物理ノード間のフローであり、仮想ノードと物理ノードとの間のフローの存在がマッピングを示すというマルチコモディティフロー問題として、当分野で周知の方法で定式化される。これは、以下の目的に関する線形最適化問題である。

20

【数16】

$$P1: \min \sum_u \sum_x \alpha_{ux} \rho_{ux} \mu_{ux} + \sum_{u,v} \sum_l \beta_{uv}^l f_{uv}^l \quad (2)$$

式中、 $f_{uv}^l$  及び  $\rho_{ux}$  は変数、 $\rho_{ux}$  はブール変数であり、仮想ノード  $x$  が物理ノード  $u$  にマップされる場合、真である。 $f_{uv}^l$  は、物理リンク  $(u, v)$  上に「流れている」仮想リンク  $l$  の帯域幅の量であり、負ではない。 $f_{uv}^l$  についての従来のフロー保護制約が適用され、従来のフロー保護制約（flow conservation constraints）は、当業者によって十分理解されている。さらに、制約

30

【数17】

$$\sum_x \rho_{ux} = 1 \quad \text{及び} \quad \sum_u \rho_{ux} \leq 1 \quad (3)$$

は、仮想ノードと物理ノードと間の1対1のマッピングを確実にし、仮想ノードは、上述したように分散される。CPU及び帯域幅に関して新しいVIによって消費される全リソースは、残っている物理リソースの量に供され、すなわち、

【数18】

$$\sum_x \rho_{ux} \mu_x \leq rCPU_u \quad \text{及び} \quad \sum_l f_{uv}^l \leq rBW_{uv} \quad (4)$$

40

であり、 $\mu_x$  は、仮想ノード  $x$  によって必要とされるCPU容量である。入力  $\mu_x$  及び  $f_{uv}^l$  はそれぞれ、リソースがリースされる時、データセンターのオペレータに対して、CPU及び帯域幅当たりの純コスト（マイナス収入）を表す。これらは、後述する価格付けポリシーから導出される。

【0085】

[0068]バックアップノード及びフェールオーバー帯域幅の追加は、信頼性ポリシー、すなわちバックアップノードがプールされるかどうか依存する。そうでない場合、問題は、新しいバックアップノード及び帯域幅を含んでP1を解くほど単純である。そうでなければ、マッピング変数  $\rho_{ux}$  への追加の制約は、新しいVIの仮想ノードと現在の仮想ノード

50



ドとの間の重複がないことを確実にするために、P 1 に追加され、すなわち、すべての占有された  $u$  について、 $u_x = 0$  である。

【0086】

[0069] P 1 に対する実行可能なソリューションがある場合、新しい要求にのみ対応することができる。次いで、細い両方向矢印を介してリソース割振りエンジン 700 にリンクされているコンポーネント、すなわち V I マップ 707、ホットスワップマップ 708、及びアカウンティングコンポーネント 706 が P 1 からのソリューションによって更新される。そうでなければ、単に不十分な物理リソースのために、更新が拒絶される。

【0087】

[0070] このメイン制御コンポーネントが障害許容力を有するようになる簡単な戦略は、複数のインスタンスにわたって同じ要求を実行することである。より効果的な方法は、いくつかの要求を処理しているが、低い  $rCPU_u$  及び  $rBW_{uv}$  の値を使用して競合状態を防止するインスタンスを複数有することである。しかし、過度の拒否のリスクがある。

価格付けポリシー

【0088】

[0071] 価格付けポリシーは、P 1 の入力  $u_x$  及び  $^1 u_v$  に影響するリソースの価格を指定する。ここで使用する価格付け戦略に固定する必要はなく、むしろ、価格付けモジュール 705 が設けられ、できる限り一般的である。特に、動的な価格付けがサポートされ、この価格付けは、需要を抑えることができ、リソースのより効率的な使用率をもたらすことができる。ある期間にわたる物理リソースアカウンティングモジュール 706 からの入力及びリソース割振りエンジン 705 からのフィードバックにより、価格付けモジュール 705 は、信頼性保証、物理リソースのタイプ（リンク、ノード）、受容率、及びリース期間の面において仮想 CPU 及び帯域幅の価格を動的に設定することができる。

仮想インフラストラクチャマップ及びホットスワップマップ

【0089】

[0072] V I マップ 707 は、認められるすべての V I、及び仮想エンティティのその物理リソースに対するマッピング、すなわち、その物理サーバ及び確保された CPU の量に対する仮想ノードのマップ、及び物理パス及びそのパスに沿って確保された帯域幅の量に対する仮想リンクのマップを記録する。さらに、V I が使用するバックアップノードのプールも格納される。

【0090】

[0073] ホットスワップマップ 708 は、バックアップノードのすべての現在のプール及びそれぞれの残りのサポート容量を記録する。この情報は、V I マップと共に、入ってくる新しい V I がバックアップノードの既存のプールを使用することができるか、又は新しい V I のための別の新しいプールを作ることができるかどうかをリソース割振りエンジン 700 が決定するのを助ける。一実施形態において、新しい V I のためのマッピングのソリューションがいったん取得されると、リソース割振りエンジン 700 はこれらの 2 つのマップに書き込む。

【0091】

[0074] 一実施形態において、これら 2 つのコンポーネントに対する障害許容力のある戦略は、物理リソースアカウンティングモジュール 706 のものと同じである。その理由は、これらのコンポーネントがデータベースのようなコンポーネントだからである。

リソースの解放

【0092】

[0075] V I によって使用されるリソースは、リースの終了時に解放されなければならない。競合状態を防止するために、リソース割振りエンジン 700 による同期ロックが解除されるまで、リソース解放モジュール 709 は、一時的にこれらのリソースを保持するガベージコレクタとして働く。このコンポーネントが故障する場合、2 つのマップに関する簡単なチェック、及び物理リソースのアカウンティング検証がこのコンポーネントを回復する。

10

20

30

40

50

## 同期及び回復機構

## 【 0 0 9 3 】

[0076]これら2つの機構は、分散された方法で機能しているあらゆる物理ノードにおけるローカルなサービスである。一実施形態において、ノード間の同期は、物理ノードのハイパーバイザで管理され、物理ノード間の監視は、ハートビート、同期信号、又は当業者に周知である他の分散監視方法を介することができる。故障が検出されると、回復手順が開始し、コントロールアーキテクチャで進行中の動作すべての代わりにする。ホットスワップノードは、均一のランダム化によって分散された方法で、各V Iの仮想の隣接者によって選択され、結合は任意に中断される。

## 【 0 0 9 4 】

[0077]したがって、データセンターにおいてホストされる仮想インフラストラクチャにおける信頼性保証を自律的に管理することができるフォールトトレラントアーキテクチャが開示される。ここで、信頼性は、仮想バックアップノードのプール及び確保されたフェールオーバー帯域幅によって保証される。アイドルのCPU容量を保護するために、バックアップがプールされ、帯域幅に対するトレードオフが定義される。バックアップを含めて、すべての仮想エンティティの物理リソースは、線形最適化フレームワークを介して割り振られる。データセンターのリソース使用率を追跡し、考慮に入れる他のコンポーネントも定義される。個々のコンポーネントは、個々に動作するように設計されており、障害に対する許容力を確実にするための方策を有する。

[リソースの割り振り：混合整数計画問題(MIXED INTEGER PROGRAMMING PROBLEM)]

## 【 0 0 9 5 】

[0078]V Iリソース割り振り問題は、マルチコモディティフロー問題(MCF)に類似する混合整数計画問題として定式化することができる。ノード間の帯域幅の要求は、フローとしてモデル化される。物理ノードと仮想ノードとの間のマッピングは、余分な「マッピング」エッジを追加し、当分野で周知の方法で、フロー問題を解決する際に、仮想ノード当たりこうした1つのエッジだけが使用されることを確実にすることによって構築される。

## 【 0 0 9 6 】

[0079]一実施形態において、MCFは、V Iノード及びリンクを物理的なインフラストラクチャにマップするために使用される。しかし、MCFは、(i)バックアップリンクが可能な限り重複することができ、(ii)バックアップノードのマッピングは、好適な1組の物理ノードに限定されるように制約する。アルゴリズム1は、を取得し、V I、及びそのバックアップノード、並びにリンクを物理的なインフラストラクチャにマップして、信頼性rを保証するための手順を列挙する。

## 【 0 0 9 7 】

[0080]図9は、帯域幅予約がどのように重複し得るかの一例を示す。図9を参照すると、左は、2つの冗長ノード(黒)を有する2ノード仮想トポロジである。ノードcは、クリティカルノードである。したがって、ノードxを冗長ノードにリンクしている1単位の帯域幅が確保されていなければならない、それは結果としてリンクDE上の2単位の予約となり得る。しかし、冗長ノードは常にノードcの1つのインスタンスしか引き継ぐことができないので、過度の予約が生じる。本明細書において使用されるMCFの制約形式において、任意のトポロジ上のこれらの重複を最大にする試みがある。

## 【 0 0 9 8 】

[0081]上記のように、既存のV Iのバックアップノードを、入ってくる新しいV Iと共有しながら、確実に不変にすることは価値がある。リソース割り振り手順が以下に提供される。行6~14は、そのバックアップノードを共有する適切なV Iを貪欲に探す。これらのV I候補を、「サポート容量」に関して配列することができる。例えば、k=3のバックアップノードは、r=99.99%のために、8~21の仮想ノードの間をサポートすることができる。8ノードのV Iは、20ノードのV Iより多くのサポート容量を有して

10

20

30

40

50

おり、したがって好ましい。この配列は、すでに共有されているV Iを考慮に入れなければならない。行1 1は、そのバックアップノードが好適な物理的位置に限定されるV Iを埋め込むことを試みる。共有が可能でない場合、バックアップノードは行1 5のように $N^P$ におけるどこからでも選択される。

【表2】

アルゴリズム1: リソース割付け手順の一例

```

1: procedure ALLOCATE(VirInf, PhyInf, r)
2:    $n \leftarrow \lfloor \text{VirInf}.N^V \rfloor$ 
3:   Compute k from n, r, p. > inverse of (1)
4:   Compute L from VirInf given k
5:   VirInf.augment ( $N^K$ , L)
6:   for all  $\zeta$  in PhyInf.getVIs()do > orderd
7:     if  $|\zeta.N^K| \neq k$  then
8:       continue
9:     end if
10:     $\Phi \leftarrow \text{phyLoc}(\zeta.N^K)$ 
11:    if MCF_OL(VirInf, PhyInf,  $\Phi$ ,  $\zeta$ ) = True then
12:      return True
13:    end if
14:  end for
15:  return MCF_OL(VirInf, PhyInf,  $N^P$ , NULL)
16: end procedure

```

【0099】

[0082] MCF問題は、以下の通りに定義される。マッピングのための拡張エッジの組を $R^P$ によって示し、

【0100】

$$R^P = \{ (u, x), (x, u) \mid x \in N^V \cap N^K, u \in [x] \} \quad (7)$$

となり、この場合、各エッジは無限の帯域幅を有する。 $[x] \cap N^P$ は、仮想ノード $x$ をホストすることができる物理ノードの組である。 $x$ がバックアップノードであり、別のV Iのバックアップと共有されることになっている場合、 $[x]$ はアルゴリズム1で定義された $R^P$ に等しい。図10は、この拡張構造の一例を示す。ノード $a$ 及び $b$ は、別のV Iの冗長ノードである。これらのノードは、物理ノードのC及びDにあり、ノード $x$ 及び $y$ を有する新しいV Iの冗長ノードであり得る。同じV Iのノードが同じ物理ノードにホストされない場合があるので、 $[a] = \{C\}$ 、 $[b] = \{D\}$ 、及び $[x] = [y] = \{A, B, E\}$ である。

【0101】

[0083] 3つの組は、以下の通りに定義される。

【0102】

$$N^A = N^P \cap N^V \cap N^K \quad (8)$$

【0103】

$$E^A = E^P \cap R^P \quad (9)$$

【0104】

$$C^K = C^V \cap \{ x \mid x \in N^V, (c, x) \in E^V, c \in C^V \} \quad (10)$$

この場合、 $N^A$ は、すべての仮想、物理、及び冗長ノードの組であり、 $E^A$ は物理及びマッピングエッジの組であり、 $C^K$ は、冗長ノードが $L^1$ においてリンクされるノードの組である。

【0105】

[0084]一実施形態において、仮想ノードとバックアップノードとの間の帯域幅予約は、フローとしてモデル化される。これらのフローによって使用される帯域幅の量は、MCF問題に対する変数である。一実施形態において、4種類のフローがある。

【0106】

10

20

30

40

50

2つの仮想ノードの間のフロー。x, y ∈ N<sup>V</sup>。リンク(i, j) ∈ E<sup>A</sup>において使用する帯域幅の量は、f<sup>x,y</sup>[i,j]によって示される。

【0107】

L<sup>1</sup>:冗長ノードa ∈ N<sup>K</sup>と仮想ノードy ∈ C<sup>K</sup>との間のフロー。冗長ノードaが何らかの仮想ノードxを引き継ぐ場合を除いて、これらのフローのいずれかにおける実際の帯域幅は、ゼロである。こうした回復が生じるときに、リンク(i, j) ∈ E<sup>A</sup>において使用される帯域幅の量は、

【数19】

$$f_{L^1}^{xy}[ij].$$

10

によって示される。これによって、冗長なフロー間の重複のモデル化が可能になる。

【0108】

冗長ノードN<sup>K</sup>と仮想ノードx ∈ C<sup>K</sup>との間のリンクにおける集計フロー。これは、リンク(i, j)上での重複後確保された帯域幅の実際の量を反映する。これは、f<sub>0</sub><sup>x</sup>[i,j]で示される。

【0109】

L<sup>2</sup>:2つの冗長ノードa, b ∈ N<sup>K</sup>間のフロー。リンク(i, j) ∈ E<sup>A</sup>において使用される帯域幅の量は、

【数20】

$$f_{L^2}^{ab}[ij].$$

20

によって示される。

【数21】

$$f_{L^1}$$

のフローとは異なり、これらのフローは、重複しない。これは、L<sup>1</sup>リンクに対する追加の信頼性を提供しない、仮想ノードxを介してノードaとbとをリンクするパス(a, x, b)を有することの取るに足らないソリューションを回避するためである。

【0110】

[0085]—実施形態において、物理ノードと、仮想又は冗長ノードとの間の双方向のマッピングは、二値変数p<sub>i,j</sub>, (i, j) ∈ R<sup>P</sup>によってモデル化される。リンク(i, j)及び(j, i)の間を流れているフローの総量が正である場合、p<sub>i,j</sub> = 1、そうでない場合、0である。したがって、MCFに対するソリューションがp<sub>x,u</sub> = 1を提供する場合、仮想ノードxは、物理ノードu上にホストされる。

30

【0111】

[0086]MCFの目的関数は、

【数22】

$$\min \sum_{w \in N^P} \alpha_w \sum_{x \in N^V} \rho_{xw} \mu_x + \sum_{(u,v) \in E^P} \beta_{uv} \times \left[ \sum_{x \in C^K} f_o^x[uv] + \sum_{(a,b) \in L^2} f_{L^2}^{ab}[uv] + \sum_{(x,y) \in E^V} f^{xy}[uv] \right] \quad (11)$$

40

として定義され、この場合、α<sub>w</sub>及びβ<sub>u,v</sub>はそれぞれ、ノード及びリンクの重みである。これによって、割り振られる計算及び帯域幅の加重和が最低限に抑えられる。負荷バランスを達成するために、重みはそれぞれ、

【数23】

$$\frac{1}{M_w + \delta}$$

及び

【数 2 4】

$$\frac{1}{H_{uv} + \delta}$$

と設定することができる。MCF に対する制約は、以下の通りである。

マッピング制約：

【数 2 5】

$$\sum_{u \in \Phi[x]} \rho_{ux} = 1, \quad \forall x \in N^V \cup N^K \quad (12)$$

10

【数 2 6】

$$\sum_{x \in N^V \cup N^K} \rho_{xu} \leq 1, \quad \forall u \in N^P \quad (13)$$

【0 1 1 2】

$$i_j H_{j i}, \quad (i, j) \in R^P \quad (14)$$

【0 1 1 3】

$$i_j = j_i, \quad (i, j) \in R^P \quad (15)$$

制約 (12) 及び (13) は、各仮想ノードが単一の物理ノードのみにマップされ、1つの物理ノードに1つ以下の仮想ノードしかマップできないようにすることを確実にする。

20

制約 (14) 及び (15) は、実行可能なフローがリンク (i, j) 上にマップされるとき、二値変数  $i_j$  を強制的に1にし、そうでなければ0にする。

(共有を条件とした) 計算容量の制約：

【0 1 1 4】

$$\rho_{iu} \mu_u \leq M_u, \quad u \in N^P, \quad i \in N^V \cup N^K \quad (16)$$

これは、マップされた仮想及び冗長ノードが物理ノード u 上の使用可能な容量  $M_u$  を超えないことを確実にする。冗長ノード、 $a \in N^K$  について、供給される最大能力は、

【数 2 7】

$$\max_{u \in C^V} \mu_u.$$

30

である。さらに、この容量が共有の冗長ノードのものを超える場合、バランスが供給されるだけで良い。

2つの仮想ノード間の帯域幅予約のためのフロー保護制約：

【数 2 8】

$$\sum_{u: (x,u) \in R^P} [f^{xy}[xu] - f^{xy}[ux]] = \eta_{xy}, \quad \forall (x,y) \in E^V \quad (17)$$

【数 2 9】

$$\sum_{u: (u,y) \in R^P} [f^{xy}[yu] - f^{xy}[uy]] = -\eta_{xy}, \quad \forall (x,y) \in E^V \quad (18)$$

40

【数 3 0】

$$\sum_{i \in N^A} [f^{xy}[ui] - f^{xy}[iu]] = 0, \quad \forall u \in N^P, \forall (x,y) \in E^V \quad (19)$$

【0 1 1 5】

[0087]制約 (17) 及び (18) は、仮想ノード x から仮想ノード y に発せられる仮想リンク (x, y) の総帯域幅  $\eta_{xy}$  を定義する。制約 (19) は、フローが中間の物理ノードで保護されることを確実にする。すなわち、ノード u から流れている総帯域幅がそのノード内に流れ込む総帯域幅に等しい。

50

$L^1$  リンク上の帯域幅を確保するためのフロー保護及び重複制約：

【数 3 1】

$$\sum_{u:(a,u) \in R^P} \left[ f_{L^1}^{acy}[au] - f_{L^1}^{acy}[au] \right] = \eta_{cy}, \quad \forall (a,y) \in L^1, \forall c \in C^V \quad (20)$$

【数 3 2】

$$\sum_{u:(u,y) \in R^P} \left[ f_{L^1}^{acy}[yu] - f_{L^1}^{acy}[uy] \right] = -\eta_{cy}, \quad \forall (a,y) \in L^1, \forall c \in C^V \quad (21)$$

【数 3 3】

$$\sum_{i \in N^A} \left[ f_{L^1}^{acy}[ui] - f_{L^1}^{acy}[iu] \right] = 0, \quad \forall u \in N^P, \forall (c,y) \in E^V, \forall a \in N^K \quad (22)$$

【数 3 4】

$$\sum_{a \in N^K, c \in F, (c,y) \in E^V} f_{L^1}^{acy}[ij] \leq f_o^y[ij], \quad \forall (i,j) \in E^A, \forall y \in C^K, \forall F \subseteq C^V, |F| \leq k \quad (23)$$

冗長ノード a がクリティカルノード x の代わりにする仮想ノード y への各フローについて、制約 (20) ~ (22) は、(17) ~ (18) における仮想フローのものと似たフロー保護モデルを定義する。制約 (23) は、すべての a にわたって合計される代わりに、重複し得る冗長なフローを処理する。1つの冗長ノード a のみが、いつでもクリティカルノード c の代わりとなり得る。次いで、フロー

【数 3 5】

$$f_{L^1}^{acy}[ij]$$

及び

【数 3 6】

$$f_{L^1}^{bcy}[ij]$$

は、重複する可能性があり、すなわち、リンク (i, j) について、

【数 3 7】

$$\max_{a \in N^K} f_{L^1}^{acy}[ij] \leq f_o^y[ij]$$

である。しかし、重複は、フロー

【数 3 8】

$$f_{L^1}^{acy}[ij]$$

及び

【数 3 9】

$$f_{L^1}^{bdy}[ij]$$

では起こらない可能性があり、この場合、冗長ノード a は、クリティカルノード c の代わりとなり、別のノード b は、クリティカルノード d の代わりとなる。これは、最大 k 個の置き換えまで生じる。制約 (23) は、これらの関係を捕捉する。

(共有及び  $C^V$  を条件とした)  $L^2$  リンク上の帯域幅予約のためのフロー保護制約：

【数 4 0】

$$\sum_{i \in N^A} \left[ f_{L^2}^{ab}[ui] - f_{L^2}^{ab}[iu] \right] = 0, \quad \forall u \in N^P, \forall (a,b) \in L^2 \quad (24)$$

10

20

30

40

【数 4 1】

$$\sum_{u:(a,u) \in R^P} \left[ f_{L^2}^{ab}[au] - f_{L^2}^{ab}[ua] \right] = \max_{x,y \in C^V} \eta_{xy}, \quad \forall (a,b) \in L^2 \quad (25)$$

【数 4 2】

$$\sum \left[ f_{L^2}^{ab}[bu] - f_{L^2}^{ab}[ub] \right] = - \max_{x,y \in C^V} \eta_{xy}, \quad \forall (a,b) \in L^2 \quad (26)$$

2つの冗長ノード a と b との間のフロー保護制約は、(17) ~ (19) における仮想フローのものと変わらない。供給される帯域幅は、 $C^V$  のノードを相互接続する仮想リンクの最大数である。しかし、これらの制約は、2つのケースにおいてのみ必要である。

10

【0 1 1 6】

1)

【数 4 3】

$$L^2 \neq \emptyset.$$

定理 2 から、クリティカルノードを相互接続している仮想リンクがある場合、 $L^2$  についての帯域幅を供給するだけで良い。

【0 1 1 7】

2) 共有。帯域幅は、共有されることになっている  $V I$  において、すでに供給されている。帯域幅の供給が十分でない場合、これらの制約は、バランスを供給するために存在する。

20

物理リンクにおけるリンク容量制約：

【数 4 4】

$$\sum_{x \in C^K} [f_o^x[uv] + f_o^x[vu]] + \sum_{(a,b) \in L^2} [f_{L^2}^{ab}[uv] + f_{L^2}^{ab}[vu]] + \sum_{(x,y) \in E^V} [f^{xy}[uv] + f^{xy}[vu]] \leq H_{uv}, \quad \forall (u,v) \in E^P \quad (27)$$

制約 (27) は、両方向の物理リンク (u、v) 上のすべてのフローを考慮に入れる。これは、物理的に残っている帯域幅  $H_{uv}$  未満でなければならない。

30

拡張マッピングリンクにおけるリンク容量制約：

【数 4 5】

$$\sum_{x \in C^K} [f_o^x[ij] + f_o^x[ji]] + \sum_{(a,b) \in L^2} [f_{L^2}^{ab}[ij] + f_{L^2}^{ab}[ji]] + \sum_{(x,y) \in E^V} [f^{xy}[ij] + f^{xy}[ji]] \leq H_{ij} \rho_{ij}, \quad \forall (i,j) \in R^P \quad (28)$$

厳密に言えば、帯域幅  $H_{ij}$  が無限であるため、マッピングリンクにおける制約はないはずである。しかし、この制約は、マッピング制約 (14) 及び (15) と連動して、いずれかの方向のそのリンクにおける任意の正のフローがある場合、マッピング二値変数  $\rho_{ij}$  を強制的に 1 にし、そうでない場合 0 にする。

40

領域制約：

$$f^{xy}[ij] = 0, \quad i, j \in N^A, \quad (x, y) \in E^V \quad (29)$$

【数 4 6】

$$f^{xy}[ij] \geq 0, \quad \forall i, j \in N^A, \forall (x, y) \in E^V \quad (29)$$

$$f_{L^1}^{ax}[ij] \geq 0 \quad \forall i, j \in N^A, \forall (a, x) \in L^1 \quad (30)$$

$$f_{L^2}^{ab}[ij] \geq 0, \quad \forall i, j \in N^A, \forall (a, b) \in L^2 \quad (31)$$

$$f_o^x[ij] \geq 0, \quad \forall i, j \in N^A, \forall x \in C^K \quad (32)$$

$$\rho_{ij} \in \{0,1\}, \quad \forall (i, j) \in R^P \quad (33)$$

10

【 0 1 1 8 】

$$f_o^x[ij] \geq 0, \quad i, j \in N^A, \quad x \in C^K \quad (32)$$

【 0 1 1 9 】

$$\rho_{ij} \in \{0, 1\}, \quad (i, j) \in R^P \quad (33)$$

これらは、この変更された MCF 問題のすべての変数における領域制約であり、すべてのフローは、非ゼロでなければならず、マッピング変数はバイナリである。

【 0 1 2 0 】

[0088] インフラストラクチャが急速に仮想化されるにつれて、信頼性保証を仮想化インフラストラクチャに提供する必要性が増加している。上記には、仮想化層自体における信頼性保証について記載されている。冗長ノードは、ネットワークにわたって分散される仮想ノードとすることができる。このために、供給された帯域幅を完備した  $n:k$  の冗長アーキテクチャ、及び物理ネットワークにわたって仮想化されたネットワークを割り振る方法が提案される。冗長ノード及びリンクによって使用されるリソースを保護するために、これらの冗長ノードを VI にわたって共有することができ、供給の間、それらの帯域幅が重複する。

20

#### コンピュータシステムの一例

【 0 1 2 1 】

[0089] 図 11 は、本明細書において記載されている 1 つ又は複数の動作を実行することができるコンピュータシステム例のブロック図である。図 11 を参照すると、コンピュータシステム 1100 は、典型的なクライアント又はサーバコンピュータシステムを備えることができる。コンピュータシステム 1100 は、情報を伝えるための通信機構又はバス 1111、及びバス 1111 に結合されて情報を処理するためのプロセッサ 1112 を備える。プロセッサ 1112 は、それだけには限定されないが、ペンティアム (Pentium) (商標)、パワー PC (Power PC) (商標)、アルファ (Alpha) (商標) など、マイクロプロセッサを含む。

30

【 0 1 2 2 】

[0090] システム 1100 は、ランダムアクセスメモリ (RAM)、又はバス 1111 に結合されてプロセッサ 1112 によって実行される命令及び情報を格納するための他の動的な記憶デバイス 1104 (メインメモリと呼ばれる) をさらに備える。プロセッサ 1112 による命令の実行中に一時的数値変数又は他の中間情報を格納するために、メインメモリ 1104 を使用することもできる。

40

【 0 1 2 3 】

[0091] コンピュータシステム 1100 は、読み取り専用メモリ (ROM) 及び / 又はバス 1111 に結合されてプロセッサ 1112 のための静的情報及び命令を格納するための他の静的記憶デバイス 1106、並びに磁気ディスク又は光ディスクなどのデータ記憶デバイス 1107 及びその対応するディスクドライブも備える。データ記憶デバイス 1107 は、バス 1111 に結合されて情報及び命令を格納するためのものである。

【 0 1 2 4 】

[0092] コンピュータシステム 1100 は、バス 1111 に結合されて、コンピュータユ

50



ーザに対して情報を表示するために、例えばブラウン管（CRT）又は液晶式ディスプレイ（LCD）などのディスプレイデバイス1121にさらに結合することができる。英数字及び他のキーを含む英数字入力デバイス1122も、バス1111に結合されてプロセッサ1112に情報及びコマンドの選択を伝えるようにすることができる。追加のユーザ入力デバイスは、バス1111に結合されて、プロセッサ1112に指示情報及びコマンドの選択を伝える、及びディスプレイ1121上のカーソルの動きを制御するための、例えばマウス、トラックボール、トラックパッド、スタイラス、又はカーソル方向キーなどのカーソル制御1123である。

【0125】

[0093]バス1111に結合することができる他のデバイスは、ハードコピーデバイス1124であり、このデバイスを使用して、例えば紙、フィルム、又は類似のタイプの媒体に情報を表すことができる。バス1111に結合することができる他のデバイスは、電話又はハンドヘルドパームデバイスに対して通信を行うための有線/無線通信機能1125である。

10

【0126】

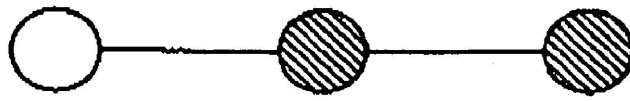
[0094]システム1100及び付随するハードウェアのコンポーネントの任意のもの又は全部を本発明において使うことができることに留意されたい。しかし、コンピュータシステムの他の構成がデバイスの一部又は全部を含むことができることを理解されよう。

【0127】

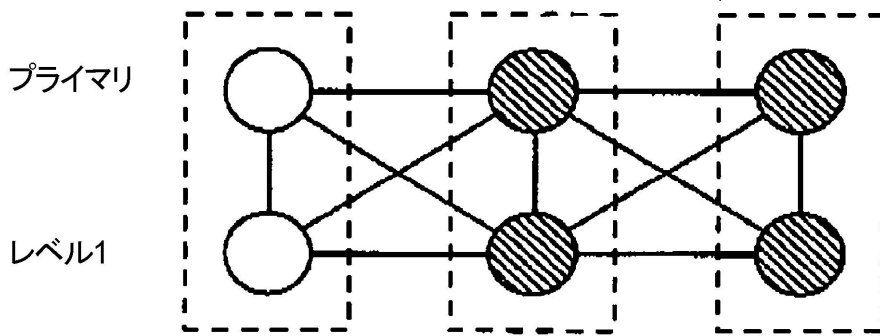
[0095]本発明の多くの変更及び修正は、間違いなく前述の説明を読んだ後に当業者にとって明らかになるが、例として示され、記載されるいかなる特定の実施形態も決して制限と見なされないことを理解されたい。したがって、様々な実施形態の詳細への言及は、本発明にとって重要であると見なされる特徴だけを詳述する請求項の範囲を制限することを目的としない。

20

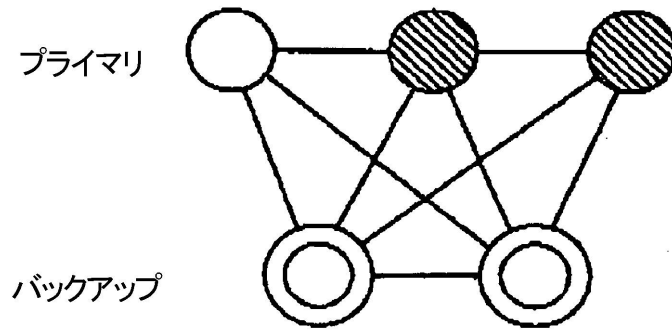
【図1】



(a) ベース仮想ネットワーク



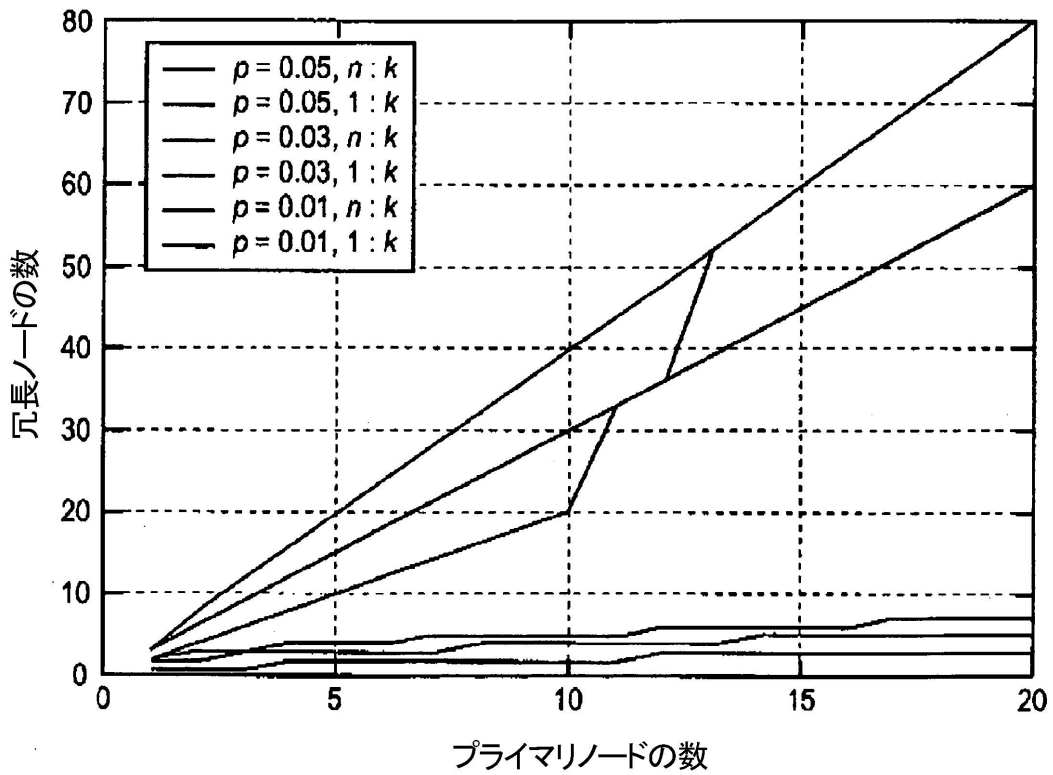
(b) シンプルな1:kの複製



(c) n:kの複製

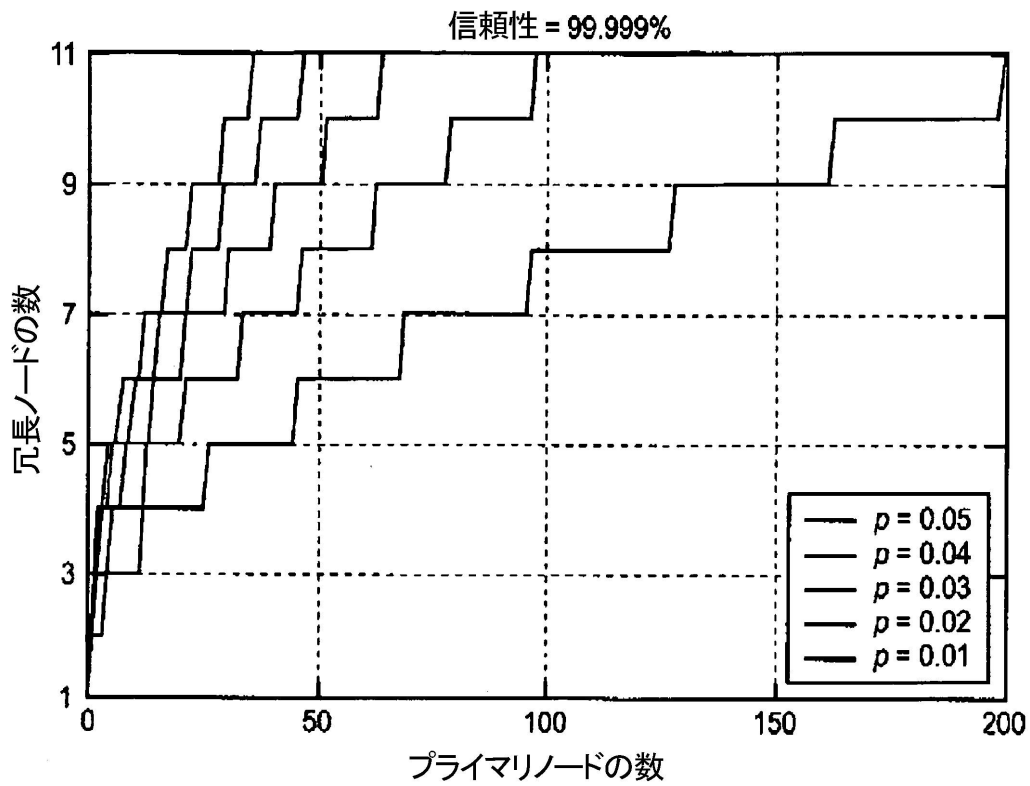
フォールトトレラントアーキテクチャの比較

【 図 2 】



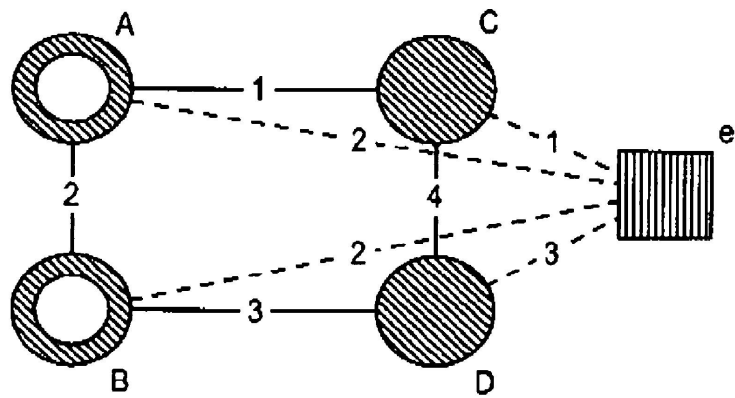
99.999%のノード信頼性のために必要な冗長ノードの数。  
 $\rho$ は、各ノードの個別の故障の確率である。

【図3】

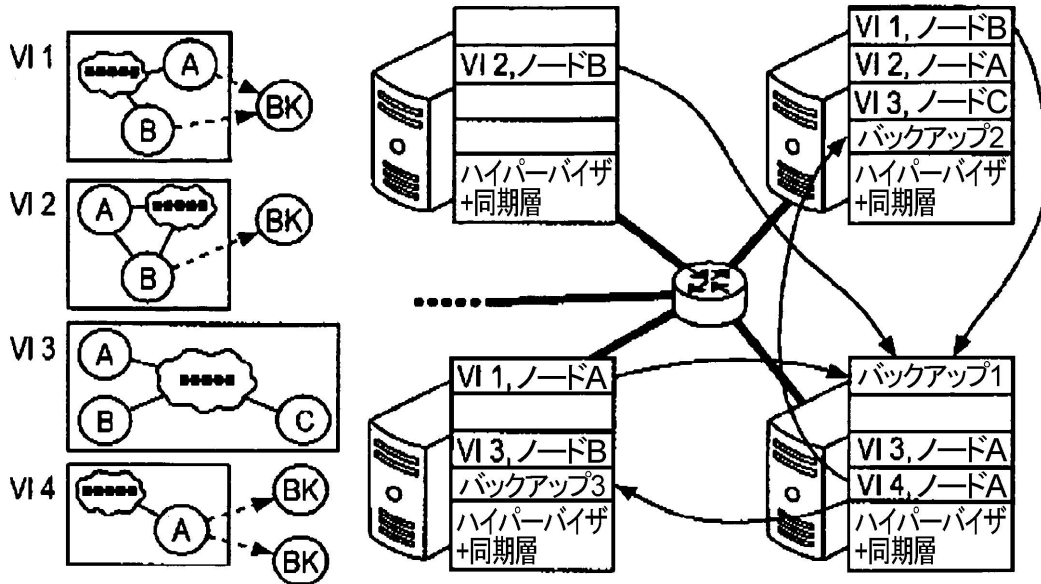


n:kの複製がサポートできるノードの数

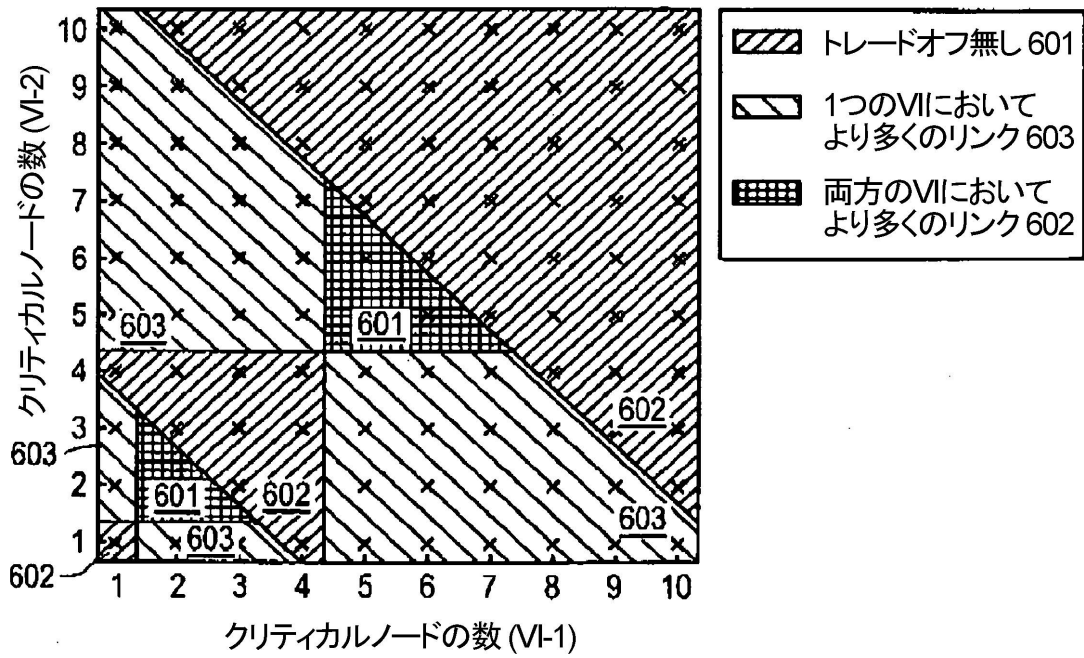
【 図 4 】



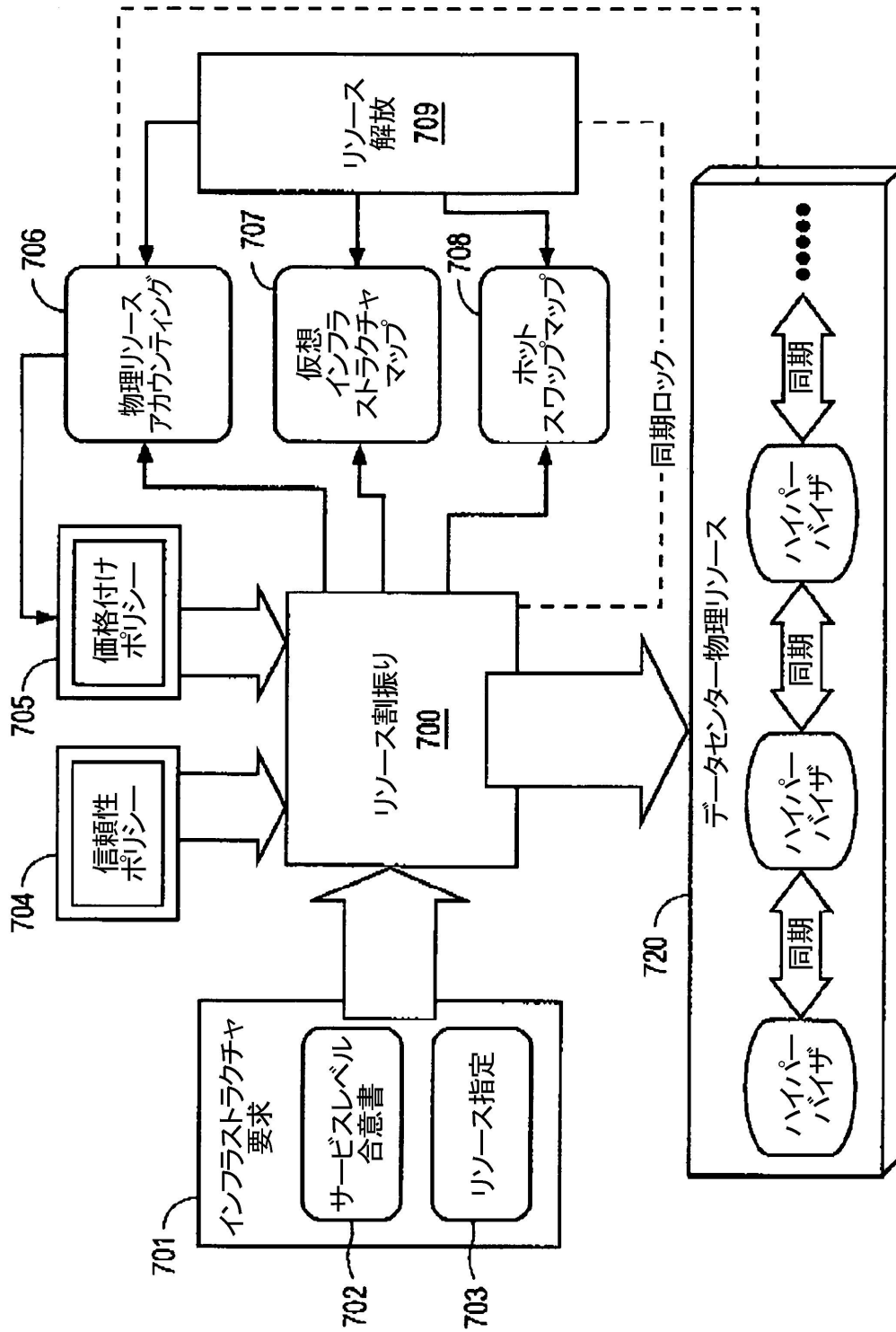
【 図 5 】



【図6】

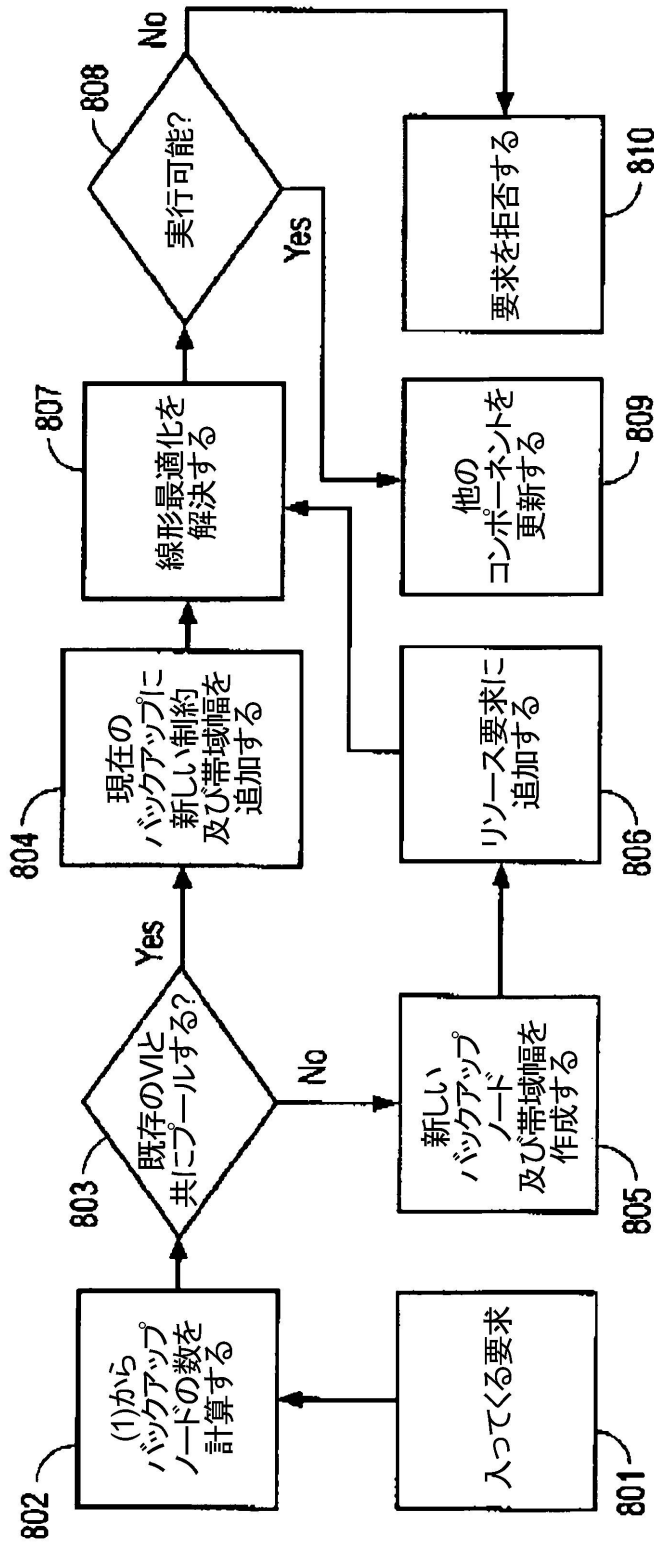


【 図 7 】

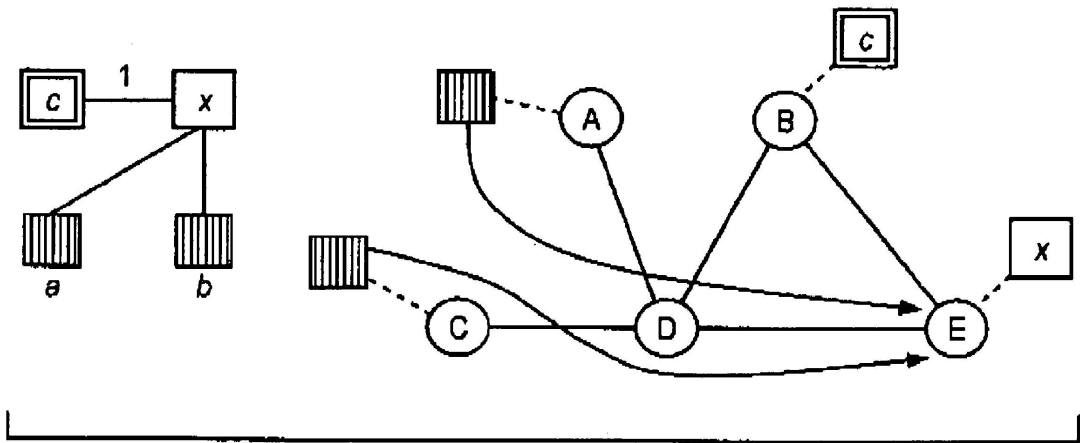




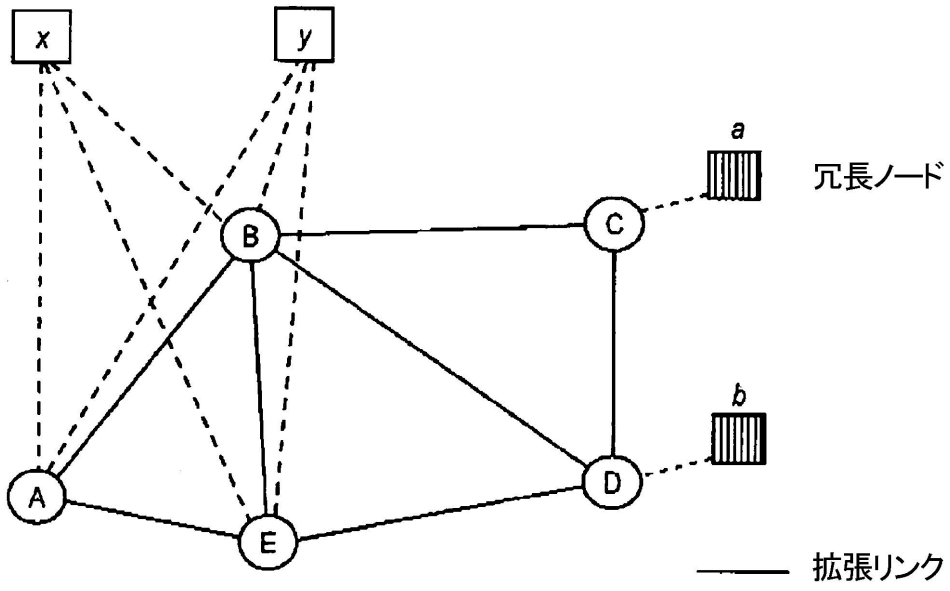
【 図 8 】



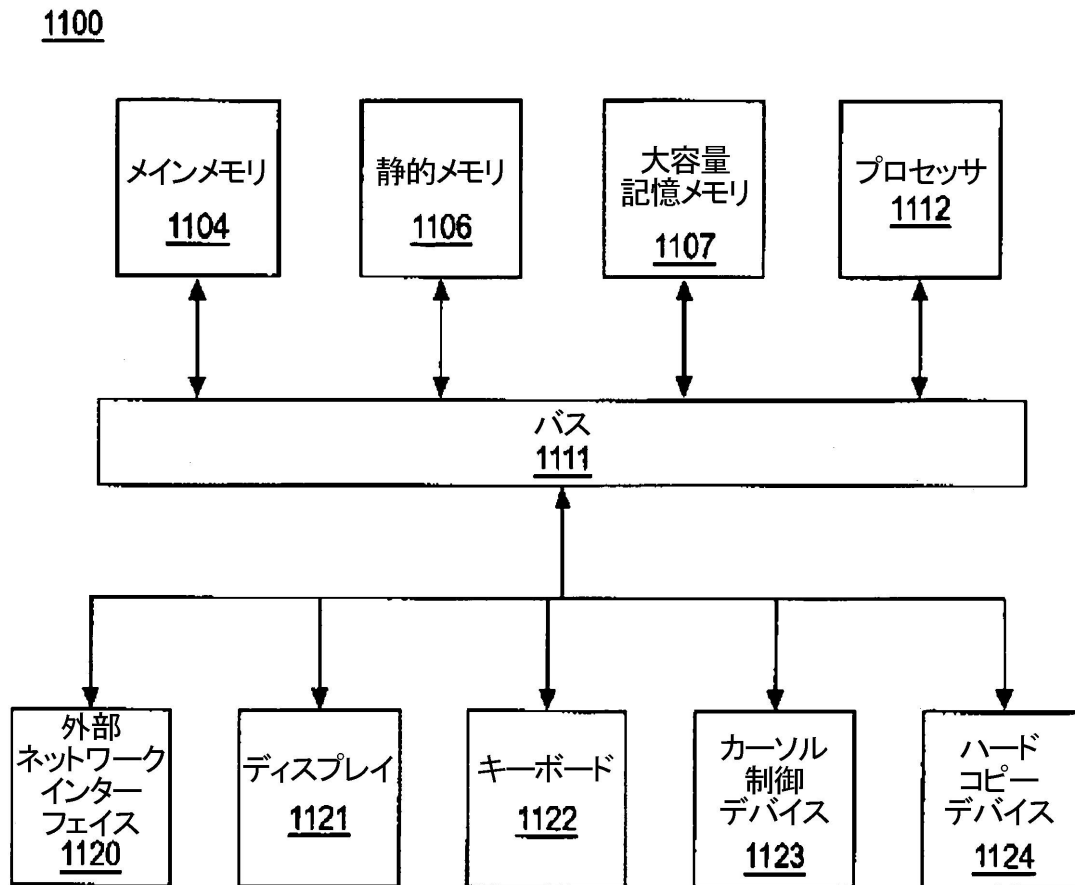
【 図 9 】



【図10】



【図11】



## フロントページの続き

- (72)発明者 ヤオ, ワイーレオン  
アメリカ合衆国, カリフォルニア州, マウンテン ビュー, アpartment 302, カ  
リフォルニア アヴェニュー 2685
- (72)発明者 ウェストファル, セドリック  
アメリカ合衆国, カリフォルニア州, サン フランシスコ, ゲレーロ ストリート 101  
9
- (72)発明者 コザット, ウラス  
アメリカ合衆国, カリフォルニア州, サンタ クララ, フローラ ヴィスタ アヴェニュー  
ナンバー349 3612

審査官 三坂 敏夫

- (56)参考文献 特開平06-028330(JP,A)  
国際公開第2009/081736(WO,A1)  
米国特許出願公開第2009/0138752(US,A1)  
Song FU, "Failure-Aware Construction and Reconfiguration of Distributed Virtual Machin  
es for High Availability Computing", CCGRID '09. 9th IEEE/ACM International Symposium  
on Cluster Computing and the Grid, 2009., 米国, IEEE, 2009年 5月21日, pages:37  
2-379

- (58)調査した分野(Int.Cl., DB名)  
G06F 11/20