

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6161312号  
(P6161312)

(45) 発行日 平成29年7月12日(2017.7.12)

(24) 登録日 平成29年6月23日(2017.6.23)

(51) Int. Cl. F I  
**G06F 9/48 (2006.01)** G06F 9/46 452F  
**G06F 9/50 (2006.01)** G06F 9/46 465D

請求項の数 10 (全 19 頁)

<p>(21) 出願番号 特願2013-23964 (P2013-23964)                  (22) 出願日 平成25年2月12日(2013.2.12)                  (65) 公開番号 特開2013-175177 (P2013-175177A)                  (43) 公開日 平成25年9月5日(2013.9.5)                  審査請求日 平成28年2月10日(2016.2.10)                  (31) 優先権主張番号 61/603,366                  (32) 優先日 平成24年2月26日(2012.2.26)                  (33) 優先権主張国 米国 (US)                  (31) 優先権主張番号 13/688,035                  (32) 優先日 平成24年11月28日(2012.11.28)                  (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 502096543                  パロ・アルト・リサーチ・センター・イン                  コーポレーテッド                  Palo Alto Research                  Center Incorporated                  アメリカ合衆国、カリフォルニア州 94                  304、パロ・アルト、コヨーテ・ヒル・                  ロード 3333                  (74) 代理人 110001210                  特許業務法人YKI国際特許事務所                  (72) 発明者 ダニエル・エイチ・グリーン                  アメリカ合衆国 カリフォルニア州 94                  087 サニーヴェイル マネ・ドライブ                  1055 ナンバー6</p>
---	--

最終頁に続く

(54) 【発明の名称】 データセンタにおけるQoS認識バランシング

(57) 【特許請求の範囲】

【請求項1】

2つの物理資源間の負荷をバランシングするためのコンピュータ実行可能な方法であって、

物理資源ごとに、

前記物理資源を共有する複数のジョブに対して資源使用モデルを確立するステップと、

前記ジョブに関連付けられたサービス品質(QoS)格付けを識別するステップと、

前記資源使用モデルと前記QoS格付けに基づいて前記物理資源に対するQoSバランスポイントを計算するステップであって、前記QoSバランスポイントは、それを上回ると実質的な資源のプロビジョニングが必要とされる、QoSの格付けを表している、ステップと、

第1の物理資源に関連付けられた第1のQoSバランスポイントと第2の物理資源に関連付けられた第2のQoSバランスポイントとの差を検出することに対応して、前記第1と前記第2のQoSバランスポイントの差が削減されるように前記第1と前記第2の物理資源間で移動される1つ以上のジョブを識別するステップと、

を含む、方法。

【請求項2】

コンピュータによって実行されると、前記コンピュータに2つの物理資源間の負荷をバランシングするための方法を実行させる命令を記憶するコンピュータ読出可能記憶媒体であって、

10

20

前記方法は、  
物理資源ごとに、

前記物理資源を共有する複数のジョブに対して資源使用モデルを確立するステップと、  
前記ジョブに関連付けられたサービス品質（QoS）格付けを識別するステップと、  
前記資源使用モデルと前記QoS格付けに基づいて前記物理資源に対するQoSバランスポイントを計算するステップであって、前記QoSバランスポイントは、それを上回ると実質的な資源のプロビジョニングが必要とされる、QoSの格付けを表している、ステップと、

第1の物理資源に関連付けられた第1のQoSバランスポイントと第2の物理資源に関連付けられた第2のQoSバランスポイントとの差を検出することに対応して、前記第1と前記第2のQoSバランスポイントの差が削減されるように前記第1と前記第2の物理資源間で移動される1つ以上のジョブを識別するステップと、

を含む、

ことよくなるコンピュータ読出可能記憶媒体。

#### 【請求項3】

2つの物理資源間で負荷をバランスングするための計算システムであって、  
前記2つの物理資源上で実行されるジョブごとに資源使用モデルを構築するように構成されている資源使用モデル構築手段と、

前記ジョブに関連付けられたQoS格付けを識別するように構成されているサービス品質（QoS）識別手段と、

資源使用モデルと前記物理資源上で実行しているジョブに関連付けられたQoS格付けに基づいて物理資源ごとにQoSバランスポイントを計算するように構成されているQoSバランスポイント計算手段であって、前記QoSバランスポイントが、それを上回ると実質的な資源のプロビジョニングが必要とされる、QoSの格付けを表しているQoSバランスポイント計算手段と、

第1の物理資源に関連付けられた第1のQoSバランスポイントと第2の物理資源に関連付けられた第2のQoSバランスポイントとの差を検出することに対応して、前記第1と前記第2のQoSバランスポイントの差が削減されるように前記第1と前記第2の物理資源間で移動される1つ以上のジョブを識別するように構成されているジョブ移行識別手段と、

を含む計算システム。

#### 【請求項4】

前記2つの物理資源に対する前記QoSバランスポイントを管理者へ提示し、  
前記管理者から、前記第1と前記第2の物理資源の間で移動する前記1つ以上のジョブを識別する入力を受信する、

ように構成されているユーザインターフェースを更に含む、請求項3に記載のシステム

#### 【請求項5】

移動される1つ以上のジョブを識別している間、前記ジョブ移行識別手段がQoSバランスを改良することができる移動の限界便益を計算するように更に構成されており、前記システムが前記計算に基づいてジョブを選択して移動を終了するように構成されているジョブ移行メカニズムを更に含む、請求項3に記載のシステム。

#### 【請求項6】

前記1つ以上のジョブの移動によって影響されるリスクの相関関係を計算するように構成されているリスク相関計算手段と、

前記リスクの相関関係が所定のしきい値を下回ることに対応して、前記識別された1つ以上のジョブの移動を終了するように構成されているジョブ移行コントローラと、

を更に含む請求項3に記載のシステム。

#### 【請求項7】

前記資源使用モデルは時変性であり、前記資源使用モデル構築手段は前記ジョブに関連

10

20

30

40

50

付けられた資源使用履歴に基づいて前記時変性の資源使用モデルを構築する、請求項 3 に記載のシステム。

【請求項 8】

2つのジョブ間の相補性レベルを測定するように構成されている相補性測定メカニズムと、

相補性レベルが所定のしきい値を上回ることに応答して、前記2つのジョブを同じ物理資源上に配置するように構成されているジョブ配置メカニズムと、  
を更に含む請求項 3 に記載のシステム。

【請求項 9】

前記相補性レベルを測定している間、前記相補性測定メカニズムは前記2つのジョブに関連付けられたピーク時の資源使用間で時間的相関関係を計算する、請求項 8 に記載のシステム。

10

【請求項 10】

1つ以上の他の互いに関連しているジョブのロケーションと、配置される前記ジョブに関連付けられるネットワーク接続要件と、前記配置されるジョブに関連付けられるセキュリティ要件の1つ以上に基づいて、特定の物理資源へジョブを配置するように構成されているジョブ配置メカニズムを更に含む、請求項 3 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

20

本明細書は、一般に、データ・センタ・オペレーションに関し、より詳細には、データ・センタ・オペレーションのサービス品質(QoS)とエネルギー効率を改良するために様々な物理資源の間でジョブをbalancingするシステムに関する。

【背景技術】

【0002】

近年の仮想化技術はデータセンタが共有環境において様々なジョブを実行することを可能にした。言い換えれば、これらの様々なジョブは、全てが単一のマシンまたは複数のマシンのクラスタによって提供することができるメモリ、中央処理装置(CPU)、および帯域幅などの同じ物理資源を共有することができる。データ・センタ・オペレーションの重要な考慮する点は様々なマシンまたは複数のマシンのクラスタの間でジョブ(または負荷)をbalancingすることである。

30

【発明の概要】

【発明が解決しようとする課題】

【0003】

データ・センタ・オペレーションにおけるジョブのbalancingは、良好に統合するジョブのグループ、即ち、マシンまたはマシンのクラスタで物理資源をより有効に活用することができるグループを検索することを含む。従来のアプローチは、各ジョブに関連付けられたQoS要件を無視する場合が多い。

【課題を解決するための手段】

【0004】

40

本発明の一実施形態は、2つの物理資源間で負荷をbalancingさせるシステムを提供する。動作中、物理資源ごとに、システムは、物理資源を共有する複数のジョブに対して資源使用モデルを確立し、ジョブに関連付けられたサービス品質(QoS)レベルを識別し、資源使用モデルとQoS格付けに基づいて物理資源に対するQoSバランスポイントを計算する。このQoSバランスポイントは、それを上回ると実質的な資源プロビジョニングが必要とされる、QoSの格付けを表している。第1の物理資源に関連付けられる第1のQoSバランスポイントと第2の物理資源に関連付けられる第2のQoSバランスポイントと差を検出することに応答して、システムは、第1と第2のQoSバランスポイントの差が削減されるように第1と第2の物理資源の間で移動される1つ以上のジョブを識別する。

50

## 【 0 0 0 5 】

本実施形態の一変形例では、移動される1つ以上のジョブを識別することは人間のオペレータによって実行される。

## 【 0 0 0 6 】

本実施形態の一変形例では、移動される1つ以上のジョブを識別することは、QoSバランスを改良することができる移動の限界便益を計算することを伴う。さらに、システムはこの計算に基づいてジョブを選択しその移動を終了させる。

## 【 0 0 0 7 】

本実施形態の一変形例では、システムは1つ以上のジョブの移動に影響されるリスクの相関関係を判定する。リスクの相関関係が所定のしきい値を下回ることに応答して、システムはこれらの移動を終了させる。

10

## 【 0 0 0 8 】

本実施形態の一変形例では、資源使用モデルは時変性であり、この時変性の資源使用モデルはこれらのジョブに関連付けられる資源使用履歴に基づいて確立される。

## 【 0 0 0 9 】

本実施形態の一変形例では、システムは2つのジョブ間の相補性レベルを測定する。所定の閾値を上回る相補性レベルに応答して、システムは同じ物理資源へ2つのジョブを配置する。

## 【 0 0 1 0 】

更なる変形例によれば、相補性レベルを測定することは、これら2つのジョブに関連付けられるピーク時の資源使用間の時間的相関関係を計算することを含む。

20

## 【 0 0 1 1 】

本実施形態の一変形例によれば、システムは、1つ以上の他の関連ジョブのロケーション、このジョブに関連付けられるネットワーク接続性要件、およびこのジョブに関連付けられるセキュリティ要件の1つ以上に基づいて、特定の物理資源へ配置されるジョブを識別する。

## 【 図面の簡単な説明 】

## 【 0 0 1 2 】

【図1】資源の関数としての確率密度関数(PDF)と例示的なジョブに対して予約された資源とを示す図である。

30

【図2】本発明の一実施形態による、2つのマシンが「バランスが崩れている」状況を示す図である。

【図3】本発明の一実施形態による、データセンタのオペレーションの負荷配置コントローラを示す図である。

【図4】本発明の一実施形態による、例示的なジョブ移行プロセスを示す流れ図である。

【図5】本発明の一実施形態による、データセンタ内のマシンの間でジョブを配置するための例示的なコンピュータシステムを示す図である。

## 【 発明を実施するための形態 】

## 【 0 0 1 3 】

以下の記載は、当業者がこれらの実施形態を活用することを可能にするために提示され、特定のアプリケーションとその要件に関して提供されている。開示されている実施形態に対する様々な変更は当業者に容易に理解され、本明細書中に定義されている一般的な原理が本明細書の精神及び範囲を逸脱しない限りに於いて他の実施形態やアプリケーションに適用可能である。よって、本発明は図示されている実施形態に限定されないが、本明細書中に開示された原理および特性を逸脱しない最も広い範囲に一致する。

40

## 【 0 0 1 4 】

本発明の実施形態は、データセンタの物理資源の間で負荷をバランスングするためのシステムを提供する。より詳細には、システムは、ジョブのQoS要件のみならずジョブの相補性とリスクの相関関係などの他の基準に基づいてグループとしてのジョブを様々な資源へ割り当てる。動作中、システムは、ジョブに必要とされるQoSを識別し、少なくとも

50

も履歴データに基づいてジョブに対して資源使用モデルを確立する。次いで、システムは、QoS格付けと現在資源に割当てられているジョブの資源使用モデルとに基づいて、資源ごとにQoSバランスポイントを計算する。QoSのバランスポイントに基づいて、自動コントローラや人間の管理者は、バランスポイントを均等化するために資源の間でジョブを移動させて、資源のグループの活用を向上させることができる。

【0015】

本明細書において、用語「物理資源」は、計算ジョブを完了するために必要とされる様々なタイプの物理的機器をさす。この用語は、処理能力、記憶領域、通信帯域幅、入/出力などを含み得る。さらに、特定の「物理資源」は、データセンターの単一マシン、マシンのクラスタ、またはすべてのマシンに言及することができる。また、用語「物理資源」と「物理的マシン」は互いに互換性がある。

10

【0016】

本明細書において、用語「ジョブ」は、共有環境における計算可能なタスクをさす。より詳細には、ジョブは、仮想マシンインスタンスまたは複数の仮想マシンインスタンスの集合であってよい。

【0017】

統計パッキング

オペレーションの負荷配置にはいくつかの重要な考慮する点がある。例えば、良好に統合するジョブのグループ、即ち、物理資源をより有効に活用することができるジョブのグループを検出するために、ジョブの負荷パターンが相補的であるか、いくつかのジョブをグルーピングすることがリスクを削減するか、およびジョブのグループがQoS要件の混合を含むかについて考える必要がある。これらの考慮する点の各々は、ジョブのグループが特定の物理資源上で一緒に実行する時の振舞い方の判断にかなり影響を与える。さらに、実際の状況において、関連する特定のジョブによっては、これらの考慮する点の1つ以上を優先することもある。

20

【0018】

仮想化環境においてジョブは物理資源を共有することができる。通常、ジョブは予約を行うことにより資源を取得し、これによって、少なくともそれぞれ自体の物理マシンを有すると同じ程度に有効に動作することを保証する。但し、予約された資源が現在ジョブによって必要とされない場合、これらの資源は他のジョブに共有される。「統計パッキング」はジョブのグループに対してまとめて予約を決定するメカニズムをいう。各ジョブにそれぞれ自体の予約をさせるというよりむしろ、このグループアプローチは仮想化環境において実行可能な更なる共有をより良好に利用できるようにする。

30

【0019】

例えば、2つだけのジョブから成るシンプルな「グループ」を考えてみると、両ジョブはそれらの資源ニーズについて実質的な不確実性を有しているが、1つジョブは上位のQoS要件を有し、1つのジョブは下位のQoS要件を有している。個別予約を行うということは、各ジョブがその不確実性に対処するために更なる資源の「バッファ」をもたなくてはならないことを意味する。しかしながら、われわれがまとめて予約決定する場合、上位のQoSジョブに必要な大規模な予約を下位のQoSジョブにも頻繁に利用できるという可能性がある。よって、下位のQoSジョブに対する資源予約をほとんど必要としない。

40

【0020】

統計パッキングをより良く定式化するために、われわれは最初にQoSレベルとジョブのQoSレベルをその資源予約にどのように連携させるかを定義しておく必要がある。用語「QoS」の汎用的意味はジョブがその目標をどの程度満たせるかをいう。いくつかのアプリケーションにおいて目標は平均的な遅延または完了したトランザクションにおいて測定される。残念ながら、異なるアプリケーションはそれらにとっては重要である異なるメトリックを有する。これらの異なるメトリックを汎用な設定に適用させることは、不十分なメトリックの影響がアプリケーション同士の間で変化することから、困難であり、ま

50

た、メトリックが標準以下の場合、その故障が不十分に書かれたアプリケーションや仮想環境における不十分な資源割当てに起因する場合、不明瞭になり得る。

【 0 0 2 1 】

多種多様なアプリケーションを介してQoSの処理を簡潔化し統合するためのゴールドスタンダードとみなされるようなより汎用なアプローチは、必要とする全ての資源を有するジョブを完全にプロビジョニングすることである。よって、アプリケーション関連メトリックに良好なパフォーマンスを分配することがジョブの役割である。完全なプロビジョニングは、通常はジョブが資源をいくつ必要とするかを明確に予測することが不可能であるので、通常は実際に必要とされる資源よりも多めに資源を予約することを要求する。たとえ余分に予約をしても不足のリスクが消えるわけではない。本明細書において、われわれは完全なプロビジョニングの至適標準が満たされないときに与えるチャンスとしてジョブによって分配されるQoSのレベルを以下に説明する。即ち、QoSレベルpを有するジョブは、測定されるインターバルの間、要求される確率pを有する資源を全て受信しない。QoSレベルpはジョブの失敗の確率と混同すべきでないことに留意されたい。最近のアプリケーションの大部分は要求される資源を100%保持していなくても十分に実行する程度に良く書かれている。例えば、重い負荷を受けているアプリケーションでも低解像度の画像を一時的に分配したり、バッチモードで実行するアプリケーションでもデッドラインまでの終了を確実にするために早めにスタートを切ったりする。しかしながら、資源不足へのレジリエンスに関してはやはりここでもアプリケーション固有の特性にすぎない。資源の汎用な管理を簡潔化するためにわれわれは至適標準が満たされない最悪のケースを想定している。混同を避けるために、完全にプロビジョニングするためにこの失敗を「不足」の状況と呼ぶ。QoSレベルpの適切な設定値（即ち、許容された「不足の確率」）は、ジョブに対する失敗の許容された確率よりも高いことに留意されたい。

【 0 0 2 2 】

負荷配置においてQoS仕様（またはQoSレベル）を使用することの重要な利点はより優れたリスク管理である。許容される「不足の確率」を特定することによって、システムは冒されるリスクをより慎重に管理することで、システムによって管理される重要かつ上位のQoSジョブのリスクが削減される。

【 0 0 2 3 】

経時的にジョブの資源の使用を監視するために、 $x_t$ で表される、直接測定された資源使用量と、 $r_t$ で表される、QoS要件を満たすために必要な資源の量と、の2つの測定値を使用することができる。仮想化される環境において、 $r_t$ は、通常は実際の使用量 $x_t$ を上回る必要な資源予約量であり。 $r_t$ が、直接的な（実際に使用された）使用量と、間接的な（例えば、予約されたバッファと可能な休止中の）資源使用量と、の両方を含むことに留意することが重要である。一般に、 $x_t$ を測定することは可能であるが、QoSが保護された使用量 $r_t$ は算定しにくい。監視ツールが、ジョブの確率論的モデル、つまり、今後起こり得る資源ニーズに対する分布 $f_t(z)$ を提供できるモデルを有している場合、 $r_t$ は累積分布とQoS仕様から算定することができる。より詳細には、累積分布

【数1】

$$\Phi(y) = \int_{-\inf}^y \phi(z), \quad (1)$$

とQoSレベルp（不足の許容された確率）と置かれた場合、

【数 2】

$$r_t = \text{inverse } \Phi(1 - p) \quad (2)$$

を得ることができる。

図 1 は資源の関数として確率密度関数 (PDF) と例示的なジョブに対して予約された資源とを示す図である。図 1 によれば、予約された資源  $r_t$  が分布曲線の残存テールにおける累積確率が許容される不足の確率を下回るほど大きいことが分かる。監視の際、 $r_t$  がジョブに対して行われた実際の予約に必ずしも一致しないことに留意されない。多くの場合、実際の予約は手動で設定され、この予約は仮想化の前にジョブが一度有していた物理資源に一致することもある。仮想化の前にも不足のリスクはあったが、通常、QoS を真剣に考慮しなくても、これらの予約は、等式 (2) で計算されるように、 $r_t$  の QoS のパフォーマンスを満たす場合もあれば満たさない場合もある。但し、監視ツールは、 $r_t$  が導出される前に  $x_t$  を最初に測定してモデリングする必要がある。

10

【0024】

$r_t$  の式が与えられる初めてわれわれは統計パッキングのアルゴリズムを定式化することができる。各ジョブの資源ニーズが分布  $(z)$  と QoS レベル  $p$  によって記述されるジョブの集合に関して、各ジョブに対して個別に計算された予約は以下のように計算され、

20

【数 3】

$$r^{(i)} = \text{inverse } \Phi^{(i)}(1 - p^{(i)}) \quad (3)$$

全てのジョブに対する全体予約は以下ようになる。

【数 4】

$$r^T = \sum_{i=1}^n r^{(i)} \quad (4)$$

30

【0025】

統計パッキングは全体予約を削減することができる。一般性を失わずに、QoS の降順、従って  $p^{(i)}$  の昇順でジョブが仕分けされると想定される。独立したジョブに関して、資源ニーズの結合された分布  $^T(z)$  は個別の分布の畳み込み  $^{(i)}(z)$  であり、以下のように書かれる。

40

【数 5】

$$\phi^T(z) = \bigotimes_{i=1}^n \phi^{(i)}(z) \quad (5)$$

【0026】

k を、予約

【数6】

$$(\hat{r}^{(k-1)} = \sum_{i=1}^{k-1} r^{(i)})$$

で表される  $k$  における部分予約（「部分」は  $s^{(k)} < r^{(k)}$  で表され、ジョブ  $k$  に対して要求される個別予約を意味する）の部分和になるような最小の索引とする。部分和と部分予約は、QoS レベル  $p^{(k)}$  におけるグループ全体のニーズを満たす程度に大きい。即ち、

【数7】

$$\hat{r}^{(k-1)} + s^{(k)} = \text{inverse } \Phi^T (1 - p^{(k)}) \quad (6)$$

である。

次いで、ジョブ  $1, 2, \dots, k$  は予約  $r^{(1)}, r^{(2)}, \dots, r^{(k-1)}$  ,  $s^{(k)}$  を得ることができ、残存ジョブは全く予約されない。予約なしではあるが、残存ジョブは、最初の  $k$  個のジョブの未使用の予約を利用することにより、それらの QoS 要件を問題なく満たすことができる。等式 (6) に示したように、統計パッキングに対する全体予約は、等式 (4) に示された全体予約から有効に削減されることは明らかである。統計パッキングアルゴリズムに可能ないくつかのばらつきがあることが理解されよう。例えば、Daniel H. Greene、Maurice Chu、Haitham Hindi、Bryan T. Preas、Nitin Parekh による「Statistical Packing of Resource Requirements in Data Centers」と題された米国特許出願公開 2010/0100877A1 と、Daniel H. Greene、Lara Crawford、Maurice Chu、John Handley「Long Term Resource Provisioning with Cascading Allocations」と題された本出願と同時に出願されている米国特許出願を参照されたい。特に、等式 (5) におけるようにジョブは必ずしも独立している必要はなく、むしろ、相関関係にあるジョブの結合した分布を学習することができる。この特許のバランシング技術は多種多様な統計パッキング方法のいずれかに基づいて行うことができる。

【0027】

他の負荷配置の考慮

前のセクションで説明した統計パッキングアルゴリズムは、ジョブのグルーピングを候補にあげ、どのグループが最良の統合を有しているかを判定し、物理資源間のジョブの移行が統合に有効に作用するロケーションを判断するように適用することができる。つまり、ジョブをパッキングしかつ予約された全体資源を判定するためにグループ統合がいかにうまく作用するかの簡単な測定が提供されている。良いパッキングとは、少ない資源の予約を行うと同時にグループの QoS 要件も満たすことである。

【0028】

しかしながら、この測定だけを使用するのではなく、負荷パターンの相補性、リスク削減ポテンシャル、資源間の QoS のバランシングなどのグループ統合に対する他の考慮する点を知ることが有用な場合もある。これらの考慮する点は、グループ統合が有効に作用する理由について人間のオペレータが洞察する力を与え、最良の統合を検索するために多種多様のグルーピングを検索する必要がないようにグルーピングの計算速度を速くするこ

10

20

30

40

50



とができる。

【 0 0 2 9 】

最も基本的な統合の考慮する点は負荷パターンの相補性である。例えば、あるジョブは常に真夜中に実行され、別のジョブは常に午後 5 時に実行される場合、これらの 2 つのジョブは相補的な負荷パターンを有しており、同じ物理資源を容易に共有することができる。相補的である負荷の場合、負荷は異なる時間で高い負荷を有する予測可能な時間的負荷パターンを有していなければならない。また、これらの負荷パターンは互いにほぼ無関係な関係であるべきである。

【 0 0 3 0 】

監視ツールに対して、ジョブの長期スケールの相関関係に関連している相補性を測定する方法を有することは有用である。実際の時変性の使用量  $x_t$  またはモデリング適用後の QoS 保護された使用量  $r_t$  のいずれかの使用量の測定のシーケンスが与えられた場合、2 つのジョブの間の相関関係は、以下のように測定することができる。

【数 8】

$$\text{corr}(z^{(1)}, z^{(2)}) = \frac{\sum_t (z_t^{(1)} - \bar{z}^{(1)}) (z_t^{(2)} - \bar{z}^{(2)})}{\sqrt{\sum_t (z_t^{(1)} - \bar{z}^{(1)})^2 \sum_t (z_t^{(2)} - \bar{z}^{(2)})^2}} \quad (7)$$

ここで、 $z_t$  は使用可能な測定値（即ち、 $x_t$  または  $r_t$ ）を表し、

【数 9】

$$\bar{z}$$

は平均使用量を表す。 $x_t$  または  $r_t$  のいずれかを使用して相補性を測定することは可能であるが、上位の QoS ジョブが含まれる場合は特に、 $r_t$  をよりしっかりと計算することが望ましい。更に、2 つのジョブが同じ資源を共有する加減を判断するためにはそれらのピーク時のニーズの相関関係を測定する方がよい場合が多い。

【数 10】

$$\max\_corr(z^{(1)}, z^{(2)}) = \frac{\max_t (z_t^{(1)} + z_t^{(2)})}{\max_u (z_u^{(1)}) + \max_v (z_v^{(2)})} \quad (8)$$

ここで、異なる索引  $t$ 、 $u$  および  $v$  はこれらの最大値が必ずしも同時に発生しないことを示すために用いられる。等式 (8) に示されているメトリックは複数のジョブの互いに関連するピーク時のニーズを測定するために一般化され得ることに留意される。2 つのジョブ間の相関するピーク時のニーズがしきい値を下回る場合、または、2 つのジョブ間の相補性レベルがしきい値を上回る場合、同じマシンへこれらの 2 つのジョブを配置することが望ましいことに留意されたい。

【 0 0 3 1 】

一部の有利なバランシングは（例えば、同じマシンまたはクラスタ内ある）同じ物理資源に相補的なジョブを配置することによって達成することができるが、この考え方のみでは通常、バランシングが達成される量に限界がある。実際、ジョブのグルーピングの方法にかかわらず、通常は、より多くの資源を必要とする平日午後などの重い負荷が掛かる時間帯がある。エネルギーの節約は、オフピークの時間帯において過剰な資源をターンオフするためにバッキングを用いることによって達成することができる。しかしながら、ピーク

10

20

30

40

50

時の資源ニーズの統合を改良するためのバランシングとパッキングは、データセンタの全体のキャパシティを判断する。相補性の利点を達成するほかに、他の考慮する点は統合を改良するとともにデータセンタのキャパシティを拡大することができる。

【0032】

いくつかのバッチジョブを除けば、データセンタの大部分のジョブは、予測不可能な資源ニーズを有している。例えば、ジョブの資源ニーズは、ウェブサイトへのビジュアリティや市場の取引のボリュームに依存する。ジョブがその必要とする資源を有することを確実にするために、通常、実際に必要とされる以上の資源を予約することが必要である。しかしながら、予測不可能な資源ニーズをカバーするために各ジョブに個別資源の余分な「バッファ」を予約させることは無駄な努力である。その代わりに、保険会社がリスクをプールすることに非常に似た方法でこれらの余分な資源をプールすることによる統合の機会が与えられる。これにはリスクが互いに無相関であることが必要とされる。例えば、保険会社にとっては大規模な地震より小規模な個別の火災に保険を掛ける方がはるかに簡単である。一般的にいえば、リスク削減の考え方として、同じ物理資源へ互いに無相関関係のリスクを有するジョブのプールを配置することを提案する。

10

【0033】

ジョブの間でのリスク削減能力を測定するには、ジョブの短期スケールの相関関係を見ていく必要がある。例えば、2つのジョブの各々が同一のQoSレベル $p$ を有し、次の時間ステップにおける予測される資源ニーズが分布 $(1)(z)$ と $(2)(z)$ によって与えられる状況を考える。これらの分布が独立している場合、ジョブの資源ニーズと一緒に考えることによってリスクを削減することができる。ジョブが通常、平均値 $\mu^{(1)}$ と $\mu^{(2)}$ と標準偏差値 $\sigma^{(1)}$ と $\sigma^{(2)}$ により分布されると仮定する。正規分布は、非現実的に良好に振舞っているが、良好な例として作用する。正規分布のテール、

20

【数11】

$$Q\left(\frac{x-\mu}{\sigma}\right) = \frac{1}{2} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}} \quad (9)$$

30

の領域に対するチャーノフ境界を用いて、QoSレベル $p$ を達成するために、これらのジョブ各々に対して要求された個別予約は、

【数12】

$$r^{(1)} = \mu^{(1)} + \sigma^{(1)} \sqrt{-2 \log(2p)} \quad (10)$$

と、

【数13】

$$r^{(2)} = \mu^{(2)} + \sigma^{(2)} \sqrt{-2 \log(2p)} \quad (11)$$

40

となる。

しかしながら、これらのジョブの結合された資源ニーズは通常平均値 $\mu^{(1)} + \mu^{(2)}$ と標準偏差値 $\sqrt{(\sigma^{(1)})^2 + (\sigma^{(2)})^2}$ を用いて分布される。

よって、適切な結合予約は、

【数 1 4】

$$r^{(T)} = \mu^{(1)} + \mu^{(2)} + \sqrt{(\sigma^{(1)})^2 + (\sigma^{(2)})^2} \sqrt{-2\log(2p)} \quad (12)$$

となる。

結合された標準偏差値は、個別の標準偏差値では直線形に成長しないので、必要とされる予約、 $r^{(T)} < r^{(1)} + r^{(2)}$ において削減されることに注意されたい。独立したジョブの結合によって保険会社が複数の独立したリスクをプールすることによってリスク削減を達成するのに非常に似たリスク削減を得ることができる。例えば、 $n$ 個のジョブを同じと結合することによって結合予約の項において $1/n$ の削減が得られる(不確実性による部分)。そこで、大まかにいえば、同じ物理資源へ独立したジョブを結合することによって $1/n$ の削減を得ることができる。

10

【0034】

ジョブがほぼ独立している重要性のお陰で、リスク削減のポテンシャルを測定するために、ジョブにおける不確実性が独立している度合を測定することを必要がある。1つにはジョブ間の短期スケールの相関関係を測定することである。

【数 1 5】

$$\text{corr}(z^{(1)}, z^{(2)}) = \frac{E\left(\left(z_t^{(1)} - \bar{z}^{(1)}\right)\left(z_t^{(2)} - \bar{z}^{(2)}\right)\right)}{E\left(z_t^{(1)} - \bar{z}^{(1)}\right)E\left(z_t^{(2)} - \bar{z}^{(2)}\right)} \quad (13)$$

20

予測が次の時間ステップにおける分布 $(1)(z)$ と $(2)(z)$ を用いて計算され、相関関係がこれらの分布の平均値に依存しないことに留意されたい。等式(7)によって記述される長期スケールの相関関係はシーケンシャルな時間ステップにわたって共に移動する平均値を明確にする。ジョブの資源ニーズが高い時、我々はジョブ間の相関関係への関心に最も高いので、相関関係のより有効な測定はジョブによって導かれる以下の「平均シフト」である。

30

【数 1 6】

$$\text{mean\_shift}(r, z^{(1)}, z^{(2)}) = E(z^{(2)} | z^{(1)} \geq r) - E(z^{(2)}) \quad (14)$$

この等式は、それ自体の予約を超えたジョブが別のジョブも余分な資源を取った可能性がどの程度であるかについていくつかの洞察を与える。等式(14)によって計算された平均シフトは、単一ジョブと複数のジョブの別のグループとの間の平均シフトへ簡単に一般化される。ジョイントガウス分布などの良好に振舞われた分布の場合において導出された平均シフトは以下に示すように相関関係に直接関係している。

40

【数 1 7】

$$\text{mean\_shift}(r, z^{(1)}, z^{(2)}) = \text{corr}(z^{(1)}, z^{(2)}) \frac{\sqrt{\frac{2}{\pi}} \sigma^{(2)} e^{-\frac{(r-\mu^{(1)})^2}{2(\sigma^{(1)})^2}}}{\text{Erfc}\left(\frac{r-\mu^{(1)}}{\sqrt{2}\sigma^{(1)}}\right)} \quad (15)$$

しかしながら、現実の分布に於いてこれらはそれほど直接関係していない可能性があり、

50

等式(14)に基づいて経験的に計算された平均シフト測定は資源割当てにさらに関係している場合がある。

【0035】

他の重要な統合の考え方は、同じ物理資源上のジョブがQoS要件の良好な混合を有しているかどうかである。不確実な資源ニーズを有しているジョブが同じ物理資源へ統合される場合、それはジョブの間での異なるQoS要件の混合を有する助けをする。言い換えれば、上位のQoS要件を有しているジョブを同じマシンへ配置するのを避けるべきであり、その逆もある。統計パッキングを用いて、上位のQoSジョブのニーズを満たすために行われる大きな予約が下位のQoSジョブのニーズも満たすことができる場合、優れた利益をもたらされる。これは上位QoSジョブの予約が必要とされることはあまりないので、これらの通常は休止中の資源でも下位のQoSジョブに十分に作用することができる。つまり、良好に混合されたQoSを有するジョブを実行する物理資源上で下位QoSジョブは上位QoSジョブの未使用の予約を有効に選り分けることができる。

10

【0036】

上位と下位のQoSジョブの混合が同じ物理資源を共有する時、グループの統合の利点が最大になるので、グループ統合の利点を高めるためには物理資源間の負荷のバランスをとるためのQoSの考え方を採用することによって非常に有利になる。前のセクションで説明したように、複数のジョブが共通の物理資源上で統合された場合、ジョブがQoS要件を満たすために必要な更なる保護的予約の一部を共有することができるので、必要とされる資源の全体量を削減することが可能である。

20

【0037】

すべての種類のジョブが資源を共有することによって利益を得ることができるが、大量の予約された資源を必要とする上位QoSジョブとそれらの下位QoS要件を満たすために「使えるだけ」ベースで休止中の資源をうまく選り分けることができる下位QoSジョブとの間では、特に良好な相乗効果が発揮される。ここで、われわれは複数のジョブを実行する物理資源ごとにQoSバランスポイントの概念を伝える。意図としては上位QoSジョブが下位QoSジョブによりバランスをとることである。QoSのバランスポイントを定量的に記述するために、われわれがQoS格付け(Q<sub>r</sub>)の概念を伝える必要がある。各ジョブは、ジョブに関連付けられたQoS要件を反映する正の実数値のQoS格付けを有している。QoS格付けを割り当てるにはさまざまな方法があり、例えば、QoS格付けは、等式(3)の個別の予約に基づいて行われてもよいし、または、これらの格付けは、ジョブのQoSとその直近の資源消費から計算された優先順位に基づいて行われてもよい。上位QoSジョブは上位QoS格付けを有し、その逆もある。QoS格付けは、ジョブの許容された不足確率である、QoSレベルpに混同されないことに留意されたい。QoS格付けはQoSレベルに逆相関する場合が多い、すなわち、より小さい許容不足確率はより上位のQoS格付けにつながる。

30

【0038】

物理資源編成のレベルごとに、QoSのバランスポイントは、等式(6)のkを検索することによって計算することができる。すなわち、ジョブのグループのバランスポイントは、統計パッキングアルゴリズム(等式(6))によって算定される索引kにおけるQoS格付けである。バランス値はQoS格付けに基づいて計算されるが、その重要な特徴は、統計パッキングがより有効に共有依存方向に遷移しているQoS順序のポイントにおけるロケーションである。QoS格付けを個別のジョブに割り当てる方法は、バランスポイントをジョブが移行される他の物理資源上のバランスポイントと比較する役割を果たす。例えば、それらの個別の予約要件に基づいているが、統合結果を反映していない個別のジョブのQoS格付けを算定するための方法は、統合結果を反映する計算されたバランスポイントを物理資源間で比較することを可能にする。

40

【0039】

QoSのバランスポイントはグループのパフォーマンスを向上させるための機会があり、そのような単純で明確な信号を発信する。例えば、クラスタ内のすべての物理マシンがほぼ同

50

じQoSバランスポイントを有している場合、そのクラスタは「インバランス（バランスが取れている）」である。一方、2つのマシンがかなり異なるバランスポイントを有している場合、2つのマシンは「アウト・オブ・バランス（バランスが取れていない）」である。したがって、ジョブの移行はバランスを改良し、ひいては、全体的なクラスタパフォーマンスを改良することができる。図2は、本発明の一実施形態による、2つのマシンが「アウト・オブ・バランス」である状況を示す図である。図2において、物理マシン（PM）クラスタ200は、各々が多数の仮想マシン（VM）またはジョブを実行する、PM202とPM204を含む。例えば、PM202は、VM206、VM208、VM210などを実行し、PM204は、VM212、VM214、VM216などを実行する。各VMはそれ自体のQoS要件を有している。例えば、VM206、208、および210に対するQoS格付けはそれぞれ、10.1、7.9、および3.9である。

10

#### 【0040】

PMで現在実行されている各ジョブのQoS格付けと資源ニーズに基づいて、システムはこのPMに対するQoSバランスポイントを計算することが可能である。一実施形態に於いて、先に説明した統計パッキングアルゴリズムはQoSバランスポイントを計算するために使用される。例えば、統計パッキングに基づいて、システムは、PM202に対してはVM208を上回るQoS格付けを有するVMを完全予約、そしてVM208に対しては部分予約を行うことによって、PM202で実行しているすべてのVMに対する資源ニーズとQoS要件が満たされることができると判断する。言い換えれば、PM202に対するQoSバランスポイントは、7.9であるVM208のQoS格付けである。先に説明した統計パッキングアルゴリズムが使用されている場合、VM210などのVM208を下回るQoS格付けを有するVMに対して予約が行われないうことに留意されたい。しかしながら、下位のQoS格付けを有するVMは上位のQoS格付けを有しているジョブに対して予約された未使用の資源をうまく選り分けることができる。同様に、システムは、PM204に対してはVM214より上位の格付けを有するVMに対して完全予約を行い、そして、VM214に対しては部分予約を行うことによって、PM204上で実行しているすべてのVMに対する資源ニーズとQoS要件を満たすことができると判断する。言い換えれば、PM204のQoSバランスポイントは2.9であるVM214のQoS格付けである。

20

#### 【0041】

統計パッキングはグループ予約を行う方法を提供する。即ち、PM202におけるすべてのVMはグループとして資源予約を行うことで、グループに対する予約された資源の量を大きく削減する。しかしながら、我々がPM202とPM204をクラスタと見なしてクラスタの全体的なパフォーマンスを向上させたい場合、我々はこれらの2つのマシン間でのQoSバランスポイントをバランシングする必要がある。

30

#### 【0042】

図2を見ると、PM202とPM204がそれぞれ大きく異なるQoSバランスポイント7.9と2.9を有していることがわかる。従って、バランスひいてはシステムのパフォーマンスを向上させるために、（矢印218で示した）PM202のバランスポイントを下げるとともに（矢印220で示した）PM204のバランスポイントを大きくすることによってバランスを改良する必要がある。QoSのバランスポイントが物理的オブジェクトの質量の中心のQoSのバランスポイントに非常に近接した直感的解釈を有していることに留意されたい。物理的なオブジェクトの場合、1つの側に質量を加算することでその側に向けて質量の中心をシフトさせる。資源の統合の場合において、QoSバランスポイントを上回る上位QoSジョブはそれらのQoSを維持するために個別の保護予約を必要とする。これとは対照的に、QoSバランスポイントを下回る下位QoSジョブは上位QoSジョブをすでに予約している予約におけるキャパシティのプールを有するQoSを分配する。このプールには限定されたキャパシティのみしか存在していないため、より多くの下位のQoSジョブを1つのマシンへ加算することはQoSのバランスポイントを下位QoSジョブへ向けてシフトさせて、最終的に、それらの一部はQoSバランスポイン

40

50

トと交差し、個別の保護的な予約を必要とする。言い換えれば、QoSバランスポイントを下回るQoS格付けを有するより多くのジョブを加算することによって、QoSバランスポイントを下方にシフトさせるとともに、下位のQoS格付けを有するジョブを除去することによってバランスポイントを上方へシフトさせることになる。よって、矢印222で示すように、PM204からPM202まで、PM204のQoSバランスポイントを下回るQoS格付けを有しているジョブ216の移行によって、PM204のQoSバランスポイントを増加させるとともにPM202のQoSのバランスポイントを低下させることにつながる。これは2つのマシンのバランスポイントを均等化し、クラスタの全体的な性能を改良する傾向がある（より少ないクラスタ内のすべてのQoS要件を満たすことが必要とされる）。PM202～PM204までジョブ208を移動させることも可能であることに留意されたい。このようなジョブの移動はPM202のQoSのバランスポイントを下げるとともにPM204のQoSのバランスポイントを上げることもできる。

10

**【0043】**

しかしながら、バランスポイントが右方向に移動しているからといって、この移行が全体的なパフォーマンスの改良につながるという保証はない。例えば、バランスポイントのシフトが過剰なため、新たな不均衡やリスク削減ポテンシャルなどの他の考慮する点が生じて自体を悪くする場合もある。実際、バランスポイントは有望な好機の合図を送る一方、パッキング計算の前後に基づいて、実際のマージナルコスト計算を用いて、実際の利益を計算して、この移行が改善点なのかどうかを判定することができる。バランスポイントの利点は、これらが有望な好機を識別することで適切な移行を検出するために必要な計算を減らすように良好な方法を提供することである。さらに、負荷監視ツールの一部として使用された場合、バランスポイントは、人間のユーザにQoSがどの程度うまく実施されるかについて洞察させる。

20

**【0044】**

多くのデータセンタは、ジョブが物理資源にどのように位置付けされるかに影響を与える多様な更なる制約を含む。典型的な制約としては、他の関連ジョブのロケーション、ネットワーク接続要件、およびセキュリティ要件が含まれる。例えば、いくつかの互いに関連するジョブが同じクラスタ内に配置され、良好な外部ネットワーク接続を有するクラスタ内に配置され、または特別に固着されたクラスタ内に配置されることが要求される可能性もある。これらの種類の制約は、ジョブが統合を改良するためにグルーピングされるときには尊重すべき制約となり得る。

30

**【0045】**

図3は、本発明の一実施形態による、データ・センタ・オペレーションの負荷配置コントローラを示す図である。図3において、負荷配置コントローラ300は、QoS識別子302と、資源使用モニタ304と、資源使用モデル構築子306と、QoSバランスポイント計算子308と、ジョブ移行コントローラ310と、リスク関連計算子312と、ユーザインターフェース314と、を含む。

**【0046】**

負荷配置コントローラ300は、クラスタに対して最高のパフォーマンスを達成するために物理マシンのクラスタ内の負荷配置をコントロールする。QoS識別子302は、クラスタ内の各ジョブに関連付けられたQoS要件を識別する。資源使用モニタ304は、全体的にランダムまたは部分的に時間的なパターンを有し得る、ジョブごとの資源ニーズを監視する。資源使用モデル構築子306は資源使用モニタ304から情報を受信しこれに応じてジョブごとに資源使用モデルを構築する。一実施形態において、資源使用モデルは、ある一定量の資源を必要とするジョブの確率を示す、資源ニーズ分布関数を含む。更なる実施形態に於いて、資源使用モデルは、ジョブの資源ニーズの時間的分布を含む。例えば、あるジョブが午前中に大量の資源を必要とする確率は高い。更なる実施形態において、資源使用モデルは、互いに関連している資源ニーズを有し得るジョブのグループに対して計算される。

40

**【0047】**

50

構築された資源モデルと特定のPM上で実行しているジョブごとのQoS要件に基づいて、QoSバランスポイント計算子308はそのPMに対してQoSバランスポイントを計算する。一実施形態に於いて、統計パッキングアルゴリズムは、QoS格付け以上のジョブに対して行われたQoS保護された予約がそのPM上の全てのジョブに対して資源ニーズおよびQoS要件を十分に満たすことができるように選択されるQoS格付けを表している、QoSバランスポイントを計算するために使用される。一実施形態に於いて、QoSバランスポイント計算子308は、クラスタ内のPMごとにQoSバランスポイントを計算する。

#### 【0048】

一実施形態に於いて、計算されたQoSバランスポイントは人間の管理者へ提示され、管理者はPM間のQoS不均衡を観察し、次いで、QoSバランスを改良することができるジョブの移行を提案することができる。これらのジョブ移行提案はジョブ移行コントローラ310へ送信され、ジョブ移行コントローラ310はどのジョブがどの方向に移行されるかを制御する。ジョブ移行コントローラ310は相補性およびリスク削減などの更なる想定を取り入れることもできる。提案された移行のいくつかに対して、限界便益の計算は統計パッキングアルゴリズムを用いて計算することができる。更なる実施形態に於いて、ジョブ移行コントローラ310は、バランス情報や恐らく他の上記したメトリックを有する人間の管理者の介入によってジョブを移行する方法を判断する。ジョブ移行コントローラ310はまた、最終のジョブ移行決定を行う前に他の考えを考慮に入れる必要もある。一実施形態に於いて、ジョブ移行コントローラはジョブの相補性を考慮し、1つのPM上で相補性の時間的パターンを有するジョブを配置する移行の提案を支持する。一実施形態に於いて、提示されたジョブ移行はQoSバランスポイント計算子308へ返送され、この計算子308は提示されたジョブ移行に基づいてPMごとに更新されたQoSバランスポイントを計算する。更新された結果がQoSバランスを改良する場合、これに応じて、ジョブ移行コントローラ310は進行してPM間でジョブを移動する。そうでない場合、この提案は放棄される。さらに、リスク相関関係計算子312はリスク相関ファクタを計算し提示されたジョブの移行に基づいてリスク削減ポテンシャルを更新する。一実施形態において、リスク削減ポテンシャルは単一マシン上のジョブの間での相関関係を測定することによって評価される。更なる実施形態に於いて、このシステムは提示されたジョブの移行によって導入された平均シフトを測定する。例えば、1つのジョブが物理マシンへ加算された場合、システムは加算されたジョブと現在マシンにある他のジョブとの間の平均シフトを測定することができる。

#### 【0049】

更新されたリスク削減ポテンシャルが所定のしきい値より大きい場合（または提案されたジョブ移行の結果としてリスクの計算された相関関係が相変わらずしきい値より小さい場合）、ジョブ移行コントローラ310は進行してジョブを移動する。更新されたリスク削減ポテンシャルがその元の値（移行前の値）に比較して大きく削減された場合、システムは提案されたジョブの移行が有益ではないと判断してこの提案を拒否する。例えば、被加算ジョブが現在マシン上にある他のジョブと強力的に互いに関連している場合（即ち、この被加算ジョブと他のジョブが同時に大きな資源を必要とする可能性が高い場合）、このような移行は全体的なシステムパフォーマンスに対して有利ではない。

#### 【0050】

図4は、本発明の一実施形態による、例示的なジョブ移行プロセスを示す流れ図である。動作中、システムはマシンのクラスタ内のジョブごとにQoS要件を識別し（動作402）、ジョブの資源使用履歴に基づいてジョブごとに資源使用モデルを構築する（動作404）。一実施形態において、資源使用モデルは資源使用確率分布関数を含む。更なる実施形態において、資源使用確率分布関数は経時的に変化する。

#### 【0051】

次に、システムはマシン上で実行するジョブに対するQoS要件と資源使用モデルに基づいてPMごとにQoSバランスポイントを計算し（動作406）、2台のマシンの間に

10

20

30

40

50

QoS不均衡が存在するかどうかを判断する(動作408)。一実施形態において、計算されたQoSのバランスポイントは管理者へ提示され、管理者はクラスタ内の任意のQoS不均衡を識別することができる。不均衡が存在する場合、管理者または自動コントローラは2台のマシン間のQoSの不均衡を改良することができるジョブの移行を提案する(動作410)。1つ以上のジョブが提案された移行に含まれることに留意されたい。

#### 【0052】

ジョブの移行に基づいて、システムは2つのマシンのためのQoSバランスポイントを再計算する(言い換えれば、潜在的な移行後の新しいバランスポイントを計算する)(動作412)。システムはバランスが改良されたかどうかを更に判断する(動作414)。バランスの改良が見られない場合、提案された移行は拒否される(動作416)。提案されたジョブの移行がQoSバランスを改良することができる場合、システムは新しいジョブ分布に対するリスク削減ポテンシャルを必要に応じて評価し(動作418)、両方のマシンに対する新しく評価されたリスク削減ポテンシャルが所定のしきい値を上回るかどうかを判断する(動作420)。所定のしきい値を超えた場合、システムは進行してジョブの移行を終了する(動作422)。しきい値を超えなかった場合、提示された移行は拒否される(動作416)。

#### 【0053】

### コンピュータシステム

図5は、本発明の一実施形態による、データセンタ内のマシン間でジョブを配置するための例示的なコンピュータシステムを示している。一実施形態において、コンピュータおよび通信システム500は、プロセッサ502、メモリ504、および記憶装置506を含む。記憶装置506はジョブ配置アプリケーション508のみならずアプリケーション510、512などの他のアプリケーションを記憶する。動作中、ジョブ配置アプリケーション508は記憶装置506からメモリ504へロードされ、次いで、プロセッサ502によって実行される。プログラムを実行している間、プロセッサ502は前述した機能を実行する。コンピュータおよび通信システム500は、任意選択のディスプレイ514、キーボード516、およびポインティングデバイス518に連結される。

#### 【0054】

本発明の実施形態は、特に、QoSの慎重な管理を最も優先すべきである状況に於いて、データセンタの負荷をバランスングするための解決を提供する。QoSの慎重な管理は、大抵の場合、不十分な資源のリスクを緩和するために余分な予約を含む。グループとして予約を行うことによって物理資源の全体的に必要なとされる予約は削減することができる。データセンタにおける負荷のバランスングの最も重要な格付けは必要とされる物理資源の量である。従って、本発明の実施形態は、グループの予約を行うために統計パッキングアルゴリズムを使用する。データセンタがどの程度良好にバランスングされているかの判定するための良好な全体的なメトリックであるだけでなく、統計パッキングアルゴリズムは、データセンタを再度バランスングするためにジョブの潜在的な移行を評価するために使用することができる。より具体的には、移行前に必要とされる物理資源を移行が有効であるかどうかを判断するためにその後必要とされる物理資源に比較することができる。負荷とリスク削減ポテンシャルの相補性などのQoSの考え方と他の考え方の組合せによってデータセンタ内のマシン間で負荷がバランスングされる目安を判断する。これらの考慮する点方はバランスングをある程度コントロールしたいとする人間のオペレータにデータセンタがどの程度有効に作用するかを洞察させる。さらに、これらの考え方は、自動アルゴリズムがパフォーマンス向上のための最も有望なりバランスング行動を探索することも可能にする。

#### 【0055】

バランスングの以上の考え方は、クラスタ内のマシン間、クラスタ間、または、データセンタ間においてさえも、ジョブを移動することを含む、データセンタにおいて複数レベルで、適用することができ、これによってQoSパフォーマンスと資源利用が改良されることに留意されたい。



【 図 1 】

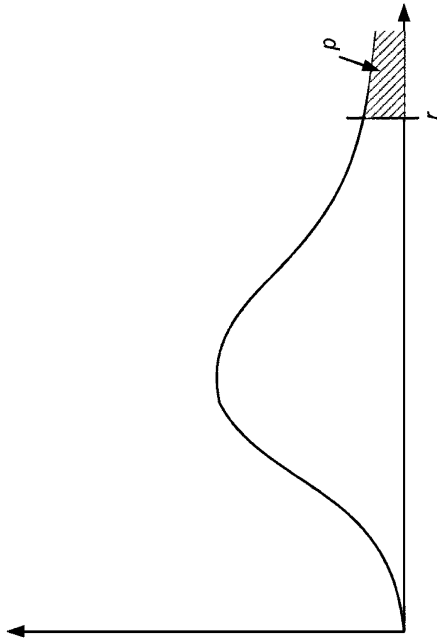
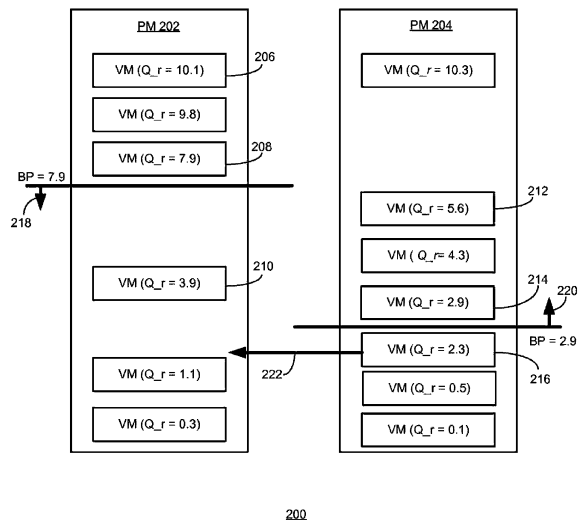


図 1

【 図 2 】



200

図 2

【 図 3 】

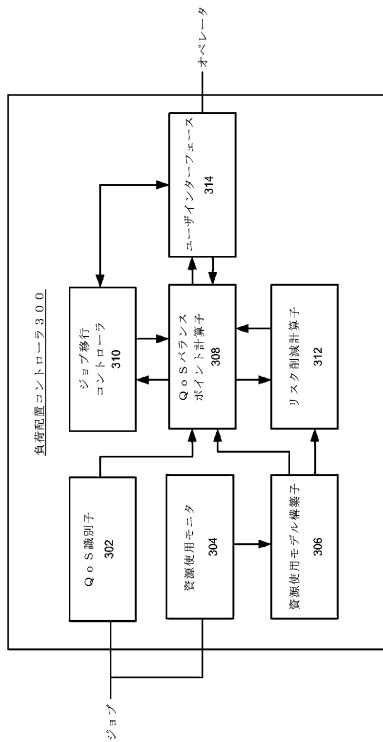


図 3

【 図 4 】

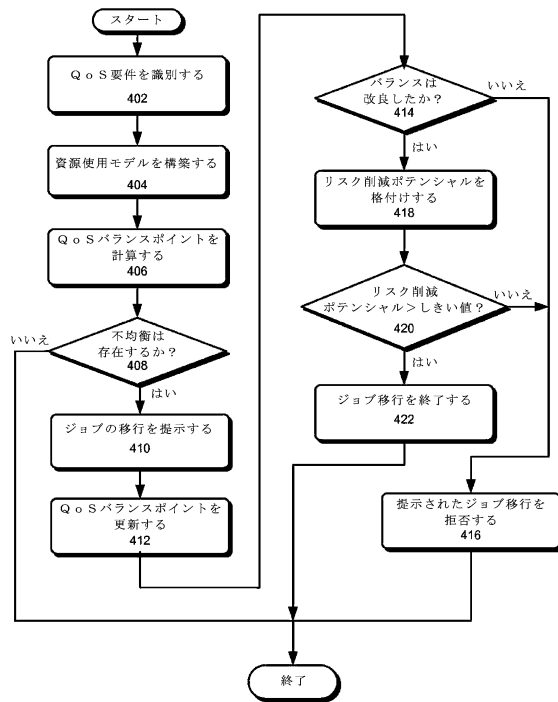


図 4

【 図 5 】

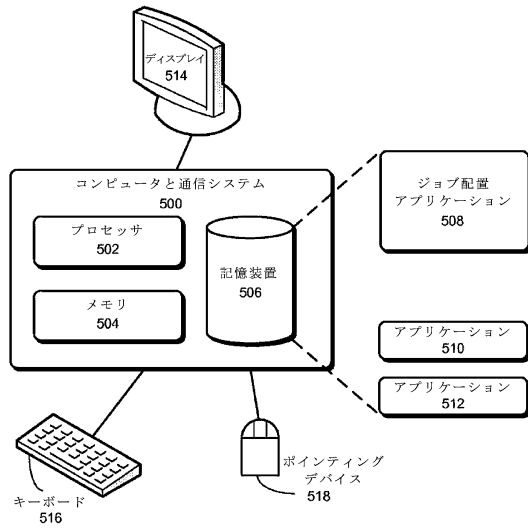


図 5

---

フロントページの続き

- (72)発明者 ララ・エス・クロフォード  
アメリカ合衆国 カリフォルニア州 94024 ロス・アルトス プレントウッド・プレイス  
720
- (72)発明者 ジョン・ハンレー  
アメリカ合衆国 カリフォルニア州 94062 エメラルド・ヒルズ パーク・ロード 246  
6

審査官 漆原 孝治

- (56)参考文献 特開2010-117760(JP,A)  
特開2009-199395(JP,A)  
特開2005-115653(JP,A)  
特開2010-244181(JP,A)  
米国特許出願公開第2011/0302578(US,A1)

- (58)調査した分野(Int.Cl., DB名)  
G06F 9/48  
G06F 9/50