

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4286703号
(P4286703)

(45) 発行日 平成21年7月1日(2009.7.1)

(24) 登録日 平成21年4月3日(2009.4.3)

(51) Int. Cl. F 1
G 0 6 F 9/50 (2006.01) G 0 6 F 9/46 4 6 2 Z

請求項の数 2 (全 34 頁)

(21) 出願番号	特願2004-105093 (P2004-105093)	(73) 特許権者	000005223
(22) 出願日	平成16年3月31日 (2004.3.31)		富士通株式会社
(65) 公開番号	特開2005-293048 (P2005-293048A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成17年10月20日 (2005.10.20)	(74) 代理人	100089118
審査請求日	平成19年3月22日 (2007.3.22)		弁理士 酒井 宏明
		(72) 発明者	園生 泰廣
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	土屋 哲
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	櫻井 英樹
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 資源計画作成プログラム

(57) 【特許請求の範囲】

【請求項1】

複数の計算機から構成される計算機クラスタが複数接続されて構成される計算機クラスタシステムが運用する複数のサービスを該複数の計算機クラスタへ配置する資源計画作成する資源計画作成プログラムであって、

運用される計算機クラスタが定められたサービスである原始サービスから該原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用の一部が移される派生サービスを、負荷予測の予測期間を所定の時間間隔で複数のタイムスロットに分割したタイムスロット毎に生成する派生サービス生成手順と、

前記派生サービス生成手順により生成された派生サービス群を前記複数の計算機クラスタへ最適配置する資源計画として、タイムスロット毎の派生サービスの配置計画であるタイムスロット配置計画を該タイムスロットの時刻順に並べた配置計画系列を、計算機クラスタ間の負荷均等度と、時間軸上で隣接するタイムスロット間のタイムスロット配置計画の類似度とを評価指標として用いて作成する計画作成手順と、

をコンピュータに実行させることを特徴とする資源計画作成プログラム。

【請求項2】

前記計画作成手順は、前記類似度と負荷均等度を評価指標とした場合に評価が最も良い配置計画系列を初期系列集合とし、該初期系列集合に含まれる配置計画系列に対して前記負荷均等度の評価がより良くなるような補正を行うことによつて前記資源計画作成することを特徴とする請求項1に記載の資源計画作成プログラム。

10

20

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、複数の計算機から構成される計算機クラスタが複数接続されて構成される計算機クラスタシステムが運用する複数のサービスの該複数の計算機クラスタへの配置計画を作成する資源計画作成プログラムに関し、特に、各計算機クラスタが過負荷になることを防ぐとともに、計算機クラスタの計算資源を効率的に使用することができる資源計画の作成プログラムに関するものである。

【背景技術】

【0002】

近年のデータセンタは、一つのデータセンタ内において、複数のサービスプロバイダ事業者から運用を委託された複数のネットワーク・サービスを同時に運用するようになってきており、それと共に、同時運用する複数のサービスの中で負荷の増加しているサービスに対しては、データセンタ内の計算資源の割当量を動的に増やし、負荷の減少しているサービスにはデータセンタ内の計算資源の割当量を動的に減らすユーティリティ方式に基づく資源割当制御を行っている（例えば、特許文献1参照）。

【0003】

しかし、データセンタで運用している一部のサービスが一時的に過負荷となり、このサービスに対する資源割当を大幅に増やす必要が生じたときに、新たに追加割当すべき遊休資源がデータセンタ内で不足し、資源追加が出来ない事態が発生する可能性がある。

【0004】

このような事態の発生を防ぐためには、データセンタでの全ての運用サービスの負荷の合計がピークになる時の必要資源量を予測し、データセンタに恒常的にピーク時必要資源量を用意することが考えられるが、この場合には、ピーク時以外の期間中は大半の資源量が遊休資源として使われないうままになり、データセンタ全体としてコストパフォーマンスや計算資源利用効率が低下してしまう。

【0005】

そこで、国際特許出願JP03/03273では、後方支援データセンタを用いる方式が提案されている。この方式では、通常、各サービスはサービスの運用を依頼したサービスプロバイダ事業者が指定した特定のデータセンタ（以下「オリジナル・センタ」という。）で恒常的に運用される（以下、このようなサービスを「原始サービス」という。）。

【0006】

そして、原始サービスが負荷の時間変動によって一時的に過負荷になり、それによってオリジナル・センタ自体もセンタ全体として過負荷状態になった場合、オリジナル・センタ上の過負荷になった原始サービスの一部を、後方支援センタと呼ぶ別異のデータセンタで処理するため、後方支援センタ上に当該原始サービスに対応する補助的サービス（以下「派生サービス」という。）を発生させて後方支援センタ上で運用を開始する。

【0007】

これによって、過負荷になった原始サービスが負担していた負荷の一部を後方支援センタ上の派生サービスが肩代わりすることができ、オリジナル・センタのセンタ全体としての過負荷状態を解消できる。

【0008】

ここで、各サービスに到着したユーザ・リクエスト量の原始サービスと派生サービスへの分配比率はネットワーク内の負荷分散装置またはリダイレクタに設定され、当該負荷分散装置等は設定された比率に基づいて、サービスに到着したユーザ・リクエスト量を原始サービスと派生サービスの間で振り分ける。

【0009】

このような従来のデータセンタ運用形態においては、オリジナル・センタは一般に複数個あり、後方支援センタは一つであることが多い。また、オリジナル・センタと後方支援センタの役割分担は明確であり、後方支援センタは専ら派生サービスの運用のためだけに

10

20

30

40

50

用いられていた。

【0010】

【特許文献1】特開2002-24192号公報

【発明の開示】

【発明が解決しようとする課題】

【0011】

しかしながら、特定のデータセンタを後方支援センタとして定め、全ての派生サービスを専らそのデータセンタでのみ運用する従来の技術には、複数のオリジナル・センタ上の多数の原始サービスから同時に一つの後方支援センタに対して派生サービスの発生要求が集中した場合、後方支援センタがセンタ全体として過負荷状態になるという問題がある。

10

【0012】

また、派生サービスの発生要求の集中がピークになる時に後方支援センタで必要となる資源量を予測し、後方支援センタに恒常的に当該ピーク時必要資源量を用意するとすると、ピーク時以外は後方支援センタの計算資源の大半は使われないこととなり、資源利用効率が低下する。

【0013】

この発明は、上述した従来技術による問題点を解消するためになされたものであり、データセンタを典型例とするような計算機資源の集合体である計算機クラスタの各々が過負荷になることを防ぐとともに、計算機クラスタの計算資源を効率的に使用することができる資源計画作成プログラムを提供することを目的とする。

20

【課題を解決するための手段】

【0014】

上述した課題を解決し、目的を達成するため、本発明は、複数の計算機から構成される計算機クラスタが複数接続されて構成される計算機クラスタシステムが運用する複数のサービスを該複数の計算機クラスタへ配置する資源計画作成プログラムであって、運用される計算機クラスタが定められたサービスである原始サービスから該原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用の一部が移される派生サービスを、負荷予測の予測期間を所定の時間間隔で複数のタイムスロットに分割したタイムスロット毎に生成する派生サービス生成手順と、前記派生サービス生成手順により生成された派生サービス群を前記複数の計算機クラスタへ最適配置する資源計画作成プログラムとして、タイムスロット毎の派生サービスの配置計画であるタイムスロット配置計画を該タイムスロットの時刻順に並べた配置計画系列を、計算機クラスタ間の負荷均等度と、時間軸上で隣接するタイムスロット間のタイムスロット配置計画の類似度とを評価指標として用いて作成する計画作成手順と、をコンピュータに実行させることを特徴とする。

30

【0015】

この発明によれば、運用される計算機クラスタが定められたサービスである原始サービスから原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用の一部が移される派生サービスを、負荷予測の予測期間を所定の時間間隔で複数のタイムスロットに分割したタイムスロット毎に生成し、生成した派生サービス群を複数の計算機クラスタへ最適配置する資源計画作成プログラムとして、タイムスロット毎の派生サービスの配置計画であるタイムスロット配置計画をタイムスロットの時刻順に並べた配置計画系列を、計算機クラスタ間の負荷均等度と、時間軸上で隣接するタイムスロット間のタイムスロット配置計画の類似度とを評価指標として用いて作成するよう構成したので、過負荷の計算機クラスタが運用するサービスの一部を計算機クラスタシステム全体で肩代わりすることができる。

40

【発明の効果】

【0016】

本発明によれば、過負荷の計算機クラスタが運用するサービスの一部を計算機クラスタシステム全体で肩代わりするので、各計算機クラスタが過負荷になることを防ぐとともに、計算機クラスタの計算資源を効率的に使用することができるという効果を奏する。

50

【発明を実施するための最良の形態】**【0017】**

以下に添付図面を参照して、この発明に係る資源計画作成プログラムの好適な実施例を詳細に説明する。なお、ここでは、本発明を複数のデータセンタの計算資源を複数のサービスに割り当てる資源計画の作成に適用した場合を中心に説明する。

【実施例】**【0018】**

まず、本実施例に係る資源計画の概念について説明する。図1は、本実施例に係る資源計画の概念を説明するための説明図である。図1の左側のグラフには、XXチケット販売会社の提供するWEBサービスに関する年間の負荷変動予測がプロットされている。

10

【0019】

このグラフ内の水平の点線で示した高さがXXチケット販売会社の自社センタの保有資源量で処理できる負荷の量を表している。このグラフの負荷変動は2箇所のピークを持ち、いずれのピークにおいても自社センタ資源で処理できる負荷量(点線の高さ)を超えてしまっている。そこで、本実施例に係る資源計画では、ピーク時に自社センタ以外のデータセンタに派生サービスを発生させて点線を超えた分の負荷を肩代わりさせる。

【0020】

ここで、左側のピークの発生期間においては、データセンタXでは資源が殆ど消費されていないのに対し、データセンタYでは大半の資源が他のサービスなどによって消費されている。従って、左側のピーク時にはデータセンタXに派生サービスを配置する。

20

【0021】

一方、右側のピークの発生期間においては、データセンタXの大半の資源が他のサービスなどによって消費されており、データセンタYでは資源が殆ど消費されていないので、データセンタYに派生サービスを配置する。

【0022】

すなわち、本実施例に係る資源計画では、複数個あるオリジナル・センタの内、運用中の原始サービスが低負荷となったことにより一時的に遊休資源が大量に余っているオリジナル・センタを選んで臨時の後方支援センタとして扱い、当該臨時の後方支援センタで、派生サービスを運用する。

【0023】

30

このような運用形態に基づくならば、オリジナル・センタと後方支援センタの役割分担は各データセンタに固有のものではなく、データセンタXが原始サービスAを恒常的に運用しており、データセンタYが原始サービスBを恒常的に運用している場合、データセンタXは、原始サービスAにとっては常にオリジナル・センタとしての役割を持ち、原始サービスBにとっては後方支援センタとしての役割を持つ場合がある。以後、ここでは、オリジナル・センタと後方支援センタの役割分担は上記のようにサービスとの相対関係で決まるものとする。

【0024】

ここで、多数のデータセンタがネットワークで繋がった環境で多数の原始サービスを運用しようとする際には、ある原始サービスが将来において過負荷になった際、どのデータセンタを派生サービス運用のための後方支援センタとして選択すればよいか問題となる。つまり、将来発生することが予想される各派生サービスを各々どのデータセンタにいつ何時に配置すべきかが問題となる。

40

【0025】

ただし、ある将来時点で発生している各派生サービスを各計算機クラスタへどのように配置するかについての最適解の探索については、計算量と計算時間を考慮する必要がある。例えば、データセンタが3個あり、ある将来時点で30個の派生サービスが同時に発生しているとすると、ある1個の派生サービスは、発生元の原始サービスのオリジナル・センタ以外の二つのデータセンタのどちらかに配置可能だから、30個の派生サービスの3個のデータセンタへの配置の仕方は2の30乗(10億以上)通りあることになる。

50

【 0 0 2 6 】

従って、これら全てを単純探索により評価するのは膨大な時間の無駄となり、派生サービス集合のデータセンタ集合への配置の最適解を発見するに当たっては、探索ステップ数ができるだけ少なく抑えられるような最適解探索手法が必要となる。

【 0 0 2 7 】

具体的には、データセンタ群と原始サービスの集合が与えられ、各原始サービスを運用するオリジナル・センタが指定された場合、派生サービスの将来発生時点を予め予測し、全てのデータセンタの計算資源利用効率が改善されるように、様々な時点で発生する各派生サービスを各々どのデータセンタで運用するか（派生サービスのセンタ配置）を最適化する必要がある（図2）。

10

【 0 0 2 8 】

また、派生サービス集合の構成や各派生サービスの処理する負荷量、各データセンタの資源消費状況などは時々刻々と変化するので、時間の経過に従って、派生サービス集合のデータセンタ群への配置の最適化を定期的に行い直す必要がある。上記の配置の最適化を定期的に行い直す際は、時間軸をタイムスロットと呼ぶ周期に分割し、各タイムスロット毎に最適な派生サービス配置を求める（図3）。

【 0 0 2 9 】

従って、最終的には、派生サービス集合のデータセンタ群への配置パターンの最適候補を一定時間周期毎（各タイムスロット毎）に求めて、それをタイムスロットの時刻順に時系列として並べた情報を出力する必要がある。この出力情報のことを以後、データセンタ群の資源計画と呼ぶ。

20

【 0 0 3 0 】

また、最適な派生サービス配置の時系列を探索するに当たって、同時発生する派生サービスの個数によっては、探索空間が巨大となり、単純探索では膨大な時間がかかるので、それを防ぐために、本実施例に係る資源計画では、以下の工夫をする。

【 0 0 3 1 】

探索対象の解は一般に複数個の指標によって評価され、それら全ての指標をなるべく最適にする解が最終的に選ばれなくてはならないが、探索空間の中から少なくとも一つの指標が最適であるような解から構成される集合（Pareto解集合）を絞り込み、そのPareto解集合の中から特定の基準によって任意の解を選び、その解を暫定解とする。

30

【 0 0 3 2 】

そして、後述するNDC入れ替え操作により、暫定解を漸近的に真の最適解に近づけてゆき、暫定解の評価値が一定の収束半径に入ったらその暫定解を最終的な最適解として出力する。

【 0 0 3 3 】

なお、データセンタが運用するサービスには、WEBサービスのようにユーザ・リクエストを受け付けてから数秒以内に応答を返さなくてはならないようなリアルタイム性を要求されるトランザクション型のサービスと、毎日、毎月といった一定の周期で決められた時刻に投入された複数のジョブ群を、予め定められたデッドライン時刻までに全て処理完了する必要があるバッチ型のサービスがある。

40

【 0 0 3 4 】

バッチ型のサービスは複数個のジョブによって構成され、一般に複数回に分けて、各々決められた時刻に、複数個ずつのジョブがまとめてデータセンタに投入される。このジョブの投入スケジュールは一日周期や1ヶ月周期等の一定周期毎に繰り返される。ジョブの投入スケジュールは、そのバッチ型サービスの運用を委託した委託元業者によって事前に指定される。

【 0 0 3 5 】

バッチ型サービスに関しては、投入されたジョブ群が指定されたデッドライン時刻までに全て処理完了するのであれば、個々のジョブは、いつ何時処理されてもよく、いわばリアルタイム性を必要としない。

50

【 0 0 3 6 】

従って、データセンタ群で運用する原始サービス集合がトランザクション型のサービスとバッチ型のサービスの2種類のサービスから構成されている場合は、トランザクション型サービスがリアルタイム性を要求されるサービスであることを考慮し、データセンタ群の持つ計算資源は、トランザクション型サービスに優先的に割り当てる。

【 0 0 3 7 】

その後、トランザクション型サービスに割り当てられなかった余剰計算資源をバッチ型サービスの処理に割り当てるような資源計画の生成を行う必要がある。従って、データセンタ群の資源計画を生成するには、まず、トランザクション型の原始サービスのみを当該データセンタ群で運用したと仮定して資源計画の生成を行い、その資源計画においてトランザクション型サービスに割り当てられていない余剰計算資源に対し、バッチ型サービスの資源計画の生成を行う。

10

【 0 0 3 8 】

また、バッチ型原始サービスの負荷変動予測に関しては、バッチ型サービスのジョブ投入スケジュールが予め既知であることから、ARIMAモデル手法のような数理的な予測アルゴリズムによってサービスの将来負荷変動を予測するのではなく、将来の負荷変動の推移を特定のジョブ・スケジューリング・アルゴリズムによって計画することが出来る。

【 0 0 3 9 】

従って、バッチ型原始サービスの集合に関する資源計画の生成では、バッチ型原始サービスの将来の負荷変動は予測するのではなく、ジョブ・スケジューリング・アルゴリズムによって事前に計画されたものを用いるという形をとる。

20

【 0 0 4 0 】

具体的なジョブ・スケジューリング・アルゴリズムは、後述するが、当該スケジューリング・アルゴリズムの特徴は、バッチ型原始サービスを構成する各ジョブを、投入時刻の遅い順にスケジューリングしてゆき、スケジューリング対象となったジョブを、当該ジョブの投入時刻からデッドライン時刻までの間の時間帯の中でデータセンタ群の余剰計算資源の最も多い時刻に当該ジョブをスケジューリングするというものである。

【 0 0 4 1 】

このようなジョブ・スケジューリング・アルゴリズムを採用する理由は、バッチ型サービスの負荷の処理は、なるべくトランザクション型サービスによる負荷が空いている時を選んで行わせることにある。

30

【 0 0 4 2 】

次に、本実施例に係る資源計画作成装置の構成について説明する。図4は、本実施例に係る資源計画作成装置の構成を示す機能ブロック図である。同図に示すように、この資源計画作成装置100は、トランザクション計画部110と、バッチ計画部120と、資源計画送信部130とを有する。

【 0 0 4 3 】

トランザクション計画部110は、トランザクション型サービスに対して資源計画を作成する処理部であり、予測部111と、派生サービス生成部112と、DC生成部113と、NDC生成部114と、最適解探索部115とを有する。

40

【 0 0 4 4 】

なお、以下のトランザクション計画部110の説明において、サービス、原始サービスおよび派生サービスについては、特に断らない限り、トランザクション型のサービスを意味する。また、トランザクション型サービスの資源計画の生成においては、タイムスロットの長さはデータセンタ群の運営業者によって、データセンタ群全体で一意的な値が手動で決定されるものとする。

【 0 0 4 5 】

予測部111は、運用すべき全ての原始サービスの負荷変動の将来予測を行う処理部である。ここで、予測の対象となる期間を以後、予測期間と呼び、その長さは、1日、1ヶ月、1年などである。

50

【 0 0 4 6 】

派生サービス生成部 1 1 2 は、予測部 1 1 1 による負荷変動予測とオリジナル・センタの保有資源量に基づいて、派生サービスの発生時刻、発生個数、各派生サービスと原始サービスとの間の負荷分配比率を求める処理部である。

【 0 0 4 7 】

派生サービスの発生条件としては以下のようなものが例として挙げられる。図 5 - 1 および図 5 - 2 は、派生サービスの発生条件を説明するための説明図 (1) および (2) である。図 5 - 1 は、データセンタ内に配置された原始サービスのリソース使用率が 9 0 % を超えたときに派生サービスを発生させる場合を示している。

【 0 0 4 8 】

また、図 5 - 2 は、三つのデータセンタ X、Y および Z のリソース使用率の平均を算出し、各データセンタのリソース使用率が平均の 1 . 2 倍以上になったときに派生サービスを発生させる場合を示している。

【 0 0 4 9 】

実際には、派生サービス生成部 1 1 2 は、各データセンタ毎に原始サービスによるリソース使用率の標準偏差を計算し、その値が一定の閾値を超えた時に派生サービスを発生させるという発生条件を採用している。

【 0 0 5 0 】

D C 生成部 1 1 3 は、予測期間を一定周期のタイムスロットに分割し、各タイムスロット毎にタイムスロット内で発生する派生サービスをどのデータセンタに配置するかを定める配置組合せパターンの候補全てを生成する処理部である。

【 0 0 5 1 】

ここで、生成される配置組合せパターンの各候補のことを D C (Deployment Candidate) と呼び、タイムスロット (t) について生成した D C の集合に属する D C のことを、タイムスロット (t) に属する D C と呼ぶ。

【 0 0 5 2 】

N D C 生成部 1 1 4 は、各タイムスロットに属する D C から V D C (Violating D C) を除外した N D C (Normal D C) を作成する処理部である。ここで、V D C とは、複数の評価指標 m_1, \dots, m_k による評価結果 e_1, \dots, e_k が所定の条件に違反する D C であり、N D C とは、それ以外の D C である。

【 0 0 5 3 】

評価指標は、評価指標を用いて派生サービス群のデータセンタ群への配置組合せの最適化をすることにより、データセンタ群全体の資源利用効率が改善されるようなものであり、例えば、負荷の均等度 (図 6 - 1)、N D C の類似度 (図 6 - 2) などがある。なお、ここでは、評価値の値が小さいほど評価が良いものとする。

【 0 0 5 4 】

図 6 - 1 は、三つのデータセンタ X、Y および Z において負荷の配置が均等でなく悪い配置と負荷の配置が均等で良い配置とを示しており、配置が良いほど負荷の均等度の値は小さいものとする。

【 0 0 5 5 】

図 6 - 2 は、四つのデータセンタ W、X、Y および Z の間での、タイムスロット (K - 1) からタイムスロット (K) に移る際の派生サービスの移動およびタイムスロット (K) からタイムスロット (K + 1) に移る際の派生サービスの移動を示しており、ここでは、時間の経過に伴うタイムスロットの遷移時におけるサービス移動個数が多いほど N D C の類似度の値は大きいものとする。

【 0 0 5 6 】

最適解探索部 1 1 5 は、Pareto 最適性理論に基づいて最適な N D C 系列を作成する処理部である。ここで、N D C 系列とは、各タイムスロットから N D C を一つずつ選んでタイムスロット順に時系列に並べたものである。すなわち、最適な N D C 系列を作成することが最適な資源計画を作成することとなる。

10

20

30

40

50

【 0 0 5 7 】

具体的には、この最適解探索部 1 1 5 は、N D C 系列毎に、それを構成する各 N D C の k 個の評価値 e_1, \dots, e_k の各々を N D C 系列全体にわたって合計して系列評価値 E_1, \dots, E_k とする。

【 0 0 5 8 】

そして、 k 個の系列評価値 E_1, \dots, E_k の中で、いずれか一つ以上の系列評価値が最小の値となるものを全て選び出し、初期系列集合とする。この初期系列集合は、N D C 系列の集合を解空間とし、 k 個の系列評価値 E_1, \dots, E_k の各々を目的関数とした時の多目的最適化問題における Pareto 解集合に相当する。

【 0 0 5 9 】

そして、求めた初期系列集合の中から系列評価関数の値が最小となるような N D C 系列を選び出し、準最適 N D C 系列とする。そして、各タイムスロット毎に、準最適 N D C 系列を構成する各 N D C を、当該 N D C と同一のタイムスロットに属する他の全ての N D C と比較し、以下の入れ替え条件が成り立つ場合には、前者の N D C を後者の N D C と入れ替える操作を行う。

【 0 0 6 0 】

ここで、 $e_i(T_p, c)$ は準最適 N D C 系列内のタイムスロット T_p に属する N D C についての e_i の値であり、 $e_i(T_p, o)$ はタイムスロット T_p に属する N D C であって準最適 N D C 系列に含まれない N D C についての e_i の値と定義する。

【 0 0 6 1 】

そして、タイムスロット T_p に属する準最適 N D C 系列内の N D C の k 個の全ての評価値について、 $e_i(T_p, c) < e_i(T_p, o)$ ($1 \leq i \leq k$) となる場合は N D C の入れ替えを行う。つまり、N D C の入れ替えの結果、 k 個の評価値の組を構成する全ての評価値に関して評価値が改善される場合は N D C の入れ替えを行う。

【 0 0 6 2 】

また、次の場合にも N D C の入れ替えを行う。 k 個の評価値の組を構成する評価値の中で、一部の評価値が入れ替えにより改善され、その他の評価値が改悪される場合、改善される評価値の改善量に入れ替え制御関数を適用した値の重み付け平均が、改悪される評価値の改悪量に入れ替え制御関数を適用した値の重み付け平均より上回るならば入れ替えを行うというものである。

【 0 0 6 3 】

つまり、整数の集合 $I = \{i | 1 \leq i \leq k\}$ が二つの部分集合 I_1 と I_2 に分割でき、 $i \in I_1$ の場合に $e_i(T_p, c) < e_i(T_p, o)$ であり、 $i \in I_2$ の場合に $e_i(T_p, c) > e_i(T_p, o)$ であるとき、以下の式 (1)

【 数 1 】

$$\sum_{i \in I_2} \phi_i \Lambda_i(\Delta e_i(T_p)) / \sum_{i \in I_2} \phi_i > \sum_{i \in I_1} \phi_i \lambda_i(\delta e_i(T_p)) / \sum_{i \in I_1} \phi_i \quad \dots (1)$$

が成り立つならば、N D C 入れ替え操作を実行する。ここで、 ϕ_i は重み係数、 $\Delta e_i(T_p) = e_i(T_p, c) - e_i(T_p, o)$ 、 $\delta e_i(T_p) = e_i(T_p, o) - e_i(T_p, c)$ 、 Λ_i および λ_i は評価値 e_i に関する入れ替え制御関数である。

【 0 0 6 4 】

そして、N D C の入れ替えによって修正された準最適 N D C 系列を系列評価関数で評価し、N D C 入れ替え操作を行う以前の準最適 N D C 系列の系列評価関数で評価した値と比較し、改善率が一定の閾値を下回ったら準最適 N D C 系列をデータセンタ群の資源計画とし、そうでなければ、N D C 入れ替えによって修正された N D C 系列を新たな準最適 N D C 系列として N D C の入れ替えを繰り返す。

【 0 0 6 5 】

ここで、図 7 ~ 図 1 2 を用いてこの最適解探索部 1 1 5 の処理の具体例について説明す

10

20

30

40

50

る。図7は、最適解探索部115による最適化手順の概略を説明するための説明図である。なお、ここでは、評価指標として負荷の均等度とNDCの類似度を用いる。

【0066】

同図に示すように、この最適解探索部115は、最適化ステップ(1)として、NDCの類似度に基づいて暫定解の初期候補群を生成し、最適化ステップ(2)として、暫定解の初期候補群を負荷の均等度とNDCの類似度を用いて評価し、評価値が最小のものを暫定解とする。

【0067】

そして、最適化ステップ(3)として、暫定解に対するNDCの入れ替えによってより良い暫定解を生成し、最適化ステップ(4)として、収束条件を満たすか否かを判定し、収束条件を満たす場合には、暫定解を最終解として出力し、収束条件を満たさない場合には、最適化ステップ(3)に戻って収束条件を満たすまでNDCの入れ替えを繰り返す。

10

【0068】

図8は、最適化手順の説明に用いるNDCテーブルの概念を説明するための説明図である。同図に示すように、NDCテーブルとは、各タイムスロットに属するNDCを負荷均等度の良い順に並べたものである。以下、このNDCテーブルを用いて各最適化ステップを説明する。

【0069】

図9は、最適化ステップ(1)を説明するための説明図である。同図に示すように、この最適化ステップでは、隣接するタイムスロット間で「NDCの類似度=0」となるNDC同士を接続して類似NDC系列を生成する。例えば、図9では、タイムスロット(0)のNDC(1)を先頭とする類似NDC系列、タイムスロット(0)のNDC(2)を先頭とする類似NDC系列などの複数の類似NDC系列が得られる。

20

【0070】

なお、NDC同士を接続する場合に、「NDCの類似度=0」となるNDCがない場合には、NDCの類似度が最も小さいNDCを選択する。この結果、生成された類似NDC系列は、NDCの類似度を評価指標とし、この評価指標の値が最も小さいNDC系列であり、これらを暫定解の初期候補群とする。また、各タイムスロットで最小の負荷均等度を持つNDCから構成されるNDC系列は、系列全体としての負荷均等度の合計が最小のNDC系列であり、これも暫定解の初期候補群に含める。

30

【0071】

図10は、最適化ステップ(2)を説明するための説明図である。同図に示すように、この最適化ステップでは、負荷均等度およびNDCの類似度を評価指標として全ての類似NDC系列を評価し、系列評価値 E_c が最小の最良類似NDC系列を暫定解とする。

【0072】

図11-1は、最適化ステップ(3)を説明するための説明図である。同図に示すように、この最適化ステップでは、暫定解の各タイムスロットにおいて、暫定解に含まれるNDC以外のNDCが所定の条件を満たす場合に、NDCを入れ替える。

【0073】

ここで、所定の条件とは、図11-2に示すように、負荷均等度およびNDCの類似度を同時に改善するNDCがある場合には、そのタイムスロットでの評価値 E_i の値が最小のものを選んで入れ替えを行う。

40

【0074】

また、どちらか一方だけが改善する場合には、以下のようにする。負荷均等度だけが改善するときは、NDCの類似度の悪化量 S を所定の関数で評価した値 $BL(S)$ よりも、負荷均等度の改善量 L の方が大きいNDCの中で $L - BL(S)$ の値が最大のものを選び、類似度だけが改善するときは、類似度の改善量 S を所定の関数で評価した値 $BL(S)$ よりも、負荷均等度の悪化量 L の方が大きいNDCの中で $BL(S) - L$ の値が最大のものを選んで入れ替えを行う。

【0075】

50

図12は、最適化ステップ(4)を説明するための説明図である。同図に示すように、この最適化ステップでは、NDC入れ替え処理を施した最良類似NDC系列の系列評価値 E_c の改善率を調べ、改善率が例えば3%以内であれば収束条件を満たしたものととして最終解を出力し、収束条件を満たさない場合には、最適化ステップ(3)に戻って、入れ替え処理を繰り返す。

【0076】

このように、NDCの類似度に基づいてNDC系列の初期候補群を生成し、初期候補群から系列評価値 E_c が最小の最良類似NDC系列を暫定解とし、NDCの入れ替えによって暫定解を改善することによって、最適なNDC系列を求めることができる。なお、ここでは、負荷均等度およびNDCの類似度の二つの評価指標を用いたが、一般には、任意の個数の評価指標を用いることができる。また、収束条件に用いる改善率も3%以外の他の値を用いることができる。

10

【0077】

バッチ計画部120は、バッチ型の原始サービス集合について資源計画の生成を行う処理部であり、負荷量算出部121と、バッチ型派生サービス生成部122と、バッチ型DC生成部123と、バッチ型NDC生成部124と、バッチ型最適解探索部125とを有する。

【0078】

なお、以下のバッチ計画部120の説明では、ジョブの投入スケジュールは1日周期で繰り返され、予測期間の長さが1日のトランザクション型サービスについての資源計画が得られていると仮定して説明する。

20

【0079】

また、あるバッチ型サービスの1個のジョブを標準的な性能のサーバ1台で処理するのに必要な時間を「ターンアラウンド時間」と呼び、この時間の長さはバッチ型サービス毎に一意に定まる。

【0080】

また、標準的な性能のサーバ1台が捌くことの出来る負荷量を「単位サーバ性能」と呼ぶ。デッドライン時刻とターンアラウンド時間は各バッチ型原始サービス毎に一意の値が定まる。また、バッチ型サービスの個々のジョブの投入時刻と予測期間の開始時刻の時間差は当該バッチ型サービスのターンアラウンド時間の整数倍とする。

30

【0081】

負荷量算出部121は、将来の負荷量を算出する処理部である。具体的には、この負荷量算出部121は、以下の手順で将来の負荷量を算出する。まず、全てのバッチ型原始サービスのターンアラウンド時間の最小公倍数を求め、この値を新たなタイムスロットの長さとする。

【0082】

そして、各データセンタ毎、各計測周期毎に、当該計測周期において当該データセンタに存在する(トランザクション型サービスに割り当てられなかった)余剰資源量で捌くことの出来る負荷量を単位サーバ性能で除算し、余剰資源量を標準性能サーバ台数の表現に換算する(以後これを、余剰サーバ台数と呼ぶ)。

40

【0083】

そして、予測期間内における各バッチ型原始サービスの将来負荷変動を0で初期化する。その後、各バッチ型原始サービス毎に、各バッチ型サービスのジョブ投入スケジュールが、「時刻 t_1 に n_1 件のジョブを投入、 t_2 に n_2 件を投入、...、 t_k に n_k 件を投入」となっている時、 $1 \leq i \leq k$ である各整数 i について以下の(S1)~(S5)を実行する。ただし $t_1 > t_2 > \dots > t_k$ とする。

【0084】

(S1) 時刻 t_i から当該バッチ型原始サービスのデッドライン時刻までの期間(これをスケジュール期間と呼ぶ)を当該バッチ型原始サービスのターンアラウンド時間の長さを持つタイム・クォンタムに分割する。

50

(S 2) 各データセンタの各タイム・クオンタム毎に余剰サーバ台数が最小となる計測周期を求め、その時の余剰サーバ台数をそのタイム・クオンタムの余剰サーバ台数とする。

【 0 0 8 5 】

(S 3) 各タイム・クオンタムにおいて、余剰サーバ台数が最大であるデータセンタを検索し、当該データセンタをそのタイム・クオンタムにおける最大センタとし、各タイム・クオンタム毎に全データセンタにわたって余剰サーバ台数を合計し、その合計値を、当該タイム・クオンタムの R 値とする。

【 0 0 8 6 】

(S 4) R 値が最大となるタイム・クオンタムを検索し、当該タイム・クオンタムにおける最大センタの余剰サーバ台数の値を 1 だけ減算し、当該タイム・クオンタムについて最大センタと R 値を求めなおすとともに、当該バッチ型原始サービスの将来負荷変動の当該タイム・クオンタムに対応する部分の負荷量に、単位サーバ性能と等しい負荷量を加算する。

10

(S 5) n_i から 1 を減算し、 $n_i > 0$ なら (S 3) に戻る。

【 0 0 8 7 】

バッチ型派生サービス生成部 1 2 2 は、負荷量算出部 1 2 1 により算出された将来負荷変動を表す負荷量に基づいてバッチ型派生サービスを発生させる処理部である。具体的には、このバッチ型派生サービス生成部 1 2 2 は、各バッチ型原始サービスについて、将来負荷変動を表す負荷量に基づいて、当該バッチ型原始サービスからバッチ型派生サービスが発生する時刻を定める。

20

【 0 0 8 8 】

その際、当該バッチ型原始サービスから適切な個数のバッチ型派生サービスを発生させ、バッチ型原始サービスの負荷量をそれらのバッチ型派生サービスの間で適切に分配するものとする。そして、各タイムスロット毎に、そのタイムスロット内で発生するバッチ型派生サービスの集合を生成する。

【 0 0 8 9 】

バッチ型 DC 生成部 1 2 3 は、各タイムスロット毎に、そのタイムスロットで発生するバッチ型派生サービス集合のデータセンタ群への配置組合せパターン全ての候補 (DC) を生成する処理部である。

【 0 0 9 0 】

バッチ型 NDC 生成部 1 2 4 は、DC から VDC を除いた NDC を作成する処理部である。すなわち、このバッチ型 NDC 生成部 1 2 4 は、各タイムスロット毎に、評価指標 m_1, \dots, m_k に基づいて、当該タイムスロットに属する各 DC を評価し、各 DC についての評価値 e_1, \dots, e_k を計算する。

30

【 0 0 9 1 】

そして、各タイムスロット毎に、そのタイムスロットに属する各 DC を評価値 e_1, \dots, e_k に基づいて NDC と VDC に分類し、VDC を当該タイムスロットに属する DC の集合から除く。

【 0 0 9 2 】

バッチ型最適解探索部 1 2 5 は、バッチ型資源計画を作成する処理部である。すなわち、このバッチ型最適解探索部 1 2 5 は、NDC 系列の集合から、トランザクション型サービスの場合と同様の手法で Pareto 解集合 (初期系列集合) を生成し、初期系列集合の中から系列評価関数の値の最小となる NDC 系列を選び、準最適 NDC 系列とする。

40

【 0 0 9 3 】

そして、準最適 NDC 系列内の各タイムスロットに対応する NDC を NDC の入れ替え操作 (トランザクション型サービスの場合と同様) によって入れ替え、NDC 入れ替え操作の前後の準最適 NDC 系列の系列評価関数の値を比較し、改善率が 3 % 以下なら準最適 NDC 系列をバッチ型資源計画とし、そうでなければ NDC の入れ替えを繰り返す。

【 0 0 9 4 】

資源計画送信部 1 3 0 は、作成した資源計画をインターネット 2 0 を介して接続される

50

n 個のデータセンタ $10_1 \sim 10_n$ に送信する処理部である。各データセンタは、この資源計画送信部 130 により送信された資源計画にしたがってサービスの運用を行う。

【0095】

次に、図3に示したトランザクション計画部 110 の処理手順について説明する。なお、以下では、トランザクション型サービスの例としてXMLベースのWEBサービスを仮定し、原始サービスは全てWEBサービスであるとする。

【0096】

また、原始サービスの負荷変動予測の対象となった予測期間を表す時間軸は一定周期のタイムスロットに分割され、各タイムスロットは原始サービスの負荷を計測する周期である計測周期に分割される。そして、予測期間に含まれるタイムスロットの全数を N_t 、予測期間内の（時間軸上の時刻の早いほうから数えて） i 番目のタイムスロットを T_i と表す。

【0097】

また、時刻 t におけるサービスの負荷量は、時刻 t に当該サービスに到着する毎秒当りのユーザ・リクエスト数とする。ここで、データセンタの集合を $C=[c_i | 1 \leq i \leq N_c]$ 、原始サービスの集合を $P=[p_i | 1 \leq i \leq N_p]$ 、タイムスロット T_i の期間内において発生している全ての派生サービスの集合を $D(T_i)=[d(T_i, j) | 1 \leq j \leq n(T_i)]$ 、DCを評価するための k 個の評価指標の組を $[m_1, m_2, \dots, m_k]$ 、 k 個の評価指標でDCを評価した k 個の評価値の組を $[e_1, e_2, \dots, e_k]$ 、NDC系列を構成する各NDCの e_i をNDC系列全体にわたって合計したものを系列評価値 E_i 、タイムスロット T_i に属するNDCの集合を $(T_i)=[(T_i, j) | 1 \leq j \leq n(T_i)]$ 、NDC系列全体の集合を S とする。

【0098】

また、あるNDC系列 $s \in S$ を構成するNDCでタイムスロット T_i に属するNDCを $s(T_i, j)$ と書き、またNDC (T_i, j) について計算した評価値 e_i の値を $e_i((T_i, j))$ と書く。従って、NDC系列の一つのインスタンス $s \in S$ は、 $s = \langle s(T_1, j_1), s(T_2, j_2), s(T_3, j_3), \dots, s(T_{N_t-1}, j_{N_t-1}), s(T_{N_t}, j_{N_t}) \rangle$ と表わされる。

【0099】

また、データセンタ c_i が原始サービス p_j のオリジナル・センタであることを $p_j = c_i$ と表し、 c_i の全資源量を使って捌くことが出来る負荷量（毎秒リクエスト数）を $R(c_i)$ と表す。

【0100】

また、ここでは、DCを評価する評価指標として以下に定義する三つの指標 m_1, m_2, m_3 を採用した。ここで、 m_1 は派生サービスを評価対象のDCに従って配置した場合のデータセンタ群全体の負荷の均等度合いを表す。データセンタ群の各データセンタで用意しておくべきピーク時必要資源量を最小化するには、予測期間全体にわたって、全てのデータセンタの間で負荷量を均等に分散すればよいからである。

【0101】

また、 m_2 は派生サービス間通信量を表す。各原始サービスがXMLベースのWEBサービスである場合、原始サービス同士がSOAPプロトコルによって互いに通信しあうケースが考えられる。このような場合、それらの原始サービスから発生した派生サービス同士もSOAP通信を行うことになる。

【0102】

従って、あるDCに従って各派生サービスを配置先データセンタに配置した際、互いに頻繁に大量のデータをSOAP通信でやり取りするような複数の派生サービスの各々が異なるデータセンタに配置されることとなった場合、データセンタ間の通信量が大幅に増大する。

【0103】

そして、通信路の帯域幅が不足することによりSOAP通信の通信遅延が増大すると、結果として各データセンタに対して不必要な負荷をかけることになる。そのため、DCの評価指標として派生サービス間通信量 m_2 を採用した。

【0104】

10

20

30

40

50

また、 m_3 はNDCの類似度を表す。タイムスロット T_i に属するNDC (T_i, j_1) の類似度の定義は次の通りである。あるNDC系列が定められたとき、NDC系列を構成するNDCで、タイムスロット T_{i-1} に対応するNDCを (T_{i-1}, j_0) 、タイムスロット T_i に対応するNDCを (T_i, j_1) 、タイムスロット T_{i+1} に対応するNDCを (T_{i+1}, j_2) としたとき、時間の経過に従ってタイムスロットが T_{i-1} から T_i に遷移する際には、派生サービス集合のデータセンタ群への配置が (T_{i-1}, j_0) から (T_i, j_1) へと変更になるに伴って、派生サービス集合 $[D(T_{i-1}) \cup D(T_i)]$ に属する幾つかの派生サービスは、あるデータセンタから別のデータセンタへ移動させられる場合がある。

【0105】

また、タイムスロットが T_i から T_{i+1} に遷移する際にも同様に派生サービス配置が (T_i, j_1) から (T_{i+1}, j_2) へと変更になるに伴い、 $[D(T_i) \cup D(T_{i+1})]$ に属する幾つかの派生サービスがデータセンタ間を移動させられる場合がある。

10

【0106】

そこで、タイムスロット T_i に対応するNDC系列内のNDC (T_i, j_1) の類似度は、タイムスロットが T_{i-1} から T_i に遷移する時に上記の理由によりデータセンタ間を移動させられる派生サービスの個数とタイムスロットが T_i から T_{i+1} に遷移する時にデータセンタ間を移動させられる派生サービスの個数を合計したものとする。

【0107】

すなわち、あるNDCの類似度が大きいということは、そのNDCに従って派生サービスを配置すると、タイムスロットの遷移に伴う派生サービスの移動回数が多くなり、その分だけデータセンタ群に対してサービス移動のための制御負荷をかけることになる。

20

【0108】

その結果、各データセンタが本来なら処理する必要のない負荷まで大量に担わされることになってしまうため、NDCの類似度をNDCの評価指標の一つとして採用する必要がある。なお、 m_3 は m_1, m_2 と異なり、NDC単体では評価できず、NDC系列が指定されて初めて評価することができる。

【0109】

図13は、図3に示したトランザクション計画部110の処理手順を示すフローチャートである。同図に示すように、このトランザクション計画部110は、まず、予測部111が、 $1 \leq i \leq N_p$ となる全ての整数 i について、 p_i の予測期間全体にわたる負荷変動予測をARIMAモデル予測などの既存の予測手法を用いて予測する(ステップS101)。なお、 p_i の負荷変動予測は p_i の負荷量の時間関数として表され、時刻 t における p_i の負荷量は $p_i(t)$ と表すことにする。

30

【0110】

そして、派生サービス生成部112が、負荷変動予測およびオリジナル・センタの資源量に基づいて派生サービスを発生させる(ステップS102)。なお、負荷変動予測に基づく派生サービスの発生させ方の詳細については後述する。そして、 $1 \leq i \leq N_t$ となる各整数 i について T_i の期間中に発生すると予測された派生サービスの集合 $D(T_i)$ を生成する(ステップS103)。

【0111】

40

そして、DC生成部113が、 $1 \leq i \leq N_t$ となる各整数 i について、 $D(T_i)$ に属する各派生サービスのデータセンタ群 C への可能な配置組合せの全ての候補(DC)を木探索アルゴリズムによって生成し、タイムスロット T_i について生成した配置組合せの候補の集合を (T_i) とする(ステップS104)。

【0112】

そして、NDC生成部114が、 $1 \leq i \leq N_t$ となる各整数 i 、 $1 \leq j \leq N_j(T_i)$ となる各整数 j の全ての組み合わせについて (T_i, j) を評価指標 m_1, m_2 で評価し、評価値 e_1, e_2 を求める(ステップS105)。なお、 (T_i, j) について e_1 および e_2 を計算する手順の詳細は後述する。

【0113】

50

そして、 $1 \leq i \leq N_t$ となる各整数 i 、 $1 \leq j \leq (T_i)$ となる各整数 j の全ての組み合わせについて、 $e_1(\gamma(T_i, j)) > 0.7$ または $e_2(\gamma(T_i, j)) > 0.8$ ならば (T_i, j) をVDCとし、集合 (T_i) から取り除く。そうでなければ、 (T_i, j) をNDCとし集合 (T_i) に残す(ステップS106)。

【0114】

そして、最適解探索部116が、NDC系列集合 S を定義し(ステップS107)、NDC系列集合 S に含まれるNDC系列で k 個の系列評価値 $E_1(s), E_2(s), E_3(s)$ のいずれか一つ以上が最小値をとるNDC系列 s を全て探し出し、初期系列集合の要素とする(ステップS108)。

【0115】

ただし、

【数2】

$$E_i(s) = \sum_{h=1}^{N_t} e_i(\gamma(T_h, j_h))$$

とする。なお、NDC系列集合 S から初期系列集合(Pareto集合)を抽出する手順の詳細は後述する。

【0116】

そして、初期系列集合に属するNDC系列を系列評価関数 $\xi(s)$ で評価し、その評価値が最小となるNDC系列を準最適NDC系列とする(ステップS109)。ただし、 $k=3$ で

【数3】

$$\xi(s) = \sum_{i=1}^k K_i E_i(s)$$

とする。また、 $K_1=1.0$ 、 $K_2=1.0$ 、 $K_3=0.1$ とした。

【0117】

そして、 $1 \leq i \leq N_t$ となる各整数 i に対応するタイムスロット T_i についてNDCの入れ替え処理を行い(ステップS110)、NDCの入れ替えを行う前の準最適NDC系列を s' 、NDC入れ替え後の準最適NDC系列を s'' として改善率 $|s' - s''| / s'$ を算出する(ステップS111)。なお、NDCの入れ替え処理の詳細については後述する。

【0118】

そして、改善率が0.03以下であるか否かを判定し(ステップS112)、改善率が0.03以下である場合には、 s'' をデータセンタ群Cの資源計画として出力し、全ての処理を終える(ステップS113)。そうでなければ、 s'' を新たな s' としてステップS110に戻る。

【0119】

このように、派生サービス生成部112が生成した派生サービスに対して最適解探索部115がPareto最適化理論に基づいて最適NDC系列を生成することによって、データセンタ群に対する最適な資源計画を作成することができる。

【0120】

次に、派生サービス生成部112による派生サービス生成の詳細について説明する。派生サービス生成部112は、全ての計測周期について、各計測周期毎に次の(S1)~(S3)を行う。

【0121】

(S1)計測周期の中心時刻 t 、 $1 \leq i \leq N_c$ となる整数 i について、 p_j 、 c_j となる全ての j について

10

20

30

40

【数4】

$$\Theta_i = \sum_j \theta_j(t)$$

を求め、集合 $[\theta_i / R(c_i) | 1 \leq i \leq N_c]$ を標本分布と見なした時の標準偏差を計算し、当該標準偏差が0.1以下であれば(S3)へ、そうでなければ(S2)へ行く。

【0122】

(S2) $1 \leq i \leq N_c$ となる各整数*i*について p_i の時刻*t*における負荷量の一部を派生サービスとして分割する。ここで、派生サービスへの負荷量の分割のしかたは次の通りである。

【0123】

各原始サービスから派生サービスが $N_c - 1$ 個ずつ発生したと仮定し、各派生サービスの仮のセンタ配置を定める。仮定した派生サービスと原始サービスの負荷量の分配比率は次のようにして求める。

【0124】

線形計画法を用いて上記の分配比率を最適化し、次の条件が満たされるようにする。各派生サービスを上記の仮のセンタ配置で配置した際、各データセンタ $c_i (1 \leq i \leq N_c)$ について、そこで運用している原始サービスと派生サービスの負荷量の合計を g_i とした時の、標本集合 $[g_i / R(c_i) | 1 \leq i \leq N_c]$ の標準偏差が0.1以下になるようにする。

【0125】

もし整数計画法によって上記の条件を充足する分割比率が得られなかった場合は、各派生サービスの仮のセンタ配置をランダムに変更して、再び整数計画法による分割比率の最適化を行うことを繰り返すことにより、各原始サービスについて派生サービスに分割する負荷量の分割比率を求める。

【0126】

ただし、 $N_c > 5$ であった場合でも一つの原始サービスから同時に発生可能な派生サービスの個数の上限は4個とする。ある派生サービスについて上記の方法で決定した分配比率が0であった場合、その派生サービスは存在しないものとして扱う。

(S3) 次の計測周期に移り、新たな計測周期について(S1)を開始する。

【0127】

次に、 (T_i, j) について e_i を計算する手順の詳細について説明する。タイムスロット T_i 内に含まれる計測周期の中心時刻の集合を $[t_h^i | 1 \leq h \leq N_m]$ とする。すなわち、タイムスロット T_i の開始時点から*h*番目の計測周期の中心時刻を t_h^i とし、一つのタイムスロット内には N_m 個の計測周期が含まれるとする。

【0128】

データセンタ $c_r (1 \leq r \leq N_c)$ に配置されている全サービスの時刻*t*における負荷量合計を $g_r(t)$ とすると、 (T_i, j) に対する e_i の値は以下のようにして求まる。

【0129】

まず、 $1 \leq r \leq N_c$ である各整数*r*についてセンタ毎負荷量合計値の時系列 $[g_r(t_h^i) | 1 \leq h \leq N_m]$ を求め、 $1 \leq r \leq N_c$ である各整数*r*について時系列 $[g_r(t_h^i) / R(c_r) | 1 \leq h \leq N_m]$ の時間軸上の平均値

【数5】

$$\mu(r, i) = \frac{\sum_{h=1}^{N_m} g_r(t_h^i) / R(c_r)}{N_m}$$

を求める。

【0130】

そして、集合 $[\mu(0, i), \mu(1, i), \mu(2, i), \dots, \mu(N_c, i)]$ を標本分布と見なした時の標準偏差を計算し、この標準偏差を (T_i, j) に対する e_i の値とする。

【0131】

10

20

30

40

50

次に、 (T_i, j) について e_2 を計算する手順の詳細について説明する。まず、評価対象 NDC の属するタイムスロット T 内の各計測周期について各データセンタ毎に、サービス間の SOAP 通信による通信量を bps の単位で見積もる。

【 0 1 3 2 】

そして、データセンタ毎に、その見積もり通信量を、当該データセンタと外部ネットワークを繋ぐ通信路の帯域幅の値 (bps 単位) で除算した値 (以下、通信路占有率と呼ぶ) を求める。

【 0 1 3 3 】

そして、計測周期毎の通信路占有率をタイムスロット T 全体にわたって平均した値を各データセンタについて計算する。当該平均値が最大となるデータセンタでの当該平均値を e_2 の値とする。

10

【 0 1 3 4 】

なお、互いに通信しあう原始サービス同士の関係が、図 1 4 に示すようなものであるとき、データセンタ c_i における計測周期が t の時の、サービス間の SOAP 通信の通信量の見積もり手順の例を以下に示す。

【 0 1 3 5 】

ここで、 e_2 による評価対象となっている NDC において配置先データセンタが指定されている派生サービスの集合を $D=[d_k | 1 \leq k \leq N_d]$ 、 d_k の配置先データセンタが c_i であることを $d_k \in c_i$ 等と表す。また、サービス X の負荷量を毎秒リクエスト数で表した値を $L(X, t)$ と表す。また、原始サービス p_j が d_k の派生元の原始サービスであることを、 $p_j \in d_k$ 等と表す。

20

【 0 1 3 6 】

まず、図 1 4 の 1 段目と 2 段目の間の通信量の見積もりで、計測周期 t の時のデータセンタ c_i における値 $B_1(c_i, t)$ を求める。まず、 $B_1(c_i, t)$ を 0 で初期化し、以下の四つの場合に分けて考える。

【 0 1 3 7 】

(1) $p_1 \in c_i$ かつ $p_2 \in c_i$ の場合

$B_1(c_i, t) = B_1(c_i, t) + |L(p_1, t) \times \omega_1 - L(p_2, t)|$ とする。

【 0 1 3 8 】

(2) $p_1 \in c_i$ であり、 $p_2 \notin c_i$ でない場合

30

$(p_2 \in d_h) \cap (d_h \in c_i)$ となる全ての添字 h の集合を H とすると

【数 6】

$$B_1(c_i, t) = B_1(c_i, t) + |L(p_1, t) \times \omega_1 - \sum_{h \in H} L(d_h, t)|$$

【 0 1 3 9 】

(3) $p_2 \in c_i$ であり、 $p_1 \notin c_i$ でない場合

$(p_1 \in d_h) \cap (d_h \in c_i)$ となる全ての添字 h の集合を H とすると

【数 7】

$$B_1(c_i, t) = B_1(c_i, t) + |L(p_2, t) - \omega_1 \times \sum_{h \in H} L(d_h, t)|$$

40

【 0 1 4 0 】

(4) $p_1 \notin c_i$ でなく、 $p_2 \notin c_i$ でもない場合

$(p_1 \in d_h) \cap (d_h \in c_i)$ となる全ての添字 h の集合を H とし、 $(p_2 \in d_k) \cap (d_k \in c_i)$ となる全ての添字 k の集合を K とすると

【数 8】

$$B_1(c_i, t) = B_1(c_i, t) + \left| \sum_{k \in K} L(d_k, t) - \omega_1 \times \sum_{h \in H} L(d_h, t) \right|$$

【 0 1 4 1 】

その次に、図 1 4 の 2 段目と 3 段目の間の通信量の見積もりで、計測周期 t の時のデータセンタ c_i における値 $B_2(c_i, t)$ を求める。まず、 $B_1(c_i, t)$ を 0 で初期化し、以下の四つの

50

場合に分けて考える。

【 0 1 4 2 】

(1) $p_2 = c_i$ かつ $p_3 = c_i$ の場合

$B_2(c_i, t) = B_2(c_i, t) + |L(p_2, t) \times \omega_2 - L(p_3, t)|$ とする。

【 0 1 4 3 】

(2) $p_2 = c_i$ であり、 $p_3 \neq c_i$ でない場合

$(p_3 = d_h) (d_h = c_i)$ となる全ての添字 h の集合を H とすると

【 数 9 】

$$B_2(c_i, t) = B_2(c_i, t) + |L(p_2, t) \times \omega_2 - \sum_{h \in H} L(d_h, t)|$$

10

【 0 1 4 4 】

(3) $p_3 = c_i$ であり、 $p_2 \neq c_i$ でない場合

$(p_2 = d_h) (d_h = c_i)$ となる全ての添字 h の集合を H とすると

【 数 1 0 】

$$B_2(c_i, t) = B_2(c_i, t) + |L(p_3, t) - \omega_2 \times \sum_{h \in H} L(d_h, t)|$$

【 0 1 4 5 】

(4) $p_2 \neq c_i$ でなく、 $p_3 \neq c_i$ でもない場合

$(p_2 = d_h) (d_h = c_i)$ となる全ての添字 h の集合を H とし、 $(p_3 = d_k) (d_k = c_i)$ となる全ての添字 k の集合を K とすると

20

【 数 1 1 】

$$B_2(c_i, t) = B_2(c_i, t) + \left| \sum_{k \in K} L(d_k, t) - \omega_2 \times \sum_{h \in H} L(d_h, t) \right|$$

【 0 1 4 6 】

ここで、 $(c_i, t) = (c_i, t) + (c_i, t)$ とすることで、計測周期 t の時のデータセンタ c_i における見積み通信量 (c_i, t) が得られる。

【 0 1 4 7 】

次に、NDC 系列集合 S から初期系列集合(Pareto集合)を抽出する手順の詳細について説明する。まず、系列評価値 $E_1(s)$ を最小にするような $s \in S$ を求める。このような s は一般に複数個存在する。

30

【 0 1 4 8 】

具体的には、 $1 \leq i \leq N_t$ となる各整数 i について (T_i) に属する全ての NDC を配列に格納して、 e_1 の値をキーとしてその配列をソートし、 (T_i) 内で e_1 の値が最小となる NDC の集合 $b(T_i)$ を生成する。そして、直積集合 $b(T_0) \times b(T_1) \times \dots \times b(T_{N_t})$ に属する各 NDC 系列を系列評価値 $E_1(s)$ を最小にするような $s \in S$ とする。

【 0 1 4 9 】

その次に、系列評価値 $E_2(s)$ を最小にするような $s \in S$ を求める。すなわち、 $1 \leq i \leq N_t$ となる各整数 i について (T_i) に属する全ての NDC を配列に格納して、 e_2 の値をキーとしてその配列をソートし、 (T_i) 内で e_2 の値が最小となる NDC の集合 $b(T_i)$ を生成する。そして、直積集合 $b(T_0) \times b(T_1) \times \dots \times b(T_{N_t})$ に属する各 NDC 系列を系列評価値 $E_2(s)$ を最小にするような $s \in S$ とする。

40

【 0 1 5 0 】

その次に、系列評価値 $E_3(s)$ を最小にするような $s \in S$ を求める。ここで時間軸全体にわたる NDC の集合 Γ を以下のように定義する。

【 数 1 2 】

$$\Gamma = \bigcup_{i=1}^{N_t} \Gamma(T_i) = \{\gamma_k \mid 1 \leq k \leq N_A\}$$

【 0 1 5 1 】

そして、 $1 \leq k \leq N_A$ となる各整数 k と $1 \leq i \leq N_t$ となる各整数 i の全ての組合せについて以

50

下の処理を行う。集合 (T_i) 内の NDC を検索して、 (T_i, j) $(T_i), 1 \leq j \leq (T_i)$ で、 k と (T_i, j) を比較した時の相対類似度が最小となるような添字 j^k の値を求める。

【0152】

ここで、相対類似度とは、二つの NDC を比較した時、二つの NDC の指定する派生サービス配置の両方に共通に含まれている派生サービスの中で、一方の NDC で指定されている配置先データセンタと他方の NDC で指定されている配置先データセンタが異なるような派生サービスの個数を表す。

【0153】

そして、 $1 \leq k \leq N_A$ となる各整数 k について以下のように表される NDC 系列 $s(k) = \langle (T_0, j^k_0), (T_1, j^k_1), \dots, (T_{N_t}, j^k_{N_t}) \rangle$ を求める。

10

【0154】

最後に、集合 $\{s(k) | 1 \leq k \leq N_A\}$ に含まれる NDC 系列 $s(k)$ の内 $E_2(s(k))$ の値が最小となるものを選び出す。

【0155】

上記のようにして $E_1(s), E_2(s), E_3(s)$ の内のいずれかが最小値となる NDC 系列を全て S の中から選び出し、Pareto 集合 (初期系列集合) の要素とする。

【0156】

次に、NDC の入れ替え処理の詳細について説明する。NDC の入れ替え処理としては以下の (S1) ~ (S3) を行う。

(S1) 集合 (T_i) 内を探索して、 (T_i, j_1) と (T_i, j_2) (ただし $j_1 \neq j_2$ とする) の間に以下の関係が成り立つような j_2 を求める。

20

$$e_1((T_i, j_1)) > e_1((T_i, j_2))$$

$$e_2((T_i, j_1)) > e_2((T_i, j_2))$$

$$e_3((T_i, j_1)) > e_3((T_i, j_2))$$

【0157】

上記の条件を満たす j_2 が見つかった場合は、 (T_i, j_1) と (T_i, j_2) を入れ替えることにより、 (T_i, j_2) を NDC 系列内のタイムスロット T_i に対応する NDC とし、(S3) に進む。上記の条件を満たす j_2 が見つからないときには (S2) に進む。

【0158】

(S2) 集合 (T_i) 内を探索して、 (T_i, j_1) と (T_i, j_2) (ただし $j_1 \neq j_2$ とする) の間に以下の関係が成り立つような j_2 を求める。

30

【数13】

$$\sum_{h \in I_2} \phi_h \Lambda_h(\Delta e_h(T_i)) / \sum_{h \in I_2} \phi_h > \sum_{h \in I_1} \phi_h \lambda_h(\delta e_h(T_i)) / \sum_{h \in I_1} \phi_h$$

【0159】

ここで、 I_1, I_2 は整数集合 $\{h | 1 \leq h \leq 3\}$ を共通部分の無い 2 つに分割した部分集合で $h \in I_1$ の時に $e_h((T_i, j_1)) > e_h((T_i, j_2))$ であり、 $h \in I_2$ の時に $e_h((T_i, j_1)) > e_h((T_i, j_2))$ であり、また、 $\Delta e_h = e_h((T_i, j_1)) - e_h((T_i, j_2))$, $\delta e_h = e_h((T_i, j_2)) - e_h((T_i, j_1))$ である。

40

【0160】

上記の条件を満たす j_2 が見つかった場合は、 (T_i, j_1) と (T_i, j_2) を入れ替えることにより、 (T_i, j_2) を NDC 系列内のタイムスロット T_i に対応する NDC とする。上記条件を満たす j_2 が見つからない場合は、タイムスロット T_i については NDC の入れ替えを行わない。ここで、 $\lambda_1, \lambda_2, \lambda_3, \lambda_1, \lambda_2, \lambda_3$ の定義及び、 $\lambda_1, \lambda_2, \lambda_3$ の値は以下の通りとする。

【0161】

$$\lambda_1(x) = x \quad \lambda_2(x) = x \quad \lambda_3(x) = 0.1 \times \exp(1.3x)$$

$$\lambda_1(x) = x \quad \lambda_2(x) = x \quad \lambda_3(x) = 0.1 \times \exp(1.3x)$$

$$\lambda_1 = 1.0 \quad \lambda_2 = 1.0 \quad \lambda_3 = 1.0$$

50

ただし、 $\exp(\)$ は自然対数の底を基数とする指数関数である。

【0162】

(S3) $i = i+1$ として次のタイムスロットに移り、(S1)に戻る。

【0163】

次に、図3に示したバッチ計画部120の処理手順について説明する。バッチ型サービスの例としては、コンビニエンス・ストアの各店舗の売り上げ情報の集計処理のための計算サービスなどが考えられる。

【0164】

以下、ここでは、バッチ型原始サービスの集合を $P=[p_i | 1 \leq i \leq N_p]$ 、 p_i の各ジョブのターンアラウンド時間を τ_i 、データセンタ c_i において、 c_i に配置されている全サービスの負荷量の合計がピークとなる時の、 c_i の必要資源量を毎秒リクエスト数で表したものを $R_{max}(c_i)$ 、データセンタ c_i の資源の内、時刻 t においてWEBサービスに割り当てられずに残っている余剰計算資源量を $R(c_i, t)$ 、単位サーバ性能の値を f とし、他の記号の定義などはトランザクション計画部110の処理手順の説明に用いた時と同様とする。

【0165】

また、ここでは、 p_i の各ジョブの投入時刻と予測期間の開始時刻との時間差は τ_i の整数倍であるとし、 τ_i はトランザクション型サービスの負荷の計測周期の整数倍であるとする。

【0166】

また、バッチ型サービスに関する資源計画の生成では、NDCを評価する指標として、WEBサービスの資源計画生成で採用した評価指標の内、 m_1 と m_3 のみを採用する。何故なら、バッチ型サービス同士がSOAP通信を行うことは通常はありえないので、サービス間のSOAP通信量に関する評価指標である m_2 は考慮する意義がないからである。

【0167】

図15は、図3に示したバッチ計画部120の処理手順を示すフローチャートである。同図に示すように、このバッチ計画部120は、まず、負荷量算出部121が、 $\tau_1, \tau_2, \dots, \tau_{N_p}$ の最小公倍数を求め、その値をタイムスロットの長さとする(ステップS201)。

【0168】

そして、時間軸をタイムスロットに分割し、 $1 \leq i \leq N_c$ となる各整数 i について、 $(c_i, t) = R(c_i, t) / f$ とし、後述するジョブ・スケジューリング・アルゴリズムにより、各バッチ型原始サービス $p_i (1 \leq i \leq N_p)$ について、将来負荷変動を計画する(ステップS202)。なお、 p_i に関する将来負荷変動は時系列データとして生成され、将来時刻 t における p_i の負荷量を (i, t) と表す。

【0169】

そして、バッチ型派生サービス生成部122が負荷量算出部121が算出した負荷量に基づいてバッチ型派生サービスを発生させる(ステップS203)。なお、バッチ型派生サービス生成部122によるバッチ型派生サービス発生の詳細については後述する。

【0170】

そして、 $1 \leq i \leq N_t$ となる各整数 i について、 T_i の期間中に発生すると予測された派生サービスの集合 $D(T_i)$ を生成する(ステップS204)。

【0171】

そして、バッチ型DC生成部123が、 $1 \leq i \leq N_t$ となる各整数 i について、 $D(T_i)$ に属する各派生サービスのデータセンタ群 C への可能な配置組合せの全ての候補(DC)を木探索アルゴリズムによって生成し(ステップS205)、タイムスロット T_i について生成した配置組合せの候補の集合を (T_i) とする。

【0172】

そして、バッチ型NDC生成部124が、 $1 \leq i \leq N_t$ となる各整数 i 、 $1 \leq j \leq |D(T_i)|$ となる各整数 j の全ての組み合わせについて (T_i, j) を評価指標 m_1 で評価し、評価値 e_1 を求める(ステップS206)。

10

20

30

40

50

【 0 1 7 3 】

そして、 $1 \leq i \leq N_t$ となる各整数 i 、 $1 \leq j \leq (T_i)$ となる各整数 j の全ての組み合わせについて、 $e_i(\gamma(T_i, j)) > 0.7$ ならば (T_i, j) をVDCとし、集合 (T_i) から取り除く。そうでなければ、 (T_i, j) をNDCとし集合 (T_i) に残す(ステップS207)。

【 0 1 7 4 】

そして、バッチ型最適解探索部125が、NDC系列集合Sを作成し(ステップS208)、NDC系列集合Sに含まれるNDC系列で k 個の系列評価値 $E_1(s), E_3(s)$ のいずれか一つ以上が最小値をとるNDC系列 s を全て探し出し、初期系列集合の要素とする(ステップS209)。

【 0 1 7 5 】

ただし、

【数14】

$$E_i(s) = \sum_{h=1}^{N_t} e_i(\gamma(T_h, j_h))$$

とする。

【 0 1 7 6 】

そして、初期系列集合に属するNDC系列を系列評価関数 $E(s)$ で評価し、その評価値が最小となるNDC系列を準最適NDC系列とする(ステップS210)。ただし、 $E(s) = K_1 E_1(s) + K_2 E_2(s)$ とする。また、 $K_1=1.0$ 、 $K_3=0.1$ とした。

【 0 1 7 7 】

そして、 $1 \leq i \leq N_t$ となる各整数 i に対応するタイムスロット T_i について後述するNDCの入れ替え処理を行う(ステップS211)。そして、NDCの入れ替えを行う前の準最適NDC系列の $()$ での評価値と、NDC入れ替え後の準最適NDC系列の $()$ での評価値を用いて改善率を算出する(ステップS212)。

【 0 1 7 8 】

そして、改善率が3%以下であるか否かを判定し(ステップS213)、改善率が3%以下なら現在の準最適NDC系列を資源計画として出力し(ステップS214)、全ての処理を終える。そうでなければステップS211に戻る。

【 0 1 7 9 】

次に、バッチ型原始サービスの将来負荷変動を計画するためのジョブ・スケジューリング・アルゴリズムについて説明する。各バッチ型原始サービス $p_i (1 \leq i \leq N_p)$ のジョブ投入スケジュールとターンアラウンド時間、デッドライン時刻が図16の様であるとき、各バッチ型原始サービス $p_i (1 \leq i \leq N_p)$ について以下の(S1)~(S9)を実行する。

【 0 1 8 0 】

ただし、予測期間の開始時刻を t_0 とすると、任意の i と j について、 $t(i, j) - t_0 = TA(i) \times N$ 、 $TD(i) - t_0 = TA(i) \times N'$ (ただし N および N' は整数)の関係が有り、 $TA(i)$ は計測周期の長さの整数倍であるものとする。また、単位サーバ性能を f とする。

【 0 1 8 1 】

(S1) $j = (i)$ とする。また、予測期間内の各時刻 t において $p_i(t) = 0$ とする。
(S2) 予測期間内の各時刻 t において、トランザクション型およびバッチ型のいかなるサービスにも割り当てられていないデータセンタ c_k 内の余剰資源量を $R(c_k, t)$ とした時の $R(c_k, t) = R(c_k, t) / f$ を求める。

【 0 1 8 2 】

(S3) 時刻 $t(i, j)$ から時刻 $TD(i)$ までの時間帯を、 $TA(i)$ の長さのタイム・クォンタムに分割し、時刻 $t(i, j)$ から開始時刻の早い順に数えて h 番目のタイム・クォンタムを $Q(i, j, h)$ とする。また $N_Q(i, j) = (TD(i) - t(i, j)) / TA(i)$ とする。

【 0 1 8 3 】

(S4) $1 \leq h \leq N_Q(i, j)$ となる各整数 h と $1 \leq k \leq N_c$ となる各整数 k の組合せについて次の

10

20

30

40

50

処理を行う。タイム・クオンタム $Q(i, j, h)$ に含まれる複数の計測周期の中で計測周期の中心時刻 t における (c_k, t) の値が最小となる計測周期を求め、当該計測周期の中心時刻 t の時の (c_k, t) の値を $Q(i, j, h)$ における余剰資源量 $(Q(i, j, h), c_k)$ とする。

【 0 1 8 4 】

(S 5) 1 $h \in N_Q(i, j)$ となる各整数 h について、次の処理を行う。 $(Q(i, j, h), c_k)$ が最大となるような添字 k の値を求め、これを $k(i, j, h)$ とし、

【数 1 5】

$$\eta(i, j, h) = \sum_{k=1}^{N_c} \eta(Q(i, j, h), c_k)$$

10

を計算する。

【 0 1 8 5 】

(S 6) (i, j, h) の値を最大にするような添字 h の値を求め、 $(c_{k(i, j, h)}, t)$ から 1 を減算し、 $R(c_{k(i, j, h)}, t)$ から f を減算する。

(S 7) $t(i, j) + TA(i) \times h$ $t = t(i, j) + TA(i) \times (h+1)$ となる各時刻 t について $i(t) = i(t) + f$ とする。

【 0 1 8 6 】

(S 8) $n(i, j)$ から 1 を減算して $n(i, j) > 0$ なら (S 5) へ戻る。

(S 9) j が 1 なら処理を終了。そうでなければ $j = j - 1$ として (S 2) に戻る。

以上の手順実行の結果得られた $i(t)$ の値がバッチ型原始サービス p_i の将来時刻 t の時の将来負荷である。

20

【 0 1 8 7 】

次に、バッチ型派生サービス生成部 1 2 2 によるバッチ型派生サービス生成の詳細について説明する。バッチ型派生サービス生成部 1 2 2 は、 $1 \leq i \leq N_c$ となる各整数 i について、以下の (S 1) および (S 2) を実行する。ここで、時刻 t におけるデータセンタ c_i の WEB サービスによる負荷量の合計を $g(c_i, t)$ と表す。

【 0 1 8 8 】

(S 1) 予測期間内の各計測周期において、次の処理を実行する。当該計測周期の中心時刻 t において、 $p_j \leq c_i$ となる全ての添字 j の集合を J とした時に

【数 1 6】

$$R_{\max}(c_i) < \sum_{j \in J} \theta(j, t) + g(c_i, t)$$

30

ならば、 $p_h \leq c_i$ となる p_h の時刻 t での負荷量 $\theta'(h, t)$ を

【数 1 7】

$$\theta'(h, t) = \left\lfloor \frac{\theta(h, t) \times (1.0 - r(i, t))}{f} \right\rfloor \times f$$

とする。

【 0 1 8 9 】

ただし、

40

【数 1 8】

$$r(i, t) = \frac{\sum_{j \in J} \theta(j, t) + g(c_i, t) - R_{\max}(c_i)}{\sum_{j \in J} \theta(j, t)}$$

である。

【 0 1 9 0 】

(S 2) 予測期間内の各計測周期において次の処理を実行する。当該計測周期の中心時刻を t とした場合、 $p_j \leq c_i$ となる全ての添字 j の集合を J とした時に

50

【数19】

$$R_{\max}(c_i) < \sum_{j \in J} \theta(j, t) + g(c_i, t)$$

ならば、 p_h c_i となる各バッチ型原始サービス p_h から時刻 t において N_c-1 個のバッチ型派生サービスを発生させ、当該各バッチ型派生サービスの時刻 t での負荷量を以下のようにする。

【数20】

$$\left[\frac{\theta(h, t) \times r(i, t)}{(N_c - 1) \times f} \right] \times f$$

10

【0191】

次に、NDCの入れ替え処理の詳細について説明する。NDCの入れ替え処理としては、以下の(S1)~(S3)を行う。

(S1) 集合 (T_i) 内を探索して、 (T_i, j_1) と (T_i, j_2) (ただし $j_1 \neq j_2$ とする)の間に以下の関係が成り立つような j_2 を求める。

$$e_1(T_i, j_1) > e_1(T_i, j_2)$$

$$e_3(T_i, j_1) > e_3(T_i, j_2)$$

【0192】

上記の条件を満たす j_2 が見つかった場合は、 (T_i, j_1) と (T_i, j_2) を入れ替えることにより、 (T_i, j_2) をNDC系列内のタイムスロット T_i に対応するNDCとし、(S3)に進む。上記の条件を満たす j_2 が見つからないときには(S2)に進む。

20

【0193】

(S2) 集合 (T_i) 内を探索して、 (T_i, j_1) と (T_i, j_2) (ただし $j_1 \neq j_2$ とする)の間に以下の式(S2)または式(S3)のいずれか一方の関係が成り立つような j_2 を求める。

【数21】

$$\Delta e_1(T_i) > \lambda(\delta e_3(T_i))$$

$$\text{ただし } e_1(\gamma_\sigma(T_i, j_1)) > e_1(\gamma(T_i, j_2))$$

$$\text{かつ } e_3(\gamma_\sigma(T_i, j_1)) < e_3(\gamma(T_i, j_2))$$

$$\text{で、 } \Delta e_1(T_i) = e_1(\gamma_\sigma(T_i, j_1)) - e_1(\gamma(T_i, j_2))$$

$$\delta e_3(T_i) = e_3(\gamma(T_i, j_2)) - e_3(\gamma_\sigma(T_i, j_1))$$

… (2)

30

【数22】

$$\lambda(\Delta e_3(T_i)) > \delta e_1(T_i)$$

$$\text{ただし、 } e_1(\gamma(T_i, j_2)) > e_1(\gamma_\sigma(T_i, j_1))$$

$$\text{かつ、 } e_3(\gamma(T_i, j_2)) < e_3(\gamma_\sigma(T_i, j_1))$$

$$\text{で、 } \Delta e_3(T_i) = e_3(\gamma_\sigma(T_i, j_1)) - e_3(\gamma(T_i, j_2))$$

$$\delta e_1(T_i) = e_1(\gamma(T_i, j_2)) - e_1(\gamma_\sigma(T_i, j_1))$$

… (3)

40

【0194】

ただし、 $(x) = 0.1 \times \exp(1.3 \times x)$ とした。上記の条件を満たす j_2 が見つかった場合は、 (T_i, j_1) と (T_i, j_2) を入れ替えることにより、 (T_i, j_2) をNDC系列内のタイムスロット T_i に対応するNDCとする。上記条件を満たす j_2 が見つからない場合は、タイムスロット T_i についてはNDCの入れ替えを行わない。

(S3) $i = i+1$ として次のタイムスロットに移り、(S1)に戻る。

【0195】

次に、本実施例に係る資源計画作成装置100による資源計画作成のシミュレーション結果について説明する。図17は、本実施例に係る資源計画作成装置100による資源計

50

画作成のシミュレーション結果を示す図である。

【0196】

同図に示すように、データセンタXのConsumerトラフィックが朝と夜の高負荷時に3分割されて、データセンタYおよびデータセンタZでも処理が行われている。また、データセンタZの商店トラフィックが朝晩のピーク時に3分割されて、データセンタXおよびデータセンタYでも処理が行われている。

【0197】

上述してきたように、本実施例では、複数のデータセンタによって運用される複数のトランザクション型サービスに対して、トランザクション計画部110の予測部111が原始サービスの負荷予測を行い、派生サービス生成部112が予測部111による負荷予測およびデータセンタの計算資源量に基づいて派生サービスを生成し、最適解探索部115が派生サービス生成部112により生成された派生サービスに対する最適な資源計画を作成する。

10

【0198】

また、複数のデータセンタによって運用される複数のバッチ型サービスに対して、負荷量算出部121が原始サービスの負荷予測を行い、バッチ計画部120のバッチ型派生サービス生成部122が負荷量算出部121による負荷予測およびトランザクション型サービスの処理に使用されない計算資源量に基づいて派生サービスを生成し、バッチ型最適解探索部125が派生サービス生成部122により生成された派生サービスに対する最適な資源計画を作成する。

20

【0199】

従って、各データセンタに対して計算資源に応じた最適な負荷割り当てを行うことができるので、データセンタの過負荷を防ぐとともに、データセンタの計算資源を効率良く使用することができる。

【0200】

また、このような資源計画の生成処理により、各派生サービスと発生元となった原始サービスとの間の負荷分配比率、および各派生サービスの発生時刻とその時の配置先データセンタが定まるので、データセンタ毎に予測期間内における原始サービス負荷と派生サービス負荷の合計量の時間変動推移が得られる。

【0201】

これを基に、各データセンタ毎に、負荷の変動に応じて派生サービスの発生を適切に行った場合の、当該データセンタのピーク時負荷量を見積もることができ、この見積もり負荷量は、派生サービスの発生を全く行わない場合や、各派生サービスの配置が不適切である場合の各データセンタのピーク時負荷量より小さい値とすることができる。

30

【0202】

また、本実施例では、資源計画作成装置100だけで負荷予測、最適資源計画の作成および資源計画のデータセンタへの送信まで行う場合について説明したが、本発明はこれに限定されるものではなく、例えば負荷予測部111を他の装置で行う場合にも同様に適用することができる。

【0203】

また、本実施例では、資源計画作成装置について説明したが、この資源計画作成装置が有する構成をソフトウェアによって実現することで、同様の機能を有する資源計画作成プログラムを得ることができる。そこで、この資源計画作成プログラムを実行するコンピュータシステムについて説明する。

40

【0204】

図18は、本実施例に係る資源計画作成プログラムを実行するコンピュータシステムを示す図である。同図に示すように、このコンピュータシステム200は、本体部201と、本体部201からの指示により表示画面202aに情報を表示するディスプレイ202と、このコンピュータシステム200に種々の情報を入力するためのキーボード203と、ディスプレイ202の表示画面202a上の任意の位置を指定するマウス204と、L

50

AN206または広域エリアネットワーク(WAN)に接続するLANインタフェースと、公衆回線207に接続するモデムとを有する。ここで、LAN206は、他のコンピュータシステム(PC)211、サーバ212、プリンタ213などとコンピュータシステム200とを接続している。

【0205】

また、図19は、図18に示した本体部201の構成を示す機能ブロック図である。同図に示すように、この本体部201は、CPU221と、RAM222と、ROM223と、ハードディスクドライブ(HDD)224と、CD-ROMドライブ225と、FDドライブ226と、I/Oインタフェース227と、LANインタフェース228と、モデム229とを有する。

10

【0206】

そして、このコンピュータシステム200において実行される資源計画作成プログラムは、フロッピディスク(FD)208、CD-ROM209、DVDディスク、光磁気ディスク、ICカードなどの可搬型記憶媒体に記憶され、これらの記憶媒体から読み出されてコンピュータシステム200にインストールされる。

【0207】

あるいは、この資源計画作成プログラムは、LANインタフェース228を介して接続されたサーバ212のデータベース、他のコンピュータシステム(PC)211のデータベースなどに記憶され、これらのデータベースから読み出されてコンピュータシステム200にインストールされる。

20

【0208】

そして、インストールされた資源計画作成プログラムは、HDD224に記憶され、RAM222、ROM223などを利用してCPU221により実行される。

【0209】

また、本実施例では、データセンタにサービスを配置する場合について説明したが、本発明はこれに限定されるものではなく、他の計算機クラスタにサービスを配置する場合にも同様に適用することができる。

【0210】

例えば、各データセンタの内部を構成する計算資源の集合が複数の区画に分割され、各区画毎に独立した運用管理がなされている場合である。この場合、各区画毎に、どの派生サービスをいつ何時配置するかについて最適化したり、どの区画で恒常的にどの原始サービスを運用するか定め、区画毎に、そこで運用する原始サービスから派生サービスを発生させる将来時点を決定するといったことも可能である。

30

【0211】

(付記1)複数の計算機から構成される計算機クラスタが複数接続されて構成される計算機クラスタシステムが運用する複数のサービスを該複数の計算機クラスタへ配置する資源計画作成する資源計画作成プログラムであって、

運用される計算機クラスタが定められたサービスである原始サービスから該原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用の一部が移される派生サービスを生成する派生サービス生成手順と、

40

前記派生サービス生成手順により生成された派生サービス群を前記複数の計算機クラスタへ配置する資源計画を所定の評価指標に基づいて作成する計画作成手順と、

をコンピュータに実行させることを特徴とする資源計画作成プログラム。

【0212】

(付記2)前記派生サービス生成手順は、前記負荷予測の予測期間を所定の時間間隔で複数のタイムスロットに分割したタイムスロット毎に派生サービスを生成し、

前記計画作成手順は、前記タイムスロット毎の派生サービスの配置計画であるタイムスロット配置計画を該タイムスロットの時刻順に並べた配置計画系列を前記資源計画として作成することを特徴とする付記1に記載の資源計画作成プログラム。

【0213】

50

(付記3) 前記計画作成手順は、派生サービス群を前記複数の計算機クラスタへ最適に配置する資源計画をPareto最適性理論に基づいて作成することを特徴とする付記2に記載の資源計画作成プログラム。

【0214】

(付記4) 前記計画作成手順は、計算機クラスタ間の負荷均等度と、時間軸上で隣接するタイムスロット間の前記タイムスロット配置計画の類似度とを評価指標として用いることを特徴とする付記3に記載の資源計画作成プログラム。

【0215】

(付記5) 前記計画作成手順は、前記類似度と負荷均等度を評価指標とした場合に評価が最も良い配置計画系列を初期系列集合とし、該初期系列集合に含まれる配置計画系列に対して前記負荷均等度の評価がより良くなるような補正を行うことによって前記資源計画を作成することを特徴とする付記4に記載の資源計画作成プログラム。

10

【0216】

(付記6) 前記計画作成手順により作成された資源計画を各計算機クラスタに送信する資源計画送信手順をさらにコンピュータに実行させることを特徴とする付記1～5のいずれか一つに記載の資源計画作成プログラム。

【0217】

(付記7) 前記計算機クラスタはデータセンタであり、前記計算機クラスタシステムは複数のデータセンタがネットワークを介して接続されたシステムであることを特徴とする付記1～6のいずれか一つに記載の資源計画作成プログラム。

20

【0218】

(付記8) 前記派生サービス生成手順は、

指定されたデッドライン時刻までに、サービスを構成する全てのジョブが処理完了されていれば、個々のジョブはいつ何時処理されても構わないことを特徴とするような非リアルタイム性のバッチ型サービスが計算機クラスタの集合で運用されている場合に、リアルタイム性を要求されるトランザクション型のサービスに対して計算機クラスタの資源を優先的に割り当て、残った余剰資源をバッチ型サービスに割り当てるために、トランザクション型サービスに対して資源計画の生成を行い、その後、割り当て対象とならなかった余剰資源についてバッチ型サービスの資源計画の生成を行う事を特徴とする付記1に記載の資源計画作成プログラム。

30

【0219】

(付記9) 前記計画作成手順は、各バッチ型原始サービスを構成する全てのジョブを、その投入時刻がデッドライン時刻に近いものから順に実行計画の対象とし、いかなるサービスによっても使用されていない余剰計算資源の量の計算機クラスタ集合全体にわたる総和が最大となる期間に当該スケジューリング対象となったジョブの実行計画を生成することを特徴とする付記1に記載の資源計画作成プログラム。

【0220】

(付記10) 複数の計算機から構成される計算機クラスタが複数接続されて構成される計算機クラスタシステムが運用する複数のサービスを該複数の計算機クラスタへ配置する資源計画を作成する資源計画作成プログラムを記録したコンピュータ読み取り可能な記録媒体であって、

40

運用される計算機クラスタが定められたサービスである原始サービスから該原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用が移される派生サービスを生成する派生サービス生成手順と、

前記派生サービス生成手順により生成された派生サービスを前記複数の計算機クラスタへ配置する資源計画を所定の評価指標に基づいて作成する計画作成手順と、

をコンピュータに実行させる資源計画作成プログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

【0221】

(付記11) 複数の計算機から構成される計算機クラスタが複数接続されて構成される計

50

算機クラスタシステムが運用する複数のサービスを該複数の計算機クラスタへ配置する資源計画を作成する資源計画作成方法であって、

運用される計算機クラスタが定められたサービスである原始サービスから該原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用の一部が移される派生サービスを生成する派生サービス生成工程と、

前記派生サービス生成工程により生成された派生サービス群を前記複数の計算機クラスタへ最適に配置する資源計画を所定の評価指標に基づいて作成する計画作成工程と、

を含んだことを特徴とする資源計画作成方法。

【0222】

(付記12) 複数の計算機から構成される計算機クラスタが複数接続されて構成される計算機クラスタシステムが運用する複数のサービスを該複数の計算機クラスタへ配置する資源計画を作成する資源計画作成装置であって、

運用される計算機クラスタが定められたサービスである原始サービスから該原始サービスの負荷予測に基づいて各計算機クラスタから他の計算機クラスタへ運用の一部が移される派生サービスを生成する派生サービス生成手段と、

前記派生サービス生成手段により生成された派生サービス群を前記複数の計算機クラスタへ最適に配置する資源計画を所定の評価指標に基づいて作成する計画作成手段と、

を備えたことを特徴とする資源計画作成装置。

【産業上の利用可能性】

【0223】

以上のように、本発明に係る資源計画作成プログラムは、複数のデータセンタへのサービス割り当てに有用であり、特に、過負荷をなくし、計算資源を有効利用する必要があるデータセンタに適している。

【図面の簡単な説明】

【0224】

【図1】本実施例に係る資源計画の概念を説明するための説明図である。

【図2】派生サービス発生と配置最適化を説明するための説明図である。

【図3】タイムスロット毎の派生サービスの配置最適化を説明するための説明図である。

【図4】本実施例に係る資源計画作成装置の構成を示す機能ブロック図である。

【図5-1】派生サービスの発生条件を説明するための説明図(1)である。

【図5-2】派生サービスの発生条件を説明するための説明図(2)である。

【図6-1】負荷の均等度合を説明するための説明図である。

【図6-2】NDCの類似度を説明するための説明図である。

【図7】最適化手順概略を説明するための説明図である。

【図8】NDCテーブルの概念を説明するための説明図である。

【図9】最適化ステップ(1)を説明するための説明図である。

【図10】最適化ステップ(2)を説明するための説明図である。

【図11-1】最適化ステップ(3)を説明するための説明図である。

【図11-2】入れ替え条件を示す図である。

【図12】最適化ステップ(4)を説明するための説明図である。

【図13】図3に示したトランザクション計画部の処理手順を示すフローチャートである。

【図14】互いにSOAP通信するサービス集合の例を示す図である。

【図15】図3に示したバッチ計画部の処理手順を示すフローチャートである。

【図16】バッチ型原始サービスのジョブ投入スケジュールの例を示す図である。

【図17】本実施例に係る資源計画作成装置100による資源計画作成のシミュレーション結果を示す図である。

【図18】本実施例に係る資源計画作成プログラムを実行するコンピュータシステムを示す図である。

【図19】図18に示した本体部の構成を示す機能ブロック図である。

10

20

30

40

50

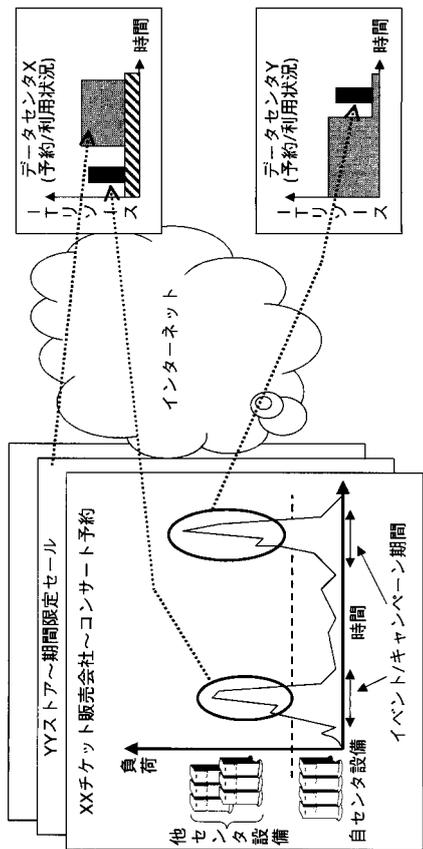
【符号の説明】

【0225】

1 0 ₁ ~ 1 0 _n	データセンタ	
2 0	インターネット	
1 0 0	資源計画作成装置	
1 1 0	トランザクション計画部	
1 1 1	予測部	
1 1 2	派生サービス生成部	
1 1 3	D C生成部	
1 1 4	N D C生成部	10
1 1 5	最適解探索部	
1 2 0	バッチ計画部	
1 2 1	負荷量算出部	
1 2 2	バッチ型派生サービス生成部	
1 2 3	バッチ型D C生成部	
1 2 4	バッチ型N D C生成部	
1 2 5	バッチ型最適解探索部	
1 3 0	資源計画送信部	
2 0 0 , 2 1 1	コンピュータシステム	
2 0 1	本体部	20
2 0 2	ディスプレイ	
2 0 2 a	表示画面	
2 0 3	キーボード	
2 0 4	マウス	
2 0 6	L A N	
2 0 7	公衆回線	
2 0 8	フロッピーディスク	
2 0 9	C D - R O M	
2 1 2	サーバ	
2 1 3	プリンタ	30
2 2 1	C P U	
2 2 2	R A M	
2 2 3	R O M	
2 2 4	ハードディスクドライブ	
2 2 5	C D - R O Mドライブ	
2 2 6	フロッピーディスクドライブ	
2 2 7	I / Oインタフェース	
2 2 8	L A Nインタフェース	
2 2 9	モデム	

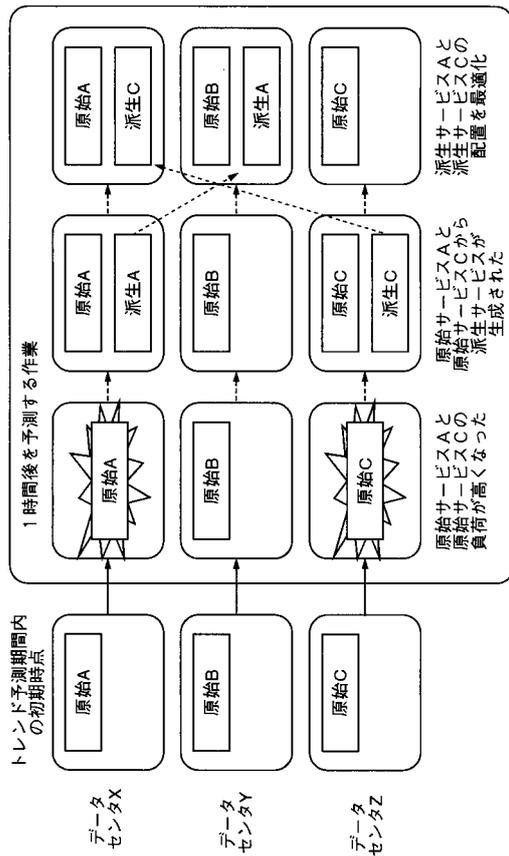
【図 1】

本実施例に係る資源計画の概念を説明するための説明図



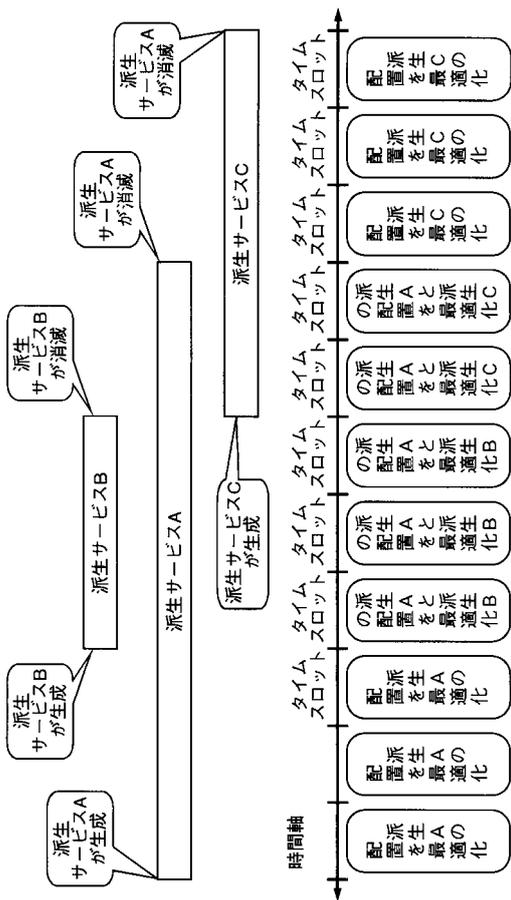
【図 2】

派生サービス発生と配置最適化を説明するための説明図



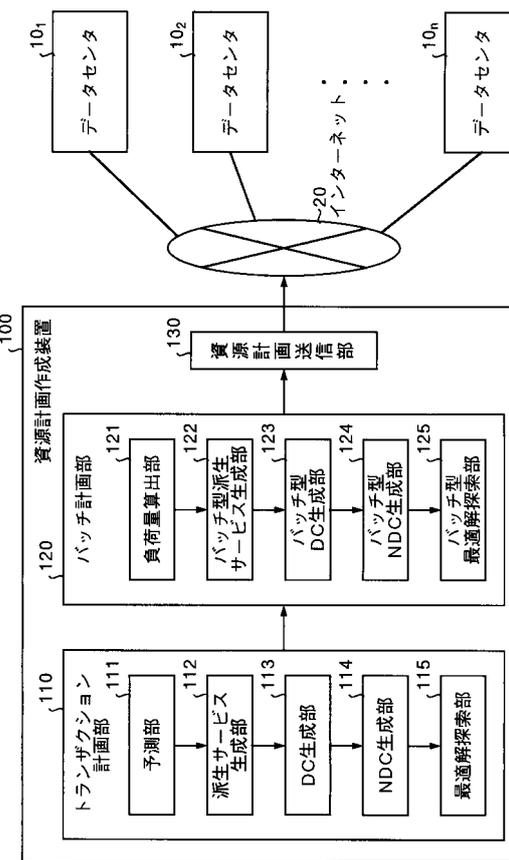
【図 3】

タイムスロット毎の派生サービスの配置最適化を説明するための説明図



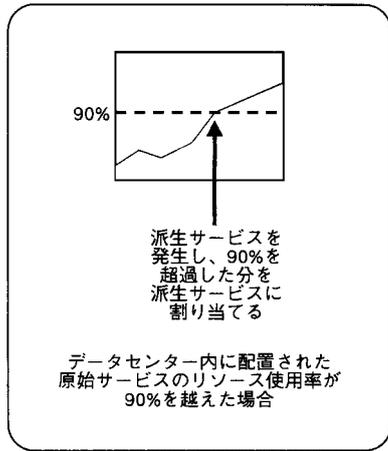
【図 4】

本実施例に係る資源計画作成装置の構成を示す機能ブロック図



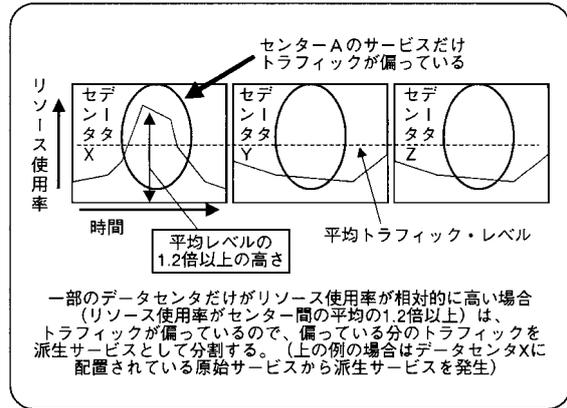
【 図 5 - 1 】

派生サービスの発生条件を説明するための説明図 (1)



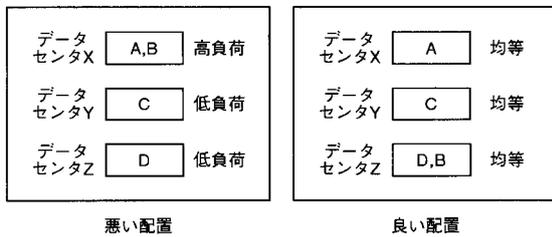
【 図 5 - 2 】

派生サービスの発生条件を説明するための説明図 (2)



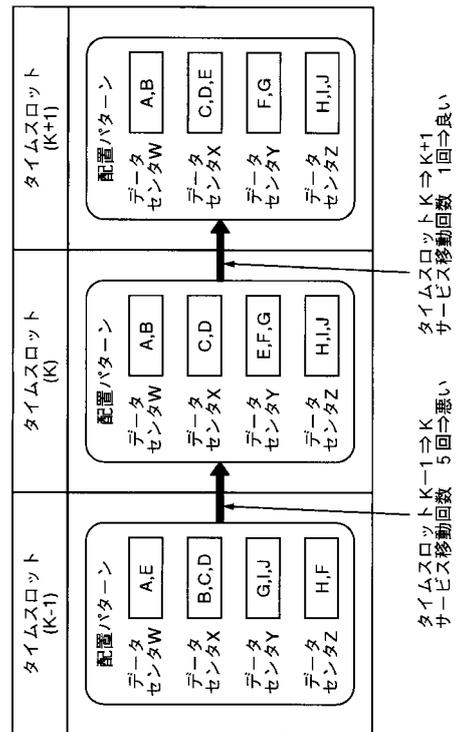
【 図 6 - 1 】

負荷の均等度合いを説明するための説明図

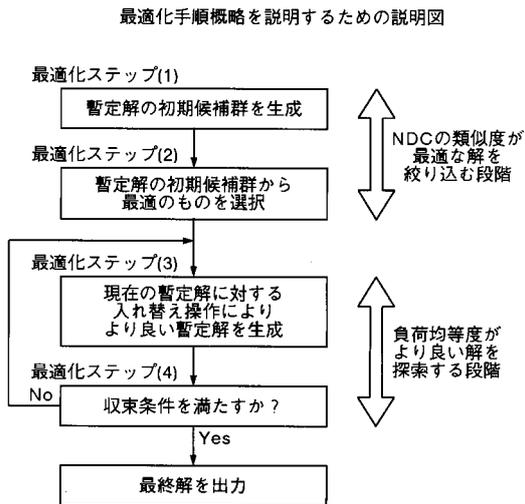


【 図 6 - 2 】

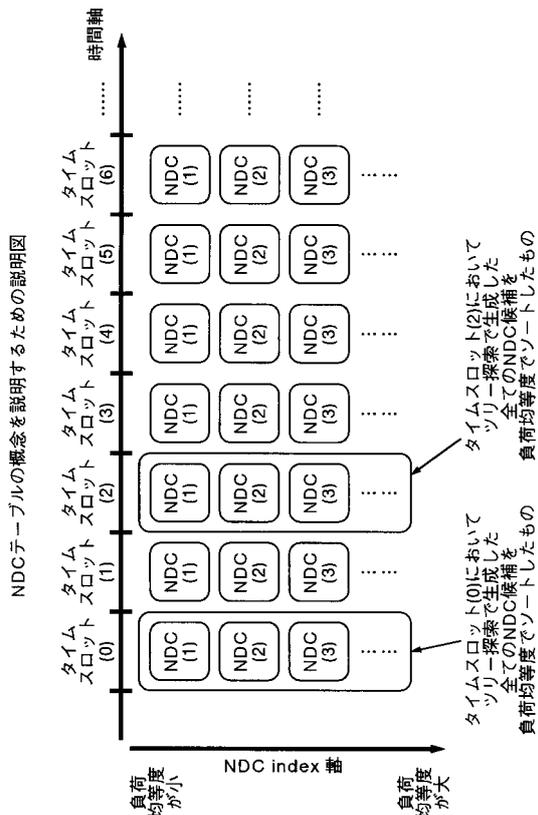
NDCの類似度を説明するための説明図



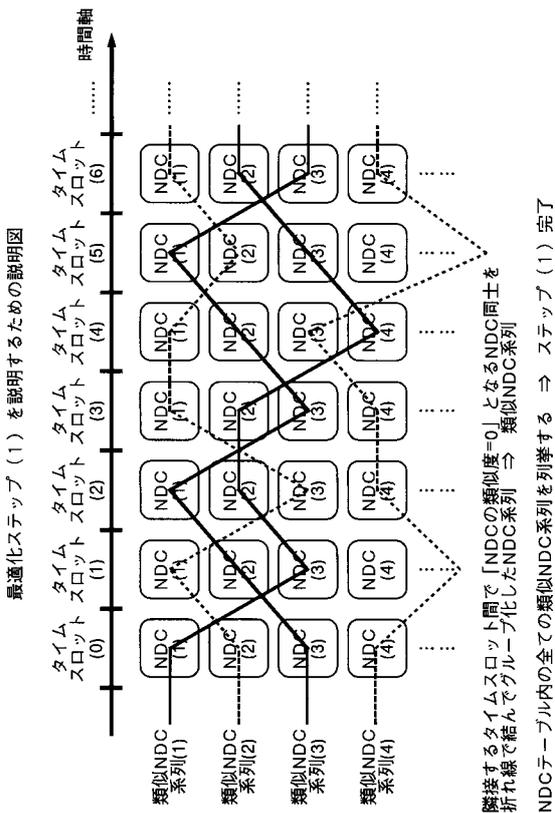
【 図 7 】



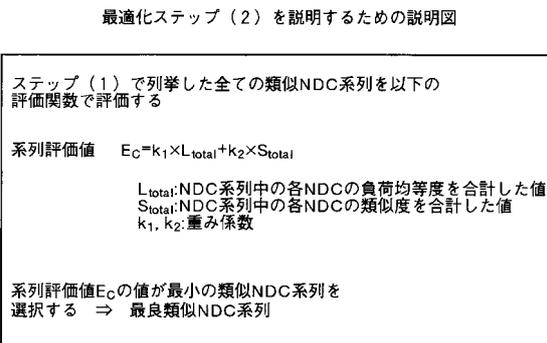
【 図 8 】



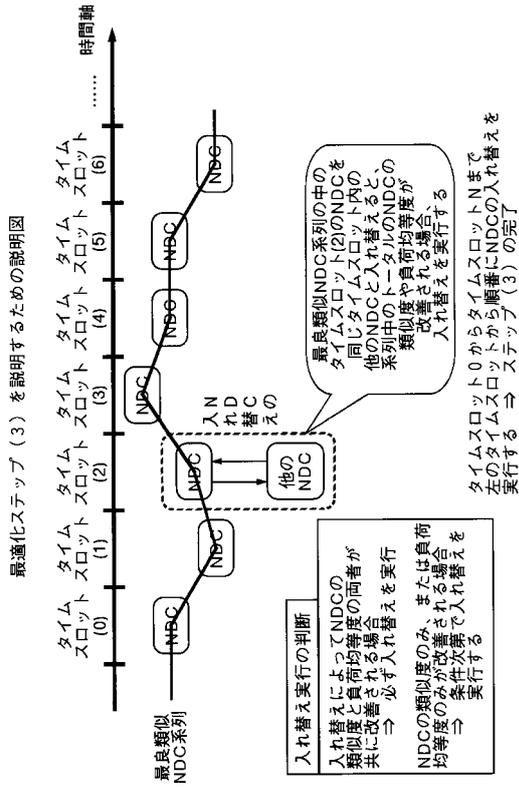
【 図 9 】



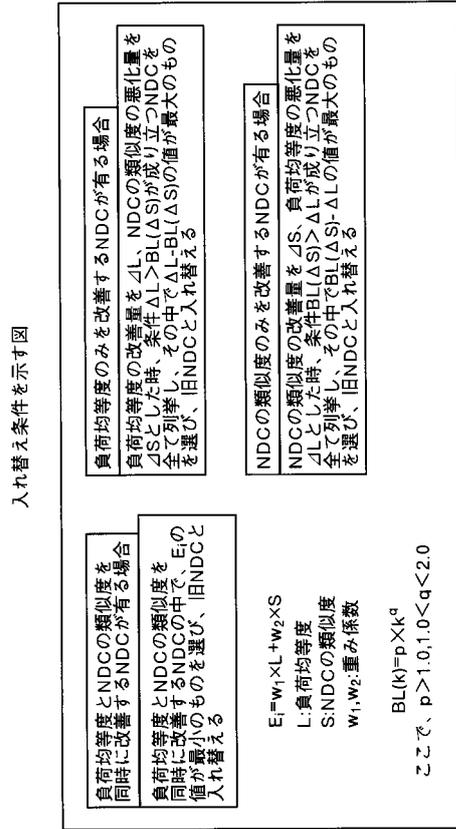
【 図 10 】



【図 1 1 - 1】

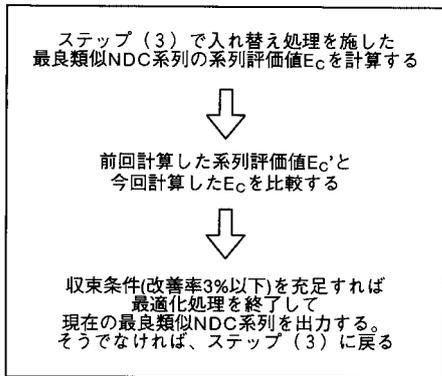


【図 1 1 - 2】



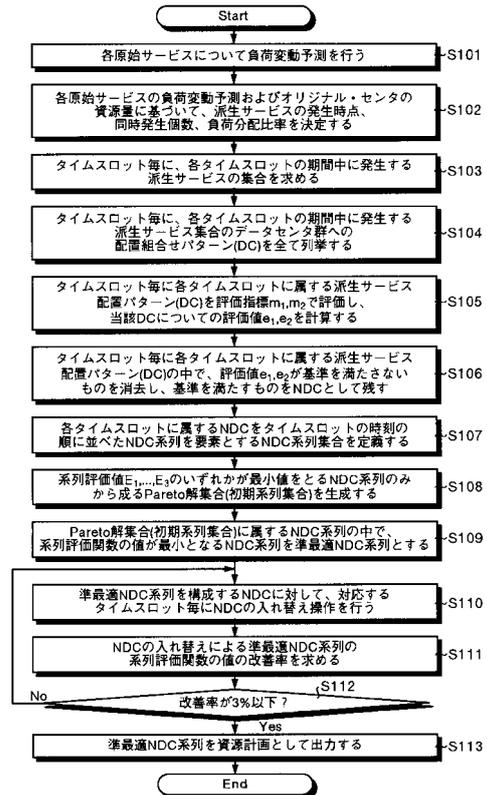
【図 1 2】

最適化ステップ (4) を説明するための説明図

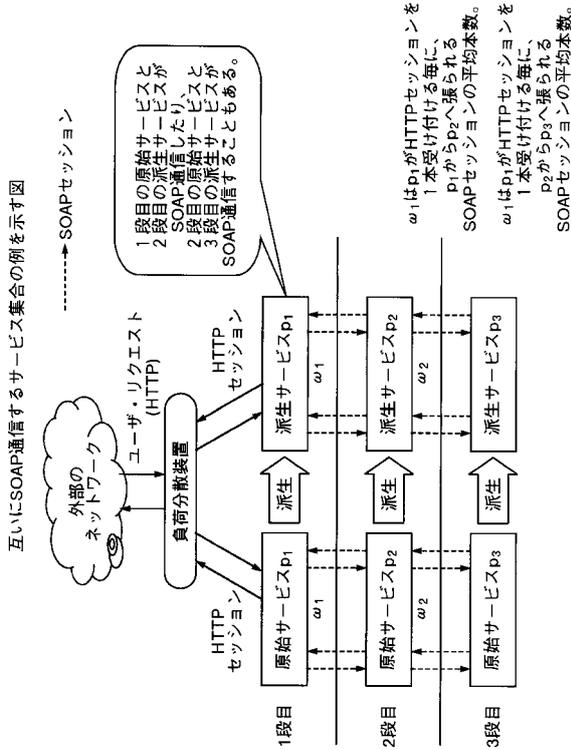


【図 1 3】

図 3 に示したトランザクション計画部の処理手順を示すフローチャート

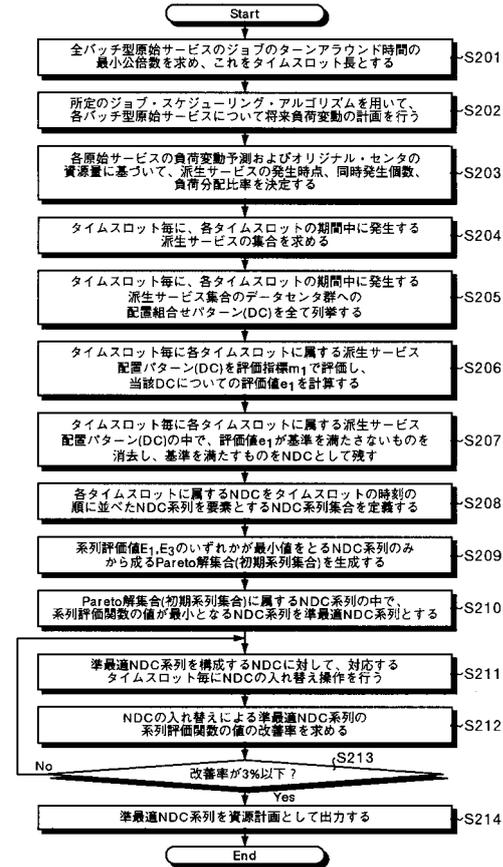


【 図 1 4 】



【 図 1 5 】

図 3 に示したバッチ計画部の処理手順を示すフローチャート



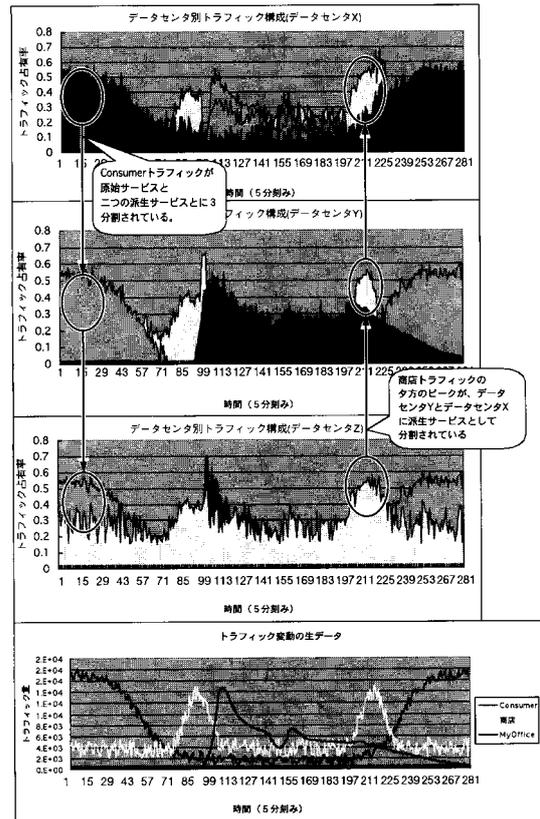
【 図 1 6 】

バッチ型原始サービスのジョブ投入スケジュールの例を示す図

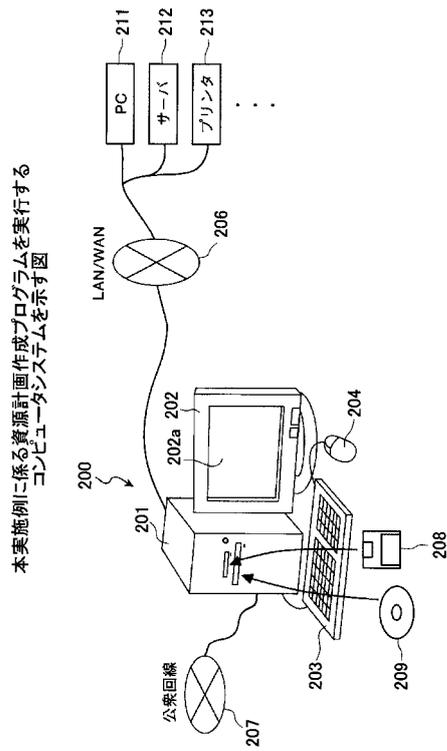
ターンアラウンド時刻 デッドライン 時刻	原始サービス p_1		原始サービス p_2		原始サービス p_p	
	時刻	投入件数	時刻	投入件数	時刻	投入件数
$t(1,1)$	$n(1,1)$	$t(2,1)$	$n(2,1)$	$t(N_p,1)$	$n(N_p,1)$	
$t(1,2)$	$n(1,2)$	$t(2,2)$	$n(2,2)$	$t(N_p,2)$	$n(N_p,2)$	
$t(1,3)$	$n(1,3)$	$t(2,3)$	$n(2,3)$	$t(N_p,3)$	$n(N_p,3)$	
$t(1,\rho(1)-1)$	$n(1,\rho(1)-1)$	$t(2,\rho(2)-1)$	$n(2,\rho(2)-1)$	$t(N_p,\rho(N_p)-1)$	$n(N_p,\rho(N_p)-1)$	
$t(1,\rho(1))$	$n(1,\rho(1))$	$t(2,\rho(2))$	$n(2,\rho(2))$	$t(N_p,\rho(N_p))$	$n(N_p,\rho(N_p))$	

【 図 1 7 】

本実施例に係る資源計画作成装置による資源計画作成のシミュレーション結果を示す図

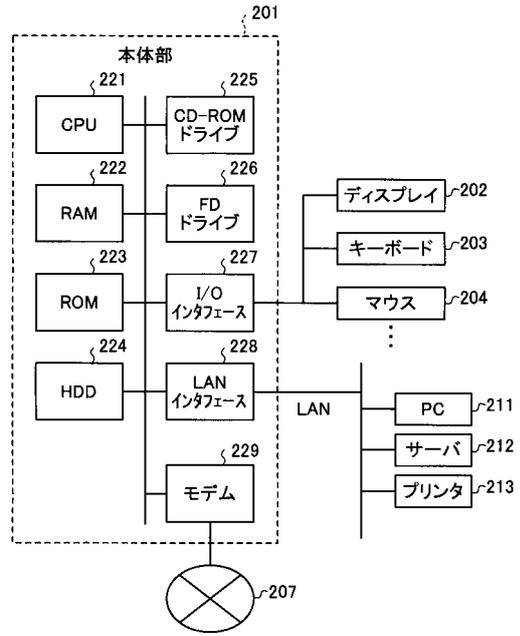


【図18】



【図19】

図18に示した本体部の構成を示す図



フロントページの続き

審査官 殿川 雅也

- (56)参考文献 特開平10-283211(JP,A)
特開2002-318791(JP,A)
特開2002-342297(JP,A)
特開2002-366373(JP,A)
特開2003-296149(JP,A)
EYMANN, T., et al., Decentralized vs. Centralized Economic Coordination of Resource Allocation in Grids, LECTURE NOTES IN COMPUTER SCIENCE, GRID COMPUTING, Springer-Verlag Berlin Heidelberg, 2004年3月6日, VOLUME 2970/2004, pp. 9 - 16
GROSU, Daniel, et al., Load Balancing in Distributed Systems: An Approach Using Cooperative Games, Parallel and Distributed Processing Symposium., Proceedings International, IPDPS'02, IEEE, 2002年4月19日, pp. 52 - 61
BUY YA, R., et al., A Case for Economy Grid Architecture for Service Oriented Grid Computing, 10th IEEE International Heterogeneous Computing Workshop (HCW), IEEE, 2001年, pp. 1 - 15
EYMANN, T., et al., The Catalaxy as a new Paradigm for the Design of Information Systems, Proceedings of the World Computer Congress of the International Federation for Information Processing, 2000, 2000年, pp. 1 - 8

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54
G06F 15/16 - 15/177