

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5244717号  
(P5244717)

(45) 発行日 平成25年7月24日 (2013. 7. 24)

(24) 登録日 平成25年4月12日 (2013. 4. 12)

(51) Int. Cl. F I  
**HO 4 L 12/713 (2013. 01)** HO 4 L 12/56 G  
**GO 6 F 13/00 (2006. 01)** GO 6 F 13/00 3 5 7 Z

請求項の数 4 (全 17 頁)

(21) 出願番号	特願2009-157717 (P2009-157717)	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成21年7月2日 (2009. 7. 2)	(74) 代理人	100080001 弁理士 筒井 大和
(65) 公開番号	特開2011-15196 (P2011-15196A)	(72) 発明者	江頭 ちひろ 神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内
(43) 公開日	平成23年1月20日 (2011. 1. 20)	(72) 発明者	小国 哲 神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内
審査請求日	平成24年2月17日 (2012. 2. 17)		

最終頁に続く

(54) 【発明の名称】 負荷割当制御方法および負荷分散システム

(57) 【特許請求の範囲】

【請求項1】

複数のクライアントからの複数の要求を、複数のサーバに割り当てる負荷分散システムにおける負荷割当制御方法であって、

前記負荷分散システムにより、前記要求による処理状況に対する前記複数のサーバのそれぞれのサービスレベルを保つための閾値を保持し、前記複数のサーバへの前記要求の割り当ての際、前記複数のサーバの優先度に従い、前記優先度の高いサーバから、前記閾値に達するまで前記要求を割り当て、前記要求に対する前記複数のサーバのそれぞれの処理状況に応じて、前記複数のサーバの電源制御を行い、前記複数のサーバへの前記要求の割り当ての際、前記複数のサーバの全ての処理状況が、前記閾値に達していた場合、前記要求を前記複数のサーバに均等に分配することを特徴とする負荷割当制御方法。

10

【請求項2】

複数のクライアントからの複数の要求を、複数のサーバに割り当てる負荷分散システムであって、

前記複数のサーバのそれぞれの、優先度、コネクション数、レスポンス時間、割り当て時刻、および状態を格納する情報テーブルと、

前記要求による処理状況に対する前記複数のサーバのそれぞれのサービスレベルを保つための閾値を保持する閾値保持部と、

前記複数のサーバの稼動状況を監視管理するサーバ管理部と、

前記複数のサーバへの前記要求の割り当ての際、前記複数のサーバの優先度に従い、前

20

記優先度の高いサーバから、前記閾値に達するまで前記要求を割り当てる転送先判定部と、

前記要求に対する前記複数のサーバのそれぞれの処理状況に応じて、前記複数のサーバの電源制御を行う電源制御部とを備え、

前記転送先判定部は、前記複数のサーバへの前記要求の割り当ての際、前記複数のサーバの全ての処理状況が、前記閾値に達していた場合、前記要求を前記複数のサーバに均等に分配することを特徴とする負荷分散システム。

【請求項 3】

複数のクライアントからの複数の要求を、複数のサーバに割り当てる負荷分散システムであって、

複数の負荷分散装置と、前記複数の負荷分散装置を管理する管理サーバとを備え、

前記管理サーバは、前記複数のサーバのそれぞれの、優先度、コネクション数、レスポンス時間、割り当て時刻、および状態を格納する情報テーブルと、前記要求による処理状況に対する前記複数のサーバのそれぞれのサービスレベルを保つための閾値を保持する閾値保持部と、前記複数のサーバの稼動状況を監視管理するサーバ管理部と、前記複数のサーバへの前記要求の割り当ての際、前記複数のサーバの前記優先度に従い、前記優先度の高いサーバから、前記閾値に達するまで前記要求を割り当てる転送先判定部と、前記複数のサーバのそれぞれの前記要求に対する処理状況に応じて、前記複数のサーバの電源制御を行う電源制御部と、前記転送先判定部での前記要求の割り当てを前記複数の負荷分散装置に指示する負荷分散装置通信部とを有することを特徴とする負荷分散システム。

【請求項 4】

請求項 3 に記載の負荷分散システムにおいて、

前記電源制御部は、前記複数のサーバの電源制御により、特定の負荷分散装置に接続されているサーバの全てを停止可能な場合、前記特定の負荷分散装置の電源を停止させることを特徴とする負荷分散システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理システムなどにおける負荷分散システムに関し、特に、負荷集中機能をもたせ、積極的に負荷量を適度に集中させることで情報処理システム内のサーバなどの稼動台数を抑制することで省電力を実現する方法に関するものである。

【背景技術】

【0002】

昨今、情報処理システムにおいてクライアント PC (パーソナルコンピュータ) からインターネットあるいは LAN (ローカルエリアネットワーク) を介してサーバに接続し、様々な情報を入手したり提供したり、情報の検索あるいは物品の販売など様々なサービスを楽しむようになっている。このようなシステムは一般に Web サーバシステムといわれている。

【0003】

これらのシステムにおいて、クライアント (クライアントは通常はパーソナルコンピュータであることが多い) からの大量のアクセスへの対応や耐障害性を維持することを目的として、クライアントからのアクセス要求を複数サーバに分配することで大量のアクセスにも対応でき、また、その内の一部が故障した場合にも、他のサーバに処理を振り分けることでシステム全体としてサービスを停止させない、といった処理をさせることがよく行われている。このような処理を行う装置が負荷分散装置と呼ばれる。

【0004】

この負荷分散装置において負荷分配は負荷分散アルゴリズムを用いて、接続クライアント台数や処理量などの負荷 (以下、負荷と呼ぶ) がシステムを構成するサーバの間で均等になるように負荷を分配する機能を提供している。

【0005】

10

20

30

40

50

一方で、地球温暖化などの環境問題への取り組みに関心が集まる中、サーバシステムの市場においても省電力化が求められるようになってきており、省電力を実現する機能を提供する技術も盛んに開発されるようになってきている。

【0006】

上記のような従来の負荷分散装置における負荷分配方式の場合、各サーバの性能が高くなくても、複数サーバで多数のクライアントからのアクセスを、均等に分散させて処理できるという利点がある。しかし、一般にこのようなWebサーバシステムではクライアントからのアクセスによる負荷が、例えば1日の内で時間帯によって大きく変動したり、あるいは、時期によって大きく変動したりすることが知られており、それでも一定のサービス品質（応答時間、同時アクセス数の上限確保等）を維持するために最大負荷の時に備えて多数のサーバを稼働させる必要があるが、一方で負荷量が少ない場合などには不必要に多数のサーバを稼働させるということになっており、省電力化のために台数を適正に増減させる技術が必要となってきた。

10

【0007】

上記課題の解決策として、特開2003-281008号公報（特許文献1）に開示された従来技術がある。

【0008】

特許文献1に開示されたものは、稼働時間が最短であるサーバに対して優先的にクライアントからの要求を分配するというものであり、なおかつクライアントからのデータ要求量とサーバシステムのデータ供給可能量とを比較し、サーバシステム全体として必要最低限のサーバだけを稼働させることを目的としたものである。

20

【0009】

特許文献1を含む従来技術では、負荷分配の方法に従来の負荷分散方式を用いて、負荷を複数サーバに分散し、その時点でのサーバシステムの負荷量を計測しながらサーバの台数を増減させていく。システム稼働中において、負荷が存在しないサーバがある場合にはサーバを停止し、稼働中サーバシステムの処理容量が不足してきたら、新たにサーバを起動する方法をとるというものである。

【先行技術文献】

【特許文献】

【0010】

【特許文献1】特開2003-281008号公報

【発明の概要】

【発明が解決しようとする課題】

【0011】

しかしながら、特許文献1に開示された方法の場合、従来の負荷分散装置に比べ電力消費を抑えることが可能ではあるが、負荷分散を行う中で、もし未使用サーバが存在すれば電源を落とし、割り当てサーバの数を少なくしていき、自然と負荷集中が起こることを待っている制御（以後、静的な制御と呼ぶ）しかできていない。

【0012】

また、この静的な制御を行う上で、ある時点で生じた負荷に対して、その時点以前に各サーバからクライアントに要求データを返信した時のサーバシステムの負荷量を用いて負荷の割当先サーバを決めるため、実際の負荷数とサーバ台数の設定にタイムラグが生じる。そのため、場合によって台数が不足、もしくは過剰になることがあり、不足する場合は恐れると常に過剰な台数のサーバを稼働させる必要がある。その結果、省電力効果については最適な制御でない状態が発生する。

40

【0013】

また、クライアントからのアクセスの未来の増減は予測できない中であって、ある時点での負荷量を見て台数を制限するといった静的な負荷集中を行う手段を提供しているが、ある特定サーバに積極的に負荷を割り当て、負荷集中の状態を自ら作り出す（以後、動的な制御と呼ぶ）ことができない。その結果、省電力化のためのサーバ稼働台数の増減を最

50

適制御できていない。

【0014】

また、従来技術において、単純なる動的な負荷集中を行うと、1台のサーバシステムが提供するサービスレベル品質を一定以上に保持することができない。

【0015】

そこで、本発明の目的は、動的な制御を行い、かつユーザが求めるサービスレベルの品質を守った上で最大の省電力効果を得ることができる負荷分散システムを提供することにある。

【0016】

本発明の前記ならびにその他の目的と新規な特徴は、本明細書の記述および添付図面から明らかになるであろう。

【課題を解決するための手段】

【0017】

本願において開示される発明のうち、代表的なものの概要を簡単に説明すれば、次の通りである。

【0018】

すなわち、代表的なものの概要は、負荷分散システムにより、クライアントからの要求に対する複数のサーバのそれぞれのサービスレベルを保つための閾値を保持し、複数のサーバへの要求の割り当ての際、複数のサーバの優先度に従い、優先度の高いサーバから、閾値に達するまで要求を割り当て、複数のサーバのそれぞれの要求に対する処理状況に応じて、複数のサーバの電源制御を行うものである。

【発明の効果】

【0019】

本願において開示される発明のうち、代表的なものによって得られる効果を簡単に説明すれば以下の通りである。

【0020】

すなわち、代表的なものによって得られる効果は、動的な制御を行い、かつユーザが求めるサービスレベルの品質を守った上で最大の省電力効果を得ることができる。

【図面の簡単な説明】

【0021】

【図1】本発明の実施の形態1に係る負荷分散システムが適用されるネットワークシステムの構成を示す構成図である。

【図2】本発明の実施の形態1に係る負荷分散システムの構成を示す構成図である。

【図3】本発明の実施の形態1に係る負荷分散システムの情報テーブルおよび割り当てテーブルの一例を示す図である。

【図4】本発明の実施の形態1に係る負荷分散システムのアルゴリズムの一例を示す図である。

【図5】本発明の実施の形態1に係る負荷分散システムの負荷集中モードの基本的な概念を説明するための説明図である。

【図6】本発明の実施の形態1に係る負荷分散システムの転送先判定部における負荷集中モードの処理を示すフローチャートである。

【図7】本発明の実施の形態1に係る負荷分散システムの転送先判定部における負荷集中モード（ラウンドロビンモード）の処理を示すフローチャートである。

【図8】本発明の実施の形態1に係る負荷分散システムの負荷集中モード（ラウンドロビンモード）の様子を説明するための情報テーブルと割り当てテーブルの一例を示す図である。

【図9】本発明の実施の形態1に係る負荷分散システムの転送先判定部における負荷集中モード（ラウンドロビンモード）の処理を示すフローチャートである。

【図10】本発明の実施の形態1に係る負荷分散システムの負荷集中モード（ラウンドロビンモード）の情報テーブルの一例を示す図である。

10

20

30

40

50

【図 1 1】本発明の実施の形態 1 に係る負荷分散システムによる効果を説明するための説明図である。

【図 1 2】本発明の実施の形態 2 に係る負荷分散システムの構成を示す構成図である。

【発明を実施するための形態】

【0022】

以下、本発明の実施の形態を図面に基づいて詳細に説明する。なお、実施の形態を説明するための全図において、同一の部材には原則として同一の符号を付し、その繰り返しの説明は省略する。

【0023】

(実施の形態 1)

図 1 により、本発明の実施の形態 1 に係る負荷分散システムが適用されるネットワークシステムの構成について説明する。図 1 は本発明の実施の形態 1 に係る負荷分散システムが適用されるネットワークシステムの構成を示す構成図である。

【0024】

図 1 において、ネットワークシステムは、複数のサーバから構成されるサーバシステム(130-1~4)が、負荷分散システム120を介して外部ネットワーク110に接続されている。

【0025】

また、この外部ネットワーク110に接続されているクライアント100とのデータパケットの送受信には必ず負荷分散システム120が仲介している。

【0026】

なお、図 1 に示す例では、クライアント台数と、サーバ台数を 4 としているが、特に 4 台には限定されないし、異なる台数であっても問題ない。

【0027】

次に、図 2 ~ 図 4 により、本発明の実施の形態 1 に係る負荷分散システムの構成について説明する。図 2 は本発明の実施の形態 1 に係る負荷分散システムの構成を示す構成図、図 3 は本発明の実施の形態 1 に係る負荷分散システムの情報テーブルおよび割り当てテーブルの一例を示す図、図 4 は本発明の実施の形態 1 に係る負荷分散システムのアルゴリズムの一例を示す図である。

【0028】

図 2 において、負荷分散システム120は、処理要求受付部121と、サーバシステムの管理を行うための手段としてサーバ管理部122と、要求転送部123と、電源制御部124と、転送先判定部125とサーバシステム全てのサーバの稼動状況や優先度などを保持する情報テーブル126と負荷分散システムが負荷を割り当て可能なサーバのアドレスを保持する割り当てテーブル127とサービスレベルアグリーメントを守るための閾値(以後アセスメントレベルと呼ぶ)を保持する閾値保持部128とを備えている。

【0029】

情報テーブル126は何らかのイベントが発生すると情報が読み出され、また、サーバの状態が変更されるもしくはされた場合にはその都度、情報テーブル126の内容も変更される。

【0030】

図 3 (a) に情報テーブル126の様子、図 3 (b) に割り当てテーブル127の様子を示す。

【0031】

アセスメントレベルは任意に設定可能で、設定内容としては接続コネクション数やCPUビジー率、I/Oビジー率、レスポンス時間などがあげられる。

【0032】

処理要求受付部121は、クライアント100からのデータ要求パケットを受け付ける。受け付けた要求パケットは転送先判定部125へ送られ、転送先判定部125において転送先サーバを選択し、要求転送部123へどのサーバへ要求を送ればよいかを指示する

10

20

30

40

50

。

【0033】

そして、要求転送部123は転送先判定部125より指示されたサーバの割り当てテーブル127のアドレスに要求を無事転送した後、情報テーブル126を更新する。

【0034】

転送先サーバの選択方法としては、負荷分散アルゴリズムを用いて従来の負荷分散方式により転送先サーバを選択する負荷分散モードと、負荷分散アルゴリズムを負荷集中型に変換することで特定のサーバに負荷を集中させる負荷集中モードがあり、某かのアルゴリズムによって制御可能である。

【0035】

図4に、既存の負荷分散装置の負荷分散アルゴリズムを参考に一般化したアルゴリズムの一覧を示す。

【0036】

サーバ管理部122は、情報テーブル126と割り当てテーブル127とサーバシステム(130-1~4)を監視管理する。

【0037】

例えば、所定期間内にサーバシステム(130-1~4)の各稼動中サーバの処理量を調査し、未使用サーバが存在し、さらに電源停止条件を満たしていると判断したならば、電源制御部に対象サーバの電源停止を指示し、情報テーブル126の内容を“停止中”と更新する。

【0038】

逆に、上記所定期間内にサーバシステム(130-1~4)の全ての稼動中サーバの負荷量が供給可能量を超えており、電源起動条件を満たしていると判断したならば、サーバシステム(130-1~4)の内、必要数の待機中サーバを優先度が高い順に起動するよう電源制御部に指示する。

【0039】

本実施の形態による負荷集中の場合、優先度の高いサーバから順に負荷を片寄せしていくため、次に電源を停止させるサーバは、稼動中サーバの中で最も優先度の低いサーバとなるため、明示的にこのサーバの負荷量を見て、未使用の状態、さらに電源停止条件を満たしていると判断したならば、対象サーバを割り当て可能対象から外し、電源制御部に対象サーバの電源停止を指示した後、対象サーバの電源停止を確認した上で情報テーブル126の内容を“停止中”と更新する、という方法も可能である。

【0040】

そして、情報テーブル126が持つ対象サーバの状態を“起動処理中”と更新する。起動処理中サーバの正常な起動を確認したら、割り当てテーブル127に対象サーバを追加し割り当て可能となった時点で、情報テーブルの内容を“稼動中”と更新する。

【0041】

また、サーバ管理部122は、サーバシステム(130-1~4)の待機中サーバの優先度の変更を行うなどして、電源の起動と停止の動作が特定のサーバに集中することを防ぐことも可能とする。

【0042】

次に、図5~図11により、本発明の実施の形態1に係る負荷分散システムの転送先判定部の負荷集中モードの処理について説明する。図5は本発明の実施の形態1に係る負荷分散システムの負荷集中モードの基本的な概念を説明するための説明図、図6は本発明の実施の形態1に係る負荷分散システムの転送先判定部における負荷集中モードの処理を示すフローチャート、図7および図9は本発明の実施の形態1に係る負荷分散システムの転送先判定部における負荷集中モード(ラウンドロビンモード)の処理を示すフローチャート、図8は本発明の実施の形態1に係る負荷分散システムの負荷集中モード(ラウンドロビンモード)の様子を説明するための情報テーブルと割り当てテーブルの一例を示す図、図10は本発明の実施の形態1に係る負荷分散システムの負荷集中モード(ラウンドロビ

10

20

30

40

50

ンモード)の情報テーブルの一例を示す図、図11は本発明の実施の形態1に係る負荷分散システムによる効果を説明するための説明図である。

【0043】

まず、負荷集中モードにおいて、図4に示した負荷分散アルゴリズム各々を負荷集中型に変換する方法について、負荷集中モードの基本的な概念を図5を用いて説明した後に、図6～図9を用いて説明する。

【0044】

本来の負荷分散と本実施の形態の負荷集中との相違点を次に示す。

【0045】

本来の負荷分散は、各サーバに負荷が均等に割り振られる。

10

【0046】

例えば、サーバが3台あり、各々#0、#1、#2とし、何らかのサーバで処理する負荷が6個順にもたらされるとすると、これらの負荷が全て等しいとするならば、3台のサーバに順に負荷1～3を割り振り、さらに負荷4～6を割り振ることで均等にすることができる。負荷1～6が全てサーバに割り振られ、処理中である様子を図5(a)に示す。

【0047】

これを本実施の形態の負荷集中モードでは、サーバに優先度を設け、優先度の高いサーバに集中的に負荷を割り振る。

【0048】

その際、負荷が特定のサーバに集中すると、そのレスポンス時間が長くなるなどしてサービスの質が落ちるため、サービスレベルアグリーメントで規定するサービスを維持できるアセスメントレベルを設定し、サーバの処理中の負荷量がこのアセスメントレベルに達するまではある特定のサーバに負荷を割り振り、アセスメントレベルに達した後は、次に優先度の高いサーバに負荷の割り振りを行う。

20

【0049】

例えば、アセスメントレベルを1台あたり4個までの負荷が許容範囲であることとし、#0、#1、#2の順に優先度が高いサーバが3台あるとする。何らかのサーバで処理する負荷が6個順にもたらされるとすると、サーバ#0には負荷1～4が、サーバ#1には負荷5～6が割り振られ、サーバ#2には何も割り振られていない状態となる(図5(b))。

30

【0050】

この状態であればサーバ#2は電源停止をただちに行うことができる。

【0051】

しかし、負荷がさらに増大していき、サーバシステムがアセスメントレベルを守れなくなった場合(これは本来、システム設計上あってはならないことではあるが、予め想定しておく必要はある)、既にこの段階でシステム設計上想定していた負荷量を超えており、できるだけ負荷を分散させた方が賢明である。

【0052】

例えば、先ほどの例に続けて負荷がかかった場合、やがて負荷量がアセスメントレベルを超えることとなる。ここから先は、負荷の増大に対してアセスメントレベルを守ることができないので、負荷分散モードに変更させることが必要となる。この状況を図5(c)に示す。すなわち、図5(c)の負荷1～12まではサーバ#0～2に4個ずつ割り振られているが、その後の負荷13～15はサーバ#0～2に均等に1個ずつを順に割り振っていくことになる。

40

【0053】

そして、その後負荷が減ってきてアセスメントレベルを守れるようになれば、また本来の負荷集中モードに戻るものとする。

【0054】

この場合のアセスメントレベルは、接続数やレスポンス時間などユーザが望むもので任意に設定可能である。

50

## 【 0 0 5 5 】

この負荷集中モードの基本的概念を基に、負荷分散アルゴリズムでの負荷集中の方法を最小コネクションモード、最速モード、ラウンドロビンモードを例にあげて説明する。

## 【 0 0 5 6 】

まず、最小コネクションモードの場合については、上記負荷集中モードの基本的概念の負荷量をコネクション数と置き換えて考える。

## 【 0 0 5 7 】

最速モードの場合には、各サーバの最新のレスポンス時間（あるいは最新の複数個のリクエストに対するレスポンス時間の平均）を見て、特定のサーバについてレスポンス時間がアセスメントレベル以下である内は負荷を割り振り、超えたら別のサーバへ割り振ると考える。

10

## 【 0 0 5 8 】

最速モードと最小コネクションモードの違いとしては、アセスメントレベルがコネクション数なのか、レスポンス時間なのかという違いである。

## 【 0 0 5 9 】

この2つのモードについて図6に示すフローチャートを用いて説明する。

## 【 0 0 6 0 】

前提条件として、新たにクライアントから接続要求があった場合、要求受付部によって受付処理が行われたリクエストは転送先判定待ち行列にスタックされるものとする。

## 【 0 0 6 1 】

20

この転送先判定待ち行列にデータが存在した場合（ステップ700）、閾値保持部128よりアセスメントレベルの読み込みと情報テーブル126にアクセスを行い読み書き可能な状態とし、テーブル情報を読み込む（ステップ701）。

## 【 0 0 6 2 】

そして、モードが負荷集中モードで、状態が稼働中の全てのサーバにおいて、負荷量がアセスメントレベルに達している場合には、負荷集中モードを負荷分散モードに切り替え（ステップ702、704、712）、モードが負荷分散モードなら、負荷分散アルゴリズムにより負荷が均等に割り振られるようにサーバを決定し（ステップ702、703）、負荷分散処理後、負荷集中モードへと戻す（ステップ706）。ここで用いる負荷分散アルゴリズムについては公知の技術であるため、ここでは詳細には述べないこととする。

30

## 【 0 0 6 3 】

稼働中のサーバにおいて、負荷量がアセスメントレベルに達していないサーバが存在する場合には、優先度が最も高いサーバを“調査サーバ”とする（ステップ704、705）。

## 【 0 0 6 4 】

そして、“調査サーバ”の負荷量とアセスメントレベルとを比較し（ステップ713）、ステップ713で、アセスメントレベルに達している場合には、優先度が次に高いサーバを“調査サーバ”として負荷量がアセスメントレベルに達しているかを確認していく（ステップ714）。

## 【 0 0 6 5 】

40

ステップ713で、“調査サーバ”の負荷量がアセスメントレベルに達していない場合には、そのサーバを割り当てサーバに決定し（ステップ707）、情報テーブル126のサーバのコネクション数欄を更新する（ステップ708）。

## 【 0 0 6 6 】

その際のアセスメントレベルとサーバの負荷量の比較とは、例えばアセスメントレベルがコネクション数の場合にはサーバの接続コネクション数の値であり、レスポンス時間の場合には、レスポンス時間の値との比較である。

## 【 0 0 6 7 】

そして、リクエスト情報を1つ転送先判定待ち行列から読み出し、転送先サーバの情報を付加して転送待ち行列にスタックする（ステップ709、710）。

50

## 【 0 0 6 8 】

そして、判定待ち行列にデータがあるかを判断し（ステップ711）、ステップ711で判定待ち行列にデータがあれば、ステップ704に戻り、ステップ711で判定待ち行列にデータが無ければステップ700に戻る。

## 【 0 0 6 9 】

要求転送部123は、この転送待ち行列に積まれているデータを取り出して処理を行う。ただし、上記では、アセスメントレベルに達していないサーバとは、新たに1つ負荷を割り振った場合にアセスメントレベルを達するか達しないかの状態にはなるが、アセスメントレベルを超える状態（つまりサービスレベルアグリーメントを満足できない状態）にはならないと仮定している。

10

## 【 0 0 7 0 】

この仮定を満足させるにはアセスメントレベルを以下のように設定しておく必要がある。すなわち、サーバに割り振った負荷量に対してサービスレベルがどのように変化するかを予め実測値や理論値を元に算出できるようにしておき（しかもこれは負荷量増大に対しサービスレベルが単純に低下する単調減少となっている必要がある）、そうでないならば、そもそも負荷量として設定している計量単位を見直すべきである）、サーバに新たに1個の負荷を割り振った場合にサービスレベルアグリーメントを満足できる、最大の負荷量をアセスメントレベルと設定する必要がある。

## 【 0 0 7 1 】

もし、例えば1個の負荷に対し負荷量の増減にばらつきがあるならば、アセスメントレベルは、あるサーバに新たに最大の負荷量をもつ1個の負荷が割り振られてもサービスレベルアグリーメントを満足する最大の負荷量とすべきであるし、もっと工夫するならば、図6に示す処理の中で、アセスメントレベルをそのときに割り振ろうとしている負荷の内容から負荷量の増分を予測し、その予測値に基づいてアセスメントレベルに達しているのかどうかを判断してもよい。

20

## 【 0 0 7 2 】

通常、負荷分散装置は処理要求の内容を解析する処理を持っている場合が多いのでこの種のインプリメントは比較的容易に可能である。また、さらに工夫するならば、サーバ毎にこのアセスメントレベルを別にしてもよい。本発明の場合、アセスメントレベルに達しているかどうかを判定するステップ（図6のステップ713）より前のステップ（図6のステップ705）で調査サーバ（すなわち、割り振り先のサーバ）が決まっているため、このような工夫を行うことができる。

30

## 【 0 0 7 3 】

次に、図7および図9のフローチャートにより、ラウンドロビンモードについて説明する。

## 【 0 0 7 4 】

本来のラウンドロビンモードは、負荷を複数台のサーバに順番に均等に割り振ることで負荷分散を達成することができるが、この考え方で負荷集中に変換させることは困難である。ラウンドロビンモードの場合には、次の2通りの考え方をを用いることとする。

## 【 0 0 7 5 】

1つ目の考え方としては、ラウンドロビンモードを単純に、負荷が発生した場合、前回割り当てたサーバの次に優先度が高いサーバに割り当てると考える、とした場合、割当先が1つしか存在しないという場合であれば、そのサーバにのみ負荷が割り当てられる。

40

## 【 0 0 7 6 】

これによって、本実施の形態では、負荷集中は以下のように考えて実現する。

## 【 0 0 7 7 】

「複数台の稼動中サーバが存在する中で、割当先はある特定の1つのサーバのみとする。負荷量がアセスメントレベルに達するまでは割当先はその特定サーバとし、アセスメントレベルに達したら、割当先を優先度が次に高いサーバとする。」

つまり、以下のように処理を行う。

50

## 【 0 0 7 8 】

( 1 ) 稼働中の全てのサーバにおいて、負荷量がアセスメントレベルに達しているならば、負荷分散モードに切り替え ( ステップ 8 0 0 、 8 0 1 、 8 0 2 、 8 0 4 、 8 1 3 ) 、負荷分散アルゴリズムに従って、つまり最後に割り当てたサーバの次に優先度の高いサーバに負荷を割り当てる ( ステップ 8 0 0 、 8 0 1 、 8 0 2 、 8 0 3 、 8 0 6 ) 。

## 【 0 0 7 9 】

( 2 ) 稼働中サーバに、アセスメントレベルに達していないサーバがある場合には、優先度が最も高いサーバを “ 調査サーバ ” と指定する ( ステップ 8 0 4 、 8 0 5 ) 。

## 【 0 0 8 0 】

( 3 ) 調査サーバの負荷量がアセスメントレベルに達しているか確認する ( ステップ 8 0 7 ) 。

10

## 【 0 0 8 1 】

( 4 ) ステップ 8 0 7 で調査サーバの負荷量がアセスメントレベルに達していなければ、調査サーバを割り当てサーバに決定し、割り当てテーブルと情報テーブルを変更し、割り当てテーブルに存在するサーバ情報の削除と調査サーバの割り当てテーブルへの追加を行い ( ステップ 8 0 8 ) 、情報テーブル 1 2 6 の “ 調査サーバ ” の負荷量を追加する ( ステップ 8 0 9 ) 。

## 【 0 0 8 2 】

( 5 ) 転送待ち行列に、転送先サーバの情報を付加したデータをスタックし ( ステップ 8 1 0 ) 、判定待ち行列にデータがあるかを判断し ( ステップ 8 1 1 ) 、ステップ 8 1 1 で判定待ち行列にデータがあれば、ステップ 8 0 4 に戻り、ステップ 8 1 1 で判定待ち行列にデータが無ければステップ 8 0 0 に戻る。

20

## 【 0 0 8 3 】

( 6 ) ステップ 8 0 7 で、調査サーバの負荷量がアセスメントレベルに達している場合には、次に優先度が高いサーバを調査サーバとしてステップ 8 0 7 に戻る ( ステップ 8 1 2 ) 。

## 【 0 0 8 4 】

例えば、アセスメントレベルがコネクション数 3 とし、情報テーブルと割り当てテーブルが図 8 ( a ) に示す状態の時点で、負荷が 1 個生起したとする。優先度順にコネクション数を確認すると、コネクション数がアセスメントレベルに達していない稼働中サーバの内、最も優先度が高いのはサーバ番号が 2 のサーバとなるため、割当先と決定する。そのため、割り当てテーブルの既存のサーバの情報を削除し、新たに対象サーバが追加され、情報テーブルと割り当てテーブルは、図 8 ( a ) に示す状態から、図 8 ( b ) に示す状態へ変更される。

30

## 【 0 0 8 5 】

また、2 つ目の考え方としては、負荷分散におけるランドロビンモードについて、ある一定期間内に割り振る負荷の数を、対象とするサーバ間で同数になるよう割り振るものと解釈し、これによって負荷集中は以下のように考えて実現する。

## 【 0 0 8 6 】

「ある一定期間内に割り振る負荷の数を、対象とするサーバの中で特定のサーバに集中して割り振る。ただし、割り振る負荷の数はアセスメントレベル以下となるようにする。」

40

つまり、以下のように処理を行う。

## 【 0 0 8 7 】

( 1 ) 稼働中の全てのサーバにおいて、負荷量がアセスメントレベルに達しているならば、負荷分散モードに切り替え ( ステップ 8 1 4 、 8 1 5 、 8 1 6 、 8 1 7 、 8 1 8 ) 、負荷分散アルゴリズムに従って、つまり最後に割り当てたサーバの次に優先度の高いサーバに負荷を割り当てる ( ステップ 8 1 4 、 8 1 5 、 8 1 7 、 8 1 9 、 8 2 2 ) 。

## 【 0 0 8 8 】

( 2 ) 稼働中サーバに、アセスメントレベルに達していないサーバがある場合には、優

50

先度が最も高いサーバを“調査サーバ”と指定する(ステップ818、820)。

【0089】

(3) 現時刻から時刻 - T の間に調査サーバに割り当てたコネクション数をカウントする(ステップ821)。

【0090】

(4) 調査サーバに割り当てたコネクション数がアセスメントレベルに達しているか確認する(ステップ823)。

【0091】

(5) ステップ823で調査サーバに割り当てたコネクション数がアセスメントレベルに達していないなら、そのサーバに割り当ててを決定し、割り当てテーブルを変更し、割り当てテーブルに割り当てた時刻と、調査サーバの情報を追加する(ステップ824)。

10

【0092】

(6) 情報テーブル内の、現時点から時刻 T 以前のデータを削除し(ステップ825)、判定待ち行列にデータがあるかを判断し(ステップ827)、ステップ827で判定待ち行列にデータがあれば、ステップ818に戻り、ステップ827で判定待ち行列にデータが無ければステップ814に戻る。

【0093】

(7) ステップ823で、調査サーバに割り当てたコネクション数がアセスメントレベルに達している場合には、次に優先度が高いサーバを調査サーバとしてステップ821に戻る(ステップ826)。

20

【0094】

情報テーブルには、例えば、図10に示すように過去に割り当てたサーバとその時刻が記録されており、一定時間(例えば1秒とか10秒とか1分とか1時間とか、アプリケーションに応じてユーザが選択する)内の割り当てた個数を算出できるようになっており、調査サーバについて、その個数がアセスメントレベル以下であるかどうかをステップ823の判定で行うようになっている。

【0095】

次に、本実施の形態の負荷分散システムによる効果について説明する。

【0096】

まず、比較例として、従来技術による省電力制御での動作を図11(a)に示す。

30

【0097】

サーバシステムに稼働サーバが3台存在し、サーバ1台あたり3個の負荷まで処理できるとする(すなわちアセスメントレベルが負荷量3としていることを意味する)。従来技術の場合、各サーバの負荷が均等になるように割り当てられるため、例えば、時刻 T1 で生じた7個の負荷に対して、各サーバが処理中である場合、この T1 の時点での負荷の割り当ては、各サーバに負荷が均等になるよう割り当てられるため、#1から順に3個、2個、2個となる。

【0098】

また、時刻 T2、T3 の時点で生起する負荷が6個、5個であった場合、サーバ2台で対応可能な負荷量が6個のため、単純に考えると2台で対応可能のためこの時点でサーバを2台に削減できることになる。

40

【0099】

しかし、負荷分散の場合、一旦ここで3台のサーバに負荷が均等となるよう割り当てられるため、例えば、時刻 T2 に生じた負荷は2個ずつ割り当てられるため、最速でも2個の処理が終了した時点での電源停止となる。

【0100】

これに対して、本実施の形態のように、アセスメントレベルに達するまでは特定サーバに負荷を割り当て、達したならば、次のサーバに割り当てるという負荷集中を行うと、図11(b)に示すように、まず、時刻 T1 で発生した7個の負荷はサーバ#1からアセス

50

メントレベルに達するまで負荷を割り当てていくため、割り当てられた負荷の個数は順に、3個、3個、1個となる。

【0101】

その後、T2で発生した負荷6個はサーバ2台で対応可能であるため、サーバ#3には既に負荷が割り当てられることはないので、時刻T2以前で電源をオフにすることができる(場合により の時点までさかのぼることも可能である)。

【0102】

なお、図11に示す簡略なモデルから、本実施の形態の省電力制御としての効果の大きさは負荷の処理時間の長さ依存していることがわかる。例えば、従来技術では、サーバを電源停止するためには、電源停止するサーバでの負荷の処理が終了しないと電源停止できないが、従来技術では、電源停止させるサーバへの負荷の割り当てが均等に割り当てられているため、本実施の形態と比べて、その負荷の個数が多くなるため、負荷の処理時間が長くなれば電源停止するまでの時間が長くなってしまふ。

10

【0103】

従って、例えば、インターネットショッピングやネットでのストリーミングデータを配信する場合のように一旦接続すると相当長い間継続する場合には、本実施の形態での負荷集中処理は効果的である。

【0104】

(実施の形態2)

実施の形態1では、負荷分散システム120に、サーバ管理部122、転送先判定部125、電源制御部124、情報テーブル126、割り当てテーブル127、閾値保持部128の機能を設けているが、実施の形態2では、複数の負荷分散装置を管理する管理サーバに設けるようにしたものである。

20

【0105】

図12により、本発明の実施の形態2に係る負荷分散システムの構成について説明する。図12は本発明の実施の形態2に係る負荷分散システムの構成を示す構成図である。

【0106】

図12において、負荷分散システムは、4台のサーバから構成されるサーバシステムを接続した負荷分散装置166が2台と、負荷分散装置166を管理する管理サーバ177から構成されている。

30

【0107】

負荷分散装置166は、処理要求受付部121と、管理サーバ177との通信を行う管理サーバ通信部162と、管理サーバ177の指示によって、負荷をサーバに転送する要求転送部123と割当先サーバの情報を保持する割り当てテーブル127とを備える。

【0108】

管理サーバ177は、各負荷分散装置166との通信を行う負荷分散装置通信部171と負荷分散装置166で受け付けた処理の配分先を判定する転送先判定部125と、全てのサーバを監視・管理するサーバ管理部122と、計算機の電源制御を行う電源制御部124と、各サーバと負荷分散装置166の情報を保持する情報テーブル126と、アセスメントレベルを保持する閾値保持部128を備える。

40

【0109】

本実施の形態の負荷分散システムとしての全体の動作は、実施の形態1と同様である。

【0110】

このように、本実施の形態では、管理サーバ177上に、サーバ管理部122、転送先判定部125、電源制御部124、情報テーブル126、割り当てテーブル127、閾値保持部128の機能を負荷分散装置166と独立してもつことで、複数の負荷分散装置166と、複数の負荷分散装置166に接続されたサーバシステムを制御可能となり、より省電力効果が得られる。

【0111】

なお、図12に示す例では、管理サーバ177による制御対象の負荷分散装置166を

50

2台としているが、2台である必要はなく、複数台の制御が可能である。

【0112】

この場合、複数の負荷分散装置166をまたいで負荷集中が可能となる。つまり、図12に示す例では、1台の負荷分散装置166に接続されているサーバシステムのみで処理可能な場合、負荷を割り当てないサーバシステム側については、負荷分散装置166ごと電源を停止することが可能となる。

【0113】

以上、本発明者によってなされた発明を実施の形態に基づき具体的に説明したが、本発明は前記実施の形態に限定されるものではなく、その要旨を逸脱しない範囲で種々変更可能であることはいうまでもない。

10

【0114】

例えば、実施の形態1、2において、クライアントからの要求の割当先であるサーバシステムは、物理サーバに限定されず、仮想サーバであってもよい。

【産業上の利用可能性】

【0115】

本発明は、情報処理システムなどにおける負荷分散システムに関し、負荷分散と共に、システム全体の省電力化が必要なシステムなどに広く適用可能である。

【符号の説明】

【0116】

100(100-1~100-4)...クライアント、110...外部ネットワーク、120...負荷分散システム、121...処理要求受付部、122...サーバ管理部、123...要求転送部、124...電源制御部、125...転送先判定部、126...情報テーブル、127...割り当てテーブル、128...閾値保持部、130(130-1~130-4)...サーバシステム、162...管理サーバ通信部、166...負荷分散装置、171...負荷分散装置通信部、177...管理サーバ。

20

【図1】

【図2】

図1

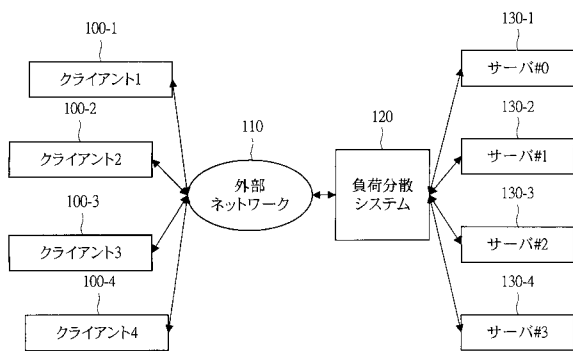
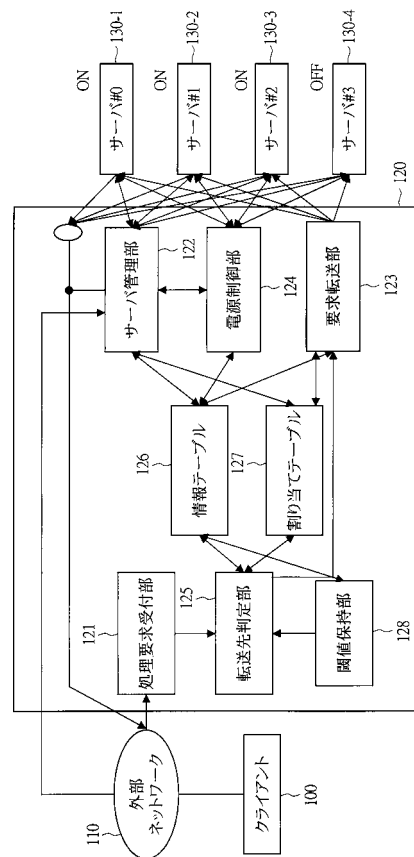


図2



【図3】

図 3

(a)

サーバ番号	優先度	状態	コネクション数	レスポンス時間	アドレス
#0	1	稼動中	3	0.0531	XXXX.XXXX.XXXX.XXXX
#1	2	稼動中	3	0.0455	XXXX.XXXX.XXXX.XXXX
#2	3	稼動中	2	0.0155	XXXX.XXXX.XXXX.XXXX
#3	4	停止中	0	0.0000	XXXX.XXXX.XXXX.XXXX

(b)

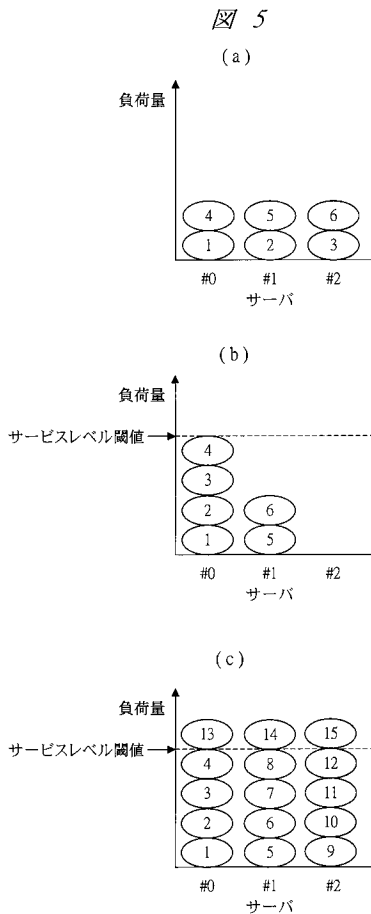
サーバ番号	アドレス	コネクション数	レスポンス時間
#0	XXXX.XXXX.XXXX.XXXX	3	0.0531
#1	XXXX.XXXX.XXXX.XXXX	3	0.0455
#2	XXXX.XXXX.XXXX.XXXX	3	0.0155

【図4】

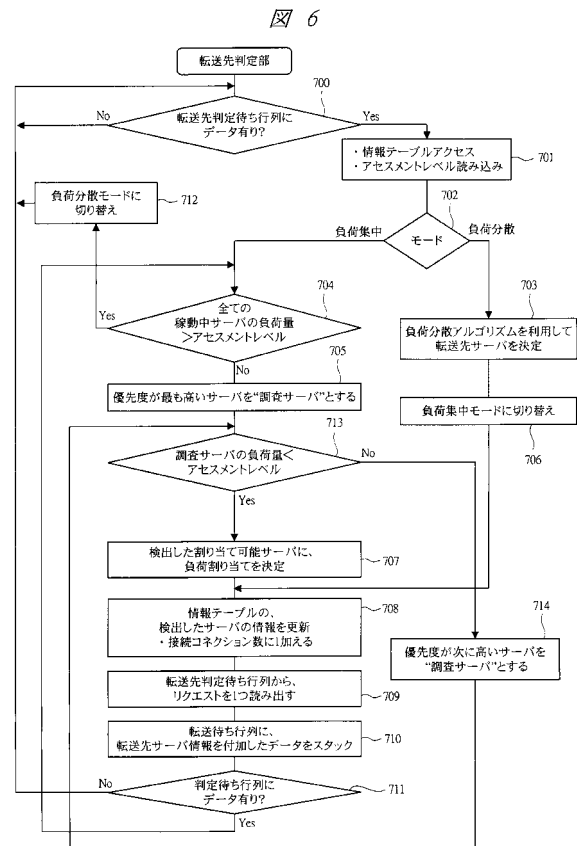
図 4

#	アルゴリズム名称	負荷分散アルゴリズムの説明
1a	最少コネクションモード	リクエストを、各サーバに接続しているTCPコネクションが最も少ないものにリクエストを割り振る。
1b	重み付き最少コネクションモード	最少コネクションモードにサーバごとに割り当ての重み付けをつけて負荷分散を実施。
2a	ラウンドロビンモード	リクエストを複数台のサーバに順番に均等に割り振る。
2b	重み付きラウンドロビンモード	ラウンドロビンにサーバごとに割り当ての重み付けをつけて負荷分散を実施。
3	最速モード	リクエストを最も応答時間が早いサーバへ優先的に割り振る。

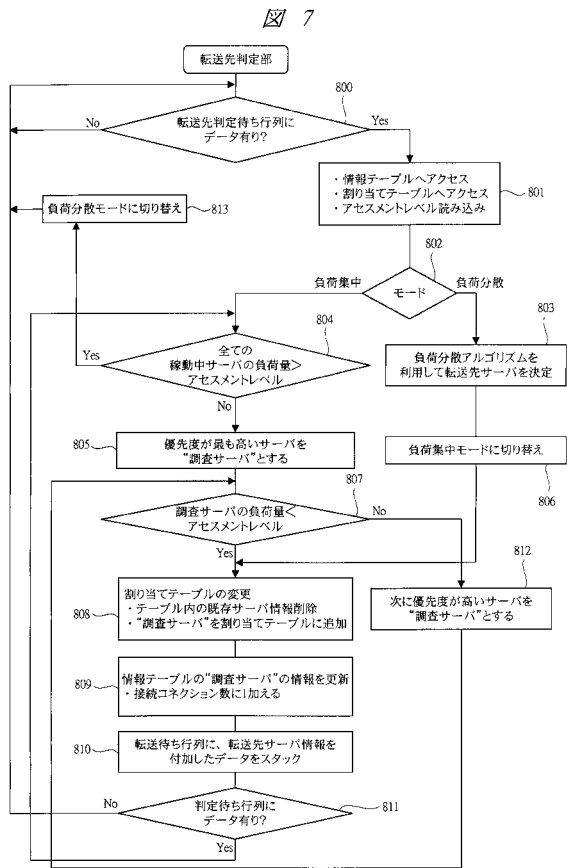
【図5】



【図6】



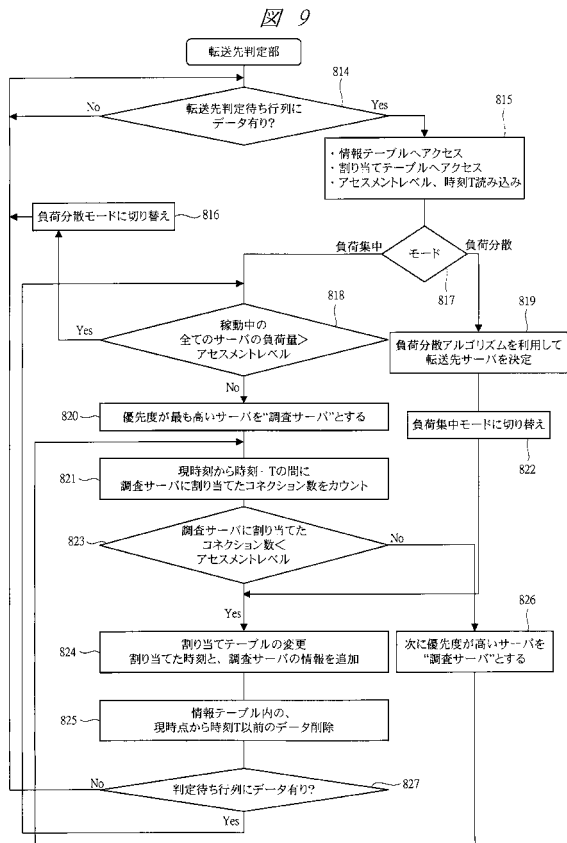
【図7】



【図8】



【図9】

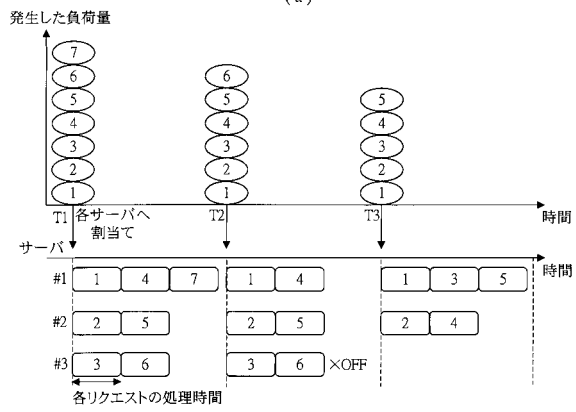


【図10】

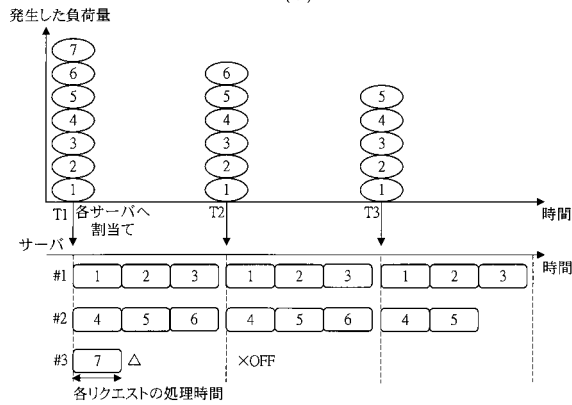
サーバ番号	時刻	状態	優先度	コネクション数	レスポンス時間	アドレス
#0	15:00	稼働中	1	1		XXXX.XXXX.XXXX.XXXX
#0	17:00	稼働中	1	1		XXXX.XXXX.XXXX.XXXX
#0	17:03	稼働中	1	1		XXXX.XXXX.XXXX.XXXX
#1	18:00	稼働中	2	1		XXXX.XXXX.XXXX.XXXX
#1	18:22	稼働中	2	1		XXXX.XXXX.XXXX.XXXX
#1	18:30	稼働中	2	1		XXXX.XXXX.XXXX.XXXX
#1	18:45	稼働中	2	1		XXXX.XXXX.XXXX.XXXX
#2	18:50	稼働中	2	1		XXXX.XXXX.XXXX.XXXX
#2	20:35	稼働中	2	1		XXXX.XXXX.XXXX.XXXX
#3	21:22	停止中	2	0		XXXX.XXXX.XXXX.XXXX
#4	21:30	停止中	2	0		XXXX.XXXX.XXXX.XXXX

【図11】

図11 (a)

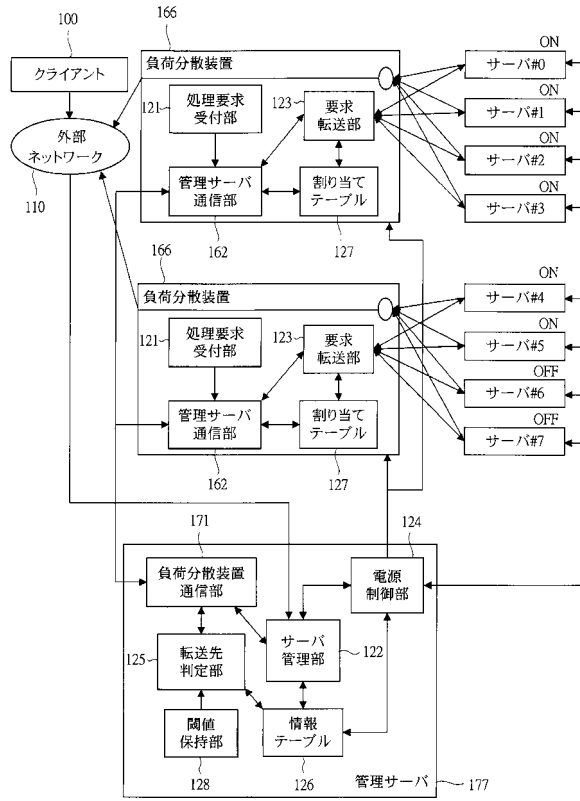


(b)



【図12】

図12



---

フロントページの続き

- (72)発明者 亀田 泰弘  
神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内
- (72)発明者 藤田 博文  
神奈川県秦野市堀山下1番地 株式会社日立製作所 エンタープライズサーバ事業部内

審査官 浦口 幸宏

- (56)参考文献 特開2008-225642(JP,A)  
特開2010-165193(JP,A)  
特開2006-227963(JP,A)  
特開2008-269249(JP,A)

- (58)調査した分野(Int.Cl., DB名)
- |      |                |
|------|----------------|
| G06F | 1/26 - 1/32    |
| G06F | 9/46 - 9/54    |
| G06F | 13/00          |
| H04L | 12/00 - 12/26  |
| H04L | 12/50 - 12/955 |