

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3987517号  
(P3987517)

(45) 発行日 平成19年10月10日(2007.10.10)

(24) 登録日 平成19年7月20日(2007.7.20)

(51) Int. Cl. F I  
G O 6 F 9/50 (2006.01) G O 6 F 9/46 4 6 2 A

請求項の数 9 (全 17 頁)

(21) 出願番号	特願2004-254427 (P2004-254427)	(73) 特許権者	000003078
(22) 出願日	平成16年9月1日(2004.9.1)		株式会社東芝
(65) 公開番号	特開2005-100387 (P2005-100387A)		東京都港区芝浦一丁目1番1号
(43) 公開日	平成17年4月14日(2005.4.14)	(74) 代理人	100058479
審査請求日	平成16年9月1日(2004.9.1)		弁理士 鈴江 武彦
(31) 優先権主張番号	特願2003-310161 (P2003-310161)	(74) 代理人	100091351
(32) 優先日	平成15年9月2日(2003.9.2)		弁理士 河野 哲
(33) 優先権主張国	日本国(JP)	(74) 代理人	100088683
			弁理士 中村 誠
		(74) 代理人	100108855
			弁理士 蔵田 昌俊
		(74) 代理人	100075672
			弁理士 峰 隆司
		(74) 代理人	100109830
			弁理士 福原 淑弘

最終頁に続く

(54) 【発明の名称】 計算機システム及びクラスタシステム用プログラム

(57) 【特許請求の範囲】

【請求項1】

2台以上の計算機を持つ計算機システムにおいて、

前記計算機が実行する複数のサービスの割当てを決定するためのサービスの優先順位とサービスに割当てする計算機の優先順位とサービス間の排他、依存関係を含む関係とサービス実行に必須の周辺機器を含むリソースの割当てと最低負荷状況の計算機にサービスを割当てするための計算機の負荷状況とのうち少なくとも一つを持つ最適配置用ポリシー情報と、サービスのスイッチオーバーの可/不可とサービスの実行可能なノードが無い場合の他サービスの停止の可/不可と前記スイッチオーバー/他サービス停止の判断基準と負荷状況の変化時の対応とのうち少なくとも一つを持つサービスの再配置用ポリシー情報とを保存するポリシー管理部と、

前記最適配置用ポリシー情報に基づいて実行するサービスを最適な計算機に割当てするサービス最適配置部と、

前記最適配置後の各計算機でのサービス負荷や計算機負荷状況によりサービスの再配置が必要か否かを判定する負荷管理部と、

前記負荷管理部によるサービスの再配置が必要との判定結果に従って、前記再配置用ポリシー情報に基づき、前記負荷管理部により再配置が必要と判定されたサービスの実行に最適な計算機を決定するサービス再配置部と、

前記再配置が必要と判定されたサービスの実行を一時的に停止させ、前記サービス配置部により決定された前記再配置が必要と判定されたサービスの実行に最適な計算機で前記

10

20

再配置が必要と判定されたサービスを起動してサービスをスイッチオーバーするサービス制御部とを具備することを特徴とする計算機システム。

【請求項 2】

前記サービス最適配置部は所望のサービスの起動時にその実行に最適な計算機を前記ポリシー管理部に保管されている最適配置用ポリシー情報を参照して決定することを特徴とする請求項 1 に記載の計算機システム。

【請求項 3】

前記サービス再配置部は計算機間における実行中のサービス配置に不均衡が発生した時にサービスの再配置の必要性を検知する手段を含み、前記検知する手段の出力により前記再配置用ポリシー情報を参照して前記サービスの再配置を行うことを特徴とする、請求項 1 に記載の計算機システム。

10

【請求項 4】

前記検知する手段は、各計算機の負荷の状況を検知する検知手段として前記負荷管理部に含まれることを特徴とする、請求項 3 に記載の計算機システム。

【請求項 5】

前記検知する手段は各計算機のノード負荷モニタを含むことを特徴とする請求項 4 に記載の計算機システム。

【請求項 6】

前記サービス再配置部は、前記各計算機の負荷状況の変化に応じてサービスの再配置の必要性を判断し、

20

当該サービスの再配置の必要性がある場合に、前記再配置用ポリシー情報に従って予備計算機の使用を含む再配置処理を実行することを特徴とする、請求項 1 に記載の計算機システム。

【請求項 7】

更に、前記計算機システムは 2 つ以上のクラスタシステムおよび各クラスタシステムが共通して使用可能なプロビジョニング計算機群を含み、

前記ポリシー管理部は、前記プロビジョニング計算機の割当て処理および切り離し処理のポリシーを指定する為の割当て、切り離し処理用ポリシー情報を保管し、

前記計算機システムは更に前記割当て、切り離し処理用ポリシー情報に従って、前記プロビジョニング計算機群から追加要求の計算機を割当てる割当て処理または余剰な計算機を切り離す切り離し処理を実行する割当て/切り離し手段を含む、請求項 1 に記載の計算機システム。

30

【請求項 8】

2 台以上の計算機が接続されて一つのクラスタシステムを実現する計算機システムによるサービス実行方法であって、

前記計算機が実行する複数のサービスの割当てを決定するためのサービスの優先順位とサービスに割当てる計算機の優先順位とサービス間の排他、依存関係を含む関係とサービス実行に必須の周辺機器を含むリソースの割当てと最低負荷状況の計算機にサービスを割当てるための計算機の負荷状況とのうち少なくとも一つを持つ最適配置用ポリシー情報と、サービスのスイッチオーバーの可/不可とサービスの実行可能なノードが無い場合の他サービスの停止の可/不可と前記スイッチオーバー/他サービス停止の判断基準と負荷状況の変化時の対応とのうち少なくとも一つを持つサービスの再配置用ポリシー情報とを保存し、

40

前記最適配置用ポリシー情報に基づいて実行するサービスを最適な計算機に割当てる処理を実行し、

前記最適配置後の各計算機でのサービス負荷や計算機負荷状況によりサービスの再配置が必要か否かを判定し、

前記負荷管理部によるサービスの再配置が必要との判定結果に従って、前記再配置用ポリシー情報に基づき、前記負荷管理部により再配置が必要と判定されたサービスの実行に最適な計算機を決定し、

前記再配置が必要と判定されたサービスの実行を一時的に停止させ、前記サービス配置

50

部により決定された前記再配置が必要と判定されたサービスの実行に最適な計算機で前記再配置が必要と判定されたサービスを起動してサービスをスイッチオーバーすることを特徴とするサービス実行方法。

【請求項 9】

2台以上の計算機が接続された計算機システムに適用し、一つのクラスタシステムを実現するためのプログラムであって、

前記計算機が実行する複数のサービスの割当てを決定するためのサービスの優先順位とサービスに割当てた計算機の優先順位とサービス間の排他、依存関係を含む関係とサービス実行に必須の周辺機器を含むリソースの割当てと最低負荷状況の計算機にサービスを割当てたための計算機の負荷状況とのうち少なくとも一つを持つ最適配置用ポリシー情報と、サービスのスイッチオーバーの可/不可とサービスの実行可能なノードが無い場合の他サービスの停止の可/不可と前記スイッチオーバー/他サービス停止の判断基準と負荷状況の変化時の対応とのうち少なくとも一つを持つサービスの再配置用ポリシー情報とを保存する手順と、

10

前記最適配置用ポリシー情報に基づいて実行するサービスを最適な計算機に割当てた処理を実行する手順と、

前記最適配置後の各計算機でのサービス負荷や計算機負荷状況によりサービスの再配置が必要か否かを判定する手順と、

前記負荷管理部によるサービスの再配置が必要との判定結果に従って、前記再配置用ポリシー情報に基づき、前記負荷管理部により再配置が必要と判定されたサービスの実行に最適な計算機を決定する手順と、

20

前記再配置が必要と判定されたサービスの実行を一時的に停止させ、前記サービス配置部により決定された前記再配置が必要と判定されたサービスの実行に最適な計算機で前記再配置が必要と判定されたサービスを起動してサービスをスイッチオーバーする手順と、を前記計算機システムに実行させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般的には、複数の計算機で構成される計算機システムに関し、特に、計算機の障害や負荷状況に応じてサービスの最適配置機能を実現するクラスタシステムの技術に関する。

30

【背景技術】

【0002】

近年、複数の計算機（例えばサーバ）で構成される計算機システムを管理し、アプリケーションプログラムを実行することによりクライアント端末（ユーザ）に提供するサービスの処理性能及び信頼性を高めるクラスタシステムと呼ばれるソフトウェア技術が開発されている。クラスタシステムは、計算機の起動時や障害発生、負荷状況の変動に対応して、計算機システム上で稼動するサービスを最適な計算機にスケジューリングなどの機能を有し、可用性の向上や負荷分散を実現する。

40

【0003】

クラスタシステムは大別して、負荷分散機能を重視した負荷分散型クラスタシステムと、フェイルオーバー機能を重視した高可用性クラスタシステムがある（例えば、非特許文献1及び非特許文献2を参照）。

【0004】

クラスタシステムは、システム運用上のルールに相当するポリシー（policy）情報の設定に基づいて、サービスを実行するのに最適な計算機を決定している。通常では、ポリシー情報はユーザの設定により変更可能である。

【0005】

また、クラスタシステムは、全ての計算機が高負荷状態で、サービスを割当てた最適な

50

計算機が無い場合には、予備用の計算機（プロビジョニング計算機）を利用して対応している。

【非特許文献1】R. Buyya, “High Performance Cluster Computing: Architecture and Systems (Volume 1&2)”, 1999年, Prentice Hall

【非特許文献2】金子哲夫、森良哉、「クラスタソフトウェア」、東芝レビュー、Vol.54 No.12(1999)、p.18-21

【発明の開示】

【発明が解決しようとする課題】

【0006】

近年では、負荷分散型クラスタシステムと高可用型クラスタシステムとが混在するクラスタシステムが開発されている。このようなシステムでは、単純に前記のポリシー情報の設定のみでサービスの最適配置（最適な計算機へのサービスの割当て）がなされると、計算機の負荷状況の変動に応じたサービスの実行が保証されない事態が発生する。具体的には、サービスの自動スイッチオーバーを実行すると、負荷の変動に伴って頻繁にスイッチオーバーが発生したり、低優先度サービスが先に実行していた時の対応が不明であったり、またサービスの実行可能な計算機が無い時には、起動されないことがあった。

【0007】

そこで、本発明の目的は、サービスの最適配置後の動的な負荷状況の変化等のサービスの実行状況に応じたサービスの再配置を確実にこなうことが可能なクラスタシステムを実現することができる。

【課題を解決するための手段】

【0008】

本発明の一観点に従った2台以上の計算機を持つ計算機システムは、前記各計算機が実行する複数のサービスの割当て処理を決定するためのポリシー情報を保存するポリシー管理部と、前記ポリシー情報に従って、各サービスを最適な計算機に割当てて処理を実行するサービス最適配置部と、前記計算機間におけるサービスの実行状況に応じて、前記サービス最適配置部により割当てられたサービスの再配置処理を前記ポリシー情報を参照して実行するサービス再配置部とを具備することを特徴とする。

【0009】

本発明の他の観点によれば、特に、負荷分散型クラスタシステムと高可用型クラスタシステムとが混在する複合クラスタシステムにおいて、動的な負荷状況の変化に応じたクラスタシステム間のサービスの最適配置を可能とする構成を有する計算機システムが構成される。

【発明の効果】

【0010】

本発明によれば、サービスの最適配置後の動的な負荷状況の変化等のサービスの実行状況に応じたサービスの再配置を確実にこなうことが可能なクラスタシステムを実現することができる。

【発明を実施するための最良の形態】

【0011】

以下図面を参照して、本発明の実施形態を説明する。

【0012】

（第1の実施形態）

図1は、第1の実施形態に関する計算機システムのシステム構成を示すブロック図である。

【0013】

計算機システムは、例えば4台の計算機C1～C4がネットワークN上で相互に接続された構成である。各計算機C1～C4はそれぞれ、オペレーティングシステム（OS-1～OS-4）31～34の制御下で動作する。なお、ここでは、1台の予備用の計算機（プロビジョニング計算機）C5がネットワークNを介して計算機システムに接続されている。更に

10

20

30

40

50

1台またはそれ以上の予備用の計算機をネットワークNに接続してもよい。

【0014】

計算機C1～計算機C4によって、クラスタシステムを構成している。このクラスタシステムでは、クラスタ制御部(CS1)10が動作する。クラスタ制御部10は、計算機C1～計算機C4にそれぞれ設けられたクラスタ制御用のプログラム(クラスタソフトウェア)(図示せず)が相互に通信を行いながら同期して一体となって動作することにより実現されるバーチャルマシンである。このため、クラスタ制御部10は計算機C1～計算機C4にまたがって存在していると考えられることができる。クラスタ制御部10は、サービス最適配置機能を実現するサービス最適配置部11、サービス再配置機能を実現するサービス再配置部12、ポリシー管理機能を実現するポリシー管理部13、負荷管理機能を実現する負荷管理部14、及びサービス制御機能を実現するサービス制御部15を有する。

10

【0015】

サービス最適配置部11は、サービスの起動が必要になった場合に、サービス実行に最適な計算機を、ポリシー管理部13により保管されているポリシー情報に従って決定する。ポリシー情報は、具体的には例えば以下のような項目(1)～(5)のポリシー(運用上のルール)を指定する。

【0016】

(1)サービスの優先順位。

【0017】

サービス毎に実行を優先する順位が付けられる。サービスの優先順位に従って、必要なリソース、即ち計算機を割当てる順序が決められる。また、優先順位の高いサービスを実行するために優先順位の低いサービスを停止させることもある。

20

【0018】

(2)サービスに割当てる計算機の優先順位。

【0019】

サービスを実行可能な計算機が複数あるときに優先的に割当てられる計算機の順位をつける。

【0020】

(3)サービス間の関係(排他,依存など)。

【0021】

同時に実行不可能なサービスを排他関係にあるサービスと言い、他のサービスが実行されている時にしか実行できないサービスを依存関係にあるサービスとする。また、同じ計算機で実行不可能なサービスをサーバ排他関係にあるサービスと言い、他のサービスが実行されている時にしか実行できないサービスをサーバ依存関係にあるサービスとする。このようなサービス相互間の関係を設定する。

30

【0022】

(4)サービス実行の必須リソース(周辺機器など)の割当て。

【0023】

サービスを実行するのに必須なリソースを設定し、そのリソースを持つ計算機以外でサービスが実行されないように設定する。

40

【0024】

(5)計算機の負荷状況(最低負荷状況の計算機に割当てる)。

【0025】

サービスを実行するときに最低負荷の計算機を選択し、そのサービスを実行しても過負荷にならない計算機を選択するなどの条件を設定する。

【0026】

サービス再配置部12は、本実施形態の要旨に関する要素であり、サービスの負荷状況の変化や計算機停止に至らない障害発生などで、サービスの計算機配置に不均衡が発生した時に、サービスの再配置をポリシー管理部13により保管されているポリシー情報に従って決定する。

50

## 【 0 0 2 7 】

この再配置に関するポリシー情報は、例えば以下のような項目（１）～（４）のポリシーを指定する。

## 【 0 0 2 8 】

（１）自サービスのスイッチオーバーの可/不可。

## 【 0 0 2 9 】

実行中のサービスを停止し、この停止したサービスを他の計算機で実行を継続させるために他の計算機に移送することをスイッチオーバーと言う。このスイッチオーバーの可・不可の設定をする。これは、予め静的に設定する場合と、クリティカルな処理を実行中に不可に設定する動的な設定をする場合とがある。

10

## 【 0 0 3 0 】

（２）サービスの実行可能なノードが無い場合に他サービス停止の可/不可。

## 【 0 0 3 1 】

あるサービスの起動時にこれを実行可能な計算機がない場合に当該サービスより優先度の低い実行中のサービスを停止してそのサービスを起動させることの可・不可を設定する。

## 【 0 0 3 2 】

この場合、停止されたサービスは他の計算機へスイッチオーバーするように設定される場合がある。これらの設定は、システム全体、サービス単位、あるいは計算機単位で設定することができる。

20

## 【 0 0 3 3 】

（３）スイッチオーバー/停止サービスの判定基準（高負荷優先/低負荷優先）。

## 【 0 0 3 4 】

判断基準の例としては、  
高負荷のサービスから先にスイッチオーバー・停止させる場合、  
低負荷のサービスから先にスイッチオーバー・停止させる場合、  
スイッチオーバー・停止させるサービスの優先順位を設定する場合  
がある。このような設定をシステム単位、計算機単位で設定することが可能である。

## 【 0 0 3 5 】

また、最後に１つだけ残ったサービスのスイッチオーバーの可/不可の設定にはそのサービスの大きさと計算機の容量との関係などを考慮して設定する必要がある。例えば、ある計算機に対しては過負荷となるサービスを同じ程度の容量の計算機にスイッチオーバーしてもやはり過負荷となり、この場合はスイッチオーバーは不可となる。

30

## 【 0 0 3 6 】

（４）負荷状況の変化時の対応。

## 【 0 0 3 7 】

計算機の負荷状況が変化したときにサービスのスイッチオーバー/停止などを実行するかどうか、を設定する。負荷状況は変化の閾値なども設定できる。

## 【 0 0 3 8 】

（４ - １）現状維持重視の場合で、サービスのスイッチオーバー/停止が起きない程度でサービスの再配置を実行する。

40

## 【 0 0 3 9 】

（４ - ２）最適配置重視の場合で、たとえサービスのスイッチオーバー/停止が起きてても最適になるようにサービスを再配置する。

## 【 0 0 4 0 】

このほかに、例えばある計算機にその停止に至らない程度の不具合が生じて一時的にその容量が低下した時なども後で述べるサービス再配置部はその必要性を検知してサービス再配置の処理が行われる。

## 【 0 0 4 1 】

これらのポリシー情報は、予めユーザにより設定が可能である。なお、再配置が決定した

50

サービスは、サービス最適配置部 1 1 によって、実行される計算機が割当てられるまで停止状態になる。

【 0 0 4 2 】

ポリシー管理部 1 3 は、サービス最適配置部 1 1 やサービス再配置部 1 2 が使用するポリシー情報を保存・管理している。

【 0 0 4 3 】

負荷管理部 1 4 は、各計算機 C 1 ~ C 4 でのサービス負荷や計算機負荷状況を判定し、この判定結果によりサービスの再配置が必要な場合には、その旨を負荷情報と共にサービス再配置部 1 2 に通知する。この通知を受けて、サービス再配置部 1 2 は、後述するようなサービスの再配置処理を実行する。

10

【 0 0 4 4 】

当該負荷情報には、各計算機 C 1 ~ C 4 の CPU、メモリ、ディスクの使用量やレスポンスタイムなどが含まれる。また、各計算機 C 1 ~ C 4 はそれぞれ、ロード負荷モニタ 2 1 ~ 2 4 を有し、それぞれの負荷状況を監視している。

【 0 0 4 5 】

( クラスタシステムの動作 )

クラスタシステム 1 0 は、ユーザが作成した並列実行型サービス及び高可用性サービスの実行を管理する。並列実行型サービスは、例えば Web サービスなどであり、一時期に複数の計算機 C 1 ~ C 4 で同時に実行可能なタイプのサービスである。並列実行型サービスが一時期に実行されるサービス数は、負荷管理部 1 4 により管理されており、高負荷になればサービス数が増大し、低負荷になればサービス数が減少する。

20

【 0 0 4 6 】

一方、ユーザが作成した高可用性サービスは、例えばデータベース検索サービスなどであり、一時期にどこか一つの計算機 ( 例えば C 2 ) でのみ実行可能なタイプのサービスである。高可用性サービスは、障害発生時にフェイルオーバや、障害予測時や高負荷時のスイッチオーバで他の計算機に移動して処理を継続するように作成されている。

【 0 0 4 7 】

例えば計算機 C 2 で実行中の高可用性サービスの負荷が急激に上昇したときに、クラスタシステム 1 0 の負荷管理部 1 4 は、当該計算機 C 2 の負荷が限界に近いと判定すると、サービスの再配置の必要をサービス再配置部 1 2 に通知する。

30

【 0 0 4 8 】

サービス再配置部 1 2 は、ポリシー管理部 1 3 で保管されているポリシー情報 ( ユーザにより設定可能 ) に従って、高可用性サービスまたは並列実行型サービスのサービス再配置処理を開始する。

【 0 0 4 9 】

具体的には、サービス再配置部 1 2 は例えば並列実行型サービスの再配置を決定すると、これを受けて、サービス制御部 1 5 が一時的に並列実行型サービスの停止を行なう。この並列実行型サービスの停止後に、サービス最適配置部 1 1 は、サービス実行に最適な計算機 ( 例えば C 1 ) を選択する。選択された計算機 ( 例えば C 1 ) 上のサービス制御部 1 5 は、並列実行型サービスを起動させることで、サービスの自動スイッチオーバを実行する。

40

【 0 0 5 0 】

以上のようなクラスタシステム 1 0 によるサービス自動スイッチオーバ機構により、動的な負荷変動に対応したサービスの最適配置が可能になる。

【 0 0 5 1 】

( サービス配置処理 )

以下図 2 のフローチャートを参照して、本実施形態のクラスタシステム 1 0 のサービス再配置処理の手順を説明する。

【 0 0 5 2 】

サービス再配置部 1 2 は、ポリシー管理部 1 3 に問い合わせを実行して、例えばユーザに

50

より設定されたポリシー情報の設定に従って再配置処理を実行する。ポリシー情報は、前述したように、例えば以下のような項目(1)～(4)のポリシーを指定する。

【0053】

(1) サービス毎のスイッチオーバーの可/不可。

【0054】

(2) サービスの実行可能なノードが無い場合に他サービス停止の可/不可。

【0055】

(3) スwitchオーバー/停止サービスの判定基準。

【0056】

(3-1) 高負荷優先/低負荷優先。

10

【0057】

(3-2) 最後のサービスのスイッチオーバーの可/不可。

【0058】

(4) 負荷状況の変化時の対応。

【0059】

(4-1) 現状維持重視の場合で、サービス停止が起きない程度で再配置。

【0060】

(4-2) 最適配置重視の場合で、サービス停止を起こしながら再配置。

【0061】

前述したように、負荷管理部14は、負荷状況の判定に応じて、サービスの再配置が必要であるか否かを判定する(ステップS1)。この判定基準としては、例えば「計算機が継続的に高負荷でサービス実行の遅延が予測される場合」や、「計算機に実行待ちになっている高負荷(予測)の高優先度サービスがある場合」等であり、サービスの再配置が必要であると判断される。

20

【0062】

以下、サービスの再配置が必要な場合の処理(ステップS1のYES)を説明する。

【0063】

サービス再配置部12は、ポリシー情報のポリシー(1及び3)に従って、サービスのスイッチオーバーや、停止可能なサービスがあるか否かを判定する(ステップS2)。判定結果が「YES」であれば、クラスタシステム10のサービス制御部15は、スイッチオーバー可能と設定されたサービスより、優先度の低いものからサービスの再配置が必要なくなるまでサービスのスイッチオーバーを実行する(ステップS3)。

30

【0064】

一方、スイッチオーバー可能なサービスが無い場合は、サービス再配置部12は、ポリシー情報のポリシー(2)に従って、強制処置が可能であるか否かを判定する(ステップS2のNO, S4)。強制処置が可能であれば、優先度の低いものからサービスの再配置が必要なくなるまでスイッチオーバーを実行する処理に移行する(ステップS4のYES, S3)。

【0065】

強制処置できない場合は、クラスタシステム10は、利用可能なプロビジョニング計算機(予備計算機)を探索し、存在する場合には当該計算機C5を追加する(ステップS4のNO, S5, S6)。ここで、追加されたプロビジョニング計算機C5は、計算機システムの負荷が低下したときに返却の指定がある場合には、当該負荷が低下したときに返却される。なお、利用可能なプロビジョニング計算機が存在しない場合には、一定時間のスリープ状態を経てリターンとなる(ステップS5のNO, S11)。

40

【0066】

次に、負荷管理部14の判定結果により、サービスの再配置が不必要である場合について説明する(ステップS1のNO)。

【0067】

サービス再配置部12は、ポリシー情報のポリシー(4-2)に従って、最適化配置重視で

50



高負荷になりつつある場合には、サービス再配置処理を実行する（ステップS7のYES，S8のYES）。そうでなければ、サービス再配置処理は終了となる（ステップS7のNO，S8のNO）。

【0068】

ここで、計算機が高負荷になりつつあるか否かの判定は、一定の間隔で平均した負荷が単調に増加していて、遠くない将来において高負荷になることが予測できるか否かで判定できる。

【0069】

さらに、サービス再配置処理を実行する場合に、サービス再配置部12は、サービスを移動した方がより最適な配置かどうかを判定し、最適な場合にはサービスのスイッチオーバーを実行する（ステップS9のYES，S10）。最適な配置であると判断できない場合には、サービス再配置処理は終了となる（ステップS9のNO）。

10

【0070】

ここで、最適な配置の判断基準は、選択された計算機で再配置するサービスを現在と同じ負荷で稼働させた場合、計算機間の負荷の状態がより平均化される場合である。また、サービスのスイッチオーバーのオーバーヘッドを加味しても、選択された計算機で処理を行う方が早いと考えられる場合などである。

【0071】

ここで、サービス再配置のポリシーとして、サービス毎にスイッチオーバーの可/不可や現状維持重視のポリシーが出来ることや、スイッチオーバーで停止しても、スイッチオーバー先の計算機で起動可能にならないかぎり実行されないことで、計算機の負荷変動に過敏に反応して、スイッチオーバーを繰り返す事を防止することが可能になる。

20

【0072】

以上要するに、本実施形態のクラスタシステムであれば、ポリシーベースで管理されたサービス再配置機能を持たせることで、動的な負荷状況の変化に応じてサービスの再配置を可能とし、かつ、ユーザの運用環境に合ったクラスタシステムの構築を容易に実現することが可能となる。

【0073】

（第2の実施形態）

図3から図5は、第2の実施形態に関する計算機システムのシステム構成及びその変化を示すブロック図である。

30

【0074】

図3に示すように、初期状態での計算機システムは、例えば5台の計算機C1～C5がネットワークN上で相互に接続された構成である。さらに、ネットワークN上には6台目の計算機C6が接続されている。当該計算機C6は、停止しており、プロビジョニング計算機（予備計算機）としてプロビジョニング計算機プール60に登録されている。

【0075】

プロビジョニング計算機プール60とは、停止している1台又は複数台の計算機をプロビジョニング計算機として登録したことを概念的に図示し総称したものである。

【0076】

40

計算機をプロビジョニング計算機としてプロビジョニング計算機プール60に登録することは、図示しないプロビジョニング計算機に関する情報（例えばプロセッサ名やMACアドレスなど）を登録情報として登録することを意味し、この登録情報によりプロビジョニング計算機プール60に登録された複数のプロビジョニング計算機を管理する。

【0077】

計算機C1～C3はそれぞれオペレーティングシステムOS（OS-1-1～OS-1-3）の制御下で稼働中である。また、計算機C4，C5はそれぞれオペレーティングシステムOS（OS-2-1,OS-2-2）の制御下で稼働中である。

【0078】

稼働中の計算機C1～C5では、プロビジョニング計算機割当て機能を実現するプロビ

50

ジョーニング計算機割当て部 3 1 と、プロビジョニング計算機切離し機能を実現するプロビジョニング計算機切離し部 3 2 と、プロビジョニングポリシー管理機能を実現するプロビジョニングポリシー管理部（以下単にポリシー管理部との略す場合がある）3 3 とが稼働している。計算機 C 1、計算機 C 2、計算機 C 3 でそれぞれプロビジョニング計算機割当て部 3 1 と、プロビジョニング計算機切離し部 3 2 と、プロビジョニングポリシー管理部 3 3 とが稼働し相互に通信を行いながら同期をとって連携することで、計算機 C 1、計算機 C 2、計算機 C 3 がクラスタシステム C S 1 を構成する。符号 3 0 は、クラスタシステム C S 1 を模式的に図示している。一方、計算機 C 4、計算機 C 5 でそれぞれプロビジョニング計算機割当て部 3 1 と、プロビジョニング計算機切離し部 3 2 と、プロビジョニングポリシー管理部 3 3 とが稼働し相互に通信を行いながら同期をとって連携することで、計算機 C 4、計算機 C 5 がクラスタシステム C S 2 を構成する。符号 4 0 は、クラスタシステム C S 2 を模式的に図示している。これらの各クラスタシステムは、相互に無関係であり、相互にサービスを関係付けたりすることは無い。

10

#### 【 0 0 7 9 】

本計算機システムには、ストレージエリアネットワーク S A N（Storage Area Network）4 5 を介して複数のストレージ装置（ディスク装置）5 0 ~ 5 7、7 0 が接続されている。

#### 【 0 0 8 0 】

本計算機システムでは、各計算機を起動するためのブートイメージを、ストレージ装置（ディスク装置）5 0 ~ 5 7 に予め記憶させ登録している。ここでブートイメージとは、計算機を起動するためのオペレーティングシステム及びこのオペレーティングシステムで実行可能なアプリケーションプログラムを含んでいる。

20

#### 【 0 0 8 1 】

各ストレージ装置 5 0 ~ 5 3 及び 5 4 ~ 5 7 には、それぞれブートイメージ O S - 1 - 1、O S - 1 - 2、O S - 1 - 3、O S - 1 - 4、O S - 2 - 1、O S - 2 - 2、O S - 2 - 3、O S - 2 - 4 が登録されている。例えば計算機 C 3 を起動させるためのブートイメージ（O S - 1 - 3）をストレージ装置 5 2 上に登録している。計算機 C 3 をこのブートイメージ（O S - 1 - 3）を用いて起動させると、当該計算機 C 3 は、O S（O S - 1 - 3）によりその動作が制御される稼働計算機となる。図 3 において、どの計算機がどのブートイメージで起動したかを矢印で図示した。

30

#### 【 0 0 8 2 】

一方、図 5 に示すように、計算機 C 3 を起動させるためのブートイメージ（O S - 2 - 4）をストレージ装置 5 7 上に登録している。計算機 C 3 をこのブートイメージ（O S - 2 - 4）を用いて起動させると、当該計算機 C 3 は、O S（O S - 2 - 4）によりその動作が制御される稼働計算機となる。図 5 において、どの計算機がどのブートイメージで起動したかを矢印で図示した。

#### 【 0 0 8 3 】

（クラスタシステムの動作）

プロビジョニング計算機割当て部 3 1 は、クラスタ制御部 3 0、4 0 で実行する計算機が必要になった場合、ポリシー管理部 3 3 を介してアクセス可能なプロビジョニングポリシーデータベース（以下ポリシー D B と略す）7 0 に蓄積されたプロビジョニングポリシー情報に従って、プロビジョニング計算機をクラスタシステムに割当てる。

40

#### 【 0 0 8 4 】

プロビジョニング計算機切離し部 3 2 は、クラスタ制御部 3 0、4 0 で実行する計算機に余剰が発生した場合、ポリシー管理部 3 3 を介してアクセス可能なポリシー D B 7 0 に従って、クラスタシステム内の計算機を切離し、プロビジョニング計算機としてプール 6 0 に登録する。

#### 【 0 0 8 5 】

ポリシー管理部 3 3 は、プロビジョニングポリシー情報（以下単にポリシー情報と略す場合がある）の設定/参照機能を提供する。当該ポリシー情報は、例えば以下のような項目（1）

50

～(4)のプロビジョニングポリシーを指定する。

【0086】

(1) クラスタシステム毎の計算機割当てレベル(優先度)。

【0087】

同時に二つ以上のクラスタシステムからプロビジョニング計算機要求が来た場合、優先的に割り当てるクラスタシステムの順位(優先度)の設定をする。クラスタシステムからの要求があったときに必要なプロビジョニングノードが無い場合に、優先度の低いクラスタシステムに割当てられた計算機を強制的に要求のあったクラスタシステムに割り当てる場合もある。

【0088】

(2) 提供計算機の返還の可/不可。

【0089】

クラスタシステムにおいて割当てられたプロビジョニング計算機をプロビジョニングプールに変換することが可能か否かを設定する。従って、この設定で不可の場合は、そのクラスタシステム内の割当て計算機数は増加する一方となる。

【0090】

(3) 提供計算機の強制返還の可/不可。

【0091】

プロビジョニングプールよりクラスタシステムへ提供されている計算機を強制的に返還させることができるか否かを設定する。すなわち、強制的に返還させてもシステムの運用に支障がないかなどを設定の際の条件とする。例えば、優先度の高いクラスタシステムより要求があったときにプロビジョニングプールに予備の計算機がない場合には、優先度の低いクラスタシステムへ強制返還の要求が行くように設定される。

【0092】

(4) 提供計算機数の指標(必須計算機数,最大計算機数,初期計算機数)。

【0093】

クラスタシステムを構成するために必要な計算機数を必須計算機数とする。クラスタシステムに割当て可能な最大の計算機数を最大計算機数とする。また、クラスタシステムの起動時に最適な割当て計算機数を初期計算機数とする。このように、クラスタシステムへ提供する計算機数を決定する際の指標を設定することができる。

【0094】

ポリシー情報は、通常では、ユーザが計算機システムの構築/保守時に、ポリシーDB70に設定される。

【0095】

図8には、図3に示したクラスタシステムにおける各計算機に登録するためのプロビジョニングDB70に登録したプロビジョニングポリシー情報の一例を示す。

【0096】

(プロビジョニング計算機割当て処理)

以下図6のフローチャートを参照して、本実施形態のプロビジョニング計算機割当て処理の手順を説明する。

【0097】

まず、図3に示すように、初期状態での計算機システムは、計算機C1～C3が稼動中であり、クラスタシステム(CS1)30が動作中である。また、計算機C4, C5が稼動中であり、クラスタシステム(CS2)40が動作中である。さらに、計算機C6は、停止しており、プロビジョニング計算機としてプール60に登録されている。

【0098】

ここで、クラスタシステム(CS2)40の負荷が増大し、2台の計算機C4, C5では処理できない状況になると、クラスタシステム(CS2)40はプロビジョニング計算機割当て部41に計算機追加を要請する(ステップS21のYES)。

【0099】

10

20

30

40

50

プロビジョニング計算機割当て部 4 1 は、プロビジョニング計算機プール 6 0 を検索し、登録されている計算機 C 6 を取り出し、要求されたクラスタシステム (CS 2) 4 0 に追加する (ステップ S 2 3 の YES, S 2 4)。ここで、プロビジョニング計算機割当て部 4 1 は、図 4 に示すように、クラスタシステム (CS 2) 4 0 に所属するブートイメージの中で、使用されていないブートイメージ (OS-2-3) をストレージ装置 5 6 から取り出し、計算機 C 6 に接続して起動させる。

**【 0 1 0 0 】**

但し、クラスタシステム (CS 2) 4 0 から、ブートイメージの満たすべき要件が、詳細に指定された場合は、その要件に合うブートイメージを検索することになる。

10

**【 0 1 0 1 】**

ところで、2つのクラスタシステム 3 0, 4 0 から同時に、計算機追加の要求がなされた場合には、プロビジョニング計算機割当て部 3 1, 4 1 は、ポリシ管理部 3 3, 4 3 を介してポリシ DB 7 0 をアクセスし、ポリシ情報に従って計算機割当てレベルの大きいクラスタ制御部を選択する (ステップ S 2 2)。そして、例えばクラスタシステム (CS 2) 4 0の方が割当てレベルが大きい場合には、プロビジョニング計算機割当て部 4 1 は、プロビジョニング計算機プール 6 0 を検索し、登録されている計算機 C 6 を優先的に割当てる (ステップ S 2 3 の YES, S 2 4)。

**【 0 1 0 2 】**

さらに、クラスタシステム (CS 2) 4 0 の負荷がさらに増大し、3台の計算機 C 4 ~ C 6 でも処理ができなくなると、クラスタ制御部 4 0 は、プロビジョニング計算機割当て部 4 1 に計算機追加を要請する。

20

**【 0 1 0 3 】**

プロビジョニング計算機割当て部 4 1 は、プロビジョニング計算機プール 6 0 には計算機が登録されていないため、前記のポリシ情報に従って強制返還可能なクラスタ制御部が存在するか否かを判断する (ステップ S 2 3 の NO, S 2 5)。存在しない場合には、一定時間のスリープ状態を経て、計算機がプール 6 0 に登録されるまで待機状態となる (ステップ S 2 5 の NO, S 2 6)。

**【 0 1 0 4 】**

一方、例えばクラスタシステム (CS 1) 3 0 が強制返還可能な場合には、プロビジョニング計算機割当て部 4 1 は、当該クラスタシステム (CS 1) 3 0 上の計算機に強制返還を要求する (ステップ S 2 5 の YES)。強制返還を要求されたクラスタシステム (CS 1) 3 0 の計算機上のプロビジョニング計算機切離し部 3 2 は、切り離し可能な計算機 (例えば C 3) を決定し、プロビジョニング計算機としてプロビジョニング計算機プール 6 0 に登録する (ステップ S 2 7)。

30

**【 0 1 0 5 】**

クラスタシステム (CS 1) 3 0 から切離された計算機 C 3 がプロビジョニング計算機プール 6 0 へ登録されると、クラスタシステム (CS 2) 4 0 のプロビジョニング計算機割当て部 4 1 は、プロビジョニング計算機プール 6 0 を検索し、登録されている計算機 C 3 を取り出して割当てる (ステップ S 2 3 の YES, S 2 4)。

40

**【 0 1 0 6 】**

プロビジョニング計算機割当て部 4 1 は、図 5 に示すように、クラスタシステム (CS 2) 4 0 に所属するブートイメージの中で、使用されていないブートイメージ (OS-2-4) をストレージ装置 5 7 から取り出し、計算機 C 3 に接続して起動させる。

**【 0 1 0 7 】**

(プロビジョニング計算機切離し処理)

次に、図 7 のフローチャートを参照して、本実施形態のプロビジョニング計算機切離し処理の手順を説明する。

**【 0 1 0 8 】**

ここでは、クラスタシステム (CS 1) 3 0 のプロビジョニング計算機切離し部 3 2 は

50

、計算機切り離し要求を受けると、ポリシー情報に従って、クラスタシステム（CS1）30上の切離し可能な計算機（ここではC3）を決定する（ステップS31のYES，S33）。

【0109】

さらに、プロビジョニング計算機切離し部32は、決定した計算機C3で稼働中のサービスにスイッチオーバー要求を出す（ステップS34）。ここで、クラスタ制御部30において、ポリシー情報に従って、切離し条件として全サービスの停止待ちの場合には、プロビジョニング計算機切離し部32は、全サービスの停止を待って、計算機C3を切離して、プロビジョニング計算機としてプロビジョニング計算機プール60に登録する（ステップS35のYES，S37，S38）。

10

【0110】

一方、切離し条件として全サービスの停止待ちではない場合には、プロビジョニング計算機切離し部32は、切離し準備として一定時間だけ待って、計算機C3を切離して、プロビジョニング計算機としてプロビジョニング計算機プール60に登録する（ステップS35のNO，S36，S38）。

【0111】

以上のように本実施形態によれば、複数のクラスタシステムからプロビジョニング計算機の追加要求があった場合に、ポリシー情報に従って、例えば強制返還が設定されているクラスタシステム（CS1）30から、相対的に計算機割当てレベルの高いクラスタシステム（CS2）40へ、計算機を切離して割当てる処理を実行できる。要するに、クラスタシステム毎にプロビジョニングポリシーを設定可能なプロビジョニング計算機の割当て/切離し機能を持つことにより、クラスタシステム間で計算機割当てレベルに基づいた最適な計算機の割当て（移動）が可能となる。このようなクラスタシステムと、例えば課金システムとを連動させることで、ネットワークサービスでの高度なSLA（service level agreement）等を実現するシステムを構築することが可能になる。

20

【0112】

この実施形態の種々の実施の態様をまとめると次のようになる。

【0113】

（1）2台以上の計算機が接続されて、2つ以上のクラスタシステムを実現する計算機システムにおいて、

30

前記各クラスタシステムが共通して使用可能な少なくとも1つのプロビジョニング計算機と、

プロビジョニング計算機の割当て処理又は切離し処理のポリシーを指定するためのポリシー情報を変更可能に保存するポリシー管理手段と、

前記ポリシー情報に従って、前記少なくとも1つのプロビジョニング計算機から追加要求の計算機を割当てる割当て処理または余剰な計算機を切離す切離し処理を実行する割当て/切離し手段と

を具備した計算機システム。

【0114】

（2）前記割当て/切離し手段は、前記ポリシー情報に従って少なくとも1つのプロビジョニング計算機として登録されている計算機、または他のクラスタシステムで使用されている計算機を、必要なクラスタシステムに割当てる（1）項に記載の計算機システム。

40

【0115】

（3）前記割当て/切離し手段は、前記ポリシー情報に従ってクラスタシステムで使用されている計算機を切離し、前記少なくとも1つのプロビジョニング計算機として登録する（1）項に記載の計算機システム。

【0116】

（4）前記ポリシー管理手段は、前記ポリシー情報を変更可能に保存するデータベースを管理し、前記各計算機からのアクセスに応じて、当該データベースから前記ポリシー情報の取出し、または設定を行なう（1）項に記載の計算機システム。

50

## 【 0 1 1 7 】

( 5 ) 2 台以上の計算機が接続された計算機システムに適用し、2 つ以上のクラスタシステムのそれぞれに含まれるプログラムであって、

変更可能なポリシー情報に従って、各クラスタシステムが共通して使用可能な少なくとも1 つのプロビジョニング計算機から追加要求の計算機を割当る処理を実行する手順と、

前記ポリシー情報に従って、各クラスタシステムが使用している少なくとも1 つのプロビジョニング計算機を切離す処理を実行する手順と

を前記計算機システムに実行させるためのプログラム。

## 【 0 1 1 8 】

なお、本発明は上記実施形態そのままに限定されるものではなく、実施段階ではその要旨を逸脱しない範囲で構成要素を変形して具体化できる。また、上記実施形態に開示されている複数の構成要素の適宜な組み合わせにより、種々の発明を形成できる。例えば、実施形態に示される全構成要素から幾つかの構成要素を削除してもよい。さらに、異なる実施形態にわたる構成要素を適宜組み合わせてもよい。

## 【 図面の簡単な説明 】

## 【 0 1 1 9 】

【 図 1 】 本発明の第 1 の実施形態に関するシステム構成を示すブロック図。

【 図 2 】 第 1 の実施形態に関するサービス再配置処理の手順を説明するためのフローチャート。

【 図 3 】 第 2 の実施形態に関するシステム構成を示すブロック図。

【 図 4 】 第 2 の実施形態に関するシステム構成の変化を示すブロック図。

【 図 5 】 第 2 の実施形態に関するシステム構成の変化を示すブロック図。

【 図 6 】 第 2 の実施形態に関するプロビジョニング計算機の割当て処理の手順を説明するためのフローチャート。

【 図 7 】 第 2 の実施形態に関するプロビジョニング計算機の切離し処理の手順を説明するためのフローチャート。

【 図 8 】 第 2 の実施形態に関するプロビジョニングポリシー情報の一例を示す図。

## 【 符号の説明 】

## 【 0 1 2 0 】

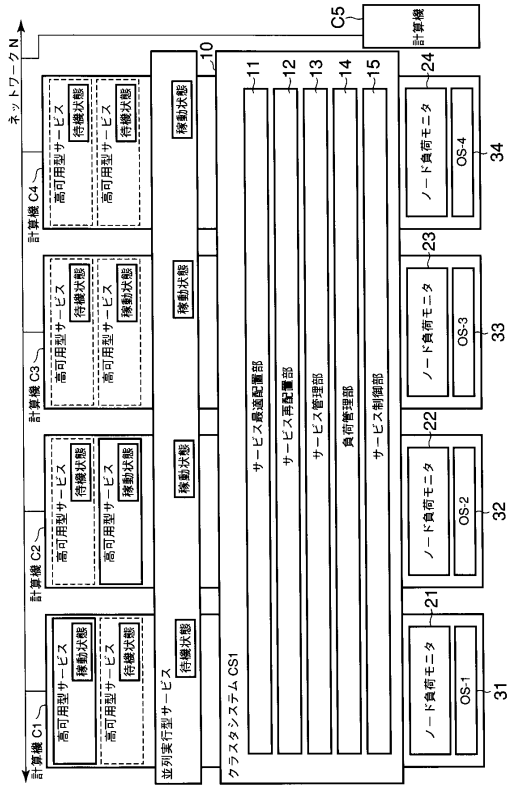
- 1 0 ... クラスタ制御部、 1 1 ... サービス最適配置部、
- 1 2 ... サービス再配置部、 1 3 ... ポリシ管理部、 1 4 ... 負荷管理部、
- 1 5 ... サービス制御部、 3 0 ... クラスタシステム C S 1 、
- 3 1 ... プロビジョニング計算機割当て部、
- 3 2 ... プロビジョニング計算機切離し部、 3 3 ... プロビジョニングポリシ管理部、
- 4 0 ... クラスタシステム C S 2 、
- 5 0 ~ 5 7 ... ストレージ装置 ( ブートイメージ登録ディスク ) 、
- 6 0 ... プロビジョニング計算機プール、
- 7 0 ... プロビジョニングポリシデータベース ( ポリシ D B ) 、 C 1 ~ C 6 ... 計算機。

10

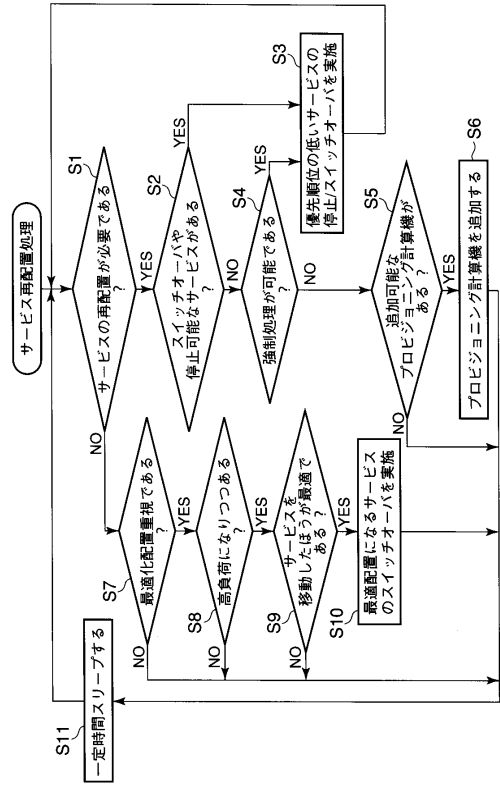
20

30

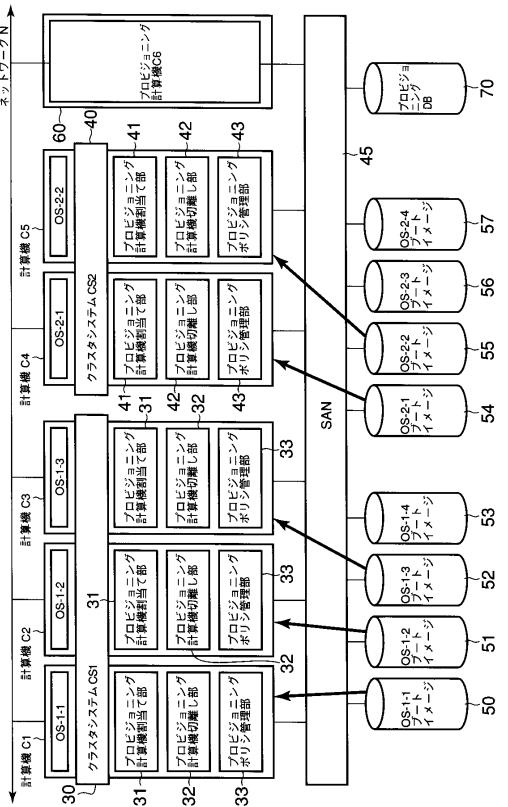
【 図 1 】



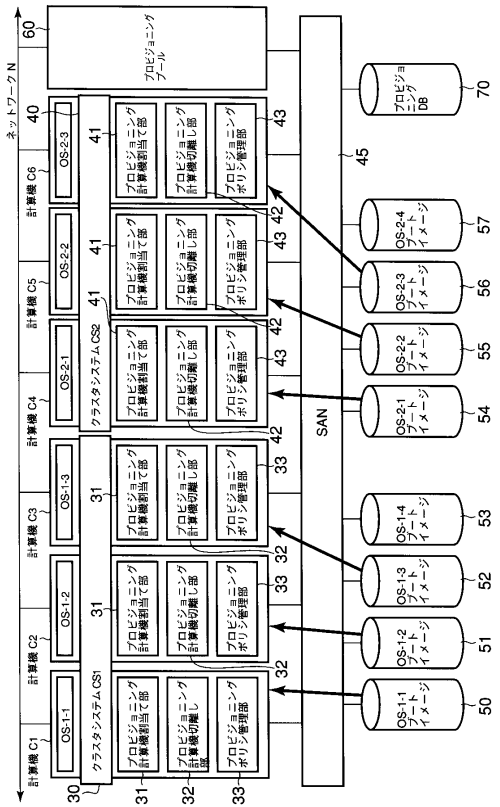
【 図 2 】



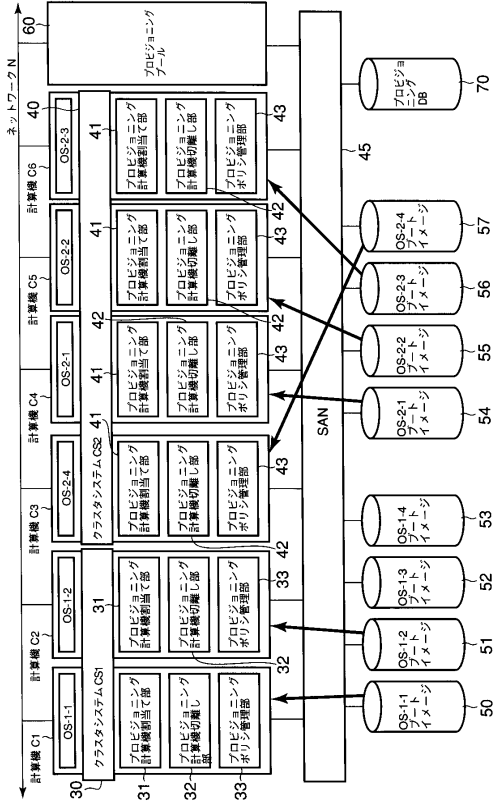
【 図 3 】



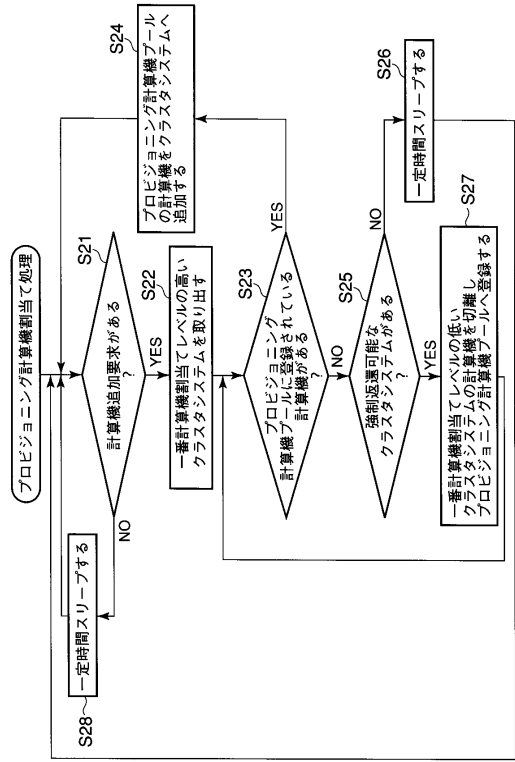
【 図 4 】



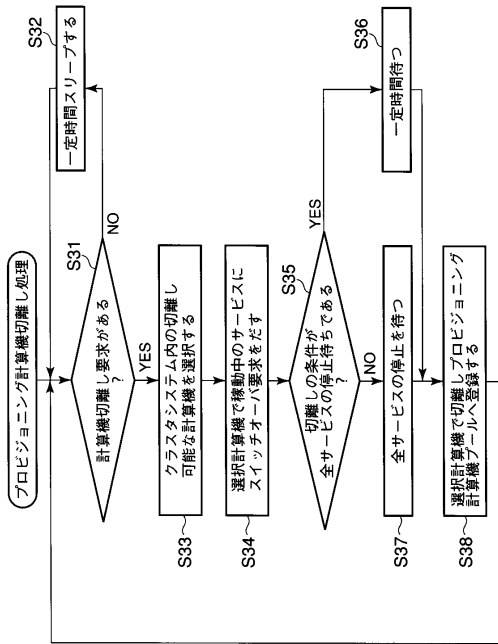
【 図 5 】



【 図 6 】



【 図 7 】



【 図 8 】

	クラスシステム CS1	クラスシステム CS2
計算割り当てレベル	2	1
提供計算機の返還	可	可
提供計算機の強制返還	可	不可
提供計算機数の指標	必須計算機数	2
	最大計算機数	4
	初期計算機数	3
		2



---

フロントページの続き

(74)代理人 100084618

弁理士 村松 貞男

(74)代理人 100092196

弁理士 橋本 良郎

(72)発明者 溝口 研一

東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

審査官 鈴木 修治

(56)参考文献 特開2000-137692(JP,A)

特開平06-068052(JP,A)

特開2003-067351(JP,A)

特開2002-108839(JP,A)

特開平08-152903(JP,A)

特開2002-007364(JP,A)

オープンソースでコストの常識を打ち破れ! 基幹システム再構築セミナー2003レポート WORKSHOP: 東芝, 日経システム構築 no. 123 SYSTEM INTEGRATION, 日本, 日経BP社, 2003年 6月25日, 第123号, 76~76

(58)調査した分野(Int.Cl., DB名)

G06F 9/46 - 9/54