

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4405511号  
(P4405511)

(45) 発行日 平成22年1月27日(2010.1.27)

(24) 登録日 平成21年11月13日(2009.11.13)

(51) Int.Cl.

F I

G 0 6 F 11/20 (2006.01)

G 0 6 F 11/20 3 1 0 F

G 0 6 F 15/177 (2006.01)

G 0 6 F 15/177 A

請求項の数 5 (全 17 頁)

(21) 出願番号	特願2006-521850 (P2006-521850)	(73) 特許権者	390009531
(86) (22) 出願日	平成16年6月30日 (2004.6.30)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公表番号	特表2007-500895 (P2007-500895A)		I N T E R N A T I O N A L B U S I N E S S M A S C H I N E S C O R P O R A T I O N
(43) 公表日	平成19年1月18日 (2007.1.18)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
(86) 国際出願番号	PCT/US2004/020906		
(87) 国際公開番号	W02005/015346		
(87) 国際公開日	平成17年2月17日 (2005.2.17)	(74) 代理人	100108501
審査請求日	平成19年6月26日 (2007.6.26)		弁理士 上野 剛史
(31) 優先権主張番号	10/604,585	(74) 代理人	100112690
(32) 優先日	平成15年7月31日 (2003.7.31)		弁理士 太佐 種一
(33) 優先権主張国	米国 (US)	(74) 代理人	100091568
早期審査対象出願			弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 複数のサービス・ポイントを有する自律コンピューティングにおける動的に構成可能な耐障害性

## (57) 【特許請求の範囲】

## 【請求項 1】

複数のプロセッサを有するシステムであって、各々が1つのサービス・ポイントを持つ少なくとも1つのクラスタを備えたシステム、を提供するように該システムを構成する方法であって、

前記システムにおける各プロセッサから他のプロセッサまでの距離を計算するステップと、

各合計距離が1つのプロセッサに関連する、複数の合計距離を計算するステップと、

前記複数の合計距離から最小合計距離を決定するステップと、

前記最小合計距離に関連するプロセッサをサービス・ポイントとして割り当てるステップと、

前記システムを複数のクラスタに区分するステップであって、

各プロセッサに関連した合計距離に従ってプロセッサを分類するステップと、

前記分類された各プロセッサを分類順に2つのクラスタの一方に割り当てるステップと、

各クラスタにおけるプロセッサに対する最小合計距離を、当該クラスタにおけるプロセッサに関連した複数の合計距離に従って決定するステップと、

各クラスタに対するサービス・ポイントとして、当該クラスタにおける最小合計距離に関連するプロセッサを割り当てるステップと、

を含むステップと、

前記 2 つのクラスタの 1 つを 2 つのクラスタに分割し、それによって前記システムを 3 つのクラスタに区分するステップと、

前記 3 つのクラスタにおけるプロセッサに関連した複数の合計距離に従って、前記 3 つのクラスタの各々におけるプロセッサに対する最小合計距離を決定するステップと、

前記最小合計距離に対応するクラスタ構成を形成するように前記 3 つのクラスタにプロセッサを割り当てるステップと、

前記 3 つのクラスタの各々に対するサービス・ポイントとして、当該クラスタにおける最小合計距離に関連するプロセッサを割り当てるステップと、  
を含む、方法。

【請求項 2】

10

前記構成する動作は、プロセッサがシステムに加えられるとき、動的に実行される、請求項 1 に記載の方法。

【請求項 3】

前記構成する動作は、プロセッサがシステムから取り除かれるとき、動的に実行される、請求項 1 に記載の方法。

【請求項 4】

前記システムの区分は、プロセッサがシステムから取り除かれるとき、動的に変更される、請求項 3 に記載の方法。

【請求項 5】

別のプロセッサをバックアップ・サービス・ポイントとして割り当てるステップを更に含む、請求項 1 に記載の方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、自律コンピューティング (autonomic computing) に関し、詳しく云えば、効率的な耐障害性の自己構成 (self-configuring) および自己修正 (self-healing) オペレーションのためのシステムにおいて、プロセッサをクラスタ化し、サービス・ポイントを割り当てるための方法に関するものである。

【背景技術】

【0002】

30

自律コンピューティングは、一般に、自己監視し、自己構成し、耐障害性であり、および自己修正するマルチプロセッサ・コンピュータ・システムを意図したものを指すが、それはかなり原理的および実用的な観点からの見方である。良好な自律コンピューティング・システムを構築する場合の 1 つの重要な考察点は、自己修正機構を拡張するためにシステム自体にそのシステムの耐障害性を組み込むことである。自己修正機構は、障害時に、システムが直ちにその障害の性質を検出し、その障害の是正を試みることを必要とするであろう。それが障害を是正できなかった場合、システムは、障害を起こしたプロセッサのタスクを 1 つまたは複数の他のプロセッサに割り振ることによって、続いて起るパフォーマンス低下を最小限にするであろう。一般的なコンピュータ・アーキテクチャでは、障害検出および管理に関するこのタスクは、プロセッサの 1 つによって、またはマスタ・プロセッサによって行われる。

40

【0003】

自律コンピューティング・システムを構成するプロセッサは地理的に広い領域に分散されることがある。更に、プロセッサは多くの種々のタイプのものがあり、多くの種々のタイプのオペレーティング・システムを稼動させ、分散ネットワークによって接続されることがある。種々のプロセッサが地理的にクラスタ状態で配列されることが多い。そのような配列は、1 つのマスタ・プロセッサにシステム全体の耐障害性を管理させることを許容するものではない。従って、幾つかのプロセッサに障害管理を行わせることは有益である。本願では、これらのプロセッサをサービス・ポイントと呼ぶことにする。

【0004】

50

図 1 には、サービス・ポイントを利用する代表的なシステムが示される。システム 1 は複数の相互接続されたプロセッサ 2 を含み、それらのプロセッサの 1 つがサービス・ポイント 10 となるように割り当てられる。一般に、サービス・ポイントは、そのプロセッサが他のプロセッサまでの最小距離を有するものとなるように選ばれる。本願では、「距離」という用語は、プロセッサ間の必要な通信時間の尺度として使用される。サービス・ポイントは、それ自身の正規のコンピューティング負荷に加えて、次のようないくつかのタスクを有する。

- ( 1 ) システム内の別の場所における障害プロセッサを検出する、
- ( 2 ) そのプロセッサのタスクを他のプロセッサに再割り振りすることによって障害プロセッサを置換する、
- ( 3 ) 他のプロセッサによって実行されるタスクを監視する、および
- ( 4 ) 最適のパフォーマンスを保証するためにシステムにおける負荷を平衡させる。

10

図 2 は、プロセッサ 3 がサービス・ポイント 10 によって障害を検出されたシステムから取り除かれてしまったこと、および、システムにおける残りのプロセッサが動作を継続していることを示す。

#### 【 0 0 0 5 】

冗長な計算を使用する耐障害性が時々使用されるが、現用の自律計算システムの自己修正および自己構成のフィーチャは、例えば、次のようないくつかの新たな問題点を提起している。

( 1 ) すべてのプロセッサ（遠隔地に設置されたものを含む）が同様のものであり且つ交換可能であるとき、自律システムの自己構成可能および自己調整可能フィーチャは非常に良好に作用する。これは、サービス・ポイントが特別のプロセッサである必要がなく、むしろ、余分な負荷を操作しているプロセッサの同じセットから選択されることを意味する。

20

( 2 ) 通常、並列およびスケーラブル・コンピュータ・アーキテクチャでは、プロセッサ・ポイントの数が一定であり、アクティブ・プロセッサの数分の一として指定することができない。しかし、少な過ぎるサービス・ポイントを持つことは自己修正機構を低速にし過ぎる。多過ぎるサービス・ポイントを持つことはシステムの全体的なパフォーマンスを低下させる。

( 3 ) 自律コンピューティング・システムは動的な環境で働くので、システム・パフォーマンスを最適化するためにクラスタ化およびサービス・ポイントの割当てを動的に最適化することが重要である。オンデマンド・コンピューティング環境では、プロセッサの合計数（従って、クラスタの構成およびサービス・ポイントの割当て）がコンピューティング負荷に応答して常に変化している。

30

#### 【 0 0 0 6 】

自律コンピューティング・システムの自己構成環境では、サービス・ポイントを事前割当てすることは一般に不可能である。従って、その状況に関する要件に従って、いずれの現用のプロセッサもサービス・ポイントになるように動的に割り当てることが可能である。一方、多過ぎるサービス・ポイントを設けることはシステムにおける大きい計算上の負荷を導くことになる。従って、サービス・ポイントの数を稼動プロセッサの数の数分の一に限定して保つことが望ましい。

40

#### 【 0 0 0 7 】

従って、現在の課題は、分散および動的環境におけるプロセッサのセット、および稼動プロセッサの合計数に対するサービス・プロセッサの最大数の比率を表す数を考慮して、サービス・ポイントおよび各サービス・ポイントがサービスするプロセッサを決定することである。

#### 【 0 0 0 8 】

クラスタ化という概念は多くの他の分野にうまく適用されている。しかし、上述の適用分野では、クラスタの数を経験的に指定することができない。余分のサービス・ポイントが常に制限されるよう、クラスタの数に上限を設定することが必要である。一定の限界を

50

伴うクラスタ化の問題点は、一般には対処し難いもの、即ち、効率的な最適の解決方法が存在しないもの、と思われている。しかし、次善の効率的な解決方法に対する要求は依然として存在する。更に詳しく云えば、システムの最適なパフォーマンス（自己構成および自己修正を含む）を保証するために、システムにおける種々のプロセッサをクラスタに動的に割り当て、各クラスタ内のサービス・ポイントを割り当てるための効果的な手順に対する要求が存在する。

【発明の開示】

【発明が解決しようとする課題】

【0009】

本発明の課題は、複数のプロセッサを有するシステムであって、各々が1つのサービス・ポイントを持つ少なくとも1つのクラスタを備えたシステム、を提供するようにそのシステムを構成する方法を提供することにより、上記の要求に対処することである。

10

【課題を解決するための手段】

【0010】

本発明によれば、これは、システムにおける各プロセッサから他のプロセッサまでの距離を計算することによって行われる。そこで、複数の合計距離が計算される。その場合、各合計距離が1つのプロセッサに関連している。それらの複数の合計距離から、最小の合計距離が決定される。1つのプロセッサが、サービス・プロセッサとなるように割り当てられる。このプロセッサは、最小の合計距離に関連するプロセッサである。

20

【0011】

本発明のもう1つの実施例によれば、本発明の方法は、システムを複数のクラスタに区分することを含む。この区分するプロセスは、各プロセッサに関連する合計距離に従ってプロセッサを分類すること、各プロセッサを2つのクラスタのうち一方に割り当てること、各クラスタにおけるプロセッサに対する最小の合計距離を、当該クラスタにおけるプロセッサに関連する複数の合計数に従って決定すること、および、各クラスタに対するサービス・プロセッサとして、当該クラスタにおける最小合計距離に関連するプロセッサを割り当てること、を含む。

【0012】

本発明の更なる実施例によれば、2つのクラスタの1つが2つのクラスタに再分割され、それによってシステムを3つのクラスタに区分する。3つのクラスタにおけるプロセッサに関連する複数の合計距離に従って、3つのクラスタの各々におけるプロセッサに対する最小合計距離が決定される。その最小合計距離に従って、プロセッサが3つのクラスタに割り当てられる。各クラスタにおける最小合計距離に関連するプロセッサがそのクラスタに対するサービス・ポイントとして割り当てられる。

30

【0013】

本発明の更なる実施例によれば、プロセッサは異なるタイプのプロセッサであってもよく、それらのタイプに従ってクラスタに割り当てられる。

【0014】

システムは、プロセッサがそのシステムに加えられるかまたはそのシステムから取り除かれるとき、動的に構成することが可能である。更に、システムの区分化は、プロセッサがシステムから取り除かれるとき、動的に変更することが可能である。

40

【0015】

各クラスタにおいて、サービス・ポイントがシステムから取り除かれる場合、そのサービス・ポイントの機能を引き受けるために他のプロセッサをバックアップ・サービス・ポイントとして割り当てることが可能である。

【0016】

本発明の他の実施例によれば、コンピュータ可読記憶媒体が提供される。その記憶媒体は、複数のプロセッサを有するシステムであって、各々が1つのサービス・ポイントを持つ少なくとも1つのクラスタを備えたシステム、を提供するようにそのシステムを構成する方法を実行するための命令を記憶している。この方法は、システムにおける各プロセッサ

50

サから他のプロセッサまでの距離を計算するステップ、各合計距離が1つのプロセッサに関連する複数の合計距離を計算するステップ、複数の合計距離から最小合計距離を決定するステップと、最小合計距離に関連するプロセッサをサービス・ポイントとして割り当てるステップを含む。

【発明を実施するための最良の形態】

【0017】

以下の説明では、システムがN個のプロセッサ  $c_1$ 、 $c_2$ 、 $\dots$ 、 $c_N$  のセットを有するものと仮定する(図3参照)。 $d(i,i)=0$  および  $d(i,j)$  が三角不等式を満足するよう、即ち、 $d(i,j)+d(j,k)$  が、 $d(i,k)$  よりも大きいかまたはそれに等しくなるよう、プロセッサ  $c_i$  および  $c_j$  の間の距離関数  $d(i,j)$  が定義される。プロセッサ間の通信の観点から、これは、j番目のプロセッサを介した通信がi番目のプロセッサおよびk番目のプロセッサの間の直接通信よりも遅くなければならないことを意味する。距離  $d(i,j)$  は、通信ネットワークの速度、データ交換およびデータ量の比率、およびプロセッサの地理的位置の関数である。

10

【0018】

N個のプロセッサのうち、サービス・ポイントとして作用するという余分な負荷を引き受けることができるプロセッサの最大分数として、分数値  $f$  が定義される。数値  $p = N * f$  は、サービス・ポイントとして作用するプロセッサの最大数を表す。従って、システムは、最大  $p$  個のクラスタを有することが可能である。

20

【0019】

本発明によれば、N個のプロセッサから成り、 $p$  個のクラスタを有するシステムにおいて、クラスタおよびサービス・ポイントを割り当てるための方法が、複雑さを増すいくつかのケースとして下記のように提供される。

【0020】

A. ケース1：同様のプロセッサの静的システム

プロセッサ  $c_1$ 、 $c_2$ 、 $\dots$ 、 $c_N$  のシステムは、 $p$  個の異なるクラスタに分けられるべきものである。この場合、各クラスタは、当該クラスタに属するすべてのプロセッサの、当該クラスタの対応するサービス・ポイントまでの最大距離の和が最小になるよう、サービス・ポイントによって識別される。例えば、図4は、 $p = 4$  個のクラスタに分けられるべき  $N = 20$  個のプロセッサを有するシステム30を示す。この例では、分数  $f$  は、 $1/5$  であり、実際には、 $f$  は一般に  $1/50$  乃至  $1/100$  の範囲の非常に小さい値である。

30

【0021】

まず、 $p = 1$  および  $p = 2$  の時の単純なケース、即ち、1個および2個のクラスタに関する方法を説明することにする。これらの2つのケースは非自明なケースの基礎を形成する。 $p = 1$  に関して、即ち、単一のクラスタに関して、次のようなアルゴリズムが使用される。

(a) アルゴリズム ( $N, 1$ )

$i = 1$  乃至  $i = N$  に対して下記の動作を行う。

$j = 1$  乃至  $N$  ( $j \neq i$ ) に対して  $d_i = \text{SUM}\{d(i,j)\}$  をセットする。

40

$i = 1$  乃至  $N$  に対して  $d_i$  が最小となるように  $i$  を出力する。

【0022】

図5は、 $N = 10$  個のプロセッサを有し、プロセッサ  $c_1$  に対して距離  $d(1,j)$  が計算されるシステムを示す。距離  $d(1,2)$ 、 $d(1,3)$ 、 $\dots$ 、 $d(1,10)$  の和がプロセッサ  $c_1$  に対する合計距離である。

【0023】

図6は、アルゴリズム ( $N, 1$ ) を示すフローチャートである。第1プロセッサに関して開始すると(ステップ301における  $i = 1$ )、そのプロセッサと他のプロセッサの各々との間の距離が計算される(ステップ302)。これらの距離の和がクラスタにおける  $i$  番目のプロセッサに対する合計距離である(ステップ303)。合計距離が最小である

50

プロセッサ（ステップ304）がそのクラスタに対するサービス・ポイントになるであろう。

【0024】

$p = 2$ 、即ち、2つのクラスタに関して、下記のアルゴリズム、即ち、アルゴリズム（ $N, 2$ ）が使用される。そのアルゴリズムでは、アルゴリズム（ $N, 1$ ）が繰返し使用される。初期ステップのように、 $N$ 個のプロセッサのセット全体が単一のクラスタとして扱われ、上記のアルゴリズム（ $N, 1$ ）が各プロセッサに対する合計距離を計算するために使用される。即ち、 $N$ 個のプロセッサの各々が当該プロセッサから他のプロセッサまでの距離の和を表す合計距離を有する。 $i$ 番目のプロセッサに対するこの距離は $d(i, 1)$ として示される。そこで、 $N$ 個のプロセッサをこの距離に従って分類することが可能である。

10

（b）アルゴリズム（ $N, 2$ ）

1. プロセッサを距離  $d(i, 1)$  によって分類する。

2. 2つのクラスタ  $L_1$  および  $L_2$  を初期化する。なお、 $L_1$  は  $N$  個のプロセッサすべてを含み、 $L_2$  は 0 個のプロセッサを含む。

2.1 アルゴリズム（ $N, 1$ ）を使って  $L_1$  に対するサービスのポイントおよびその最小距離を計算する。

2.2  $d_0$  をその距離に初期化する。

3. プロセッサ  $c_1$ 、 $c_2$ 、...、 $c_N$  をそれらの分類順に走査する。

3.1 各プロセッサ  $c_i$  に関して、それを  $L_1$  から取り除き、それを  $L_2$  に入れる

20

3.2  $L_1$  および  $L_2$  に対するサービスのポイントおよび距離  $d_{1i}$  および  $d_{2i}$  を計算する。

3.3 距離の和  $d_i = d_{1i} + d_{2i}$  をセットする。

4.  $d_i$  が最小となるよう  $i = 1$  乃至  $N$  に対する距離  $d_i$  を選定し、その距離に対するクラスタ  $L_1$  および  $L_2$  を出力する。

【0025】

図7は、 $N = 10$  個のプロセッサを有し、2つのクラスタ  $L_1$  および  $L_2$  に分けられたシステム30の概略図である。図8は、アルゴリズム（ $N, 2$ ）を図解したフローチャートを示す。上述のように、まず、アルゴリズム（ $N, 1$ ）を使用してシステムの各プロセッサに対する合計距離が計算され（ステップ321）、しかる後、この距離に従って、プロセッサが分類される（ステップ322）。クラスタのその初期割り当てでは、10個のプロセッサすべてがクラスタ  $L_1$  にあり、クラスタ  $L_2$  にはクラスタがない。従って、初期最小合計距離  $d_0$  は、アルゴリズム（ $N, 1$ ）を使用して前に計算された値と同じである（ステップ323）。各プロセッサ（最小距離で始まり、アルゴリズム（ $N, 1$ ）を使用して計算されたように距離による分類順序で生じる）がクラスタ  $L_1$  からクラスタ  $L_2$  に順次移動し（ステップ324）、各繰返しに対してアルゴリズム（ $N, 1$ ）を使用して各クラスに対する最小合計距離が計算される（ステップ325）。 $i$  番目の繰返しにおいて、最小合計距離が、クラスタ  $L_1$  に対する  $d_{1i}$  およびクラスタ  $L_2$  に対する  $d_{2i}$  として示される。そこで、結合された最小距離  $d_i = d_{1i} + d_{2i}$  が計算される（ステップ326）。 $N$  個の繰返しの1つにおいて、 $d_{1i}$  の値が最小になるであろう。この値は、対応するクラスタ割振りの際にシステムに対する合計距離として選定される（ステップ327）。

30

40

【0026】

$p$  個のクラスタの非自明なケースに関して、上記のアルゴリズム（ $N, 2$ ）が繰返し使用される。各ステージにおいて、アルゴリズム（ $N, 2$ ）を使用して、作成されたクラスタの各々が2つの部分に分けられる。そこで、合計距離関数を最小にするクラスタ割り当てが識別される。

【0027】

例えば、アルゴリズム（ $N, 2$ ）が実行されるとき、図7に示された  $N = 10$  個のプロセッサのシステムが2つのクラスタ  $L_1$  および  $L_2$  に分けられる。そこで、次のステージ

50

において、 $L_2$  はそのままにして、クラスタ  $L_1$  が 2 つのクラスタ  $L_{1.1}$  および  $L_{1.2}$  に分けられる (図 9)。次に、合計距離 ( $d_{2.1}$  として示される) が、 $L_{1.1}$ 、 $L_{2.1}$ 、および  $L_2$  のクラスタ距離の和として計算される。次に、クラスタ  $L_1$  をそのままにして、クラスタ  $L_2$  が 2 つのクラスタ  $L_{2.1}$  および  $L_{2.2}$  に分けられる (図 10)。再び、合計距離 ( $d_{2.2}$  として示される) が、 $L_1$ 、 $L_{2.1}$ 、および  $L_{2.2}$  のクラスタ距離の和として計算される。 $d_{2.1}$  および  $d_{2.2}$  の小さい方が、それに対応するクラスタ構成と共に選ばれる。従って、 $d_{2.1} < d_{2.2}$  である場合、クラスタ  $L_{1.1}$ 、 $L_{1.2}$ 、および  $L_2$  が選ばれる。再分割である  $L_{1.1}$  および  $L_{1.2}$  はそれぞれ  $L_1$  および  $L_3$  として再命名され、クラスタ  $L_3$  のこのセットに対する合計距離が  $d_{2.1}$  に等しくセットされる。

【0028】

10

同様に、いずれの繰返しステージ  $q + 1$  (但し、 $q = 2, \dots, N - 1$ ) においても、クラスタ  $L_1, \dots, L_q$  が存在する。 $q + 1$  個のクラスタの  $q$  個のセットが、他のすべてのクラスタをそのままにして、クラスタ  $L_1, \dots, L_q$  の各々を 2 つの部分  $L_{\{j.1\}}$  および  $L_{\{j.2\}}$  (但し、 $j = 1, \dots, q$ ) に一時に 1 つずつ分けることによって作成される。そこで、これらのクラスタの各々に対して、距離  $d_{\{q.1\}}, \dots, d_{\{q.q\}}$  が計算される。これらの距離のうちの最小距離が、対応するクラスタ構成と共に選ばれる。 $d_{\{m.j\}} = \min\{d_{\{q.1\}}, \dots, d_{\{q.q\}}\}$  であると仮定する。そこで、対応するクラスタ  $\{L_1, L_2, \dots, L_q\}$  が選ばれる。次に、距離  $d_{\{q+1\}}$  が  $d_{\{m.j\}}$  に等しくセットされ、 $L_m = L_{m.1}$  および  $L_{\{q+1\}} = L_{m.2}$  にセットされる。このプロセスが、それぞれが自身のサービス・ポイントを有する  $q + 1$  個のクラスタ

20

のセットを生じる。

【0029】

$p$  個のクラスタのためのアルゴリズムは次のようになる。

(c) アルゴリズム ( $N, p$ )

0. アルゴリズム ( $N, 2$ ) を使用してプロセッサのセットを 2 つのクラスタに分ける。

1.  $i = 2$  乃至  $p$  に対して下記の動作を行う。

1.1  $j = 1$  乃至  $i$  に対して下記の動作を行う。

1.1.1 アルゴリズム ( $N, 2$ ) を用いて  $L_j$  を  $L_{j.1}$  および  $L_{j.2}$  に分ける。

1.1.2 クラスタの合計距離を合算し、それを  $d_{\{i.j\}}$  と呼ぶ。

1.2  $d_{\{i+1\}} = \min\{d_{\{i.j\}}\}$  をセットする ( $j=1$  乃至  $i$ )。

30

最小値が生じた  $j$  の値を  $j = m$  にすると、

1.3  $L_m = L_{m.1}$  をセットする。

1.4  $L_{\{i+1\}} = L_{m.2}$  をセットする。

2. 結果  $L_1, \dots, L_p$  を戻す

【0030】

図 11 は、上記のアルゴリズム ( $N, p$ ) に関するフローチャートを示す。まず、アルゴリズム ( $N, 2$ ) を使用してシステムが 2 つのクラスタに分けられる (ステップ 351)。次に、繰返しが行われ、アルゴリズム ( $N, 2$ ) を用いて  $j$  番目のクラスタが 2 つに分けられ (ステップ 352)、すべてのクラスタ (2 つの再分割されたクラスタを含む) に対する合計距離が計算される (ステップ 353)。これらの合計距離の最小値がつけられ (ステップ 354)、それに対応するクラスタ構成が選ばれる (ステップ 355)。そこで、それらのクラスタが、図 9 および図 10 に関連して上述したように、再命名される。 $p$  個のクラスタのセットが見つけれられるまで、クラスタが再分割される。しかる後、クラスタは、各々が 1 つのサービス・ポイントを有する  $p$  個のクラスタのセットを有することになる。このプロセスの結果が、 $N = 20$  個のプロセッサおよび  $p = 4$  個のクラスタに関して、図 12 に概略的に示される (図 4 と比較されたい)。システム 30 は、それぞれのサービス・ポイント 311、312、313、314 を有する 4 つのクラスタ 31、32、33、34 に分けられる。

40

【0031】

B. ケース 2: 非同ープロセッサの静的システム

50

図13は、プロセッサすべてが必ずしも同一のものではないが、それらのオペレーティング・システムおよび/またはテクノロジーに従って事前にグループ分けされており、或るタイプのプロセッサに対するサービス・ポイントは同じタイプのものでなければならない、というシステム40の概略図である。システム40では、プロセッサが2つのタイプ41および42のものである。種々のタイプのプロセッサが、TCP/IPのようなプロセッサに無関係のプロトコルを実行するネットワークを介して相互に通信を行うことが可能である。例えば、インテル社のチップ上で稼動するウィンドウズ・ベースのプロセッサのセットおよびAIX(IBMコーポレーションの商標)を実行するサーバのセットを使用して、自律コンピューティング環境を設定することも可能である。その場合、インテル社のシステムに対するサービス・ポイントはインテル社のシステムでなければならない、AIXに関するサービス・ポイントはAIXベースのシステムでなければならない。従って、システム40は、1つのクラスタ内のすべてのプロセッサが同じタイプのものとなるよう、クラスタに分けられなければならない。

#### 【0032】

このケースに対する解決方法を、上記のケース1から得ることが可能である。m個の異なるタイプのプロセッサが存在し、これらの異なるタイプの各々に対するサービス・ポイントがそれ自身の種類のみから選ばれたものであることが可能であると仮定する。この場合、先ず $m < p$ であるかどうかチェックされる。但し、 $p$ はクラスタの最大数である。 $m = p$ である(実用的なシステムでは極めてありそうもないことである)場合、クラスタは、単に、プロセッサのタイプに従ってグループ分けされる。 $m < p$ に関しては、各々が同じタイプのプロセッサだけを含むm個のクラスタ $L_1, \dots, L_m$ にプロセッサをグループ分けすることによって、クラスタが初期設定される。しかる後、アルゴリズム( $N, p$ )がこれらのm個のクラスタに適用される。これらのプロセスの結果が、 $N = 20$ 、 $m = 2$ 、および $p = 4$ に関連して図14に概略的に示される。プロセッサ41はクラスタ411、412、413にグループ化され、一方、プロセッサ42はクラスタ420にグループ化される。

#### 【0033】

C. ケース3: 動的システム(プロセッサ編入・離脱システム)

このケースはケース2の拡張であり、プロセッサの数が動的に変化する。プロセッサは、障害のために、またはそれらがオンデマンド環境で動作しているために、システムに編入することまたはシステムから離脱することが可能である。更に、プロセッサの数の変化がクラスタの数または配列の変化に通じることがあり、従って、クラスタも動的に変化する。更に、プロセッサが複数のタイプのものである(従って、複数の異なるグループにある)場合、グループの数が動的に変化することもある。例えば、図15は、クラスタ50におけるプロセッサ501がシステムから取り除かれることになる状況を示す。そのクラスタは、図16に示されるように、クラスタ51として再構成される。プロセッサ501を取り除いた結果、クラスタ51のサービス・ポイントの再割当てが生じることがあり得る。

#### 【0034】

クラスタ内のプロセッサとサービス・ポイントとの間の合計距離は、システムのパフォーマンスの損失を回避するために動的に最小にされなければならない。1つまたは複数のプロセッサがシステムに加えられるかまたはシステムから取り除かれるたびに、アルゴリズム( $N, p$ )を使用することによってそのクラスタを再計算することが可能である。この方法は、計算のオーバーヘッドによって非常に高価となり、大きなパフォーマンス低下を引き起こすことがある。従って、(i)システムにプロセッサを加えること、および(ii)システムからプロセッサを取り除くこと、の2つのケースの各々に関してシステムを動的に更新することが必要である。これらの2つのケースに関する手順は後述される。アルゴリズム( $N, p$ )を使用してクラスタの最適なセットが事前に見つかっているものと仮定する。

#### 【0035】



(i) 新たなプロセッサがシステムに加えられる。この場合、すべての既存のクラスタにおけるすべてのサービス・ポイントに対比してその加えられたプロセッサがテストされる。しかる後、その新たなプロセッサは、合計距離が最小であるクラスタに加えられる。その新たなプロセッサが加えられるクラスタでは、新たなサービス・ポイントを見つけるために、アルゴリズム (N, 1) が使用される。この手順は必ずしもシステムを最適化するものではないが、非常に高速であることに留意すべきである。この手順は図 17 に概略的に示される。図 17 では、システム 55 が 2 つのクラスタ 551 および 552 を含み、新たなプロセッサ 560 を加えることが望ましい。プロセッサ 560 はクラスタ 552 に加えられるであろう。それは、システムの合計距離がこれによって最小に保たれるためである。

10

#### 【0036】

(ii) 既存のプロセッサがシステムから取り除かれる。この場合、どちらのクラスタからプロセッサが取り除かれるかが先ず決定される。この状況に関しては、次の 2 つのケースが存在する。

(a) 取り除かれたプロセッサがサービス・ポイントから最も遠いものではない場合、それはそのクラスタに関する最大距離に影響を与えない。その場合、そのプロセッサは取り除かれ、このクラスタに関してアルゴリズム (N, 1) を使用することにより、新たなサービス・ポイントが再計算される。この状況が図 18 に示される。そこでは、システム 57 がクラスタ 571 および 572 を含み、プロセッサ 574 が取り除かれるべきものである。その結果、クラスタ構成に関する変化は生じない。

20

(b) 取り除かれたプロセッサがサービス・ポイントから最も遠いものである場合、その取り除きはクラスタにおけるサービス・ポイントからの最大距離に影響を与える。その場合、プロセッサが先ずクラスタから取り除かれる。一般に、システムは、最適なパフォーマンスを維持するために再平衡化されなければならない。プロセッサがクラスタから取り除かれるとき、他のクラスタから当該クラスタに他のプロセッサを組み入れることによって、システムを更に効率的にすることが可能であろう。この状況が図 19 に示される。そこでは、プロセッサ 577 がシステムから取り除かれている (図 18 と比較されたい)。その結果、システムの合計距離を最小にするように、新たなクラスタ 578 および 579 が形成される。これを達成するために、クラスタの最終的セットを形成するためのアルゴリズム (N, p) の各ステップにおいて、前に形成されたクラスタが、アルゴリズム (N, 2) を使って分割される。従って、クラスタを形成するプロセス全体をバイナリ・ツリーで表現することが可能である。最終的なクラスタはすべてバイナリ・ツリーのリーフ・ノードである。

30

#### 【0037】

プロセッサがクラスタから取り除かれるとき、そのプロセッサが取り除かれるクラスタの他の同胞 (sibling) ノードが考察される。その同胞ノードにおけるプロセッサは、プロセッサを同胞ノードから当該クラスタに移動させることがシステムの全体的な距離を最小にするかどうかに関して調べられる。システム全体の全体的な距離が影響を受けない場合、アクションが取られる必要はない。しかし、プロセッサを移動させることによって全体的な距離が減少する場合、プロセッサが同胞ノードから取り除かれて当該クラスタに組み入れられ、従って全体的な距離が再び最小にされる。同胞ノードがリーフ・ノードまたは既存のクラスタである場合、更なるアクションを取る必要はない。そうでない場合、プロセッサの取り除き後、そのチルドレン・ノードを平衡化するためにアルゴリズム (N, 2) が使用される。

40

#### 【0038】

上記の調整は本来ローカルのことであり、非常に高速で行うことが可能である。しかし、それらはシステム全体を総体的に平行化するものではない。従って、多くの動的な調整後、システムはかなり平衡を失ったものになることがあり、時々アルゴリズム (N, p) を使用してシステムを最適化し直すことが必要であろう。

#### 【0039】

50

#### D．ケース４：サービス・ポイントに対するバックアップ

システムのフェールセーフ動作を保証するために、本発明は、各クラスタにデュアル・サービス・ポイントを提供するための方法も含む。この方法では、バックアップ・サービス・ポイントがアクティブなサービス・ポイント内に記憶された情報のミラー情報を保持し、更新時にサービス・ポイントが障害を生じた場合、アクティブなサービス・ポイントとして機能し始める。それらのサービス・ポイントがシステムの普通のプロセッサとは全く異なることは明らかであろう。従って、サービス・ポイントが障害を生じ、システムから取り除き得ることは可能である。この場合のシステム障害を防ぐために、サービス・ポイントのサービス情報が同じクラスタにおける他のプロセッサに記憶される。例えば、図２０に示されたシステム６０は４つのクラスタ６１、６２、６３、および６４を含む。クラスタ６１のサービス・ポイントはプロセッサ６０１を含み、プロセッサ６０２はバックアップ・プロセッサである。アクティブ・サービス・プロセッサ６０１における障害の場合、バックアップ・サービス・ポイント６０２がサービス・ポイントの機能を引き受ける。そこで、アルゴリズム（ $N, 1$ ）を、必要な場合にクラスタを平衡化し直すために使用することが可能である。

#### 【００４０】

#### E．動的な数のサービス・ポイント

更に一般的な問題点は、サービス・ポイントの数  $p$  に関する最大の許容可能な制限も変更されるという状況である。これは、かなりの数のプロセッサが一時にシステムに加えられるかまたはシステムから取り除かれるときに起こり得る。これは、分数値  $f$  ( $p = N * f$  に従って、サービス・ポイントの最大の許容可能な数  $p$  を与える) が変更される場合にも起こり得る。次のような考慮すべき２つのケースが存在する。

(i) 最大の許容可能なサービス・ポイントの新たな数  $p_1$  が  $p$  よりも大きい。この場合、アルゴリズム（ $N, p$ ）は現在のクラスタ数と共に使用されるが、そのアルゴリズムは、新しいクラスタ数  $p_1$  を用いて実行される。これは、ケース１において説明された解決方法と同じであり、第１ステップにおいて、 $p$  個のクラスタの各々が１回分割され、一方、他のクラスタはそのままにされる。各ケースにおいて、最も近いサービス・ポイントまでの距離の総和が計算され、その総和を最小にするクラスタ構成が選ばれる。このプロセスは、 $p + 1$ 、 $\dots$ 、 $p_1$  に対して繰り返し実行される。

(ii) 最大の許容可能なサービス・ポイントの新たな数  $p_2$  が  $p$  よりも小さい。この場合、 $p$  個のクラスタの各々を生成した命令を再コールすること（ケース１の上記説明参照）が必要である。クラスタは、 $p_2$  個のクラスタだけが残るまで、逆順に再結合される。

#### 【００４１】

本発明を特定の実施例によって説明したが、当業者にとっては多くの代替、修正、および変更が明らかであることは、上記の説明から鑑みて明白である。従って、本発明は、本発明および「特許請求の範囲」の技術的範囲内にあるそのような代替、修正、および変更をすべて包含することを意図するものである。

#### 【図面の簡単な説明】

#### 【００４２】

【図１】サービス・ポイントを有するシステムにおけるプロセッサの代表的な配列を示す概略図である。

【図２】障害のあるプロセッサがシステムから取り除かれる後の図１のシステムを示す概略図である。

【図３】プロセッサ間の距離が本発明に従って決定される場合の複数のプロセッサを有するシステムの概略図である。

【図４】本発明の実施例に従って、サービス・ポイントを有するクラスタに構成されるべき複数のプロセッサを有するシステムの概略図である。

【図５】本発明に従って、システムにおける１つのプロセッサから他のプロセッサまでの合計距離を計算するための手順を示す概略図である。

【図６】本発明の実施例に従ってアルゴリズム（ $N, 1$ ）におけるステップを詳細に示す

フローチャートである。

【図 7】 2つのクラスタに分割されたシステムを示す概略図である。

【図 8】 本発明の実施例に従ってアルゴリズム (N, 2) におけるステップを詳細に示すフローチャートである。

【図 9】 本発明の実施例に従って、クラスタが再分割されるシステムを示す概略図である。

【図 10】 システムにおけるクラスタの別の再分割を示す概略図である。

【図 11】 本発明の実施例に従ってアルゴリズム (N, p) におけるステップを詳細に示すフローチャートである。

【図 12】 図 4 のシステムにおいてクラスタを構成し、サービス・ポイントを割り当てるためのプロセスの結果を示す概略図である。

10

【図 13】 本発明の更なる実施例に従って、サービス・ポイントを有するクラスタに構成されるべき異なるタイプの複数のプロセッサを有するシステムを示す概略図である。

【図 14】 図 13 のシステムにおいてクラスタを構成し、サービス・ポイントを割り当てるためのプロセスの結果を示す概略図である。

【図 15】 複数のプロセッサを有し、1つのプロセッサがシステムから取り除かれるべきシステムを示す概略図である。

【図 16】 プロセッサが取り除かれた後の図 15 のシステムを示す概略図である。

【図 17】 複数のプロセッサを有し、1つのプロセッサがシステムに加えられるべきシステムを示す概略図である。

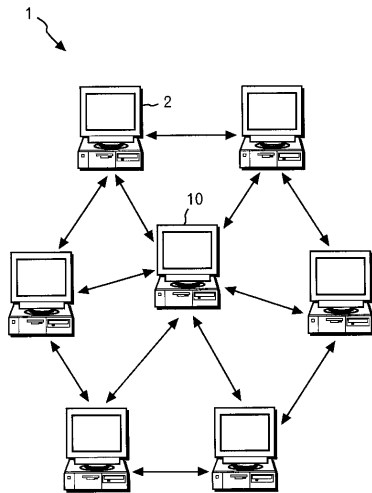
20

【図 18】 2つのクラスタを有し、1つのプロセッサがシステムから取り除かれるべきシステムを示す概略図である。

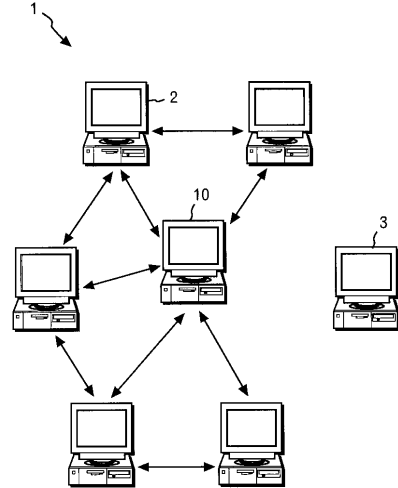
【図 19】 クラスタが再構成された場合、プロセッサが取り除かれた後の図 18 のシステムを複数のプロセッサを有し、1つのプロセッサがシステムから取り除かれるべきシステムを示す概略図である。

【図 20】 クラスタのサービス・ポイントがシステムから取り除かれるべき場合、各々がサービス・ポイントを有するクラスタに構成された複数のプロセッサを有するシステムを示す概略図である。

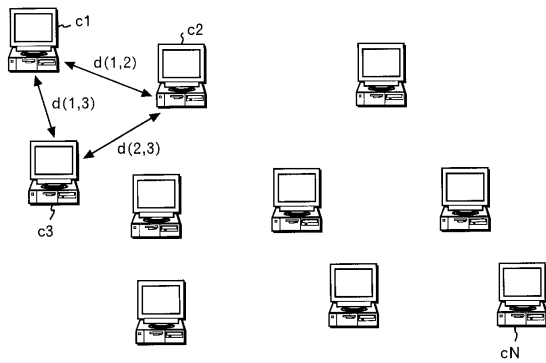
【図 1】



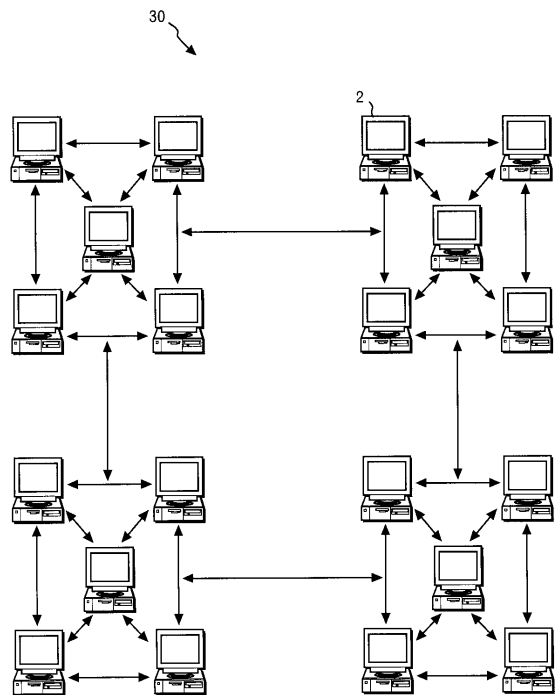
【図 2】



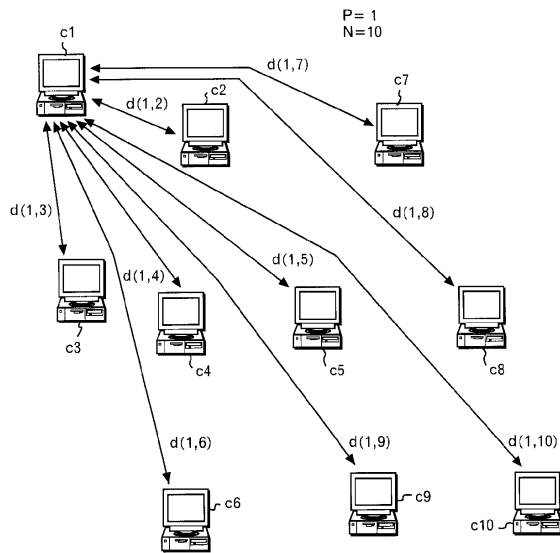
【図 3】



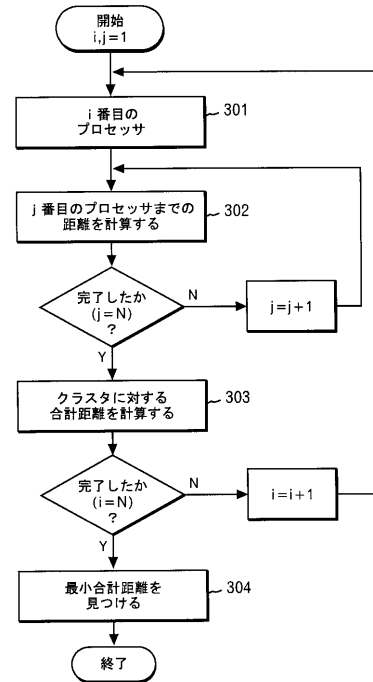
【図 4】



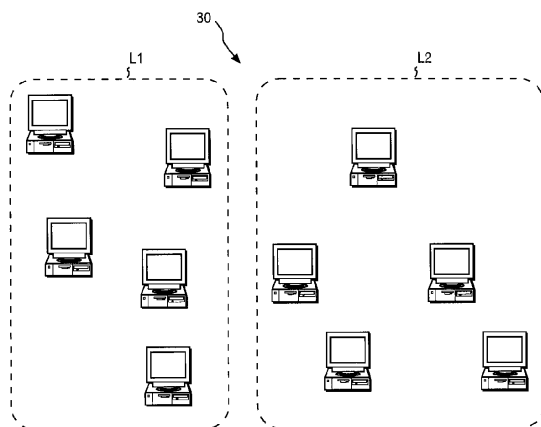
【図 5】



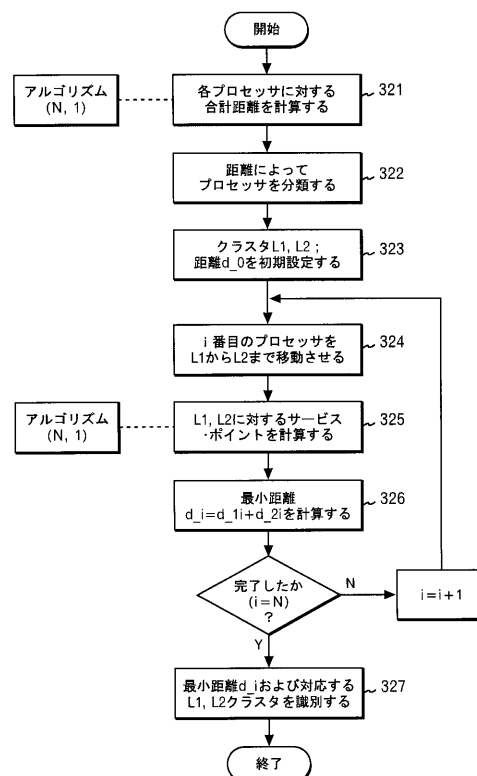
【図 6】



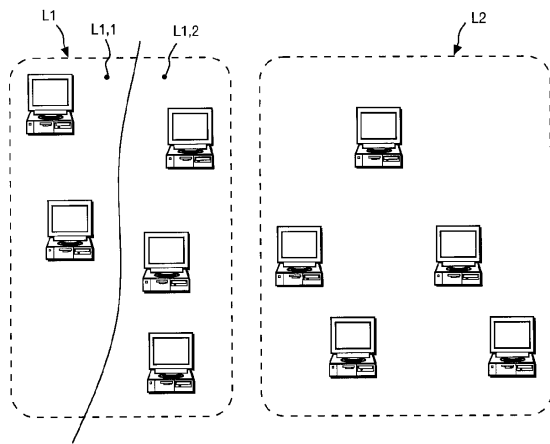
【図 7】



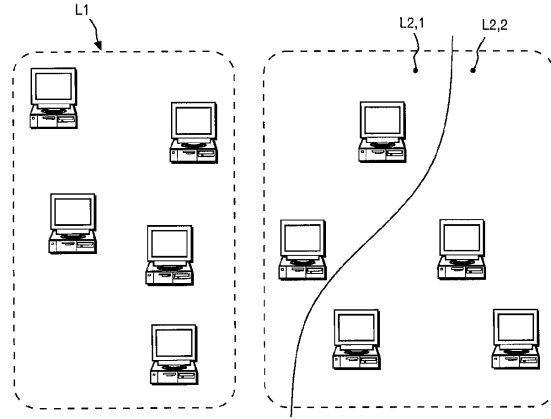
【図 8】



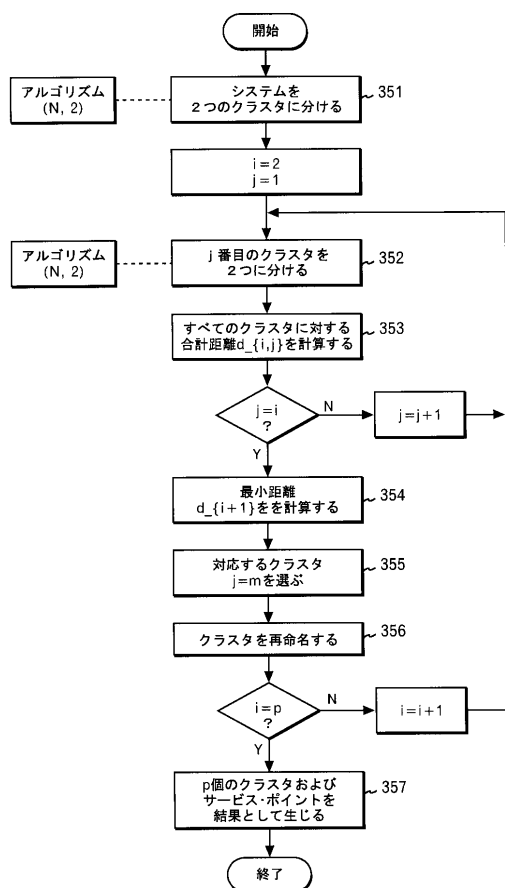
【図 9】



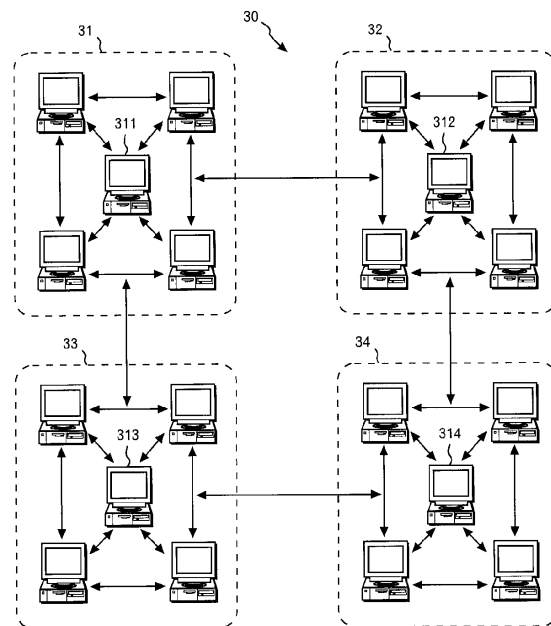
【図 10】



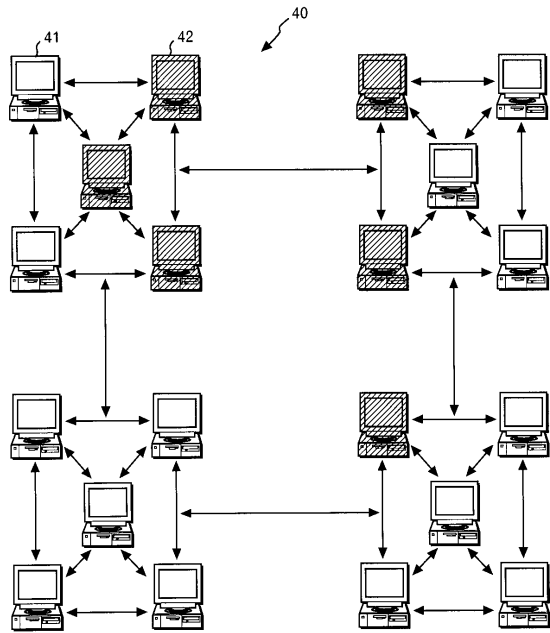
【図 11】



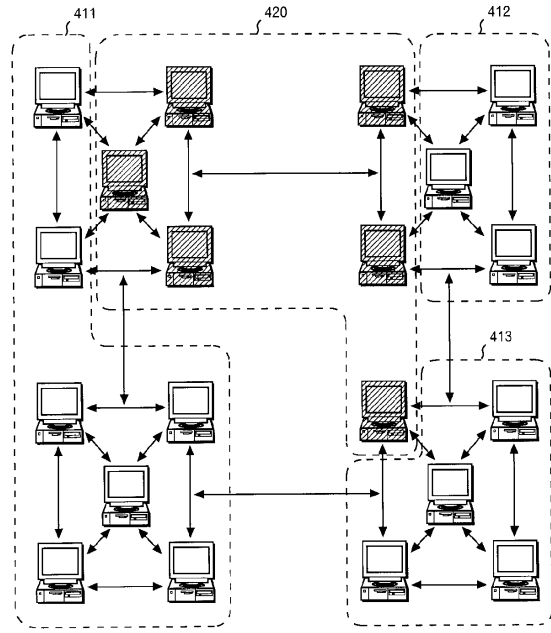
【図 12】



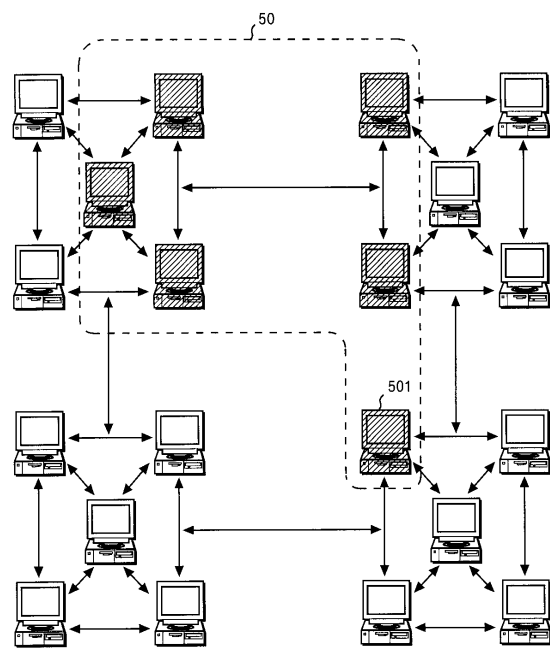
【図 13】



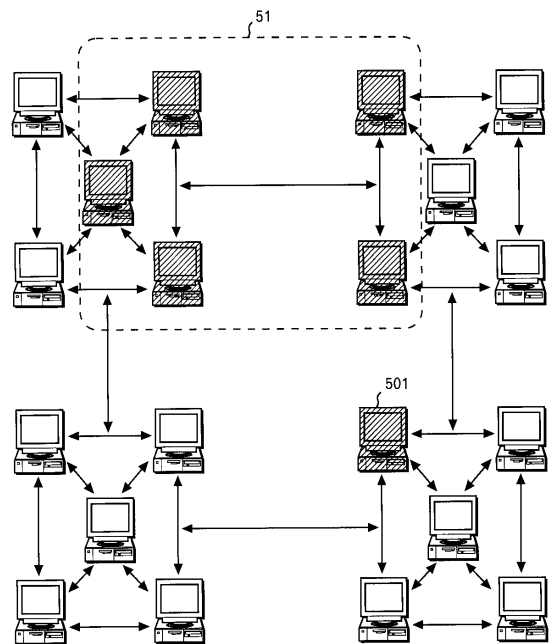
【図 14】



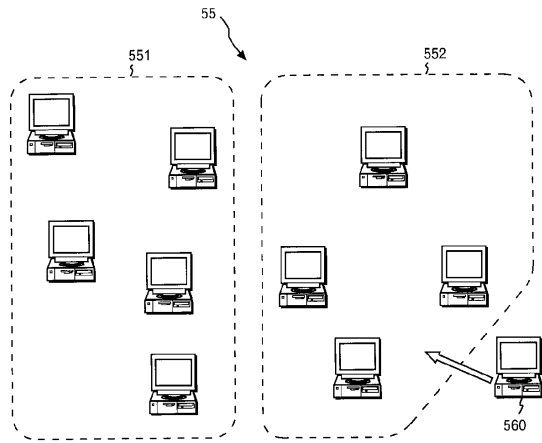
【図 15】



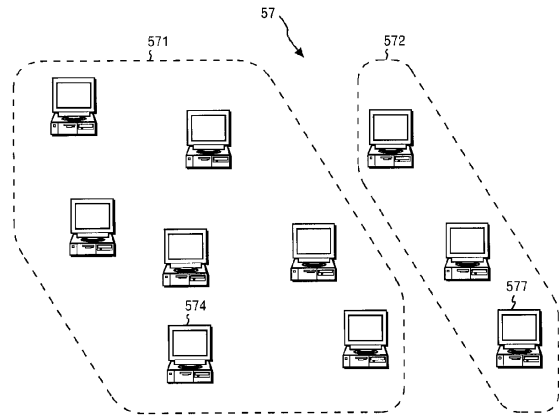
【図 16】



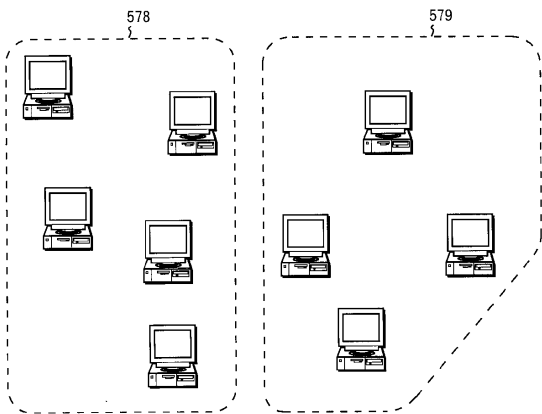
【図 17】



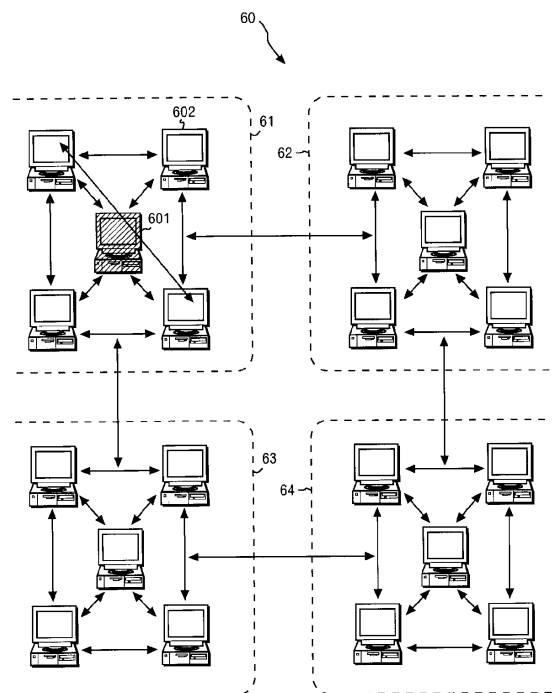
【図 18】



【図 19】



【図 20】





---

フロントページの続き

(74)代理人 100086243

弁理士 坂口 博

(72)発明者 ムクハージー、マハラージャ

アメリカ合衆国 1 2 5 9 0、ニューヨーク州ワピンジャーズ・フォールス、ブラザーズ・ロード  
1 2

審査官 久保 正典

(56)参考文献 米国特許第 0 6 4 0 0 6 9 2 ( U S , B 1 )

特開平 1 0 - 0 4 0 2 2 7 ( J P , A )

特開 2 0 0 3 - 0 3 0 1 6 5 ( J P , A )

特開平 0 3 - 2 3 2 0 5 5 ( J P , A )

Liu J et al. , Distributed Distance Measurement for Large-Scale Networks , Computer Networks , Elsevier Science , 2 0 0 3 年 2 月 5 日 , vol.41 , no.2 , p177-192

Tillett, J. et al. , Cluster-head identification in ad hoc sensor networks using particle swarm optimization , Personal Wireless Communications, 2002 IEEE International Conference on , IEEE , 2 0 0 2 年 1 2 月 1 5 日 , p201-205

(58)調査した分野(Int.Cl. , D B 名)

G06F11/16-11/20

G06F15/177