

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4243376号  
(P4243376)

(45) 発行日 平成21年3月25日(2009.3.25)

(24) 登録日 平成21年1月9日(2009.1.9)

(51) Int. Cl. F I  
**G06F 17/30 (2006.01)**  
 G06F 17/30 210D  
 G06F 17/30 419A  
 G06F 17/30 220B

請求項の数 7 (全 18 頁)

(21) 出願番号	特願平11-17644	(73) 特許権者	596170170
(22) 出願日	平成11年1月26日(1999.1.26)		ゼロックス コーポレイション
(65) 公開番号	特開平11-316768		XEROX CORPORATION
(43) 公開日	平成11年11月16日(1999.11.16)		アメリカ合衆国、コネチカット州 068
審査請求日	平成18年1月23日(2006.1.23)		56、ノーウォーク、ピーオーボックス
(31) 優先権主張番号	09/013, 668		4505、グローバー・アヴェニュー 4
(32) 優先日	平成10年1月26日(1998.1.26)		5
(33) 優先権主張国	米国 (US)	(74) 代理人	100075258
			弁理士 吉田 研二
		(74) 代理人	100096976
			弁理士 石田 純
		(72) 発明者	クレイグ デー シルバーステイン
			アメリカ合衆国 カリフォルニア州 スタ
			ンフォード サンタ フェ アベニュー
			817

最終頁に続く

(54) 【発明の名称】 任意のコーパスサブセットをほぼ一定時間でクラスタ化するための方法および装置

(57) 【特許請求の範囲】

【請求項1】

関連する複数のドキュメントがメタドキュメントにまとめられ、このメタドキュメントが階層的に関連づけられて電子的に記憶されるコーパスを処理して、事前に識別された関心に関連するドキュメントをクラスタ化する方法であって、

プロセッサが、

ユーザの入力に応じて、少なくとも一つの初期メタドキュメントであって下位のメタドキュメントを含む初期メタドキュメントを選択する初期選択ステップと、

選択された初期メタドキュメントを階層的関連に基づいて下層側に拡張し、次の階層の複数メタドキュメントからなるフォーカスセットを選択する拡張ステップと、

前記フォーカスセット内のメタドキュメントの良好度を検査して最悪メタドキュメントを選択する最悪選択ステップと、

前記選択された最悪メタドキュメントをその下層側の子孫のメタドキュメントに置き換える置き換えステップと、

関心のあるドキュメントを含まない子孫のメタドキュメントを除去する除去ステップと

、前記フォーカスセットのメタドキュメントの数が所定数になるまで、前記最悪選択ステップと、置き換えステップと、除去ステップを繰り返すステップと、

を実行し、

得られたフォーカスセット内のメタドキュメントをクラスタとして、所定数にクラスタ

10

20

化された電子的に記憶されたドキュメントのコーパスを得ることを特徴とする方法。

【請求項 2】

請求項 1 に記載の方法において、さらに、  
前記メタドキュメントの要約を確定するステップと、  
前記要約をユーザに提示するステップと、  
を含むことを特徴とする方法。

【請求項 3】

請求項 2 に記載の方法において、前記要約は、  
各メタドキュメントにおいて最も頻繁に現れる固定数の話題のワードと、  
各メタドキュメント内の少なくとも一つの典型的なドキュメントの名称と、  
を含むことを特徴とする方法。

10

【請求項 4】

請求項 1 に記載の方法において、前記拡張ステップは、さらに、選択されたメタドキュメント内の関心のあるドキュメントの数がカットオフ値を超えるかを決定するステップを含むことを特徴とする方法。

【請求項 5】

請求項 4 に記載の方法において、前記選択されたメタドキュメント内の関心のあるドキュメントの数がカットオフ値未満である場合、前記ドキュメントは別のドキュメントセットに加えられることを特徴とする方法。

【請求項 6】

請求項 5 に記載の方法において、前記拡張ステップは、さらに、前記別のドキュメントセットを前記次のメタドキュメントに加えるステップを含むことを特徴とする方法。

20

【請求項 7】

関連する複数のドキュメントがメタドキュメントにまとめられ、このメタドキュメントが階層的に関連づけられて電子的に記憶されるコーパスを処理して、事前に識別された関心に関連するドキュメントをクラスタ化する装置であって、

ユーザの入力に応じて、少なくとも一つの初期メタドキュメントであって下位のメタドキュメントを含む初期メタドキュメントを選択する初期選択手段と、

選択された初期メタドキュメントを階層的関連に基づいて下層側に拡張し、次の階層の複数メタドキュメントからなるフォーカスセットを選択する拡張手段と、

30

前記フォーカスセット内のメタドキュメントの良好度を検査して最悪メタドキュメントを選択する最悪選択手段、

前記選択された最悪メタドキュメントをその下層側の子孫のメタドキュメントに置き換える置き換え手段と、

関心のあるドキュメントを含まない子孫のメタドキュメントを除去する除去手段と、

前記フォーカスセットのメタドキュメントの数が所定数になるまで、前記最悪選択手段と、置き換え手段と、除去手段の処理を繰り返させる繰り返し手段と、

を含み、

得られたフォーカスセット内のメタドキュメントをクラスタとして、所定数にクラスタ化された電子的に記憶されたドキュメントのコーパスを得ることを特徴とする装置。

40

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、電子ドキュメントをほぼ一定時間でクラスタ化するための方法および装置に関する。特に、本発明は、電子ドキュメントの大きなコーパス（記録されたテキストの集合）をほとんど一定時間で、それに比べて非常に小さなクラスタのセットに分割することを目的とする。

【0002】

【従来の技術】

ドキュメントブラウジングは、大きなテキストコレクションにアクセスするために使用さ

50

れる有力なツールである。ブラウジングは、クエリー（質問）がないため検索と識別され、余りにも一般的であるかまたは余りにも漠然としているためのいずれかによって、幾つかの検索言語によって有効に表現できない情報ニーズに対して、都合よく作動する。たとえば、ユーザが、関心のある話題を適切な言語で記述することに不慣れである場合、またはワードの特定の選択にコミットする（明言する）ことを希望しない場合がある。実際、ユーザは、特定のものは何も要求しないが、コレクションの一般的な（全般的な）情報コンテンツ（内容）を探求したい場合がある。

#### 【0003】

この状況において、情報アクセスシステムは有用である。情報アクセスシステムは、コレクションのコンテンツを提示すること、および、ユーザがコンテンツの幾つかの話題に固有のサブセットに関心を集中することの両者をナビゲーションできるコレクションのアウトラインを含む。このようなブラウジングシステムは、P e d e r s e nらによる米国特許第5,442,778号（分散/集合：Scatter/Gather）および米国特許第5,483,650号に開示されており、これらの各特許は引用により本願に援用する。

10

#### 【0004】

分散/集合においては、注意は、常に、ドキュメントのフォーカスセット、特にユーザの関心を引く可能性のあるサブセットに向けられる。最初に、フォーカスセットは、ドキュメントコレクション全体であるかもしれない。フォーカスセット内のドキュメントは、ドキュメントの少数の話題固有の（トピックに密着したtopic-coherent）サブセット、またはドキュメントのクラスタにクラスタ化される。用語「クラスタ化（clustering）」および「分散（scattering）」は、同義語として使用される。したがって、フォーカスセット内のドキュメントは、クラスタに分散される、ということができる。

20

#### 【0005】

分散/集合においては、クラスタ要約（サマリー）が作成（develop）され、ユーザに提示される。クラスタ要約は、通常、フォーカスセットのドキュメントの輪郭（アウトライン）を示すコンテンツの表（テーブル）である。クラスタ要約は、各クラスタ内のドキュメントから自動的に決定される示唆に富むテキストを含む。各クラスタ要約は、2種類の情報を含む。すなわち、クラスタのドキュメントに最も頻繁に出現する話題の（topical）ワード、およびクラスタ内の数個の典型的なドキュメントの名称である。要約は、クラスタプロファイルに基づく。このクラスタプロファイルはクラスタのドキュメントに出現するワードを反映する。

30

#### 【0006】

そこで、ユーザは、最も関心を持たせるように見える複数のクラスタを識別し、選択する。選択されたクラスタは、一緒に集合され、新しい小さいフォーカスセットを形成する。すなわち、新しいフォーカスセットは、選択されたクラスタ内のドキュメントの合併（ユニオン：union）である。ユーザがドキュメントに個別にアクセスすること、またはクエリーに基づく探索方法を使用することを求めるまで、この処理が所望の回数反復される。

#### 【0007】

分散/集合は、必ずしも独立型情報アクセスツールである必要はない。むしろ、分散/集合は、論理探索または類似性探索などの探索方法と連携して使用することができる。類似する例は、リファレンスブック（reference book）であり、リファレンスブックは二つのアクセス方法を提供する。一つは、ブラウジングのための前にある目次（a table of contents）であり、他は、さらに直接的な探索のための最後にある索引（index）である。分類/集合は、必ずしも特定のドキュメントを探索するために使用されとは限らない。むしろ、クラスタ要約に存在する用語（ボキャブラリー）を公開することによって、分散/集合は、相補的な探索方法を補助する。たとえば、クラスタプロファイルを、類似性探索においてコレクション全体に対するクエリーとして使用することができる。逆に言えば、分散/集合を使用し、過剰の多数のドキュメントを検索するワードを基礎とするクエリーの結果を編成することができる。

40

#### 【0008】

50

図9は、1990年8月のNew York Times News Serviceに掲載された約5,000記事のテキストコレクションに適用される分散/集合方法を示す図である。図9においては、分散/集合方法を一層簡単に提示するために、実際のクラスタ要約の代わりに単独のワード文字(ラベル)が示されている。

【0009】

図9に示す例においては、ユーザの情報ニーズは、1990年に発生したことを一般的に決定することである。特定の話題記述は全く存在しないため、この情報ニーズを効果的に表現するワードに基づくクエリーを構成することは困難である。ユーザは、一般的な話題、たとえば、「国際的事件」を考えるが、この話題記述は、国際事件に関する記事は、通常、これらのワードを使用しないため、有効ではない。

10

【0010】

分散/集合によって、ある用語を提供することが強制されるのではなく、ユーザは、クラスタ要約のセット、すなわちコレクションの輪郭を提供される。ユーザニーズは、関心のある話題に関連する可能性があると考えられるクラスタを選択する。図9に示す分散/集合処理においては、その月の主要な新聞記事は、最初の分散からすぐに明らかになり、イラクのクエート侵入およびドイツ再統合問題である。これによって、ユーザは、国際問題に焦点を絞るようになり、「イラク」、「ドイツ」、および「石油」クラスタを選択する。これらの3クラスタは一緒に集合され、より小さなフォーカスセットを形成する。

【0011】

次に、この比較的小さなフォーカスセットは、クラスタ化、すなわち分散され、減少されたコレクションを包含する8個の新しいクラスタを形成する。減少されたコレクションは、記事のサブセットのみを含むので、これらの新クラスタによって、元の8個のクラスタより細かいレベルの詳細が明らかになる。イラク侵入に関する記事および石油記事の一部は、米国軍展開、石油市場に対するイラク侵入の影響、およびクエートにおける人質についてのクラスタに分離される。

20

【0012】

ユーザが、これらの主要な新聞記事を適切に理解するが、世界の他の部分で何が発生したかを見つけないと思う場合、ユーザは、たとえば、「パキスタン」クラスタを選択、-このクラスタも他の外国政治新聞記事を含む-およびアフリカに関する記事を含むクラスタを選択することができる。これらのクラスタを分散することによって、多数の特定の国際状況ならびに多方面にわたる国際記事の小さなコレクションが明らかになる。このようにして、ユーザは、パキスタンにおける政変およびトリニダードで発生した人質について知る。これらの記事は、別の状態では、その月の一層重要な記事の中に埋没してしまうものである。

30

【0013】

図10は、分散/集合の操作を示す図である。図10に示す例においては、テキストコレクション(またはフォーカスセット)20は、グロリエ(Grolier)の百科事典のオンライン版である。フォーカスセット内の2,700,000記事は、それぞれ、独立のドキュメントとして処理される。図10に示す例においては、ユーザは、宇宙開発における女性の役割を調査することに関心がある。この情報ニーズを正式の(formal)クエリーによって表現しようとするのではなく、ユーザは、代わりに、クラスタの記述から、関心のある話題に関連すると考えられる多数のトップレベルのクラスタ22A~22Iを提供される。次に、ユーザは、軍事経過(ヒストリー)クラスタ22A、科学および産業検出子22Cおよびアメリカ社会クラスタ22Hを選択し、グロリエの事典から得られる記事の指示されたサブセットの減少されたコーパス(またはフォーカスセット)24を形成する。

40

【0014】

次に、減少されたコーパスは、浮動によって(on the fly)もう一度クラスタ化され、減少されたコーパス24を対象にする新しいクラスタのセット26A~26Jを生成する。減少されたコーパスはグロリエの事典の記事のサブセットを含むため、これらの新クラスタ

50

タは、トップレベルクラスタ22A~22Iより細かいレベルの詳細である。ユーザは、再度、関心のあるクラスタを選択する。この場合、選択されたクラスタは、軍用機クラスタ26E、工業技術クラスタ26G、および物理クラスタ26Hである。再度、さらに減少されたコーパス28が形成され、再クラスタ化される。最終セットのクラスタ30A~30Fは、軍用機クラスタ30A、アポロ計画クラスタ30B、航空宇宙産業クラスタ30C、天候クラスタ30D、天文学クラスタ30E、および民間航空機クラスタ30Fを含む。この段階において、クラスタは、十分に小さく、記事名称の網羅的なリストを通じて直接に精読することができる。関心のある少なくとも一つの記事が見出されると仮定すると、ユーザは、同じクラスタ内に類似の特性の記事をさらに見出すこと、またはことによると探し当てた記事またはクラスタ記述の用語集(ボキャブラリー)に基づいて方向を

10

【0015】

【発明が解決しようとする課題】

ドキュメントクラスタ化に関する以前の成果は、線形時間(linear-time)法、たとえば、分散/集合および米国特許第5,483,650号に記載の線形時間法を含み、この方法によれば、クラスタ化のために要する時間は僅か数分に減少される。これは、広範囲のワードに基づくクエリーを使用し、中程度の大きさのコレクションを探索するのに十分な速度である。たとえば、毎秒およそ3000ドキュメントの速度を、分散/集合を使用し、サンマイクロシステムズ(Sun Microsystems)のSPARCSTATION2上において、実現できる。しかし、線形時間クラスタ化でさえも、非常に大きなドキュメントコレクションの対話型ブラウジングを支援するためには遅すぎる。このことは、約750,000のドキュメントを含むテキスト検索評価のために、分散/集合をTIPSTERコレクション、DARPA標準に適用する場合を考慮すれば、特に、明らかである。毎秒3000ドキュメントの速度において、これは、分散するために4時間以上を必要とし、対話型にとっては長すぎると考えられる。したがって、ドキュメントをクラスタ化するために、一層迅速なさらに効率的な方法を見出すことが必要とされる。

20

【0016】

本発明は、顧客対応可能な時間/精度トレードオフを持つコーパスサブセットをほぼ一定時間でクラスタ化するための方法および装置を提供するものである。

【0017】

本発明は、基礎的なブラウジング方法、たとえば、分散/集合に使用することが可能であり、大きなドキュメントコレクションを関連のあるドキュメントのクラスタに効率的に分割するほぼ一定時間でクラスタ化するための方法も提供するものである。

30

【0018】

【課題を解決するための手段】

本発明による再クラスタ化方法および装置においては、入力は、全体のドキュメントの複数のメタドキュメントへのクラスタ化であり、複数のメタドキュメントから「最悪」メタドキュメントが選択される。「最悪」メタドキュメントは、その子のメタドキュメントによって置換され、関心のあるドキュメントを含まないこれらの子は除去される(pruned)。次に、残りのメタドキュメントは一緒に集合され再クラスタ化される。ユーザが所望の程度の特定性を得るまで、この処理が反復される。

40

【0019】

このクラスタ化方法は従来の方法より速く、この方法においては、クラスタは本来の資質(in their own right)でドキュメントとして処理され、既存の階層(hierarchy)を使用しクラスタの新しいセットを生成する。すなわち、本発明による再クラスタ化方法および装置においては、クラスタは、大きな個別ドキュメントであるかのように、クラスタ化する必要があるメタドキュメントとして処理され、クラスタ化される。したがって、ファンアウトkを有するクラスタ階層の場合、本発明による再クラスタ化方法および装置は、最小のクラスタから開始し、各クラスタをそのk個の子によって置換する。親クラスタは検査され、最悪クラスタが除去される。すなわち、「最悪」親クラスタは、そのk個の子に

50

よって置換される。

【 0 0 2 0 】

本発明のこれらおよび他の特徴および利益は、以下の好適な実施形態に関する詳細記述に記載され、明らかとなる。

【 0 0 2 1 】

以下、本発明を添付図面を参照して詳細に述べる。図面において、同じ符号は、同じ構成要素を示す。

【 0 0 2 2 】

【発明の実施の形態】

図 1 は、本発明による再クラスタ化システム 1 0 の一実施形態を示すブロック図である。システム 1 0 は、プロセッサ 1 1、ROM 1 2、RAM 1 3、不揮発性メモリ 1 4、コーパス入力 1 5、ユーザ入力装置 1 6、ディスプレイ装置 1 7、および出力装置 1 8 を備える。

10

【 0 0 2 3 】

ブラウジング手順を実行する前に、ドキュメントコーパスがコーパス入力 1 5 から入力される。次に、ドキュメントコーパスは、プロセッサ 1 1 によって分割される。分割手順の結果は、ディスプレイ装置 1 7 に表示される。操作者は、ユーザ入力装置 1 6、たとえば、マウス、キーボード、タッチスクリーン、スタイラス、またはこれらの要素の組合せなどを使用し、コマンドおよびデータを入力することができる。ユーザは、ドキュメントのハードコピーのみでなくクラスタダイジェスト要約（サマリー）のプリント出力も出力装置 1 8、たとえばプリンタに出力することができる。

20

【 0 0 2 4 】

従来は、プロセッサ 1 1 によって、ドキュメントの初期順序付け（initial ordering）が準備される。初期順序付けは、たとえば、分散 / 集合に記載の分別法を使用して準備される。プロセッサ 1 1 によって、コーパスの最初の順序付けの要約も決定され、この要約はディスプレイ装置 1 7 に表示、または出力装置 1 8 によってユーザに出力することができる。この要約は、たとえば、分散 / 集合に記載されているクラスタダイジェスト法を使用し、決定することができる。

【 0 0 2 5 】

ユーザからユーザ入力装置 1 6 を経由して適切な命令を受領後、プロセッサ 1 1 は、コーパスのさらに進んだ順序付けを実行することができる。このさらに進んだ順序付けは、たとえば、分散 / 集合に記載されているバックショット（buckshot：大きめの散弾）法を使用し、形成される。次に、このステップの所望の数の反復が実行され、コーパスがさらに狭くされる。結局、個別のドキュメントが検査され、または幾つかの有向探索ツールが限定コーパスに適用される場合がある。

30

【 0 0 2 6 】

図 2 は、本発明による再クラスタ化の一実施形態の輪郭（アウトライン）を示す図である。処理は、ステップ S 1 0 0 において開始され、ステップ S 2 0 0 に続く。ステップ S 2 0 0 において、ユーザは、全ドキュメントコレクションの一部またはコーパスの一部を表現するドキュメントセットを選択する。後のステップにおける反復のために、フォーカスセットはメタドキュメントを含み、メタドキュメントは、それぞれ、コレクションの一部のみを表現する（代表する）。メタドキュメントセット中のメタドキュメントの数は、ほぼ所定の最大数に等しく、最大数は、たとえば、5 0 0 または 1 0 0 0 とすることができる。次に、ステップ S 3 0 0 において、最初のメタドキュメントセットは、プロセッサ 1 1 によって選択され、クラスタ化される。好適には、メタドキュメントクラスタの所定数は、1 0 である。一般に、必要とされることは、新メタドキュメントの所定数は、その後のメタドキュメントの所定最大数より小さいことが必要であるということのみである。メタドキュメントを選択し、クラスタ化する処理は、図 3 および図 4 に関連して、以下に述べる。次に、制御はステップ S 4 0 0 に続く。

40

【 0 0 2 7 】

50

ステップS 4 0 0において、新メタドキュメントは、要約されて利用できる形式になる。次に、ステップS 5 0 0において、たとえば、ディスプレイ装置17または出力装置18を使用し、ユーザに提示される。次に、処理は、ステップS 6 0 0に続き、ステップS 6 0 0において、処理は停止する。

**【0028】**

メタドキュメントセットはクラスタ階層Hを有し、クラスタ階層Hは、k個の子のファンアウト (fan-out) およびルートノードrを有する。階層は、クラスタのツリー構造であり、クラスタはノードと呼ばれ、ノード1のk個の子の合併はノード1自体と同じドキュメントを有するように、ノードはメタドキュメントを表現する。ドキュメントのセットSは、クラスタ化ルーチンに入力される。この処理の結果、k個のクラスタのセットとなり、このクラスタはS中のドキュメントを正確に含む。

10

**【0029】**

図3は、図2のメタドキュメント選択およびクラスタ化ステップS 3 0 0の第1実施形態のさらに詳細な輪郭を示す図である。ステップS 3 0 0から始まり、制御はステップS 3 2 0に進む。ステップS 3 2 0において、収集する必要があるノードの最大数Mが設定される。次に、ステップS 3 3 0において、初期フォーカスセットTが、階層Hのルートノードrとして設定される。次に、ルートノードは、そのk個の子によって直ちに置換される。次に、制御はステップS 3 4 0に進む。

**【0030】**

ステップS 3 4 0 ~ S 3 6 0において、ある方法において「良好」である関心のあるノードがクラスタ階層中に見出される。ノードの良好度を決定する方法について、以下に詳細に述べる。

20

**【0031】**

ステップS 3 4 0において、フォーカスセットTのk個のノードは検査され、「最悪」ノードがピックアップされる。「最悪」ノードは、以下の述べる「良好度」検査によって決定される。次に、ステップS 3 5 0において、「最悪」ノードは除去され、そのk個の子のノードによって置換され、子は関心のあるドキュメントを含む。関心のあるドキュメントを含まない子は含まれず、効果的に除去される (pruned)。

**【0032】**

次に、ステップS 3 6 0において、制御ルーチンは、フォーカスセットTが収集する必要があるノードの最大数Mに等しいノード数またはそれより大きいノード数を有するかを決定する。フォーカスセットTのノード数が収集する必要がある最大ノード数M未満である場合、制御はステップS 3 4 0に跳び戻る。そうではなく、フォーカスセットのノード数が少なくともMに等しい場合、制御は、ステップS 3 7 0に続く。

30

**【0033】**

ステップS 3 7 0において、フォーカスセットTはクラスタ化され、クラスタPのセットが得られる。次に、ステップS 3 8 0において、クラスタPのこのセットの各ノードは、クラスタ内の、Sにおいては存在しなかったドキュメントを削除するために、関心のあるドキュメント $I_S(n)$ によって置換される。次に、制御はステップS 3 9 0に続き、ここで制御はステップS 4 0 0に戻る。

40

**【0034】**

前述したクラスタ化ステップにおいて、見出されたM個のノードは、線形時間クラスタ化方法を使用し、クラスタ化される。選択されるノード数が限定される限り、これによって、一定時間 (constant-time) のクラスタ化が与えられる。

**【0035】**

クラスタ階層のノードの数は大きい場合があるため、すべてのノードを検査して「良好」ノードを見出すことはできない。その代わりに、クラスタ階層は、トップからファンアウトする。階層Hのルートノードから始まり、ルートノードは、直ちにその子によって置換される。得られるセットのk個のノードは検査され、「最悪」ノードがピックアップされる。「最悪」ノードは除去され、そのk個の子によって置換される。この処理は、今、検討中の2

50

$k - 1$  のノードについて反復される。実際は、すべての  $k$  個の子ノードは、必ずしも含まれない。むしろ、子ノードのサブセットのみが、検討される。  $M$  個のノードが収集されると、処理は停止される。

【 0 0 3 6 】

この時点において、共通集合（積集合）テーブル  $I_S$  が生成される。任意のノード  $n$  に対して、そのノードの共通集合  $I_S(n)$  は、  $S \cap n$  におけるドキュメントのセットである。すなわち、  $I_S(n)$  は、ドキュメントセット  $S$  とノード  $n$  に含まれるドキュメント間の共通集合である。したがって、共通集合テーブル  $I_S$  によって、ドキュメントセット  $S$  およびノード  $n$  の両者に含まれる関心のあるドキュメントのみが、提供される。  $I_S$  は、  $|S| \log(n)$  時間内に作成される。共通集合テーブル  $I_S$  を使用し、結果として得られる各ノードが、共通集合  $I_S(n)$  によって置換され、ドキュメントセット  $S$  に存在しない、クラスタ中のドキュメントが削除される。得られるノードは、クラスタ化され、  $k$  個のクラスタとなり、各ノードはなお単独の実体（エンティティ：entity）として処理される。

10

【 0 0 3 7 】

任意のノード  $n$  に対する  $S$  および  $n$  の共通集合を求めるために、ドキュメントを処理し、ドキュメントを含む階層  $H$  のすべてのノードを戻す関数を使用される。この関数は、ドキュメントセット  $S$  に従属せず、階層  $H$  が決定されると同時に決定されることができる。階層  $H$  は、一定の  $k$  のファンアウトを有するので、階層  $H$  は深度  $\log n$  を有し、したがって、各ドキュメントは、  $\log n$  ノードにある。

20

【 0 0 3 8 】

$I_S$  を求めるために、テーブルが構成され、ノードによって索引される。テーブルの各項目は、原始状態においては、空である。ドキュメントセット  $S$  の各ドキュメントに対して、事前に計算された前述した関数を使用し、どのノードがドキュメントを包含するかを見出す。次に、ドキュメントセットをこのような各ノードに対するテーブル項目に追加する。理論上は、一定時間内に任意のサイズの空テーブルを構成することが可能であるが、実際には、明白な線形時間（リニア・タイム）アルゴリズムは極めて迅速である。テーブル更新は、ドキュメント当たり時間  $\log n$ 、または全体で時間  $|S| \log n$  を要する。得られるテーブルは、  $I_S$ 、すなわち必要とされる共通集合計算ツールである。

30

【 0 0 3 9 】

図 4 は、図 2 のメタドキュメント選択およびクラスタ化ステップ  $S 3 0 0$  の第 2 実施形態のさらに詳細な輪郭を示す図であり、如何にして、追加されるカットオフ値を有する任意のデータセットに対するノードが決定されるかを示す。ステップ  $S 3 0 0$  において開始され、制御はステップ  $S 1 3 0 5$  に続く。次に、ステップ  $S 1 3 0 5$  において、カットオフ値が、  $c$  に設定され、その結果、  $c$  未満のドキュメントを含むノードは単独のドキュメントノードによって置換されることができる。再度、収集する必要があるノードの最大数  $M$  も、設定される。次に、ステップ  $S 1 3 1 5$  において、初期フォーカスセット  $T$  が、階層  $H$  のルートノードとして設定される。次に、制御はステップ  $S 1 3 2 0$  に続く。

【 0 0 4 0 】

ステップ  $S 1 3 2 0$  において、小さいドキュメントセット  $E$  は、ゼロに設定される。次に、ステップ  $S 1 3 2 5$  において、フォーカスセットの  $k$  ノードが検査され、「最悪」ノードがピックアップされる。次に、ステップ  $S 1 3 3 0$  において、「最悪」ノードは、照合され、そのノードが、カットオフ値  $c$  未満の数のドキュメントを含むかまたはそれに等しい数のドキュメントを含むかが、決定される。ノードのドキュメントの数がカットオフ値  $c$  未満である場合、制御は、ステップ  $S 1 3 3 5$  に続く。そうではなく、選択されるノードがカットオフ値  $c$  未満の数のドキュメントを含まない場合、制御は、ステップ  $S 1 3 4 0$  に跳ぶ。

40

【 0 0 4 1 】

ステップ  $S 1 3 3 5$  において、ノード内のその数のドキュメントが、小さなドキュメントセット  $E$  に加えられる。次に、制御は、ステップ  $S 1 3 4 5$  に跳ぶ。ステップ  $S 1 3 4 0$  に

50

において、関心のあるドキュメントを含むノードの子がフォーカスセットTに加えらる。関心のあるドキュメントを含まない子は包含されず、効果的に、「除去される：pruned」。次に、制御は、ステップS 1 3 4 5に続く。

【0042】

ステップS 1 3 4 5において、フォーカスセットは、照合され、フォーカスセットが収集する必要がある最大数Mに達しているかが決定される。収集する必要があるノードの最大数Mに達している場合、制御は、ステップA 1 3 5 0に続く。その他の場合は、制御は、ステップS 1 3 2 5に跳び戻り、次の最悪ノードを見出す。

【0043】

次に、ステップS 1 3 5 0において、小さいドキュメントセットEがフォーカスセットT 10  
に加えらる。次に、ステップS 1 3 5 5において、フォーカスセットTは、クラスタ化され、クラスタPのセットが得らる。次に、ステップS 1 3 6 0において、各ノードPは、関心のあるドキュメント $I_s(n)$ によって、置換される。次に、制御は、ステップS 1 3 6 5に続く。ステップS 1 3 6 5において、制御は、図2のステップS 4 0 0に戻る。

【0044】

このように、追加されるカットオフ値を有する任意のデータセットに対して、ノードがドキュメントセットSから得られる数個のドキュメントのみを含む場合、これらのドキュメントは、ノードを拡張する時間を消費する代わりに別のセットEに追加される。

【0045】

ノードをその子によって置換する場合、「空」の子、すなわちドキュメントセットSに 20  
いかなるドキュメントも含まない子は、明白に回避することができる。「単集合（Singleton）」子、すなわち、ドキュメントSから得られる一つのドキュメントのみしか包含しない子も、特別に取り扱うことができる。ノード内に一つのドキュメントしか存在しない場合は、ノード全体が包含されない。ドキュメントが簡単に取り出され、それ自体がノードとして処理される。これは、適切な終端子孫（リーフディセendent：leaf descendent）によって子ノードを置換することと等価である。一般に、カットオフ値c未満のドキュメントを包含するノードは、c個の単独ドキュメントノードによって置換することができる。一定の数のノードのみが検査されるため、この方法によって生成される新しいノードの数も一定である。 30

【0046】

如何にして多数のノードが拡張されるかにcの値が影響を及ぼすようにすることは望ましくないので、単独のドキュメントノードは、通常のノードと別に数えられる。すなわち、単独のドキュメントノードをフォーカスセットT内に保持するのではなく、単独のドキュメントノードは、別のセットEに移動される。この処理は、フォーカスセットTが所定のサイズに達するまで続く。 $|E|$ は、定数によって限定されるので、この値は実行時間の解析に影響を及ぼさない。

【0047】

たとえば、図3のステップS 3 4 0および図4のステップS 1 3 2 5において、「最悪」ノードを決定するために使用される幾つかの「良好度」検査がある。使用することができる一つの「良好度」検査は、適合度検査または割合（RATIO）検査である。ノードが 40  
包含する大部分のドキュメントもドキュメントセットSから得られる関心のあるドキュメントである場合、ノードは、「良好」である。

【0048】

たとえば、nは、dドキュメントを有する場合、nの良好度は、下式によって表される。

【0049】

【数1】

$$g = |I_s(n)| / d$$

関数 $f(S, T)$ によって、フォーカスセットT内の最低の良好度を有するノードは返される。この関数は、僅かしか一致しないノード、すなわち一致しない子を有する可能性の 50

あるノードに有利であるので、この良好度検査は、結果として、広い範囲の除去 (pruned) となり、結果が改善される。他方、かなり良好な割合を有する大きなノードは、絶対値の項に多数の非一致ドキュメントを含む場合でも、フォーカスセット T 内にそのまま留まる。

【0050】

一つの大きなノードが、ドキュメントセット S 内に多数のドキュメントを包含する場合、割合検査は、このノードに有利である。このことはクラスタ化の場合に問題となる場合があり、その理由は、クラスタ化方法は、ノード内のドキュメントすべてを単独の実体として処理し、不均衡なクラスタサイズとなる可能性があるためである。このような大きなノードの拡張は、良好度値に重みを付けることによって促進される。たとえば、ノード n は、d ドキュメントを有する場合、ノード n の加重良好度  $g'$  は、下式によって表される。

【0051】

【数2】

$$g' = |I_S(n)| / d$$

この場合、ドキュメントセット S 内に多数のドキュメントを有することは、良好な割合の保証にはならない。実際に、比較的少数のドキュメント d を有することが、一層有利である。このことによって、出力ノードは、すべて、ドキュメントセット S から得られるほぼ等しい数のドキュメントを有することを保証することが容易になる。

【0052】

良好度を決定する他の手法は、情報理論による測定を使用する。ノードの子がノード自体より多くのドキュメントセット S に関する情報をコード化する場合、そのノードは、その子によって置換される良い候補である。このことは、親における一致は、子の間に不均一に分散され、その結果、劣悪な子は除去され、良好な子が維持されることを暗に示す。

【0053】

たとえば、ノード n が、サイズ d を有する場合、ノード  $n_i$  は、ノード n の子であり、サイズ  $d_i$  を有する。ノード n 内の情報  $I(n)$  は、下式で表される。

【0054】

【数3】

$$I(n) = - ( |I_S(n)| / d ) \cdot \log_2 ( |I_S(n)| / d )$$

ノード n に対する情報ゲイン  $G(n)$  は、下式で表される。

【0055】

【数4】

$$G(n) = I(n) - \{ ( |d_i| / |d| ) \cdot I(n_i) \}$$

ここで、 $i$  は、 $i$  についてのサメンションである。

【0056】

ノード n に対する適切な良好度測定は、 $G(n)$  によって与えられる。関数  $f(S, T)$  によって、フォーカスセット T 内の最高の情報ゲインを有するノードは返される。このことは、その子によって置換されることにより最も利益が得られるノードがピックアップされるという利点を有する。不都合なことに、これらの一致が子の間に均一に分散される場合、このことは、僅かな一致しか有しない大きなノードを無視することになる。

【0057】

本発明においては、非所定数の個別ドキュメントの代わりに、所定数のメタドキュメントが、クラスタ化または分散のための手順において使用される。メタドキュメントは、ツリー、たとえば、図5から図8までのツリーなどのメタドキュメントから得られる降順の複数の個別ドキュメントを表現する。

【0058】

図5から図8までの以下の討議の場合、本発明に従って、たとえば前述した割合検査などの幾つかの「良好度」検査の一つを使用し、「最悪」メタドキュメントを選択することができる。しかし、討議を容易にするために、図5から図8までにおいて、「最悪」メタドキュメントは、最低数の関心のあるドキュメントを有するメタドキュメントを選択するこ

10

20

30

40

50

とによって簡単に選択されるものとする。

【0059】

図5において、ツリー81のノード82～86は、個別ドキュメント、たとえば、ドキュメント88などのコレクションを表現するメタドキュメントである。たとえば、図5において、ノード89は3個の子、ドキュメント88a、88b、88cを有する内部ノードである。内部ノード89も、内部ノード84の子であり、内部ノード84自体はルートノード82の子である。ルートノード82は、ドキュメントコレクション全体を表現するメタドキュメントである。メタドキュメント83～86はメタドキュメント82から直接に得られる子である。さらに、メタドキュメント89のレベル87は、メタドキュメント83～86から直接に得られる子である。最後に、個別ドキュメント88、すなわちツリーの葉は、メタドキュメント87から直接に得られる子である。ツリー81は、説明上、非常に簡単にしてある。実際には、大きなコーパスは非常に多数の個別ドキュメントおよび便利に示す必要があるメタドキュメントのレベルを有する。

10

【0060】

一例として、10,000のドキュメントをクラスタ化し、10の話題に関連するグループ、すなわちクラスタとする場合を考える。この例の場合、同じ10,000ドキュメントの、たとえば500クラスタへの原型のクラスタ化は、既に利用可能である。互いに極端に類似しているドキュメントは、通常、同じクラスタに現れるので、500のクラスタの内の所定のクラスタのすべてのドキュメントは、所望の10のクラスタの内の同じクラスタに同様に出現するものとする。言い換えれば、細粒度クラスタ化において一緒にクラスタ化されるほど十分に類似しているドキュメントは、粗粒度クラスタ化において、一緒にクラスタ化されることになる。これは、米国特許第5,483,650号に開示されているクラスタリファインメント(refinement)仮説である。

20

【0061】

本発明は、既存のクラスタをメタドキュメントとして処理し、このメタドキュメントは全体としてコーパス全体の圧縮表現を形成する。すべての個別ドキュメントを直接にクラスタ化する代わりに、本発明は、すべての個別ドキュメントを表現するメタドキュメントをクラスタ化する。前述した例において、10,000の個別ドキュメントをクラスタ化する代わりに、本発明によれば、500のメタドキュメントをクラスタ化することができる。クラスタ洗練仮説によれば、メタドキュメントクラスタ化および個別ドキュメントクラスタ化は、同様な結果を生成する。

30

【0062】

たとえば、ステップS340からS360までの第1反復の場合、図5のフォーカスセット100は、ドキュメントコレクション全体を表現するルートノードすなわちメタドキュメント82のみを含む。当然、第1反復中は、このメタドキュメント82は、フォーカスセットTの唯一のメタドキュメントであるので、ステップS340において選択される。ステップS350において、メタドキュメント82は、その直接の子孫、すなわち子であるメタドキュメント83～86に拡張される。次に、これらの子メタドキュメント83～86を使用し、図6に示すように、フォーカスセット100において、メタドキュメント82を置換する。このようにして、フォーカスセット100は、子孫のメタドキュメント83～86を含む。

40

【0063】

次に、ステップS340が、図6のフォーカスセット100に関して反復される。フォーカスセット100内のメタドキュメントの数がステップS360における所定の最大数未満である限り、クラスタ化処理はステップS340～S360を経由して循環を継続する。メタドキュメント83～86の内、メタドキュメント84は、最低数の個別ドキュメント88を表現する。すなわち、メタドキュメント84は、6個の個別ドキュメントを表現し、一方、メタドキュメント83、85、および86は、それぞれ、7、8、および9個の個別ドキュメントを表現する。したがって、メタドキュメント84は、図7に示すように、選択され、その子孫、すなわち孫、メタドキュメント89～92に拡張される。しか

50

し、孫メタドキュメント 90 および 91 は、関心のあるドキュメントを含まないので除去される。したがって、フォーカスセット 100 は、今度は、メタドキュメント 83、85 ~ 86、89、および 92 を含む。

【0064】

所定の最大数のメタドキュメントが、ステップ S360 において、まだ実現されない場合、ステップ S340 が、図 7 に示すフォーカスセット 100 に関して反復される。最低数の個別ドキュメントを表現する子メタドキュメント 83 が、ステップ S350 において選択され、図 8 に示すように、その子孫、すなわち孫、メタドキュメント 87、および 93 ~ 95 に拡張される。しかし、メタドキュメント 95 は、関心のあるドキュメントを包含しないので、メタドキュメント 95 は除去される。したがって、フォーカスセット 100 は、ここで、子孫メタドキュメント 85 ~ 87、89、および 92 ~ 94 を包含する。

10

【0065】

図 2、図 3 または図 4、および図 5 に輪郭を示す処理は、フォーカスセット内のメタドキュメントの数が所定の最大数未満である限り継続される。所定の最大数が十分に高い場合、フォーカスセットは、実質上、個別ドキュメントを含む。その場合、ステップ S360 によって、メタドキュメントおよび個別ドキュメントの全数が所定の最大数未満であるかが決定される。しかし、この状況は、通常発生せず、特に、処理の僅かしかない第 1 反復中には発生しない。

【0066】

図 1 に示すように、再クラスタ化システム 10 は、好適には、プログラム式汎用コンピュータ上において実現される。しかし、再クラスタ化システム 10 は、専用コンピュータ、プログラム式マイクロプロセッサまたはマイクロコントローラおよび周辺一体型回路構成要素、ASIC または他の一体型回路、デジタル信号プロセッサ、有線（ハードワイヤード：hardwired）電子または論理回路たとえば個別要素（ディスクリートエレメント：discrete element）回路、PLD、PLA、FPGA、PAL などのプログラマブル論理装置、などによっても実現することができる。一般に、図 2 から図 5 に示す流れ図を実行することができる有限状態機械（finite state machine）を実現できるいかなる装置を使用しても、再クラスタ化システム 10 を実現することができる。

20

【0067】

以上、特定の実施形態について述べたが、多数の代替方法、変形、および異形は当業者には明らかであることは、明白である。したがって、前述した本発明の好適な実施形態は、説明を目的とするものであり、これに限定されるものではない。特許請求の範囲によって規定される本発明の思想および範囲を離脱することなく、種々の変化を実施し得る。

30

[付記]

[付記 1] 電子的に記憶されるドキュメントのコーパスを処理し、一つ以上の事前に識別された関心のあるドキュメントをクラスタ化する方法であって、

複数のドキュメントを代表する少なくとも一つの初期メタドキュメントを含むフォーカスセットを拡張し、複数の次のメタドキュメントとするステップであって、それぞれの次のメタドキュメントは前記初期メタドキュメントのサブセットであるステップと、

前記フォーカスセット内のメタドキュメントを選択するステップと、

40

前記選択されたメタドキュメントを子孫のメタドキュメントに拡張するステップと、

少なくとも一つの関心のあるドキュメントを含まない子孫のメタドキュメントを除去するステップと、

前記次のメタドキュメントの数が少なくとも所定の最大数に等しくなるまで、前記選択および拡張ステップを反復するステップと、

を含む拡張ステップと、

前記次のメタドキュメントをクラスタ化し、所定数のクラスタとするステップと、を含むことを特徴とする方法。

[付記 2] 付記 1 に記載の方法において、クラスタの前記所定数は、前記所定最大数未満であることを特徴とする方法。

50

[付記3] 付記1に記載の方法において、前記少なくとも一つの初期メタドキュメントは、ドキュメントのコーパス全体を代表する単独のメタドキュメントであることを特徴とする方法。

[付記4] 付記1に記載の方法において、前記所定最大数は、前記拡張および選択ステップが、共に、時間制約内に完了するように決定されることを特徴とする方法。

[付記5] 付記1に記載の方法において、さらに、前記新メタドキュメントの要約を確定するステップと、前記要約をユーザに提示するステップと、を含むことを特徴とする方法。

[付記6] 付記5に記載の方法において、前記要約は、各新メタドキュメントにおいて最も頻繁に現れる固定数の話題のワードと、各新メタドキュメント内の少なくとも一つの典型的なドキュメントの名称と、を含むことを特徴とする方法。

10

[付記7] 付記1に記載の方法において、前記クラスタ化ステップは、各メタドキュメントによって表現されるドキュメントの数に関係なく、多くても、所定量の時間を要することを特徴とする方法。

[付記8] 付記1に記載の方法において、前記拡張ステップは、さらに、選択されたメタドキュメント内の関心のあるドキュメントの数がカットオフ値を超えるかを決定するステップを含むことを特徴とする方法。

[付記9] 付記8に記載の方法において、前記選択されたメタドキュメント内の関心のあるドキュメントの数がカットオフ値未満である場合、前記ドキュメントは別のドキュメントセットに加えられることを特徴とする方法。

20

[付記10] 付記9に記載の方法において、前記拡張ステップは、さらに、前記クラスタ化ステップにおいてクラスタ化するために、前記別のドキュメントセットを前記次のメタドキュメントに加えるステップを含むことを特徴とする方法。

[付記11] 電子的に記憶されるドキュメントのコーパスを処理し、少なくとも一つの事前に識別された関心のあるドキュメントをクラスタ化する装置であって、

複数のドキュメントを代表する少なくとも一つの初期メタドキュメントを含むフォーカスセットを拡張し、複数の次のメタドキュメントとする拡張手段であって、それぞれの次のメタドキュメントは前記少なくとも一つの初期メタドキュメントのサブセットである拡張手段と、

30

前記フォーカスセット内のメタドキュメントを選択するための選択手段であって、選択されたメタドキュメントは拡張手段によってその子孫のメタドキュメントに拡張される選択手段と、

少なくとも一つの関心のあるドキュメントを含まない子孫のメタドキュメントを除去するための除去手段と、

を含む拡張手段と、

前記次のメタドキュメントをクラスタ化し、所定数のクラスタとするためのクラスタ化手段と、

を備え、

40

前記拡張手段は、前記次のメタドキュメントの数が所定の最大数に少なくとも等しくなるまで、前記フォーカスセットを拡張することを特徴とする装置。

[付記12] 付記11に記載の装置において、新メタドキュメントの前記所定数は、前記所定最大数未満であることを特徴とする装置。

[付記13] 付記11に記載の装置において、前記少なくとも一つの初期メタドキュメントは、ドキュメントのコーパス全体を代表する単独のメタドキュメントであることを特徴とする装置。

[付記14] 付記11に記載の装置において、前記所定最大数は、前記クラスタ化手段が前記次のメタドキュメントのクラスタ化を時間制約内に完了するように決定されることを特徴とする装置。

50

【付記 15】 付記 11 に記載の装置において、さらに、  
新メタドキュメントの要約を確定し、前記要約をユーザに提示する要約手段を備えることを特徴とする装置。

【付記 16】 付記 15 に記載の装置において、前記要約は、  
各新メタドキュメントにおいて最も頻繁に現れる固定数の話題のワードと、  
各新メタドキュメント内の少なくとも一つの典型的なドキュメントの名称と、  
を含むことを特徴とする装置。

【付記 17】 付記 11 に記載の装置において、前記クラスタ化手段は、各メタドキュメントによって表現されるドキュメントの数に関係なく、多くても、所定量の時間を要することを特徴とする方法。

10

【付記 18】 付記 11 に記載の装置において、前記拡張手段は、前記メタドキュメント内の関心のあるドキュメントの数がカットオフ値を超えるかを決定することを特徴とする装置。

【付記 19】 付記 18 に記載の装置において、前記メタドキュメント内の関心のあるドキュメントの数がカットオフ値未満である場合、前記ドキュメントは別のドキュメントセットに加えられることを特徴とする装置。

【付記 20】 付記 19 に記載の装置において、前記拡張手段は、クラスタ化手段によってクラスタ化するために、前記別のドキュメントセットを前記次のメタドキュメントに加えることを特徴とする装置。

#### 【図面の簡単な説明】

20

【図 1】 本発明による装置の一実施形態を示すブロック図である。

【図 2】 本発明による再クラスタ化方法の一実施形態の輪郭を示す流れ図である。

【図 3】 図 2 のメタドキュメント拡張ステップの第 1 実施形態の輪郭をより詳細に示す流れ図である。

【図 4】 図 2 のメタドキュメント拡張ステップの第 2 実施形態の輪郭をより詳細に示す流れ図である。

【図 5】 本発明の好適な実施形態によるフォーカスセットのツリーおよび変化するコンテンツを示す図である。

【図 6】 本発明の好適な実施形態によるフォーカスセットのツリーおよび変化するコンテンツを示す図である。

30

【図 7】 本発明の好適な実施形態によるフォーカスセットのツリーおよび変化するコンテンツを示す図である。

【図 8】 本発明の好適な実施形態によるフォーカスセットのツリーおよび変化するコンテンツを示す図である。

【図 9】 分散 / 集合手順を広く示す図である。

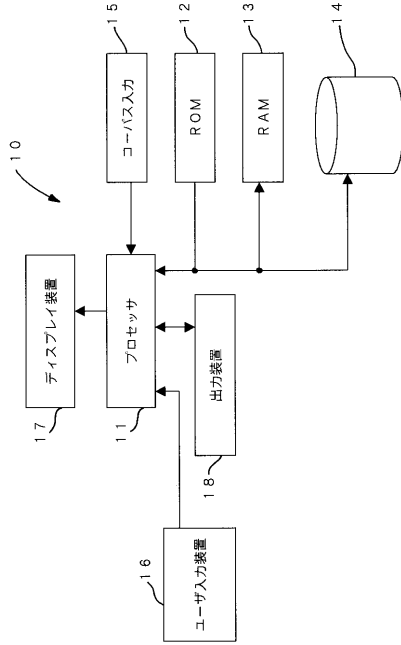
【図 10】 従来の分散 / 集合ドキュメントブラウジング法を、ドキュメントの特定のコーパスに適用する場合を示す図である。

#### 【符号の説明】

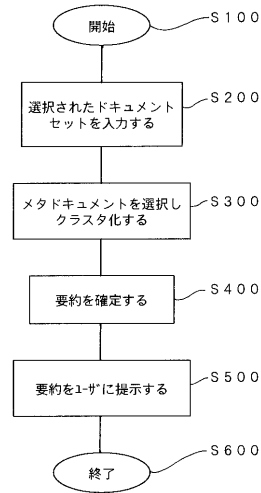
10 再クラスタ化システム、11 プロセッサ、12 ROM、13 RAM、14 不揮発性メモリ、15 コーパス入力、16 ユーザ入力装置、17 ディスプレイ装置、18 出力装置、81 ツリー、82 ルートノード(メタドキュメント)、83~87、89~95 ノード(メタドキュメント)、88 ドキュメント、100 フォーカスセット。

40

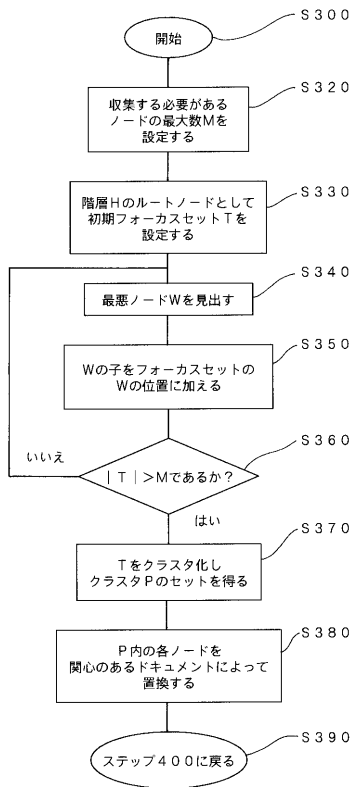
【図1】



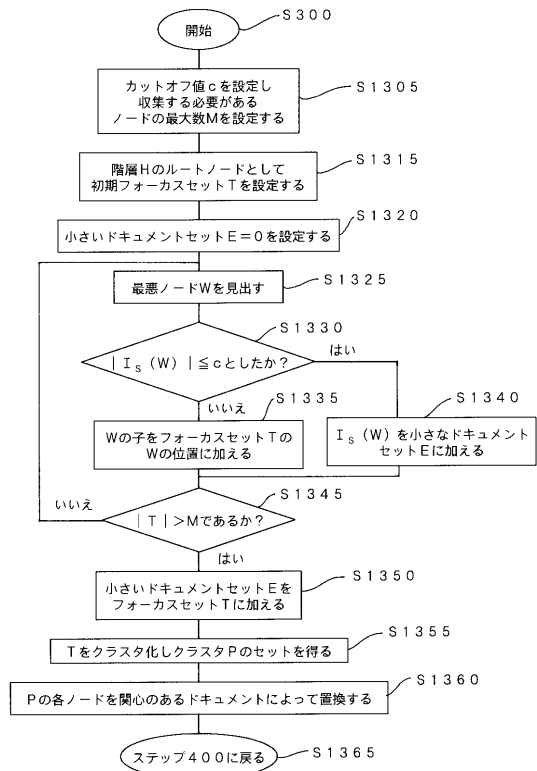
【図2】



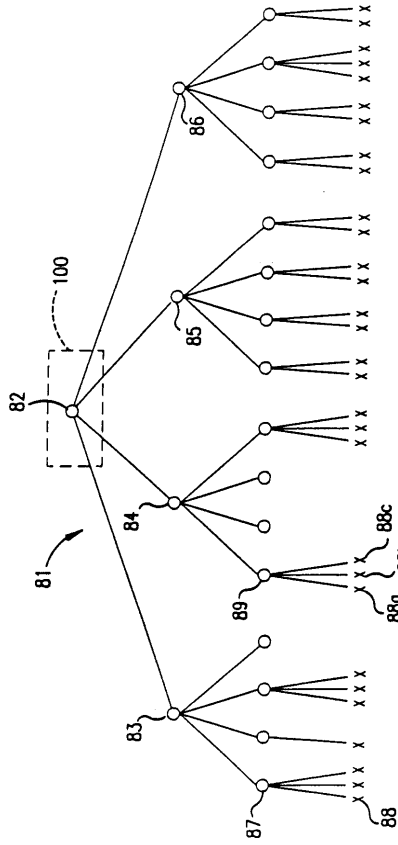
【図3】



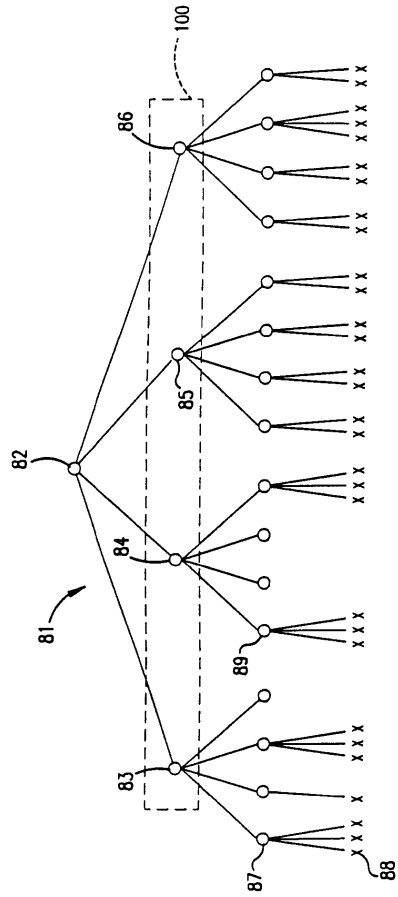
【図4】



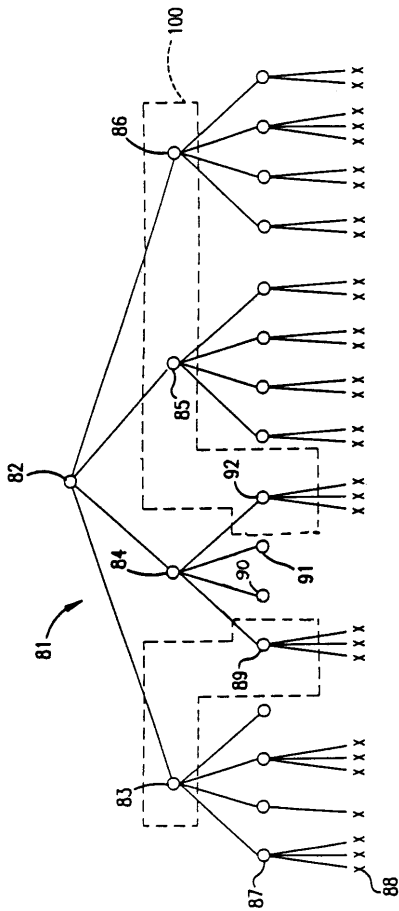
【 図 5 】



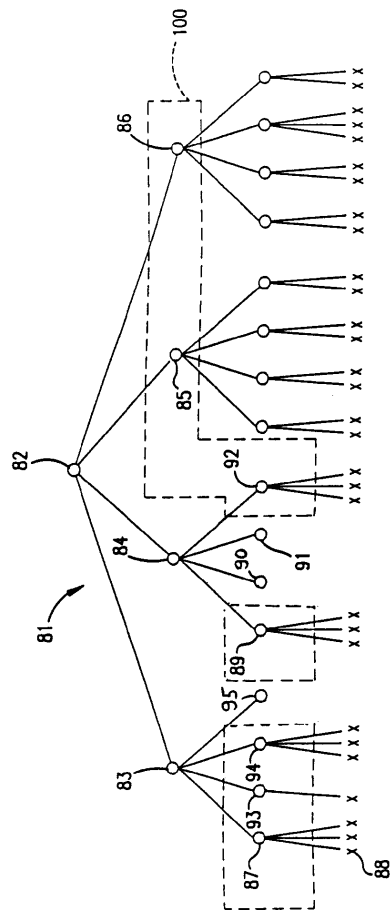
【 図 6 】



【 図 7 】

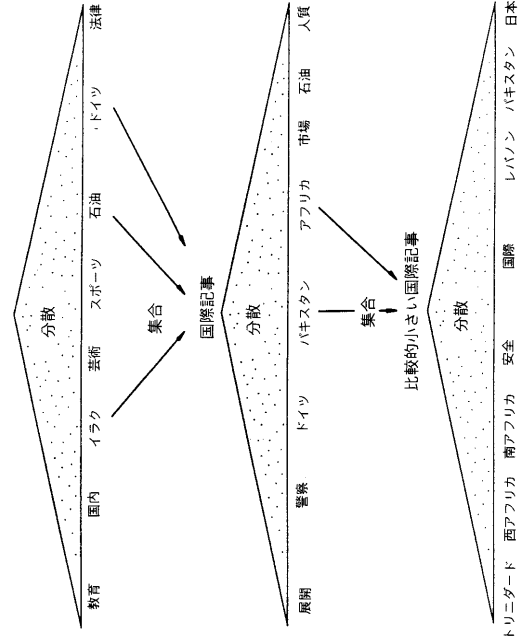


【 図 8 】

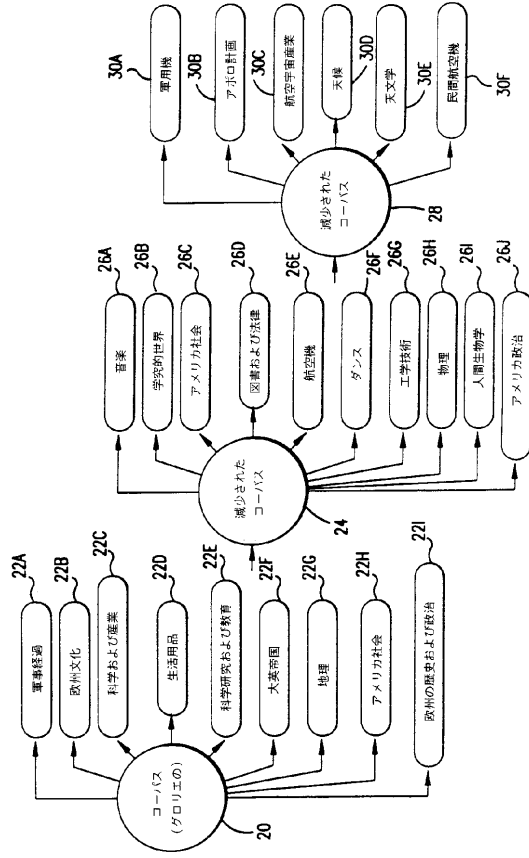


【図9】

ニューヨークタイムズニュースサービス、1990年8月



【図10】



---

フロントページの続き

審査官 辻本 泰隆

(56)参考文献 特開平08-153121(JP,A)  
特開平05-225256(JP,A)  
特開平8-235198(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G06F 17/30