

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2010-501096  
(P2010-501096A)

(43) 公表日 平成22年1月14日(2010.1.14)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 414A	5B075
	G06F 17/30 210D	
	G06F 17/30 419A	

審査請求 未請求 予備審査請求 未請求 (全 26 頁)

(21) 出願番号 特願2009-524708 (P2009-524708)  
 (86) (22) 出願日 平成19年8月16日 (2007. 8. 16)  
 (85) 翻訳文提出日 平成21年2月16日 (2009. 2. 16)  
 (86) 国際出願番号 PCT/US2007/018417  
 (87) 国際公開番号 W02008/021561  
 (87) 国際公開日 平成20年2月21日 (2008. 2. 21)  
 (31) 優先権主張番号 11/465, 026  
 (32) 優先日 平成18年8月16日 (2006. 8. 16)  
 (33) 優先権主張国 米国 (US)

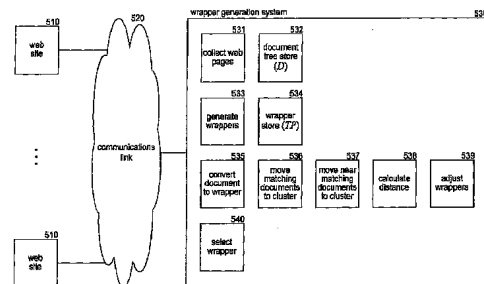
(71) 出願人 500046438  
 マイクロソフト コーポレーション  
 アメリカ合衆国 ワシントン州 9805  
 2-6399 レッドモンド ワン マイ  
 クロソフト ウェイ  
 (74) 代理人 100077481  
 弁理士 谷 義一  
 (74) 代理人 100088915  
 弁理士 阿部 和夫  
 (72) 発明者 チーロン ウェン  
 アメリカ合衆国 98052 ワシントン  
 州 レッドモンド ワン マイクロソフト  
 ウェイ マイクロソフト コーポレーシ  
 ョン インターナショナル パテンツ内

最終頁に続く

(54) 【発明の名称】 ラッパー生成およびテンプレート検出の協同最適化

(57) 【要約】

テンプレート検出およびラッパー生成を協同して最適化することによって階層的に編成された文書のラッパーを生成する方法およびシステムを提供する。ラッパー生成システムは、文書ツリーのクラスタを識別することと、クラスタのラッパーツリーを生成することとによって、類似するテンプレートを有する文書のラッパーを生成する。ラッパーツリーは、クラスタのテンプレートと一致する文書のラッパーを定義する。ラッパー生成システムは、初期文書ツリーに基づいてクラスタのラッパーツリーを生成することによって、文書ツリーをクラスタ化する。ラッパー生成システムは、次に、任意の他の文書ツリーがクラスタのラッパーツリーと一致またはほぼ一致するかどうかを繰り返して判定し、そうである場合には、その文書ツリーをクラスタに追加し、適当にラッパーツリーを調整し、その結果、新たに追加された文書ツリーを含むすべての文書ツリーがラッパーツリーと一致するようにする。



## 【特許請求の範囲】

## 【請求項 1】

階層的に編成された文書のラッパーを生成するコンピューティング装置での方法であって、各文書は、文書ツリーを有し、前記方法は、

文書ツリーのラッパーツリーを作成すること（535）と、

前記ラッパーツリーへのその距離がしきい値以内である文書ツリーを選択すること（537）と、

前記文書ツリーに基づいて前記ラッパーツリーを調整すること（539）と

を含み、前記ラッパーは、前記調整されたラッパーツリーに基づくことを特徴とする方法。

10

## 【請求項 2】

前記ラッパーツリーの前記調整は、前記選択された文書ツリーのラッパーツリーおよび前記作成されたラッパーツリーをマージすることを含むことを特徴とする請求項 1 に記載の方法。

## 【請求項 3】

文書ツリーのラッパーツリーの前記作成は、前記文書ツリー内のノードの連続するサブフォレストを組み合わせることを含むことを特徴とする請求項 2 に記載の方法。

## 【請求項 4】

複数の文書ツリーについて前記選択および調整を実行することを含むことを特徴とする請求項 1 に記載の方法。

20

## 【請求項 5】

前記ラッパーツリーは、前記ラッパーツリーがそれに関して作成される前記文書ツリーおよび前記ラッパーツリーがそれに関して調整される前記文書ツリーを含む文書ツリーのクラスタに関することを特徴とする請求項 1 に記載の方法。

## 【請求項 6】

前記ラッパーツリーが生成された後に、文書ツリーのもう 1 つのクラスタについてもう 1 つのラッパーツリーを生成することを含むことを特徴とする請求項 5 に記載の方法。

## 【請求項 7】

文書ツリーとラッパーツリーとの間の距離は、文書ノードと一致しないラッパーノードの個数およびラッパーノードと一致しない文書ノードの個数に基づくことを特徴とする請求項 1 に記載の方法。

30

## 【請求項 8】

前記距離は、前記文書ツリーおよび前記ラッパーツリーの重さに基づいて正規化されることを特徴とする請求項 7 に記載の方法。

## 【請求項 9】

複数のラッパーツリーが生成された時に、文書ツリーと前記ラッパーツリーとの間の距離に基づいて前記文書ツリーからデータを抽出するのに使用すべきラッパーを識別することを特徴とする請求項 1 に記載の方法。

## 【請求項 10】

前記しきい値は、適応式であることを特徴とする請求項 1 に記載の方法。

40

## 【請求項 11】

前記しきい値は、増やされたしきい値のゆえに選択される文書に基づく前記ラッパーツリーの前記調整が前記ラッパーツリーの有効性を下げるまで増やされることを特徴とする請求項 10 に記載の方法。

## 【請求項 12】

階層的に編成された文書とラッパーツリーとの間の類似性を判定するコンピューティングシステムであって、前記文書は、文書ツリーを有し、前記システムは、

前記文書ツリーのノードを前記ラッパーツリーのノードに位置合せするコンポーネント（1101）と、

位置合せされないノードの個数からメトリックを生成するコンポーネント（1102～

50

1106)であって、前記メトリックは、前記文書ツリーと前記ラッパーツリーとの間の類似性を示す、コンポーネント(1102~1106)とを含むことを特徴とするコンピューティングシステム。

【請求項13】

前記メトリックは、前記ラッパーツリーのノードと位置合せされない前記文書ツリーのノードの個数および前記文書ツリーのノードと位置合せされない前記ラッパーツリーのノードの個数に基づくことを特徴とする請求項12に記載のコンピューティングシステム。

【請求項14】

位置合せされないノードの個数は、前記ツリーの重さに基づいて正規化されることを特徴とする請求項13に記載のコンピューティングシステム。

【請求項15】

前記メトリックは、次式

【数1】

$$\Psi(T_w, T_d) = \left( \frac{C_w(T_w, T_d)}{W(T_w)} + \frac{C_d(T_w, T_d)}{W(T_d)} \right) / 2$$

によって表されることを特徴とする請求項12に記載のコンピューティングシステム。

【請求項16】

ラッパーツリーごとに、

以前に選択されたことがない文書ツリーを選択すること(603)と、

前記選択された文書ツリーの前記ラッパーツリーを作成すること(604)と、

前記ラッパーからのその距離がしきい値未満である選択されていない文書ツリーが存在する時に、前記文書ツリーを選択し(606、607)、前記選択された文書ツリーに基づいて前記ラッパーツリーを調整することと

を含む方法によって、文書ツリーのラッパーツリーを生成するためにコンピューティングシステムを制御する命令を含むことを特徴とするコンピュータ可読媒体。

【請求項17】

前記距離が0である時に、前記ラッパーツリーを調整せずに前記文書ツリーを選択することを特徴とする請求項16に記載のコンピュータ可読媒体。

【請求項18】

ラッパーツリーの前記選択された文書ツリーは、文書ツリーのクラスタを形成することを特徴とする請求項17に記載のコンピュータ可読媒体。

【請求項19】

前記距離は、前記文書ツリーと前記ラッパーツリーとの間の位置合せされないノードの個数に基づくことを特徴とする請求項16に記載のコンピュータ可読媒体。

【請求項20】

文書の文書ツリーからのその距離が最小である前記ラッパーツリーを識別することと、前記識別されたラッパーツリーを前記文書からデータを抽出するためのテンプレートとして使用することとによって、前記文書からデータを抽出することを含むことを特徴とする請求項16に記載のコンピュータ可読媒体。

【発明の詳細な説明】

【背景技術】

【0001】

ワールドワイドウェブ(「ウェブ」)は、ウェブページを介してアクセス可能な膨大な量の情報を提供する。ウェブページは、スタティックコンテンツまたはダイナミックコンテンツを含むことができる。スタティックコンテンツは、一般に、ウェブページの多数のアクセスにまたがって同一のままであることができる情報を指す。ダイナミックコンテンツは、一般に、ウェブデータベース内に格納され、検索要求に回答してウェブページに追加される情報を指す。ダイナミックコンテンツは、ディープウェブまたは隠しウェブと呼

10

20

30

40

50

ばれてきたものを表す。

【0002】

多くの検索エンジンサービスは、ユーザがウェブのスタティックコンテンツを検索することを可能にする。ユーザが、検索用語を含む検索要求またはクエリをサブミットした後に、検索エンジンサービスは、これらの検索用語に関連する可能性があるウェブページを識別する。これらのウェブページが、検索結果である。関連するウェブページをすばやく識別するために、検索エンジンサービスは、ウェブページへのキーワードのマッピングを維持する場合がある。このマッピングは、各ウェブページのキーワードを識別するためにウェブを「クロール」することによって生成することができる。ウェブをクロールするために、検索エンジンサービスは、ルートウェブページを介してアクセス可能なすべてのウェブページを識別するためにルートウェブページのリストを使用する場合がある。特定のウェブページのキーワードは、見出しとしての単語の識別、ウェブページのメタデータ内で供給される単語、強調表示される単語など、さまざまな周知の情報検索技法を使用して識別することができる。

10

【0003】

しかし、これらの検索エンジンサービスは、一般に、非クロール可能コンテンツとも考えられるダイナミックコンテンツの検索を提供しない。多くのウェブページは、構造化されたソース（たとえば、リレーショナルデータベース）から生成されたダイナミックページを含む。そのようなダイナミックコンテンツを含むウェブページが生成される時に、基礎になる構造化されたソースの構造化データは、構造化されない形または半構造化された形でウェブページ内においてエンコードされる。そのようなダイナミックコンテンツの検索に関連する1つの問題は、ウェブページから対応する構造化されたソースのスキーマを識別することが難しいことである。スキーマは、基礎になる構造化ソースに格納された情報または属性を定義する。この問題のゆえに、そのようなダイナミックコンテンツを有するウェブページの照会は、しばしば、不満足な結果をもたらす。

20

【0004】

ウェブページのダイナミックコンテンツのスキーマを識別する試みが行われ、その結果、検索を容易にするためにより構造化されたフォーマットにコンテンツを変換できるようになってきた。構造化されたフォーマットでのウェブページからの情報およびその編成の抽出は、「ラッパー」と呼ばれるプログラムによって実行される。ウェブサイトのウェブページ用のラッパーを手作業で生成するのは、時間がかかる可能性がある。したがって、ダイナミックコンテンツを提供する数千個のウェブサイトの数百万個のウェブページのラッパーを手作業で生成するのは、非実用的である。

30

【0005】

いくつかの自動ラッパー「誘導」システムまたは自動ラッパー生成システムが開発されてきた。ラッパー誘導は、ウェブページのダイナミックコンテンツのスキーマを学習することと、そのウェブページからデータを抽出し、抽出されたデータをスキーマによって識別される構造化フォーマットで格納するラッパーを生成することというプロセスである。これらの自動ラッパー誘導システムは、ラッパーの表現力に対して有効性をトレードオフするものである。有効性とは、ラッパー誘導プロセスには使用されないが、同一の「テンプレート」を共有するウェブページからコンテンツを抽出する際にラッパーがどれほど正確であるかを指す。ラッパー誘導システムは、ウェブページのトレーニングセットを使用してテンプレートに関してラッパーを生成する。次に、ラッパーは、同一のテンプレートを共有するウェブページからデータを抽出するのに使用される。表現力は、ラッパーのテンプレートによって識別されるラッパーによって処理できるウェブページの範囲を指す。ラッパーをより表現力のあるものにするために、ラッパー誘導システムは、一般に、ラッパーにワイルドカード（たとえば、「\*」）を導入し、その結果、より多くのウェブページがラッパーの範囲に含まれるようにする。しかし、一般に、ラッパーの表現力が高まるほど、その有効性が下がり、逆も同様である。

40

【0006】

50

有効性と表現力との間の受け入れられるトレードオフを提供するために、通常のラッパー誘導システムは、トレーニングウェブページを、ウェブページ上のダイナミックコンテンツの編成を表すテンプレートに従ってクラスタ化する。したがって、類似する編成を用いる（たとえば、同一のテンプレートを有する）ウェブページは、一緒にクラスタ化される。これらのラッパー誘導システムは、あるクラスタ内のウェブページのラッパーを自動的に生成することができる。クラスタのウェブページが類似するので、そのようなラッパーは、表現力を高めながら許容できる有効性を達成するために、限られたワイルドカードを使用する可能性がある。

【先行技術文献】

【非特許文献】

【0007】

【非特許文献1】Liu, B., Grossman, R., and Zhai, Y., "Mining Data Records in Web Pages," Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2003, pp. 601-606

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかし、そのような通常のラッパー誘導システムによって生成されるラッパーの正確さは、大部分、同一のテンプレートを有するウェブページを正しくクラスタ化することの正確さに依存する。一部のラッパー誘導システムは、単純に、ウェブページのURLの間の類似性に基づいてウェブページをクラスタ化する。このクラスタ化の単純な手法は、ウェブサイトが、そのウェブサイトの同一のサブディレクトリ内の同一のテンプレートを使用するウェブページを格納する場合には適当である。その場合に、それらのURLは、サブディレクトリの位置を示すために同一のプレフィックスを有する。しかし、多くのウェブサイトは、ウェブページのURLを定義する時に、はるかにより複雑な手法を使用する。その結果、類似するURLを有するウェブページが、非常に異なるテンプレートを有する可能性があり、非常に異なるURLを有するウェブページが、非常に類似するテンプレートを有する場合がある。したがって、ウェブページの編成の類似性に基づいてウェブページを正確にクラスタ化することは、非常に難しい可能性があり、有効性と表現力との間の許容できないトレードオフを有するラッパーがもたらされる。

【課題を解決するための手段】

【0009】

テンプレート検出とラッパー生成とを協同して最適化することによって階層的に編成された文書のラッパーを生成する方法およびシステムを提供する。ラッパー生成システムは、文書のテンプレートを検出すると同時に、文書のラッパーを動的に生成する。ラッパー生成システムは、ラベル付き葉ノードを有する文書ツリーとして文書を表す。ラッパー生成システムは、文書ツリーのクラスタを識別することと、クラスタのラッパーツリーを生成することとによって、類似するテンプレートを有する文書のラッパーを生成する。ラッパーツリーは、クラスタのテンプレートと一致する文書のラッパーを定義する。ラッパー生成システムは、初期文書ツリーに基づいてクラスタのラッパーツリーを生成することによって文書ツリーをクラスタ化する。次に、ラッパー生成システムは、他の文書ツリーがそのクラスタのラッパーツリーと一致またはほぼ一致するかどうかを繰り返して判定し、そうである場合には、その文書ツリーをクラスタに追加し、新たに追加された文書ツリーを含むすべての文書ツリーがラッパーツリーと一致するようにするために、適当にラッパーツリーを調整する。ラッパーツリーと一致するかほぼ一致する文書ツリーがもうない時には、そのクラスタのラッパーツリーの生成が完了する。次に、ラッパー生成システムは、このプロセスを繰り返して、文書ツリーの新しいクラスタを形成し、そのラッパーツリーを生成する。次に、ラッパー生成システムは、これらのラッパーツリーを使用して、ラッパーを定義する。新しい文書のデータを抽出するために、新しい文書の文書ツリーが生成され、最もよく一致するラッパーツリーが識別され、その識別されたラッパーツリー

10

20

30

40

50

のラッパーが、データを抽出するのに使用される。

【0010】

この「課題を解決するための手段」は、下の「発明を実施するための形態」でさらに説明する概念の選択物を単純化された形で導入するために提供されたものである。この「課題を解決するための手段」は、請求される主題の主要な特徴または本質的特徴を識別することを意図されたものではなく、請求される主題の範囲を判定する際の助けとして使用されることを意図されたものでもない。

【図面の簡単な説明】

【0011】

【図1】一実施形態のラッパー生成システムの動作の高水準概要を示すブロック図である

10

【図2】文書ツリーのラッパーツリーへの変換を示す図である。

【図3】一実施形態でのラッパーツリーのマージを示す図である。

【図4】ラッパーツリーと文書ツリーとの位置合わせを示す図である。

【図5】一実施形態でのラッパー生成システムのコンポーネントを示すブロック図である

【図6】一実施形態でのラッパー生成システムのラッパー生成コンポーネントの処理を示す流れ図である。

【図7】一実施形態でのラッパー生成システムのラッパーへの文書変換コンポーネントの処理を示す流れ図である。

20

【図8】一実施形態でのラッパー生成システムのクラスタへの一致する文書の移動コンポーネントの処理を示す流れ図である。

【図9】一実施形態でのラッパー生成システムのクラスタへのほぼ一致する文書の移動コンポーネントの処理を示す流れ図である。

【図10】一実施形態でのラッパー生成システムのほぼ一致する文書のチェックコンポーネントの処理を示す流れ図である。

【図11】一実施形態でのラッパー生成システムの距離計算コンポーネントの処理を示す流れ図である。

【図12】一実施形態でのラッパー生成システムのラッパー調整コンポーネントの処理を示す流れ図である。

30

【図13】一実施形態でのラッパー生成システムのラッパー位置合わせコンポーネントの処理を示す流れ図である。

【発明を実施するための形態】

【0012】

テンプレート検出およびラッパー生成を協同して最適化することによって階層的に編成された文書のラッパーを生成する方法およびシステムを提供する。一実施形態で、ラッパー生成システムは、文書のテンプレートを検出すると同時に文書のラッパーを動的に生成する。ラッパー生成システムは、ウェブページなど、それぞれが文書ツリーと称するその階層のツリー構造によって表される、階層的に編成された文書のコレクションを与えられる。たとえば、ウェブページが、ドキュメントオブジェクトモデル(「DOM」)に準拠する場合に、文書ツリーは、DOM文書のタグに対応するノードを含む。ラッパー生成システムは、文書ツリーの葉ノードのラベルをも与えられる。これらのラベルは、基礎になる構造化データの識別子、フィールド、または属性に対応する。たとえば、自動車に関する情報を含むウェブページは、「メーカー」、「モデル」、「年式」、「色」、「価格」などとしてラベル付けされた葉ノードを有する可能性がある。ラッパー生成システムは、コレクションの文書ツリーを選択し、その文書ツリーに対応するラッパーツリーを生成することによって、類似するテンプレートを有する文書のラッパーを生成する。ラッパーツリーは、検出されたテンプレートと一致する文書のラッパーを定義する。ラッパー生成システムは、コレクションの任意の他の文書ツリーがラッパーツリーとほぼ一致するかどうかを判定する。そうである場合には、ラッパー生成システムは、文書ツリーが、ラッパー

40

50

を生成するのに使用された文書ツリーと同一のテンプレートを共有すると考える。ラッパー生成システムは、ほぼ一致する文書ツリーについてラッパーツリーを調整し、その結果、両方の文書ツリーがそのラッパーツリーと一致するようにする。文書ツリーが、たまたまラッパーツリーと正確に一致する場合には、ラッパー生成システムは、ラッパーツリーを調整する必要がない。次に、ラッパー生成システムは、コレクションの任意の他の文書ツリーが調整されたラッパーツリーとほぼ一致するかどうかの判定およびラッパーツリーの調整を、調整されたラッパーツリーとほぼ一致するコレクション内の文書ツリーがなくなるまで繰り返す。ラッパーツリーの生成および調整に使用される文書ツリーを、類似するテンプレートを有し、ラッパーツリーと一致する文書ツリーの「クラスタ」と称する。次に、ラッパー生成システムは、新しいクラスタを形成するための文書ツリーの選択、その文書ツリーのラッパーツリーの生成、およびクラスタに追加されるほぼ一致するツリーに関するラッパーツリーの動的調整というプロセスを繰り返す。この形で、ラッパー生成システムは、動的に生成されるラッパーツリーに基づいてテンプレートを検出し、動的に識別されたテンプレートに基づいてラッパーツリーを生成する。次に、ラッパー生成システムは、ラッパーツリーを使用してラッパーを定義する。

10

20

30

40

50

**【0013】**

一実施形態で、ラッパー生成システムは、距離メトリックを使用して、文書ツリーがラッパーツリーと一致し、またはほぼ一致するかどうかを判定する。ラッパーツリーは、類似するテンプレートを有する文書ツリーを定義するラッパーノードの階層であり、テンプレートを有する文書ツリーからデータを抽出するのに使用されるラッパーを表す。ラッパーツリーのラッパーノードは、ラッパーのノードと一致する文書ツリーの対応する文書ノードを定義するタグ、ラベル、またはワイルドカードを含むことができる。文書ツリーとラッパーツリーとの間の距離を判定するために、ラッパー生成システムは、文書ノードをラッパーノードに位置合せする。次に、ラッパー生成システムは、ラッパーノードと位置合せされないツリーノードの個数と、ツリーノードと位置合せされないラッパーノードの個数とをカウントする。ラッパー生成システムは、これらの位置合せされないノードの個数に基づいて距離メトリックを計算する。ラッパー生成は、距離メトリックを正規化し、その結果、同一の距離しきい値を、そのサイズに関わりなく文書ツリーとラッパーツリーとの間のほぼ一致を定義するのに使用できるようにすることもできる。たとえば、5つの位置合せされないノードを有する文書ツリーは、100個のノードを有するラッパーツリーとほぼ一致するが、10個のノードのみを有するラッパーツリーとほぼ一致はしないものとすることができる。

**【0014】**

一実施形態で、ラッパー生成システムは、固定しきい値または適応しきい値のいずれかを使用して、文書ツリーがラッパーツリーとほぼ一致するかどうかを判定することができる。固定しきい値を使用する時には、ラッパー生成システムは、ラッパー生成プロセス全体を通じて、文書ツリーがラッパーツリーとほぼ一致するかどうかを判定するのに同一のしきい値を使用する。適応しきい値を使用する時には、ラッパー生成システムは、当初に、文書ツリーをクラスタ化する時に小さいしきい値を使用することができる。クラスタのラッパーツリーとほぼ一致する文書ツリーがもうない場合に、ラッパー生成システムは、しきい値を増やし、その後、増やされたしきい値を使用し、ラッパーツリーを適当に調整することによって、ラッパーツリーとほぼ一致する文書ツリーをクラスタに追加することができる。その後、ラッパー生成システムは、増やされたしきい値に基づいて調整されたラッパーツリーが、以前のしきい値を用いるラッパーツリーよりよく動作するかどうかをテストする。ラッパー生成システムは、新たに調整されたラッパーツリーに基づくラッパーを使用してクラスタの文書ツリーからデータを抽出することによって、性能をテストすることができる。この性能が、大幅によりよいものではない場合には、ラッパー生成システムは、増やされたしきい値のゆえに追加された文書ツリーを除外するように最終的なクラスタをセットすることと、最終的なラッパーツリーに、増やされたしきい値のゆえに調整される前のラッパーツリーをセットすることとによって、増やされたしきい値の影響を

ロールバックする。しかし、調整されたラッパーツリーの性能の方がよい場合には、ラッパー生成システムは、もう一度しきい値を増やし、ほぼ一致する文書ツリーをクラスタに追加する。ラッパー生成システムは、増やされたしきい値に基づいて調整されたラッパーツリーが以前のしきい値を用いるラッパーツリーよりよくまたは大幅によりよく動作しなくなるまで、このプロセスを継続する。

#### 【0015】

図1は、一実施形態のラッパー生成システムの動作の高水準概要を示すブロック図である。このラッパー生成システムは、ウェブページなどの文書のトレーニングセットを与えられる101。このラッパー生成システムは、次に、ウェブページを解析して102、文書ツリーを生成し、その葉ノードのラベルを入力する。このラッパー生成システムは、テンプレートを協同して検出し103、動的に調整されるラッパーツリーに基づいて文書ツリーをクラスタ化することによってラッパーを生成する103。次に、このラッパー生成システムは、ラッパーツリーに基づいて定義されたラッパーを出力する104。次に、このラッパー生成システムは、ラッパーを使用してウェブページからデータを抽出することができる。このラッパー生成システムは、新しいウェブページを受け取る105時に、新しいウェブページを解析して106、文書ツリーを生成する。次に、このラッパー生成システムは、距離メトリックに基づいて、生成された文書ツリーに最も近いラッパーツリーを選択し107、選択されたラッパーツリーに対応するラッパーを使用してツリーからデータを抽出する108。一実施形態で、文書は、XMLフォーマットで表すことができる。

#### 【0016】

ラッパー生成システムは、「括弧」ノードを含めることができる各ラッパーノードに割り当てられた記号を有する修正DOMツリー(modified DOM tree)としてラッパーツリーを表す。ラッパー生成システムは、ラッパーノードの記号を使用して、文書ツリーとラッパーツリーとの間の距離を計算する時に、文書ノードおよびラッパーノードが位置合せされているかどうかを判定する。ラッパー生成システムは、ラッパーノードの記号S( )が、1、整数N(N ≥ 2)、またはワイルドカードすなわち?、+、および\*のうちの1つになるように定義する。1の記号は、ラッパーノードが1つのツリーノードだけと一致することができることを意味する。Nの記号は、ラッパーノードがN個の連続するツリーノードだけと一致することができることを意味する。?の記号は、ラッパーノードが0個または1個のツリーノードだけと一致することができることを意味する。+の記号は、ラッパーノードが連続するN個のツリーノードと一致することができる、N ≥ 1であることを意味する。\*の記号は、ラッパーノードが0個または連続するN個のツリーノードと一致することができる、N ≥ 1であることを意味する。?または\*の記号を有するラッパーノードは、ツリーノードと一致しない可能性があるため、「ソフト」ノードと考えられる。他のすべてのラッパーノードは、「ハード」ノードと考えられる。ラッパーツリーには、「ブロック」ノードと称する特殊なノードを含めることができる。括弧ノードは、タグを有しておらず、括弧の対のように振る舞い、したがって、葉ノードになることはできない。ラッパーツリーの他のすべてのノードは、「タグ」ノードと呼ばれる。

#### 【0017】

一実施形態で、ラッパー生成システムは、文書ツリーとラッパーツリーとの間の距離メトリックを、文書ツリーおよびラッパーツリーの重さに基づいて正規化する。ラッパー生成システムは、ツリーノードのツリーノード重さW( )を、をルートとするサブツリーのノード数と等しくなるように定義する。そのルートノードがであるツリーT<sub>d</sub>の文書ツリー重さW(T<sub>d</sub>)は、W( )である。ラッパー生成システムは、ラッパーノードのラッパーノード重さW( )を、そのラッパーノードがソフトノードである時には0になり、そのラッパーノードがハードノードである時にはその子ノードの重さの合計になり、そのラッパーノードがハード葉タグノードである時には1になり、そのラッパーノードがハード非葉タグノードである時には1にその子ノードの重さの合計を加えたもの



になるように定義する。ラッパーツリー重さは、ラッパーツリー  $T_w$  の  $W(T_w)$  は、 $T_w$  が文書ツリー  $T_d$  によって生成される時には  $W(T_d)$  であり、ラッパーツリー  $T_w$  がラッパーツリー

【0018】

【数1】

$$T_{w_1}$$

【0019】

および

【0020】

【数2】

$$T_{w_2}$$

【0021】

の組合せによって生成される時には

【0022】

【数3】

$$\max(W(T_{w_1}), W(T_{w_2}))$$

【0023】

である。

【0024】

ラッパー生成システムは、ラッパーツリーを生成するのに使用される文書ツリーの個数として「ラッパーレベル」を定義する。レベル1ラッパーツリーは、文書ツリーをラッパーツリーに変換することによって生成されるラッパーツリーである。文書ツリーのラッパーツリーへの変換は、次の式によって表される。

$$T_d \quad T_w$$

ここで、 $T_d$  は、文書  $d$  の文書ツリーを表し、 $T_w$  は、ラッパー  $w$  のラッパーツリーを表す。定義により、レベル1ラッパーの重さは、それがそこから生成された文書ツリーの重さである。文書ツリーをラッパーツリーに変換する時に、ラッパー生成システムは、反復パターン組合せアルゴリズムを実行して、 $T_w$  を  $T_d$  より簡潔にする。一実施形態で、ラッパー生成システムは、参照によって本明細書に組み込まれている非特許文献1に記載のアルゴリズムに類似するアルゴリズムを使用する。このアルゴリズムが、文書ノードの同一の連続するサブツリーを検出する場合に、このアルゴリズムは、それらをラッパーツリーの1つのラッパーノードにマージし、その記号に同一のサブツリーの個数をセットする。このアルゴリズムが、文書ノードの同一の連続するサブフォレストを識別する場合に、このアルゴリズムは、ラッパーツリー内のブロックノードの下の1つのサブフォレストとしてそれらをマージし、その記号に同一のサブフォレストの個数をセットする。このアルゴリズムは、文書ツリーをラッパーツリーに変換する時に、葉ノードのラベルを検討する。図2に、文書ツリーのラッパーツリーへの変換を示す。文書ツリー210が、ラッパーツリー220に変換される。文書ツリーのルートノードAは、BノードおよびCノードの反復サブフォレストを含むので、ラッパー生成システムは、括弧ノードXをラッパーツリーに追加し、その記号に2をセットする。各Bノードは、連続するノード  $D_1$  を含む（添字は、ノードのラベルを表す）ので、ラッパー生成システムは、これらのノードをラッパーツリー内で組み合わせ、その記号に2をセットする。

【0025】

ラッパー生成システムは、2つの低水準ラッパーツリーの位置合わせに基づいて高水準ラッパーツリーを生成する。ラッパー生成システムは、トップダウンの順序でレイヤごとに2つのラッパーツリー

【0026】

10

20

30

40

50

【数 4】

 $T_{w_1}$ 

【0 0 2 7】

および

【0 0 2 8】

【数 5】

 $T_{w_2}$ 

【0 0 2 9】

を位置合せする。ルートノードから同一の深さにあるノードは、同一レイヤに含まれ、ブロックノードは、1つのレイヤに含まれるとは考えられない。ラッパー生成システムは、同一レイヤ内のノードだけを位置合せする。ラッパー生成システムは、ラッパーノード

10

【0 0 3 0】

【数 6】

 $\sigma_{w_1}$ 

および

【0 0 3 1】

【数 7】

 $\sigma_{w_2}$ 

が次の条件を満足する時に、これらのノードが一致すると考える。

【0 0 3 2】

【数 8】

 $\sigma_{w_1}$ 

【0 0 3 3】

および

【0 0 3 4】

【数 9】

 $\sigma_{w_2}$ 

【0 0 3 5】

が、両方とも非葉ノードであるか、両方とも葉ノードであり、

【0 0 3 6】

【数 10】

$$T(\sigma_{w_1}) = T(\sigma_{w_2}),$$

【0 0 3 7】

であり、

【0 0 3 8】

【数 11】

 $\sigma_{w_1}$ 

【0 0 3 9】

および

【0 0 4 0】

20

30

40

【数 1 2】

$\sigma_{w_2}$

【0 0 4 1】

が両方とも葉ノードである場合に

【0 0 4 2】

【数 1 3】

$$L(\sigma_{w_1}) = L(\sigma_{w_2})$$

【0 0 4 3】

である。

ここで、 $T(\quad)$  は、ノード のタグを表し、 $L(\quad)$  は、ノード のラベルを表す。各レイヤで、ラッパー生成システムは、位置合わせ関数

【0 0 4 4】

【数 1 6】

$$A(F_{w_1}, F_{w_2})$$

【0 0 4 5】

を呼び出すことによって表されるように、サブツリーのアレイ

【0 0 4 6】

【数 1 4】

$F_{w_1}$

【0 0 4 7】

と

【0 0 4 8】

【数 1 5】

$F_{w_2}$

【0 0 4 9】

との間のシーケンス位置合わせを実行する。ラッパー生成システムは、最小コスト位置合わせを入手するために動的計画法を使用する。

【0 0 5 0】

【数 1 7】

$F_{w_1}$

【0 0 5 1】

および

【0 0 5 2】

【数 1 8】

$F_{w_2}$

【0 0 5 3】

内のすべての不一致のルートノードは、その重みをコストとして

【0 0 5 4】

【数 1 9】

$$A(F_{w_1}, F_{w_2})$$

【0 0 5 5】

に寄与する。非葉ノードである一致するノードの対

10

20

30

40

50

【 0 0 5 6 】

【 数 2 0 】

$\sigma_{w_2}$

【 0 0 5 7 】

および

【 0 0 5 8 】

【 数 2 1 】

$\sigma_{w_2}$

【 0 0 5 9 】

について、ラッパー生成システムは、

【 0 0 6 0 】

【 数 2 2 】

$A(\text{childF}_{w_1}, \text{childF}_{w_2})$

【 0 0 6 1 】

を再帰的に呼び出し、ここで、

【 0 0 6 2 】

【 数 2 3 】

$\text{childF}_{w_1}$

【 0 0 6 3 】

および

【 0 0 6 4 】

【 数 2 4 】

$\text{childF}_{w_2}$

【 0 0 6 5 】

は、

【 0 0 6 6 】

【 数 2 5 】

$\sigma_{w_2}$

【 0 0 6 7 】

および

【 0 0 6 8 】

【 数 2 6 】

$\sigma_{w_2}$

【 0 0 6 9 】

の子ノードをルートとするサブツリーからなるサブフォレストである。ラッパー生成システムは、関数

【 0 0 7 0 】

【 数 2 7 】

$A(\text{childF}_{w_1}, \text{childF}_{w_2})$

【 0 0 7 1 】

を呼び出すことによって計算されたコストを

【 0 0 7 2 】

10

20

30

40

【数 2 8】

$$A(F_{w_1}, F_{w_2})$$

【0073】

のコストに加算する。ラッパー生成システムは、ラッパーノードをトップダウンの再帰的な形で位置合せするので、2つのラッパー内のノードの両方がルートノードであるか、その親ノードが互いに位置合せされている場合に限って、それらのノードの位置合わせを試みる。

【0074】

図3は、一実施形態でのラッパーツリーのマージを示す図である。ラッパーツリー310および320がマージされて、ラッパーツリー330を形成する。影付きのノードDおよびGは、他方のラッパーツリー内に一致するノードを有しない。その結果、マージされたラッパーツリーの対応するラッパーノードは、ソフトノードである。ラッパー生成システムは、当初に、これらのラッパーツリーのルートノードを渡して、 $A(A, A)$ として位置合わせ関数を呼び出す。この関数は、第2レイヤのラッパーノードを渡して $A(B(C^3DE^*), BC^3E)$ としてそれ自体を再帰的に呼び出す。ワイルドカード?はソフトなので、この関数は、 $A(BC^3DE^*, BC^3E)$ および $A(B, BC^3E)$ としてそれ自体を再帰的に呼び出す。前者が、より低いコスト(すなわち、よりよい一致)をもたらすので、この関数は、前者を位置合わせとして選択する。関数 $A(B(C^3DE^*), BC^3E)$ の動的計画法プロセス中に、関数 $A(F^2, FG^+)$ が、 $A(BC^3DE^*, BC^3E)$ と $A(B, BC^3E)$ との両方によって再帰的に呼び出されて、この2つの解のコストを計算する。位置合わせが、2つのラッパーツリーの間の最適解を入手した後に、ラッパー生成システムは、次の記号生成関数Fを使用して新しいラッパーツリーを構築する。

$F(1, NULL)$	= ?		$F(? , N)$	= *
$F(? , NULL)$	= ?		$F(? , +)$	= *
$F(n, NULL)$	= *		$F(1, *)$	= *
$F(+, NULL)$	= *		$F(N, *)$	= *
$F(*, NULL)$	= *		$F(? , *)$	= *
$F(1, 1)$	= 1		$F(+, *)$	= *
$F(N, N)$	= N		$F(1, N)$	= +
$F(+, +)$	= +		$F(N, +)$	= +
$F(? , ?)$	= ?		$F(1, +)$	= +
$F(*, *)$	= *		$F(N_1, N_2)$	= +
$F(1, ?)$	= ?			

(1)

ここで、NULLは、ラッパーノードの不一致を表す。たとえば、 $F(1, NULL)$ は、その記号が1であるラッパーノードの不一致を示す。

【0075】

ラッパー生成システムは、次を除いて、ラッパーツリーが位置合せされる方法に似た形で文書ツリーおよびラッパーツリーを位置合せする。ラッパーツリーは、1対1の形で位置合せされるが、その記号が+または\*であるラッパーノード(タグノードのみ)を、複数の文書ノードと位置合わせすることができる。また、2つのラッパーノードを位置合せできるかどうかを判定する時には、ラッパー生成システムは、ノードのラベルを考慮に入れない。ラッパーツリー $T_w$ と文書ツリー $T_d$ との間の位置合わせについて、ラッパー生成システムは、不一致の文書ノードによって寄与される総コストを示すために $C_d(T_w, T_d)$ を使用する。 $C_d$ は、ラッパーノードと一致しない文書ツリー $T_d$ のノードのカウントを表し、 $C_w$ は、文書ノードと一致しないラッパーツリー $T_w$ のノードを表す。

【0076】

10

20

30

40

【数 29】

$$T_{w1} + T_{w2} \rightarrow T_{w3}$$

【0077】

の場合には、すべての文書ツリー  $T_d$  について、

【0078】

【数 30】

$$T_{w3}$$

【0079】

と  $T_d$  との間の位置合わせは、少なくとも

【0080】

【数 31】

$$T_{w1}$$

【0081】

と  $T_d$  との間ならびに

【0082】

【数 32】

$$T_{w2}$$

【0083】

と  $T_d$  との間と同数の一致する対を作る。というのは、

【0084】

【数 33】

$$T_{w3}$$

が、少なくとも

【0085】

【数 34】

$$T_{w1}$$

【0086】

および

【0087】

【数 35】

$$T_{w2}$$

【0088】

と同数のワイルドカードを有するからである。また、

【0089】

【数 36】

$$T_{w3}$$

【0090】

に現れるが

【0091】

【数 37】

$$T_{w1}$$

10

20

30

40

50

【 0 0 9 2 】

と

【 0 0 9 3 】

【 数 3 8 】

 $T_{w2}$ 

【 0 0 9 4 】

との両方には存在しないラッパーノードは、ソフトノードであり、したがって、コストには寄与しない。したがって、すべての文書ツリー  $T_d$  について、次の条件が満足される。

【 0 0 9 5 】

【 数 3 9 】

$$\begin{aligned} C_w(T_{w_3}, T_d) &\leq C_w(T_{w_1}, T_d) \\ C_d(T_{w_3}, T_d) &\leq C_d(T_{w_1}, T_d) \quad (i=1,2) \end{aligned} \quad (2)$$

【 0 0 9 6 】

ラッパー生成システムは、次の式によってラッパーツリー  $T_w$  と文書ツリー  $T_d$  との間の距離を定義する。

【 0 0 9 7 】

【 数 4 0 】

$$\Psi(T_w, T_d) = \left( \frac{C_w(T_w, T_d)}{W(T_w)} + \frac{C_d(T_w, T_d)}{W(T_w)} \right) / 2 \quad (3)$$

【 0 0 9 8 】

この式は、すべてのラッパーツリー  $T_w$  および文書ツリー  $T_d$  について、 $0 \leq \Psi(T_w, T_d) \leq 1$  であるという特性を有する。また、すべての文書ツリー  $T_d$  について、次の条件が満足される。

【 0 0 9 9 】

【 数 4 1 】

$$\Psi(T_{w_3}, T_d) \leq \Psi(T_{w_1}, T_d) \quad (i=1,2) \quad (4)$$

【 0 1 0 0 】

図 4 は、ラッパーツリーと文書ツリーとの位置合わせを示す図である。ラッパーノードおよび文書ノードは、破線によって示されているように位置合せされる。

【 0 1 0 1 】

一実施形態で、ラッパー生成システムは、文書ツリーがラッパーツリーとほぼ一致するかどうかを判定する時に、適応しきい値を使用する。ラッパー生成システムは、当初に、小さいしきい値から開始する。文書ツリーが現在のしきい値内にない時には、ラッパー生成システムは、小さい量だけしきい値を増やす。ラッパー生成システムは、増やされたしきい値がラッパーツリーによって表されるラッパーの性能を大きくは高めなくなるまで、クラスタ化を繰り返す。ラッパー生成システムは、対応するクラスタ化された文書ツリーをテストすることによって、増やされたしきい値の下で生成されたラッパーツリー  $T_w'$  の性能を評価する。ラッパー生成システムは、次に、前のラッパーツリーおよび現在のラッパーツリーの精度  $p$ 、リコール  $r$ 、および  $F1 = f$  を計算する。ラッパー生成システムは、次の式によって、大きい改善を表すことができる。

【 0 1 0 2 】

10

20

30

40

【数 4 2】

$$\begin{aligned} f_{r_c} &> f_{r_c} \\ r_c - r_c &> \lambda \end{aligned} \quad (5)$$

【0103】

ここで、 $\lambda$  は、しきい値増加が生成されるラッパーのリコールの大きい改善につながることを保証するのに使用される小さい値（たとえば、0.005）である。ラッパー生成システムは、初期しきい値と1との間になるように定義される停止値に達する時に、しきい値の増加を終了することもできる。

10

【0104】

図5は、一実施形態でのラッパー生成システムのコンポーネントを示すブロック図である。ラッパー生成システム530は、通信リンク520を介してウェブサイト510に接続される。ラッパー生成システムは、ウェブページ収集コンポーネント531および文書ツリーストア532を含む。ウェブページ収集コンポーネントは、さまざまなウェブサイトをクロールして、ラッパー生成用のトレーニングデータとしてウェブページを収集する。ウェブページ収集コンポーネントは、各ウェブページのドキュメントオブジェクトモデル（「DOM」）表現を生成し、その表現を文書ツリーとして文書ツリーストア内に格納する。ラッパー生成システムは、ユーザが文書ツリーの葉ノードに手作業でラベルを付けるためのユーザインターフェースをも提供することができる。ラッパー生成システムは、これらのラベルを文書ツリーストア内に格納する。したがって、ラッパー生成用のトレーニングデータは、葉ノードのラベルと一緒に文書ツリーを含む。ラッパー生成システムは、ラッパー生成コンポーネント533およびラッパーストア534をも含む。ラッパー生成コンポーネントは、類似するテンプレートを共有するウェブページをクラスタ化し、ツリーストアのトレーニングデータのラッパーツリーを生成し、そのラッパーツリーをラッパーストアに格納する。ラッパー生成コンポーネントは、類似するテンプレートを共有するウェブページを動的にクラスタ化し、クラスタ化中にラッパーツリーを動的に調整する。したがって、クラスタのウェブページは、それらがクラスタのラッパーツリーとどれほどよく一致するかに基づいて選択され、ラッパーツリーは、クラスタの選択されたウェブページに基づいて調整される。ラッパー生成コンポーネントは、ラッパーへの文書変換コンポーネント535、クラスタへの一致する文書の移動コンポーネント536、クラスタへのほぼ一致する文書の移動コンポーネント537、距離計算コンポーネント538、およびラッパー調整コンポーネント539を呼び出す。ラッパーへの文書変換コンポーネントは、トレーニングデータの文書ツリーをラッパーツリーに変換する。クラスタへの一致する文書の移動コンポーネントは、クラスタの現在のラッパーツリーと一致するトレーニングデータの文書ツリーを識別し、これらの識別された文書ツリーをクラスタに移動する。クラスタへのほぼ一致する文書の移動コンポーネントは、クラスタの現在のラッパーツリーとほぼ一致するトレーニングデータの文書ツリーを識別し、これらの識別された文書ツリーをクラスタに移動し、これらの識別された文書ツリーに基づいてクラスタのラッパーツリーを調整する。距離計算コンポーネントは、文書ツリーとラッパーツリーとの間の距離を計算して、文書ツリーがどれほどよく一致するかを判定する。ラッパー調整コンポーネントは、新しい文書ツリーがクラスタに追加される時のクラスタのラッパーツリーの動的調整をもたらすために、2つのラッパーツリーをマージする。ラッパー生成コンポーネントがウェブページをクラスタ化した後に、ラッパー生成コンポーネントは、各クラスタの動的に生成されたラッパーツリーをラッパーストアに格納する。ラッパー生成システムは、ラッパー選択コンポーネント540をも含む。ラッパー選択コンポーネントは、ウェブページを受け取り、そのウェブページの文書ツリーを生成し、そのツリーに最も近いラッパーストアのラッパーツリーを識別し、識別されたラッパーツリーのラッパーを使用してウェブページからデータを抽出する。

20

30

40

【0105】

50



ラッパー生成システムを実施できるコンピューティング装置は、中央処理装置、メモリ、入力装置（たとえば、キーボードおよびポインティング装置）、出力装置（たとえば、ディスプレイ装置）、およびストレージ装置（たとえば、ディスクドライブ）を含むことができる。物理メモリおよびストレージ装置は、ラッパー生成システムを実施する命令およびデータ構造を含むことができるコンピュータ可読媒体である。さらに、このデータ構造および命令を、格納するか、通信リンク上の信号などのデータ伝送媒体を介して伝送することができる。インターネット、ローカルエリアネットワーク、広域ネットワーク、ポイントツーポイントダイヤルアップネットワーク、セル電話網など、さまざまな通信リンクを使用して、システムのコンポーネントを接続することができる。

【0106】

10

ラッパー生成システムの実施形態は、パーソナルコンピュータ、サーバコンピュータ、マルチプロセッサシステム、マイクロプロセッサベースのシステム、ネットワークPC、ミニコンピュータ、メインフレームコンピュータ、上記のシステムまたは装置のいずれかを含む分散コンピューティング環境などを含むさまざまなオペレーティング環境で実施し、使用することができる。ユーザコンピューティング装置は、セル電話機、携帯情報端末、スマートホン、パーソナルコンピュータ、プログラマブル消費者エレクトロニクス、デジタルカメラなどを含むことができる。

【0107】

20

ラッパー生成システムは、1つまたは複数のコンピュータまたは他の装置によって実行される、プログラムモジュールなどのコンピュータ実行可能命令の全体的な文脈で説明することができる。一般に、プログラムモジュールは、特定のタスクを実行するか特定の抽象データ型を実施する、ルーチン、プログラム、オブジェクト、コンポーネント、データ構造などを含む。通常、プログラムモジュールの機能性は、さまざまな実施形態で望み通りに組み合わせるか分散させることができる。ラッパー生成システムによって使用される文書には、ウェブページ、XMLベースの文書、HTMLベースの文書など、すべての階層的に編成された文書を含めることができる。

【0108】

30

40

50

図6は、一実施形態でのラッパー生成システムのラッパー生成コンポーネントの処理を示す流れ図である。このコンポーネントは、類似するテンプレートを有するウェブページを動的にクラスタ化し、ウェブページの各クラスタのラッパーツリーを動的に生成する。このコンポーネントは、当初に、トレーニングコレクションD内にトレーニングデータのすべての文書ツリーを有する状態で開始する。ブロック601~609では、このコンポーネントは、ウェブページのクラスタを識別し、クラスタのラッパーツリーを動的に調整しながらループする。判断ブロック601で、トレーニングコレクションが空の場合には、このコンポーネントは完了し、そうでない場合には、このコンポーネントはブロック602で継続する。ブロック602では、このコンポーネントは、新しいクラスタTPを作成する。ブロック603で、このコンポーネントは、コレクションからの文書ツリー $T_d$ をクラスタの最初の文書ツリーとして選択する。一実施形態で、このコンポーネントは、文書ツリーをランダムに選択する。ブロック604で、このコンポーネントは、ラッパーへの文書変換コンポーネントを呼び出して、選択された文書ツリー $T_d$ を新しいクラスタTPの初期ラッパーツリー $T_w$ に変換する。ブロック605で、このコンポーネントは、コレクションDからの選択された文書ツリー $T_d$ を新しいクラスタTPに移動する。ブロック606で、このコンポーネントは、クラスタへの一致する文書の移動コンポーネントを呼び出して、初期ラッパーツリーと一致する文書ツリーをコレクションDから新しいクラスタTPに移動する。ブロック607で、このコンポーネントは、クラスタへのほぼ一致する文書の移動コンポーネントを呼び出して、調整されたラッパーツリーとほぼ一致する文書ツリーをコレクションDから新しいクラスタTPに移動し、移動された文書に基づいてラッパーツリーを調整する。ブロック608で、このコンポーネントは、クラスタTPをクラスタのコレクションRに追加する。ブロック609で、このコンポーネントは、ラッパーツリー $T_w$ をラッパーツリーのコレクションWに追加する。その後、このコンポ

ーメントは、ブロック601にループして、クラスタの識別を継続する。

【0109】

図7は、一実施形態でのラッパー生成システムのラッパーへの文書変換コンポーネントの処理を示す流れ図である。このコンポーネントは、文書ツリーのルートノードを渡され、ラッパーツリーを生成するためにこのコンポーネント自体を再帰的に呼び出す。このコンポーネントは、再帰的であるものとして図示されているが、当業者は、このコンポーネントを、その代わりに非再帰的な形で実施できることを了解するであろう。判断ブロック701で、渡されたノードが葉ノードである場合には、このコンポーネントはリターンし、そうでない場合には、このコンポーネントはブロック702で継続する。ブロック702~704では、このコンポーネントは、渡されたノードの各子ノードを選択し、このコンポーネントを再帰的に呼び出しながらループする。ブロック702で、このコンポーネントは、渡されたノードの次の子ノードを選択する。判断ブロック703で、すべての子ノードを既に選択し終えている場合に、このコンポーネントは、ブロック705で継続し、そうでない場合には、このコンポーネントは、ブロック704で継続する。ブロック704で、このコンポーネントは、選択されたノードを渡して、ラッパーへの文書変換コンポーネントを再帰的に呼び出し、その後、ブロック702にループして、次の子ノードを選択する。ブロック705で、このコンポーネントは、渡されたノードの連続するサブツリーを組み合わせる。ブロック706で、このコンポーネントは、渡されたノードの連続するサブフォレストを組み合わせる。その後、このコンポーネントは、リターンする。

10

【0110】

図8は、一実施形態でのラッパー生成システムのクラスタへの一致する文書の移動コンポーネントの処理を示す流れ図である。このコンポーネントは、ラッパーツリー $T_w$ と一致するコレクションDのすべての文書ツリー $T_d$ をクラスタTPに移動する。ブロック801で、このコンポーネントは、コレクションDの次の文書ツリー $T_d$ を選択する。判断ブロック802で、すべての文書ツリーが既に選択済みである場合には、このコンポーネントは、リターンし、そうでない場合には、このコンポーネントは、ブロック803で継続する。ブロック803では、このコンポーネントは、距離計算コンポーネントを呼び出して、選択された文書ツリー $T_d$ とラッパーツリー $T_w$ との間の距離を計算する。判断ブロック804で、距離が0である場合には、選択された文書ツリーは、ラッパーツリーと一致し、このコンポーネントは、ブロック805で継続し、そうでない場合には、このコンポーネントは、ブロック801にループして、次の文書ツリーを選択する。ブロック805で、このコンポーネントは、選択された文書ツリー $T_d$ をクラスタTPに移動し、その後、ブロック801にループして、次の文書ツリーを選択する。

20

30

【0111】

図9は、一実施形態でのラッパー生成システムのクラスタへのほぼ一致する文書の移動コンポーネントの処理を示す流れ図である。このコンポーネントは、ラッパーツリー $T_w$ とほぼ一致するコレクションDの文書ツリーをクラスタTPに移動し、ラッパーツリーを動的に調整する。このコンポーネントは、コレクション内のどの文書も、調整されたラッパーツリーとほぼ一致しなくなるまでこの処理を繰り返す。ブロック901~906では、このコンポーネントは、ほぼ一致する文書ツリーをクラスタに移動しながらループする。ブロック901で、このコンポーネントは、文書ツリー $T_d$ がラッパーツリー $T_w$ とほぼ一致するかどうかを判定するために、ほぼ一致する文書のチェックコンポーネントを呼び出す。判断ブロック902で、ほぼ一致する文書ツリーが見つかった場合には、このコンポーネントは、ブロック903で継続し、そうでない場合には、このコンポーネントは、ブロック907で継続する。ブロック903で、このコンポーネントは、ラッパーへの文書の変換コンポーネントを呼び出して、ほぼ一致する文書ツリーをラッパーツリー $T_w'$ に変換する。ブロック904で、このコンポーネントは、ラッパー調整コンポーネントを呼び出して、クラスタに追加される文書ツリーから生成されたラッパーツリー $T_w'$ に基づいて、ラッパーツリー $T_w$ を調整する。ブロック905で、このコンポーネントは、文書ツリー $T_d$ をコレクションDからクラスタTPに移動する。ブロック906で、この

40

50

コンポーネントは、クラスタへの一致する文書の移動コンポーネントを呼び出して、調整されたラッパーツリーと一致するコレクションのすべての文書ツリーをクラスタに移動する。その後、このコンポーネントは、ブロック 901 にループして、調整されたラッパーツリーとほぼ一致するさらなる文書ツリーをチェックする。

#### 【0112】

ブロック 907 ~ 911 で、このコンポーネントは、クラスタへの文書ツリーの移動に関する適応しきい値を実施する。このコンポーネントが固定しきい値を使用した場合には、このコンポーネントは、ブロック 907 で継続するのではなく、リターンする。ブロック 907 では、このコンポーネントは、現在のしきい値が、前のしきい値を使用して生成されたラッパーツリーより改善されたラッパーツリーをもたらしたかどうかを判定する。判断ブロック 908 で、改善が大きい場合には、このコンポーネントは、ブロック 909 でしきい値を増やし、ブロック 901 にループして、増やされたしきい値に基づいて文書ツリーをクラスタに移動する。しかし、改善が大きくはない場合には、このコンポーネントは、ブロック 910 で継続する。判断ブロック 910 では、改善がある場合に、このコンポーネントは、リターンし、そうでない場合には、このコンポーネントは、ブロック 911 で継続する。ブロック 911 で、このコンポーネントは、改善をもたらさなかった現在のしきい値に関する文書ツリー移動の影響をロールバックし、その後、リターンする。

#### 【0113】

図 10 は、一実施形態でのラッパー生成システムのほぼ一致する文書のチェックコンポーネントの処理を示す流れ図である。このコンポーネントは、ほぼ一致である文書ツリーが見つかるか、一致する文書ツリーがないと判定されるまで、コレクション D の文書ツリーとラッパーツリー  $T_w$  との間の距離をチェックしながらループする。ブロック 1001 で、このコンポーネントは、コレクション D の次の文書ツリー  $T_d$  を選択する。判断ブロック 1002 で、コレクションの文書ツリーのすべてが既に選択済みである場合には、ほぼ一致である文書ツリーはなく、このコンポーネントは、ほぼ一致がないことを示してリターンし、そうでない場合には、このコンポーネントは、ブロック 1003 で継続する。ブロック 1003 で、このコンポーネントは、ラッパーツリー  $T_w$  および文書ツリー  $T_d$  を渡して距離計算コンポーネントを呼び出して、文書ツリーとラッパーツリーとの間の距離を計算する。判断ブロック 1004 で、距離がしきい値未満である場合に、このコンポーネントは、ほぼ一致を示して文書ツリーを返し、そうでない場合には、このコンポーネントは、ブロック 1001 にループして、コレクションの次の文書ツリーを選択する。代替案では、このコンポーネントは、最も近いほぼ一致するが同一ではない文書ツリー、最も遠いほぼ一致する文書ツリー、またはランダムに選択されたほぼ一致する文書ツリーを返すことができる。

#### 【0114】

図 11 は、一実施形態でのラッパー生成システムの距離計算コンポーネントの処理を示す流れ図である。このコンポーネントは、ラッパーツリー  $T_w$  および文書ツリー  $T_d$  を渡され、そのラッパーツリーとその文書ツリーとの間の距離を計算する。ブロック 1101 で、このコンポーネントは、文書ツリーをラッパーツリーに位置合せする。ブロック 1102 で、このコンポーネントは、ラッパーツリーの位置合せされないノードの個数  $C_w$  をカウントする。ブロック 1103 で、このコンポーネントは、ラッパーツリーの重さ  $W(T_w)$  を計算する。ブロック 1104 で、このコンポーネントは、文書ツリーの位置合せされないノードの個数  $C_d$  をカウントする。ブロック 1105 で、このコンポーネントは、文書ツリーの重さ  $W(T_d)$  を計算する。ブロック 1106 で、このコンポーネントは、式 5 を使用して距離を計算する。その後、このコンポーネントは、リターンする。

#### 【0115】

図 12 は、一実施形態でのラッパー生成システムのラッパー調整コンポーネントの処理を示す流れ図である。このコンポーネントは、ラッパーツリーの対を渡され、これらを単一のラッパーツリーにマージする。ブロック 1201 で、このコンポーネントは、ラッパー位置合わせコンポーネントを呼び出して、ラッパーツリーのノードを位置合せする。ラ

10

20

30

40

50

ッパ-位置合わせコンポーネントは、ラッパ-ツリーのサブツリーのすべての可能な位置合わせをテストする再帰コンポーネントである。このコンポーネントは、動的計画法技法を使用して、以前にテストされた解の再テストを防ぐことができる。ブロック1202で、このコンポーネントは、ラッパ-ツリーのルートノードの次の解を選択する。判断ブロック1203で、すべての解が既に選択済みである場合には、このコンポーネントは、ブロック1205で継続し、そうでない場合には、このコンポーネントは、ブロック1204で継続する。ブロック1204で、このコンポーネントは、解のコストを集計し、その後、ブロック1202にループして、次の解を選択する。ブロック1205で、このコンポーネントは、最小のコストを有する解を選択し、その後、リターンする。

【0116】

図13は、一実施形態でのラッパ-生成システムのラッパ-位置合わせコンポーネントの処理を示す流れ図である。このコンポーネントは、すべての可能な解のコストを判定するために再帰的に呼び出される。判断ブロック1301で、両方のラッパ-ツリーのレイヤがまだある場合には、このコンポーネントは、ブロック1302で継続し、そうでない場合には、このコンポーネントは、リターンする。ブロック1302~1305では、このコンポーネントは、新しい解を選択し、ラッパ-位置合わせコンポーネントを再帰的に呼び出しながらループする。ブロック1302で、このコンポーネントは、現在のレイヤの次の解を選択する。判断ブロック1303で、すべての解が既に選択済みである場合には、このコンポーネントは、リターンし、そうでない場合には、このコンポーネントは、ブロック1304で継続する。ブロック1304で、このコンポーネントは、ラッパ-ツリーの次のレイヤのノードを位置合せするために、ラッパ-位置合わせコンポーネントを再帰的に呼び出す。ブロック1305で、このコンポーネントは、レイヤの選択された解のコストをセットし、その後、ブロック1302にループして、次の解を選択する。

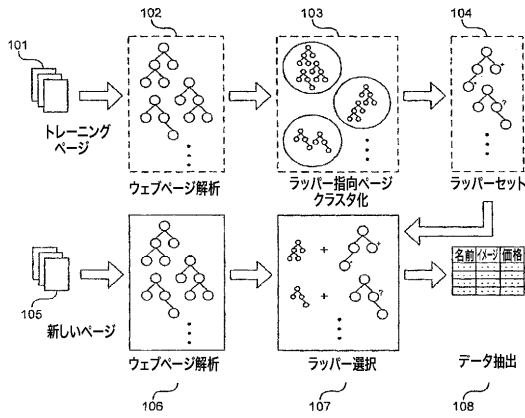
【0117】

本主題を、構造的特徴および/または方法論的動作に固有の言葉で説明してきたが、添付の特許請求の範囲で定義される本主題が、必ずしも上で説明した特定の特徴または動作に限定されないことを理解されたい。そうではなく、上で説明した特定の特徴および動作は、特許請求の範囲を実施する例の形として開示されたものである。したがって、本発明は、添付の特許請求の範囲によるものを除いて限定されない。

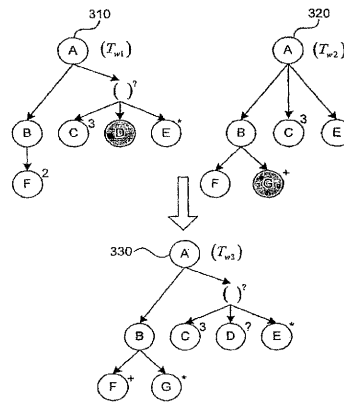
10

20

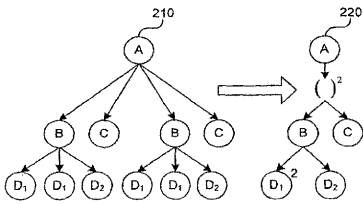
【 図 1 】



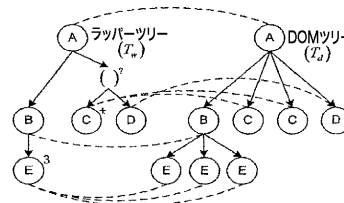
【 図 3 】



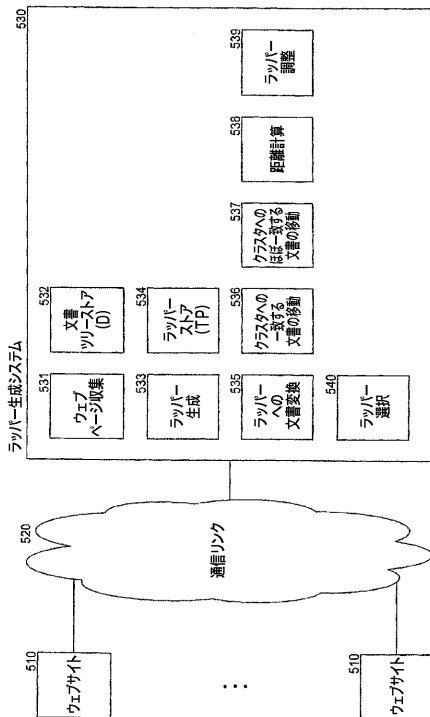
【 図 2 】



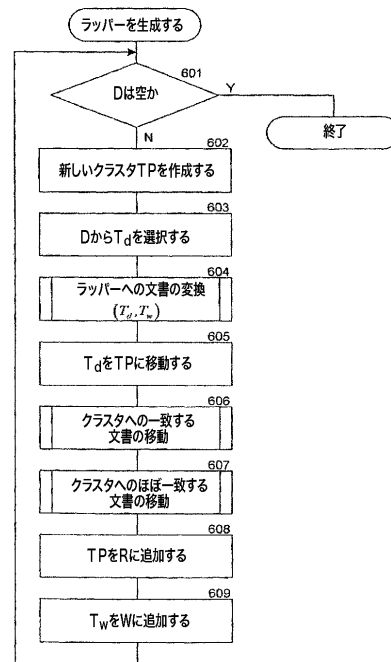
【 図 4 】



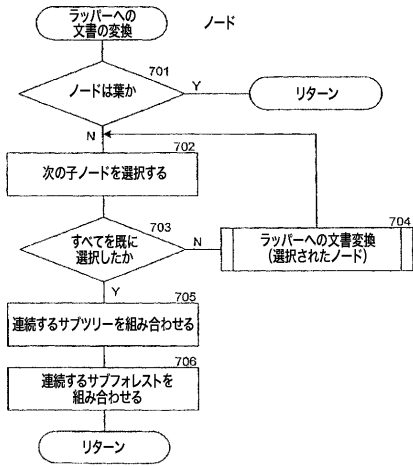
【 図 5 】



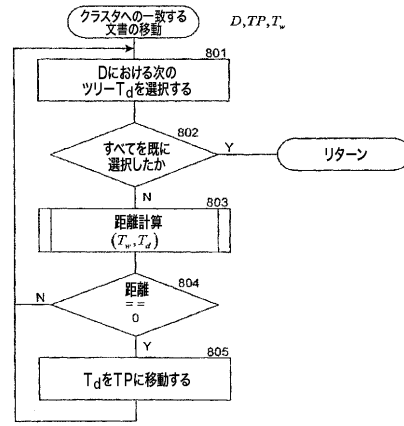
【 図 6 】



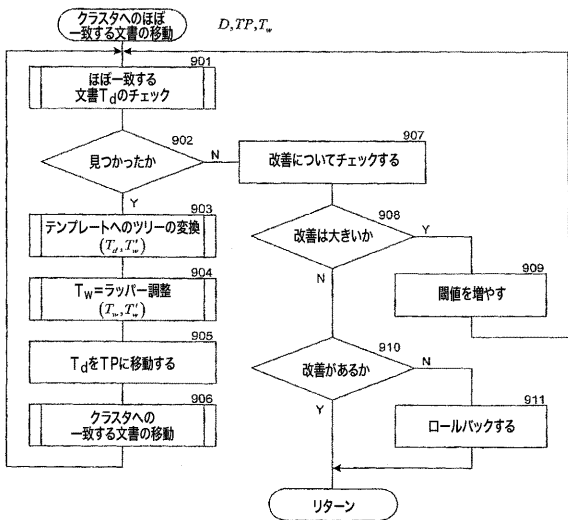
【 図 7 】



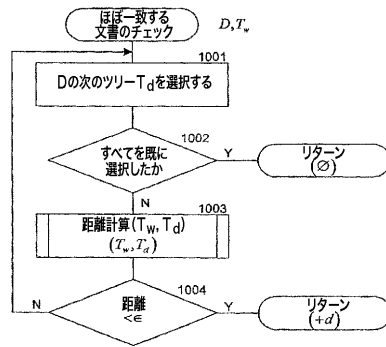
【 図 8 】



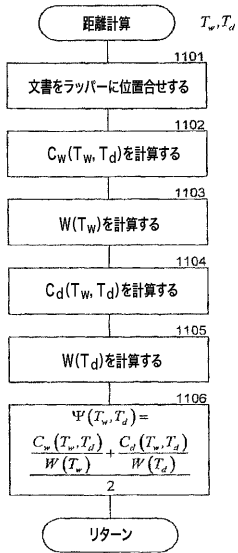
【 図 9 】



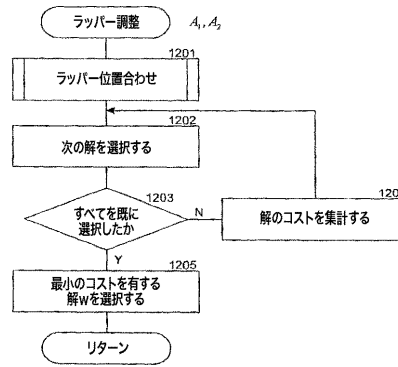
【 図 10 】



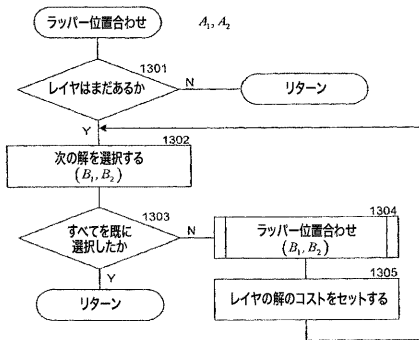
【 図 1 1 】





【 図 1 2 】



【 図 1 3 】



## 【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. <b>PCT/US2007/018417</b>
<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
<i>G06F 17/00(2006.01)i, G06F 17/24(2006.01)i</i>		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) IPC8 G06F 17/00		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Korean Utility models and applications for Utility models since 1975 Japanese Utility models and applications for Utility models since 1975		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) ekipass "wrapper, tree, web-page, cluster, parsing, query, extract, hierarchy, document, node, metadata"		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y A	WO 2005/072072 A2 (KIM, SUN KWON) 11 AUG. 2005 See abstract; figures 2,6-8,11. claims 1-6,24,29,35.	1-11,16-20 12-15
Y A	WO 2006/036376 A1 (GOOGLE INC.) 06 APR. 2006 See abstract; figures 7-9b. claims 1,5-8,18,19,24,26,27,30-32,42.	1-11,16-20 12-15
A	US 2004/0068697 A1 (GEORGES HARIK et al.) 08 APR. 2004 See abstract; figures 22-25. claims 1,21,41,61.	1-20
A	KR 1020030069639 A (LEE, UI BUM) 27 AUG. 2003 See abstract; figure 3. claims 1,2.	1-20
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 01 FEBRUARY 2008 (01.02.2008)		Date of mailing of the international search report <b>01 FEBRUARY 2008 (01.02.2008)</b>
Name and mailing address of the ISA/KR  Korean Intellectual Property Office 920 Dunsan-dong, Seo-gu, Daejeon 302-701, Republic of Korea Facsimile No. 82-42-472-7140		Authorized officer KIM, Jung Jin Telephone No. 82-42-481-5962 



**INTERNATIONAL SEARCH REPORT**  
Information on patent family members

International application No.

**PCT/US2007/018417**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
W02005072072A2	11.08.2005	CN1906615A KR2005077681A US2007110047A1 W02005072072A3	31.01.2007 03.08.2005 17.06.2007 06.10.2005
W02006036376A1	06.04.2006	AU2005290154A1 CA2581713A1 EP01800226A1 KR1020070058685A US2006074907AA	06.04.2006 06.04.2006 27.06.2007 08.06.2007 06.04.2006
US20040068697A1	08.04.2004	US2007208772AA US7231393BA	06.09.2007 12.06.2007
KR1020030069639A	27.08.2003	None	

## フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, MT, NL, PL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, SV, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW

(72)発明者 ミン ワン

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マイクロソフト コーポレーション インターナショナル パテント内

(72)発明者 ルイホワ ソン

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マイクロソフト コーポレーション インターナショナル パテント内

(72)発明者 ウェイ - イン マ

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マイクロソフト コーポレーション インターナショナル パテント内

(72)発明者 シュイ ゼン

アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マイクロソフト コーポレーション インターナショナル パテント内

Fターム(参考) 5B075 ND35 NK43 NR06