

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-205060
(P2010-205060A)

(43) 公開日 平成22年9月16日(2010.9.16)

(51) Int.Cl. F I テーマコード(参考)
G06F 17/30 (2006.01) G06F 17/30 210A 5B075
 G06F 17/30 170B

審査請求 未請求 請求項の数 10 O L (全 16 頁)

(21) 出願番号 特願2009-50950 (P2009-50950)
 (22) 出願日 平成21年3月4日(2009.3.4)

(71) 出願人 000155469
 株式会社野村総合研究所
 東京都千代田区丸の内一丁目6番5号
 (74) 代理人 100080001
 弁理士 筒井 大和
 (74) 代理人 100093023
 弁理士 小塚 善高
 (74) 代理人 100117008
 弁理士 筒井 章子
 (72) 発明者 竹原 一彰
 東京都千代田区丸の内一丁目6番5号 株
 式会社野村総合研究所内
 Fターム(参考) 5B075 ND08 NK02 NK31 PQ74

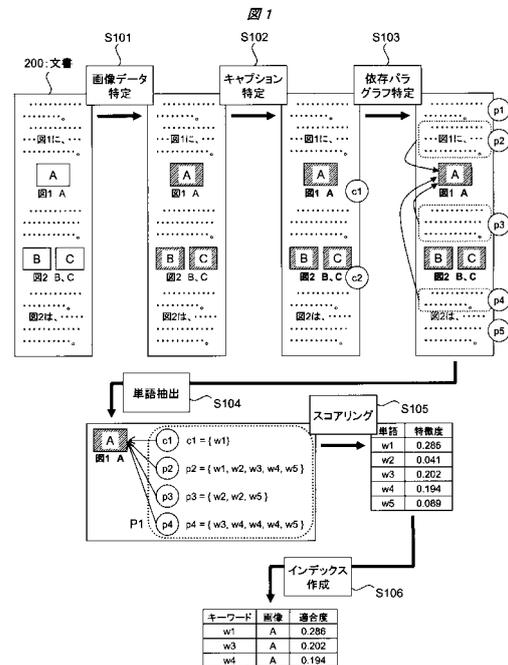
(54) 【発明の名称】 文書内画像検索方法および文書内画像検索システム

(57) 【要約】

【課題】 文書から抽出した画像に対して、キーワードを文脈情報を考慮して効率的に精度良く付与することにより、キーワードを利用した高精度で効率の良い画像の検索を可能とする文書内画像検索方法を提供する。

【解決手段】 文書の中から画像の位置を特定し、画像のデータを抽出する処理(S101)と、画像のキャプション領域を特定し、画像名とキャプションとを抽出する処理(S102)と、画像について記述している依存パラグラフを特定する処理(S103)と、依存パラグラフから単語を抽出する処理(S104)と、各単語について依存パラグラフ内での特徴度をスコアリングする処理(S105)と、特徴度が上位の単語をキーワードとして抽出し、インデックステーブルに格納する処理(S106)とを実行し、指定された検索語に基づいてインデックステーブル内のキーワードを検索し、一致するキーワードが付与された画像を出力する処理を実行する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

コンピュータシステムにより、画像と文字列が混在する文書から抽出した画像に対してキーワードを付与し、ユーザから指定された検索語に基づいて前記キーワードを検索し、一致する前記キーワードが付与された画像を出力する文書内画像検索方法であって、

前記コンピュータシステムは、

前記文書を解析して前記文書の中から画像の位置を特定し、前記画像のデータを抽出して格納する画像データ特定処理と、

前記画像データ特定処理で特定した前記画像について、前記文書を解析して前記画像のキャプション領域を特定し、前記キャプション領域から画像名とキャプションとを抽出して前記画像と対応付けて格納するキャプション特定処理と、

前記画像データ特定処理で特定した前記画像について、前記文書を解析して前記文書中で前記画像について記述しているパラグラフである依存パラグラフを特定する依存パラグラフ特定処理と、

前記依存パラグラフ特定処理で特定した前記依存パラグラフから単語を抽出する単語抽出処理と、

前記単語抽出処理で抽出した前記各単語について、前記依存パラグラフ内での特徴度を所定の方法によりスコアリングするスコアリング処理と、

前記スコアリング処理で算出した前記各単語の特徴度が上位の所定の前記単語を前記キーワードとして抽出し、抽出した前記キーワードを対象の前記画像のインデックスとし、その前記特徴度を対象の前記画像に対する適合度として、インデックステーブルに格納するインデックス作成処理とを実行し、

前記ユーザによって指定された前記検索語に基づいて、前記インデックステーブル内の前記キーワードを検索し、一致する前記キーワードが付与された前記画像を出力する画像検索処理を実行することを特徴とする文書内画像検索方法。

【請求項 2】

請求項 1 に記載の文書内画像検索方法において、

前記依存パラグラフ特定処理では、

前記画像データ特定処理で特定した前記画像の前記画像名によって前記文書をサーチし、前記画像名の文字列が最初に出現したパラグラフから、次の画像の画像名の文字列が出現するパラグラフの直前のパラグラフまでを、対象の前記画像についての前記依存パラグラフとして特定することを特徴とする文書内画像検索方法。

【請求項 3】

請求項 1 または 2 に記載の文書内画像検索方法において、

前記スコアリング処理では、

前記単語抽出処理で抽出した前記各単語について、対象の前記画像の前記依存パラグラフ内での前記単語の出現頻度と、前記文書内の全ての画像の前記依存パラグラフの中での前記単語が出現する前記依存パラグラフの数とに基づいて、前記特徴度をスコアリングすることを特徴とする文書内画像検索方法。

【請求項 4】

請求項 3 に記載の文書内画像検索方法において、

前記スコアリング処理では、

前記単語抽出処理で抽出した前記各単語について、対象の前記画像の前記依存パラグラフ内での前記単語の出現頻度を、対象の前記画像の前記依存パラグラフ内での前記単語の出現位置に基づいて重み付けして算出することを特徴とする文書内画像検索方法。

【請求項 5】

請求項 1 ~ 4 のいずれか 1 項に記載の文書内画像検索方法において、

前記画像検索処理では、

前記検索語に基づいて前記画像を出力する際に、前記インデックステーブルから、前記検索語に一致する前記キーワードと対応する前記画像との前記適合度を取得し、前記適合

10

20

30

40

50

度に応じて前記画像の出力方法を制御することを特徴とする文書内画像検索方法。

【請求項 6】

画像検索サーバおよび前記画像検索サーバに接続されたクライアント端末を有し、画像と文字列が混在する文書から抽出した画像に対してキーワードを付与し、ユーザから指定された検索語に基づいて前記キーワードを検索し、一致する前記キーワードが付与された画像を出力する文書内画像検索システムであって、

前記画像検索サーバは、

前記文書の中から画像の位置を特定し、前記画像のデータを抽出して格納する画像データ特定部と、

前記画像データ特定部で特定された前記画像について、前記画像のキャプション領域を特定し、前記キャプション領域から画像名とキャプションとを抽出して前記画像と対応付けて格納するキャプション特定部と、

前記画像データ特定部で特定された前記画像について、前記文書中で前記画像について記述しているパラグラフである依存パラグラフを特定する依存パラグラフ特定部と、

前記依存パラグラフ特定部で特定された前記依存パラグラフから単語を抽出する単語抽出部と、

前記単語抽出部で抽出された前記各単語について、前記依存パラグラフ内での特徴度を所定の方法によりスコアリングするスコアリング部と、

前記スコアリング部で算出された前記各単語の特徴度が上位の所定の前記単語を前記キーワードとして抽出し、抽出した前記キーワードを対象の前記画像のインデックスとし、その前記特徴度を対象の前記画像に対する適合度として、インデックステーブルに格納するインデックス作成部と、

前記クライアント端末を利用して前記ユーザによって指定された前記検索語に対して、前記インデックステーブルから前記検索語と一致する前記キーワードに対応する前記画像を取得する検索処理部と、

前記クライアント端末上に表示させる、前記検索語の入力や検索結果の出力のための画面を生成するユーザインタフェース部とを有することを特徴とする文書内画像検索システム。

【請求項 7】

請求項 6 に記載の文書内画像検索システムにおいて、

前記依存パラグラフ特定部は、

前記画像データ特定部で特定された前記画像の前記画像名によって前記文書をサーチし、前記画像名の文字列が最初に出現したパラグラフから、次の画像の画像名の文字列が出現するパラグラフの直前のパラグラフまでを、対象の前記画像についての前記依存パラグラフとして特定することを特徴とする文書内画像検索システム。

【請求項 8】

請求項 6 または 7 に記載の文書内画像検索システムにおいて、

前記スコアリング部は、

前記単語抽出部で抽出された前記各単語について、対象の前記画像の前記依存パラグラフ内での前記単語の出現頻度と、前記文書内の全ての画像の前記依存パラグラフの中での前記単語が出現する前記依存パラグラフの数とに基づいて、前記特徴度をスコアリングすることを特徴とする文書内画像検索システム。

【請求項 9】

請求項 8 に記載の文書内画像検索システムにおいて、

前記スコアリング部は、

前記単語抽出部で抽出された前記各単語について、対象の前記画像の前記依存パラグラフ内での前記単語の出現頻度を、対象の前記画像の前記依存パラグラフ内での前記単語の出現位置に基づいて重み付けして算出することを特徴とする文書内画像検索システム。

【請求項 10】

請求項 6 ~ 9 のいずれか 1 項に記載の文書内画像検索システムにおいて、

10

20

30

40

50

前記検索処理部は、

前記ユーザによって指定された前記検索語に対して、前記インデックステーブルから前記検索語と一致する前記キーワードに対応する前記画像を前記適合度と合わせて取得し、

前記ユーザインタフェース部は、

前記検索結果の出力のための画面を生成する際に、前記適合度に応じて検索結果の前記画像の出力方法を制御することを特徴とする文書内画像検索システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像と文字列が混在する文書内の画像を検索する技術に関し、特に、画像の内容を表すキーワードを指定することにより画像を検索する文書内画像検索方法および文書内画像検索システムに適用して有効な技術に関するものである。

10

【背景技術】

【0002】

近年、IT技術の進展により、従来は紙などの物理的な媒体によって保存されていた文書等を含む大量の情報が電子化されて保存されるようになってきている。さらに、これらの情報に対してコンピュータを利用して、例えば検索エンジンや検索システム等によって検索して所望の情報を取得し、情報を有効活用するということが行われている。

【0003】

電子化された文書群に対して、検索語を指定し、文書内のテキスト（もしくは文書の内容を表すキーワードやタグ等）に検索語と一致する文字列を含む文書や、その文書内における位置などを検索することは広く一般的に行われている。一方、電子化された文書には、テキスト情報以外に図や表、写真などの画像も含まれる。この文書内に含まれる画像についてもテキストと同様に検索語を指定することにより検索したいという要望がある。この場合、画像データ自体は文字情報を含まないため、画像に対してその内容を表す文字情報を何らかの手段で付与する必要がある。

20

【0004】

これに対して、例えば、特開平8-202731号公報（特許文献1）には、スキャナにより入力した文書を、画像分離手段により文字領域と画像領域とに分離し、文字領域から文字認識手段により文字列を認識し、認識された文字列から単語分離手段により単語を抽出し、抽出した単語の文書内での出現頻度に基づいて入力画像に付加するキーワードを判定し、キーワードと入力画像をデータベースに登録することにより、オペレータの手を介さずに入力画像にキーワードを付加する技術が開示されている。

30

【0005】

また、例えば、特開平11-25113号公報（特許文献2）には、画像および文字列が混在した文書から画像を抽出して画像DBに格納する際に、文書中から画像について記述した文字列（キーテキスト）を自動的に抽出して画像に関連付けて格納し、入力された検索語に基づいてキーテキストを検索することによって該当する画像を得る技術が開示されている。

40

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開平8-202731号公報

【特許文献2】特開平11-25113号公報

【発明の概要】

【発明が解決しようとする課題】

【0007】

特許文献1に記載されたような画像に対するキーワードの付与方法では、キーワードの抽出に際して、文書中での対象の画像に対する言及などの文脈情報を考慮したものとなっていないため、特に、文書中の画像が複数になった場合には、画像の内容を表した適切な

50

キーワードを付与することができず、画像検索の際の適合率が低くなるという問題が生じる。

【0008】

一方、特許文献2に記載されたようなキーテキストの付与方法では、画像について記述した文字列をキーテキストとするため、文脈情報を考慮したキーテキストを画像に付与することができる。また、キーキャプションを使用すればノイズが少ない検索が可能であり、また、キーページを使用すれば広範囲の検索を行うことが可能である。

【0009】

しかし、逆に、キーキャプションを使用すれば漏れが大きくなり、また、キーページを使用すればノイズが大きくなるため適合率が低くなるという問題が生じる。さらに、画像と関連付けて格納するキーテキストの情報（特にキーページ）がキーワードの場合と比べて格段に大きくなるという問題や、検索時に検索対象のキーテキストの種別を多く指定するほど検索処理に時間を要するという問題を生じる。

10

【0010】

そこで本発明の目的は、文書から抽出した画像に対して、画像の内容を表すキーワードを文脈情報を考慮して効率的に精度良く付与することにより、キーワードを利用した高精度で効率の良い画像の検索を可能とする文書内画像検索方法および文書内画像検索システムを提供することにある。本発明の前記ならびにその他の目的と新規な特徴は、本明細書の記述および添付図面から明らかになるであろう。

【課題を解決するための手段】

20

【0011】

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

【0012】

本発明の代表的な実施の形態による文書内画像検索方法は、文書を解析して前記文書の中から画像の位置を特定し、前記画像のデータを抽出して格納する画像データ特定処理と、前記画像データ特定処理で特定した前記画像について、前記文書を解析して前記画像のキャプション領域を特定し、前記キャプション領域から画像名とキャプションとを抽出して前記画像と対応付けて格納するキャプション特定処理と、前記画像データ特定処理で特定した前記画像について、前記文書を解析して前記文書中で前記画像について記述しているパラグラフである依存パラグラフを特定する依存パラグラフ特定処理と、前記依存パラグラフ特定処理で特定した前記依存パラグラフから単語を抽出する単語抽出処理と、前記単語抽出処理で抽出した前記各単語について、前記依存パラグラフ内での特徴度を所定の方法によりスコアリングするスコアリング処理と、前記スコアリング処理で算出した前記各単語の特徴度が上位の所定の前記単語を前記キーワードとして抽出し、抽出した前記キーワードを対象の前記画像のインデックスとし、その前記特徴度を対象の前記画像に対する適合度として、インデックステーブルに格納するインデックス作成処理とを実行し、ユーザによって指定された前記検索語に基づいて、前記インデックステーブル内の前記キーワードを検索し、一致する前記キーワードが付与された前記画像を出力する画像検索処理を実行することを特徴とするものである。

30

40

【発明の効果】

【0013】

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

【0014】

本発明の代表的な実施の形態によれば、文書から抽出した画像に対して、画像に関連するキーワードを文脈情報を考慮して効率的に精度良く付与することが可能となり、キーワードを利用した高精度で効率の良い画像の検索が可能となる。

【図面の簡単な説明】

【0015】

50

【図 1】本発明の実施の形態 1 における、画像抽出部でのインデックス作成処理の例について説明する図である。

【図 2】本発明の実施の形態 1 である文書内画像検索システムの構成例の概要を示す図である。

【図 3】本発明の実施の形態 1 における、画像情報およびインデックステーブルのデータ構成の例を示した図である。

【図 4】本発明の実施の形態 1 における、依存パラグラフに含まれる各単語の T F × I P F 値を算出した例を示した図である。

【図 5】本発明の実施の形態 1 における、画像を検索する際にクライアント端末に表示されるユーザインタフェースの例を示した図である。

【図 6】本発明の実施の形態 2 における、依存パラグラフに含まれる各単語の T F × I P F 値を、単語の出現位置に応じて重み付けして算出した例を示した図である。

【発明を実施するための形態】

【0016】

以下、本発明の実施の形態を図面に基づいて詳細に説明する。なお、実施の形態を説明するための全図において、同一部には原則として同一の符号を付し、その繰り返しの説明は省略する。

【0017】

< 実施の形態 1 >

本発明の実施の形態 1 である文書内画像検索システムは、画像と文字列が混在した文書から画像を抽出し、文書中で当該画像について記述しているパラグラフ（依存パラグラフ）内の単語から、特徴度のスコアリングにより上位のものをキーワードとして抽出し、当該キーワードを当該画像に関連するキーワードとして付与してインデックスを作成する。このとき、当該キーワードの特徴度を当該画像に対する適合度とする。また、ユーザにより画像を検索するための検索語が入力されると、インデックスに基づいて検索語と一致するキーワードに対応する画像を取得して画面表示により出力する。このとき、対応する画像が複数ある場合は、適合度に応じて優先付けして画面表示する。

【0018】

[システム構成]

図 2 は、本発明の実施の形態 1 である文書内画像検索システムの構成例の概要を示す図である。文書内画像検索システム 1 は、例えば、コンピュータシステムによる画像検索サーバ 100 とデータベース、および、インターネットや社内 LAN 等のネットワーク 500 を介して画像検索サーバ 100 に接続された、PC 等のクライアント端末 400 から構成される。また、画像検索サーバ 100 は、データベースもしくはファイル等により、画像とテキストが混在する複数の文書 200 を保持している。

【0019】

画像検索サーバ 100 は、例えば、画像抽出部 110 および画像検索部 120 を有する。また、データベースとして、画像情報 310 およびインデックステーブル 320 を有する。これらのデータベースは、画像検索サーバ 100 が直接保持してもよいし、アクセス可能な他のデータベースサーバに保持する構成としてもよい。画像抽出部 110 は、文書 200 内の画像を抽出し、キーワードを付与してインデックスを作成する処理を行い、例えば、画像データ特定部 111、キャプション特定部 112、依存パラグラフ特定部 113、単語抽出部 114、スコアリング部 115、およびインデックス作成部 116 を有する。

【0020】

画像データ特定部 111 は、各文書 200 を解析して文書 200 中の画像の位置を特定し、当該画像のデータを抽出して、画像情報 310 に格納する。キャプション特定部 112 は、画像データ特定部 111 によって特定された各画像について、文書 200 を解析して対象の画像のキャプション領域を特定し、当該キャプション領域から画像名とキャプションとを抽出して、画像情報 310 の対象の画像のエントリに格納する。依存パラグラフ

10

20

30

40

50

特定部 1 1 3 は、画像データ特定部 1 1 1 によって特定された各画像について、文書 2 0 0 を解析して後述する依存パラグラフを特定し、画像情報 3 1 0 の対象の画像のエントリに格納する。

【 0 0 2 1 】

単語抽出部 1 1 4 は、依存パラグラフ特定部 1 1 3 によって特定された依存パラグラフから自然言語処理により単語（複合名詞）を抽出する。スコアリング部 1 1 5 は、単語抽出部 1 1 4 によって抽出された各単語について、依存パラグラフ内での特徴度を後述する方法によりスコアリングする。インデックス作成部 1 1 6 は、スコアリング部 1 1 5 によって算出された各単語の特徴度が上位の所定の単語をキーワードとして抽出し、抽出したキーワードを対象の画像のインデックスとし、その特徴度を対象の画像に対する適合度として、インデックステーブル 3 2 0 に格納する。

10

【 0 0 2 2 】

画像検索部 1 2 0 は、クライアント端末 4 0 0 を利用してユーザによって指定された検索語に基づいて、インデックステーブル 3 2 0 内のキーワードを検索し、一致するキーワードが付与された画像を出力する画像検索処理を行い、例えば、検索処理部 1 2 1 およびユーザインタフェース部 1 2 2 を有する。検索処理部 1 2 0 は、ユーザによって指定された検索語に対して、インデックステーブル 3 2 0 から検索語と一致するキーワードに対応する画像を適合度と合わせて取得する。ユーザインタフェース部 1 2 2 は、クライアント端末 4 0 0 上に表示させる、検索語の入力や検索結果の出力のための画面を生成する。

【 0 0 2 3 】

画像抽出部 1 1 0 および画像検索部 1 2 0 の各部は、ソフトウェアプログラムとして実現され、例えば、図示しない Web サーバプログラム上で稼働するアプリケーションプログラムとして実現される。また、ユーザインタフェース部 1 2 2 では、例えば、HTML (HyperText Markup Language) によって画面を生成し、図示しない Web サーバプログラムを介して、クライアント端末 4 0 0 上の図示しない Web ブラウザによって表示させる。

20

【 0 0 2 4 】

文書 2 0 0 は、例えば、ワードプロセッサ等のアプリケーションプログラムで作成された画像を含むテキスト文書や、HTML 等のタグ文書など、画像抽出部 1 1 0 により画像と文字列の認識が可能である電子化された文書であれば取り扱うことが可能である。なお、紙媒体の文書であっても、例えば、特許文献 1、2 等に記載されているように、スキャナによって紙媒体の文書を読み取り、読み取ったデータに基づいて文字領域と画像領域とを識別し、文字領域については OCR (Optical Character Reader) 等により文字認識を行うことによって文書 2 0 0 として取り込むことが可能である。

30

【 0 0 2 5 】

[データ構成]

図 3 は、画像情報 3 1 0 およびインデックステーブル 3 2 0 のデータ構成の例を示した図である。画像情報 3 1 0 は、例えば、画像 ID 3 1 1、画像データ 3 1 2、文書名 3 1 3、位置 3 1 4、画像名 3 1 5、キャプション 3 1 6、および依存パラグラフ 3 1 7 の項目を有し、文書 2 0 0 から抽出された画像に関する情報を保持する。

40

【 0 0 2 6 】

画像 ID 3 1 1 は、文書内画像検索システム 1 内で対象の画像を一意に特定するために付与される ID である。画像データ 3 1 2 は、文書 2 0 0 から抽出された画像のバイナリデータである。文書名 3 1 3 および位置 3 1 4 は、対象の画像が含まれる文書 2 0 0 の文書名および文書 2 0 0 内の位置（行数）である。画像名 3 1 5 およびキャプション 3 1 6 は、対象の画像の画像名およびキャプションである。依存パラグラフ 3 1 7 は、対象の画像の依存パラグラフの文字列である。

【 0 0 2 7 】

インデックステーブル 3 2 0 は、例えば、キーワード 3 2 1、画像 ID 3 2 2、および適合度 3 2 3 の項目を有し、ユーザから指定された検索語によって画像検索部 1 2 0 にお

50

いて画像を検索する際に利用するインデックスを保持する。キーワード321は、文書200から抽出した各画像に対して画像抽出部110での処理によって付与されたキーワードである。画像ID322は、対象のキーワードが付与された画像のIDである。適合度323は、対象のキーワードの対象の画像に対する適合度を示すスコアである。なお、インデックステーブル320は、データベースに限らずファイル形式であってもよい。また、画像情報310およびインデックステーブル320の各項目は上記のものに限らず、他の項目を有していてもよい。

【0028】

[インデックス作成処理]

図1は、本実施の形態の画像抽出部110でのインデックス作成処理の例について説明する図である。まず、画像データ特定部111により、対象の文書200を解析して文書200の中から画像の位置を特定し、当該画像データを抽出して画像情報310に格納する、画像データ特定処理を行う(ステップS101)。画像の位置の特定については、例えば、特許文献1や特許文献2に記載されているような方法をとることができる。図1では、画像A、B、Cの3つの画像を特定した場合の例を示している。

10

【0029】

なお、抽出した画像にはIDを付与し、画像のバイナリデータと合わせて、画像情報310の画像ID311および画像データ312にそれぞれ格納する。また、当該画像が含まれる文書200の文書名および文書200内の位置(行数)を、画像情報310の文書名313および位置314にそれぞれ格納する。

20

【0030】

次に、キャプション特定部112により、ステップS101で特定した画像について画像のキャプション領域を特定し、キャプション領域から画像名とキャプションとを抽出して画像情報310の該当の画像のエントリに格納する、キャプション特定処理を行う(ステップS102)。ここで、キャプション領域とは、図や表などの画像についての短い説明が記載された領域であり、例えば、「図1」や「表2」などの画像名と、「インデックス作成処理の例について説明する図」などの画像に対して付与された文字列であるキャプションから構成される。

【0031】

キャプション領域の特定については、例えば、特許文献2に記載されているような方法をとることができる。ここで、例えば、学术论文などの文書では、一般的に画像が図である場合にはキャプション領域は画像の下部に配置され、画像が表である場合には画像の上部に配置される。従って、画像の上部および下部の所定の小領域をキャプション領域として特定する。キャプション領域から画像名とキャプションを特定する際には、例えば、キャプション領域内の文字列から画像名に相当する文字列を判定するための正規表現を用いて画像名を特定し、その後続く1文をキャプションとして特定する方法をとることができる。

30

【0032】

画像名に相当する文字列を判定するための正規表現としては、例えば、「図¥d*」、「表¥d*」、「図表¥d*」、「グラフ¥d*」などを用いることができる。これらの正規表現は、予め定義してファイル等に保持しておく。図1では、画像Aについては「図1 A」というキャプション領域(画像名「図1」、キャプションc1「A」)、画像B、Cについては「図2 B、C」というキャプション領域(画像名「図2」、キャプションc2「B、C」)を特定した場合の例を示している。なお、抽出した画像名およびキャプションは、画像情報310の該当の画像のエントリの画像名315およびキャプション316にそれぞれ格納する。

40

【0033】

次に、依存パラグラフ特定部113により、ステップS101で特定した画像について、文書200中で当該画像について記述しているパラグラフである依存パラグラフを特定する、依存パラグラフ特定処理を行う(ステップS103)。依存パラグラフの特定につ

50

いては、例えば、ステップ S 1 0 2 で特定した画像の画像名によって文書 2 0 0 をサーチし、画像名の文字列が最初に出現したパラグラフから、次の画像の画像名の文字列が出現するパラグラフの直前のパラグラフまでを、対象の画像についての依存パラグラフとして特定する。

【 0 0 3 4 】

図 1 では、パラグラフ p 1 ~ p 5 のうち、画像 A についての依存パラグラフ P 1 として、画像 A の画像名である「図 1」が最初に出現するパラグラフ p 2 から、次の画像である画像 B、C の画像名である「図 2」が出現するパラグラフ p 5 の直前のパラグラフ p 4 までを特定した場合の例を示している。なお、抽出した依存パラグラフ内の文字列は、画像情報 3 1 0 の該当の画像のエントリの依存パラグラフ 3 1 7 に格納する。

10

【 0 0 3 5 】

ここで、実際は、パラグラフ p 5 以降にも画像 A（「図 1」）についての記述がされているパラグラフが存在する場合も想定される。しかし、これらのパラグラフについては、特定するのに多くの処理を要するのに比して、その記述内容と画像 A との直接の関連度はそれほど高くない場合が多く、これらのパラグラフから取得されるキーワードの画像 A との適合度は低い場合が多い。また、記述内容と画像 A との関連度がある場合であっても、記述内容が依存パラグラフと同じような内容である等により、適合度が高いキーワードを独自に抽出できるケースはそれほど多くない。従って、本実施の形態では、上述したように、次の画像名が出現するまでのパラグラフを依存パラグラフとすることで、効率良く十分な精度のキーワードが抽出できる依存パラグラフの特定を可能とする。

20

【 0 0 3 6 】

なお、ステップ S 1 0 2 において画像名とキャプションが特定できなかった場合（画像にキャプション領域がない場合や、キャプション領域を有していても正規表現と一致する画像名がない場合など）は、ステップ S 1 0 3 以降の処理は行わず、キーワードを付与しないようにしてもよいし、例えば、特許文献 2 に記載されているような方法やその他の方法により、依存パラグラフに相当するパラグラフを特定するようにしてもよい。

【 0 0 3 7 】

次に、単語抽出部 1 1 4 により、ステップ S 1 0 3 で特定した依存パラグラフから自然言語処理によって単語（複合名詞）を抽出する、単語抽出処理を行う（ステップ S 1 0 4）。ここでは、例えば、一般的な形態素解析により依存パラグラフから複合名詞を抽出する。なお、ステップ S 1 0 2 で特定したキャプションは、処理の便宜上、例えば、依存パラグラフの 0 段落目（先頭）に相当するものとして依存パラグラフに含めるものとし、同様に形態素解析を行って複合名詞を抽出する。

30

【 0 0 3 8 】

図 1 では、画像 A（「図 1」）について、依存パラグラフ P 1（キャプション c 1、およびパラグラフ p 2 ~ p 4）からそれぞれ、w 1 ~ w 5 の各単語（複合名詞）を抽出した場合の例を示している。ここで、例えばパラグラフ p 3 で単語 w 2 が 2 つ抽出されているのは、単語 w 2 がパラグラフ p 3 で 2 回出現していることを示している。

【 0 0 3 9 】

次に、スコアリング部 1 1 5 により、ステップ S 1 0 4 で抽出した各単語について、依存パラグラフ内での特徴度を所定の方法によりスコアリングする、スコアリング処理を行う（ステップ S 1 0 5）。ここでは、各単語について、後述する T F × I P F 値（Term Frequency × Inversed Paragraph Frequency）を算出して特徴度とする。図 1 では、単語 w 1 ~ w 5 について、それぞれ T F × I P F 値を算出して特徴度とした場合の例を示している。

40

【 0 0 4 0 】

次に、画像抽出部 1 1 0 により、ステップ S 1 0 5 で算出した各単語の特徴度が上位の所定の単語をキーワードとして抽出し、抽出したキーワードを対象の画像のインデックスとし、その特徴度を対象の画像に対する適合度として、インデックステーブル 3 2 0 に格納する、インデックス作成処理を行う（ステップ S 1 0 6）。

50

【 0 0 4 1 】

図 1 では、単語 $w_1 \sim w_5$ のうち、例えば、特徴度が平均値以上である単語 w_1 、 w_3 、 w_4 の 3 つをキーワードとして抽出し、画像 A のインデックスとした場合の例を示している。なお、ここでは平均値以上の特徴度を有する単語をキーワードとして抽出しているが、例えば、特徴度が上位から所定の順位のものまでを抽出するなど他の方法であってもよい。また、キャプションに含まれる単語は、画像に直接的に言及しているということから、特徴度のスコアに関わりなくキーワードとして抽出するようにしてもよい。

【 0 0 4 2 】

抽出したキーワードと対応する画像の ID、およびその適合度は、インデックステーブル 3 2 0 のキーワード 3 2 1、画像 ID 3 2 2、および適合度 3 2 3 にそれぞれ格納する。以上の処理により、文書 2 0 0 群から画像を抽出し、抽出した画像に対してキーワードを効率的に精度良く付与して、適合度と合わせてインデックス化したインデックステーブル 3 2 0 を生成することができる。

【 0 0 4 3 】

[スコアリングとキーワード抽出]

以下では、スコアリング部 1 1 5 におけるスコアリング処理（ステップ S 1 0 5）、および、インデックス作成部 1 1 6 におけるインデックス作成処理（ステップ S 1 0 6）について説明する。スコアリング処理（ステップ S 1 0 5）では、図 1 の単語抽出処理（ステップ S 1 0 4）にて抽出された依存パラグラフ内の各単語について、特徴度として $TF \times IPF$ 値を算出してスコアリングする。 $TF \times IPF$ 値とは、 TF (Term Frequency) 値と IPF (Inverse Paragraph Frequency) 値の積である。

【 0 0 4 4 】

TF 値および IDF (Inversed Document Frequency) 値を用いてある文書中の特徴的な単語（重要とみなされる単語）を抽出することは一般的に行われている。本実施の形態のスコアリング処理でもこの手法を適用して特徴度を算出するが、本実施の形態では、 TF 値および IDF 値の算出時における単位である「文書 (Document)」を依存パラグラフ P とした、 TF 値および IPF 値を用いて特徴度を算出する。なお、特徴度の算出手法はこれに限るものではなく、単語毎に数値として画像との適合度を評価することが可能な手法であれば利用することができる。

【 0 0 4 5 】

本実施の形態のスコアリング処理において、 TF 値は、依存パラグラフ P 内における各単語（複合名詞）の出現頻度であり、この値が大きいほど当該単語は依存パラグラフ P （すなわち対応する画像）の特徴をよく表しているものと考えられる。ある依存パラグラフ P_j における単語 w_i の TF 値は、例えば、依存パラグラフ P_j 内の単語 w_i の出現頻度を、依存パラグラフ P_j において出現する延べ単語数で正規化して以下の式で表される。

【 0 0 4 6 】

【 数 1 】

$$tf_{w_i} = \frac{|w_i|}{\sum_n |w_n|}$$

tf_{w_i} : 単語 w_i の TF (Term Frequency) 値

$|w_i|$: 単語 w_i の依存パラグラフ P_j 内での出現頻度

$\sum_n |w_n|$: 依存パラグラフ P_j に出現する延べ単語数

【 0 0 4 7 】

一方、 TF 値が大きい単語であっても、他の画像についての依存パラグラフ P にも頻繁

10

20

30

40

50

に出現する単語は、特定の依存パラグラフ P の特徴を表す単語ではない一般的な単語である場合が多い。ここで、IPF 値は、対象の単語が出現する依存パラグラフ P の数の逆数であり、この値が大きいほどこの単語が出現する依存パラグラフ P の数が少ない。すなわち、この単語は特定の依存パラグラフ P の特徴をよく表しているものと考えられる。ある単語 w_i の IPF 値は、例えば、単語 w_i が出現する依存パラグラフ P の数の逆数を、対象の文書 200 内の全ての依存パラグラフ P の数で正規化して以下の式で表される。

【0048】

【数2】

$$ipf_{w_i} = \log \frac{|P|}{|w_i \in P|}$$

10

ipf_{w_i} : 単語 w_i の IPF (Inversed Paragraph Frequency) 値

$|P|$: 依存パラグラフ P の総数

$|w_i \in P|$: 単語 w_i を含む依存パラグラフ P の数

【0049】

20

上記の TF 値と IPF 値の両者の値が大きい単語 w_i が、対象の依存パラグラフ P (すなわち対応する画像) の特徴を真によく表していると考えられるため、TF 値と IPF 値の積である TF × IPF 値を算出して、これを単語 w_i の特徴度のスコアとする。この TF × IPF 値を、依存パラグラフ P 内の各単語について算出する。TF × IPF 値が大きい単語は、対象の画像の内容をよく表しており、キーワードとしての適合度が高いものと考えられる。なお、上記の TF 値、IPF 値の算出式については一例であり、精度や処理時間などに応じて正規化や対数計算の式などを適当なものにすることができる。

【0050】

図4は、図1の例に示した画像 A について、その依存パラグラフ P1 に含まれる各単語の TF × IPF 値を算出した例を示した図である。画像 A についての依存パラグラフ P1 (キャプション c1 およびパラグラフ p2 ~ p4) に含まれる単語 $w_1 \sim w_5$ について、依存パラグラフ P1 内での出現頻度に基づいて数1により算出した TF 値と、出現した依存パラグラフ P の数に基づいて数2により算出した IPF 値、および TF × IPF 値のスコア (特徴度) が示されている。この特徴度に基づいて、インデックス作成処理 (ステップ S106) では、例えば、各単語の特徴度がその平均値 (0.162) 以上である単語 w_1 、 w_3 、 w_4 の3つをキーワードとして抽出する。これにより、画像 A の内容をよく表した精度の高いキーワードを抽出することができる。

30

【0051】

[画像検索処理]

以下では、ユーザがクライアント端末 400 を利用して文書 200 群に含まれる画像を検索する際のユーザインタフェースおよび画像検索部 120 での画像検索処理について説明する。図5は、画像を検索する際にクライアント端末 400 に表示されるユーザインタフェースの例を示した図である。当該画面は、上述したように、画像検索部 120 のユーザインタフェース部 122 によって、例えば、HTML によって生成され、図示しない Web サーバプログラムを介して、クライアント端末 400 上の図示しない Web ブラウザによって表示される。

40

【0052】

図5に示した画面の上部には、例えば、ユーザが画像を検索するための検索語を指定することができるフィールドを有する。当該フィールドにはユーザが検索語を複数指定することも可能である。なお、本実施の形態では、画像検索サーバ 100 のインデックスステー

50

ブル320にキーワードのリストを有しているため、これを参照することにより、ユーザが検索語を入力している途中であっても、途中まで入力された文字列に一致するキーワードの候補を「単語候補」のフィールドに表示することが可能である。ユーザは、表示された候補の中から所望のキーワードをマウスによるクリック等で選択して、検索語として確定させることができる。

【0053】

このキーワードの候補を表示する処理は、Google（登録商標）等のWebサイトで一般的に行われているように、例えば、当該画面コンテンツにAjax（Asynchronous JavaScript（登録商標）+XML）等を利用したモジュールを組み込み、当該モジュールが、ユーザが入力した検索語の文字列を取得して画像検索サーバ100に非同期で送信し、画像検索サーバ100では、画像検索部120の検索処理部121によりインデックステーブル320を検索することによって、入力された文字列を先頭を含むキーワード321のリストを取得してクライアント端末400に送信し、クライアント端末400によって「単語候補」のフィールドに表示することで実現することができる。

10

【0054】

図5では、ユーザが検索語として「20」まで入力した時点で、これに該当するキーワード（検索語の候補）として、「2010年」、「2060年」、「20世紀」の単語をそれぞれインデックステーブル320から取得して「単語候補」フィールドに表示した場合の例を示している。これにより、ユーザが検索語としてキーワードを指定する際の労力を大幅に低減させることができる。

20

【0055】

検索語が確定すると、確定した複数の検索語のAND条件で、検索語に一致するキーワードが付与された文書200群内の画像を「画像一覧」のフィールドに表示する。ここでは、例えば、上述のようなAjax等を利用したモジュールや、検索ボタン等の押下に伴う処理によって、確定した複数の検索語を画像検索サーバ100に送信する。

【0056】

画像検索サーバ100は、検索処理部121によりインデックステーブル320を検索し、受信した検索語のAND条件により該当する画像ID322を取得する。さらに、画像情報310から、対応する画像ID311のエントリの画像データ312や、文書名313、位置314、画像名315、キャプション316等の他の情報を取得してクライアント端末400に送信し、クライアント端末400によって「画像一覧」フィールドに画像データや他の情報を表示する。

30

【0057】

なお、複数の画像を表示する際に、例えば、画像検索サーバ100からクライアント端末400に送信する画像の検索結果の情報に、キーワードと画像との適合度の情報をインデックステーブル320の適合度323から取得して追加することができる。これにより、クライアント端末400では、例えば、画像のキーワードに対する適合度の値（複数のキーワードに対応する場合はその合計）が大きいものを、表示順序を上位にしたり、視覚的に目立つようにしたりなど優先的に表示し、画像とキーワードとの適合度に応じて出力方法を柔軟に制御することが可能となる。

40

【0058】

また、表示する画像の適合度の閾値をユーザにより設定できるようにしておき、閾値未満の適合度の画像は表示しない（もしくは画像を検索する際の対象から除外する）ようにしてもよい。例えば、依存パラグラフPが十分な長さを有しておらず短い場合や、依存パラグラフP内のどの単語もあまり特徴的ではなく、各単語のTF×IPF値が近似する（TF×IPF値の分散が小さい）場合などは、TF×IPF値が小さくなる傾向が高い。この場合、これらの単語はキーワードとしての精度が低いため、閾値を調整することによって対応する画像が表示されないようにすることができる。

【0059】

以上のように、本実施の形態の文書内画像検索システム1によれば、文書200から抽

50

出した画像に対して、当該画像について記述している依存パラグラフPを特定することで、画像に関連するキーワードを文脈情報を考慮して効率的に精度良く付与することが可能となり、キーワードを利用した高精度で効率の良い画像の検索が可能となる。また、各キーワードと画像の組合せに対してスコア（適合度）を有するため、適合度に応じて検索結果の画像の表示順序等の出力方法を制御することによってユーザの利便性を高めることが可能となる。

【0060】

<実施の形態2>

本発明の実施の形態2である文書内画像検索システムは、上述した実施の形態1の文書内画像検索システム1において、スコアリング部115での特徴度のスコアリング処理（ステップS105）で、依存パラグラフP内の各単語の特徴度をスコアリングする際に、単語の出現位置の情報に基づいて重み付けを行うことによって、抽出するキーワードの精度をより高くすることを可能とするものである。なお、スコアリング部115以外の他の構成や処理内容は、実施の形態1で説明したものと同様であるため、再度の説明は省略する。

10

【0061】

図6は、図1の例に示した画像Aについて、その依存パラグラフP1に含まれる各単語のTF×IPF値を、単語の出現位置に応じて重み付けして算出した例を示した図である。まず、依存パラグラフP1内で出現する各単語（w1～w5）を、依存パラグラフP1内で出現した行に応じて図6の中段の表に示すように集計する。このとき、例えば、キャプション中の単語は0行目に出現したものとし、依存パラグラフP1内の各パラグラフ（p2～p4）を連結して1行目からカウントするものとする。

20

【0062】

ここで、行の値をx、重み付け値をyとした重み付け関数 $y = f(x)$ を利用して各行での重み付け値を算出する。図6では、0行目で重み付け値が1であり、20行目で0となる、傾きマイナス0.05の一次関数によって重み付け値を算出している。これは、キャプションを始めとして、依存パラグラフP内の先頭に近い位置で出現した単語ほど、対応する画像に対して直接的に言及している場合が多いことを考慮した重み付け関数である。

【0063】

この重み付け関数によれば、キャプションに含まれる単語を無条件に抽出するという処理を行わなくても、これらの単語には自動的に大きい重み付け値を付与することができるため、キーワードとして抽出されるようにすることができる。なお、重み付け関数は、図6に示したものに限らず、例えば、対象の画像の出現行を中心とした正規分布曲線を有する確率密度関数など、単語の出現位置による画像との適合度のモデルに基づいて種々のものを用いることができる。

30

【0064】

この重み付け値に基づいて、各単語の出現頻度の値を図6の下段の表に示すように補正する。例えば、単語w1は、0行目（重み付け値1.00）で1回、1行目（重み付け値0.95）で1回出現しているため、重み補正後の出現頻度は、

$$1 \times 1.00 + 1 \times 0.95 = 1.95$$

となる。また、単語w2は、1行目（重み付け値0.95）で1回、5行目（重み付け値0.75）で1回、6行目（重み付け値0.70）で1回出現しているため、重み補正後の出現頻度は、

40

$$1 \times 0.95 + 1 \times 0.75 + 1 \times 0.70 = 2.40$$

となる。以下、単語w3～w5についても同様に算出する。

【0065】

以上のように算出された重み補正後の出現頻度に基づいて、上述した数1により重み補正後のTF値を算出し、実施の形態1の場合と同様に数2により算出したIPF値と乗算することで、重み補正後のTF×IPF値を算出する。図6の例では、重み補正後のTF

50

× I P F 値の平均値は 0 . 1 7 4 となり、平均値以上の単語をキーワードとして抽出すると、単語 w 1、w 3 の 2 つとなる。実施の形態 1 の場合と比較して単語 w 4 がキーワードとして抽出されなくなっているが、これは、単語 w 4 は依存パラグラフ P 1 の中で後半部分に多く出てきていることから、上述の処理により特徴度（画像との適合度）が相対的に低いものと判断されるためである。

【 0 0 6 6 】

以上のように、本実施の形態の文書内画像検索システム 1 によれば、依存パラグラフ P 内の各単語の特徴度をスコアリングする際に、単語の出現位置を変数とした重み付け関数を利用して出現頻度に重み付けを行うことによって、単語の出現位置による画像についての言及の程度の違いを考慮して特徴度をスコアリングする。これにより、抽出するキーワードの精度をより高くし、画像検索の際のノイズを低減することが可能となる。

10

【 0 0 6 7 】

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

【 産業上の利用可能性 】

【 0 0 6 8 】

本発明は、画像の内容を表すキーワードを指定することにより文書内の画像を検索する文書内画像検索方法および文書内画像検索システムに利用可能である。

【 符号の説明 】

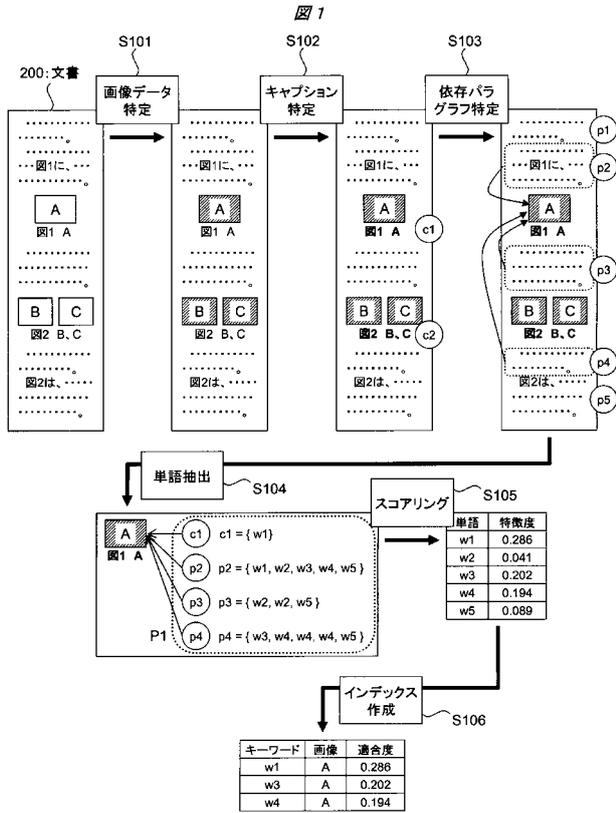
20

【 0 0 6 9 】

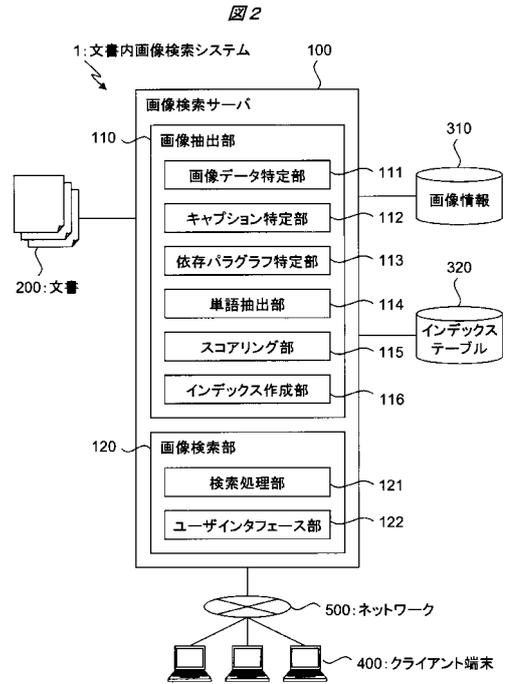
1 ... 文書内画像検索システム、
 1 0 0 ... 画像検索サーバ、 1 1 0 ... 画像抽出部、 1 1 1 ... 画像データ特定部、 1 1 2 ... キャプション特定部、 1 1 3 ... 依存パラグラフ特定部、 1 1 4 ... 単語抽出部、 1 1 5 ... スコアリング部、 1 1 6 ... インデックス作成部、 1 2 0 ... 画像検索部、 1 2 1 ... 検索処理部、 1 2 2 ... ユーザインタフェース部、
 2 0 0 ... 文書、
 3 1 0 ... 画像情報、 3 1 1 ... 画像 I D、 3 1 2 ... 画像データ、 3 1 3 ... 文書名、 3 1 4 ... 位置、 3 1 5 ... 画像名、 3 1 6 ... キャプション、 3 1 7 ... 依存パラグラフ、 3 2 0 ... インデックステーブル、 3 2 1 ... キーワード、 3 2 2 ... 画像 I D、 3 2 3 ... 適合度、
 4 0 0 ... クライアント端末、
 5 0 0 ... ネットワーク。

30

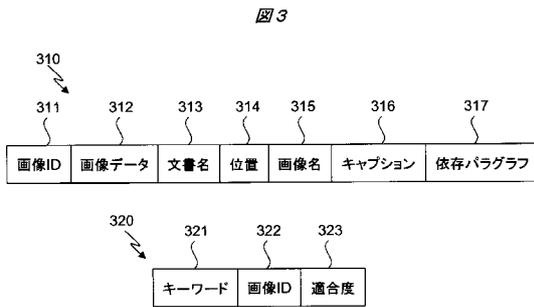
【 図 1 】



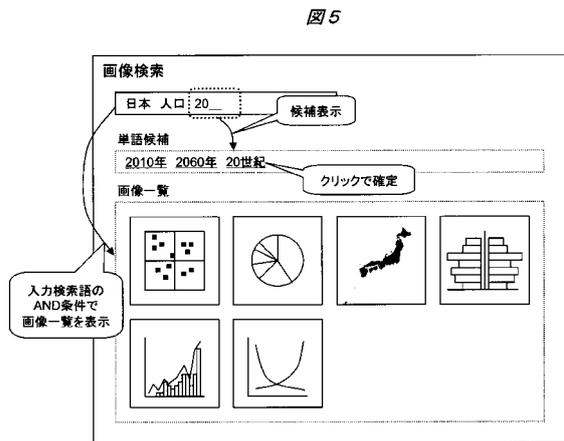
【 図 2 】



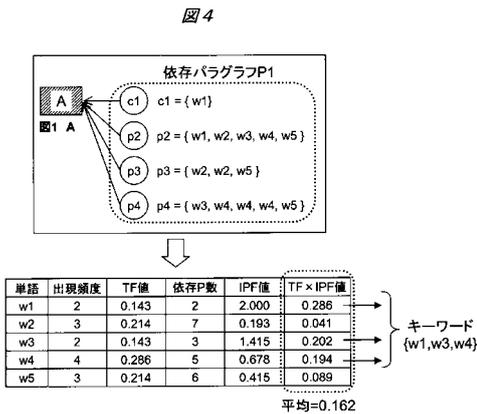
【 図 3 】



【 図 5 】



【 図 4 】



【 図 6 】

