

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5979650号
(P5979650)

(45) 発行日 平成28年8月24日 (2016. 8. 24)

(24) 登録日 平成28年8月5日 (2016. 8. 5)

(51) Int. Cl. F I
G O 6 F 1 7 / 2 7 (2 0 0 6 . 0 1) G O 6 F 1 7 / 2 7 6 7 0

請求項の数 20 (全 25 頁)

<p>(21) 出願番号 特願2014-152580 (P2014-152580)</p> <p>(22) 出願日 平成26年7月28日 (2014. 7. 28)</p> <p>(65) 公開番号 特開2016-31572 (P2016-31572A)</p> <p>(43) 公開日 平成28年3月7日 (2016. 3. 7)</p> <p>審査請求日 平成28年1月5日 (2016. 1. 5)</p> <p>早期審査対象出願</p>	<p>(73) 特許権者 390009531 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 INTERNATIONAL BUSINESS MACHINES CORPORATION 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 New Orchard Road, Armonk, New York 10504, United States of America</p> <p>(74) 代理人 100108501 弁理士 上野 剛史</p> <p style="text-align: right;">最終頁に続く</p>
--	--

(54) 【発明の名称】用語を適切な粒度で分割する方法、並びに、用語を適切な粒度で分割するためのコンピュータ及びそのコンピュータ・プログラム

(57) 【特許請求の範囲】

【請求項1】

用語を適切な粒度で分割する方法であって、コンピュータが抽出手段と分割手段とを備えており、前記方法は、

前記抽出手段が、コンテンツを格納したメモリ又は記憶装置から読み取り、構文解析により、粒度を規定する構成要素を前記コンテンツから抽出するステップを実行し、ここで、前記構成要素は、少なくとも1つの名詞又は記号を含む1又は複数の単語列であり、

前記分割手段が、前記用語がその一部に少なくとも1つの前記構成要素を含む場合に、前記用語を前記構成要素がある位置で分割し、当該分割した後の用語を当該分割した後の用語を入れるリストに格納するステップを実行し、

前記分割するステップが、

前記用語が当該用語の末尾から最長一致する前記構成要素を含む場合に、前記用語を前記末尾から最長一致する構成要素がある位置で分割するステップと、

前記用語から前記末尾から最長一致する前記構成要素を除いた後の用語が当該除いた後の用語の先頭から最長一致する前記構成要素を含む場合に、前記除いた後の用語を前記先頭から最長一致する構成要素がある位置で分割するステップと

を含む、

前記方法。

【請求項2】

前記用語を前記末尾から最長一致する前記構成要素がある位置で分割するステップが、

前記末尾から最長一致する前記構成要素を前記用語の主要語として保存するステップを含む、請求項 1 に記載の方法。

【請求項 3】

前記除いた後の用語を前記先頭から最長一致する前記構成要素がある位置で分割するステップが、

前記先頭から最長一致する前記構成要素を前記用語の第 1 の修飾語として保存するステップ

をさらに含む、請求項 1 又は 2 に記載の方法。

【請求項 4】

前記除いた後の用語を前記先頭から最長一致する前記構成要素がある位置で分割するステップが、

前記先頭から最長一致する前記構成要素以外の部分を第 2 の修飾語として保存するステップ

を含む、請求項 3 に記載の方法。

【請求項 5】

前記構成要素を抽出するステップが、

前記コンテンツ中のテキストそれぞれに前記構文解析を適用して、文節を抽出するステップと、

前記抽出した文節のうちの名詞又は記号を含む文節から前記構成要素となりうる部分を抽出するステップと

を含む、請求項 1 ~ 4 のいずれか一項に記載の方法。

【請求項 6】

前記構成要素を抽出するステップが、

前記コンテンツから、前記構成要素を抽出する対象のテキストを切り出すステップ

をさらに含み、

前記文節を抽出するステップが、前記切り出したテキストそれぞれに前記構文解析を適用して行われる、請求項 5 に記載の方法。

【請求項 7】

前記構成要素を抽出するステップが、

前記切り出したテキストを事前定義した文字がある場所で分割するステップ

をさらに含み、

前記文節を抽出するステップが、前記分割したテキストそれぞれに前記構文解析を適用して行われる、請求項 6 に記載の方法。

【請求項 8】

前記用語が用語リスト中の用語であり、

前記構成要素を抽出するステップが、

前記構成要素となりうる部分から前記用語リスト中にある同じ用語を削除し、当該削除した残りを前記構成要素とするステップ

をさらに含む、請求項 5 ~ 7 のいずれか一項に記載の方法。

【請求項 9】

前記分割するステップが、

前記分割するステップに従い分割した分割回数と、予め設定された分割回数を規定する分割パラメータとを比較し、前記分割回数が前記分割パラメータよりも少ないことに応じて、前記用語を前記構成要素がある位置でさらに分割するステップ

を含む、請求項 1 ~ 8 のいずれか一項に記載の方法。

【請求項 10】

前記用語が用語リスト中の用語である、請求項 1 ~ 7 及び 9 のいずれか一項に記載の方法。

【請求項 11】

前記用語が、前記コンテンツ中の所定の長さよりも長い用語である、請求項 1 ~ 7 及び

10

20

30

40

50

10のいずれか一項に記載の方法。

【請求項12】

前記用語が名詞、記号又はそれらの組み合わせを含む単語列である、請求項1～11のいずれか一項に記載の方法。

【請求項13】

前記用語が複合名詞である、請求項1～11のいずれか一項に記載の方法。

【請求項14】

前記構成要素が、少なくとも1つの名詞又は記号を含む1又は複数の単語列である、請求項1～13のいずれか一項に記載の方法。

【請求項15】

用語を適切な粒度で分割するためのコンピュータであって、
コンテンツを格納するメモリ又は記憶装置と、
前記メモリ又は記憶装置からコンテンツを読み取り、構文解析により、粒度を規定する構成要素を前記コンテンツから抽出する抽出手段であって、前記構成要素は、少なくとも1つの名詞又は記号を含む1又は複数の単語列である、前記抽出手段と、

前記用語がその一部に少なくとも1つの前記構成要素を含む場合に、前記用語を前記構成要素がある位置で分割し、当該分割した後の用語を当該分割した後の用語を入れるリストに格納する分割手段と

を備えており、

前記分割手段が、

前記用語が当該用語の末尾から最長一致する前記構成要素を含む場合に、前記用語を前記末尾から最長一致する構成要素がある位置で分割すること、

前記用語から前記末尾から最長一致する前記構成要素を除いた後の用語が当該除いた後の用語の先頭から最長一致する前記構成要素を含む場合に、前記除いた後の用語を前記先頭から最長一致する構成要素がある位置で分割すること

を実行する、前記コンピュータ。

【請求項16】

前記分割手段が、前記末尾から最長一致する前記構成要素を前記用語の主要語として保存することを実行する、請求項15に記載のコンピュータ。

【請求項17】

前記分割手段が、前記先頭から最長一致する前記構成要素を前記用語の第1の修飾語として保存することを実行する、請求項15又は16に記載のコンピュータ。

【請求項18】

前記分割手段が、前記先頭から最長一致する前記構成要素以外の部分を第2の修飾語として保存することを実行する、請求項17に記載のコンピュータ。

【請求項19】

前記抽出手段が、

前記コンテンツ中のテキストそれぞれに前記構文解析を適用して、文節を抽出すること、

前記抽出した文節のうちの名詞又は記号を含む文節から前記構成要素となりうる部分を抽出すること

を実行する、請求項15～18のいずれか一項に記載のコンピュータ。

【請求項20】

用語を適切な粒度で分割するためのコンピュータ・プログラムであって、コンピュータに、請求項1～14のいずれか一項に記載の方法の各ステップを実行させる、前記コンピュータ・プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、用語の分割技法に関する。特に、本発明は、用語を適切な粒度で分割する

10

20

30

40

50

技法に関する。

【背景技術】

【0002】

名詞や名詞相当の接辞が複数接続して、複数の単語（例えば、2～6語の単語）からなる複合名詞が限りなく作り出される。

【0003】

システム開発において作成される用語集は、上記複合名詞を含む。しかしながら、複合名詞はその意味が一見して不明なことが多い。

【0004】

特に、金融機関のシステム用の用語集は、例えば「定期預金毎月振替限度額」（漢字からなる複合名詞である）や「基準価額型金信解約合計金額」（漢字からなる複合名詞である）などの複合名詞を含みうる。

10

【0005】

また、英語においても、語が複数連なって名詞句を形成する。例えば、「Beneficiary right seller's business security deposit」や「Financial instruments intermediary service」である。

【0006】

複合名詞に形態素解析技術を使用して、複合名詞を分割する技術が知られている（例えば、下記非特許文献1～2を参照）。しかしながら、形態素解析技術では、形態素解析器が持つシステム辞書及び文法に基づいて複合名詞が分割されるために、必ずしも望ましい結果が得られていない。

20

【0007】

形態素解析器が、例えば上記「定期預金毎月振替限度額」（漢字からなる複合名詞である）を形態素解析技術を使用して分割すると、「定期+預金+毎月+振替+限度額」（漢字からなる単語である）のように、単語ごとに細分化して分割してしまう。また、形態素解析器が、上記「基準価額型金信解約合計金額」（漢字からなる複合名詞である）を形態素解析技術を使用して分割すると、「基準+価額+型+金+信+解約+合計+金額」（漢字からなる単語である）のように、本来一語である「金信」が漢字（すなわち、一つの単語であり、「金銭信託」（漢字である）の略語である）一文字に細分化して分割してしまう。

30

【0008】

形態素解析器が、例えば上記「business security deposit」を形態素解析技術を使用して分割すると、「"business security" + "deposit"」若しくは「"business" + "security deposit"」に分割するのか、又は、分割せずに「"business security deposit"」のままにするのかを判定することは難しい。

【0009】

下記特許文献1～10は、文章の解析やキーワードの抽出を記載する。

【先行技術文献】

【特許文献】

【0010】

40

【特許文献1】特開2007-257390号公報

【特許文献2】特開平10-207890号公報

【特許文献3】特開平7-85101号公報

【特許文献4】特開2007-264718号公報

【特許文献5】特開平8-305695号公報

【特許文献6】特開2001-325284号公報

【特許文献7】特開2010-204866号公報

【特許文献8】特開2011-96245号公報

【特許文献9】特開2008-140359号公報

【特許文献10】特開2012-234512号公報

50

【非特許文献】

【0011】

【非特許文献1】太田悟 等、「規則・用例融合型の日本語複合名詞構造解析法」、言語処理学会 第3回年次大会 発表論文集、313～316頁1997年3月、<URL:http://www.anlp.jp/proceedings/annual_meeting/2003/pdf_dir/C6-2.pdf> から入手可能

【非特許文献2】高橋允彦 等、「構造化規則を用いた日本語複合名詞解析」、言語処理学会 第9回年次大会 発表論文集、541～544頁、2003年3月、<URL:http://www.anlp.jp/proceedings/annual_meeting/2003/pdf_dir/C6-2.pdf> から入手可能

【非特許文献3】宮崎正弘 等、「構造化チャートパーザを用いた日本語複合名詞構造解析器」、言語処理学会、2008年、<URL:http://www.languetech.co.jp/out/08nlp-miyazaki.pdf> から入手可能

10

【発明の概要】

【発明が解決しようとする課題】

【0012】

複合名詞は限りなく作り出される故に、当該複合名詞の全てを辞書に登録することは不可能である。その為に、複合名詞を分割することが試みられている。

【0013】

しかしながら、上記した通り、作り出される複合名詞の数は膨大であるので、人が、用語集に登録される全ての用語が複合名詞であるかどうかをチェックし、複合名詞である場合に当該複合名詞を分割することは現実的でない。

20

【0014】

また、人手によって複合名詞を分割する場合には、複合名詞の分割の仕方が作業者に依存してしまい、例えば何が主要語又は修飾語であるのかの判断が作業者の主観によって異なる。

【0015】

また、主要語又は修飾語のリストは、事前に定義されていない場合がほとんどであり、また事前に定義されていたとしてもそのリストは不完全である。

【0016】

また、例えばシステム開発において作成される用語集に登録される複合名詞の数は数百にものぼり、特に、巨大な業務システム開発において作成される用語集に登録される用語の数は1000を超える場合がある。

30

【0017】

さらに、特定のプロジェクトの用語辞書に含まれる用語を、単語ごと又は一文字まででなく、幾つかの単語のまとまりとして分割したいという要望がある。例えば上記「定期預金毎月振替限度額」（漢字からなる複合名詞である）を、「定期預金＋毎月＋振替限度額」（いずれも漢字からなる単語である）に分割したいという要望がある。

【0018】

そこで、本発明は、特定のプロジェクトが独自に持つ粒度で用語を分割する技法を提供することを目的とする。

【0019】

また、本発明は、作業者の主観に依存することなしに、用語を分割する技法を提供することを目的とする。

40

【0020】

さらに、本発明は、主要語又は修飾語のリストを必要とすることなしに、用語を分割する技法を提供することを目的とする。

【課題を解決するための手段】

【0021】

本発明は、用語を適切な粒度で分割する技法を提供する。当該技法は、用語を適切な粒度で分割する方法、並びに、用語を適切な粒度で分割するためのコンピュータ、そのコンピュータ・プログラム及びコンピュータ・プログラム製品を包含しうる。

50

【 0 0 2 2 】

本発明に従う第 1 の態様において、用語を適切な粒度で分割する方法は、コンピュータが、

(A) 構文解析により、コンテンツから粒度を規定する構成要素 (element word) を抽出するステップと、

(B) 上記用語がその一部に少なくとも 1 つの上記構成要素を含む場合に、上記用語を上記構成要素がある位置で分割するステップと

を実行することを含む。

【 0 0 2 3 】

本発明の一つの実施態様において、(A) 上記構成要素を抽出するステップが、

(A - 3) 上記コンテンツ中のテキストそれぞれに上記構文解析を適用して、文節を抽出するステップと、

(A - 4) 上記抽出した文節のうちの名詞又は記号を含む文節から上記構成要素となりうる部分を抽出するステップと

を含みうる。

【 0 0 2 4 】

本発明の一つの実施態様において、(A) 上記構成要素を抽出するステップが、

(A - 1) 上記コンテンツから、上記構成要素を抽出する対象のテキストを切り出すステップ

をさらに含み、

上記文節を抽出するステップが、上記切り出したテキストそれぞれに上記構文解析を適用して行われうる。

【 0 0 2 5 】

本発明の一つの実施態様において、(A) 上記構成要素を抽出するステップが、

(A - 2) 上記切り出したテキストを事前定義した文字がある場所で分割するステップをさらに含み、

上記文節を抽出するステップが、上記分割したテキストそれぞれに上記構文解析を適用して行われうる。

【 0 0 2 6 】

本発明の一つの実施態様において、上記用語が用語リスト中の用語であり、

(A) 上記構成要素を抽出するステップが、

(A - 5) 上記構成要素となりうる部分のうちから上記用語リスト中にある用語を削除し、当該削除した残りを上記構成要素とするステップ

をさらに含みうる。

【 0 0 2 7 】

本発明の一つの実施態様において、(B) 上記分割するステップが、

(B - 1) 上記用語が当該用語の末尾から最長一致する上記構成要素 (第 1 の構成要素) を含む場合に、上記用語を上記末尾から最長一致する上記構成要素 (第 1 の構成要素) がある位置で分割するステップ

を含みうる。

【 0 0 2 8 】

本発明の一つの実施態様において、(B - 1) 上記用語を上記末尾から最長一致する上記構成要素 (第 1 の構成要素) がある位置で分割するステップが、

上記末尾から最長一致する上記構成要素 (第 1 の構成要素) を上記用語の主要語として保存するステップ

を含みうる。

【 0 0 2 9 】

本発明の一つの実施態様において、(B) 上記分割するステップが、

(B - 2) 上記用語から上記末尾から最長一致する上記構成要素 (第 1 の構成要素) を除いた後の用語が当該除いた後の用語の先頭から最長一致する上記構成要素 (第 2 の構成

10

20

30

40

50

要素)を含む場合に、上記除いた後の用語を上記先頭から最長一致する上記構成要素(第2の構成要素)がある位置で分割するステップを含みうる。

【0030】

本発明の一つの実施態様において、(B-2)上記除いた後の用語を上記先頭から最長一致する上記構成要素(第2の構成要素)がある位置で分割するステップが、

上記先頭から最長一致する上記構成要素(第2の構成要素)を上記用語の第1の修飾語として保存するステップ

をさらに含みうる。

【0031】

本発明の一つの実施態様において、(B-2)上記除いた後の用語を上記先頭から最長一致する上記構成要素(第2の構成要素)がある位置で分割するステップが、

上記先頭から最長一致する上記構成要素(第2の構成要素)以外の部分を第2の修飾語として保存するステップ

を含みうる。

【0032】

本発明の一つの実施態様において、(B)上記分割するステップが、

予め設定された分割回数を規定する分割パラメータに従って、上記用語を上記構成要素がある位置で分割するステップ

を含みうる。

【0033】

本発明に従う第2の態様において、用語を適切な粒度で分割するためのコンピュータは、

構文解析により、粒度を規定する構成要素をコンテンツから抽出する抽出手段と、

上記用語がその一部に少なくとも1つの上記構成要素を含む場合に、上記用語を上記構成要素がある位置で分割する分割手段と

を備えている。

【0034】

本発明の一つの実施態様において、上記抽出手段が、上記コンテンツ中のテキストそれぞれに上記構文解析を適用して、文節を抽出し、上記抽出した文節のうちの名詞又は記号を含む文節から上記構成要素となりうる部分を抽出しうる。

【0035】

本発明の一つの実施態様において、上記抽出手段が、さらに、上記コンテンツから、上記構成要素を抽出する対象のテキストを切り出し、当該切り出したテキストそれぞれに上記構文解析を適用して、上記文節を抽出しうる。

【0036】

本発明の一つの実施態様において、上記抽出手段が、さらに、上記切り出したテキストを事前定義した文字がある場所で分割し、当該分割したテキストそれぞれに上記構文解析を適用して、上記文節を抽出しうる。

【0037】

本発明の一つの実施態様において、上記用語が用語リスト中の用語であり、上記抽出手段が、上記構成要素となりうる部分のうちから上記用語リスト中にある用語を削除し、当該削除した残りを上記構成要素としうる。

【0038】

本発明の一つの実施態様において、上記分割手段が、上記用語が当該用語の末尾から最長一致する上記構成要素(第1の構成要素)を含む場合に、上記用語を上記末尾から最長一致する上記構成要素(第1の構成要素)がある位置で分割しうる。

【0039】

本発明の一つの実施態様において、上記分割手段が、上記用語を上記末尾から最長一致する上記構成要素(第1の構成要素)がある位置で分割し、上記末尾から最長一致する上

10

20

30

40

50

記構成要素（第1の構成要素）を上記用語の主要語として保存しうる。

【0040】

本発明の一つの実施態様において、上記分割手段が、上記用語から上記末尾から最長一致する上記構成要素（第1の構成要素）を除いた後の用語が当該除いた後の用語の先頭から最長一致する上記構成要素（第2の構成要素）を含む場合に、上記除いた後の用語を上記先頭から最長一致する上記構成要素（第2の構成要素）がある位置で分割しうる。

【0041】

本発明の一つの実施態様において、上記分割手段が、上記除いた後の用語を上記先頭から最長一致する上記構成要素（第2の構成要素）がある位置で分割し、上記先頭から最長一致する上記構成要素（第2の構成要素）を上記用語の第1の修飾語として保存しうる。

10

【0042】

本発明の一つの実施態様において、上記分割手段が、上記除いた後の用語を上記先頭から最長一致する上記構成要素（第2の構成要素）がある位置で分割し、上記先頭から最長一致する上記構成要素（第2の構成要素）以外の部分を第2の修飾語として保存しうる。

【0043】

本発明の一つの実施態様において、上記分割手段が、予め設定された分割回数を規定する分割パラメータに従って、上記用語を上記構成要素がある位置で分割しうる。

【0044】

また、本発明に従う第3の態様において、コンピュータ・プログラム及びコンピュータ・プログラム製品は、上記コンピュータに、本発明に従う第1の態様に記載の用語を適切な粒度で分割する方法の各ステップを実行させる。

20

【0045】

本発明の実施態様に従うコンピュータ・プログラムはそれぞれ、一つ又は複数のフレキシブル・ディスク、MO、CD-ROM、DVD、BD、ハードディスク装置、USBに接続可能なメモリ媒体、ROM、MRAM、RAM等の任意のコンピュータ読み取り可能な記録媒体に格納することができる。当該コンピュータ・プログラムは、記録媒体への格納のために、通信回線で接続する他のデータ処理システム、例えばコンピュータからダウンロードしたり、又は他の記録媒体から複製したりすることができる。また、本発明の実施態様に従うコンピュータ・プログラムは、圧縮し、又は複数に分割して、単一又は複数の記録媒体に格納することもできる。また、様々な形態で、本発明の実施態様に従うコンピュータ・プログラム製品を提供することも勿論可能であることにも留意されたい。本発明の実施態様に従うコンピュータ・プログラム製品は、例えば、上記コンピュータ・プログラムを記録した記憶媒体、又は、上記コンピュータ・プログラムを伝送する伝送媒体を包含しうる。

30

【0046】

本発明の上記概要は、本発明の必要な特徴の全てを列挙したものではなく、これらの構成要素のコンビネーション又はサブコンビネーションもまた、本発明となりうることに留意すべきである。

【0047】

本発明の実施態様において使用されるコンピュータの各ハードウェア構成要素を、複数のマシンと組み合わせ、それらに機能を配分し実施する等の種々の変更は当業者によって容易に想定され得ることは勿論である。それらの変更は、当然に本発明の思想に包含される概念である。ただし、これらの構成要素は例示であり、そのすべての構成要素が本発明の必須構成要素となるわけではない。

40

【0048】

また、本発明は、ハードウェア、ソフトウェア、又は、ハードウェア及びソフトウェアの組み合わせとして実現可能である。ハードウェアとソフトウェアとの組み合わせによる実行において、上記コンピュータ・プログラムをインストールされたコンピュータにおける当該プログラムの実行が典型的な例として挙げられる。かかる場合、当該コンピュータ・プログラムが当該コンピュータのメモリにロードされて実行されることにより、当該コ

50

ンピュータ・プログラムは、当該コンピュータを制御し、本発明にかかる処理を実行させる。当該コンピュータ・プログラムは、任意の言語、コード、又は、表記によって表現可能な命令群から構成されうる。そのような命令群は、当該コンピュータが特定の機能を直接的に、又は、1. 他の言語、コード若しくは表記への変換及び、2. 他の媒体への複製、のいずれか一方若しくは双方が行われた後に、実行することを可能にするものである。

【発明の効果】

【0049】

本発明の実施態様に従うと、コンテンツ（例えば、マニュアル、業務手順書）から抽出される構成要素に従う粒度で、用語（例えば、用語辞書中の用語）を分割することが可能になる。コンテンツ又はコンテンツの技術分野が異なれば当該コンテンツから抽出される構成要素は異なる。従って、本発明の実施態様に従い、コンテンツ又はコンテンツの技術分野に従って粒度が動的に変わることから、コンテンツ又はコンテンツの技術分野に適した粒度で用語が分割されることが可能になる。

10

【0050】

また、本発明の実施態様に従うと、作業者の主観に依存することなしに、コンテンツから抽出される構成要素に従う粒度で、用語が分割される。

【0051】

さらに、本発明の実施態様に従うと、主要語又は修飾語のリストを必要とすることなしに、コンテンツから抽出される構成要素に従う粒度で、用語が分割される。

【図面の簡単な説明】

20

【0052】

【図1】本発明の実施態様において使用されうるコンピュータの一例を示した図である。

【図2A】本発明の実施態様に従い、用語リスト中の用語（英語である）を、コンテンツから抽出された粒度を規定する構成要素がある位置で分割する例を示す。

【図2B】本発明の実施態様に従い、用語リスト中の用語（漢字を含む）を、コンテンツから抽出された粒度を規定する構成要素がある位置で分割する例を示す。

【図3A】本発明の実施態様に従い、粒度を規定する構成要素をコンテンツから抽出する処理の為にフローチャートを示す。

【図3B】本発明の実施態様に従い、用語を構成要素がある位置で分割する処理の為にフローチャートを示す。

30

【図4】図1に従うハードウェア構成を好ましくは備えており、図3A及び図3Bそれぞれに示すフローチャートに従って本発明の実施態様を実施するコンピュータの機能ブロック図の一例を示す図である。

【発明を実施するための形態】

【0053】

本発明の実施形態を、以下に図面に従って説明する。以下の図を通して、特に断らない限り、同一の符号は同一の対象を指す。本発明の実施形態は、本発明の好適な態様を説明するためのものであり、本発明の範囲をここで示すものに限定する意図はないことを理解されたい。

【0054】

40

図1は、本発明の実施態様において使用されうるコンピュータの一例を示した図である。

【0055】

本発明の実施態様に従うコンピュータは、1又は複数のコンピュータから構成されうる。

【0056】

図1は、本発明の実施態様において使用されうるコンピュータを実現するためのハードウェア構成の一例を示した図である。

【0057】

コンピュータ(101)は例えば、コンピュータ(例えば、デスクトップ・コンピュー

50

タ、ノートブック・コンピュータ、ウルトラブック・コンピュータ、サーバ・コンピュータ)でありうる。

【0058】

コンピュータ(101)は、CPU(102)とメイン・メモリ(103)とを備えており、これらはバス(104)に接続されている。CPU(102)は好ましくは、32ビット又は64ビットのアーキテクチャに基づくものである。当該CPU(102)は例えば、インテル社のCore(商標 i)シリーズ、Core(商標) 2シリーズ、Atom(商標)シリーズ、Xeon(登録商標)シリーズ、Pentium(登録商標)シリーズ若しくはCeleron(登録商標)シリーズ、AMD(Advanced Micro Devices)社のAシリーズ、Phenom(商標)シリーズ、Athlon(商標)シリーズ、Turion(商標)シリーズ若しくはSempron(商標)、又は、国際ナショナル・ビジネス・マシーンス・コーポレーションのPower(商標)シリーズでありうる。

10

【0059】

バス(104)には、ディスプレイ・コントローラ(105)を介して、ディスプレイ(106)、例えば液晶ディスプレイ(LCD)が接続されうる。また、液晶ディスプレイ(LCD)は例えば、タッチパネル・ディスプレイ又はフローティング・タッチ・ディスプレイであってもよい。ディスプレイ(106)は、コンピュータ(101)上で動作中のソフトウェア、例えば本発明の実施態様に従うコンピュータ・プログラムが稼働することによって表示される情報(例えば、用語リスト中の用語、コンテンツ、構成要素、又は分割された用語)を、適当なグラフィック・インタフェースで表示するために使用されうる。

20

【0060】

バス(104)には任意的に、例えばSATA又はIDEコントローラ(107)を介して、記憶装置(108)、例えばハードディスク又はソリッド・ステート・ドライブに接続されうる。

【0061】

バス(104)には任意的に、例えばSATA又はIDEコントローラ(107)を介して、記憶装置(108)、ドライブ(109)、例えばCD、DVD又はBDドライブが接続されうる。

30

【0062】

バス(104)には、周辺装置コントローラ(110)を介して、例えばキーボード・マウス・コントローラ又はUSBバスを介して、任意的に、キーボード(111)及びマウス(112)が接続されうる。

【0063】

記憶装置(108)には、オペレーティング・システム、Windows(登録商標)OS、UNIX(登録商標)、Linux(登録商標)(例えば、Red Hat(登録商標)、Debian(登録商標))、Mac OS(登録商標)、及びJ2EEなどのJava(登録商標)処理環境、Java(登録商標)アプリケーション、Java(登録商標)仮想マシン(VM)、Java(登録商標)実行時(JIT)コンパイラを提供するプログラム、本発明の実施態様に従うコンピュータ・プログラム、及びその他のプログラム、並びにデータ(例えば、用語リスト、コンテンツ)が、メイン・メモリ(103)にロード可能なように記憶されうる。

40

【0064】

記憶装置(108)は、コンピュータ(101)内に内蔵されていてもよく、当該コンピュータ(101)がアクセス可能なようにケーブル(例えば、USBケーブル又はLANケーブル)を介して接続されていてもよく、又は、当該コンピュータ(101)がアクセス可能なように有線又は無線ネットワークを介して接続されていてもよい。

【0065】

ドライブ(109)は、必要に応じて、例えばCD-ROM、DVD-ROM又はBD

50

- ROMからプログラム、例えばオペレーティング・システム又はアプリケーションを記憶装置(108)にインストールするために使用されうる。

【0066】

通信インタフェース(114)は、例えばイーサネット(登録商標)・プロトコルに従う。通信インタフェース(114)は、通信コントローラ(113)を介してバス(104)に接続され、コンピュータ(101)を通信回線(115)に有線又は無線接続する役割を担い、コンピュータ(101)のオペレーティング・システムの通信機能のTCP/IP通信プロトコルに対して、ネットワーク・インタフェース層を提供する。通信回線は例えば、有線LAN接続規格に基づく有線LAN環境、又は無線LAN接続規格に基づく無線LAN環境、例えばIEEE802.11a/b/g/nなどのWi-Fi無線LAN環境、若しくは携帯電話網環境(例えば、3G、又は4G(LTEを含む)環境)でありうる。

10

【0067】

コンピュータ(101)は、通信回線(115)を介して例えば他の装置(例えば、コンピュータ又はネットワーク・アタッチト・ストレージ)からのデータを受信し、記憶装置(108)上に格納しうる。

【0068】

図2A及び図2Bはそれぞれ、本発明の実施態様に従い、粒度を規定する構成要素をコンテンツから抽出し、そして用語リスト中の用語を、当該抽出された構成要素がある位置で分割する例を示す。

20

【0069】

図2Aは、コンテンツ及び用語リストが英語である場合の例を示す。

【0070】

ユーザは、用語を分割する為の粒度を規定する構成要素を抽出する為のコンテンツ(201)及び、分割対象の用語を含む用語リスト(202)を用意する。コンテンツ(201)及び用語リスト(202)それぞれの内容は、図2Aに記載の通りである。コンテンツ(201)はコンピュータ分野の業務手順書であり、及び用語リスト(202)もまたコンピュータ分野の用語リストであるとする。

【0071】

コンピュータ(101)は、大別して、粒度を規定する構成要素をコンテンツから抽出する工程、及び、用語リスト中の用語を、当該抽出された構成要素がある位置で分割する工程を含む。

30

【0072】

(粒度を規定する構成要素をコンテンツから抽出する工程)

【0073】

コンピュータ(101)は、コンテンツ(201)及び用語リスト(202)を入力として受け取り、例えばメモリ(103)又は記憶装置(108)に格納する。

【0074】

コンピュータ(101)は、コンテンツ(201)から、構成要素を抽出する対象のテキストを切り出す。例えば、コンピュータ(101)は、コンテンツ(201)から、例えば変更履歴やコメント又は注釈を削除して、例えば本文のテキストを切り出しうる。

40

【0075】

コンピュータ(101)は、上記切り出したテキストを、事前定義した文字がある場所(すなわち、事前定義した文字がある場所の前後)で分割する。事前定義した文字は例えば、広義の句読点でありうる。広義の句読点は例えば、狭義の句読点(句点、読点)、疑問符、感嘆符、省略符、括弧(例えば、丸括弧、鉤括弧、角括弧、波括弧、亀甲括弧、山括弧、若しくは、隅付き括弧)、又は、その他文章に使う様々な記号を含みうる。コンピュータ(101)は例えば、コンテンツ(201)中の「- AAA」、「- BBB」及び「- ZUR」の記号 - の前後でテキストを分割する。

【0076】

50

コンピュータ(101)は、上記事前定義した文字で分割したテキストそれぞれに、当業者に知られている任意の構文解析技術を適用して、文節を抽出する。

【0077】

コンピュータ(101)は、上記抽出した文節のうちの名詞又は記号を含む文節から構成要素となりうる部分(すなわち、構成要素の候補)(203)を抽出する。名詞は、所謂文法上の名詞に分類される文字でありうる。記号は、自然言語処理において、辞書中に存在しない単語である未知語や省略語を含みうる。構成要素は、少なくとも1つの名詞又は記号を含む1又は複数の単語列でありうる。構成要素の候補(203)の内容は、図2Aに記載の通りである。構成要素の候補(203)中、「AAA」、「BBB」、「PPP」、「QQQ」及び「RRR」はいずれも、名詞又は記号を含む文節である。また、構成要素の候補(203)中、「ZUR」は固有名詞であり、「EOF」は「End-Of-File」の省略形である。

10

【0078】

コンピュータ(101)は、構成要素の候補(203)が、用語リスト(202)中にある用語を含むかを判断する。当該用語は、名詞、記号又はそれらの組み合わせを含む単語列でありうる。また、当該用語は例えば、複合名詞でありうる。コンピュータ(101)は、構成要素の候補(203)が用語リスト(202)中にある用語「ZUR EOF mark」を含む為に、構成要素の候補(203)から構成要素「ZUR EOF mark」を削除し、当該構成要素「ZUR EOF mark」を削除した残りを、粒度を規定する構成要素(204)とする。

【0079】

(用語リスト中の用語を、当該抽出された構成要素がある位置で分割する工程)

20

【0080】

コンピュータ(101)は、用語リスト(202)から一つの利用語を取り出し、当該取り出した用語が、当該用語の末尾から最長一致する構成要素(204)を含むかを判断する。コンピュータ(101)は、当該取り出した用語が、当該用語の末尾から最長一致する構成要素(204)を含むことに応じて、当該用語をその末尾から最長一致する構成要素がある位置で分割する。そして、コンピュータ(101)は、上記末尾から最長一致する構成要素を分割した用語を、分割した後の用語を入れるリストL(205)に格納する。

【0081】

コンピュータ(101)は、用語リスト(202)中の「Beneficiary right seller's business security deposit」を取り出し、当該「Beneficiary right seller's business security deposit」が、その末尾から最長一致する構成要素(204)を含むかを判断する。コンピュータ(101)は、当該「Beneficiary right seller's business security deposit」が、その末尾から最長一致する構成要素「business security deposit」を含んでいることに応じて、「Beneficiary right seller's business security deposit」をその末尾から最長一致する構成要素「business security deposit」がある位置(すなわち、「business security deposit」の直前)で分割し、すなわち、「Beneficiary right seller's」と「business security deposit」とに分割する。そして、コンピュータ(101)は、上記末尾から最長一致する構成要素を分割した用語「business security deposit」を、上記リストL(205)に格納する。従って、「Beneficiary right seller's business security deposit」は、コンテンツ(201)中の構成要素「business security deposit」の粒度に従って分割されている。

30

40

【0082】

同様に、コンピュータ(101)は、用語リスト(202)中の「Financial instruments intermediary service」を取り出し、当該「Financial instruments intermediary service」が、その末尾から最長一致する構成要素(204)を含むかを判断する。コンピュータ(101)は、当該「Financial instruments intermediary service」が、その末尾から最長一致する構成要素「intermediary service」を含んでいることに応じて、「Financial instruments intermediary service」をその末尾から最長一致する構成要素「intermediary service」がある位置(すなわち、「intermediary service」の直前)で分割し

50

、すなわち、「Financial instruments」と「intermediary service」とに分割する。そして、コンピュータ（101）は、上記末尾から最長一致する構成要素を分割した用語「intermediary service」を、上記リストL（205）に格納する。従って、「Financial instruments intermediary service」は、コンテンツ（201）中の構成要素「intermediary service」の粒度に従って分割されている。

【0083】

同様に、コンピュータ（101）は、用語リスト（202）中の「ZUR EOF mark」を取り出し、当該「ZUR EOF mark」が、その末尾から最長一致する構成要素（204）を含むか判断する。コンピュータ（101）は、当該「ZUR EOF mark」が、その末尾から最長一致する構成要素「mark」を含んでいることに応じて、「ZUR EOF mark」をその末尾から最長一致する構成要素「mark」がある位置（すなわち、「mark」の直前）で分割し、すなわち、「ZUR EOF」と「mark」とに分割する。そして、コンピュータ（101）は、上記末尾から最長一致する構成要素を分割した用語「mark」を、上記リストL（205）に格納する。従って、「ZUR EOF mark」は、コンテンツ（201）中の構成要素「mark」の粒度に従って分割されている。

10

【0084】

次に、コンピュータ（101）は、上記取り出した用語から上記末尾から最長一致する構成要素を分割した用語を除いた後の用語が、当該除いた後の用語の先頭から最長一致する構成要素（204）を含むかを判断する。コンピュータ（101）は、上記取り出した用語から上記末尾から最長一致する構成要素を分割した用語を除いた後の用語が当該除いた後の用語の先頭から最長一致する構成要素を含むことに応じて、当該除いた後の用語をその先頭から最長一致する構成要素がある位置で分割する。そして、コンピュータ（101）は、上記先頭から最長一致する構成要素を分割した用語を、上記リストL（205）に格納する。

20

【0085】

コンピュータ（101）は、用語リスト（202）中の「Beneficiary right seller's business security deposit」中の「business security deposit」を除いた用語「Beneficiary right seller's」がその先頭から最長一致する構成要素（204）を含むか判断する。コンピュータ（101）は、当該「business security deposit」がその先頭から最長一致する構成要素（204）を含まないことに応じて、当該「business security deposit」を、上記リストL（205）に格納して、分割処理を終了する。

30

【0086】

同様に、コンピュータ（101）は、用語リスト（202）中の「Financial instruments intermediary service」中の「intermediary service」を除いた用語「Financial instruments」がその先頭から最長一致する構成要素（204）を含むか判断する。コンピュータ（101）は、当該「Financial instruments」がその先頭から最長一致する構成要素（204）を含まないことに応じて、当該「Financial instruments」をそのまま、上記リストL（205）に格納して、分割処理を終了する。

【0087】

同様に、コンピュータ（101）は、用語リスト（202）中の「ZUR EOF mark」中の「mark」を除いた用語「ZUR EOF」がその先頭から最長一致する構成要素（204）を含むか判断する。当該除いた用語「ZUR EOF」はその先頭から最長一致する構成要素「ZUR EOF」を含んでいる。しかしながら、両者は同一であるから分割出来ない。従って、コンピュータ（101）は、当該「ZUR EOF」をそのまま、上記リストL（205）に格納して、分割処理を終了する。

40

【0088】

従って、上記分割処理後のリストL（205）は、「Beneficiary right seller's」、「business security deposit」、「Financial instruments」、「intermediary service」、「ZUR EOF」、及び「mark」を含む。

【0089】

50

また、コンピュータ(101)は、上記分割処理後のリストL(205)を、図2Aに示すように、末尾から最長一致する最初の構成要素を主要語として、及び、先頭から最長一致する最初の構成要素を修飾語1として、並びに、次に、末尾から最長一致する構成要素を修飾語2(ある場合)として、表示装置(106)上に表示しうる。

【0090】

代替的には、コンピュータ(101)は、上記分割処理後のリストL(205)を、下記のように、分割した箇所を示す記号、例えば、|を入れて、表示装置(106)上に表示しうる。例えば、以下の通りである。

Beneficiary right seller's | business security deposit
Financial instruments | intermediary service
ZUR EOF | mark

10

【0091】

用語リスト(202)中の用語「Beneficiary right seller's business security deposit」、「Financial instruments intermediary service」、及び「ZUR EOF mark」それぞれは、従来技術に従う形態素解析器に従うと、単語は通常空白によってわかち書きされている為に、単語ごとに分割される。一方、本願発明の態様に従うと、上記分割処理後のリストL(205)に示すように、コンテンツ(201)から抽出された構成要素(204)の粒度に従い用語が分割される。

【0092】

図2Bは、コンテンツ及び用語リストが日本語(漢字を含む)である場合の例を示す。

20

【0093】

ユーザは、用語を分割する為の粒度を規定する構成要素を抽出する為のコンテンツ(211)及び、分割対象の用語を含む用語リスト(212)を用意する。コンテンツ(211)及び用語リスト(212)それぞれの内容は、図2Bに記載の通りである。コンテンツ(211)は金融分野の業務手順書であり、及び用語リスト(212)もまた金融分野の用語リストであるとする。

【0094】

コンピュータ(101)は、大別して、粒度を規定する構成要素をコンテンツから抽出する工程、及び、用語リスト中の用語を、当該抽出された構成要素がある位置で分割する工程を含む。

30

【0095】

(粒度を規定する構成要素をコンテンツから抽出する工程)

【0096】

コンピュータ(101)は、コンテンツ(211)及び用語リスト(212)を入力として受け取り、例えばメモリ(103)又は記憶装置(108)に格納する。

【0097】

コンピュータ(101)は、コンテンツ(211)から、構成要素を抽出する対象のテキストを切り出す。当該テキストの切り出しとは、上記において述べた通りである。

【0098】

コンピュータ(101)は、上記切り出したテキストを、事前定義した文字がある場所(すなわち、事前定義した文字がある場所の前後)で分割する。事前定義した文字とは、上記において述べた通りである。コンピュータ(101)は例えば、コンテンツ(211)中の「(金信)」の括弧書き (の前後、及び括弧書き) の前後でテキストを分割する。

40

【0099】

コンピュータ(101)は、上記事前定義した文字で分割したテキストそれぞれに、当業者に知られている任意の構文解析技術を適用して、文節を抽出する。

【0100】

コンピュータ(101)は、上記抽出した文節のうちの名詞又は記号を含む文節から構成要素となりうる部分(すなわち、構成要素の候補)(213)を抽出する。構成要素の

50

候補(213)の内容は、図2Bに記載の通りである。

【0101】

コンピュータ(101)は、構成要素の候補(213)が、用語リスト(212)中にある用語を含むかを判断する。コンピュータ(101)は、構成要素の候補(213)が用語リスト(212)中にある用語「金信期日後収益金税額」、「延滞元金額」及び「補正計算元本額」を含む為に、構成要素の候補(213)から構成要素「金信期日後収益金税額」、「延滞元金額」及び「補正計算元本額」を削除し、これらの構成要素を削除した残りを、粒度を規定する構成要素(214)とする。

【0102】

(用語リスト中の用語を、当該抽出された構成要素がある位置で分割する工程)

10

【0103】

コンピュータ(101)は、用語リスト(212)から一つの利用語を取り出し、当該取り出した用語が、当該用語の末尾から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該取り出した用語が、当該用語の末尾から最長一致する構成要素(214)を含むことに応じて、当該用語をその末尾から最長一致する構成要素がある位置で分割する。そして、コンピュータ(101)は、上記末尾から最長一致する構成要素を分割した用語を、分割した後の用語を入れるリストL(215)に格納する。

【0104】

コンピュータ(101)は、用語リスト(212)中の「金信期日後収益金税額」を取り出し、当該「金信期日後収益金税額」が、その末尾から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該「金信期日後収益金税額」が、その末尾から最長一致する構成要素「税額」を含んでいることに応じて、「金信期日後収益金税額」をその末尾から最長一致する構成要素「税額」がある位置(すなわち、「税額」の直前)で分割し、すなわち、「金信期日後収益金」と「税額」とに分割する。そして、コンピュータ(101)は、上記末尾から最長一致する構成要素を分割した用語「税額」を、上記リストL(215)に格納する。従って、「金信期日後収益金税額」は、コンテンツ(211)中の構成要素「税額」の粒度に従って分割されている。

20

【0105】

同様に、コンピュータ(101)は、用語リスト(212)中の「延滞元金額」を取り出し、当該「延滞元金額」が、その末尾から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該「延滞元金額」が、その末尾から最長一致する構成要素「元金額」を含んでいることに応じて、「延滞元金額」をその末尾から最長一致する構成要素「元金額」がある位置(すなわち、「元金額」の直前)で分割し、すなわち、「延滞」と「元金額」とに分割する。そして、コンピュータ(101)は、上記末尾から最長一致する構成要素を分割した用語「元金額」を、上記リストL(215)に格納する。従って、「延滞元金額」は、コンテンツ(211)中の構成要素「元金額」の粒度に従って分割されている。

30

【0106】

同様に、コンピュータ(101)は、用語リスト(212)中の「補正計算元本額」を取り出し、当該「補正計算元本額」が、その末尾から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該「補正計算元本額」が、その末尾から最長一致する構成要素「元本額」を含んでいることに応じて、「補正計算元本額」をその末尾から最長一致する構成要素「元本額」がある位置(すなわち、「元本額」の直前)で分割し、すなわち、「補正計算」と「元本額」とに分割する。そして、コンピュータ(101)は、上記末尾から最長一致する構成要素を分割した用語「元本額」を、上記リストL(215)に格納する。従って、「補正計算元本額」は、コンテンツ(211)中の構成要素「元本額」の粒度に従って分割されている。

40

【0107】

次に、コンピュータ(101)は、上記取り出した用語から上記末尾から最長一致する

50

構成要素を分割した用語を除いた後の用語が、当該除いた後の用語の先頭から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、上記取り出した用語から上記末尾から最長一致する構成要素を分割した用語を除いた後の用語が当該除いた後の用語の先頭から最長一致する上記構成要素を含むことに応じて、当該除いた後の用語をその先頭から最長一致する構成要素がある位置で分割する。そして、コンピュータ(101)は、上記先頭から最長一致する構成要素を分割した用語を、上記リストL(215)に格納する。

【0108】

コンピュータ(101)は、用語リスト(212)中の「金信期日後収益金税額」中の「税額」を除いた用語「金信期日後収益金」の先頭から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該「金信期日後収益金」がその先頭から最長一致する構成要素「金信」を含んでいることに応じて、「金信期日後収益金」をその先頭から最長一致する構成要素「金信」がある位置(すなわち、「金信」の直後)で分割し、すなわち、「金信」と「期日後収益金」とに分割する。そして、コンピュータ(101)は、上記先頭から最長一致する構成要素を分割した用語「金信」を、上記リストL(215)に格納する。従って、「金信期日後収益金」は、コンテンツ(211)中の構成要素「金信」の粒度に従って分割されている。

10

【0109】

同様に、コンピュータ(101)は、用語リスト(212)中の「延滞元金額」中の「元金額」を除いた用語「延滞」の先頭から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該「延滞」がその先頭から最長一致する構成要素(214)を含まないことに応じて、当該「延滞」をそのまま、上記リストL(215)に格納して、分割処理を終了する。

20

【0110】

同様に、コンピュータ(101)は、用語リスト(212)中の「補正計算元本額」中の「元本額」を除いた用語「補正計算」の先頭から最長一致する構成要素(214)を含むかを判断する。コンピュータ(101)は、当該「補正計算」がその先頭から最長一致する構成要素「補正計算」を含むが、両者は同一である為に分割できない。そこで、コンピュータ(101)は、当該「補正計算」をそのまま、上記リストL(215)に格納して、分割処理を終了する。

30

【0111】

次に、コンピュータ(101)は、上記「金信期日後収益金税額」のうち、残った用語「期日後収益金」が、当該用語の末尾から最長一致する構成要素(214)を含むかを判断する。当該残った用語「期日後収益金」はその先頭から最長一致する構成要素「期日後収益金」を含んでいる。しかしながら、両者は同一であるから分割出来ない。従って、コンピュータ(101)は、当該「期日後収益金」をそのまま、上記リストL(215)に格納して、分割処理を終了する。

【0112】

従って、上記分割処理後のリストL(215)は、「金信」、「期日後収益金」、「税額」、「延滞」、「元金額」、「補正計算」、及び「元本額」を含む。

40

【0113】

また、コンピュータ(101)は、上記分割処理後のリストL(215)を、図2Bに示すように、末尾から最長一致する最初の構成要素を主要語として、及び、先頭から最長一致する最初の構成要素を修飾語1として、並びに、次に、末尾から最長一致する構成要素を修飾語2(ある場合)として、表示装置(106)上に表示しうる。

【0114】

代替的には、コンピュータ(101)は、上記分割処理後のリストL(215)を、下記のように、分割した箇所を示す記号、例えば、|を入れて、表示装置(106)上に表示しうる。例えば、以下の通りである。

金信 | 期日後収益金 | 税額

50

延滞 | 元金額
 補正計算 | 元本額

【0115】

用語リスト(212)中の用語「金信期日後収益金税額」、「延滞元金額」、及び「補正計算元本額」それぞれは、従来技術に従う形態素解析器に従うと、例えば「金信期日後収益金税額」、「延滞元金額」、及び「補正計算元本額」というように、従来技術に従う形態素解析器が持つシステム辞書及び文法に基づいて複合名詞が分割される。従って、必ずしも、望ましい箇所での用語が分割されていない。一方、本願発明の態様に従うと、上記分割処理後のリストL(215)に示すように、コンテンツ(211)から抽出された構成要素(214)の粒度に従い用語が分割される。

10

【0116】

図2A及び図2Bそれぞれに示す本願発明の実施態様に従い、コンテンツから抽出された構成要素の粒度に従い用語を分割することによって得られた語は、下記(1)~(3)それぞれに示す場面において利用されうる。

【0117】

(1)主要語又は例えば、修飾語1若しくは修飾語2で、分割された語をソートする。当該語をソートをすることによって、ユーザは、例えば、統一すべき類似の又は同義の語を発見したり、追加すべき語を洗い出したり、又は短い主要語のみからなる語を発見したりすることが用意になる。

20

【0118】

(2)主要語の抽象度が高い場合、ユーザは、当該主要語を複数の具体的な用語に置き換えることが容易になる。抽象度が高い主要語は例えば、「実行金額」である。

【0119】

(3)ユーザは、主要語が非常に長い用語を発見することが容易になる。このことは、ユーザが、当該用語が仕様書に説明がない用語だと判定することを容易にしたり、又は、関係者間で共通の理解を持つ為に、当該用語の説明の為の記述を仕様書に追加することを可能にする。例えば、上記分割の結果、コンテンツ(211)から抽出された構成要素(214)中に該当する構成要素がなく、「分割可能上限金額」及び「利息算出対象元金」が主要語として得られたとする。このような場合には、コンピュータ(101)は、当該長い主要語を持つ用語「分配可能上限金額」及び「利息算出対象元金」として、抽出する。

30

【0120】

図3A及び図3Bは、本発明の実施態様に従い、粒度を規定する構成要素をコンテンツから抽出し、そして用語リスト中の用語を、当該抽出された構成要素がある位置で分割する処理の為にフローチャートを示す。

【0121】

図3Aは、本発明の実施態様に従い、粒度を規定する構成要素をコンテンツから抽出する処理の為にフローチャートを示す。

【0122】

ステップ301において、コンピュータ(101)は、コンテンツから上記構成要素を抽出する処理を開始する。

40

【0123】

ステップ302において、コンピュータ(101)は、粒度を規定する構成要素をそこから抽出する為にコンテンツを、コンテンツを記録した記録媒体(331)から読み取る。コンテンツは例えば、索引を作成することが必要な文書、例えばビジネス文書(例えば、仕様書、業務手順書、ビジネスツール定義書)でありうるがこれらに限定されるものではない。コンテンツが例えば、仕様書である場合には、例えばデータ項目やビジネス・プロセスが説明とともに記載されている。コンテンツが属する技術分野や適用分野が異なると、それから抽出される構成要素も異なってくる。すなわち、コンテンツが変われば、用語を分割する為に粒度も変わってくる。

50

【 0 1 2 4 】

引き続き、ステップ 3 0 2 において、コンピュータ (1 0 1) は、任意的に、当該読み取ったコンテンツから、上記構成要素を抽出する対象のテキストを切り出しうる。例えば、コンピュータ (1 0 1) は、コンテンツ (3 3 1) から、例えば変更履歴やコメント又は注釈を削除して、例えば本文のテキストを切り出しうる。

【 0 1 2 5 】

ステップ 3 0 3 において、コンピュータ (1 0 1) は、ステップ 3 0 2 で切り出したテキストを、事前定義した文字がある場所 (すなわち、事前定義した文字がある場所の前後) で分割する。事前定義した文字は例えば、広義の句読点でありうる。広義の句読点は例えば、狭義の句読点 (句点、読点)、疑問符、感嘆符、省略符、括弧 (例えば、丸括弧、鉤括弧、角括弧、波括弧、亀甲括弧、山括弧、若しくは、隅付き括弧)、又は、その他文章に使う様々な記号を含みうる。例えば、テキストが、「Confirm the item number (check whether the order is valid) in this process.」 (英語である) である場合には、コンピュータ (1 0 1) は、「Confirm the item number | (| check whether the order is valid |) | in this process」 (| は、分割する位置を示す) と、丸括弧の前後で分割する。同様に、例えば、テキストが、「品物番号を確認 (オーダーと同じか) する」 (日本語である) である場合には、コンピュータ (1 0 1) は、「品物番号を確認 | (| オーダーと同じか |) | する」 (日本語である) (| は、分割する位置を示す) と、丸括弧の前後で分割する。

【 0 1 2 6 】

ステップ 3 0 4 において、コンピュータ (1 0 1) は、ステップ 3 0 3 で分割したテキストそれぞれに、当業者に知られている任意の構文解析技術を適用して、文節を抽出する。例えば、テキストが、「In the calculation / of business security deposit / , PPP / , QQQ / and RRR / are used.」 (英語である) である場合には、コンピュータ (1 0 1) は、「In the calculation / of business security deposit / , PPP / , QQQ / and RRR / are used.」 (/ は、文節を抽出する区切りを示す) として文節を抽出する。同様に、例えば、テキストが、「中途解約受取金額は預金期間から決まる」 (日本語である) である場合には、コンピュータ (1 0 1) は、「中途解約受取金額は / 預金期間から / 決まる」 (日本語である) (/ は、文節を抽出する区切りを示す) として文節を抽出する。

【 0 1 2 7 】

ステップ 3 0 5 において、コンピュータ (1 0 1) は、ステップ 3 0 4 で抽出した文節のうち名詞又は記号を含む文節から構成要素となりうる部分を抽出する。コンピュータ (1 0 1) は、当該抽出した部分を、構成要素の候補としてリストしうる。名詞は、所謂文法上の名詞に分類される文字でありうる。記号は、自然言語処理において、辞書中に存在しない単語である未知語や省略語を含みうる。構成要素は、少なくとも 1 つの名詞又は記号を含む 1 又は複数の単語列でありうる。構成要素となりうる部分を抽出することは例えば、英語の場合、冠詞を除くことを含む。また、構成要素となりうる部分を抽出することは、当該抽出された部分について、複数形を単数形に変換したり、大文字を小文字に変形したり、旧字体を新字体に変形したりするような変形処理をすることを含みうる。

【 0 1 2 8 】

ステップ 3 0 6 において、コンピュータ (1 0 1) は、任意的に、ステップ 3 0 5 で抽出した部分に、当該部分の直前の文字列を補完するかどうかを判断する。例えば、文字列が専門用語の場合には、任意の構文解析技術を使用すると、一つの単語が複数に分割される場合がある。そこで、ステップ 3 0 5 で抽出した部分に、当該部分の直前の文字列を補完することによって、コンピュータ (1 0 1) は、複数に分割された一つの用語を、本来の一つの用語になるように直前の文字列で補完しうる。例えば、上記抽出した部分が「File:Open / Menu」 (英語である) (/ は、区切りを示す) である場合、「File:Open」はスペースがなく且つ記号「:」が挿入されているために、当該「File:Open」は未知語として検出され、及び、「Menu」は名詞として別々に検出される。従って、上記抽出した部分

10

20

30

40

50

「File:Open / Menu」は、「Menu」の直前に「File:Open」を補って、「File:Open Menu」と本来の一つのまとまりのある用語にする。例えば、上記抽出した部分が、「採 / 番」（漢字である）（ / は、区切りを示す）である場合には、「番」の直前に「採」を補って、「採番」（漢字である；「採番」とは、データ管理のために、それぞれのデータに固有の番号を与えることを意味する）と本来の一つの用語にする。コンピュータ（101）は、上記直前の文字列を補完することに応じて、処理をステップ307に進める。一方、コンピュータ（101）は、上記直前の文字列を補完しないことに応じて、処理をステップ308に進める。

【0129】

ステップ307において、コンピュータ（101）は、ステップ305で抽出した部分に、当該部分の直前の文字列を補完する。そして、コンピュータ（101）は、当該補完した文字列で、構成要素の候補の上記リストを更新しうる。引き続き、コンピュータ（101）は、ステップ306に戻り、さらに補完する必要があるかどうか判断しうる。

【0130】

ステップ308において、コンピュータ（101）は、用語を例えばメモリ（103）中に読み取る。用語は、名詞、記号（未知語や省略語を含む）又はそれらの組み合わせを含む単語列でありうる。当該用語の読み取りは、例えば、用語を格納した用語リストを、当該用語リストを記録した記録媒体（332）から読み取ることによって行われうる。また、代替的には、コンピュータ（101）は、ユーザによって指定された用語を入力として読み取りうる。さらに、代替的には、コンピュータ（101）は、コンテンツ（331）中の所定の長さよりも長い（例えば、構成要素の平均文字長よりも長い、又は例えば、10文字よりも長い）用語を、分割対象の用語として読み取りうる。引き続き、コンピュータ（101）は、当該用語リスト中に、ステップ305で抽出した、構成要素となりうる部分と同じ用語があるかを判断する。用語リスト中に上記構成要素となりうる部分と同じ用語がある場合には、当該構成要素となりうる部分は、用語リスト中の用語を分割する為の構成要素となり得ないからである。コンピュータ（101）は、用語リスト中に上記構成要素となりうる部分と同じ用語がある場合には、処理をステップ309に進める。一方、コンピュータ（101）は、用語リスト中に上記構成要素となりうる部分と同じ用語がない場合には、処理をステップ310に進める。

【0131】

ステップ309において、コンピュータ（101）は、用語リスト中に上記構成要素となりうる部分と同じ用語がある場合には、当該同じ用語を、上記構成要素の候補の上記リストから削除する。用語リストは例えば、データ項目やビジネス・プロセスが記載されている。用語リストは例えば、 $Kwds = \{k_1, k_2, k_3, \dots, k_n\}$ で表されうる。

【0132】

ステップ310において、コンピュータ（101）は、上記構成要素の候補の上記リストを、コンテンツから粒度を規定する構成要素として、例えば構成要素を格納する記録媒体（333）に格納する。

【0133】

ステップ311において、コンピュータ（101）は、コンテンツから上記構成要素を抽出する処理を終了する。

【0134】

図3Bは、本発明の実施態様に従い、用語を構成要素がある位置で分割する処理の為のフローチャートを示す。

【0135】

ステップ321において、コンピュータ（101）は、用語を、ステップ310で作成した構成要素がある位置で分割する処理を開始する。

【0136】

ステップ322において、コンピュータ（101）は、用語を例えば用語リスト（332）から一つ取り出す。そして、コンピュータ（101）は、当該取り出した一つの利用語

10

20

30

40

50

が、当該用語の末尾から最長一致する構成要素（記憶媒体（333）に格納されている）を含むかを判断する。コンピュータ（101）は、当該取り出した用語が、当該用語の末尾から最長一致する構成要素（204）を含むことに応じて、当該用語をその末尾から最長一致する構成要素がある位置で分割する。そして、コンピュータ（101）は、上記末尾から最長一致する構成要素を分割した用語を、分割した後の用語を入れるリストL（334）に格納する。コンピュータ（101）は、上記末尾から最長一致する構成要素を分割した用語を、主要語としてリストL（334）に格納しうる。

【0137】

ステップ323において、コンピュータ（101）は、ステップ322での上記分割ができたことに応じて、用語の分割回数dを1つ増加し（d++）、処理をステップ324に進める。一方、コンピュータ（101）は、ステップ322での上記分割ができなかったことに応じて、処理をステップ328に進める。

10

【0138】

ステップ324において、コンピュータ（101）は、分割回数dと用語の分割回数を規定する分割パラメータQとを比較する。コンピュータ（101）は、分割回数dが分割パラメータQよりも少ないことに応じて、さらに分割処理をする為に、処理をステップ325に進める。一方、コンピュータ（101）は、分割回数dが分割パラメータQ以上であることに応じて、これ以上分割処理をしない為に、処理をステップ328に進める。

【0139】

ステップ325において、コンピュータ（101）は、ステップ322で取り出した用語から上記末尾から最長一致する構成要素を分割した用語を除いた後の用語が、当該除いた後の用語の先頭から最長一致する構成要素を含むかを判断する。コンピュータ（101）は、上記取り出した用語から上記末尾から最長一致する構成要素を分割した用語を除いた後の用語が当該除いた後の用語の先頭から最長一致する構成要素を含むことに応じて、当該除いた後の用語をその先頭から最長一致する構成要素がある位置で分割する。そして、コンピュータ（101）は、上記先頭から最長一致する構成要素を分割した用語を、分割した後の用語を入れるリストL（334）に格納する。コンピュータ（101）は、上記先頭から最長一致する構成要素を分割した用語を、修飾語1としてリストL（334）に格納しうる。

20

【0140】

ステップ326において、コンピュータ（101）は、ステップ325での上記分割ができたことに応じて、用語の分割回数dを1つ増加し（d++）、処理をステップ327に進める。一方、コンピュータ（101）は、ステップ325での上記分割ができなかったことに応じて、処理をステップ328に進める。

30

【0141】

ステップ327において、コンピュータ（101）は、分割回数dと用語の分割回数を規定する分割パラメータQとを比較する。コンピュータ（101）は、分割回数dが分割パラメータQよりも少ないことに応じて、さらに分割処理をする為に、処理をステップ322に戻す。コンピュータ（101）は、引き続き、末尾から最長一致する構成要素があるか、そして、先頭が最長一致する構成要素があるかを繰り返して行いうる。一方、コンピュータ（101）は、分割回数dが分割パラメータQ以上であることに応じて、これ以上分割処理をしない為に、処理をステップ328に進める。

40

【0142】

ステップ328において、コンピュータ（101）は、分割した語と、分割した語の残りがあある場合には当該残りとを分割後の語を格納する用語リストL（334）に格納する。

【0143】

ステップ329において、コンピュータ（101）は、任意的に、用語リストLの内容（上記分割した語と、分割した語の残りがあある場合には当該残り分割後の語）を、例えば表示装置（106）上に表示しうる。

50

【 0 1 4 4 】

ステップ 3 3 0 において、コンピュータ (1 0 1) は、用語を、上記構成要素がある位置で分割する処理を終了する。

【 0 1 4 5 】

図 4 は、図 1 に従うハードウェア構成を好ましくは備えており、図 3 A 及び図 3 B それぞれに示すフローチャートに従って本発明の実施態様を実施するコンピュータの機能ブロック図の一例を示す図である。

【 0 1 4 6 】

コンピュータ (4 0 1) は、図 1 A に示すコンピュータ (1 0 1) に示されている構成、例えば CPU (1 0 2)、メイン・メモリ (1 0 3)、記憶装置 (1 0 8)、及びディスク (1 0 8) を備えている。

10

【 0 1 4 7 】

コンピュータ (4 0 1) は、抽出手段 (4 1 1)、分割手段 (4 1 2)、及び表示手段 (4 1 3) を備えている。

【 0 1 4 8 】

抽出手段 (4 1 1) は、構文解析により、粒度を規定する構成要素をコンテンツから抽出する。

【 0 1 4 9 】

また、抽出手段 (4 1 1) は、上記コンテンツ中のテキストそれぞれに上記構文解析を適用して、文節を抽出し、上記抽出した文節のうちの名詞又は記号を含む文節から上記構成要素となりうる部分を抽出する。

20

【 0 1 5 0 】

また、抽出手段 (4 1 1) は、上記コンテンツから、上記構成要素を抽出する対象のテキストを切り出し、当該切り出したテキストそれぞれに上記構文解析を適用して、上記文節を抽出する。

【 0 1 5 1 】

また、抽出手段 (4 1 1) は、上記切り出したテキストを事前定義した文字がある場所で分割し、当該分割したテキストそれぞれに上記構文解析を適用して、上記文節を抽出する。

【 0 1 5 2 】

また、抽出手段 (4 1 1) は、上記用語が用語リスト中の用語である場合に、上記構成要素となりうる部分のうちから上記用語リスト中にある用語を削除し、当該削除した残りを上記構成要素とする。

30

【 0 1 5 3 】

抽出手段 (4 1 1) は、図 3 A に記載の各ステップを実行しうる。

【 0 1 5 4 】

分割手段 (4 1 2) は、上記用語がその一部に少なくとも 1 つの上記構成要素を含む場合に、上記用語を上記構成要素がある位置で分割する。

【 0 1 5 5 】

また、分割手段 (4 1 2) は、上記用語が当該用語の末尾から最長一致する上記構成要素 (第 1 の構成要素) を含む場合に、上記用語を上記末尾から最長一致する上記構成要素 (第 1 の構成要素) がある位置で分割する。

40

【 0 1 5 6 】

また、分割手段 (4 1 2) は、上記用語を上記末尾から最長一致する上記構成要素 (第 1 の構成要素) がある位置で分割し、上記末尾から最長一致する上記構成要素 (第 1 の構成要素) を上記用語の主要語として保存する。

【 0 1 5 7 】

また、分割手段 (4 1 2) は、上記用語から上記末尾から最長一致する上記構成要素 (第 1 の構成要素) を除いた後の用語が当該除いた後の用語の先頭から最長一致する上記構成要素 (第 2 の構成要素) を含む場合に、上記除いた後の用語を上記先頭から最長一致す

50

る上記構成要素（第2の構成要素）がある位置で分割する。

【0158】

また、分割手段（412）は、上記除いた後の用語を上記先頭から最長一致する上記構成要素（第2の構成要素）がある位置で分割し、上記先頭から最長一致する上記構成要素（第2の構成要素）を上記用語の第1の修飾語として保存する。

【0159】

また、分割手段（412）は、上記除いた後の用語を上記先頭から最長一致する上記構成要素（第2の構成要素）がある位置で分割し、上記先頭から最長一致する上記構成要素（第2の構成要素）以外の部分を第2の修飾語として保存する。

【0160】

また、分割手段（412）は、予め設定された分割回数を規定する分割パラメータに従って、上記用語を上記構成要素がある位置で分割する。

【0161】

分割手段（412）は、図3Bに記載のステップ322～328を実行しうる。

【0162】

表示手段（413）は、用語リストLの内容を、例えば表示装置（106）上に表示する。

【0163】

表示手段（413）は、図3Bに記載のステップ329を実行しうる。

【0164】

本発明の実施態様に従うと、上記した通り、コンテンツから抽出される構成要素に従う粒度で、用語を分割することが可能である。従って、コンテンツが属する技術分野や適用分野が異なると、それから抽出される構成要素も異なってくる為に、用語を分割する為の構成要素の粒度も異なってくる。このようにして分割された語は、下記（1）及び（2）それぞれに示す場面において利用されうる。

【0165】

（1）例えば、システムの大規模改修やシステム統合により新システムに移行するという場面では、例えば以前のシステムの設計で用いられていた用語を見直す必要がある。例えば、新システムで使わない用語を削除したり、新システムで新規に検討しなくてはならない用語（例えば、新しいビジネス・プロセス関連の用語）を洗い出したり、曖昧に使用されていた用語を見直したり、同じ意味で異なる用語が異なるシステム間で使用されている場合に、当該異なる用語を統一したりする必要がある。このような場合に、本発明の実施態様に従い、コンテンツから抽出した構成要素で用語を分割することによって当該用語が適切な粒度で分割される為に当該分割された結果に基づいて上記用語の見直しが可能となる点で、本発明の実施態様に従う上記分割は有用である。

【0166】

（2）用語の見直しをするに際して、当該用語がデータ項目の場合には、当該用語がデータベースのカラム名になっていたり、プログラム中の変数となっていたりする為に、構造ルールを設けるがある。このような場合に、本発明の実施態様に従いコンテンツから抽出した構成要素で用語を分割することによって当該用語が適切な粒度で分割される為に当該分割された結果に基づいて、上記構造ルールを設けることが可能となる、当該構造ルールを設けることで、用語の構造を共通理解でき、さらには用語が意味する概念についての理解が一意に決まることが期待でき、また、不明確な用語を作成する必要がないという点で、本発明の実施態様に従う上記分割は有用である。

10

20

30

40

【 図 1 】

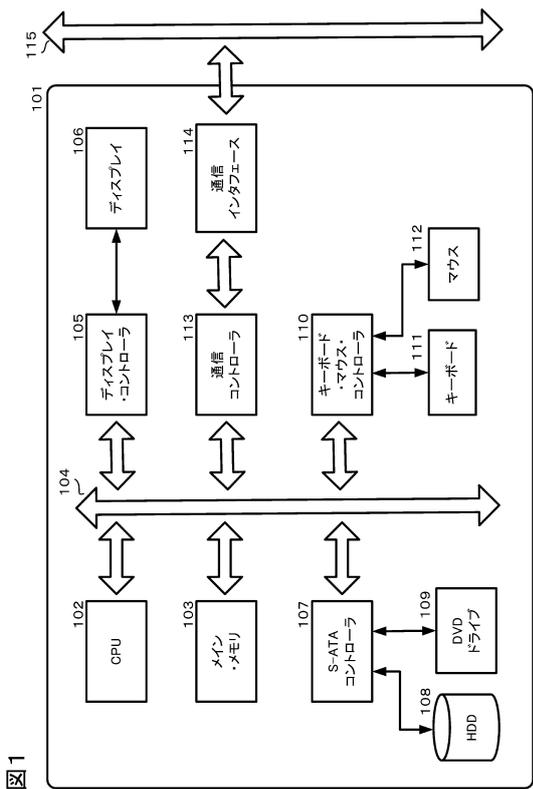


図 1

【 図 2 A 】

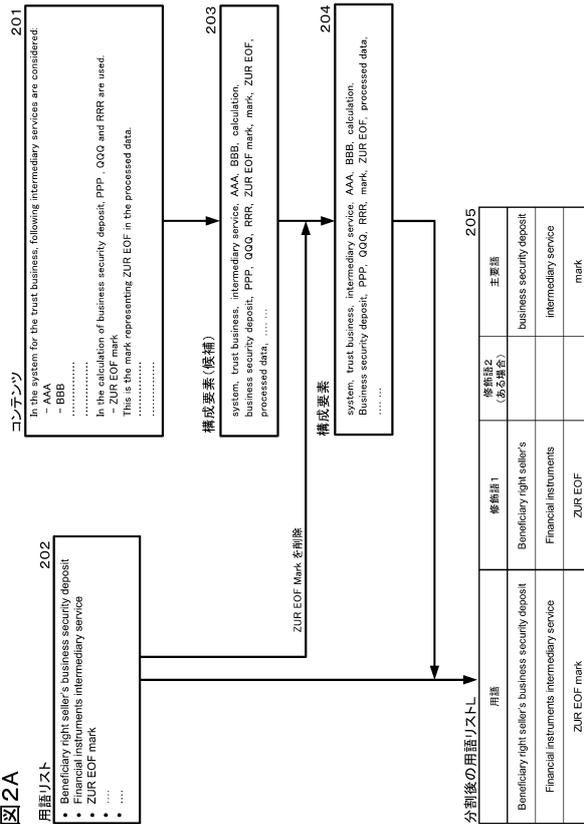


図 2A

【 図 2 B 】

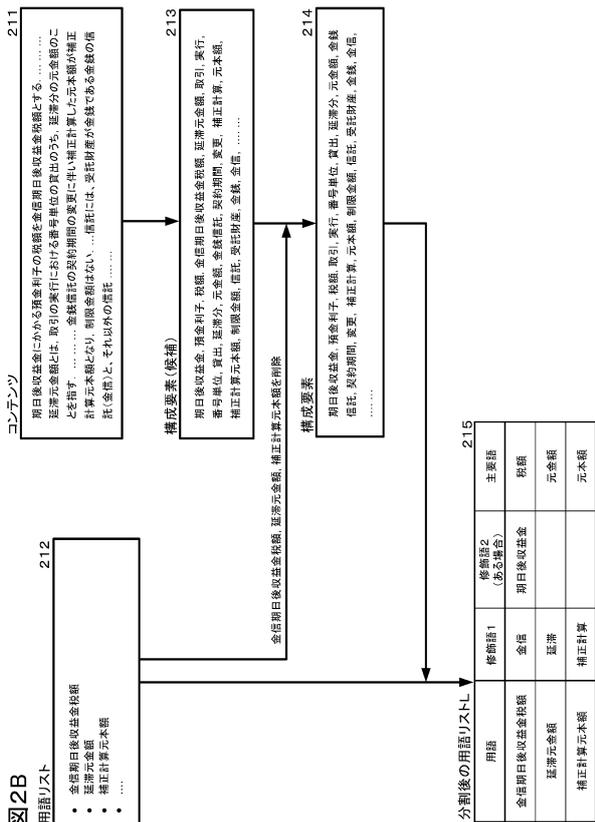


図 2B

【 図 3 A 】

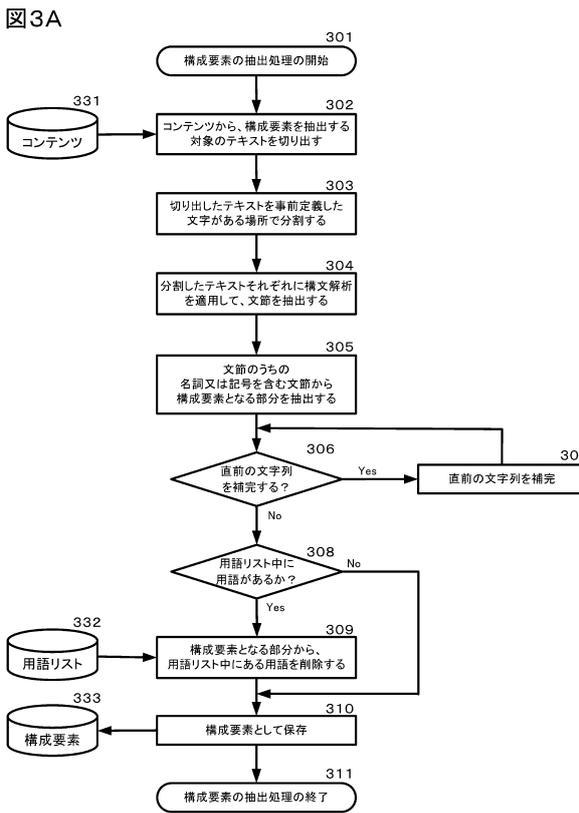
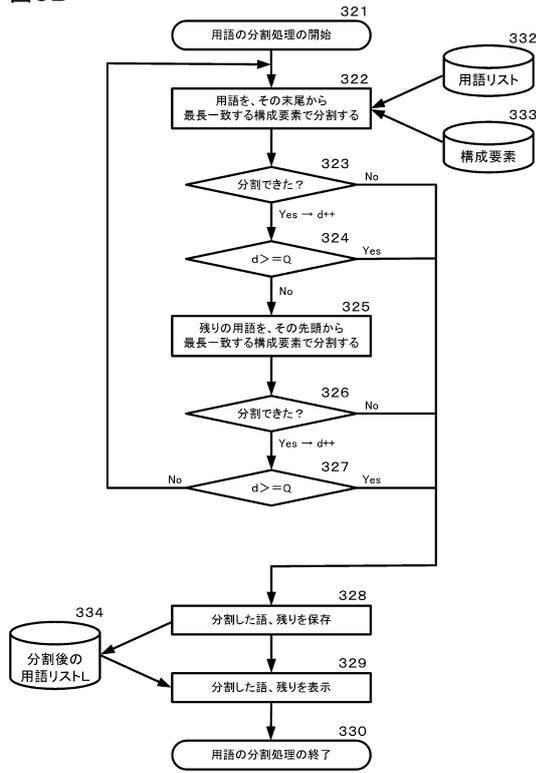


図 3A

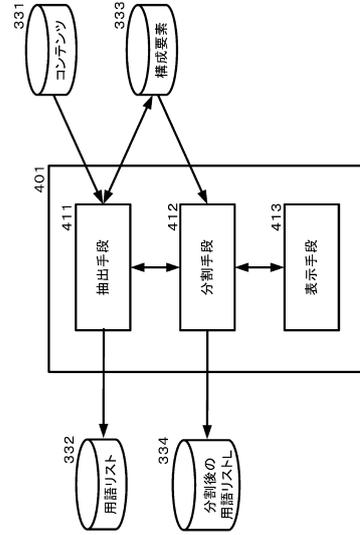
【図3B】

図3B



【図4】

図4



フロントページの続き

- (74)代理人 100112690
弁理士 太佐 種一
- (72)発明者 竹内 広宣
東京都江東区豊洲五丁目6番52号 NBF豊洲チャンネルフロント 日本アイ・ピー・エム株式会
社 東京基礎研究所内
- (72)発明者 中村 大賀
東京都江東区豊洲五丁目6番52号 NBF豊洲チャンネルフロント 日本アイ・ピー・エム株式会
社 東京基礎研究所内
- (72)発明者 十河 維
東京都中央区日本橋箱崎町19番地21 日本アイ・ピー・エム株式会社内

審査官 長 由紀子

- (56)参考文献 特開2011-096245(JP,A)
特開平10-260824(JP,A)
特開2003-015869(JP,A)
特開平04-215182(JP,A)
特開平06-231188(JP,A)
特開平07-065008(JP,A)
特開2002-259370(JP,A)
特開2005-293582(JP,A)
高田 京児, ビジネスの視点でデータを整理 内部統制にも対応できる統合データベース設計の
研究, DB Magazine, 日本, 株式会社翔泳社, 2009年 2月 1日, 第18巻第10
号, p.98-110
- (58)調査した分野(Int.Cl., DB名)
G06F 17/20-28