

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2021年3月4日(04.03.2021)



(10) 国際公開番号  
**WO 2021/038887 A1**

- (51) 国際特許分類:  
G06F 16/383 (2019.01) G06F 16/90 (2019.01)  
G06F 16/31 (2019.01) G06F 16/93 (2019.01)
- (21) 国際出願番号: PCT/JP2019/034306
- (22) 国際出願日: 2019年8月30日(30.08.2019)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (71) 出願人: 富士通株式会社 (FUJITSU LIMITED) [JP/JP]; 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 Kanagawa (JP).
- (72) 発明者: 馬場 謙介 (BABA, Kensuke); 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP). 野呂 智哉

(NORO, Tomoya); 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP). 福田 茂紀 (FUKUTA, Shigeki); 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP). 大倉 清司 (OKURA, Seiji); 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP).

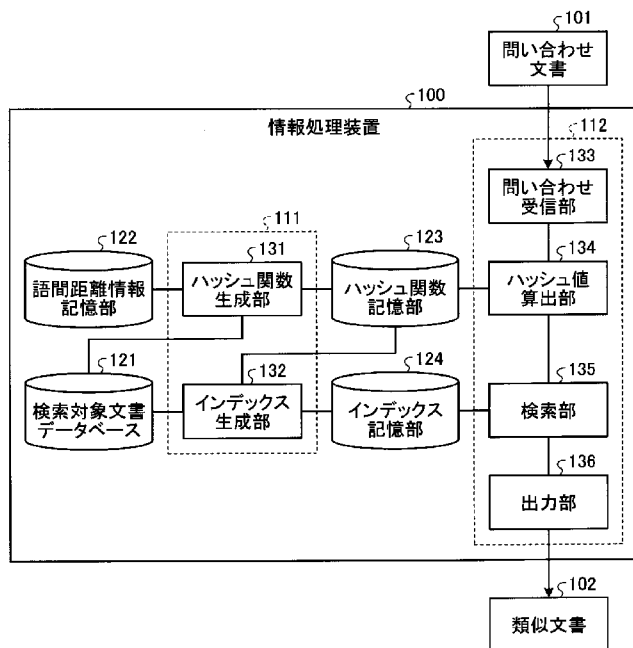
(74) 代理人: 特許業務法人酒井国際特許事務所 (SAKAI INTERNATIONAL PATENT OFFICE); 〒1000013 東京都千代田区霞が関3丁目8番1号 虎の門三井ビルディング Tokyo (JP).

(81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ,

(54) Title: SIMILAR DOCUMENT RETRIEVAL METHOD, SIMILAR DOCUMENT RETRIEVAL PROGRAM, SIMILAR DOCUMENT RETRIEVAL DEVICE, INDEX INFORMATION CREATION METHOD, INDEX INFORMATION CREATION PROGRAM, AND INDEX INFORMATION CREATION DEVICE

(54) 発明の名称: 類似文書検索方法、類似文書検索プログラム、類似文書検索装置、索引情報作成方法、索引情報作成プログラムおよび索引情報作成装置

[図1]



- 100... INFORMATION PROCESSING DEVICE
- 101... QUERY DOCUMENT
- 102... SIMILAR DOCUMENT
- 102... RETRIEVAL DOCUMENT
- 121... RETRIEVAL DOCUMENT DATABASE
- 122... WORD INTERVAL DISTANCE INFORMATION STORAGE UNIT
- 123... HASH FUNCTION STORAGE UNIT
- 124... INDEX STORAGE UNIT
- 131... HASH FUNCTION GENERATION UNIT
- 132... INDEX GENERATION UNIT
- 133... QUERY RECEIVING UNIT
- 134... HASH VALUE COMPUTATION UNIT
- 135... RETRIEVAL UNIT
- 136... OUTPUT UNIT

(57) Abstract: In a similar document retrieval method according to an embodiment, a computer executes a generation process, a computation process, and a retrieval process. The generation process generates a hash function for allocating a value to each word included in a set of words in a retrievable document on the basis of the set of words and word

WO 2021/038887 A1

BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

- (84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類 :

- 一 国際調査報告 (条約第21条(3))

---

interval information indicating word meaning closeness, said hash function allocating a closer value with reference to a specific word in the order of closeness to said word. The computation process computes summary information for each of a plurality of retrievable documents on the basis of the generated hash function and computes summary information for an inputted document on the basis of the generated hash function. The retrieval process retrieves a document similar to the inputted document among the plurality of retrievable documents on the basis of a comparison between the computed summary information for the retrievable documents and the computed summary information for the inputted document.

(57) 要約 : 実施形態の類似文書検索方法は、生成する処理と、算出する処理と、検索する処理とをコンピュータが実行する。生成する処理は、検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成する。算出する処理は、生成したハッシュ関数に基づいて複数の検索対象文書それぞれの要約情報を算出し、生成したハッシュ関数に基づいて入力文書の要約情報を算出する。検索する処理は、算出した検索対象文書の要約情報と、入力文書の要約情報との間の比較に基づいて、複数の検索対象文書の中から入力文書に類似する文書を検索する。

## 明 細 書

発明の名称：

類似文書検索方法、類似文書検索プログラム、類似文書検索装置、索引情報作成方法、索引情報作成プログラムおよび索引情報作成装置

### 技術分野

[0001] 本発明の実施形態は、類似文書検索方法、類似文書検索プログラム、類似文書検索装置、索引情報作成方法、索引情報作成プログラムおよび索引情報作成装置に関する。

### 背景技術

[0002] 従来、コンピュータによる自然言語処理の一つとして、データベースに記憶された文書の中から入力文書に類似する文書を検索する検索処理がある。例えば、問い合わせ文書のサンプルと当該サンプルに対応する返答文書とをデータベースに登録しておく。そして、入力された問い合わせ文書に類似するサンプルを検索し、当該類似するサンプルに対応する返答文書を出力するチャットボット等の対話インタフェースを構築することが考えられる。

[0003] データベースに登録されたサンプルの中から入力文書に類似するサンプルを検索する方法としては、2つの文書の間で出現する単語の共通度を評価する方法がある。例えば、以下のような検索方法が考えられる。ある文書に含まれる単語集合から1つのハッシュ値を算出するハッシュ関数（Min-hash関数と言うことがある）を複数個定義しておく。各ハッシュ関数は、異なる単語に対して異なる値を対応付けた対応関係を持ち、ある単語集合に含まれる単語に対応する値のうち最小の値をハッシュ値として出力する。データベースに登録されたサンプルに対して、この複数のハッシュ関数を用いて算出された複数のハッシュ値を列挙したベクトルを予め生成しておく。そして、入力文書に含まれる単語集合と上記の複数のハッシュ関数から同様にベクトルを算出し、データベースに登録されたサンプルのベクトルと近似するものを検索する。

## 先行技術文献

### 特許文献

[0004] 特許文献1：特開2018-173909号公報

### 発明の概要

#### 発明が解決しようとする課題

[0005] しかしながら、上記の従来技術では、入力文書が短い、文書に用いられる表現が多様でサンプルと一致する表現が少ないなどの場合に、単語の共通度が低く評価され、検索精度が低下するという問題がある。例えば、チャットボットの入力文書は、主に話し言葉であり、1文が短く表現も多様であることから、入力文書と、サンプルとの間で共通の単語が出現する確率が全体的に低くなる。この結果、入力文書に対するサンプルの検索精度が低くなりやすい。

[0006] 1つの側面では、類似する文書の検索精度を向上させる類似文書検索方法、類似文書検索プログラム、類似文書検索装置、索引情報作成方法、索引情報作成プログラムおよび索引情報作成装置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

[0007] 1つの案では、類似文書検索方法は、生成する処理と、算出する処理と、検索する処理とをコンピュータが実行する。生成する処理は、検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成する。算出する処理は、生成したハッシュ関数に基づいて複数の検索対象文書それぞれの要約情報を算出し、生成したハッシュ関数に基づいて入力文書の要約情報を算出する。検索する処理は、算出した検索対象文書の要約情報と、入力文書の要約情報との間の比較に基づいて、複数の検索対象文書の中から入力文書に類似する文書を検索する。

#### 発明の効果

[0008] 1つの側面では、類似する文書の検索精度が向上する。

### 図面の簡単な説明

[0009] [図1]図1は、実施形態にかかる情報処理装置の機能構成例を示すブロック図である。

[図2]図2は、実施形態にかかる情報処理装置の処理フローの一例を示す説明図である。

[図3]図3は、ハッシュ関数の生成処理の一例を示すフローチャートである。

[図4]図4は、ハッシュ関数を説明する説明図である。

[図5]図5は、語間の類似度を説明する説明図である。

[図6]図6は、Min hashによるハッシュ値を説明する説明図である。

[図7]図7は、実施形態にかかる情報処理装置の動作の概要を説明する説明図である。

[図8]図8は、検索対象文書の絞り込みを説明する説明図である。

[図9]図9は、操作画面の表示例を示す説明図である。

[図10]図10は、プログラムを実行するコンピュータの一例を示す図である。

### 発明を実施するための形態

[0010] 以下、図面を参照して、実施形態にかかる類似文書検索方法、類似文書検索プログラム、類似文書検索装置、索引情報作成方法、索引情報作成プログラムおよび索引情報作成装置を説明する。実施形態において同一の機能を有する構成には同一の符号を付し、重複する説明は省略する。なお、以下の実施形態で説明する類似文書検索方法、類似文書検索プログラム、類似文書検索装置、索引情報作成方法、索引情報作成プログラムおよび索引情報作成装置は、一例を示すに過ぎず、実施形態を限定するものではない。また、以下の各実施形態は、矛盾しない範囲内で適宜組みあわせてもよい。

[0011] [概要]

本実施形態では、データベースに登録されたサンプル（以下、検索対象文書とも呼ぶ）の中から入力文書（以下、問い合わせ文書とも呼ぶ）に類似す

る文書を検索する情報処理装置を例示する。

[0012] この情報処理装置では、先ず事前処理を行い、次いで検索処理を行う。事前処理では、ハッシュ関数を用いて検索対象文書それぞれについてハッシュ値を算出し、算出したハッシュ値より検索対象文書を検索するための木構造などの索引構造（例えば探索木）を作る。

[0013] 検索処理では、ハッシュ関数を用いて問い合わせ文書のハッシュ値を算出する。次いで、情報処理装置は、問い合わせ文書のハッシュ値と、検索対象文書それぞれのハッシュ値とを比較し、索引構造で示された検索対象文書それぞれのハッシュ値の中から近いものを探す。次いで、情報処理装置は、問い合わせ文書のハッシュ値と近いハッシュ値の検索対象文書を類似する文書の検索結果とする。

[0014] ハッシュ値を算出するハッシュ関数については、データベースに登録された検索対象文書に含まれる単語を抽出して得られた単語の集合をもとに、複数個定義しておく。具体的には、情報処理装置は、 $W$ を単語の集合、 $h$ を $W$ から  $\{0, 1, \dots, |W| - 1\}$  へのすべての単射の集合（複数のハッシュ関数）としてハッシュ関数を複数生成しておく。

[0015] 情報処理装置は、検索対象文書および問い合わせ文書について次の処理を行うことで、複数のハッシュ値によるベクトルを求めて検索対象文書および問い合わせ文書をハッシュ関数で要約した要約情報を得る。

- ・  $h$  からハッシュ関数をランダムに選択する。
- ・ 検索対象文書に含まれる単語を抽出して単語の集合を取得する。
- ・ 選択した  $h$  により各単語から得る整数の中で最小のものをハッシュ値とする。
- ・ 関数のランダムな選択を複数回行うことで複数のハッシュ値を得る。

[0016] 互いに類似する文書同士は、共通語の出現割合（ジャカル係数）が高くなる。 $h$  からのハッシュ関数の選択について、ハッシュ値が一致する確率はジャカル係数に一致する。よって、類似する文書を求める際の文書のハッシュ値同士の比較において、ベクトルの各要素の一致の割合からジャカル

係数を確率的に計算でき、ハッシュ値のハミング距離（例えば不一致の数）が文書間の近さ（類似度）を反映している。

[0017] また、 $n$ を検索対象文書の数、 $m$ を検索対象文書ごとの単語数（平均値）、 $k$ をハッシュ関数の数とする。すると、検索対象文書のベクトルを計算する計算量は $O(kmn)$ となる。1つの問い合わせ文書のベクトルを計算する計算量は $O(km)$ となり、このベクトルを用いた近傍探索の計算コストは $O(\log(n))$ となる。ハッシュ関数は予めランダムに $k$ 個生成される。 $k$ を $O(\log(n))$ とすると、異なる検索対象文書から同一のハッシュ値が算出される衝突確率を十分に小さくすることができる。

[0018] しかしながら、ハミング距離による類似度の検証では、ベクトルの各要素の差に意味はなく、要素が不一致の場合（共通語の出現がない場合）、類似度は0となる。このため、共通語の出現が少ない場合、類似検索の精度が低くなりやすい。

[0019] そこで、本実施形態では、情報処理装置は、単語間の語意の近さを示す語間距離情報をもとに、所定の単語を基準として語意の近い順に近い値を割り当てるハッシュ関数を生成する。具体的には、情報処理装置は、検索対象文書に含まれる単語を抽出して単語の集合し、ランダムに基準とする単語を選択する。次いで、情報処理装置は、語間距離情報をもとに、基準とする単語に対して語意の近い順に単語の集合に含まれる単語をソートし、ソート順に一意的な値（例えばソート順に大きくなる整数値）の割り当てを行う。情報処理装置は、上記の処理を繰り返すことで、複数のハッシュ関数を生成する。

[0020] 例えば、単語の集合が{猫、ご飯、…にゃんこ、えさ}であり、基準とする単語を{猫}とする場合、情報処理装置は、{猫}を基準として語意の近い順にソートする。これにより、{猫、にゃんこ、…ご飯、えさ}が得られる。このようにソートした単語の集合について、情報処理装置は、{猫=0、にゃんこ=1、…ご飯=5、えさ=6}などのように、ソート順に整数値を割り当てる。

[0021] このようにして生成されたハッシュ関数によるハッシュ値は、所定の単語

に対する語意の距離（語間とも呼ぶ）に応じた値となる。このため、ハッシュ値間のハミング距離だけでなく、ユークリッド距離による類似度の検証が可能となる。したがって、共通語の出現がなく、ハッシュ値が不一致（要素が不一致の場合）であっても、ユークリッド距離により類似度を検証することができ、類似検索の精度を向上させることができる。

[0022] [構成例]

図1は、実施形態にかかる情報処理装置の機能構成例を示すブロック図である。図1に示すように、情報処理装置100は、入力された問い合わせ文書101について、検索対象文書データベース121に格納された検索対象文書の中から類似する類似文書102を求める処理を行う装置である。情報処理装置100としては、例えばパーソナルコンピュータ等を適用できる。

[0023] 情報処理装置100は、索引構造を作成する事前処理を行うインデックス生成モジュール111と、検索処理を行う検索モジュール112と、処理に関する各種データを格納する記憶部（検索対象文書データベース121、語間距離情報記憶部122、ハッシュ関数記憶部123およびインデックス記憶部124）とを有する。すなわち、情報処理装置100は、類似文書検索装置および索引情報作成装置の一例である。

[0024] 検索対象文書データベース121は、問い合わせ文書101に対して検索対象となる検索対象文書が登録されたデータベースである。検索対象文書データベース121における検索対象文書は、事前に登録されていてもよいし、情報処理装置100を使用するユーザとのチャットボット等の対話インタフェースにおける対話を通じて追加されてもよいし、ネットワークから自動的に収集されてもよい。

[0025] 語間距離情報記憶部122は、単語それぞれについて、他の単語との語意の近さ（語間）を表す語間距離情報を格納する。具体的には、語間距離情報としては、単語それぞれ（ $v$ ）と、他の単語（ $w$ ）との語意の距離（語間）を表す関数 $d(d, w)$ などがある。すなわち、語間距離情報は、単語の語意の近さを示す語間情報の一例である。

- [0026] ハッシュ関数記憶部 1 2 3 は、インデックス生成モジュール 1 1 1 が生成した異なる複数のハッシュ関数を記憶する。各ハッシュ関数は、検索対象文書に出現し得る単語それぞれに対して一意な整数（所定の単語に対して語意の近い順）を対応付ける対応関係をもち、単語集合を受け付けて1つのハッシュ値を出力する。また、異なるハッシュ関数は異なる対応関係をもつ。
- [0027] インデックス記憶部 1 2 4 は、問い合わせ文書 1 0 1 に類似する検索対象文書を検索するための索引構造を記憶する。索引構造は、複数のハッシュ関数を用いて検索対象文書から算出されたベクトル（要約情報）に基づいてインデックス生成部 1 3 2 が生成したものであり、検索対象文書それぞれの要約情報を探索する索引情報の一例である。
- [0028] 索引構造としては、例えば探索木を適用できる。探索木は、木構造に接続された複数のノード（葉ノードと、葉ノードに至る各ノード）を含む。探索木の葉ノードは、検索対象文書を指し示す。例えば、葉ノードは、検索対象文書のベクトルと当該検索対象文書を識別する識別情報（例えば文書ID）とを含む。ただし、葉ノードがベクトルを含まなくてもよい。
- [0029] 葉ノード以外の各ノードには2つの子ノードが接続されている。葉ノード以外の各ノードは、ベクトルの中の特定の次元に対する閾値をもつ。入力されたベクトルの中の特定の次元のハッシュ値が閾値以上である場合は右子ノードに進み、特定の次元のハッシュ値が閾値未満である場合は左子ノードに進む。このようにして、探索木をルートノードから葉ノードに向かって辿ることで、所定のベクトルに近い検索対象文書を効率的に検索することができる。
- [0030] インデックス生成モジュール 1 1 1 は、ハッシュ関数生成部 1 3 1 と、インデックス生成部 1 3 2 とを有する。
- [0031] ハッシュ関数生成部 1 3 1 は、ハッシュ関数生成部 1 3 1 に記憶された検索対象文書と、問い合わせ受信部 1 3 3 に記憶された語間距離情報とに基づいて、複数のハッシュ関数を生成し、生成した複数のハッシュ関数をハッシュ関数記憶部 1 2 3 に格納する。

[0032] 具体的には、ハッシュ関数生成部131は、検索対象文書に含まれる単語を抽出して単語の集合し、ランダムに基準とする単語を選択する。次いで、ハッシュ関数生成部131は、語間距離情報記憶部122の語間距離情報を参照し、基準とする単語に対して語意の近い順に単語の集合に含まれる単語を整列（ソート）する。次いで、ハッシュ関数生成部131は、単語の集合に含まれる単語について、整列順に一意的な値（例えばソート順に大きくなる整数値）を割り当ててハッシュ関数を生成する。ハッシュ関数生成部131は、上記のハッシュ関数を生成する処理を繰り返すことで、複数のハッシュ関数を生成する。

[0033] 例えば、語間距離情報として単語それぞれ（ $v$ ）と、他の単語（ $w$ ）との語意の距離（語間）を表す関数  $d(v, w)$  が与えられているものとする。この関数  $d(v, w)$  については、語間の類似度や単語のベクトル表現などを参照して予め作成することができる。また、関数  $d(v, w)$  は、次のとおりである。

・任意の  $v, w \in W$  について、 $0 \leq d(v, w)$  である。

・任意の  $w \in W$  について、 $d(w, w) = 0$  である。

[0034] また、ハッシュ関数については、任意の  $u, v \in W$  について、 $h(u) > h(v) \Leftrightarrow d(u, w) > d(v, w)$  となる  $w \in W$  があるものとする。

[0035] ハッシュ関数生成部131は、 $W$  から  $w$  をランダムに選択し、すべての  $v \in W$  を  $d(v, w)$  の小さい順にソートする。次いで、ハッシュ関数生成部131は、ソートされた語  $w, v_1, v_2, \dots$  に整数  $0, 1, 2, \dots$  を割り当てる。なお、ハッシュ関数生成部131は、整数を割り当てる代わりに  $d(v, w)$  の値をそのまま使ってもよい（重複がないものとする）。

[0036] インデックス生成部132は、検索対象文書データベース121に記憶された検索対象文書とハッシュ関数記憶部123に記憶されたハッシュ関数に基づいて、索引構造を生成し、生成した索引構造をインデックス記憶部124に格納する。

[0037] 具体的には、インデックス生成部132は、検索対象文書ごとに単語集合

を抽出し、抽出した単語集合を複数のハッシュ関数それぞれに入力して、ハッシュ値のベクトル、すなわち検索対象文書の要約情報を算出する。次いで、インデックス生成部132は、複数の検索対象文書に対応する複数のベクトルを効率的に検索できるように、索引構造を生成する。例えば、インデックス生成部132は、ベクトルの中の1つの次元に着目し、ベクトルの集合が二分割されるように当該次元のハッシュ値の閾値を決定することを繰り返すことで、探索木を生成する。このとき、インデックス生成部132は、探索木の葉ノードにはできる限り単一のベクトルが対応付けられるように中間ノードを生成する。

[0038] 検索モジュール112は、問い合わせ受信部133と、ハッシュ値算出部134と、検索部135と、出力部136とを有する。

[0039] 問い合わせ受信部133は、問い合わせ文書101を受信する。問い合わせ受信部133は、ユーザから文字列として入力された問い合わせ文書101を受信してもよいし、ユーザが口頭で発した問い合わせ発話の音声信号を文字列に変換してもよい。また、問い合わせ受信部133は、他の情報処理装置から文字列または音声信号を受信してもよい。

[0040] ハッシュ値算出部134は、ハッシュ関数記憶部123に記憶された複数のハッシュ関数に基づいて、問い合わせ文書101に対応するベクトル、すなわち問い合わせ文書101の要約情報を生成する。具体的には、ハッシュ値算出部134は、問い合わせ文書101から単語集合を抽出し、抽出した単語集合を複数のハッシュ関数それぞれに入力して、ハッシュ値のベクトルを算出する。

[0041] 検索部135は、インデックス記憶部124に記憶された索引構造と問い合わせ文書101のベクトルに基づいて、近傍探索により問い合わせ文書101に最も類似する検索対象文書を検索する。具体的には、問い合わせ文書101に最も類似する検索対象文書は、ベクトル同士を比較したときにハッシュ値が一致する次元が最も多いものである。

[0042] 例えば、検索部135は、問い合わせ文書101のベクトルの中の特定の

次元のハッシュ値と閾値とを比較しながら、インデックス記憶部124に記憶された探索木をルートノードから葉ノードに向かって辿り、特定の葉ノードに到達する。検索部135は、到達した葉ノードに対応する検索対象文書を選択する。

[0043] なお、到達した葉ノードに2以上の検索対象文書が対応付けられている場合、すなわち、ハッシュ値が一致する次元数が同じであり、探索木では検索対象文書を1つに絞り込めない場合、検索部135は、ハッシュ値同士を比較してユークリッド距離を求める。次いで、検索部135は、ユークリッド距離がより近いものを最も類似する検索対象文書とする。

[0044] 出力部136は、検索された検索対象文書を類似文書102として出力する。例えば、出力部136は、類似文書102の文字列をディスプレイ等に表示してもよいし、類似文書102を音声信号に変換してスピーカにより音声を再生してもよい。また、出力部136は、他の情報処理装置に類似文書102の文字列または音声信号を送信してもよい。

[0045] また、出力部136は、類似文書102を出力する代わりに、検索された検索対象文書に対して検索対象文書データベース121において予め紐付けられた処理を実施してもよい。具体的には、検索対象文書データベース121には、検索対象文書ごとに、所定の処理（例えばスケジュール登録、メール送信）が登録されているものとする。出力部136は、検索された検索対象文書に紐付けられた処理を検索対象文書データベース121より読み出して実行することで、問い合わせ文書101に対応した処理を行うことが可能となる。

[0046] [動作例]

図2は、実施形態にかかる情報処理装置の処理フローの一例を示す説明図である。図2に示すように、情報処理装置100は、索引構造を作成する事前処理（S1）と、問い合わせ文書101に対する類似文書102を検索して出力する検索処理（S2）とを行う。

[0047] 先ず、事前処理（S1）について説明する。事前処理（S1）では、先ず

検索対象文書データベース121より検索対象文書が読み出され、ハッシュ関数生成部131に入力される(S11)。

[0048] ハッシュ関数生成部131では、問い合わせ受信部133より語間の距離の入力を受け付け(S13)、入力された検索対象文書と、語間の距離とに基づき、複数のハッシュ関数123aを生成する(S12)。

[0049] 図3は、ハッシュ関数の生成処理の一例を示すフローチャートである。図3に示すように、ハッシュ関数生成部131は、検索対象文書の入力を受け付け(S31)、検索対象文書に含まれる単語(出現語)を抽出することで(S32)、単語の集合(語集合)を取得する(S33)。

[0050] 次に、ハッシュ関数生成部131は、S34~S39の処理を生成するハッシュ関数の数であるk回繰り返すことで、複数のハッシュ関数123aを生成する。

[0051] 具体的には、ハッシュ関数生成部131は、語集合の中からランダムに1つの単語を選択し(S35)、語間の距離を示す語間距離情報の入力を受け付け(S36)、選択した単語と他の単語との距離を語間距離情報より参照する(S37)。次に、ハッシュ関数生成部131は、選択した単語との距離の近い順に語集合の単語をソートし(S38)、整列した各単語に整列順に大きくなるような整数値をハッシュ値として割り当てる(S39)。

[0052] 次に、ハッシュ関数生成部131は、S34~S39の処理をk回繰り返して得られた複数のハッシュ関数123aを出力し、インデックス生成部132に格納する(S40)。

[0053] 図4は、ハッシュ関数123aを説明する説明図である。図4に示すように、ハッシュ関数123aにおける $h_1, h_2, \dots, h_2, \dots$ が1つのハッシュ関数である。

[0054] 例えば、 $h_1$ は、(猫)を基準の単語としており、語集合における各単語について(猫)に対する語間の距離に応じた整数値が割り当てられている。また、 $h_2$ は、(ごはん)を基準の単語としており、語集合における各単語について(ごはん)に対する語間の距離に応じた整数値が割り当てられている。

また、 $h_3$ は、（にゃんこ）を基準の単語としており、語集合における各単語について（にゃんこ）に対する語間の距離に応じた整数値が割り当てられている。また、 $h_4$ は、（えさ）を基準の単語としており、語集合における各単語について（えさ）に対する語間の距離に応じた整数値が割り当てられている。また、 $h_5$ は、（花）を基準の単語としており、語集合における各単語について（花）に対する語間の距離に応じた整数値が割り当てられている。また、 $h_6$ は、（水）を基準の単語としており、語集合における各単語について（水）に対する語間の距離に応じた整数値が割り当てられている。

[0055] ここで、図4に例示したハッシュ関数1 2 3 aにより次の文書A～文書Cのハッシュ値を求める場合を例示する。

文書A：「猫にごはんをやる」

文書B：「にゃんこにえさをやる」

文書C：「花に水をやる」

[0056] 文書Aの語集合は{猫、ご飯}、文書Bの語集合は{にゃんこ、えさ}、文書Cの語集合は{花、水}である。よって、図4に例示したハッシュ関数1 2 3 aによる文書Aのハッシュ値 $H_A$ は、 $H_A = 0 0 1 1 3 3$ となる。同様に、文書Bのハッシュ値 $H_B$ は、 $H_B = 1 1 0 0 2 2$ となる。また、文書Cのハッシュ値 $H_C$ は、 $H_C = 4 2 4 2 0 0$ となる。

[0057] ここで、文書A～Cのハッシュ値を比較してハミング距離を計算すると、 $H_A, H_B : 6$ 、 $H_A, H_C : 6$ 、 $H_B, H_C : 6$ となる。また、文書A～Cのハッシュ値を比較してユークリッド距離を計算すると、 $H_A, H_B : 1$ 、 $H_A, H_C : 6.9$ 、 $H_B, H_C : 6.2$ となる。このように、ハミング距離では類似度の検証が困難な場合（共通語が含まれていない場合）でも、ユークリッド距離により互いに類似する文書（図示例では文書A、B）を検証することができる。

[0058] 図2に戻り、S 1 2に次いで、インデックス生成部1 3 2は、入力された検索対象文書と、生成したハッシュ関数1 2 3 aに基づいて、索引構造1 2 4 aを生成し（S 1 4）、生成した索引構造1 2 4 aをインデックス記憶部

124に格納する。

[0059] 検索処理(S2)では、問い合わせ受信部133が受信した問い合わせ文書101がハッシュ値算出部134に入力される(S21)。ハッシュ値算出部134は、ハッシュ関数記憶部123に記憶された複数のハッシュ関数に基づいて、入力された問い合わせ文書101のハッシュ値を複数生成し(S22)、問い合わせ文書101に対応するベクトルを得る。

[0060] 次に、検索部135は、インデックス記憶部124に記憶された索引構造における検索対象文書それぞれのベクトルと、問い合わせ文書101のベクトルのハッシュ値を照合し(S23)、問い合わせ文書101に最も類似する検索対象文書を検索する。出力部136は、検索部135が検索した類似文書102を出力する(S24)。

[0061] 図5は、語間の類似度を説明する説明図である。具体的には、図5は、語間の類似度を示す高次元の空間における語W1~W6それぞれの配置を俯瞰した図である。図5の語W1~W6それぞれは、文書に含まれる単語を示す。ここで、語W1~W3は、例えば「猫」などの単語について類似するものであり、点線で示すクラスターを形成している。同様に、語W4~W6は、例えば「犬」などの単語について類似するものであり、語W1~W3とは別のクラスターを形成している。

[0062] 図5に示すように、語間の類似度を示す高次元の空間において、単純な射影では、類似度をうまく評価することは困難である。例えば、軸A1への直交射影では、類似する単語(語W1~W3または語W4~W6)の値が近くなる。しかしながら、軸A1とは異なる軸A2への直交射影では、互いに類似しない単語(例えば語W1と語W4)が互いに類似する単語(例えば語W1と語W3)よりも近くなる場合がある。

[0063] 本実施形態では、ランダムに基準とする単語を選択し、Min hashによるハッシュ値を用いることから、基準とする単語(基準点)からの距離による射影を用いている。

[0064] 図6は、Min hashによるハッシュ値を説明する説明図である。こ

ここで、語W1は「猫」、語W2は「にゃんこ」、語W3は「キヤット」、語W4は「犬」、語W5は「鼠」、語W6は「ドッグ」、語W7は「わんこ」であるものとする。また、基準点は語W1であるものとする。

[0065] 図6に示すように、基準点である語W1をもとに、ハッシュ関数123aでは、1:「猫」、2:「にゃんこ」、3:「キヤット」、4:「犬」、5:「鼠」、6:「ドッグ」、7:「わんこ」の値が割り当てられる。ここで、基準点からの距離ならば、基準点に近いほど類似度の大小関係が保たれている(例えば1:「猫」、2:「にゃんこ」、3:「キヤット」など)。しかしながら、基準点からの距離が遠いところでは(例えば5:「鼠」、6:「ドッグ」)、類似しない語の値が近くなる場合がある。

[0066] 本実施形態では、語集合のハッシュ値の中の最小値を用いる(Min Hash)ので、類似度の大小関係が保たれた、語の類似度を表現できる射影として適切なものとなっている。

[0067] 図7は、実施形態にかかる情報処理装置100の動作の概要を説明する説明図である。図7に示すように、事前処理(S1)において、情報処理装置100は、語間距離情報による語間の類似度をもとに、検索対象文書データベース121の検索対象文書121aそれぞれに含まれる単語の集合について、所定の単語を基準として語意の近い順に近い値を割り当てるハッシュ関数を複数生成する。次いで、情報処理装置100は、生成した複数のハッシュ関数を用いて検索対象文書121aそれぞれをMin Hashによる変換を行い、算出されたベクトルを検索するための索引構造を生成してインデックス記憶部124に格納する。

[0068] 検索処理(S2)において、情報処理装置100は、入力された問い合わせ文書101について、S1で生成した同じ複数のハッシュ関数を用いてMin Hashによる変換を行う。次いで、情報処理装置100は、問い合わせ文書101より算出したハッシュ値のベクトルと、インデックス記憶部124に格納されたハッシュ値のベクトルとを比較することで、問い合わせ文書101に最も類似する検索対象文書121aを検索する。

- [0069] 図示例では、「会議を調整したい」とする問い合わせ文書101に対し、検索対象文書121aの中で「スケジュールの調整」と、「打ち合わせの調整」とが「調整」という単語においてハッシュ値が一致することから、ハミング距離が近いものとなる。ここで、問い合わせ文書101に含まれる「会議」に対し、「スケジュール」よりも「会議」の方がより語意が近く、ユークリッド距離が近いものとなる。したがって、「会議を調整したい」とする問い合わせ文書101に対しては、「打ち合わせを調整したい」とする類似文書102が得られることとなる。
- [0070] なお、出力部136は、検索された検索対象文書について、最も類似する1つの類似文書102を出力してもよいし、ハミング距離およびユークリッド距離により得られた類似度が所定の閾値以上である複数の類似文書102を出力してもよい。
- [0071] 図8は、検索対象文書の絞り込みを説明する説明図である。具体的には、図8は、類似度の閾値と検索対象文書121aのヒット数との関係例を示すグラフである。
- [0072] グラフ10は、問い合わせ文書101と検索対象文書121aの類似度の閾値と、類似度が閾値より大きい検索対象文書121aの数（ヒット数）との間の関係を示す。類似度は、問い合わせ発話のベクトルと検索対象発話のベクトルの間でハッシュ値が一致する次元の数に、ハッシュ値同士のユークリッド距離を合わせたものである。
- [0073] (A) 関連語を考慮せずにベクトルを算出する方法では、問い合わせ文書101と各検索対象文書121aの類似度が全体として低く算出される。よって、類似度の閾値とヒット数との間の関係は曲線11のようになる。すなわち、類似度の閾値を低く設定してもヒット数が少なくなり、類似する検索対象文書121aの検索漏れが多くなる。
- [0074] (B) 関連語を同一視してベクトルを算出する方法では、問い合わせ文書101と各検索対象文書121aの類似度が全体として高く算出される。よって、類似度の閾値とヒット数との間の関係は曲線12のようになる。すな

わち、類似度の閾値を高く設定してもヒット数が多くなり、類似する検索対象文書 1 2 1 a を効率的に絞り込むことが難しい。

[0075] (C) 本実施形態の方法では、問い合わせ文書 1 0 1 と各検索対象文書 1 2 1 a の類似度がユークリッド距離を加味したものとなる。よって、類似度の閾値とヒット数との間の関係は曲線 1 3 のように連続的なものとなる。その結果、問い合わせ文書 1 0 1 に類似する検索対象文書 1 2 1 a を効率的に絞り込むことができる。

[0076] 図 9 は、操作画面の表示例を示す説明図である。操作画面 2 0 は、例えばチャットボット等の対話インタフェースにおいて情報処理装置 1 0 0 のユーザに提示することでユーザからの各種操作を受け付ける画面である。操作画面 2 0 において、表示領域 2 1 は処理結果などを表示する領域であり、入力領域 2 2 は文書などの入力を行う領域である。例えば、情報処理装置 1 0 0 は、入力領域 2 2 に入力された問い合わせ文書 1 0 1 に対して類似文書 1 0 2 の検索を行い、検索結果に応じた出力を表示領域 2 1 に行う。

[0077] 例えば、情報処理装置 1 0 0 は、「スケジュール調整したい」とする入力文書 2 1 a を問い合わせ文書 1 0 1 として類似文書 1 0 2 の検索を行い、検索結果 2 1 b を表示領域 2 1 に表示する。具体的には、「スケジュール調整したい」とする問い合わせ文書 1 0 1 に対し、検索対象文書 1 2 1 a の中から類似文書 1 0 2 として得られた「スケジュールの調整」に対応するスケジュール登録の処理が実行される。

[0078] 情報処理装置 1 0 0 では、図 9 の例におけるチャットボット等の対話インタフェースのように、問い合わせ文書 1 0 1 が主に話し言葉であり、1 文が短く表現も多様である場合であっても、適切に類似文書 1 0 2 の検索を行うことができる。

[0079] [効果]

以上のように、情報処理装置 1 0 0 は、ハッシュ関数生成部 1 3 1 と、インデックス生成部 1 3 2 と、ハッシュ値算出部 1 3 4 と、検索部 1 3 5 とを有する。ハッシュ関数生成部 1 3 1 は、検索対象文書 1 2 1 a に含まれる単

語の集合と、単語の語意の近さを示す語間距離情報とに基づき、単語の集合に含まれる単語それぞれに対して、所定の単語を基準として語意の近い順に近い値を割り当てるハッシュ関数を生成する。インデックス生成部132は、生成したハッシュ関数に基づいて複数の検索対象文書121aそれぞれの要約情報を算出する。ハッシュ値算出部134は、生成したハッシュ関数に基づいて入力文書（問い合わせ文書101）の要約情報を算出する。検索部135は、算出した検索対象文書121aの要約情報と、問い合わせ文書101の要約情報との間の比較に基づいて、複数の検索対象文書121aの中から入力文書に類似する文書を検索する。

[0080] このため、情報処理装置100において、生成したハッシュ関数に基づいて算出した検索対象文書121aの要約情報と、問い合わせ文書101の要約情報とに含まれるハッシュ値は、所定の単語に対する語意の距離に応じた値となる。したがって、情報処理装置100では、検索対象文書121aの要約情報と、問い合わせ文書101の要約情報との間の比較において、例えば互いの要約情報の間におけるハミング距離だけでなく、ユークリッド距離を用いた類似度合いの検証を行うことができ、類似する文書の検索精度を向上させることができる。

[0081] また、インデックス生成部132およびハッシュ値算出部134のそれぞれは、生成したハッシュ関数において検索対象文書121aまたは問い合わせ文書101に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する。このように、情報処理装置100では、Min-hash関数により検索対象文書121aまたは問い合わせ文書101の要約情報を算出できる。

[0082] また、ハッシュ関数生成部131は、所定の単語を選び直してハッシュ関数を生成する処理を繰り返すことでハッシュ関数を複数生成する。そして、インデックス生成部132およびハッシュ値算出部134のそれぞれは、生成した複数のハッシュ関数によって検索対象文書121aまたは問い合わせ文書101に含まれる単語の集合から算出される複数のハッシュ値を含むべ

クトルを要約情報として算出する。これにより、情報処理装置100では、検索対象文書121aの要約情報と、問い合わせ文書101の要約情報との間の比較において、複数のハッシュ値を列挙したベクトルの比較により、例えばハミング距離、ユークリッド距離などを求めて類似度合いを検証することができる。

[0083] また、情報処理装置100のインデックス生成部132は、算出した複数の検索対象文書121aそれぞれの要約情報と、検索対象文書121aとを対応付けた索引情報を生成する。検索部135は、問い合わせ文書101の要約情報と、索引情報において検索対象文書121aと対応付けられた要約情報との間の比較を行う。情報処理装置100では、索引情報を生成しておくことで、複数の検索対象文書121aの中から問い合わせ文書101に類似する文書の検索を高速に行うことができる。

[0084] [その他]

なお、図示した各装置の各構成要素は、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。

[0085] 例えば、本実施形態では、インデックス生成モジュール111と、検索モジュール112とを有する情報処理装置100を例示したが、インデックス生成モジュール111と、検索モジュール112とはそれぞれ異なる情報処理装置が有していてもよい。すなわち、事前処理(S1)と、検索処理(S2)とは、それぞれ異なる情報処理装置で実施してもよい。

[0086] また、情報処理装置100で行われる各種処理機能は、CPU（またはMPU、MCU (Micro Controller Unit) 等のマイクロ・コンピュータ) 上で、その全部または任意の一部を実行するようにしてもよい。また、各種処理機能は、CPU（またはMPU、MCU等のマイクロ・コンピュータ) で解析実行されるプログラム上、またはワイヤードロジックによるハードウエ

ア上で、その全部または任意の一部を実行するようにしてもよいことは言うまでもない。また、情報処理装置100で行われる各種処理機能は、クラウドコンピューティングにより、複数のコンピュータが協働して実行してもよい。

[0087] ところで、上記の実施形態で説明した各種の処理は、予め用意されたプログラムをコンピュータで実行することで実現できる。そこで、以下では、上記の実施例と同様の機能を有するプログラムを実行するコンピュータ（ハードウェア）の一例を説明する。図10は、プログラムを実行するコンピュータの一例を示す図である。

[0088] 図10に示すように、コンピュータ1は、各種演算処理を実行するCPU201と、データ入力を受け付ける入力装置202と、モニタ203と、スピーカ204とを有する。また、コンピュータ1は、記憶媒体からプログラム等を読み取る媒体読取装置205と、各種装置と接続するためのインタフェース装置206と、有線または無線により外部機器と通信接続するための通信装置207とを有する。また、コンピュータ1は、各種情報を一時記憶するRAM208と、ハードディスク装置209とを有する。また、コンピュータ1内の各部（201～209）は、バス210に接続される。

[0089] ハードディスク装置209には、上記の実施形態で説明した各種の処理を実行するためのプログラム211が記憶される。また、ハードディスク装置209には、プログラム211が参照する各種データ212（例えば語間距離情報記憶部122、検索対象文書データベース121、ハッシュ関数記憶部123およびインデックス記憶部124の情報）が記憶される。入力装置202は、例えば、コンピュータ1の操作者から操作情報の入力を受け付ける。モニタ203は、例えば、操作者が操作する各種画面を表示する。インタフェース装置206は、例えば印刷装置等が接続される。通信装置207は、LAN（Local Area Network）等の通信ネットワークと接続され、通信ネットワークを介した外部機器との間で各種情報をやりとりする。

[0090] CPU201は、ハードディスク装置209に記憶されたプログラム21

1を読み出して、RAM 208に展開して実行することで、ハッシュ関数生成部131、インデックス生成部132、問い合わせ受信部133、ハッシュ値算出部134、検索部135および出力部136に関する各種の処理を行う。なお、プログラム211は、ハードディスク装置209に記憶されていなくてもよい。例えば、コンピュータ1が読み取り可能な記憶媒体に記憶されたプログラム211を、コンピュータ1が読み出して実行するようにしてもよい。コンピュータ1が読み取り可能な記憶媒体は、例えば、CD-ROMやDVDディスク、USB (Universal Serial Bus) メモリ等の可搬型記録媒体、フラッシュメモリ等の半導体メモリ、ハードディスクドライブ等が対応する。また、公衆回線、インターネット、LAN等に接続された装置にこのプログラム211を記憶させておき、コンピュータ1がこれらからプログラムを読み出して実行するようにしてもよい。

## 符号の説明

- [0091] 1…コンピュータ
- 10…グラフ
- 11～13…曲線
- 20…操作画面
- 21…表示領域
- 21a…入力文書
- 21b…検索結果
- 22…入力領域
- 100…情報処理装置
- 101…問い合わせ文書
- 102…類似文書
- 111…インデックス生成モジュール
- 112…検索モジュール
- 121…検索対象文書データベース
- 121a…検索対象文書

1 2 2 …語間距離情報記憶部  
1 2 3 …ハッシュ関数記憶部  
1 2 3 a …ハッシュ関数  
1 2 4 …インデックス記憶部  
1 2 4 a …索引構造  
1 3 1 …ハッシュ関数生成部  
1 3 2 …インデックス生成部  
1 3 3 …問い合わせ受信部  
1 3 4 …ハッシュ値算出部  
1 3 5 …検索部  
1 3 6 …出力部  
2 0 1 …CPU  
2 0 2 …入力装置  
2 0 3 …モニタ  
2 0 4 …スピーカ  
2 0 5 …媒体読取装置  
2 0 6 …インタフェース装置  
2 0 7 …通信装置  
2 0 8 …RAM  
2 0 9 …ハードディスク装置  
2 1 0 …バス  
2 1 1 …プログラム  
2 1 2 …各種データ  
A 1 ～ A 2 …軸  
W 1 ～ W 7 …語

## 請求の範囲

- [請求項1] 検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、前記単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成し、
- 生成した前記ハッシュ関数に基づいて複数の前記検索対象文書それぞれの要約情報を算出し、
- 生成した前記ハッシュ関数に基づいて入力文書の要約情報を算出し、
- 算出した前記検索対象文書の要約情報と、前記入力文書の要約情報との間の比較に基づいて、複数の前記検索対象文書の中から前記入力文書に類似する文書を検索する、
- 処理をコンピュータが実行することを特徴とする類似文書検索方法。
- [請求項2] 前記算出する処理のそれぞれは、生成した前記ハッシュ関数において前記検索対象文書または前記入力文書に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する、
- ことを特徴とする請求項1に記載の類似文書検索方法。
- [請求項3] 前記生成する処理は、前記所定の単語を選び直して前記ハッシュ関数を生成する処理を繰り返すことで前記ハッシュ関数を複数生成し、
- 前記算出する処理のそれぞれは、生成した複数の前記ハッシュ関数によって前記検索対象文書または前記入力文書に含まれる単語の集合から算出される複数のハッシュ値を含むベクトルを要約情報として算出する、
- ことを特徴とする請求項2に記載の類似文書検索方法。
- [請求項4] 算出した複数の前記検索対象文書それぞれの要約情報と、前記検索対象文書とを対応付けた索引情報を生成する処理をさらにコンピュー

タが実行し、

前記検索する処理は、前記入力文書の要約情報と、前記索引情報において前記検索対象文書と対応付けられた要約情報との間の比較を行う、

ことを特徴とする請求項 1 に記載の類似文書検索方法。

[請求項5]

検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、前記単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成し、

生成した前記ハッシュ関数に基づいて複数の前記検索対象文書それぞれの要約情報を算出し、

生成した前記ハッシュ関数に基づいて入力文書の要約情報を算出し、

算出した前記検索対象文書の要約情報と、前記入力文書の要約情報との間の比較に基づいて、複数の前記検索対象文書の中から前記入力文書に類似する文書を検索する、

処理をコンピュータに実行させることを特徴とする類似文書検索プログラム。

[請求項6]

前記算出する処理のそれぞれは、生成した前記ハッシュ関数において前記検索対象文書または前記入力文書に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する、

ことを特徴とする請求項 5 に記載の類似文書検索プログラム。

[請求項7]

前記生成する処理は、前記所定の単語を選び直して前記ハッシュ関数を生成する処理を繰り返すことで前記ハッシュ関数を複数生成し、

前記算出する処理のそれぞれは、生成した複数の前記ハッシュ関数によって前記検索対象文書または前記入力文書に含まれる単語の集合から算出される複数のハッシュ値を含むベクトルを要約情報として算

出する、

ことを特徴とする請求項6に記載の類似文書検索プログラム。

[請求項8]

算出した複数の前記検索対象文書それぞれの要約情報と、前記検索対象文書とを対応付けた索引情報を生成する処理をさらにコンピュータが実行し、

前記検索する処理は、前記入力文書の要約情報と、前記索引情報において前記検索対象文書と対応付けられた要約情報との間の比較を行う、

ことを特徴とする請求項5に記載の類似文書検索プログラム。

[請求項9]

検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、前記単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成するハッシュ関数生成部と、

生成した前記ハッシュ関数に基づいて複数の前記検索対象文書それぞれの要約情報を算出する第1の算出部と、

生成した前記ハッシュ関数に基づいて入力文書の要約情報を算出する第2の算出部と、

算出した前記検索対象文書の要約情報と、前記入力文書の要約情報との間の比較に基づいて、複数の前記検索対象文書の中から前記入力文書に類似する文書を検索する検索部と、

を有することを特徴とする類似文書検索装置。

[請求項10]

前記第1の算出部および前記第2の算出部のそれぞれは、生成した前記ハッシュ関数において前記検索対象文書または前記入力文書に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する、

ことを特徴とする請求項9に記載の類似文書検索装置。

[請求項11]

前記ハッシュ関数生成部は、前記所定の単語を選び直して前記ハッシュ関数を生成する処理を繰り返すことで前記ハッシュ関数を複数生

成し、

前記第1の算出部および前記第2の算出部のそれぞれは、生成した複数の前記ハッシュ関数によって前記検索対象文書または前記入力文書に含まれる単語の集合から算出される複数のハッシュ値を含むベクトルを要約情報として算出する、

ことを特徴とする請求項10に記載の類似文書検索装置。

[請求項12]

算出した複数の前記検索対象文書それぞれの要約情報と、前記検索対象文書とを対応付けた索引情報を生成する索引情報生成部をさらに有し、

前記検索部は、前記入力文書の要約情報と、前記索引情報において前記検索対象文書と対応付けられた要約情報との間の比較を行う、

ことを特徴とする請求項9に記載の類似文書検索装置。

[請求項13]

検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、前記単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成し、

生成した前記ハッシュ関数に基づいて複数の前記検索対象文書それぞれの要約情報を算出し、

算出した複数の前記検索対象文書それぞれの要約情報を探索する索引情報を生成する、

処理をコンピュータが実行することを特徴とする索引情報作成方法。

[請求項14]

前記算出する処理は、生成した前記ハッシュ関数において前記検索対象文書に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する、

ことを特徴とする請求項13に記載の索引情報作成方法。

[請求項15]

前記生成する処理は、前記所定の単語を選び直して前記ハッシュ関数を生成する処理を繰り返すことで前記ハッシュ関数を複数生成し、

前記算出する処理は、生成した複数の前記ハッシュ関数によって前記検索対象文書に含まれる単語の集合から算出される複数のハッシュ値を含むベクトルを要約情報として算出する、

ことを特徴とする請求項 1 4 に記載の索引情報作成方法。

[請求項16] 検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、前記単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成し、

生成した前記ハッシュ関数に基づいて複数の前記検索対象文書それぞれの要約情報を算出し、

算出した複数の前記検索対象文書それぞれの要約情報を探索する索引情報を生成する、

処理をコンピュータに実行させることを特徴とする索引情報作成プログラム。

[請求項17] 前記算出する処理は、生成した前記ハッシュ関数において前記検索対象文書に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する、

ことを特徴とする請求項 1 6 に記載の索引情報作成プログラム。

[請求項18] 前記生成する処理は、前記所定の単語を選び直して前記ハッシュ関数を生成する処理を繰り返すことで前記ハッシュ関数を複数生成し、

前記算出する処理は、生成した複数の前記ハッシュ関数によって前記検索対象文書に含まれる単語の集合から算出される複数のハッシュ値を含むベクトルを要約情報として算出する、

ことを特徴とする請求項 1 7 に記載の索引情報作成プログラム。

[請求項19] 検索対象文書に含まれる単語の集合と、単語の語意の近さを示す語間情報とに基づき、前記単語の集合に含まれる単語それぞれに対して、所定の単語を基準として当該単語に対して語意の近い順に近い値を割り当てるハッシュ関数を生成するハッシュ関数生成部と、

生成した前記ハッシュ関数に基づいて複数の前記検索対象文書それぞれの要約情報を算出する算出部と、

算出した複数の前記検索対象文書それぞれの要約情報を探索する索引情報を生成する索引情報生成部と、

を有することを特徴とする索引情報作成装置。

[請求項20] 前記算出部は、生成した前記ハッシュ関数において前記検索対象文書に含まれる単語の集合の単語それぞれに割り当てられた値の中で最小の値をハッシュ値として算出する、

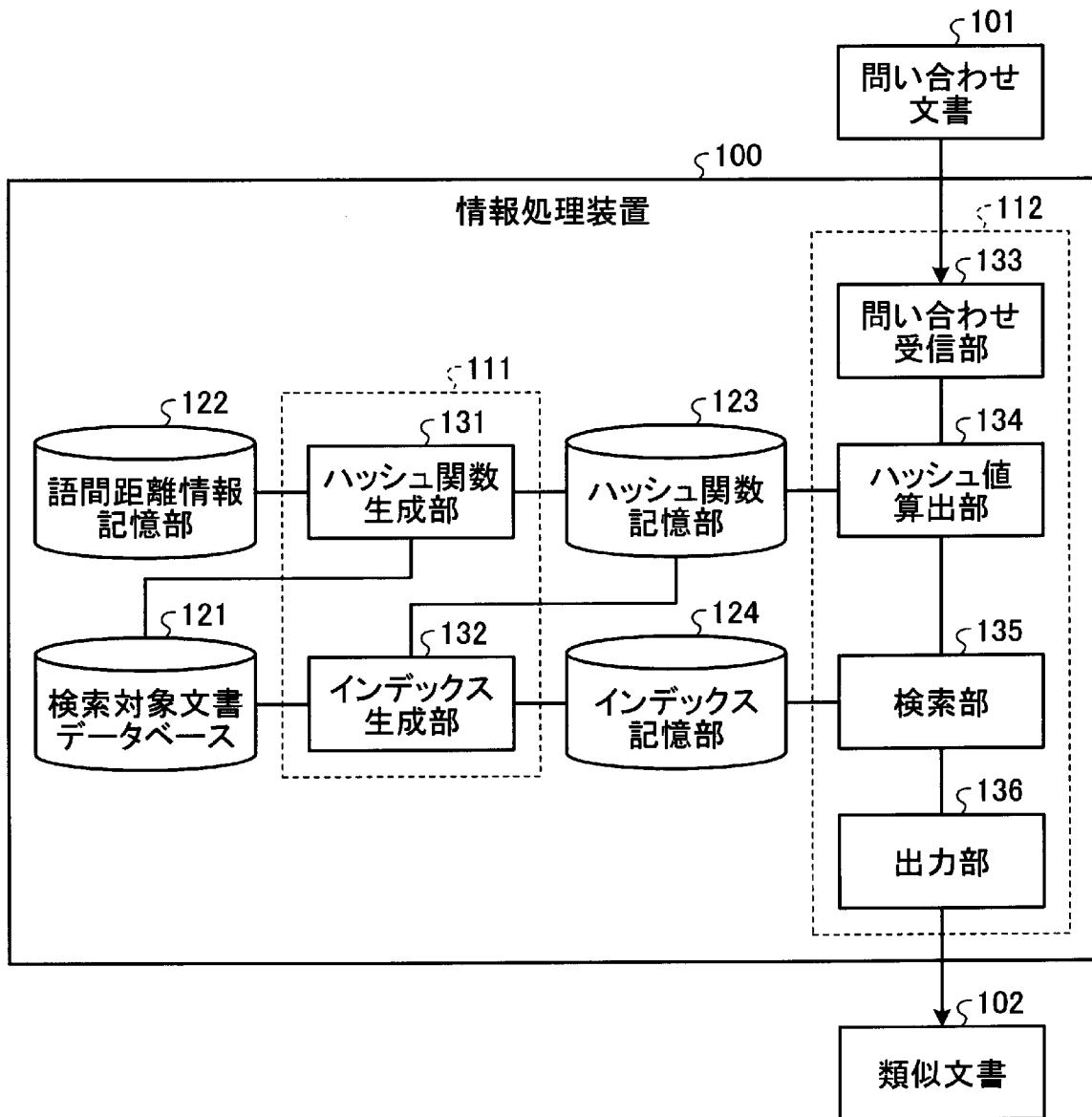
ことを特徴とする請求項19に記載の索引情報作成装置。

[請求項21] 前記ハッシュ関数生成部は、前記所定の単語を選び直して前記ハッシュ関数を生成する処理を繰り返すことで前記ハッシュ関数を複数生成し、

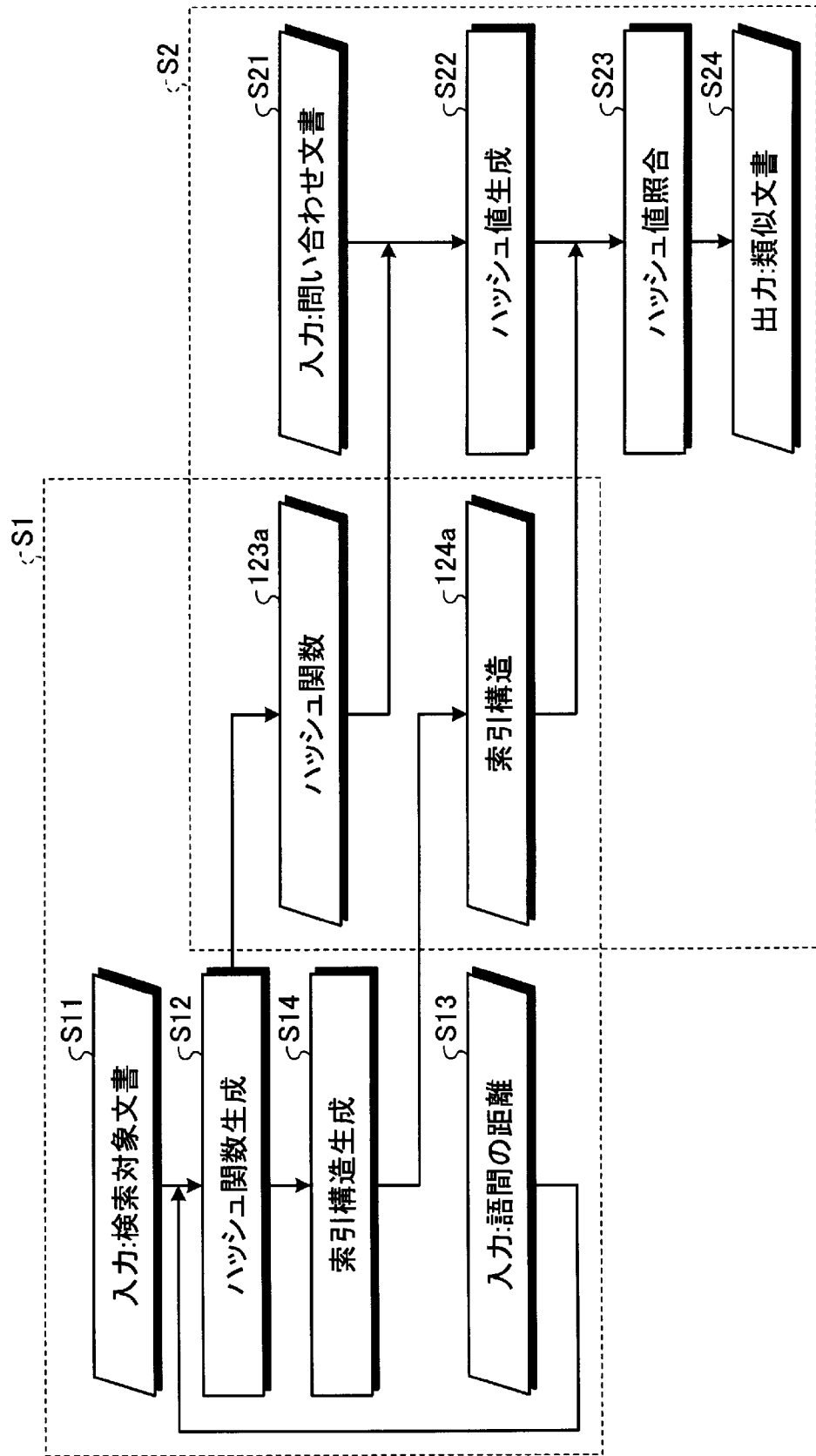
前記算出部は、生成した複数の前記ハッシュ関数によって前記検索対象文書に含まれる単語の集合から算出される複数のハッシュ値を含むベクトルを要約情報として算出する、

ことを特徴とする請求項20に記載の索引情報作成装置。

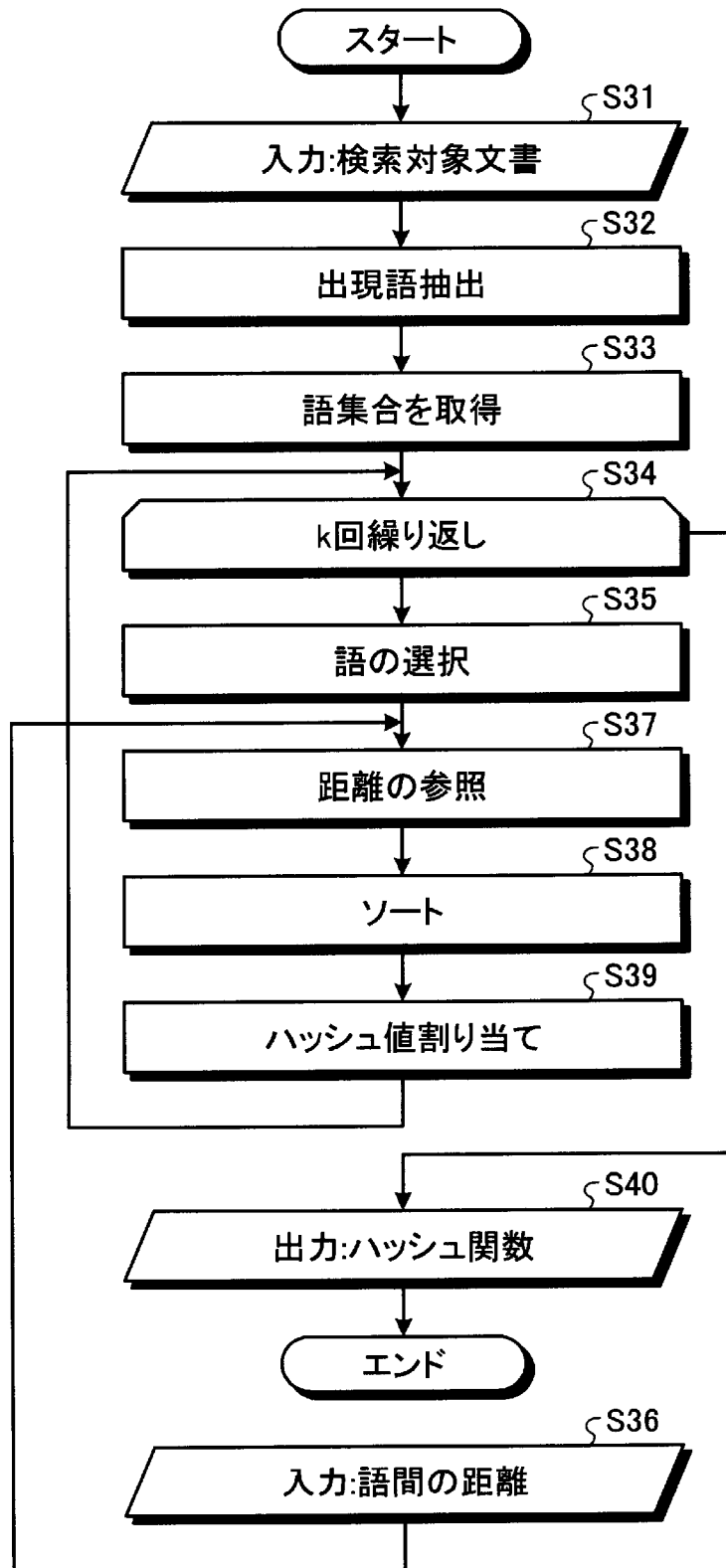
[図1]



[図2]

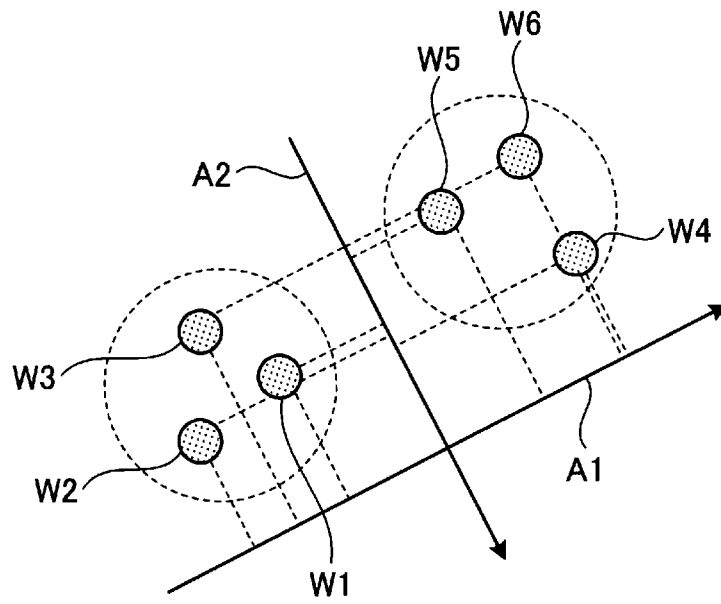


[図3]

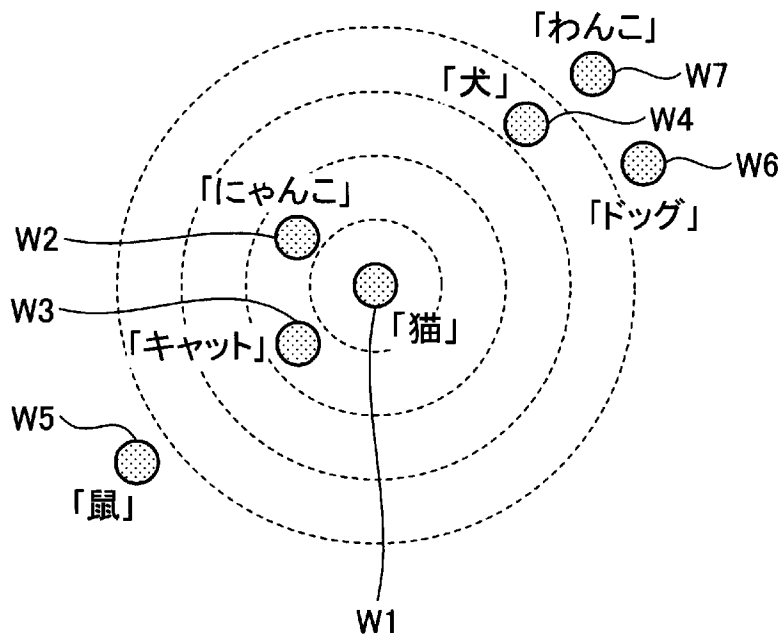




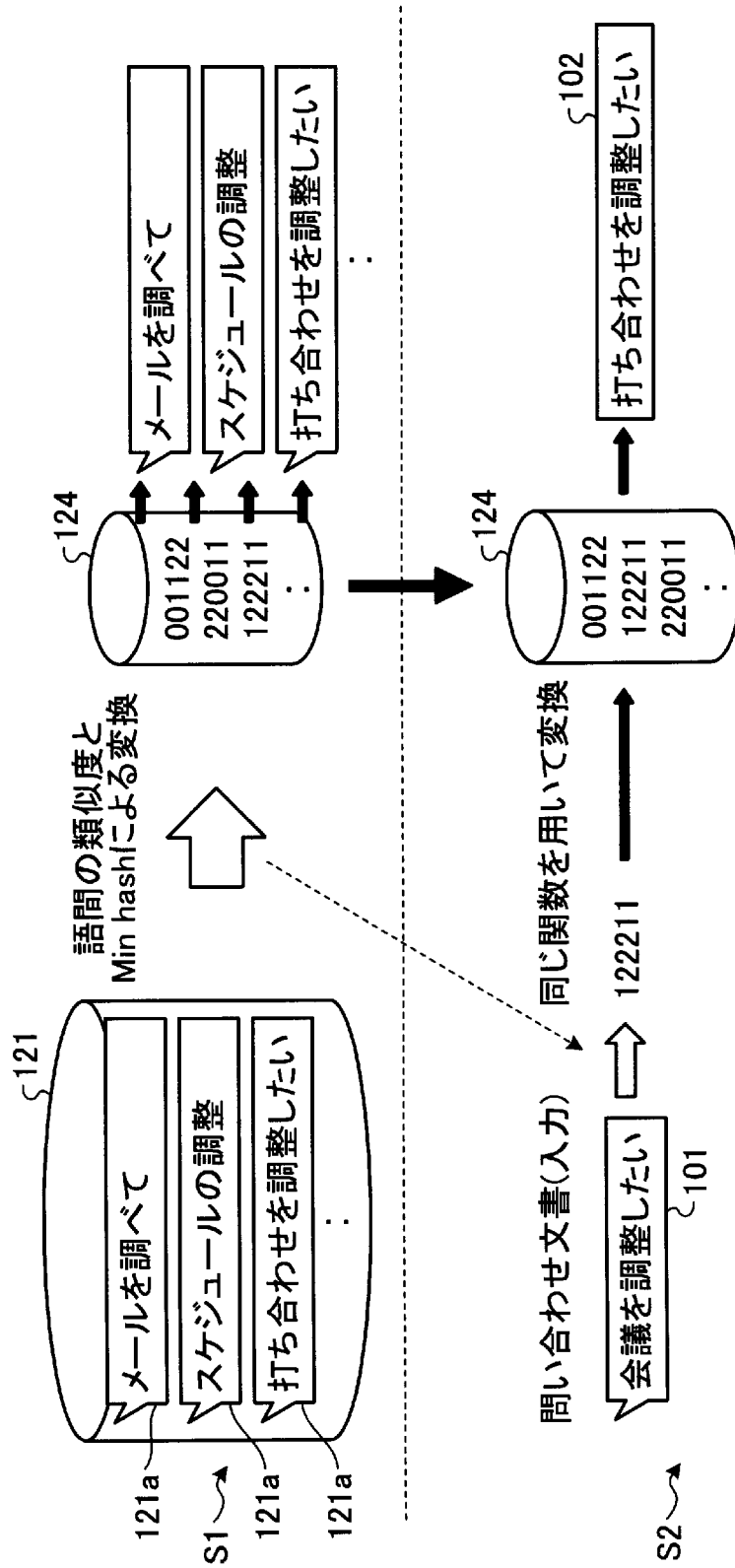
[図5]



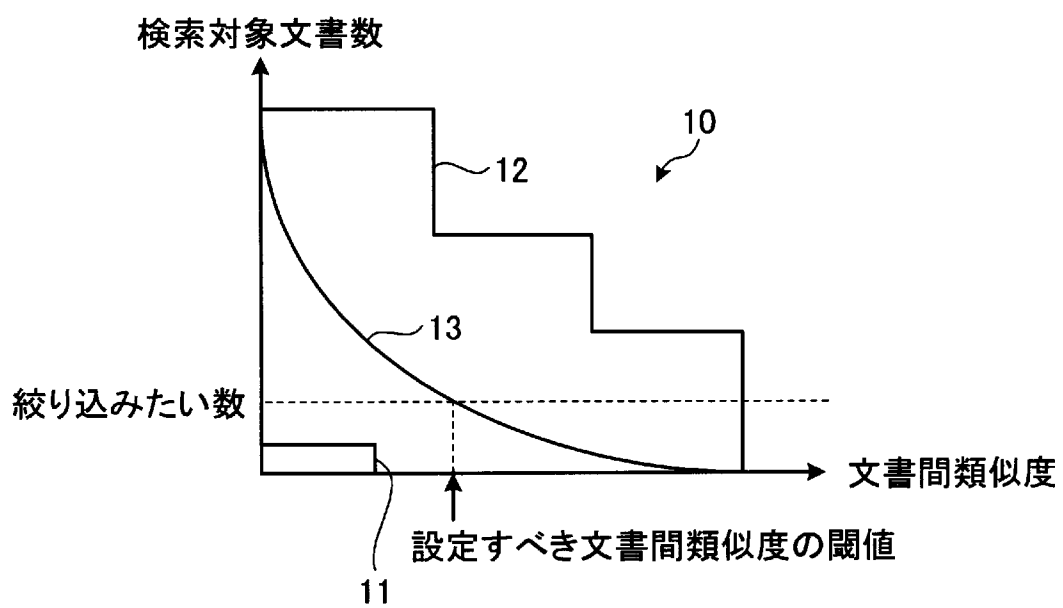
[図6]



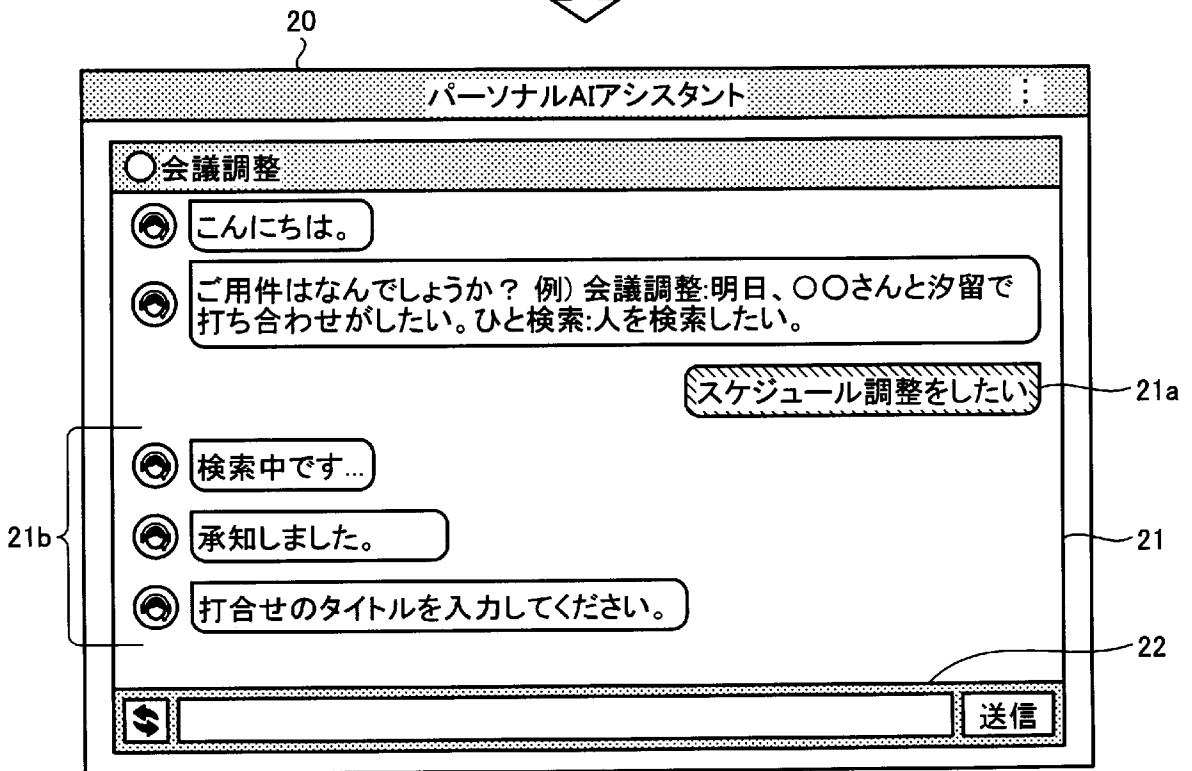
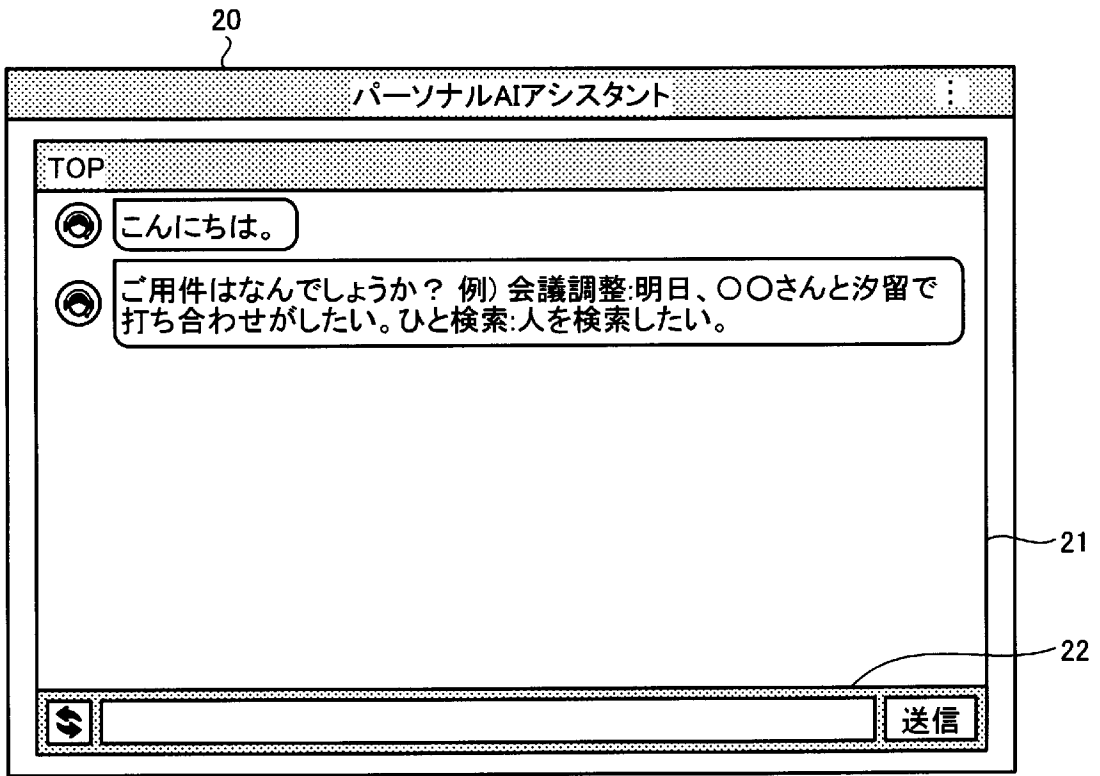
[図7]



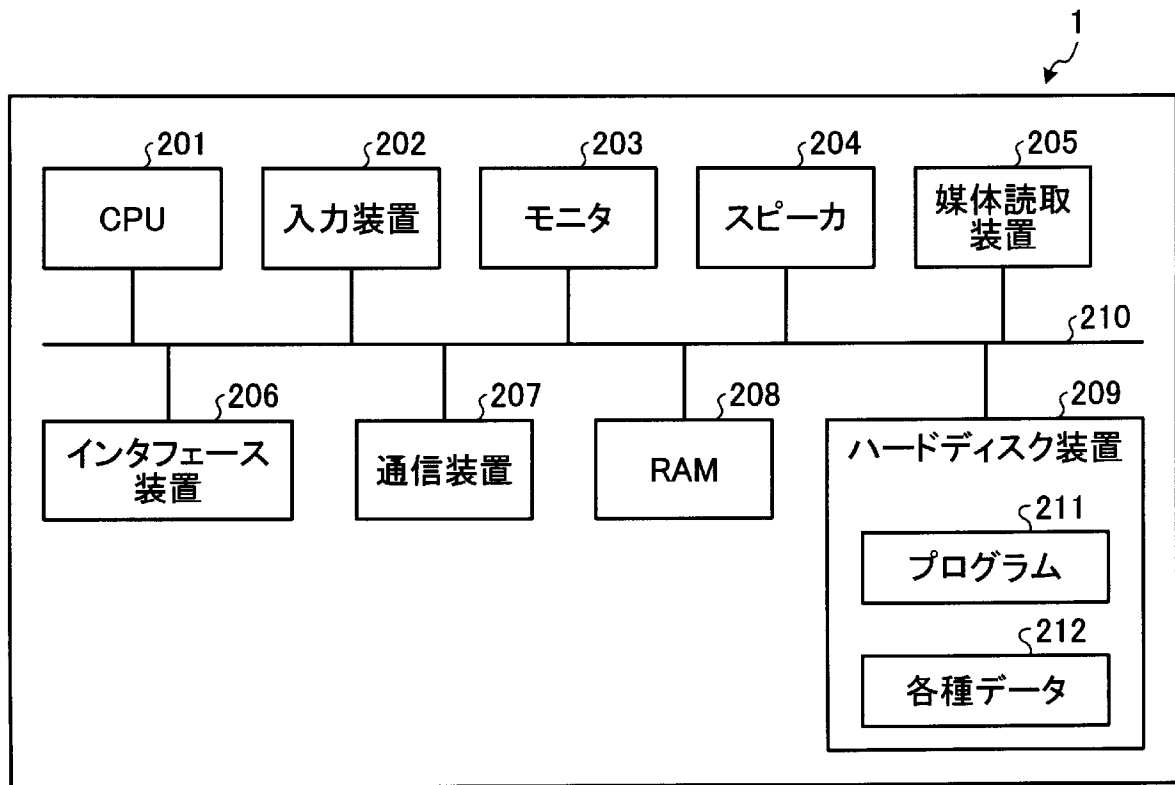
[図8]



[図9]



[図10]



**INTERNATIONAL SEARCH REPORT**

International application No.

PCT/JP2019/034306

**A. CLASSIFICATION OF SUBJECT MATTER**

Int.Cl. G06F16/383(2019.01)i, G06F16/31(2019.01)i, G06F16/90(2019.01)i,  
G06F16/93(2019.01)i

According to International Patent Classification (IPC) or to both national classification and IPC

**B. FIELDS SEARCHED**

Minimum documentation searched (classification system followed by classification symbols)

Int.Cl. G06F16/383, G06F16/31, G06F16/90, G06F16/93

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Published examined utility model applications of Japan	1922-1996
Published unexamined utility model applications of Japan	1971-2019
Registered utility model specifications of Japan	1996-2019
Published registered utility model applications of Japan	1994-2019

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	JP 2010-267108 A (NIPPON TELEGRAPH AND TELEPHONE CORP.) 25 November 2010, entire text, all drawings (Family: none)	1-21
A	JP 2015-201042 A (NIPPON TELEGRAPH AND TELEPHONE CORP.) 12 November 2015, entire text, all drawings (Family: none)	1-21
A	US 2011/0087669 A1 (STRATIFY, INC.) 14 April 2011, entire text, all drawings (Family: none)	1-21

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents:	"I" 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
"A" document defining the general state of the art which is not considered to be of particular relevance	"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
"E" earlier application or patent but published on or after the international filing date	"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
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&" document member of the same patent family
"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	

Date of the actual completion of the international search  
08 October 2019 (08.10.2019)

Date of mailing of the international search report  
15 October 2019 (15.10.2019)

Name and mailing address of the ISA/  
Japan Patent Office  
3-4-3, Kasumigaseki, Chiyoda-ku,  
Tokyo 100-8915, Japan

Authorized officer  
  
Telephone No.

**INTERNATIONAL SEARCH REPORT**

International application No.

PCT/JP2019/034306

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	CN 106156154 A (ALIBABA GROUP HOLDING LTD.) 23 November 2016, entire text, all drawings (Family: none)	1-21
A	CN 107784110 A (BEIJING RUN TECHNOLOGY CO., LTD.) 09 March 2018, entire text, all drawings (Family: none)	1-21

A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int.Cl. G06F16/383(2019.01)i, G06F16/31(2019.01)i, G06F16/90(2019.01)i, G06F16/93(2019.01)i

B. 調査を行った分野

調査を行った最小限資料 (国際特許分類 (IPC))

Int.Cl. G06F16/383, G06F16/31, G06F16/90, G06F16/93

最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公報	1922-1996年
日本国公開実用新案公報	1971-2019年
日本国実用新案登録公報	1996-2019年
日本国登録実用新案公報	1994-2019年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	JP 2010-267108 A (日本電信電話株式会社) 2010.11.25, 全文、全図 (ファミリーなし)	1-21
A	JP 2015-201042 A (日本電信電話株式会社) 2015.11.12, 全文、全図 (ファミリーなし)	1-21
A	US 2011/0087669 A1 (STRATIFY, INC) 2011.04.14, 全文、全図 (ファミリーなし)	1-21

☑ C欄の続きにも文献が列挙されている。

☐ パテントファミリーに関する別紙を参照。

\* 引用文献のカテゴリー

「A」特に関連のある文献ではなく、一般的技術水準を示すもの  
 「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの  
 「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)  
 「O」口頭による開示、使用、展示等に言及する文献  
 「P」国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献

「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの  
 「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの  
 「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの  
 「&」同一パテントファミリー文献

国際調査を完了した日

08.10.2019

国際調査報告の発送日

15.10.2019

国際調査機関の名称及びあて先

日本国特許庁 (ISA/J P)  
 郵便番号 100-8915  
 東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

松尾 真人

5N

8384

電話番号 03-3581-1101 内線 3586

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	CN 106156154 A (ALIBABA GROUP HOLDING LTD) 2016. 11. 23, 全文、全図 (ファミリーなし)	1-21
A	CN 107784110 A (BEIJING RUN TECHNOLOGY CO LTD) 2018. 03. 09, 全文、全図 (ファミリーなし)	1-21