

(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
2023年12月7日(07.12.2023)



(10) 国際公開番号

WO 2023/233633 A1

- (51) 国際特許分類:
G06F 40/166 (2020.01) G06F 40/253 (2020.01)
G06F 40/216 (2020.01) G06F 40/44 (2020.01)
G06F 40/232 (2020.01)
- (21) 国際出願番号: PCT/JP2022/022525
- (22) 国際出願日: 2022年6月2日(02.06.2022)
- (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: 片岡 正弘 (KATAOKA, Masahiro);
〒2118588 神奈川県川崎市中原区上小田中4

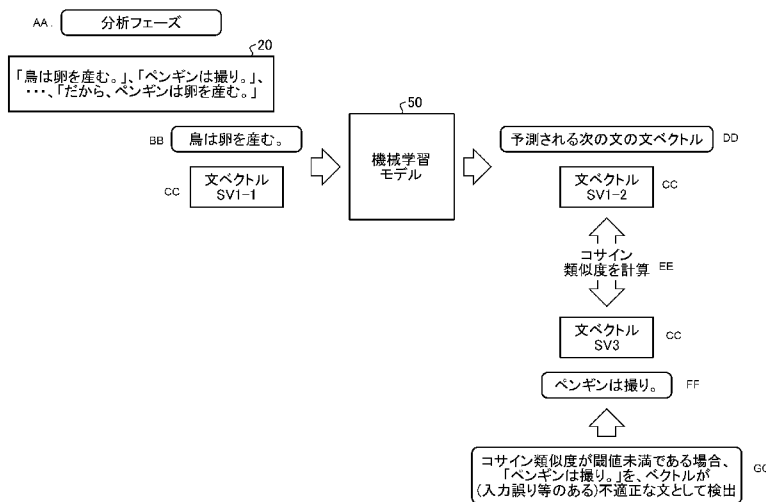
丁目1番1号 富士通株式会社内 Kanagawa (JP). 松村 量 (MATSUMURA, Ryo); 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP). 尾上 聡 (ONOE, Satoshi); 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP).

(74) 代理人: 弁理士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, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN,

(54) Title: INFORMATION PROCESSING PROGRAM, INFORMATION PROCESSING METHOD, AND INFORMATION PROCESSING DEVICE

(54) 発明の名称: 情報処理プログラム、情報処理方法および情報処理装置



20 "Bird lays egg", "Penguin takes photo.", "Therefore, penguin lays egg."
50 Machine-learning model
AA Analysis phase
BB Bird lays egg
CC Sentence vectors
DD Sentence vectors of predicted next sentence
EE Compute cosine similarity
FF Penguin takes photo.
GG When cosine similarity is smaller than threshold, detect "Penguin takes photo." as sentence of improper vectors (having input error, etc.)

(57) Abstract: This information processing device calculates vectors of a plurality of sentences having a relationship in front and back sentences in a plurality of continuous sentences. The information processing device generates a machine-learning model which predicts sentence vectors of a sentence that is input after a certain sentence when the sentence vectors of the certain sentence are input to the machine-learning model by inputting, in order, the plurality of sentence vectors to the machine-learning model and training the machine-learning model. The information processing device cal-



WO 2023/233633 A1

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, IQ, IR, IS, IT, JM, 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, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, WS, 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, SC, 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))

calculates vectors of a first sentence and vectors of a second sentence that follows the first sentence. The information processing device calculates vectors of a sentence that is predicted to follow the first sentence by inputting the vectors of the first sentence to the machine-learning model, and determines whether the vectors of the second sentence are proper.

(57) 要約 : 情報処理装置は、連続する複数の文であって、前後の文に関係性を有する複数の文のベクトルをそれぞれ算出する。情報処理装置は、複数の文のベクトルを順番に機械学習モデルに入力して訓練することで、機械学習モデルにある文のベクトルを入力した際に、ある文の次に入力される文の文ベクトルを予測する機械学習モデルを生成する。情報処理装置は、第1文のベクトルと、第1文に続く第2文のベクトルとを算出する。情報処理装置は、第1文のベクトルを機械学習モデルに入力することで第1文に続くと予測される文のベクトルを算出し、第2文のベクトルが適正であるか否かを判定する。

明 細 書

発明の名称：

情報処理プログラム、情報処理方法および情報処理装置

技術分野

[0001] 本発明は、情報処理プログラム等に関する。

背景技術

[0002] 近年、文のベクトルを算出し、算出したベクトルを利用して、他言語への翻訳、データベースの検索等の各種処理を実行するサービスが提供されている。しかし、ユーザに指定される文自体に入力誤り等が存在すると、文のベクトルを精度よく算出することができず、翻訳、検索等の処理に誤りが生じる場合がある。

[0003] たとえば、適正な文「その機能が特徴である」と、入力誤りの文「その昨日が特徴である」とは、相互に大きく意味の異なる文となり、各文のベクトルも大きく異なる。

[0004] 文の入力誤りを修正する従来技術として、修正履歴から、入力誤りと、その修正文とのペアのデータセットを用いて、学習モデルを訓練しておき、訓練した学習モデルに対象となる文を入力することで、対象となる文の入力誤りを修正する従来技術がある。

先行技術文献

特許文献

[0005] 特許文献1：特開2019-101993号公報

非特許文献

[0006] 非特許文献1：三木一弘、他“BERTを用いた英文空所補充問題の一解法”岡山大学工学部情報系学科、DEIM2020 G2-4 (day1 p47)

非特許文献2：田中佑、他“Wikipediaの修正履歴を用いた日本語入力誤りデータセットの構築”京都大学 大学院情報学科研究科、言語処理学会、第26回年次大会、発表論文集、2020年3月

発明の概要

発明が解決しようとする課題

[0007] しかしながら、上述した従来技術では、文章の一部の単語などがマスクされた空所補充する技術であり、複数の単語で構成される文の空所を補充する単語の精度は高いものの、複数の文で構成される文章の空所を補充する文に関する高精度化の記述は少なく、かつ、入力誤りを含む文を検出するものではなかった。また、従来技術では、誤字、脱字等の入力誤りを修正できるものの、誤変換の入力誤りを正しく修正できないケースが多かった。

[0008] 1つの側面では、本発明は、複数の文で構成される文章の空所を補充する文の推定や、入力誤りを含む文を検出することができる情報処理プログラム、情報処理方法および情報処理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

[0009] 第1の案では、コンピュータに次の処理を実行させる。コンピュータは、連続する複数の文であって、前後の文に係り性を有する複数の文のベクトルをそれぞれ算出する。コンピュータは、複数の文のベクトルを順番に機械学習モデルに入力して訓練することで、機械学習モデルにある文のベクトルを入力した際に、ある文の次に入力される文の文ベクトルを予測する機械学習モデルを生成する。コンピュータは、第1文のベクトルと、第1文に続く第2文のベクトルとを算出する。コンピュータは、第1文のベクトルを前記機械学習モデルに入力することで第1文に続くと予測される文のベクトルを算出し、第2文のベクトルが適正であるか否かを判定する。

発明の効果

[0010] 複数の文で構成される文章の空所を補充する文の推定や、入力誤りを含む文を検出することができる。

図面の簡単な説明

[0011] [図1]図1は、本実施例に係る情報処理装置の学習フェーズの処理を説明するための図である。

[図2]図 2 は、本実施例に係る情報処理装置の分析フェーズの処理を説明するための図である。

[図3]図 3 は、本実施例に係る情報処理装置の構成を示す機能ブロック図である。

[図4]図 4 は、単語ベクトル辞書のデータ構造の一例を示す図である。

[図5A]図 5 A は、文ベクトルを算出する処理を説明するための図（1）である。

[図5B]図 5 B は、文ベクトルを算出する処理を説明するための図（2）である。

[図6]図 6 は、文転置インデックスを生成する処理を説明するための図である。

[図7]図 7 は、本実施例に係る情報処理装置の学習フェーズの処理手順を示すフローチャートである。

[図8]図 8 は、本実施例に係る情報処理装置の分析フェーズの処理手順を示すフローチャートである。

[図9]図 9 は、情報処理装置のその他の処理を説明するための図（1）である。

[図10]図 1 0 は、情報処理装置のその他の処理を説明するための図（2）である。

[図11]図 1 1 は、実施例の情報処理装置と同様の機能を実現するコンピュータのハードウェア構成の一例を示す図である。

発明を実施するための形態

[0012] 以下に、本願の開示する情報処理プログラム、情報処理方法および情報処理装置の実施例を図面に基づいて詳細に説明する。なお、この実施例によりこの発明が限定されるものではない。

実施例

[0013] 本実施例に係る情報処理装置の処理について説明する。情報処理装置は、学習フェーズの処理を実行した後に、分析フェーズの処理を実行する。図 1

は、本実施例に係る情報処理装置の学習フェーズの処理を説明するための図である。

- [0014] 学習フェーズにおいて、情報処理装置は、教師データ141に含まれる複数の文章を用いて、機械学習モデル50の学習を実行する（機械学習モデル50を訓練する）。機械学習モデル50は、BERT（Pre-training of Deep Bidirectional Transformers for Language Understanding）、Next Sentence Prediction、Transformers等のNN（Neural Network）である。
- [0015] 教師データ141に含まれる文章には、複数の文が含まれる。複数の文は、前後の文に所定の関係性を有する。各文は、帰納法や演繹法の三段論法等に基づいて予め設定された文である。
- [0016] たとえば、文章10aには、先頭から順に、文「鳥は卵を産む。」、文「ペンギンは鳥。」、・・・、文「だから、ペンギンは卵を産む。」が含まれる。文章10bには、先頭から順に、文「鳥は卵から生まれる。」、文「鳩は鳥の仲間である。」、・・・、文「従って、鳩は卵から生まれる。」が含まれる。
- [0017] 情報処理装置は、文章10a、10b、その他の文章に含まれる各文の文ベクトルを算出する。たとえば、情報処理装置は、文に対して形態素解析を実行して単語に分割し、各単語のベクトルを積算することで、文ベクトルを算出する。
- [0018] 文章10aの文「鳥は卵を産む。」の文ベクトルを「SV1-1」とする。文「ペンギンは鳥」の文ベクトルを「SV1-2」とする。文「だから、ペンギンは卵を産む。」の文ベクトルを「SV1-3」とする。
- [0019] 文章10bの文「鳥は卵から生まれる。」の文ベクトルを「SV2-1」とする。文「鳩は鳥の仲間である。」の文ベクトルを「SV2-2」とする。文「従って、鳩は卵から生まれる。」の文ベクトルを「SV2-3」とする。
- [0020] 情報処理装置は、文章に含まれる先頭の文のベクトルから順番に、機械学習モデル50に入力する処理を繰り返し実行する。たとえば、情報処理装置

は、文ベクトル「SV1-1」、「SV1-2」、・・・、「SV1-3」の順に、機械学習モデル50に文ベクトルを入力する。情報処理装置は、文ベクトル「SV2-1」、「SV2-2」、・・・、「SV2-3」の順に、機械学習モデル50に文ベクトルを入力する。

[0021] 情報処理装置が、上記の学習フェーズの処理を実行することで、ある第1文の文ベクトルが入力された場合に、第1文の次の第2文の文ベクトルを予測する機械学習モデル50が生成される。

[0022] 図2は、本実施例に係る情報処理装置の分析フェーズの処理を説明するための図である。分析フェーズにおいて、情報処理装置は、訓練済みの機械学習モデル50を用いて、処理対象の文章に含まれる文ベクトルを算出し、コサイン類似度などにより、不適正な文を検出する。

[0023] 図2の説明では、入力誤り等を含む処理対象の文章を文章20とする。文章20には、先頭から順に、文「鳥は卵を産む。」、文「ペンギンは撮り。」、・・・、文「だから、ペンギンは卵を産む。」により構成される。文「ペンギンは撮り。」は、教師データ141の文章10aに含まれる正しい文「ペンギンは鳥。」に対して、単語「鳥」の同音異義語「撮り」の入力誤りを含む文である。

[0024] 情報処理装置は、文「鳥は卵を産む。」の文ベクトル「SV1-1」を算出し、算出した文ベクトル「SV1-1」を、機械学習モデル50に入力することで、文「鳥は卵を産む。」の次の文の文ベクトルを予測する。図2に示す例では、機械学習モデル50によって、文「鳥は卵を産む。」の次の文の文ベクトルとして、「SV1-2」が予測されている。

[0025] 情報処理装置は、文章20に含まれる文であって、文「鳥は卵を産む。」の次の文「ペンギンは撮り」の文ベクトル「SV3」を算出する。

[0026] 情報処理装置は、機械学習モデル50によって予測された次の文の文ベクトル「SV1-2」と、文章20に含まれる文であって、文「鳥は卵を産む。」の次の文「ペンギンは撮り」の文ベクトル「SV3」とのコサイン類似度を算出する。

- [0027] 情報処理装置は、文章10aに含まれる文であって、文「鳥は卵を産む。」の次の文「ペンギンは鳥。」は、コサイン類似度が閾値未満の場合として、正しい（以下、適正な、と表記する）文であると判定する。一方、情報処理装置は、文章20に含まれる文であって、文「鳥は卵を産む。」の次の文「ペンギンは撮り。」は、コサイン類似度が閾値未満の場合として、入力誤り等を含む不適正な文であると判定する。
- [0028] 上記のように、情報処理装置は、教師データ141に含まれる文章の各文のベクトルを順番に機械学習モデル50に入力することで、ある第1文の文ベクトルが入力された場合に、第1文の次の第2文の文ベクトルを予測する機械学習モデル50を生成する。情報処理装置は、生成した機械学習モデルに、処理対象の文章の文の文ベクトルを入力し、次の文の文ベクトルを予測し、予測した文ベクトルを基にして、処理対象の文章から、入力誤りのある文を検出する。すなわち、処理対象の文章に含まれる各文から、入力誤り等を含み、不適正な文ベクトルを持つ文を検出することができる。
- [0029] なお、情報処理装置は、図2の処理において、文「ペンギンは撮り」が不適正な文ベクトルの文であると判定した場合に、機械学習モデル50によって予測された文ベクトルSV1-2を基にして、適正な文ベクトルの文「ペンギンは鳥。」をDB (Data Base) などから検索して、正しい修正候補として表示装置に出力（以下、適正化と表記する）してもよい。
- [0030] 更に、情報処理装置は、単語単位のベクトルの順番を学習した他の機械学習モデルに、不適正な文ベクトルを検出した文「ペンギンは撮り」を構成する複数の単語「ペンギン」、「は」、「撮り。」の各単語ベクトルを算出し、乖離した単語「撮り。」の入力誤り等を適正化してもよい。
- [0031] 次に、図1及び図2で説明した処理を実行する情報処理装置の構成例について説明する。図3は、本実施例に係る情報処理装置の構成を示す機能ブロック図である。図3に示すように、情報処理装置100は、通信部110、入力部120、表示部130、記憶部140、制御部150を有する。
- [0032] 通信部110は、有線又は無線で外部装置等に接続され、外部装置等との

間で情報の送受信を行う。たとえば、通信部110は、NIC (Network Interface Card) 等によって実現される。通信部110は、図示しないネットワークに接続されていてもよい。

[0033] 入力部120は、各種の情報を、情報処理装置100に入力する入力装置である。入力部120は、キーボードやマウス、タッチパネル等に対応する。たとえば、ユーザは、入力部120を操作して、文章のデータ等を入力してもよい。

[0034] 表示部130は、制御部150から出力される情報を表示する表示装置である。表示部130は、液晶ディスプレイ、有機EL (Electro Luminescence) ディスプレイ、タッチパネル等に対応する。たとえば、入力誤りのある文が、表示部130に表示される。

[0035] 記憶部140は、機械学習モデル50、教師データ141、単語ベクトル辞書142を有する。記憶部140は、たとえば、RAM (Random Access Memory)、フラッシュメモリ (Flash Memory) 等の半導体メモリ素子、または、ハードディスク、光ディスク等の記憶装置によって実現される。

[0036] 機械学習モデル50は、図1で説明したBERT、Next Sentence Prediction、Transformers等のNN等である。

[0037] 教師データ141は、図1で説明した教師データ141である。教師データ141に含まれる文章には、複数の文が含まれる。複数の文は、前後の文に所定の関係性を有する。各文は、帰納法や演繹法の三段論法等に基づいて予め設定された文である。

[0038] 単語ベクトル辞書142は、単語に割り当てられた符号、単語ベクトルを定義するテーブルである。図4は、単語ベクトル辞書のデータ構造の一例を示す図である。図4に示すように、この単語ベクトル辞書142は、符号、単語、単語ベクトル(1)~(7)を有する。符号は、単語に割り当てられる符号(Code)である。単語は、文字列に含まれる単語である。単語ベクトル(1)~(7)は、単語に割り当てられたベクトルである。単語ベクトルの第n成分を単語ベクトル(n)と表記する(n=1~7)。

- [0039] DB 143は、様々な文章を有する。文章には複数の文が含まれ、各文には複数の単語が含まれる。DB 143は、教師データ141に含まれる文章を有していてもよい。
- [0040] 文転置インデックス144は、文ベクトルと、位置ポインタとを対応付ける。位置ポインタは、文ベクトルに対応する文が存在するDB 143の位置を示す。
- [0041] 図3の説明に戻る。制御部150は、前処理部151と、学習部152と、分析部153とを有する。制御部150は、たとえば、CPU (Central Processing Unit) やMPU (Micro Processing Unit) により実現される。また、制御部150は、例えばASIC (Application Specific Integrated Circuit) やFPGA (Field Programmable Gate Array) 等の集積回路により実行されてもよい。
- [0042] 前処理部151は、各種の前処理を実行する。たとえば、前処理部151は、DB 143から未処理に文を取得し、文の文ベクトルを算出する。前処理部151は、算出した文ベクトルと、文ベクトルに対応する文の位置ポインタとの関係を、文転置インデックス144に設定する。
- [0043] 前処理部151が、文の文ベクトルを算出する処理の一例について説明する。図5A及び図5Bは、文ベクトルを算出する処理を説明するための図である。ここでは、文「馬は人参が好きです。」の文ベクトルを算出する場合について説明する。前処理部151は、文「馬は人参が好きです。」に対して、形態素解析を実行することで、複数の単語に分解する。分解した各単語には、「△ (スペース)」を付与する。たとえば、文1「馬は人参が好きです。」は、「馬△」、「は△」、「人参△」、「が△」、「好き△」、「です△」、「。△」に分割する。
- [0044] 前処理部151は、分割した各単語と、単語ベクトル辞書45とを比較することで、各単語に対応する符号を特定し、単語と置き換える。たとえば、各単語「馬△」、「は△」、「人参△」、「が△」、「好き△」、「です△」、「。△」は、それぞれ、「C1」、「C2」、「C3」、「C4」、「

C 5」、「C 6」、「C 7」に置き換えられる。

[0045] 図5 Bの説明に移行する。前処理部151は、単語ベクトル辞書45と、各符号とを基にして、符号に割り当てられた単語ベクトル(1)～(7)を特定する。たとえば、符号「C 1」の単語ベクトル(1)～(7)は、 $wv_{1-1} \sim 1-7$ とする。符号「C 2」の単語ベクトル(1)～(7)は、 $wv_{2-1} \sim 2-7$ とする。符号「C 3」の単語ベクトル(1)～(7)は、 $wv_{3-1} \sim 3-7$ とする。

[0046] 符号「C 4」の単語ベクトル(1)～(7)は、 $wv_{4-1} \sim 4-7$ とする。符号「C 5」の単語ベクトル(1)～(7)は、 $wv_{5-1} \sim 5-7$ とする。符号「C 6」の単語ベクトル(1)～(7)は、 $wv_{6-1} \sim 6-7$ とする。符号「C 7」の単語ベクトル(1)～(7)は、 $wv_{7-1} \sim 7-7$ とする。

[0047] 前処理部151は、要素毎に単語ベクトルを積算することで、文の文ベクトルSV1を算出する。たとえば、前処理部151は、各単語ベクトル(1)となる $wv_{1-1} \sim 7-1$ を積算することで、文ベクトルSV1の第1成分「SV1-1」を算出する。前処理部151は、各単語ベクトル(2)となる $wv_{1-2} \sim 7-2$ を積算することで、文ベクトルSV1の第2成分「SV1-2」を算出する。各単語ベクトル(3)となる $wv_{1-3} \sim 7-3$ を積算することで、文ベクトルSV1の第3成分「SV1-3」を算出する。

[0048] 前処理部151は、各単語ベクトル(4)となる $wv_{1-4} \sim 7-4$ を積算することで、文ベクトルSV1の第4成分「SV1-4」を算出する。前処理部151は、各単語ベクトル(5)となる $wv_{1-5} \sim 7-5$ を積算することで、文ベクトルSV1の第5成分「SV1-5」を算出する。前処理部151は、各単語ベクトル(6)となる $wv_{1-6} \sim 7-6$ を積算することで、文ベクトルSV1の第6成分「SV1-6」を算出する。前処理部151は、各単語ベクトル(7)となる $wv_{1-7} \sim 7-7$ を積算することで、文ベクトルSV1の第7成分「SV1-7」を算出する。

- [0049] 前処理部151は、DB143に含まれる他の文章の各文についても、上記処理を繰り返し実行することで、各文の文ベクトルを算出する。
- [0050] 前処理部151は、算出した各文の文ベクトルと、DB143の位置ポイントとを対応付けることで、文転置インデックス144を生成する。なお、前処理部151は、図6に示すようなデータ構造の文転置インデックス144を生成してもよい。図6は、文転置インデックスを生成する処理を説明するための図である。図6に示すように、前処理部151は、文ベクトルと、複数のレコードポイントと、複数の位置ポイントとを対応付け、各レコードポイント、位置ポイントを、DB143の各文に対応付けてもよい。
- [0051] 図3の説明に戻る。学習部152は、図1で説明した学習フェーズの処理を実行することで、ある第1文の文ベクトルが入力された場合に、第1文の次の第2文の文ベクトルを予測する機械学習モデル50を生成する。
- [0052] たとえば、学習部152は、教師データ141の文章に含まれる各文の文ベクトルを算出し、算出した文ベクトルを順番に機械学習モデル50に入力することで、機械学習モデル50の学習を実行する。学習部152のその他の処理は、図1で説明した処理と同様である。学習部152が、文の文ベクトルを算出する処理は、前処理部151が文の文ベクトルを算出する処理と同様である。
- [0053] 分析部153は、図2で説明した分析フェーズの処理を実行することで、処理対象の文章に含まれる文から、文ベクトルが不適正な文を検出する。
- [0054] たとえば、分析部153は、処理対象の文章20を受け付けた場合に、文章20に含まれる文の文ベクトルを算出する。分析部153は、文章20に含まれる句点「。」を基にして、文章20に含まれる文を特定する。分析部153が、文の文ベクトルを算出する処理は、前処理部151が文の文ベクトルを算出する処理と同様である。文章20の先頭からn番目の文の文ベクトル「SV_n」と表記する（n=0~M）。
- [0055] 分析部153は、文ベクトルSV_nを、訓練済みの機械学習モデル50に入力して、文章20の先頭からn+1番目の文の文ベクトルSV_{n+1}'を

予測する。分析部153は、機械学習モデル50を用いて予測した文ベクトル SV_{n+1}' と、文のベクトル SV_{n+1} とのコサイン類似度を算出する。

[0056] 分析部153は、文ベクトル SV_{n+1}' と、文ベクトル SV_{n+1} とのコサイン類似度が閾値以上の場合には、先頭から $n+1$ 番目の文が適正な文であると判定する。一方、分析部153は、文ベクトル SV_{n+1}' と、文ベクトル SV_{n+1} とのコサイン類似度が閾値以上の場合には、先頭から $n+1$ 番目の文が、文ベクトルの不適正な文であると判定する。

[0057] 分析部153は、文ベクトルの不適正な文であると判定した場合に、文ベクトル SV_{n+1}' と、文転置インデックス144とを比較して、文ベクトル SV_{n+1}' に対応する文の位置ポイントを特定する。分析部153は、位置ポイントを基にして、文ベクトル SV_{n+1}' に対応する文を、DB143から検索する。分析部153は、文ベクトルの不適正な文と、検索した文とを対応付けて、表示部130に表示させる。

[0058] 分析部153は、文ベクトルの不適正な文と、検索した文とを単語単位に比較して、文ベクトルの不適正な文から、入力誤りの単語を検出し、検出した単語を表示させてもよい。

[0059] 次に、本実施例に係る情報処理装置100の処理手順の一例について説明する。図7は、本実施例に係る情報処理装置の学習フェーズの処理手順を示すフローチャートである。図7に示すように、情報処理装置100の学習部152は、教師データ141から、未選択の文章を選択する（ステップS101）。

[0060] 学習部152は、選択した文章に含まれる各文の文ベクトルをそれぞれ算出し、文ベクトルとDBのレコードと文の位置とを対応付けた文転置インデックスを生成する（ステップS102）。学習部152は、選択した文章に含まれる先頭の文の文ベクトルから順に、機械学習モデル50に入力することで、学習を実行する（ステップS103）。

[0061] 学習部152は、学習を継続する場合には（ステップS104, Yes）

、ステップS101に移行する。一方、学習部152は、学習を継続しない場合には（ステップS104, No）、学習フェーズの処理を終了する。

[0062] 図8は、本実施例に係る情報処理装置の分析フェーズの処理手順を示すフローチャートである。図8に示すように、情報処理装置100の分析部153は、処理対象の文章の入力を受け付ける（ステップS201）。

[0063] 分析部153は、入力された文章に含まれる各文の文ベクトルをそれぞれ算出する（ステップS202）。分析部153は、 n を初期値に設定する（ステップS203）。

[0064] 分析部153は、文章に含まれる複数の文のうち、 n 番目の文の文ベクトル SV_n を機械学習モデル50に入力し、 $n+1$ 番目の文の文ベクトル SV_{n+1}' を予測する（ステップS204）。

[0065] 分析部153は、文章に含まれる複数の文のうち、 $n+1$ 番目の文の文ベクトル SV_{n+1} と、予測した文の文ベクトル SV_{n+1}' とのコサイン類似度を算出する（ステップS205）。

[0066] 分析部153は、コサイン類似度が閾値以上である場合には（ステップS206, Yes）、ステップS210に移行する。

[0067] 一方、分析部153は、コサイン類似度が閾値以上でない場合には（ステップS206, No）、 $n+1$ 番目の文を、文ベクトルが不適正な文として検出する（ステップS207）。分析部153は、予測した文ベクトル SV_{n+1}' と、文転置インデックス144とを基にして、文ベクトル SV_{n+1}' に対応する文をDB143から検出する（ステップS208）。

[0068] 分析部153は、文ベクトルが不適正な文と、DB143から検出した文とを、表示部130に表示する（ステップS209）。

[0069] ステップS210以降の処理について説明する。分析部153は、 n が L 以上である場合には（ステップS210, Yes）、処理を終了する。 L は、処理対象の文章に含まれる文の数である。分析部153は、 n が L 以上でない場合には（ステップS210, No）、 n に1を加算した値によって、 n を更新し（ステップS211）、ステップS204に移行する。

- [0070] 次に、本実施例に係る情報処理装置100の効果について説明する。情報処理装置100は、教師データ141に含まれる文章の各文のベクトルを順番に機械学習モデル50に入力することで、ある第1文の文ベクトルが入力された場合に、第1文の次の第2文の文ベクトルを予測する機械学習モデル50を生成する。情報処理装置100は、生成した機械学習モデル50に、処理対象の文章の文の文ベクトルを入力し、次の文の文ベクトルを予測し、予測した文ベクトルを基にして、処理対象の文章から、不適正な文ベクトルを持つ文を検出する。また、その不適正な文から入力誤り等の単語を適正化することができる。
- [0071] 情報処理装置100は、機械学習モデル50によって予測された次の文の文ベクトルと、処理対象の文章に含まれる文の次の文の文ベクトルとのコサイン類似度を基にして、文ベクトルの不適正な文を検出し、入力誤り等を適正化する。これによって、計算コストを抑えて、文ベクトルの不適正な文を検出し、入力誤り等を適正化することができる。
- [0072] 情報処理装置100は、帰納法または演繹法に基づいて並び順が決定された複数の文のベクトルを順番に機械学習モデルに入力して訓練する。これによって、帰納法または演繹法に基づいた対象文の次の文を予測することができる。
- [0073] 情報処理装置100は、修正対象の文であると判定された場合に、機械学習モデル50に予測されたベクトルを基にして、修正後の文を検索する。これによって、修正後の文を通知することができる。
- [0074] なお、上述した情報処理装置100の処理内容は一例であり、情報処理装置100は、その他の処理を実行してもよい。以下では、情報処理装置100のその他の処理について説明する。
- [0075] 図9及び図10は、情報処理装置のその他の処理を説明するための図である。上記の情報処理装置100は、機械学習モデル50に、三段論法に基づく文のベクトルの順番を学習させていたが、文のベクトルの代わりに、タンパク質の配列であり、単語に相当する複数のアミノ酸配列で構成されるタン

パク質一次構造のベクトルの順番を学習させてもよい。以下の説明では、タンパク質の連続アミノ酸配列を「基本構造」と、タンパク質一次構造を「一次構造」と表記する。

- [0076] 図9について説明する。学習フェーズにおいて、情報処理装置100は、教師データ241に含まれる複数のタンパク質の配列20a, 20bを用いて、機械学習モデル50の学習を実行する。
- [0077] たとえば、配列20aには、一次構造「 α 一次構造」、「 β 一次構造」、 \dots 、「 γ 一次構造」が含まれる。配列20bには、一次構造「 Δ 一次構造」、「 ε 一次構造」、 \dots 、「 ζ 一次構造」が含まれる。
- [0078] 情報処理装置100は、基本構造と、ベクトルとを対応付けたタンパク質の基本構造のベクトル辞書を用いて、各一次構造のベクトルを特定する。たとえば、複数の基本構造で構成される一次構造「 α 一次構造」のベクトルを「V20-1」、一次構造「 β 一次構造」のベクトルを「V20-2」、一次構造「 γ 一次構造」のベクトルを「V20-3」とする。一次構造「 Δ 一次構造」のベクトルを「V21-1」、一次構造「 ε 一次構造」のベクトルを「V21-2」、一次構造「 ζ 一次構造」のベクトルを「V21-3」とする。各一次構造のベクトルは、その一次構造を構成する複数の基本構造の各基本構造のベクトルをもとに算出される。
- [0079] 情報処理装置100は、タンパク質の配列に含まれる先頭の一次構造のベクトルから順番に、機械学習モデル50に入力する処理を繰り返し実行する。たとえば、情報処理装置は、ベクトル「V20-1」、「V20-2」、 \dots 、「V20-3」の順に、機械学習モデル50にベクトルを入力する。情報処理装置は、ベクトル「V21-1」、「V21-2」、 \dots 、「V21-3」の順に、機械学習モデル50にベクトルを入力する。
- [0080] 情報処理装置100が、上記の学習フェーズの処理を実行することで、ある一次構造のベクトルが入力された場合に、ある一次構造の次の一次構造のベクトルを予測する機械学習モデル50が生成される。
- [0081] 図10について説明する。分析フェーズにおいて、処理対象のタンパク質

の配列を、配列 25 とする。配列 25 には、先頭から順に、一次構造「 α 一次構造」、「 η 一次構造」、 \dots 、「 γ 一次構造」が含まれる。

[0082] 情報処理装置 100 は、一次構造「 α 一次構造」のベクトル「V20-1」を算出し、算出したベクトル「V20-1」を、機械学習モデル 50 に入力することで、一次構造「 α 一次構造」の次の一次構造のベクトルを予測する。図 10 に示す例では、機械学習モデル 50 によって、一次構造「 α 一次構造」の次の一次構造のベクトルとして、「V20-2」が予測されている。

[0083] 情報処理装置 100 は、配列 25 に含まれる一次構造であって、一次構造「 α 一次構造」の次の「 η 一次構造」のベクトル「V22」を算出する。

[0084] 情報処理装置 100 は、機械学習モデル 50 によって予測された次の一次構造のベクトル「V20-2」と、配列 25 に含まれる一次構造であって、基本構造「 α 一次構造」の次の「 η 一次構造」のベクトル「V22」とのコサイン類似度を算出する。

[0085] 情報処理装置は、コサイン類似度が閾値以上の場合、配列 25 に含まれる一次構造であって、一次構造「 α 一次構造」の次の「 η 一次構造」が正しい一次構造であると判定する。一方、情報処理装置は、コサイン類似度が閾値未満の場合、配列 25 に含まれる一次構造であって、一次構造「 α 一次構造」の次の「 η 一次構造」を不適正な一次構造であると判定し、一次構造「 η 一次構造」に含まれる基本構造の突然変異等を適正化する。

[0086] 図 9 及び図 10 に示した処理を、情報処理装置 100 が実行することで、タンパク質の配列に含まれる複数の一次構造から、不適正な一次構造ベクトルを持つ一次構造を検出し、突然変異等のある基本構造を適正化することができる。これにより、複数のタンパク質一次構造で構成される受容体に発生した突然変異等（SNPs が代表例である）を持つタンパク質一次構造を検出することができる。さらに、受容体を構成する多数のタンパク質一次構造と、受容体に結合する単一または複数のタンパク質一次構造を結合順に機械学習することで、受容体に結合するリガンドのタンパク質一次構造のベクトル

ルを予測することができる。これにより、既に、バイオ医薬品として製品化されたリガンドと類似し、優れた薬効を持ち、副反応が抑制された、新しいタンパク質一次構造のベクトルを持つリガンドの改良を支援することができる。

[0087] 次に、上記実施例に示した情報処理装置100と同様の機能を実現するコンピュータのハードウェア構成の一例について説明する。図11は、実施例の情報処理装置と同様の機能を実現するコンピュータのハードウェア構成の一例を示す図である。

[0088] 図11に示すように、コンピュータ300は、各種演算処理を実行するCPU301と、ユーザからのデータの入力を受け付ける入力装置302と、ディスプレイ303とを有する。また、コンピュータ300は、有線または無線ネットワークを介して、外部装置等との間でデータの授受を行う通信装置304と、インタフェース装置305とを有する。また、コンピュータ300は、各種情報を一時記憶するRAM306と、ハードディスク装置307とを有する。そして、各装置301~307は、バス308に接続される。

[0089] ハードディスク装置307は、前処理プログラム307a、学習プログラム307b、分析プログラム307cを有する。また、CPU301は、各プログラム307a~307cを読み出してRAM306に展開する。

[0090] 前処理プログラム307aは、前処理プロセス306aとして機能する。学習プログラム307bは、学習プロセス306bとして機能する。分析プログラム307cは、分析プロセス306cとして機能する。

[0091] 前処理プロセス306aの処理は、前処理部151の処理に対応する。学習プロセス306bの処理は、学習部152の処理に対応する。分析プロセス306cの処理は、分析部153の処理に対応する。

[0092] なお、各プログラム307a~307cについては、必ずしも最初からハードディスク装置307に記憶させておかなくても良い。例えば、コンピュータ300に挿入されるフレキシブルディスク(FD)、CD-ROM、D

V D、光磁気ディスク、I Cカードなどの「可搬用の物理媒体」に各プログラムを記憶させておく。そして、コンピュータ300が各プログラム307 a~307 cを読み出して実行するようにしてもよい。

符号の説明

[0093]	50	機械学習モデル
	100	情報処理装置
	110	通信部
	120	入力部
	130	表示部
	140	記憶部
	141	教師データ
	142	単語ベクトル辞書
	143	D B
	144	文転置インデックス
	150	制御部
	151	前処理部
	152	学習部
	153	分析部

請求の範囲

- [請求項1] 連続する複数の文であって、前後の文に係性を有する前記複数の文のベクトルをそれぞれ算出し、
- 前記複数の文のベクトルを順番に機械学習モデルに入力して訓練することで、前記機械学習モデルにある文のベクトルを入力した際に、前記ある文の次に入力される文の文ベクトルを予測する前記機械学習モデルを生成し、
- 第1文のベクトルと、前記第1文に続く第2文のベクトルとを算出し、
- 前記第1文のベクトルを前記機械学習モデルに入力することで前記第1文に続くと予測される文のベクトルを算出し、前記第2文のベクトルが適正であるか否かを判定する
- 処理をコンピュータに実行させることを特徴とする情報処理プログラム。
- [請求項2] 前記判定する処理は、前記第1文のベクトルを、前記機械学習モデルに入力することで予測されるベクトルと、前記第2文のベクトルとのコサイン類似度を基にして、前記第2文のベクトルが適正であるか否かを判定することを特徴とする請求項1に記載の情報処理プログラム。
- [請求項3] 前記連続する複数の文は、帰納法または演繹法に基づいて並び順が決定された複数の文であり、前記機械学習モデルを生成する処理は、帰納法または演繹法に基づいて並び順が決定された複数の文のベクトルを順番に機械学習モデルに入力して訓練することを特徴とする請求項2に記載の情報処理プログラム。
- [請求項4] 前記第2文のベクトルが不適正であると判定された場合に、前記第1文のベクトルを前記機械学習モデルに入力することで前記第1文に続くと予測される文のベクトルを算出し、算出したベクトルと類似した文を検索し適正な文の候補として提示するために、算出した前記第

1文に続くと予測される文のベクトルを基にして、適正な文を推薦する処理を更にコンピュータに実行させることを特徴とする請求項1に記載の情報処理プログラム。

[請求項5] 連続する複数の文であって、前後の文に関係性を有する前記複数の文のベクトルをそれぞれ算出し、

前記複数の文のベクトルを順番に機械学習モデルに入力して訓練することで、前記機械学習モデルにある文のベクトルを入力した際に、前記ある文の次に入力される文の文ベクトルを予測する前記機械学習モデルを生成し、

第1文のベクトルと、前記第1文に続く第2文のベクトルとを算出し、

前記第1文のベクトルを前記機械学習モデルに入力することで前記第1文に続くと予測される文のベクトルを算出し、前記第2文のベクトルが適正であるか否かを判定する

処理をコンピュータが実行することを特徴とする情報処理方法。

[請求項6] 前記判定する処理は、前記第1文のベクトルを、前記機械学習モデルに入力することで予測されるベクトルと、前記第2文のベクトルとのコサイン類似度を基にして、前記第2文のベクトルが適正であるか否かを判定することを特徴とする請求項5に記載の情報処理方法。

[請求項7] 前記連続する複数の文は、帰納法または演繹法に基づいて並び順が決定された複数の文であり、前記機械学習モデルを生成する処理は、帰納法または演繹法に基づいて並び順が決定された複数の文のベクトルを順番に機械学習モデルに入力して訓練することを特徴とする請求項6に記載の情報処理方法。

[請求項8] 前記第2文のベクトルが不適正であると判定された場合に、前記第1文のベクトルを前記機械学習モデルに入力することで前記第1文に続くと予測される文のベクトルを算出し、算出したベクトルと類似した文を検索し適正な文の候補として提示するために、算出した前記第

1文に続く予測される文のベクトルを基にして、適正な文を推薦する処理を更にコンピュータに実行させることを特徴とする請求項6に記載の情報処理方法。

[請求項9] 連続する複数の文であって、前後の文に係性を有する前記複数の文のベクトルをそれぞれ算出し、

前記複数の文のベクトルを順番に機械学習モデルに入力して訓練することで、前記機械学習モデルにある文のベクトルを入力した際に、前記ある文の次に入力される文の文ベクトルを予測する前記機械学習モデルを生成し、

第1文のベクトルと、前記第1文に続く第2文のベクトルとを算出し、

前記第1文のベクトルを前記機械学習モデルに入力することで前記第1文に続く予測される文のベクトルを算出し、前記第2文のベクトルが適正であるか否かを判定する

処理を実行する制御部を有することを特徴とする情報処理装置。

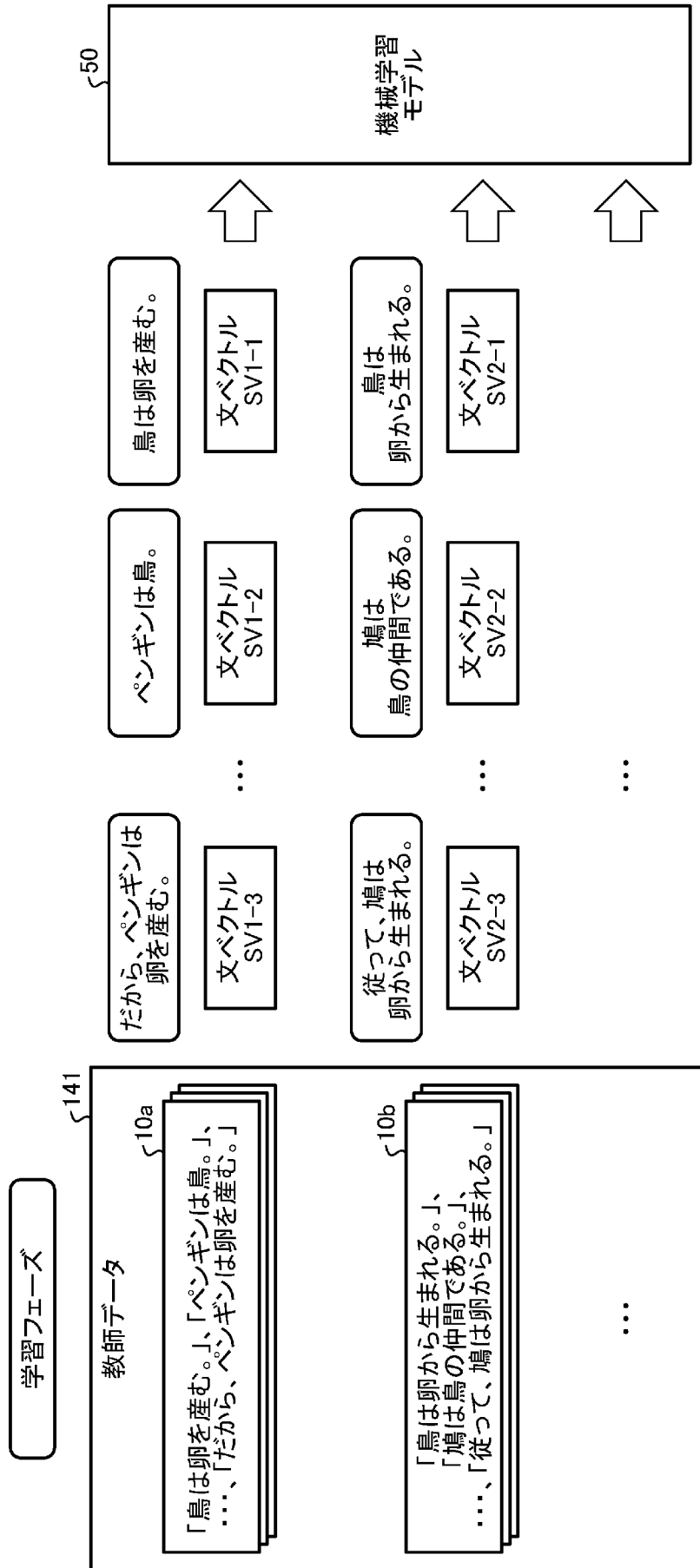
[請求項10] 前記制御部が実行する前記判定する処理は、前記第1文のベクトルを、前記機械学習モデルに入力することで予測されるベクトルと、前記第2文のベクトルとのコサイン類似度を基にして、前記第2文のベクトルが適正であるか否かを判定することを特徴とする請求項9に記載の情報処理装置。

[請求項11] 前記連続する複数の文は、帰納法または演繹法に基づいて並び順が決定された複数の文であり、前記制御部が実行する前記機械学習モデルを生成する処理は、帰納法または演繹法に基づいて並び順が決定された複数の文のベクトルを順番に機械学習モデルに入力して訓練することを特徴とする請求項10に記載の情報処理装置。

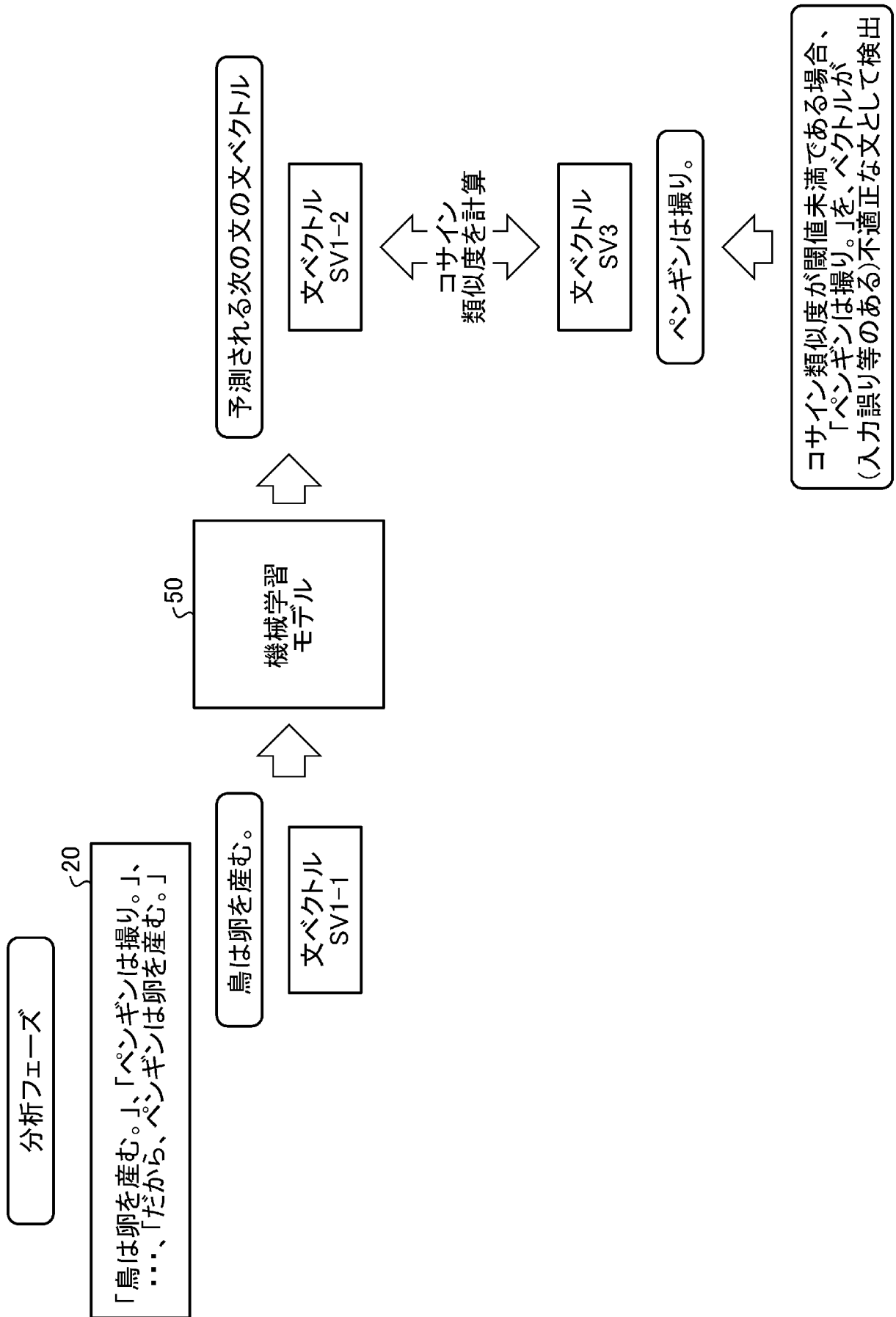
[請求項12] 前記制御部は、前記第2文のベクトルが不適正であると判定された場合に、前記第1文のベクトルを前記機械学習モデルに入力することで前記第1文に続く予測される文のベクトルを算出し、算出したベ

クトルと類似した文を検索し適正な文の候補として提示するために、算出した前記第1文に続くと予測される文のベクトルを基にして、適正な文を推薦する処理を更にコンピュータに実行させることを特徴とする請求項9に記載の情報処理装置。

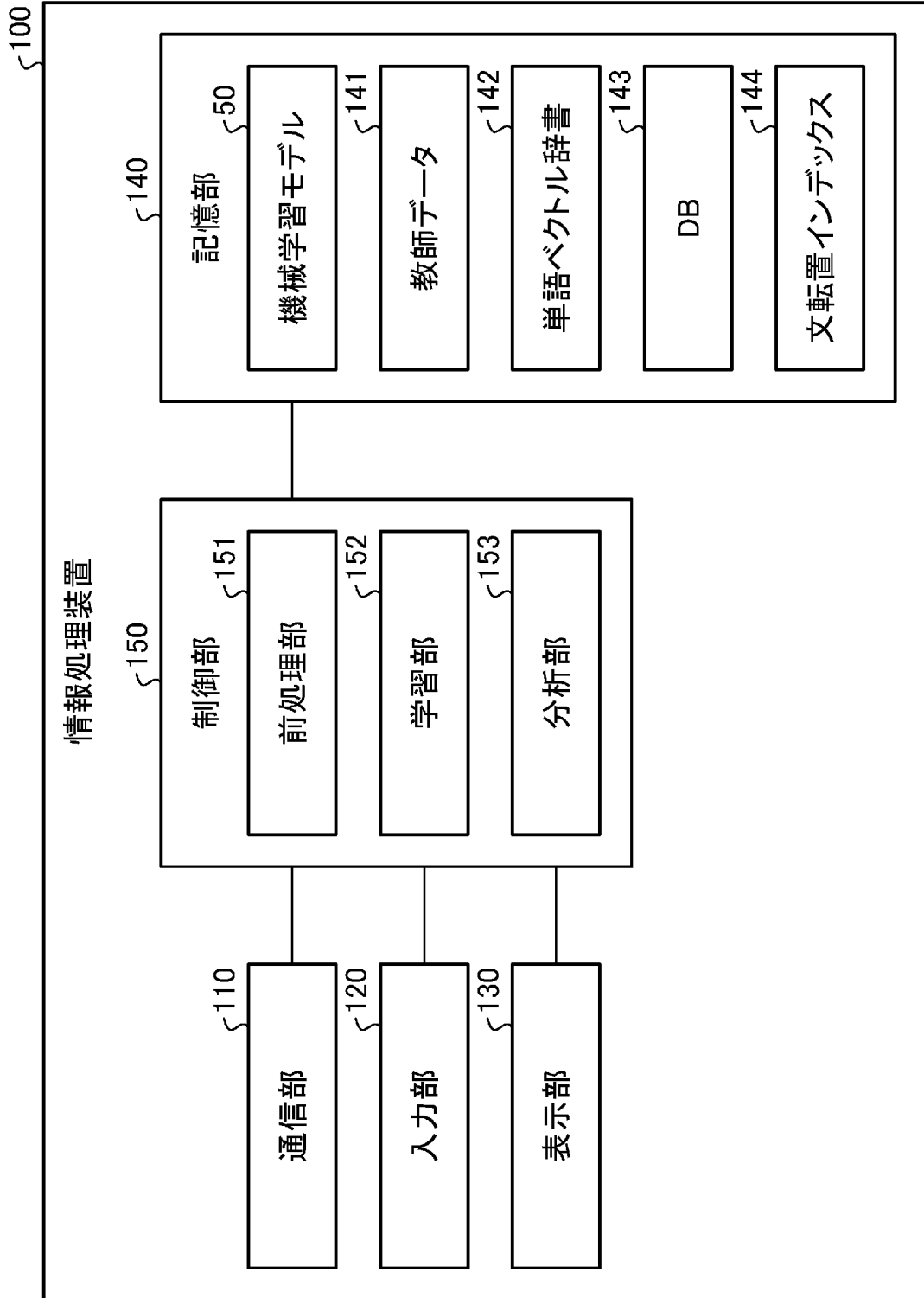
[図1]



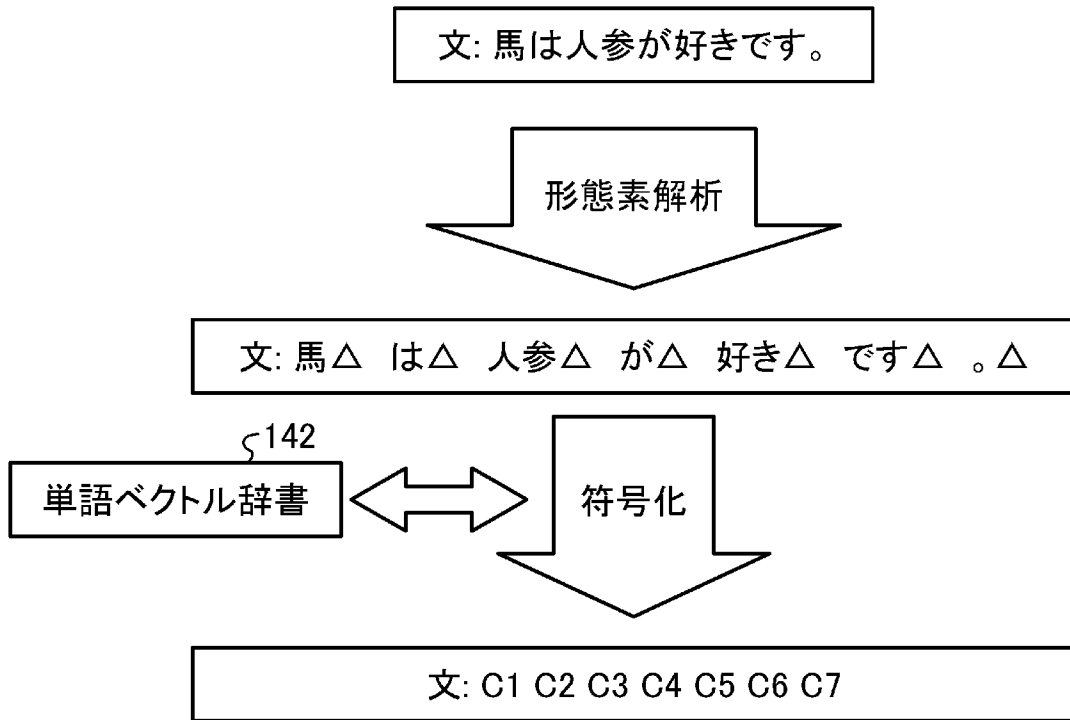
[図2]



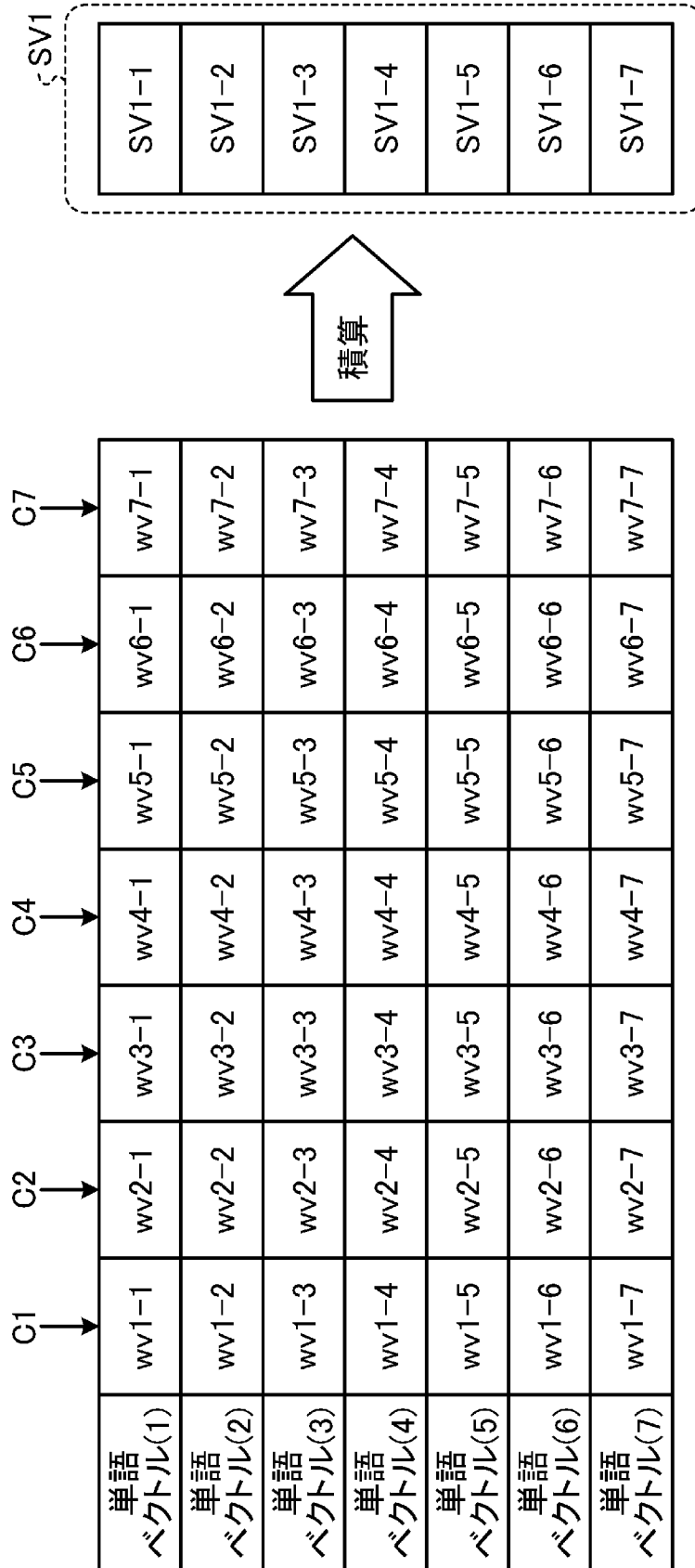
[図3]



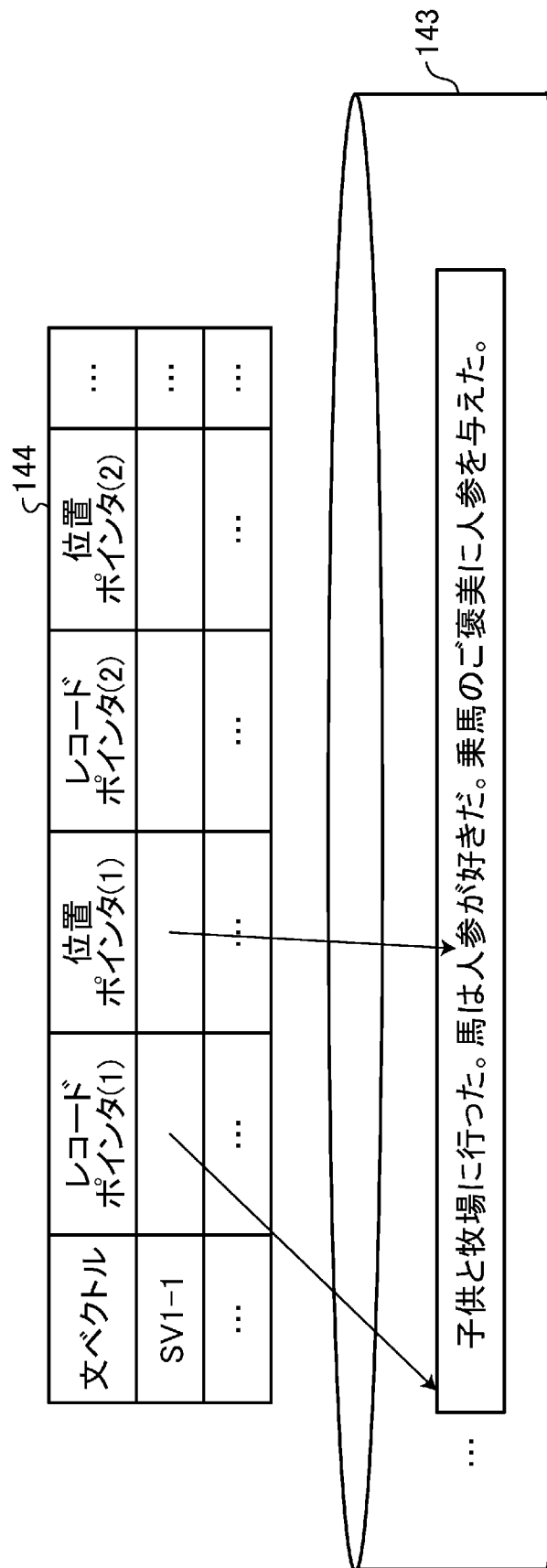
[図5A]



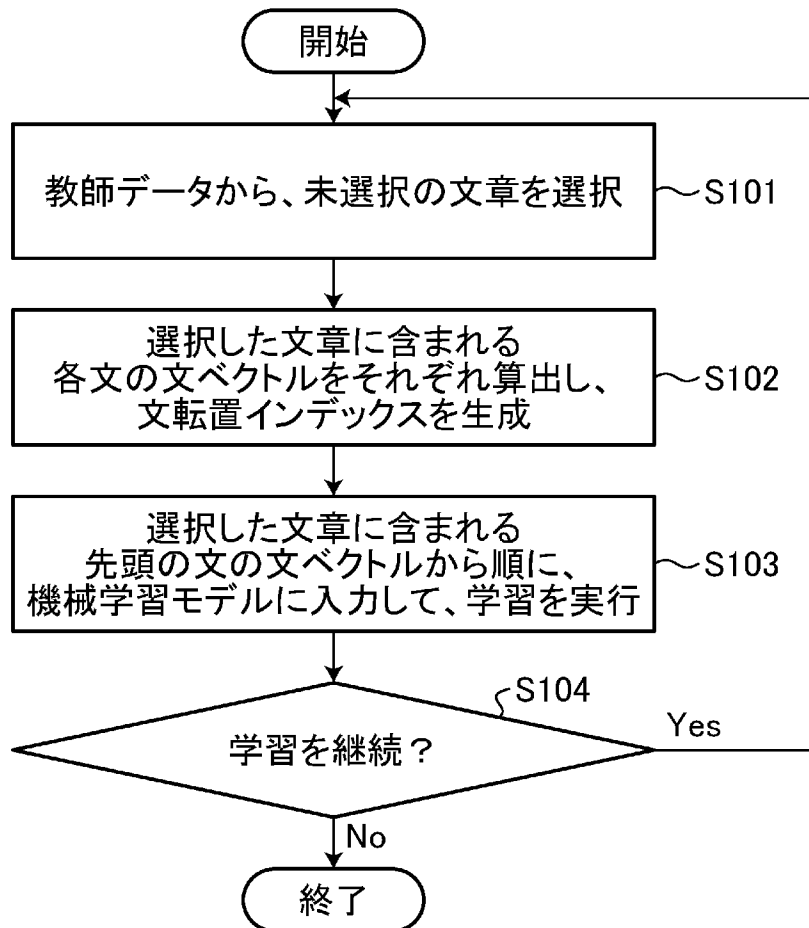
[図5B]



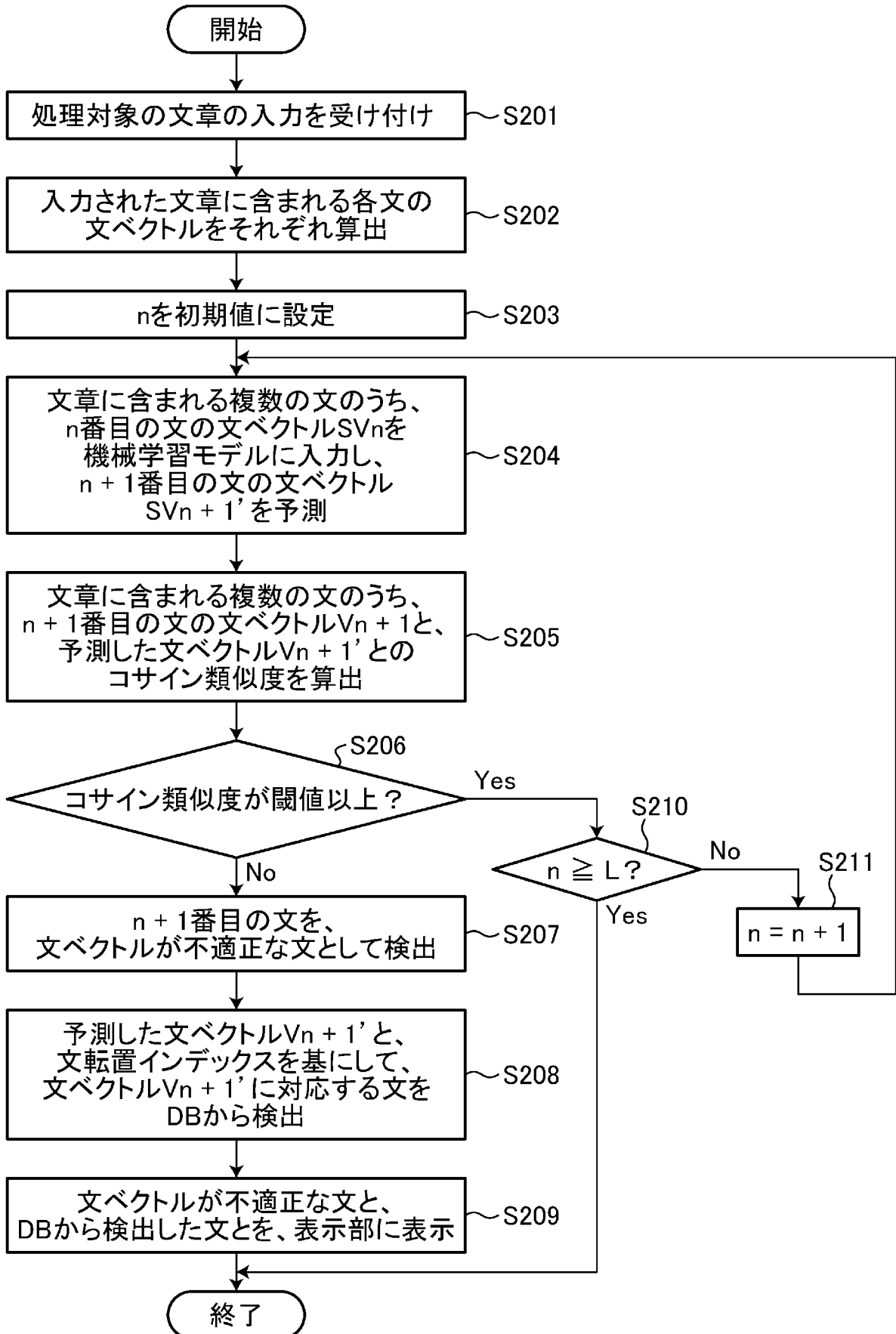
[図6]



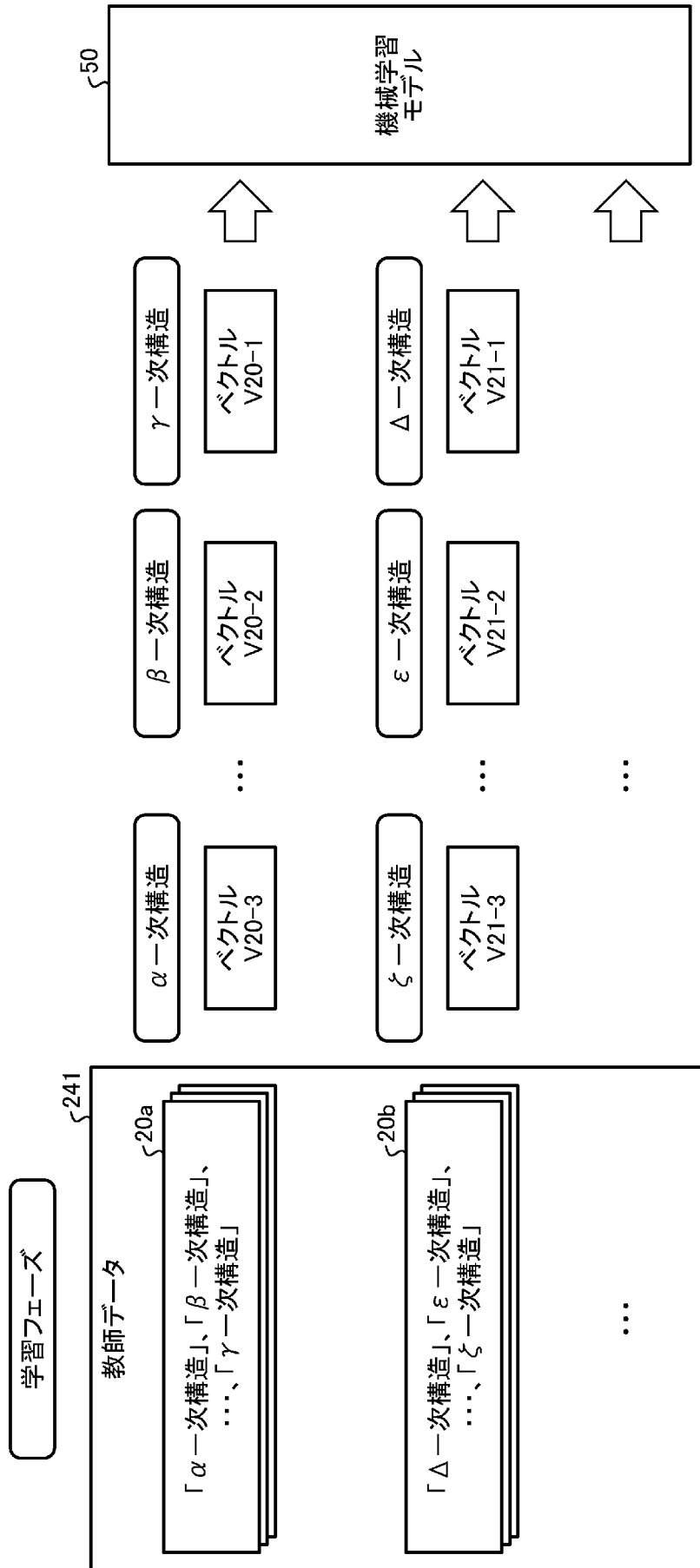
[図7]



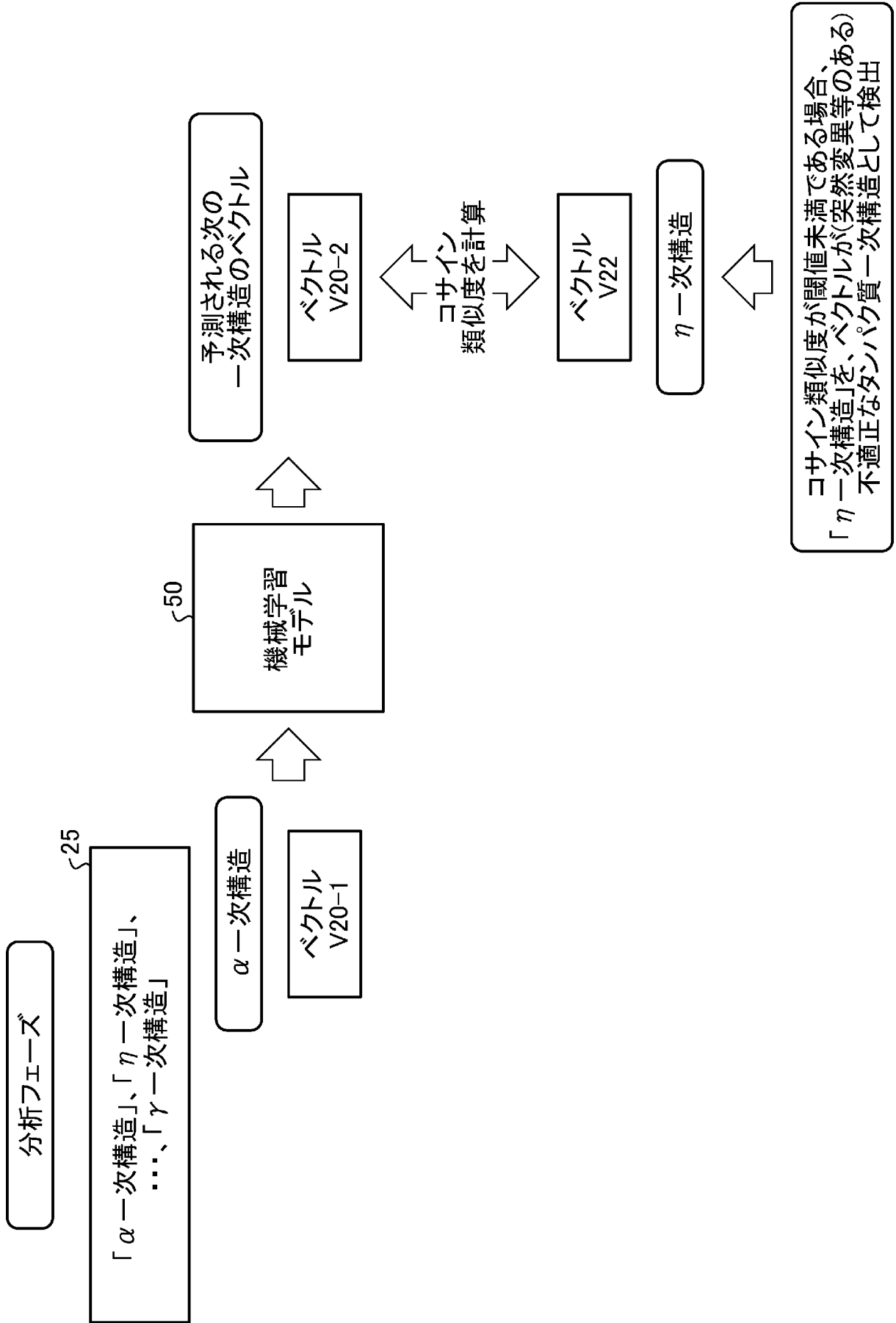
[図8]



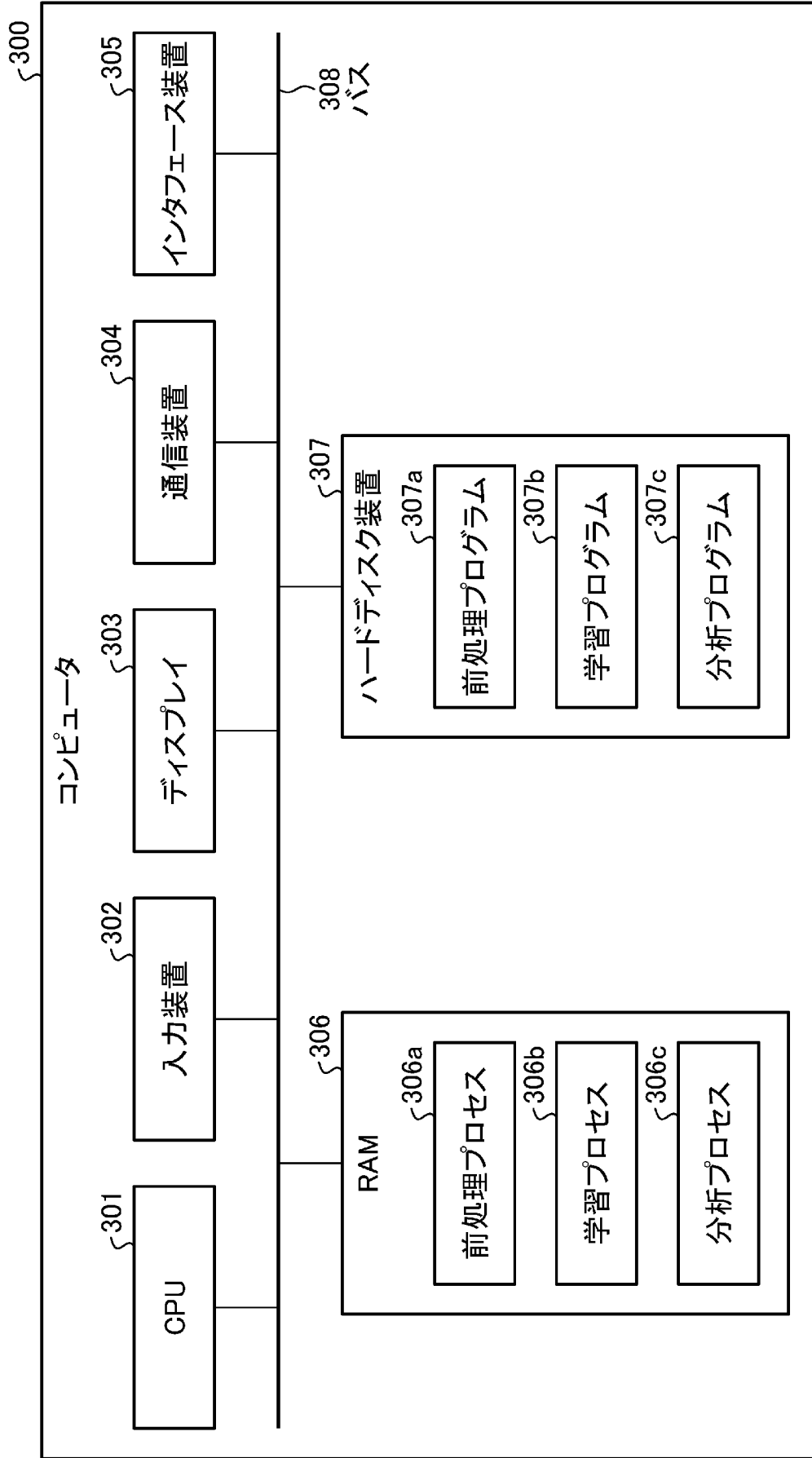
[図9]



[図10]



[図11]



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2022/022525

A. CLASSIFICATION OF SUBJECT MATTER		
<i>G06F 40/166</i> (2020.01)i; <i>G06F 40/216</i> (2020.01)i; <i>G06F 40/232</i> (2020.01)i; <i>G06F 40/253</i> (2020.01)i; <i>G06F 40/44</i> (2020.01)i FI: G06F40/166; G06F40/216; G06F40/232; G06F40/253; G06F40/44		
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) G06F40/00-40/58		
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 1971-2022 Registered utility model specifications of Japan 1996-2022 Published registered utility model applications of Japan 1994-2022		
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	CN 111539199 A (CHINA MOBILE (HANGZHOU) INFORMATION TECHNOLOGY CO., LTD.) 14 August 2020 (2020-08-14) entire text, all drawings	1-12
A	JP 2019-016140 A (ASAHI SHIMBUN PUBLISHING) 31 January 2019 (2019-01-31) entire text, all drawings	1-12
A	WO 2021/124490 A1 (FUJITSU LIMITED) 24 June 2021 (2021-06-24) entire text, all drawings	1-12
A	JP 2021-089696 A (IB RES KK) 10 June 2021 (2021-06-10) entire text, all drawings	1-12
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 11 July 2022		Date of mailing of the international search report 02 August 2022
Name and mailing address of the ISA/JP Japan Patent Office (ISA/JP) 3-4-3 Kasumigaseki, Chiyoda-ku, Tokyo 100-8915 Japan		Authorized officer Telephone No.

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/JP2022/022525

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
CN	111539199	A	14 August 2020	(Family: none)	
JP	2019-016140	A	31 January 2019	(Family: none)	
WO	2021/124490	A1	24 June 2021	(Family: none)	
JP	2021-089696	A	10 June 2021	(Family: none)	

A. 発明の属する分野の分類（国際特許分類（IPC）） G06F 40/166(2020.01)i; G06F 40/216(2020.01)i; G06F 40/232(2020.01)i; G06F 40/253(2020.01)i; G06F 40/44(2020.01)i FI: G06F40/166; G06F40/216; G06F40/232; G06F40/253; G06F40/44		
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）） G06F40/00-40/58 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 1971-2022年 日本国実用新案登録公報 1996-2022年 日本国登録実用新案公報 1994-2022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	CN 111539199 A (CHINA MOBILE (HANGZHOU) INFORMATION TECHNOLOGY CO., LTD.) 14.08.2020 (2020-08-14) 全文, 全図	1-12
A	JP 2019-016140 A (株式会社朝日新聞社) 31.01.2019 (2019-01-31) 全文, 全図	1-12
A	WO 2021/124490 A1 (富士通株式会社) 24.06.2021 (2021-06-24) 全文, 全図	1-12
A	JP 2021-089696 A (アイピーリサーチ株式会社) 10.06.2021 (2021-06-10) 全文, 全図	1-12
<input type="checkbox"/> C欄の続きにも文献が列挙されている。 <input checked="" type="checkbox"/> 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、一般的な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抵触するものではなく、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	11.07.2022	国際調査報告の発送日 02.08.2022
名称及びあて先 日本国特許庁(ISA/JP) 〒100-8915 日本国 東京都千代田区霞が関三丁目4番3号	権限のある職員（特許庁審査官） 木村 大吾 5N 4807 電話番号 03-3581-1101 内線 3586	

国際調査報告
パテントファミリーに関する情報

国際出願番号

PCT/JP2022/022525

引用文献	公表日	パテントファミリー文献	公表日
CN 111539199 A	14.08.2020	(ファミリーなし)	
JP 2019-016140 A	31.01.2019	(ファミリーなし)	
WO 2021/124490 A1	24.06.2021	(ファミリーなし)	
JP 2021-089696 A	10.06.2021	(ファミリーなし)	