

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)特許出願公表番号

特表2023-541106
(P2023-541106A)

(43)公表日 令和5年9月28日(2023.9.28)

(51)Int.Cl. F I テーマコード(参考)
 G 0 6 F 16/90 (2019.01) G 0 6 F 16/90 1 0 0 5 B 1 7 5
 G 0 6 F 16/33 (2019.01) G 0 6 F 16/33

審査請求 有 予備審査請求 未請求 (全 26 頁)

(21)出願番号	特願2023-507493(P2023-507493)	(71)出願人	520412811 エルアンドティー テクノロジー サービス シズ リミテッド インド国, チェンナイ 6 0 0 0 8 9, ラ ーマプラーム, マウント プーナマリー ロード, ディーエルエフ アイティー エ スイーゼット パーク, ブロック 3, 2 番 フロア, 1 / 1 2 4
(86)(22)出願日	令和4年3月18日(2022.3.18)	(74)代理人	100147485 弁理士 杉村 憲司
(85)翻訳文提出日	令和5年3月14日(2023.3.14)	(74)代理人	230118913 弁護士 杉村 光嗣
(86)国際出願番号	PCT/IB2022/052469	(74)代理人	100192924 弁理士 石井 裕充
(87)国際公開番号	WO2022/269369		
(87)国際公開日	令和4年12月29日(2022.12.29)		
(31)優先権主張番号	202141028709		
(32)優先日	令和3年6月25日(2021.6.25)		
(33)優先権主張国・地域又は機関	インド(IN)		

最終頁に続く

(54)【発明の名称】 文書の関連性を調査するためのシステム及び方法

(57)【要約】

本明細書で開示されるのは、文書の関連性を調査するためのシステム102及び方法である。システム102は、ユーザからの要求に基づいて、1つ又は複数のデータソース210から文書を抽出する。次いで、システム102は、ユーザから、ユーザ意図情報108及びユーザクエリ110を取得する。次いで、システム102は、各文書の関連性レベルを決定するために、ユーザ意図情報108に関して各文書を解析する。関連性レベルは、ランキングスコアの形式で示される。システムは、文書をランク付けし、それらのスコアの順に文書をユーザに表示する。システム102はまた、文書からの重要な抜粋をハイライトし、一つ一つの文書に関してユーザによって送信された1つ又は複数のユーザクエリ110への1つ又は複数の応答222を提供する。受信された応答に基づいて、ユーザは、システムをさらに訓練するためのフィードバックを提供し、それによって、より良い精度を達成する。

【選択図】 図4

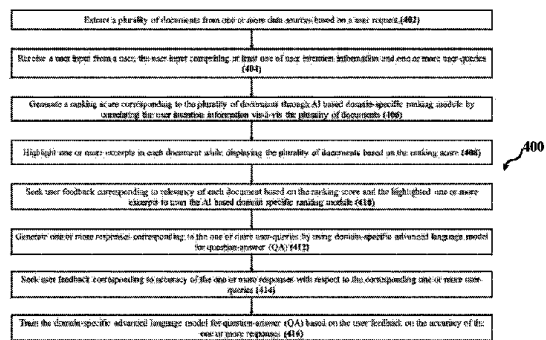


Figure 4

【特許請求の範囲】

【請求項1】

文書の関連性を調査するための方法であって、ユーザ要求に基づいて、1つ又は複数のデータソース（210）から複数の文書を抽出する（402）ことと、ユーザからユーザ入力を受信する（404）ことであって、前記ユーザ入力は、ユーザ意図情報（108）及び1つ又は複数のユーザクエリ（110）のうちの少なくとも1つを含む、前記受信する（404）ことと、前記複数の文書に対して前記ユーザ意図情報（108）を相関させることによって、ドメイン固有ランキングモジュールを介して前記複数の文書の各々のランキングスコアを生成する（406）ことであって、前記ランキングスコアは、前記ユーザ意図情報（108）に関して前記複数の文書の各々の関連性レベルを示す、前記生成する（406）ことと、前記ランキングスコアに基づいて前記複数の文書を表示する間、各文書における1つ又は複数の抜粋をハイライトする（408）ことと、前記ドメイン固有ランキングモジュールを訓練するために、前記ランキングスコアと前記ハイライトされた1つ又は複数の抜粋とに基づいて、各文書の関連性に対応するユーザフィードバックを探索する（410）こととを含む、方法。

10

【請求項2】

質疑応答（QA）のためのドメイン固有の高度な言語モデルを使用することによって、前記1つ又は複数のユーザクエリ（110）に対応する1つ又は複数の応答（222）を生成する（412）ことと、前記対応する1つ又は複数のユーザクエリ（110）に関して前記1つ又は複数の応答の精度に対応するユーザフィードバックを探索する（414）ことと、前記1つ又は複数の応答（222）の前記精度への前記ユーザフィードバックに基づいて、前記質疑応答（QA）のための高度な言語モデルを訓練する（416）こととをさらに含む、請求項1に記載の方法。

20

【請求項3】

前記ユーザ要求は、前記複数の文書に関連するキーワードのセット（104）、又は、前記複数の文書に係る複数の固有ID（106）を含む、請求項1に記載の方法。

30

【請求項4】

前記複数の文書に対して前記ユーザ意図情報を相関させることは、文ペアの、第1文に対応する第1の依存木と、第2文に対応する第2の依存木とを生成する（302、304）ことと、ドメイン固有の高度な言語モデルを採用することによって、前記第1文に対応する第1のトークンのセットと、前記第2文に対応する第2のトークンのセットとを符号化する（306）ことであって、前記第1のトークンのセットは、前記第1文の単語のセットに対応し、前記第2のトークンのセットは、前記第2文の単語のセットに対応する、前記符号化する（306）ことと、前記第1文の単語のセットと前記第2文の前記単語のセットとの間の依存関係を確立するために、前記第1文の前記単語のセットに対応する第1のアテンションマトリクスと、前記第2文の前記単語のセットに対応する第2のアテンションマトリクスとを生成する（308）ことと、相互アテンションを採用することによって、前記第1文及び前記第2文に対応する意味表現を生成する（310）ことと、前記生成された埋め込みが意味的に関連して、意味的に意味のある文の埋め込みを生成するために類似性尺度で区別されることができるようするために、特定の損失関数を使用することによってネットワーク内の複数の重みを更新することによって、前記第1文の前記単語のセットと、前記第2文の前記単語のセットとを微調整する（312）ことと、少なくとも1つの類似性測定技術を採用することによって、前記微調整された第1文及び

40

50

前記第2文のランキングスコアを生成する(314)こととを含む、請求項1に記載の方法。

【請求項5】

各文書における前記1つ又は複数の抜粋をハイライトする(408)ことは、前記ユーザ意図情報(108)に関して各文書内の各段落のランキングスコアを生成することと、

各文書から、前記ユーザ意図情報(108)に関して最も高いランキングスコアを有する少なくとも1つの段落を選択することと

によって、各文書内の少なくとも1つの重要な段落をハイライトすることを含み、

前記各文書における前記1つ又は複数の抜粋をハイライトする(408)ことは、前記少なくとも1つの重要な段落の前記複数の単語のアテンション重みと前記ユーザ意図情報とを計算することと、

前記複数の単語の各々にスコアを割り当てることと、

最大スコアに基づいて、前記少なくとも1つの重要な段落から開始点及び終了点を取得することと

によって、前記少なくとも1つの重要な段落の少なくとも1つの重要な文をハイライトすることを含む、

請求項1に記載の方法。

【請求項6】

文書の関連性を調査するためのシステムであって、

ユーザ要求に基づいて、1つ又は複数のデータソース(210)から複数の文書を抽出するように構成されている抽出ユニット(231)と、

ユーザからユーザ入力を受信するように構成されている受信ユニット(232)であって、前記ユーザ入力は、ユーザ意図情報(108)及び1つ又は複数のユーザクエリ(110)のうちの少なくとも1つを含む、前記受信ユニット(232)と、

前記複数の文書に対して前記ユーザ意図情報(108)を相関させることによって、ドメイン固有ランキングモジュール(300)を介して前記複数の文書に対応するランキングスコアを生成するように構成されているスコア生成ユニット(233)であって、各ランキングスコアは、前記ユーザ意図情報(108)に関して文書の関連性レベルを示す、前記スコア生成ユニット(233)と、

前記ランキングスコアに基づいて前記複数の文書を表示する間、各文書における1つ又は複数の抜粋をハイライトするように構成されているハイライトユニット(234)と、

前記ランキングスコアと前記ハイライトされた1つ又は複数の抜粋とに基づいて、各文書の関連性に対応するユーザフィードバックを探索するように構成されているフィードバックユニット(236)と、

前記ユーザフィードバックに基づいて、前記ドメイン固有ランキングモジュール(300)を訓練するように構成されている訓練ユニット(237)と

を含む、システム。

【請求項7】

質疑応答(QA)のための高度な言語モデルを使用することによって、前記1つ又は複数のユーザクエリ(110)に対応する1つ又は複数の応答(222)を生成するように構成されている応答生成ユニット(235)をさらに含み、

前記フィードバックユニット(236)は、前記対応する1つ又は複数のユーザクエリ(110)に関して前記1つ又は複数の応答(222)の精度に対応するユーザフィードバックを探索するようにさらに構成されており、

前記訓練ユニット(237)は、前記1つ又は複数の応答(222)の前記精度に関する前記ユーザフィードバックに基づいて、前記質疑応答(QA)のための高度な言語モデルを訓練するようにさらに構成されている、請求項7に記載のシステム。

【請求項8】

前記ユーザ要求は、前記複数の文書に関連するキーワードのセット(104)、又は、前

記複数の文書に係る複数の固有ID（106）を含む、請求項7に記載のシステム。

【請求項9】

前記複数の文書に対して前記ユーザ意図情報を相関させるために、前記スコア生成ユニット（233）は、

文ペアの、第1文に対応する第1の依存木と第2文に対応する第2の依存木とを生成し、ドメイン固有の高度な言語モデルを採用することによって、前記第1文に対応する第1のトークンのセットと、前記第2文に対応する第2のトークンのセットとを符号化することであって、前記第1のトークンのセットは、前記第1文の単語のセットに対応し、前記第2のトークンのセットは、前記第2文の単語のセットに対応する、前記符号化し、前記第1文の単語のセットと前記第2文の前記単語セットとの間の依存関係を確立するために、前記第1文の前記単語のセットに対応する第1のアテンションマトリクスと、前記第2文の前記単語のセットに対応する第2のアテンションマトリクスとを生成し、相互アテンションを採用することによって、前記第1文及び前記第2文に対応する意味表現を生成し、

前記生成された埋め込みが意味的に関連して、意味的に意味のある文の埋め込みを生成するために類似性尺度で区別されることができるようになるために、特定の損失関数を使用することによって前記ネットワーク内の複数の重みを更新することによって、前記第1文の前記単語のセットと、前記第2文の前記単語のセットとを微調整し、少なくとも1つの類似性測定技術を採用することによって、前記微調整された第1文及び前記第2文のランキングスコアを生成する

ようにさらに構成されている、請求項7に記載のシステム。

【請求項10】

各文書における前記1つ又は複数の抜粋をハイライトするために、前記ハイライトユニット（234）は、

前記ユーザ意図情報（108）に関して各文書内の各段落のランキングスコアを生成することと、

各文書から、前記ユーザ意図情報（108）に関して最も高いランキングスコアを有する少なくとも1つの段落を選択すること

によって、各文書内の少なくとも1つの重要な段落をハイライトするようにさらに構成されており、

前記ハイライトユニット（234）は、

前記少なくとも1つの重要な段落の複数の単語のアテンション重みと前記ユーザ意図情報（108）とを計算することと、

前記複数の単語の各々にスコアを割り当てることと、

前記スコアに基づいて、前記少なくとも1つの重要な段落から開始点及び終了点を取得することと

によって、前記少なくとも1つの重要な段落の少なくとも1つの重要な文をハイライトするようにさらに構成されている、請求項7に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、データ解析の分野に関する。より具体的には、文書の関連性を調査するように、様々なソースから抽出された文書を解析することに関する。

【背景技術】

【0002】

以下の説明は、本開示を理解するのに有用であり得る情報を含む。本明細書で提供される任意の情報が、先行技術であること、又は、本請求項に関連すること、又は、具体的に或いは黙示的に参照される任意の刊行物が先行技術であることを認めるものではない。

【0003】

調査ベースの研究のための報告書を作成することは、多数の文書を手動で抽出し、解析す

10

20

30

40

50

ることを必要とし、したがって、非常に退屈な作業である。個人は、関連文書を見つけるために長い時間を費やさなければならない、関連文書を見つけた後、文書を研究に含めるか否かを決定するために、一つ一つの文書を研究しなければならない。例えば、個人は、医療デバイスの臨床評価を記載する臨床評価報告書（CER）を準備したい。このために、まず、個人は、既存の文献、臨床経験、臨床試験、又はその3つの任意の組み合わせから臨床データを識別しなければならない。次いで、個人は、収集されたデータに基づいて、データの関連性、適用性、品質、及び意義を評価し、最終的に、CERに結論を明示しなければならない。さらに、これらの2つのステップは、作成される報告書の品質に不可欠であるので、これらのステップは、細心の注意を払って実行されなければならない。したがって、手動で実行されるとき、プロセス全体は、非常に時間がかかり、精神的に疲弊する。

10

【0004】

したがって、データの識別及び評価と、文書の関連性を調査することとに関わる手動労力を軽減する方法及びシステムの必要性がある。

【発明の概要】

【0005】

本開示は、先行技術の1つ又は複数の欠点を克服し、本開示を通して議論される追加の利点を提供する。追加の特徴及び利点は、本開示の技術を通して実現される。本開示の他の実施形態及び態様は、本明細書において詳細に説明され、請求された開示の一部とみなされる。

20

【0006】

本開示の非限定的な一実施形態では、文書の関連性を調査するための方法が開示される。方法は、ユーザ要求に基づいて、1つ又は複数のデータソースから複数の文書を抽出することを含む。方法は、ユーザからユーザ入力を受信することをさらに含み、ユーザ入力は、ユーザ意図情報及び1つ又は複数のユーザクエリのうちの少なくとも1つを含む。方法は、複数の文書に対してユーザ意図情報を相関させることによって、ドメイン固有ランキングモジュールを介して複数の文書の各々のランキングスコアを生成することをさらに含み、ランキングスコアは、ユーザ意図情報に関して複数の文書の各々の関連性レベルを示す。方法は、ランキングスコアに基づいて複数の文書を表示する間、各文書における1つ又は複数の抜粋をハイライトすることをさらに含む。方法は、ドメイン固有ランキングモジュールを訓練するために、ランキングスコアとハイライトされた1つ又は複数の抜粋とに基づいて、各文書の関連性に対応するユーザフィードバックを探索することをさらに含む。

30

【0007】

本開示のさらに別の非限定的な実施形態では、文書の関連性を調査するためのシステムが開示される。システムは、ユーザ要求に基づいて、1つ又は複数のデータソースから複数の文書を抽出するように構成されている抽出ユニットを含む。システムは、ユーザからユーザ入力を受信するように構成されている受信ユニットをさらに含み、ユーザ入力は、ユーザ意図情報及び1つ又は複数のユーザクエリのうちの少なくとも1つを含む。システムは、複数の文書に対してユーザ意図情報を相関させることによって、ドメイン固有ランキングモジュールを介して複数の文書に対応するランキングスコアを生成するように構成されているスコア生成ユニットをさらに含み、各ランキングスコアは、ユーザ意図情報に関して文書の関連性レベルを示す。システムは、ランキングスコアに基づいて複数の文書を表示する間、各文書における1つ又は複数の抜粋をハイライトするように構成されているハイライトユニットをさらに含む。システムは、ランキングスコアとハイライトされた1つ又は複数の抜粋とに基づいて、各文書の関連性に対応するユーザフィードバックを探索するように構成されているフィードバックユニットをさらに含む。システムは、ユーザフィードバックに基づいて、ドメイン固有ランキングモジュールを訓練するように構成されている訓練ユニットをさらに含む。

40

【0008】

50

前述の概要は、例示的にすぎず、任意の方法においても限定的であることを意図していない。上述された例示的な態様、実施形態、及び特徴に加えて、さらなる態様、実施形態、及び特徴は、図面及び以下の詳細な説明を参照することによって明らかになるであろう。

【図面の簡単な説明】

【0009】

本開示の実施形態自体、ならびに好ましい使用態様、さらなるその目的及び利点は、添付の図面と共に読むとき、例示的な一実施形態の以下の詳細な説明を参照することによって、最も良く理解されるであろう。ここで、1つ又は複数の実施形態は、例としてのみ、添付の図面を参照しながら記載される。

【図1A】図1Aは、本開示の一実施形態による、文書の関連性を調査するためのシステムの環境100を説明する。

【図1B】図1Bは、本開示の一実施形態による、文書の関連性を調査するためのシステムの環境100を説明する。

【図1C】図1Cは、本開示の一実施形態による、文書の関連性を調査するためのシステムの環境100を説明する。

【図2】図2は、本開示の一実施形態による、文書の関連性を調査するためのシステムのブロックダイアグラム200を説明する。

【図3】図3は、本開示の一実施形態による、ドメイン固有ランキングモジュールを説明する例示的な一実施形態300を説明する。

【図3A】図3Aは、本開示の一実施形態による、第1文に対応する依存木300Aを説明する。

【図3B】図3Bは、本開示の一実施形態による、第2文に対応する依存木300Bを説明する。

【図3C】図3Cは、本開示の一実施形態による、第1文に対応するアテンションマトリクス300Cを説明する。

【図3D】図3Dは、本開示の一実施形態による、第2文に対応するアテンションマトリクス300Dを説明する。

【図4】図4は、本開示の一実施形態による、文書の関連性を調査するための方法のフローチャート400を説明する。

【0010】

図は、説明の目的のみのために本開示の実施形態を描いている。当業者は、以下の説明から、本明細書で記載される開示の原理から逸脱することなく、本明細書で説明される構造及び方法の代替的な実施形態が採用され得ることを容易に認識するであろう。

【発明を実施するための形態】

【0011】

上記は、以下に続く本開示の詳細な説明がより良く理解され得るために、本開示の特徴及び技術的利点を大まかに説明した。開示された着想及び具体的な実施形態が、本開示の同じ目的を遂行するための他の構造を修正し、或いは設計するための基礎として容易に利用され得ることを当業者によって理解されるべきである。

【0012】

本開示の特徴であると考えられる新規な特徴は、その構成と動作方法との両方に関して、さらなる目的及び利点と共に、添付の図と関連して考慮されるとき、以下の説明からより良く理解されるであろう。しかしながら、図の各々は、例示及び説明の目的のみのために提供され、本開示の限定の定義として意図されないことを明示的に理解されるべきである。

【0013】

本明細書で開示されるのは、文書の関連性を調査するためのシステム及び方法である。調査報告書、例えば、評価報告書、文献調査報告書、定性的報告書、及び定量的報告書を作成する作業は、非常に面倒で時間のかかる作業である。個人は、個人の研究に関連するデータを収集するために長い時間を費やさなければならない。その後、個人は、収集された

10

20

30

40

50

データを解析し、データの関連性、適用性、品質、及び意義を調査しなければならない。しかしながら、データを収集する最初のステップと、データを解析するその後のステップとは、非常に面倒であり、手動で行われるとき、完了するのにしばしば数日かかり得る。データを収集する作業と、データを解析する作業とは、作成される報告書の品質に関して非常に不可欠であるので、これらの作業は、細心の注意を払って実行されることが必要となる。作業を手動で実行するとき、個人が一つ一つの文書を細心の注意を払って解析することは人間的に可能でないので、常に改善の余地がある。

【0014】

本開示は、文書を収集し、解析する作業においてユーザを支援するシステムを提供する。システムは、ユーザからの要求に基づいて、データソースから文書を抽出する。抽出された文書は、ユーザが行いたい検索に大まかに基づいている。次いで、システムは、ユーザが、抽出された文書に関してシステムによって回答されることを希望する特定のクエリと共に、研究を実行する背後にあるユーザの意図を受信する。クエリは、ドメイン固有であり、ユーザのニーズ及び研究ドメインに基づいて、ユーザによってカスタマイズされ得る。例えば、生物医学ドメインにおいて働いているユーザは、ナノマテリアルドメイン又は輸送ドメインにおいて働いているユーザと比較して、異なるクエリを有することとなる。次いで、システムは、ユーザの意図に関して各文書を解析し、各文書の関連性レベルを決定する。関連性レベルは、ランキングスコアの形式で示されてもよい。システムは、文書をランク付けし、ランキングスコアの順に文書をユーザに表示する。これは、ユーザが、関連性の順に文書を閲覧することを可能にする。さらに、システムはまた、文書からの重要な抜粋をハイライトし、ユーザが、文書全体を読むことの代わりにハイライトされた抜粋のみに集中することを助ける。さらに、システムはまた、一つ一つの文書に関してユーザによって送信されたクエリへの応答を提供する。このようにして、ユーザが、ユーザの研究の文書を含めるか、或いは除外するかを決定することが比較的容易になる。その関連性を調査するために文書全体を研究する必要がないので、そのような決定に達するまでに要する時間は、大幅に短縮される。システムの詳細な動作及び説明は、後続の段落において記載される。

【0015】

本開示の一実施形態によれば、図1A-図1Cは、文書の関連性を調査するためのシステムの例示的な環境100を示す。図1Aの例示的な環境100が、ユーザが生物医学ドメインにおける臨床評価報告書(CER)を作成したいことを考慮して説明されていることを当業者によって留意されなければならない。さらに、当業者は、システム102がまた、図1Aに示される以外の、様々な環境において実装されてもよいことを理解し得る。

【0016】

本開示の一実施形態によれば、例示的な環境100は、システム102のブロックダイアグラム200を示す図2と共に説明される。本開示は、システム102がサーバ上で実装されることを考慮して説明されるが、システム102は、ラップトップコンピュータ、デスクトップコンピュータ、ノートブック、ワークステーション、メインフレームコンピュータ、サーバ、ネットワークサーバ、クラウドベースコンピューティング環境などの、様々なコンピューティングシステムにおいてツールとして実装されてもよいことを理解され得る。

【0017】

一実施形態では、システム102は、I/Oインターフェース202と、プロセッサ204と、メモリ206と、ユニット208とを含んでもよい。メモリ206は、プロセッサ204及びユニット208に通信可能に結合されてもよい。さらに、メモリ206は、ユーザ意図情報108と、1つ又は複数のクエリ110と、1つ又は複数の応答222と、1つ又は複数のツール224とを格納してもよい。格納された量の各々の意義及び使用は、後続の段落において説明される。プロセッサ204は、1つ又は複数のマイクロプロセッサ、マイクロコンピュータ、マイクロコントローラ、デジタル信号プロセッサ、中央処理装置、ステートマシン、論理回路、及び/又は操作命令に基づいて信号を操作する任意の

デバイスとして実装されてもよい。他の機能の中で、プロセッサ204は、メモリ206に格納されたコンピュータ可読命令をフェッチし、実行するように構成されている。I/Oインタフェース202は、様々なソフトウェア及びハードウェアインタフェース、例えば、ウェブインタフェース、グラフィカルユーザインタフェースなどを含んでもよい。I/Oインタフェース202は、システム102が、ウェブサーバ及び外部データサーバ（図示せず）などの、他のコンピューティングデバイスと通信することを可能にしてもよい。I/Oインタフェース202は、例えば、LAN、ケーブルなどの有線ネットワークと、WLAN、セルラー、又は衛星などの無線ネットワークとを含む、多種多様なネットワーク及びプロトコルタイプ内の複数の通信を促進してもよい。I/Oインタフェース202は、多くのデバイスを互いに、或いは別のサーバに接続するための1つ又は複数のポートを含んでもよい。

10

【0018】

一実装形態では、ユニット208は、抽出ユニット231と、受信ユニット232と、スコア生成ユニット233と、ハイライトユニット234と、応答生成ユニット235と、フィードバックユニット236と、訓練ユニット237と、決定ユニット238とを含んでもよい。本開示の実施形態によれば、これらのユニット231-238は、システム102の様々な動作を実行するための、プロセッサ、マイクロプロセッサ、マイクロコントローラ、特定用途向け集積回路のようなハードウェアコンポーネントを含んでもよい。本開示の様々な実施形態によれば、プロセッサ204は、ユニット231-238の全ての機能を実行してもよいことを理解されなければならない。

20

【0019】

図1Aを参照すると、環境100は、ユーザ要求に基づいて、1つ又は複数のデータソース210から複数の文書を抽出するシステム102を示す。図2によれば、文書の抽出は、抽出ユニット231によって実行される。一実施形態では、ユーザ要求は、複数の文書に係るキーワードのセット104を含む。例えば、ユーザが作成したいCERの主題が、「非セミノーマ精巣癌を有する患者を管理するための補助としてのヒト血清中のアクセス α -フェトプロテイン (AFP) の定量決定」に関連する報告書である場合、ユーザは、「AFP」、「ヒト」、「精巣癌」を含み得るキーワードのセット104を提供する。キーワードのセットに基づく抽出ユニット231は、複数の文書を抽出する。別の実施形態では、ユーザ要求は、ユーザが、複数の文書に係る複数の固有ID106を提供することを含む。複数の固有IDは、複数の固有ID106に対応する複数の文書を抽出する抽出ユニット231に提供される。図1Aに描かれた例示的な環境100によれば、抽出された複数の文書は、250である。

30

【0020】

複数の文書が抽出されると、受信ユニット232は、ユーザ入力を受信する。一実施形態では、ユーザ入力は、自然言語を使用して、CERを作成する背後にある意図を示すユーザ意図情報108を含む。例示的な環境100によれば、ユーザ意図情報108は、「アクセス α -フェトプロテイン (AFP) アッセイは、非セミノーマ精巣癌を有する患者の管理における補助としてのヒト血清中のAFPの定量決定のためのアクセスイノムアッセイシステムでの使用のための常磁性粒子、化学発光イノムアッセイである。胎児の神経管閉鎖障害 (ONTD) の検出を補助するための母体血清及び羊水」のように読んでもよい。図1A及び図1Cに示されるように、ユーザ入力は、複数の文書に係る1つ又は複数のユーザクエリ110をさらに含む。例えば、1つ又は複数のユーザクエリ110は、「文書の意図は何ですか?」、「デバイスは何のために使用されますか?」、「研究の参加者の種類は何ですか?」などの質問を含んでもよい。1つ又は複数のユーザクエリ110は、ユーザのニーズに基づいてカスタマイズ可能であることを当業者によって留意されなければならない。すなわち、1つ又は複数のユーザクエリ110は、ユーザが行っている研究の種類及び研究が属するドメインに基づいて編集され、或いは起草され得る。例えば、例示的な環境100によれば、研究は、生物医学ドメインであり、1つ又は複数のクエリ110は、生物医学ドメインに関連する。しかしながら、研究がナノ材料のドメインであ

40

50

るユーザは、上述されるものとは異なるクエリのセットを有してもよく、ユーザのニーズに基づいてクエリをカスタマイズし得る。

【0021】

ユーザ意図情報108に関して複数の文書の関連性を調査するために、スコア生成ユニット233は、複数の文書に対してユーザ意図情報108を相関させ、各文書のランキングスコアを生成する。ランキングスコアは、ユーザ意図情報108に関して各文書の関連性レベルを示す。図1Aに描かれた、例示的な環境100によれば、ユーザ意図情報108は、250個の抽出された文書の各々に対して相関される。

【0022】

図3に説明されるように、複数の文書に対してユーザ意図情報108を相関させるために、スコア生成ユニット233は、ドメイン固有ランキングモジュール300を採用する。ドメイン固有ランキングモジュール300は、1つ又は複数のツール224のうちの1つとしてメモリ206に格納される。ドメイン固有ランキングモジュール300の詳細は、後続の段落において一例を用いて説明される。

10

【0023】

ドメイン固有ランキングモジュール300のステップ302において、文ペアは、入力される。文ペアは、ユーザ意図情報108に由来する第1文と、複数の文書のうちの1つに由来する第2文とを含む。例えば、文ペアは、以下を含んでもよい。

第1文—AFP is detected in patients with NSTC (AFPは、NSTCの患者において検出される)

第2文—Elevated AFP revealed mixed germ cell tumours (上昇されたAFPは、混合性胚細胞腫瘍を明らかにした)。

20

【0024】

ドメイン固有ランキングモジュール300のステップ304において、第1文に対応する第1の依存木300Aと第2文に対応する第2の依存木300Bとは、生成される。図3A及び図3Bに説明されるように、第1及び第2の依存木300A、300Bは、第1文及び第2文の文法構造を表し、その統語構造を明らかにする。

【0025】

ドメイン固有ランキングモジュール300のステップ306において、第1文及び第2文の意味を変更しないように、ドメイン固有の高度な言語モデルが採用され、第1文に対応する第1のトークン (又は単語) のセットと、第2文に対応する第2のトークン (又は単語) のセットとをドメイン固有の方法で符号化する。第1文及び第2文のためのドメイン固有の高度な言語モデルの動作は、それぞれ、表1及び表2に説明される。

30

【0026】

【表1】

表1：第1文のためのドメイン固有変換

第1文	AFP	Is	detected	In	Patients	With	NSTC	.
トークン	'afp'	'is'	'detected'	'in'	'patients'	'with'	'nstc'	'.'
符号化	28634	165	2490	121	568	190	3281	205

40

【0027】

50

【表 2】

表 2 : 第 2 文のためのドメイン固有変換

第 1 文	Elevated	AFP	Revealed	Mixed	Germ	Cell	tumours	.
トークン	'elevated',	'afp',	'revealed',	'mixed',	'germ',	'cell',	'tumours',	'.'
符号化	5161	2863 4	2861	4055	2927	377	11571	20 5

10

【0028】

しかしながら、ドメイン固有の高度な言語モデルの代わりに、通常言語モデルが使用される場合、第 1 文のトークンは、['a', '##f', '##p', 'is', 'detected', 'in', 'patients', 'with', 'n', '##s', '##t', '##c', '.'] のように生成されることとなり、第 2 文のトークンは、['elevated', 'a', '##f', '##p', 'revealed', 'mixed', 'ge', '##rm', 'cell', 'tu', '##mour', '##s', '.'] のように生成されることとなる。したがって、第 1 文及び第 2 文の意味は、失われることとなる。

20

【0029】

ドメイン固有ランキングモジュール 300 のステップ 308 において、第 1 文の単語のセットと第 2 文の単語のセットとの間に明示的な依存関係を確立するために、第 1 文の単語のセットに対応する第 1 のアテンションマトリクスと、第 2 文の単語のセットに対応する第 2 のアテンションマトリクスとは、作成される。第 1 文の単語のセットと第 2 文の単語のセットとに対応するアテンションマトリクスは、それぞれ、図 3 C と図 3 D とに描かれており、ハイライトされたセルは、アテンション重みに関してである。

【0030】

ドメイン固有ランキングモジュール 300 のステップ 310 において、相互アテンションの技術は、第 1 文及び第 2 文に対応する意味表現を生成するために採用される。相互アテンションは、単語間の類似性が、単語の位置決め起因して、或いは同義語又は類義語の使用起因して見過ごされることがないように、異なる位置での単語の同義語及び類義語を考慮する。例えば、第 1 文 “AFP is detected in patients with NSTC” はまた、“NSTC patients are detected with AFP levels (NSTC の患者は、AFP レベルと共に検出される)” として表現されることができ、第 2 文 “Elevated AFP revealed mixed germ cell tumours” はまた、“Mixed germ cell tumours are detected with high AFP (混合性胚細胞腫瘍は、高い AFP と共に検出される)” として表現されることができ。

30

40

【0031】

ドメイン固有ランキングモジュール 300 のステップ 312 において、第 1 文の単語のセットと第 2 文の単語のセットとは、意味的に意味のある文表現又は埋め込みを生成するために微調整される。微調整は、損失関数を最小化するために、第 1 文及び第 2 文の単語のセットを、共有パラメータを有するデュアルブランチネットワークに提供することによって達成される。共有パラメータを有するこのデュアルブランチネットワークは、類似度フィードバックのラベル付き文ペアを有するモデルを提供し、文の意味表現間の損失を最小化することによって、損失関数を最小化するために予め訓練される。

50

【0032】

ドメイン固有ランキングモジュール300のステップ314において、ランキングスコアは、コサイン類似度、マンハッタン距離類似度、ユークリッド距離類似度などの、類似度測定技術を採用することによって、微調整された第1文及び第2文のために生成される。しかしながら、上述されたもの以外の類似度測定技術はまた、使用されてもよいことを当業者によって留意され得る。

【0033】

ステップ302-314において記載された手順が繰り返し実行され、ユーザ意図情報108にしたがって、複数の文書の各々に関してランキングスコアを決定する。

【0034】

複数の文書のランキングスコアが生成されると、図1Bに描かれているように、複数の文書は、そのランキングスコアの順に表示される。一実施形態では、最も高いランキングスコアを有する文書は、一番上に表示されてもよい。しかしながら、別の実施形態では、最も高いランキングスコアを有する文書は、一番下に表示されてもよい。さらに、図1Bに示されるように、各文書は、その固有ID及びそのランキングスコアと共に表示される。図1Bに描かれたランキングスコアは、100点中であるが、ランキングスコアは、他の形式において生成され、或いは表現され得ることは、当業者によって留意され得る。

【0035】

さらに、表示された各文書に関して、ハイライトユニット234は、ユーザ意図情報108の観点から重要である1つ又は複数の抜粋をハイライトする。1つ又は複数の抜粋は、1つ又は複数の段落及び/又は1つ又は複数の文を含んでもよいことを当業者によって理解され得る。1つ又は複数の重要な段落をハイライトするために、ドメイン固有ランキングモジュール300は、ハイライトユニット234によって採用され、ユーザ意図情報108と最も高い関連性を有する文書内の段落を識別する。さらに、重要な段落の1つ又は複数の重要な文をハイライトするために、重要な段落の文のそれぞれのトークン（又は単語）のアテンション重みとユーザ意図情報108とは、スコア生成ユニット233によって計算される。アテンション重みに基づいて、各トークンは、スコアを割り当てられる。スコアに基づいて、ハイライトユニット234によって重要な文としてハイライトされる最も関連性のある段落の重要部分の開始及び終了は、取得される。

【0036】

1つ又は複数の抜粋のハイライトは、ユーザが、ユーザ意図情報108においてユーザによって定義されたユーザの意図する目的に最も関連する文書のセクションを容易に識別することを可能にする。

【0037】

さらに、図1Cに示されるように、システム102は、ユーザが、示されるように応答タブを選択することによって、ユーザによって提供された1つ又は複数のユーザクエリ110への1つ又は複数の応答を閲覧することを可能にする。図1Bでは、1つ又は複数のユーザクエリ110への応答は、システム102のメモリ206に1つ又は複数のツール224として格納された質疑応答(QA)のためのドメイン固有の高度な言語モデルを採用する応答生成ユニット235によって提供される。生成された1つ又は複数の応答222は、メモリ206に格納される。さらに、生成された1つ又は複数の応答に基づいて、フィードバックユニット236は、ユーザが、1つ又は複数の応答にフィードバックを提供することを可能にする。例えば、図1Aの例示的環境100に示されるように、ユーザが提供された応答に満足しない場合、或いは、応答生成ユニットが応答を提供することができなかった場合、ユーザは、応答生成ユニット235によって提供された応答を検証するオプション、又は応答を変更するオプションを提供される。フィードバックに基づいて、訓練ユニット237は、質疑応答(QA)のための高度な言語モデルを再訓練する。

【0038】

図1Aに示されるように、システム102はさらに、決定タブを提供することによって、ユーザが1つ又は複数の関連文書を選択することを可能にする。決定タブは、決定ユニッ

10

20

30

40

50

ト238によって制御され、特定の文書を「含む」或いは「除外する」のいずれかのオプションをユーザに提供する。しかしながら、ユーザが決定できない場合、決定ユニット238は、「決定できない」オプションをユーザに提供することによって、ユーザが特定の文書を再検討することを可能にする。

【0039】

したがって、その様々なコンポーネントの相互作用を介したシステム102は、閲覧者が最小限の時間及び労力で関連文書を識別することをより容易にし、それによって、ユーザエクスペリエンスを向上させる。さらに、ドメイン固有ランキングモジュール300が様々な技術の組み合わせを採用するので、複数の文書の各々のランキングスコアが生成される方法は、非常に広範囲且つ正確である。

10

【0040】

図4は、本開示の一実施形態による、複数の文書の関連性を調査するための方法400を描いている。

【0041】

図4に説明されるように、方法400は、製造工場における情報を管理するための方法を説明する1つ又は複数のブロックを含む。方法400は、コンピュータ実行可能命令の一般的なコンテキストで記載されてもよい。概して、コンピュータ実行可能命令は、特定の機能を実行する、或いは特定の抽象データ型を実装する、ルーチン、プログラム、オブジェクト、コンポーネント、データ構造、手順、モジュール、及び機能を含んでもよい。

【0042】

方法400が記載される順序は、限定として解釈されることを意図していなく、記載された方法ブロックの任意の数は、方法を実装するために任意の順序で組み合わせられてもよい。さらに、個々のブロックは、記載された主題の本旨及び範囲から逸脱することなく、方法から削除されてもよい。

20

【0043】

ブロック402において、方法400は、ユーザ要求に基づいて、1つ又は複数のデータソース210から複数の文書を抽出することを含んでもよい。

【0044】

ブロック404において、方法400は、ユーザからユーザ入力を受信することを含んでもよく、ユーザ入力は、ユーザ意図情報108及び1つ又は複数のユーザークエリ110のうちの少なくとも1つを含む。

30

【0045】

ブロック406において、方法400は、複数の文書に対してユーザ意図情報108を相関させることによって、複数の文書の各々のランキングスコアを生成することを含んでもよい。

【0046】

ブロック408において、方法400は、ランキングスコアに基づいて複数の文書を表示する間、各文書における1つ又は複数の抜粋をハイライトすることを含んでもよい。

【0047】

ブロック410において、方法400は、ドメイン固有ランキングモジュールを訓練するために、ランキングスコアとハイライトされた1つ又は複数の抜粋とに基づいて、各文書の関連性に対応するユーザフィードバックを探索することを含んでもよい。

40

【0048】

ブロック412において、方法400は、質疑応答(QA)のためのドメイン固有の高度な言語モデルを使用することによって、1つ又は複数のユーザークエリ110に対応する1つ又は複数の応答222を生成することを含んでもよい。

【0049】

ブロック414において、方法400は、対応する1つ又は複数のユーザークエリ110に関して1つ又は複数の応答222の精度に対応するユーザフィードバックを探索することを含んでもよい。

50

【0050】

ブロック416において、方法400は、1つ又は複数の応答222の精度に関するユーザフィードバックに基づいて、質疑応答(QA)ためのドメイン固有の高度な言語モデルを訓練することを含んでもよい。

【0051】

互いに通信するいくつかのコンポーネントを有する一実施形態の説明は、全てのそのようなコンポーネントが必要とされることを意味しない。それどころか、様々なオプションのコンポーネントは、本開示の多種多様な可能性のある実施形態を説明するために記載される。

【0052】

単一のデバイス又は物品が本明細書に記載されているとき、2つ以上のデバイス／物品（それらが協働するか否か）が、単一のデバイス／物品の代わりに使用されてもよいことは明らかであろう。同様に、2つ以上のデバイス又は物品（それらが協働するか否か）が本明細書に記載されている場合、単一のデバイス／物品が、2つ以上のデバイス又は物品の代わりに使用されてもよいこと、又は、異なる数のデバイス／物品が、示された数のデバイス又はプログラムに代えて使用されてもよいことは明らかであろう。代替的に、デバイスの機能性及び／又は特徴は、そのような機能性／特徴を有するものとして明示的に記載されていない1つ又は複数の他のデバイスによって具現化されてもよい。したがって、本開示の他の実施形態は、デバイス自体を含む必要はない。

【0053】

最後に、本明細書で使用される言語は、可読性及び説明目的のために主として選択されており、発明的主題を画成し、或いは制限するために選択されないことがある。したがって、本開示の範囲は、この詳細な説明によってではなく、むしろ、本明細書に基づく出願に関して発行される任意の請求項によって限定されることを意図している。したがって、本開示の実施形態は、以下の特許請求の範囲に規定される本開示の範囲を例示するが、限定しないことを意図している。

【0054】

本明細書では、様々な態様及び実施形態が開示されているが、他の態様及び実施形態は、当業者に明らかであろう。本明細書に開示された様々な態様及び実施形態は、例示の目的であり、限定することを意図しておらず、真の範囲及び本旨は、以下の特許請求の範囲によって示される。

10

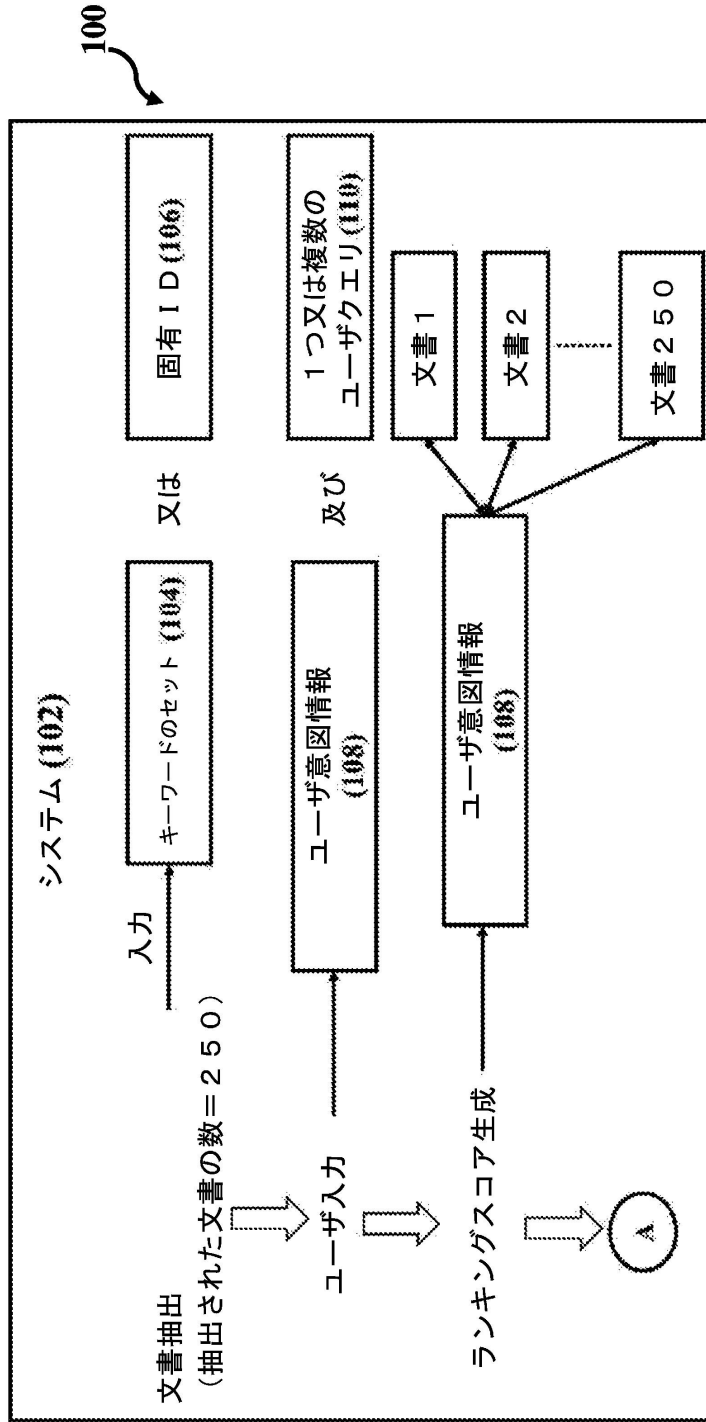
20

30

40

50

【図1A】



10

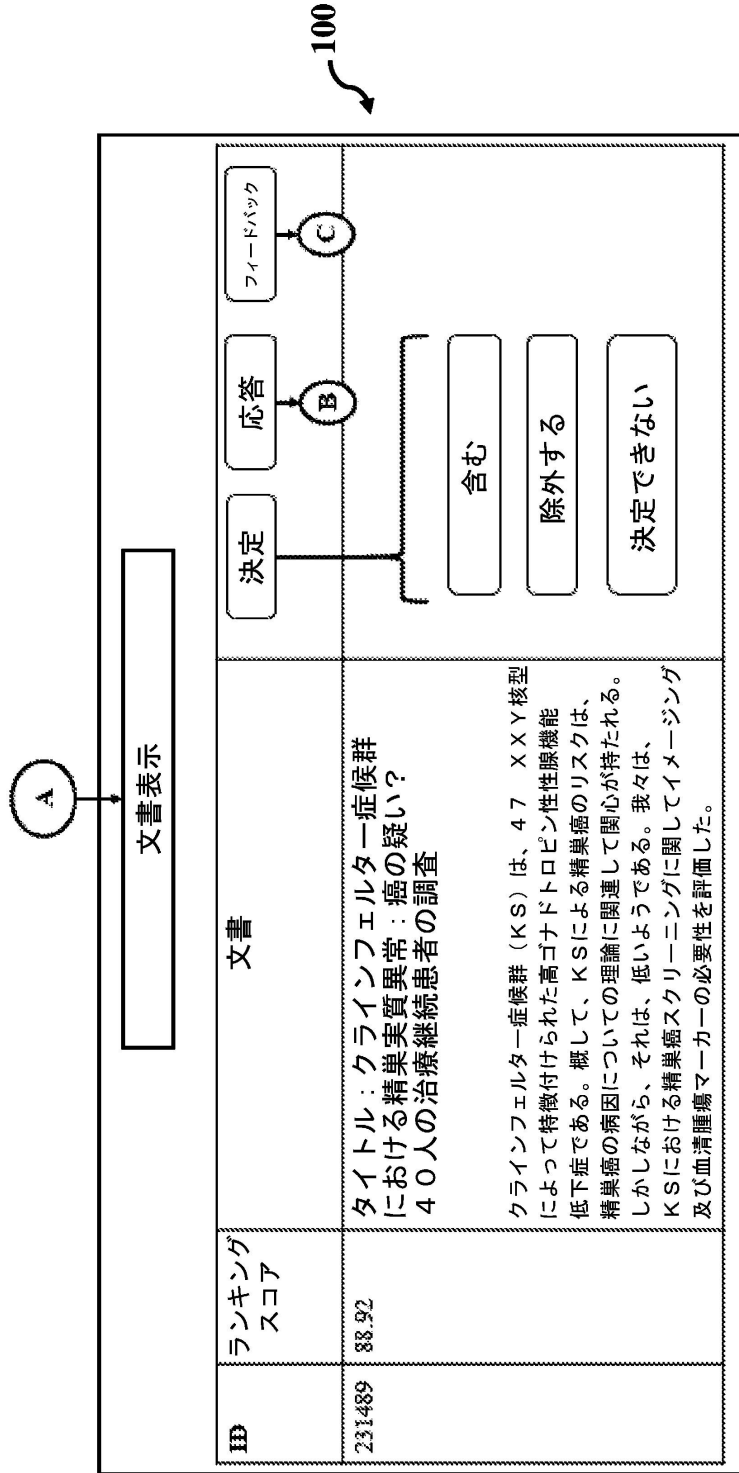
20

30

40

50

【図 1 B】



B

10

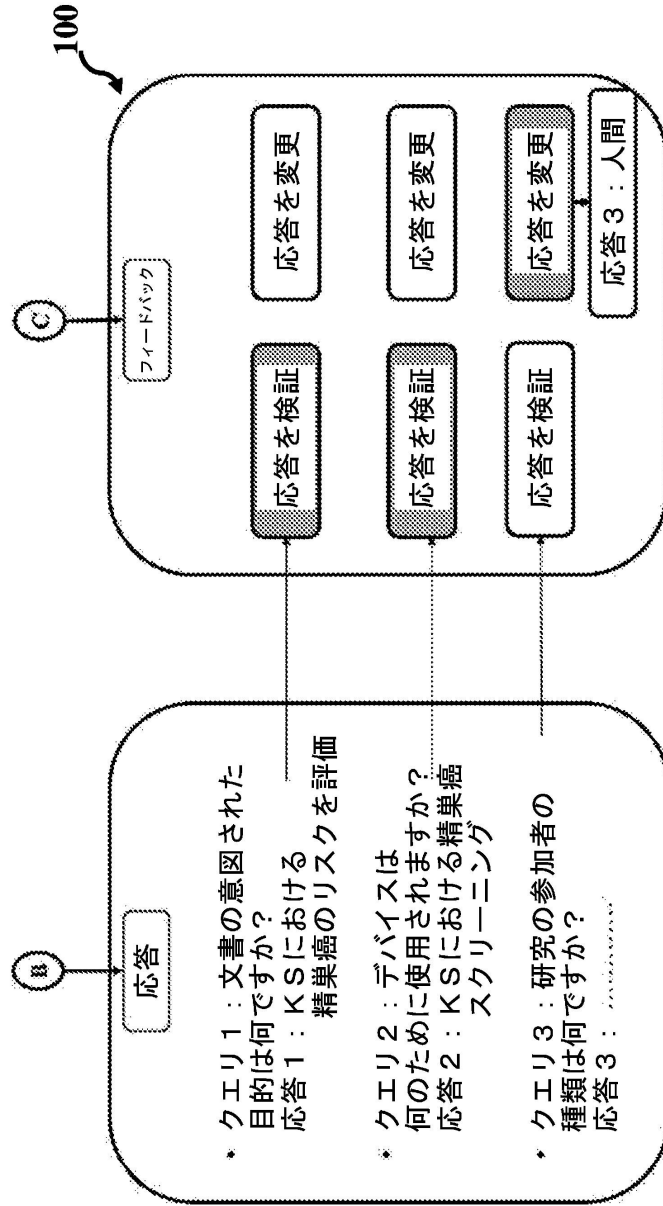
20

30

40

50

【図1C】



10

20

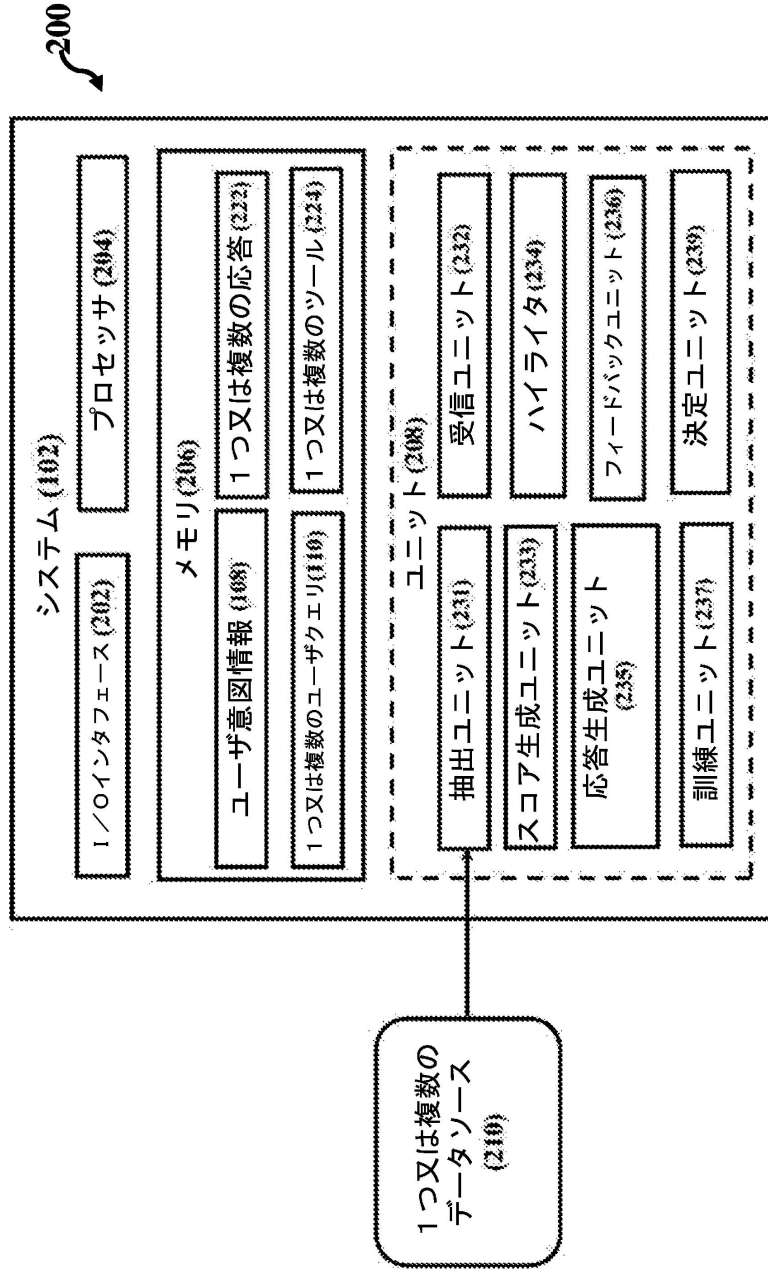
C

30

40

50

【図2】



10

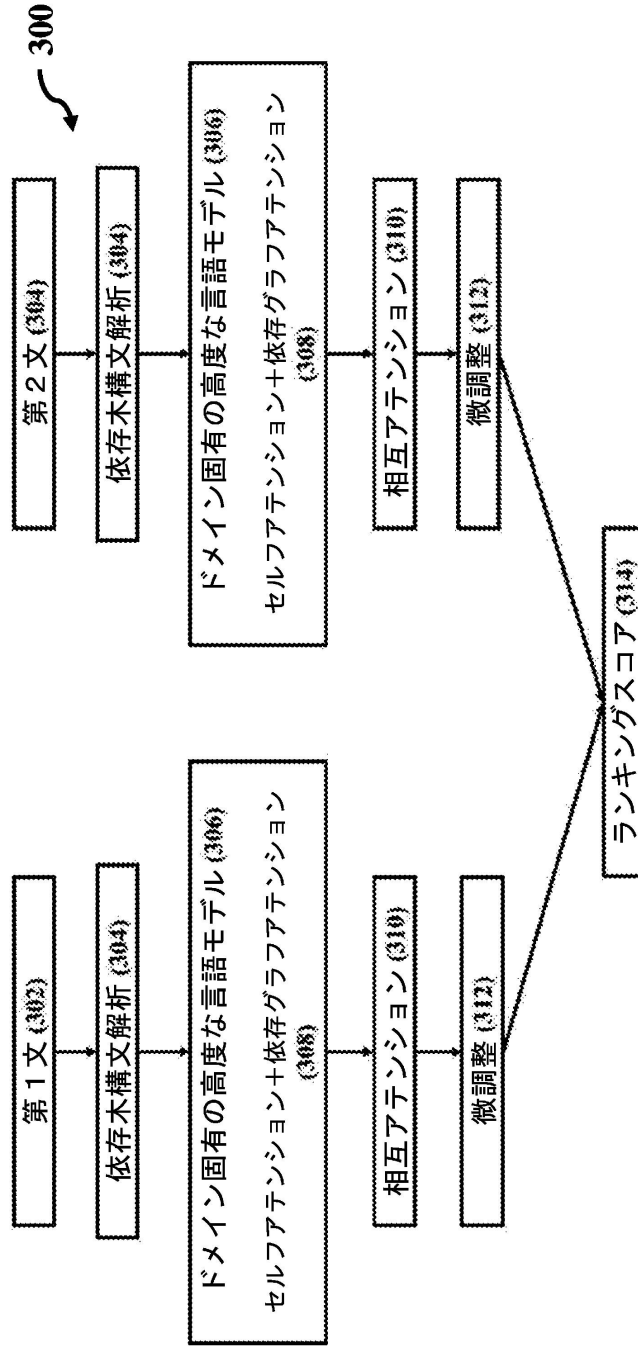
20

30

40

50

【図 3】



10

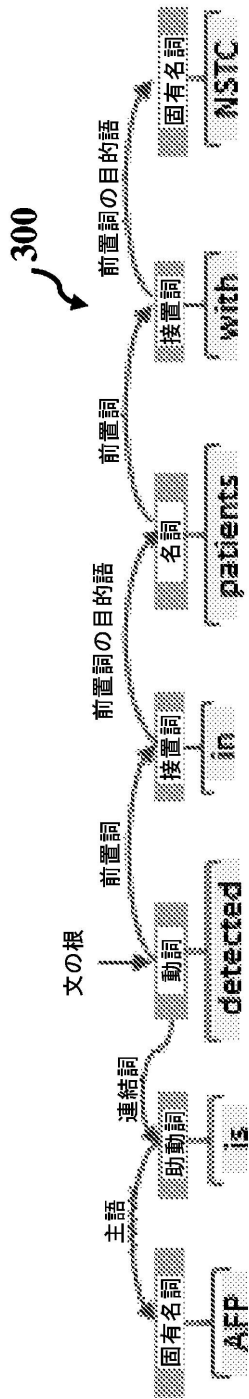
20

30

40

50

【図 3 A】



A

10

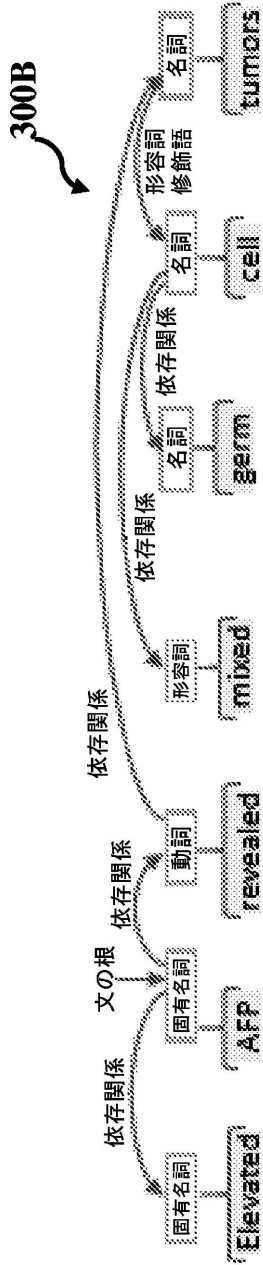
20

30

40

50

【図 3 B】



B

10

20

30

40

50

【図 3 C】

300C

	AFP	is	detected	in	patients	with	NSTC
AFP	■						
is		■					
detected			■				
in				■			
patients					■		
with						■	
NSTC							■

Figure 3C

10

20

30

40

50

【図 3 D】

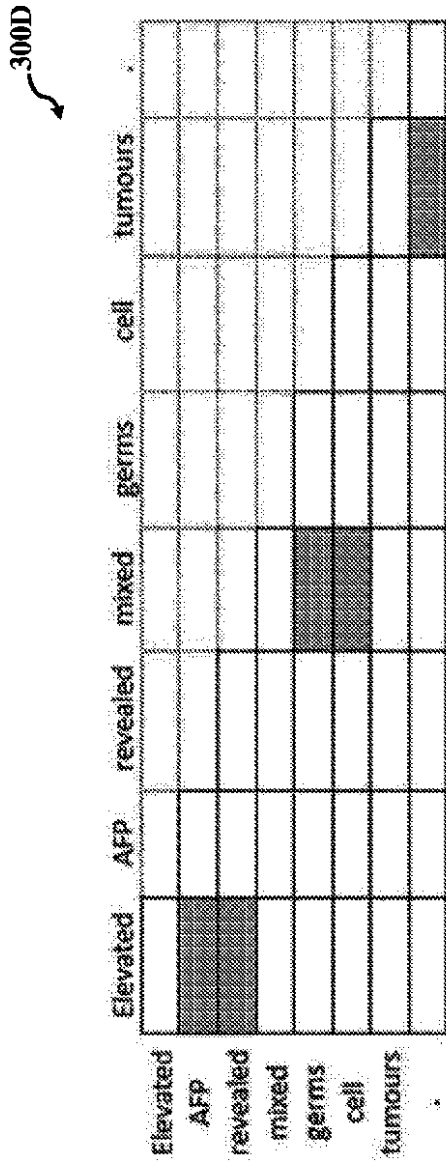


Figure 3D

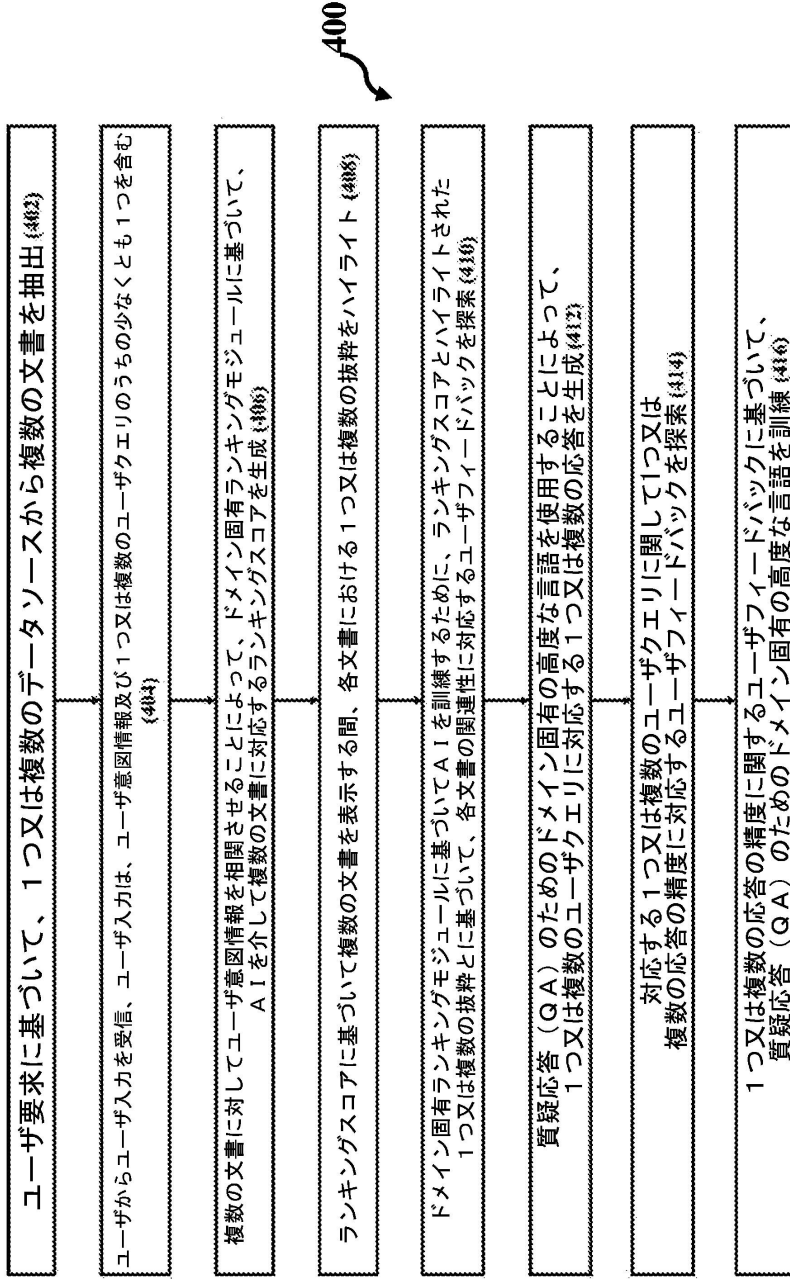
10

20

30

40

50



10

20

30

40

50

【国際調査報告】

INTERNATIONAL SEARCH REPORT		International application No. PCT/IB2022/052469
A. CLASSIFICATION OF SUBJECT MATTER G06F16/783,G06F16/93,G06F7/00 Version=2022.01 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) G06F Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) Databases - PatSeer, IPO Internal Database Keywords - relevancy document, highlight excerpts, rank, score		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	KR20110027729A (BELENZON SHARON, ET AL.) 16 March 2011 (16-03-2011) {Whole Document specially Figures 3-6,11 with description}	1-10
Y	US9449080B1 (ZHANG GUANGSHENG) 20 September 2016 (20-09-2016) {Whole Document specially Columns 1-2}	1-10
<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: "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 "A" document defining the general state of the art which is not considered to be of particular relevance "D" document cited by the applicant in the international application "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 "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) "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 "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 "&" document member of the same patent family		
Date of the actual completion of the international search 30-06-2022		Date of mailing of the international search report 30-06-2022
Name and mailing address of the ISA/ Indian Patent Office Plot No.32, Sector 14,Dwarka,New Delhi-110075 Facsimile No.		Authorized officer Saket Kumar Gupta Telephone No. +91-1125300200

Form PCT/ISA/210 (second sheet) (July 2019)

10

20

30

40

50

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/IB2022/052469

Citation	Pub.Date	Family	Pub.Date
KR 20110027729 A	16-03-2011	CA 2727963 A1	30-12-2009
		EP 2321772 A2	18-05-2011
		JP 2011525673 A	22-09-2011
		US 2011093449 A1	21-04-2011
		WO 2009156987 A2	30-12-2009

10

20

30

40

50

 フロントページの続き

(81)指定国・地域 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(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), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, 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, 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

(72)発明者 アンキット マルヴィヤ
 インド国 460330 マディヤ プラデーシュ ベトウル チチョリ ディストリクト ベトウル ロード ナンバー1289 ニア グル シャハーブ カレヅジ

(72)発明者 マダスダン シン
 インド国 560100 カルナータカ バンガロール エレクトロニック シティー-1 ニーラドリ ロード ファースト クロス アジュメラ ストーン パーク ビー-603

(72)発明者 ムリダル バララマン
 インド国 560037 カルナータカ バンガロール ドッダーネクンディ チナパナハリ メイン ロード エスヴィーエス パームス 2 ナンバー206

(72)発明者 プラカー スリヴァスタヴァ
 インド国 211002 ウッター プラデーシュ プラヤーグラージ タゴール タウン ジャワハル ラール ネルー ロード ナンバー 4-シー

Fターム(参考) 5B175 DA01 EA01 HA01 JB02