

(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2023-507881  
(P2023-507881A)

(43)公表日 令和5年2月28日(2023.2.28)

(51)Int.Cl. F I  
G 0 6 F 40/279 (2020.01) G 0 6 F 40/279

テーマコード(参考)  
5 B 0 9 1

審査請求 未請求 予備審査請求 未請求 (全 31 頁)

(21)出願番号 特願2022-525481(P2022-525481)  
 (86)(22)出願日 令和2年12月30日(2020.12.30)  
 (85)翻訳文提出日 令和4年6月20日(2022.6.20)  
 (86)国際出願番号 PCT/IB2020/062535  
 (87)国際公開番号 WO2021/137166  
 (87)国際公開日 令和3年7月8日(2021.7.8)  
 (31)優先権主張番号 201941054421  
 (32)優先日 令和1年12月30日(2019.12.30)  
 (33)優先権主張国・地域又は機関  
 インド(IN)

(71)出願人 520412811  
 エルアンドティー テクノロジー サービ  
 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  
 (74)代理人 100093861  
 弁理士 大賀 眞司  
 (74)代理人 100129218  
 弁理士 百本 宏之

最終頁に続く

(54)【発明の名称】ドメインベースのテキスト抽出方法およびシステム

(57)【要約】

本開示は、入力ファイルのコンテンツから情報を抽出するための方法およびシステムに関する。本方法は、入力ファイルからのテキストデータを識別することと、複数のテキストエンティティから関連テキストエンティティを識別するためにユーザからテキスト入力を受け取ることと、テキスト入力に対応する検索パターンを自動的に生成することと含んでよい。本方法は、複数のテキストエンティティの各々に関連付けられたパターンを決定することと、テキスト入力に対応する検索パターンを複数のテキストエンティティに関連付けられたパターンと対応付けることとをさらに含んでよい。本方法は、対応付けに基づいて複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンを識別することと、1つまたは複数の合致パターンに対応する関連テキストエンティティを複数のテキストエンティティから抽出することとをさらに含んでよい。

【選択図】図2

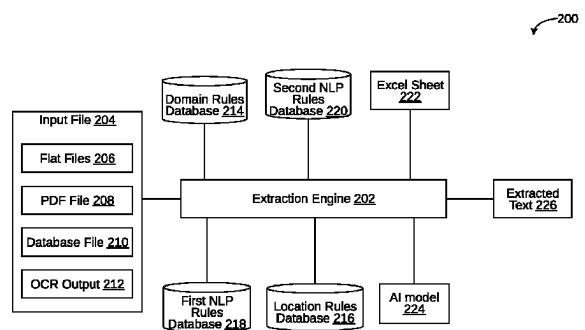


FIG. 2

## 【特許請求の範囲】

## 【請求項1】

入力ファイルのコンテンツから情報を抽出する方法であって、  
前記入力ファイルからテキストデータを識別することであって、識別された前記テキストデータは複数のテキストエンティティを含む、識別することと、  
前記複数のテキストエンティティから関連テキストエンティティを識別するために、ユーザからテキスト入力を受け取ることと、  
前記テキスト入力に対応する検索パターンを自動的に生成することと、  
前記複数のテキストエンティティの各々に関連付けられたパターンを決定することと、  
前記テキスト入力に対応する前記検索パターンを、前記複数のテキストエンティティに関連付けられたパターンと対応付けることと、  
前記対応付けに基づいて、前記複数のテキストエンティティに関連付けられた前記パターンから1つまたは複数の合致パターンを識別することと、  
前記1つまたは複数の合致パターンに対応する関連テキストエンティティを前記複数のテキストエンティティから抽出することと  
を含む方法。

10

## 【請求項2】

抽出された前記関連テキストエンティティを用いて表現を生成することを含む、請求項1に記載の方法。

## 【請求項3】

前記入力ファイルは、画像ファイルおよびポータブル・ドキュメント・フォーマット（PDF）ファイルのうち少なくとも1つを含む、請求項1に記載の方法。

20

## 【請求項4】

前記複数のテキストエンティティの各々に関連付けられたドメインベースの名前を決定することをさらに含み、  
前記ドメインベースの名前は、前記複数のテキストエンティティの各々をデータデポジトリに記憶されたドメインベースの名前のリストと対応付けることに基づいて決定される、  
請求項1に記載の方法。

## 【請求項5】

前記入力ファイルに関連付けられたタイプを識別することと、  
前記入力ファイルに関連付けられた前記タイプに基づいて、前記入力ファイルにおける1つまたは複数の関連テキストエンティティの位置を決定することと、  
前記位置に基づいて、前記入力ファイルから前記1つまたは複数の関連テキストエンティティを抽出することと  
をさらに含む、請求項1に記載の方法。

30

## 【請求項6】

前記位置は、(x, y)座標に基づくものであり、前記位置は、前記データリポジトリに記憶された、前記入力ファイルに関連付けられたタイプに対応するテンプレートに基づいて予め決定される、請求項5に記載の方法。

## 【請求項7】

前記複数のテキストエンティティに関連付けられた品詞（POS）を識別することと、  
POSが名詞として識別される1つまたは複数のテキストエンティティを抽出することと  
をさらに含み、  
前記POSは、名詞、代名詞、動詞、副詞、形容詞、接続詞、前置詞、および間投詞のうちの1つである、  
請求項1に記載の方法。

40

## 【請求項8】

POSが識別されない1つまたは複数のユニークなテキストエンティティを識別することと、  
前記1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティ

50

ティタイプを決定することと

をさらに含み、

前記エンティティタイプは、値エンティティおよびパターンエンティティのうちの1つであり、各値エンティティは関連付けられた値を有し、各パターンエンティティは関連付けられたパターンを有する、

請求項7に記載の方法。

【請求項9】

各値エンティティについて、

前記値エンティティの前記関連付けられた値と同じ関連付けられた値を有する1つまたは複数のテキストエンティティを自動的に識別すること

をさらに含む、請求項8に記載の方法。

【請求項10】

各パターンエンティティについて、

前記パターンエンティティに関連付けられたパターンに対応する検索パターンを自動的に生成することと、

前記パターンエンティティに対応する前記検索パターンを、前記複数のテキストエンティティに関連付けられたパターンと対応付けることと、

前記対応付けに基づいて、前記複数のテキストエンティティに関連付けられた前記パターンから1つまたは複数の合致パターンを決定することと、

前記1つまたは複数の合致パターンに対応する合致テキストエンティティを前記複数のテキストエンティティから抽出することと

をさらに含む、請求項8に記載の方法。

【請求項11】

分類ベースの機械学習（ML）モデルを用いて、履歴データに基づいて推奨を生成することを含む、請求項1に記載の方法。

【請求項12】

入力ファイルのコンテンツから情報を抽出するためのシステムであって、

プロセッサと、

前記プロセッサに通信可能に結合されたメモリと

を含み、

前記メモリは、プロセッサ実行可能命令を記憶し、前記プロセッサ実行可能命令は、前記プロセッサにより実行されたとき、前記プロセッサに、

前記入力ファイルからテキストデータを識別することであって、識別された前記テキストデータは複数のテキストエンティティを含む、識別すること、

前記複数のテキストエンティティから関連テキストエンティティを識別するために、ユーザからテキスト入力を受け取ること、

前記テキスト入力に対応する検索パターンを自動的に生成すること、

前記複数のテキストエンティティの各々に関連付けられたパターンを決定すること、

前記テキスト入力に対応する前記検索パターンを、前記複数のテキストエンティティに関連付けられたパターンと対応付けること、

前記対応付けに基づいて、前記複数のテキストエンティティに関連付けられた前記パターンから1つまたは複数の合致パターンを識別すること、および、

前記1つまたは複数の合致パターンに対応する関連テキストエンティティを前記複数のテキストエンティティから抽出することを行わせる、

システム。

【請求項13】

前記入力ファイルは、画像ファイルおよびポータブル・ドキュメント・フォーマット（PDF）ファイルのうちの少なくとも1つを含む、請求項12に記載のシステム。

【請求項14】

前記プロセッサ実行可能命令は、前記プロセッサにより実行されたとき、前記プロセッサ

10

20

30

40

50

に、  
 ドメインベースの抽出であって、前記ドメインベースの抽出を行うことは、  
 前記複数のテキストエンティティの各々に関連付けられたドメインベースの名前を決定することであって、前記ドメインベースの名前は、前記複数のテキストエンティティの各々をデータデポジトリに記憶されたドメインベースの名前のリストと対応付けることに基づいて決定される、決定すること  
 を含む、ドメインベースの抽出、  
 位置ベースの抽出であって、前記位置ベースの抽出を行うことは、  
 前記入力ファイルに関連付けられたタイプを識別すること、  
 前記入力ファイルに関連付けられた前記タイプに基づいて、前記入力ファイルにおける1つまたは複数の関連テキストエンティティの位置を決定すること、および、  
 前記位置に基づいて、前記入力ファイルから前記1つまたは複数の関連テキストエンティティを抽出することであって、前記位置は、前記データリポジトリに記憶された、前記入力ファイルに関連付けられたタイプに対応するテンプレートに基づいて予め決定される、抽出すること、  
 を含む、位置ベースの抽出、  
 品詞（POS）ベースの抽出であって、前記品詞（POS）ベースの抽出を行うことは、  
 前記複数のテキストエンティティに関連付けられたPOSを識別することであって、前記POSは、名詞、代名詞、動詞、副詞、形容詞、接続詞、前置詞、および間投詞のうちの1つである、識別すること、および、  
 POSが名詞として識別される1つまたは複数のテキストエンティティを抽出すること、  
 POSが識別されない1つまたは複数のユニークなテキストエンティティを識別すること  
 、  
 前記1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティティタイプを決定することであって、前記エンティティタイプは、値エンティティおよびパターンエンティティのうちの1つであり、各値エンティティは関連付けられた値を有し、各パターンエンティティは関連付けられたパターンを有する、決定すること、  
 各値エンティティについて、  
 前記値エンティティの前記関連付けられた値と同じ関連付けられた値を有する1つまたは複数のテキストエンティティを自動的に識別すること、  
 各パターンエンティティについて、  
 前記パターンエンティティに関連付けられたパターンに対応する検索パターンを自動的に生成すること、  
 前記パターンエンティティに対応する前記検索パターンを、前記複数のテキストエンティティに関連付けられたパターンと対応付けること、  
 前記対応付けに基づいて、前記複数のテキストエンティティに関連付けられた前記パターンから1つまたは複数の合致パターンを決定すること、および、  
 前記1つまたは複数の合致パターンに対応する合致テキストエンティティを前記複数のテキストエンティティから抽出すること  
 を含む、品詞（POS）ベースの抽出、ならびに、  
 機械学習（ML）ベースの抽出であって、前記MLベースの抽出を行うことは、  
 分類ベースのMLモデルを用いて、履歴データに基づいて推奨を生成すること  
 を含む、機械学習（ML）ベースの抽出のうちの少なくとも1つをさらに行わせる、  
 請求項12に記載のシステム。

【請求項15】

入力ファイルからテキストデータを識別することであって、識別された前記テキストデータは複数のテキストエンティティを含む、識別することと、  
 前記複数のテキストエンティティから関連テキストエンティティを識別するために、ユーザからテキスト入力を受け取ることと、  
 前記テキスト入力に対応する検索パターンを自動的に生成することと、

前記複数のテキストエンティティの各々に関連付けられたパターンを決定することと、  
 前記テキスト入力に対応する前記検索パターンを、前記複数のテキストエンティティに  
 関連付けられたパターンと対応付けることと、  
 前記対応付けに基づいて、前記複数のテキストエンティティに関連付けられた前記パター  
 ンから1つまたは複数の合致パターンを識別することと、  
 前記1つまたは複数の合致パターンに対応する関連テキストエンティティを前記複数のテ  
 キストエンティティから抽出することと  
 を含むステップを1つまたは複数のプロセッサを含むコンピュータに行わせるコンピュ  
 ータ実行可能命令のセットが記憶された非一時的コンピュータ可読記憶媒体。

【請求項16】

前記コンピュータ実行可能命令のセットは、前記1つまたは複数のプロセッサを含む前記  
 コンピュータに、

ドメインベースの抽出であって、前記ドメインベースの抽出を行うことは、  
 前記複数のテキストエンティティの各々に関連付けられたドメインベースの名前を決定す  
 ることであって、前記ドメインベースの名前は、前記複数のテキストエンティティの各々  
 をデータデポジトリに記憶されたドメインベースの名前のリストと対応付けることに基  
 づいて決定される、決定すること

を含む、ドメインベースの抽出、

位置ベースの抽出であって、前記位置ベースの抽出を行うことは、

前記入力ファイルに関連付けられたタイプを識別すること、

前記入力ファイルに関連付けられた前記タイプに基づいて、前記入力ファイルにおける1  
 つまたは複数の関連テキストエンティティの位置を決定すること、および、

前記位置に基づいて、前記入力ファイルから前記1つまたは複数の関連テキストエンティ  
 ティを抽出することであって、前記位置は、前記データリポジトリに記憶された、前記入  
 力ファイルに関連付けられたタイプに対応するテンプレートに基づいて予め決定される、  
 抽出すること、

を含む、位置ベースの抽出、

品詞（POS）ベースの抽出であって、前記品詞（POS）ベースの抽出を行うことは、  
 前記複数のテキストエンティティに関連付けられたPOSを識別することであって、前記P  
 OSは、名詞、代名詞、動詞、副詞、形容詞、接続詞、前置詞、および間投詞のうちの1  
 つである、識別すること、および、

POSが名詞として識別される1つまたは複数のテキストエンティティを抽出すること、  
 POSが識別されない1つまたは複数のユニークなテキストエンティティを識別すること

、  
 前記1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティ  
 ティタイプを決定することであって、前記エンティティタイプは、値エンティティおよび  
 パターンエンティティのうちの1つであり、各値エンティティは関連付けられた値を有し  
 、各パターンエンティティは関連付けられたパターンを有する、決定すること、

各値エンティティについて、

前記値エンティティの前記関連付けられた値と同じ関連付けられた値を有する1つまたは  
 複数のテキストエンティティを自動的に識別すること、

各パターンエンティティについて、

前記パターンエンティティに関連付けられたパターンに対応する検索パターンを自動的に  
 生成すること、

前記パターンエンティティに対応する前記検索パターンを、前記複数のテキストエンティ  
 ティに関連付けられたパターンと対応付けること、

前記対応付けに基づいて、前記複数のテキストエンティティに関連付けられた前記パター  
 ンから1つまたは複数の合致パターンを決定すること、および、

前記1つまたは複数の合致パターンに対応する合致テキストエンティティを前記複数のテ  
 キストエンティティから抽出すること

を含む、品詞（POS）ベースの抽出、ならびに、機械学習（ML）ベースの抽出であって、前記MLベースの抽出を行うことは、分類ベースのMLモデルを用いて、履歴データに基づいて推奨を生成することを含む、機械学習（ML）ベースの抽出のうちの少なくとも1つを行わせる、請求項15に記載の非一時的コンピュータ可読記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、概してデータ抽出に関し、より詳細には、1つまたは複数のデータ抽出アプローチを用いて入力ファイルのコンテンツからテキスト情報を抽出するための方法およびシステムに関する。

10

【背景技術】

【0002】

最近、テキスト抽出技法が重要性を帯びてきている。例えば、光学文字認識（Optical Character Recognition; OCR）などの抽出技法は、ユーザが画像またはポータブル・ドキュメント・フォーマット（Portable Document Format; PDF）ファイルなどのファイルからテキストデータを抽出することを可能とし得る。さらに、関連情報を抽出し、抽出された関連情報を用いて表現を生成することが望ましい場合がある。

【発明の概要】

【発明が解決しようとする課題】

20

【0003】

いくつかの利用可能な技法によれば、入力ファイルがタグを含む場合、または識別されたテキストにおいてパターンを識別することができる場合に、入力ファイルからテキスト情報を抽出し、抽出された情報を用いて表現を生成することが可能であり得る。しかしながら、タグが存在しない場合があるまたはパターンを識別することができない複雑な文書においては、関連情報を抽出したり、表現を生成したりすることが難しい。

【課題を解決するための手段】

【0004】

本発明による方法およびシステムは、入力ファイルのコンテンツから情報を抽出するものであって、入力ファイルからテキストデータを識別することと、複数のテキストエンティティから関連テキストエンティティを識別するために、ユーザからテキスト入力を受け取ることと、テキスト入力に対応する検索パターンを自動的に生成することと、複数のテキストエンティティの各々に関連付けられたパターンを決定することと、テキスト入力に対応する前記検索パターンを、前記複数のテキストエンティティに関連付けられたパターンと対応付けることと、対応付けに基づいて、複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンを識別することと、1つまたは複数の合致パターンに対応する関連テキストエンティティを複数のテキストエンティティから抽出することを含む。識別されたテキストデータは、複数のテキストエンティティを含む。

30

【0005】

その他、本願が開示する課題とその解決方法は、発明を実施するための形態の欄、および図面の記載によって明らかにされる。

40

【0006】

本開示に組み込まれてその一部を構成する添付の図面は、例示的实施形態を示し、説明と共に開示の原理を解説する役割を果たす。

【図面の簡単な説明】

【0007】

【図1】本開示のいくつかの実施形態に係る、入力ファイルのコンテンツから情報を抽出するための例示的システムの機能ブロック図である。

【図2】本開示のいくつかの実施形態に係る、自然言語処理（Natural Language Processing; NLP）メタデータ抽出フレームワークの機能ブロック図である。

50

【図3】本開示のいくつかの実施形態に係る、スペルチェックシステムのブロック図である。

【図4】本開示のいくつかの実施形態に係る、推奨システムのブロック図である。

【図5】本開示のいくつかの実施形態に係る、メタデータ更新システムのブロック図である。

【図6】本開示の様々な実施形態に係る、入力ファイルのコンテンツから情報を抽出する方法のフローチャートである。

【図7】本開示の様々な実施形態に係る、入力ファイルのコンテンツから情報を抽出する方法のフローチャートである。

【図8】本開示の様々な実施形態に係る、入力ファイルのコンテンツから情報を抽出する方法のフローチャートである。

【図9】本開示のいくつかの実施形態に係る、情報を抽出する対象であるノイジーなエンティティを有する例示的な入力ファイルのスナップショットである。

【図10】本開示の様々な実施形態に係る、情報を抽出する対象である別の例示的な入力ファイルのスナップショットである。

【発明を実施するための形態】

【0008】

添付の図面を参照して、例示的な実施形態を説明する。都合の良い場合には、同じまたは同様の部分を示すために、図面の全体を通じて同じ参照番号が用いられる。開示の原理の例および特徴を本明細書で説明するが、開示の実施形態の趣旨および範囲から逸脱しない限りにおいて、修正、改変、および他の実装が可能である。以下の詳細な説明は単に例示的なものとみなされることが意図されており、真の範囲および趣旨は、以下の特許請求の範囲によって示される。

【0009】

ここで図1を参照すると、本開示のいくつかの実施形態に係る、入力ファイルのコンテンツから情報を抽出するための例示的なシステム100の機能ブロック図が示されている。システム100は、テキスト識別デバイス102を含んでよい。いくつかの実施形態において、テキスト識別デバイス102は、入力ファイルのコンテンツからテキストデータを識別してよい。例えば、入力ファイルは、画像またはポータブル・ドキュメント・フォーマット (Portable Document Format; PDF) ファイルであってよい。識別されたテキストデータは、複数のテキストエンティティを含んでよい。したがって、いくつかの実施形態において、テキスト識別デバイス102は、入力ファイルからテキストデータを識別するために、光学文字認識 (Optical Character Recognition; OCR) 技法を用いてよい。代替的な実施形態において、テキスト識別デバイス102は、入力ファイルからテキストデータを識別するために、当技術分野において知られている任意の他の技法を用いてよい。

【0010】

システム100は、テキスト識別デバイス102と通信可能な情報抽出デバイス104をさらに含んでよい。例えば、情報抽出デバイス104は、テキスト識別デバイス102によって入力ファイルから識別された複数のテキストエンティティから、関連情報を抽出するように構成されてよい。いくつかの実施形態において、情報を抽出するべく、情報抽出デバイス104は、複数のテキストエンティティから関連テキストエンティティを識別するためにユーザからテキスト入力を受け取り、テキスト入力に対応する検索パターンを自動的に生成してよい。情報抽出デバイス104はさらに、複数のテキストエンティティの各々に関連付けられたパターンを決定し、テキスト入力に対応する検索パターンを複数のテキストエンティティに関連付けられたパターンと対応付けてよい。情報抽出デバイス104はさらに、対応付けに基づいて複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンを識別し、1つまたは複数の合致パターンに対応する関連テキストエンティティを複数のテキストエンティティから抽出してよい。

【0011】

システム100は、ディスプレイ114をさらに含んでよい。いくつかの実施形態において、情報抽出デバイス104は、1つまたは複数のプロセッサ110およびコンピュータ可読媒体（例えばメモリ）112を含んでよい。コンピュータ可読記憶媒体112は、1つまたは複数のプロセッサ110により実行されたとき、本開示の態様により、入力ファイルのコンテンツに基づいて表現を生成することを1つまたは複数のプロセッサ110に行わせる命令を記憶してよい。コンピュータ可読記憶媒体112はまた、システム100によって取り込まれ、処理され、かつ／または必要とされ得る様々なデータ（例えば、識別されたテキストデータ、値エンティティ、データ、パターンエンティティデータなど）を記憶してよい。システム100は、ディスプレイ114を介してアクセス可能なユーザインターフェース116を介してユーザとやり取りしてよい。システム100はまた、様々なデータを送るまたは受け取るために、通信ネットワーク108を介して1つまたは複数の外部デバイス106とやり取りしてよい。外部デバイス106は、リモートサーバ、デジタルデバイス、または別のコンピューティングシステムを含んでよいが、これらに限定されないものであってよい。

#### 【0012】

ここで図2を参照すると、本開示のいくつかの実施形態に係る自然言語処理（Natural Language Processing; NLP）メタデータ抽出フレームワーク200の機能ブロック図が示されている。NLPメタデータ抽出フレームワーク200は、システム100に、特に情報抽出デバイス104に実装されてよい。NLPメタデータ抽出フレームワーク200は、入力ファイルのコンテンツから情報を抽出するための様々な機能を行う様々なモジュールを含んでよい。いくつかの実施形態において、抽出エンジン202が、入力ファイル204のコンテンツからテキストデータを抽出／識別してよい。入力ファイル204は、フラットファイル206、またはPDFファイル208、またはデータベースファイル210、またはOCR出力212であってよい。

#### 【0013】

抽出エンジン202はさらに、ドメインルールデータベース214から1つまたは複数のドメインルールを受け取ってよい。さらに、いくつかの実施形態において、ドメインルールデータベース214は、（従来の手続き型コードによるものではなく）「if-then」ルールに基づくものであってよい。ドメインルールは、複数のテキストエンティティの各々に関連付けられたドメインベースの名前を決定するために用いられてよいことに留意されたい。ドメインベースの名前は、複数のテキストエンティティの各々を、データデポジトリ（図1では不図示）に記憶されたドメインベースの名前のリストと対応付けることに基づいて決定されてよい。例えば、識別されたデータが、ある国のある都市のPINコードに属するテキストエンティティを含む場合、ドメインベースのデータベースは、例えばPINコードをデータリポジトリに記憶されたドメインベースの名前のリストと対応付けることにより、ドメイン名、すなわちそのPINコードを有する都市の名前を決定してよい。

#### 【0014】

抽出エンジン202はさらに、位置ルールデータベース216から1つまたは複数の位置ルールを受け取ってよい。いくつかの実施形態において、位置ルールは、OCR出力に基づいて、入力ファイルにおけるテキストデータからのエンティティに関連付けられた位置を決定するために用いられてよいことに留意されたい。位置は、入力ファイルに関連付けられたタイプを最初に識別し、入力ファイルに関連付けられたタイプに基づいて入力ファイルにおける1つまたは複数の関連テキストエンティティの位置を決定したときに、決定されてよい。関連テキストエンティティの位置が決定されると、関連テキストエンティティが位置に基づいて入力ファイルから抽出されてよい。1つまたは複数の位置ルールは、表を含む入力データ、すなわち表形式で配置された入力データの場合にはより関連性が高いものであり得るが、1つまたは複数の位置ルールは、任意の他のタイプの入力データ、例えばフリーテキストを有する文書にも同様に適用可能であってよいことに留意されたい。

10

20

30

40

50

## 【0015】

例として、入力ファイルのタイプは、「A a d h a r」カードのものであってよい。「A a d h a r」カードのような文書においては、データが特定の形式で存在していてよいことをさらに理解されたい。したがって、関連テキストエンティティの位置は、データリポジトリに記憶された、入力ファイルに関連付けられたタイプに対応するテンプレートに基づいて予め決定されてよい。例えば、(対象の人物の)テキストエンティティ「名前」および「住所」の位置は、文書の中程であってよく、この情報は、システム100のデータベースにおいて取得可能であり記憶されていてよい。特に、位置は、(x, y)座標に基づくものであってよい。上記の例において、テキストエンティティ「名前」および「住所」の位置は、(x, y)座標(5, 10)～(15, 15)によって画定される領域内であってよい。したがって、文書のタイプが「A a d h a r」カードとして識別されると、システム100は、上記位置にアクセスして、その位置から関連テキストエンティティ(すなわち「名前」および「住所」)を抽出してよい。

10

## 【0016】

抽出エンジン202はさらに、第1のNLPルールデータベース218から1つまたは複数の第1のNLPルールを受け取ってよい。第1のNLPルールデータベース218は、正解値/属性について提供されるNLPルールに基づいて正規表現(Regular Expression; Reg ex)を作成してよいことに留意されたい。当業者には理解されるように、Reg exは、検索パターンを記述するための特別な文字列であってよい。

## 【0017】

抽出エンジン202はさらに、第2のNLPルールデータベース220から1つまたは複数の第2のNLPルールを受け取ってよい。いくつかの実施形態において、第2のNLPルールデータベース220は、正解値/属性について品詞(part of speech; POS)を定義してよいことに留意されたい。例えば、複数のテキストエンティティに関連付けられたPOSが識別されてよい。当業者には理解されるように、POSは、名詞、代名詞、動詞、副詞、形容詞、接続詞、前置詞、および間投詞のうちの1つであってよい。POSを識別すると、POSが名詞である複数のテキストエンティティから1つまたは複数のユニークなテキストエンティティが選択されてよい。POSが名詞として識別された1つまたは複数の(選択された)ユニークなテキストエンティティは、関連データとして抽出されてよい。

20

30

## 【0018】

例えば、給与明細文書(入力データ)は、対象の人物の名前を含む複数のテキストエンティティを含んでよい。名前を抽出することが、抽出エンジン202の目的であってよい。さらに、名前は、名詞のPOSカテゴリに属するものであることが理解されよう。したがって、POSが名詞であるテキストエンティティが選択されることになる。したがって、必要とされる名前が、選択されたエンティティから抽出されることになる。

## 【0019】

いくつかの実施形態において、複数のテキストエンティティに関連付けられたPOSを識別する試行が行われると、POSが識別されない1つまたは複数のユニークなテキストエンティティが識別されてよい。1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティティタイプが決定されてよい。エンティティタイプは、値エンティティまたはパターンエンティティであってよい。各値エンティティは関連付けられた値を有してよく、各パターンエンティティは関連付けられたパターンを有する。

40

## 【0020】

例えば、識別されたテキストエンティティは、電子メールIDまたは油田名(例えばそれぞれ「j o h n @ g m a i l . c o m」または「A B - 1 2 3」)のようなコードネームを示すテキスト文字の組み合わせから構成されていてよい。したがって、そのようなテキストエンティティについては、POSが識別されない場合がある。したがって、そのようなテキストエンティティは、値エンティティまたはパターンエンティティのいずれかとして識別されてよい。したがって、電子メールID(「j o h n @ g m a i l . c o m」)

50

は値エンティティとして決定されてよく、値は電子メールIDである。これは、電子メールIDのパターンが通常は典型的かつ既知のパターンであり、全ての電子メールIDが同じまたは類似のパターンを共有し得るためであり得る。一方、油田名に対応する「AB-123」のようなテキストエンティティは、パターンエンティティとして決定されてよい。これは、油田名が同じパターンを共有し得るが、そのようなパターンは一般的に知られていない場合がある、すなわちNLPメタデータ抽出フレームワーク200がそのようなパターンを予め記憶していない場合があるためである。値エンティティを識別すると、値エンティティの関連付けられた値と同じ関連付けられた値を有する入力ファイルにおけるテキストエンティティが、自動的に識別されてよい。換言すると、全ての電子メールIDテキストエンティティが、識別および抽出されてよい。

10

【0021】

パターンエンティティを識別すると、パターンエンティティに関連付けられたパターン (Regex) に対応する検索パターンが、自動的に生成されてよい。さらに、入力ファイルにおける複数のテキストエンティティの各々に関連付けられたパターンが決定されてよい。パターンエンティティに対応する検索パターン (Regex) は、複数のテキストエンティティに関連付けられたパターンと対応付けられてよく、対応付けに基づいて、複数のテキストエンティティに関連付けられたパターンからの1つまたは複数の合致パターンが決定されてよい。1つまたは複数の合致パターンに対応する合致テキストエンティティが、複数のテキストエンティティから抽出されてよい。

【0022】

例えば、油田に対応するテキストエンティティ「AB-123」について、検索パターン (Regex) 「apd」 (ここで、「a」は1つまたは複数のアルファベットを意味し、「p」は1つまたは複数の句読点を意味し、「d」は1つまたは複数の数字を意味する) が自動的に生成されてよい。さらに、入力ファイルにおける他のテキストエンティティに関連付けられたパターンが決定され、検索パターン「apd」が他のテキストエンティティに関連付けられたパターンと対応付けられてよく、合致パターンが識別されてよい。その後、合致パターンに対応する全てのテキストエンティティが抽出される。これにより、入力ファイルにおける全ての油田名が抽出されてよい。

20

【0023】

いくつかの実施形態において、入力ファイルからのテキストデータが識別されると、第1のNLPルールは、複数のテキストエンティティから関連テキストエンティティを識別するためにユーザからテキスト入力を受け取り、テキスト入力に対応する検索パターンを自動的に生成することを行わせてよい。第1のNLPルールはさらに、複数のテキストエンティティの各々に関連付けられたパターンを決定し、テキスト入力に対応する検索パターンを複数のテキストエンティティに関連付けられたパターンと対応付けることを行わせてよい。第1のNLPルールはさらに、対応付けに基づいて複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンを識別し、1つまたは複数の合致パターンに対応する関連テキストエンティティを複数のテキストエンティティから抽出することを行わせてよい。

30

【0024】

例えば、ユーザは、入力ファイルから抽出された複数のテキストエンティティから関連テキストエンティティを識別するために、テキスト入力「AB-123」を提供してよい。検索パターン (Regex) 「apd」が、入力されたテキストエンティティ「AB-123」について識別されてよい。さらに、入力ファイルにおける他のテキストエンティティに関連付けられたパターンが決定されてよく、検索パターン (Regex) 「apd」が他のテキストエンティティに関連付けられたパターンと対応付けられてよく、合致パターンが識別される。その後、合致パターンに対応する全てのテキストエンティティが抽出されてよい。これにより、入力ファイルにおける全ての油田名が抽出されてよい。

40

【0025】

いくつかの実施形態において、ユーザからのテキスト入力は、「Microsoft E

50

x c e l（登録商標）」シートにおけるエントリの形態で受け取られてよいことに留意されたい。テキスト入力に対応する検索パターン（R e g e x）を自動的に生成し得る1つまたは複数のルールが定義されてよい。その後、1つまたは複数のルールは、入力ファイルの複数のテキストエンティティの各々に関連付けられたパターンを決定し、テキスト入力に対応する検索パターンを複数のテキストエンティティに関連付けられたパターンと対応付け、対応付けに基づいて複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンを識別し、1つまたは複数の合致パターンに対応する関連テキストエンティティを複数のテキストエンティティから抽出してよい。

#### 【0026】

したがって、いくつかの実施形態において、「M i c r o s o f t E x c e l（登録商標）」データは、1つまたは複数の「M i c r o s o f t E x c e l（登録商標）」シート222から受け取られてよい。様々なルール（ドメインルール、NLPルール、位置ルール等）が適用されることで、抽出されたデータ226が得られてよい。いくつかの実施形態において、「M i c r o s o f t E x c e l（登録商標）」マクロが抽出に用いられてよい。さらに、いくつかの実施形態において、同マクロがP y t h o n（言語；登録商標）に埋め込まれてよい。ここで、抽出されたデータ226は、複数のテキストエンティティから抽出された関連テキストエンティティ、または、様々なルールを用いて抽出された関連テキストエンティティに基づいて生成され得る表現を指すものであってよい。いくつかの追加的な実施形態において、上記のアプローチ、すなわちドメインベース、位置ベース、P O Sベース、およびR e g e xベースのアプローチがテキスト抽出を提供することに失敗した場合、機械学習（Machine Learning; ML）ベースのアプローチが用いられてよい。MLベースのアプローチには、人工知能（A I）モデル（すなわちMLモデル）224が用いられてよい。

#### 【0027】

ここで図3を参照すると、本開示のいくつかの実施形態に係るスペルチェックシステム300のブロック図が示されている。いくつかの実施形態において、抽出モジュール302（テキスト識別デバイス102に対応する）により抽出されたテキストエンティティが、スペルチェックシステム300により受け取られてよい。抽出モジュール302により生成されたこれらの表現は、スペルチェックシステム300への入力テキスト304として働いてよい。入力テキスト304は、前処理モジュール306により前処理されてよい。

#### 【0028】

前処理が行われると、コサイン類似度モジュール308により、前処理されたデータに対してコサイン類似度解析が行われてよい。一例として、コサイン類似度モジュール308は、閾値を用いてコサイン類似度解析を行ってよい。最小編集距離モジュール310が、コサイン類似度モジュール308から受け取られたデータに対して最小編集距離解析を行ってよい。いくつかの実施形態において、最小編集距離解析は、フィルタリング済みの単語に対して行われてよい。

#### 【0029】

最大コサイン類似度モジュール312が、最小編集距離モジュール310から受け取られたデータに対して最大コサイン類似度解析を行ってよい。最大コサイン類似度解析は、最小編集距離単語内の単語に基づくものであってよいことに留意されたい。置換モジュール314は、最大コサイン類似度モジュール312により行われた解析に基づいて、不正確な単語を修正された単語で置換してよい。

#### 【0030】

ここで図4を参照すると、本開示のいくつかの実施形態に係る推奨システム400のブロック図が示されている。例えば、推奨システム400は、エンティティタイプ、すなわち値エンティティまたはパターンエンティティを識別するために用いられてよい。いくつかの実施形態において、推奨システム400は、分類ベースの機械学習（ML）モデル402を用いてよい。推奨システム400は、入力データ（NLP抽出・スペルチェックデータ）406を受け取ってよい。入力データは、NLPおよびスペルチェックが行われたデ

ータであってよいことに留意されたい。推奨システム400は、設定（config）ファイル408を含んでよい。configファイル408は、入力データについて値が推奨され得るか、またはパターンが推奨され得るかを決定することをトリガしてよい。

#### 【0031】

推奨システム400は、入力データをMLモデル402に供給してよい。MLモデル402は、アーカイブデータベース404からの履歴データを用いてよい。履歴データに基づいて、MLモデル402は、予測データ410を提供してもよく、または推奨データ412を提供してもよい。したがって、推奨値414が予測データ410に基づいて生成されてもよく、または推奨パターン416が推奨データ412に基づいて生成されてもよい。

#### 【0032】

ここで図5を参照すると、本開示のいくつかの実施形態に係るメタデータ更新システム500のブロック図が示されている。いくつかの実施形態において、メタデータ更新システム500は、NLP抽出済みデータ502を受け取ってよい。メタデータ更新システム500は、システム日時モジュール504を含んでよい。メタデータ更新システム500は、信頼度モジュール506および確率モジュール508をさらに含んでよい。信頼度モジュール506は、抽出されたデータおよびNLPルールにより生成された表現の正確度に関連付けられた信頼度スコアを決定してよい。確率モジュール508は、信頼度スコアに基づいて、抽出されたデータおよびNLPルールにより生成された表現の精度に関連付けられた正確度スコアを決定してよい。換言すると、確率モジュール508は、抽出されたデータまたはNLPルールを用いて生成された表現がどれだけ正確／有用であるかについての確率を決定してよい。

#### 【0033】

確率モジュール508により決定される確率が十分に高い（すなわち閾値よりも大きい）場合、抽出されたデータおよびNLPルールを用いて生成された表現、すなわちルールベース値514が提供されてよい。一方、確率モジュール508により決定される確率が十分に高くない（すなわち閾値よりも小さい）場合、抽出されたデータおよびMLモデルを用いて生成された表現、すなわち機械学習（ML）値516が提供されてよい。前述の通り、ML値は、アーカイブデータベース510に記憶され得る履歴データを用いてMLモデルにより生成されてよい。抽出されたデータおよび生成された表現は、最終データベース512に記憶されてよい。

#### 【0034】

ここで図6を参照すると、本開示の一実施形態に係る、入力ファイルのコンテンツから情報を抽出する方法600のフローチャートが示されている。理解されるように、方法600は、テキスト抽出のRegexベースの方法を提供してよい。ステップ602において、入力ファイルからテキストデータが識別されてよい。識別されたテキストデータは、複数のテキストエンティティを含んでよい。いくつかの例において、入力ファイルは、画像ファイルおよびポータブル・ドキュメント・フォーマット（PDF）ファイルのうちの少なくとも1つを含んでよいことに留意されたい。

#### 【0035】

ステップ604において、複数のテキストエンティティから関連テキストエンティティを識別するために、テキスト入力ユーザから受け取られてよい。例えば、ユーザは、テキストを「Microsoft Excel（登録商標）」シートにおけるエントリとして提供してよい。ステップ606において、テキスト入力に対応する検索パターン（Regex）が、自動的に生成されてよい。ステップ608において、複数のテキストエンティティの各々に関連付けられたパターンが決定されてよい。ステップ610において、テキスト入力に対応する検索パターンが、複数のテキストエンティティに関連付けられたパターンと対応付けられてよい。ステップ612において、複数のテキストエンティティに関連付けられたパターンからの1つまたは複数の合致パターンが、対応付けに基づいて識別されてよい。ステップ614において、1つまたは複数の合致パターンに対応する関連テキストエンティティが、複数のテキストエンティティから抽出されてよい。

10

20

30

40

50

## 【0036】

加えて、いくつかの実施形態においては、複数のテキストエンティティから1つまたは複数の合致パターンに対応する関連テキストエンティティを抽出すると、抽出された関連テキストエンティティを用いて、表現が生成されてよい。

## 【0037】

さらに、例えば方法600では関連エンティティを抽出することができない場合に、関連エンティティを抽出するために機械学習（ML）モデルが用いられてよい。換言すると、データ（関連テキストエンティティ）を抽出することができない場合、MLベースのアプローチが開始してよい。

## 【0038】

ここで図7を参照すると、本開示の別の実施形態に係る、入力ファイルのコンテンツから情報を抽出する方法700のフローチャートが示されている。ステップ702において、入力ファイルからテキストデータが識別されてよい。識別されたテキストデータは、複数のテキストエンティティを含んでよい。入力ファイルは、画像ファイルおよびポータブル・ドキュメント・フォーマット（PDF）ファイルのうち少なくとも1つを含んでよいことに留意されたい。

## 【0039】

ステップ704において、複数のテキストエンティティの各々に関連付けられたドメインベースの名前が決定されてよい。ドメインベースの名前は、複数のテキストエンティティの各々を、データデポジトリに記憶されたドメインベースの名前のリストと対応付けることに基づいて決定されてよい。ステップ706において、入力ファイルに関連付けられたタイプが識別されてよい。ステップ708において、入力ファイルに関連付けられたタイプに基づいて、入力ファイルにおける1つまたは複数の関連テキストエンティティの位置が決定されてよい。ステップ710において、位置に基づいて、1つまたは複数の関連テキストエンティティが入力ファイルから抽出されてよい。

## 【0040】

ステップ712において、複数のテキストエンティティに関連付けられた品詞（POS）が識別されてよい。POSは、名詞、代名詞、動詞、副詞、形容詞、接続詞、前置詞、および間投詞のうちの一つであってよいことを理解されたい。ステップ714において、POSが名詞として識別される1つまたは複数のテキストエンティティが決定されてよい。

## 【0041】

ステップ716において、POSが識別されない1つまたは複数のユニークなテキストエンティティが選択されてよい。ステップ718において、1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティティタイプが決定されてよい。エンティティタイプは、値エンティティおよびパターンエンティティのうちの一つであってよい。さらに、各値エンティティは関連付けられた値を有してよく、各パターンエンティティは関連付けられたパターンを有する。図4に関連して説明したように、1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティティタイプを決定するために、（MLベースの）推奨システム400が用いられてよい。

## 【0042】

ステップ720において、各値エンティティについて、値エンティティの関連付けられた値と同じ関連付けられた値を有する1つまたは複数のテキストエンティティが、自動的に識別されてよい。各パターンエンティティについて、ステップ722～728が行われてよい。したがって、ステップ722において、パターンエンティティに関連付けられたパターンに対応する検索パターン（Regex）が、自動的に生成されてよい。ステップ724において、パターンエンティティに対応する検索パターンが、複数のテキストエンティティに関連付けられたパターンと対応付けられてよい。ステップ726において、対応付けに基づいて、複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンが決定されてよい。ステップ728において、1つまたは複数の合致パターンに対応する合致テキストエンティティが、複数のテキストエンティティから抽出

10

20

30

40

50

されてよい。

#### 【0043】

方法700は、方法600と組み合わせて行われてよいことに留意されたい。例示的シナリオにおいて、方法600は、方法700のステップ710が完了した後、方法700のステップ712が始まる前に開始してよい。

#### 【0044】

ここで図8を参照すると、本開示の別の実施形態に係る、入力ファイルのコンテンツから情報を抽出する方法800のフローチャートが示されている。ステップ802において、入力ファイルのコンテンツからテキストデータが抽出されてよい。抽出されたテキストデータは、複数のテキストエンティティを含んでよい。

10

#### 【0045】

ステップ804において、複数のテキストエンティティの各々に関連付けられたドメインベースの名前が決定されてよい。ドメインベースの名前は、複数のテキストエンティティの各々を、データデポジトリに記憶されたドメインベースの名前のリストと対応付けることに基づいて決定されてよい。例えば、PINコードが取得可能な場合、例えばPINコードをナレッジベースに記憶された（PINコードおよび関連付けられた都市名の）リストと対応付けることにより、ドメインベースのアプローチを用いてPINコードから都市が決定されてよい。同様に、都市名（例えばデリー）が識別された場合、国名（例えばインド）が決定されてよい。いくつかの実施形態においては、それに応じて、第2の複数のエンティティの各々にタグが割り当てられてよい。

20

#### 【0046】

ステップ806において、ドメインベースの名前を決定することに成功したか否かを確認するためのチェックが行われてよい。ドメインベースの名前を決定することに成功した場合、本方法はステップ822（「Yes」の経路）に進んでよく、このとき、識別されたテキストエンティティが抽出されてよい。さらに、抽出されたテキストエンティティを用いて、表現が生成されてよい。ステップ806において、ドメインベースの名前が決定されない場合、本方法はステップ808（「No」の経路）に進んでよい。換言すると、ドメインベースの名前を識別することができない場合、本方法は、次の代替的アプローチの実行に進んでよい。

30

#### 【0047】

ステップ808において、複数のテキストエンティティの各々の位置に基づいて、所定のフィールドの値が決定されてよい。位置は、（xおよびy）座標の形態であってよい。一例として、位置は、データリポジトリに記憶されたテンプレートに基づいて予め決定されてよい。例えば、いくつかのシナリオにおいて、特定の名前の値がその名前の横または特定の名前ヘッダの下に記載されていてよく、これは（xおよびy）座標に基づいて識別可能であってよい。例えば、「Adharカード」、または「PANカード」、または運転免許証において、エンティティ、例えばユーザの名前が特定の位置において取得可能であってよく、そこから名前を識別することができる。位置の技法は、表の場合に特に有用であり得る。これは、表におけるデータは構造化された形式で利用可能であり、エンティティの位置が容易かつ正確に特定され得るためである。したがって、入力ファイルに関連付けられたタイプ、すなわち入力ファイルが「Adharカード」であるか、または「PANカード」であるか、または運転免許証であるかが、まず識別されてよい。さらに、入力ファイルに関連付けられたタイプに基づいて、入力ファイルにおける1つまたは複数の関連テキストエンティティの位置が決定されてよい。

40

#### 【0048】

ステップ810において、入力ファイルにおける1つまたは複数の関連テキストエンティティの位置を決定することに成功したか否かを確認するためのチェックが行われてよい。位置を決定することに成功した場合、方法800は再びステップ822（「Yes」の経路）に進んでよく、このとき、入力ファイルからの1つまたは複数の関連テキストエンティティが位置に基づいて抽出されてよく、抽出されたテキストエンティティを用いて表現

50

が生成されてよい。位置が決定されないと認められる場合、ステップ810において、方法800はステップ812に進んでよい。

#### 【0049】

ステップ812において、関連テキストエンティティを抽出するために、入力ファイルに対してRegexベースのアプローチが試行されてよい。この目的で、複数のテキストエンティティから関連テキストエンティティを識別するためにユーザからテキスト入力を受け取られてよく、テキスト入力に対応する検索パターンが自動的に生成されてよい。換言すると、Regexが生成されてよい。さらに、複数のテキストエンティティの各々に関連付けられたパターンが生成されてよく、テキスト入力に対応する検索パターン (Regex) が、複数のテキストエンティティに関連付けられたパターンと対応付けられてよく、対応付けに基づいて、複数のテキストエンティティに関連付けられたパターンからの1つまたは複数の合致パターンが識別されてよい。いくつかの実施形態において、方法800は、Regexに対応する、抽出されたテキストデータから関連エンティティを決定する確率を算出することを含んでよい。

10

#### 【0050】

ステップ814において、Regexベースのアプローチを適用することに成功したかを決定するためのチェックが行われてよい。Regexベースのアプローチを適用することに成功した場合、方法800はステップ822 (「Yes」の経路) に進んでよく、このとき、1つまたは複数の合致パターンに対応する関連テキストエンティティが複数のテキストエンティティから抽出されてよく、抽出されたテキストエンティティを用いて表現が生成されてよい。Regexベースのアプローチが成功しない場合、方法800はステップ816 (「No」の経路) に進んでよい。

20

#### 【0051】

ステップ816において、品詞 (POS) ベースのアプローチが試行されてよい。例えば、複数のテキストエンティティの各々に関連付けられたPOSが決定されてよい。換言すると、ステップ816において、エンティティがいずれのPOSと関連付けられ得るか、例えばPOSが名詞、副詞、形容詞等であり得るかを決定するための試行が行われてよい。さらに、ステップ816において、1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティティタイプを決定するための試行が行われてよい。この目的で、POSが識別されない1つまたは複数のユニークなテキストエンティティが選択されてよい。例えば、アルファベット、数字、句読点等の組み合わせを有するユニークなエンティティ (上記で述べた電子メールIDまたは油田名など) については、POSが識別されない場合がある。そのようなテキストエンティティについて、1つまたは複数のユニークなテキストエンティティの各々に関連付けられたエンティティタイプが決定されてよい。エンティティタイプは、値エンティティおよびパターンエンティティのうちの1つであってよい。各値エンティティは関連付けられた値を有してよく、各パターンエンティティは関連付けられたパターンを有してよい。さらに、各値エンティティについて、値エンティティの関連付けられた値と同じ関連付けられた値 (例えば電子メールID) を有する1つまたは複数のテキストエンティティを自動的に識別するための試行が行われてよい。

30

40

#### 【0052】

ステップ818において、試行が成功したか (すなわちPOSベースのアプローチが成功したか) を決定するためのチェックが行われてよい。試行が成功したと認められる場合、方法800はステップ822 (「Yes」の経路) に進んでよく、このとき、値エンティティの関連付けられた値と同じ関連付けられた値を有する1つまたは複数のテキストエンティティ (関連テキストエンティティ) が抽出されてよく、抽出されたテキストエンティティを用いて表現が生成されてよい。

#### 【0053】

さらに、各パターンエンティティについて、パターンエンティティに関連付けられたパターンに対応する検索パターンを自動的に生成し、パターンエンティティに対応する検索パ

50

ターンを複数のテキストエンティティに関連付けられたパターンと対応付け、対応付けに基づいて複数のテキストエンティティに関連付けられたパターンから1つまたは複数の合致パターンを決定するための試行が行われてよい。ステップ818において、試行が成功したか（すなわちPOSベースのアプローチが成功したか）を決定するためのチェックが行われてよい。試行が成功したと認められる場合、方法800はステップ822（「Yes」の経路）に進んでよく、このとき、1つまたは複数の合致パターンに対応する合致テキストエンティティが複数のテキストエンティティから抽出されてよく、抽出されたテキストエンティティを用いて表現が生成されてよい。ステップ818において、POSベースのアプローチが成功しないと認められる場合、方法800はステップ820（「No」の経路）に進んでよい。

10

## 【0054】

ステップ820において、機械学習（ML）モデルを用いて複数のテキストエンティティから関連テキストエンティティ（情報）を抽出するための試行が行われてよい。MLモデルは、関連テキストエンティティのタイプを抽出するようにまず訓練されてよいことが理解されよう。したがって、MLモデルの訓練に基づいて、MLベースの分類が、関連テキストエンティティを抽出するように入力ファイルに適用されてよい。したがって、ステップ822において、MLベースの分類に基づいて関連テキストエンティティが抽出されてよく、抽出されたテキストエンティティを用いて表現が生成されてよい。

## 【0055】

いくつかの実施形態において、様々なステップ（すなわちステップ804、808、812、816、および820）は、上記で説明した順序で行われてもよく、または任意の他の順序で行われてもよいことに留意されたい。

20

## 【0056】

さらに、下記のシナリオにおいてNLPベースのアプローチが用いられてよいことに留意されたい。

(i) 検索キー（すなわちユーザからのテキスト入力）が利用可能であり、かつそれに関連付けられたテキスト抽出（すなわち関連エンティティ）が入力ファイルにおいて予期される場合。

(ii) 検索キー（すなわちユーザからのテキスト入力）が利用可能でないが、それに関連付けられたテキスト抽出（すなわち関連エンティティ）が入力ファイルにおいて予期される場合。

30

(iii) 検索キー（すなわちユーザからのテキスト入力）が利用可能であるが、それに関連付けられたテキスト抽出（すなわち関連エンティティ）が入力ファイルにおいて予期されない場合。

## 【0057】

さらに、各々のアプローチにおいて信頼度スコア（確率スコア）が算出されてよいことに留意されたい。例えば、「ドメインベース」のアプローチについては、検索キー（すなわちユーザからのテキスト入力）が入力ファイルにおける残りのテキストエンティティと厳密に合致する場合に、最大の確率スコア（例えば1）が得られてよい。さらに、「位置ベース」のアプローチでは、確率スコア（ $p$ ）は、検索キーの合致の確率（ $w$ ）および位置

40

( $l$ )の識別、すなわち $p = w * l$ に基づいてよい。Regexベースのアプローチについては、位置識別についての確率は0または1のいずれかであってよく、抽出の確率スコア（ $e$ ）は1または0であり得る。さらに、抽出される値は複数（ $n$ ）であり得、それらの値のうちの1つが選択されてよい。ゆえに、選択される値の確率は「 $1/n$ 」であってよく、信頼度スコアは $e * (1/n)$ であってよい。「位置ベース」および「Regexベース」のアプローチの両方が用いられる場合、各アプローチに50%の重みが与えられてよい。ゆえに、信頼度スコアは $(w * l) / 2 + (e * (1/n)) / 2$ であってよい。MLベースのアプローチについては、確率スコアが混同行列（正解率）に基づいて導出されてよい。

## 【0058】

50

さらに、以下のアプローチ、すなわち純粋な単語合致アプローチ、ルールベースのアプローチ（すなわちドメインベース、位置ベース、POSベース、およびRegexベースのアプローチのうちの一つまたは複数）、およびMLベースのアプローチが用いられる場合、複合確率スコア（P）が算出されてよい。

したがって、

複合確率スコア（P）＝（P1＊0.3）＊（P3＊0.7）、または

複合確率スコア（P）＝（P1＊0.3）＊（P2＊0.7）

式中、

P1＝純粋な単語合致アプローチの確率スコア、

P2＝ルールベースのアプローチの確率スコア、および

P3＝MLベースのアプローチの確率スコア。

【0059】

<ケースシナリオ1>

ここで図9を参照すると、ノイジーなエンティティ902が入力ファイル900に存在する例示的な入力ファイル900のスナップショットが示されている。例えば、エンティティの第1の部分は汚れまたはかすれ等により判読できない場合があり、エンティティの第2の部分は「United」という単語を含む場合がある。したがって、実際のエンティティは「United Airlines（登録商標）」または「United States」等である可能性がある。このエンティティを決定するために、一つまたは複数のアプローチ（すなわちドメインベース、位置ベース、POSベース、Regexベース、またはMLベース）が実際のエンティティの決定に用いられてよい。例えば、ドメインベース、位置ベース、POSベース、およびRegexベースのアプローチが抽出を提供することに失敗した場合、MLベースのアプローチが開始してよく、分類ベースの機械学習（ML）モデルを用いて履歴データに基づいてテキスト抽出を提供してよい。MLモデルは、決定論的モデルおよび／または確率論的モデルであってよい。

【0060】

<ケースシナリオ2>

ここで図10を参照すると、入力ファイル「Local Order」1000のスナップショットが示されており、そこから様々なテキストエンティティを抽出することが必要とされ得る。図10に示すように、入力ファイル「Local Order」1000は表形式構造である。

【0061】

例えば、属性「RO Number」についての値を抽出するために、RO dateの属性を対応付けることにより、「RO Number」がドメインディクショナリデータベース（ルックアップテーブル）1002から抽出されてよい。

【0062】

属性Part Numberを抽出するために、入力ファイル1000に存在するキー「Part Number」が識別され、またはユーザ入力として受け取られてよい。その後、位置ベースのアプローチにより、データベースに記憶されたテンプレート情報を用いて、Part Numberの可能性のある位置が決定されてよい。例えば、位置ベースのアプローチにより、Part Numberが入力ファイル1000の「右」側に存在し得ることが示唆されてよい。したがって、必要とされるPart Number「32145643」が、入力ファイル「Local Order」1000の表形式構造の右のセル座標から抽出されてよい。

【0063】

入力ファイル1000に存在する電子メールIDの属性を抽出するために、Regexベースのアプローチが用いられてよい。したがって、電子メールの属性に関連付けられたキー（Regex）が識別され、またはデータベースから受け取られてよい。したがって、Regexは、電子メールID「ajay.thakur@gmail.com」に対応する「apapapa」（「a」は一つまたは複数のアルファベットを意味し、「p」は

10

20

30

40

50

1つまたは複数の句読点を意味し、「d」は1つまたは複数の数字を意味する)であってよい。したがって、合致パターン (Regex) を有するテキストエンティティ、すなわち電子メールIDのテキストエンティティが識別されてよく、値「ajay. thakur@gmail.com」が入力ファイル1000の表形式構造の対応するセルから抽出されてよい。

【0064】

入力ファイル1000から属性Nameを抽出するために、POSベースのアプローチ用いられてよい。例えば、まず属性 (Name) について存在するキーが入力ファイル1000において識別されてよい。その後、例えば固有表現抽出 (Named Entity Recognizer; NER) タガーを用いて、入力ファイル1000における識別されたテキストエンティティのPOSが決定されてよい。属性Nameについてデータベースにおいて指定されるNERタガーが検索されてよい。したがって、名詞のPOSを有する属性Name («Ajay Thakur») に対応するテキストエンティティが抽出されてよい。したがって、抽出されたテキストエンティティ「Ajay Thakur」を用いる表現が生成されてよい。

【0065】

入力ファイルのコンテンツから情報を抽出するための1つまたは複数のテキスト抽出技法が、上記に開示されている。上記の技法は、従来の技法に対して1つまたは複数の利点を提供する。例えば、これらの技法は、テキスト抽出を行うための様々なアプローチ、すなわちドメインベース、位置ベース、POSベース、Regexベースのアプローチ、およびMLベースのアプローチを提供する。したがって、単一のアプローチ、または任意の順序で適用される複数のアプローチを用いて、テキスト抽出を行うことができる。さらに、これらの技法は、抽出された情報を用いて表現を生成することを可能とする。さらに、これらの技法は、タグが存在しない、または抽出された情報においてパターンを識別することができない複雑な文書においても、情報を抽出し、抽出された情報を用いて表現を生成することを可能とする。

【0066】

本開示と整合する実施形態の実装において、1つまたは複数のコンピュータ可読記憶媒体が利用されてよい。コンピュータ可読記憶媒体とは、プロセッサにより読み取り可能な情報またはデータが記憶され得る任意のタイプの物理メモリを指す。よって、コンピュータ可読記憶媒体は、本明細書に記載の実施形態と整合するステップまたは段階をプロセッサに実行させるための命令を含む、1つまたは複数のプロセッサが実行するための命令を記憶してよい。「コンピュータ可読媒体」という用語は、有形の物品を含み、搬送波および過渡信号を除外する、すなわち非一時的なものであると理解されるべきである。例としては、ランダムアクセスメモリ (RAM)、リードオンリメモリ (ROM)、揮発性メモリ、不揮発性メモリ、ハードドライブ、CD-ROM、DVD、フラッシュドライブ、ディスク、および任意の他の既知の物理記憶媒体が挙げられる。

【0067】

開示および例は単に例示的なものとみなされることが意図されており、開示の実施形態の真の範囲および趣旨は、以下の特許請求の範囲によって示される。

10

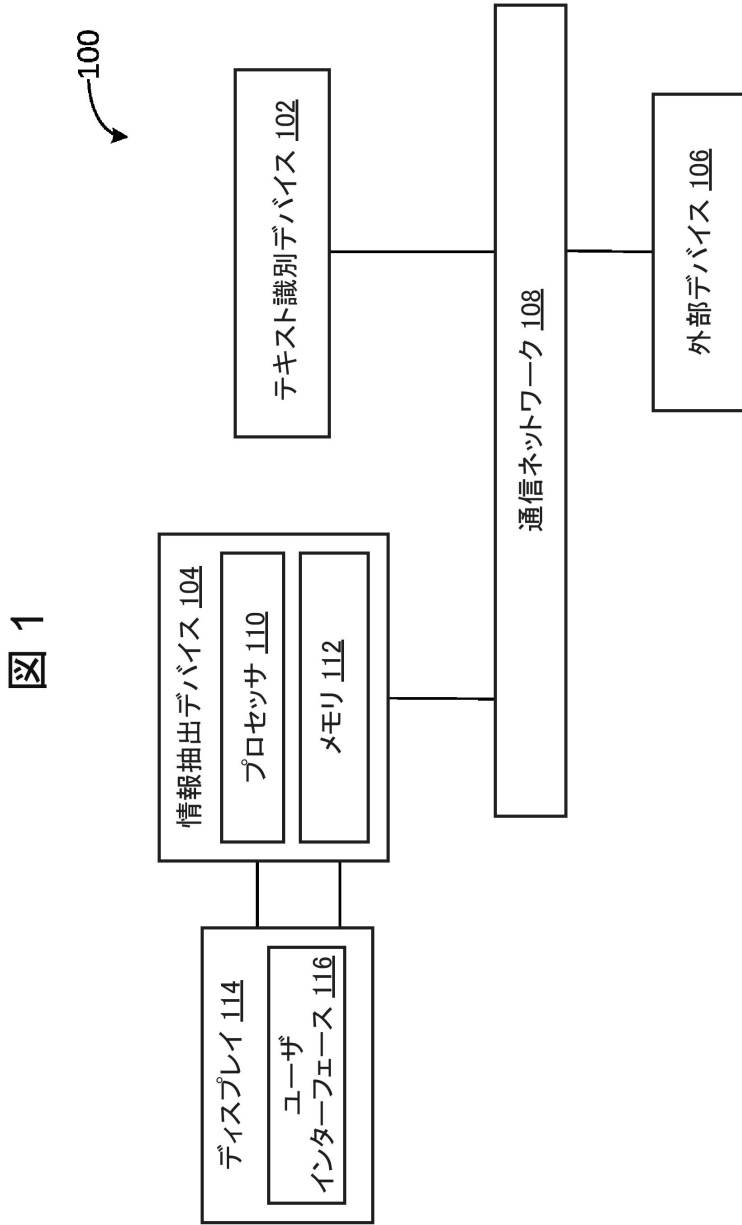
20

30

40

50

【図1】



10

20

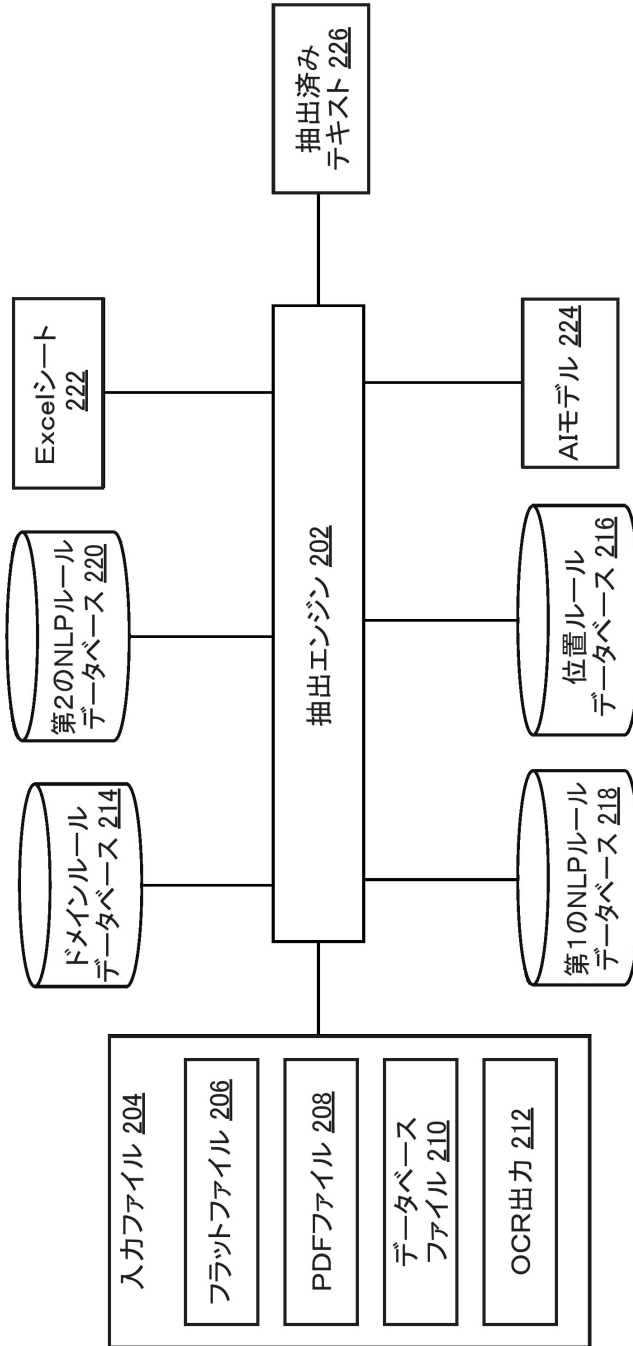
30

40

50

【図2】

図2  
200



10

20

30

40

50

【図 3】

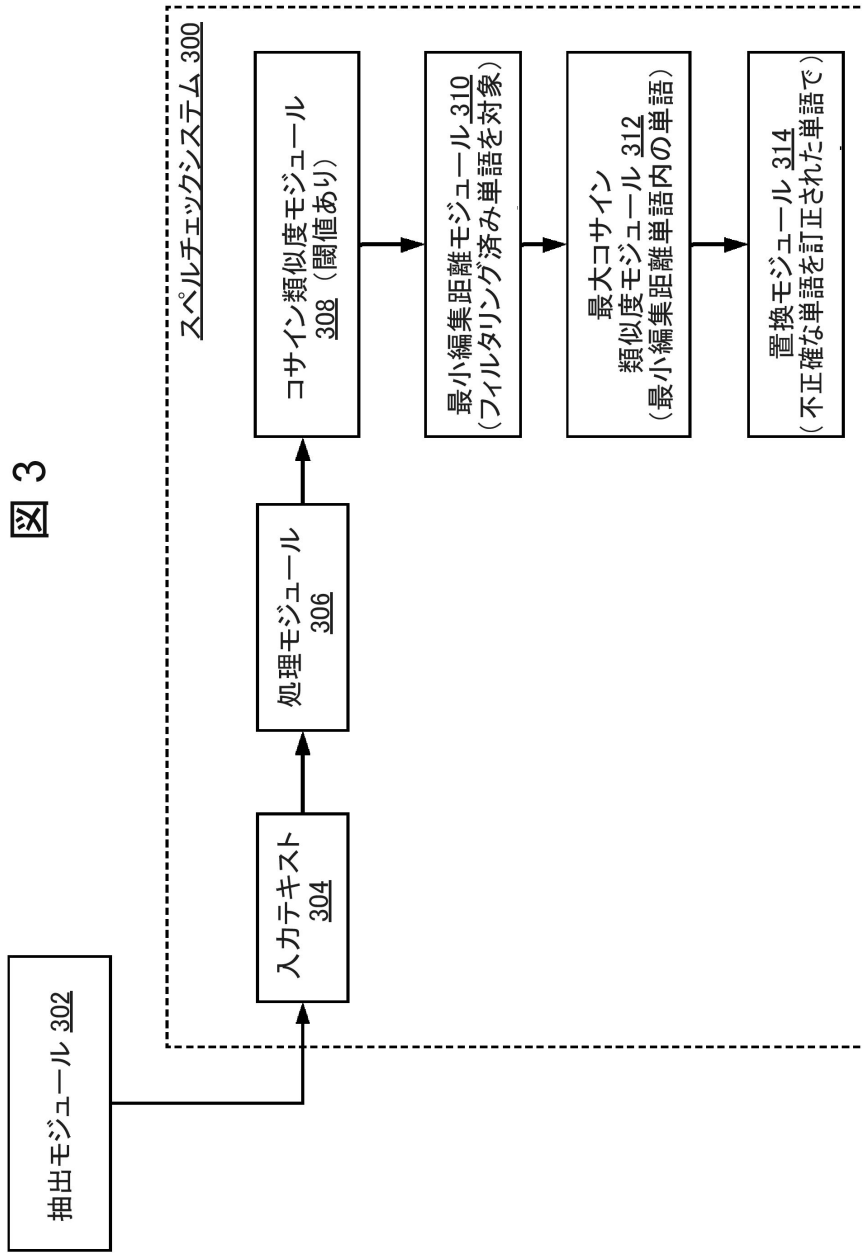


図 3

10

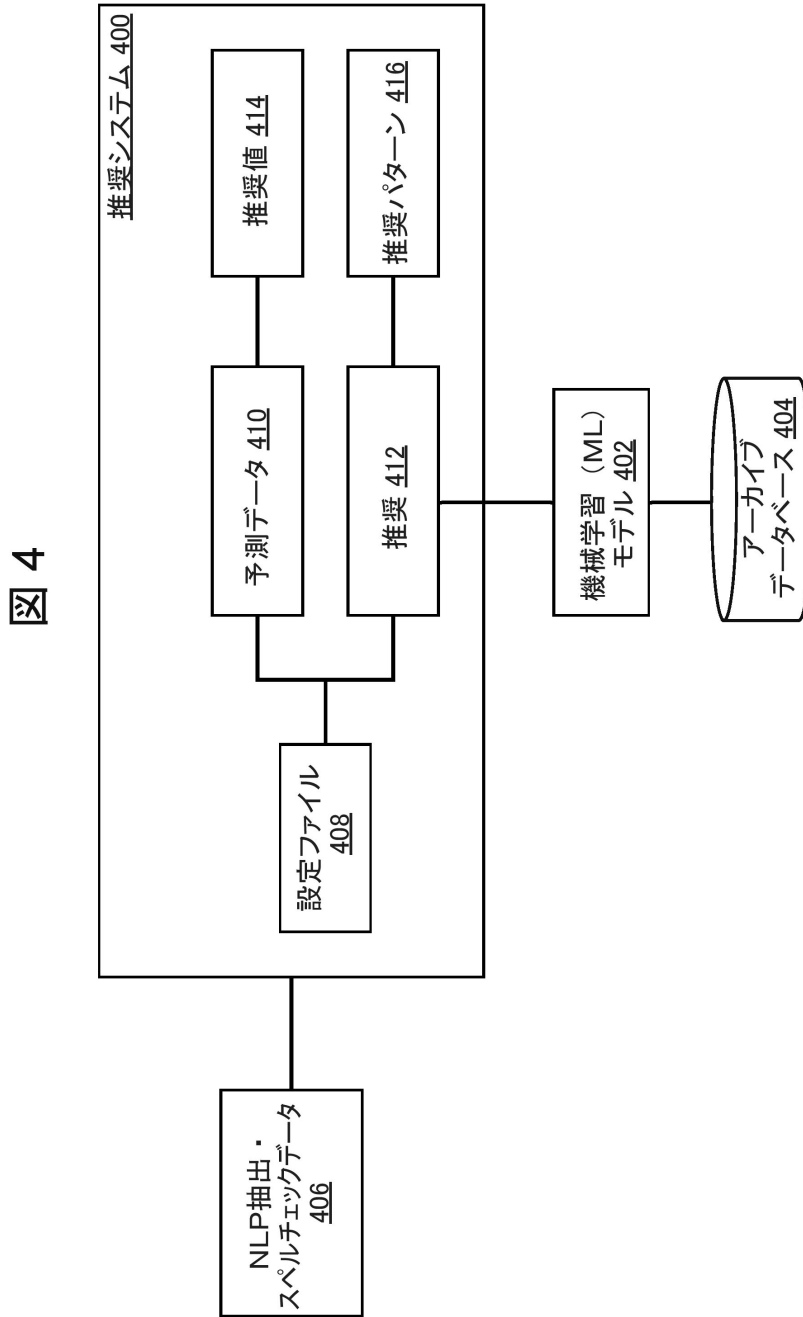
20

30

40

50

【図4】



10

20

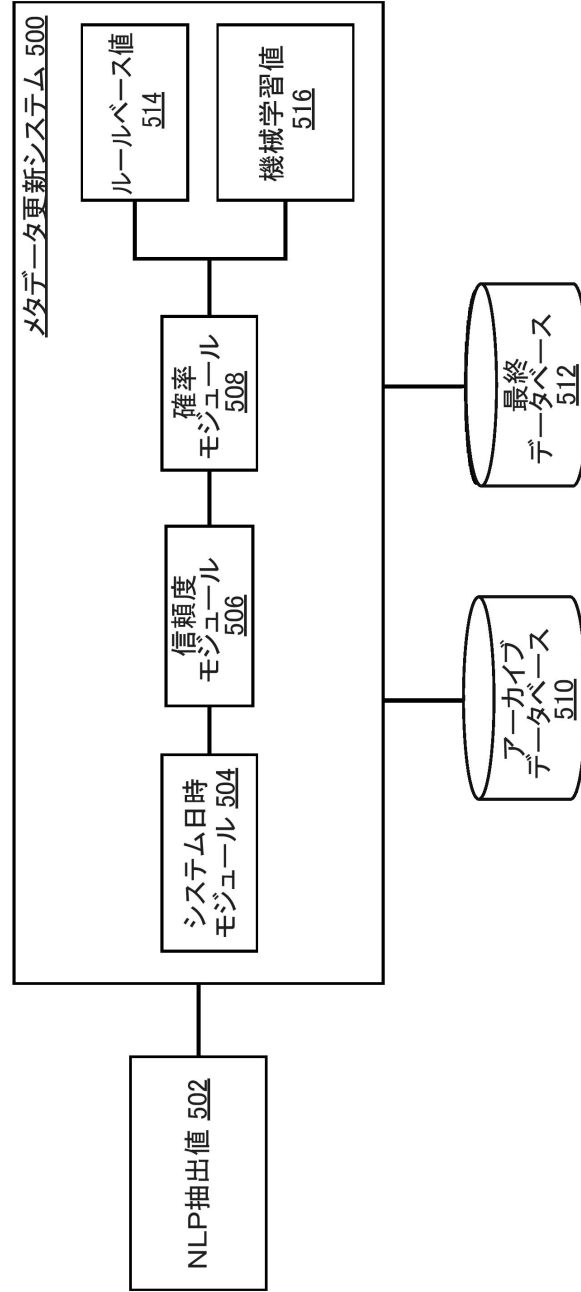
30

40

50

【図5】

図5



10

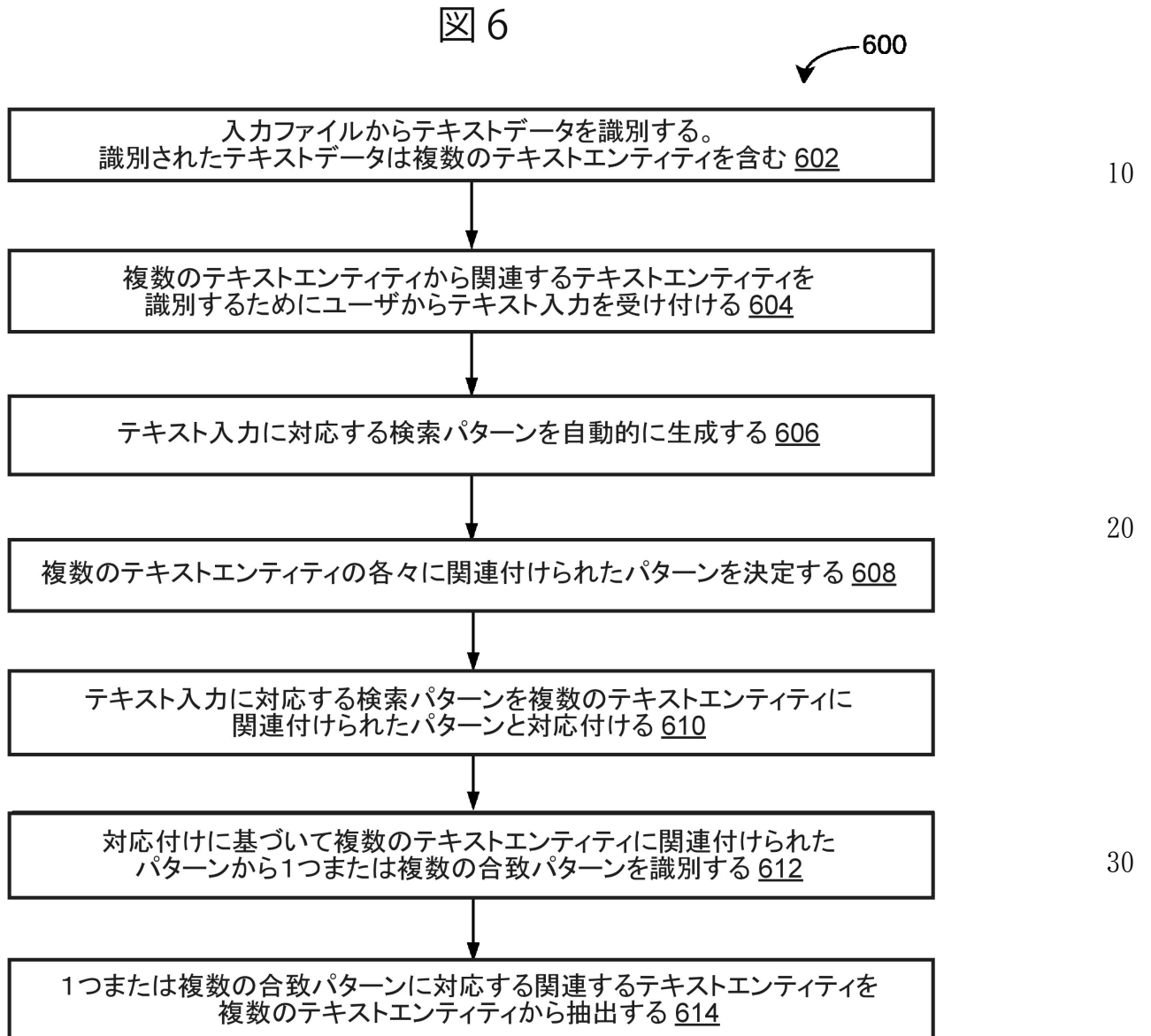
20

30

40

50

【図6】



10

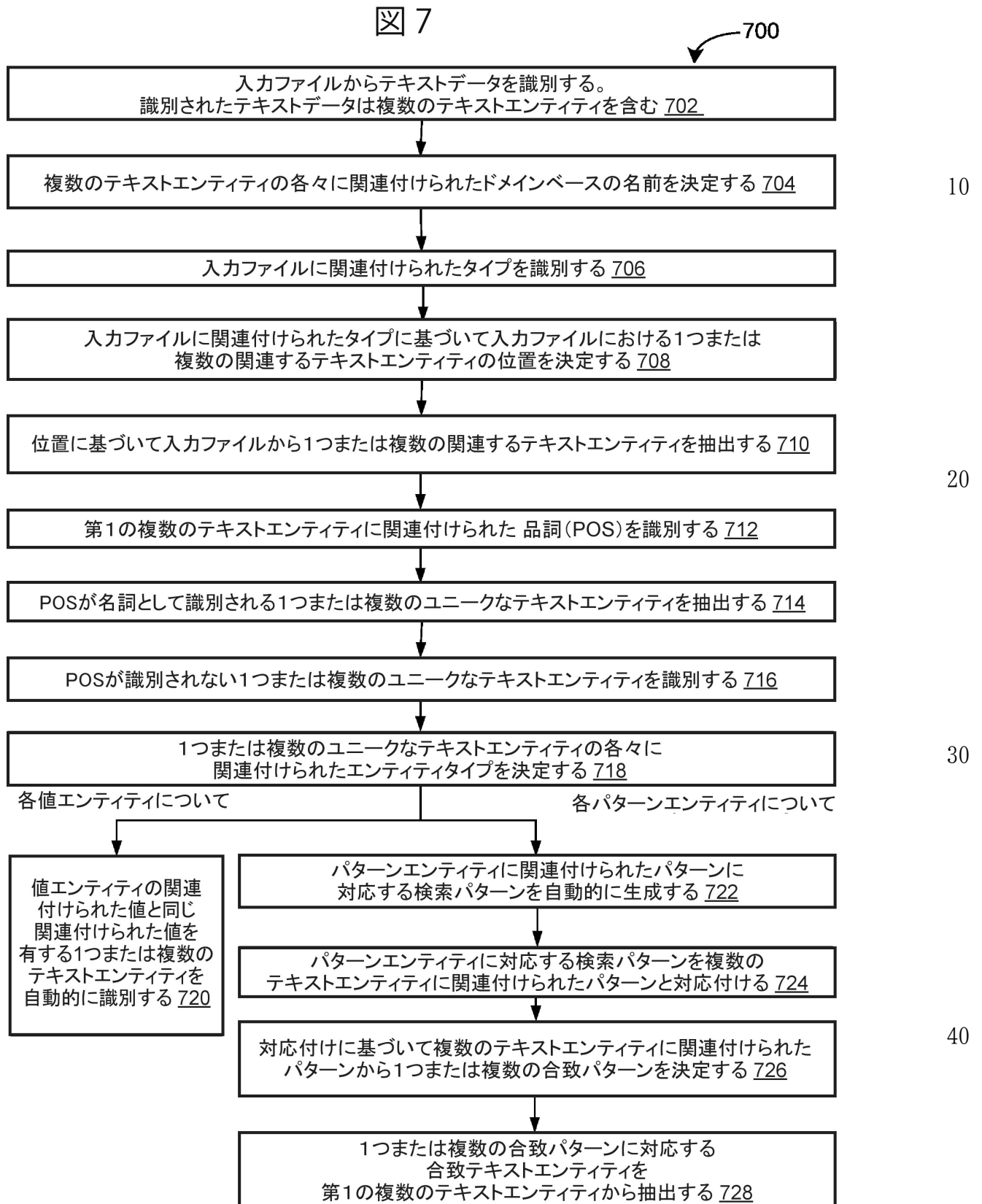
20

30

40

50

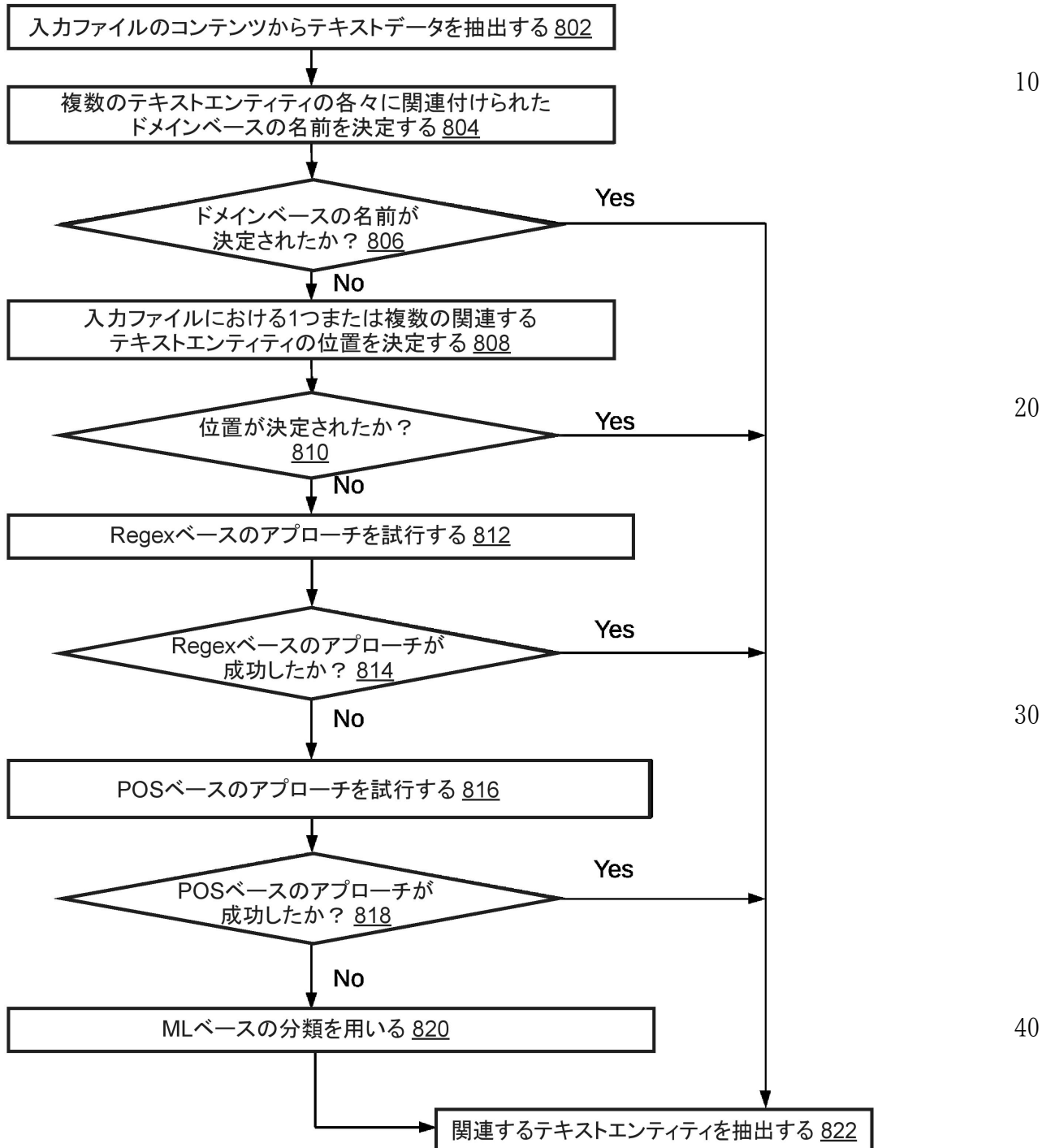
【図7】



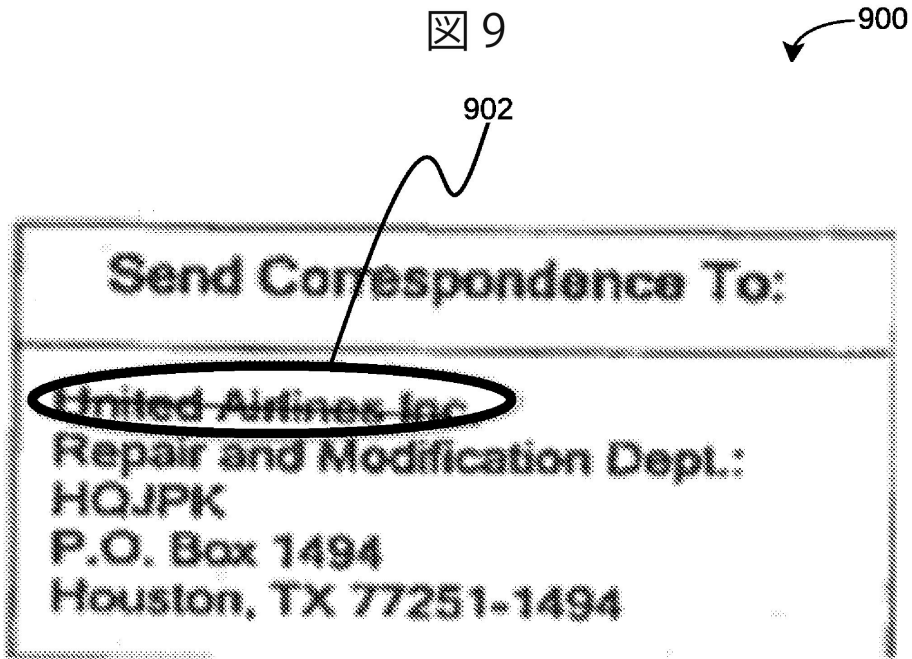
【図 8】

図 8

800



【図 9】



10

20

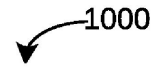
30

40

50

【図 10】

図 10



10

**Local Purchase Order**

R.O. Number		256789214200	Account ID:	125673
Date Required		12/02/2020	Special Instruction:	None
R.O. Date		14/02/2020	Warranty:	None
Mode of Payment	Serial No.	234761	Customer Info	
Online	Part No.	32145643	Name: Ajay Thakur	
	IMEI No.	987345823421	Email: ajay.thakur@gmail.com	
	Tracking Tag:	1345276	Phone No. 8567434651 Fax No: 3214-292-2123	

20

This document constitutes an agreement between the vendor and the buyer. See terms and conditions of this purchase listed on the reverse side

Item No.	Specification	Unit	Quantity	Unit Price	Total Value
1	Mobile Display Capacitive Touch Screen-PS123Y1	1	1	5500	5500
2					
3					
4					
				Tax	50.85
				Total Value	5550.85

30

40

50

## 【国際調査報告】

INTERNATIONAL SEARCH REPORT		International application No. PCT/IB2020/062535
A. CLASSIFICATION OF SUBJECT MATTER G06F16/00,G06K9/46,G06K9/72,G06F40/00 Version=2021.01 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, G06K 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: TotalPatent One, IPO Internal Database Keywords:extract, data, pattern, document, search		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 20190286898 A1 (FIRST AMERICAN FINANCIAL CORPORATION) 19 September 2019 (19/09/2019) (abstract, paragraphs [0005]-[0007], [0029]-[0037], [0043]-[0048], figures 2, 4 and 5	1-16
<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 21-04-2021		Date of mailing of the international search report 21-04-2021
Name and mailing address of the ISA/ Indian Patent Office Plot No.32, Sector 14,Dwarka,New Delhi-110075 Facsimile No.		Authorized officer Prashant Singh 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/IB2020/062535

Citation	Pub.Date	Family	Pub.Date
US 2019286898 A1	19-09-2019	AU 2015203150 A1	21-01-2016
		CA 2895917 A1	30-12-2015

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, 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 シン, マドウスダン  
 インド 560100 バンガロール, エレクトロニック シティー—1, ファースト クロス, アジメラ ストーン パーク, ハウス ナンバー ビー—603

(72)発明者 ハルダー, カウシク  
 インド 700010 コルカタ, ベレガタ, 152—ラジャ ラジェンドラ ラル ミトラ ロード, アイディアアル アパートメント, ブロッカー—エー, フラット—4エー

(72)発明者 パナパリ ベンカタ ラメッシュ ラユル, ニルマル  
 インド 531035 ヴィシャーカパトナム, サババラム マンダル, ビハインド ガバメント ジュニア カレッジ, ゴボルナガー, ハウス ナンバー 6—37

(72)発明者 ゴーシュ ダスティダル, アリトラ  
 インド 700082 コルカタ, ハリデブプール, 181 パナマリ パネルジー ロード, フラット ナンバー 06

(72)発明者 シャ, アジェイ  
 インド 736182, アリプルドゥアル, ジャイガオン ディストリクト, エヌ・エス ロード, ビサイド ポスト オフィス, シャンプー シャ内, ハウス ナンバー 298

Fターム(参考) 5B091 AA15