

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-173522
(P2018-173522A)

(43) 公開日 平成30年11月8日(2018.11.8)

(51) Int.Cl.		F I		テーマコード (参考)
G 1 0 L 15/22 (2006.01)		G 1 0 L 15/22	2 0 0 V	
G 1 0 L 15/10 (2006.01)		G 1 0 L 15/10	5 0 0 Z	

審査請求 未請求 請求項の数 9 O L (全 12 頁)

(21) 出願番号 特願2017-71168 (P2017-71168)
(22) 出願日 平成29年3月31日 (2017. 3. 31)

(71) 出願人 000002897
大日本印刷株式会社
東京都新宿区市谷加賀町一丁目1番1号
(74) 代理人 100096091
弁理士 井上 誠一
(72) 発明者 松本 征二
東京都新宿区市谷加賀町一丁目1番1号
大日本印刷株式会社内

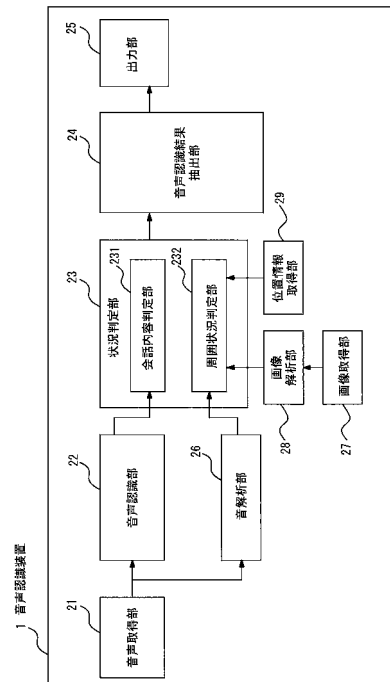
(54) 【発明の名称】 音声認識装置、音声認識方法、及びプログラム

(57) 【要約】

【課題】 状況に適した音声認識結果を得ることが可能な音声認識装置、音声認識方法、及びプログラムを提供する。

【解決手段】 音声認識装置 1 は、会話の音声データを取得し、音声認識部 2 2 により音声データをテキストに変換し、音声認識結果として複数の変換候補を求める。また、状況判定部 2 3 は、会話の内容のジャンルや、周囲の音の特徴等を解析したりすることにより、コンテキストや周囲の状況（場所、シーン、会話の目的）等を求め、会話の状況を判定する。音声認識結果抽出部 2 4 は、複数の変換候補の中から会話の状況に適した変換候補を抽出し、出力する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

音声データを取得する音声取得手段と、
取得した音声データを認識し、音声認識結果として複数の変換候補を求める音声認識手段と、
音声データ取得時の状況を判定する状況判定手段と、
前記音声認識手段により求めた複数の変換候補の中から前記状況判定手段により判定した状況に適した変換候補を抽出する抽出手段と、
を備えることを特徴とする音声認識装置。

【請求項 2】

状況と関連付けられた単語を記憶する記憶手段を備え、
前記状況判定手段は、前記単語を音声データから抽出し、前記単語に基づいて前記状況を判定することを特徴とする請求項 1 に記載の音声認識装置。

【請求項 3】

前記状況判定手段は、更に、取得される周囲の音に基づいて前記状況を判定することを特徴とする請求項 2 に記載の音声認識装置。

【請求項 4】

前記状況判定手段は、更に、取得される画像に基づいて前記状況を判定することを特徴とする請求項 2 または請求項 3 に記載の音声認識装置。

【請求項 5】

前記状況判定手段は、更に、取得される位置情報に基づいて前記状況を判定することを特徴とする請求項 2 から請求項 4 のいずれかに記載の音声認識装置。

【請求項 6】

前記状況判定手段により判定した会話の状況の履歴を記憶する記憶手段を更に備え、
前記抽出手段は、前記記憶手段に記憶された会話の状況の履歴に基づき、前記音声認識手段により求めた複数の変換候補の中から会話の状況に適した変換候補を抽出することを特徴とする請求項 1 から請求項 5 のいずれかに記載の音声認識装置。

【請求項 7】

前記抽出手段は、設定された尤度に基づいて抽出候補に優先付けを行うことを特徴とする請求項 1 から請求項 6 のいずれかに記載の音声認識装置。

【請求項 8】

コンピュータが、
音声データを取得するステップと、
取得した音声データを認識し、音声認識結果として複数の変換候補を求めるステップと、
音声データ取得時の状況を判定するステップと、
前記複数の変換候補の中から前記状況に適した変換候補を抽出するステップと、
を含むことを特徴とする音声認識方法。

【請求項 9】

コンピュータを、
音声データを取得する音声取得手段、
取得した音声データを認識し、音声認識結果として複数の変換候補を求める音声認識手段、
音声データ取得時の状況を判定する状況判定手段、
前記音声認識手段により求めた複数の変換候補の中から前記状況判定手段により判定した状況に適した変換候補を抽出する抽出手段、
として機能させるためのプログラム。

【発明の詳細な説明】**【技術分野】**

【0001】

本発明は、音声認識装置、音声認識方法、及びプログラムに関し、詳細には、音声認識精度を向上するための技術に関する。

【背景技術】

【0002】

従来より、スマートフォンやカーナビゲーションシステム等のユーザインターフェースとして音声入力を用いるものが利用されている。また、AI (Artificial Intelligence ; 人工知能) や対話型ロボットが普及し始め、機器と人が音声によってコミュニケーションをとる機会が増加している。そのため、精度のよい音声認識システムが要望されている。例えば、特許文献1には、ユーザが発話する言葉に含まれるユーザの意図や感情を推定することで、ユーザとの間の対話をより円滑にする機能を有する音声対話装置について記載されている。

10

【0003】

また、従来の音声認識処理では、音声認識の結果、複数の変換候補がある場合にはシステムが第1候補を自動的に選択する方法等が一般的である。例えば、「おすすめのたいけんはありますか」という音声が入力された際の音声認識結果としては、「No. 1 : おすすめの大剣は何ですか」、「No. 2 : おすすめの体験は何ですか」、「No. 3 : お勧めの体験はなんですか」、「No. 4 : お奨めの体験は何ですか」、「No. 5 : おすすめの大剣はなんですか」のような複数の変換候補が得られる。この場合、観光のシーンであればNo. 2、3、4が正しい変換候補となるが、ゲームについての会話中であれば、No. 1、あるいは5が正しい変換となる。したがって適切な音声の文書変換を行うには、シーンや会話の流れ等の状況を把握し、それに応じた候補を出力する必要がある。

20

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2006-313287号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

例えば、上述の特許文献1には、ユーザの感情や生理状態等を音声、画像、生理信号等の非言語情報を用いて入力音声の意図(質問、詰問、疑い)を推定し判断する技術について開示されている。しかしながら、特許文献1は、ユーザの感情や生理状態によるユーザ情報のみから対話内容を推定するため、簡単な応答文など2者択一形式の認識には対応できるものの、前述のように多義的であって複雑な会話内容に対しては適切な認識を行うことが困難である。

30

【0006】

本発明は、このような課題に鑑みてなされたもので、会話のシーンやコンテキストといった状況に適した音声認識結果を得ることが可能な音声認識装置、音声認識方法、及びプログラム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】

40

【0007】

前述した課題を解決するため第1の発明は、音声データを取得する音声取得手段と、取得した音声データを認識し、音声認識結果として複数の変換候補を求める音声認識手段と、音声データ取得時の状況を判定する状況判定手段と、前記音声認識手段により求めた複数の変換候補の中から前記状況判定手段により判定した状況に適した変換候補を抽出する抽出手段と、を備えることを特徴とする音声認識装置である。

【0008】

第1の発明によれば、会話の音声データを認識し、音声認識結果として複数の変換候補を求め、会話の状況(シーンやコンテキスト)を判定し、複数の変換候補の中から会話の状況に適した変換候補を抽出する。これにより、シーンやコンテキストといった会話の状

50

況に適した音声認識結果を得ることが可能となり、音声認識精度の高い音声認識装置を提供できる。

【0009】

第1の発明において、前記状況判定手段は、会話の内容に基づいて前記状況を求めることが望ましい。これにより、会話の内容そのものから状況を推定し、適切な音声認識結果を抽出可能となる。また、前記状況判定手段は、更に、取得される周囲の音に基づいて前記状況を求めてもよい。これにより、会話が行われている状況を周囲の音の特徴から取得できるため、より正確に音声認識結果を得ることができる。

【0010】

また、前記状況判定手段は、更に、取得される画像に基づいて前記状況を求めてもよい。更に、前記状況判定手段は、更に、取得される位置情報に基づいて前記状況を求めてもよい。音の特徴のみならず、画像や位置情報等からも会話の状況を求めることで、状況を正確に求めることが可能となり、これにより正確な音声認識結果を得ることが可能となる。

10

【0011】

また、前記状況判定手段により判定した会話の状況の履歴を記憶する記憶手段を更に備え、前記抽出手段は、前記記憶手段に記憶された会話の状況の履歴に基づき、前記音声認識手段により求めた複数の変換候補の中から会話の状況に適した変換候補を抽出することが望ましい。これにより、会話の状況の履歴から、話の流れを認識することが可能となり、コンテキストに適した音声認識結果をより正確に得ることが可能となる。

20

【0012】

第2の発明は、コンピュータが、会話の音声データを取得するステップと、取得した音声データを認識し、音声認識結果として複数の変換候補を求めるステップと、会話の状況を判定するステップと、前記複数の変換候補の中から前記会話の状況に適した変換候補を抽出するステップと、を含むことを特徴とする音声認識方法である。

【0013】

第2の発明によれば、コンピュータは、会話の音声データを認識し、音声認識結果として複数の変換候補を求め、会話の状況を判定し、複数の変換候補の中から会話の状況に適した変換候補を抽出する。これにより、シーンやコンテキストといった会話の状況に適した音声認識結果を得ることが可能となり、音声認識精度を向上させることが可能となる。

30

【0014】

第3の発明は、コンピュータを、会話の音声データを取得する音声取得手段、取得した音声データを認識し、音声認識結果として複数の変換候補を求める音声認識手段、会話の状況を判定する状況判定手段、前記音声認識手段により求めた複数の変換候補の中から前記状況判定手段により判定した会話の状況に適した変換候補を抽出する抽出手段、として機能させるためのプログラムである。

【0015】

第3の発明により、コンピュータを第1の発明の音声認識装置として機能させることが可能となる。

【発明の効果】

40

【0016】

本発明により、シーンやコンテキストといった状況に適した音声認識結果を得ることが可能な音声認識装置、音声認識方法、及びプログラムを提供できる。

【図面の簡単な説明】

【0017】

【図1】音声認識装置1の構成図

【図2】音声認識装置1として機能するコンピュータ10の構成図

【図3】音声認識装置1が実行する音声認識処理の流れを示すフローチャート

【図4】音声認識の変換候補と会話のジャンルとを関連づけたデータであるジャンルデータ5の例

50

【図5】会話内容による状況判定結果、及び音による状況判定結果の具体例

【図6】本発明に係る音声認識装置1を利用した音声認識システム100の例

【発明を実施するための形態】

【0018】

以下、図面に基づいて本発明の好適な実施形態について詳細に説明する。

【0019】

図1は、本発明に係る音声認識装置1の構成を示す図である。音声認識装置1は、音声取得部21、音声認識部22、状況判定部23、音声認識結果抽出部24、出力部25、及び音解析部26を備える。また、これらの構成に加え、画像取得部27、画像解析部28、及び位置情報取得部29を備えてもよい。

10

【0020】

図2は、音声認識装置1として機能させるコンピュータ10の構成例を示す図である。図2に示すように、コンピュータ10は、制御部11、記憶部12、メディア入出力部13、周辺機器I/F部14、入力部15、表示部16、通信制御部17、マイク18等がバス19を介して接続されて構成される。コンピュータ10を音声認識装置1として機能させる場合、コンピュータ10の制御部11は、図1に示す音声認識装置1の各部（音声取得部21、音声認識部22、状況判定部23、音声認識結果抽出部24、出力部25、音解析部26、画像取得部27、画像解析部28、及び位置情報取得部29）の機能を記述したプログラムを実行する。

【0021】

制御部11は、CPU（Central Processing Unit）、ROM（Read Only Memory）、RAM（Random Access Memory）等により構成される。

20

CPUは、記憶部12、ROM、記録媒体等に格納されるプログラムをRAM上のワークメモリ領域に呼び出して実行し、バス19を介して接続された各部を駆動制御する。ROMは、コンピュータ10のブートプログラムやBIOS等のプログラム、データ等を恒久的に保持する。RAMは、ロードしたプログラムやデータを一時的に保持するとともに、制御部11が各種処理を行うために使用するワークエリアを備える。制御部11は、上記プログラムを読み出して実行することにより、図1に示す音声認識装置1の各部（音声取得部21、音声認識部22、状況判定部23、音声認識結果抽出部24、出力部25、及び音解析部26等）として機能する。すなわち、制御部11はマイク18または通信制御部17等から入力された音声データを取得し、取得した音声データについて後述する音声認識処理（図3参照）を実行する。

30

【0022】

記憶部12は、例えば、ハードディスクドライブ等の記憶装置である。記憶部12には制御部11が実行するプログラムや、プログラム実行に必要なデータ、オペレーティングシステム等が格納されている。これらのプログラムコードは、制御部11により必要に応じて読み出されてRAMに移され、CPUに読み出されて実行される。

【0023】

メディア入出力部13は、例えば、CD、DVD、MO等の各種記録媒体（メディア）のドライブ装置であり、メディアに対してデータの入出力（書込み／読み出し）を行う。

40

【0024】

周辺機器I/F（インタフェース）部14は、周辺機器を接続させるためのポートであり、周辺機器I/F部14を介して周辺機器とのデータの送受信を行う。周辺機器I/F部14は、USB等で構成されており、通常複数の周辺機器I/Fを有する。周辺機器との接続形態は有線、無線を問わない。

【0025】

入力部15は、例えば、キーボード、マウス等のポインティング・デバイス、テンキー等の入力装置であり、入力されたデータを制御部11へ出力する。

表示部16は、例えば液晶パネル、CRTモニタ等のディスプレイ装置と、ディスプレ

50

イ装置と連携して表示処理を実行するための論理回路（ビデオアダプタ等）で構成され、制御部 11 の制御により入力された表示情報をディスプレイ装置上に表示させる。なお、入力部 15 及び表示部 16 は、表示画面にタッチパネル等の入力装置を一体的に設けたタッチパネルディスプレイとしてもよい。

【0026】

通信制御部 17 は、通信制御装置、通信ポート等を有し、ネットワーク 3 等との通信を制御する。

マイク 18 は、音声を収集し、音声データとして制御部 11 に入力する。

バス 19 は、各装置間の制御信号、データ信号等の授受を媒介する経路である。

【0027】

図 1 を参照して本発明に係る音声認識装置 1 の機能構成を説明する。

音声取得部 21 は、会話の音声データを取得する。会話の音声データは、音声認識装置 1 がマイク 18 を備えるものであれば、マイク 18 から入力された音声データでもよいし、通信制御部 17 及びネットワーク 3 を介して音声認識装置 1 と通信接続された機器とから入力されたものでもよい。

【0028】

音声認識部 22 は、取得した音声データのユーザの発話の内容を音声認識し、音声認識結果として 1 または複数の変換候補を求める。音声認識部 22 は、発話の音声データと語とを対応付けた発話辞書や、音響モデル、言語モデル等の音声認識用データを有し、これらの音声認識用データを用いて、発話の音響や言語を解析し、発話の内容をテキストに変換する音声認識処理を行う。音声認識結果であるテキストは、状況判定部 23 の会話内容判定部 231 に出力される。

【0029】

状況判定部 23 は、会話の状況を判定する。会話の状況とは、具体的には、会話の内容（ジャンル）、及び会話が行われている場所や目的等の周囲状況（シーン）である。図 1 に示すように、状況判定部 23 は、会話内容を判定するための会話内容判定部 231 と、周囲状況を判定するための周囲状況判定部 232 とを有する。

【0030】

会話内容判定部 231 は、会話の文に含まれる単語を解析することにより、会話の内容を求める。ここで求める会話の内容とは、話のジャンルまたは目的等である。ジャンルとは、「観光」、「ゲーム」、「飲食」、「映画」、「学校」、「医療」、...等のように、何についての会話であるかを示す分類である。目的とは、「接客」や「雑談」等のように会話がどのような目的で行われているかを示す分類である。会話内容判定部 231 は、例えば、単語とジャンルとを関連付けたデータをジャンルデータ 5（図 4 参照）として記憶部 12 に予め記憶しており、このジャンルデータ 5 を参照することにより会話の内容（ジャンル）を判定する。各単語は複数のジャンルに跨って含まれていてもよい。会話内容判定部 231 は、音声認識結果として得られる 1 または複数の文に含まれる単語から、会話内容の候補を求める。また、判定対象とする文だけでなく、それより前に入力された音声データから認識された文（音声認識結果）を判定対象に含むようにすることが望ましい。これにより、コンテキスト（文脈）を考慮して会話の内容を求めることができる。

【0031】

周囲状況判定部 232 は、会話の音声データが入力されたときの周囲の音データの特徴に基づいて場所等の周囲状況を求める。音解析部 26 は、音声取得部 21 により取得した音声データから周囲の音データを抽出し、この周囲の音データの特徴を抽出し、周囲状況判定部 232 に出力する。周囲状況判定部 232 は、抽出した音データの特徴と状況とを関連付けたデータを音解析用データとして記憶部 12 に予め記憶しており、この音解析用データに基づいて会話の周囲状況を判定する。例えば、「レストラン」の音解析用データには、食器等の音やテーブルでの会話、接客の音等の特徴が含まれる。また「アミューズメントパーク」の音解析用データには、歓声やアトラクションの音等の特徴が含まれる。周囲状況判定部 232 は、音データの特徴から 1 または複数の周囲状況の候補を求める。

10

20

30

40

50

【 0 0 3 2 】

なお、周囲状況判定部 2 3 2 は、音解析のみならず、画像や位置情報に基づいて周囲状況を判定してもよい。具体的には、図 1 に示すように画像取得部 2 7 により会話中の様子や場所を撮影した画像（映像または静止画）等を解析する画像解析部 2 8 を備え、画像解析部 2 8 によって会話の場所や目的等、周囲状況を判定してもよい。また、GPS（Global Positioning System）等の位置情報取得部 2 9 を更に備え、周囲状況判定部 2 3 2 は、位置情報及び地図データ等に基づいて会話の場所（店舗や施設）等を求めることにより周囲状況を求めてもよい。

【 0 0 3 3 】

音声認識結果抽出部 2 4 は、音声認識部 2 2 により求めた複数の変換候補の中から状況判定部 2 3 により判定した会話の状況に適した変換候補を抽出する。変換候補の抽出については後述する。

【 0 0 3 4 】

出力部 2 5 は、音声認識結果抽出部 2 4 により抽出した変換候補（テキスト）を出力する。出力は、表示部 1 6 への表示や、制御部 1 1 への通知、ネットワーク 3 を介した通信接続先への送信等、当該音声認識装置 1 に接続された各種機器に対する制御情報としての送信等も含むものとする。

【 0 0 3 5 】

次に、図 3 を参照して、音声認識装置 1 が実行する音声認識処理について説明する。

制御部 1 1 は、記憶部 1 2 から図 3 に示す音声認識処理に関するプログラム及びデータを読み出し、このプログラム及びデータに基づいて処理を実行する。

【 0 0 3 6 】

まず制御部 1 1（音声取得部 2 1）は、会話の音声データを取得する（ステップ S 1 0 1）。音声データは、マイク 1 8 から入力されたものでもよいし、通信制御部 1 7 及びネットワーク 3 を介して音声認識装置 1 と通信接続された機器から入力されたものでもよい。制御部 1 1（音声認識部 2 2）は、取得した音声データについて音声認識を行う（ステップ S 1 0 2）。ステップ S 1 0 2 では、制御部 1 1（音声認識部 2 2）は、音声データに含まれる会話の音声を認識し、テキストに変換する処理を行う。制御部 1 1（音声認識部 2 2）は、音声認識処理の結果、1 または複数の変換候補を得る。複数の変換候補がある場合に、ステップ S 1 0 3 ~ ステップ S 1 0 4 の処理により会話の状況を判定する。

【 0 0 3 7 】

制御部 1 1（状況判定部 2 3 の会話内容判定部 2 3 1）は、会話の状況として、会話の内容（ジャンル等）を判定する（ステップ S 1 0 3）。制御部 1 1（会話内容判定部 2 3 1）は、ステップ S 1 0 2 の音声認識の結果（変換候補）に含まれる語の意味を解析することにより、会話の内容を求める。ここで求める会話の内容とは、会話のジャンルまたは目的等である。会話内容判定部 2 3 1 は、例えば、記憶部 1 2 に予め記憶されているジャンルデータ 5 を参照することにより会話の内容（ジャンル）を判定する。

【 0 0 3 8 】

ジャンルデータ 5 は、図 4 に示すように、単語の読み（音声認識結果）について 1 または複数の変換候補となる語と、その語のジャンルとを関連付けたデータである。例えば、音声認識結果「たいけん」の変換候補は、「大剣」と「体験」等があり、変換候補「大剣」のジャンルは「RPG（ゲーム）」、変換候補「体験」のジャンルは「観光」である。このように、ひとつの音声認識結果について 1 または複数の変換候補と各変換候補に応じたジャンルが格納されている。各変換候補について複数のジャンルが関連づけられていてもよい。会話内容判定部 2 3 1 は、ジャンルデータ 5 を参照することにより、音声認識結果について、1 または複数の会話内容の候補（ジャンル候補）を求める。例えば、音声認識結果「おすすめのたいけんはありますか」であれば、「たいけん」という語が含まれるため、ジャンル候補として、「RPG（ゲーム）」と「観光」が求められる。

【 0 0 3 9 】

次に、制御部 1 1（状況判定部 2 3 の周囲状況判定部 2 3 2）は、周囲の状況を判定す

10

20

30

40

50

る（ステップS104）。制御部11（周囲状況判定部232）は、会話の音声データが入力されたときの周囲の音の特徴を解析し、音の特徴に基づいて場所等の周囲状況を求める。例えば、「レストラン」で収録された音には、食器等の音や接客の音等の特徴が含まれている。周囲状況判定部232は、音の特徴と状況とを関連付けたデータを音特徴データとして記憶部12に予め記憶しており（不図示）、この音特徴データに基づいて会話の周囲状況を判定するようにしてもよい。制御部11（周囲状況判定部232）は、1または複数の周囲状況の候補を求める。なお、周囲状況は、場所に限定されず、「接客」、「授業」、「雑談」等のように、会話の目的等としてもよい。制御部11は、音の特徴解析による周囲状況の判定結果として、例えば、「観光案内所」、「接客」等を得る。

【0040】

なお、制御部11（周囲状況判定部232）は、周囲の音の特徴のみならず、画像や位置情報に基づいて周囲状況を判定してもよい。具体的には、画像取得部27（カメラ等）により会話中の様子を撮影した映像（画像）等を取得し、解析する画像解析部28を備え、画像解析部28によって会話の音声データが入力されたときの会話の場所や目的等、周囲状況を求めてもよい。また、GPS等の位置情報を取得し、位置情報及び予め記憶されている地図データに基づいて会話の場所（店舗や施設）等を求めることにより周囲状況を求めてもよい。

【0041】

制御部11（音声認識結果抽出部24）は、ステップS102で得た音声認識結果の複数の変換候補のうち、ステップS103及びステップS104において求めた会話の状況（会話内容（ジャンル）及び周囲状況）に適した変換候補を抽出する（ステップS105）。例えば、音声認識結果が「おすすめのたいけんはありますか」の場合、この文に含まれる「たいけん」の語には、「大剣」と「体験」の変換候補がある。ステップS103で会話の内容が「RPG（ゲーム）」、「観光」、...と判定され、ステップS104で周囲の状況が「観光」、「接客」、...と判定された場合、制御部11は会話内容のジャンルと周囲状況とをマッチングし、尤度の高いジャンルの語を抽出する。上記例では、会話の状況として「観光」が尤もらしいと判定されるため、変換候補「体験」を選択し、入力音声の音声認識結果として「おすすめの体験はありますか」を得る。

【0042】

制御部11は、ステップS103及びステップS104で判定した状況（上記例では、「観光」）を状況履歴データとして時間情報（音声データの入力時刻等）と関連付けて記憶部12に保存する（ステップS106）。

【0043】

制御部11（出力部25）は、ステップS105で抽出した音声認識結果を出力する。出力は、表示部16への表示や、制御部11への通知、ネットワーク3を介した通信接続先への送信等、当該音声認識装置1に接続された各種機器への制御信号の送信等も含むものとする。

【0044】

音声認識結果を出力すると、入力された音声データに対する音声認識処理を終了する。

【0045】

なお、上述の音声認識処理において、ステップS106で保存した履歴に基づき、会話内容を判定するようにしてもよい。すなわち、ステップS103において、前の文までの会話の状況の履歴が保存されている場合は、制御部11は、前の文までの会話の状況から会話の内容（ジャンル）を絞り込んでよい。

【0046】

例えば、図5に示すように、17時50分に「いらっしゃいませ」、「何かお探しですか」という会話の音声が入力され、17時50分における音声（「いらっしゃいませ」、「何かお探しですか」）の状況として、「店」、「ファミレス」、「ドラッグストア」等の会話ジャンルが求められるものとする。会話ジャンルの各候補にはそれぞれ尤度が付与されているものとする。例えば、語「いらっしゃいませ」に対する会話ジャンル「店」の

10

20

30

40

50

尤度は「1.0」であり、「ファミレス」の尤度は「0.9」であり、「ドラッグストア」の尤度は「0.9」、...等である。尤度は、例えば語とジャンルとを対応付けたジャンルデータ5に予め付与されているものとする。また、17時50分における音解析による状況判定結果として、「観光案内所」と「受付」が尤度とともに求められる。例えば、「観光案内所」の尤度は「0.8」、「受付」の尤度は「0.5」のように求められるものとする。この音解析による状況判定結果の尤度は、入力された音と予め記憶されている音特徴データとの一致度等から付与するものとするればよい。

【0047】

次に、17時56分に「 レジャー施設はどこですか」、「ここから5分の場所にあります」という会話の音声が入力される。制御部11は、17時56分における会話のジャンルを「遊園地」、「観光」等と判定する。それぞれの尤度は「遊園地」が「0.9」、「観光」が「0.8」とする。また17時56分における音解析による状況判定結果として、「観光案内所」、「店頭」、「接客」を得る。「観光案内所」の尤度は「0.8」、「店頭」の尤度は「0.8」、「接客」の尤度は「0.7」とする。

10

【0048】

その後、18時00分に、処理対象である「おすすめのたいけんはありますか」という音声が入力されるものとする。制御部11は、18時00分における会話のジャンルを「RPG(ゲーム)」、「観光」等と判定する。それぞれの尤度として「RPG(ゲーム)」は「0.5」、「観光」は「0.2」を得るものとする。また18時00分における音解析による状況判定結果として、「観光案内所」、「接客」を得る。「観光案内所」の尤度は「0.8」、「接客」の尤度は「0.7」とする。

20

【0049】

「たいけん」の各変換候補(「大剣」、「体験」)の会話内容に基づく尤度は、RPG(ゲーム)は「0.5」、観光は「0.2」であるが、音による状況判定では、RPG(ゲーム)という候補はなし(尤度「0」)、観光(観光案内所)は尤度「0.8」である。これらを併せると、

「RPG(ゲーム)」の尤度 = 会話内容「0.5」 + 音判定「0」 = 0.5

「観光」の尤度 = 会話内容「0.2」 + 音判定「0.8」 = 1.0

となる。

【0050】

したがって、状況としては「RPG(ゲーム)」よりも「観光」の尤度が高く適切である。従って、制御部11は、「たいけん」の変換結果(音声認識結果)としては「体験」が適していると判断する。このように、音声認識や音判定の履歴を遡って状況判定に利用すれば、コンテキスト(文脈、会話の流れ)を考慮した音声認識結果を得ることが可能となる。例えば最近の「RPG(ゲーム)」、「観光」のコンテキストの履歴を見ると「観光」が多く出現するので、この値(=重み)を音声認識や音判定の尤度に掛け合わせて足したものを比較し判定することもできる。

30

「RPG(ゲーム)」の尤度 = (会話内容「0.5」 + 音判定「0」) × ゲームのコンテキストの重み「0」 = 0

「観光」の尤度 = (会話内容「0.2」 + 音判定「0.8」) × 観光のコンテキストの重み「0.7」 = 0.7

40

【0051】

以上説明したように、本実施の形態の音声認識装置1は、会話の音声データを認識し、音声認識結果として複数の変換候補を求め、音声認識の結果のみならず周囲音等を考慮して会話の状況を判定することにより、複数の変換候補の中から会話の状況に適した候補を抽出する。これにより、会話の状況に適した音声認識結果を得ることが可能となり、音声認識精度を向上できる。

【0052】

なお、本発明の音声認識装置1をスマートフォン2やタブレット等の通信機器や、インターネット等のネットワーク3に接続されたPC(Personal Computer)7に適用する場

50

合において、図 6 に示す音声認識システム 100 のように、スマートフォン 2 等からアクセス可能なサーバに本発明に係る音声認識装置 1 の各機能部（音声取得部 21、音声認識部 22、状況判定部 23、音声認識結果抽出部 24、出力部 25、音解析部 26 等）を備える構成としてもよい。すなわち、スマートフォン 2、PC 7 等はマイク 18 から入力された会話の音声データをネットワーク 3 を介して音声認識装置 1（サーバ）に送信すると、音声認識装置 1 は、図 3 に示す音声認識処理を実行し、音声認識結果を音声入力元のスマートフォン 2 等に返すものとしてもよい。

【0053】

また、本発明に係る音声認識装置 1 は、対話型ロボット 6 に適用してもよい。この場合、本発明に係る音声認識装置 1 の各機能（音声取得部 21（マイク）、音声認識部 22、状況判定部 23、音声認識結果抽出部 24、出力部 25、音解析部 26、画像取得部 27（カメラ）、画像解析部 28、位置情報取得部 29）をロボット 6 が備える構成とする。或いは、対話型ロボット 6 の音声取得部 21（マイク）、画像取得部 27（カメラ）から入力された音声や画像を、サーバ（音声認識装置 1）に送信し、サーバ（音声認識装置 1）は、図 3 に示す音声認識処理を実行し、音声認識結果を音声入力元の対話型ロボット 6 に返すものとしてもよい。

10

【0054】

その他、本発明に係る音声認識装置 1 は、カーナビゲーションシステム等の各種情報機器や家電等に適用することも可能である。

【0055】

以上、添付図面を参照して、本発明に係る音声認識装置等の好適な実施形態について説明したが、本発明は係る例に限定されない。当業者であれば、本願で開示した技術的思想の範疇内において、各種の変更例または修正例に想到し得ることは明らかであり、それらについても当然に本発明の技術的範囲に属するものと了解される。

20

【符号の説明】

【0056】

- 1 …… 音声認識装置
- 10 …… コンピュータ
- 11 …… 制御部
- 12 …… 記憶部
- 13 …… メディア入出力部
- 14 …… 周辺機器 I / F 部
- 15 …… 入力部
- 16 …… 表示部
- 17 …… 通信制御部
- 18 …… マイク
- 19 …… バス
- 21 …… 音声取得部
- 22 …… 音声認識部
- 23 …… 状況判定部
- 23 1 …… 会話内容判定部
- 23 2 …… 周囲状況判定部
- 24 …… 音声認識結果抽出部
- 25 …… 出力部
- 26 …… 音解析部
- 27 …… 画像取得部
- 28 …… 画像解析部
- 29 …… 位置情報取得部
- 3 …… ネットワーク
- 5 …… ジャンルデータ

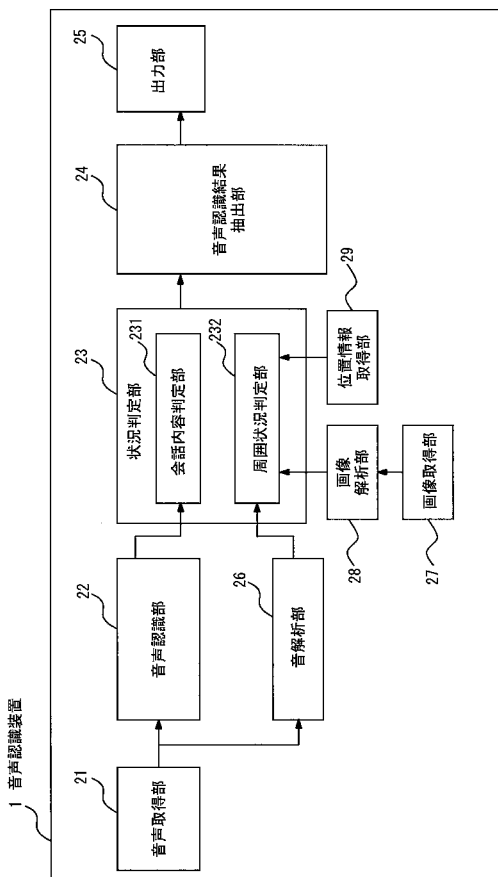
30

40

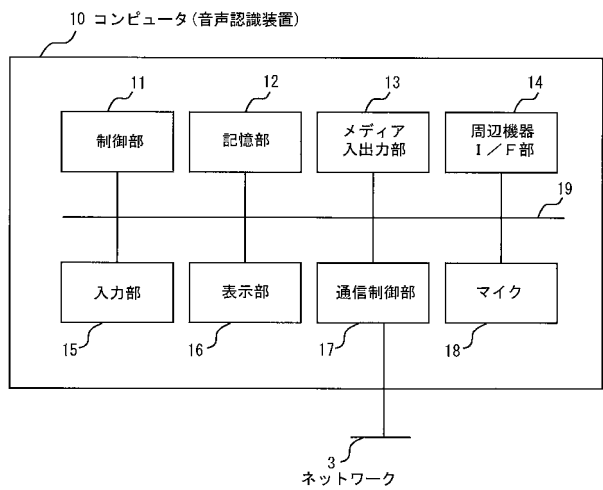
50

- 6 対話型ロボット
- 7 P C
- 1 0 0 音声認識システム

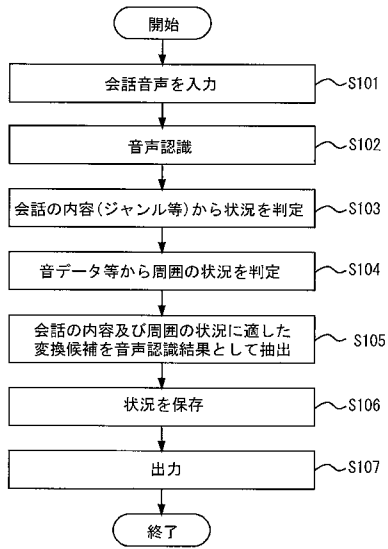
【 図 1 】



【 図 2 】



【 図 3 】

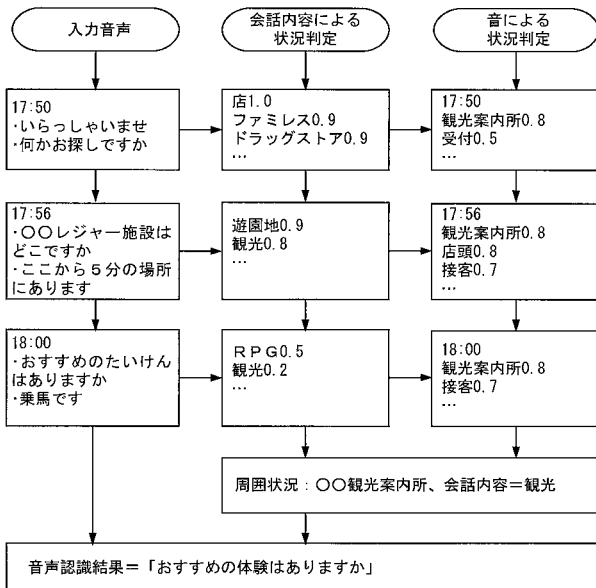


【 図 4 】

5 ジャンルデータ

音声認識結果	変換候補	ジャンル
たいけん	大剣	RPG (ゲーム)、...
	体験	観光、...
	...	
...		

【 図 5 】



【 図 6 】

