

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7691027号
(P7691027)

(45)発行日 令和7年6月11日(2025.6.11)

(24)登録日 令和7年6月3日(2025.6.3)

(51)国際特許分類 F I
 G 1 0 L 21/007 (2013.01) G 1 0 L 21/007
 G 1 0 L 13/02 (2013.01) G 1 0 L 13/02 1 1 0 Z
 G 1 0 L 15/20 (2006.01) G 1 0 L 15/20 3 0 0

請求項の数 10 (全28頁)

(21)出願番号	特願2024-504041(P2024-504041)	(73)特許権者	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(86)(22)出願日	令和4年3月1日(2022.3.1)	(74)代理人	100104765 弁理士 江上 達夫
(86)国際出願番号	PCT/JP2022/008597	(74)代理人	100107331 弁理士 中村 聡延
(87)国際公開番号	WO2023/166557	(74)代理人	100131015 弁理士 三輪 浩誉
(87)国際公開日	令和5年9月7日(2023.9.7)	(72)発明者	カク レイ 東京都港区芝五丁目7番1号 日本電気株式会社内
審査請求日	令和6年8月6日(2024.8.6)	(72)発明者	山本 仁 東京都港区芝五丁目7番1号 日本電気株式会社内

最終頁に続く

(54)【発明の名称】 音声認識システム、音声認識方法、及び記録媒体

(57)【特許請求の範囲】

【請求項1】

話者が発話したリアル発話データを取得する発話データ取得手段と、
 前記リアル発話データをテキストデータに変換するテキスト変換手段と、
 前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成する音声合成手段と、
 前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成する変換モデル生成手段と、
 前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、
 を備える音声認識システム。

10

【請求項2】

前記変換モデル生成手段は、前記入力音声と、前記音声認識手段の認識結果と、を用いて前記変換モデルのパラメータを調整する、
 請求項1に記載の音声認識システム。

【請求項3】

前記対応合成音声を含むデータを用いて音声認識モデルを生成する音声認識モデル生成手段を更に備え、
 前記音声認識手段は、前記音声認識モデルを用いて音声認識する、
 請求項1又は2に記載の音声認識システム。

【請求項4】

20

前記音声認識モデル生成手段は、前記変換モデルを用いて変換された前記合成音声と、前記音声認識手段の認識結果と、を用いて前記音声認識モデルのパラメータを調整する、請求項 3 に記載の音声認識システム。

【請求項 5】

前記話者の属性を示す属性情報を取得する属性取得手段を更に備え、前記音声合成手段は、前記属性情報を用いて音声合成を行うことで前記対応合成音声を生成する、

請求項 1 から 4 のいずれか一項に記載の音声認識システム。

【請求項 6】

所定の条件ごとに前記リアル発話データを記憶する複数のリアル発話音声コーパスを更に備え、

前記発話データ取得手段は、前記複数のリアル発話音声コーパスから 1 つを選択して前記リアル発話データを取得する、

請求項 1 から 5 のいずれか一項に記載の音声認識システム。

【請求項 7】

前記テキストデータ及び前記対応合成音声の少なくとも一方にノイズを付与するノイズ付与手段を更に備える、

請求項 1 から 6 のいずれか一項に記載の音声認識システム。

【請求項 8】

手話データを取得する手話データ取得手段と、前記手話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記手話データに対応する対応合成音声を生成する音声合成手段と、

前記手話データ及び前記対応合成音声を用いて、入力される手話を合成音声に変換する変換モデルを生成する変換モデル生成手段と、

前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える音声認識システム。

【請求項 9】

少なくとも 1 つのコンピュータによって、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成し、

前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成し、

前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法。

【請求項 10】

少なくとも 1 つのコンピュータに、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成し、

前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成し、

前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法を実行させる コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

この開示は、音声認識システム、音声認識方法、及び記録媒体の技術分野に関する。

【背景技術】

【0002】

この種のシステムとして、合成音声を生成するものが知られている。例えば特許文献1では、音声の声色を表す特徴量を学習済みの変換モデルによって変換するなどして、合成音声を生成することが開示されている。特許文献2では、音声認識結果として取得されたテキストデータからターゲット言語の文を生成し、そのターゲット言語の文から合成音声を生成することが開示されている。

【0003】

その他の関連する技術として、例えば特許文献3では、学習用コーパスを用いて音声変換モデルの学習を行うことが開示されている。

10

【先行技術文献】

【特許文献】

【0004】

【文献】国際公開第2021/033685号

【文献】国際公開第2014/010450号

【文献】特開2020-166224号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

この開示は、先行技術文献に開示された技術を改善することを目的とする。

20

【課題を解決するための手段】

【0006】

この開示の音声認識システムの一の様子は、話者が発話したリアル発話データを取得する発話データ取得手段と、前記リアル発話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成する音声合成手段と、前記リアル発話データ及び前記対応合成音声をを用いて、入力音声を合成音声に変換する変換モデルを生成する変換モデル生成手段と、前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える。

30

【0007】

この開示の音声認識システムの一の様子は、手話データを取得する手話データ取得手段と、前記手話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記手話データに対応する対応合成音声を生成する音声合成手段と、前記手話データ及び前記対応合成音声をを用いて、入力される手話を合成音声に変換する変換モデルを生成する変換モデル生成手段と、前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える

【0008】

この開示の音声認識方法の一の様子は、少なくとも1つのコンピュータによって、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成し、前記リアル発話データ及び前記対応合成音声をを用いて、入力音声を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する。

40

【0009】

この開示の記録媒体の一の様子は、少なくとも1つのコンピュータに、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成し、前記リアル発話データ及び前記対応合成音声をを用いて、入力音声を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識す

50

る、音声認識方法を実行させるコンピュータプログラムが記録されている。

【図面の簡単な説明】

【0010】

【図1】第1実施形態に係る音声認識システムのハードウェア構成を示すブロック図である。

【図2】第1実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【図3】第1実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【図4】第1実施形態に係る音声認識システムによる音声認識動作の流れを示すフローチャートである。

10

【図5】第2実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【図6】第2実施形態に係る音声認識システムによる変換モデル学習動作の流れを示すフローチャートである。

【図7】第3実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【図8】第3実施形態に係る音声認識システムによる音声認識モデル生成動作の流れを示すフローチャートである。

【図9】第4実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【図10】第4実施形態に係る音声認識システムによる音声認識モデル学習動作の流れを示すフローチャートである。

【図11】第5実施形態に係る音声認識システムの機能的構成を示すブロック図である。

20

【図12】第5実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【図13】第6実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【図14】第6実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【図15】第7実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【図16】第7実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【図17】第7実施形態の変形例に係る音声認識システムの機能的構成を示すブロック図である。

30

【図18】第7実施形態の変形例に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【図19】第8実施形態の変形例に係る音声認識システムの機能的構成を示すブロック図である。

【図20】第8実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【図21】第8実施形態に係る音声認識システムによる音声認識動作の流れを示すフローチャートである。

【発明を実施するための形態】

【0011】

40

以下、図面を参照しながら、音声認識システム、音声認識方法、及び記録媒体の実施形態について説明する。

【0012】

<第1実施形態>

第1実施形態に係る音声認識システムについて、図1から図4を参照して説明する。

【0013】

(ハードウェア構成)

まず、図1を参照しながら、第1実施形態に係る音声認識システムのハードウェア構成について説明する。図1は、第1実施形態に係る音声認識システムのハードウェア構成を示すブロック図である。

50

【0014】

図1に示すように、第1実施形態に係る音声認識システム10は、プロセッサ11と、RAM(Random Access Memory)12と、ROM(Read Only Memory)13と、記憶装置14とを備えている。音声認識システム10は更に、入力装置15と、出力装置16と、を備えていてもよい。上述したプロセッサ11と、RAM12と、ROM13と、記憶装置14と、入力装置15と、出力装置16とは、データバス17を介して接続されている。

【0015】

プロセッサ11は、コンピュータプログラムを読み込む。例えば、プロセッサ11は、RAM12、ROM13及び記憶装置14のうちの少なくとも一つが記憶しているコンピュータプログラムを読み込むように構成されている。或いは、プロセッサ11は、コンピュータで読み取り可能な記録媒体が記憶しているコンピュータプログラムを、図示しない記録媒体読み取り装置を用いて読み込んでよい。プロセッサ11は、ネットワークインタフェースを介して、音声認識システム10の外部に配置される不図示の装置からコンピュータプログラムを取得してもよい(つまり、読み込んでよい)。プロセッサ11は、読み込んだコンピュータプログラムを実行することで、RAM12、記憶装置14、入力装置15及び出力装置16を制御する。本実施形態では特に、プロセッサ11が読み込んだコンピュータプログラムを実行すると、プロセッサ11内には、音声認識を行うための機能ブロックが実現される。即ち、プロセッサ11は、音声認識システム10における各制御を実行するコントローラとして機能してよい。

【0016】

プロセッサ11は、例えばCPU(Central Processing Unit)、GPU(Graphics Processing Unit)、FPGA(field-programmable gate array)、DSP(Demand-Side Platform)、ASIC(Application Specific Integrated Circuit)として構成されてよい。プロセッサ11は、これらのうち一つで構成されてもよいし、複数を並列で用いるように構成されてもよい。

【0017】

RAM12は、プロセッサ11が実行するコンピュータプログラムを一時的に記憶する。RAM12は、プロセッサ11がコンピュータプログラムを実行している際にプロセッサ11が一時的に使用するデータを一時的に記憶する。RAM12は、例えば、D-RAM(Dynamic Random Access Memory)や、SRAM(Static Random Access Memory)であってよい。また、RAM12に代えて、他の種類の揮発性メモリが用いられてもよい。

【0018】

ROM13は、プロセッサ11が実行するコンピュータプログラムを記憶する。ROM13は、その他に固定的なデータを記憶していてもよい。ROM13は、例えば、P-ROM(Programmable Read Only Memory)や、EPROM(Erasable Read Only Memory)であってよい。また、ROM13に代えて、他の種類の不揮発性メモリが用いられてもよい。

【0019】

記憶装置14は、音声認識システム10が長期的に保存するデータを記憶する。記憶装置14は、プロセッサ11の一時記憶装置として動作してもよい。記憶装置14は、例えば、ハードディスク装置、光磁気ディスク装置、SSD(Solid State Drive)及びディスクアレイ装置のうちの少なくとも一つを含んでいてもよい。

【0020】

入力装置15は、音声認識システム10のユーザからの入力指示を受け取る装置である。入力装置15は、例えば、キーボード、マウス及びタッチパネルのうちの少なくとも一つを含んでいてもよい。入力装置15は、スマートフォンやタブレット等の携帯端末として構成されていてもよい。入力装置15は、例えばマイクを含む音声入力可能な装置で

10

20

30

40

50

あってもよい。

【0021】

出力装置16は、音声認識システム10に関する情報を外部に対して出力する装置である。例えば、出力装置16は、音声認識システム10に関する情報を表示可能な表示装置（例えば、ディスプレイ）であってもよい。また、出力装置16は、音声認識システム10に関する情報を音声出力可能なスピーカ等であってもよい。出力装置16は、スマートフォンやタブレット等の携帯端末として構成されていてもよい。また、出力装置16は、画像以外の形式で情報を出力する装置であってもよい。例えば、出力装置16は、音声認識システム10に関する情報を音声で出力するスピーカであってもよい。

【0022】

なお、図1では、複数の装置を含んで構成される音声認識システム10の例を挙げたが、これらの全部又は一部の機能を、1つの装置（音声認識装置）として実現してもよい。その場合、音声認識装置は、例えば上述したプロセッサ11、RAM12、ROM13のみを備えて構成され、その他の構成要素（即ち、記憶装置14、入力装置15、出力装置16）については、音声認識装置に接続される外部の装置が備えるようにしてもよい。また、音声認識装置は、一部の演算機能を外部の装置（例えば、外部サーバやクラウド等）によって実現するものであってもよい。

【0023】

（機能的構成）

次に、図2を参照しながら、第1実施形態に係る音声認識システム10の機能的構成について説明する。図2は、第1実施形態に係る音声認識システムの機能的構成を示すブロック図である。

【0024】

図2に示すように、第1実施形態に係る音声認識システム10は、その機能を実現するための構成要素として、発話データ取得部110と、テキスト変換部120と、音声合成部130と、変換モデル生成部140と、音声変換部210と、音声認識部220と、を備えて構成されている。発話データ取得部110、テキスト変換部120、音声合成部130、変換モデル生成部140、音声変換部210、音声認識部220の各々は、例えば上述したプロセッサ11（図1参照）によって実現される処理ブロックであってよい。

【0025】

発話データ取得部110は、話者が発話したリアル発話データを取得可能に構成されている。リアル発話データは、音声データ（例えば、波形データ）であってもよい。リアル発話データは、例えば複数のリアル発話データを蓄積するデータベース（リアル発話音声コーパス）から取得されてよい。発話データ取得部110で取得されたリアル発話データは、テキスト変換部120及び変換モデル生成部140に出力される構成となっている。

【0026】

テキスト変換部120は、発話データ取得部110で取得されたリアル発話データをテキストデータに変換可能に構成されている。即ち、テキスト変換部120は、音声データをテキスト変換する処理を実行可能に構成されている。なお、テキスト変換の具体的な手法については、既存の技術が適宜採用されてよい。テキスト変換部120で変換されたテキストデータ（即ち、リアル発話データに対応するテキストデータ）は、音声合成部130に出力される構成となっている。

【0027】

音声合成部130は、テキスト変換部120で変換されたテキストデータを音声合成することで、リアル発話データに対応する対応合成音声を生成可能に構成されている。なお、音声合成の具体的な手法については、既存の技術を適宜採用することができる。音声合成部130で生成された対応合成音声は、変換モデル生成部140に出力される構成となっている。なお、対応合成音声は、複数の対応合成を蓄積可能なデータベース（合成音声コーパス）に蓄積されてから、変換モデル生成部140に出力されてもよい。

【0028】

10

20

30

40

50

変換モデル生成部 140 は、発話データ取得部 110 で取得されたリアル発話データと、音声合成部 130 で合成された対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成可能に構成されている。変換モデルは、例えば、話者が発話した入力音声（即ち、人間の音声）を、合成音声（即ち、機械的な音声）に近づくように変換する。変換モデル生成部 140 は、例えば GAN (Generative Adversarial Network: 敵対的生成ネットワーク) を用いて、変換モデルを生成するように構成されてよい。変換モデル生成部 140 で生成された変換モデルは、音声変換部 210 に出力される構成となっている。

【0029】

音声変換部 210 は、変換モデル生成部 140 で生成された変換モデルを用いて、入力音声を合成音声に変換可能に構成されている。音声変換部 210 に入力される入力音声は、例えばマイク等を用いて入力される音声であってよい。音声変換部 210 で変換された合成音声は、音声認識部 220 に出力される構成となっている。

【0030】

音声認識部 220 は、音声変換部 210 で変換された合成音声を音声認識することが可能に構成されている。即ち、音声認識部 220 は、合成音声をテキスト化する処理を実行可能に構成されている。音声認識部 220 は、合成音声の音声認識結果を出力可能に構成されてよい。なお、音声認識結果の利用方法については特に限定されない。

【0031】

(変換モデル生成動作)

次に、図 3 を参照しながら、第 1 実施形態に係る音声認識システム 10 による変換モデルを生成する際の動作（以下、適宜「変換モデル生成動作」と称する）の流れについて説明する。図 3 は、第 1 実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

【0032】

図 3 に示すように、第 1 実施形態に係る音声認識システム 10 による変換モデル生成動作が開始されると、まず発話データ取得部 110 が、リアル発話データを取得する（ステップ S101）。そして、テキスト変換部 120 が、発話データ取得部 110 で取得されたリアル発話データをテキストデータに変換する（ステップ S102）。

【0033】

続いて、音声合成部 130 が、テキスト変換部 120 で変換されたテキストデータを音声合成し、リアル発話データに対応する対応合成音声を生成する（ステップ S103）。そして、変換モデル生成部 140 が、発話データ取得部 110 で取得されたリアル発話データ及び音声合成部 130 で生成された対応合成音声に基づいて、変換モデルを生成する（ステップ S104）。その後、変換モデル生成部 140 は、生成した変換モデルを音声変換部 210 に出力する（ステップ S105）。

【0034】

(変換認識動作)

次に、図 4 を参照しながら、第 1 実施形態に係る音声認識システム 10 による音声認識を行う際の動作（以下、適宜「音声認識動作」と称する）の流れについて説明する。図 3 は、第 1 実施形態に係る音声認識システムによる音声認識動作の流れを示すフローチャートである。

【0035】

図 4 に示すように、第 1 実施形態に係る音声認識システム 10 による音声認識動作が開始されると、まず音声変換部 210 が入力音声を取得する（ステップ S151）。そして、音声変換部 210 は、変換モデル生成部 140 で生成された変換モデルを読み込む（ステップ S152）。その後、音声変換部 210 は、読み込んだ変換モデルを用いて音声変換を行い、入力音声を合成音声に変換する（ステップ S153）。

【0036】

続いて、音声認識部 220 は、音声認識モデル（即ち、音声認識をするためのモデル）

10

20

30

40

50

を読み込む（ステップS154）。そして、音声認識部220は、読み込んだ音声認識モデルを用いて、音声変換部210で合成された合成音声を音声認識する（ステップS155）。その後、音声認識部220は、音声認識結果を出力する（ステップS156）。

【0037】

（技術的効果）

次に、第1実施形態に係る音声認識システム10によって得られる技術的効果について説明する。

【0038】

図1から図4で説明したように、第1実施形態に係る音声認識システム10では、変換モデルを生成する際に、リアル発話データ及びリアル発話データに対応する対応合成音声
10
が用いられる。そして特に、リアル発話データに対応する対応合成音声は、リアル発話データをテキスト変換し、テキストデータを音声合成することで生成される。このようにすれば、リアル発話データと、それに対応する合成音声と、の両方を用意する必要がなくなる（即ち、リアル発話データのみ用意すれば、対応合成音声を生成できる）ため、変換モデルを生成するのに要するコストを抑制することができる。その結果、低コストで認識精度の高い音声認識を実現することが可能となる。

【0039】

<第2実施形態>

第2実施形態に係る音声認識システム10について、図5及び図6を参照して説明する。なお、第2実施形態は、上述した第1実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第1実施形態と同一であってよい。このため、以下では、すでに説明した第1実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。
20

【0040】

（機能的構成）

まず、図5を参照しながら、第2実施形態に係る音声認識システム10の機能的構成について説明する。図5は、第2実施形態に係る音声認識システムの機能的構成を示すブロック図である。なお、図5では、図2で示した構成要素と同様の要素に同一の符号を付している。

【0041】

図5に示すように、第2実施形態に係る音声認識システム10は、その機能を実現するための構成要素として、発話データ取得部110と、テキスト変換部120と、音声合成部130と、変換モデル生成部140と、音声変換部210と、音声認識部220と、を備えて構成されている。そして第2実施形態では特に、変換モデル生成部140に、音声変換部210に入力される入力音声及び音声認識部220による認識結果が入力される構成となっている。第2実施形態に係る変換モデル生成部140は、音声変換部210に入力される入力音声及び音声認識部220による認識結果に基づいて、変換モデルの学習を実行可能に構成されている。
30

【0042】

（変換モデル学習動作）

次に、図6を参照しながら、第2実施形態に係る音声認識システム10による変換モデルを学習する際の動作（以下、適宜「変換モデル学習動作」と称する）の流れについて説明する。図6は、第2実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。
40

【0043】

図6に示すように、第2実施形態に係る音声認識システム10による変換モデル学習動作が開始されると、まず変換モデル生成部140が、音声変換部210に入力される入力音声を取得する（ステップS201）。そして、変換モデル生成部140は更に、その入力音声が入力された際の音声認識結果（即ち、図4に示すステップS156で出力される音声認識結果）を取得する（ステップS202）。
50

【 0 0 4 4 】

続いて、変換モデル生成部 1 4 0 は、取得した入力音声及び音声認識結果に基づいて、変換モデルを学習する（ステップ S 2 0 3）。この際、変換モデル生成部 1 4 0 は、すでに生成していた変換モデルのパラメータ調整を行ってよい。その後、変換モデル生成部 1 4 0 は、学習した変換モデルを音声変換部 2 1 0 に出力する（ステップ S 2 0 4）。

【 0 0 4 5 】

（技術的効果）

次に、第 2 実施形態に係る音声認識システム 1 0 によって得られる技術的効果について説明する。

【 0 0 4 6 】

図 5 及び図 6 で説明したように、第 2 実施形態に係る音声認識システム 1 0 では、入力音声及び音声認識結果に基づいて変換モデルが学習される。このようにすれば、入力音声を実際にどのように音声認識されるかを考慮して学習が行われるため、より適切な音声変換が行えるように変換モデルを学習できる。具体的には、音声変換した合成音声を用いて行う音声認識の精度が向上するように、変換モデルを学習できる。

【 0 0 4 7 】

< 第 3 実施形態 >

第 3 実施形態に係る音声認識システム 1 0 について、図 7 及び図 8 を参照して説明する。なお、第 3 実施形態は、上述した第 1 及び第 2 実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第 1 及び第 2 実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

【 0 0 4 8 】

（機能的構成）

まず、図 7 を参照しながら、第 3 実施形態に係る音声認識システム 1 0 の機能的構成について説明する。図 7 は、第 3 実施形態に係る音声認識システムの機能的構成を示すブロック図である。なお、図 7 では、図 2 で示した構成要素と同様の要素に同一の符号を付している。

【 0 0 4 9 】

図 7 に示すように、第 3 実施形態に係る音声認識システム 1 0 は、その機能を実現するための構成要素として、発話データ取得部 1 1 0 と、テキスト変換部 1 2 0 と、音声合成部 1 3 0 と、変換モデル生成部 1 4 0 と、音声変換部 2 1 0 と、音声認識部 2 2 0 と、音声認識モデル生成部 3 1 0 と、を備えて構成されている。即ち、第 3 実施形態に係る音声認識システム 1 0 は、第 1 実施形態の構成（図 2 参照）に加えて、音声認識モデル生成部 3 1 0 を更に備えている。なお、音声認識モデル生成部 3 1 0 は、例えば上述したプロセッサ 1 1（図 1 参照）によって実現される処理ブロックであってよい。

【 0 0 5 0 】

音声認識モデル生成部 3 1 0 は、入力音声を合成音声に変換する音声認識モデルを生成可能に構成されている。具体的には、音声認識モデル生成部 3 1 0 は、音声合成手段で生成された対応合成音声を用いて、音声認識モデルを生成可能に構成されている。なお、音声認識モデルは、対応合成音声と、それ以外の合成音声とを用いて、音声認識モデルを生成してもよい。音声認識モデル生成部 3 1 0 は、音声合成部 1 3 0 から直接対応合成音声を取得するよう構成されてもよいし、音声合成手段で生成された対応合成音声を複数記憶する合成音声コーパスから対応合成音声を取得するように構成されてもよい。音声認識モデル生成部 3 1 0 で生成された音声認識モデルは、音声認識部 2 2 0 に出力される構成となっている。

【 0 0 5 1 】

（音声認識モデル生成動作）

次に、図 8 を参照しながら、第 3 実施形態に係る音声認識システム 1 0 による音声認識モデルを生成する際の動作（以下、適宜「音声認識モデル生成動作」と称する）の流れに

10

20

30

40

50

ついて説明する。図 8 は、第 3 実施形態に係る音声認識システムによる音声認識モデル生成動作の流れを示すフローチャートである。

【 0 0 5 2 】

図 8 に示すように、第 3 実施形態に係る音声認識システム 1 0 による音声認識モデル生成動作が開始されると、まず音声認識モデル生成部 3 1 0 が、音声合成部 1 3 0 で生成された対応合成音声を取得する（ステップ S 3 0 1 ）。

【 0 0 5 3 】

続いて、音声認識モデル生成部 3 1 0 は、取得した対応合成音声を用いて音声認識モデルを生成する（ステップ S 3 0 2 ）。その後、音声認識モデル生成部 3 1 0 は、生成した音声認識モデルを音声認識部 2 2 0 に出力する（ステップ S 3 0 3 ）。

10

【 0 0 5 4 】

（技術的効果）

次に、第 3 実施形態に係る音声認識システム 1 0 によって得られる技術的効果について説明する。

【 0 0 5 5 】

図 7 及び図 8 で説明したように、第 3 実施形態に係る音声認識システム 1 0 では、対応合成音声を用いて音声認識モデルが生成される。このようにすれば、音声認識モデルを生成するための合成音声を別途用意する必要がない（即ち、音声変換モデルを生成するために用いた対応合成音声を利用できる）ため、効率的に音声認識モデルを生成することが可能である。

20

【 0 0 5 6 】

< 第 4 実施形態 >

第 4 実施形態に係る音声認識システム 1 0 について、図 9 及び図 1 0 を参照して説明する。なお、第 4 実施形態は、上述した第 3 実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第 1 から第 3 実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

【 0 0 5 7 】

（機能的構成）

まず、図 9 を参照しながら、第 4 実施形態に係る音声認識システム 1 0 の機能的構成について説明する。図 9 は、第 4 実施形態に係る音声認識システムの機能的構成を示すブロック図である。なお、図 9 では、図 7 で示した構成要素と同様の要素に同一の符号を付している。

30

【 0 0 5 8 】

図 9 に示すように、第 4 実施形態に係る音声認識システム 1 0 は、その機能を実現するための構成要素として、発話データ取得部 1 1 0 と、テキスト変換部 1 2 0 と、音声合成部 1 3 0 と、変換モデル生成部 1 4 0 と、音声変換部 2 1 0 と、音声認識部 2 2 0 と、音声認識モデル生成部 3 1 0 と、を備えて構成されている。そして第 4 実施形態では特に、音声認識モデル生成部 3 1 0 に、音声変換部 2 1 0 で変換された合成音声及び音声認識部 2 2 0 による認識結果が入力される構成となっている。第 4 実施形態に係る音声認識モデル生成部 3 1 0 は、音声変換部 2 1 0 で変換された合成音声及び音声認識部 2 2 0 による認識結果に基づいて、音声認識モデルの学習を実行可能に構成されている。

40

【 0 0 5 9 】

（音声認識モデル学習動作）

次に、図 1 0 を参照しながら、第 4 実施形態に係る音声認識システム 1 0 による音声認識モデルを学習する際の動作（以下、適宜「音声認識モデル学習動作」と称する）の流れについて説明する。図 1 0 は、第 3 実施形態に係る音声認識システムによる音声認識モデル学習動作の流れを示すフローチャートである。

【 0 0 6 0 】

図 1 0 に示すように、第 4 実施形態に係る音声認識システム 1 0 による音声認識モデル

50

学習動作が開始されると、まず音声認識モデル生成部 310 が、音声変換部 210 で変換された合成音声（即ち、音声認識部 220 に入力される合成音声）を取得する（ステップ S401）。そして、音声認識モデル生成部 310 は更に、その合成音声の音声認識結果（即ち、図 4 に示すステップ S156 で出力される音声認識結果）を取得する（ステップ S402）。

【0061】

続いて、音声認識モデル生成部 310 は、取得した合成音声及び音声認識結果に基づいて、音声認識モデルを学習する（ステップ S403）。この際、音声認識モデル生成部 310 は、すでに生成していた変換モデルのパラメータ調整を行ってよい。その後、音声認識モデル生成部 310 は、学習した音声認識モデルを音声変換部 210 に出力する（ステップ S404）。

10

【0062】

（技術的効果）

次に、第 4 実施形態に係る音声認識システム 10 によって得られる技術的効果について説明する。

【0063】

図 9 及び図 10 で説明したように、第 4 実施形態に係る音声認識システム 10 では、合成音声及び音声認識結果に基づいて変換モデルが学習される。このようにすれば、合成音声が実際にどのように音声認識されるかを考慮して学習が行われるため、より適切な音声認識が行えるように音声認識モデルを学習できる。具体的には、音声認識の精度が向上するように、音声認識モデルを学習できる。

20

【0064】

＜第 5 実施形態＞

第 5 実施形態に係る音声認識システム 10 について、図 11 及び図 12 を参照して説明する。なお、第 5 実施形態は、上述した第 1 から第 4 実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第 1 から第 4 実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

【0065】

（機能的構成）

まず、図 11 を参照しながら、第 5 実施形態に係る音声認識システム 10 の機能的構成について説明する。図 11 は、第 5 実施形態に係る音声認識システムの機能的構成を示すブロック図である。なお、図 11 では、図 2 で示した構成要素と同様の要素に同一の符号を付している。

30

【0066】

図 11 に示すように、第 5 実施形態に係る音声認識システム 10 は、その機能を実現するための構成要素として、発話データ取得部 110 と、テキスト変換部 120 と、音声合成部 130 と、変換モデル生成部 140 と、属性情報取得部 150 と、音声変換部 210 と、音声認識部 220 と、を備えて構成されている。即ち、第 5 実施形態に係る音声認識システム 10 は、第 1 実施形態の構成（図 2 参照）に加えて、属性情報取得部 150 を更に備えている。なお、属性情報取得部 150 は、例えば上述したプロセッサ 11（図 1 参照）によって実現される処理ブロックであってよい。

40

【0067】

属性情報取得部 150 は、リアル発話データの話者に関する属性情報を取得可能に構成されている。属性情報は、例えば話者の性別、年齢、職業等に関する情報を含んでよい。属性情報取得部 150 は、例えば話者が保有する端末や ID カード等から属性情報を取得可能に構成されてよい。或いは、属性情報取得部 150 は、話者が入力した属性情報を取得するように構成されてよい。属性情報取得部 150 で取得された属性情報は、音声合成部 130 に出力される構成になっている。属性情報は、リアル発話データに紐付けた状態でリアル発話音声コーパスに記憶されてもよい。この場合、属性情報は、リアル発話

50

音声コーパスから音声合成部 130 に出力されるように構成されればよい。

【0068】

(変換モデル生成動作)

次に、図12を参照しながら、第5実施形態に係る音声認識システム10による変換モデル生成動作の流れについて説明する。図12は、第5実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。なお、図12では、図3に示した処理と同様の処理に同一の符号を付している。

【0069】

図12に示すように、第5実施形態に係る音声認識システム10による変換モデル生成動作が開始されると、まず発話データ取得部110が、リアル発話データを取得する(ステップS101)。そして、属性情報取得部150が、リアル発話データの話者に関する属性情報を取得する(ステップS501)。なお、ステップS101とS102の処理は相前後して実行されてもよいし、同時に並行して実行されてもよい。

【0070】

続いて、テキスト変換部120が、発話データ取得部110で取得されたリアル発話データをテキストデータに変換する(ステップS102)。その後、音声合成部130が、テキスト変換部120で変換されたテキストデータを音声合成し、リアル発話データに対応する対応合成音声を生成するが、本実施形態では特に、属性情報も用いて音声合成を行う(ステップS502)。例えば、音声合成部130は、リアル発話データの話者の性別や年齢、職業等を考慮した音声合成を行ってよい。

【0071】

続いて、変換モデル生成部140が、発話データ取得部110で取得されたリアル発話データ及び音声合成部130で生成された対応合成音声(ここでは、属性情報に基づいて音声合成された合成音声)に基づいて、変換モデルを生成する(ステップS104)。なお、変換モデル生成部140に入力されるリアル発話データ及び対応合成音声の組には、属性情報が付与されていてよい。その場合、変換モデル生成部140は、属性情報も考慮して、変換モデルを生成してよい。その後、変換モデル生成部140は、生成した変換モデルを音声変換部210に出力する(ステップS105)。

【0072】

(技術的効果)

次に、第5実施形態に係る音声認識システム10によって得られる技術的効果について説明する。

【0073】

図11及び図12で説明したように、第5実施形態に係る音声認識システム10では、話者の属性情報を用いて対応合成音声生成される。このようにすれば、話者の属性が考慮された状態で対応合成音声生成されるため、より適切な音声変換モデルを生成することが可能となる。また、上述した第3実施形態のように、対応合成音声を用いて音声認識モデルを生成する場合(図7及び図8参照)も、属性が考慮された対応合成音声を用いることで、より適切な音声認識モデルを生成することが可能となる。

【0074】

<第6実施形態>

第6実施形態に係る音声認識システム10について、図13及び図14を参照して説明する。なお、第6実施形態は、上述した第1から第5実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第1から第5実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

【0075】

(機能的構成)

まず、図13を参照しながら、第6実施形態に係る音声認識システム10の機能的構成について説明する。図13は、第6実施形態に係る音声認識システムの機能的構成を示す

10

20

30

40

50

ブロック図である。なお、図 13 では、図 11 で示した構成要素と同様の要素に同一の符号を付している。

【0076】

図 13 に示すように、第 6 実施形態に係る音声認識システム 10 は、その機能を実現するための構成要素として、複数のリアル発話音声コーパス 105 a、105 b、及び 105 c（以下、適宜まとめて「リアル発話音声コーパス 105」と称する）と、発話データ取得部 110 と、テキスト変換部 120 と、音声合成部 130 と、変換モデル生成部 140 と、音声変換部 210 と、音声認識部 220 と、を備えて構成されている。即ち、第 6 実施形態に係る音声認識システム 10 は、第 1 実施形態の構成（図 2 参照）に加えて、複数のリアル発話音声コーパス 105 を更に備えている。なお、複数のリアル発話音声コーパス 105 は、例えば上述した記憶装置 14（図 1 参照）によって構成されてよい。

10

【0077】

複数のリアル発話音声コーパス 105 は、リアル発話データを所定の条件ごとに記憶している。ここでの「所定の条件」は、例えばリアル発話データを分類するために設定される条件である。例えば、複数のリアル発話音声コーパス 105 の各々は、分野別にリアル発話データを記憶するものであってよい。この場合、リアル発話音声コーパス 105 a が法律の分野に関するリアル発話データを記憶し、リアル発話音声コーパス 105 b が科学の分野に関するリアル発話データを記憶し、リアル発話音声コーパス 105 c が医療の分野に関するリアル発話データを記憶するように構成されてよい。なお、ここでは説明の便宜上 3 つのリアル発話音声コーパス 105 を図示しているが、リアル発話音声コーパス 105 の数は特に限定されるものではない。

20

【0078】

第 6 実施形態に係る発話データ取得部 110 は、上述した複数のリアル発話音声コーパス 105 から 1 つを選択してリアル発話データを取得可能に構成されている。なお、ここで選択されたリアル発話音声コーパス 105 に関する情報（具体的には、所定の条件に関する情報）は、リアル発話データと共に変換モデル生成部 140 に出力されてよい。そして、変換モデル生成部 140 は、変換モデルを生成する際に選択されたリアル発話音声コーパス 105 に関する情報を用いてもよい。また、上述した第 3 実施形態のように、音声認識モデルを生成する構成では、選択されたリアル発話音声コーパス 105 に関する情報が、音声認識モデル生成部 310 に出力されてもよい。そして、音声認識モデル生成部 310 は、音声認識モデルを生成する際に選択されたリアル発話音声コーパス 105 に関する情報を用いてもよい。

30

【0079】

（変換モデル生成動作）

次に、図 14 を参照しながら、第 6 実施形態に係る音声認識システム 10 による変換モデル生成動作の流れについて説明する。図 14 は、第 6 実施形態に係る音声認識システム 10 による変換モデル生成動作の流れを示すフローチャートである。なお、図 14 では、図 12 に示した処理と同様の処理に同一の符号を付している。

【0080】

図 14 に示すように、第 6 実施形態に係る音声認識システム 10 による変換モデル生成動作が開始されると、まず発話データ取得部 110 が、複数のリアル発話音声コーパス 105 の中から、発話データを取得するコーパスを選択する（ステップ S601）。そして、発話データ取得部 110 は、選択したリアル発話音声コーパスから、リアル発話データを取得する（ステップ S602）。

40

【0081】

続いて、テキスト変換部 120 が、発話データ取得部 110 で取得されたリアル発話データをテキストデータに変換する（ステップ S102）。そして、音声合成部 130 が、テキスト変換部 120 で変換されたテキストデータを音声合成し、リアル発話データに対応する対応合成音声を生成する（ステップ S103）。

【0082】

50

続いて、変換モデル生成部 140 が、発話データ取得部 110 で取得されたリアル発話データ及び音声合成部 130 で生成された対応合成音声に基づいて、変換モデルを生成するが、本実施形態では特に、選択されたリアル発話音声コーパスに関する情報も用いられる(ステップ S606)。その後、変換モデル生成部 140 は、生成した変換モデルを音声変換部 210 に出力する(ステップ S105)。

【0083】

(技術的効果)

次に、第 6 実施形態に係る音声認識システム 10 によって得られる技術的効果について説明する。

【0084】

図 13 及び図 14 で説明したように、第 6 実施形態に係る音声認識システム 10 では、変換モデルを生成する際に、リアル発話データを取得する際に選択したリアル発話音声コーパス 105 に関する情報が用いられる。このようにすれば、リアル発話データの分類に用いられた所定の条件(例えば、分野)が考慮されることになるため、より適切な変換モデルを生成することが可能となる。

【0085】

<第 7 実施形態>

第 7 実施形態に係る音声認識システム 10 について、図 15 及び図 16 を参照して説明する。なお、第 7 実施形態は、上述した第 1 から第 6 実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第 1 から第 6 実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

【0086】

(機能的構成)

まず、図 15 を参照しながら、第 7 実施形態に係る音声認識システム 10 の機能的構成について説明する。図 15 は、第 7 実施形態に係る音声認識システムの機能的構成を示すブロック図である。なお、図 15 では、図 2 で示した構成要素と同様の要素に同一の符号を付している。

【0087】

図 15 に示すように、第 7 実施形態に係る音声認識システム 10 は、その機能を実現するための構成要素として、発話データ取得部 110 と、テキスト変換部 120 と、音声合成部 130 と、変換モデル生成部 140 と、ノイズ付与部 160 と、音声変換部 210 と、音声認識部 220 と、を備えて構成されている。即ち、第 7 実施形態に係る音声認識システム 10 は、第 1 実施形態の構成(図 2 参照)に加えて、ノイズ付与部 160 を更に備えている。なお、ノイズ付与部 160 は、例えば上述したプロセッサ 11 (図 1 参照)によって実現される処理ブロックであってよい。

【0088】

ノイズ付与部 160 は、テキスト変換部 120 で生成されるテキストデータにノイズを付与可能に構成されている。ノイズ付与部 160 は、例えば、テキスト変換前のリアル発話データにノイズを付与することで、テキストデータにノイズが付与されるようにしてもよいし、テキスト変換後のテキストデータにノイズを付与するようにしてもよい。或いは、ノイズ付与部 160 は、テキスト変換部 120 がリアル発話データをテキスト変換する際にノイズを付与するようにしてもよい。ノイズ付与部 160 は、予め設定されたノイズを付与するようにしてもよいし、ランダムに設定したノイズを付与するようにしてもよい。

【0089】

(変換モデル生成動作)

次に、図 16 を参照しながら、第 7 実施形態に係る音声認識システム 10 による変換モデル生成動作の流れについて説明する。図 16 は、第 7 実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。なお、図 16 では、図 3 に示した処理と同様の処理に同一の符号を付している。

10

20

30

40

50

【 0 0 9 0 】

図 1 6 に示すように、第 7 実施形態に係る音声認識システム 1 0 による変換モデル生成動作が開始されると、まず発話データ取得部 1 1 0 が、リアル発話データを取得する（ステップ S 1 0 1）。ここで本実施形態では特に、ノイズ付与部 1 6 0 がテキスト変換部 1 2 0 にノイズ情報を出力する（ステップ S 7 0 1）。そして、テキスト変換部 1 2 0 は、発話データ取得部 1 1 0 で取得されたリアル発話データを、ノイズが付与されたテキストデータに変換する（ステップ S 7 0 2）。

【 0 0 9 1 】

続いて、音声合成部 1 3 0 が、テキスト変換部 1 2 0 で変換されたテキストデータ（ここでは、ノイズが付与されたテキストデータ）を音声合成し、リアル発話データに対応する対応合成音声を生成する（ステップ S 1 0 3）。そして、変換モデル生成部 1 4 0 が、発話データ取得部 1 1 0 で取得されたリアル発話データ及び音声合成部 1 3 0 で生成された対応合成音声に基づいて、変換モデルを生成する（ステップ S 1 0 4）。その後、変換モデル生成部 1 4 0 は、生成した変換モデルを音声変換部 2 1 0 に出力する（ステップ S 1 0 5）。

10

【 0 0 9 2 】

（技術的効果）

次に、第 7 実施形態に係る音声認識システム 1 0 によって得られる技術的効果について説明する。

【 0 0 9 3 】

図 1 5 及び図 1 6 で説明したように、第 7 実施形態に係る音声認識システム 1 0 では、リアル発話データが、ノイズが付与されたテキストデータに変換される。このようにすれば、ノイズを含むデータを用いて変換モデルが生成されることになるため、ノイズに強い変換モデル（例えば、入力音声にノイズが含まれていても適切に音声変換できる変換モデル）を生成することが可能である。

20

【 0 0 9 4 】

＜第 7 実施形態の変形例＞

第 7 実施形態の変形例に係る音声認識システム 1 0 について、図 1 7 及び図 1 8 を参照して説明する。なお、第 7 実施形態の変形例は、上述した第 7 実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第 1 から第 7 実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

30

【 0 0 9 5 】

（機能的構成）

まず、図 1 7 を参照しながら、第 7 実施形態の変形例に係る音声認識システム 1 0 の機能的構成について説明する。図 1 7 は、第 7 実施形態の変形例に係る音声認識システムの機能的構成を示すブロック図である。なお、図 1 7 では、図 1 5 で示した構成要素と同様の要素に同一の符号を付している。

【 0 0 9 6 】

図 1 7 に示すように、第 7 実施形態の変形例に係る音声認識システム 1 0 は、その機能を実現するための構成要素として、発話データ取得部 1 1 0 と、テキスト変換部 1 2 0 と、音声合成部 1 3 0 と、変換モデル生成部 1 4 0 と、ノイズ付与部 1 6 0 と、音声変換部 2 1 0 と、音声認識部 2 2 0 と、を備えて構成されている。ただし、第 7 実施形態の変形例に係る音声認識システム 1 0 では、ノイズ付与部 1 6 0 が、音声合成部 1 3 0 にノイズ情報を出力可能に構成されている。即ち、第 7 実施形態の変形例では、音声合成部 1 3 0 による音声合成の際にノイズが付与される構成となっている。

40

【 0 0 9 7 】

（変換モデル生成動作）

次に、図 1 8 を参照しながら、第 7 実施形態の変形例に係る音声認識システム 1 0 による変換モデル生成動作の流れについて説明する。図 1 8 は、第 7 実施形態の変形例に係る

50

音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。なお、図 18 では、図 16 に示した処理と同様の処理に同一の符号を付している。

【0098】

図 18 に示すように、第 7 実施形態の変形例に係る音声認識システム 10 による変換モデル生成動作が開始されると、まず発話データ取得部 110 が、リアル発話データを取得する（ステップ S101）。そして、テキスト変換部 120 が、発話データ取得部 110 で取得されたリアル発話データをテキストデータに変換する（ステップ S102）。

【0099】

続いて、本実施形態では特に、ノイズ付与部 160 が音声合成部 130 にノイズ情報を出力する（ステップ S751）。そして、音声合成部 130 は、テキスト変換部 120 で変換されたテキストデータを音声合成し、ノイズが付与された対応合成音声を生

10

【0100】

成する（ステップ S752）。続いて、変換モデル生成部 140 が、発話データ取得部 110 で取得されたリアル発話データ及び音声合成部 130 で生成された対応合成音声（ここでは、ノイズが付与された対応合成音声）に基づいて、変換モデルを生成する（ステップ S104）。その後、変換モデル生成部 140 は、生成した変換モデルを音声変換部 210 に出力する（ステップ S105）。

【0101】

（技術的効果）

次に、第 7 実施形態の変形例に係る音声認識システム 10 によって得られる技術的効果について説明する。

20

【0102】

図 17 及び図 18 で説明したように、第 7 実施形態の変形例に係る音声認識システム 10 では、ノイズが付与された対応合成音声が生

【0103】

< 第 8 実施形態 >

第 8 実施形態に係る音声認識システム 10 について、図 19 から図 21 を参照して説明する。なお、第 8 実施形態は、上述した第 1 から第 7 実施形態と一部の構成及び動作が異なるのみであり、その他の部分については第 1 から第 7 実施形態と同一であってよい。このため、以下では、すでに説明した各実施形態と異なる部分について詳細に説明し、その他の重複する部分については適宜説明を省略するものとする。

30

【0104】

（機能的構成）

まず、図 19 を参照しながら、第 8 実施形態に係る音声認識システム 10 の機能的構成について説明する。図 19 は、第 8 実施形態に係る音声認識システムの機能的構成を示すブロック図である

40

【0105】

図 19 に示すように、第 8 実施形態に係る音声認識システム 10 は、その機能を実現するための構成要素として、手話データ取得部 410 と、テキスト変換部 420 と、音声合成部 430 と、変換モデル生成部 440 と、音声変換部 510 と、音声認識部 520 と、を備えて構成されている。手話データ取得部 410、テキスト変換部 420、音声合成部 430、変換モデル生成部 440、音声変換部 510、音声認識部 520 の各々は、例えば上述したプロセッサ 11（図 1 参照）によって実現される処理ブロックであってよい。

【0106】

手話データ取得部 410 は、手話発話データを取得可能に構成されている。手話データは、例えば手話の動画データであってよい。手話データは、例えば複数の手話データを蓄

50

積するデータベース（手話コーパス）から取得されてよい。手話データ取得部 4 1 0 で取得された手話データは、テキスト変換部 1 2 0 及び変換モデル生成部 1 4 0 に出力される構成となっている。

【 0 1 0 7 】

テキスト変換部 4 2 0 は、手話データ取得部 4 1 0 で取得された手話データをテキストデータに変換可能に構成されている。即ち、テキスト変換部 4 2 0 は、手話データに含まれる手話の内容をテキスト変換する処理を実行可能に構成されている。なお、テキスト変換の具体的な手法については、既存の技術が適宜採用されてよい。テキスト変換部 4 2 0 で変換されたテキストデータ（即ち、手話データに対応するテキストデータ）は、音声合成部 4 3 0 に出力される構成となっている。

10

【 0 1 0 8 】

音声合成部 4 3 0 は、テキスト変換部 4 2 0 で変化されたテキストデータを音声合成することで、手話データに対応する対応合成音声を生成可能に構成されている。なお、音声合成の具体的な手法については、既存の技術を適宜採用することができる。音声合成部 4 3 0 で生成された対応合成音声は、変換モデル生成部 4 4 0 に出力される構成となっている。なお、対応合成音声は、複数の対応合成を蓄積可能なデータベース（合成音声コーパス）に蓄積されてから、変換モデル生成部 4 4 0 に出力されてもよい。

【 0 1 0 9 】

変換モデル生成部 4 4 0 は、手話データ取得部 4 1 0 で取得された手話データと、音声合成部 4 3 0 で合成された対応合成音声を用いて、入力手話を合成音声に変換する変換モデルを生成可能に構成されている。変換モデルは、例えば、入力される入力手話（例えば、手話の動画）を、合成音声（即ち、機械的な音声）に変換する。変換モデル生成部 4 4 0 は、例えば G A N を用いて、変換モデルを生成するように構成されてよい。変換モデル生成部 4 4 0 で生成された変換モデルは、音声変換部 5 1 0 に出力される構成となっている。

20

【 0 1 1 0 】

音声変換部 5 1 0 は、変換モデル生成部 4 4 0 で生成された変換モデルを用いて、入力手話を合成音声に変換可能に構成されている。音声変換部 5 1 0 に入力される入力手話は、例えばカメラ等を用いて入力される動画であってよい。音声変換部 5 1 0 で変換された合成音声は、音声認識部 5 2 0 に出力される構成となっている。

30

【 0 1 1 1 】

音声認識部 5 2 0 は、音声変換部 5 1 0 で変換された合成音声を音声認識することが可能に構成されている。即ち、音声認識部 5 2 0 は、合成音声をテキスト化する処理を実行可能に構成されている。音声認識部 5 2 0 は、合成音声の音声認識結果を出力可能に構成されてよい。なお、音声認識結果の利用方法については特に限定されない。

【 0 1 1 2 】

（変換モデル生成動作）

次に、図 2 0 を参照しながら、第 8 実施形態に係る音声認識システム 1 0 による変換モデル生成動作の流れについて説明する。図 2 0 は、第 8 実施形態に係る音声認識システムによる変換モデル生成動作の流れを示すフローチャートである。

40

【 0 1 1 3 】

図 2 0 に示すように、第 8 実施形態に係る音声認識システム 1 0 による変換モデル生成動作が開始されると、まず手話データ取得部 4 1 0 が、手話データを取得する（ステップ S 8 0 1）。そして、テキスト変換部 4 2 0 が、手話データ取得部 4 1 0 で取得された手話データをテキストデータに変換する（ステップ S 8 0 2）。

【 0 1 1 4 】

続いて、音声合成部 4 3 0 が、テキスト変換部 4 2 0 で変換されたテキストデータを音声合成し、手話データに対応する対応合成音声を生成する（ステップ S 4 0 3）。そして、変換モデル生成部 1 4 0 が、手話データ取得部 4 1 0 で取得された手話データ及び音声合成部 4 3 0 で生成された対応合成音声に基づいて、変換モデルを生成する（ステップ S

50

804)。その後、変換モデル生成部440は、生成した変換モデルを音声変換部510に出力する(ステップS805)。

【0115】

(変換認識動作)

次に、図21を参照しながら、第8実施形態に係る音声認識システム10による音声認識動作の流れについて説明する。図21は、第8実施形態に係る音声認識システムによる音声認識動作の流れを示すフローチャートである。

【0116】

図21に示すように、第1実施形態に係る音声認識システム10による音声認識動作が開始されると、まず音声変換部510が入力手話を取得する(ステップS851)。そして、音声変換部510は、変換モデル生成部440で生成された変換モデルを読み込む(ステップS852)。その後、音声変換部210は、読み込んだ変換モデルを用いて音声変換を行い、入力手話を合成音声に変換する(ステップS853)。

10

【0117】

続いて、音声認識部520は、音声認識モデルを読み込む(ステップS854)。そして、音声認識部520は、読み込んだ音声認識モデルを用いて、音声変換部510で合成された合成音声を音声認識する(ステップS855)。その後、音声認識部520は、音声認識結果を出力する(ステップS856)。

【0118】

(技術的効果)

次に、第8実施形態に係る音声認識システム10によって得られる技術的効果について説明する。

20

【0119】

図19から図21で説明したように、第8実施形態に係る音声認識システム10では、変換モデルを生成する際に、手話データ及び手話データに対応する対応合成音声がいられる。そして特に、手話データに対応する対応合成音声は、手話データをテキスト変換し、テキストデータを音声合成することで生成される。このようにすれば、手話データと、それに対応する合成音声と、の両方を用意する必要がなくなる(即ち、手話データのみ用意すれば、対応合成音声を生成できる)ため、変換モデルを生成するのに要するコストを抑制することができる。その結果、低コストで認識精度の高い音声認識を実現することが可能となる。

30

【0120】

上述した各実施形態の機能を実現するように該実施形態の構成を動作させるプログラムを記録媒体に記録させ、該記録媒体に記録されたプログラムをコードとして読み出し、コンピュータにおいて実行する処理方法も各実施形態の範疇に含まれる。すなわち、コンピュータ読取可能な記録媒体も各実施形態の範囲に含まれる。また、上述のプログラムが記録された記録媒体はもちろん、そのプログラム自体も各実施形態に含まれる。

【0121】

記録媒体としては例えばフロッピー(登録商標)ディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、磁気テープ、不揮発性メモリカード、ROMを用いることができる。また該記録媒体に記録されたプログラム単体で処理を実行しているものに限らず、他のソフトウェア、拡張ボードの機能と共同して、OS上で動作して処理を実行するものも各実施形態の範疇に含まれる。更に、プログラム自体がサーバに記憶され、ユーザ端末にサーバからプログラムの一部または全てをダウンロード可能なようにしてもよい。

40

【0122】

<付記>

以上説明した実施形態に関して、更に以下の付記のようにも記載されうるが、以下には限られない。

【0123】

50

(付記 1)

付記 1 に記載の音声認識システムは、話者が発話したリアル発話データを取得する発話データ取得手段と、前記リアル発話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を生成する音声合成手段と、前記リアル発話データ及び前記対応合成音声をを用いて、入力音声を合成音声に変換する変換モデルを生成する変換モデル生成手段と、前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える音声認識システムである。

【 0 1 2 4 】

(付記 2)

付記 2 に記載の音声認識システムは、前記変換モデル生成手段は、前記入力音声と、前記音声認識手段の認識結果と、を用いて前記変換モデルのパラメータを調整する、付記 1 に記載の音声認識システムである。

【 0 1 2 5 】

(付記 3)

付記 3 に記載の音声認識システムは、前記対応合成音声を含むデータを用いて音声認識モデルを生成する音声認識モデル生成手段を更に備え、前記音声認識手段は、前記音声認識モデルを用いて音声認識する、付記 1 又は 2 に記載の音声認識システムである。

【 0 1 2 6 】

(付記 4)

付記 4 に記載の音声認識システムは、前記音声認識モデル生成手段は、前記変換モデルを用いて変換された前記合成音声と、前記音声認識手段の認識結果と、を用いて前記音声認識モデルのパラメータを調整する、付記 3 に記載の音声認識システムである。

【 0 1 2 7 】

(付記 5)

付記 5 に記載の音声認識システムは、前記話者の属性を示す属性情報を取得する属性取得手段を更に備え、前記音声合成手段は、前記属性情報を用いて音声合成を行うことで前記対応合成音声を生成する、付記 1 から 4 のいずれか一項に記載の音声認識システムである。

【 0 1 2 8 】

(付記 6)

付記 6 に記載の音声認識システムは、所定の条件ごとに前記リアル発話データを記憶する複数のリアル発話音声コーパスを更に備え、前記発話データ取得手段は、前記複数のリアル発話音声コーパスから 1 つを選択して前記リアル発話データを取得する、付記 1 から 5 のいずれか一項に記載の音声認識システムである。

【 0 1 2 9 】

(付記 7)

付記 7 に記載の音声認識システムは、前記テキストデータ及び前記対応合成音声の少なくとも一方にノイズを付与するノイズ付与手段を更に備える、付記 1 から 6 のいずれか一項に記載の音声認識システムである。

【 0 1 3 0 】

(付記 8)

付記 8 に記載の音声認識システムは、手話データを取得する手話データ取得手段と、前記手話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記手話データに対応する対応合成音声を生成する音声合成手段と、前記手話データ及び前記対応合成音声をを用いて、入力される手話を合成音声に変換する変換モデルを生成する変換モデル生成手段と、前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える音声認識システムである。

【 0 1 3 1 】

(付記 9)

10

20

30

40

50

付記 9 に記載の音声認識方法は、少なくとも 1 つのコンピュータによって、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を作成し、前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法である。

【 0 1 3 2 】

(付記 1 0)

付記 1 0 に記載の記録媒体は、少なくとも 1 つのコンピュータに、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を作成し、前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法を実行させるコンピュータプログラムが記録された記録媒体である。

10

【 0 1 3 3 】

(付記 1 1)

付記 1 1 に記載のコンピュータプログラムは、少なくとも 1 つのコンピュータに、話者が発話したリアル発話データを取得し、前記リアル発話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を作成し、前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法を実行させるコンピュータプログラムである。

20

【 0 1 3 4 】

(付記 1 2)

付記 1 2 に記載の音声認識装置は、話者が発話したリアル発話データを取得する発話データ取得手段と、前記リアル発話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記リアル発話データに対応する対応合成音声を作成する音声合成手段と、前記リアル発話データ及び前記対応合成音声を用いて、入力音声を合成音声に変換する変換モデルを生成する変換モデル生成手段と、前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える音声認識装置である。

30

【 0 1 3 5 】

(付記 1 3)

付記 1 3 に記載の音声認識方法は、少なくとも 1 つのコンピュータによって、手話データを取得し、前記手話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記手話データに対応する対応合成音声を作成し、前記手話データ及び前記対応合成音声を用いて、入力される手話を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法である。

【 0 1 3 6 】

(付記 1 4)

付記 1 4 に記載の記録媒体は、少なくとも 1 つのコンピュータに、手話データを取得し、前記手話データをテキストデータに変換し、前記テキストデータを用いた音声合成により、前記手話データに対応する対応合成音声を作成し、前記手話データ及び前記対応合成音声を用いて、入力される手話を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法を実行させるコンピュータプログラムが記録された記録媒体である。

40

【 0 1 3 7 】

(付記 1 5)

付記 1 5 に記載のコンピュータプログラムは、少なくとも 1 つのコンピュータに、手話データを取得し、前記手話データをテキストデータに変換し、前記テキストデータを用い

50

た音声合成により、前記手話データに対応する対応合成音声を生成し、前記手話データ及び前記対応合成音声を用いて、入力される手話を合成音声に変換する変換モデルを生成し、前記変換モデルを用いて変換された前記合成音声を音声認識する、音声認識方法を実行させるコンピュータプログラムである。

【 0 1 3 8 】

(付記 1 6)

付記 1 6 に記載の音声認識装置は、手話データを取得する手話データ取得手段と、前記手話データをテキストデータに変換するテキスト変換手段と、前記テキストデータを用いた音声合成により、前記手話データに対応する対応合成音声を生成する音声合成手段と、前記手話データ及び前記対応合成音声を用いて、入力される手話を合成音声に変換する変換モデルを生成する変換モデル生成手段と、前記変換モデルを用いて変換された前記合成音声を音声認識する音声認識手段と、を備える音声認識装置である。

10

【 0 1 3 9 】

この開示は、請求の範囲及び明細書全体から読み取ることのできる発明の要旨又は思想に反しない範囲で適宜変更可能であり、そのような変更を伴う音声認識システム、音声認識方法、及び記録媒体もまたこの開示の技術思想に含まれる。

【符号の説明】

【 0 1 4 0 】

1 0 音声認識システム

1 1 プロセッサ

1 4 記憶装置

1 0 5 リアル発話音声コーパス

1 1 0 発話データ取得部

1 2 0 テキスト変換部

1 3 0 音声合成部

1 4 0 変換モデル生成部

1 5 0 属性情報取得部

1 6 0 ノイズ付与部

2 1 0 音声変換部

2 2 0 音声認識部

3 1 0 音声認識モデル生成部

4 1 0 手話データ取得部

4 2 0 テキスト変換部

4 3 0 音声合成部

4 4 0 変換モデル生成部

5 1 0 音声変換部

5 2 0 音声認識部

20

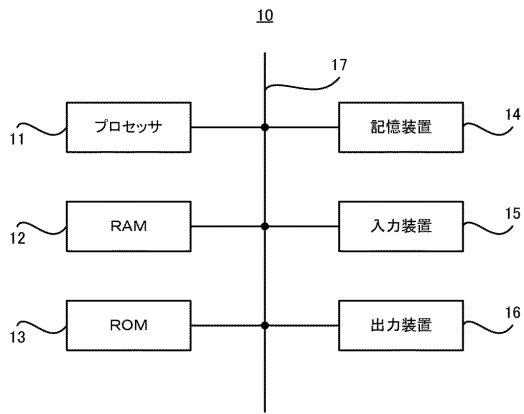
30

40

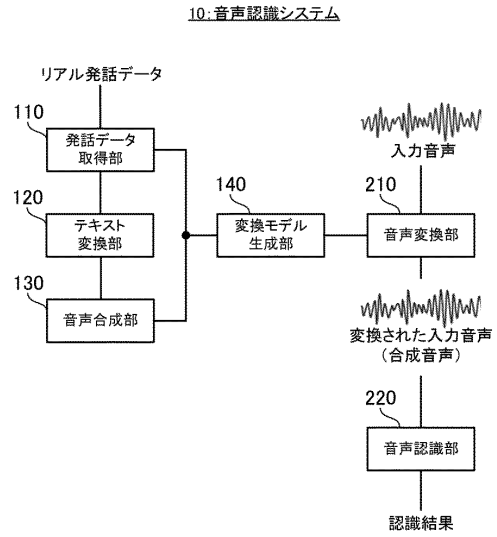
50

【図面】

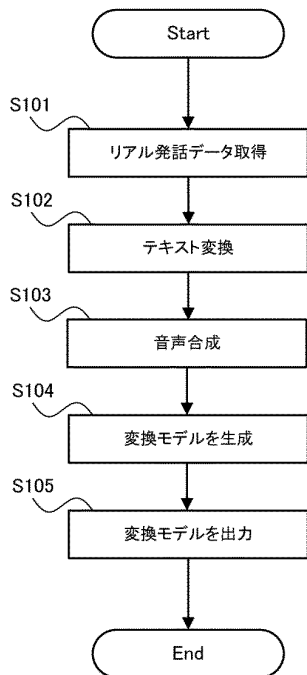
【図 1】



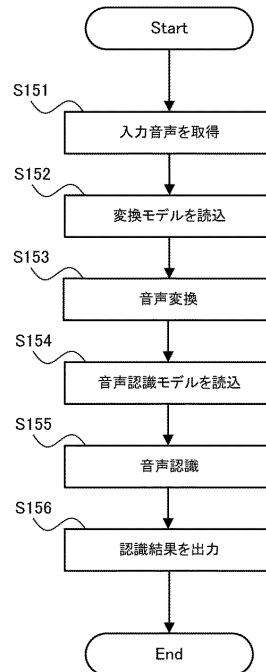
【図 2】



【図 3】



【図 4】



10

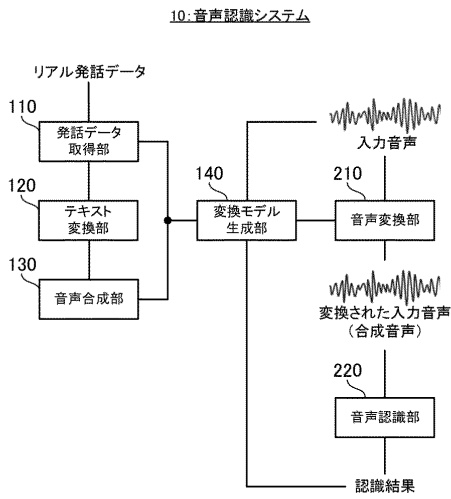
20

30

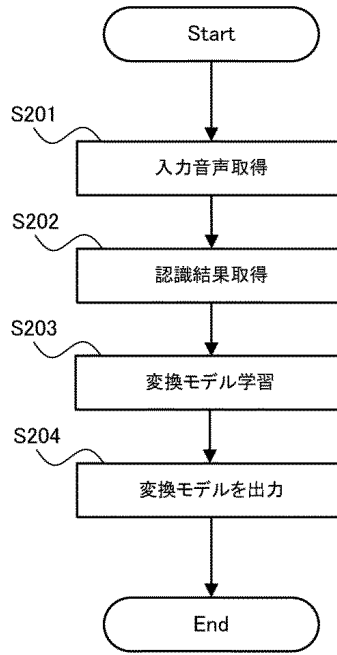
40

50

【 図 5 】



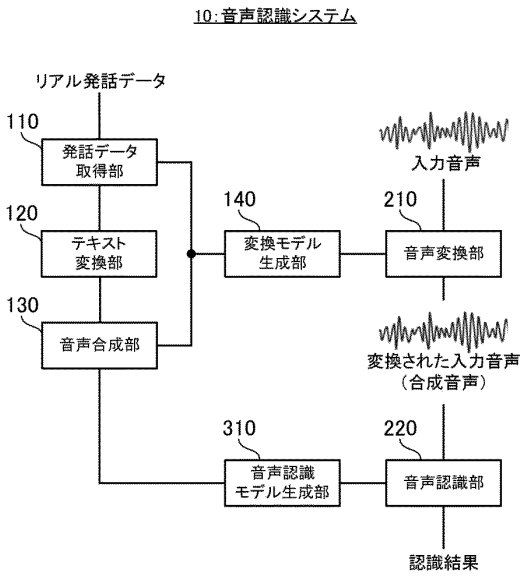
【 図 6 】



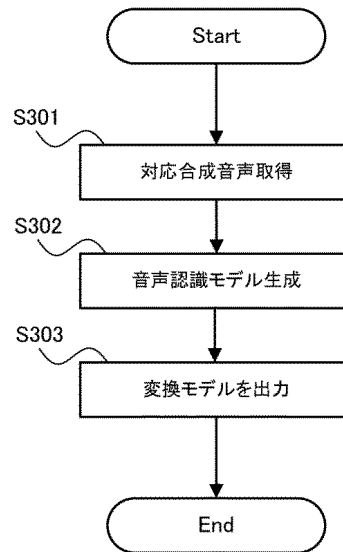
10

20

【 図 7 】



【 図 8 】

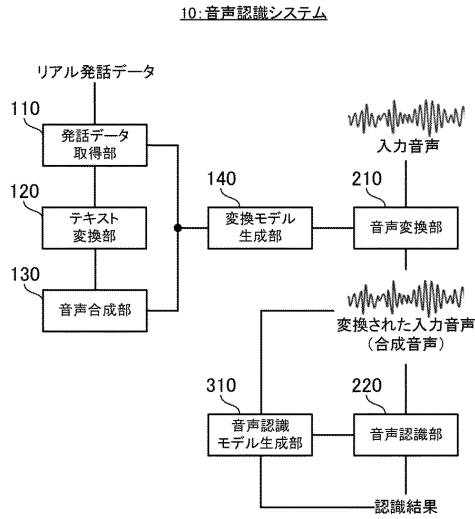


30

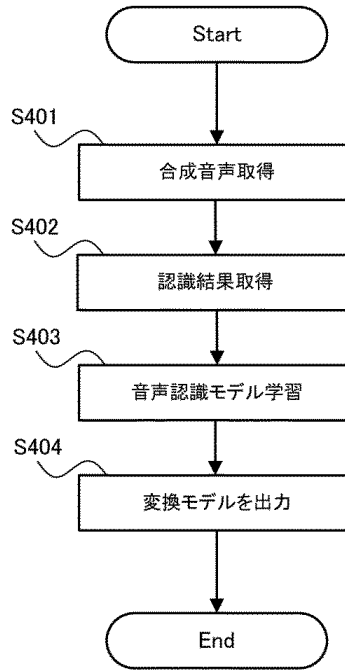
40

50

【 図 9 】



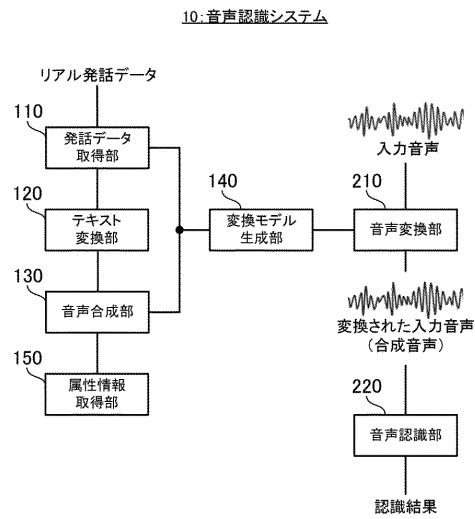
【 図 10 】



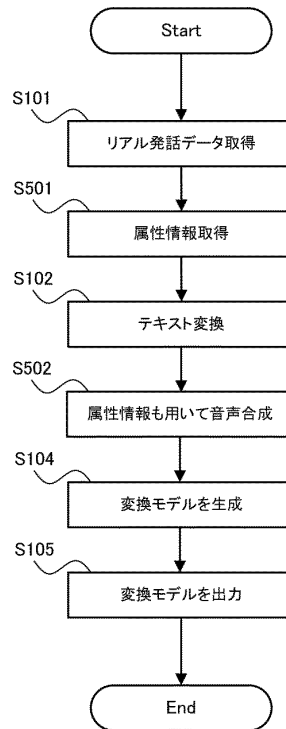
10

20

【 図 11 】



【 図 12 】

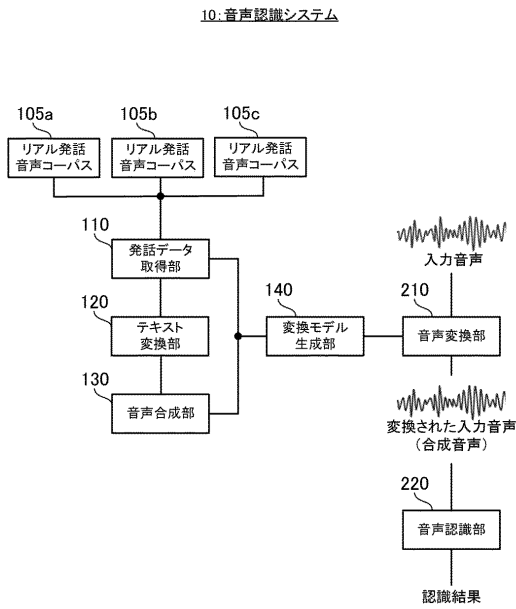


30

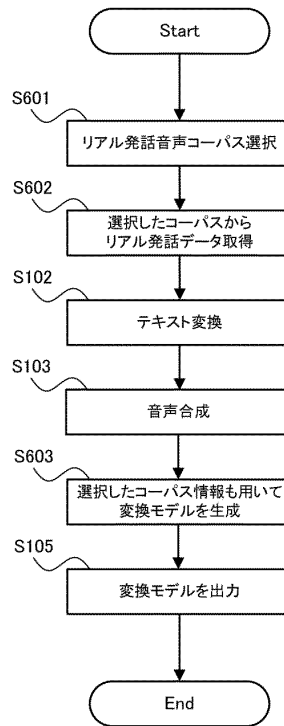
40

50

【 図 1 3 】



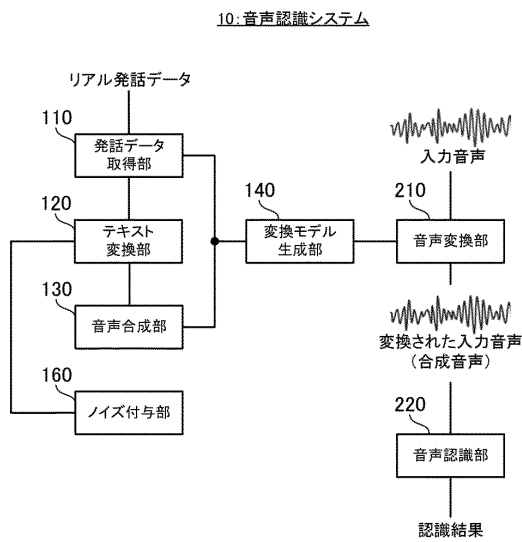
【 図 1 4 】



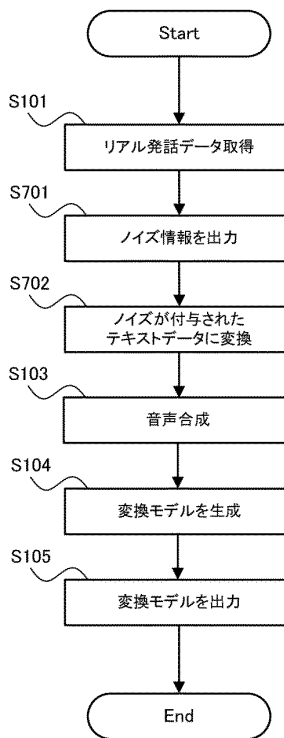
10

20

【 図 1 5 】



【 図 1 6 】

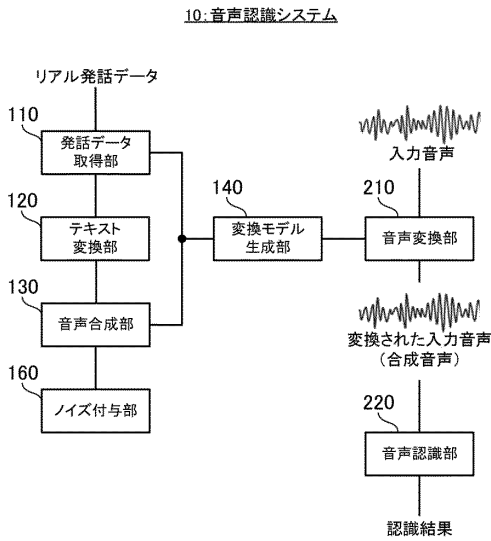


30

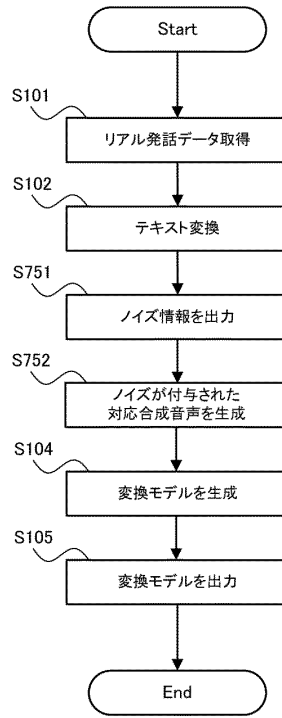
40

50

【図 17】



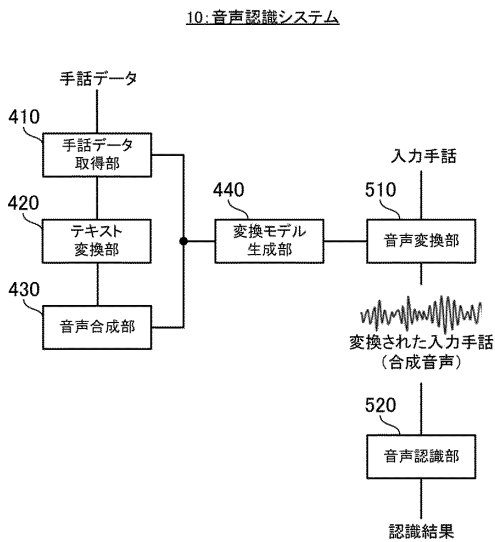
【図 18】



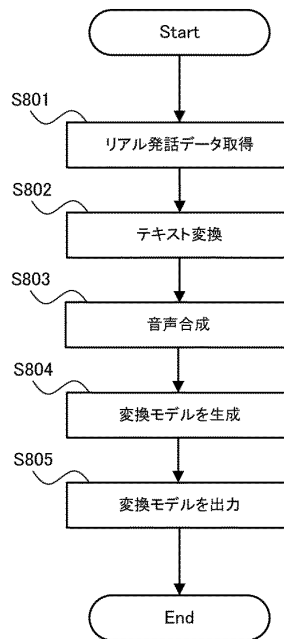
10

20

【図 19】



【図 20】

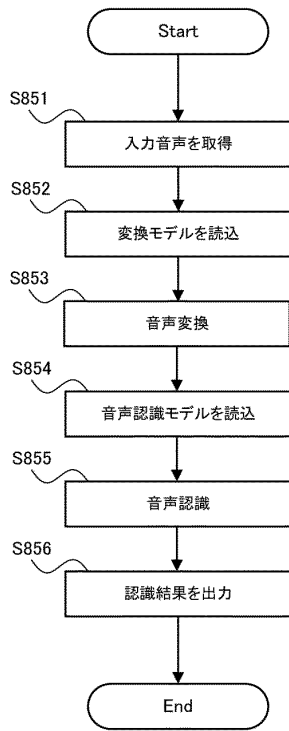


30

40

50

【図 21】



10

20

30

40

50

フロントページの続き

審査官 大野 弘

(56)参考文献 特開 2019 - 008120 (JP, A)

特表 2003 - 522978 (JP, A)

(58)調査した分野 (Int.Cl., DB名)

G10L 21/007

G10L 13/02

G10L 15/20