

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6903613号  
(P6903613)

(45) 発行日 令和3年7月14日(2021.7.14)

(24) 登録日 令和3年6月25日(2021.6.25)

(51) Int.Cl.

F I

G 1 0 L 15/06 (2013.01)

G 1 0 L 15/06 3 0 0 Y

請求項の数 13 (全 14 頁)

<p>(21) 出願番号 特願2018-168708 (P2018-168708)</p> <p>(22) 出願日 平成30年9月10日 (2018.9.10)</p> <p>(65) 公開番号 特開2020-42130 (P2020-42130A)</p> <p>(43) 公開日 令和2年3月19日 (2020.3.19)</p> <p>審査請求日 令和2年7月1日 (2020.7.1)</p>	<p>(73) 特許権者 000003078 株式会社東芝 東京都港区芝浦一丁目1番1号</p> <p>(74) 代理人 110002147 特許業務法人酒井国際特許事務所</p> <p>(72) 発明者 丁 寧 東京都港区芝浦一丁目1番1号 株式会社 東芝内</p> <p>審査官 中村 天真</p>
--	---

最終頁に続く

(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) 【特許請求の範囲】

【請求項1】

言語モデルと第1音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する生成部と、

前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付与する重みを前記個数に応じて決定する決定部と、

前記重みに基づいて前記音声データを選択する選択部と、

前記選択部により選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成する適応部と、

を備える音声認識装置。

10

【請求項2】

言語モデルと第1音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する生成部と、

前記ラベルを用いて、前記音声データに含まれる音声フレームと、前記音声データに含まれる非音声フレームとの比率を計算する計算部と、

前記音声フレームの比率が所定の選択範囲以内である音声データを選択する選択部と、

前記選択部により選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成する適応部と、

を備える音声認識装置。

【請求項3】

20

言語モデルと第 1 音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する生成部と、

前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付与する重みを前記個数に応じて決定する決定部と、

前記ラベルを用いて、前記音声データに含まれる音声フレームと、前記音声データに含まれる非音声フレームとの比率を計算する計算部と、

前記重みと、前記音声フレームの比率とに基づいて、前記音声データを選択する選択部と、

前記選択部により選択された音声データを用いて、前記第 1 音響モデルを適応させることにより、第 2 音響モデルを生成する適応部と、

を備える音声認識装置。

10

【請求項 4】

前記決定部は、前記個数が多いほど、前記重みを小さく決定する、請求項 1 に記載の音声認識装置。

【請求項 5】

前記選択部は、前記重みが閾値より大きいか否かを判定し、前記閾値よりも大きい重みが付与された音声データを選択する、

請求項 1 に記載の音声認識装置。

【請求項 6】

前記言語モデルと前記第 2 音響モデルとを用いて、前記音声データの音声認識を行う認識部、

を更に備える請求項 1 に記載の音声認識装置。

20

【請求項 7】

前記言語モデル及び前記第 1 音響モデルは、音声の言語的な特徴、及び、音声の音響的な特徴の両方を区別せずにモデル化された 1 つの音声認識ベースモデルにより表され、

前記適応部は、前記選択部により選択された音声データを用いて、前記音声認識ベースモデルを適応させる、

請求項 1 に記載の音声認識装置。

【請求項 8】

言語モデルと第 1 音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成するステップと、

前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付与する重みを前記個数に応じて決定するステップと、

前記重みに基づいて前記音声データを選択するステップと、

前記選択するステップにより選択された音声データを用いて、前記第 1 音響モデルを適応させることにより、第 2 音響モデルを生成するステップと、

を含む音声認識方法。

30

【請求項 9】

言語モデルと第 1 音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成するステップと、

前記ラベルを用いて、前記音声データに含まれる音声フレームと、前記音声データに含まれる非音声フレームとの比率を計算するステップと、

前記音声フレームの比率が所定の選択範囲以内である音声データを選択する選択部と、

前記選択するステップにより選択された音声データを用いて、前記第 1 音響モデルを適応させることにより、第 2 音響モデルを生成するステップと、

を含む音声認識方法。

40

【請求項 10】

言語モデルと第 1 音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成するステップと、

前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付

50

与する重みを前記個数に応じて決定するステップと、

前記ラベルを用いて、前記音声データに含まれる音声フレームと、前記音声データに含まれる非音声フレームとの比率を計算するステップと、

前記重みと、前記音声フレームの比率とに基づいて、前記音声データを選択するステップと、

前記選択するステップにより選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成するステップと、

を含む音声認識方法。

【請求項11】

コンピュータを、

言語モデルと第1音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する生成部と、

前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付与する重みを前記個数に応じて決定する決定部と、

前記重みに基づいて前記音声データを選択する選択部と、

前記選択部により選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成する適応部、

として機能させるためのプログラム。

【請求項12】

コンピュータを、

言語モデルと第1音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する生成部と、

前記ラベルを用いて、前記音声データに含まれる音声フレームと、前記音声データに含まれる非音声フレームとの比率を計算する計算部と、

前記音声フレームの比率が所定の選択範囲以内である音声データを選択する選択部と、

前記選択部により選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成する適応部、

として機能させるためのプログラム。

【請求項13】

コンピュータを、

言語モデルと第1音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する生成部と、

前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付与する重みを前記個数に応じて決定する決定部と、

前記ラベルを用いて、前記音声データに含まれる音声フレームと、前記音声データに含まれる非音声フレームとの比率を計算する計算部と、

前記重みと、前記音声フレームの比率とに基づいて、前記音声データを選択する選択部と、

前記選択部により選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成する適応部、

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】

本発明の実施形態は音声認識装置、音声認識方法及びプログラムに関する。

【背景技術】

【0002】

音響モデルと言語モデルとを用いて音声データを認識し、音声データに含まれる発話のテキストを出力する音声認識技術が従来から知られている。音響モデルは予め大量（例えば数百時間以上）のデータを用いて学習される。しかし、どのような条件で用いても高い

10

20

30

40

50

認識率（例えば85%以上）が得られるような音響モデルを学習することは困難である。例えば、クリーンな環境で収録された音声データを用いて学習された音響モデルが用いられた場合、残響が大きい会議室での認識率が劣化してしまう。認識率の劣化を防ぐ有効な方法の一つとして、音響モデルの適応がある。

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特許第5852550号公報

【発明の概要】

【発明が解決しようとする課題】

10

【0004】

しかしながら、従来の技術では、音響モデルの適応を行った場合、悪影響も生じていた。例えば、同じ内容の発話が繰り返された場合、音響モデルの適応によって、この発話を認識しやすくなるが、他の発話を認識しにくくなる。また例えば、音声データには音声及び非音声の両方が含まれているが、非音声の部分が多い場合、音響モデルの適応によって、非音声の認識結果が出やすくなり、音声の認識結果が出にくくなる。本発明が解決しようとする課題は、音響モデルの適応による悪影響を抑制できる音声認識装置、音声認識方法及びプログラムを提供することである。

【課題を解決するための手段】

【0005】

20

実施形態の音声認識装置は、生成部と決定部と選択部と適応部とを備える。生成部は、言語モデルと第1音響モデルとを用いて音声データを認識し、前記音声データに含まれる発話を識別するラベルを生成する。決定部は、前記ラベルを用いて、同じ発話を含む音声データの個数を特定し、前記音声データに付与する重みを前記個数に応じて決定する。選択部は、前記重みに基づいて前記音声データを選択する。適応部は、前記選択部により選択された音声データを用いて、前記第1音響モデルを適応させることにより、第2音響モデルを生成する。

【図面の簡単な説明】

【0006】

【図1】第1実施形態の音声認識装置の機能構成の例を示すブロック図。

30

【図2】第1実施形態のラベル情報の例を示す図。

【図3】第1実施形態の音声認識装置の動作方法の例を示すフローチャート。

【図4】第2実施形態の音声認識装置の機能構成の例を示すブロック図。

【図5】第2実施形態の音声データの例を示す図。

【図6】第3実施形態の音声認識装置の機能構成の例を示すブロック図。

【図7】第4実施形態の音声認識装置の機能構成の例を示すブロック図。

【図8】第5実施形態の音声認識装置の機能構成の例を示すブロック図。

【図9】第1乃至第5実施形態の音声認識装置のハードウェア構成の例を示す図。

【発明を実施するための形態】

【0007】

40

以下に添付図面を参照して、音声認識装置、音声認識方法及びプログラムの実施形態を詳細に説明する。

【0008】

はじめに、音響モデルの適応について説明する。音響モデルの適応は、学習された音響モデルに基づき、適応データを用いて当該音響モデルを再学習することにより行われる。以下、はじめに学習された音響モデルをベース音響モデル（第1音響モデル）と呼び、適応させた音響モデルを適応音響モデル（第2音響モデル）と呼ぶ。

【0009】

音響モデルを適応させる方法は、おおむね二種類ある（教師あり適応及び教師なし適応）。教師あり適応では、音声データ、及び、音声データの正解ラベルの両方を含む適応デ

50

ータが使用される。教師なし適応では、音声データのみを含む適応データが使用される（正解ラベルがない。）。

【 0 0 1 0 】

教師あり適応は正解ラベルがあるため適応に対して良いが、書き起こしなどによって正解ラベルを作成する必要があるため、コストが高い。

【 0 0 1 1 】

一方、教師なし適応は正解ラベルの作成が要らないためコストが低い。教師なし適応では、音声データを認識し、音声認識結果をラベルとして用いる。音声認識結果の誤りは適応に悪影響を及ぼす可能性があるため、基本的には、音声認識精度は高いほどよい。従来の教師なし適応方法では、言語モデルとベース音響モデルとを用いて音声を認識し、ラベル、信頼度及び音響尤度を出力する。従来の教師なし適応方法では、信頼度がより高く、かつ、音響尤度がより小さい音声データを選択して、音響モデルの適応を行う。

【 0 0 1 2 】

（第1実施形態）

はじめに、第1実施形態の音声認識装置10の機能構成の例について説明する。

【 0 0 1 3 】

[ 機能構成の例 ]

図1は第1実施形態の音声認識装置10の機能構成の例を示す図である。第1実施形態の音声認識装置10は、生成部1、決定部2、選択部3及び適応部4を備える。音声認識装置10の一部又は全ての機能は、ソフトウェア（プログラム）で実現されても良いし、ハードウェアで実現されても良い。

【 0 0 1 4 】

また、第1実施形態の音声認識装置10は、言語モデル101、ベース音響モデル102及び適応音響モデル103を記憶する。言語モデル101は、音声の言語的な特徴をモデル化したデータである。ベース音響モデル102及び適応音響モデル103は、音声の音響的な特徴をモデル化したデータである。ベース音響モデル102は、はじめに学習されたデータである。適応音響モデル103は、適応データを用いてベース音響モデル102を再学習することにより得られたデータである。なお、言語モデル101、ベース音響モデル102及び適応音響モデル103を記憶する記憶部は、外部の装置に備えられていてもよい。

【 0 0 1 5 】

生成部1は、言語モデル101とベース音響モデル102とを用いて音声データを認識し、ラベルを生成する。音声データは、例えば発話毎に区切られたデータである。ラベルは、音声データの音声認識結果から変換されたデータである。ラベルは、音声データに含まれる発話を識別する情報である。

【 0 0 1 6 】

決定部2は、ラベルを用いて、同じ発話を含む音声データの個数を特定し、音声データに付与する重みを、当該個数に応じて決定する。

【 0 0 1 7 】

ラベル及び当該ラベルの個数は、例えば図2に示すラベル情報として、音声認識装置10に記憶される。

【 0 0 1 8 】

図2は第1実施形態のラベル情報の例を示す図である。第1実施形態のラベル情報は、音声データ、音声認識結果、ラベル、カウント数及び重みを含む。

【 0 0 1 9 】

音声認識結果は、音声データの認識結果である。図2の例では、ラベルは、音声認識結果をひらがなに変換したデータである。なお、ラベルは、ひらがなに限らずローマ字等でもよい。

【 0 0 2 0 】

カウント数は、ラベルの個数を示す。例えば、発話 - 1、発話 - 3 及び発話 - 5 のラベ

10

20

30

40

50

ルは同じである。発話 - 1 のラベル生成時には、当該ラベルのカウント数は 1 となる。発話 - 3 のラベル生成時には、当該ラベルのカウント数は 2 となる。発話 - 5 のラベル生成時には、当該ラベルのカウント数は 3 となる。

【 0 0 2 1 】

重みは、ラベルの重みを示す。図 2 の例では、ラベルのカウント数が大きいほど、当該ラベルの重みは小さくなる。

【 0 0 2 2 】

生成部 1 は、例えば下記式 ( 1 ) により、ラベルの重みを決定する。

【 0 0 2 3 】

$$\mu = e^{-x} \cdot \dots (1)$$

10

【 0 0 2 4 】

ここで、 $\mu$  は重みであり、 $x$  はカウント数である。図 2 の例では、式 ( 1 ) により重みが決定されている。例えば発話 - 1、発話 - 3 及び発話 - 5 のラベルの重みは、それぞれ 1 . 0 0、0 . 3 7、0 . 1 4 である。発話 - 2、発話 - 4 及び発話 - 6 のラベルの重みは、1 . 0 0 である。

【 0 0 2 5 】

なお、重みを決定する式は、上述の式 ( 1 ) に限られず、他の減少関数でも良い。

【 0 0 2 6 】

図 1 に戻り、選択部 3 は、生成部 1 により生成されたラベル情報に含まれる重みに基づいて、適応データとして使用する音声データ ( 発話 ) を選択する。適応データの中で同じ内容の発話が複数存在する場合、適応によって、同じ内容の発話の事後確率が高くなって、当該発話の認識がしやすくなる。一方、この場合、他の発話の事後確率が低くなるため、他の発話を認識しにくくなる。

20

【 0 0 2 7 】

したがって、選択部 3 は、各発話の重みと重み閾値とを比較し、重み閾値より大きい発話を適応データとして選択する。これにより、適応データを使用して生成された適応音響モデル 1 0 3 を使用して、音声認識をする場合の悪影響を抑制することができる。

【 0 0 2 8 】

重み閾値は、例えば下記式 ( 2 ) により決定される。

【 0 0 2 9 】

$$= e^{-n} \cdot \dots (2)$$

30

【 0 0 3 0 】

ここで、 $\mu$  は重み閾値であり、 $\alpha$  は発話係数であり、 $n$  は全発話数である。つまり、同じ内容の発話については、カウント数  $x$  が全発話数  $n$  の  $\alpha$  倍より小さい場合 (  $x < \alpha n$  )、適応データとして選択される。

【 0 0 3 1 】

発話係数  $\alpha$  は、例えば 0 . 2 である。図 2 の例では、全発話数  $n$  は 6 であるため、重み閾値  $\mu$  は 0 . 8 2 になる。発話 - 1、発話 - 2、発話 - 3 及び発話 - 5 の重みは、重み閾値  $\mu$  より大きいため、選択部 3 により適応データとして選択される。一方、発話 - 4 及び発話 - 6 の重みは、重み閾値  $\mu$  より小さいため、選択部 3 により適応データとして選択されない。

40

【 0 0 3 2 】

なお、第 1 実施形態の説明では、発話係数が 0 . 2 の場合について説明したが、必要に応じて、発話係数を 1 以下の他の数値を設定しても良い。また、全発話数  $n$  の比率  $\alpha$  ではなく、絶対発話数 ( 全発話数  $n$  ) に基づいて、重み閾値  $\mu$  を決定してもよい。この場合、上述の式 ( 2 ) の  $\alpha n$  を  $n$  に変更すればよい。

【 0 0 3 3 】

適応部 4 は、選択部 3 により選択された適応データを用いて、ベース音響モデル 1 0 2 を適応させることにより、適応音響モデル 1 0 3 を生成する。具体的には、ベース音響モデル 1 0 2 の適応は、ベース音響モデル 1 0 2 のパラメータを、適応データを用いて最適

50

化することにより行われる。ベース音響モデル102を適応させる方法は、例えばDNN (Deep Neural Network)、CNN (Convolutional Neural Network)、及び、RNN (Recurrent Neural Network)などを使用する方法がある。適応音響モデル103は、音声認識装置10の外部の記憶部に記憶されるようにしてもよい。

#### 【0034】

##### [動作方法の例]

図3は第1実施形態の音声認識装置10の動作方法の例を示すフローチャートである。はじめに、生成部1が、言語モデル101とベース音響モデル102とを用いて音声データを認識する(ステップS1)。次に、生成部1が、ステップS1の処理により認識された音声データに含まれる発話を識別するラベルを生成する(ステップS2)。

10

#### 【0035】

次に、決定部2が、ラベルを用いて、同じ発話を含む音声データの個数を特定し、当該音声データに付与する重みを当該個数に応じて決定する(ステップS3)。次に、選択部3が、適応データとして使用する音声データを、重みに基づいて選択する(ステップS4)。次に、適応部4が、選択部3により選択された音声データ(適応データ)を用いて、ベース音響モデル102を適応させることにより、適応音響モデル103を生成する(ステップS5)。

#### 【0036】

以上説明したように、第1実施形態の音声認識装置10では、生成部1が、言語モデル101とベース音響モデル102(第1音響モデル)とを用いて音声データを認識し、当該音声データに含まれる発話を識別するラベルを生成する。決定部2が、ラベルを用いて、同じ発話を含む音声データの個数を特定し、当該音声データに付与する重みを当該個数に応じて決定する。選択部3が、重みに基づいて音声データを選択する。そして、適応部4が、選択部3により選択された音声データ(適応データ)を用いて、ベース音響モデル102(第1音響モデル)を適応させることにより、適応音響モデル103(第2音響モデル)を生成する。

20

#### 【0037】

これにより第1実施形態の音声認識装置10によれば、音響モデルを適応させた場合に生じる音声認識に与える悪影響を抑制することができる。

30

#### 【0038】

##### (第2実施形態)

次に第2実施形態について説明する。第2実施形態の説明では、第1実施形態と同様の説明については省略する。

#### 【0039】

適応データに含まれる非音声の部分が多いほど、当該適応データを用いた適応によって、非音声の確率が高くなる(音声の確率が低くなる)ため、音声の認識結果が非音声になることが多くなる。一方、適応データに含まれる非音声の部分が少ないほど、当該適応データを用いた適応によって、非音声の確率が低くなる(音声の確率が高くなる)ため、非音声の認識結果が音声になることが多くなる。

40

#### 【0040】

教師あり学習の場合、音声データから手動で発話ごとに切り出すため、非音声の部分のデータ量を制御できる。一方、教師なしの学習の場合、基本的にVAD (voice activity detection)等の音声区間検出処理により、自動的に発話を切り出すため、非音声の部分のデータ量の制御が困難である。

#### 【0041】

第2実施形態では、音声データに含まれる音声(または非音声)の部分が多い場合でも、適応の悪影響を抑制できる構成について説明する。

#### 【0042】

##### [機能構成の例]

50

図4は第2実施形態の音声認識装置10-2の機能構成の例を示すブロック図である。第2実施形態の音声認識装置10-2は、生成部1、選択部3-2、適応部4及び計算部5を備える。生成部1及び適応部4の説明は、第1実施形態と同様なので省略する。

【0043】

計算部5は、生成部1により生成されたラベルを用いて、音声データに含まれる音声フレームと、当該音声データに含まれる非音声フレームとの比率を計算する。

【0044】

図5は第2実施形態の音声データの例を示す図である。図5の例では、音声データに含まれるフレームの数が20である場合を示す。1、2、18、19及び20番目のフレームは、非音声フレームの一例である。なお、silは、silenceの略である。3～17番目のフレームは、音声フレームである。図5の音声データに含まれる発話は、「おはようございます」であり、当該発話のラベルも「おはようございます」である。

【0045】

計算部5は、フレームごとの音素を表すために、生成したラベルを用いてアライメントを行う。発音の長さによって二つ以上のフレームに、一つの音素が対応することもある。図5の例では、例えば、フレーム4及び5に対応する音素は同じになる。

【0046】

計算部5は、音声フレームと非音声フレームとの比率を計算する。図5の例では、音声フレームの比率は $15 / 20 = 0.75$ である。非音声フレームの比率は $5 / 20 = 0.25$ である。

【0047】

図4に戻り、選択部3-2は、音声フレームの比率が所定の選択範囲以内である音声データを、適応データとして選択する。所定の選択範囲は、例えば0.3～0.9である。図5の例では、音声フレームの比率は0.75であるので、当該音声フレームは選択部3-2により適応データとして選択される。

【0048】

所定の選択範囲は、適応の目的に応じて設定すれば良い。音声認識装置10-2から、できるだけ音声の認識結果を出力したい場合、所定の選択範囲として、値がより高い区間の範囲を使用する(例えば、0.4～1.0)。一方、音声データに背景雑音が入っているため、音声認識装置10-2から、できるだけ背景雑音の認識結果を出力したくない場合、所定の選択範囲として、値がより低い区間の範囲を使用する(例えば、0.0～0.5)。

【0049】

以上説明したように、第2実施形態の音声認識装置10-2によれば、例えば非音声フレームの比率が高い(例えば0.7以上)音声データが含まれている場合でも、選択部3-2により、当該音声データが選択されない。これにより、適応音声モデル103を使用した音声認識結果への悪影響を抑制できる。

【0050】

(第3実施形態)

次に第3実施形態について説明する。第3実施形態の説明では、第1及び第2実施形態と同様の説明については省略する。第3実施形態では、第1及び第2実施形態を組み合わせる場合の動作について説明する。

【0051】

[機能構成の例]

図6は第3実施形態の音声認識装置10-3の機能構成の例を示すブロック図である。第3実施形態の音声認識装置10-3は、生成部1、決定部2、選択部3-3、適応部4及び計算部5を備える。生成部1、決定部2及び適応部4の説明は、第1実施形態と同様なので省略する。計算部5の説明は、第2実施形態と同様なので省略する。

【0052】

第1実施形態の選択部3による適応データの選択方法を選択方法Aとし、第2実施形態

10

20

30

40

50

の選択部 3 - 2 による適応データの選択方法を選択方法 B とする。選択方法 A 及び B は独立である。そのため、選択方法 A 及び B の組み合わせによって、適応データとして使用する音声データを選択することが可能である。

【 0 0 5 3 】

選択部 3 - 3 は、決定部 2 により決定された重みと、計算部 5 により計算された音声フレームの比率とに基づいて、適応データとして使用する音声データを選択する。具体的には、選択部 3 - 3 は、例えば選択方法 A により適応データ候補を選択し、次に、適応データ候補から選択方法 B により適応データを選択する。また例えば、選択部 3 - 3 は、選択方法 B により適応データ候補を選択し、次に、適応データ候補から選択方法 A により適応データを選択する。

10

【 0 0 5 4 】

これにより第 3 実施形態の音声認識装置 1 0 - 3 によれば、第 1 及び第 2 実施形態の効果を得ることができる。

【 0 0 5 5 】

( 第 4 実施形態 )

次に第 4 実施形態について説明する。第 4 実施形態の説明では、第 1 実施形態と同様の説明については省略する。第 4 実施形態では、適応音響モデル 1 0 3 を使用して、音声認識をする構成について説明する。

【 0 0 5 6 】

[ 機能構成の例 ]

図 7 は第 4 実施形態の音声認識装置 1 0 - 4 の機能構成の例を示す図である。第 4 実施形態の音声認識装置 1 0 - 4 は、生成部 1、決定部 2、選択部 3、適応部 4 及び認識部 6 を備える。生成部 1、決定部 2、選択部 3 及び適応部 4 の説明は、第 1 実施形態と同様なので省略する。

20

【 0 0 5 7 】

認識部 6 は、言語モデル 1 0 1 及び適応音響モデル 1 0 3 を用いて、音声データの音声認識を行う。例えば、適応データが取得された環境と類似する環境で取得された音声データの音声認識をする場合、適応音響モデル 1 0 3 のパラメータは、ベース音響モデル 1 0 2 のパラメータより好ましい。また例えば、適応データに含まれる発話の話者と類似する話者（または同じ話者）の音声データの音声認識をする場合、適応音響モデル 1 0 3 のパラメータは、ベース音響モデル 1 0 2 のパラメータより好ましい。そのため、適応音響モデル 1 0 3 を用いて音声認識を行う場合、より高い音声認識精度が得られる。

30

【 0 0 5 8 】

( 第 5 実施形態 )

次に第 5 実施形態について説明する。第 5 実施形態の説明では、第 1 実施形態と同様の説明については省略する。第 1 実施形態では、言語モデル 1 0 1 及びベース音響モデル 1 0 2 の 2 種類のモデルを用いて適応を行っていた。第 5 実施形態では、言語モデル 1 0 1 及びベース音響モデル 1 0 2 を区別せずに、End - t o - End の音声認識方法により、適応を行う場合の構成について説明する。

【 0 0 5 9 】

[ 機能構成の例 ]

図 8 は第 5 実施形態の音声認識装置 1 0 - 5 の機能構成の例を示す図である。第 5 実施形態の音声認識装置 1 0 - 5 は、生成部 1 - 2、決定部 2、選択部 3 及び適応部 4 - 2 を備える。決定部 2 及び選択部 3 の説明は、第 1 実施形態と同様なので省略する。

40

【 0 0 6 0 】

第 5 実施形態の音声認識装置 1 0 - 5 は、音声認識ベースモデル 1 0 4 及び音声認識適応モデル 1 0 5 を記憶する。音声認識ベースモデル 1 0 4 は、音声の言語的な特徴、及び、音声の音響的な特徴の両方を区別せずにモデル化したデータである。

【 0 0 6 1 】

生成部 1 - 2 は、音声認識ベースモデル 1 0 4 を用いて音声データを認識し、ラベルを

50

生成する。第5実施形態では、音声認識ベースモデル104が、言語モデル101及びベース音響モデル102の役割を果たす。ラベルの生成方法の説明は、第1実施形態と同じなので省略する。

【0062】

適応部4-2は、選択部3により選択された適応データを用いて、音声認識ベースモデル104を適応させることにより、音声認識適応モデル105を生成する。具体的には、音声認識ベースモデル104の適応は、音声認識ベースモデル104のパラメータを、適応データを用いて最適化することにより行われる。音声認識ベースモデル104を適応させる方法は、例えばDNN、CNN及びRNN(Recurrent Neural Network)などを使用する方法がある。音声認識適応モデル105は、音声認識装置10の外部の記憶部に記憶されるようにしてもよい。

10

【0063】

最後に、第1乃至第5実施形態の音声認識装置10(10-2, 10-3, 10-4, 10-5)のハードウェア構成の例について説明する。

【0064】

[ハードウェア構成の例]

図9は第1乃至第5実施形態の音声認識装置10(10-2, 10-3, 10-4, 10-5)のハードウェア構成の例を示す図である。以下では、第1実施形態の音声認識装置10の場合を例にして説明する。なお、第2乃至第5実施形態の音声認識装置10-2(10-3, 10-4, 10-5)のハードウェア構成も、第1実施形態の音声認識装置10のハードウェア構成と同様である。

20

【0065】

第1実施形態の音声認識装置10は、制御装置301、主記憶装置302、補助記憶装置303、表示装置304、入力装置305及び通信装置306を備える。制御装置301、主記憶装置302、補助記憶装置303、表示装置304、入力装置305及び通信装置306は、バス310を介して接続されている。

【0066】

制御装置301は、補助記憶装置303から主記憶装置302に読み出されたプログラムを実行する。主記憶装置302は、ROM(Read Only Memory)、及び、RAM(Random Access Memory)等のメモリである。補助記憶装置303は、HDD(Hard Disk Drive)、及び、メモリカード等である。

30

【0067】

表示装置304は表示情報を表示する。表示装置304は、例えば液晶ディスプレイ等である。入力装置305は、音声認識装置10を操作するためのインタフェースである。入力装置305は、例えばキーボードやマウス等である。音声認識装置10がスマートフォン及びタブレット型端末等のスマートデバイスの場合、表示装置304及び入力装置305は、例えばタッチパネルである。通信装置306は、他の装置と通信するためのインタフェースである。

【0068】

第1実施形態の音声認識装置10で実行されるプログラムは、インストール可能な形式又は実行可能な形式のファイルでCD-ROM、メモリカード、CD-R及びDVD(Digital Versatile Disc)等のコンピュータで読み取り可能な記憶媒体に記録されてコンピュータ・プログラム・プロダクトとして提供される。

40

【0069】

また第1実施形態の音声認識装置10で実行されるプログラムを、インターネット等のネットワークに接続されたコンピュータ上に格納し、ネットワーク経由でダウンロードさせることにより提供するように構成してもよい。また第1実施形態の音声認識装置10で実行されるプログラムをダウンロードさせずにインターネット等のネットワーク経由で提供するように構成してもよい。

50

## 【 0 0 7 0 】

また第1実施形態の音声認識装置10のプログラムを、ROM等に予め組み込んで提供するように構成してもよい。

## 【 0 0 7 1 】

第1実施形態の音声認識装置10で実行されるプログラムは、上述の機能ブロックのうち、プログラムによっても実現可能な機能ブロックを含むモジュール構成となっている。当該各機能ブロックは、実際のハードウェアとしては、制御装置301が記憶媒体からプログラムを読み出して実行することにより、上記各機能ブロックが主記憶装置302上にロードされる。すなわち上記各機能ブロックは主記憶装置302上に生成される。

## 【 0 0 7 2 】

なお上述した各機能ブロックの一部又は全部をソフトウェアにより実現せずに、IC (Integrated Circuit)等のハードウェアにより実現してもよい。

## 【 0 0 7 3 】

また複数のプロセッサを用いて各機能を実現する場合、各プロセッサは、各機能のうち1つを実現してもよいし、各機能のうち2以上を実現してもよい。

## 【 0 0 7 4 】

また第1実施形態の音声認識装置10の動作形態は任意でよい。第1実施形態の音声認識装置10を、例えばネットワーク上のクラウドシステムとして動作させてもよい。

## 【 0 0 7 5 】

本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

## 【 符号の説明 】

## 【 0 0 7 6 】

- 1 生成部
- 2 決定部
- 3 選択部
- 4 適応部
- 5 計算部
- 6 認識部
- 101 言語モデル
- 102 ベース音響モデル
- 103 適応音響モデル
- 104 音声認識ベースモデル
- 105 音声認識適応モデル
- 301 制御装置
- 302 主記憶装置
- 303 補助記憶装置
- 304 表示装置
- 305 入力装置
- 306 通信装置
- 310 バス

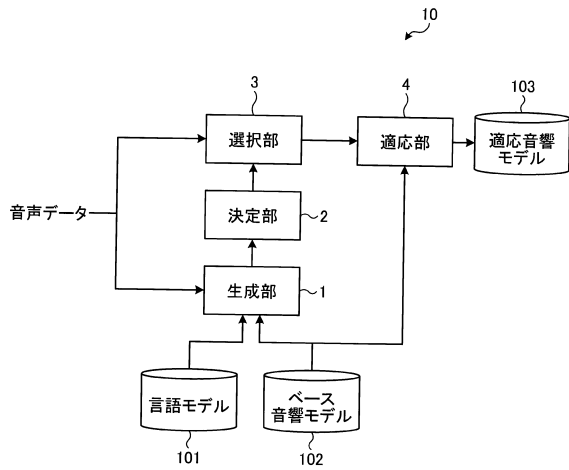
10

20

30

40

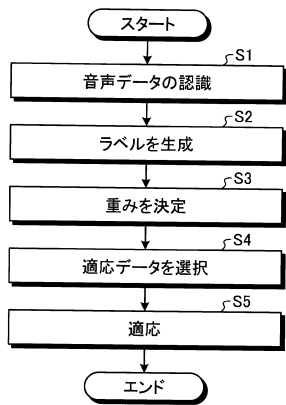
【図1】



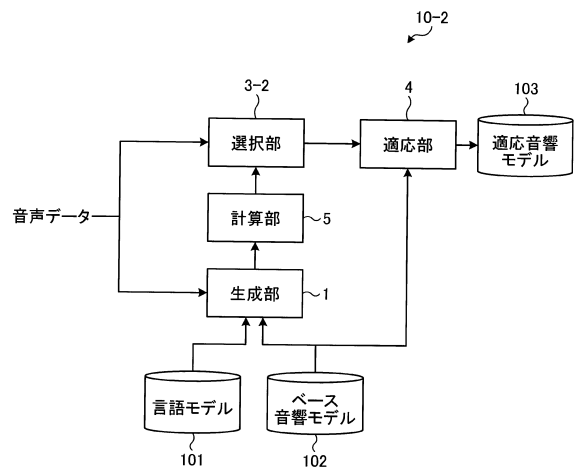
【図2】

音声データ	音声認識結果	ラベル	カウント数	重み
発話-1	あのー	あのー	1	1.00
発話-2	先週の進捗ですが	せんしゅうのしんちよくですが	1	1.00
発話-3	あのー	あのー	2	0.37
発話-4	いろんなトランプがあったので	いろんなとらぶらがあつたので	1	1.00
発話-5	あのー	あのー	3	0.14
発話-6	うまく進んでない	うまくすすんでない	1	1.00

【図3】



【図4】





---

フロントページの続き

(56)参考文献 特開2014-102345(JP,A)  
特開2002-156994(JP,A)  
特開2000-305589(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G10L 15/00-15/34  
IEEE Xplore