

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6845489号
(P6845489)

(45) 発行日 令和3年3月17日(2021.3.17)

(24) 登録日 令和3年3月2日(2021.3.2)

(51) Int.Cl. F I
G 1 0 L 17/00 (2013.01) G 1 0 L 17/00 2 0 0 C

請求項の数 10 (全 16 頁)

<p>(21) 出願番号 特願2019-504164 (P2019-504164) (86) (22) 出願日 平成29年3月7日(2017.3.7) (86) 国際出願番号 PCT/JP2017/008979 (87) 国際公開番号 W02018/163279 (87) 国際公開日 平成30年9月13日(2018.9.13) 審査請求日 令和1年5月21日(2019.5.21)</p>	<p>(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 (74) 代理人 100077838 弁理士 池田 憲保 (74) 代理人 100129023 弁理士 佐々木 敬 (72) 発明者 山本 仁 東京都港区芝五丁目7番1号 日本電気株式会社社内 (72) 発明者 越仲 孝文 東京都港区芝五丁目7番1号 日本電気株式会社社内 審査官 大野 弘</p>
--	--

最終頁に続く

(54) 【発明の名称】 音声処理装置、音声処理方法、および音声処理プログラム

(57) 【特許請求の範囲】

【請求項 1】

音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出する貢献度推定手段と、

前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する話者特徴算出手段とを備える、音声処理装置。

【請求項 2】

前記音声信号に含まれる音の種類比率を表す音声統計量を算出する音声統計量算出手段をさらに備え、

前記話者特徴算出手段は、前記音声信号の前記音声統計量と、前記音声信号の前記貢献度とに基づいて、前記認識特徴量を算出する、請求項 1 に記載の音声処理装置。

【請求項 3】

前記貢献度推定手段は、前記音声信号の前記貢献度として、

前記音声信号の一部が音声か否かを識別して算出した音声らしさを表す値、前記音声信号の一部が話者認識に正解する音声か否かを識別して算出した話者認識の正解しやすさを表す値、前記音声信号の一部が話者認識誤りを起こす音声か否かを識別して算出した話者認識の誤りやすさを表す値の少なくともいずれかひとつを算出する、請求項 1 または 2 に記載の音声処理装置。

【請求項 4】

前記貢献度推定手段は、
ニューラルネットワークを用いて前記音声信号の前記貢献度を算出する、請求項 3 に記載の音声処理装置。

【請求項 5】

前記話者特徴算出手段は、
前記認識特徴量として i-vector を算出する、請求項 3 または 4 に記載の音声処理装置。

【請求項 6】

前記認識特徴量に基づいて前記属性情報を認識する属性認識手段を備える、請求項 1 ~ 5 のいずれか 1 項に記載の音声処理装置。

10

【請求項 7】

前記特定の属性情報は、
音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、音声信号から推定される話者の性格の少なくともいずれか 1 つを表す情報である、請求項 1 ~ 6 のいずれか 1 項に記載の音声処理装置。

【請求項 8】

音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出し、

前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する、音声処理方法。

20

【請求項 9】

前記音声信号に含まれる音の種類比率を表す音声統計量をさらに算出し、
前記音声信号の前記音声統計量と、前記音声信号の前記貢献度とに基づいて、前記認識特徴量を算出する、請求項 8 に記載の音声処理方法。

【請求項 10】

コンピュータに、
音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出する処理と、

前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する処理とを実行させる、音声処理プログラム。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声処理装置、音声処理方法、および音声処理プログラムに関する。

【背景技術】

【0002】

音声信号から、音声を発した話者を特定するための個人性を表す話者特徴を算出する音声処理装置が知られている。また、この話者特徴を用いて、音声信号を発した話者を推定する話者認識装置が知られている。

40

【0003】

この種の音声処理装置を用いる話者認識装置は、話者を特定するために、第 1 の音声信号から抽出した第 1 の話者特徴と、第 2 の音声信号から抽出した第 2 の話者特徴との類似度を評価する。そして、話者認識装置は、類似度の評価結果に基づいて 2 つの音声信号の話者が同一か否かを判定する。

【0004】

非特許文献 1 には、音声信号から話者特徴を抽出する技術が記載されている。非特許文献 1 に記載の話者特徴抽出技術は、音声モデルを用いて音声信号の音声統計量を算出する

50

。そして、非特許文献1に記載の話者特徴抽出技術は、因子分析技術に基づいてその音声統計量を処理し、所定の要素数で表現される話者特徴ベクトルとして算出する。すなわち、非特許文献1においては、話者特徴ベクトルを話者の個人性を表す話者特徴として利用する。

【先行技術文献】

【非特許文献】

【0005】

【非特許文献1】Najim Dehak, Patrick Kenny, Reda Dehak, Pierre Dumouchel, and Pierre Ouellet, "Front-End Factor Analysis for Speaker Verification," IEEE Transaction on Audio, Speech and Language Processing, Vol. 19, No. 4, pp. 788-798, 2011.

10

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかしながら、非特許文献1に記載の技術には、抽出した話者特徴を用いる話者認識の精度が十分でないという問題があった。

【0007】

非特許文献1に記載の技術は、話者特徴抽出装置に入力された音声信号に対して所定の統計処理を行い、話者特徴ベクトルを算出する。具体的には、非特許文献1に記載の技術は、話者特徴抽出装置に入力された音声信号の全体に対して一律の統計処理を行うことにより、話者特徴ベクトルを算出している。そのため、非特許文献1に記載の技術は、音声信号の部分区間に、話者の個人性を算出する元として適切ではない信号が含まれている場合であっても、音声信号の全体から話者特徴ベクトルを算出してしまうので、話者認識の精度を損なうおそれがある。具体的には、音声信号の部分区間に、例えば、話者の不明瞭な発声、話者の咳や笑い声などの話し声とは異なる音、雑音などが混入している場合に、話者認識の精度を損なうおそれがある。

20

【0008】

本発明は、上記問題に鑑みてなされたものであり、その目的は、話者認識の精度をより高めた音声処理装置、音声処理方法、および音声処理プログラムを提供することにある。

【課題を解決するための手段】

30

【0009】

本発明の第1の態様の音声処理装置は、音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出する貢献度推定手段と、前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する話者特徴算出手段とを備える。

【0010】

本発明の第2の態様の音声処理方法は、音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出し、前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する。

40

【0011】

本発明の第3の態様の音声処理プログラムは、コンピュータに、音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出する処理と、前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する処理とを実行させる。

【発明の効果】

【0012】

本発明によれば、話者認識の精度をより高めた音声処理装置、音声処理方法、およびプログラムを提供することができる。

【図面の簡単な説明】

50

【 0 0 1 3 】

【 図 1 】 本発明の第 1 の実施形態に係る音声処理装置の構成を示すブロック図である。

【 図 2 】 本発明の第 1 の実施形態に係る音声処理装置の動作の流れを示すフローチャートである。

【 図 3 】 本発明の第 2 の実施形態に係る音声処理装置の構成を示すブロック図である。

【 図 4 】 本発明の第 2 の実施形態に係る音声処理装置の動作の流れを示すフローチャートである。

【 図 5 】 本発明の第 3 の実施形態に係る音声処理装置の構成を示すブロック図である。

【 図 6 】 本発明のその他の実施形態に係る音声処理装置の構成を示すブロック図である。

【 発明を実施するための形態 】

10

【 0 0 1 4 】

以下、音声処理装置等および話者特徴抽出装置の実施形態について、図面を参照して説明する。なお、実施形態において同じ符号を付した構成要素は同様の動作を行うので、再度の説明を省略する場合がある。

【 0 0 1 5 】

< 第 1 の実施形態 >

図 1 は、本発明の第 1 の実施形態に係る音声処理装置の構成を示すブロック図である。

【 0 0 1 6 】

音声処理装置 1 0 0 は、貢献度推定部 1 1 と、話者特徴算出部 1 2 とを備える。

【 0 0 1 7 】

20

貢献度推定部 1 1 は、外部から音声を表す音声信号を受け取る。また、貢献度推定部 1 1 は、受けた音声信号に基づき、その音声信号の部分区間の品質の程度を数値で表した貢献度を算出する。

【 0 0 1 8 】

話者特徴算出部 1 2 は、貢献度推定部 1 1 が算出した音声信号の部分区間の貢献度を、その部分区間の重みとして用いて、音声信号から特定の属性情報を認識するための認識特徴量を算出する。

【 0 0 1 9 】

ここで、特定の属性情報とは、音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、および音声信号から推定される話者の性格等を示す情報である。

30

【 0 0 2 0 】

図 2 を参照し、音声処理装置 1 0 0 の動作の流れについて説明する。図 2 は本発明の第 1 の実施形態に係る音声処理装置の動作の流れを示すフローチャートである。

【 0 0 2 1 】

まず、貢献度推定部 1 1 は、外部から受けた音声信号に基づいて、音声信号の部分区間の貢献度を算出する（ステップ S 1 0 1）。次いで、貢献度推定部 1 1 は、算出した音声信号の部分区間の貢献度を話者特徴算出部 1 2 に出力する。

【 0 0 2 2 】

次いで、話者特徴算出部 1 2 は、貢献度推定部 1 1 から受けた貢献度に基づいて、認識特徴量を算出する（ステップ S 1 0 2）。

40

【 0 0 2 3 】

< 第 2 の実施形態 >

図 3 は、第 2 の実施形態における音声処理装置 2 0 0 のブロック図である。音声処理装置 2 0 0 は、貢献度推定部 1 1、話者特徴算出部 1 2、音声区間検出部 2 1、および音声統計量算出部 2 2 を備える。また、音声処理装置 2 0 0 は、さらに、貢献度記憶部 2 3 および貢献度学習部 2 4 を備えてもよい。

【 0 0 2 4 】

音声区間検出部 2 1 は、外部から音声信号を受け取る。また、音声区間検出部 2 1 は、受け取った音声信号に含まれる音声区間を検出して区分化する。この時、音声区間検出部

50

21は、音声信号を一定の長さに区分化してもよいし、異なる長さに区分化してもよい。例えば、音声区間検出部21は、音声信号のうち音量が一定時間継続して所定値より小さい区間を無音と判定し、その区間の前後を異なる音声区間と判定して区分化してもよい。そして、音声区間検出部21は、区分化した結果（音声区間検出部21の処理結果）である区分化音声信号を、貢献度推定部11および音声統計量算出部22に出力する。ここで、音声信号の受け取りとは、例えば、外部の装置または他の処理装置からの音声信号の受信、または他のプログラムからの音声信号処理の処理結果の引き渡しのことである。また、出力とは、例えば、外部の装置や他の処理装置への送信、または他のプログラムへの音声区間検出部21の処理結果の引き渡しのことである。

【0025】

音声統計量算出部22は、音声区間検出部21から区分化音声信号を受け取る。音声統計量算出部22は、受け取った区分化音声信号に基づいて、該区分化音声信号に含まれる音の種類を表す音声統計量を算出する。ここで、音の種類とは、例えば、言語により定まる音素や単語、音声信号を類似度に基づいてクラスタリングして得られる音のグループである。そして、音声統計量算出部22は、音声統計量を話者特徴算出部12に出力する。以降、ある音声信号に対して算出された音声統計量を、該音声信号の音声統計量と呼ぶ。

【0026】

音声統計量算出部22が、音声統計量を算出する方法の一例について説明する。具体的には、音声統計量算出部22は、音声区間検出部21から受け取った区分化音声信号に基づいて、該区分化音声信号を周波数分析処理した計算結果で表現される音響特徴を算出し、算出した結果を出力する。例えば、音声統計量算出部22は、音声区間検出部21から受け取った区分化音声信号を、短時間フレーム時系列に変換する。そして、音声統計量算出部22は、短時間フレーム時系列のそれぞれのフレームを周波数分析し、その処理結果を音響特徴として出力する。この場合、音声統計量算出部22は、例えば、短時間フレーム時系列として、25ミリ秒区間のフレームを10ミリ秒ごとに生成する。音声統計量算出部22は、例えば、周波数分析結果である音響特徴として、高速フーリエ変換処理（Fast Fourier Transform; FFT）およびフィルタバンク処理によって得られた周波数フィルタバンク特徴や、さらに加えて離散コサイン変換処理を施して得られたメル周波数ケプストラム係数（Mel-Frequency Cepstrum Coefficients; MFCC）特徴などを算出する。

【0027】

そして、音声統計量算出部22は、音響特徴の時系列と、音響特徴と音の種類との対応関係を格納する音声モデルを用いて、音の種類を表す数値情報の時系列を算出する。音声統計量算出部22は、例えば、音声モデルがガウス混合モデル（Gaussian Mixture Model; GMM）である場合、ガウス混合モデルが有する各要素分布の平均、分散、および混合係数に基づいて、各要素分布の事後確率を算出する。ここで、各要素分布の事後確率は、音声信号に含まれる音の種類それぞれの出現度である。また、音声統計量算出部22は、例えば、音声モデルがニューラルネットワーク（Neural Network）である場合、音響特徴と、ニューラルネットワークが有する重み係数に基づいて、音声信号に含まれる音の種類の出現度を算出する。

【0028】

貢献度記憶部23は、1つ以上の貢献度推定器を記憶する。貢献度推定器は、音声信号を信号の品質によって複数の種類に仕分けるよう動作するように構成されるものである。貢献度推定器は、例えば、音声信号の品質を表す数値情報を出力する。信号の品質の種類とは、例えば、音声・非音声・無音である。また、信号の品質の種類とは、例えば、話者認識に正解する音声・話者認識に誤りを起こす音声である。

【0029】

具体的には、貢献度記憶部23は、貢献度推定器が保有するパラメタを記憶する。貢献度記憶部23は、例えば、貢献度推定器がニューラルネットワークである場合、それを構成するノードの数やノード間の接続重み係数などの一式をパラメタとして記憶する。

10

20

30

40

50

【0030】

なお、図3では、貢献度記憶部23が音声処理装置200内に内蔵されることを例に説明を行ったが、本発明はこれに限定されるものではない。貢献度記憶部23は、音声処理装置200の外部に設けられた記憶装置で実現されるものであってもよい。

【0031】

貢献度推定部11は、音声区間検出部21から区分化音声信号を受け取る。貢献度推定部11は、貢献度記憶部23に記憶されている貢献度推定器を用いて、区分化音声信号の品質を表す数値情報を算出する。貢献度推定部11は、音声統計量算出部22と同様に、区分化音声信号を短時間フレーム時系列に変換し、それぞれのフレームの音響特徴を算出し、音響特徴の時系列を算出する。続いて、貢献度推定部11は、各フレームの音響特徴と貢献度推定器のパラメタとを用いて、各フレームの品質を表す数値を算出する。以降、ある音声信号に対して算出された信号の品質を表す数値のことを音声信号の貢献度と呼ぶ。

10

【0032】

具体的には、貢献度推定部11は、例えば、貢献度推定器がニューラルネットワークである場合、音響特徴と、ニューラルネットワークが有する重み係数とに基づいて、音響特徴の貢献度を算出する。例えば、貢献度推定器がニューラルネットワークであり、その出力層が、2つの信号の品質の種類「話者認識に正解する信号」と「話者認識誤りを起こす信号」とに相当するものであるとする。このとき、貢献度推定器は、音響特徴が話者認識に正解する信号である確率と、音響特徴が話者認識誤りを起こす信号である確率とを算出し、貢献度として、例えば、「話者認識に正解する信号」である確率を出力する。また、貢献度推定部11は、話者認識を実行する前に、音声信号の部分区間が音声か否かを識別して音声である確率を算出してもよい。

20

【0033】

話者特徴算出部12は、音声統計量算出部22が出力した音声統計量および貢献度推定部11が出力した貢献度を受け取る。話者特徴算出部12は、音声統計量および貢献度を用いて、音声信号から特定の属性情報を認識するための認識特徴量を算出する。

【0034】

話者特徴算出部12が音声信号xの認識特徴量としてi-vectorに基づく特徴ベクトルF(x)を算出する方法の一例について説明する。なお、話者特徴算出部12が算出する特徴ベクトルF(x)は、音声信号xに対して所定の演算を施して算出できるベクトルであればよく、i-vectorはその一例である。

30

【0035】

話者特徴算出部12は、音声統計量算出部22から、音声信号xの統計量の情報として、例えば、短時間フレームごとに算出された音響事後確率 $P_t(x)$ および音響特徴 $A_t(x)$ ($t = \{1 \dots T\}$ 、 T は1以上の自然数)とを受け取る。また、話者特徴算出部12は、貢献度推定部11から、音声信号xの貢献度の情報として、例えば、短時間フレームごとに算出された貢献度 $C_t(x)$ を受け取る。話者特徴算出部12は、以下の式(1)のように、音響事後確率 $P_t(x)$ の各要素に対して、貢献度 $C_t(x)$ をかけて、その結果を $Q_t(x)$ として算出する。

40

【0036】

【数1】

$$Q_{t,c}(x) = C_t(x)P_{t,c}(x) \cdot \cdot \cdot \cdot (1)$$

【0037】

話者特徴算出部12は、貢献度によって重みづけされた音響事後確率 $Q_t(x)$ および音響特徴 $A_t(x)$ を用いて、以下の式(2)に基づいて音声信号xの0次統計量 $S_0(x)$ を算出し、式(3)に基づいて1次統計量 $S_1(x)$ を算出する。

【0038】

【数 2】

$$S_0(x) = \begin{pmatrix} S_{0,1}I_D & \cdots & 0_D \\ \vdots & \ddots & \vdots \\ 0_D & \cdots & S_{0,c}I_D \end{pmatrix}, \quad S_{0,c} = \sum_{t=1}^T Q_{t,c}(x) \cdot \cdots \cdot (2)$$

【0039】

【数 3】

$$S_1(x) = (S_{1,1}, S_{1,2}, \dots, S_{1,c})^T, \quad S_{1,c} = \sum_{t=1}^T Q_{t,c}(x)(A_t(x) - m_c) \cdot \cdots \cdot (3) \quad 10$$

【0040】

話者特徴算出部 12 は、続いて、以下の式 (4) に基づいて音声信号 x の i -vector である $F(x)$ を算出する。

【0041】

【数 4】

$$F(x) = (I + T^T \Sigma^{-1} S_0(x) T)^{-1} T^T \Sigma^{-1} S_1(x) \cdot \cdots \cdot (4)$$

【0042】

20

式 (1) ~ 式 (4) において、 C は統計量 $S_0(x)$ および $S_1(x)$ の要素数、 D は音響特徴 $A_t(x)$ の要素数 (次元数)、 m_c は音響特徴空間における c 番目の領域の音響特徴の平均ベクトル、 I は単位行列、 0 は零行列を表す。 T は i -vector 計算用のパラメタであり、 Σ は音響特徴空間における音響特徴の共分散行列である。

【0043】

話者特徴算出部 12 が上述の手順で特徴ベクトル $F(x)$ を算出する際に、音声信号 x のすべての時刻 t ($t = \{1 \dots T\}$ 、 T は 1 以上の自然数) において、その貢献度 $C_t(x)$ が 1 であれば、非特許文献 1 に記載の i -vector 算出手順と等価である。本実施形態において、話者特徴算出部 12 は、貢献度推定部 11 が音声信号 x の時刻 t に応じて推定した貢献度 $C_t(x)$ を用いることにより、非特許文献 1 に記載の i -vector とは異なる特徴ベクトル $F(x)$ を算出できる。 30

【0044】

このように、音声処理装置 200 において、話者特徴算出部 12 が、音声信号 x に対して、該音声信号の各部分区間の品質に応じた貢献度 $C_t(x)$ を用いて特徴ベクトル $F(x)$ を算出することにより、音声信号の品質に応じた特徴ベクトルを出力することができる。

【0045】

貢献度学習部 24 は、訓練用音声信号を用いて貢献度記憶部 23 が記憶できる貢献度推定器を学習する。貢献度学習部 24 は、例えば、貢献度推定器がニューラルネットワークである場合、それを構成するノード間の接続重み係数などのパラメタを、一般的な最適化基準に従って最適化する。貢献度学習部 24 が使用する訓練用音声信号は、複数の音声信号を集めたものであり、それぞれの音声信号は、貢献度推定部 11 が出力する信号の品質の種類の内いずれかと対応付けられたものである。 40

【0046】

以下では、入力が音響特徴であり、出力が「話者認識に正解する音声」および「話者認識に誤りを起こす音声」の 2 種類の信号の品質である貢献度推定器を貢献度学習部 24 が学習する方法の一例を説明する。

【0047】

(a) まず、貢献度学習部 24 は、話者ラベル付きの複数の音声信号を用いて、音声信号の話者ラベルを識別することのできる識別器を学習する。(b) 次に、貢献度学習部 2 50

4 は、話者ラベル付きの複数の音声信号のそれぞれを、短時間フレームごとに算出した音響特徴の時系列に変換し、(a)で学習した識別器を用いて、各フレームの話者ラベルを識別する。(c)次に、貢献度学習部24は、識別された各フレームの話者ラベルのうち、事前に付与された話者ラベルと、識別器が識別した話者ラベルが同一であるフレームを「話者認識に正解する音声」、そうでないフレームを「話者認識に誤りを起こす音声」とする。(d)そして、貢献度学習部24は、「話者認識に正解する音声」および「話者認識に誤りを起こす音声」を訓練用音声信号として、貢献度推定器を学習する。

【0048】

以上述べたように、本実施形態に係る音声処理装置200において、貢献度推定部11は、音声信号の部分区間に応じた品質を表す指標として、音声信号の貢献度を算出できる。また、話者特徴算出部12は、音声信号の音響統計量と貢献度とに基づいて特徴ベクトルを算出する。これにより、音声信号に対して、音声信号の各部分区間の品質を考慮した特徴ベクトルを出力できる。すなわち、本実施形態にかかる音声処理装置200は、話者認識の精度を高めるのに適した話者特徴を算出できる。

10

【0049】

なお、本実施形態に係る音声処理装置200における貢献度記憶部23は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。また、貢献度記憶部23に貢献度推定器が記憶される過程は特に限定されない。例えば、記録媒体を介して貢献度推定器が貢献度記憶部23に記憶されてもよいし、通信回線等を介して送信された貢献度推定器が貢献度記憶部23に記憶されてもよい。または、入力デバイスを介して入力された貢献度推定器が貢献度記憶部23で記憶されてもよい。

20

【0050】

(第2の実施形態の動作)

次に、第2の実施形態における音声処理装置200の動作について、図4のフローチャートを用いて説明する。図4は、音声処理装置200の動作の一例を示すフローチャートである。

【0051】

音声処理装置200は、外部から1つ以上の音声信号を受け取り、音声区間検出部21に提供する。具体的には、音声区間検出部21は、受け取った音声信号を区分化し、区分化音声信号を貢献度推定部11および音声統計量算出部22に出力する(ステップS201)。

30

【0052】

音声統計量算出部22は、受け取った1つ以上の区分化音声信号それぞれについて、短時間フレーム分析処理を行い、音響特徴と音声統計量の時系列を算出する(ステップS202)。

【0053】

貢献度推定部11は、受け取った1つ以上の区分化音声信号のそれぞれについて、短時間分析フレーム処理を行い、貢献度の時系列を算出する(ステップS203)。

【0054】

話者特徴算出部12は、受け取った1つ以上の音響特徴・音声統計量・貢献度の時系列に基づいて、話者認識特徴量を算出して出力する。(ステップS204)。音声処理装置200は、外部からの音声信号の受理が終了したら、一連の処理を終了する。

40

【0055】

(第2の実施形態の効果)

以上、説明したように、本実施形態にかかる音声処理装置200によれば、音声処理装置200が算出した話者特徴を用いる話者認識の精度を高めることができる。なぜならば、音声処理装置200は、貢献度推定部11が音声信号の品質を貢献度として算出し、話者特徴算出部12が貢献度を考慮した特徴ベクトルを算出することで、音声信号の品質の高い部分区間に重きを置いた特徴ベクトルを出力するからである。

【0056】

50

このように、本実施形態に係る音声処理装置200は、音声信号に対して、各部分区間の品質に応じた貢献度を考慮した特徴ベクトルを算出する。これにより、音声信号の部分区間に、話者の不明瞭な発声、話者の咳や笑い声などの話し声とは異なる音、雑音などが混入している場合にも、話者認識に適した認識特徴量を求めることができる。

【0057】

<第3の実施形態>

図5は、本発明の第3の実施形態に係る、音声処理装置の構成の一例を示すブロック図である。

【0058】

図5に示すように、音声処理装置300は、貢献度推定部11と、話者特徴算出部12と、属性認識部13とを備える。音声処理装置300は、属性情報を認識することのできる音声処理装置である。

【0059】

貢献度推定部11および話者特徴算出部12については、第1および第2の実施形態と同様なので説明は省略する。

【0060】

属性認識部13は、話者特徴算出部12から属性情報を認識するための認識特徴量を受け取る。属性認識部13は、認識特徴量に基づいて、音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、音声信号から推定される発話者の性格等を認識する。具体的には、属性認識部13は、例えば、認識特徴量を比較するための比較用音声データを格納する記憶装置(図示しない)を参照する。この場合、属性認識部13は、認識特徴量と、比較用音声データの類似の度合い等を算出することで、属性情報を認識することができる。

【0061】

<第3の実施形態の具体例>

次に、本発明の第3の実施形態に係る音声処理装置300の具体的な応用例について説明する。

【0062】

本発明の第3の実施形態に係る音声処理装置300が算出した話者特徴は、音声信号の話者を推定する話者認識に利用可能である。例えば、第1の音声信号から算出した第1の話者特徴と、第2の音声信号から算出した第2の話者特徴とから、2つの話者特徴の類似性を現す指標として、コサイン類似度を算出する。例えば、話者照合することを目的とする場合は、前記の類似度に基づく照合可否の判定情報を出力してもよい。また、話者識別することを目的とする場合は、第1の音声信号に対して複数の第2の音声信号を用意して各々の類似度を求め、値の大きい組を出力してもよい。

【0063】

本発明の第3の実施形態に係る音声処理装置300は、音声信号から特定の属性情報を認識するための認識特徴量を算出する特徴算出装置の一例である。音声処理装置300は、特定の属性が音声信号を発した話者であるとき、話者特徴抽出装置として利用可能である。また、音声処理装置300は、例えば文発話の音声信号に対して、当該話者特徴を用いて推定した話者情報に基づいて、当該話者の話し方の特徴に適応化する機構を備える音声認識装置の一部としても利用可能である。また、ここで、話者を示す情報は、話者の性別を示す情報や、話者の年齢あるいは年齢層を示す情報であってもよい。

【0064】

本発明の第3の実施形態に係る音声処理装置300は、特定の属性を音声信号が伝える言語(音声信号を構成する言語)を示す情報とするとき、言語特徴算出装置として利用可能である。また、音声処理装置300は、例えば文発話の音声信号に対して、当該言語特徴を用いて推定した言語情報に基づいて、翻訳する言語を選択する機構を備える音声翻訳装置の一部としても利用可能である。

【0065】

10

20

30

40

50

本発明の第3の実施形態に係る音声処理装置300は、特定の属性が話者の発話時の感情を示す情報であるとき、感情特徴算出装置として利用可能である。また、音声処理装置300は、例えば蓄積された多数の発話の音声信号に対して、当該感情特徴を用いて推定した感情情報に基づいて、特定の感情に対応する音声信号を特定する機構を備える音声検索装置や音声表示装置の一部としても利用可能である。この感情情報には、例えば、感情表現を示す情報、発話者の性格を示す情報等が含まれる。

【0066】

以上のように、本実施形態における特定の属性情報は、音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、音声信号から推定される発話者の性格、の少なくともいずれか一つを表す情報である。

10

【0067】

(ハードウェア構成についての説明)

以上、実施形態を用いて本発明を説明したが、本発明は、上記実施形態に限定されるものではない。本発明の構成や詳細には、本発明の範囲内で当業者が理解しうる様々な変更をすることができる。すなわち、本発明は、以上の実施形態に限定されることなく、種々の変更が可能であり、それらも本発明の範囲内に包含されるものであることは言うまでもない。

【0068】

以上のように、本発明の一態様における音声処理装置等は、音声信号の品質を考慮した特徴ベクトルを抽出し話者認識の精度を高めることができるという効果を有しており、音声処理装置等および話者認識装置として有用である。なお、本発明において使用者に関する情報を取得、利用する場合は、これを適法に行うものとする。

20

【0069】

<その他の実施形態>

音声処理装置は、ハードウェアによって実現してもよいし、ソフトウェアによって実現してもよい。また、音声処理装置は、ハードウェアとソフトウェアの組み合わせによって実現してもよい。

【0070】

図6は、音声処理装置を構成する情報処理装置(コンピュータ)の一例を示すブロック図である。

30

【0071】

図6に示すように、情報処理装置400は、制御部(CPU: Central Processing Unit)410と、記憶部420と、ROM(Read Only Memory)430と、RAM(Random Access Memory)440と、通信インターフェース450と、ユーザインターフェース460とを備えている。

【0072】

制御部(CPU)410は、記憶部420またはROM430に格納されたプログラムをRAM440に展開して実行することで、音声処理装置および話者認識装置の各種の機能を実現することができる。また、制御部(CPU)410は、データ等を一時的に格納できる内部バッファを備えていてもよい。

40

【0073】

記憶部420は、各種のデータを保持できる大容量の記憶媒体であって、HDD(Hard Disc Drive)、およびSSD(Solid State Drive)等の記憶媒体で実現することができる。また、記憶部420は、情報処理装置400が通信インターフェース450を介して通信ネットワークと接続されている場合には、通信ネットワーク上に存在するクラウドストレージであってもよい。また、記憶部420は、制御部(CPU)410が読み取り可能なプログラムを保持していてもよい。

【0074】

ROM430は、記憶部420と比べると小容量なフラッシュメモリ等で構成できる不揮発性の記憶装置である。また、ROM430は、制御部(CPU)410が読み取り可

50

能なプログラムを保持していてもよい。なお、制御部（CPU）410が読み取り可能なプログラムは、記憶部420およびROM430の少なくとも一方が保持していればよい。

【0075】

なお、制御部（CPU）410が読み取り可能なプログラムは、コンピュータが読み取り可能な様々な記憶媒体に非一時的に格納した状態で、情報処理装置400に供給してもよい。このような記憶媒体は、例えば、磁気テープ、磁気ディスク、光磁気ディスク、CD-ROM、CD-R、CD-R/W、半導体メモリである。

【0076】

RAM440は、DRAM（Dynamic Random Access Memory）及びSRAM（Static Random Access Memory）等の半導体メモリであり、データ等を一時的に格納する内部バッファとして用いることができる。

10

【0077】

通信インターフェース450は、有線または無線を介して、情報処理装置400と、通信ネットワークとを接続するインターフェースである。

【0078】

ユーザインターフェース460は、例えば、ディスプレイ等の表示部、およびキーボード、マウス、タッチパネル等の入力部である。

【0079】

上記の実施の形態の一部又は全部は、以下の付記のようにも記載され得るが以下には限られない。

20

【0080】

[付記1]

音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出する貢献度推定手段と、

前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する話者特徴算出手段とを備える、音声処理装置。

【0081】

[付記2]

30

前記音声信号に含まれる音の種類比率を表す音声統計量を算出する音声統計量算出手段をさらに備え、

前記話者特徴算出手段は、前記音声信号の前記音声統計量と、前記音声信号の前記貢献度とに基づいて、前記認識特徴量を算出する、付記1に記載の音声処理装置。

【0082】

[付記3]

前記貢献度推定手段は、前記音声信号の前記貢献度として、

前記音声信号の一部が音声か否かを識別して算出した音声らしさを表す値、前記音声信号の一部が話者認識に正解する音声か否かを識別して算出した話者認識の正解しやすさを表す値、前記音声信号の一部が話者認識誤りを起こす音声か否かを識別して算出した話者認識の誤りやすさを表す値の少なくともいずれかひとつを算出する、付記1または2に記載の音声処理装置。

40

【0083】

[付記4]

前記貢献度推定手段は、

ニューラルネットワークを用いて前記音声信号の前記貢献度を算出する、付記3に記載の音声処理装置。

【0084】

[付記5]

前記話者特徴算出手段は、

50

前記認識特微量として i -vector を算出する、付記 3 または 4 に記載の音声処理装置。

【 0 0 8 5 】

[付記 6]

前記認識特微量に基づいて前記属性情報を認識する属性認識手段を備える、付記 1 ~ 5 のいずれか 1 つに記載の音声処理装置。

【 0 0 8 6 】

[付記 7]

前記特定の属性情報は、

音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、音声信号から推定される話者の性格の少なくともいずれか 1 つを表す情報である、付記 1 ~ 6 のいずれか 1 つに記載の音声処理装置。

10

【 0 0 8 7 】

[付記 8]

音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出し、

前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特微量を算出する、音声処理方法。

【 0 0 8 8 】

[付記 9]

前記音声信号に含まれる音の種類比率を表す音声統計量をさらに算出し、

20

前記音声信号の前記音声統計量と、前記音声信号の前記貢献度とに基づいて、前記認識特微量を算出する、付記 8 に記載の音声処理方法。

【 0 0 8 9 】

[付記 1 0]

前記音声信号の前記貢献度として、

前記音声信号の一部が音声か否かを識別して算出した音声らしさを表す値、前記音声信号の一部が話者認識に正解する音声か否かを識別して算出した話者認識の正解しやすさを表す値、前記音声信号の一部が話者認識誤りを起こす音声か否かを識別して算出した話者認識の誤りやすさを表す値の少なくともいずれかひとつを算出する、付記 8 または 9 に記載の音声処理方法。

30

【 0 0 9 0 】

[付記 1 1]

ニューラルネットワークを用いて前記音声信号の前記貢献度を算出する、付記 1 0 に記載の音声処理方法。

【 0 0 9 1 】

[付記 1 2]

前記認識特微量として i -vector を算出する、付記 1 0 または 1 1 に記載の音声処理方法。

【 0 0 9 2 】

[付記 1 3]

前記認識特微量に基づいて前記属性情報を認識する、付記 8 ~ 1 2 のいずれか 1 つに記載の音声処理方法。

40

【 0 0 9 3 】

[付記 1 4]

前記特定の属性情報は、

音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、音声信号から推定される話者の性格の少なくともいずれか 1 つを表す情報である、付記 8 ~ 1 3 のいずれか 1 つに記載の音声処理方法。

【 0 0 9 4 】

[付記 1 5]

50

コンピュータに、
 音声を表す音声信号に基づき、前記音声信号の部分区間の品質の程度を表す貢献度を算出する処理と、

前記音声信号の前記部分区間の前記貢献度を、前記部分区間の重みとして用いて、前記音声信号から特定の属性情報を認識するための認識特徴量を算出する処理とを実行させる、音声処理プログラム。

【0095】

[付記16]

前記コンピュータに、

前記音声信号に含まれる音の種類比率を表す音声統計量をさらに算出する処理と、
 前記音声信号の前記音声統計量と、前記音声信号の前記貢献度とに基づいて、前記認識特徴量を算出する処理とを実行させる、付記15に記載の音声処理方法。

10

【0096】

[付記17]

前記コンピュータに、

前記音声信号の前記貢献度として、

前記音声信号の一部が音声か否かを識別して算出した音声らしさを表す値、前記音声信号の一部が話者認識に正解する音声か否かを識別して算出した話者認識の正解しやすさを表す値、および前記音声信号の一部が話者認識誤りを起こす音声か否かを識別して算出した話者認識の誤りやすさを表す値の少なくともいずれかひとつを算出する処理を実行させる、付記15または16に記載の音声処理プログラム。

20

【0097】

[付記18]

前記コンピュータに、

ニューラルネットワークを用いて前記音声信号の前記貢献度を算出する処理を実行させる、付記17に記載の音声処理プログラム。

【0098】

[付記19]

前記コンピュータに、

前記認識特徴量として i-vector を算出する処理を実行させる、付記17または18に記載の音声処理プログラム。

30

【0099】

[付記20]

前記コンピュータに、

前記認識特徴量に基づいて前記属性情報を認識する処理を実行させる、付記15～19のいずれか1つに記載の音声処理プログラム。

【0100】

[付記21]

前記特定の属性情報は、

音声信号を発した話者、音声信号を構成する言語、音声信号に含まれる感情表現、音声信号から推定される話者の性格の少なくともいずれか1つを表す情報である、付記15～20のいずれか1つに記載の音声処理プログラム。

40

【符号の説明】

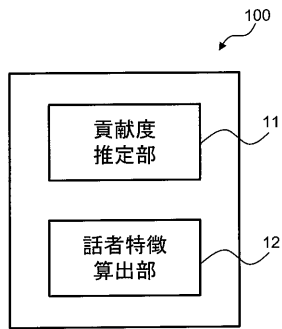
【0101】

- 11・・・貢献度推定部
- 12・・・話者特徴算出部
- 13・・・属性認識部
- 21・・・音声区間検出部
- 22・・・音声統計量算出部
- 23・・・貢献度記憶部

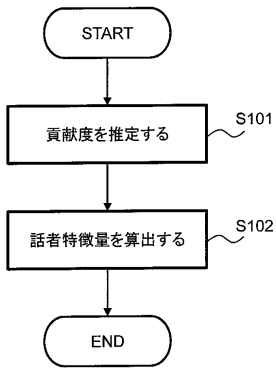
50

- 2 4 . . . 貢献度学習部
- 1 0 0 , 2 0 0 , 3 0 0 . . . 音声処理装置
- 4 0 0 . . . 情報処理装置
- 4 1 0 . . . 制御部 (C P U)
- 4 2 0 . . . 記憶部
- 4 3 0 . . . R O M
- 4 4 0 . . . R A M
- 4 5 0 . . . 通信インターフェース
- 4 6 0 . . . ユーザーインターフェース

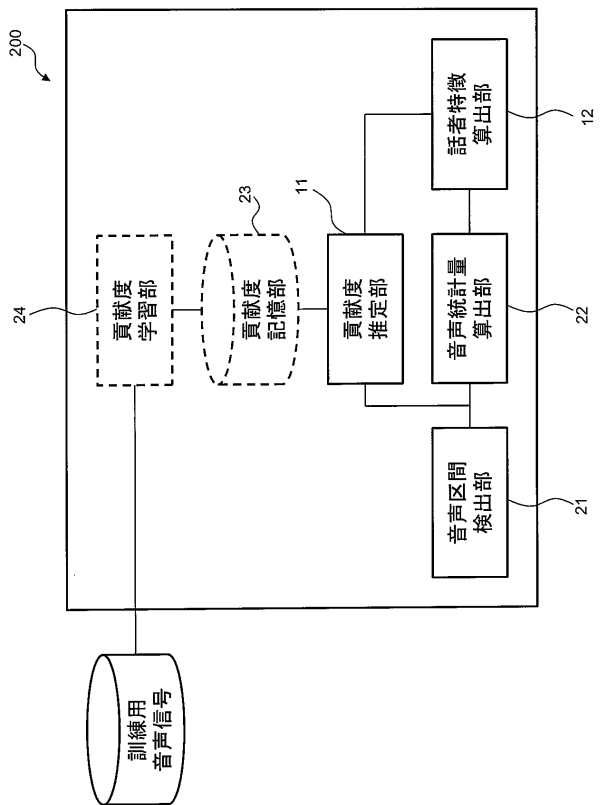
【図1】



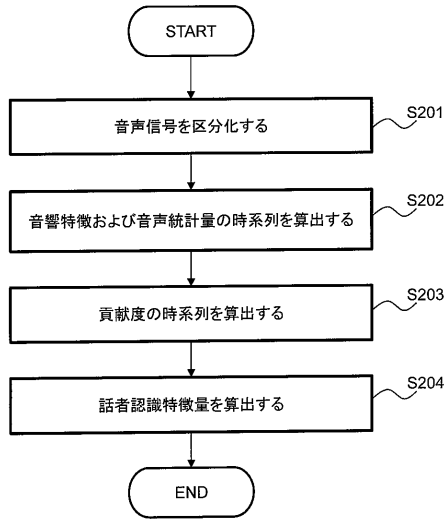
【図2】



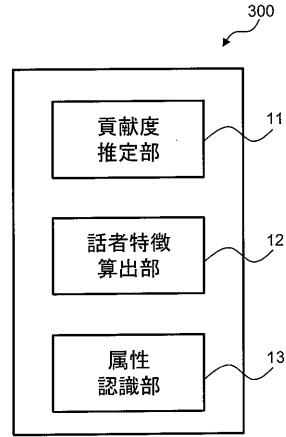
【図3】



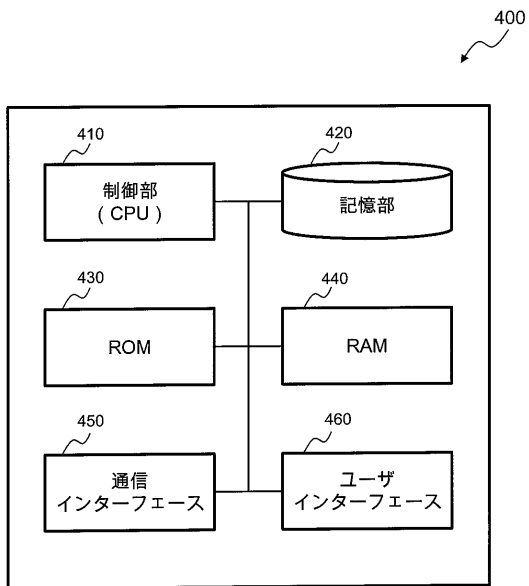
【図4】



【図5】



【図6】



フロントページの続き

(56)参考文献 特開2015-184378(JP,A)
特開2004-341340(JP,A)

(58)調査した分野(Int.Cl., DB名)
G10L 17/00