

(19) 日本国特許庁(JP)

再公表特許(A1)

(11) 国際公開番号

W02018/159612

発行日 令和2年1月9日 (2020.1.9)

(43) 国際公開日 平成30年9月7日 (2018.9.7)

(51) Int.Cl.

G10L 21/007 (2013.01)

F I

G10L 21/007

テーマコード (参考)

審査請求 未請求 予備審査請求 有 (全 28 頁)

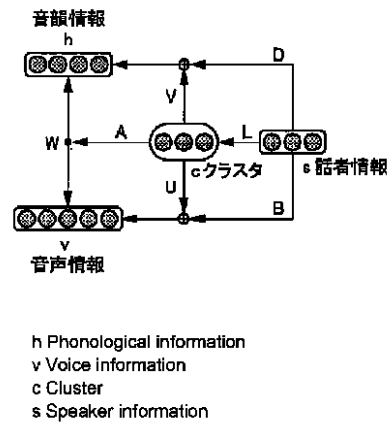
<p>出願番号 特願2019-503021 (P2019-503021)</p> <p>(21) 国際出願番号 PCT/JP2018/007268</p> <p>(22) 国際出願日 平成30年2月27日 (2018.2.27)</p> <p>(31) 優先権主張番号 特願2017-36109 (P2017-36109)</p> <p>(32) 優先日 平成29年2月28日 (2017.2.28)</p> <p>(33) 優先権主張国・地域又は機関 日本国 (JP)</p>	<p>(71) 出願人 504133110 国立大学法人電気通信大学 東京都調布市調布ヶ丘一丁目5番地1</p> <p>(74) 代理人 110000925 特許業務法人信友国際特許事務所</p> <p>(72) 発明者 中鹿 亘 東京都調布市調布ヶ丘一丁目5番地1 国立大学法人電気通信大学内</p>
---	---

最終頁に続く

(54) 【発明の名称】 声質変換装置、声質変換方法およびプログラム

(57) 【要約】

パラメータ学習ユニットとパラメータ記憶ユニットと声質変換処理ユニットとを備える。パラメータ学習ユニットは、入力データを表現する可視素子と、潜在的な情報を表現した隠れ素子との間に結合重みが存在すると仮定した制限ボルツマンマシンによる確率モデルを用意する。その確率モデルとして、固有の適応行列を持つ複数個の話者クラスタを定義し、それぞれの話者について、複数個の話者クラスタへの重みを推定して、パラメータを決定する。パラメータ記憶ユニットは、パラメータを記憶する。声質変換処理ユニットは、パラメータ記憶ユニットが記憶したパラメータと目標話者の話者情報とに基づいて、入力話者の音声に基づく音声情報の声質変換処理を行う。



【特許請求の範囲】

【請求項 1】

入力話者の音声を目標話者の音声に声質変換する声質変換装置であって、
学習用の音声に基づく音声情報およびその音声情報に対応する話者情報から、声質変換のためのパラメータを決定するパラメータ学習ユニットと、

前記パラメータ学習ユニットが決定したパラメータを記憶するパラメータ記憶ユニットと、

前記パラメータ記憶ユニットが記憶したパラメータと前記目標話者の話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う声質変換処理ユニットとを備え、

前記パラメータ学習ユニットは、音声に基づく音声情報、音声情報に対応する話者情報および音声中の音韻を表す音韻情報のそれぞれを変数とすることで、前記音声情報、前記話者情報および前記音韻情報のそれぞれの間の結合エネルギーの関係性を前記パラメータによって表す確率モデルを取得し、前記確率モデルとして、固有の適応行列を持つ複数個の話者クラスタを定義するようにした

声質変換装置。

【請求項 2】

さらに、前記パラメータ記憶ユニットが記憶したパラメータを前記入力話者の音声に適応して、適応後のパラメータを得る適応ユニットを備え、

前記パラメータ記憶ユニットは、前記適応ユニットで適応後のパラメータを記憶し、前記声質変換処理ユニットは、適応後のパラメータと前記目標話者の話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う

請求項 1 に記載の声質変換装置。

【請求項 3】

前記パラメータ学習ユニットと前記適応ユニットは共通の演算処理部で構成され、

前記学習用の音声に基づいてパラメータを決定する処理と、前記入力話者の音声に基づいて適応後のパラメータを得る処理を、前記共通の演算処理部で行うようにした

請求項 2 に記載の声質変換装置。

【請求項 4】

前記パラメータ学習ユニットが学習する際には、複数のクラスタが最も離れるように学習し、学習した複数のクラスタの間で、話者クラスタへの重みの位置を設定する

請求項 1 に記載の声質変換装置。

【請求項 5】

前記声質変換処理ユニットは、前記パラメータから前記目標話者の話者情報を得、得られた話者情報から前記目標話者の音声情報を得るようにした

請求項 1 に記載の声質変換装置。

【請求項 6】

音声情報の特徴量 $v = [v_1, \dots, v_I] \in \mathbb{R}^I$ と、音韻情報の特徴量 $h = [h_1, \dots, h_J] \in \{0, 1\}^J$ 、 $\sum_j h_j = 1$ との間に、話者情報の特徴量 $s = [s_1, \dots, s_R] \in \{0, 1\}^R$ 、 $\sum_r s_r = 1$ に依存した双方な結合重み $W \in \mathbb{R}^{I \times J}$ が存在すると仮定したとき、前記話者クラスタとして、話者クラスタ $c \in \mathbb{R}^K$ を導入し、話者クラスタ c を、

$$c \triangleq Ls$$

(但し、 $L \in \mathbb{R}^{K \times R} = [l_1 \dots l_R]$ の各列ベクトル l_r は、それぞれの話者クラスタへの重みを表す非負パラメータであり、 $\|l_r\|_1 = 1$ 、 l_r の制約を課す) と表現し、音響情報の特徴量のクラスタ依存項のバイアスパラメータを $U \in \mathbb{R}^{I \times K}$ 、音韻情報の特徴量のクラスタ依存項のバイアスパラメータを $V \in \mathbb{R}^{J \times K}$ 、として、話者非依存項、クラスタ依存項、および話者依存項のそれぞれを、

10

20

30

40

$$\tilde{W} \triangleq A \circ_{\frac{1}{3}} cW$$

$$\tilde{b} \triangleq b + Uc + Bs$$

$$\tilde{d} \triangleq d + Vc + Ds$$

として示す

請求項 1 に記載の声質変換装置。

【請求項 7】

入力話者の音声を目標話者の音声に声質変換する声質変換方法であって、
音声に基づく音声情報、音声情報に対応する話者情報および音声中の音韻を表す音韻情報のそれぞれを変数とすることで、前記音声情報、前記話者情報および前記音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを用意し、その確率モデルとして、固有の適応行列を持つ複数個の話者クラスタを定義し、それぞれの話者について、前記複数個の話者クラスタへの重みを推定して、学習用の音声についての前記パラメータを決定するパラメータ学習ステップと、

前記パラメータ学習ステップで得られたパラメータ、又は当該パラメータを前記入力話者の音声に適応した適応後のパラメータと、前記目標話者の前記話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う声質変換処理ステップとを含む、声質変換方法。

【請求項 8】

音声に基づく音声情報、音声情報に対応する話者情報および音声中の音韻を表す音韻情報のそれぞれを変数とすることで、前記音声情報、前記話者情報および前記音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを用意し、その確率モデルとして、固有の適応行列を持つ複数個の話者クラスタを定義し、それぞれの話者について、前記複数個の話者クラスタへの重みを推定して、学習用の音声についての前記パラメータを決定して記憶するパラメータ学習ステップと、

前記パラメータ学習ステップで得られたパラメータ、又は当該パラメータを入力話者の音声に適応した適応後のパラメータと、目標話者の前記話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う声質変換処理ステップと、
をコンピュータに実行させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は任意話者声質変換を可能とする声質変換装置、声質変換方法およびプログラムに関する。

【背景技術】

【0002】

従来、入力話者音声の音韻情報を保存したまま、話者性に関する情報のみを出力話者のものへ変換させる技術である声質変換の分野では、モデルの学習時において、入力話者と出力話者の同一発話内容による音声対であるパラレルデータを使用するパラレル声質変換が主流であった。

パラレル声質変換としては、GMM (Gaussian Mixture Model) に基づく手法、NMF (Non-negative Matrix Factorization) に基づく手法、DNN (Deep Neural Network) に基づく手法など、様々な統計的アプローチが提案されている (特許文献 1 参照)。パラレル声質変換では、パラレル制約のおかげで比較的高い精度が得られる反面、学習データとしては入力話者と出力話者の発話内容を一致させる必要があるため、利便性が損なわれてしまうという問題があった。

【0003】

これに対して、モデルの学習時に上述のパラレルデータを使用しない非パラレル声質変

10

20

30

40

50

換が注目を浴びている。非パラレル声質変換は、パラレル声質変換に比べて精度面で劣るものの自由発話を用いて学習を行うことができるため利便性や実用性は高い。非特許文献 1 には、入力話者の音声と出力話者の音声を用いて事前に個々のパラメータを学習しておくことで、学習データに含まれる話者を入力話者または目標話者とする声質変換を可能とする技術が記載されている。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2008 - 58696 号公報

【非特許文献】

【0005】

【非特許文献 1】T. Nakashika, T. Takiguchi, and Y. Ariki: "Parallel-Data-Free, Many-To-Many Voice Conversion Using an Adaptive Restricted Boltzmann Machine," Proceedings of Machine Learning in Spoken Language Processing (MLSLP) 2015, 6 pages, 2015.

【発明の概要】

【発明が解決しようとする課題】

【0006】

非特許文献 1 に記載の技術は、統計的な非パラレル声質変換アプローチとして、制限ボルツマンマシン (Restricted Boltzmann Machine: 以下 R B M と称する) を適用した、適応型 R B M (A R B M) に基づく声質変換に基づく声質変換を行う。このアプローチでは、複数の話者による音声データから自動的にそれぞれの話者固有の適応行列と、音響特徴量 (メルケプストラム) から話者に依存しない潜在特徴 (以下、これらを潜在的な音韻または単に音韻と呼ぶ) への射影行列を同時に推定する。これにより、入力話者の音声および入力話者の適応行列から計算した潜在的な音韻と、目標話者の適応行列を用いて音響特徴量を計算することで目標話者に近い音声を得るようにしている。

【0007】

一度学習によって潜在的な音韻を得るための射影行列が推定されれば、新たな入力話者・目標話者に対してそれぞれの適応行列のみを推定 (このステップを適応と呼ぶ) することで変換が可能となる。しかし、話者固有の適応行列は音響特徴量の二乗個のパラメータを含むため、音響特徴量の次元数や話者数が増えるほどパラメータ数が膨大となり、学習コストが掛かってしまう。そして、適応時に必要となるデータ数が多くなり、事前に学習していない話者のその場での変換が困難となってしまふといった問題が発生する。また、声質変換を利用する場面では、その場で音声を収録し、即座に変換を行いたいケースが考えられるが、従来技術では、即座に変換することは困難であった。

【0008】

本発明はかかる点に鑑み、各話者の発話について少ないデータ数で簡単に声質変換が可能な声質変換装置、声質変換方法およ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】

【0009】

上記課題を解決するため、本発明の声質変換装置は、入力話者の音声を目標話者の音声に声質変換する声質変換装置であって、パラメータ学習ユニットとパラメータ記憶ユニットと声質変換処理ユニットとを備える。

パラメータ学習ユニットは、学習用の音声に基づく音声情報およびその音声情報に対応する話者情報から、声質変換のためのパラメータを決定する。

パラメータ記憶ユニットは、パラメータ学習ユニットが決定したパラメータを記憶する。

声質変換処理ユニットは、パラメータ記憶ユニットが記憶したパラメータと目標話者の話者情報とに基づいて、入力話者の音声に基づく音声情報の声質変換処理を行う。

ここで、パラメータ学習ユニットは、音声に基づく音声情報、音声情報に対応する話者

10

20

30

40

50

情報および音声の音韻を表す音韻情報のそれぞれを変数とすることで、音声情報、話者情報および音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを取得し、確率モデルとして、固有の適応行列を持つ複数の話者クラスタを定義するようにした。

【0010】

また、本発明の声質変換方法は、入力話者の音声を目標話者の音声に声質変換する方法であって、パラメータ学習ステップと声質変換処理ステップとを含む。

パラメータ学習ステップは、音声に基づく音声情報、音声情報に対応する話者情報および音声の音韻を表す音韻情報のそれぞれを変数とすることで、音声情報、話者情報および音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを用意する。そして、その確率モデルとして、固有の適応行列を持つ複数の話者クラスタを定義し、それぞれの話者について、複数の話者クラスタへの重みを推定して、学習用の音声についてのパラメータを決定する。

声質変換処理ステップは、パラメータ学習ステップで得られたパラメータ、又は当該パラメータを入力話者の音声に適応した適応後のパラメータと、目標話者の話者情報とに基づいて、入力話者の音声に基づく音声情報の声質変換処理を行う。

【0011】

また本発明のプログラムは、上述した声質変換方法のパラメータ学習ステップと声質変換処理ステップとをコンピュータに実行させるものである。

【0012】

本発明によれば、話者クラスタにより目標話者を設定することができるため、従来よりも非常に少ないデータ数で、入力話者音声を目標話者音声に声質変換できるようになる。

【図面の簡単な説明】

【0013】

【図1】本発明の一実施の形態例に係る声質変換装置の構成例(例1)を示すブロック図である。

【図2】本発明の一実施の形態例に係る声質変換装置の構成例(例2)を示すブロック図である。

【図3】声質変換装置のハードウェア構成例を示すブロック図である。

【図4】従来の確率モデルを模式的に示す説明図である。

【図5】声質変換装置のパラメータ推定部が備える確率モデルを模式的に示す説明図である。

【図6】本発明の一実施の形態例に係る処理全体の流れを示すフローチャートである。

【図7】図6のステップS3の学習の詳細例を示すフローチャートである。

【図8】図6のステップS4の適応の詳細例を示すフローチャートである。

【図9】図6のステップS8の声質変換の詳細例を示すフローチャートである。

【図10】本発明の一実施形態によるクラスタの重み分布の例を示す説明図である。

【図11】声質変換装置のパラメータ推定部が備える確率モデルの別の例を示す説明図である。

【発明を実施するための形態】

【0014】

以下、本発明の好適な一実施形態例について説明する。

【0015】

[1. 構成]

図1は、本発明の一実施形態例にかかる声質変換装置の構成例(例1)を示す図である。図1においてPC等により構成される声質変換装置1は、事前に、学習用音声信号と学習用音声信号に対応する話者の情報(対応話者情報)に基づいて学習を行っておくことで、任意の話者による変換用音声信号(適応話者音声信号)を、目標話者の声質に変換し、変換済み音声信号として出力する。

学習用音声信号は、予め記録された音声データに基づく音声信号でもよく、また、マイ

10

20

30

40

50

クロフォン等により話者が話す音声（音波）を直接電気信号に変換したものでよい。また、対応話者情報は、ある学習用音声信号と他の学習用音声信号とが同じ話者による音声信号か異なる話者による音声信号かを区別できるものであればよい。

【0016】

声質変換装置1は、パラメータ学習ユニット11と声質変換処理ユニット12とパラメータ記憶ユニット13とを備える。パラメータ学習ユニット11は、学習用音声信号と対応話者情報とに基づいた学習処理により声質変換のためのパラメータを決定する。パラメータ学習ユニット11が決定したパラメータは、パラメータ記憶ユニット13に記憶される。パラメータ記憶ユニット13に記憶されたパラメータは、適応処理によって、パラメータ学習ユニット11が入力話者の適応後のパラメータに変換する。声質変換処理ユニット12は、上述の学習処理および適応処理によりパラメータが決定された後、決定されたパラメータと目標とする話者の情報（目標話者情報）とに基づいて変換用音声信号の声質を目標話者の声質に変換し、変換済み音声信号として出力する。なお、パラメータ学習ユニット11が学習処理と適応処理の双方を行うのは一例であり、後述する図2に示すように、パラメータ学習ユニット11と別に適応ユニット14を備えるようにしてもよい。

10

【0017】

パラメータ学習ユニット11は、音声信号取得部111と前処理部112と話者情報取得部113とパラメータ推定部114を備える。音声信号取得部111は、前処理部112に接続され、前処理部112および話者情報取得部113は、それぞれパラメータ推定部114に接続される。

20

【0018】

音声信号取得部111は、接続された外部機器から学習用音声信号を取得するものであり、例えば、マウスやキーボード等の図示しない入力部からのユーザの操作に基づいて学習用音声信号が取得される。また、音声信号取得部111は、接続される不図示のマイクロフォンから、話者の発話をリアルタイムに取り込むようにしてもよい。なお、以下の説明では、パラメータ学習ユニット11が学習用音声信号を取得してパラメータを得る処理を述べるが、パラメータ学習ユニット11が適応話者音声信号に適応したパラメータを得る適応処理時にも、各処理部は同様の処理が行われる。適応処理の詳細については後述するが、適応処理時には、学習処理でパラメータ記憶ユニット13に記憶されたパラメータを、適応話者音声信号に適応したパラメータとする適応化処理が行われる。

30

前処理部112は、音声信号取得部111で取得された学習用音声信号を単位時間ごと（以下、フレームという）に切り出し、MFC C（Mel-Frequency Cepstrum Coefficient s：メル周波数ケプストラム係数）やメルケプストラム特徴量などのフレームごとの音声信号のスペクトル特徴量を計算した後、正規化を行うことで学習用音声情報を生成する。

【0019】

対応話者情報取得部113は、音声信号取得部111による学習用音声信号の取得に紐付けられた対応話者情報を取得する。対応話者情報は、ある学習用音声信号の話者と他の学習用音声信号の話者とを区別できるものであればよく、例えば、図示しない入力部からのユーザの入力によって取得される。また、複数の学習用音声信号のそれぞれについて互いに話者が異なることが明らかであれば、学習用音声信号の取得に際して対応話者情報取得部113が自動で対応話者情報を付与してもよい。例えば、パラメータ学習ユニット11が10人の話し声の学習を行うと仮定すると、対応話者情報取得部113は、音声信号取得部111に入力中の学習用音声信号が10人の内のどの話者の話し声の音声信号であるかを区別する情報（対応話者情報）を、自動的にまたはユーザからの入力により取得する。なお、ここで話し声の学習を行う人数を10人としたのは、あくまでも一例である。パラメータ学習ユニット11は、最低でも2人の音声が入力されれば学習が可能であるが、人数が多い方がより精度の高い学習ができることになる。

40

【0020】

パラメータ推定部114は、音声情報推定部1141と話者情報推定部1142と音韻情報推定部1143とによって構成されるRBM（制限ボルツマンマシン）を適用した、

50

適応型 R B M (A R B M) の確率モデルを持ち、学習用音声信号に基づいてパラメータの推定を行う。パラメータ推定部 1 1 4 が学習処理によって推定したパラメータは、パラメータ記憶ユニット 1 3 に記憶される。この学習処理で得たパラメータは、適応話者の音声信号がパラメータ学習ユニット 1 1 に入力されたとき、パラメータ記憶ユニット 1 3 からパラメータ学習ユニット 1 1 に読み出され、そのときの適応話者の音声信号に適応したパラメータとされる。

【 0 0 2 1 】

パラメータ推定部 1 1 4 がパラメータを推定する際に適用される本実施形態例の確率モデルでは、各推定部 1 1 4 1 , 1 1 4 2 , 1 1 4 3 が持つ音声情報、話者情報、および音韻情報の他に、話者の特徴から得た複数の話者クラスタの情報を持つ。すなわち、パラメータ推定部 1 1 4 は、この話者クラスタを計算する話者クラスタ計算部 1 1 4 4 を有する。さらに、本実施形態例の確率モデルでは、各情報のそれぞれの間の結合エネルギーの関係性を表すパラメータを持つ。なお、以下の説明では、本実施形態例の確率モデルを、話者クラスタ適応型 R B M と称する。話者クラスタ適応型 R B M の詳細については後述する。

10

【 0 0 2 2 】

音声情報推定部 1 1 4 1 は、音韻情報および話者情報ならびに各種パラメータを用いて音声情報を取得する。ここで、音声情報とは、それぞれの話者の音声信号の音響ベクトル (スペクトル特徴量やケプストラム特徴量など) を意味する。

【 0 0 2 3 】

話者情報推定部 1 1 4 2 は、音声情報および音韻情報ならびに各種パラメータを用いて話者情報を推定する。ここで、話者情報とは、話者を特定するための情報であり、それぞれの話者の音声を持つ音響ベクトル情報である。すなわち、この話者情報 (話者ベクトル) は、同じ話者の音声信号に対しては全て共通であり、異なる話者の音声信号に対しては互いに異なるような、音声信号の発話者を特定させるベクトルを意味している。

20

【 0 0 2 4 】

音韻情報推定部 1 1 4 3 は、音声情報および話者情報ならびに各種パラメータにより音韻情報を推定する。ここで音韻情報とは、音声情報に含まれる情報の中から、学習を行う全ての話者に共通となる情報である。例えば、入力した学習用音声信号が、「こんにちは」と発話した音声の信号であるとき、この音声信号から得られる音韻情報は、その「こんにちは」と発話した言葉の情報に相当する。但し、本実施の形態例での音韻情報は、言葉に相当する情報であっても、いわゆるテキストの情報ではなく、言語の種類に限定されない音韻の情報であり、どのような言語で話者が話した場合にも共通となる、音声信号の中で潜在的に含まれる、話者情報以外の情報を表すベクトルである。

30

【 0 0 2 5 】

話者クラスタ計算部 1 1 4 4 は、入力中の学習用音声信号から得た話者情報に対応したクラスタを計算する。すなわち、パラメータ推定部 1 1 4 が備える話者クラスタ適応型 R B M は、話者情報を示すクラスタを複数持ち、話者クラスタ計算部 1 1 4 4 は、入力中の学習用音声信号から得た話者情報に対応するクラスタを計算する。

【 0 0 2 6 】

また、パラメータ推定部 1 1 4 が備える話者クラスタ適応型 R B M は、音声情報、話者情報、音韻情報および話者クラスタの情報を持つだけでなく、各情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表すようにしている。

40

【 0 0 2 7 】

声質変換処理ユニット 1 2 は、音声信号取得部 1 2 1 と前処理部 1 2 2 と話者情報設定部 1 2 3 と声質変換部 1 2 4 と後処理部 1 2 5 と音声信号出力部 1 2 6 とを備える。音声信号入力 1 2 1、前処理部 1 2 2、声質変換部 1 2 4、後処理部 1 2 5 および音声信号出力部 1 2 6 は順次接続され、声質変換部 1 2 4 には、更にパラメータ学習ユニット 1 1 のパラメータ推定部 1 1 4 が接続される。

【 0 0 2 8 】

50

音声信号取得部 1 2 1 は、変換用音声信号を取得し、前処理部 1 2 2 は、変換用音声信号に基づき変換用音声情報を生成する。本実施の形態例では、音声信号取得部 1 2 1 が取得する変換用音声信号は、任意の話者による変換用音声信号でよい。

音声信号取得部 1 2 1 および前処理部 1 2 2 は、上述したパラメータ学習ユニット 1 1 の音声信号取得部 1 1 1 および前処理部 1 1 2 の構成と同じであり、別途設置することなくこれらを兼用してもよい。

【 0 0 2 9 】

話者情報設定部 1 2 3 は、声質変換先である目標話者を設定し目標話者情報を出力する。話者情報設定部 1 2 3 で設定される目標話者は、ここでは、パラメータ学習ユニット 1 1 のパラメータ推定部 1 1 4 が事前に学習処理して話者情報を取得した話者の中から選ばれる。話者情報設定部 1 2 3 は、例えば、図示しないディスプレイ等に表示された複数の目標話者の選択肢（パラメータ推定部 1 1 4 が事前に学習処理した話者の一覧など）からユーザが図示しない入力部によって 1 つの目標話者を選択するものであってもよく、また、その際に、図示しないスピーカにより目標話者の音声を確認できるようにしてもよい。

10

【 0 0 3 0 】

声質変換部 1 2 4 は、目標話者情報に基づいて変換用音声情報に声質変換を施し、変換済み音声情報を出力する。声質変換部 1 2 4 は、音声情報設定部 1 2 4 1、話者情報設定部 1 2 4 2、音韻情報設定部 1 2 4 3、および話者クラスタ計算部 1 2 4 4 を持つ。この音声情報設定部 1 2 4 1、話者情報設定部 1 2 4 2、音韻情報設定部 1 2 4 3、および話者クラスタ計算部 1 2 4 4 は、上述のパラメータ推定部 1 1 4 において、話者クラスタ適応型 R B M の確率モデルが持つ音声情報推定部 1 1 4 1、話者情報推定部 1 1 4 2、音韻情報推定部 1 1 4 3、および話者クラスタ計算部 1 1 4 4 と同等の機能を持つ。

20

【 0 0 3 1 】

すなわち、音声情報設定部 1 2 4 1、話者情報設定部 1 2 4 2 および音韻情報設定部 1 2 4 3 には、それぞれ音声情報、話者情報および音韻情報が設定されるが、音韻情報設定部 1 2 4 3 に設定される音韻情報は、前処理部 1 2 2 から供給される音声情報に基づいて得た情報である。一方、話者情報設定部 1 2 4 2 に設定される話者情報は、パラメータ学習ユニット 1 1 内の話者情報推定部 1 1 4 2 での推定結果から取得した目標話者についての話者情報（話者ベクトル）である。音声情報設定部 1 2 4 1 に設定される音声情報は、これら話者情報設定部 1 2 4 2 および音韻情報設定部 1 2 4 3 に設定された話者情報および音韻情報と各種パラメータとから得られる。話者クラスタ計算部 1 2 4 4 は、目標話者の話者クラスタ情報を計算する。

30

なお、図 1 では声質変換部 1 2 4 を設ける構成を示したが、声質変換部 1 2 4 を別途設置することなく、パラメータ推定部 1 1 4 の各種パラメータを固定することで、パラメータ推定部 1 1 4 が声質変換の処理を実行する構成としてもよい。

【 0 0 3 2 】

後処理部 1 2 5 は、声質変換部 1 2 4 で得られた変換済み音声情報に逆正規化処理を施し、更に逆 F F T 処理することでスペクトル情報をフレームごとの音声信号へ戻した後に結合し、変換済み音声信号を生成する。

音声信号出力部 1 2 6 は、接続される外部機器に対して変換済み音声信号を出力する。接続される外部機器としては、例えば、スピーカなどが挙げられる。

40

【 0 0 3 3 】

図 2 は、本発明の一実施形態例にかかる声質変換装置の別の構成例（例 2）を示す図である。

図 2 に示す声質変換装置 1 は、適応話者音声信号によりパラメータの適応処理を行う適応ユニット 1 4 を備える点が、図 1 に示す声質変換装置 1 と異なる。すなわち、図 1 に示す声質変換装置 1 では、パラメータ学習ユニット 1 1 が、学習処理と適応処理の双方を行うようにしたのに対して、図 2 に示す声質変換装置 1 では、適応ユニット 1 4 が適応処理を行うようにした点が異なる。

【 0 0 3 4 】

50

適応ユニット 1 4 は、音声信号取得部 1 4 1 と前処理部 1 4 2 と適応話者情報取得部 1 4 3 とパラメータ推定部 1 4 4 を備える。音声信号取得部 1 4 1 は、適応話者音声信号を取得し、取得した音声信号を前処理部 1 4 2 に出力する。前処理部 1 4 2 は、音声信号の前処理を行って適応用音声情報を得、得られた適応用音声情報をパラメータ推定部 1 4 4 に供給する。適応話者情報取得部 1 4 3 は、適応話者についての話者情報を取得し、取得した適応話者情報をパラメータ推定部 1 4 4 に供給する。

パラメータ推定部 1 4 4 は、音声情報推定部 1 4 4 1 と話者情報推定部 1 4 4 2 と音韻情報推定部 1 4 4 3 と話者クラスタ計算部 1 4 4 4 を有し、音声情報、話者情報、音韻情報、および話者クラスタの情報を保持する。

【 0 0 3 5 】

適応ユニット 1 4 で得られた適用後のパラメータは、パラメータ記憶ユニット 1 3 に記憶した後、声質変換処理ユニット 1 2 に供給される。あるいは、適応ユニット 1 4 で得られた適用後のパラメータを、直接、声質変換処理ユニット 1 2 に供給するようにしてもよい。

図 2 に示す声質変換装置 1 のその他の部分については、図 1 に示す声質変換装置 1 と同様に構成する。

【 0 0 3 6 】

図 3 は、声質変換装置 1 のハードウェア構成例を示す図である。ここでは、声質変換装置 1 をコンピュータ (P C) で構成した例を示す。

図 3 に示すように、声質変換装置 1 は、バス 1 0 7 を介して相互に接続された C P U (中央制御ユニット : Central Processing Unit) 1 0 1、R O M (Read Only Memory) 1 0 2、R A M (Random Access Memory) 1 0 3、H D D (Hard Disk Drive) / S S D (Solid State Drive) 1 0 4、接続 I / F (Interface) 1 0 5、通信 I / F 1 0 6 を備える。C P U 1 0 1 は、R A M 1 0 3 をワークエリアとして R O M 1 0 2 または H D D / S S D 1 0 4 等に格納されたプログラムを実行することで、声質変換装置 1 の動作を統括的に制御する。接続 I / F 1 0 5 は、声質変換装置 1 に接続される機器とのインターフェースである。通信 I / F は、ネットワークを介して他の情報処理機器と通信を行うためのインターフェースである。

【 0 0 3 7 】

学習用音声信号、変換用音声信号、および変換済み音声信号の入出力および設定は、接続 I / F 1 0 5 または通信 I / F 1 0 6 を介して行われる。パラメータ記憶ユニット 1 3 でのパラメータの記憶は、R A M 1 0 3 または H D D / S S D 1 0 4 により行われる。図 1 で説明した声質変換装置 1 の機能は、C P U 1 0 1 において所定のプログラムが実行されることで実現される。プログラムは、記録媒体を経由して取得してもよく、ネットワークを経由して取得してもよく、R O M に組み込んで使用してもよい。また、一般的なコンピュータとプログラムの組合せでなく、A S I C (Application Specific Integrated Circuit) や F P G A (Field Programmable Gate Array) などの論理回路を組むことで、声質変換装置 1 の構成を実現するためのハードウェア構成にしてもよい。

【 0 0 3 8 】

[2 . 話者クラスタ適応型 R B M の定義]

次に、パラメータ推定部 1 1 3 および符号化部 1 2 3 が持つ確率モデルである、話者クラスタ適応型 R B M について説明する。

まず、本発明に適用される話者クラスタ適応型 R B M を説明する前に、既に提案した確率モデルである、適応型 R B M について説明する。

図 4 は、適応型 R B M のグラフ構造を模式的に示す図である。

適応型 R B M の確率モデルは、音声情報 v 、話者情報 s および音韻情報 h と、それぞれの情報の結合エネルギーの関係性を示すパラメータを持つ。ここでは、音響 (メルケプストラム) 情報の特徴量 $v = [v_1, \dots, v_I] \quad R^I$ と、音韻情報の特徴量 $h = [h_1, \dots, h_J] \quad \{ 0, 1 \}^J$ 、 $h_j = 1$ との間に、話者特徴量 $s = [s_1, \dots, s_R] \quad \{ 0, 1 \}^R$ 、 $r s_r = 1$ に依存した双方向な結合重み $W \quad R^I \times J$ が

10

20

30

40

50

存在すると仮定したとき、適応型 R B M の確率モデルは、次の [数 1] 式 ~ [数 3] 式で示される条件付き確率密度関数で示される。

【 0 0 3 9 】

【 数 1 】

$$p(\mathbf{v}, \mathbf{h} | \mathbf{s}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h} | \mathbf{s})}$$

【 0 0 4 0 】

【 数 2 】

$$E(\mathbf{v}, \mathbf{h} | \mathbf{s}) = \frac{1}{2} \left\| \frac{\mathbf{v} - \tilde{\mathbf{b}}}{\sigma} \right\|_2^2 - \tilde{\mathbf{d}}^\top \mathbf{h} - \left(\frac{\mathbf{v}}{\sigma^2} \right)^\top \tilde{\mathbf{W}} \mathbf{h}$$

10

【 0 0 4 1 】

【 数 3 】

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h} | \mathbf{s})}$$

【 0 0 4 2 】

但し、 \mathbf{R}^I は音響特徴量の偏差を表すパラメータであり、 $\mathbf{b} \in \mathbf{R}^I$ および $\mathbf{d} \in \mathbf{R}^J$ はそれぞれ話者特徴量 \mathbf{s} に依存した音響特徴量、音韻特徴量のバイアスを表す。式の中の記号の上に付けられた「 \sim 」は、該当する情報が話者に依存した情報であることを示す。なお、明細書の中では、表記上の制約のため、「 \sim 」を記号の上に付与できないので、例えば $W(\sim)$ のように、記号の後に括弧で示す。「 \wedge 」などの、記号の上に付与して示す他の記号についても、同様に表記する。

20

また、[数 2] 式の右辺の括弧および「 \cdot^2 」は、それぞれ要素ごとの除算、要素ごとの二乗を表す。話者依存の項 $W(\sim)$ 、 $\mathbf{b}(\sim)$ 、 $\mathbf{d}(\sim)$ は、話者非依存パラメータと話者依存パラメータを用いて、下記の [数 4] 式 ~ [数 6] 式のように定義される。

【 0 0 4 3 】

【 数 4 】

$$\tilde{\mathbf{W}} \triangleq \sum_r \mathbf{A}_r \mathbf{s}_r \mathbf{W} = \mathbf{A} \circ_3^1 \mathbf{s} \mathbf{W}$$

【 0 0 4 4 】

【 数 5 】

$$\tilde{\mathbf{b}} \triangleq \mathbf{b} + \sum_r \mathbf{b}_r \mathbf{s}_r = \mathbf{b} + \mathbf{B} \mathbf{s}$$

【 0 0 4 5 】

【 数 6 】

$$\tilde{\mathbf{d}} \triangleq \mathbf{d} + \sum_r \mathbf{d}_r \mathbf{s}_r = \mathbf{d} + \mathbf{D} \mathbf{s}$$

40

【 0 0 4 6 】

ここで、 $\mathbf{W} \in \mathbf{R}^{I \times J}$ 、 $\mathbf{b} \in \mathbf{R}^I$ 、 $\mathbf{d} \in \mathbf{R}^J$ は話者非依存パラメータを表し、 $\mathbf{A}_r \in \mathbf{R}^{I \times I}$ ($\mathbf{A} = \{ \mathbf{A}_r \}_{r=1}^R$)、 $\mathbf{b}_r \in \mathbf{R}^I$ ($\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R]$)、 $\mathbf{d}_r \in \mathbf{R}^J$ ($\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_R]$) は、話者 r に依存したパラメータを表す。また、 \mathbf{i}^j は左テンソルのモード i 、右テンソルのモード j に沿った内積演算を表す。

【 0 0 4 7 】

ここでは、音響特徴量はクリーン音声のメルケプストラムとし、発話者の違いによるバ

50

ラメータ変動は、話者特徴量 s によって規定される話者依存項（[数4]式，[数5]式，[数6]式）で吸収する。したがって、音韻特徴量は話者に依存しないいずれかの要素のみがアクティブとなる観測不可能な特徴量である、音韻の情報が含まれることになる。

【0048】

このように、適応型 R B M によって音響特徴量と音韻特徴量を得ることができるが、適応型 R B M では、話者依存パラメータの数は $(I^2 R)$ に比例し、音響特徴量の二乗 (I^2) が比較的大きいため、話者数が増加するほど推定するパラメータ数が膨大となり、計算に要するコストが増加してしまう。また、ある話者 r の適応時においても、推定すべきパラメータ数が $(I^2 + I + J)$ となり、過学習を避けるために相応に多くのデータを必要とする問題があった。

【0049】

ここで、本発明では、これらの問題を解決するために、話者クラスタ適応型 R B M を適用する。

図5は、話者クラスタ適応型 R B M のグラフ構造を模式的に示す図である。

話者クラスタ適応型 R B M の確率モデルは、音声情報 v 、話者情報 s および音韻情報 h と、それぞれの情報の結合エネルギーの関係性を示すパラメータの他に、話者クラスタ $c \in R^K$ を持つ。話者クラスタ c は、次の[数7]式と恒等的に表現される。

【0050】

【数7】

$$c \triangleq Ls$$

【0051】

但し、 $L \in R^{K \times R} = [l_1 \dots l_R]$ の各列ベクトル l_r は、それぞれの話者クラスタへの重みを表す非負パラメータであり、 $\|l_r\|_1 = 1$ 、 l_r の制約を課す。

先に説明した適応型 R B M (図4)では、話者ごとに適応行列を用意したが、本発明の話者クラスタ適応型 R B M ではクラスタごとに適応行列を用意する。また、音響特徴量、音韻特徴量のバイアスは、話者非依存項、クラスタ依存項、話者依存項の加算で表現される。すなわち、話者依存の項 $W(-)$ 、 $b(-)$ 、 $d(-)$ は、下記の[数8]式～[数10]式のように定義される。

【0052】

【数8】

$$\tilde{W} \triangleq A o_3^1 c W$$

【0053】

【数9】

$$\tilde{b} \triangleq b + Uc + Bs$$

【0054】

【数10】

$$\tilde{d} \triangleq d + Vc + Ds$$

【0055】

ここで、音響情報の特徴量のクラスタ依存項のバイアスパラメータを $U \in R^{I \times K}$ 、音韻情報の特徴量のクラスタ依存項のバイアスパラメータを $V \in R^{J \times K}$ とする。

[数8]式で示される $A = \{A_k\}_{k=1}^K$ と、先に説明した適応型 R B M での[数4]式における A を比較すると、適応型 R B M では $(I^2 R)$ 個のパラメータが含まれていたのに対して、話者クラスタ適応型 R B M では $(I^2 K)$ 個となり、大幅にパラメータ数を削減することができる。例えば、一例としては、 $R = 58$ 、 $I = 32$ 、 $K = 8$ に設定し

10

20

30

40

50

た場合、先に説明した適応型 R B M ではパラメータ数 5 9 3 9 2 個になるが、話者クラスタ適応型 R B M では 8 1 9 2 個になり、大幅にパラメータ数を削減できる。

【 0 0 5 6 】

また、先に説明した適応型 R B M では、話者一人につき $I^2 + I + J$ ($= 1 0 7 2$) 個のパラメータ ($H = 1 6$ の場合) であったのに対して、話者クラスタ適応型 R B M では、話者一人につき $K + I + J$ ($= 5 6$) 個のパラメータでよい。したがって、話者クラスタ適応型 R B M によると、大幅にパラメータ数を削減することができ、少ないデータで適応が可能になる。

【 0 0 5 7 】

話者クラスタ適応型 R B M においても、条件付き確率 $p(v, h | s)$ を、先に説明した [数 1] 式 ~ [数 3] 式で定義する。このとき、条件付き確率 $p(v | h, s)$, $p(h | v, s)$ は、それぞれ次の [数 1 1] 式および [数 1 2] 式に示すようになる。

【 0 0 5 8 】

【 数 1 1 】

$$p(v|h, s) = \mathcal{N}(v|\tilde{b} + \tilde{W}h, \sigma^2)$$

【 0 0 5 9 】

【 数 1 2 】

$$p(h|v, s) = \mathcal{B}(h|f(\tilde{d} + \tilde{W}^\top \frac{v}{\sigma^2}))$$

【 0 0 6 0 】

但し、[数 1 1] 式の右辺の $\mathcal{N}(\cdot)$ は次元独立の多変量正規分布、[数 1 2] 式の右辺の $\mathcal{B}(\cdot)$ は多次元ベルヌーイ分布、 $f(\cdot)$ は要素ごとのsoftmax関数を表す。

音韻特徴量 h は既知であり、ある話者 r の音響特徴量の平均ベクトル μ_r を考えると、[数 1 1] 式より、平均ベクトルは [数 1 3] 式に示すようになる。

【 0 0 6 1 】

【 数 1 3 】

$$\begin{aligned} \mu_r &= \mathbf{b} + \mathbf{b}_r + \mathbf{U}\lambda_r + \mathcal{A} \circ_3^1 \lambda_r \mathbf{W}h \\ &= \mathbf{M}\lambda'_r + \mathbf{b}_r \end{aligned}$$

【 0 0 6 2 】

但し、 $\mathbf{r} = [\mathbf{r}^\top \ 1]^\top$ は、 \mathbf{r} の拡張ベクトルであり、 $\mathbf{M} = [\mu_1, \dots, \mu_{K+1}]$ の各列ベクトルは、[数 1 4] 式で定義される。

【 0 0 6 3 】

【 数 1 4 】

$$\mu_k = \begin{cases} \mathbf{u}_k + \mathbf{A}_k \mathbf{W}h & (k = 1, \dots, K) \\ \mathbf{b} & (k = K + 1) \end{cases}$$

【 0 0 6 4 】

本発明の一実施形態例による話者クラスタ適応型 R B M では、話者依存項 \mathbf{b}_r が存在し、話者非依存平均ベクトル μ_k が [数 1 4] 式のように構造化される特徴を持つ。また、潜在的な音韻特徴量を陽に確率変数として定義している。

【 0 0 6 5 】

また、本発明の一実施形態例による話者クラスタ適応型 R B M では、話者非依存パラメータと話者クラスタ重みを同時に推定することができる。すなわち、 R 人の話者による N フレームの音声データ $\{v_n | s_n\}_{n=1}^N$ に対する対数尤度 ([数 1 5] 式) を最大

10

20

30

40

50

化するように、確率的勾配法を用いて全てのパラメータ $\theta = \{W, U, V, A, L, B, D, b, d, \dots\}$ を同時に更新し推定することが可能である。ここでは、それぞれのパラメータの勾配は省略する。

【0066】

【数15】

$$\mathcal{L} = \log \prod_n p(\mathbf{v}_n | \mathbf{s}_n) = \sum_n \log \sum_h p(\mathbf{v}_n, \mathbf{h}_n | \mathbf{s}_n)$$

【0067】

各勾配には計算困難なモデルに対する期待値が出現するが、通常のRBMの確率モデルと同様に、CD法(Contrastive Divergence法)を用いることで、効率よく近似することができる。

また、クラスタ重みの非負条件を満たすために、 $w_r = e^{z_r}$ と置き換えて、 z_r でパラメータ更新を行う。クラスタ重みはパラメータ更新後、 $\|w_r\|_1 = 1$ を満たすように正規化する。

さらに、モデルの学習が行われれば、音韻特徴量およびクラスタの形成が完了したとみなし、新たな話者 r について、 $\theta_r = \{w_r, b_r, d_r\}$ のみを更新し推定し、他のパラメータは固定する。

【0068】

この話者クラスタ適応型RBMを声質変換に適用する際には、ある入力話者の音声の音響特徴量 $v^{(i)}$ および話者特徴量 $s^{(i)}$ 、目標話者の話者特徴量 $s^{(o)}$ が与えられたとき、最も確率の高い音響特徴量 $v^{(o)}$ が目標話者の音響特徴量であるとして、[数16]式に示すように定式化される。

【0069】

【数16】

$$\begin{aligned} \hat{v}^{(o)} &\triangleq \operatorname{argmax}_v p(v | v^{(i)}, s^{(i)}, s^{(o)}) \\ &\simeq \operatorname{argmax}_v p(v | \hat{h}, s^{(o)}) \\ &= \mathbf{b} + \mathbf{B}s^{(o)} + \mathbf{U}Ls^{(o)} + \mathcal{A} \circ_3^1 Ls^{(o)} \mathbf{W}\hat{h} \end{aligned}$$

30

【0070】

但し、 $h^{(\wedge)}$ は、入力話者の音響特徴量および話者特徴量が与えられたときの h の条件付き期待値であり、[数17]式で表される。

【0071】

【数17】

$$\begin{aligned} \hat{h} &\triangleq \mathbb{E}[h | v^{(i)}, s^{(i)}] \\ &= \mathbf{f}(d + \mathbf{V}Ls^{(i)} + \mathbf{D}s^{(i)} + \mathbf{W}^\top (\mathcal{A} \circ_3^1 Ls^{(i)})^\top \frac{v^{(i)}}{\sigma^2}) \end{aligned}$$

40

【0072】

[3. 声質変換動作]

図6は、本発明の実施形態例による声質変換処理動作を示すフローチャートである。図6に示すように、パラメータ学習処理として、声質変換装置1のパラメータ学習ユニット11の音声信号取得部111と話者情報取得部113とは、図示しない入力部によるユーザの指示に基づいて学習用音声信号とその対応話者情報とをそれぞれ取得する(ステップS1)。

前処理部112は、音声信号取得部111が取得した学習用音声信号からパラメータ推定部114に供給する学習用音声情報を生成する(ステップS2)。ここでは、例えば学

50

習用音声信号をフレームごと（例えば、5 m s e c ごと）に切り出し、切り出された学習用音声信号に F F T 処理などを施すことでスペクトル特徴量（例えば、M F C C やメルケプストラム特徴量）を算出する。そして、算出したスペクトル特徴量の正規化処理（例えば、各次元の平均と分散を用いて正規化）を行うことで学習用音声情報 v を生成する。

生成された学習用音声情報 v は、話者情報取得部 1 1 3 によって取得された対応話者情報 s とともにパラメータ推定部 1 1 4 へ出力される。

【 0 0 7 3 】

パラメータ推定部 1 1 4 は、話者クラスタ適応型 R B M の学習処理を行う（ステップ S 3）。ここでは、学習用話者情報 s に対応した話者クラスタ c と、学習用音声情報 v を用いて各種パラメータの推定のための学習を行う。

10

【 0 0 7 4 】

次に、ステップ S 3 の詳細について、図 7 を参照して説明する。まず、図 7 に示すように、話者クラスタ適応型 R B M の確率モデルにおいて、全パラメータに任意の値を入力し（ステップ S 1 1）、音声情報推定部 1 1 4 1 に取得した学習用音声情報 v を入力し、話者情報推定部 1 1 4 2 に取得した対応話者情報 s を入力する（ステップ S 1 2）。

そして、話者情報推定部 1 1 4 2 が取得した対応話者情報 s から、話者クラスタ計算部 1 1 4 4 が話者クラスタ c を計算し、その計算した話者クラスタ c と、音声情報推定部 1 1 4 1 に取得した学習用音声情報 v を入力とする（ステップ S 1 3）。

【 0 0 7 5 】

次に、ステップ S 1 3 で入力された話者クラスタ c と学習用音声情報 v とを用いて音韻情報 h の条件付き確率密度関数を決定し、その確率密度関数に基づいて音韻情報 h をサンプリングする（ステップ S 1 4）。ここで「サンプリングする」とは、条件付き確率密度関数に従うデータをランダムに 1 つ生成することをいい、以下、同じ意味で用いる。

20

【 0 0 7 6 】

さらに、ステップ S 1 4 でサンプルされた音韻情報 h と話者クラスタ c とを用いて音声情報 v の条件付き確率密度関数を決定し、その確率密度関数に基づいて学習用音声情報 v をサンプリングする（ステップ S 1 5）。

【 0 0 7 7 】

次に、ステップ S 1 4 でサンプルされた音韻情報 h と、ステップ S 1 5 でサンプルされた学習用音声情報 v とを用いて音韻情報 h の条件付き確率密度関数を決定し、その確率密度関数に基づいて音韻情報 h を再サンプリングする（ステップ S 1 6）。

30

【 0 0 7 8 】

そして、上述の [数 1 5] 式で示される対数尤度 L をそれぞれのパラメータで偏微分し、勾配法により全パラメータを更新する（ステップ S 1 7）。具体的には、確率的勾配法が用いられ、サンプルされた学習用音声情報 v 、音韻情報 h 、および対応話者情報 s を用いてモデルに対する期待値を近似計算することができる。

【 0 0 7 9 】

全パラメータを更新した後、所定の終了条件を満たしていれば（ステップ S 1 8 の Y E S）、次のステップに進み、満たしていなければ（ステップ S 1 8 の N O）ステップ S 1 1 に戻り、以降の各ステップを繰り返す（ステップ S 1 8）。なお、所定の終了条件としては、例えば、これら一連のステップの繰り返し数が挙げられる。

40

【 0 0 8 0 】

再び、図 6 に戻り、説明を続ける。パラメータ推定部 1 1 4 は、上述の一連のステップにより推定されたパラメータを学習により決定されたパラメータとして、パラメータ記憶ユニット 1 3 に記憶する。そして、その記憶したパラメータを、入力した適応話者音声信号に基づいて、適応後のパラメータとする適用処理を行う。この適用処理で得られた適応後のパラメータを、声質変換ユニット 1 2 の声質変換部 1 2 4 へ引き渡す（ステップ S 4）。

【 0 0 8 1 】

次に、ステップ S 4 での適応処理の詳細について、図 8 を参照して説明する。まず、図

50

8に示すように、話者固有パラメータとして任意の値を入力し（ステップS21）、音声情報推定部1441に取得した適応話者音声情報 v を入力し、話者情報推定部1442に取得した適応話者情報 s を入力する（ステップS22）。

そして、話者情報推定部1442が取得した適応話者情報 s から、話者クラスタ計算部1444が話者クラスタ c を計算し、その計算した話者クラスタ c と、音声情報推定部1441に取得した適応話者音声情報 v を入力とする（ステップS23）。

【0082】

次に、ステップS23で入力された話者クラスタ c と適応話者音声情報 v とを用いて音韻情報 h の条件付き確率密度関数を決定し、その確率密度関数に基づいて音韻情報 h をサンプルする（ステップS24）。

さらに、ステップS24でサンプルされた音韻情報 h と話者クラスタ c とを用いて音声情報 v の条件付き確率密度関数を決定し、その確率密度関数に基づいて適応話者音声情報 v をサンプルする（ステップS25）。

【0083】

次に、ステップS24でサンプルされた音韻情報 h と、ステップS25でサンプルされた適応話者音声情報 v とを用いて音韻情報 h の条件付き確率密度関数を決定し、その確率密度関数に基づいて音韻情報 h を再サンプルする（ステップS26）。

【0084】

そして、上述の[数15]式で示される対数尤度 L をそれぞれのパラメータで偏微分し、勾配法により適応話者に固有のパラメータを更新する（ステップS27）。

【0085】

適応話者に固有のパラメータを更新した後、所定の終了条件を満たしていれば（ステップS28のYES）、次のステップに進み、満たしていなければ（ステップS28のNO）ステップS21に戻り、以降の各ステップを繰り返す（ステップS28）。

【0086】

再び、図6に戻り、説明を続ける。

声質変換処理として、ユーザは、図示しない入力部を操作して声質変換ユニット12の話者情報設定部123において声質変換の目標となる目標話者の情報 $s(o)$ を設定する（ステップS5）。そして、音声信号取得部121により変換用音声信号を取得する（ステップS6）。

前処理部122は、パラメータ学習処理の場合と同じく変換用音声信号に基づいて音声情報を生成し、話者情報取得部123によって取得された対応話者情報 s とともに声質変換部124へ出力される（ステップS7）。

声質変換部124は、話者クラスタ適応型RBMを適用して、適応話者の音声を目標話者の音声に変換する声質変換を行う（ステップS8）。

【0087】

次に、ステップS8の詳細について、図9を参照して説明する。まず、図9に示すように、話者クラスタ適応型RBMの確率モデルにおいて、決定された全パラメータを入力し（ステップS31）、音声情報設定部1241に音声情報 v を入力し、話者情報設定部1242に入力話者情報 s を入力し、話者クラスタ計算部1244が入力話者の話者クラスタ c を計算する（ステップS32）。

そして、ステップS32で計算された話者クラスタ c と音声情報 v とを用いて、音韻情報 h を推定する（ステップS33）。

【0088】

次に、声質変換部124は、パラメータ学習処理で学習済みの目標話者の話者情報 s を取得し、話者クラスタ計算部1244が目標話者の話者クラスタ c を計算する（ステップS34）。そして、ステップS34で計算された目標話者の話者クラスタ c とステップS33で推定した音韻情報 h とを用いて、音声情報設定部1241が変換済み音声情報 v を推定する（ステップS35）。推定された変換済み音声情報 $v(o)$ は、後処理部125へ出力される。

10

20

30

40

50

【 0 0 8 9 】

再び、図 6 に戻り、説明を続ける。後処理部 1 2 5 は、変換済み音声情報 v を用いて変換済み音声信号を生成する（ステップ S 9）。具体的には、正規化されている変換済み音声信号 v に非正規化処理（ステップ S 2 で説明した正規化処理に用いる関数の逆関数を施す処理）を施し、非正規化処理のなされたスペクトル特徴量を逆変換することでフレームごとの変換済み音声信号を生成し、これらフレームごとの変換済み音声信号を時刻順に結合することで変換済み音声信号を生成する。

後処理部 1 2 5 により生成された変換済み音声信号は、音声信号出力部 1 2 6 より外部へ出力される（ステップ S 1 0）。変換済み音声信号を外部に接続されたスピーカで再生することにより、目標話者の音声に変換された入力音声を聞くことができる。

10

【 0 0 9 0 】

[4 . 評価実験例]

次に、本発明による話者クラスタ適応型 R B M の効果を実証するため、声質変換実験を行った例について説明する。

確率モデルの学習には日本音響学会研究用連続音声データベース（ASJ-JIPDEC）の中からランダムに $R = 8 ; 1 6 ; 5 8$ 名の話者を選び、40 センテンスの音声データを用いた。学習話者の評価には、男性 1 名（ECL0001）を入力話者、女性 1 名（ECL1003）を目標話者とし、学習データとは別の 1 0 センテンスの音声データを用いた。確率モデルの適応には、学習時に含まれない女性話者（ECL1004）、男性話者（ECL0002）をそれぞれ入力話者、目標話者とし、適応データのセンテンス数を 0 . 2 から 4 0 まで変えて評価を行った。適応話者の評価についても適応データに含まれない 1 0 センテンスの音声データを用いた。分析合成ツール（WORLD：URL <http://ml.cs.yamanashi.ac.jp/world/index.html>）によって得られたスペクトルから計算した 3 2 次元のメルケプストラムを入力特徴量に用いた（ $I = 3 2$ ）。また、潜在音韻特徴量の数を $J = 8 ; 1 6 ; 2 4$ 、クラスタの数を $K = 2 ; 3 ; 4 ; 6 ; 8$ とし、最も高い精度となるものを採用した。学習率 0 : 0 1、モーメント係数 0 : 9、バッチサイズ 1 0 0 × R、繰り返し回数 1 0 0 の確率的勾配法を用いて確率モデルを学習した。

20

声質変換の精度を測る指標として、以下の [数 1 8] 式で定義される M D I R (mel-cepstral distortion improvement ratio) の平均値を用いた。

【 0 0 9 1 】

30

【 数 1 8 】

$$MDIR[dB] = \frac{10\sqrt{2}}{\ln 10} (\|v_o - v_i\|_2 - \|v_o - \hat{v}_o\|_2)$$

【 0 0 9 2 】

ここで、 v_o 、 v_i 、 v_o (^) は、それぞれ、入力話者とアライメントをとった目標話者音声のメルケプストラム特徴量、同アライメントをとった入力話者音声のメルケプストラム特徴量、 v_i に対して声質変換を施した音声のメルケプストラム特徴量を示す。M D I R は改善率を表し、値が大きいほど高い変換精度を示す。

まず、 $K = 2 ; R = 8$ および $K = 3 ; R = 1 6$ としたとき、推定された各話者のクラスタ重み w_r の分布を図 1 0 A および図 1 0 B に示す。図 1 0 A の例は、 $K = 2$ であり、男性のクラスタ (Cluster 1) と女性のクラスタ (Cluster 2) との 2 つのクラスタが自動的に形成されている。図 1 0 B の例は、 $K = 3$ であり、男性のクラスタ (Cluster 1) と女性のクラスタ (Cluster 2) の他に、さらに男女が混ざった別のクラスタ (Cluster 3) が自動的に形成されている。この図 1 0 A および図 1 0 B において、各学習者の話者クラスタの位置 R 1 1 ~ R 1 8 および R 2 1 ~ R 3 0 を示し、印で示す音声は男性の音声であり、×印で示す音声は女性の音声である。

40

【 0 0 9 3 】

図 1 0 A および図 1 0 B から分かるように、印で示す男性の音声は、(Cluster 1) に近い位置 (クラスタ重み) になり、×印で示す女性の音声は、(Cluster 2) に近い位置に

50

学習されており、性別の教師を与えていないにも関わらず、男性のクラスタ(Cluster 1)と女性のクラスタ(Cluster 2)が自動的に形成されていることが分かる。また、図10Aおよび図10Bに示すように、学習データでは、二つのクラスタが最も離れるように学習されている。すなわち、互いに最も離れている話者ペアが、それぞれのクラスタ(Cluster 1及びCluster 2)と重なる位置に設定されている。そして、各クラスタが最も離れるように学習した複数のクラスタの間で、話者クラスタへの重みの位置を設定する。このように複数のクラスタが最も離れるように学習する性質は、各クラスタ(代表話者)を内分する点を自由に調節することで任意の声へ変換する際、調節の幅が広くなり好ましい。

【0094】

次に、本発明による話者クラスタ適応型RBMによる確率モデル(CABと示す)と、従来の非パラレル声質変換手法である適応型RBM(ARBMと示す)の学習話者の変換精度を比較した例を、[表1]に示す。ここでは、学習人数が8人、16人、58人の例を示し、値が高いほど精度が高いことを示す。

【0095】

【表1】

# persons	8	16	58
ARBM	3.70	2.64	3.02
CAB	3.21	3.06	3.23

【0096】

従来の適応型RBM(ARBM)では、話者数の少ない場合には高い精度を示すが、話者数を増加させると精度が低下することが分かる。一方、話者ごとのパラメータ数を抑えた話者クラスタ適応型RBMによる確率モデル(CAB)では、話者数を増加させても精度に変化はあまり見られない。

[表2]は、本発明による話者クラスタ適応型RBMによる確率モデルと、従来の適応型RBM(ARBM)による確率モデルとの、センテンス数による変換精度を比較した例である。

【0097】

【表2】

# sent.	0.2	0.5	1	10	40
ARBM	2.48	3.25	3.21	3.41	3.45
CAB	3.14	3.54	3.63	3.60	3.58

【0098】

[表2]から明らかのように、適応に用いるセンテンス数が1以下のとき、従来モデルでは精度の低下が見られるが、話者クラスタ適応型RBMによる確率モデル(CAB)では、0.5センテンス程度で、10センテンス以上の場合と同等のパフォーマンスが得られる。

【0099】

以上、本発明によれば、話者情報から話者クラスタを取得して、その話者クラスタを使って確率モデルを得るようにしたので、従来よりも非常に少ないデータ数で、入力話者音声を目話者音声に声質変換できるようになる。

【0100】

[5. 変形例]

なお、ここまで説明した実施形態例では、目標話者の音声情報vと音韻情報nとを得る処理として、図5の話者クラスタ適応型RBMのグラフ構造に示すように、話者クラスタ

10

20

30

40

50

c が持つパラメータ A , V , U から、演算で目標話者の音声情報 v と音韻情報 n を得るようにした。

これに対して、図 1 1 に示すように、話者クラスタ c が持つパラメータ A , V , U から、目標話者の話者情報 s を得、得られた話者情報 s を使って、話者に依存したパラメータ D , A , B を得た後、これらのパラメータ D , A , B から、目標話者の音声情報 v と音韻情報 n を得るようにしてもよい。話者に依存したパラメータ D , A , B から、目標話者の音声情報 v と音韻情報 n を得る処理については、例えば図 4 の適応型 R B M のグラフ構造で説明した処理が適用可能である。

この図 1 1 に示すように、話者クラスタ c を使って目標話者の話者情報 s を得た後、目標話者の音声情報 v と音韻情報 n を得るようにすることも、図 5 の例と同様に、適切な目標話者の音声情報 v と音韻情報 n を得ることができる。この図 1 1 に示す処理を行う場合には、目標話者の音声情報 v と音韻情報 n が、目標話者の話者情報 s から得られるため、それぞれの情報の精度が向上する効果を有する。但し、データ量については、図 5 の例よりも増加する。

【 0 1 0 1 】

また、ここまで説明した実施形態例では、学習用の音声信号による学習で、声質変換のためのパラメータを学習処理した後、適応話者音声信号の入力で、パラメータを適応話者音声信号に適応した後、適応されたパラメータを使って、目標話者の音声信号に声質変換するようにした。このようにすることで、事前に学習されていない音声信号（適応話者音声信号）を、目標話者の音声信号に声質変換することができる。これに対して、適応話者音声信号の入力を省略して、学習用の音声信号で得たパラメータを使って、学習用の音声信号を目標話者の音声信号に声質変換してもよい。

この場合には、声質変換装置 1 は、例えば図 1 に示す構成として、パラメータ学習ユニット 1 1 での学習で得られたパラメータをパラメータ記憶ユニット 1 3 が記憶し、声質変換処理ユニット 1 2 は、パラメータ記憶ユニット 1 3 が記憶したパラメータを適用して、入力音声を目話者の音声に変換処理すればよい。

【 0 1 0 2 】

また、ここまで説明した実施形態例では、学習を行う入力音声（入力話者の音声）や適応を行う入力音声として、人間の話し声の音声処理する例について説明したが、実施形態例で説明した各情報を得る学習が可能であれば、学習用や適応を行う音声信号（入力信号）として、人間の話し声以外の様々な音とし、その音声信号を学習又は適応するようにしてもよい。例えば、サイレンの音や動物の鳴き声などのような音を学習又は適応するようにしてもよい。

【 符号の説明 】

【 0 1 0 3 】

1 . . . 声質変換装置、 1 1 . . . パラメータ学習ユニット、 1 2 . . . 声質変換処理ユニット、 1 3 . . . パラメータ記憶ユニット、 1 4 . . . 適応ユニット、 1 0 1 . . . CPU、 1 0 2 . . . ROM、 1 0 3 . . . RAM、 1 0 4 . . . HDD / SDD、 1 0 5 . . . 接続 I / F、 1 0 6 . . . 通信 I / F、 1 1 1 , 1 2 1 , 1 4 1 . . . 音声信号取得部、 1 1 2 , 1 2 2 , 1 4 2 . . . 前処理部、 1 1 3 . . . 対応話者情報取得部、 1 1 4 , 1 4 4 . . . パラメータ推定部、 1 1 4 1 , 1 4 4 1 . . . 音声情報推定部、 1 1 4 2 , 1 4 4 2 . . . 話者情報推定部、 1 1 4 3 , 1 4 4 3 . . . 音韻情報推定部、 1 1 4 4 , 1 4 4 4 . . . 話者クラスタ計算部、 1 2 3 . . . 話者情報設定部、 1 2 4 . . . 声質変換部、 1 2 4 1 . . . 音声情報設定部、 1 2 4 2 . . . 話者情報設定部、 1 2 4 3 . . . 音韻情報設定部、 1 2 4 4 . . . 話者クラスタ計算部、 1 2 5 . . . 後処理部、 1 2 5 . . . 音声信号出力部

10

20

30

40

【図1】

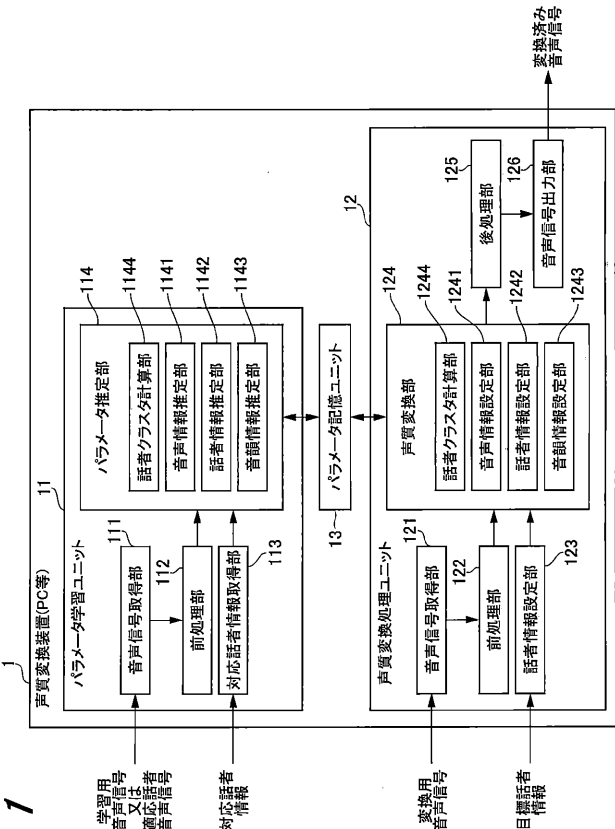
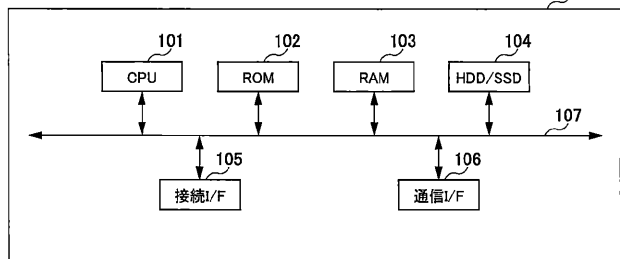


FIG. 1

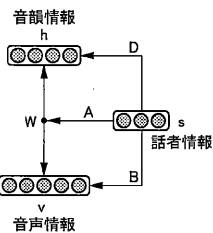
【図3】

FIG. 3



【図4】

FIG. 4



【図2】

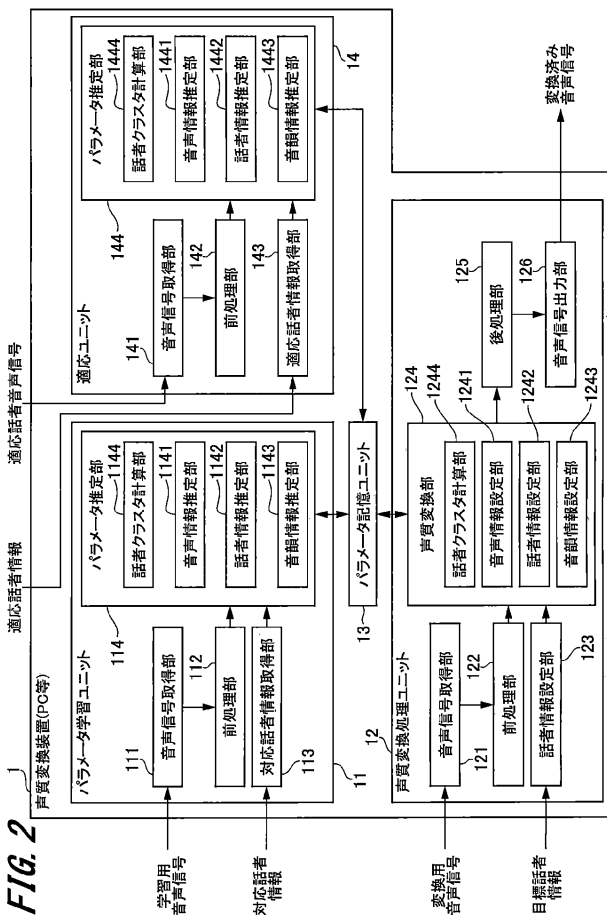
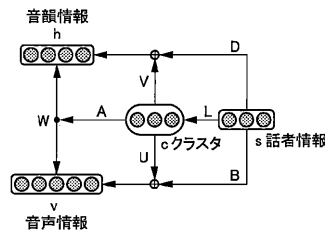


FIG. 2

【図5】

FIG. 5

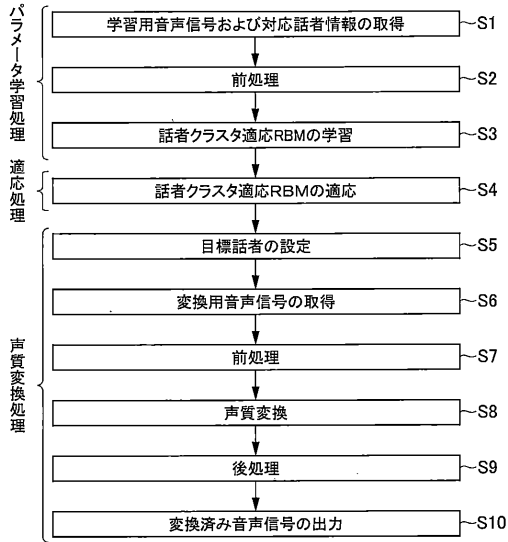


【図4】

FIG. 4

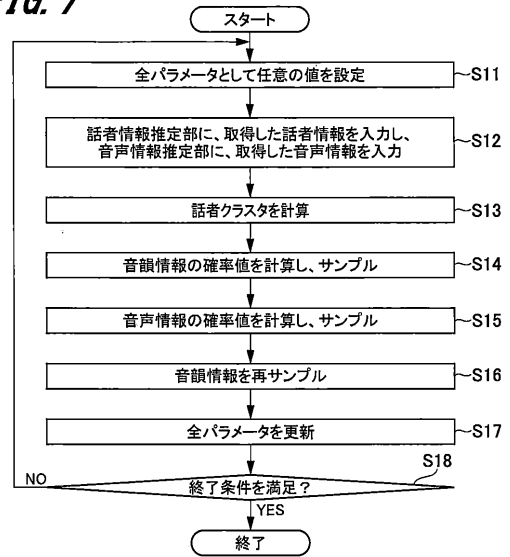
【 図 6 】

FIG. 6



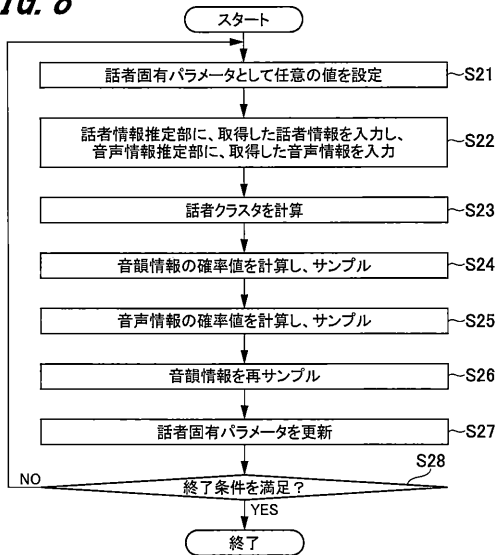
【 図 7 】

FIG. 7



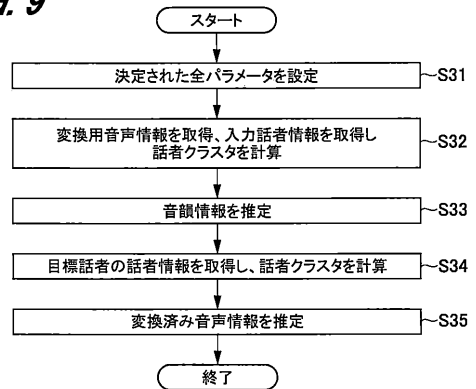
【 図 8 】

FIG. 8



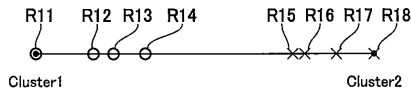
【 図 9 】

FIG. 9



【図 1 0】

FIG. 10A



【図 1 1】

FIG. 11

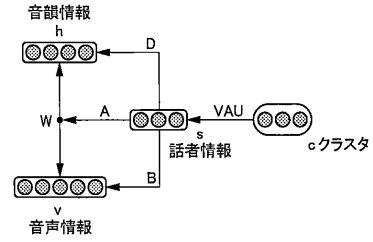
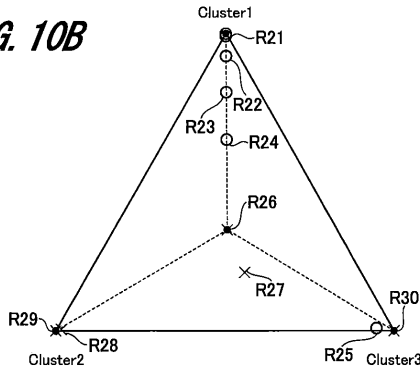


FIG. 10B



【手続補正書】

【提出日】平成30年10月22日(2018.10.22)

【手続補正 1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項 1】

入力話者の音声を目標話者の音声に声質変換する声質変換装置であって、

学習用の音声に基づく音声情報およびその音声情報に対応する話者情報から、声質変換のためのパラメータを決定するパラメータ学習ユニットと、

前記パラメータ学習ユニットが決定したパラメータを記憶するパラメータ記憶ユニットと、

前記パラメータ記憶ユニットが記憶したパラメータと前記目標話者の話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う声質変換処理ユニットとを備え、

前記パラメータ学習ユニットは、音声に基づく音声情報、音声情報に対応する話者情報および音声内の音韻を表す音韻情報のそれぞれを変数とすることで、前記音声情報、前記話者情報および前記音韻情報のそれぞれの間の結合エネルギーの関係性を前記パラメータによって表す確率モデルを取得し、前記確率モデルとして、固有の適応行列を持つ複数の話者クラスタを定義するようにし、

前記声質変換処理ユニットは、前記パラメータから前記目標話者の話者情報を得、得られた話者情報から前記目標話者の音声情報を得るようにした

声質変換装置。

【請求項 2】

さらに、前記パラメータ記憶ユニットが記憶したパラメータを前記入力話者の音声に適應して、適應後のパラメータを得る適應ユニットを備え、

前記パラメータ記憶ユニットは、前記適應ユニットで適應後のパラメータを記憶し、前記声質変換処理ユニットは、適應後のパラメータと前記目標話者の話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う

請求項 1 に記載の声質変換装置。

【請求項 3】

前記パラメータ学習ユニットと前記適應ユニットは共通の演算処理部で構成され、

前記学習用の音声に基づいてパラメータを決定する処理と、前記入力話者の音声に基づいて適應後のパラメータを得る処理を、前記共通の演算処理部で行うようにした

請求項 2 に記載の声質変換装置。

【請求項 4】

前記パラメータ学習ユニットが学習する際には、複数のクラスタが最も離れるように学習し、学習した複数のクラスタの間で、話者クラスタへの重みの位置を設定する

請求項 1 に記載の声質変換装置。

【請求項 5】

(削除)

【請求項 6】

音声情報の特徴量 $v = [v_1, \dots, v_I] \in \mathbb{R}^I$ と、音韻情報の特徴量 $h = [h_1, \dots, h_J] \in \{0, 1\}^J$ 、 $\sum_j h_j = 1$ との間に、話者情報の特徴量 $s = [s_1, \dots, s_R] \in \{0, 1\}^R$ 、 $\sum_r s_r = 1$ に依存した双方な結合重み $W \in \mathbb{R}^{I \times J}$ が存在すると仮定したとき、前記話者クラスタとして、話者クラスタ $c \in \mathbb{R}^K$ を導入し、話者クラスタ c を、

$$c \triangleq Ls$$

(但し、 $L \in \mathbb{R}^{K \times R} = [l_1 \dots l_R]$ の各列ベクトル l_r は、それぞれの話者クラスタへの重みを表す非負パラメータであり、 $\|l_r\|_1 = 1$ 、 l_r の制約を課す) と表現し、音響情報の特徴量のクラスタ依存項のバイアスパラメータを $U \in \mathbb{R}^{I \times K}$ 、音韻情報の特徴量のクラスタ依存項のバイアスパラメータを $V \in \mathbb{R}^{J \times K}$ 、として、話者非依存項、クラスタ依存項、および話者依存項のそれぞれを、

$$\tilde{W} \triangleq A \circ_3 cW$$

$$\tilde{b} \triangleq b + Uc + Bs$$

$$\tilde{d} \triangleq d + Vc + Ds$$

として示す

請求項 1 に記載の声質変換装置。

【請求項 7】

入力話者の音声を目話者の音声に声質変換する声質変換方法であって、

音声に基づく音声情報、音声情報に対応する話者情報および音声中の音韻を表す音韻情報のそれぞれを変数とすることで、前記音声情報、前記話者情報および前記音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを用意し、その確率モデルとして、固有の適應行列を持つ複数の話者クラスタを定義し、それぞれの話者について、前記複数の話者クラスタへの重みを推定して、学習用の音声についての前記パラメータを決定するパラメータ学習ステップと、

前記パラメータ学習ステップで得られたパラメータ、又は当該パラメータを前記入力話者の音声に適應した適應後のパラメータと、前記目標話者の前記話者情報とに基づいて、

前記入力話者の音声に基づく前記音声情報の声質変換処理を行う声質変換処理ステップとを含み、

前記声質変換処理ステップでの声質変換処理では、前記パラメータから前記目標話者の話者情報を得、得られた話者情報から前記目標話者の音声情報を得るようにした声質変換方法。

【請求項 8】

音声に基づく音声情報、音声情報に対応する話者情報および音声の中の音韻を表す音韻情報のそれぞれを変数とすることで、前記音声情報、前記話者情報および前記音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを用意し、その確率モデルとして、固有の適応行列を持つ複数個の話者クラスタを定義し、それぞれの話者について、前記複数個の話者クラスタへの重みを推定して、学習用の音声についての前記パラメータを決定して記憶するパラメータ学習ステップと、

前記パラメータ学習ステップで得られたパラメータ、又は当該パラメータを入力話者の音声に適応した適応後のパラメータと、目標話者の前記話者情報とに基づいて、前記入力話者の音声に基づく前記音声情報の声質変換処理を行う声質変換処理ステップと、をコンピュータに実行させるプログラムであり、

前記声質変換処理ステップでの声質変換処理は、前記パラメータから前記目標話者の話者情報を得、得られた話者情報から前記目標話者の音声情報を得るようにしたプログラム。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】0003

【補正方法】変更

【補正の内容】

【0003】

声質変換が可能な声質変換装置、声質変換方法およ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

[0009]

上記課題を解決するため、本発明の声質変換装置は、入力話者の音声を目話者の音声に声質変換する声質変換装置であって、パラメータ学習ユニットとパラメータ記憶ユニットと声質変換処理ユニットとを備える。

パラメータ学習ユニットは、学習用の音声に基づく音声情報およびその音声情報に対応する話者情報から、声質変換のためのパラメータを決定する。

パラメータ記憶ユニットは、パラメータ学習ユニットが決定したパラメータを記憶する。

声質変換処理ユニットは、パラメータ記憶ユニットが記憶したパラメータと目標話者の話者情報とに基づいて、入力話者の音声に基づく音声情報の声質変換処理を行う。

ここで、パラメータ学習ユニットは、音声に基づく音声情報、音声情報に対応する話者情報および音声の中の音韻を表す音韻情報のそれぞれを変数とすることで、音声情報、話者情報および音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデルを取得し、確率モデルとして、固有の適応行列を持つ複数個の話者クラスタを定義するようにし、声質変換処理ユニットは、パラメータから目標話者の話者情報を得、得られた話者情報から目標話者の音声情報を得るようにした。

[0010]

また、本発明の声質変換方法は、入力話者の音声を目話者の音声に声質変換する方法であって、パラメータ学習ステップと声質変換処理ステップとを含む。

パラメータ学習ステップは、音声に基づく音声情報、音声情報に対応する話者情報および音声の中の音韻を表す音韻情報のそれぞれを変数とすることで、音声情報、話者情報および音韻情報のそれぞれの間の結合エネルギーの関係性をパラメータによって表す確率モデル

ルを用意する。そして、その確率モデ

【手続補正 3】

【補正対象書類名】明細書

【補正対象項目名】0004

【補正方法】変更

【補正の内容】

【0004】

ルとして、固有の適応行列を持つ複数個の話者クラスタを定義し、それぞれの話者について、複数個の話者クラスタへの重みを推定して、学習用の音声についてのパラメータを決定する。

声質変換処理ステップは、パラメータ学習ステップで得られたパラメータ、又は当該パラメータを入力話者の音声に適応した適応後のパラメータと、目標話者の話者情報とに基づいて、入力話者の音声に基づく音声情報の声質変換処理を行う。声質変換処理ステップでの声質変換処理では、パラメータから目標話者の話者情報を得、得られた話者情報から目標話者の音声情報を得るようにした。

[0011]

また本発明のプログラムは、上述した声質変換方法のパラメータ学習ステップと声質変換処理ステップとをコンピュータに実行させるものである。

[0012]

本発明によれば、話者クラスタにより目標話者を設定することができるため、従来よりも非常に少ないデータ数で、入力話者音声を目標話者音声に声質変換できるようになる。

図面の簡単な説明

[0013]

[図1] 本発明の一実施の形態例に係る声質変換装置の構成例(例1)を示すブロック図である。

[図2] 本発明の一実施の形態例に係る声質変換装置の構成例(例2)を示すブロック図である。

[図3] 声質変換装置のハードウェア構成例を示すブロック図である。

[図4] 従来の確率モデルを模式的に示す説明図である。

[図5] 声質変換装置のパラメータ推定部が備える確率モデルを模式的に示す説明図である。

[図6] 本発明の一実施の形態例に係る処理全体の流れを示すフローチャートである。

[図7] 図6のステップS3の学習の詳細例を示すフローチャートである。

[図8] 図6のステップS4の適応の詳細例を示すフローチャートである。

[図9] 図6のステップS8の声質変換の詳細例を示すフローチャートである。

[図10] 本発明の一実施形態によるクラスタの重み分布の例を示す説明図である。

[図11] 声質変換装置のパラメータ推定部が備える確率モデルの別の例を示す説明図である。

発明を実施するための形態

【国際調査報告】

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2018/007268

A. CLASSIFICATION OF SUBJECT MATTER Int.Cl. G10L21/007(2013.01)i		
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) Int.Cl. G10L21/007		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Published examined utility model applications of Japan	1922-1996	
Published unexamined utility model applications of Japan	1971-2018	
Registered utility model specifications of Japan	1996-2018	
Published registered utility model applications of Japan	1994-2018	
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) IEEE Xplore		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X A	中鹿亘、滝口 哲也、制約付き Three-Way Restricted Boltzmann Machine を用いた音響・音韻・話者情報の同時モデリング、電子情報通信学会技術研究報告 vol. 115, no. 346, 25 November 2015, vol. 115, no. 346, pp. 7-12, (NAKASHIKA, Toru, TAKIGUCHI, Tetsuya, "Simultaneous Modelling of Acoustic, Phonetic, Speaker Features Using Improved Three-Way Restricted Boltzmann Machine IEICE technical report", IEICE technical report)	1-3, 7-8 4-6
A	JP 2016-29779 A (KDDI CORPORATION) 03 March 2016, paragraph [0005] (Family: none)	1-8
P, X	WO 2017/146073 A1 (THE UNIVERSITY OF ELECTRO-COMMUNICATIONS) 31 August 2017, claim 1 (Family: none)	1-3, 7-8
<input type="checkbox"/> Further documents are listed in the continuation of Box C.		<input type="checkbox"/> See patent family annex.
* Special categories of cited documents:		
"A" document defining the general state of the art which is not considered to be of particular relevance	"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	
"E" earlier application or patent but published on or after the international filing date	"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	
"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	"&" document member of the same patent family	
"P" document published prior to the international filing date but later than the priority date claimed		
Date of the actual completion of the international search 07 May 2018 (07.05.2018)	Date of mailing of the international search report 22 May 2018 (22.05.2018)	
Name and mailing address of the ISA/ Japan Patent Office 3-4-3, Kasumigaseki, Chiyoda-ku, Tokyo 100-8915, Japan	Authorized officer Telephone No.	

国際調査報告		国際出願番号 PCT/J P 2 0 1 8 / 0 0 7 2 6 8									
A. 発明の属する分野の分類 (国際特許分類 (IPC)) Int.Cl. G10L21/007(2013.01)i											
B. 調査を行った分野 調査を行った最小限資料 (国際特許分類 (IPC)) Int.Cl. G10L21/007											
最小限資料以外の資料で調査を行った分野に含まれるもの <table border="0"> <tr> <td>日本国実用新案公報</td> <td>1922-1996年</td> </tr> <tr> <td>日本国公開実用新案公報</td> <td>1971-2018年</td> </tr> <tr> <td>日本国実用新案登録公報</td> <td>1996-2018年</td> </tr> <tr> <td>日本国登録実用新案公報</td> <td>1994-2018年</td> </tr> </table>				日本国実用新案公報	1922-1996年	日本国公開実用新案公報	1971-2018年	日本国実用新案登録公報	1996-2018年	日本国登録実用新案公報	1994-2018年
日本国実用新案公報	1922-1996年										
日本国公開実用新案公報	1971-2018年										
日本国実用新案登録公報	1996-2018年										
日本国登録実用新案公報	1994-2018年										
国際調査で使用した電子データベース (データベースの名称、調査に使用した用語) IEEE Xplore											
C. 関連すると認められる文献											
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号									
X A	中鹿 亘、滝口 哲也, 制約付き Three-Way Restricted Boltzmann Machine を用いた音響・音韻・話者情報の同時モデリング, 電子情報通信学会技術研究報告 Vol. 115 No. 346, 2015.11.25, 第115巻、第346号, p.7-12	1-3, 7-8 4-6									
A	JP 2016-29779 A (KDD I 株式会社) 2016.03.03, [0005] (ファミリーなし)	1-8									
<input checked="" type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。											
* 引用文献のカテゴリー 「A」特に関連のある文献ではなく、一般的技術水準を示すもの 「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す) 「O」口頭による開示、使用、展示等に言及する文献 「P」国際出願日前で、かつ優先権の主張の基礎となる出願		の日の後に公表された文献 「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」同一パテントファミリー文献									
国際調査を完了した日 07.05.2018		国際調査報告の発送日 22.05.2018									
国際調査機関の名称及びあて先 日本国特許庁 (ISA/J P) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号		特許庁審査官 (権限のある職員) 冨澤 直樹	5Z 4188								
		電話番号 03-3581-1101	内線 3591								

国際調査報告		国際出願番号 PCT/JP2018/007268
C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
P, X	WO 2017/146073 A1 (国立大学法人電気通信大学) 2017.08.31, [請求項1] (ファミリーなし)	1-3, 7-8

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, 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, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT

(注) この公表は、国際事務局(WIPO)により国際公開された公報を基に作成したものである。なおこの公表に係る日本語特許出願(日本語実用新案登録出願)の国際公開の効果は、特許法第184条の10第1項(実用新案法第48条の13第2項)により生ずるものであり、本掲載とは関係ありません。