

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6783475号  
(P6783475)

(45) 発行日 令和2年11月18日(2020.11.18)

(24) 登録日 令和2年10月26日(2020.10.26)

(51) Int.Cl. F I  
G 1 O L 21/007 (2013.01) G 1 O L 21/007

請求項の数 4 (全 15 頁)

<p>(21) 出願番号 特願2018-501721 (P2018-501721)                  (86) (22) 出願日 平成29年2月22日 (2017. 2. 22)                  (86) 国際出願番号 PCT/JP2017/006478                  (87) 国際公開番号 W02017/146073                  (87) 国際公開日 平成29年8月31日 (2017. 8. 31)                  審査請求日 令和1年12月25日 (2019. 12. 25)                  (31) 優先権主張番号 特願2016-32488 (P2016-32488)                  (32) 優先日 平成28年2月23日 (2016. 2. 23)                  (33) 優先権主張国・地域又は機関                  日本国 (JP)</p>	<p>(73) 特許権者 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                  (74) 代理人 110000925                  特許業務法人信友国際特許事務所                  (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                  (72) 発明者 南 泰浩                  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                  審査官 岩田 淳</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) 【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、  
 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、前記音声情報および前記音韻情報に対応する前記話者情報を前記確率モデルに順次入力することで、前記パラメータを学習により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により決定された前記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、目標話者の音韻情報を推定し、その推定した音韻情報を使って、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、  
 を備える声質変換装置。

【請求項2】

前記パラメータは、前記音声情報と前記音韻情報との関係性の度合いを表すM、前記音韻情報と前記話者情報との関係性の度合いを表すV、前記話者情報と前記音声情報との関係性の度合いを表すU、前記話者情報によって決定される射影行列集合A、前記音声情報のバイアスb、前記音韻情報のバイアスc、および前記音声情報の偏差の7つのパラメータからなり、

これら7つのパラメータは、前記音声情報をv、前記音韻情報をh、前記話者情報をsとすることで、以下の(A)式~(D)式によって関係付けられ、

$$\begin{aligned}
 E(\mathbf{v}, \mathbf{h}, \mathbf{s}) &= \frac{1}{2} \mathbf{v}^T \bar{\mathbf{v}} - \mathbf{b}^T \bar{\mathbf{v}} - \mathbf{c}^T \mathbf{h} - \mathbf{h}^T \mathbf{V} \mathbf{s} - \mathbf{s}^T \mathbf{U} \bar{\mathbf{v}} - \bar{\mathbf{v}}^T \mathbf{A}_s \mathbf{M} \mathbf{h}, \quad \dots\dots (\text{A}) \\
 p(\mathbf{v} | \mathbf{h}, \mathbf{s}) &= \mathcal{N}(\mathbf{v} | \mathbf{b} + \mathbf{U}^T \mathbf{s} + \mathbf{A}_s \mathbf{M} \mathbf{h}, \sigma^2) \quad \dots\dots (\text{B}) \\
 p(\mathbf{h} | \mathbf{s}, \mathbf{v}) &= \mathcal{B}(\mathbf{h} | \mathbf{f}(\mathbf{c} + \mathbf{V} \mathbf{s} + \mathbf{M}^T \mathbf{A}_s^T \bar{\mathbf{v}})) \quad \dots\dots (\text{C}) \\
 p(\mathbf{s} | \mathbf{v}, \mathbf{h}) &= \mathcal{B}(\mathbf{s} | \mathbf{f}(\mathbf{U} \bar{\mathbf{v}} + \mathbf{V}^T \mathbf{h} + [\bar{\mathbf{v}}^T \mathbf{A}_k] \mathbf{M} \mathbf{h})) \quad \dots\dots (\text{D})
 \end{aligned}$$

さらに、前記声質変換処理ユニットが音韻情報を推定する際には、前記7つのパラメータの内の少なくとも、前記音声情報と前記音韻情報との関係性の度合いを表すM、前記音韻情報と前記話者情報との関係性の度合いを表すV、前記話者情報と前記音声情報との関係性の度合いを表すU、および前記音韻情報のバイアスcを使った式で推定するようにした

10

請求項1に記載の声質変換装置。

【請求項3】

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に、前記音声情報および前記音声情報に対応する前記話者情報を前記確率モデルに順次入力することで、前記パラメータを学習により決定するパラメータ学習ステップと、

20

前記パラメータ学習ステップにより決定された前記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、目標話者の音韻情報を推定し、その推定した音韻情報を使って、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4】

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に、前記音声情報および前記音声情報に対応する前記話者情報を前記確率モデルに順次入力することで、前記パラメータを学習により決定するパラメータ学習ステップと、

30

前記パラメータ学習ステップにより決定された前記パラメータと目標話者の前記話者情報とに基づいて、目標話者の音韻情報を推定し、その推定した音韻情報を使って、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】

40

【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

50

う。

【0003】

これに対して、モデルの学習時に上述の平行データを使用しない非平行声質変換が注目を浴びている。非平行声質変換は、平行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ものである。

【先行技術文献】

【特許文献】

10

【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.

20

【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では、平行データを必要とする平行声質変換と比較して、平行データを必要としない分利便性や実用性が高いが、事前に入力話者の音声を学習させておく必要があるという問題がある。また、変換時において事前に入力話者を指定する必要があり、入力話者の音声を問わず特定話者の音声を出力したいという要求を満たすことはできないという問題があった。

【0007】

本発明は、上記従来の問題点に鑑み提案されたものであり、その目的とするところは、予め入力話者を特定しなくとも目標話者の声質へ声質変換を可能とすることにある。

30

【課題を解決するための手段】

【0008】

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、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、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、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し、音声情報および音声情報に対応する話者情報を確率モデルに順次入力することで、パラメータを学習により決定する。

40

声質変換処理ユニットは、パラメータ学習ユニットにより決定されたパラメータと目標話者の話者情報とに基づいて、目標話者の音韻情報を推定し、その推定した音韻情報を使って、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。

【発明の効果】

【0009】

本発明によれば、話者を考慮しつつ音声のみから音韻を推定することができるため、入力話者を特定しなくとも目標話者への声質変換が可能となる。

【図面の簡単な説明】

【0010】

50

【図 1】本発明の一実施形態にかかる声質変換装置の構成例を示すブロック図である。

【図 2】図 1 のパラメータ推定部が備える確率モデル Three - Way R B M (Restricted Boltzmann machine) を模式的に示す図である。

【図 3】図 1 の声質変換装置のハードウェア構成例を示す図である。

【図 4】実施形態の処理例を示すフローチャートである

【図 5】図 4 の前処理の詳細例を示すフローチャートである。

【図 6】図 4 の確率モデル 3 W R B M による学習の詳細例を示すフローチャートである。

【図 7】図 4 の声質変換の詳細例を示すフローチャートである。

【図 8】図 4 の後処理の詳細例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】

10

【 0 0 1 1 】

以下、本発明の好適な実施形態について説明する。

【 0 0 1 2 】

< 構成 >

図 1 は本発明の一実施形態にかかる声質変換装置の構成例を示す図である。図 1 において P C 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を目標とする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、マイク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。

20

【 0 0 1 3 】

声質変換装置 1 は、パラメータ学習ユニット 1 1 と声質変換処理ユニット 1 2 とを備える。パラメータ学習ユニット 1 1 は、学習用音声信号と対応話者情報とに基づいて学習により声質変換のためのパラメータを決定するものである。また、声質変換処理ユニット 1 2 は、上述の学習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ものである。

30

【 0 0 1 4 】

パラメータ学習ユニット 1 1 は、音声信号取得部 1 1 1 と前処理部 1 1 2 と話者情報取得部 1 1 3 とパラメータ推定部 1 1 4 を備える。音声信号取得部 1 1 1 は、前処理部 1 1 2 に接続され、前処理部 1 1 2 および話者情報取得部 1 1 3 は、それぞれパラメータ推定部 1 1 4 に接続される。

【 0 0 1 5 】

音声信号取得部 1 1 1 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 1 1 1 は、マイクロフォンに接続され、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。

40

前処理部 1 1 2 は、音声信号取得部 1 1 1 が取得し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、M F C C (Mel-Frequency Cepstrum Coefficients : 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。

【 0 0 1 6 】

対応話者情報取得部 1 1 3 は、音声信号取得部 1 1 1 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話者情報取得部

50

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が入力中かを区別する情報(対応話者情報)を、ユーザの入力又は自動的に取得する。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。

#### 【0017】

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確率モデルThree-Way RBM(3WRBM)を備える。

音声情報推定部1141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。音声情報は、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)である。

話者情報推定部1142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。話者情報は、話者を特定するための情報であり、それぞれの話者の音響が持つ話者ベクトルの情報である。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である。

音韻情報推定部1143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。音韻情報は、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。

また、パラメータ推定部114が備える確率モデル3WRBMとしては、各推定部1141, 1142, 1143が推定した3つの情報(音声情報、話者情報、音韻情報)を持つことになるが、確率モデル3WRBMでは、これら音声情報、話者情報、音韻情報を持つだけでなく、3つの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。

これら音声情報推定部1141、話者情報推定部1142および音韻情報推定部1143、音声情報、話者情報および音韻情報、各種パラメータ並びに確率モデル3WRBMについての詳細については後述する。

#### 【0018】

声質変換処理ユニット12は、音声信号取得部121と前処理部122と話者情報設定部123と声質変換部124と後処理部125と音声信号出力部126とを備える。音声信号入力121、前処理部122、声質変換部124、後処理部125および音声信号出力部126は順次接続され、声質変換部124には、更にパラメータ学習ユニット11のパラメータ推定部114が接続される。

#### 【0019】

音声信号取得部121は、変換用音声信号を取得し、前処理部122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121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。つまり、事前に学習がされていない話者の話し声が、音声信号取得部121に供給される。

音声信号取得部121および前処理部122は、上述したパラメータ学習ユニット11の音声信号取得部111および前処理部112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。

#### 【0020】

話者情報設定部123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ものである。話者情報設定部123で設定する目標話者は、ここでは、パラメータ学習ユ

10

20

30

40

50

ユニット 11 のパラメータ推定部 114 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 123 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 114 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。

【0021】

声質変換部 124 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 124 は、音声情報設定部 1241、話者情報設定部 1242 および音韻情報設定部 1243 を持つ。この音声情報設定部 1241、話者情報設定部 1242 および音韻情報設定部 1243 は、上述のパラメータ推定部 114 において、確率モデル 3WRBM が持つ音声情報推定部 1141、話者情報推定部 1142 および音韻情報推定部 1143 と同等の機能を持つ。すなわち、音声情報設定部 1241、話者情報設定部 1242 および音韻情報設定部 1243 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 1243 に設定される音韻情報は、前処理部 122 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 1242 に設定される話者情報は、パラメータ学習ユニット 11 内の話者情報推定部 1142 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 1241 に設定される音声情報は、これら話者情報設定部 1242 および音韻情報設定部 1243 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。

なお、図 1 では声質変換部 124 を設ける構成を示したが、声質変換部 124 を別途設置することなく、パラメータ推定部 114 の各種パラメータを固定することで、パラメータ推定部 114 が声質変換の処理を実行する構成としてもよい。

【0022】

後処理部 125 は、声質変換部 124 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 FFT 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 126 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。

【0023】

図 2 はパラメータ推定部 114 の備える確率モデル 3WRBM を模式的に示す図である。確率モデル 3WRBM は、上述のとおり、音声情報推定部 1141、話者情報推定部 1142 および音韻情報推定部 1143 を備え、これらが音声情報  $v$ 、話者情報  $s$  および音韻情報  $h$  のそれぞれを変数とする以下の 3 変数同時確率密度関数の (1) 式で表現される。なお、話者情報  $s$  と音韻情報  $h$  はバイナリベクトルであり、諸要素がオン（アクティブ）になっている状態を 1 で表す。

【0024】

【数 1】

$$p(\mathbf{v}, \mathbf{h}, \mathbf{s}) = \frac{1}{N} e^{-E(\mathbf{v}, \mathbf{h}, \mathbf{s})} \quad \dots\dots(1)$$

$$\mathbf{v} = [v_1, \dots, v_D] \in \mathbb{R}^D$$

$$\mathbf{s} = [s_1, \dots, s_R] \in \{0, 1\}^R, \sum_k s_k = 1$$

$$\mathbf{h} = [h_1, \dots, h_H] \in \{0, 1\}^H, \sum_j h_j = 1$$

【0025】

ここで、(1) 式の  $E$  は音声モデリングのためのエネルギー関数であり、 $N$  は正規化項

10

20

30

40

50

である。ここでエネルギー関数  $E$  は、以下の (2) ~ (5) 式に示されるように、音声情報と音韻情報との関係性の度合いを表す  $M$ 、音韻情報と話者情報との関係性の度合いを表す  $V$ 、話者情報と音声情報との関係性の度合いを表す  $U$ 、更に  $M$  を線形変換する、話者情報  $s$  によって決定される射影行列集合  $A$ 、音声情報のバイアス  $b$ 、音韻情報のバイアス  $c$ 、音声情報の偏差  $\sigma^2$  の 7 つのパラメータ ( $\theta = \{ M, A, U, V, b, c, \sigma^2 \}$ ) によって関係付けられる。

【0026】

【数2】

$$E(v, h, s) = \frac{1}{2} v^T \bar{v} - b^T \bar{v} - c^T h - h^T V s - s^T U \bar{v} - \bar{v}^T A_s M h, \quad \dots\dots(2)$$

10

【0027】

ただし、 $A_s = \sum_k A_k s_k$ 、 $M = [m_1, \dots, m_H]$  とし、便宜上  $A = \{ A_k \}$  とする。また、 $\bar{v}$  は、 $v$  を要素ごとにパラメータ  $\sigma^2$  で除算したベクトルを表す。なお、本明細書中に示す「 $\bar{v}$ 」の「 $\bar{\cdot}$ 」は、上述の (2) 式に示すように、本来は「 $\cdot$ 」が「 $v$ 」の上に付加されるものであるが、本明細書では記載上の制約から「 $\bar{v}$ 」と記載することとする。なお、 $\bar{s}$ 、 $\bar{h}$  の「 $\bar{\cdot}$ 」、および  $h^\wedge$  の「 $^\wedge$ 」も、本来は文字の上に付加されるものであるが同様の理由により、明細書中では上述のとおり記載している。

20

このときそれぞれの条件付き確率は、以下の (3) ~ (5) 式となる。

【0028】

【数3】

$$p(v|h, s) = \mathcal{N}(v | b + U^T s + A_s M h, \sigma^2) \quad \dots\dots(3)$$

$$p(h|s, v) = \mathcal{B}(h | f(c + V s + M^T A_s^T \bar{v})) \quad \dots\dots(4)$$

$$p(s|v, h) = \mathcal{B}(s | f(U \bar{v} + V^T h + [\bar{v}^T A_k] M h)) \quad \dots\dots(5)$$

【0029】

ここで  $N$  は次元独立の多変量正規分布、 $B$  は多次元ベルヌーイ分布、 $f$  は要素ごとの  $\text{softmax}$  関数を表す。

30

上述の (1) ~ (5) 式において、 $R$  人の話者による  $T$  フレームの音声情報に対する対数尤度を最大化するように各種パラメータを推定する。なお、各種パラメータ推定の詳細は後述する。

【0030】

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。図3に示すように、声質変換装置1は、バス107を介して相互に接続されたCPU(Central Processing Unit)101、ROM(Read Only Memory)102、RAM(Random Access Memory)103、HDD(Hard Disk Drive)/SSD(Solid State Drive)104、接続I/F(Interface)105、通信I/F106を備える。CPU101は、RAM103をワークエリアとしてROM102またはHDD/SSD104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/F105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。

40

音声信号の入出力ならびに話者情報の入力および設定は、接続I/F105または通信I/F106を介して行われる。図1で説明した声質変換装置1の機能は、CPU101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、ROM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、ASIC(Application Specific Integrated Circuit)やFPGA(Field Programmable Gate Array)

50

)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。

【0031】

<動作>

図4は、上述の実施形態の処理例を示すフローチャートである。図4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)。

以下、ステップS2の詳細については、図5を参照して説明する。図5に示すように、前処理部112は、学習用音声信号をフレームごと(例えば、5msごと)に切り出し(ステップS21)、切り出された学習用音声信号にFFT処理などを施すことでスペクトル特徴量(例えば、MFCCやメルケプストラム特徴量)を算出する(ステップS22)。そして、ステップS22で得られ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(ステップS23)。

学習用音声情報 $v$ は、話者情報取得部113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114へ出力される。

【0032】

パラメータ推定部114は、確率モデル3WRBMにおいて、学習用音声情報 $v$ と対応話者情報 $s$ を用いて各種パラメータ( $M$ 、 $V$ 、 $U$ 、 $A$ 、 $b$ 、 $c$ 、)の推定のための学習を行う(ステップS3)。

$R$ 人( $R \geq 2$ )の話者による $T$ フレームの音声データ(学習用音声情報と対応話者情報との組) $X = \{v_t, s_t\}^T_{t=1}$ に対する対数尤度 $L$ 、以下(6)式を最大化するように各種パラメータ $M$ 、 $V$ 、 $U$ 、 $A$ 、 $b$ 、 $c$ 、を推定する。なお、 $t$ は時刻 $t$ を表し、 $v_t$ 、 $s_t$ 、 $h_t$ はそれぞれ時刻 $t$ における音声情報、話者情報、音韻情報を表す。

【0033】

【数4】

$$\mathcal{L} = \log p(\mathbf{X}) = \sum_t \log \sum_h p(v_t, h_t, s_t) \dots\dots\dots (6)$$

【0034】

次に、ステップS3の詳細について、図6を参照して説明する。まず、図6に示すように、確率モデル3WRBMにおいて、各種パラメータ $M$ 、 $V$ 、 $U$ 、 $A$ 、 $b$ 、 $c$ 、に任意の値を入力し(ステップS31)、音声情報推定部1141に学習用音声情報 $v$ を入力し、話者情報推定部1142に対応話者情報 $s$ を入力する(ステップS32)。

そして、上述の(4)式により、学習用音声情報 $v$ と対応話者情報 $s$ 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S33)。ここで「サンプルする」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。

【0035】

次に、サンプルされた音韻情報 $h$ と上述の学習用音声情報 $v$ とを用いて上述の(5)式により対応話者情報 $s$ 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話者情報 $\tilde{s}$ をサンプルする。そして、サンプルされた音韻情報 $h$ とサンプルされた対応話者情報 $\tilde{s}$ とを用いて上述の(3)式により学習用音声情報 $v$ 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 $\tilde{v}$ をサンプルする(ステップS34)。

次に、上述のステップS34でサンプルされた対応話者情報 $\tilde{s}$ と学習用音声情報 $\tilde{v}$ 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

10

20

30

40

50



韻情報  $\tilde{h}$  を再サンプルする (ステップ S 3 5)。

【 0 0 3 6 】

そして、上述の ( 6 ) 式で示される対数尤度  $\mathcal{L}$  をそれぞれのパラメータで偏微分し、勾配法により各種パラメータを更新する (ステップ S 3 6)。具体的には、確率的勾配法が用いられ、対数尤度  $\mathcal{L}$  をそれぞれのパラメータで偏微分した以下の ( 7 ) ~ ( 1 3 ) 式が用いられる。ここで、各偏微分項右辺の  $\langle \cdot \rangle_{\text{data}}$  はそれぞれのデータに対する期待値を表し、 $\langle \cdot \rangle_{\text{model}}$  は、モデルの期待値を表している。モデルに対する期待値は項数が膨大となり計算困難だが、CD (Contrastive Divergence) 法を適用し、上述のとおりサンプルされた学習用音声情報  $\tilde{v}$ 、対応話者情報  $\tilde{s}$ 、音韻情報  $\tilde{h}$  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。

10

【 0 0 3 7 】

【数 5】

$$\frac{\partial \mathcal{L}}{\partial \mathbf{M}} = \left\langle \sum_k \mathbf{A}_k^T \tilde{\mathbf{v}} \tilde{\mathbf{h}}^T \mathbf{s}_k \right\rangle_{\text{data}} - \left\langle \sum_k \mathbf{A}_k^T \tilde{\mathbf{v}} \tilde{\mathbf{h}}^T \mathbf{s}_k \right\rangle_{\text{model}} \quad \dots\dots (7)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}_k} = \langle \tilde{\mathbf{v}} \tilde{\mathbf{h}}^T \mathbf{s}_k \mathbf{M}^T \rangle_{\text{data}} - \langle \tilde{\mathbf{v}} \tilde{\mathbf{h}}^T \mathbf{s}_k \mathbf{M}^T \rangle_{\text{model}} \quad \dots\dots (8)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}} = \langle \tilde{\mathbf{s}} \tilde{\mathbf{v}}^T \rangle_{\text{data}} - \langle \tilde{\mathbf{s}} \tilde{\mathbf{v}}^T \rangle_{\text{model}} \quad \dots\dots (9)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{V}} = \langle \tilde{\mathbf{h}} \tilde{\mathbf{s}}^T \rangle_{\text{data}} - \langle \tilde{\mathbf{h}} \tilde{\mathbf{s}}^T \rangle_{\text{model}} \quad \dots\dots (10)$$

20

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = \langle \tilde{\mathbf{v}} \rangle_{\text{data}} - \langle \tilde{\mathbf{v}} \rangle_{\text{model}} \quad \dots\dots (11)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{c}} = \langle \tilde{\mathbf{h}} \rangle_{\text{data}} - \langle \tilde{\mathbf{h}} \rangle_{\text{model}} \quad \dots\dots (12)$$

$$\frac{\partial \mathcal{L}}{\partial \sigma} = \frac{1}{\sigma^3} \circ \left( \langle \mathbf{v} \circ \mathbf{v} - 2\mathbf{v} \circ (\mathbf{b} + \mathbf{U}^T \mathbf{s} + \mathbf{A}_s \mathbf{M} \mathbf{h}) \rangle_{\text{data}} \right. \\ \left. - \langle \mathbf{v} \circ \mathbf{v} - 2\mathbf{v} \circ (\mathbf{b} + \mathbf{U}^T \mathbf{s} + \mathbf{A}_s \mathbf{M} \mathbf{h}) \rangle_{\text{model}} \right), \quad \dots\dots (13)$$

【 0 0 3 8 】

30

各種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 (YES)、次のステップに進み、満たしていなければ (NO) ステップ S 3 2 に戻り、以降の各ステップを繰り返す (ステップ S 3 7)。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。

なお、学習処理として、一度各種パラメータを決定したあと、新たに別の人のパラメータを追加する場合には、一部の式で示すパラメータのみを更新するようにしてもよい。例えば、[数 5] で示す ( 7 ) 式 ~ ( 1 3 ) 式の中で、( 8 ) 式、( 9 ) 式、および ( 1 0 ) 式により、新たに得た学習音声で当該パラメータを更新する。( 7 ) 式、( 1 1 ) 式、( 1 2 ) 式、および ( 1 3 ) 式で得られるパラメータについては、既に学習済みのパラメータを更新せずにそのまま使用してもよく、また、他のパラメータと同様にパラメータ

40

【 0 0 3 9 】

再び、図 4 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 4 0 】

次に、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 1 2 の話者情報設定部 1 2 3 において声質変換の目標となる目標話者の情報  $s^{(\circ)}$  を設定する (ステップ S 5)。そして、音声信号取得部 1 2 1 により変換用音声信号を取得す

50

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変換用音声情報 $v^{(i)}$ を生成し、上述の対応する目標話者情報 $s^{(o)}$ とともに声質変換部124へ出力する（ステップS7）。なお、変換用音声信号 $v^{(i)}$ の生成は、上述のステップS2（図5のステップS21～S23）と同様の手順で行われる。

【0041】

声質変換処理部124は、目標話者情報 $s^{(o)}$ に基づいて変換用音声情報 $v^{(i)}$ から変換済み音声情報 $v^{(o)}$ を生成する（ステップS8）。

ステップS8の詳細は図7に示されている。以下、図7を参照してステップS8について具体的に説明する。まず、確率モデル3WRBMにおいてパラメータ学習ユニット11のパラメータ推定部114から取得した各種パラメータを設定する（ステップS81）。そして、前処理部122から変換音声情報を取得し（ステップS82）、以下の(14)式に入力することで音韻情報 $\hat{h}$ を推定する（ステップS83）。

続いて、話者情報設定部123での設定に基づいて、パラメータ学習処理で学習済みの目標話者の話者情報 $s^{(o)}$ を設定する（ステップS84）。なお、以下の(14)式の三行目、分母に用いられる $h'$ 、 $s'$ は、分子に用いられる $h$ 、 $s$ と計算上区別するために用いられるものであり、その意味は $h$ 、 $s$ と同じである。

【0042】

【数6】

$$\begin{aligned} \hat{h} &\triangleq \mathbb{E}[h|v^{(i)}] && \dots\dots(14) \\ &= [p(h_j = 1|v^{(i)})] \\ &= \left[ \frac{\sum_s p(v^{(i)}, h_j = 1, s)}{\sum_{h'} \sum_{s'} p(v^{(i)}, h', s')} \right] \\ &= f(\mathbf{c} + \mathbf{g}(\mathbf{V} + \bar{v}^{(i)\top} \mathbf{U}^\top + \mathbf{M}^\top [\mathbf{A}_k^\top \bar{v}^{(i)}])), \end{aligned}$$

【0043】

そして、算出された音韻情報 $\hat{h}$ を用いて、以下の(15)式により変換済み音声情報 $v^{(o)}$ を推定する（ステップS85）。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。

【0044】

【数7】

$$\begin{aligned} \hat{v}^{(o)} &\triangleq \operatorname{argmax}_{v^{(o)}} p(v^{(o)}|v^{(i)}, s^{(o)}) && \dots\dots(15) \\ &= \operatorname{argmax}_{v^{(o)}} \sum_h p(h|v^{(i)}, s^{(o)}) p(v^{(o)}|h, v^{(i)}, s^{(o)}) \\ &\simeq \operatorname{argmax}_{v^{(o)}} p(\hat{h}|v^{(i)}, s^{(o)}) p(v^{(o)}|\hat{h}, v^{(i)}, s^{(o)}) \\ &= \operatorname{argmax}_{v^{(o)}} p(v^{(o)}|\hat{h}, s^{(o)}) \\ &= \mathbf{b} + \mathbf{U}_{\alpha}^\top + \mathbf{A}_o \mathbf{M} \hat{h}, \end{aligned}$$

【0045】

図4に戻り、後処理部125は、変換済み音声情報 $v^{(o)}$ を用いて変換済み音声信号を生成する（ステップS9）。具体的には、図8に示すように、正規化されている変換済み音声信号 $v^{(o)}$ に非正規化処理（上述の正規化処理に用いる関数の逆関数を施す処理）を施し（ステップS91）、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（ステップS92）、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（ステップS93）。

図4に示すように、後処理部125により生成された変換済み音声信号は、音声信号出

10

20

30

40

50

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。

【 0 0 4 6 】

以上、本発明によれば、確率モデル 3 W R B M により話者情報を考慮しながら音声情報のみから音韻情報を推定することができるため、声質変換の際、入力話者を指定しなくとも目標話者への声質変換が可能となり、また、入力話者の音声学習時において学習のために用意されていない音声であったとしても目標話者の声質へ変換することが可能となる。

【 0 0 4 7 】

< 実験例 >

本発明の効果を実証するため、[ 1 ] 従来の非パラレル声質変換と本発明との変換精度を比較する実験と、[ 2 ] 本発明による話者非指定型と話者指定型の変換精度を比較する実験を行った。

実験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男性 2 7 名、女性 3 1 名の計 5 8 名の話者を選び、5 発話分の音声データを学習に用いるとともに、他の 1 0 発話分の音声データを評価に用いた。スペクトル特徴量としては、3 2 次元のメルケプストラム特徴量を用いた。また、音韻情報の次元数を 1 6 とした。評価尺度には客観評価基準である M D I R (mel-distortion improvement ratio) を用いた。

以下、( 1 6 ) 式は、実験に用いた M D I R を示す式であり、数値が大きいほど高い精度を表す。学習率 0 . 0 1、モーメント係数 0 . 9、バッチサイズ 1 0 0、繰り返し回数 5 0 の確率的勾配法を用いてモデルを学習した。

【 0 0 4 8 】

【 数 8 】

$$MDIR[dB] = \frac{10\sqrt{2}}{\ln 10} (\|v^{(o)} - v^{(t)}\|^2 - \|v^{(o)} - \hat{v}^{(o)}\|^2) \dots\dots(16)$$

【 0 0 4 9 】

【 表 1 】

Method	ARBM	SATBM	Proposed
MDIR [dB]	2.11	2.66	<b>3.07</b>

【 0 0 5 0 】

【 表 2 】

	MDIR [dB]
Correct speaker specified	3.07
Different speaker specified	2.79
Arbitrary source approach	3.03

【 0 0 5 1 】

[ 実験結果 ]

まず、本発明による 3 W R B M による声質変換と、従来の非パラレル声質変換法である A R B M (Adaptive Restricted Boltzmann Machine) 及び S A T B M (Speaker Adaptive Trainable Boltzmann Machine) と比較した。上述の [ 表 1 ] に示すように、本発明による手法で最も高い精度が得られたことが分かる。

次に、本発明で述べた 3 W R B M において、話者非指定型と、話者指定型の変換精度を比較した。実験結果を上述の [ 表 2 ] に示す。本発明において、話者非指定型 (arbitrary source approach) は入力話者を指定していないにもかかわらず、正しい入力話者を指

10

20

30

40

50

定した場合 (correct speaker specified) と遜色ない結果が得られた。なお、正しくない入力話者を指定した場合 (different speaker specified)、精度が下がることを確認した。

【0052】

<変形例>

なお、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)として、人間の話し声の音声を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、学習用の音声信号(入力信号)として、人間の話し声以外の様々な音として、その音声信号を学習してもよい。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してもよい。

10

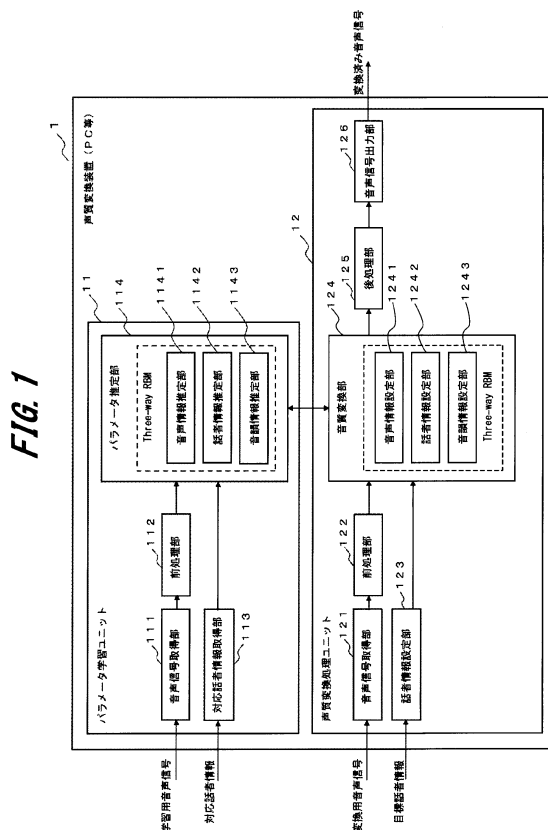
【符号の説明】

【0053】

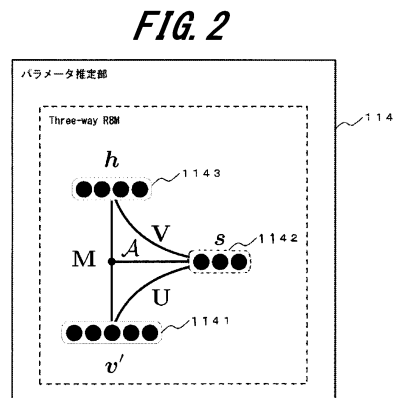
1・・・音質変換装置、11・・・パラメータ学習ユニット、12・・・音質変換処理ユニット、101・・・CPU、102・・・ROM、103・・・RAM、104・・・HDD/SDD、105・・・接続I/F、106・・・通信I/F、111,121・・・音声信号取得部、112,122・・・前処理部、113・・・対応話者情報取得部、114・・・パラメータ推定部、1141・・・音声情報推定部、1142・・・話者情報推定部、1143・・・音韻情報推定部、123・・・話者情報設定部、124・・・声質変換部、1241・・・音声情報設定部、1242・・・話者情報設定部、1243・・・音韻情報設定部、125・・・後処理部、125・・・音声信号出力部

20

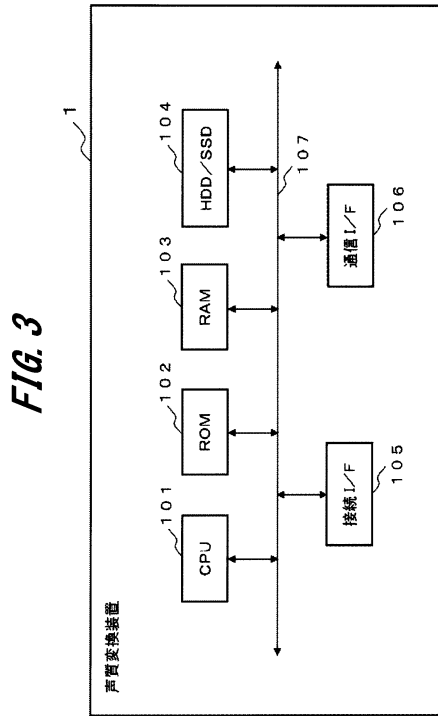
【図1】



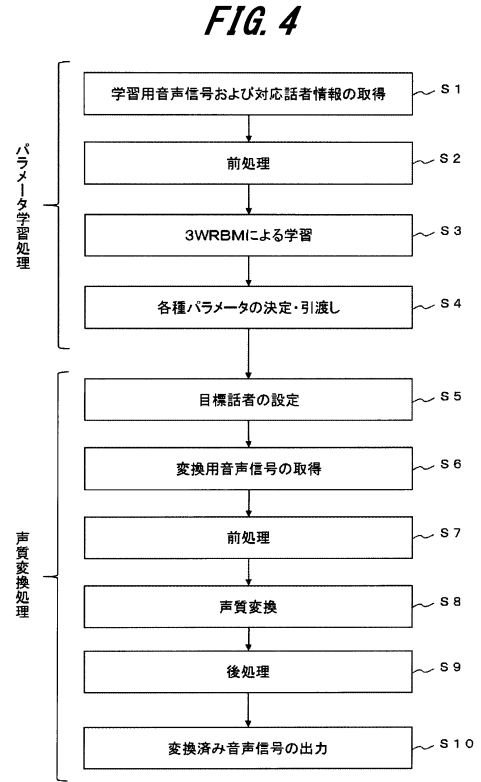
【図2】



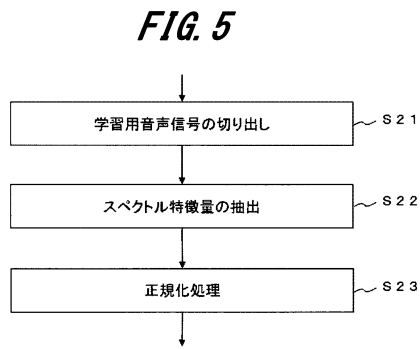
【図3】



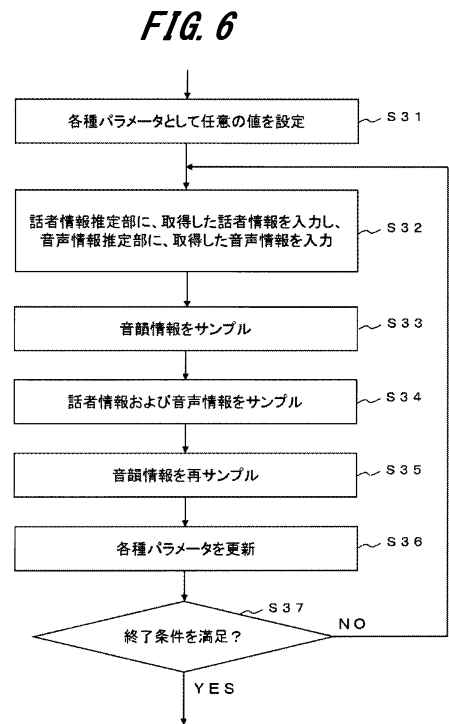
【図4】



【図5】

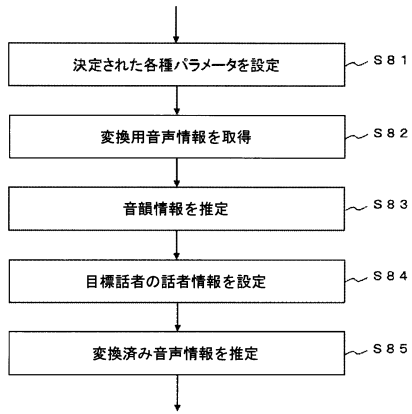


【図6】



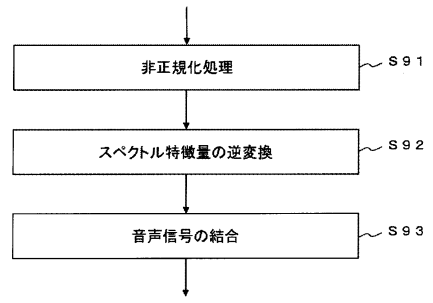
【図7】

**FIG. 7**



【図8】

**FIG. 8**



---

フロントページの続き

- (56)参考文献 特開2015-040903(JP,A)  
特開2010-014913(JP,A)  
中鹿 亘,外1名,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,情報処理学会 研究報告  
音声言語情報処理(SLP), [online], 日本, 情報処理学会, 2015年12月 3日, 2015-SLP-  
109, SP2015-71, 第1-6頁

(58)調査した分野(Int.Cl., DB名)

G10L 13/00 - 13/10  
G10L 19/00 - 99/00