

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5194197号  
(P5194197)

(45) 発行日 平成25年5月8日(2013.5.8)

(24) 登録日 平成25年2月8日(2013.2.8)

|                          |                       |
|--------------------------|-----------------------|
| (51) Int.Cl.             | F I                   |
| G 1 0 L 21/007 (2013.01) | G 1 0 L 21/04 1 2 0 D |
| G 1 0 L 25/15 (2013.01)  | G 1 0 L 11/00 1 0 1 D |
| G 1 0 L 25/75 (2013.01)  | G 1 0 L 11/00 5 0 1   |

請求項の数 17 (全 39 頁)

|               |                              |           |                     |
|---------------|------------------------------|-----------|---------------------|
| (21) 出願番号     | 特願2012-551826 (P2012-551826) | (73) 特許権者 | 000005821           |
| (86) (22) 出願日 | 平成24年7月12日 (2012.7.12)       |           | パナソニック株式会社          |
| (86) 国際出願番号   | PCT/JP2012/004517            |           | 大阪府門真市大字門真1006番地    |
| (87) 国際公開番号   | W02013/008471                | (74) 代理人  | 100109210           |
| (87) 国際公開日    | 平成25年1月17日 (2013.1.17)       |           | 弁理士 新居 広守           |
| 審査請求日         | 平成24年11月14日 (2012.11.14)     | (72) 発明者  | 釜井 孝浩               |
| (31) 優先権主張番号  | 特願2011-156042 (P2011-156042) |           | 日本国大阪府門真市大字門真1006番地 |
| (32) 優先日      | 平成23年7月14日 (2011.7.14)       |           | パナソニック株式会社内         |
| (33) 優先権主張国   | 日本国(JP)                      | (72) 発明者  | 廣瀬 良文               |
| 早期審査対象出願      |                              |           | 日本国大阪府門真市大字門真1006番地 |
|               |                              |           | パナソニック株式会社内         |
|               |                              | 審査官       | 山下 剛史               |

最終頁に続く

(54) 【発明の名称】 声質変換システム、声質変換装置及びその方法、声道情報生成装置及びその方法

(57) 【特許請求の範囲】

【請求項1】

声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換システムであって、

互いに種類が異なる複数の母音の音声を受け付ける母音受付部と、

前記母音受付部によって受け付けられた複数の母音の音声を分析することにより、前記母音の種類毎に、第1声道形状情報を生成する分析部と、

前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合部と、

入力音声の声道形状情報及び音源情報を取得し、前記入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換し、変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成部とを備える

声質変換システム。

【請求項2】

前記混合部は、

前記母音の種類毎に生成された複数の第1声道形状情報を平均することにより、1つの平均声道形状情報を算出する平均声道情報算出部と、

前記母音受付部によって受け付けられた母音の種類毎に、当該母音の第1声道形状情報と前記平均声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合声道情報生成部とを備える

請求項1に記載の声質変換システム。

【請求項3】

平均声道情報算出部は、前記複数の第1声道形状情報を重み付き算術平均することにより、前記平均声道形状情報を算出する

請求項2に記載の声質変換システム。

【請求項4】

前記混合部は、前記入力音声に含まれる母音の局所的発話速度が大きいほど、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報が前記母音の種類毎に生成された複数の第1声道形状情報の平均に近づくように、前記第2声道形状情報を生成する

請求項1～3のいずれか1項に記載の声質変換システム。

【請求項5】

前記混合部は、母音の種類に応じて設定された混合比率を用いて、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合する

請求項1～4のいずれか1項に記載の声質変換システム。

【請求項6】

前記混合部は、ユーザーによって設定された混合比率を用いて、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合する

請求項1～5のいずれか1項に記載の声質変換システム。

【請求項7】

前記混合部は、前記入力音声の言語種類に応じて設定された混合比率を用いて、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合する

請求項1～6のいずれか1項に記載の声質変換システム。

【請求項8】

前記声質変換システムは、さらに、  
前記入力音声の声道形状情報及び音源情報が記憶されている入力音声記憶部を備え、  
前記合成部は、前記入力音声記憶部から、前記入力音声の声道形状情報及び音源情報を取得する

請求項1～7のいずれか1項に記載の声質変換システム。

【請求項9】

入力音声の声質を変換する際に用いられる、声道の形状を示す声道形状情報を生成する声道情報生成装置であって、

互いに種類が異なる複数の母音の音声を分析することにより、前記母音の種類毎に、第1声道形状情報を生成する分析部と、

前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合部とを備える

声道情報生成装置。

【請求項10】

さらに、

前記母音の種類毎に、前記第2声道形状情報を用いて合成音を生成する合成部と、

前記合成音を音声として出力する出力部とを備える

請求項9に記載の声道情報生成装置。

【請求項11】

声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換装置であっ

10

20

30

40

50

て、

母音の種類毎に、当該母音の第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより生成された第2声道形状情報を記憶している母音声道情報記憶部と、

入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換し、変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成部とを備える

声質変換装置。

【請求項12】

声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換方法であって、

互いに種類が異なる複数の母音の音声を受け付ける母音受付ステップと、

前記母音受付ステップにおいて受け付けられた複数の母音の音声を分析することにより、前記母音の種類毎に第1声道形状情報を生成する分析ステップと、

前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合ステップと、

入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換する変換ステップと、

変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成ステップとを含む

声質変換方法。

【請求項13】

入力音声の声質を変換する際に用いられる、声道の形状を示す声道形状情報を生成する声道情報生成方法であって、

互いに種類が異なる複数の母音の音声を分析することにより、前記母音の種類毎に第1声道形状情報を生成する分析ステップと、

前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合ステップとを含む

声道情報生成方法。

【請求項14】

声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換方法であって、

入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の第1声道形状情報及び前記入力音声に含まれる母音と異なる種類の母音の第1声道形状情報を混合することにより生成された、前記入力音声に含まれる母音と同じ種類の母音の第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換する変換ステップと、

変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成ステップとを含む

声質変換方法。

【請求項15】

請求項12に記載の声質変換方法をコンピュータに実行させるためのプログラム。

【請求項16】

請求項13に記載の声道情報生成方法をコンピュータに実行させるためのプログラム。

【請求項17】

請求項14に記載の声質変換方法をコンピュータに実行させるためのプログラム。

10

20

30

40

50

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、声質変換技術に関する。

## 【背景技術】

## 【0002】

従来の声質変換技術としては、互いに異なる2つの話し方（例えば感情）で発声された同一内容の音声の対を大量に用意し、それらから2つの話し方の間の変換規則を学習する技術がある（例えば、特許文献1参照）。特許文献1に記載の声質変換技術では、学習モデルに基づいて無感情音声から感情音声への変換を行うことができる。

10

## 【0003】

特許文献2記載の声質変換技術では孤立発声された少量の母音から特徴量を抽出することによって目的の音声への変換を実現している。

## 【先行技術文献】

## 【特許文献】

## 【0004】

【特許文献1】特開平7-72900号公報

【特許文献2】国際公開第2008/142836号

## 【発明の概要】

## 【発明が解決しようとする課題】

20

## 【0005】

しかしながら、上記の声質変換技術では、入力音声を滑らかで自然な音声に変換することができない場合がある。

## 【0006】

そこで、本発明は、入力音声を滑らかで自然な音声に変換することができる声質変換システムを提供する。

## 【課題を解決するための手段】

## 【0007】

本発明の一態様に係る声質変換システムは、声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換システムであって、互いに種類が異なる複数の母音の音声を受け付ける母音受付部と、前記母音受付部によって受け付けられた複数の母音の音声を分析することにより、前記母音の種類毎に、第1声道形状情報を生成する分析部と、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合部と、入力音声の声道形状情報及び音源情報を取得し、前記入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換し、変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成部とを備える。

30

## 【0008】

なお、これらの全般的または具体的な態様は、システム、方法、集積回路、コンピュータプログラムまたはコンピュータ読み取り可能なCD-ROM (Compact Disc Read Only Memory) などの記録媒体で実現されてもよく、システム、方法、集積回路、コンピュータプログラムおよび記録媒体の任意な組み合わせで実現されても良い。

40

## 【発明の効果】

## 【0009】

本発明の一態様に係る声質変換システムによれば、入力音声を滑らかで自然な音声に変換することができる。

## 【図面の簡単な説明】

50

【 0 0 1 0 】

【図 1】図 1 は、母音のスペクトル包絡の一例を示す模式図である。

【図 2 A】図 2 A は、孤立母音の第 1 及び第 2 フォルマント周波数の分布を示す図である。

【図 2 B】図 2 B は、文中母音の第 1 及び第 2 フォルマント周波数の分布を示す図である。

【図 3】図 3 は、人間の声道についての音響管モデルを示す図である。

【図 4 A】図 4 A は、孤立母音と平均声道形状情報との関係を示す図である。

【図 4 B】図 4 B は、文中母音と平均声道形状情報との関係を示す図である。

【図 5 A】図 5 A は、孤立母音の第 1 及び第 2 フォルマント周波数の平均を示す図である。

【図 5 B】図 5 B は、文中母音の第 1 及び第 2 フォルマント周波数の平均を示す図である。

【図 6】図 6 は、文中母音の  $F_1 - F_2$  平均、孤立母音の  $F_1 - F_2$  平均、及び平均声道形状情報の各々と、複数の文中母音の第 1 及び第 2 フォルマント周波数との二乗平均平方根誤差を示す図である。

【図 7】図 7 は、 $F_1 - F_2$  平面における各孤立母音の位置を平均声道形状情報の位置に向かって移動させたときの効果を説明するための図である。

【図 8】図 8 は、実施の形態 1 における声質変換システムの構成図である。

【図 9】図 9 は、実施の形態 1 における分析部の詳細な構成の一例を示す図である。

【図 10】図 10 は、実施の形態 1 における合成部の詳細な構成の一例を示す図である。

【図 11 A】図 11 A は、実施の形態 1 における声質変換システムの処理動作を示すフローチャートである。

【図 11 B】図 11 B は、実施の形態 1 における声質変換システムの処理動作を示すフローチャートである。

【図 12】図 12 は、実施の形態 1 における声質変換システムの処理動作を示すフローチャートである。

【図 13 A】図 13 A は、日本語の入力音声の声質を変換したときの実験結果を示す図である。

【図 13 B】図 13 B は、英語の入力音声の声質を変換したときの実験結果を示す図である。

【図 14】図 14 は、 $F_1 - F_2$  平面に英語の 13 母音を配置した図である。

【図 15】図 15 は、実施の形態 1 における母音受付部の一例を示す図である。

【図 16】図 16 は、全ての孤立母音の第 1 及び第 2 フォルマント周波数を比率  $q$  で移動させた場合に  $F_1 - F_2$  平面上で形成される多角形を示す図である。

【図 17】図 17 は、声道長変換比率  $r$  で声道断面積関数を伸縮する変換方法について説明するための図である。

【図 18】図 18 は、声道長変換比率  $r$  で声道断面積関数を伸縮する変換方法について説明するための図である。

【図 19】図 19 は、声道長変換比率  $r$  で声道断面積関数を伸縮する変換方法について説明するための図である。

【図 20】図 20 は、実施の形態 2 における声質変換システムの構成図である。

【図 21】図 21 は、実施の形態 2 における声道情報生成装置が出力する各母音の音声を説明するための図である。

【図 22】図 22 は、実施の形態 3 における声質変換システムの構成図である。

【図 23】図 23 は、他の実施の形態に係る声質変換システムの構成図である。

【図 24】図 24 は、特許文献 1 における声質変換装置の構成図である。

【図 25】図 25 は、特許文献 2 における声質変換装置の構成図である。

【発明を実施するための形態】

【 0 0 1 1 】

10

20

30

40

50

(本発明の基礎となった知見)

機器やインタフェースにおいて音声出力機能は、操作方法や機器の状態をユーザーに知らせるなどの重要な役割を担っている。また、情報機器においては、音声出力機能は、ネットワークを介して取得したテキスト情報などを読み上げる機能としても用いられる。

【0012】

さらに最近では、機器が擬人化されるとともに特徴的な声を出力することが求められる場合も増えている。例えば、人は、人型ロボットに人格を感じるため、人型ロボットが単調な合成音声で話したときには違和感を覚えることが多い。

【0013】

また、有名人やアニメのキャラクターの声で好きな言葉を喋らせる事ができるサービスが登場している。このようなサービスを提供するためのアプリケーションでは、話す内容以上に声の特徴がニーズの中心となっている。

【0014】

このように、音声出力機能への要求は、かつての明瞭性あるいは正確性から、声の種類が選べること、あるいは好みの声に変化させられることへと広がっている。

【0015】

さて、このような音声出力機能を実現する手段としては、人が話した声を録音して再生する録音再生方式と、テキストや発音記号から音声波形を生成する音声合成方式とがある。録音再生方式は、音が良いのが長所であるが、記憶容量が大きくなることと状況に応じて発話させる内容が変えられないことが短所である。

【0016】

一方、音声合成方式は、テキストで発話内容を変えることができるので記憶容量の増大は避けられるが、音質やイントネーションの自然さという点において録音再生方式には及ばない。したがって、メッセージの種類が少ない場合は録音再生方式が選ばれ、多い場合は音声合成方式が選ばれることが多い。

【0017】

ところが、いずれの方式を用いても、声の種類は予め用意した種類に限られる。すなわち、男性と女性など2種類の声を使いたい場合は、両方の声を録音しておくか両方の声の音声合成部を用意する必要があり、機器のコストや開発のコストが増大する。まして、好みの声に調整したり変えたりすることは不可能である。

【0018】

そこで、声の特徴を別の話者の声の特徴に近似させる声質変換技術の要求が高まっている。

【0019】

上述したように、従来の声質変換技術としては、互いに異なる2つの話し方(例えば感情)で発声された同一内容の音声の対を大量に用意し、それらから2つの話し方の間の変換規則を学習する技術がある(例えば、特許文献1参照)。

【0020】

図24は、特許文献1に記載の声質変換装置の構成図である。

【0021】

この図に示す声質変換装置は、音響的分析部2002と、スペクトルのDP(Dynamic Programming)マッチング部2004と、各音素の時間長伸縮部2006と、ニューラルネットワーク部2008とを備える。

【0022】

ニューラルネットワーク部2008は、無感情な音声の音響的特徴パラメータを、感情を伴った音声の音響的特徴パラメータに変換するための学習を行う。その後、学習済みの当該ニューラルネットワーク部2008を用いて無感情な音声に感情が付与される。

【0023】

スペクトルのDPマッチング部2004は、音響的分析部2002で抽出された特徴パラメータのうち、スペクトルの特徴パラメータについて、無感情の音声と感情を伴った音

10

20

30

40

50

声との間の類似度を時々刻々調べる。そして、スペクトルのDPマッチング部2004は、同一の音素毎の時間的な対応をとることによって、無感情音声に対する感情音声の音素毎の時間的な伸縮率を求める。

【0024】

各音素の時間長伸縮部2006は、スペクトルのDPマッチング部2004で得られた音素毎の時間的な伸縮率に応じて、感情音声の特徴パラメータの時系列を時間的に正規化して無感情音声の特徴パラメータの時系列に合うようにする。

【0025】

ニューラルネットワーク部2008は、学習時においては、時々刻々と入力層に与えられる無感情音声の音響的特徴パラメータと出力層に与えられる感情音声の音響的特徴パラメータとの違いを学習する。

10

【0026】

また、ニューラルネットワーク部2008は、感情の付与時においては、学習時に決定されたネットワーク内部の重み係数を用いて、時々刻々と入力層に与えられる無感情音声の音響的特徴パラメータから感情音声の音響的特徴パラメータを推定する計算を行なう。以上により、声質変換装置は、学習モデルに基づいて無感情音声から感情音声への変換を行う。

【0027】

しかしながら、特許文献1の技術では、予め決められた学習用文章と同一の内容の文章の音声を、目標とする感情を伴った発声で収録する必要がある。したがって、話者変換に用いる場合は、目標とする話者(目標話者)に予め決められた学習用文章を全て発話してもらう必要がある。したがって、目標話者に対する負担が大きくなることという課題がある。

20

【0028】

そこで、目標話者の発声負担が少なくなる技術として、少量の音声から目標話者の特徴量を抽出して用いる技術が提案されている(例えば、特許文献2参照)。

【0029】

図25は、特許文献2に記載の声質変換装置の構成図である。

【0030】

この図に示す声質変換装置は、入力音声の母音の声道情報を入力された変換比率で目標話者の母音の声道情報に変換することにより、入力音声の声質を変換する。ここで、声質変換装置は、目標母音声道情報保持部2101と、変換比率入力部2102と、母音変換部2103と、子音声道情報保持部2104と、子音選択部2105と、子音変形部2106と、合成部2107とを備える。

30

【0031】

目標母音声道情報保持部2101には、目標話者が発声した代表的な母音から抽出された目標母音声道情報が保持されている。母音変換部2103は、入力音声の母音区間の声道情報を、目標母音声道情報を用いて変換する。

【0032】

この時、母音変換部2103は、変換比率入力部2102から与えられた変換比率に基づいて、入力音声の母音区間の声道情報と目標母音声道情報とを混合する。子音選択部2105は、前後の母音との接続性を考慮して子音声道情報保持部2104から子音の声道情報を選択する。そして、子音変形部2106は、選択された子音の声道情報を、前後の母音になめらかに繋がるように変形する。合成部2107は、入力音声の音源情報と、母音変換部2103、子音選択部2105及び子音変形部2106により変形された声道情報とを用いて、合成音を生成する。

40

【0033】

しかしながら、特許文献2の技術では目標音声の声道情報として孤立発声された母音の声道情報を用いているので、変換された音声は滑らかさに欠け、ぎこちない印象となる。これは、別々に発声した母音の特徴と、文として連続して発声された音声の中の母音の特徴

50

との間に違いがあることに起因する。したがって、日常会話などの音声を対象に声質変換を行うと自然性の低下が著しくなる。

【0034】

以上説明したように、従来の声質変換技術では、少量の目標音声のサンプルを用いて入力音声の声質を変換する場合に、滑らかで自然な音声に変換することができなかった。すなわち、特許文献1の技術では、多量の同一内容の発声音声対から変換規則を学習する必要があるために、目標話者による大量の発声が必要になるという課題があった。一方、特許文献2の技術では、目標話者による母音の音声の入力のみで声質変換が可能であるという利点を有するが、利用できる音声特徴量が孤立発声された母音のものであるために生成される音声の自然性が低いという課題があった。

10

【0035】

このような課題を鑑みて、本願発明者らが見出した知見を以下に説明する。

【0036】

孤立して発声された音声 (discrete utterance speech) に含まれる母音は、文章として発声された音声に含まれる母音と異なる特徴を有する。例えば、「あ (a)」のみ発声したときの母音である「a」は、日本語の「こんにちは / k o N n i c h i w a /」に含まれる文末の「a」とは、異なる特徴を有する。また、「え (e)」のみ発声した時の母音である「e」は、英語の「Hello」に含まれる「e」とは、異なる特徴を有する。

20

【0037】

以下、孤立して発声することを「孤立発声」とも表記し、文章として連続して発声することを「連続発声」又は「文発声」とも表記する。また、孤立発声された母音を「孤立母音」とも表記し、文章として連続発声された母音を「文中母音」とも表記する。本願発明者らは、鋭意研究を行った結果、孤立発声の母音と文発声の母音の違いに関する新たな知見を見出した。以下、説明する。

【0038】

図1は、母音のスペクトル包絡の一例を示す模式図である。図1において、縦軸はパワーを示し、横軸は周波数を示す。図1に示すように、母音のスペクトルは複数のピークを有する。この複数のピークは、声道の共振に対応する。最も小さい周波数のピークは、第1フォルマントと呼ばれる。2番目に小さい周波数のピークは、第2フォルマントと呼ばれる。それぞれのピークの位置に対応する周波数 (中心周波数) を、それぞれ第1フォルマント周波数、第2フォルマント周波数と呼ぶ。母音の種類は、主に、第1フォルマント周波数と第2フォルマント周波数との関係で決まる。

30

【0039】

図2Aは、孤立母音の第1及び第2フォルマント周波数の分布を示す。図2Bは、文中母音の第1及び第2フォルマント周波数の分布を示す。図2A及び図2Bにおいて、横軸は第1フォルマント周波数を示し、縦軸は第2フォルマント周波数を示す。図2A及び図2Bに示す第1及び第2フォルマント周波数で定義された2次元平面をF1 - F2平面と呼ぶ。

40

【0040】

具体的には、図2Aは、ある話者が日本語の5母音を孤立発声したときの母音の第1及び第2フォルマント周波数を示す。また、図2Bは、同じ話者が日本語の文章を連続発声したときの母音の第1及び第2フォルマント周波数を示す。図2A及び図2Bにおいて、/ a / / i / / u / / e / / o / の5つの母音は、異なる記号で示されている。

【0041】

図2Aに示すように、5つの孤立母音を結ぶ点線の形状は、五角形となる。また、/ a / / i / / u / / e / / o / の5つの孤立母音は、F1 - F2平面において互いに離れて配置される。これは、/ a / / i / / u / / e / / o / の5つの孤立母音は、それぞれ異なる特徴を有することを意味する。例えば / a / と / i / の孤立母音は、/ a / と / o / の孤立母音よりも大きく離れていることが分かる。

50



## 【 0 0 4 2 】

しかし、図 2 B に示すように、5 つの文中母音は、F 1 - F 2 平面において互いの位置が近付いている。つまり、図 2 B に示す文中母音の位置は、図 2 A に示す孤立母音の位置よりも五角形の中心又は重心に近付いている。

## 【 0 0 4 3 】

文中母音では、その母音の前後の音素又は子音との調音が行なわれる。そのために、それぞれの文中母音に発声の怠け (reduction of articulation) が生じる。このため、文章として連続発声されたときの個々の母音は曖昧な発音になる。ただし、文章全体を通して音声は、なめらかで自然に聞こえる。

## 【 0 0 4 4 】

逆に、孤立母音と同じように、1 つ 1 つの文中母音をはっきりと発音された場合、調音運動が不自然になる。その結果、文章全体を通して音声は、滑らかではなく、ぎこちなく聞こえる。したがって、連続音声を合成する際には、発声の怠けを模擬する音声を用いることが重要である。

## 【 0 0 4 5 】

発声の怠けを実現するためには、文発声の音声から母音特徴量を抽出しても良い。しかし、そのためには多くの文発声の音声を用意する必要があるため、実用上使いやすさが大きく損なわれる。さらに、文中母音は、前後の音韻の影響を強く受ける。前後の音韻 (音韻環境) が近い母音を用いなければ、自然性が損なわれた音声となる。このため、膨大な量の文発声の音声が必要となる。例えば、数十文程度の文発声の音声では、必要十分な量とはならない。

## 【 0 0 4 6 】

本願発明者らは、( 1 ) 少量の音声を用意すれば良いという利便性を活かすために、孤立母音の特徴量を取得し、( 2 ) 発声の怠けを模擬するために、F 1 - F 2 平面において孤立母音によって形成される五角形を縮小する方向に孤立母音の特徴量を移動させるという知見を見出した。この知見に基づき、具体的な方法を説明する。

## 【 0 0 4 7 】

1 つ目の方法は、F 1 - F 2 平面において各母音を五角形の重心に向けて移動させる方法である。ここで、F 1 - F 2 平面上での第 i 母音の位置ベクトル  $\mathbf{b}$  を式 ( 1 ) のように定義する。

## 【 0 0 4 8 】

【数 1】

$$\mathbf{b}_i = [f1_i \quad f2_i] \quad (1)$$

## 【 0 0 4 9 】

ここで、 $f 1_i$  は、第 i 母音の第 1 フォルマント周波数を示し、 $f 2_i$  は、第 i 母音の第 2 フォルマント周波数を示す。i は母音の種類を表すインデックスである。5 母音の場合は、1 i 5 となる。

## 【 0 0 5 0 】

重心  $\mathbf{g}$  は、下記の式 ( 2 ) で表される。

## 【 0 0 5 1 】

【数 2】

$$\mathbf{g} = \frac{1}{N} \sum_{i=1}^N \mathbf{b}_i \quad (2)$$

## 【 0 0 5 2 】

ここで、N は母音の種類の数である。すなわち、重心  $\mathbf{g}$  は、母音の位置ベクトルの算術平均である。続いて、第 i 母音の位置ベクトルを下記の式 ( 3 ) のように変換する。

【 0 0 5 3 】

【 数 3 】

$$\hat{b}_i = ag + (1-a)b \quad (3)$$

【 0 0 5 4 】

ここで、 $a$ は、0から1の間の値であり、母音の位置ベクトル $b$ を重心 $g$ に近付ける度合いを表す曖昧化度合い係数である。曖昧化度合い係数 $a$ が1に近いほど、全ての母音は重心 $g$ に近づく。その結果、母音の位置ベクトル $b$ の違いも小さくなる。言い換えれば、図2Aに示す $F_1 - F_2$ 平面上において、各母音の音響的特徴が曖昧になる。

10

【 0 0 5 5 】

上記の考え方により、母音の曖昧化ができる。しかし、フォルマント周波数を直接変更することは、問題がある。図2Aには、第1フォルマント周波数と第2フォルマント周波数のみが示されている。しかし、孤立母音と文中母音とでは、第1及び第2フォルマント周波数だけではなく、他の物理量も異なっている。他の物理量は、例えば、第2フォルマント周波数よりも高次のフォルマント周波数又は各フォルマントのバンド幅などである。したがって、例えば、母音の第2フォルマント周波数のみをより高い周波数に変化させたとき、第2フォルマント周波数が第3フォルマント周波数に接近しすぎることが考えられる。

【 0 0 5 6 】

その結果、スペクトル包絡において異常に鋭いピークが現れ、合成フィルタが発振する、又は合成音の振幅が異常に大きくなる可能性がある。このような場合は、正常な音声を合成することができない。

20

【 0 0 5 7 】

音声の声質を変換する場合、音声の特徴を表す複数のパラメータがバランスを保った状態で変化しなければ、変換後の音声が妥当ではない音になってしまう。したがって、第1フォルマント周波数及び第2フォルマント周波数の2つのパラメータだけを変化させた場合、複数のパラメータのバランスが崩れ、著しく音質が劣化する。

【 0 0 5 8 】

この課題を解決するために、本願発明者らは、フォルマント周波数を直接変更するのではなく、声道形状を変形させることで母音を曖昧化する方法を見出した。

30

【 0 0 5 9 】

( 声道断面積関数 )

声道形状を示す情報(以下、「声道形状情報」という)としては、例えば、声道断面積関数がある。図3は、人間の声道についての音響管モデルを示す。人間の声道とは、声帯から口唇までの空間である。

【 0 0 6 0 】

図3の(a)において、縦軸は断面積の大きさを示し、横軸は音響管のセクション番号を示す。ここで、音響管のセクション番号とは、声道の中の位置を示す。横軸の左端は、口唇(Lip)の位置に対応し、横軸の右端は、声門(glottis)の位置に対応する。

40

【 0 0 6 1 】

図3の(a)に示す音響管モデルは、複数の円形の音響管が縦続接続されている。声道の断面積を、各セクションの音響管の断面積として、声道形状を模擬している。ここで、声道の長さ方向の位置と、その位置に対応する断面積の大きさとの関係を声道断面積関数と呼ぶ。

【 0 0 6 2 】

声道の断面積は、LPC分析に基づくPARCOR係数と一意に対応することが知られている。下記の式(4)により、PARCOR係数を、声道の断面積に変換できる。以下、PARCOR係数 $k_i$ を、声道形状情報の一例として説明する。ただし、声道形状情報

50

は、PARCOR係数に限定されるものではなく、PARCOR係数に等価なLSP (Line Spectrum Pairs) やLPCなどであっても良い。また、上述の音響管モデルにおける音響管の間の反射係数とPARCOR係数とは、符号が反転していることが違うだけである。このため、声道形状情報として反射係数が用いられても良い。

【0063】

【数4】

$$\frac{A_i}{A_{i+1}} = \frac{1-k_i}{1+k_i} \quad (4)$$

10

【0064】

ここで、 $A_i$ は、図3の(b)に示す第*i*区間の音響管の断面積であり、 $k_i$ は、第*i*番目と第*i+1*番目との境界のPARCOR係数(反射係数)である。

【0065】

PARCOR係数は、LPC分析により分析された線形予測係数 $\hat{a}_i$ を用いて算出することができる。具体的には、PARCOR係数は、Levinson-Durbin-Itakuraアルゴリズムを用いることにより算出される。なお、PARCOR係数は次の特徴を有する。

- ・線形予測係数は分析次数*p*に依存するが、PARCOR係数は分析の次数に依存しない。
- ・低次係数の値の変動はスペクトルへの影響が大きく、高次になるにつれて値の変動がスペクトルに与える影響が小さくなる。
- ・高次係数の値の変動のスペクトルへの影響は全周波数帯域に渡って平坦なものである。

20

【0066】

なお、声道形状情報は、必ずしも声道の断面積を示す情報である必要はなく、声道の各セクションの容積を示す情報であっても良い。

【0067】

(声道形状の変形)

次に、声道形状の変形について説明する。上述のように、声道の形状は、式(4)に示すPARCOR係数から求められる。ここでは、声道形状を変形するために、複数の声道形状情報を混合する。具体的には、複数の声道断面積関数の加重平均を求める代わりに、複数のPARCOR係数ベクトルの加重平均を求める。第*i*母音のPARCOR係数ベクトルは、式(5)で表される。

30

【0068】

【数5】

$$\mathbf{k}_i = (k_1^i \quad k_2^i \quad \cdots \quad k_M^i) \quad (5)$$

【0069】

複数の母音のPARCOR係数ベクトルの加重平均は式(6)で表される。

40

【0070】

【数6】

$$\bar{\mathbf{k}} = \sum_i w_i \mathbf{k}_i$$

$$\sum_i w_i = 1 \quad (6)$$

【0071】

ここで $w_i$ は、重み係数である。混合したい母音の声道形状情報が2つの場合、重み係

50

数は、2つの声道形状情報の混合比に対応する。

【0072】

(声道形状情報の曖昧化)

次に、母音を曖昧化するために複数の母音の声道形状情報を混合する手順を説明する。

【0073】

まず、N個の種類の母音の平均声道形状情報を式(7)で求める。つまり、各母音の声道形状情報が示す値(ここではPARCOR係数)の算術平均を算出することにより、平均声道形状情報を生成する。

【0074】

【数7】

$$\bar{\mathbf{k}} = \frac{1}{N} \sum_{i=0}^{N-1} \mathbf{k}_i \quad (7)$$

10

【0075】

次に、第i母音の曖昧化度合い係数aを用いて、第i母音の声道形状情報を曖昧化後の声道形状情報に変換する。すなわち、各母音の声道形状情報が示す値を平均声道形状情報が示す値に近付けることにより、曖昧化後の各母音の声道形状情報を生成する。つまり、第i母音の声道形状情報と他の母音の声道形状情報とを混合して、曖昧化後の声道形状情報を生成する。

20

【0076】

【数8】

$$\hat{\mathbf{k}}_i = a\bar{\mathbf{k}} + (1-a)\mathbf{k}_i \quad (8)$$

$\mathbf{k}_i$  : 曖昧化前の母音の声道形状情報、 $\hat{\mathbf{k}}_i$  : 曖昧化後の母音の声道形状情報

【0077】

このようにして生成された曖昧化後の母音の声道形状情報を用いて音声の合成を行うことで、音質を劣化させず、発声の怠けを再現することができる。

【0078】

以下に、実際に実験を行った結果について説明する。

30

【0079】

図4Aは、孤立母音と平均声道形状情報との関係を示す。また、図4Bは、文中母音と平均声道形状情報との関係を示す。図4A及び図4Bにおいて、平均声道形状情報は、式(7)に従って、図2Aに示す孤立母音の情報を用いて求めたなお、図4A及び図4Bに示す星印は、平均声道形状情報を用いて合成された母音の第1及び第2フォルマント周波数を示す。

【0080】

図4Aにおいて、平均声道形状情報は、5つの母音によって形成される五角形の重心近傍に位置する。図4Bにおいて、平均声道形状情報は、文中母音が分布する領域の中心近傍に位置する。

40

【0081】

図5Aは、孤立母音(図2Aに示す15個の母音)の第1及び第2フォルマント周波数の平均を示す。また、図5Bは、文中母音(図2Bに示す95個の母音)の第1及び第2フォルマント周波数の平均を示す。なお、以下において、第1及び第2フォルマント周波数の平均をF1 - F2平均とも呼ぶ。

【0082】

図5A及び図5Bにおいて、第1フォルマント周波数及び第2フォルマント周波数の平均は、破線で示されている。また、図5A及び図5Bには、図4A及び図4Bに示した平均声道形状情報も星印で示されている。

50

## 【 0 0 8 3 】

式(7)を用いて求めた図4Aに示す平均声道形状情報の位置は、図5Aに示す孤立母音のF1 - F2平均の位置よりも、図5Bに示す文中母音のF1 - F2の平均の位置に近い。したがって、式(7)及び式(8)を用いて求めた平均声道形状情報は、孤立母音のF1 - F2の平均よりも、実際の発声の急げに近似している。以下に、具体的な座標値を用いて説明する。

## 【 0 0 8 4 】

図6は、文中母音のF1 - F2平均、孤立母音のF1 - F2平均、及び平均声道形状情報の各々と、複数の文中母音の第1及び第2フォルマント周波数との二乗平均平方根誤差(RMSE: root mean square error)を示す。

10

## 【 0 0 8 5 】

図6に示すように、平均声道形状情報のRMSEは、孤立母音のF1 - F2平均のRMSEよりも、文中母音のF1 - F2平均のRMSEに近い。ただし、RMSEが近いことだけが、音声の自然さに貢献するとは言えないが、発声の急げの近似度合いを表す指標として見ることはできる。

## 【 0 0 8 6 】

次に、図7は、式(8)を用いて、F1 - F2平面における各孤立母音の位置を平均声道形状情報の位置に向かって移動させたときの効果を説明するための図である。図7において、大きい白丸はa = 0の場合の各母音の位置、小さい白丸はa = 1の場合の各母音の位置すなわち平均声道形状における位置を表し、黒い点はaを0.1刻みで大きくしていった場合の各母音の位置を表している。全ての母音が孤立母音の位置から平均声道形状における母音の位置に向かって連続的に移動している。このように、声道形状情報を混合して声道形状を変形することにより、第1及び第2フォルマント周波数は平均化及び曖昧化が可能であることが分かった。

20

## 【 0 0 8 7 】

そこで、本発明の一態様に係る声質変換システムは、声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換システムであって、互いに種類が異なる複数の母音の音声を受け付ける母音受付部と、前記母音受付部によって受け付けられた複数の母音の音声を分析することにより、前記母音の種類毎に、第1声道形状情報を生成する分析部と、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合部と、入力音声の声道形状情報及び音源情報を取得し、前記入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換し、変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成部とを備える。

30

## 【 0 0 8 8 】

この構成によれば、母音の種類毎に、複数の第1声道形状情報を混合して第2声道形状情報を生成することができる。つまり、少量の音声のサンプルから母音の種類毎に第2声道形状情報を生成することができる。このように母音の種類毎に生成された第2声道形状情報は、曖昧化された母音の声道形状情報に相当する。したがって、第2声道形状情報を用いて入力音声の声質を変換することにより、入力音声を滑らかで自然な音声に変換することが可能となる。

40

## 【 0 0 8 9 】

また例えば、前記混合部は、前記母音の種類毎に生成された複数の第1声道形状情報を平均することにより、1つの平均声道形状情報を算出する平均声道情報算出部と、前記母音受付部によって受け付けられた母音の種類毎に、当該母音の第1声道形状情報と前記平均声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合声道情報生成部とを備えても良い。

## 【 0 0 9 0 】

50

この構成によれば、第2声道形状情報を平均声道形状情報に容易に近付けることが可能となる。

【0091】

また例えば、平均声道情報算出部は、前記複数の第1声道形状情報を重み付き算術平均することにより、前記平均声道形状情報を算出しても良い。

【0092】

この構成によれば、複数の第1声道形状情報の重み付き算術平均を平均声道形状情報として算出することができる。したがって、例えば、目標話者の発声の怠けの特徴に応じて第1声道形状情報に重み付けすることにより、入力音声をより滑らかで自然な目標話者の音声に変換することも可能となる。

10

【0093】

また例えば、前記混合部は、前記入力音声に含まれる母音の局所的発話速度が大きいほど、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報が前記母音の種類毎に生成された複数の第1声道形状情報の平均に近付くように、前記第2声道形状情報を生成しても良い。

【0094】

この構成によれば、入力音声に含まれる母音の局所的発話速度に応じて複数の第1声道形状情報の混合比率を設定することができる。文中母音の曖昧化度合いは、局所的発話速度に依存する。したがって、入力音声をより滑らかで自然な音声に変換することが可能となる。

20

【0095】

また例えば、前記混合部は、母音の種類に応じて設定された混合比率を用いて、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合しても良い。

【0096】

この構成によれば、母音の種類に応じて、複数の第1声道形状情報の混合比率を設定することができる。文中母音の曖昧化度合いは、母音の種類に依存する。したがって、入力音声をより滑らかで自然な音声に変換することが可能となる。

【0097】

また例えば、前記混合部は、ユーザーによって設定された混合比率を用いて、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合しても良い。

30

【0098】

この構成によれば、複数の母音の曖昧化度合いを、ユーザーの好みにあわせて設定することができる。

【0099】

また例えば、前記混合部は、前記入力音声の言語種類に応じて設定された混合比率を用いて、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合しても良い。

【0100】

この構成によれば、入力音声の言語種類に応じて、複数の第1声道形状情報の混合比率を設定することができる。文中母音の曖昧化度合いは、入力音声の言語種類に依存する。したがって、各言語にふさわしい曖昧化度合いを設定することができる。

40

【0101】

また例えば、前記声質変換システムは、さらに、前記入力音声の声道形状情報及び音源情報が記憶されている入力音声記憶部を備え、前記合成部は、前記入力音声記憶部から、前記入力音声の声道形状情報及び音源情報を取得しても良い。

【0102】

本発明の一態様に係る声道情報生成装置は、入力音声の声質を変換する際に用いられる、声道の形状を示す声道形状情報を生成する声道情報生成装置であって、互いに種類が異

50

なる複数の母音の音声进行分析することにより、前記母音の種類毎に、第1声道形状情報を生成する分析部と、前記母音の種類毎に、当該母音の前記第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する混合部とを備える。

【0103】

この構成によれば、母音の種類毎に、複数の第1声道形状情報を混合して第2声道形状情報を生成することができる。つまり、少量の音声のサンプルから母音の種類毎に第2声道形状情報を生成することができる。このように母音の種類毎に生成された第2声道形状情報は、曖昧化された母音の声道形状情報に相当する。したがって、第2声道形状情報が声質変換装置に出力されれば、声質変換装置は、第2声道形状情報を用いて入力音声を滑らかに自然な音声に変換することができる。

10

【0104】

また例えば、さらに、前記母音の種類毎に、前記第2声道形状情報を用いて合成音を生成する合成部と、前記合成音を音声として出力する出力部とを備えても良い。

【0105】

この構成によれば、母音の種類毎に第2声道形状情報を用いて生成された合成音を音声として出力することができる。したがって、従来の声質変換装置を用いて、入力音声を滑らかに自然な音声に変換することができる。

【0106】

本発明の一態様に係る声質変換装置は、声道の形状を示す声道形状情報を用いて入力音声の声質を変換する声質変換装置であって、母音の種類毎に、当該母音の第1声道形状情報と、当該母音と異なる種類の母音の前記第1声道形状情報とを混合することにより生成された第2声道形状情報を記憶している母音声道情報記憶部と、入力音声に含まれる母音の声道形状情報と、前記入力音声に含まれる母音と同じ種類の母音の前記第2声道形状情報とを混合することにより、前記入力音声の声道形状情報を変換し、変換後の前記入力音声の声道形状情報と前記入力音声の音源情報とを用いて合成音を生成することにより、前記入力音声の声質を変換する合成部とを備える。

20

【0107】

この構成によれば、上記声質変換システムと同様の効果を奏することができる。

【0108】

なお、これらの全般的または具体的な態様は、方法、集積回路、コンピュータプログラムまたはコンピュータ読み取り可能なCD-ROMなどの記録媒体で実現されてもよく、方法、集積回路、コンピュータプログラムおよび記録媒体の任意な組み合わせで実現されても良い。

30

【0109】

以下本発明の実施の形態について、図面を参照しながら説明する。

【0110】

なお、以下で説明する実施の形態は、いずれも本発明の一具体例を示す。以下の実施の形態で示される数値、形状、材料、構成要素、構成要素の配置位置及び接続形態、ステップ、ステップの順序などは、一例であり、本発明を限定する主旨ではない。また、以下の実施の形態における構成要素のうち、最上位概念を示す独立請求項に記載されていない構成要素については、任意の構成要素として説明される。

40

【0111】

(実施の形態1)

図8は、実施の形態1における声質変換システム100の構成図である。

【0112】

声質変換システム100は、声道の形状を示す声道形状情報を用いて入力音声の声質を変換する。図8に示すように、声質変換システム100は、入力音声記憶部101と、母音受付部102と、分析部103と、第1母音声道情報記憶部104と、混合部105と、第2母音声道情報記憶部107、合成部108と、出力部109と、混合比率入力部1

50

10と、変換比率入力部111とを備える。それぞれの構成要素は、有線又は無線で接続されており、互いに情報を送受信する。以下、各構成要素について、説明する。

【0113】

(入力音声記憶部101)

入力音声記憶部101は、入力音声情報と、入力音声情報と対応付けられた付属情報とを記憶している。入力音声情報とは、変換対象となる入力音声に関する情報である。具体的には、入力音声情報は、複数の音素で構成される音声の情報である。例えば、ある歌手が歌った音声等を予め録音しておくことにより、入力音声情報が準備される。より具体的には、入力音声記憶部101は、入力音声情報を声道情報と音源情報とに分離した形式で記憶している。

10

【0114】

付属情報は、入力音声において音素の境界を示す時間の情報と、音素の種類の情報とを含む。

【0115】

(母音受付部102)

母音受付部102は、母音の音声を受け付ける。本実施の形態では、母音受付部102は、入力音声と同じ言語の母音の音声であって、互いに種類が異なる複数の母音の音声を受け付ける。互いに種類が異なる複数の母音の音声とは、複数の異なる種類の母音を含んでいれば良く、同じ種類の複数の母音を含んでも良い。

20

【0116】

母音受付部102は、分析部103に、母音の音声に対応する電気信号である母音の音響信号を送信する。

【0117】

母音受付部102は、例えば、話者が発した音声を受け付ける場合は、マイクロホンを有する。母音受付部102は、例えば、予め電気信号に変換されている音響信号を受け付ける場合、オーディオ回路及びアナログデジタル変換器を有する。母音受付部102は、例えば、予め音響信号がデジタルデータに変換された音響データを受け付ける場合、データ読出器を有する。

【0118】

なお、母音受付部102は、表示部を備えても良い。表示部は、目標話者に発声させたい単母音又は文章と、発声タイミングとを表示する。

30

【0119】

また、母音受付部102が受け付ける音声は、孤立発声された母音であっても良い。例えば、母音受付部102は、代表的な母音の音響信号を受け付けても良い。代表的な母音は、言語により異なる。例えば、日本語の代表的な母音とは、/a//i//u//e//o/の5種類の母音である。英語の代表的な母音は、以下に国際音声記号(International Phonetic Alphabet)で示す13種類の母音である。

【0120】

【数9】

[i][u][ɪ][ʊ][e][o][ə][ɛ][ʌ][ɔ][æ][ɑ][ɒ]

40

【0121】

母音受付部102は、例えば日本語の母音の音声を受け付ける場合は、/a//i//u//e//o/の5種類の母音を目標話者に孤立発声(すなわち各母音の間を開けて発声)させることで、母音の音声を受け付ける。このように話者に母音を孤立発声してもらうことにより、分析部103は、パワー情報を用いて母音区間を切り出すことが可能となる。

【0122】

50



ただし、母音受付部 102 は、必ずしも孤立発声された母音の音声を受け付ける必要はない。母音受付部 102 は、文章として連続発声された母音を受け付けても良い。例えば話者が緊張していて意識的にはっきりとした発声が行われた場合は、文章として連続発声された母音も、孤立発声された母音に近い音声になることがある。母音受付部 102 が文発声の母音を受け付ける場合は、例えば 5 母音を含む文章（例えば「本日は晴天なり」など）を話者に発声させれば良い。この場合、分析部 103 は、HMM (Hidden - Markov - Model) などを用いた音素自動セグメンテーション技術によって母音区間を切り出すことができる。

#### 【0123】

(分析部 103)

分析部 103 は、母音受付部 102 から母音の音響信号を受け付ける。分析部 103 は、母音受付部 102 で受け付けられた母音の音響信号に対して、付属情報を付与する。さらに、分析部 103 は、例えば LPC (Linear Predictive Coding) 分析や ARX (Auto-regressive Exogenous) 分析などの分析方法を用いて各母音の音響信号を分析することにより、各母音の音響信号を声道情報と音源情報とに分離する。

#### 【0124】

声道情報には、母音が発声された時の声道の形状を示す声道形状情報が含まれる。分析部 103 によって分離された声道情報に含まれる声道形状情報を第 1 声道形状情報と呼ぶ。つまり、分析部 103 は、母音受付部 102 によって受け付けられた複数の母音の音声を分析することにより、母音の種類毎に、第 1 声道形状情報を生成する。

#### 【0125】

第 1 声道形状情報の例としては、上述の LPC の他に、PARCOR 係数、PARCOR 係数と等価な LSP (Line Spectrum Pairs) などがある。また、音響管モデルにおける音響管の間の反射係数と PARCOR 係数との関係は、符号が反転していることのみである。このため、反射係数そのものを第 1 声道形状情報として用いても良い。

#### 【0126】

付属情報は、各母音の種類 (/a / /i / など) と、母音区間中心の時刻とを含む。分析部 103 は、第 1 母音声道情報記憶部 104 に、母音の種類毎に、少なくとも母音の第 1 声道形状情報を格納する。

#### 【0127】

次に、母音の第 1 声道形状情報の生成方法の一例を説明する。

#### 【0128】

図 9 は、実施の形態 1 における分析部 103 の詳細な構成の一例を示す。分析部 103 は、母音安定区間抽出部 1031 と母音声道情報作成部 1032 とを備える。

#### 【0129】

母音安定区間抽出部 1031 は、入力された母音を含む音声から孤立母音の区間 (母音区間) を抽出することにより母音区間中心の時刻を算出する。母音区間の抽出方法は特に限定される必要はない。例えば、母音安定区間抽出部 1031 は、パワーが一定以上の区間を安定区間とし、当該安定区間を母音区間として抽出するようにしても良い。

#### 【0130】

母音声道情報作成部 1032 は、母音安定区間抽出部 1031 により抽出された孤立母音の母音区間中心に対して、母音の声道形状情報を作成する。例えば、母音声道情報作成部 1032 は、上述の PARCOR 係数を第 1 声道形状情報として算出する。母音声道情報作成部 1032 は、第 1 母音声道情報記憶部 104 に、母音の第 1 声道形状情報を格納する。

#### 【0131】

(第 1 母音声道情報記憶部 104)

第 1 母音声道情報記憶部 104 は、母音の種類毎に、少なくとも母音の第 1 声道形状情

10

20

30

40

50

報を記憶している。つまり、第1母音声道情報記憶部104は、分析部103によって母音の種類毎に生成された複数の第1声道形状情報を記憶している。

【0132】

(混合部105)

混合部105は、母音の種類毎に、当該母音の第1声道形状情報と、当該母音と異なる種類の母音の第1声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する。具体的には、混合部105は、母音の種類毎に、当該母音の第2声道形状情報が当該母音の第1声道形状情報よりも平均声道形状情報に近づくように、当該母音の第2声道形状情報を生成する。このように生成される第2声道形状情報は、曖昧化された声道形状情報に相当する。

10

【0133】

なお、平均声道形状情報とは、母音の種類毎に生成された複数の第1声道形状情報の平均である。また、複数の声道形状情報を混合するとは、複数の声道形状情報の各々が示す値又はベクトルを重み付け加算することを意味する。

【0134】

ここで、混合部105の詳細な構成の一例を説明する。混合部105は、例えば、平均声道情報算出部1051と混合声道情報生成部1052とを備える。

【0135】

(平均声道情報算出部1051)

平均声道情報算出部1051は、第1母音声道情報記憶部104に記憶された複数の第1声道形状情報を取得する。平均声道情報算出部1051は、取得した複数の第1声道形状情報を平均することにより、1つの平均声道形状情報を算出する。具体的な処理については、後述する。平均声道情報算出部1051は、混合声道情報生成部1052に平均声道形状情報を送信する。

20

【0136】

(混合声道情報生成部1052)

混合声道情報生成部1052は、平均声道情報算出部1051から平均声道形状情報を受信する。また、混合声道情報生成部1052は、第1母音声道情報記憶部104に記憶された複数の第1声道形状情報を取得する。

【0137】

そして、混合声道情報生成部1052は、母音受付部102によって受け付けられた母音の種類毎に、当該母音の第1声道形状情報と平均声道形状情報とを混合することにより、当該母音の第2声道形状情報を生成する。具体的には、混合声道情報生成部1052は、母音の種類毎に、第1声道形状情報を平均声道形状情報に近付ける処理を行うことにより、第2声道形状情報を生成する。

30

【0138】

第1声道形状情報と平均声道形状情報との混合比率は、母音の曖昧化度合いに応じて設定されれば良い。本実施の形態では、混合比率は、式(8)における曖昧化度合い係数 $a$ に相当する。つまり、混合比率は、値が大きいほど曖昧化度合いが高くなる。混合声道情報生成部1052は、混合比率入力部110から入力された混合比率を用いて、第1声道形状情報と平均声道形状情報とを混合する。

40

【0139】

なお、混合声道情報生成部1052は、予め記憶されている混合比率を用いて、第1声道形状情報と平均声道形状情報とを混合しても良い。この場合、声質変換システム100は、必ずしも混合比率入力部110を備える必要はない。

【0140】

ある種類の母音の第2声道形状情報を平均声道形状情報に近付けた場合、その種類の母音の第2声道形状情報は、他の種類の母音の第2声道形状情報に近づく。すなわち、第2声道形状情報が平均声道形状情報により近づくように混合比率が設定されれば、混合声道情報生成部1052は、より曖昧化された第2声道形状情報を生成することができる。こ

50

のようなより曖昧化された第2声道形状情報を用いて生成された合成音は、滑舌が悪い音声となる。例えば、幼児の声に入力音声の声質を変換するときには、このように第2声道形状情報が平均声道形状情報に近づくように混合比率が設定されることが有効である。

【0141】

また、第2声道形状情報を平均声道形状情報にあまり近付けない場合、第2声道形状情報は、孤立母音の声道形状情報に近くなる。例えば、口を大きく開けてはっきり調音する傾向にある歌声に入力音声の声質を変換するときには、このように第2声道形状情報が平均声道形状情報にあまり近付かないように混合比率が設定されることが適している。

【0142】

混合声道情報生成部1052は、第2母音声道情報記憶部107に、母音の種類毎の第2声道形状情報を格納する。

10

【0143】

(第2母音声道情報記憶部107)

第2母音声道情報記憶部107は、母音の種類別に、第2声道形状情報を記憶している。つまり、第2母音声道情報記憶部107は、混合部105によって母音の種類毎に生成された複数の第2声道形状情報を記憶している。

【0144】

(合成部108)

合成部108は、入力音声記憶部101に記憶されている入力音声情報を取得する。また、合成部108は、第2母音声道情報記憶部107に記憶されている母音の種類毎の第2声道形状情報を取得する。

20

【0145】

そして、合成部108は、入力音声情報に含まれる母音の声道形状情報と、入力音声情報に含まれる母音と同じ種類の母音の第2声道形状情報とを混合することにより、入力音声の声道形状情報を変換する。その後、合成部108は、入力音声の変換後の声道形状情報と、入力音声記憶部101に記憶されている入力音声の音源情報とを用いて合成音を生成することにより、入力音声の声質を変換する。

【0146】

具体的には、合成部108は、変換比率入力部111から入力された変換比率を混合比率として用いて、入力音声情報に含まれる母音の声道形状情報と、当該母音と同じ種類の母音の第2声道形状情報とを混合する。この変換比率は、入力音声を変化させる度合いに応じて設定されれば良い。

30

【0147】

なお、合成部108は、予め記憶されている変換比率を用いて、入力音声情報に含まれる母音の声道形状情報と当該母音と同じ種類の母音の第2声道形状情報とを混合しても良い。この場合、声質変換システム100は、必ずしも変換比率入力部111を備える必要はない。

【0148】

合成部108は、このように生成された合成音の信号を出力部109に送信する。

【0149】

ここで、合成部108の詳細な構成の一例を説明する。なお、以下に説明する合成部108の詳細な構成は、特許文献2と同様の構成である。

40

【0150】

図10は、実施の形態1における合成部108の詳細な構成の一例を示す。合成部108は、母音変換部1081と、子音選択部1082と、子音声道情報記憶部1083と、子音変形部1084と、音声合成部1085とを備える。

【0151】

母音変換部1081は、入力音声記憶部101から、音素境界付き声道情報と音源情報とを取得する。

【0152】

50

音素境界付き声道情報は、入力音声の声道情報に、入力音声に対応する音素情報と各音素の時間長の情報とが付された情報である。母音変換部1081は、母音区間ごとに該当する母音の第2声道形状情報を第2母音声道情報記憶部107から読み出す。そして、母音変換部1081は、母音区間の声道形状情報と読み出した第2声道形状情報とを混合することにより、入力音声の母音部の声質変換を行なう。この時の変換割合は、変換比率入力部111から入力された変換比率に基づく。

【0153】

子音選択部1082は、前後の母音との接続性を考慮して子音声道情報記憶部1083から子音の声道情報を選択する。そして、子音変形部1084は、選択された子音の声道情報を、前後の母音になめらかに繋がるように変形する。音声合成部1085は、入力音声の音源情報と、母音変換部1081、子音選択部1082及び子音変形部1084により変形された声道情報とを用いて、合成音を生成する。

10

【0154】

このように、特許文献2における目標母音声道情報を第2声道形状情報に置き換えて声質変換が実行される。

【0155】

(出力部109)

出力部109は、合成部108から合成音信号を受信する。出力部109は、合成音信号を合成音として出力する。出力部109は、例えば、スピーカで構成される。

【0156】

(混合比率入力部110)

混合比率入力部110は、混合声道情報生成部1052で用いる混合比率を受け付ける。混合比率入力部110は、混合声道情報生成部1052に、受け付けた混合比率を送信する。

20

【0157】

(変換比率入力部111)

変換比率入力部111は、合成部108で用いる変換比率を受け付ける。変換比率入力部111は、合成部108に、受け付けた変換比率を送信する。

【0158】

次に、以上のように構成された声質変換システム100における各種動作について説明する。

30

【0159】

図11A、図11B及び図12は、実施の形態1における声質変換システム100の処理動作を示すフローチャートである。

【0160】

具体的には、図11Aは、声質変換システム100において母音の音声を受け付けてから第2声道形状情報を生成するまでの処理の流れを示す。また、図11Bは、図11Aに示す第2声道形状情報生成処理(S600)の詳細を示す。また、図12は、実施の形態1において入力音声の声質を変換する処理の流れを示す。

【0161】

(ステップS100)

母音受付部102は、目標話者が発声した母音が含まれる音声を受け付ける。母音が含まれる音声とは、例えば、日本語の場合、日本語の5母音を「アー、イー、ウー、エー、オー」と発声したときの音声である。各母音の間隔は、500ms程度であれば良い。

40

【0162】

(ステップS200)

分析部103は、母音受付部102が受け付けた音声に含まれる1つの母音の声道形状情報を第1声道形状情報として生成する。

【0163】

(ステップS300)

50

分析部 103 は、生成された第 1 声道形状情報を、第 1 母音声道情報記憶部 104 に格納する。

【0164】

(ステップ S400)

分析部 103 は、母音受付部 102 が受け付けた音声に含まれる全ての種類の母音について、第 1 声道形状情報が生成されたか否かを判定する。例えば、分析部 103 は、母音受付部 102 が受け付けた音声に含まれる母音の種類情報を取得する。さらに、分析部 103 は、取得した母音の種類情報を参照して、音声に含まれる全ての種類の母音の第 1 声道形状情報が第 1 母音声道情報記憶部 104 に記憶されているか否かを判定する。ここで、全ての種類の母音の第 1 声道形状情報が第 1 母音声道情報記憶部 104 に記憶されている場合に、分析部 103 は、完了と判断する。一方、いずれかの種類の母音の第 1 声道形状情報が記憶されていない場合には、分析部 103 は、ステップ S200 の処理を行う。

10

【0165】

(ステップ S500)

平均声道情報算出部 1051 は、第 1 母音声道情報記憶部 104 に記憶されている全ての種類の母音の第 1 声道形状情報を用いて、1 つの平均声道形状情報を算出する。

【0166】

(ステップ S600)

混合声道情報生成部 1052 は、ステップ S100 で受け付けられた音声に含まれる母音の種類毎に、平均声道形状情報と、第 1 母音声道情報記憶部 104 に記憶されている第 1 声道形状情報とを用いて、第 2 声道形状情報を生成する。

20

【0167】

ここで、図 11B を用いて、ステップ S600 の詳細を説明する。

【0168】

(ステップ S601)

混合声道情報生成部 1052 は、第 1 母音声道情報記憶部 104 に記憶されている 1 つの母音の第 1 声道形状情報に平均声道形状情報を混合することによって、当該母音の第 2 声道形状情報を生成する。

【0169】

(ステップ S602)

混合声道情報生成部 1052 は、第 2 母音声道情報記憶部 107 に、ステップ S601 で生成された第 2 声道形状情報を格納する。

30

【0170】

(ステップ S603)

混合声道情報生成部 1052 は、ステップ S100 で受け付けられた音声に含まれる全ての種類の母音について、ステップ S602 の処理が行われたか否かを判定する。例えば、混合声道情報生成部 1052 は、母音受付部 102 が受け付けた音声に含まれる母音の種類情報を取得する。そして、混合声道情報生成部 1052 は、取得した母音の種類情報を参照して、音声に含まれる全ての種類の母音の第 2 声道形状情報が第 2 母音声道情報記憶部 107 に記憶されているか否かを判定する。

40

【0171】

ここで、全ての種類の母音の第 2 声道形状情報が第 2 母音声道情報記憶部 107 に記憶されている場合に、混合声道情報生成部 1052 は、完了と判断する。一方、いずれかの種類の母音の第 2 声道形状情報が第 2 母音声道情報記憶部 107 に記憶されていない場合には、混合声道情報生成部 1052 は、ステップ S601 の処理を行う。

【0172】

次に、このように母音の種類毎に生成された第 2 声道形状情報を用いて入力音声の声質を変換する処理について図 12 を用いて説明する。

【0173】

(ステップ S800)

50

合成部108は、第2母音声道情報記憶部107に記憶されている第2声道形状情報を用いて、入力音声記憶部101に記憶されている入力音声の声道形状情報を変換する。具体的には、合成部108は、入力音声に含まれる母音の声道形状情報を、入力音声に含まれる母音と同じ種類の母音の第2声道形状情報と混合することにより、入力音声の声道形状情報を変換する。

【0174】

(ステップS900)

合成部108は、ステップS800で変換された入力音声の声道形状情報と、入力音声記憶部101に記憶されている入力音声の音源情報とを用いて、合成音を生成する。これにより、入力音声の声質が変換された合成音が生成される。つまり、声質変換システム100は、入力音声の特徴を変化させることができる。

10

【0175】

(実験結果)

次に、実際に入力音声の声質を変換する実験を行って効果を確認した結果について説明する。図13Aは、日本語の入力音声の声質を変換したときの実験結果を示す。ここでは、入力音声は、ある女性話者によって文発声された音声である。また、目標話者は、入力音声を発声した女性話者とは別の女性話者である。図13Aには、その目標話者が孤立発声した母音に基づいて入力音声の声質が変換された結果が示されている。

【0176】

図13Aの(a)は、従来技術で声質変換されたスペクトログラムを示す。図13Aの(b)は、本実施の形態における声質変換システム100により声質変換されたスペクトログラムを示す。本実験では、式(8)における曖昧化度合い係数 $a$ (混合比率)として、「0.3」を用いた。

20

【0177】

また、発話内容は、日本語の「ねえご隠居さん、昔から鶴は千年、亀は万年なんてことを言いますね」( / ne e go i N kyo sa N , mu ka shi ka ra , tsu ru wa se N ne N , ka me wa ma N ne N na N te ko to o i i ma su ne / 、 “ Hi daddy . They say crane lives longer than a thousand years , and tortoise lives longer than ten thousand years , don't they ? ” ) である。

30

【0178】

図13Aの(b)では、(a)と比べて、全体に時間方向のフォルマント軌跡が滑らかになっており、連続発声として自然性が改善している。特に、図13Aにおいて白線で囲んだ部分は、(a)と(b)との間で顕著な違いが見られる。

【0179】

図13Bは、英語の入力音声の声質を変換したときの実験結果を示す。具体的には、図13Bの(a)は、従来技術で声質変換されたスペクトログラムを示す。図13Bの(b)は、本実施の形態における声質変換システム100によって声質変換されたスペクトログラムを示す。

40

【0180】

図13Bにおいて、入力音声の話者と目標話者とは、図13Aと同様である。また、曖昧化度合い係数 $a$ も図13Aと同様である。

【0181】

発話内容は、英語の “ Work hard today . ” である。なお、英語の発話内容は、カタカナの「ワークハードトゥデイ」という文字列に置き換えられ、日本語の音素で合成音が生成されている。

【0182】

声質変換後の音声の韻律(すなわちイントネーションパターン)は、入力音声の韻律と

50

同じになるため、日本語の音素を用いて声質変換しても、声質変換後の音声はある程度英語らしく聞こえる。しかし、英語の母音は日本語に比べて数が多いため、日本語の代表的な母音だけでは、英語の母音を表現できないという問題がある。

【0183】

そこで、本実施の形態に示す技術で母音を曖昧化することによって、日本語らしさを低減し、結果として英語音声として自然さを増すことが可能となる。特に、以下にIPAで示す、曖昧母音である *schwa* は、日本語の5母音とは全く異なり、F1 - F2平面において日本語の5母音によって形成される五角形の重心付近に位置するために、本実施の形態による曖昧化の効果が大きい。

【0184】

【数10】

[ə]

【0185】

図13Bにおいて、特に白線で囲んだ部分は、(a)と(b)との間で顕著な違いが見られる。時刻1.2秒においては、第1及び第2フォルマント周波数だけではなく、第3フォルマント周波数にも違いが現れていることが分かる。実際に合成音を聞いた印象では、(a)はカタカナをそのまま話しているような感じであり、(b)は英語として受け入れやすい感じであった。また、(a)は英語を話すときに力を入れて調音している感じであり、(b)はリラックスして話している感じであった。

【0186】

ところで、発声の怠けは発話速度によって変化する。ゆっくり発話するときには、各母音は、孤立母音と同様に正確に調音される。この特徴は、歌を歌う場合などに顕著に現れる。入力音声が歌声の場合は、声質変換システム100は、孤立母音をそのまま用いて声質変換を行っても、違和感のない合成音を生成することが可能である。

【0187】

一方で、会話調の話し方で早く発話するときには、顎や舌などの調音器官の動きが発話速度に追いつかないために、発声の怠けが大きくなる。そこで、曖昧化度合い(混合比率)は、当該音韻周辺の局所的発話速度に応じて設定されても良い。つまり、混合部105は、入力音声に含まれる母音の局所的発話速度が大きいほど、入力音声に含まれる母音と同じ種類の母音の第2声道形状情報が平均声道形状情報に近づくように、第2声道形状情報を生成しても良い。これにより、入力音声をより滑らかで自然な音声に変換することが可能となる。

【0188】

具体的には、式(8)の曖昧化度合い係数  $a$  (混合比率) は、例えば次の式(9)のように局所的発話速度  $r$  (単位は1秒あたりの音素数など) の関数として設定されれば良い。

【0189】

【数11】

$$a = a_0 + h(r - r_0) \quad (9)$$

【0190】

ここで、 $a_0$  は基準の曖昧化度合いを表す値であり、 $r_0$  は基準の発話速度(単位は  $r$  と同じ)である。また、 $h$  は所定の値であり、 $r$  によって  $a$  を変化させる感度である。

【0191】

なお、文中母音は、F1 - F2平面において、孤立母音よりも多角形の内側に移動するが、その度合いは母音によって異なる。例えば図4A及び図4Bにおいて、/o/ は比較的变化が少ないが、/a/ は少数の外れ値を除いて大きく内側に移動している。また、/i/ も多くが特定の方向に移動しているが、/u/ は移動する方向もまちまちである。

【0192】

10

20

30

40

50

そこで、母音によって曖昧化度合い（混合比率）を変えることも有効と考えられる。つまり、混合部 105 は、母音の種類に応じて設定された混合比率を用いて、母音の種類毎に、当該母音の第 1 声道形状情報と、当該母音と異なる種類の母音の第 1 声道形状情報とを混合しても良い。この場合、/o/ の曖昧化度合いを小さく、/a/ の曖昧化度合いを大きくしても良い。また /i/ は曖昧化度合いを大きく、/u/ はどちらの方向に移動させれば良いか分からないために曖昧化度合いを小さくしても良い。これらの傾向は個人によって異なる可能性があるので、目標話者が誰であるかによって曖昧化度合いを変えても良い。

【0193】

もちろん、曖昧化度合いは、ユーザーの好みによって変えられても構わない。この場合、ユーザーは、混合比率入力部 110 を介して、母音の種類毎に、好みの曖昧化度合いを示す混合比率を入力すれば良い。つまり、混合部 105 は、ユーザーによって設定された混合比率を用いて、母音の種類毎に、当該母音の第 1 声道形状情報と、当該母音と異なる種類の母音の第 1 声道形状情報とを混合しても良い。

10

【0194】

また、平均声道情報算出部 1051 は、式(7)に示すように、複数の第 1 声道形状情報の算術平均（相加平均）を算出することにより、平均声道形状情報を算出したが、必ずしも式(7)のように平均声道形状情報を算出する必要はない。例えば、平均声道情報算出部 1051 は、式(6)の重み係数  $w_i$  を不均一にして、平均声道形状情報を算出しても良い。

20

【0195】

つまり、平均声道形状情報は、互いに種類が異なる複数の母音の第 1 声道形状情報の重み付き算術平均であっても構わない。例えば、個人ごとに発声の怠けの特徴を調べて、その個人の発声の怠けを近似するように重み係数の調整を行なうことは効果的である。例えば、目標話者の発声の怠けの特徴に応じて第 1 声道形状情報に重み付けすることにより、入力音声より滑らかで自然な目標話者の音声に変換することも可能となる。

【0196】

また、平均声道情報算出部 1051 は、式(7)のような相加平均ではなく、相乗平均や調和平均を平均声道形状情報として算出しても構わない。具体的には、式(10)のように PARCOR 係数の平均ベクトルを表すと、平均声道情報算出部 1051 は、式(11)のように、複数の母音の第 1 声道形状情報の相乗平均を平均声道形状情報として算出しても良い。また、平均声道情報算出部 1051 は、式(12)のように、複数の母音の第 1 声道形状情報の調和平均を平均声道形状情報として算出しても良い。

30

【0197】

【数12】

$$\bar{\mathbf{k}} = (\bar{k}_1 \quad \bar{k}_2 \quad \dots \quad \bar{k}_M) \quad (10)$$

【0198】

【数13】

40

$$\bar{k}_m = \sqrt[N]{\prod_{i=1}^N k_m^i} = \sqrt[N]{k_m^1 k_m^2 \dots k_m^N} \quad (11)$$

【0199】



【数 1 4】

$$\bar{k}_m = \frac{N}{\sum_{i=1}^N \frac{1}{k_m^i}} = \frac{N}{\frac{1}{k_m^1} + \frac{1}{k_m^2} + \dots + \frac{1}{k_m^N}} \quad (12)$$

【0200】

要するに、複数の母音の第1声道形状情報の平均は、各母音の第1声道形状情報と混合されたときに、F1 - F2平面における母音の分布範囲が縮小されるように算出されれば良い。

10

【0201】

例えば日本語の5母音 / a /、/ i /、/ u /、/ e /、/ o / の場合、式(7)や式(11)、式(12)のような平均声道形状を求めることは必ずしも必要ではない。例えば、ある母音と別の母音を混合することによってその母音を五角形の重心に近づける操作が行なわれても良い。例えば母音 / a / のあいまい化を行う場合、/ a / とは別の種類の母音を少なくとも2つ選び、選ばれた2つの母音を用いて所定の重みで混合を行っても良い。F1 - F2平面上で5母音が形成する五角形が凸五角形(全ての内角の大きさが二直角より小さい五角形)であれば、/ a / と他の任意の2つの母音を混合して作られた母音は必ずこの五角形の内側に位置する。多くの場合、日本語の5母音形成する五角形は凸五角形であり、この方法によって母音を曖昧化できる。

20

【0202】

また、上述したように英語には日本語よりも母音の数が多いため、F1 - F2平面において母音間の距離が小さい傾向にある。この傾向は言語によって異なるので、曖昧化度合い係数は、言語に応じて設定されることが望ましい。つまり、混合部105は、入力音声の言語種類に応じて定められた混合比率を用いて、母音の種類毎に、当該母音の第1声道形状情報と、当該母音と異なる種類の母音の第1声道形状情報とを混合しても良い。これにより、各言語にふさわしい曖昧化度合いを設定することができ、入力音声をより滑らかで自然な音声に変換することが可能となる。

【0203】

英語の母音種類は日本語よりも多いため、F1 - F2平面での多角形は日本語の多角形よりも複雑である。図14は、F1 - F2平面に英語の13母音を配置した図である。なお、図14は、「Ghonim, A., Smith, J. and Wolfe, J. (2007) "The sounds of world English", <http://www.phys.unsw.edu.au/swe>」から引用した。英語では母音のみを発声することは難しいので、[h]と[d]で挟まれた仮想的な単語で母音が表されている。13母音を全て加算平均して求めた平均声道形状と各母音を混合した場合、各母音が重心に近づく方向に移動するため曖昧化される。

30

【0204】

しかし、日本語の場合に述べたように、全ての母音を用いて平均声道形状を求めることは必ずしも必要ではない。図14の配置を用いると、“heed”、“haired”、“had”、“hard”、“hod”、“howd”、“whod”を用いて凸多角形を構成することができる。この多角形の辺に近い母音は日本語と同様に、当該母音をそれとは別の少なくとも2母音を選び混合することで曖昧化が可能である。一方、多角形の内部に位置する母音(図では“heard”)については、それらがもともと曖昧な音であるためにそのまま利用する。

40

【0205】

このように、本実施の形態における声質変換システム100によれば、少量の母音を入力するだけで滑らかな文発声の音声を生成することができる。さらに、日本語母音を用いて英語の音声を生成することができるなど、飛躍的に柔軟な声質変換が可能になる。

【0206】

50

つまり、本実施の形態における声質変換システム100によれば、母音の種類毎に、複数の第1声道形状情報を混合して第2声道形状情報を生成することができる。つまり、少量の音声のサンプルから母音の種類毎に第2声道形状情報を生成することができる。このように母音の種類毎に生成された第2声道形状情報は、曖昧化された母音の声道形状情報に相当する。したがって、第2声道形状情報を用いて入力音声の声質を変換することにより、入力音声を滑らかで自然な音声に変換することが可能となる。

#### 【0207】

なお、母音受付部102は、前述したとおり典型的にはマイクロホンを有するが、さらに、ユーザーに発声内容とタイミングとを指示するための表示装置(prompter)を有することが望ましい。具体例としては、図15に示すように、母音受付部102は、  
10  
マイクロホン1021と、マイクロホン1021の近傍に配置された液晶ディスプレイなどの表示部1022とから構成されても良い。この場合、表示部1022は、目標話者に発声させる内容1023(この場合は母音)とタイミング1024とを表示すれば良い。

#### 【0208】

なお、本実施の形態では、混合部105は、平均声道形状情報を算出していたが、必ずしも平均声道形状情報を算出する必要はない。例えば、混合部105は、母音の種類毎に、当該母音の第1声道形状情報と、当該母音とは異なる種類の母音の声道形状情報とを所定の混合比率で混合することにより、当該母音の第2声道形状情報を生成すれば良い。このとき、所定の混合比率は、第2声道形状情報が第1声道形状情報よりも平均声道形状情報に近づくように設定されれば良い。  
20

#### 【0209】

つまり、混合部105は、F1-F2平面上で母音間の距離が近づくように第2声道形状情報が生成されれば、どのように複数の第1声道形状情報が混合されても構わない。例えば、混合部105は、入力音声においてある母音から別の母音に遷移する時に声道形状情報が急峻に変わらないように第2声道形状情報を生成しても良い。つまり、混合部105は、入力音声に含まれる母音の並びに適応して混合比率を変化させながら、入力音声に含まれる母音と同じ種類の母音の第1声道形状情報と、入力音声に含まれる母音と異なる種類の母音の第1声道形状情報とを混合しても良い。その結果、第2声道形状情報から得られる母音のF1-F2平面における位置は、同じ種類の母音であっても、多角形領域内で動くことになる。これは、PARCOR係数の時系列を移動平均法などにより平滑化することで実現可能である。  
30

#### 【0210】

(実施の形態1の変形例)

次に、実施の形態1の変形例について説明する。

#### 【0211】

実施の形態1では、母音受付部102は、当該言語における代表的な全ての種類の母音(日本語では5母音)を受け付けていたが、本変形例では、母音受付部102は、必ずしも全ての種類の母音を受け付ける必要はない。本変形例では、実施の形態1よりも少ない種類の母音で声質変換を実現する。以下、その方法について説明する。

#### 【0212】

母音の種類は第1フォルマント周波数と第2フォルマント周波数とで特徴付けられるが、それらの値は個人によって異なっている。それでも、同一の母音と知覚される理由を説明するモデルとして、第1フォルマント周波数と第2フォルマント周波数との比によって母音が特徴付けられるとみなしたモデルがある。ここで、第*i*母音の第1フォルマント周波数 $f_{1i}$ 及び第2フォルマント周波数 $f_{2i}$ からなるベクトル $v_i$ を式(13)で表すとし、第1フォルマント周波数と第2フォルマント周波数との比を保ったままベクトル $v_i$ を移動したベクトル $v_i'$ を式(14)で表すとする。  
40

#### 【0213】

【数 1 5】

$$\mathbf{v}_i = [f1_i \quad f2_i] \quad (13)$$

【0 2 1 4】

【数 1 6】

$$\mathbf{v}'_i = q\mathbf{v}_i = q[f1_i \quad f2_i] = [qf1_i \quad qf2_i] \quad (14)$$

【0 2 1 5】

q はベクトル  $\mathbf{v}_i$  とベクトル  $\mathbf{v}'_i$  との比率である。上述のモデルに基づけば、比率 q の値を変化させてもベクトル  $\mathbf{v}_i$  とベクトル  $\mathbf{v}'_i$  とは同じ母音として知覚される。

10

【0 2 1 6】

このように、全ての孤立母音の第 1 及び第 2 フォルマント周波数を比率 q で移動した場合、F 1 - F 2 平面上で母音の第 1 及び第 2 フォルマント周波数によって形成される多角形は、図 1 6 に示すように互いに相似となる。図 1 6 では、元の多角形 A と、 $q > 1$  の時の多角形 B と、 $q < 1$  の時の多角形 C 及び D とが表されている。

【0 2 1 7】

このように第 1 フォルマント周波数  $f 1_i$  と第 2 フォルマント周波数  $f 2_i$  との比を保ったまま声道形状を変形する方法としては、声道の長さを変更するという方法がある。声道長を  $1 / q$  倍にすれば、全てのフォルマントの周波数が q 倍になる。そこで、まず声道長変換比率  $r = 1 / q$  を求め、次に声道長変換比率 r で声道断面積関数を伸縮するような変換を行なう。

20

【0 2 1 8】

まず、声道長変換比率 q を求める方法について説明する。

【0 2 1 9】

P A R C O R 係数は、分析次数が十分高ければ高次の係数になるに従って絶対値が小さくなる傾向にある。特に、声帯の位置に相当するセクション番号以上の次数では小さな値が続く。そこで、高次の係数から順に低い次数へと値を検査し、絶対値がある閾値を超えたところを声帯位置とみなし、その次数 k を記憶しておく。この方法により、あらかじめ用意された母音から取り出した k を  $k_a$ 、入力された母音から取り出した k を  $k_b$  とすれば、声道長変換比率 r は、式 ( 1 5 ) のように計算することができる。

30

【0 2 2 0】

【数 1 7】

$$r = \frac{kb}{ka} \quad (15)$$

【0 2 2 1】

次に、声道長変換比率 r で声道断面積関数を伸縮する変換方法について説明する。

40

【0 2 2 2】

図 1 7 は、ある母音の声道断面積関数を示す。横軸は、口唇から声帯へ向かっての距離をセクション番号で表す。縦軸は、声道断面積を表す。破線は、声道断面積をスプライン関数などにより内挿して連続値にしたものである。

【0 2 2 3】

連続値になった声道断面積関数を新たなセクション間隔  $1 / r$  でサンプリングし ( 図 1 8 )、サンプリングされた値を元のセクション間隔で配置しなおす ( 図 1 9 )。図 1 9 の例では、声道末端部分 ( 声帯側 ) に余剰セクションが生まれるが ( 図 1 9 の網掛け部分 )、余剰セクションの部分は一定の断面積にしておく。これは、声道長を超えるセクションでは P A R C O R 係数の絶対値が非常に小さい値になるからである。つまり、P A R C O

50

R係数の符号を反転したものはセクション間の反射係数であり、反射係数が0であると言うことはセクション間の断面積に差がないことを意味するからである。

【0224】

上記の例では、声道長を短くする場合 ( $r < 1$ ) の変換方法を示した。一方、声道長を長くする場合 ( $r > 1$ ) は、声道末端部分 (声帯側) には収まりきらないセクションが生まれるが、これらのセクションの値は捨てる。捨てるPARCOR係数の絶対値が小さくなるように、元々の分析次数を高めにとっておくが良い。例えばサンプリング周波数10kHzの音声に対して通常のPARCOR分析では次数を10前後にするが、20などの高い値にしておけば良い。

【0225】

このような方法で、入力された単一の母音と、あらかじめ用意された母音から、全ての母音の声道形状情報を推定することが可能である。つまり、母音受付部102は、全ての種類の母音を受け付ける必要がなくなる。

【0226】

(実施の形態2)

次に、実施の形態2について説明する。

【0227】

本実施の形態では、声質変換システムが2つの装置によって構成される点が、実施の形態1における声質変換システムと異なる。以下において、実施の形態1と異なる点を中心に説明する。

【0228】

図20は、実施の形態2における声質変換システム200の構成図である。図20において、図8と同じ機能を有する構成要素については同じ符号を用い、適宜説明を省略する。

【0229】

図20に示すように、声質変換システム200は、声道情報生成装置201と声質変換装置202とを備える。

【0230】

声道情報生成装置201は、入力音声の声質を変換する際に用いられる、声道の形状を示す第2声道形状情報を生成する。声道情報生成装置201は、母音受付部102と、分析部103と、第1母音声道情報記憶部104と、混合部105と、混合比率入力部110と、第2母音声道情報記憶部107と、合成部108aと、出力部109とを備える。

【0231】

合成部108aは、母音の種類毎に、第2母音声道情報記憶部107に記憶されている第2声道形状情報を用いて合成音を生成する。そして、合成部108aは、生成した合成音の信号を出力部109に送信する。声道情報生成装置201の出力部109は、母音の種類毎に生成された合成音の信号を音声として出力する。

【0232】

図21は、実施の形態2における声道情報生成装置201が出力する母音の音声を説明するための図である。図21では、声道情報生成装置201の母音受付部102によって受け付けられる複数の母音の音声によりF1 - F2平面に形成される五角形を実線で表わす。また、声道情報生成装置201の出力部109によって母音の種類毎に出力される音声によりF1 - F2平面に形成される五角形を破線で表わす。

【0233】

図21から明らかなように、声道情報生成装置201の出力部109は、曖昧化された母音の音声を出力する。

【0234】

声質変換装置202は、声道形状情報を用いて入力音声の声質を変換する。声質変換装置202は、母音受付部102と、分析部103と、第1母音声道情報記憶部104と、入力音声記憶部101と、合成部108bと、変換比率入力部111と、出力部109と

10

20

30

40

50

を備える。この声質変換装置 202 は、図 25 に示す特許文献 2 の声質変換装置と同様の構成である。

【0235】

合成部 108b は、第 1 母音声道情報記憶部 104 に記憶されている第 1 声道形状情報を用いて、入力音声の声質を変換する。ただし、本実施の形態では、声質変換装置 202 の母音受付部 102 は、声道情報生成装置 201 によって曖昧化された母音の音声を受け付けている。つまり、声質変換装置 202 の第 1 母音声道情報記憶部 104 に記憶されている第 1 声道形状情報は、実施の形態 1 における第 2 声道形状情報に相当する。したがって、声質変換装置 202 の出力部 109 は、実施の形態 1 と同様の音声を出力する。

【0236】

以上のように、本実施の形態における声質変換システム 200 によれば、声道情報生成装置 201 と声質変換装置 202 との 2 つの装置によって構成することができる。そして、声質変換装置 202 は、従来の声質変換装置と同様の構成にすることができる。つまり、本実施の形態における声質変換システム 200 によれば、実施の形態 1 と同様の効果を、従来の声質変換装置を用いて実現することが可能となる。

【0237】

(実施の形態 3)

次に、実施の形態 3 について説明する。

【0238】

本実施の形態では、声質変換システムが 2 つの装置によって構成される点が、実施の形態 1 における声質変換システムと異なる。以下において、実施の形態 1 と異なる点を中心に説明する。

【0239】

図 22 は、実施の形態 3 における声質変換システム 300 の構成図である。図 22 において、図 8 と同じ機能を有する構成要素については同じ符号を用い、適宜説明を省略する。

【0240】

図 22 に示すように、声質変換システム 300 は、声道情報生成装置 301 と声質変換装置 302 とを備える。

【0241】

声道情報生成装置 301 は、第 1 母音声道情報記憶部 104 と、混合部 105 と、混合比率入力部 110 とを備える。声質変換装置 302 は、入力音声記憶部 101 と、母音受付部 102 と、分析部 103 と、合成部 108 と、出力部 109 と、変換比率入力部 111 と、母音声道情報記憶部 303 と、母音声道情報入出力切替部 304 とを備える。

【0242】

母音声道情報入出力切替部 304 は、第 1 のモード又は第 2 のモードで動作する。具体的には、母音声道情報入出力切替部 304 は、第 1 のモードでは、母音声道情報記憶部 303 に記憶されている第 1 声道形状情報を第 1 母音声道情報記憶部 104 に出力する。一方、母音声道情報入出力切替部 304 は、第 2 のモードでは、混合部 105 から出力された第 2 声道形状情報を、母音声道情報記憶部 303 に格納する。

【0243】

母音声道情報記憶部 303 には、第 1 声道形状情報及び第 2 声道形状情報が格納される。つまり、母音声道情報記憶部 303 は、実施の形態 1 における第 1 母音声道情報記憶部 104 及び第 2 母音声道情報記憶部 107 に相当する。

【0244】

以上、本実施の形態における声質変換システムによれば、母音を曖昧化する機能を有する声道情報生成装置 301 を独立した装置として構成することができる。そして、声道情報生成装置 301 は、マイクロホンなどが不要であるので、コンピュータソフトウェアとして実現することができる。したがって、声道情報生成装置 301 は、声質変換装置 302 の性能を高めるために後付けするソフトウェア（いわゆるプラグイン）として提供する

10

20

30

40

50

ことができる。

【0245】

また、声道情報生成装置301は、サーバアプリケーションとして実現することもできる。この場合、声道情報生成装置301は、ネットワークを介して声質変換装置302と接続されれば良い。

【0246】

以上、本発明の一態様に係る声質変換システム、声質変換装置、及び声道情報生成装置について、実施の形態に基づいて説明したが、本発明は、これらの実施の形態に限定されるものではない。本発明の趣旨を逸脱しない限り、当業者が思いつく各種変形を本実施の形態に施したものの、あるいは異なる実施の形態における構成要素を組み合わせて構築される形態も、本発明の範囲内に含まれる。

10

【0247】

例えば、上記実施の形態1～3において、声質変換システムは、複数の構成要素を備えていたが、必ずしもそれらの構成要素のすべてを備える必要はない。例えば、声質変換システムは、図23に示すように構成されても良い。

【0248】

図23は、他の実施の形態に係る声質変換システム400の構成図である。なお、図23において、図8と同様の構成要素については、同一の符号を付し、適宜説明を省略する。

【0249】

図23に示す声質変換システム400は、声道情報生成装置401と声質変換装置402とを備える。なお、図23において、図8と同様の構成要素については、同一の符号を付し、説明を省略する。

20

【0250】

図23に示す声質変換システム400は、分析部103及び混合部105を有する声道情報生成装置401と、第2母音声道情報記憶部107及び合成部108を有する声質変換装置402とを備える。なお、声質変換システム400は、必ずしも第2母音声道情報記憶部107を備える必要はない。

【0251】

声質変換システム400は、このように構成されても、曖昧化された声道形状情報である第2声道形状情報を用いて入力音声の声質を変換することができるので、実施の形態1における声質変換システム100と同様の効果を奏することができる。

30

【0252】

また、上記各実施の形態における声質変換システム、声質変換装置、又は声道情報生成装置が備える構成要素の一部又は全部は、1個のシステムLSI(Large Scale Integration:大規模集積回路)から構成されているとしても良い。

【0253】

システムLSIは、複数の構成要素を1個のチップ上に集積して製造された超多機能LSIであり、具体的には、マイクロプロセッサ、ROM(Read Only Memory)、RAM(Random Access Memory)などを含んで構成されるコンピュータシステムである。前記ROMには、コンピュータプログラムが記憶されている。前記マイクロプロセッサが、前記コンピュータプログラムに従って動作することにより、システムLSIは、その機能を達成する。

40

【0254】

なお、ここでは、システムLSIとしたが、集積度の違いにより、IC、LSI、スーパーLSI、ウルトラLSIと称されることもある。また、集積回路化の手法はLSIに限るものではなく、専用回路又は汎用プロセッサで実現しても良い。LSI製造後に、プログラムすることが可能なFPGA(Field Programmable Gate Array)、あるいはLSI内部の回路セルの接続や設定を再構成可能なリプログラマブル・プロセッサを利用しても良い。

50

## 【 0 2 5 5 】

さらには、半導体技術の進歩又は派生する別技術により L S I に置き換わる集積回路化の技術が登場すれば、当然、その技術を用いて機能ブロックの集積化を行っても良い。バイオ技術の適用等が可能性としてありえる。

## 【 0 2 5 6 】

また、本発明の一態様は、このような特徴的な構成要素を備える声質変換システム、声質変換装置、又は声道情報生成装置だけでなく、声質変換システム、声質変換装置、又は声道情報生成装置に含まれる特徴的な処理部をステップとする声質変換方法又は声道情報生成方法であっても良い。また、本発明の一態様は、声質変換方法又は声道情報生成方法に含まれる特徴的な各ステップをコンピュータに実行させるコンピュータプログラムであ

10

## 【産業上の利用可能性】

## 【 0 2 5 7 】

本発明の一態様に係る声質変換システムは、音声加工ツール、ゲーム、家電製品等の音声ガイド、ロボットの音声出力等として有用である。また、ある人の声を別の人の声に変換する用途ではなくとも、テキスト音声合成の出力を滑らかに聞きやすい印象にするための用途にも応用できる。

20

## 【符号の説明】

## 【 0 2 5 8 】

1 0 0、2 0 0、3 0 0、4 0 0 声質変換システム

1 0 1 入力音声記憶部

1 0 2 母音受付部

1 0 3 分析部

1 0 4 第1母音声道情報記憶部

1 0 5 混合部

1 0 7 第2母音声道情報記憶部

1 0 8、1 0 8 a、1 0 8 b 合成部

1 0 9 出力部

1 1 0 混合比率入力部

1 1 1 変換比率入力部

2 0 1、3 0 1、4 0 1 声道情報生成装置

2 0 2、3 0 2、4 0 2 声質変換装置

3 0 3 母音声道情報記憶部

3 0 4 母音声道情報入出力切替部

1 0 2 1 マイクロホン

1 0 2 2 表示部

1 0 3 1 母音安定区間抽出部

1 0 3 2 母音声道情報作成部

1 0 5 1 平均声道情報算出部

1 0 5 2 混合声道情報生成部

1 0 8 1 母音変換部

1 0 8 2 子音選択部

1 0 8 3 声道情報記憶部

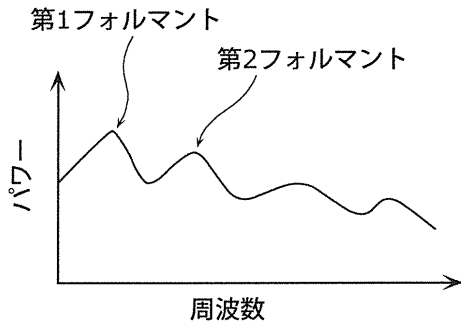
1 0 8 4 子音変形部

1 0 8 5 音声合成部

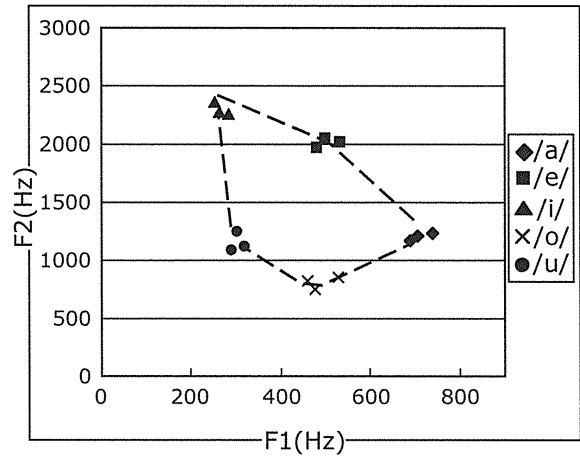
30

40

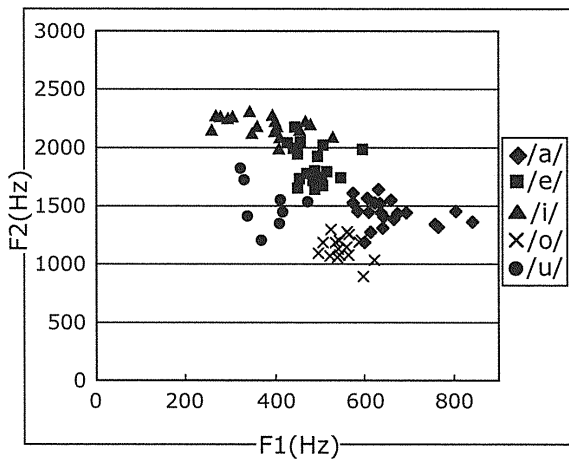
【図1】



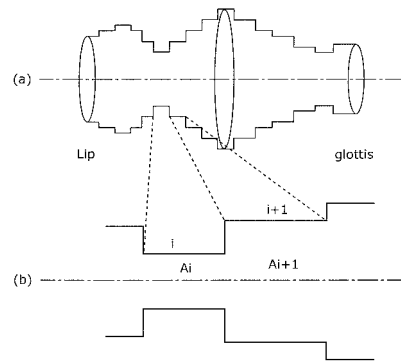
【図2A】



【図2B】

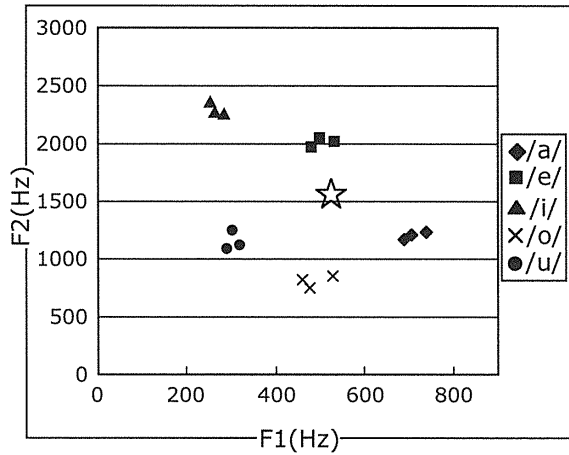


【図3】



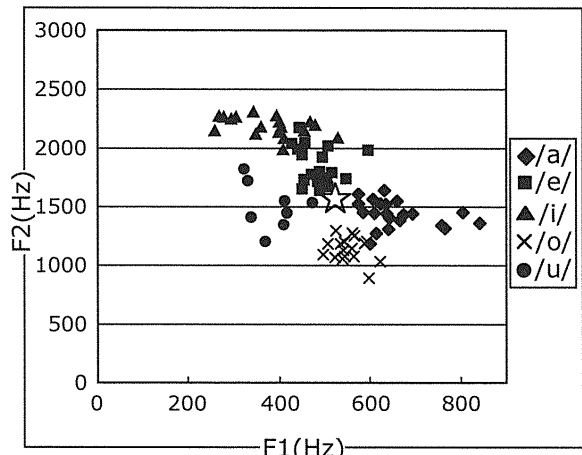


【図 4 A】



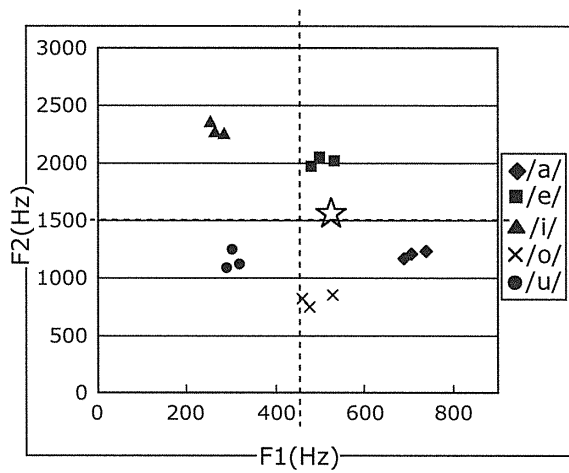
☆ 平均声道形状情報

【図 4 B】

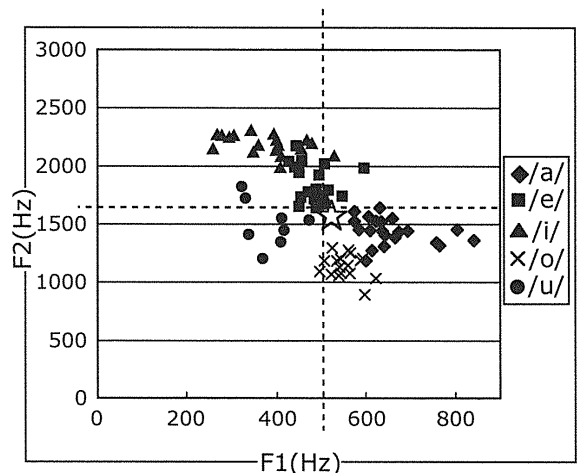


☆ 平均声道形状情報

【図 5 A】



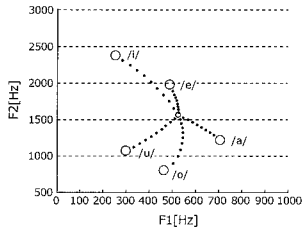
【図 5 B】



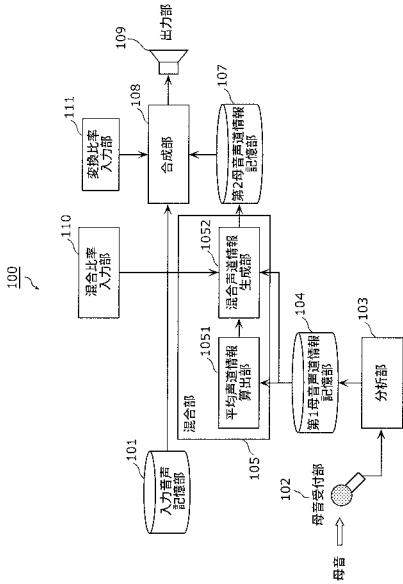
【図 6】

|               | 平方根二乗誤差(RMSE) |
|---------------|---------------|
| 文中母音のF1-F2平均  | 399.5         |
| 孤立母音のF1-F2平均  | 422.5         |
| 孤立母音の平均声道形状情報 | 404.0         |

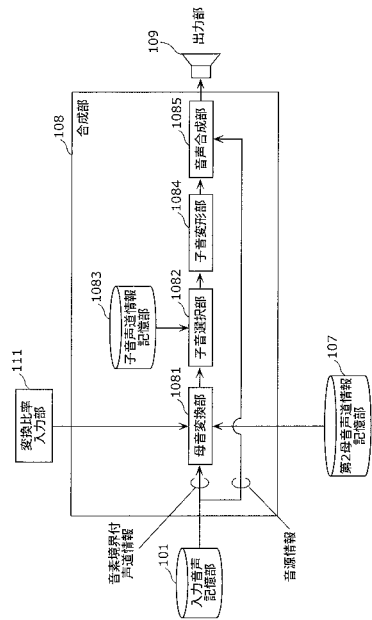
【図7】



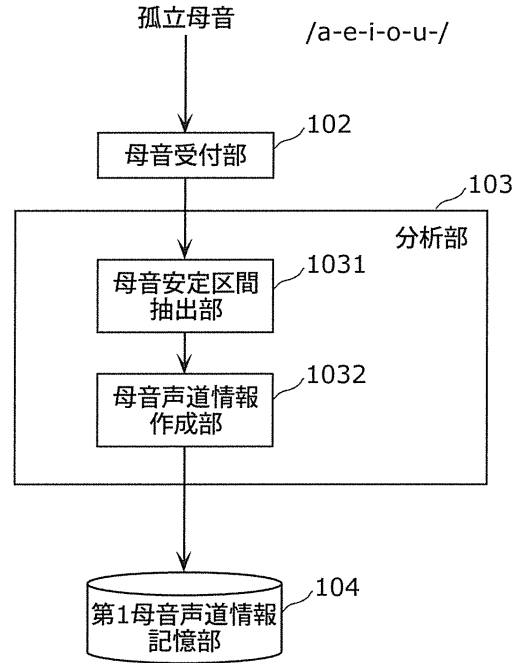
【図8】



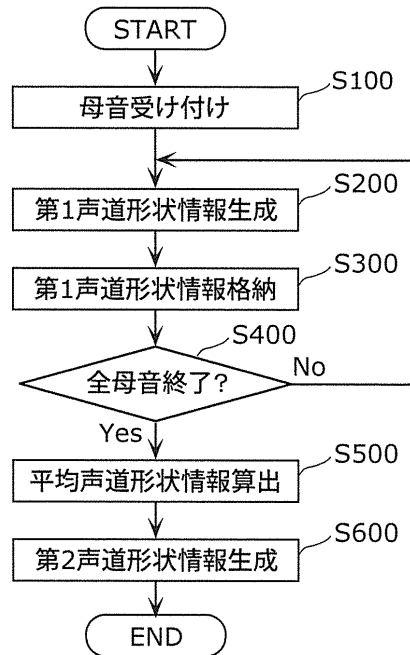
【図10】



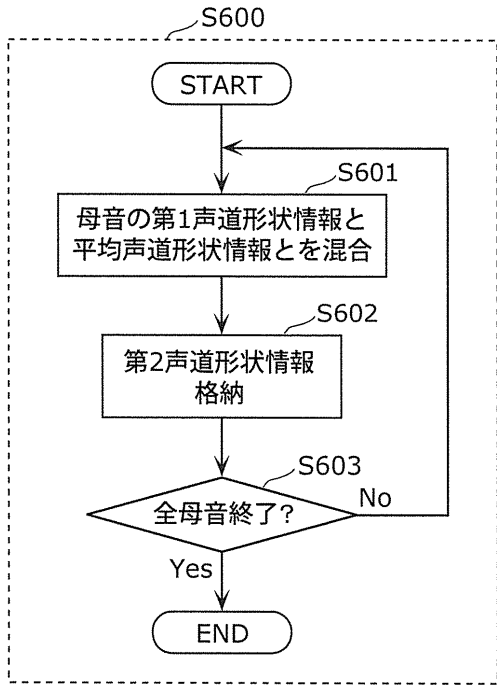
【図9】



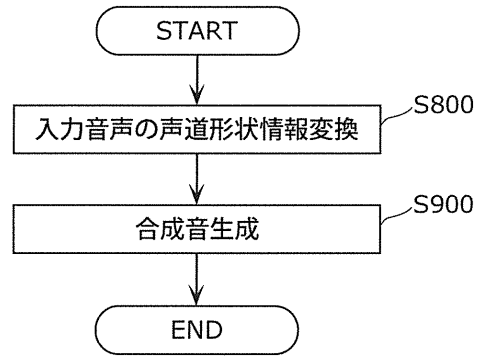
【図11A】



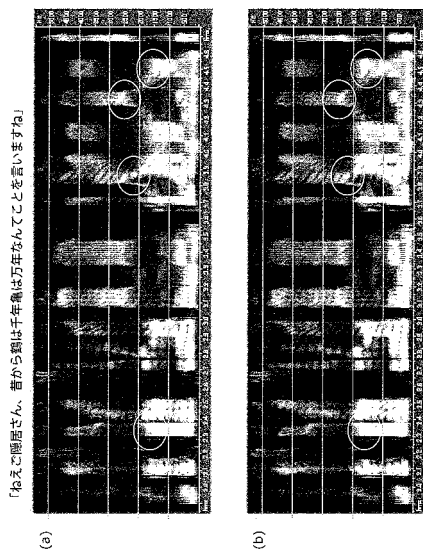
【図 1 1 B】



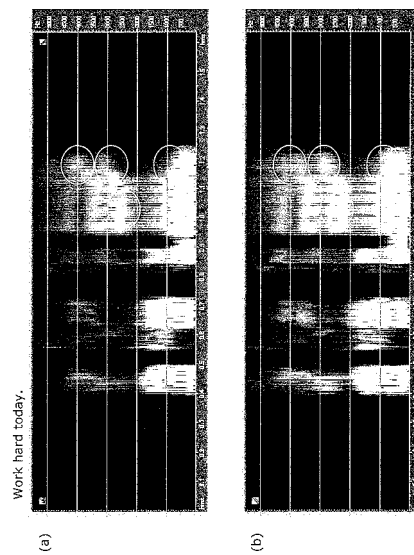
【図 1 2】



【図 1 3 A】

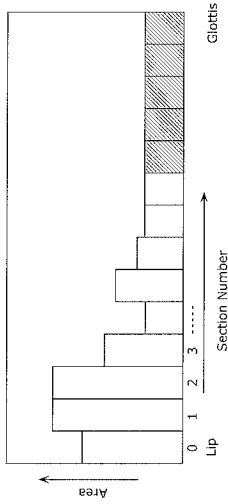


【図 1 3 B】

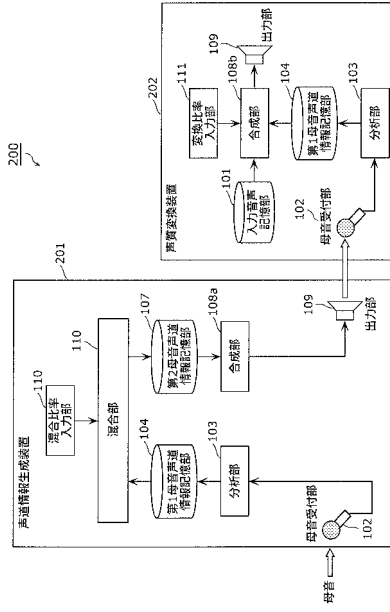




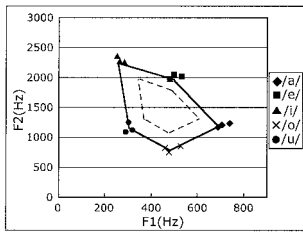
【 図 19 】



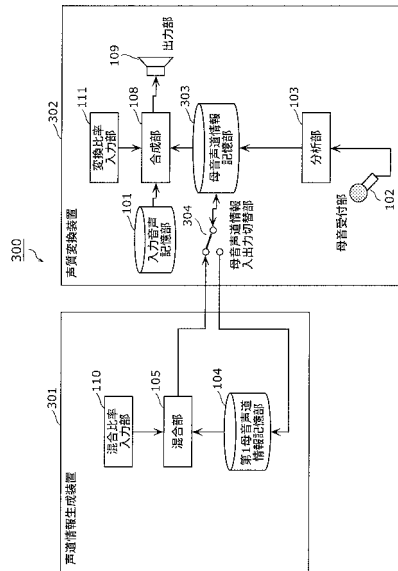
【 図 20 】



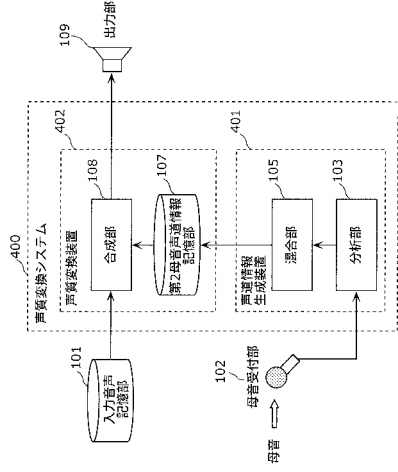
【 図 21 】



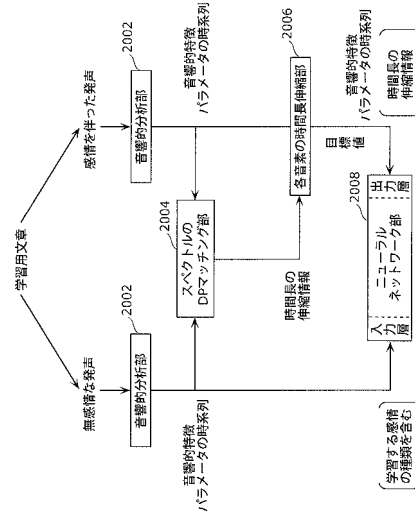
【 図 22 】



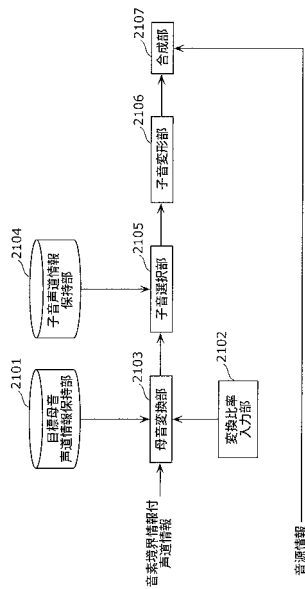
【図 23】



【図 24】



【図 25】



---

フロントページの続き

- (56)参考文献 国際公開第2008/148547(WO, A1)  
特開2007-50143(JP, A)  
特開2001-282300(JP, A)  
特開2006-330343(JP, A)  
国際公開第2008/142836(WO, A1)  
国際公開第2010/035438(WO, A1)

- (58)調査した分野(Int.Cl., DB名)  
G10L 21/00-25/93