

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4829477号  
(P4829477)

(45) 発行日 平成23年12月7日(2011.12.7)

(24) 登録日 平成23年9月22日(2011.9.22)

(51) Int.Cl.		F I	
<b>G 1 0 L</b>	<b>13/08</b>	<b>(2006.01)</b>	G 1 0 L 13/08 1 3 0 B
<b>G 1 0 L</b>	<b>11/00</b>	<b>(2006.01)</b>	G 1 0 L 11/00 1 0 1
<b>G 1 0 L</b>	<b>13/06</b>	<b>(2006.01)</b>	G 1 0 L 13/06 2 0 0
<b>G 1 0 L</b>	<b>21/04</b>	<b>(2006.01)</b>	G 1 0 L 21/04 1 2 0 A

請求項の数 27 (全 34 頁)

(21) 出願番号	特願2004-79079 (P2004-79079)	(73) 特許権者	000004237
(22) 出願日	平成16年3月18日 (2004.3.18)		日本電気株式会社
(65) 公開番号	特開2005-266349 (P2005-266349A)		東京都港区芝五丁目7番1号
(43) 公開日	平成17年9月29日 (2005.9.29)	(74) 代理人	100080816
審査請求日	平成19年1月15日 (2007.1.15)		弁理士 加藤 朝道
審判番号	不服2010-7652 (P2010-7652/J1)	(72) 発明者	三井 康行
審判請求日	平成22年4月12日 (2010.4.12)		東京都港区芝五丁目7番1号 日本電気株式会社内
			合議体
			審判長 板橋 通孝
			審判官 加藤 恵一
			審判官 溝本 安展

最終頁に続く

(54) 【発明の名称】 声質変換装置および声質変換方法ならびに声質変換プログラム

(57) 【特許請求の範囲】

【請求項1】

変換目標となる話者の音声（「目標話者音声」という）を入力する目標話者音声入力部と、

音声合成用のデータを記憶する音声合成用データ記憶部と、

前記入力された目標話者音声と同一又は類似の発声内容を記述する発音記号列を入力する発音記号入力部と、

前記入力された発音記号列にしたがって、前記音声合成用データ記憶部に記憶されている音声合成用のデータに基づき、音声を合成して出力する音声合成部と、

前記音声合成部から出力される音声信号（「合成音」という）の特徴パラメータを抽出する合成音特徴パラメータ抽出部と、

前記目標話者音声の特徴パラメータを抽出する目標話者音声特徴パラメータ抽出部と、

前記合成音特徴パラメータ抽出部からの前記合成音の特徴パラメータと、目標話者音声特徴パラメータ抽出部からの前記目標話者音声の特徴パラメータとを入力し、特徴パラメータを表す空間において、前記合成音の第1の部分と、前記目標話者音声の第2の部分との時間軸上の対応関係を求め、前記合成音の前記第1の部分におけるスペクトル形状を、前記目標話者音声の前記第2の部分におけるスペクトル形状に変換する変換関数の同定を行う変換関数生成部と、

を備えている、ことを特徴とする声質変換装置。

【請求項2】

10

20

変換対象の音声信号を入力し、前記変換関数生成部で同定された前記変換関数を用いて、前記変換対象となる音声信号を、前記目標話者音声の声質を持つ音声信号に変換して出力する声質変換部をさらに備えている、ことを特徴とする請求項 1 記載の声質変換装置。

【請求項 3】

前記目標話者音声入力部により入力された前記目標話者音声を入力して音声認識し、前記発音記号列を生成して前記音声合成部へ出力する音声認識部をさらに備えている、ことを特徴とする請求項 1 記載の声質変換装置。

【請求項 4】

前記音声合成部で前記合成音が生成される際に用いられた、少なくとも時刻情報に対応づいた音素列を含むセグメンテーション情報を記憶し、前記変換関数生成部へ出力する第 1 のセグメンテーション情報記憶部をさらに備えている、ことを特徴とする請求項 1 ~ 3 のいずれかーに記載の声質変換装置。

10

【請求項 5】

前記音声合成部で作成される前記合成音の、少なくとも時刻情報に対応づいた音素列を含むセグメンテーション情報を入力して、前記音声合成部へ出力するセグメンテーション情報入力部をさらに備えている、ことを特徴とする請求項 1 ~ 3 のいずれかーに記載の声質変換装置。

【請求項 6】

前記音声認識部で認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含むセグメンテーション情報を記憶し、前記変換関数生成部および前記音声合成部の少なくとも一方へ出力する第 2 のセグメンテーション情報記憶部を備えている、ことを特徴とする請求項 3 又は 4 記載の声質変換装置。

20

【請求項 7】

前記音声合成用データ記憶部に記憶される前記音声合成用のデータとして、互いに異なる複数の話者の音声データを含むことを特徴とする請求項 1 ~ 3 のいずれかーに記載の声質変換装置。

【請求項 8】

請求項 2 記載の声質変換装置の前記声質変換部から出力される変換後の信号を音声合成用データとして記憶するデータ記憶部と、

発声内容を記述する記号列を入力する記号列入力部と、

30

前記入力された記号列にしたがって前記音声合成用データ記憶部内の音声合成用データから音声合成して出力する音声合成出力部と、

を備えている、ことを特徴とする合成音生成装置。

【請求項 9】

請求項 1 記載の声質変換装置における前記合成音の特徴パラメータと前記特徴パラメータに対応して同定された変換関数とを対応付けて記憶する変換データ記憶部と、

被変換音声信号を入力する音声入力部と、

前記入力された被変換音声信号の特徴パラメータを求め、求めた特徴パラメータに対応する前記変換データ記憶部内の前記合成音の特徴パラメータを探索し、前記合成音の特徴パラメータに対応する前記変換関数を読み出し、前記変換関数によって前記被変換音声信号を変換して出力する音声変換出力部と、

40

を備えている、ことを特徴とする合成音生成装置。

【請求項 10】

声質変換装置により声質を変換する方法であって、

声質の変換目標となる話者の音声（「目標話者音声」という）を入力するステップと、前記目標話者音声の発声内容を記述する発音記号列を入力し、記憶部に予め記憶されている音声合成用のデータを用いて、前記発音記号列から、合成音を作成するステップと、前記合成音を分析し、前記合成音の特徴パラメータを抽出するステップと、

前記目標話者音声进行分析し、前記目標話者音声の特徴パラメータを抽出するステップと

50

前記目標話者音声の特徴パラメータと前記合成音の特徴パラメータとの時間軸上の対応付けを行うステップと、

対応付けがなされた前記目標話者音声及び前記合成音の特徴パラメータに基づき、前記合成音のスペクトル形状を、前記目標話者音声のスペクトル形状に、変換するための変換関数を生成するステップと、

を含む、ことを特徴とする声質変換方法。

【請求項 1 1】

前記変換関数を用いて、変換対象となる音声信号を前記目標話者音声の声質を持つ音声信号に変換するステップをさらに含む、ことを特徴とする請求項 1 0 記載の声質変換方法。

10

【請求項 1 2】

入力される前記目標話者音声を音声認識するステップを含み、

前記音声認識の結果から前記発音記号列を生成する、ことを特徴とする請求項 1 0 記載の声質変換方法。

【請求項 1 3】

前記合成音の少なくとも時刻情報に対応づいた音素列を含むセグメンテーション情報を入力して、前記セグメンテーション情報に基づき前記合成音を生成する、ことを特徴とする請求項 1 0 ~ 1 2 のいずれか一に記載の声質変換方法。

【請求項 1 4】

前記合成音の生成時に用いられる、少なくとも時刻情報に対応づいた音素列を含む第 1 のセグメンテーション情報によって、前記時間軸上の対応付けを行うステップを含む、ことを特徴とする請求項 1 0 ~ 1 2 のいずれか一に記載の声質変換方法。

20

【請求項 1 5】

前記音声認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含む第 2 のセグメンテーション情報によって、前記時間軸上の対応付けを行うステップを含む、ことを特徴とする請求項 1 0 ~ 1 2、1 4 のいずれか一に記載の声質変換方法。

【請求項 1 6】

前記音声認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含む第 2 のセグメンテーション情報に基づき、前記合成音を生成するステップを含む、ことを特徴とする請求項 1 0 ~ 1 2、1 4 のいずれか一に記載の声質変換方法。

30

【請求項 1 7】

前記音声認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含む第 2 のセグメンテーション情報によって、前記時間軸上の対応付けを行い、前記第 2 のセグメンテーション情報に基づき、前記合成音を生成するステップを含む、ことを特徴とする請求項 1 0 ~ 1 2、1 4 のいずれか一に記載の声質変換方法。

【請求項 1 8】

声質変換装置により声質を変換する方法であって、

声質の変換目標となる話者の音声（「目標話者音声」という）を入力するステップと、  
入力される前記目標話者音声を音声認識するステップと、

前記音声認識の結果から前記目標話者音声の発声内容を記述する発音記号列を生成するステップと、

40

を含み、さらに、

( a ) 記憶部に予め記憶されている音声合成用のデータを用いて、前記発音記号列から、合成音を作成するステップと、

( b ) 前記合成音と前記目標話者音声とを分析し、それぞれの特徴パラメータを抽出するステップと、

( c ) 前記目標話者音声の特徴パラメータと前記合成音の特徴パラメータとの時間軸上の対応付けを行うステップと、

( d ) 対応付けされた 2 つの前記特徴パラメータに基づいて、前記合成音のスペクトル形状を前記目標話者音声のスペクトル形状に変換する変換関数を生成するステップと、

50

( e ) 前記生成された変換関数を用いて変換対象となる音声の声質を変換するステップと、

( f ) 前記声質の変換結果を、前記音声合成用のデータとして前記記憶部に格納するステップと、

を所定の収束条件に至るまで繰り返すことを特徴とする声質変換方法。

【請求項 19】

声質変換装置を構成するコンピュータに、

声質の変換目標となる話者の音声(「目標話者音声」という)を入力する処理と、

前記目標話者音声の発声内容を記述する発音記号列を入力し、記憶部に予め記憶されている音声合成用のデータを用いて、前記発音記号列から、合成音を作成する処理と、

前記合成音と前記目標話者音声とを分析し、それぞれの特徴パラメータを抽出する処理と、

前記目標話者音声の特徴パラメータと前記合成音の特徴パラメータとの時間軸上の対応付けを行う処理と、

対応付けがなされた前記目標話者音声及び前記合成音の特徴パラメータに基づき、前記合成音のスペクトル形状を、前記目標話者音声のスペクトル形状に、変換するための変換関数を生成する処理と、

を実行させるプログラム。

【請求項 20】

請求項 19 記載のプログラムにおいて、

前記変換関数を用いて、変換対象となる音声信号を前記目標話者音声の声質を持つ音声信号に変換する処理を、さらに前記コンピュータに実行させるプログラム。

【請求項 21】

請求項 19 記載のプログラムにおいて、

入力される前記目標話者音声を音声認識する処理と、

前記音声認識の結果から前記発音記号列を生成する処理とを、前記コンピュータに実行させるプログラム。

【請求項 22】

請求項 19 ~ 21 のいずれかーに記載のプログラムにおいて、

前記合成音の少なくとも時刻情報に対応づいた音素列を含むセグメンテーション情報を入力して、前記セグメンテーション情報に基づき前記合成音を生成する処理を、前記コンピュータに実行させるプログラム。

【請求項 23】

請求項 19 ~ 21 のいずれかーに記載のプログラムにおいて、

前記合成音の生成時に用いられる、少なくとも時刻情報に対応づいた音素列を含む第 1 のセグメンテーション情報によって、前記時間軸上の対応付けを行う処理を、前記コンピュータに実行させるプログラム。

【請求項 24】

請求項 19 ~ 21、23 のいずれかーに記載のプログラムにおいて、

前記音声認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含む第 2 のセグメンテーション情報によって、前記時間軸上の対応付けを行う処理を、前記コンピュータに実行させるプログラム。

【請求項 25】

請求項 19 ~ 21、23 のいずれかーに記載のプログラムにおいて、

前記音声認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含む第 2 のセグメンテーション情報に基づき、前記合成音を生成する処理を、前記コンピュータに実行させるプログラム。

【請求項 26】

請求項 19 ~ 21、23 のいずれかーに記載のプログラムにおいて、

前記音声認識された目標話者音声の、少なくとも時刻情報に対応づいた音素列を含む第

10

20

30

40

50

2のセグメンテーション情報によって、前記時間軸上の対応付けを行い、前記第2のセグメンテーション情報に基づき前記合成音を生成する処理を、前記コンピュータに実行させるプログラム。

【請求項27】

声質変換装置を構成するコンピュータに、  
 声質の変換目標となる話者の音声（「目標話者音声」という）を入力する処理と、  
 入力される前記目標話者音声を音声認識する処理と、  
 前記音声認識の結果から前記目標話者音声の発声内容を記述する発音記号列を生成する処理と、を実行させ、さらに、

（a）記憶部に予め記憶されている音声合成用のデータを用いて、前記発音記号列から合成音を作成する処理と、

（b）前記合成音と前記目標話者音声とを分析し、それぞれの特徴パラメータを抽出する処理と、

（c）前記目標話者音声の特徴パラメータと前記合成音の特徴パラメータとの時間軸上の対応付けを行う処理と、

（d）対応付けされた2つの前記特徴パラメータに基づいて、前記合成音のスペクトル形状を前記目標話者音声のスペクトル形状に変換する変換関数を生成する処理と、

（e）前記生成された変換関数を用いて変換対象となる音声の声質を変換する処理と、

（f）前記声質の変換結果を、前記音声合成用のデータとして前記記憶部に格納する処理と、

を所定の収束条件に至るまで繰り返す、

処理を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、声質変換装置および声質変換方法ならびに声質変換プログラムに関する。

【背景技術】

【0002】

従来から、ある話者が発声した音声を、別の話者の声質を持つ音声へと変換する、声質変換技術についての研究がなされてきた。例えば、特許文献1において、話者Aの声質を話者Bの声質へと変換する技術が開示されている。特許文献1における声質変換法を図38に示す。この声質変換法では、話者Aの音声をLPC分析101によって分析し、話者Aのコードブック103を用い、ベクトル量子化105によって量子化する。また、話者Bの音声をLPC分析102によって分析し、話者Bのコードブック104を用い、ベクトル量子化106によって量子化する。それぞれ量子化されたデータは、時間軸の整合をDPマッチングによる対応付け107によって対応付け（写像）がなされ、写像を元に、ヒストグラムを用いた変換コードブック作成108によってスペクトル変換コードブック109を作成する。話者Aの声質がスペクトル変換コードブック109を用いて声質変換される。

【0003】

さらに近年では、一人ないし複数の話者の音声データから生成ならびに蓄積された音声合成用データベースを用いて、入力されたテキストの内容を音声として出力する音声合成装置において、所望の声質を持つ音声を合成するための声質変換技術も研究されている。この方法の利点は、ある声質を持つ合成音を生成する際に、その声質を持つ話者の音声データをその都度録音してデータベースを生成ならびに蓄積する必要がなく、一人ないし複数の標準話者データベースを予め蓄積しておけば、そのデータベースを目的の声質を持つように声質変換することで、所望の声質を持つ合成音が生じることができるという点にある。

【0004】

例えば、特許文献2において、予め記憶してある音声データベースと目標話者の音声との間の写像コードブックを作成し、これを変換関数として変換された音声データベースを

10

20

30

40

50

テキスト音声合成用データベースとして用いる技術が開示されている。特許文献2における声質変換法を図39に示す。複数の登録話者の音響特徴パラメータを含む音声データベースを予め記憶する。この声質変換装置は、複数の登録話者の音響特徴パラメータを含む音声データベース203とそのコードブック204を予め記憶しておく。入力された目標話者の少なくとも1単語の音声信号に基づいて、声質変換をすべき目標話者に最も近い話者を、複数の登録話者の中から選択する選択手段201と、選択手段201によって選択された話者の音響空間と目標話者の音響空間との間の差分を計算することにより、選択された話者から目標話者への写像コードブック205を計算する生成手段202と、入力された音声合成すべき文字列に基づいて、音声データベース203に記憶された選択された話者の音声の音響特徴パラメータを上記選択された話者のコードブックを用いて量子化し、選択された話者のコードブックと写像コードブックの対応関係に基づいて文字列に対応する目標話者の音声信号の音響特徴パラメータを生成する写像処理手段206と、写像処理手段206によって生成された目標話者の音声信号の音響特徴パラメータに基づいて、文字列に対応する目標話者の音声信号を発生して出力する音声合成手段207とを備えている。また、生成手段202は、移動ベクトル場平滑化法を用いて、選択された話者から目標話者への写像コードブックを計算するようにしている。

10

【0005】

なお、音声信号処理、音声認識の一般的な技術については、非特許文献1、非特許文献2において解説されている。

【0006】

20

【特許文献1】特開平1-97997号公報（図4）

【特許文献2】特開平8-248994号公報（図1）

【非特許文献1】古井貞熙著、「デジタル音声処理」、東海大学出版会、1985年

【非特許文献2】中川聖一著、「確率モデルによる音声認識」、電子情報通信学会、1988年

【発明の開示】

【発明が解決しようとする課題】

【0007】

特許文献1に開示されている声質変換方法では、変換元話者Aから変換目標話者Bへと声質変換する場合に、AとBが全く同一内容の発声をしている原音の音声データ（あるいは特徴量データ）を用いて変換関数を作成する必要があった。このため、この方法で声質変換を実施するためには、話者Aと話者Bの全く同一内容の音声データがその都度必要となってしまう、自由な声質変換ができなかった。

30

【0008】

また、特許文献2に開示されている声質変換および音声合成方法では、目標話者の音声の音響空間と音声データベースの音響空間との間の差分を取って移動ベクトル場平滑化法を用いて写像コードブックを求めるため、データベースの情報量が目標話者の情報量に比べ圧倒的に多い場合、非常に疎な対応付けを平滑化して全体の対応付けとする必要があった。このため、対応付けの精度が低く、変換後の音声の音質が劣化する、あるいは目標音声の音質との類似度が低いという問題点が生じた。

40

【0009】

したがって、本発明の主たる目的は、変換目標話者の発声内容に依存しない自由な声質変換を実現する装置および方法ならびにプログラムを提供することにある。

【0010】

本発明の他の目的は、所望の音質へ高精度に変換できる変換関数の同定を可能とする装置および方法ならびにプログラムを提供することにある。

【0011】

本発明のさらに他の目的は、変換目標話者音声と同一あるいは類似の発声内容を持つ合成音の自動生成を可能とする装置および方法ならびにプログラムを提供することにある。

【課題を解決するための手段】

50

## 【0012】

前記目的を達成するために、本発明に係る声質変換装置は、第1のAspectによれば、変換目標となる話者の音声（「目標話者音声」という）を入力する目標話者音声入力部と、音声合成用のデータを記憶する音声合成用データ記憶部と、入力された目標話者音声と同一又は類似の発声内容を記述する発音記号列を入力する発音記号入力部と、を備える。また、入力された発音記号列にしたがって、音声合成用データ記憶部に記憶されている音声合成用のデータに基づき、音声を合成して出力する音声合成部と、音声合成部から出力される音声信号（「合成音」という）の特徴パラメータを抽出する合成音特徴パラメータ抽出部と、目標話者音声の特徴パラメータを抽出する目標話者音声特徴パラメータ抽出部と、を備える。さらに、合成音特徴パラメータ抽出部からの合成音の特徴パラメータと、目標話者音声特徴パラメータ抽出部からの目標話者音声の特徴パラメータとを入力し、特徴パラメータを表す空間において、合成音の第1の部分と、目標話者音声の第2の部分との時間軸上の対応関係を求め、合成音の第1の部分におけるスペクトル形状を、目標話者音声の第2の部分におけるスペクトル形状に変換する変換関数の同定を行う変換関数生成部と、を備える構成とされる。

10

## 【0013】

また、本発明に係る声質変換方法は、第2のAspectによれば、声質変換装置により声質を変換する方法である。声質変換方法は、声質の変換目標となる話者の音声（「目標話者音声」という）を入力するステップと、目標話者音声の発声内容を記述する発音記号列を入力し、記憶部に予め記憶されている音声合成用のデータを用いて、発音記号列から、合成音を作成するステップと、を含む。また、合成音と目標話者音声とを分析し、それぞれの特徴パラメータを抽出するステップと、目標話者音声の特徴パラメータと合成音の特徴パラメータとの時間軸上の対応付けを行うステップと、対応付けがなされた目標話者音声及び合成音の特徴パラメータに基づき、合成音のスペクトル形状を、目標話者音声のスペクトル形状に、変換するための変換関数を生成するステップと、を含む。

20

## 【0014】

さらに、本発明に係る声質変換方法は、第3のAspectによれば、声質変換装置により声質を変換する方法である。声質変換方法は、声質の変換目標となる話者の音声（「目標話者音声」という）を入力するステップと、入力される目標話者音声を音声認識するステップと、音声認識の結果から目標話者音声の発声内容を記述する発音記号列を生成するステップと、を含む。さらに、(a)記憶部に予め記憶されている音声合成用のデータを用いて、発音記号列から、合成音を作成するステップと、(b)合成音と目標話者音声とを分析し、それぞれの特徴パラメータを抽出するステップと、(c)目標話者音声の特徴パラメータと合成音の特徴パラメータとの時間軸上の対応付けを行うステップと、(d)対応付けされた2つの特徴パラメータに基づいて、合成音のスペクトル形状を目標話者音声のスペクトル形状に変換する変換関数を生成するステップと、(e)生成された変換関数を用いて変換対象となる音声の声質を変換するステップと、(f)声質の変換結果を、音声合成用のデータとして記憶部に格納するステップと、を収束条件に至るまで繰り返す。

30

## 【0015】

また、本発明に係る声質変換プログラムは、第4のAspectによれば、声質変換装置を構成するコンピュータに実行させるプログラムである。このプログラムは、声質の変換目標となる話者の音声（「目標話者音声」という）を入力する処理と、目標話者音声の発声内容を記述する発音記号列を入力し、記憶部に予め記憶されている音声合成用のデータを用いて、発音記号列から、合成音を作成する処理と、を実行させる。また、合成音と目標話者音声とを分析し、それぞれの特徴パラメータを抽出する処理と、目標話者音声の特徴パラメータと合成音の特徴パラメータとの時間軸上の対応付けを行う処理と、対応付けがなされた目標話者音声及び合成音の特徴パラメータに基づき、合成音のスペクトル形状を、目標話者音声のスペクトル形状に、変換するための変換関数を生成する処理と、を実行させる。

40

50

## 【 0 0 1 6 】

さらに、本発明に係る声質変換プログラムは、第5のAspectによれば、声質変換装置を構成するコンピュータに実行させるプログラムである。このプログラムは、声質の変換目標となる話者の音声（「目標話者音声」という）を入力する処理と、入力される目標話者音声を音声認識する処理と、音声認識の結果から目標話者音声の発声内容を記述する発音記号列を生成する処理と、を実行させる。さらに、（a）記憶部に予め記憶されている音声合成用のデータを用いて、発音記号列から合成音を作成する処理と、（b）合成音と目標話者音声とを分析し、それぞれの特徴パラメータを抽出する処理と、（c）目標話者音声の特徴パラメータと合成音の特徴パラメータとの時間軸上の対応付けを行う処理と、（d）対応付けされた2つの特徴パラメータに基づいて、合成音のスペクトル形状を目標話者音声のスペクトル形状に変換する変換関数を生成する処理と、（e）生成された変換関数を用いて変換対象となる音声の声質を変換する処理と、（f）声質の変換結果を、音声合成用のデータとして記憶部に格納する処理と、を収束条件に至るまで繰り返す処理を実行させる。

10

## 【発明の効果】

## 【 0 0 1 7 】

本発明によれば、目標話者音声と同一あるいは類似の発声内容を持つ変換元音声を予め用意することなく、目標話者音声の声質への変換を行うことができる。したがって、使用者の負担が軽減する。この理由は、目標話者音声と同一あるいは類似の発声内容を表す発音記号列を入力することで、目標話者音声と同一あるいは類似の発声内容を持つ合成音を生成することができるためである。

20

## 【 0 0 1 8 】

また、学習用音声の対応付けおよび声質変換を行うための変換関数の同定が高精度にできる。したがって、従来に比べて所望の声質に近い声質を持つ高音質な音声を得られる。この理由の1つは、合成音作成に使用されるデータを目標話者音声の情報から推定することによって、目標話者音声への変換関数を同定しやすいように合成音を作成することができるためである。もう1つの理由は、合成音作成時に使用されたデータを記憶しておくことによって、このデータを音声から抽出する特徴パラメータに加えて変換関数生成に用いることができるためである。

30

## 【 0 0 1 9 】

さらに、目標話者音声と同一あるいは類似の発声内容を表す発音記号列を予め用意することなく、目標話者音声と同一あるいは類似の発声内容の合成音が作成できる。したがって、処理が自動化し、処理速度が向上する。この理由は、目標話者音声を音声認識することによって、目標話者音声と同一あるいは類似の発声内容を表す発音記号列を自動的に生成することができるためである。

## 【発明を実施するための最良の形態】

## 【 0 0 2 0 】

## [第1の実施形態]

図1に、本発明の第1の実施形態に係る声質変換装置のブロック図を示す。声質変換装置は、目標話者音声入力部11と、発音記号列入力部12と、音声合成用データ記憶部13と、音声合成部14と、目標話者音声特徴パラメータ抽出部15と、合成音特徴パラメータ抽出部16と、変換関数生成部17と、声質変換部18とを備えている。

40

## 【 0 0 2 1 】

目標話者音声入力部11は、目標とする声質を持つ音声データを入力する。発音記号列入力部12は、作成したい音声の発声内容が記述されている発音記号列を入力する。この時、入力される発音記号列は、目標話者音声入力部11から入力される目標話者音声と同一あるいは類似の発声内容の音声を合成するように記述してあるものとする。音声合成用データ記憶部13は、音声合成部14で用いる、音声や音節等の単位、時間長情報等のデータを記憶している。音声合成部14は、発音記号列に対応するデータにしたがって、音声合成用データ記憶部13に蓄えられたデータから合成用のデータを算出し、これを用い

50

て音声を作成して出力する。

【0022】

また、目標話者音声特徴パラメータ抽出部15は、目標話者音声入力部11から入力された音声データに対しスペクトル分析を施して、特徴パラメータを抽出する。合成音特徴パラメータ抽出部16は、音声合成部14から入力された音声データに対しスペクトル分析を施して、特徴パラメータを抽出する。ここで、目標話者音声および合成音から抽出する特徴パラメータとしては、例えば非特許文献1に記載してあるようなLPC (linear predictive Coding) 係数、フォルマント周波数ならびにバンド帯域幅、韻律データ等の内、少なくとも1つを含む。変換関数生成部17は、目標話者音声から抽出された特徴パラメータおよび合成音から抽出された特徴パラメータによって、合成音の声質から目標話者音声の声質へと変換する関数を同定する。声質変換部18は、同定された変換関数を用いて、変換対象である被変換入力信号を変換して変換後出力信号を出力する。

10

【0023】

次に、本発明の第1の実施形態に係る声質変換装置の動作を図1および図13を参照して説明する。図13は、本発明の第1の実施形態に係る声質変換装置の動作を示すフローチャートである。まず、ユーザが望む声質を持つ話者の音声データを目標音声として用意し、目標話者音声入力部11に入力する(ステップS11)。次に、目標話者音声入力部11に入力した目標話者音声と同一あるいは類似の発声内容を持つように発音記号列を作成し、発音記号列入力部12に入力する。入力された発音記号列に対応するデータにしたがって音声合成用データ記憶部13から合成用のデータを算出し、音声合成部14において合成音を生成する。この時、目標話者音声と合成音の発声の時間長は、ずれていても構わない。音声合成用データと発音記号との対応付けの規則は、発音記号が指定されると音声合成用データ記憶部13に記憶されている音声素片データや時間長情報等の音声合成用データから当該発音記号に対応した音声合成用データを算出するように予め作成されている(ステップS12)。

20

【0024】

また、目標話者音声特徴パラメータ抽出部15において、目標話者音声に対し分析を行い、目標話者音声の特徴パラメータを抽出する。また、合成音特徴パラメータ抽出部16において、合成音に対し分析を行い、合成音の特徴パラメータを抽出する(ステップS13)。合成音と目標音声の時間軸のずれを修正するために、目標話者音声から抽出したパラメータと合成音から抽出した特徴パラメータを用いて、合成音と目標音声の間で時間軸の対応付けを行う。対応付けの方法としては、DPマッチングによる時間軸伸縮等が考えられる(ステップS14)。全学習データ(目標話者音声と合成音の組)の時間軸の対応付けが行われた後、この時間軸対応付け済みの特徴パラメータを用いて、変換関数生成部17で合成音の声質から目標話者音声の声質へと変換する関数を作成する(ステップS15)。

30

【0025】

さらに、被変換入力信号のどの部分にどのような変換関数が適用されるかを決定するために、被変換入力信号(素片データ)と作成された変換関数との対応付けを行う(ステップS16)。被変換入力信号に対応付けられた変換関数を用いて、被変換入力信号を変換する(ステップS17)。変換されたデータを変換後出力信号として出力する(ステップS18)。

40

【0026】

本発明の第1の実施形態によれば、目標話者音声と同一あるいは類似の発声内容の合成音を音声合成で生成してから分析するため、目標話者音声の発声内容にかかわらず、高精度の対応付けを行ったうえで変換関数を生成することが可能となる。このため、変換目標話者の発声内容に依存しない自由な声質変換を実現することができる。したがって、目標話者音声の発声内容に合わせて、その都度変換元音声を入力する必要がなくなるので、使用者の負担を軽減し、処理を迅速化することができる。

【0027】

50

[ 第 2 の実施形態 ]

図 2 に、本発明の第 2 の実施形態に係る声質変換装置のブロック図を示す。図 2 において、図 1 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。第 2 の実施形態に係る声質変換装置は、第 1 の実施形態に係る声質変換装置に、さらに音声認識部 19 を備えている。音声認識部 19 は、目標話者音声入力部 11a で入力された目標話者音声を音声認識し、目標話者音声の発声内容を発音記号とともに発音記号列入力部 12a へ出力する。

【 0028 】

次に、本発明の第 2 の実施形態に係る声質変換装置の動作を図 2 および図 14 を参照して説明する。ただし、図 14 において、ステップ S21 および S23 ~ S28 は、それぞれ図 13 の S11 および S13 ~ S18 と同様の動作をするため、新たに追加したステップの動作のみを説明する。ステップ S21 で目標話者音声入力部 11a から入力された目標話者音声を音声認識部 19 で音声認識して (ステップ S221) 発音記号列を生成し (ステップ S222)、生成された発音記号列を用いて合成音を作成する (ステップ S223)。

10

【 0029 】

本発明の第 2 の実施形態によれば、目標話者音声の音声認識結果を合成音生成のための発音記号列として入力するため、ユーザがテキスト等で発音記号列を入力する必要が無く、処理の自動化を図ることが可能となる。本発明は、目標話者音声として多くの文を用いて対応付けの高精度化を図る際に有効である。

20

【 0030 】

[ 第 3 の実施形態 ]

本発明の第 3 の実施形態に係る声質変換装置は、第 1 の実施形態に係る声質変換装置の構成に加え、音声合成部で使用したセグメンテーション情報を記憶し、変換関数生成部に出力するセグメンテーション情報記憶部により構成される。このような構成を採用し、変換関数生成部でセグメンテーション情報を用いることにより、所望の声質へ高精度に変換できる変換関数の同定が可能である。なお、以下の説明において、セグメンテーション情報とは、時刻情報に対応づいた音素 (音素ラベル) を主とし、ピッチパターンあるいはパワー等の韻律情報を付加したものでよい。

【 0031 】

図 3 に、本発明の第 3 の実施形態に係る声質変換装置のブロック図を示す。図 3 において、目標話者音声入力部 11 と、発音記号列入力部 12 と、音声合成用データ記憶部 13 と、目標話者音声特徴パラメータ抽出部 15 と、合成音特徴パラメータ抽出部 16 と、声質変換部 18 は、図 1 と同等であり、その説明を省略する。本発明の第 3 の実施形態に係る声質変換装置は、さらに合成音セグメンテーション情報記憶部 20 を備えている。合成音セグメンテーション情報記憶部 20 は、音声合成部 14a で発音記号列に従って音声を合成する際に用いられるセグメンテーション情報を記憶して、変換関数生成部 17a へと出力する。

30

【 0032 】

次に、本発明の第 3 の実施形態に係る声質変換装置の動作を図 3 および図 15 を参照して説明する。ただし、図 15 において、図 13 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S12a では、音声合成部 14 で音声を合成する際に用いられたセグメンテーション情報 23a を合成音セグメンテーション情報記憶部 20 に記憶し、変換関数生成部 17a へ出力する。変換関数生成部 17a に入力されたセグメンテーション情報 23a は、時間軸の対応付け (ステップ S14a)、および変換関数の作成 (ステップ S15a) で用いられる。

40

【 0033 】

本発明の第 3 の実施形態によれば、合成音を生成する際に利用するセグメンテーション情報を変換関数生成時に活用するため、合成音特徴パラメータと目標話者音声特徴パラメ

50

ータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

【 0 0 3 4 】

[ 第 4 の実施形態 ]

図 4 に、本発明の第 4 の実施形態に係る声質変換装置のブロック図を示す。図 4 において、図 2 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 4 の実施形態に係る声質変換装置は、図 2 に対し、さらに合成音セグメンテーション情報記憶部 20 を備えている。なお、合成音セグメンテーション情報記憶部 20 は、図 3 で説明したものと同様である。

【 0 0 3 5 】

次に、本発明の第 4 の実施形態に係る声質変換装置の動作を図 4 および図 16 を参照して説明する。ただし、図 16 において、図 13 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S223a では、音声合成部 14a で音声を合成する際に用いられたセグメンテーション情報 23b を出力し、変換関数生成部 17a へ入力する。入力されたセグメンテーション情報 23b は、時間軸の対応付け (ステップ S24a)、および変換関数の作成 (ステップ S25b) で用いられる。

【 0 0 3 6 】

本発明の第 4 の実施形態によれば、目標話者音声を音声合成して発音記号列を自動生成し、これを用いて合成音を生成する際に利用するセグメンテーション情報を変換関数生成時に活用するため、処理の自動化を図りつつ、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

【 0 0 3 7 】

[ 第 5 の実施形態 ]

図 5 に、本発明の第 5 の実施形態に係る声質変換装置のブロック図を示す。図 5 において、図 1 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 5 の実施形態に係る声質変換装置は、さらにセグメンテーション情報入力部 22 を備えている。セグメンテーション情報入力部 22 は、目標話者音声入力部 11 で入力された目標話者音声を不図示の手段で分析してセグメンテーション情報 23c を抽出し、音声合成部 14b に出力する。

【 0 0 3 8 】

次に、本発明の第 5 の実施形態に係る声質変換装置の動作を図 5 および図 17 を参照して説明する。ただし、図 17 において、図 13 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S12b では、目標話者音声入力部 11 で入力された目標話者音声に基づいたセグメンテーション情報 23c を入力する。入力されたセグメンテーション情報 23c は、音声合成部 14b に出力され、発音記号列入力部 12 で入力された発音記号列の発声内容を持ち、発声の時間長が目標話者音声と同一である音声が発音合成部 14b において合成される。

【 0 0 3 9 】

本発明の第 5 の実施形態によれば、目標話者音声に基づいたセグメンテーション情報を入力して合成音を生成するため、生成された合成音の時間長情報は当該目標話者音声の時間長情報と等しくなり、変換関数生成部での時間軸の対応付け処理を不必要とする、あるいは簡素化することが可能となる。

【 0 0 4 0 】

[ 第 6 の実施形態 ]

図 6 に、本発明の第 6 の実施形態に係る声質変換装置のブロック図を示す。図 6 において、図 2 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 6 の実施形態に係る声質変換装置は、さらにセグメンテーション情報入力部 22 を備えている。セグメンテーション情報入力部 22 は、目標話者音声入力部 11a で入力された目標話者音声を不図示の手段で分析してセグメンテーション情報 23c を抽出し、音声合成部 14b に出力する。

## 【 0 0 4 1 】

次に、本発明の第 6 の実施形態に係る声質変換装置の動作を図 6 および図 1 8 を参照して説明する。ただし、図 1 8 において、図 1 3 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S 2 2 2 a では、目標話者音声入力部 1 1 a で入力された目標話者音声に基づいたセグメンテーション情報 2 3 c を入力する。入力されたセグメンテーション情報 2 3 c は音声合成部 1 4 b に出力され、発音記号列入力部 1 2 a で入力された発音記号列の発声内容を持ち、発声の時間長が目標話者音声と同一である音声合成される。

## 【 0 0 4 2 】

本発明の第 6 の実施形態によれば、目標話者音声に基づいたセグメンテーション情報を 10  
入力して合成音を生成するため、生成された合成音の時間長情報は当該目標話者音声の時間長情報と等しくなり、変換関数生成部での時間軸の対応付け処理を不必要とする、あるいは簡素化することが可能となる。さらに、本発明では、音声認識によって目標話者音声から発音記号列を生成しているため、処理の自動化を図りつつ、発音記号列と目標話者音声から抽出される時間長情報との整合性が取りやすくすることが可能となる。

## 【 0 0 4 3 】

## [ 第 7 の実施形態 ]

図 7 に、本発明の第 7 の実施形態に係る声質変換装置のブロック図を示す。図 7 において、図 2 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 7 の実施形態に係る声質変換装置は、さらに目標話者音声セグメンテーション情報記憶部 2 1 を備えている。目標話者音声セグメンテーション情報記憶部 2 1 は、目標話者音声入力部 1 1 a で入力された目標話者音声のセグメンテーション情報を記憶して、変換関数生成部 1 7 b に出力する。 20

## 【 0 0 4 4 】

次に、本発明の第 7 の実施形態に係る声質変換装置の動作を図 7 および図 1 9 を参照して説明する。ただし、図 1 9 において、図 1 4 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S 2 2 1 a では、目標話者音声信号を音声認識部 1 9 で音声認識した結果、抽出されるセグメンテーション情報 2 3 d を出力する。セグメンテーション情報 2 3 d は、変換関数生成部 1 7 b へ出力され、時間軸の対応付け (ステップ S 2 4 a)、および変換関数の作成 (ステップ S 2 5 a) で用いられる。 30

## 【 0 0 4 5 】

本発明の第 7 の実施形態によれば、目標話者音声を音声認識する際に同時に目標話者音声のセグメンテーション情報を抽出して変換関数生成部に入力するため、目標話者音声の特徴パラメータに加えてセグメンテーション情報を変換関数生成に用いることができ、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

## 【 0 0 4 6 】

## [ 第 8 の実施形態 ]

図 8 に、本発明の第 8 の実施形態に係る声質変換装置のブロック図を示す。図 8 において、図 2 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 8 の実施形態に係る声質変換装置は、さらに目標話者音声セグメンテーション情報記憶部 2 1 a を備えている。目標話者音声セグメンテーション情報記憶部 2 1 a は、目標話者音声入力部 1 1 a で入力された目標話者音声のセグメンテーション情報を記憶して、音声合成部 1 4 b に出力する。 40

## 【 0 0 4 7 】

次に、本発明の第 8 の実施形態に係る声質変換装置の動作を図 8 および図 2 0 を参照して説明する。ただし、図 2 0 において、図 1 4 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S 2 2 1 a では、目標話者を音声認識部 1 9 で音声認識した結果抽出されるセグメンテーシ 50

オン情報 2 3 e を出力する。ステップ S 2 2 3 b において、音声合成部 1 4 b は、セグメンテーション情報 2 3 e を入力し、発音記号列入力部 1 2 a で入力された発音記号列の発声内容を持ち、目標話者音声のセグメンテーション情報 2 3 e に基づいた音声を合成する。

【 0 0 4 8 】

本発明の第 8 の実施形態によれば、目標話者音声を音声認識する際に同時に同音声のセグメンテーション情報を抽出して音声合成部に入力するため、目標話者音声に基づいたセグメンテーション情報を用いて合成音を生成することができ、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

10

【 0 0 4 9 】

[ 第 9 の実施形態 ]

図 9 に、本発明の第 9 の実施形態に係る声質変換装置のブロック図を示す。図 9 において、図 2 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 9 の実施形態に係る声質変換装置は、さらに目標話者音声セグメンテーション情報記憶部 2 1 b を備えている。目標話者音声セグメンテーション情報記憶部 2 1 b は、目標話者音声入力部 1 1 a で入力された目標話者音声のセグメンテーション情報を記憶して、音声合成部 1 4 b および変換関数生成部 1 7 b に出力する。

【 0 0 5 0 】

次に、本発明の第 9 の実施形態に係る声質変換装置の動作を図 9 および図 2 1 を参照して説明する。ただし、図 2 1 において、図 1 4 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S 2 2 1 a では、目標話者を音声認識部 1 9 で音声認識した結果抽出されるセグメンテーション情報 2 3 f を出力する。セグメンテーション情報 2 3 f は、音声合成部 1 4 b へ出力され、発音記号列入力部 1 2 a で入力された発音記号列の発声内容を持ち、目標話者音声のセグメンテーション情報 2 3 f に基づいた音声を合成する。さらに、セグメンテーション情報 2 3 f は、変換関数生成部 1 7 b にも出力され、時間軸の対応付け（ステップ S 2 4 a ）、および変換関数の作成（ステップ S 2 5 a ）で用いられる。

20

【 0 0 5 1 】

本発明の第 9 の実施形態によれば、目標話者音声を音声認識する際に同時に同音声のセグメンテーション情報を抽出して音声合成部および変換関数生成部に入力するため、目標話者音声に基づいたセグメンテーション情報を用いて合成音を生成することができ、かつ、目標話者音声の特徴パラメータに加えてセグメンテーション情報を変換関数生成に用いることができる。したがって、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

30

【 0 0 5 2 】

[ 第 1 0 の実施形態 ]

図 1 0 に、本発明の第 1 0 の実施形態に係る声質変換装置のブロック図を示す。図 1 0 において、図 2 と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第 1 0 の実施形態に係る声質変換装置は、さらに合成音セグメンテーション情報記憶部 2 0 と目標話者音声セグメンテーション情報記憶部 2 1 を備えている。合成音セグメンテーション情報記憶部 2 0 は、図 3 で説明した合成音セグメンテーション情報記憶部 2 0 と同じである。また目標話者音声セグメンテーション情報記憶部 2 1 は、図 7 で説明した目標話者音声セグメンテーション情報記憶部 2 1 と同様のものである。

40

【 0 0 5 3 】

次に、本発明の第 1 0 の実施形態に係る声質変換装置の動作を図 1 0 および図 2 2 を参照して説明する。ただし、図 2 2 において、図 1 3 の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップ S 2 2 1 a では、目標話者を音声認識部 1 9 で音声認識した結果抽出されるセグメンテーション情報 2 3 d を出力する。セグメンテーション情報 2 3 d は、変換関数生成部 1 7

50

cに出力される。また、ステップS223aでは、音声合成部14aで音声を合成する際に用いられたセグメンテーション情報23gを記憶し、変換関数生成部17cへ出力される。変換関数生成部17cに入力されたセグメンテーション情報23fおよび23gは時間軸の対応付け(ステップS24b)、および変換関数の作成(ステップS25b)で用いられる。

【0054】

本発明の第10の実施形態によれば、変換関数生成部において合成音と目標話者音声両方のセグメンテーション情報を利用できるため、合成音および目標話者音声の特徴パラメータに加えてセグメンテーション情報を変換関数生成に用いることができ、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図

10

【0055】

[第11の実施形態]

図11に、本発明の第11の実施形態に係る声質変換装置のブロック図を示す。図11において、図2と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第11の実施形態に係る声質変換装置は、さらに合成音セグメンテーション情報記憶部20と目標話者音声セグメンテーション情報記憶部21aを備えている。合成音セグメンテーション情報記憶部20および目標話者音声セグメンテーション情報記憶部21aは、図3の合成音セグメンテーション情報記憶部20および図8の目標話者音声セグメンテーション情報記憶部21aと同様のものである。

20

【0056】

次に、本発明の第11の実施形態に係る声質変換装置の動作を図11および図23を参照して説明する。ただし、図23において、図20の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップS223cでは、音声合成部14cで音声を合成する際に用いられたセグメンテーション情報23gを記憶し、変換関数生成部17aへ入力する。入力されたセグメンテーション情報23gは時間軸の対応付け(ステップS24a)、および変換関数の作成(ステップS25a)で用いられる。

【0057】

本発明の第11の実施形態によれば、目標話者音声を音声認識する際に同時に同音声のセグメンテーション情報を抽出して音声合成部に入力するため、目標話者音声に基づいたセグメンテーション情報を用いて合成音を生成することができる。また、合成音を生成する際に利用するセグメンテーション情報を変換関数生成時に活用するため、合成音の特徴パラメータに加えてセグメンテーション情報を変換関数生成に用いることができる。したがって、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

30

【0058】

[第12の実施形態]

図12に、本発明の第12の実施形態に係る声質変換装置のブロック図を示す。図12において、図2と同一の符号は、同一部あるいは相当部を示し、その説明を省略する。本発明の第12の実施形態に係る声質変換装置は、さらに合成音セグメンテーション情報記憶部20と目標話者音声セグメンテーション情報記憶部21bを備えている。合成音セグメンテーション情報記憶部20および目標話者音声セグメンテーション情報記憶部21bは、それぞれ図3の合成音セグメンテーション情報記憶部20および図9の目標話者音声セグメンテーション情報記憶部21bと同様のものである。

40

【0059】

次に、本発明の第12の実施形態に係る声質変換装置の動作を図12および図24を参照して説明する。ただし、図24において、図21の同一の符号は、同一あるいは同等の処理を行い、その説明を省略し、新たに追加したステップの動作のみを説明する。ステップS223cでは、音声合成部14cで音声を合成する際に用いられたセグメンテーシ

50

ン情報 2 3 g を記憶し、変換関数生成部 1 7 c へ入力する。入力されたセグメンテーション情報 2 3 g は時間軸の対応付け（ステップ S 2 4 b）、および変換関数の作成（ステップ S 2 5 b）で用いられる。

#### 【 0 0 6 0 】

本発明の第 1 2 の実施形態によれば、目標話者音声を音声認識する際に同時に同音声のセグメンテーション情報を抽出して音声合成部および変換関数生成部に入力するため、目標話者音声に基づいたセグメンテーション情報を用いて合成音を生成することができる。また、目標話者音声の特徴パラメータに加えてセグメンテーション情報を変換関数生成に用いることができる。また、変換関数生成部において合成音と目標話者音声両方のセグメンテーション情報を利用できるため、合成音および目標話者音声の特徴パラメータに加えてセグメンテーション情報を変換関数生成に用いることができる。したがって、合成音特徴パラメータと目標話者音声特徴パラメータとを対応付ける際の処理の簡素化、高精度化を図ることが可能となる。

#### 【 実施例 1 】

#### 【 0 0 6 1 】

図 2 5 は、本発明の第 1 の実施例に係る声質変換装置のブロック図である。声質変換装置は、目標話者音声入力部 1 1、発音記号列入力部 1 2、音声合成用データ記憶部 1 3、音声合成部 1 4、L P C 分析部 2 5、2 6、3 0、D P マッチング部 2 7、スペクトル形状変換部 2 8、変換関数導出部 2 9、対応フレーム探索部 3 1、スペクトル変換部 3 2 を備える。なお、図 2 5 において、図 1 と同一の符号は、同一物あるいは相当物を示し、特に記載無き場合、その説明を省略する。

#### 【 0 0 6 2 】

本実施例では、あらかじめ少なくとも 1 人以上の話者から音声合成用のデータベースが作成されており、音声合成用データ記憶部 1 3 に保存されているものとする。データベースは、音声素片、継続時間長、ピッチパターン等のデータを含んでいる。ただし、必ずしもこれらの全てのデータをデータベース内に記憶しておく必要はなく、これらのうち 1 つないし 2 つのデータだけでもよい。

#### 【 0 0 6 3 】

今、目標話者音声の発声内容 C が明らかになっており、予めこの発声内容と同一あるいは類似の発声内容を持つ合成音を生成できるような発音記号列が発音記号列入力部 1 2 から入力され用意されているとする。この発音記号列に従って、音声合成部 1 4 は、音声合成用データ内の素片の声質を持った発声の合成音 B を合成する。目標話者音声入力部 1 1 から入力される目標話者音声としては、操作開始時点で録音装置に発声する場合、予め保存してあった音声データを用いる場合等が考えられる。また、目標話者音声も複数あっても構わない。

#### 【 0 0 6 4 】

次に、L P C 分析部 2 5 は、目標話者音声入力部 1 1 から入力された、話者 A の発声内容 C の音声（以降では、目標 A とする）を例えば 2 0 m s e c といった長さの分析フレームごとに L P C 分析し、目標 A の L P C 係数 3 5 a を抽出する。また、L P C 分析部 2 6 は、音声合成部 1 4 で作成した合成音 B の発声内容 C の音声（以降では、合成音 B とする）を同様に L P C 分析し、合成音 B の L P C 係数 3 5 b を抽出する。本実施例では分析方法（特徴パラメータ抽出）として、L P C 分析を説明するが、他に、L S P（line spectrum pair）分析、P A R C O R（partial auto-correlation）分析等のスペクトル分析や零交叉計数法、自己相関法等のピッチ抽出法、およびこれら複数の組み合わせによる分析等が考えられる。また、抽出する特徴パラメータとして L P C 係数のほかに、自己相関係数、ケプストラム係数等や、ピッチ周波数等の韻律パラメータを用いることもできる。

#### 【 0 0 6 5 】

次に、声質変換の主要部となる、D P マッチング部 2 7、スペクトル形状変換部 2 8、変換関数導出部 2 9、対応フレーム探索部 3 1、スペクトル変換部 3 2 における処理の流れを説明する。図 2 6 は、声質変換の主要部における音声データの処理の流れを説明する

図である。LPC分析された目標話者音声（目標A）と合成音Bとは、DPマッチング部27によってDPマッチングが行われて対応関係が求められる。求められた対応関係を用いて、合成音Bから目標Aへの変換関数が、スペクトル形状変換部28、および変換関数導出部29により生成される。

【0066】

変換対象となる音声素片は、対応フレーム探索部31によって合成音素片との対応付けがなされる。この場合、対応付けは、例えば変換対象となる音声素片と合成音素片とのそれぞれのLPC係数での特徴空間で距離が小さくなるものを選択することでなされる。変換対象となる音声素片は、対応付けがされた合成音素片に対応する先に求めた変換関数によってスペクトル変換がなされ、目標の声質を持つ音声素片が生成されることとなる。以下に、各部についてより詳しく説明する。

10

【0067】

DP (dynamic programming) マッチング部27は、目標AのLPC係数35aと合成音BのLPC係数35bとを用いて、目標Aと合成音Bの時間軸を合わせるために、DP (dynamic programming) マッチングによる時間軸伸縮を行う。これにより、目標Aと合成音Bとの分析フレームごとの対応が作成される。この時、DPマッチングで用いる特徴量空間内距離を表す尺度としては、差分ベクトルの二乗和、WLR (weighted likelihood ratio) 尺度、最尤スペクトル距離、LPCケプストラム距離等を用いることができる。

【0068】

図27は、分析フレームごとのDPマッチングを模式的に表した図である。目標話者音声（目標A）のフレームをA1、A2、A3、A4、A5、・・・Am、・・・とし、合成音BのフレームをB1、B2、B3、B4、B5、・・・Bn、・・・とする。目標話者音声（目標A）と合成音Bとをフレーム毎にLPC分析し、その上で、例えば、フレームA1とフレームB1との特徴量空間内距離、フレームA1とフレームB2との特徴量空間内距離、フレームA2とフレームB1との特徴量空間内距離、の中で最も距離の短い対応関係を選択する。ここでは、フレームA2とフレームB1との特徴量空間内距離が選択されたとすると、次には、例えば、フレームA3とフレームB1との特徴量空間内距離、フレームA2とフレームB2との特徴量空間内距離、フレームA3とフレームB2との特徴量空間内距離、の中で最も距離の短い対応を選択する。このようにして、特徴量空間内距離の最小の距離の対応を順次求めていき、目標話者音声（目標A）と合成音BとのDPマッチングが行われる。

20

30

【0069】

次に、スペクトル形状変換部28は、分析フレームごとに対応付けられた目標AのLPC係数35aと合成音BのLPC係数35bとを用いて、合成音Bの声質が目標Aの声質へと変換されるような変換関数を同定する。合成音Bから目標Aへと声質変換されるような変換を実現するには、合成音Bの分析フレームごとの周波数特性が目標Aの周波数特性とできるだけ等しくなるように、合成音Bの周波数領域のデータを変換関数で変換すればよい。これを分析フレームごとに考えると、合成音Bの1つの分析フレームBn (n番目の分析フレームBnとする)の周波数特性が、この分析フレームに時間軸上で対応付けられた目標Aの分析フレームAm (m番目の分析フレームAmとする)の周波数特性に変換されるような変換関数を同定すればよいということになる。そこで、分析フレームをフーリエ変換して波形の周波数領域の形状(スペクトル包絡)を求め、周波数軸上でDPマッチング等による伸縮を行い、合成音Bから目標Aへの変換関数を求める。すなわち、合成音Bの分析フレームBnのスペクトル形状を、目標Aの分析フレームAmのスペクトル形状に合致させるような変換を行う。

40

【0070】

さらに、スペクトル形状変換について説明する。図28は、周波数軸上でのDPマッチングを模式的に表した図である。目標話者音声（目標A）と合成音Bのスペクトル包絡をそれぞれSA、SBとする。周波数軸上でスペクトル包絡SAとスペクトル包絡SBとの

50

対応関係をDPマッチングによって求める。対応関係を表すパス $P_0$ 、 $\dots$ 、 $P_i$ 、 $\dots$ 、 $P_n$ が変換関数に相当し、合成音Bのスペクトルをパス $P_0$ 、 $\dots$ 、 $P_i$ 、 $\dots$ 、 $P_n$ によって非線型に写像することで目標話者音声(目標A)のスペクトルが得られることとなる。なお、予めスペクトル包絡 $S_A$ 、 $S_B$ に対し、高域強調あるいは低域強調等の前処理を行い、前処理がなされたスペクトル包絡に対してDPマッチングを行うようにしてもよい。

【0071】

以上で説明したように、変換関数同定においては、LPC分析部26で合成音Bのn番目の分析フレーム $B_n$ をLPC分析した結果から、スペクトル包絡 $S_{B_n}$ を導出する。同様に、 $B_n$ に時間軸上で対応付けられた目標Aの分析フレーム $A_m$ のスペクトル包絡 $S_{A_m}$ を同定する。その上で、 $S_{B_n}$ が $S_{A_m}$ に変換されるような変換関数を作成する。このようにして、合成音Bのn番目の分析フレーム $B_n$ から、時間軸上で対応付けられた目標Aの分析フレーム $A_m$ への変換関数が求まるので、これを当該フレーム全てに適用し、合成音A全体に対する変換関数を求めることができる。

10

【0072】

なお、以上の説明では、目標Aと合成音Bといった1文だけを学習データとして用いて変換関数をフレームごとに求める例を示したが、学習データをさらに増やした場合、目標Aと合成音BのLPC係数を特徴量空間内で、例えば音素ごとにクラスタリングし、クラスタの代表点ごとに変換関数を同定することも可能である。また、GMM(Gaussian Mixture model)等の確率密度分布で特徴量空間を表現して、密度分布状態で対応付ける方法も考えられる。

20

【0073】

図29は、特徴量空間内でクラスタリングされた音素間の変換関数を模式的に示す図である。合成音Bの特徴パラメータ空間内と目標Aの特徴パラメータ空間内において、例えば音素(a、i、u、e、o、s)がそれぞれクラスタに分類され、合成音Bの特徴量空間内のクラスタ中の代表点を変換関数によって目標Aの特徴量空間内のクラスタ中の代表点に変換される。

【0074】

次に、変換対象である音声合成用素片データと変換関数との対応付けを行う。そのために、LPC分析部30は、あらかじめ音声合成用の素片信号をLPC分析して、素片データの分析フレームごとのLPC係数を求めておく。

30

【0075】

対応フレーム探索部31は、LPC分析部30が出力する素片データの分析フレーム毎のLPC係数に対して合成音Bのフレーム毎のLPC係数35b中の特徴量空間内距離が最も近いものを対応フレームとして求める。

【0076】

変換関数導出部29は、スペクトル形状変換部28によって求めた合成音Bのフレームごとの変換関数を素片データの変換関数とする。この時、学習データを増やした場合は、図29で説明した変換関数の同定時と同様に、フレームごとではなくクラスタごとに変換関数を設定するといったことも可能である。

【0077】

スペクトル変換部32は、変換関数導出部29において求めた素片に対する変換関数の中から、先の対応フレームに対する変換関数を選択し、選択された変換関数を用いて素片信号(素片データ)のスペクトルを変換する。これにより、目標Aの声質を持つ変換後素片信号(合成音を作成できる素片データ)を得ることができる。

40

【0078】

さらに、この素片データを、音声合成用データ記憶部13に保存されている変換前の音声合成用データベース内の素片データと差し替えることで、任意のテキストに対して目標Aの声質を持つ合成音を作成できるデータベースが完成する。

【0079】

なお、第1の実施例では、合成音を作成するための発音記号列が予め分かっているもの

50

としたが、目標話者音声の発声内容を漢字かな混じり文等のテキストデータで表して形態素解析等を用いて分析し、その結果によって発音記号列を生成することも可能である。

【 0 0 8 0 】

また、被変換対象である入力信号は、上記の素片データ以外にも、上記以外の素片データを使うこともできる。例えば、音声合成用データ記憶部 1 3 に保存されている素片データが旅行用の会話のデータであり、入力信号となる素片データが一般の会話用のデータである場合などがある。

【 0 0 8 1 】

さらに、上記の素片データで作成された合成音、上記以外の素片データで作成された合成音、ユーザの発声による音声等を被変換対象である入力信号とすることもできる。

10

【 0 0 8 2 】

さらに、第 1 の実施例の変形として、変換関数生成部内の合成音と目標話者音声の時間軸対応付け処理を行わないで変換関数の同定を行うという方法も考えられる。この方法を実現できる理由を以下に示す。本発明では、合成音と目標話者音声の発声内容を同一あるいは類似のものにして処理を行っているため、合成音および目標話者音声の音素情報は同一の部分が多く、特徴量空間内の特徴パラメータの分布が互いに類似している。このため、例えば、音素ごとにクラスタリングして変換関数を求める場合、合成音と目標話者音声の時間軸対応を取らなくても、変換関数を生成することができる。この方法を用いれば、時間軸対応付け処理を行わないので、処理速度の大幅な向上を図ることが可能である。

【 実施例 2 】

20

【 0 0 8 3 】

図 3 0 は、本発明の第 2 の実施例に係る声質変換装置のブロック図である。声質変換装置は、目標話者音声入力部 1 1、発音記号列入力部 1 2、音声合成用データ記憶部 1 3、音声合成部 1 4 a、音声認識部 1 9、合成音セグメンテーション情報記憶部 2 0、目標話者音声セグメンテーション情報記憶部 2 1、L P C 分析部 2 5、2 6、3 0、D P マッチング部 2 7 a、スペクトル形状変換部 2 8、変換関数導出部 2 9、対応フレーム探索部 3 1、スペクトル変換部 3 2 を備える。なお、図 3 0 において、図 2 5 と同一の符号は、同一物あるいは相当物を示し、特に記載無き場合、その説明を省略する。

【 0 0 8 4 】

音声認識部 1 9 は、入力された目標話者音声信号に対し音声認識を行うことによって、目標話者音声信号と同一あるいは類似の発声内容を持つ合成音を生成するための発音記号列を生成し、発音記号列入力部 1 2 に出力する。今、発声内容 C を持つ音声为目标話者音声であるとすると、これを音声認識することによって、発声内容 C を持つ発音記号列を自動的に生成することができる。このため、発音記号列が予め明らかでない場合でも、目標話者音声と同一あるいは類似の発声内容を持つ合成音を生成することができ、目標話者音声のみを入力してしまえば、第 1 の実施例と同様の処理を自動的に行うことができる。この時の音声認識の方法としては、例えば非特許文献 2 にあるような H M M ( hidden Markov model ) による音声認識法等がある。

30

【 0 0 8 5 】

合成音セグメンテーション情報記憶部 2 0 は、音声合成部 1 4 a において音声を合成する際に音声合成用データ内から算出されて用いられるセグメンテーション情報を、変換関数の生成時に活用するために D P マッチング部 2 7 a に出力する。今、発声内容 C を持つ音声为目标話者音声であるとすると、これを音声認識部 1 9 において音声認識することによって、発声内容 C を持つ発音記号列を自動的に生成し発音記号列入力部 1 2 に出力される。発音記号列が音声合成部 1 4 a に入力されたとすると、音声合成部 1 4 a は、この発音記号列にしたがって音声合成用データ記憶部 1 3 から発声内容 C に対応したセグメンテーション情報を算出する。合成音セグメンテーション情報記憶部 2 0 は、このセグメンテーション情報を記憶しておき、セグメンテーション情報 2 3 i として D P マッチング部 2 7 a に出力し、変換関数生成の際に利用する。例えば、目標話者音声を分析して抽出された特徴パラメータの時間変化に対して、セグメンテーション情報と合成音を分析して抽出

40

50

された特徴パラメータの時間変化を対応付けることで、目標話者音声と合成音の時間軸対応付けを簡素化、高精度化することができる。セグメンテーション情報としては、図31に示すような、フレーム番号と音素が対応付けられた表（例えば、フレーム番号1、2、・・・10、11、・・・50、51、・・・に対しそれぞれ音素「i」、「i」、・・・「i」、「e」、・・・「u」、「s」、・・・が対応している）を用いる。その他の利用法としては、目標話者音声と合成音の組が複数文であった場合、音素セグメンテーション情報を用いてクラスタリングするといった方法も考えられる。

【0086】

目標話者音声セグメンテーション情報記憶部21は、音声認識部19が目標話者音声信号を音声認識する際に、出力される目標話者音声のセグメンテーション情報を入力する。今、発声内容Cを持つ目標話者音声を目話者音声入力部11に入力したとする。音声認識部19は、この目標話者音声を音声認識して、発声内容Cを記述する発音記号列を生成して発音記号列入力部12に出力する一方で、音声認識の結果であるセグメンテーション情報を目標話者音声セグメンテーション情報記憶部21に出力し、目標話者音声セグメンテーション情報記憶部21は、セグメンテーション情報を記憶しておく。このセグメンテーション情報23hを、DPマッチング部27aに出力して合成音と目標話者音声との対応付け処理に利用する。利用法としては、合成音を分析して抽出された特徴パラメータの時間変化に対して、セグメンテーション情報と目標話者音声を分析して抽出された特徴パラメータの時間変化を対応付けることで、目標話者音声と合成音の時間軸対応付けを簡素化、高精度化することができる。なお、セグメンテーション情報としては、図31で示した表と同様のものを用いる。

【0087】

DPマッチング部27aは、実施例1で説明したDPマッチング部27と同様に、目標AのLPC係数35aと合成音BのLPC係数35bとを用いて、目標Aと合成音Bの時間軸を合わせるために、DPマッチングによる時間軸伸縮を行う。さらに、DPマッチング部27aは、合成音セグメンテーション情報記憶部20が出力するセグメンテーション情報23hと、目標話者音声セグメンテーション情報記憶部21が出力するセグメンテーション情報23iとを用いてDPマッチングによる時間軸対応付けの簡素化、高精度化を図っている。

【0088】

次にセグメンテーション情報を用いたDPマッチングによる時間軸対応付けの一例について説明する。図32は、セグメンテーション情報を用いたDPマッチングによる時間軸対応付けの第1の例を示す図である。縦軸方向に分析された目標A（音素「a」「s」「u」）、横軸方向に分析された合成音B（音素「a」「s」「u」）が配置されている。セグメンテーション情報23hとして、目標Aの音素「a」と「s」の間に音素境界P1が、目標Aの音素「s」と「u」の間に音素境界Q1が音声合成部14aにおいて設定されている。また、セグメンテーション情報23iとして、合成音Bの音素「a」と「s」の間に音素境界P2が、目標Aの音素「s」と「u」の間に音素境界Q2が音声認識部19において設定されている。この時、音素境界P1と音素境界P2との交点を拘束点Pとし、音素境界Q1と音素境界Q2との交点を拘束点Qとする。DPマッチングを音素「a」「s」「u」に対して順次行い、DPパスを求めていく。この時、DPパスが拘束点P、拘束点Qを通るように制約を課してDPパスを求めようとする。すなわち、拘束点を通るようにDPパスを決定することにより、制約条件付のDPマッチングを行い、音素境界同士が対応付けられるようにする。なお、拘束点は、必ずしも通らなくてもよく、拘束点の近傍を通るような緩やかな制約を課してもよい。

【0089】

セグメンテーション情報は、ある発声内容に対して、各音素の開始時刻と終了時刻、および音素のラベル等が記述された情報である。セグメンテーション情報によって、音素のラベルとその音素の開始時刻、終了時刻が示されるために、目標Aあるいは合成音Bにおける音素の境界が明確に示されることとなる。したがって、目標Aと合成音Bとの各フレ

10

20

30

40

50

ーム間の時間軸上の対応付け（DPパス）を求める際に、セグメンテーション情報を用いて制約を付けることで、DPマッチングによる対応付けが音素境界付近であいまいになった場合であっても、精度の高い対応付けを実現することができることとなる。

【0090】

次にセグメンテーション情報を用いたDPマッチングによる時間軸対応付けの他の例について説明する。図33は、セグメンテーション情報を用いたDPマッチングによる時間軸対応付けの第2の例を示す図である。目標Aと合成音Bは、図32で説明したと同様に配置されている。ただし、図33において、図32と異なる点は、セグメンテーション情報23hが入力されていない。すなわち、目標話者音声セグメンテーション情報記憶部21が存在せず、目標Aに音素境界がないことである。この場合、通常のDPマッチングによってDPパスを求めた上で、DPパスが合成音Bの音素境界P1を通る点に対応する目標Aの推定音素境界P3を求める。また、DPパスが合成音Bの音素境界Q1を通る点に対応する目標Aの推定音素境界Q3を求める。すなわち、制約条件なしでDPマッチングを行った後に、セグメンテーション情報がある音声の音素境界と対応付いた箇所を、セグメンテーション情報がない音声の音素境界として推定することができる。通常、DPマッチングを行っただけでは、音素境界を判定することができないため、音素境界が推定できるこの方法は、変換関数を選択する際などにより有効な手段となる。

10

【0091】

以上の説明では、セグメンテーション情報をDPマッチングに適用する例について説明したが、他にセグメンテーション情報を、スペクトル変換部32において被変換入力信号（素片信号）を変換する際に適用することも可能である。すなわち、被変換入力信号の各フレームがどの変換関数で変換されるかを判定する際に、そのフレームのセグメンテーション情報に付するラベル情報を用いれば、どの集合（例えば、音素毎のクラスタ等）に属するかを容易に判別することができる。この様子を図34に示す。図34において、ラベル情報（「i」「i」「e」「u」）は、特徴パラメータ空間中の各クラスタに対応付けがなされ、ラベル情報から直接的に特徴量空間内のクラスタ中の代表点に変換することができる。

20

【実施例3】

【0092】

次に声質変換方法を繰り返し行い精度を高めて行く例について説明する。図35は、本発明の第3の実施例に係る声質変換方法を表すフローチャート図である。図35において、ステップS31～S39は、それぞれ図14のステップS21、S221、S222、S223、S23、S24、S25、S26、S27と同等の処理を行うステップであり、その説明を省略する。ステップS40において、変換後の信号を音声合成用のデータとして音声合成用データ記憶部13に登録する。

30

【0093】

ステップS40において、繰り返しが所定の収束条件に達したか否かを判定する。達していなければ、ステップS34に戻り、変換後の音声合成用データを用いて合成音を生成する。達していればステップS42で一連の処理が終了する。収束条件としては、例えば、目標話者音声と合成音とのLPC係数空間内距離が一定値以下になった場合、もしくはこの一定値以下の状態が一定回数継続した場合がある。また、スペクトルやパワー等の前回との差分値の合計値が一定値以下になった場合、もしくはこの一定値以下の状態が一定回数継続した場合を収束条件としてもよい。さらに繰り返し回数が一定回数繰り返した場合等を収束条件としてもよい。

40

【0094】

第3の実施例では、第1の実施例のように、素片を被変換入力信号として音声合成用データの声質を変換する場合には、変換後の音声合成用データを用いて再度目標話者音声と同一あるいは類似の発声内容の合成音を生成し、実施例1と同様の方法で変換関数の生成を行い、素片を変換するという処理を複数回繰り返す。さらに、第2の実施例を組み合わせることで、最初に目標話者音声を入力してしまえば、繰り返しの処理も全て自動

50

で行うことができ、変換精度（声質の類似度）を高めていくことが可能となる。なお、図 35 では、セグメンテーション情報を用いない例を示したが、勿論セグメンテーション情報を用いてもよい。

【0095】

なお、以上説明した実施例 1～3 において、音声合成用データ記憶部 13 に複数人数分の音声データからなる音声合成用データを記憶しておき、入力された目標話者の声質によって使用する音声合成用データを選択することができるようにすることも可能である。例えば、目標話者音声が男声であった場合は男声の素片データを用い、女声であった場合は女声の素片データを用いるといった方法が考えられる。この方法を用いれば、極端な変換を避けて変換処理による音質の劣化を少なくすることが可能となる。

【実施例 4】

【0096】

次に、本発明の声質変換装置、声質変換方法あるいは声質変換プログラムによって生成された合成音データを用いて合成音を生成する装置について説明する。図 36 は、本発明の第 4 の実施例に係る合成音生成装置を表すブロック図である。図 36 において、記号列入力部 41、音声合成出力部 42、データ記憶部 43 は、それぞれ図 1 の発音記号列入力部 12、音声合成部 14、音声合成用データ記憶部 13 と同等の機能を有するものである。ただし、音声合成出力部 42 は、生成された合成音を出力する機能を持つ。また、データ記憶部 43 は、先に説明した声質変換装置、声質変換方法あるいは声質変換プログラムによって生成される変換後出力信号を合成音データとして蓄えるものである。音声合成出力部 42 は、記号列入力部 41 に入力される記号列に基づいてデータ記憶部 43 から読み出した合成音データを用いて合成音を生成して出力する。出力される合成音は、変換後出力信号、すなわち目標話者の声質に変換されたデータに基づいて生成されるので、目標話者の声質を備えた合成音となる。

【実施例 5】

【0097】

次に、本発明の声質変換装置、声質変換方法あるいは声質変換プログラムによって生成された合成音データを用いて合成音を生成する他の装置について説明する。図 37 は、本発明の第 5 の実施例に係る合成音生成装置を表すブロック図である。図 37 において、音声入力部 51 は、ユーザの発する音声信号を入力する。変換データ記憶部 53 は、先に説明した声質変換装置、声質変換方法あるいは声質変換プログラムによって生成される、分析フレーム毎の LPC 係数と対応する変換関数とを記憶している。音声変換出力部 52 は、音声入力部 51 から出力される音声信号をフレーム毎に、例えば LPC 分析し、分析されたフレームにパラメータ空間距離の最も近いフレームを変換データ記憶部 53 内において探索して対応する変換関数を読み出し、この変換関数によって音声信号のスペクトルを変換して出力する。出力される合成音は、目標話者の声質に変換されたデータに基づいて生成されるので、目標話者の声質を備えた合成音となる。

【産業上の利用可能性】

【0098】

本発明によれば、電話機やトランシーバー等の通信機器で、自由に声質を変えるといた用途に適用できる。また、パーソナルコンピュータや携帯電話等で、電子メールやチャットのテキストを読み上げる際の合成音の声質をユーザの望む声質にするといった用途にも適用できる。さらに、アニメーションの音声の録音や外国映画の吹き替え等をテキスト音声合成で行う場合に、登場人物に合ったキャラクタ音声の声質を生成するといった用途にも適用できる。

【図面の簡単な説明】

【0099】

【図 1】本発明の第 1 の実施形態に係る声質変換装置の構成を示すブロック図である。

【図 2】本発明の第 2 の実施形態に係る声質変換装置の構成を示すブロック図である。

【図 3】本発明の第 3 の実施形態に係る声質変換装置の構成を示すブロック図である。

10

20

30

40

50

- 【図 4】本発明の第 4 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 5】本発明の第 5 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 6】本発明の第 6 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 7】本発明の第 7 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 8】本発明の第 8 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 9】本発明の第 9 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 10】本発明の第 10 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 11】本発明の第 11 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 12】本発明の第 12 の実施形態に係る声質変換装置の構成を示すブロック図である。
- 【図 13】本発明の第 1 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 14】本発明の第 2 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 15】本発明の第 3 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 16】本発明の第 4 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 17】本発明の第 5 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 18】本発明の第 6 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 19】本発明の第 7 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 20】本発明の第 8 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 21】本発明の第 9 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 22】本発明の第 10 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 23】本発明の第 11 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 24】本発明の第 12 の実施形態に係る声質変換装置の動作を示すフローチャートである。
- 【図 25】本発明の第 1 の実施例に係る声質変換装置の構成を示すブロック図である。
- 【図 26】声質変換の主要部における音声データの処理の流れを説明する図である。
- 【図 27】分析フレームごとの DP マッチングを模式的に表した図である。
- 【図 28】周波数軸上での DP マッチングを模式的に表した図である。
- 【図 29】特徴量空間内でクラスタリングされた音素間の変換関数を模式的に示す図である。
- 【図 30】本発明の第 2 の実施例に係る声質変換装置のブロック図である。
- 【図 31】フレーム番号と音素との対応付けを示す図である。
- 【図 32】セグメンテーション情報を用いた DP マッチングによる時間軸対応付けの第 1 の例を示す図である。
- 【図 33】セグメンテーション情報を用いた DP マッチングによる時間軸対応付けの第 2 の例を示す図である。
- 【図 34】ラベル情報とクラスタとの対応付けを示す図である。
- 【図 35】本発明の第 3 の実施例に係る声質変換方法を表すフローチャート図である。

10

20

30

40

50

【図36】本発明の第4の実施例に係る合成音生成装置を表すブロック図である。

【図37】本発明の第5の実施例に係る合成音生成装置を表すブロック図である。

【図38】第1の従来例の声質変換方法を表すブロック図である。

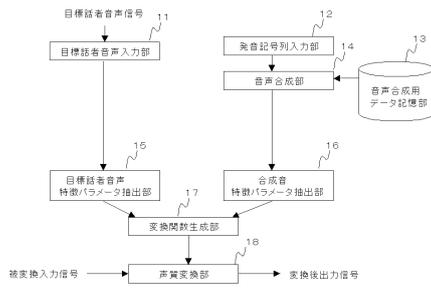
【図39】第2の従来例の声質変換方法を表すブロック図である。

【符号の説明】

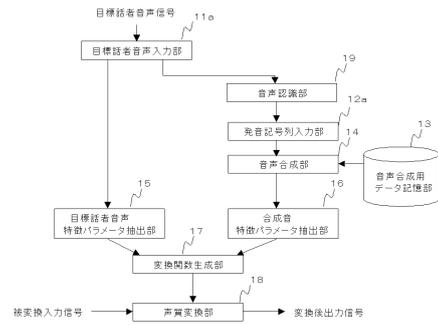
【0100】

11、11a	目標話者音声入力部	
12、12a	発音記号列入力部	
13	音声合成用データ記憶部	
14、14a、14b、14c	音声合成部	10
15	目標話者音声特徴パラメータ抽出部	
16	合成音特徴パラメータ抽出部	
17、17a、17b、17c	変換関数生成部	
18	声質変換部	
19	音声認識部	
20	合成音セグメンテーション情報記憶部	
21、21a、21b	目標話者音声セグメンテーション情報記憶部	
22	セグメンテーション情報入力部	
23a、23b、23c、23d、23e、23f、23g、23h、23i	セグメンテーション情報	20
25、26、30	LPC分析部	
27、27a	DPマッチング部	
28	スペクトル形状変換部	
29	変換関数導出部	
31	対応フレーム探索部	
32	スペクトル変換部	
35a、35b	LPC係数	
41	記号列入力部	
42	音声合成出力部	
43	データ記憶部	30
51	音声入力部	
52	音声変換出力部	
53	変換データ記憶部	

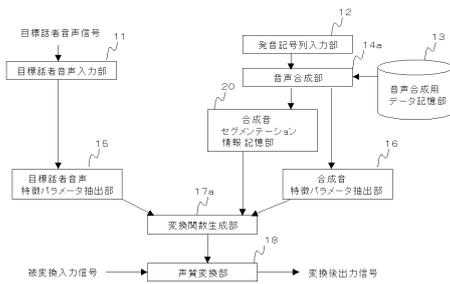
【図1】



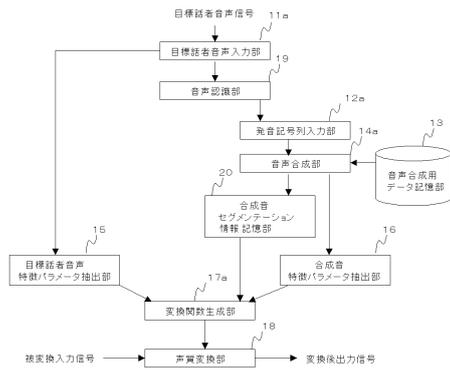
【図2】



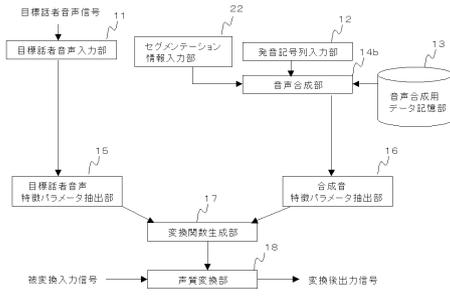
【図3】



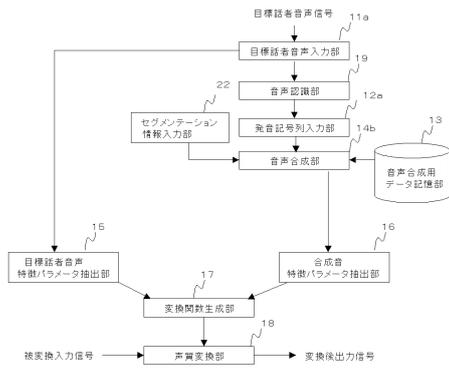
【図4】



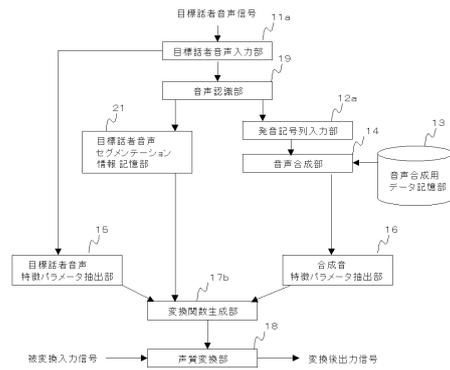
【図5】



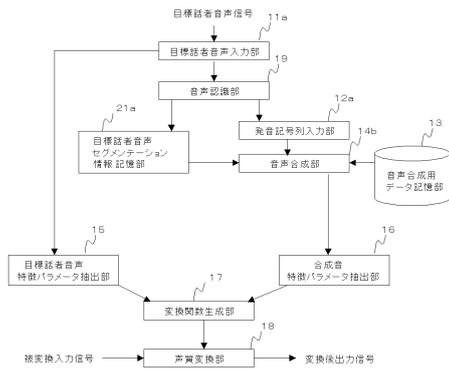
【図6】



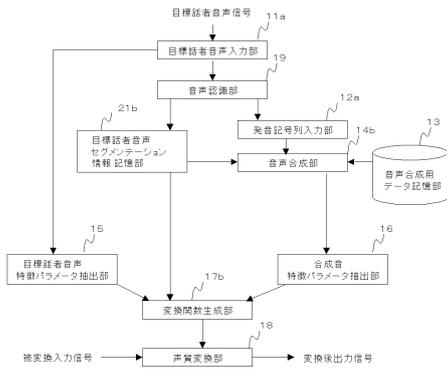
【図7】



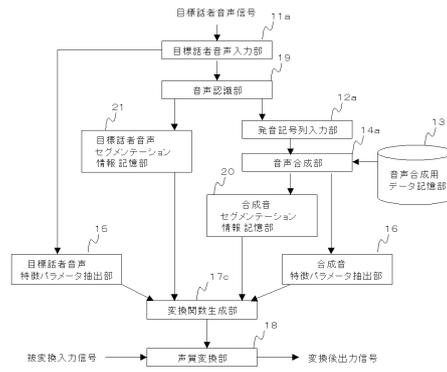
【図8】



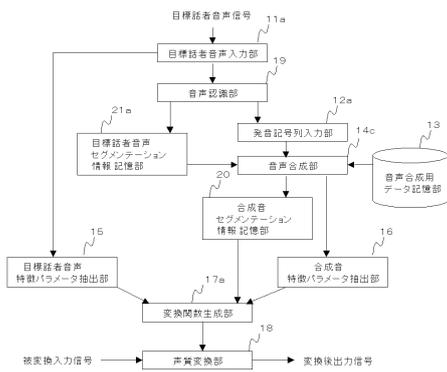
【図 9】



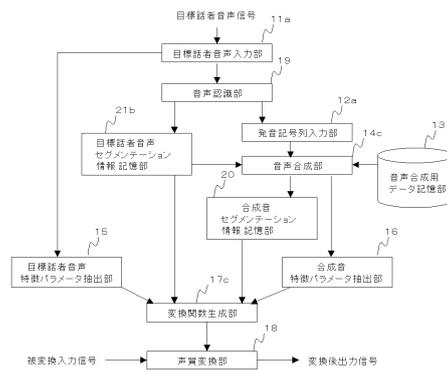
【図 10】



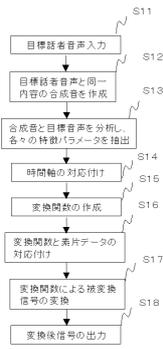
【図 11】



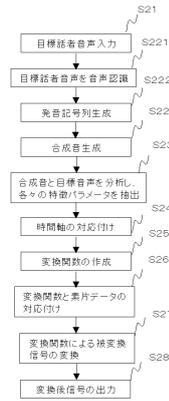
【図 12】



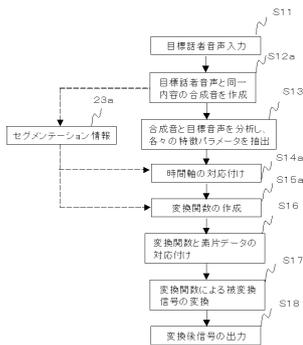
【図 13】



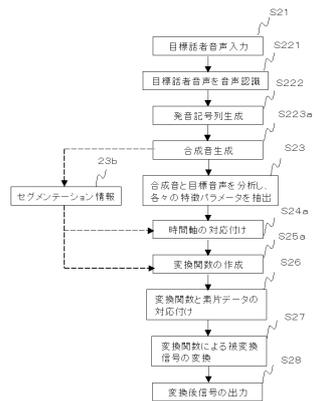
【図 14】



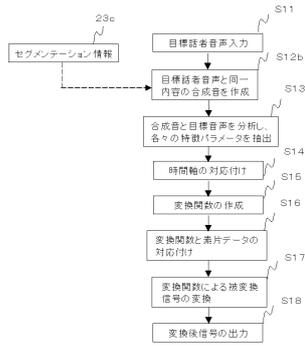
【図 15】



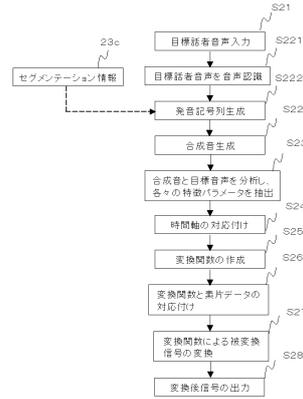
【図 16】



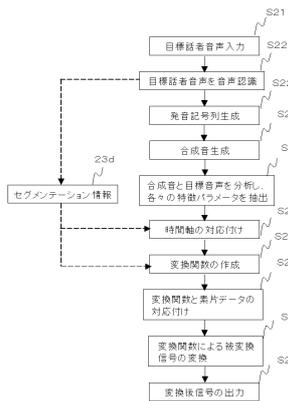
【図 17】



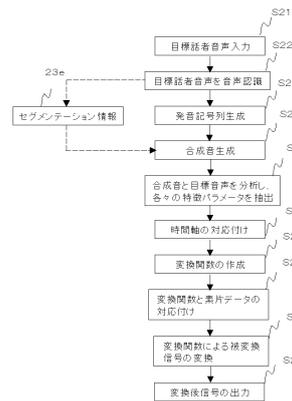
【図 18】



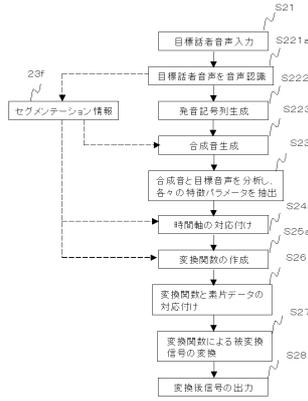
【図 19】



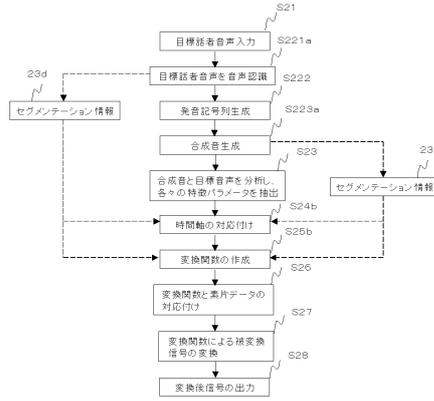
【図 20】



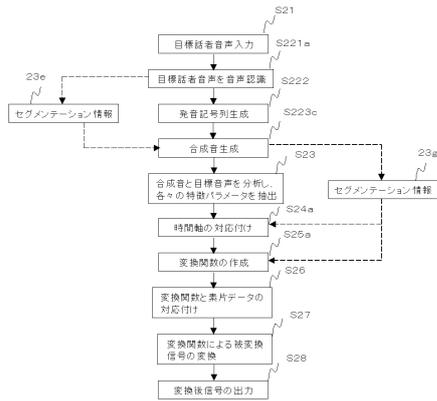
【図 2 1】



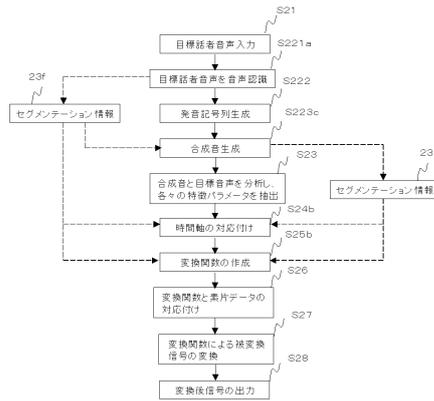
【図 2 2】



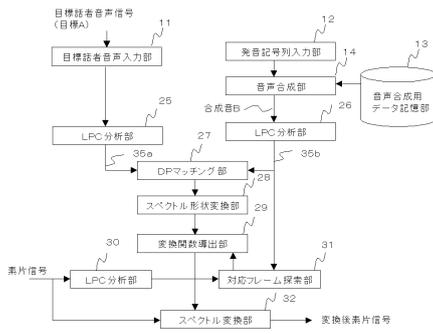
【図 2 3】



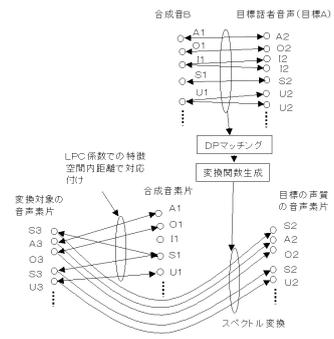
【図 2 4】



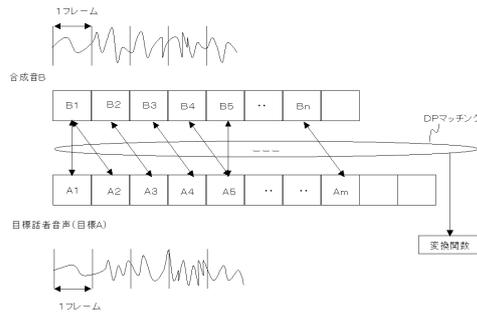
【図25】



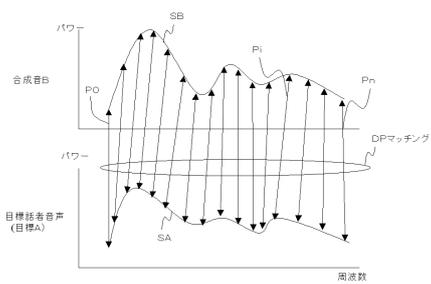
【図26】



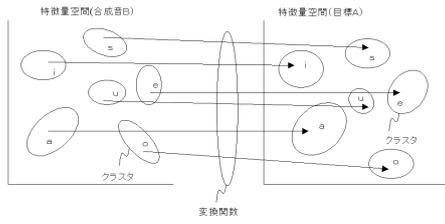
【図27】



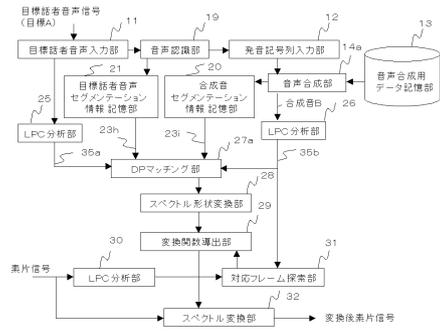
【図28】



【図 29】



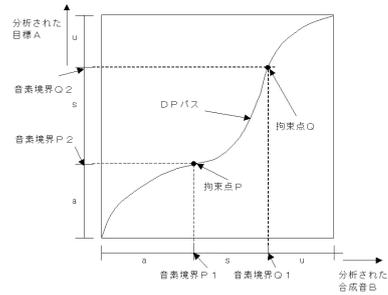
【図 30】



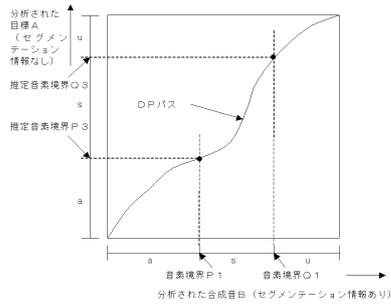
【図 31】

フレーム番号	1	2	...	10	11	...	50	51	...
自音	i	i	...	i	e	...	u	s	...

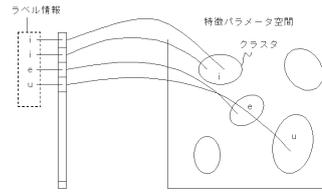
【図 32】



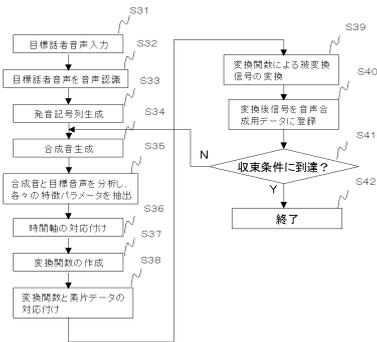
【図 33】



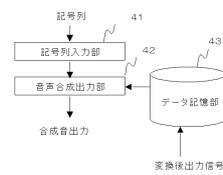
【図 34】



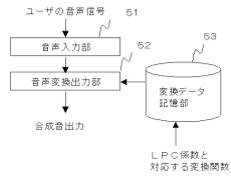
【図 35】



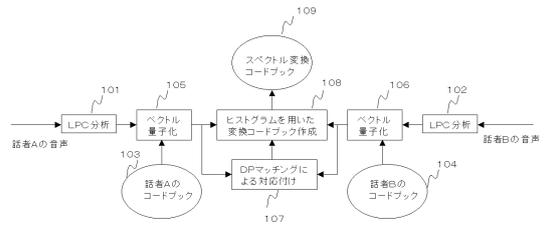
【図 36】



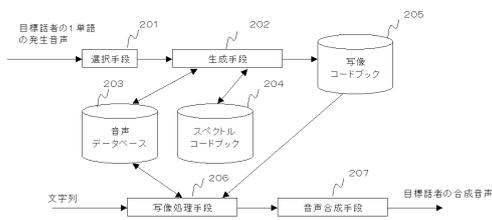
【図 37】



【図 38】



【図 39】



---

フロントページの続き

- (56)参考文献 特開2002-215199(JP,A)  
特開2000-29487(JP,A)  
特開平10-301599(JP,A)  
特開平09-179576(JP,A)  
特開平09-244694(JP,A)  
特開平09-305197(JP,A)  
特許第2880508(JP,B2)  
特開2001-034280(JP,A)