

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-4997
(P2018-4997A)

(43) 公開日 平成30年1月11日(2018.1.11)

(51) Int.Cl.		F I		テーマコード (参考)
G 1 0 L 13/06 (2013.01)		G 1 0 L 13/06	1 3 0	
G 1 0 L 13/06 (2013.01)		G 1 0 L 13/08	1 5 0 A	

審査請求 未請求 請求項の数 4 O L (全 13 頁)

(21) 出願番号	特願2016-132586 (P2016-132586)	(71) 出願人	000004352 日本放送協会 東京都渋谷区神南2丁目2番1号
(22) 出願日	平成28年7月4日(2016.7.4)	(71) 出願人	591053926 一般財団法人NHKエンジニアリングシステム 東京都世田谷区砧一丁目10番11号
		(74) 代理人	100121119 弁理士 花村 泰伸
		(72) 発明者	尾上 和穂 東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内
		(72) 発明者	齋藤 礼子 東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内

最終頁に続く

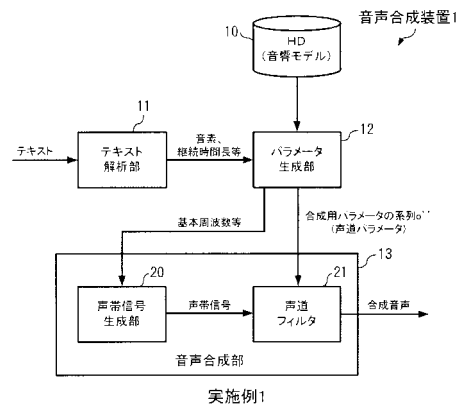
(54) 【発明の名称】 音声合成装置及びプログラム

(57) 【要約】

【課題】合成音声は自然の音に近くなるように、合成用パラメータを生成する。

【解決手段】音声合成装置1のパラメータ生成部12は、テキストの音素列について、音素の各状態に対応するパラメータの系列oを音響モデルから生成する。パラメータ生成部12は、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列oが観測される音響モデルの確率分布の尤度が最大となるパラメータの系列o'を生成し、さらに、合成用パラメータの系列o'に対し、MFCC等の分散に基づいた正規乱数(ガウスノイズ)を付加し、合成用パラメータの系列o''を生成する。音声合成部13は、ピッチの基本周波数等に基づいて、テキストに対応する声帯信号を生成し、声帯信号に対し、合成用パラメータの系列o''を用いて声道フィルタによるフィルタ処理を施すことで、合成音声を生成する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

予め学習した音響モデルを用いて、テキストに対応した合成音声を生成する音声合成装置において、

前記音響モデルを用いて、テキストに対応した合成用パラメータを生成するパラメータ生成部と、

前記パラメータ生成部により生成された前記合成用パラメータを用いて、前記テキストに対応した声帯信号に対し、声道フィルタのフィルタ処理を施して合成音声を生成する音声合成部と、を備え、

前記パラメータ生成部は、

前記テキストを構成する音素の各状態に対応する特徴量を、前記音響モデルから読み出し、前記特徴量の確率分布の尤度が最大となるパラメータを生成するパラメータ生成手段と、

前記テキストを構成する音素の各状態に対応する特徴量の分散を、前記音響モデルから読み出し、前記パラメータ生成手段により生成された前記パラメータに、前記分散に基づいた正規乱数を付加し、前記合成用パラメータを生成する分散付加手段と、を備えたことを特徴とする音声合成装置。

【請求項 2】

請求項 1 に記載の音声合成装置において、

前記テキストを構成する音素の各状態に対応する特徴量を、MFCC、MFCCの一次回帰係数及びMFCCの二次回帰係数とする、ことを特徴とする音声合成装置。

【請求項 3】

請求項 1 または 2 に記載の音声合成装置において、

前記パラメータ生成部に備えた前記分散付加手段に代わる新たな分散付加手段は、

前記テキストを構成する音素の各状態に対応する特徴量の分散を、前記音響モデルから読み出し、前記パラメータ生成手段により生成された前記パラメータに、前記分散に基づいた正規乱数を付加し、複数の合成用パラメータを生成し、

前記音声合成部に代わる新たな音声合成部は、

前記新たな分散付加手段により生成された前記複数の合成用パラメータを用いて、前記テキストに対応した声帯信号に対して声道フィルタのフィルタ処理を施し、前記複数の合成用パラメータのそれぞれに対応する複数の合成音声を生成し、前記複数の合成音声を平均化する、ことを特徴とする音声合成装置。

【請求項 4】

コンピュータを、請求項 1 から 3 までのいずれか一項に記載の音声合成装置として機能させるためのプログラム。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、統計的音響モデルを用いて、音声合成のための合成用パラメータを生成する音声合成装置及びプログラムに関する。

【背景技術】**【0002】**

従来、テキストを分析して言語情報を生成し、言語情報からテキストの文に対応する音声波形を生成する音声合成手法が知られている。この音声合成手法は、大きく 2 種類に分類することができ、一方が波形接続音声合成手法であり、他方が統計的音声合成手法である。

【0003】

波形接続音声合成手法は、得られた音韻の系列に従って音声素片を選択し、韻律情報に従って音声素片のピッチ及び継続時間長を変形して接続することで、合成音声を生成する。これに対し、統計的音声合成手法は、予め学習しておいた統計量に基づいて、最適な合

10

20

30

40

50

成用パラメータを生成することで、合成音声を生成する。

【0004】

具体的には、統計的音声合成手法を用いる音声合成装置は、入力したテキストを音素に変換し、予め学習しておいた各音素の状態毎の統計量を有する音響モデルを用いて、確率分布の尤度が最も高い合成用パラメータを生成する（例えば、非特許文献1を参照）。

【0005】

ここで、音響モデルを、音素の状態の系列を q 、音響モデルから観測されるパラメータの系列を o 、音素の状態の系列 q と音響モデルにおいて、パラメータの系列 o が観測される確率分布を P とすると、合成用パラメータは、以下の式にて推定される。

【数1】

$$o' = \arg \max_o P(o|q, \lambda) \quad \dots (1)$$

前記数式(1)は、音素の状態の系列 q と音響モデルにおいて、パラメータの系列 o が観測される確率分布 P の尤度が最大となるパラメータの系列 o を、合成用パラメータの系列 o' として求めることを表している。

【0006】

音響モデルを生成する際には、大量の学習用音声から、音素の状態毎にMFCC (Mel-Frequency Cepstrum Coefficients:メル周波数ケプストラム係数)、その一次回帰係数及び二次回帰係数を含む特徴量が学習される。そして、音響モデルには、これらのパラメータの平均及び分散、及び、音素の状態間の遷移確率、並びに、出力確率が格納される。このように、合成用パラメータの系列は、MFCCの一次回帰係数及び二次回帰係数が考慮されて生成されるため、発話文章全体で平均をトレースするような滑らか値となる。

【0007】

図6は、従来技術において、音響モデルを用いて生成された合成用パラメータの系列を示す図である。横軸は時間 t を示し、縦軸は、MFCC c 及びMFCCの一次回帰係数 c の値を示す。

【0008】

図6に示すように、従来の統計的音声合成手法により生成される合成用パラメータの系列 o' (c , c)は、音素を構成する複数の状態について、実線で表す平均を基準にして、滑らかに変化する値となる。四角の点線は分散の範囲を示す。

【0009】

しかしながら、このような統計的音声合成手法により生成される合成用パラメータの系列 o' は、図6に示したとおり、平均を基準にして滑らかに変化する値となっており、自然な音声から得られる揺らぎを持つパラメータの振る舞いとは一致しない。つまり、この合成用パラメータの系列 o' により生成される合成音声は、不自然な音になってしまう。

【0010】

一方で、音素間で合成用パラメータの値が不連続とならないように、滑らかに変化する自然な合成用パラメータの系列を生成する手法が開示されている（例えば、特許文献1を参照）。この手法は、言語レベル（音素、音節、単語等）を単位とする言語区間のスペクトルパラメータを算出し、複数の言語区間のそれぞれに対するスペクトルパラメータを、言語情報に基づいて複数のクラスタに分類する。そして、同一クラスタに属する複数のスペクトルパラメータの特徴を示す音響モデルをクラスタ毎に生成する。合成用パラメータの系列を生成する際には、テキストの言語情報に基づいて、クラスタ毎の音響モデルからテキストに応じた音響モデルを選択する。

【先行技術文献】

【特許文献】

【0011】

【特許文献1】特開2010-237323号公報

10

20

30

40

50

【非特許文献】

【0012】

【非特許文献1】Keiichi Tokuda, Heiga Zen, “Fundamentals and recent advances in HMM-based speech synthesis” [online], [平成28年5月20日検索], インターネット<URL: <http://hts.sp.nitech.ac.jp/?Tutorial>>

【発明の概要】

【発明が解決しようとする課題】

【0013】

前述の特許文献1の手法は、言語区間に対応したクラスタ毎の音響モデルから、テキストに応じた音響モデルを選択することで、音素間で不連続点のない滑らかに変化する合成用パラメータの系列を生成することができる。

10

【0014】

しかしながら、人間の発声には毎回揺らぎが存在し、実際の音声信号から得られる特徴量は、滑らかに変化する特性を有さない。このため、特許文献1の手法により生成される合成用パラメータの系列は、発話文章全体で平均をトレースするような滑らかな値となることに変わりがなく、前述の非特許文献1の手法と同様に、合成音声の音は、実際の人間の発声に特有な揺らぎが表現できず、不自然であるという問題があった。

【0015】

そこで、本発明は前記課題を解決するためになされたものであり、その目的は、自然の音に近い合成音声を生成可能な音声合成装置及びプログラムを提供することにある。

20

【課題を解決するための手段】

【0016】

前記課題を解決するために、請求項1の音声合成装置は、予め学習した音響モデルを用いて、テキストに対応した合成音声を生成する音声合成装置において、前記音響モデルを用いて、テキストに対応した合成用パラメータを生成するパラメータ生成部と、前記パラメータ生成部により生成された前記合成用パラメータを用いて、前記テキストに対応した声帯信号に対し、声道フィルタのフィルタ処理を施して合成音声を生成する音声合成部と、を備え、前記パラメータ生成部が、前記テキストを構成する音素の各状態に対応する特徴量を、前記音響モデルから読み出し、前記特徴量の確率分布の尤度が最大となるパラメータを生成するパラメータ生成手段と、前記テキストを構成する音素の各状態に対応する特徴量の分散を、前記音響モデルから読み出し、前記パラメータ生成手段により生成された前記パラメータに、前記分散に基づいた正規乱数を付加し、前記合成用パラメータを生成する分散付加手段と、を備えたことを特徴とする。

30

【0017】

また、請求項2の音声合成装置は、請求項1に記載の音声合成装置において、前記テキストを構成する音素の各状態に対応する特徴量を、MFCC、MFCCの一次回帰係数及びMFCCの二次回帰係数とする、ことを特徴とする。

【0018】

また、請求項3の音声合成装置は、請求項1または2に記載の音声合成装置において、前記パラメータ生成部に備えた前記分散付加手段に代わる新たな分散付加手段が、前記テキストを構成する音素の各状態に対応する特徴量の分散を、前記音響モデルから読み出し、前記パラメータ生成手段により生成された前記パラメータに、前記分散に基づいた正規乱数を付加し、複数の合成用パラメータを生成し、前記音声合成部に代わる新たな音声合成部が、前記新たな分散付加手段により生成された前記複数の合成用パラメータを用いて、前記テキストに対応した声帯信号に対して声道フィルタのフィルタ処理を施し、前記複数の合成用パラメータのそれぞれに対応する複数の合成音声を生成し、前記複数の合成音声を平均化する、ことを特徴とする。

40

【0019】

さらに、請求項4のプログラムは、コンピュータを、請求項1から3までのいずれか一項に記載の音声合成装置として機能させることを特徴とする。

50

【発明の効果】

【0020】

以上のように、本発明によれば、自然の音に近い合成音声を生成することが可能となる。

【図面の簡単な説明】

【0021】

【図1】実施例1の音声合成装置の構成例を示すブロック図である。

【図2】パラメータ生成部において、合成用パラメータの系列を生成するための構成例を示すブロック図である。

【図3】音響モデルを用いて生成された合成用パラメータの系列 θ の例を示す図である。

10

【図4】音響モデルのパラメータの系列 θ 、合成用パラメータの系列 θ' 及び合成用パラメータの系列 θ'' の例を示す図である。

【図5】実施例2の音声合成装置の構成例を示すブロック図である。

【図6】従来技術において、音響モデルを用いて生成された合成用パラメータの系列を示す図である。

【発明を実施するための形態】

【0022】

以下、本発明を実施するための形態について図面を用いて詳細に説明する。本発明は、合成音声の合成用パラメータを生成する際に、予め学習しておいた音響モデルの統計量である分散を用いて、揺らぎを付加した合成用パラメータを生成することを特徴とする。

20

【0023】

つまり、本発明では、音響モデルの統計量である分散に基づいた正規乱数（ガウスノイズ）を、人間の発声の揺らぎとみなし、この揺らぎを合成用パラメータへ反映する。これにより、合成用パラメータを用いて生成される合成音声は、人間の発声の音に近くなる。つまり、自然の音に近い合成音声を生成することが可能となる。

【0024】

（音声合成装置／実施例1）

まず、実施例1の音声合成装置について説明する。図1は、実施例1の音声合成装置の構成例を示すブロック図である。この音声合成装置1は、HD10、テキスト解析部11、パラメータ生成部12及び音声合成部13を備えている。

30

【0025】

HD10には、予め学習しておいた音響モデル（統計的音響モデル、HMM（隠れマルコフモデル））が格納されている。この音響モデルは、音素を構成する複数の状態（音素の始まり、中間及び終わりの各状態）に対するガウス分布のパラメータ、及び、各状態間の遷移確率、並びに、各状態における出力確率により構成される。ここで、ガウス分布のパラメータとは、例えば、人間の音声知覚の特徴を考慮した声道特性を表す特徴量であるメル周波数ケプストラム係数（MFCC）、このMFCCの一次回帰係数及びMFCCの二次回帰係数からなるスペクトルパラメータの平均及び分散から構成される。

40

【0026】

図示しない音響モデル学習部は、音響モデルを以下の手順により生成する。すなわち、音響モデル学習部は、学習対象の音声信号に対し、当該音声信号を構成する各フレームのMFCCを算出し、音素区間等の複数のMFCCをベクトル化したスペクトルパラメータを算出する。そして、音響モデル学習部は、複数のスペクトルパラメータを近似するガウス分布のパラメータ、及び、各状態間の遷移確率、並びに、各状態における出力確率を算出する。これにより、音素を構成する複数の状態のそれぞれについてのパラメータからなる音響モデルを生成する。

【0027】

テキスト解析部11は、合成音声の生成対象であるテキストを入力し、テキストに対して形態素解析等の処理を行う。これにより、テキストの音素列、音素の開始時間及び終了

50

時間、アクセントの有無、有声音及び無声音の区別情報等を生成する。そして、テキスト解析部 11 は、音素毎に、その開始時間及び終了時間から継続時間長を算出する。テキスト解析部 11 は、テキストの音素列、音素毎の継続時間長、アクセントの有無、有声音及び無声音の区別情報等をパラメータ生成部 12 に出力する。

【0028】

パラメータ生成部 12 は、テキスト解析部 11 からテキストの音素列、音素毎の継続時間長、アクセントの有無、有声音及び無声音の区別情報等を入力する。そして、パラメータ生成部 12 は、まず、音素列の音素の各状態に対応するパラメータの系列 o を、HD10 の音響モデルから生成する。つまり、パラメータ生成部 12 は、音素列の音素の各状態に対応するそれぞれのパラメータを HD10 の音響モデルから読み出し、パラメータの系列 o を生成する。次に、パラメータ生成部 12 は、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列 o が観測される確率分布 P の尤度が最大となり、かつ MFCC 等の分散が反映された合成用パラメータの系列 o' を生成する。パラメータ生成部 12 は、合成用パラメータの系列 o' を音声合成部 13 に出力する。

10

【0029】

また、パラメータ生成部 12 は、音素列の各音素について、その継続時間長及びアクセントの有無等の情報に基づいて、ピッチの基本周波数を生成する。そして、パラメータ生成部 12 は、音素列の各音素についてのピッチの基本周波数、有声音及び無声音の区別情報等を音声合成部 13 に出力する。

20

【0030】

図 2 は、パラメータ生成部 12 において、合成用パラメータの系列を生成するための構成例を示すブロック図である。このパラメータ生成部 12 は、パラメータ生成手段 30 及び分散付加手段 31 を備えている。

【0031】

パラメータ生成手段 30 は、テキスト解析部 11 から、テキストの音素列及び音素の継続時間長を入力し、HD10 の MFCC の一次回帰係数 及び MFCC の二次回帰係数 からなるスペクトルパラメータから生成した音響モデルから、音素の各状態に対応するパラメータの系列 o を生成する。そして、パラメータ生成手段 30 は、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列 o が観測される音響モデルの確率分布 P の尤度が最大となる合成用パラメータの系列 o' を生成する。そして、パラメータ生成手段 30 は、合成用パラメータの系列 o' を分散付加手段 31 に出力する。

30

【0032】

具体的には、パラメータ生成手段 30 は、音響モデルを M 、音素の状態の系列を q 、音響モデル M から観測されるパラメータの系列を o 、音素の状態の系列 q と音響モデル M において、パラメータの系列 o が観測される音響モデルの確率分布を P として、音素の各状態について、前記数式 (1) により、音素の状態の系列 q と音響モデル M において、パラメータの系列 o が観測される音響モデル M の確率分布 P の尤度が最大となるパラメータの系列 o を合成用パラメータの系列 o' として算出する。

【0033】

これにより、合成用パラメータの系列 o' は、MFCC の一次回帰係数 及び二次回帰係数 が考慮されて生成されるため、発話文章全体で平均をトレースするような滑らかな値となる。

40

【0034】

分散付加手段 31 は、パラメータ生成手段 30 から合成用パラメータの系列 o' を入力すると共に、HD10 の音響モデルから、音素列の音素の各状態について合成用パラメータの系列 o' に対応する MFCC 等の分散を読み出す。

【0035】

分散付加手段 31 は、音素の各状態の合成用パラメータの系列 o' に対し、合成用パラメータの系列 o' に対応する分散に基づいた正規乱数 (ガウスノイズ) を付加し、合成用

50

パラメータの系列 o'' を生成する。そして、分散付加手段 31 は、合成用パラメータの系列 o'' を音声合成部 13 に出力する。

【0036】

この分散に基づいた正規乱数（ガウスノイズ）は、人間の発声の揺らぎを表現するものであり、分散付加手段 31 により生成される合成用パラメータの系列 o'' は、人間の発声の揺らぎの成分が付加されたパラメータの系列となる。

【0037】

具体的には、分散付加手段 31 は、MFCC 等の分散を σ^2 、分散 σ^2 に基づいたガウスノイズを $N(0, \sigma^2)$ として、合成用パラメータの系列 o' に、分散 σ^2 に基づいたガウスノイズ $N(0, \sigma^2)$ を加算することで、合成用パラメータの系列 o'' を求める。つまり、分散付加手段 31 は、合成用パラメータの系列 o'' を、以下の数式（2）にて算出する。

[数2]

$$o'' = o' + N(0, \sigma^2) \quad \dots (2)$$

【0038】

図3は、音響モデルを用いて生成された合成用パラメータの系列 o'' の例を示す図である。横軸は時間 t を示し、縦軸はMFCC c を示す。図3には、音響モデルのパラメータの系列 o （MFCC c の平均）、合成用パラメータの系列 o' 及び合成用パラメータの系列 o'' の特性を示す。四角の点線は分散を示す。

【0039】

図4は、音響モデルのパラメータの系列 o （MFCC c の平均）、合成用パラメータの系列 o' 及び合成用パラメータの系列 o'' の例を示す図であり、図3に示す分散の領域 σ^2 を拡大した図である。図4（1）は音響モデルのパラメータの系列 o （MFCC c の平均）の例であり、図4（2）は合成用パラメータの系列 o' の例であり、図4（3）は合成用パラメータの系列 o'' の例である。

【0040】

図4（1）は、音響モデルのパラメータの系列 o （MFCC c の平均）が、音素の状態毎に一定値であることを示す。図4（2）は、合成用パラメータの系列 o' が、図6に示した合成用パラメータの系列 o' の特性と同じであることを示す。この合成用パラメータの系列 o' は、MFCC の一次回帰係数 β_1 及び二次回帰係数 β_2 が考慮されて生成されるため、音響モデルのパラメータの系列 o である MFCC c の平均を基準にして、滑らかに変化する値となる。この合成用パラメータの系列 o' は、図2に示したパラメータ生成手段 30 により生成される。

【0041】

図4（3）は、合成用パラメータの系列 o'' が、図4（2）に示した合成用パラメータの系列 o' に対し、人間の発声の揺らぎ成分を表現する、MFCC 等の分散に基づいた正規乱数（ガウスノイズ）が付加された特性となることを示す。この合成用パラメータの系列 o'' は、図2に示した分散付加手段 31 により生成される。

【0042】

このように、パラメータ生成部 12 は、音素の各状態について、MFCC 等の分散を反映した合成用パラメータの系列 o'' 、すなわち、音素の状態の系列と音響モデルにおいて、パラメータの系列 o が観測される音響モデルの確率分布 P の尤度を最大とし、かつ MFCC 等の分散を反映した合成用パラメータの系列 o'' を生成する。この合成用パラメータの系列 o'' は、MFCC 等の分散を考慮した揺らぎが与えられ、後述する声道フィルタ 21 は、これを声道パラメータとして用いる。

【0043】

図1に戻って、音声合成部 13 は、パラメータ生成部 12 から、音素の各状態におけるピッチの基本周波数、有声音及び無声音の区別情報等、並びに合成用パラメータの系列 o'' を入力する。そして、音声合成部 13 は、基本周波数、有声音及び無声音の区別情報等に基づいて、テキストに対応する声帯信号を生成し、声帯信号に対し、合成用パラメータの系列 o'' を用いて声道フィルタによるフィルタ処理を施すことで、合成音声を生成する

10

20

30

40

50

。音声合成部 13 は、生成した合成音声を出力する。

【0044】

図 1 に示すように、音声合成部 13 は、声帯信号生成部 20 及び声道フィルタ 21 を備えている。声帯信号生成部 20 は、パラメータ生成部 12 から、音素列の各音素についてのピッチの基本周波数、有声音及び無声音の区別情報等を入力する。そして、声帯信号生成部 20 は、対象区間が有声音である場合、基本周波数の逆数であるピッチ周期のパルス列を声帯信号として生成し、対象区間が無声音である場合、白色雑音を声帯信号として生成する。声帯信号生成部 20 は、生成した声帯信号を声道フィルタ 21 に出力する。

【0045】

声道フィルタ 21 は、声帯信号生成部 20 から声帯信号を入力すると共に、パラメータ生成部 12 から合成用パラメータの系列 σ を入力し、声帯信号に対し、合成用パラメータの系列 σ を用いて声道フィルタによるフィルタ処理を施すことで、合成音声を生成する。声道フィルタ 21 は、生成した合成音声を出力する。尚、声道フィルタ 21 による合成音声の生成手法は既知であり、非特許文献 1 等を参照されたい。

【0046】

以上のように、実施例 1 の音声合成装置 1 によれば、パラメータ生成部 12 は、テキストの音素列について、音素の各状態に対応するパラメータの系列 σ を音響モデルから生成する。そして、パラメータ生成部 12 は、前記数式 (1) にて、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列 σ が観測される音響モデルの確率分布の尤度が最大となる合成用パラメータの系列 σ' を生成する。さらに、パラメータ生成部 12 は、前記数式 (2) にて、合成用パラメータの系列 σ' に対し、MFC C 等の分散に基づいた正規乱数 (ガウスノイズ) を付加し、合成用パラメータの系列 σ を生成する。

【0047】

音声合成部 13 は、ピッチの基本周波数等に基づいて、テキストに対応する声帯信号を生成し、声帯信号に対し、合成用パラメータの系列 σ を用いて声道フィルタによるフィルタ処理を施すことで、合成音声を生成する。

【0048】

このように、音響モデルの統計量である分散を用いて、人間の発声の揺らぎを表現する成分を合成用パラメータの系列 σ に反映するようにした。これにより、音響モデルの学習時に用いる音声信号の分布に適合した合成用パラメータの系列 σ を生成することができる。つまり、人間の音声のパラメータに近い合成用パラメータの系列 σ である声道パラメータを生成することができるため、この声道パラメータを用いて生成される合成音声は、実際に人間が発声する音に近くなる。すなわち、自然の音に近い合成音声を生成することができ、人間の発話に近い統計的合成音声の提供が可能となる。

【0049】

(音声合成装置 / 実施例 2)

次に、実施例 2 の音声合成装置について説明する。図 5 は、実施例 2 の音声合成装置の構成例を示すブロック図である。この音声合成装置 2 は、HD 10、テキスト解析部 11、パラメータ生成部 14 及び音声合成部 15 を備えている。

【0050】

図 1 に示した実施例 1 の音声合成装置 1 と、図 5 に示す実施例 2 の音声合成装置 2 とを比較すると、両音声合成装置 1, 2 は、同じテキスト解析部 11 を備えている点で共通する。一方、実施例 2 の音声合成装置 2 は、実施例 1 の音声合成装置 1 に備えたパラメータ生成部 12 及び音声合成部 13 とは異なるパラメータ生成部 14 及び音声合成部 15 を備えている点で、実施例 1 の音声合成装置 1 と相違する。

【0051】

テキスト解析部 11 は、図 1 に示したテキスト解析部 11 と同一であるから、ここでは説明を省略する。

【0052】

10

20

30

40

50

パラメータ生成部 14 は、テキスト解析部 11 からテキストの音素列、音素毎の継続時間長、アクセントの有無、有声音及び無声音の区別情報等を入力する。そして、パラメータ生成部 12 は、音素列の音素の各状態に対応するパラメータの系列 o を、HD10 の音響モデルから生成する。

【0053】

パラメータ生成部 14 は、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列 o が観測される音響モデルの確率分布の尤度が最大となり、かつ MFCC 等の分散が反映された複数の合成用パラメータの系列 o を生成し、複数の合成用パラメータの系列 o を音声合成部 15 に出力する。具体的には、パラメータ生成部 14 は、前記数式 (1) にて、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列 o が観測される音響モデルの確率分布の尤度が最大となる合成用パラメータの系列 o' を生成する。さらに、パラメータ生成部 14 は、前記数式 (2) にて、合成用パラメータの系列 o' に対し、MFCC 等の分散に基づいた正規乱数 (ガウスノイズ) を付加し、複数の合成用パラメータの系列 o を生成する。

10

【0054】

また、パラメータ生成部 14 は、図 1 に示したパラメータ生成部 12 と同様に、ピッチの基本周波数を生成し、音素列の各音素についてのピッチの基本周波数等を音声合成部 15 に出力する。

【0055】

例えば、パラメータ生成部 14 は、合成用パラメータの系列 o' に対し、MFCC 等の分散に基づいた正規乱数 (ガウスノイズ) を付加することで、異なる合成用パラメータの系列 o_1 , o_2 , o_3 を生成する。ここで、合成用パラメータの系列 o_1 , o_2 , o_3 は、同一の合成用パラメータの系列 o' に、MFCC 等の分散に基づく異なる正規乱数 (ガウスノイズ) を付加して生成する。尚、MFCC 等の分散を用いるから、合成用パラメータの系列 o_1 , o_2 , o_3 は、生成処理毎に異なる値となる。

20

【0056】

音声合成部 15 は、パラメータ生成部 14 から、音素列の各音素についてのピッチの基本周波数等、及び複数の合成用パラメータの系列 o_1 , o_2 , o_3 をそれぞれ入力する。そして、音声合成部 15 は、ピッチの基本周波数等に基づいて、テキストに対応する声帯信号を生成し、声帯信号に対し、複数の合成用パラメータの系列 o_1 , o_2 , o_3 のそれぞれを用いて声道フィルタによるフィルタ処理を施すことで、複数の合成音声 w_1 , w_2 , w_3 を生成する。尚、複数の合成用パラメータの系列 o_1 , o_2 , o_3 のそれぞれを用いて声道フィルタによるフィルタ処理を施すことで生成された複数の合成音声 w_1 , w_2 , w_3 は、時間が揃っている。音声合成部 15 は、生成した複数の合成音声 w_1 , w_2 , w_3 の時間波形を平均化し、平均化した時間波形を合成音声として出力する。

30

【0057】

図 5 に示すように、音声合成部 15 は、声帯信号生成部 20、声道フィルタ 22 及び平均化部 23 を備えている。声帯信号生成部 20 は、図 1 に示した声帯信号生成部 20 と同一であるため、ここでは説明を省略する。

40

【0058】

声道フィルタ 22 は、声帯信号生成部 20 から声帯信号を入力すると共に、パラメータ生成部 14 から複数の合成用パラメータの系列 o_1 , o_2 , o_3 をそれぞれ入力する。そして、声道フィルタ 22 は、声帯信号に対し、複数の合成用パラメータの系列 o_1 , o_2 , o_3 のそれぞれを用いて声道フィルタによるフィルタ処理を施すことで、複数の合成音声 w_1 , w_2 , w_3 を生成する。声道フィルタ 22 は、生成した複数の合成音声 w_1 , w_2 , w_3 を平均化部 23 に出力する。尚、声道フィルタ 22 による合成音声の生成手法は既知であるため、ここでは詳細な説明を省略する。

【0059】

声道フィルタ 22 は、第 1 の合成音声 w_1 、第 2 の合成音声 w_2 及び第 3 の合成音声 w_3

50

3におけるそれぞれの時間波形を平均化部23に出力する。ここで、第1の合成音声 w_1 、第2の合成音声 w_2 及び第3の合成音声 w_3 は、同じテキストに対応した音声である。一方、合成用パラメータの系列 o_1 、 o_2 、 o_3 が異なるため、異なる時間波形となる。

【0060】

平均化部23は、声道フィルタ22から複数の合成音声 w_1 、 w_2 、 w_3 の時間波形を入力し、複数の合成音声 w_1 、 w_2 、 w_3 の時間波形を平均化し、平均化した時間波形を合成音声として出力する。

【0061】

以上のように、実施例2の音声合成装置2によれば、パラメータ生成部14は、テキストの音素列について、音素の各状態に対応するパラメータの系列 o を音響モデルから生成する。そして、パラメータ生成部14は、前記数式(1)及び前記数式(2)にて、音素の各状態について、音素の状態の系列と音響モデルにおいて、パラメータの系列 o が観測される音響モデルの確率分布の尤度が最大となり、かつMFC C等の分散が反映された複数の合成用パラメータの系列 o_1 、 o_2 、 o_3 を生成する。

10

【0062】

音声合成部15は、ピッチ基本周波数等に基づいて、テキストに対応する声帯信号を生成し、声帯信号に対し、複数の合成用パラメータの系列 o_1 、 o_2 、 o_3 を用いて声道フィルタによるフィルタ処理を施すことで、複数の合成音声 w_1 、 w_2 、 w_3 を生成する。そして、音声合成部15は、複数の合成音声 w_1 、 w_2 、 w_3 の時間波形を平均化する。

20

【0063】

このように、音響モデルの統計量である分散を用いて、人間の発声の揺らぎを表現する成分を複数の合成用パラメータの系列 o_1 、 o_2 、 o_3 に反映するようにした。これにより、これらの複数の合成用パラメータの系列 o_1 、 o_2 、 o_3 を用いて生成される合成音声は、人間の発声の音に一層近くなる。すなわち、自然の音に一層近い合成音声を生成することができる。

【0064】

尚、本発明の実施例1、2による音声合成装置1、2のハードウェア構成としては、通常のコンピュータを使用することができる。音声合成装置1、2は、CPU、RAM等の揮発性の記憶媒体、ROM等の不揮発性の記憶媒体、及びインターフェース等を備えたコンピュータによって構成できる。音声合成装置1に備えたHD10、テキスト解析部11、パラメータ生成部12及び音声合成部13の各機能は、これらの機能を記述したプログラムをCPUに実行させることによりそれぞれ実現できる。また、音声合成装置2に備えたHD10、テキスト解析部11、パラメータ生成部14及び音声合成部15の各機能は、これらの機能を記述したプログラムをCPUに実行させることによりそれぞれ実現できる。

30

【0065】

これらのプログラムは、磁気ディスク(フロッピー(登録商標)ディスク、ハードディスク等)、光ディスク(CD-ROM、DVD等)、半導体メモリ等の記憶媒体に格納して頒布することもでき、ネットワークを介して送受信することもできる。

40

【0066】

以上、実施例1、2を挙げて本発明を説明したが、本発明は前記実施例1、2に限定されるものではなく、その技術思想を逸脱しない範囲で種々変形可能である。例えば、前記実施例1において、音声合成装置1のパラメータ生成部12は、音響モデルとして学習されたMFC C等の統計量を用いて、MFC C等の分散が反映された合成用パラメータの系列 o を生成するようにした。

【0067】

この場合、パラメータ生成部12は、MFC Cの統計量を用いて、MFC Cの分散が反映された合成用パラメータの系列 o を生成するようにしてもよい。具体的には、パラメ

50

ータ生成部 12 のパラメータ生成手段 30 は、前記数式 (1) にて、音素の状態の系列と、MFCC を用いて学習した音響モデルにおいて、パラメータの系列 θ が観測される確率分布の尤度が最大となる合成用パラメータの系列 θ' を生成する。そして、分散付加手段 31 は、前記数式 (2) にて、合成用パラメータの系列 θ' に対し、MFCC を用いて学習した音響モデルにおける MFCC の分散に基づいた正規乱数 (ガウスノイズ) を付加し、合成用パラメータの系列 θ'' を生成する。

【0068】

また、パラメータ生成部 12 は、音響モデルとして学習された MFCC 及び MFCC の一次回帰係数 の統計量を用いて、MFCC の分散及び MFCC の一次回帰係数 の分散が反映された合成用パラメータの系列 θ'' を生成するようにしてもよい。具体的には、パラメータ生成部 12 のパラメータ生成手段 30 は、前記数式 (1) にて、音素の状態の系列と、MFCC 及び MFCC の一次回帰係数 を用いて学習した音響モデルにおいて、パラメータの系列 θ が観測される確率分布の尤度が最大となる合成用パラメータの系列 θ' を生成する。そして、分散付加手段 31 は、前記数式 (2) にて、合成用パラメータの系列 θ' に対し、MFCC 及び MFCC の一次回帰係数 を用いて学習した音響モデルにおける MFCC の分散及び MFCC の一次回帰係数 の分散に基づいた正規乱数 (ガウスノイズ) を付加し、合成用パラメータの系列 θ'' を生成する。

10

【0069】

また、パラメータ生成部 12 は、音響モデルとして学習された MFCC、MFCC の一次回帰係数 及び MFCC の二次回帰係数 の統計量を用いて、MFCC の分散、MFCC の一次回帰係数 の分散及び MFCC の二次回帰係数 の分散が反映された合成用パラメータの系列 θ'' を生成するようにしてもよい。具体的には、パラメータ生成部 12 のパラメータ生成手段 30 は、前記数式 (1) にて、音素の状態の系列と、MFCC、MFCC の一次回帰係数 及び MFCC の二次回帰係数 を用いて学習した音響モデルにおいて、パラメータの系列 θ が観測される確率分布の尤度が最大となる合成用パラメータの系列 θ' を生成する。そして、分散付加手段 31 は、前記数式 (2) にて、合成用パラメータの系列 θ' に対し、MFCC、MFCC の一次回帰係数 及び MFCC の二次回帰係数 を用いて学習した音響モデルにおける MFCC の分散、MFCC の一次回帰係数 の分散及び MFCC の二次回帰係数 の分散に基づいた正規乱数 (ガウスノイズ) を付加し、合成用パラメータの系列 θ'' を生成する。実施例 2 についても同様である。

20

30

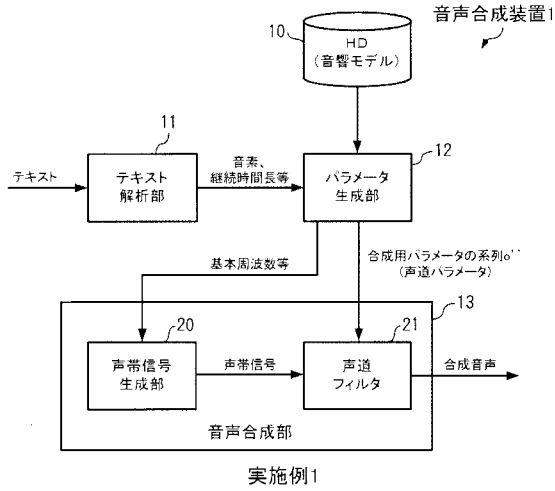
【符号の説明】

【0070】

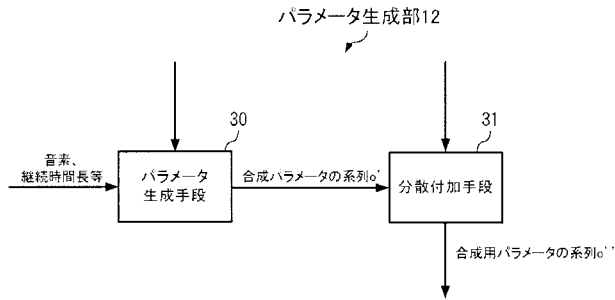
- 1, 2 音声合成装置
- 10 HD (ハードディスク)
- 11 テキスト解析部
- 12, 14 パラメータ生成部
- 13, 15 音声合成部
- 20 声帯信号生成部
- 21, 22 声道フィルタ
- 23 平均化部
- 30 パラメータ生成手段
- 31 分散付加手段

40

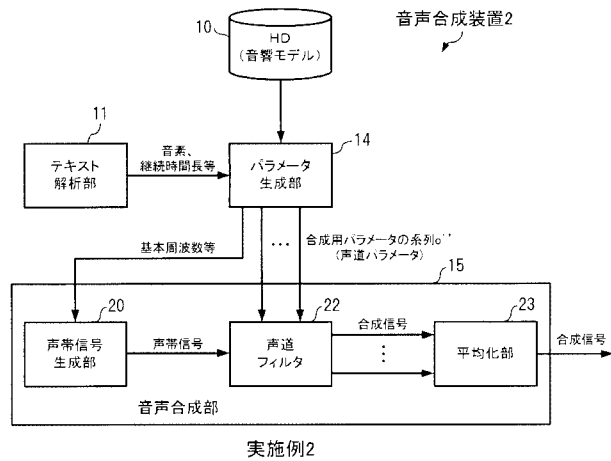
【 図 1 】



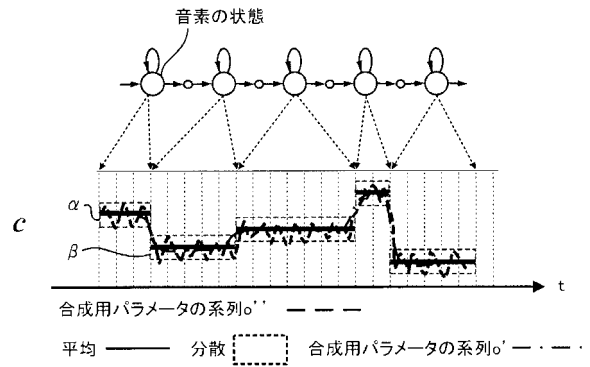
【 図 2 】



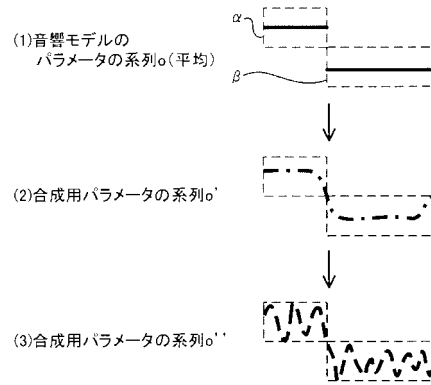
【 図 5 】



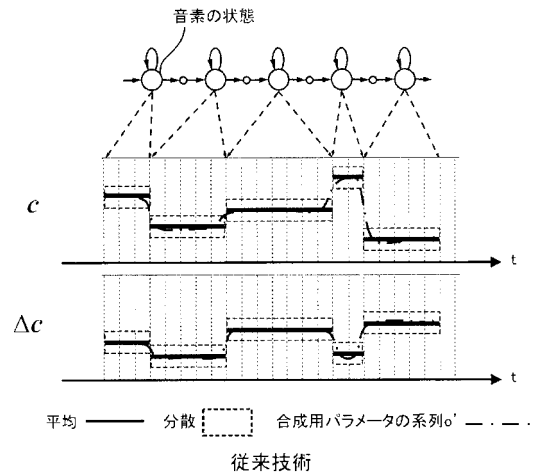
【 図 3 】



【 図 4 】



【 図 6 】



フロントページの続き

(72)発明者 清山 信正

東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内

(72)発明者 今井 篤

東京都世田谷区砧一丁目10番11号 日本放送協会放送技術研究所内

(72)発明者 都木 徹

東京都世田谷区砧一丁目10番11号 一般財団法人NHKエンジニアリングシステム内