

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4054507号  
(P4054507)

(45) 発行日 平成20年2月27日(2008.2.27)

(24) 登録日 平成19年12月14日(2007.12.14)

(51) Int.Cl.

F I

G 0 6 F 3/16 (2006.01)  
 G 1 0 L 13/00 (2006.01)  
 G 1 0 L 13/08 (2006.01)  
 G 1 0 L 13/06 (2006.01)

G 0 6 F 3/16 3 3 O K  
 G 1 0 L 3/00 E  
 G 1 0 L 3/00 H  
 G 1 0 L 5/04 F

請求項の数 11 (全 16 頁)

(21) 出願番号 特願2000-99534 (P2000-99534)  
 (22) 出願日 平成12年3月31日(2000.3.31)  
 (65) 公開番号 特開2001-282282 (P2001-282282A)  
 (43) 公開日 平成13年10月12日(2001.10.12)  
 審査請求日 平成16年12月10日(2004.12.10)  
 審判番号 不服2005-14360 (P2005-14360/J1)  
 審判請求日 平成17年7月27日(2005.7.27)

早期審査対象出願

(73) 特許権者 000001007  
 キヤノン株式会社  
 東京都大田区下丸子3丁目30番2号  
 (74) 代理人 100076428  
 弁理士 大塚 康德  
 (74) 代理人 100112508  
 弁理士 高柳 司郎  
 (74) 代理人 100115071  
 弁理士 大塚 康弘  
 (74) 代理人 100116894  
 弁理士 木村 秀二  
 (72) 発明者 深田 俊明  
 東京都大田区下丸子3丁目30番2号 キ  
 ヤノン株式会社内

最終頁に続く

(54) 【発明の名称】 音声情報処理方法および装置および記憶媒体

(57) 【特許請求の範囲】

【請求項1】

音韻系列を受信する受信工程と、

基本周波数の時間変化を多項式セグメントモデルによってモデル化したセグメントピッチパターンモデルに基づいて、前記音韻系列を構成する各音韻の基本周波数を生成する生成工程と、

前記生成工程で生成された前記各音韻の基本周波数に基づいて音声合成する音声合成工程と、

を有することを特徴とする音声情報処理方法。

【請求項2】

前記セグメントピッチパターンモデルは、音素、音節、単語の少なくともいずれかを単位としたモデルであることを特徴とする請求項1に記載の音声情報処理方法。

【請求項3】

前記セグメントピッチパターンモデルは、アクセント型、モーラ数、モーラ位置、品詞の少なくとも1つを考慮したモデルであることを特徴とする請求項1に記載の音声情報処理方法。

【請求項4】

前記セグメントピッチパターンモデルは、単一混合分布、多混合分布の少なくともいずれかによってモデリングされたモデルであることを特徴とする請求項1に記載の音声情報処理方法。

**【請求項 5】**

前記セグメントピッチパターンモデルは、アクセント句、単語、フレーズ、文の少なくともいずれかからなる単位ごとに正規化されたモデルであることを特徴とする請求項 1 に記載の音声情報処理方法。

**【請求項 6】**

請求項 1 乃至 5 のいずれか 1 項に記載の音声情報処理方法を実行するプログラムを記憶したことを特徴とする、コンピュータにより読取り可能な記憶媒体。

**【請求項 7】**

音韻系列を受信する受信手段と、

基本周波数の時間変化を多項式セグメントモデルによってモデル化したセグメントピッチパターンモデルに基づいて、前記音韻系列を構成する各音韻の基本周波数を生成する生成手段と、

前記生成手段により設定された前記各音韻の基本周波数に基づいて音声を合成する音声合成手段と、

を有することを特徴とする音声情報処理装置。

**【請求項 8】**

前記セグメントピッチパターンモデルは、音素、音節、単語の少なくともいずれかを単位としたモデルであることを特徴とする請求項 7 に記載の音声情報処理装置。

**【請求項 9】**

前記セグメントピッチパターンモデルは、アクセント型、モーラ数、モーラ位置、品詞の少なくとも 1 つを考慮したモデルであることを特徴とする請求項 7 に記載の音声情報処理装置。

**【請求項 10】**

前記セグメントピッチパターンモデルは、単一混合分布、多混合分布の少なくともいずれかによってモデリングされたモデルであることを特徴とする請求項 7 に記載の音声情報処理装置。

**【請求項 11】**

前記セグメントピッチパターンモデルは、アクセント句、単語、フレーズ、文の少なくともいずれかからなる単位ごとに正規化されたモデルであることを特徴とする請求項 7 に記載の音声情報処理装置。

**【発明の詳細な説明】****【0001】****【発明の属する技術分野】**

本発明は、音声合成或いは音声認識に際して実施される所定のセグメント単位での時系列の基本周波数（ピッチパターン）を設定する音声情報処理方法及びその装置、及び、前記音声合成方法を実施するプログラムを記憶した、コンピュータにより読取り可能な記憶媒体に関するものである。

**【0002】****【従来の技術】**

近年、任意の文字系列を音韻系列に変換し、その音韻系列を所定の音声規則合成方式に従って合成音声に変換する音声合成装置が開発されている。

**【0003】****【発明が解決しようとする課題】**

しかしながら、従来の音声合成装置から出力される合成音声は、人間が発声する自然音声と比較すると不自然で機械的なものであった。この原因の一つとして、例えば「おんせい」という文字系列を構成する音韻系列「o, X, s, e, i」において、各音韻のアクセントやイントネーションを生成する韻律生成規則の精度が挙げられる。精度が悪い場合、音韻系列に対して十分なピッチパターンが生成されないため、合成される音声は不自然で機械的なものとなる。

**【0004】**

本発明は上記従来例に鑑みてなされたもので、所定単位の音韻の基本周波数の時間変化をモデル化することにより、自然なイントネーションを与える音声合成を行うことができる音声情報処理方法及び装置を提供することを目的とする。

【 0 0 0 5 】

又本発明の目的は、所定単位の音韻の基本周波数の時間変化をモデル化することにより、このモデル化した情報を用いて高精度に音声認識ができる音声情報処理方法及び装置を提供することにある。

【 0 0 0 6 】

【課題を解決するための手段】

上記目的を達成するために本発明の音声情報処理方法は以下のような工程を備える。即ち、

音韻系列を受信する受信工程と、

基本周波数の時間変化を多項式セグメントモデルによってモデル化したセグメントピッチパターンモデルに基づいて、前記音韻系列を構成する各音韻の基本周波数を生成する生成工程と、

前記生成工程で生成された前記各音韻の基本周波数に基づいて音声合成する音声合成工程とを有することを特徴とする。

【 0 0 0 7 】

上記目的を達成するために本発明の音声情報処理方法は以下のような工程を備える。即ち、

音声を受信する受信工程と、前記音声の特徴パラメータを抽出する抽出工程と、セグメントピッチパターンモデルに基づいて、前記特徴パラメータを認識する音声認識工程と、を有することを特徴とする。

【 0 0 0 8 】

上記目的を達成するために本発明の音声情報処理装置は以下のような構成を備える。即ち、

音韻系列を受信する受信手段と、

基本周波数の時間変化を多項式セグメントモデルによってモデル化したセグメントピッチパターンモデルに基づいて、前記音韻系列を構成する各音韻の基本周波数を生成する生成手段と、

前記生成手段により設定された前記各音韻の基本周波数に基づいて音声合成する音声合成手段とを有することを特徴とする。

【 0 0 0 9 】

上記目的を達成するために本発明の音声情報処理装置は以下のような構成を備える。即ち、

音声を受信する受信手段と、前記音声の特徴パラメータを抽出する抽出手段と、セグメントピッチパターンモデルに基づいて、前記特徴パラメータを認識する音声認識手段と、を有することを特徴とする。

【 0 0 1 0 】

【発明の実施の形態】

【 0 0 1 1 】

本発明の実施の形態における多項式セグメントモデルの概要は以下の通りである。Lフレーム長のD次元の観測ベクトル $\{y_1, \dots, y_L\}$   $y_t = [y_{t,1}, y_{t,2}, \dots, y_{t,D}]$ を  $L \times D$ の行列で表現した

【 0 0 1 2 】

【数 1】

10

20

30

40

$$Y = \begin{bmatrix} y_{1,1} & y_{1,2} & \cdots & y_{1,D} \\ y_{2,1} & y_{2,2} & \cdots & y_{2,D} \\ \vdots & \vdots & & \vdots \\ y_{L,1} & y_{L,2} & \cdots & y_{L,D} \end{bmatrix} \quad (1)$$

を R 次の多項式セグメントモデルによって、

【 0 0 1 3 】

【 数 2 】

10

$$Y = ZB + E, \quad (2)$$

... 式 ( 2 )

と表す。ここで、Z は  $L \times (R + 1)$  のデザインマトリクスであり、

【 0 0 1 4 】

【 数 3 】

$$Z = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 1 & \frac{1}{L-1} & \cdots & \left(\frac{1}{L-1}\right)^R \\ \vdots & \vdots & & \vdots \\ 1 & \frac{t-1}{L-1} & \cdots & \left(\frac{t-1}{L-1}\right)^R \\ \vdots & \vdots & & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix}, \quad (3)$$

20

30

... 式 ( 3 )

と表される。また、B は  $(R + 1) \times D$  のパラメータ系列行列

【 0 0 1 5 】

【 数 4 】

$$B = \begin{bmatrix} b_{0,1} & b_{0,2} & \cdots & b_{0,D} \\ b_{1,1} & b_{1,2} & \cdots & b_{1,D} \\ \vdots & \vdots & & \vdots \\ b_{R,1} & b_{R,2} & \cdots & b_{R,D} \end{bmatrix}, \quad (4)$$

40

... 式 ( 4 )

であり、E は  $L \times D$  の予測誤差行列

【 0 0 1 6 】

【 数 5 】

$$E = \begin{bmatrix} e_{1,1} & e_{1,2} & \dots & e_{1,D} \\ e_{2,1} & e_{2,2} & \dots & e_{2,D} \\ \vdots & \vdots & & \vdots \\ e_{L,1} & e_{L,2} & \dots & e_{L,D} \end{bmatrix} \quad (5)$$

...式(5)

である。デザインマトリクスZによって異なる長さのセグメントを“0”から“1”の間に正規化することができる。

10

【0017】

セグメントYがラベルaによって生成されるときに尤度は次のように表される。

【0018】

【数6】

$$P(Y|a) = \prod_{t=1}^L f(y_t). \quad (6)$$

【0019】

20

...式(6)

上式(6)において、 $f(y_t)$ は、ラベルaに対する特徴ベクトル $y_t$ の尤度であり、次式によって与えられる。

【0020】

【数7】

$$f(y_t) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_a|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (y_t - z_t B_a)^T \Sigma_a^{-1} (y_t - z_t B_a) \right\} \quad (7)$$

30

【0021】

...式(7)

ここで、 $B_a$ と  $a$ はラベルaを表す単一ガウスセグメントモデルのパラメータである。上式において、 $z_t$ は、

【0022】

【数8】

$$z_t = \left[ 1, \frac{t-1}{L-1}, \dots, \left( \frac{t-1}{L-1} \right)^R \right] \quad (8)$$

40

...式(8)

と与えられる。いま、ラベルaに対して、K個のセグメント $Y_1, Y_2, \dots, Y_K$ がある場合に、モデルパラメータ $B_a$ 及び  $a$ を求めたいとする。このとき、 $B_a$ 及び  $a$ に対するこれらのセグメントの確率は、

【0023】

【数9】

$$\begin{aligned}
 P(Y_1, Y_2, \dots, Y_K | B_a, \Sigma_a) &= \prod_{k=1}^K P(Y_k | B_a, \Sigma_a) \\
 &= \prod_{k=1}^K \prod_{t=1}^L f(y_{k,t})
 \end{aligned} \tag{9}$$

...式(9)

と与えられる。これより、上式の確率を最大化する  $B_a$ 、 $a$ を求めることによりモデルパラメータが求まる。これらの推定値は、

【0024】

【数10】

$$\bar{B}_a = \left[ \sum_{k=1}^K Z_k^T Z_k \right]^{-1} \left[ \sum_{k=1}^K Z_k^T Y_k \right], \tag{10}$$

...式(10)

【0025】

【数11】

20

$$\bar{\Sigma}_a = \frac{\sum_{k=1}^K (Y_k - Z_k \bar{B}_a)^T (Y_k - Z_k \bar{B}_a)}{\sum_{k=1}^K L_k} \tag{11}$$

...式(11)

として得ることができる。

30

【0026】

このように、セグメントピッチパターンの時間変化を多項式によってモデリングすることによって、セグメントピッチパターンの時系列間の相関を考慮することが可能になり、前記従来例の問題点が解決できる。

【0027】

以下、添付図面を参照して本発明の好適な実施の形態を詳細に説明する。

【0028】

[実施の形態1]

図1は、本発明の実施の形態1に係る音声合成装置の構成を示すブロック図である。

【0029】

40

図1において、101はCPUで、ROM102に記憶された制御プログラム、或いは外部記憶装置104からRAM103にロードされた制御プログラムに従って、本実施の形態の音声合成装置における各種制御を行う。ROM102は、各種パラメータやCPU101が実行する制御プログラムなどを格納している。RAM103は、CPU101による各種制御の実行時に作業領域を提供するとともに、CPU101により実行される制御プログラムを記憶する。104はハードディスク、フロッピーディスク、CD-ROM等の外部記憶装置で、この外部記憶装置がハードディスクの場合には、CD-ROMやフロッピーディスク等からインストールされた各種プログラムが記憶されている。105は入力部で、キーボード、マウス等のポインティングデバイスを有している。又、この入力部105は、例えば通信回線等を介してインターネット等からのデータを入力しても良い。

50

106は液晶やCRT等の表示部で、CPU101の制御により各種データの表示を行う。107はスピーカで、音声信号(電気信号)を可聴音である音声に変換して出力する。108は上記各部を接続するバスである。109は音声合成・認識ユニットである。

#### 【0030】

図2は、本実施の形態1に係る音声合成・認識ユニット109の動作を示すフローチャートである。以下に示される各ステップは、ROM102に格納された制御プログラム、或いは外部記憶装置104からRAM103にロードされた制御プログラムをCPU101が実行することによって実現される。

#### 【0031】

まずステップS201で、漢字かな混じりの日本語テキストデータ、又は他の言語のテキストデータが入力部105から入力されるとステップS202に進み、この入力されたテキストデータを、言語解析辞書201を用いて解析し、入力テキストデータに対する音韻系列(読み)やアクセントなどの情報を抽出する。次にステップS203に進み、これらの情報を用いて、ステップS202で求めた音韻系列を構成する各音韻の継続時間長、基本周波数(セグメントピッチパターンともいう)、パワー等のプロソディ(韻律情報ともいう)を生成する。この際、セグメントピッチパターンはピッチパターンモデル202を用いて決定され、また継続時間長、パワー等は韻律制御モデル203を用いて決定される。

#### 【0032】

次にステップS204に進み、ステップS202で解析して抽出された音韻系列、及びステップS203で生成されたプロソディに基づいて、音声素片辞書204から、その音韻系列に対応する合成音声を生成するための音声素片(波形もしくは特徴パラメータ)を複数個選択する。次にステップS205に進み、それら選択された音声素片を用いて合成音声信号を生成し、ステップS206において、その生成された合成音声信号に基づいて音声をスピーカ107から出力する。最後にステップS207において、入力されたテキストデータに対する処理が全て終了したか否かの判断を行い、終了していない場合はステップS201に戻り、前述の処理が続けられる。

#### 【0033】

図3は、図2のステップS203のプロソディ生成処理で使用した上述の多項式セグメントモデルに基づくセグメントピッチパターンモデルの作成手順を示すフローチャートである。

#### 【0034】

このセグメントピッチパターンモデルを作成するためには、まずステップS301で、複数の学習サンプルを有する音声ファイル301を用いて、所定単位の音韻系列の基本周波数(ピッチパターン)を抽出する。この基本周波数の抽出において、有声・無声の判別結果、ピッチマーク等の情報を使用する場合には、基本周波数抽出に必要な情報を格納したサイド情報ファイル302も併せて利用する。

#### 【0035】

次に、ステップS302に進み、所定単位の音韻系列を構成する音素、音節、単語などを単位とした音韻の時間情報が付与されたラベルファイル303を用いて、音韻系列のピッチパターンをセグメント単位に分割する。そして最後にステップS303に進み、同一カテゴリに属するセグメント毎に、前述の式(10)及び式(11)を用いてセグメントピッチパターンモデルのモデルパラメータを計算する。

#### 【0036】

以下、具体例を挙げて本実施の形態1に係る処理手順を、図3乃至図9を参照して説明する。

#### 【0037】

図4は、サイド情報ファイル302に記憶された音韻系列「音声(oNsee)」に関するサイド情報の一例を示す図、図5は、図4の有声区間(o,N,e,e)に対する基本周波数の一例を示す図、図6はラベルファイル303に記憶された音韻系列「音声(oNsee)」に関する情

10

20

30

40

50

報の一例を示す図、図 7 は図 5 のピッチパターンをモデル化した図、図 8 は音韻系列「アクセント(akuseNto)」に対するピッチパターンの一例を示す図、そして図 9 はラベルファイル 3 0 3 に記憶された音韻系列「アクセント(akuseNto)」に関する情報の一例を示す図である。

#### 【 0 0 3 8 】

いま音韻系列「音声(oNsee)」のサイド情報ファイル 3 0 2 が図 4 で与えられるとする。図 4 では、各音素(o,N,s,e,e)の開始時刻、終了時刻、及び有声か、無声かを示すフラグがセットされている。尚、「pau」はポーズを示す。このとき、ステップ S 3 0 1 の基本周波数抽出処理では、図 4 の有声区間(o,N,e,e)を検出し、それらの基本周波数を図 5 のように抽出する。次にラベルファイル 3 0 3 が図 6 のように与えられるとき、開始時刻および終了時刻の情報から、有声音の音素区間をステップ S 3 0 2 においてセグメントに分割（この場合は各音素に分割）する。次にステップ S 3 0 3 に進み、例えば、R 次（R = 1：直線）のセグメントモデルによって図 5 に示すピッチパターンの各セグメントピッチパターンをモデル化すると図 7 のように表される。

#### 【 0 0 3 9 】

また、音韻系列「アクセント(akuseNto)」のピッチパターンが図 8 のように抽出されたとする。また、このときのラベルファイル 3 0 3 が図 9 で与えられるとする。このときユニット 1 0 9 は、図 5 及び図 8 に示される 2 つのピッチパターンを用いて、同じ音韻・言語環境に属するセグメントを検出し、それらをモデリングして 1 つのセグメントピッチパターンモデルを生成する。いま、音韻・言語環境として、モーラ位置とアクセント型を選ぶと、「音声」の第 1 モーラの“o”（図 6）及び「アクセント」の第 1 モーラの“a”（図 9）は共にアクセント型が“1”であるため、それらを 1 つのセグメントピッチパターンとしてモデリングする（第 2、第 3、第 4 モーラも同様）。

#### 【 0 0 4 0 】

上述のようにしてモデリングされたセグメントピッチパターンモデルのモデルパラメータを、ピッチパターンモデル 2 0 2 に保持することによって、ステップ S 2 0 3 のプロソディ生成処理では、音韻系列（ $p = \{ p_a, \dots, p_J \}$ ）に対する音韻・言語環境と継続時間長モデル 2 0 3 から得られる各音韻の継続時間長（ $d = \{ dp_1, \dots, dp_J \}$ ）に基づいて、各音韻のセグメントピッチパターン  $Y_{pj}$  を、

$$Y_{pj} = Z dp_j B_{pj} \quad \dots \text{式 ( 1 2 )}$$

として生成することができる。ここで、 $Z dp_j$  は  $dp_j$  フレームのデザインマトリクス、 $B_{pj}$  は音韻  $p_j$  の音韻・言語環境に対応するセグメントピッチパターンモデルのモデルパラメータである。

#### 【 0 0 4 1 】

以上説明したように本実施の形態 1 によれば、セグメントピッチパターン時系列の相関を考慮した多項式セグメントモデルに基づいて、各セグメントピッチパターンをモデリングし、このモデルを用いて所定単位の音韻系列を構成する各音韻のピッチパターンを設定することにより、自然なイントネーションを与える音声合成して出力できるという効果がある。

#### 【 0 0 4 2 】

##### [ 実施の形態 2 ]

上述の実施の形態 1 では、モデル化したセグメントピッチパターンモデルを用いて音声合成する例について説明したが、この実施の形態 2 では、セグメントピッチパターンモデルを用いて音声認識する例について説明する。本実施の形態 2 に係るハードウェア構成は図 1 と同様のものを用いることができる。ここで、入力部 1 0 5 はマイクロフォンである。

#### 【 0 0 4 3 】

図 1 0 は、本発明の実施の形態 2 に係る音声合成・認識ユニット 1 0 9 の動作を示すフローチャートである。以下に示される各ステップは、ROM 1 0 2 に格納された制御プログラムあるいは外部記憶装置 1 0 4 から RAM 1 0 3 にロードされた制御プログラムを CPU 1 0 1 が実行することによって実現される。



## 【 0 0 4 4 】

まずステップ S 4 0 1 で、マイクロフォンなどを備える入力部 1 0 5 から音声波形が入力される。次ステップ S 4 0 2 に進み、その入力された音声波形の特徴パラメータの抽出が行われ、広く用いられているケプストラムなどの周波数特徴量の時系列  $O_a(t)$  に加え、基本周波数やその回帰パラメータなどのピッチに関する特徴量の時系列  $O_p(t)$  を抽出する。

## 【 0 0 4 5 】

次にステップ S 4 0 3 に進み、言語モデル 4 0 1 ( 単語認識の場合は不要 )、上述のセグメントピッチパターンモデルを保持する音響・ピッチパターンモデル 4 0 2、認識辞書 4 0 3 を用いて、ステップ S 4 0 2 で得られた特徴パラメータの尤度を最大とする音声認識結果を探索処理によって求める。次にステップ S 4 0 4 に進み、表示部 1 0 6 への画面表示、或いはスピーカ 1 0 7 による音声出力などの所望の手段によって音声認識結果を出力する。最後にステップ S 4 0 5 において、入力部 1 0 5 からの音声入力終了したか否かの判断を行い、終了していない場合はステップ S 4 0 1 に戻り、前述した処理を実行する。

10

## 【 0 0 4 6 】

いま、音響特徴量ベクトル  $O_a$  に対する単語仮説  $W$  の対数音響尤度を  $P_a(O_a | W)$ 、ピッチ特徴量ベクトル  $O_p$  に対する単語仮説  $W$  の対数ピッチ尤度を  $P_p(O_p | W)$  とし、単語仮説  $W$  の対数言語尤度を  $P_l(W)$  とすると、ステップ S 4 0 3 の探索処理で得られる認識結果  $\sim W$  は、

$$\sim W = \operatorname{argmax}\{w_a P_a(O_a | W) + w_p P_p(O_p | W) + w_l P_l(W)\}$$

20

(  $w$   $W$  )

... 式 ( 1 3 )

で表される。ここで、 $w_a$ ,  $w_p$ ,  $w_l$  は、それぞれ対数音響尤度、対数ピッチ尤度、対数言語尤度に対する重み係数である。ここで、対数音響尤度は HMM ( 隠れマルコフモデル )、対数言語尤度は単語  $n$ -gram に基づく方法など、従来広く用いられている方法によって求めることができる。また、対数ピッチ尤度は、上記式 ( 9 ) を用いて求めることができる。

## 【 0 0 4 7 】

## [ 実施の形態 3 ]

上記実施の形態 1 及び 2 では、上述の式 ( 7 ) に表されるように、セグメントピッチパターンを単一混合ガウス分布によって作成していたが、本実施の形態 3 では、これを多混合ガウス分布によってピッチパターンをモデル化する。

30

## 【 0 0 4 8 】

このとき、 $f(y_t)$  は以下のように表される。

## 【 0 0 4 9 】

## 【 数 1 2 】

$$f(y_t) = \sum_{m=1}^M w_m f_m(y_t). \quad (14)$$

40

## 【 0 0 5 0 】

... 式 ( 1 4 )

ここで、

## 【 0 0 5 1 】

## 【 数 1 3 】

$$f_m(y_t) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_m|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (y_t - z_t B_m)^T \Sigma_m^{-1} (y_t - z_t B_m) \right\}, \quad (15)$$

...式(15)

であり、式(14)における $w_m$ は $m$ 番目の混合分布における重みであり、 $w_m = 1$  ( $m=1 \sim M$ )を満たす。このとき、式(15)におけるモデルパラメータ $B_m$ ,  $\Sigma_m$ ,  $w_m$ はクラスタリング法、もしくはEM(Expectation-Maximization)法によって求めることができる。

10

このようにして得られる多混合ガウス分布によるピッチパターンモデルを用いれば、上記実施の形態2における音声認識装置の性能を向上させることが可能となる。

【0052】

[実施の形態4]

上記実施の形態1では、基本周波数の絶対値から直接セグメントピッチパターンモデルを作成し、このモデルを用いて音声合成におけるピッチパターンの設定を行っていたが、一般にピッチパターンはコンテキストや話者による変動が大きいため、ピッチパターンを抽出する際に、アクセント句、単語、フレーズ(呼気段落)、文などの所望の発話単位(発話もひとまとまりとして処理できる単位)ごとに基本周波数の最大値や最小値などを抽出し、これらの値を利用することによってピッチパターンを正規化し、この正規化されたピッチパターンを用いて、セグメントピッチパターンのモデルを作成するようにしても良い。

20

【0053】

図8に示すピッチパターンを基本周波数の最大値で正規化したときのピッチパターンの例を図11に示す。このように、正規化したピッチパターンからピッチパターンモデルを作成することにより、よりコンテキストなどの変動を大きく吸収した高精度なピッチパターンモデルが作成できる。

【0054】

但し、このピッチパターンモデルを用いて音声合成装置におけるピッチパターンを生成する場合、正規化に用いたパラメータ(図11の場合は最大値)を推定する必要があるが、これは、音韻・言語コンテキストを要因とした線形もしくは非線形モデルなどの公知の方法によって求めることが可能である。

30

【0055】

[実施の形態5]

上記実施の形態では、音素という比較的時間的に短い音韻単位を用いてセグメントピッチパターンをモデル化していたが、本発明はこれに限らず、例えば単語やアクセントといった比較的長い音韻単位に対してモデル化することも可能である。この場合、基本周波数が存在しない無声音の区間をセグメントピッチパターンのモデリングから除外する必要があるが、これは上記式(3)のデザインマトリクスにおいて、次式のように無声音の区間の行を“0”と置くことにより、無声音区間を除外してセグメントピッチパターンをモデル化することができる。

40

【0056】

【数14】

$$\mathbf{Z} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 1 & \frac{1}{L-1} & \cdots & \left(\frac{1}{L-1}\right)^R \\ \vdots & \vdots & & \vdots \\ 1 & \frac{t_s-1}{L-1} & \cdots & \left(\frac{t_s-1}{L-1}\right)^R \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 0 \\ 1 & \frac{t_e-1}{L-1} & \cdots & \left(\frac{t_e-1}{L-1}\right)^R \\ \vdots & \vdots & & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \quad (16)$$

10

【 0 0 5 7 】

...式 ( 1 6 )

このようにして、図 5 に示される 1 単語のピッチパターンをセグメントピッチパターンとして多項式セグメントモデルによってモデリングすることにより、図 1 2 に示されるような、無声区間を含むピッチパターンモデルを得ることができる。

20

【 0 0 5 8 】

なお、上記各実施の形態における構成は本発明の一実施の形態を示したものであり、各種変形が可能である。この変形例を示せば以下の通りである。

【 0 0 5 9 】

実施の形態 1 では、モーラ位置およびアクセント型を音韻・言語環境として考慮してセグメントピッチパターンモデルを作成したが、モーラ数や品詞など他の環境を用いてもよい。また、本発明は日本語以外の言語にも適用可能である。

【 0 0 6 0 】

また前述の実施の形態 1 では、回帰次数 1 (  $R = 1$  ) によってモデリングする例を示したが、 $R$  は 0 以上 ( ただし、 $R < L$  ) の任意の整数値を用いてモデリングしても良い。

30

【 0 0 6 1 】

又前述の実施の形態 2 では、ワンパスの音声認識手法を用いた音声認識装置における例を示したが、従来の音声認識手法を用いて  $N$  ベスト (  $N_{\text{best}}$  ) もしくは単語 ( 音素 ) グラフなどによる認識候補に対して、セグメントピッチパターンモデルによって得られる対数ピッチ尤度を用いて認識結果をリスクアリングする、マルチパス探索の音声認識手法に基づく音声認識装置に対しても適用可能である。

【 0 0 6 2 】

又前述の実施の形態 4 では、基本周波数の最大値によってピッチパターンの正規化処理を行ったが、本発明はこれに限定されるものでなく、例えば最小値を用いた正規化処理や最大値と最小値の差で与えられるダイナミックレンジを用いた正規化処理など他の正規化処理を用いてもよい。

40

【 0 0 6 3 】

また本発明の目的は、前述した実施の形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体を、システムあるいは装置に供給し、そのシステムあるいは装置のコンピュータ ( または CPU や MPU ) が記憶媒体に格納されたプログラムコードを読み出し実行することによっても達成される。

【 0 0 6 4 】

この場合、記憶媒体から読出されたプログラムコード自体が前述した実施の形態の機能を

50

実現することになり、そのプログラムコードを記憶した記憶媒体は本発明を構成することになる。プログラムコードを供給するための記憶媒体としては、例えば、フロッピディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、DVD、磁気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

【0065】

また、コンピュータが読出したプログラムコードを実行することにより、前述した実施の形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているOS（オペレーティングシステム）などが実際の処理の一部または全部を行い、その処理によって前述した実施の形態の機能が実現される場合も含まれる。

【0066】

更に、記憶媒体から読出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行い、その処理によって前述した実施の形態の機能が実現される場合も含まれる。

【0067】

以上説明したように本実施の形態によれば、セグメントピッチパターン系列の相関を考慮して各セグメントピッチパターンを統計的にモデリングすることによって、高精度に所定単位の音韻系列のピッチパターンをモデル化することができるようになり、音声合成装置におけるイントネーション生成の自然性の向上、もしくは基本周波数を特徴量にもつ音声認識装置における認識性能の向上が可能になるという効果がある。

【0068】

【発明の効果】

以上説明したように本発明によれば、所定単位の音韻の基本周波数の時間変化をモデル化することにより、自然なイントネーションを与える音声合成を行うことができる。

【0069】

又本発明によれば、所定単位の音韻の基本周波数の時間変化をモデル化することにより、このモデル化した情報を用いて高精度に音声認識ができるという効果がある。

【図面の簡単な説明】

【図1】本発明の実施の形態に係る音声合成装置（音声認識装置）のハードウェア構成を示したブロック図である。

【図2】本実施の形態に係る音声合成装置における音声合成の処理手順を示したフローチャートである。

【図3】図2のステップS203における多項式セグメントモデルに基づくセグメントピッチパターンモデルの作成手順を示したフローチャートである。

【図4】本発明の実施の形態に係るサイド情報ファイルに記憶された「音声(oNsee)」に関するサイド情報の一例を示す図である。

【図5】本発明の実施の形態に係る「音声」という単語発声に対するピッチパターンの一例を示す図である。

【図6】本発明の実施の形態に係るラベルファイルに記憶された「音声(oNsee)」に関する情報の一例を示す図である。

【図7】図5のピッチパターンを図6に示す音素セグメントごとに回帰次数1のセグメントモデルによってモデリングした場合のピッチパターンの一例を示す図である。

【図8】本発明の実施の形態に係る「アクセント」という単語発声に対するピッチパターンの一例を示す図である。

【図9】本発明の実施の形態に係るラベルファイルに記憶された「アクセント(akuseNto)」に関する情報の一例を示す図である。

【図10】本発明の実施の形態2に係る音声認識装置における音声認識の処理手順を示したフローチャートである。

【図11】本発明の実施の形態4に係る、図8に示すピッチパターンを基本周波数の最大

10

20

30

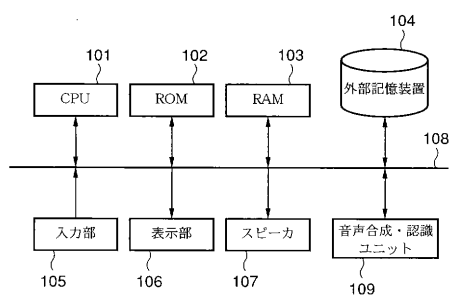
40

50

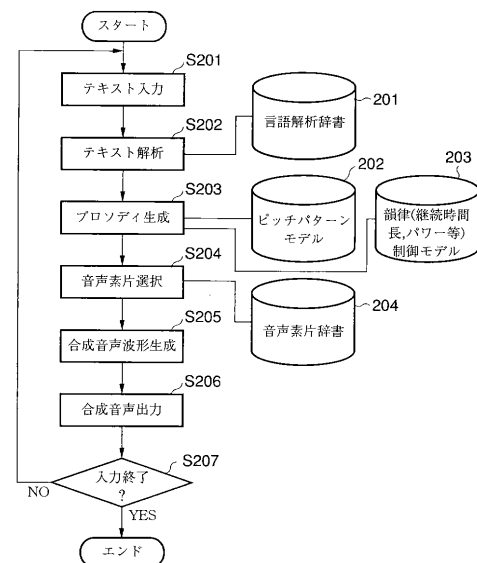
値で正規化したときのピッチパターンの一例を示す図である。

【図 1 2】本発明の実施の形態 5 に係る、図 5 に示すピッチパターンを単語全体の有声音部分のピッチパターンに対して、多項式セグメントモデルによってモデリングした場合のピッチパターンの一例を示す図である。

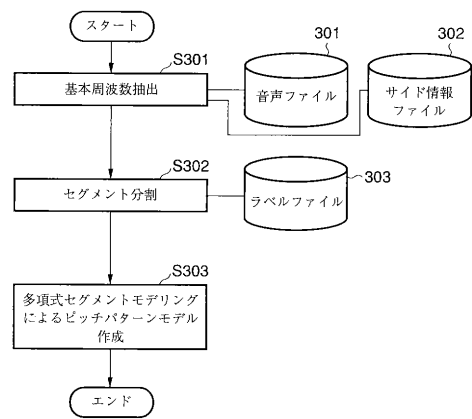
【図 1】



【図 2】



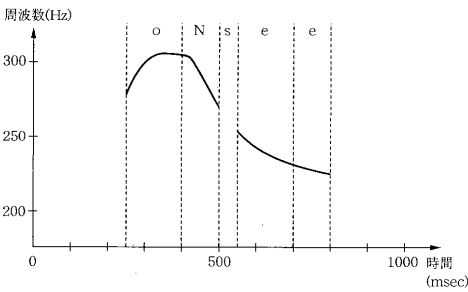
【図 3】



【図 4】

音声ファイル番号	開始時刻 (msec)	終了時刻 (msec)	音素	無声(0)/有声(1)
1	0	250	pau	0
1	250	400	o	1
1	400	500	N	1
1	500	550	s	0
1	550	700	e	1
1	700	800	e	1
1	800	1000	pau	0

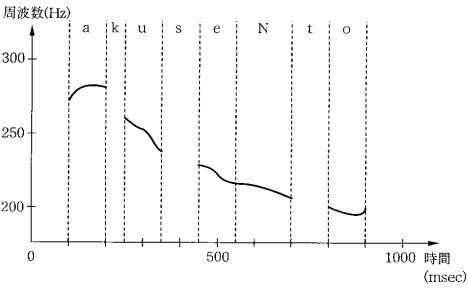
【図 5】



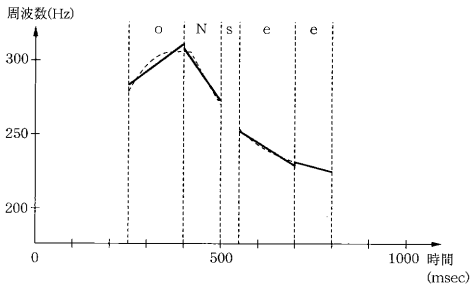
【図 6】

音声ファイル番号	開始時刻 (msec)	終了時刻 (msec)	音素	モーラ位置	アクセント型
1	0	250	pau	-1	-1
1	250	400	o	1	1
1	400	500	N	2	1
1	500	550	s	3	1
1	550	700	e	3	1
1	700	800	e	4	1
1	800	1000	pau	-1	-1

【図 8】



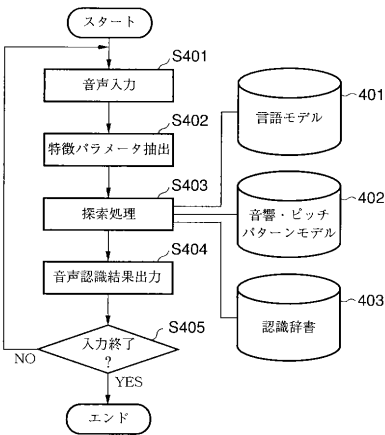
【図 7】



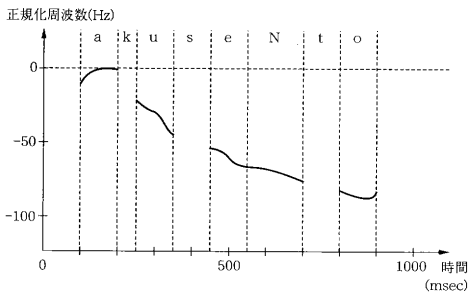
【図 9】

音声ファイル番号	開始時刻 (msec)	終了時刻 (msec)	音素	モーラ位置	アクセント型
2	0	100	pau	-1	-1
2	100	200	a	1	1
2	200	250	k	2	1
2	250	350	u	2	1
2	350	450	s	3	1
2	450	550	e	3	1
2	550	700	N	4	1
2	700	800	t	5	1
2	800	900	o	5	1
2	900	1000	pau	-1	-1

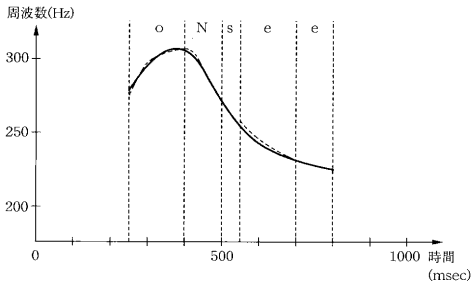
【図 10】



【図 11】



【図 12】



---

フロントページの続き

合議体

審判長 板橋 通孝

審判官 脇岡 剛

審判官 松永 稔

(56)参考文献 特開平 1 1 - 0 9 5 7 8 3 ( J P , A )

特開平 0 2 - 1 9 7 8 9 7 ( J P , A )

特開平 0 8 - 1 2 3 4 6 9 ( J P , A )

特開平 1 0 - 1 4 9 1 8 9 ( J P , A )

深田俊明ほか, HMM統計情報に基づく単語ビッチパターン生成, 日本音響学会平成 6 年度春季  
研究発表会講演論文集 - I -, 日本音響学会, 平成 6 年 3 月 2 3 日, p . 2 2 9 - 2 3 0

深田俊明ほか, 混合分布セグメントモデルのためのモデルパラメータ推定法, 電子情報通信学会  
技術研究報告, 電子情報通信学会, 1 9 9 6 . 0 6 . 1 6 , V o l . 9 6 , N o . 9 3 ( S P 9  
6 2 0 - 3 2 ), p . 3 1 - 3 8

(58)調査した分野(Int.Cl. , D B 名)

G06F 3/16

G10L 3/00

G10L 5/04