

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4793776号
(P4793776)

(45) 発行日 平成23年10月12日(2011.10.12)

(24) 登録日 平成23年8月5日(2011.8.5)

(51) Int.Cl. F I
G 1 0 L 13/08 (2006.01) G 1 0 L 13/08 1 3 0 E
 G 1 0 L 13/08 1 2 7 C

請求項の数 3 外国語出願 (全 17 頁)

<p>(21) 出願番号 特願2005-98067 (P2005-98067) (22) 出願日 平成17年3月30日 (2005. 3. 30) (65) 公開番号 特開2006-276660 (P2006-276660A) (43) 公開日 平成18年10月12日 (2006.10.12) 審査請求日 平成20年2月28日 (2008. 2. 28)</p> <p>(出願人による申告) 平成16年度独立行政法人情報通信研究機構、研究テーマ「大規模コーパスベース音声対話翻訳技術の研究開発」に関する委託研究、産業活力再生特別措置法第30条の適用を受ける特許出願</p>	<p>(73) 特許権者 393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2 (74) 代理人 100099933 弁理士 清水 敏 (72) 発明者 ジンフ・ニ 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内 (72) 発明者 河井 恒 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内</p> <p>審査官 山下 剛史</p> <p style="text-align: right;">最終頁に続く</p>
---	--

(54) 【発明の名称】 イントネーションの変化の特徴を声調の変形により表す方法及びそのコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

イントネーションの変化の特徴を声調の変形により表す方法であって、

話者の、個々の音節から得た語の声調の各々について、基本周波数(F0)ターゲットに関する参考値の所定の組を準備するステップを含み、前記F0ターゲットの参考値の組は、対応する語の声調を特徴づけるものであり、

前記話者のサンプル音声データ中の各音節についてF0ターゲット値を抽出するステップと、

前記サンプル音声データ中の各音節の前記F0ターゲット値の各々について、その音節の語の声調に関する参考値から前記F0ターゲット値への変化の度合いを表す所定の第1のパラメータを計算するステップとをさらに含み、

前記準備するステップは、

語の声調の各々について前記話者による複数個の個々の音節を録音するステップと、

それぞれの語の声調に従って、録音された個々の音節のF0ターゲット値を抽出するステップと、

語の声調の各々について、語の声調を特徴づけるF0ターゲットの各々のF0ターゲット値を平均して前記参考値を求めるステップとを含む、イントネーションの変化の特徴を声調の変形により表す方法。

【請求項2】

所定の第2のパラメータの分布が、当該所定の第2のパラメータの所定の基準値の両側で

つりあうように、前記所定の第1のパラメータを前記所定の第2のパラメータに正規化するステップをさらに含む、請求項1に記載の方法。

【請求項3】

コンピュータ上で実行されると、請求項1又は請求項2に記載の全てのステップを当該コンピュータに行わせる、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は話し言葉の処理に関し、特に、話し言葉でのイントネーションの変化を測定して所望のイントネーションの音声を合成することに関する。 10

【背景技術】

【0002】

中国語の基本周波数(F0)の輪郭(一般的な意味でのイントネーション)は、語の複数の声調、及び平叙文と疑問文との対比を表すような実際のイントネーション(語の声調を除く)を明らかにするものである。伝統的に第一声、第二声、第三声、第四声と呼ばれ(声調1から4)、その各々が他と区別される独自の特徴を持った4つの語の声調と、このような顕著な特徴のない中立声調(声調0)とがある。

【0003】

声調の種類は中国語の音節を直接に構成する要素である。例えば、「ma」は声調の種類によって以下の5つの異なる意味を持つ。 20

【0004】

【数1】

ma 1 (妈 : 母)
 ma 2 (麻 : 麻痺した)
 ma 3 (马 : 馬)
 ma 4 (骂 : 呪う)
 ma 0 (吗 : 疑問の不変化詞)

このために重要な問題が生じる。テキスト-トゥ-スピーチ(text-to-speech: TTS)合成においてイントネーションを合成する際に、語の声調と実際のイントネーションとの相互作用をどのように明らかにするか、ということである。これはTTSを会話システムに適用する際に非常に重要である。会話システムでは例えば、疑問、メッセージの確認、及び感情が、人間によって、通常は音節のイントネーション(すなわち語の声調)と区別され、さらに通常の平叙文とも区別されるイントネーションのパターンで実現される[非特許文献1参照]。 30

【0005】

これに対してとり得る解決策はおそらく、F0輪郭をアクセントと句の成分とに分解するフジサキのモデルであろう[非特許文献2参照]。イントネーションの変化をアクセントと句の成分との両者に分配してもよいが、モデルのパラメータ数は限られている。実際のイントネーションが語の声調に及ぼす影響に対処するため、言語学者は一般に音節[非特許文献3]または句[非特許文献4]のレベルでのピッチ範囲の変化に注目する。 40

【非特許文献1】G. コチャンスキー及びC. シー、「ソフトテンプレートを用いた韻律学モデリング」音声コミュニケーション、第39巻、pp. 311-352、2003年(G. Kochanski and C. Shih, "Prosody modeling with soft templates," Speech Communication, Vol. 39, pp. 311-352, 2003.)

【非特許文献2】H. フジサキ及びK. ヒロセ、「日本語宣言文における音声基本周波数輪郭の分析」日本音響学会誌、第5巻、第4号、pp. 233-242、1984年(H. Fujisaki and K. Hirose, "Analysis of voice fundamental frequency contours for declarative sentences of Japanese," J. Acoust. Soc. Japan, Vol. 50

5, No.4, pp. 233-242, 1984.)

【非特許文献3】J. シェン、「北京方言における声調とイントネーションのピッチ範囲」、実験的音声学における調査報告書、T. リン及びL. J. ワン編、北京大学出版局、pp. 73 - 130、1985年(中国語)(J. Shen, "Pitch range of tone and intonation in Beijing dialect," in Working papers in experimental phonetics, ed. by T. Lin and L. J. Wang, Beijing Univ. Press, pp. 73-130, 1985. (in Chinese))

【非特許文献4】Z. ウー、「標準中国語のためのイントネーション分析の新方法：文中の句輪郭の周波数転位処理」話し言葉の分析、知覚及び処理、G. ファンら編、pp. 255 - 268、1996年(Z. Wu, "A new method of intonation analysis for standard Chinese: frequency transposition processing of phrasal contours in a sentence," Analysis, perception and processing of spoken language, ed. by G. Fant, et al, pp. 255-268, 1996.)

【非特許文献5】Y. R. チャオ、中国語話し言葉の文法。パークレー、カリフォルニア大学出版局、1968年(Y. R. Chao, A grammar of spoken Chinese. Berkeley, University of California Press, 1968.)

【非特許文献6】P. クラトチヴィル、北京語のイントネーション、イントネーションシステム、20ヶ国語の調査内、D. ハースト及びA. D. クリスト編、ケンブリッジ大学出版局、417 - 431、1998年(P. Kratochvil, Intonation in Beijing Chinese, in Intonation systems, a survey of twenty languages, ed. by D. Hirst and A. D. Cristo, Cambridge Uni. Press, 417-431, 1998.)

【非特許文献7】J. ニ及びK. ヒロセ、「標準中国語文の基本周波数輪郭の機能的モデリングの実験的評価」ISCSLP2000、北京、pp. 319 - 322、2000年(J. Ni and K. Hirose, "Experimental evaluation of a functional modeling of fundamental frequency contours of standard Chinese sentences," ISCSLP2000, Beijing, pp. 319-322, 2000.)

【非特許文献8】J. ニ及びH. カワイ、「ピッチ範囲が中国語の声調とイントネーションパターンを固定する」音声韻律学2004、奈良、pp. 95 - 98、2004年(J. Ni and H. Kawai, "Pitch targets anchor Chinese tone and intonation patterns," Speech Prosody 2004, Nara, pp. 95-98, 2004.)

【非特許文献9】J. ニ及びH. カワイ、「パラメトリックモデリング及び合成による分析ベースのパターンマッチングを通じた声調特徴量の抽出」ICASSP2003、pp. 72 - 75、2003年(J. Ni and H. Kawai, "Tone feature extraction through parametric modeling and analysis-by-synthesis-based pattern matching," ICASSP2003, pp. 72-75, 2003)

【非特許文献10】J. ニ及びH. カワイ、「関数モデル及びその評価による中国語基本周波数輪郭の骨格化」TAL2004、pp. 151 - 154、北京、2004年(J. Ni and H. Kawai, "Skeletonising Chinese fundamental frequency contours with a functional model and its evaluation," TAL2004, pp. 151-154, Beijing, 2004.)

【非特許文献11】J. トゥハート、R. コリナー及びC. コーエン、イントネーションの知覚的研究：音声のメロディに対する実験的、音声学的アプローチ、ケンブリッジ大学出版局、1990年(J. 'tHart, R. Collier and A. Cohen, A perceptual study of intonation: an experimental-phonetic approach to speech melody, Cambridge University Press, 1990.)

【発明の開示】

【発明が解決しようとする課題】

【0006】

このようなアプローチの限界は、測定されたピッチ範囲が多少とも語の声調の影響を含んでいることである。さらに、もしある発話中の語の声調がたまたま全て声調1であった

10

20

30

40

50

場合、ピッチ範囲の計算ができなくなる。というのも、声調 1 は高音域レベルの特性を有し、ピッチ範囲を推定するのに基準として利用可能な低音域の特徴がないからである。

【0007】

この発明は、このイントネーションの変化を測定するという問題に別の方向から取組み、分離された個々の音節からの参考値の内部での声調変化を含む、声調の種類への依存性と、F0輪郭の起伏とを分解する際に生じる困難さを避けるようにする。

【0008】

従って、この発明の目的の1つは、自然な条件下で、音声の基にあるイントネーションの変化を測定可能な方法を提供することである。

【0009】

この発明の別の目的は、語の声調に影響されることなく、音声の基にあるイントネーションの変化を測定可能な方法を提供することである。

【課題を解決するための手段】

【0010】

この発明の第1の局面に従えば、イントネーションの種類を声調の変形により特徴づける方法は、話者の個々の音節から得た語の声調の各々について、基本周波数(F0)ターゲットに関する参考値の所定の組を準備するステップを含み、F0ターゲットの参考値の組は対応する語の声調を特徴づけるものであり、話者のサンプル音声データ中の各音節についてF0ターゲット値を抽出するステップと、サンプル音声データ中の各音節のF0ターゲット値の各々について、その音節の語の声調に関する参考値から当該F0ターゲット値への変化の度合いを表す所定の第1のパラメータを計算するステップとをさらに含み、前記準備するステップは、語の声調の各々について前記話者による複数個の個々の音節を録音するステップと、それぞれの語の声調に従って、録音された個々の音節のF0ターゲット値を抽出するステップと、語の声調の各々について、語の声調を特徴づけるF0ターゲットの各々のF0ターゲット値を平均して前記参考値を求めるステップとを含む。

【0012】

より好ましくは、この方法は、所定の第2のパラメータの分布が所定の第2のパラメータの所定の基準値の両側でつりあうように、前記所定の第1のパラメータを所定の第2のパラメータに正規化するステップをさらに含む。

【0013】

この発明の第2の局面は、コンピュータ上で実行されると、上記したいずれかの全てのステップをコンピュータに行わせる、コンピュータプログラムに関する。

【発明を実施するための最良の形態】

【0014】

A. 方法の概観

A. 1 変形

非特許文献7で扱われている、機能モデルで構築された変形は、さまざまな声域でのF0輪郭を時空間と呼ばれる正規化された空間にマッピングすることを可能にする。ここで、 f_0 はヘルツ表示のF0を表すものとし、 λ は(正規化された周波数)でのF0を表すものとする。 f_0 と λ との間の変形は以下の式で表される。

【0015】

【数2】

$$\frac{\ln f_0 - \ln f_{0_b}}{\ln f_{0_t} - \ln f_{0_b}} = \frac{A(\lambda, \xi) - A(\lambda_b, \xi)}{A(\lambda_t, \xi) - A(\lambda_b, \xi)}, \quad (1)$$

ここでA(,)は単純な共振システム内での振幅 - 周波数応答を表す。

【0016】

【数3】

$$A(\lambda, \zeta) = \frac{1}{\sqrt{(1 - (1 - 2\zeta^2)\lambda)^2 + 4\zeta^2(1 - 2\zeta^2)\lambda}}, \lambda \geq 1. \quad (2)$$

は共振システムの減衰比を表す。物理的には、減衰比は共振システム中の粘性抵抗の等価物を表す。他のモデルパラメータは以下を示す。

【0017】

[f_{0b} , f_{0t}]: 声域の最高周波数と最低周波数

[b , t]: で表した声域の最高周波数と最低周波数

声域 [f_{0b} , f_{0t}] は話者に依存する。実際には、対象となる話者の発話の周波数範囲として測定することができる。ほとんどの場合、 t と b とはそれぞれ 1 及び 2 に固定できる。

【0018】

と t が与えられると、 f_0 は上述の変換で直接計算できる。便宜上、 $T_{f_0}(\)$ における b から f_0 への変形を示すものとする。

【0019】

$$f_0 = T_{f_0}(b, t) \quad (3)$$

他方で、 b (又は t) は、 f_0 と t (又は b) が与えられれば、反復処理によって決定することもできる。 $T_{f_0}(\)$ が b での f_0 から t への変形を表すものとする。 f_0 が大きくなるほど、 b で表した値は小さくなる。

【0020】

$$b = T_{f_0}(f_0, t) \quad (4)$$

さらに、 $T_{f_0}(\)$ が b から f_0 への変形のための t を表すものとする。

【0021】

$$t = T_{f_0}(b, f_0) \quad (5)$$

A.2 声調の変形

この変換により、以下の b で示すように、 $[f_{0b}, f_{0t}]$ 内での f_{0_1} から f_{0_2} への変化を測定する方法が提供される。

【0022】

$$b = T_{f_0}(T_{f_0}(f_{0_1}, f_{0_0}), f_{0_2}) \quad (6)$$

ここで f_{0_0} は、 f_{0_1} 及び f_{0_2} をともに b 値にマッピングするときの b の基準値である。好ましくは、 f_{0_0} は 0.156 に固定される。

【0023】

f_{0_1} 及び f_{0_2} 間の一対一のマッピングを保証するために、 b は $(0, 0.7]$ の集合に属していなければならない。これにより、以下の $f_{0_1} = T_{f_0}(b_i, f_{0_0})$ という条件下で図1に見られるように、個々の b_i について、 f_{0_1} 及び f_{0_2} 間での制約が導かれる。

【0024】

$$f_{0_2} = T_{f_0}(T_{f_0}(f_{0_1}, f_{0_0}), f_{0_0}) \quad (7)$$

b_i が基準の f_{0_0} ($= 0.156$) から遠ざかるにつれて、 f_{0_1} は非線形にかつ単調に f_{0_2} へと変化し、その範囲は領域 $[1, 2]$ の両端において急激に狭くなる。

【0025】

b_i を f_{0_0} の両側でつりあわせるため、正規化された減衰比 ζ_n を $\zeta_n \in [-1, 1]$ とし、次のように定義する。

【0026】

【数4】

$$\zeta_n = \begin{cases} (\zeta - \zeta_0) / (0.7 - \zeta_0), & \text{for } \zeta \geq \zeta_0, \\ (\zeta - \zeta_0) / \zeta_0, & \text{for } \zeta < \zeta_0. \end{cases} \quad (8)$$

この方法を拡張して、語の声調及びピッチアクセント等の、2個の F_0 ターゲットのシ

ーケンス間の変化を測定することが可能である。ある声調の中でのすべてのF0ターゲットは、同じ f_0 における f_0 による相対量として表される。この方法を2個の声調間の変化を測定するために用いる利点は、声調内の内部変化が見え、このため、実際の声調の変化を測定可能となることである。

【0027】

図2から図4はこの声調変形をマンダリン語の声調に適用した例を示す。図2(a)は4個の語の声調(ボックス30に示すように、声調1から声調4を同じ時間軸上で重ねたもの)を6回繰返した様子を示し、図2(b)は $f_n = 0$ を示し、これはターゲット声調変化がない、基準となる語の声調を表す。図3(b)に示すように、 f_n が2秒間に0から-1まで線形に変化すると、図2(a)の声調のシーケンスは図3(a)に示すものへと変化する。 f_n は図4(b)の太線に対応し、図2(a)の声調シーケンスは図4(a)に示す太線へと変化する。確かに、声域の非常に高い/低い領域ではピッチ範囲が狭くなる現象が実際の発声でよく見られる。

10

【0028】

A.3 イントネーションの変化測定

音節のイントネーションは声調と呼ばれる。音節と一致する時間-F0輪郭は声調パターンとして知られている。チャオ(Chao)の声調理論[非特許文献5を参照されたい。]に従って、4つの語の声調を4個の声調パターンとして表し、さらにこれを、図5に示すようないくつかの選択されたF0ターゲットにより表す。各声調は主要ターゲットによって特徴づけられる[非特許文献6を参照されたい。]。図5では主要ターゲットを黒丸で示す。

20

【0029】

F0輪郭で明示される声調の変化は、基となる語の声調を特定の態様で変更したものである[非特許文献6を参照]。F0輪郭は、F0ターゲットのシーケンスで信頼性をもって表すことができ、F0ターゲットの数と種類とは、声調パターンに従い、基となる語の声調から決定できる[非特許文献8を参照]。従って、声調変形を用いてF0輪郭から声調の変化を測定するアルゴリズムは、基本的に以下のステップを含む。

【0030】

・初期化：話者による個々の音節から測定された平均のF0ターゲットに従って、4つの声調パターンについてF0ターゲットの基準値(参考値)を決定する。

30

【0031】

・ステップ1：図5の声調パターンに従って、F0輪郭からF0ターゲット(観測値)を抽出する。F0輪郭からF0ターゲットを推定するためのアルゴリズムを、非特許文献9及び10に記載のとおり利用することができ、これによってまず声調特徴を抽出し、その後これをF0ターゲットに変換する。

【0032】

・ステップ2：声調パターンについて対 (f_{0i}, \hat{f}_{0i}) を作成する。ここで、 f_{0i} はi番目のF0ターゲットの観測値を表し、 \hat{f}_{0i} (「f」の前の「^」記号は本来fの上部に表記すべきものである。)はその参考値を表す。声調0については、このF0ターゲットの参考値は、単に先行する声調での最後のF0ターゲットの参考値をとるものとする。

40

【0033】

・ステップ3： $f_i = T(T(\hat{f}_{0i}, f_{00}), f_{0i})$ 、及び f_n を計算する。ただし、 $i = 1, \dots, N$ (F0ターゲットの数)とする。これがイントネーションの変化の特徴を表している。

【0034】

図6は、(a) f_n (丸)により特徴が表されたイントネーションパターンの推定に用いられたF0ターゲット対と、(b)対応する発話データで得られたF0輪郭のためのF0ターゲット対との、参考値(三角)と観測値(丸)とをプロットしている。線P0P4は $f_n = -1.045t + 0.686$ を示し、線P5P7は $f_n = -0.809t + 1.$

50

198を示す。

【0035】

B. 実施例の説明

B. 1 構造

B. 1. 1 機能ブロック

図7はこの発明の一実施例に従った音声合成システム40を示すブロック図である。図7を参照して、音声合成システム40は、所定の話者の基準発話のための記憶装置50と、話者のサンプル発話を記憶するための記憶装置52と、基準発話の声調の各々に対する基準F0ターゲットを抽出し、さらに記憶装置52に記憶されたサンプル発話の各々について、イントネーション変化を示す正規化された減衰比 α_n のシーケンスを抽出するためのイントネーション抽出モジュール54とを含む。

10

【0036】

音声合成システム40はさらに、基準発話の基準F0ターゲットを記憶するための記憶装置56と、 α_n のシーケンスを記憶するための記憶装置58とを含む。減衰比 α_n のシーケンスは、サンプル発話のイントネーション変化の特徴を表すものである。従って、ユーザは、記憶装置58に記憶された α_n のシーケンスを利用して、所望のイントネーションを指定することができる。

【0037】

音声合成システム40はさらに、合成すべき入力テキスト62と関連付けられたイントネーション情報60を受け、入力テキスト62中の音節の各々についてF0を合成するためのF0シンセサイザ64と、入力されたテキスト62とF0シンセサイザ64から出力されたF0とに従って音声信号を合成するための音声シンセサイザ66とを含む。

20

【0038】

イントネーション抽出モジュール54は、記憶装置50内の基準発話の音節の各々からF0ターゲットを抽出し、抽出されたf0ターゲットを記憶装置56に記憶するための第1のターゲット抽出モジュール80と、記憶装置52内のサンプル発話の音節の各々からF0ターゲットを抽出するための第2のターゲット抽出モジュール82と、第2のターゲット抽出モジュール82から出力されたF0ターゲットの各々について、減衰比 α_n を計算し、 α_n のシーケンスを記憶装置58に出力するための α_n 計算モジュール84とを含む。

30

【0039】

F0シンセサイザ64は、イントネーション情報内の α_n のシーケンスから α_n を計算する α_n 計算モジュール90と、以下の式に従って、入力テキスト62の各々の音節の f_{0_i} を計算し、計算された f_{0_i} を音声シンセサイザ66に出力するためのF0計算モジュール90とを含む。

【0040】

$$f_{0_i} = T_{f_0} (T (f_{0_i}, f_{0_0}), \alpha_n) \quad (9)$$

B. 1. 2 コンピュータによる実現

図7に示されたモジュールは、この実施例ではコンピュータソフトウェアで実現される。図8は第1のターゲット抽出モジュール80を実現するコンピュータプログラムの制御構造を示す。図8を参照して、プログラムはステップ100で始まり、基準発話に見出される声調1～声調4の各々について、ステップ102～120が繰返される。

40

【0041】

ステップ102で、変数SUMがゼロに初期化される。

【0042】

ステップ110で、基準発話内の、関心のある声調データの全てについて、ステップ112～116が繰返される。ステップ114で、音節の音声データからF0ターゲットが抽出される。抽出されたF0はステップ116でSUMに加えられる。

【0043】

ステップ112から116が関心のある声調の音節全てに対し繰返された後、ステップ

50

118でSUMの平均を求める。ステップ120で、この平均が、対象の声調と関連付けた上でメモリに記憶される。

【0044】

この処理の終わりには、声調1～声調4の平均F0がメモリに記憶されていることになる。

【0045】

図9は図7に示す第2のターゲット抽出モジュール82及び f_n 計算モジュール84を実現するコンピュータプログラムの制御構造を示す。図9を参照して、ステップ140で、記憶装置52に記憶されたサンプル発話の全てについてF0輪郭が計算される。ステップ142で、入力テキスト62(図7を参照)の全ての音節について、ステップ144から152が繰返される。

10

【0046】

この繰返しでは、まず、処理中の音節の声調のF0ターゲットが抽出される。抽出されたi番目のF0ターゲットを f_{0_i} 、 $1 \leq i \leq N$ (発話中のターゲットの数)とする。

【0047】

ステップ146で、ステップ144で抽出された f_{0_i} が音節の声調パターンの \hat{f}_{0_i} と対にされる。ここで \hat{f}_{0_i} は f_{0_i} の参考値を表す。声調0については、そのF0ターゲットの参考値は単に、先行する声調の最後のF0ターゲットの参考値をとるだけである。

20

【0048】

ステップ148で、 ζ_i が以下の式に従って計算される。

【0049】

$$\zeta_i = T(\hat{f}_{0_i}, f_{0_i}) \quad (10)$$

ステップ150で、正規化された ζ_{n_i} ($1 \leq i \leq N$)が以下の式に従って計算される。

【0050】

【数5】

$$\zeta_{ni} = \begin{cases} (\zeta_i - \zeta_0) / (0.7 - \zeta_0), & \text{for } \zeta_i \geq \zeta_0 \\ (\zeta_i - \zeta_0) / \zeta_0 & \text{for } \zeta_i < \zeta_0 \end{cases} \quad (11)$$

30

ステップ152で、結果 ζ_{n_i} が記憶装置58に記憶される(図7を参照)。

【0051】

記憶装置52に記憶されているサンプル発話の音節全てについて上述の処理を繰返した後、ユーザは正規化された ζ_{n_i} を用いればどのようなイントネーションも記述できる。従って、イントネーション情報60は ζ_{n_i} のシーケンスの形で準備することができる。

【0052】

この実施例では、図7に示すF0シンセサイザ64もまたコンピュータソフトウェアで実現される。このコンピュータプログラムの制御構造を図10に示す。

40

【0053】

図10を参照して、F0シンセサイザ64が起動されると、まずイントネーション情報60内のイントネーションデータ ζ_{n_i} を読み出す。次に、ステップ172で、入力テキスト62の音節全てについてステップ174から178を繰返す。ここで ζ_{n_i} ($1 \leq i \leq N$)はイントネーション情報60の正規化された減衰率のシーケンスとする。

【0054】

ステップ174で、式(11)の逆関数に従って、 ζ_{n_i} から ζ_i を計算する。

【0055】

ステップ176で、i番目の音節(声調)のF0ターゲット f_{0_i} が以下の式に従って計算される。

50

【0056】

$$f_{0i} = T_{f0} (T(\wedge f_{0i}, 0), i) \quad (12)$$

ここで $\wedge f_{0i}$ は基準発話から抽出された参考値(F0ターゲット)を表し、 0 は定数(好ましくは、 0 は0.156)を表す。

【0057】

ステップ178で、このようにして計算された f_{0i} がメモリに記憶される。

【0058】

入力テキスト62の全ての音節について、ステップ174から178が繰返された後、イントネーション情報60によりイントネーションパターンが指定された入力テキスト62中の声調のシーケンスのF0ターゲットとして、 f_{0i} のシーケンスがステップ180で出力される。

10

【0059】

B.1.3 コンピュータハードウェア

図11は上述のコンピュータプログラムを実行するこの実施例のコンピュータシステム330の外観を示し、図12はこのシステム330をブロック図で示す。

【0060】

図11を参照して、このコンピュータシステム330は、FD(フレキシブルディスク)ドライブ352およびCD-ROM(コンパクトディスク読出専用メモリ)ドライブ350を有するコンピュータ340と、キーボード346と、マウス348と、モニタ342と、一対のスピーカ372と、マイクロフォン370と、を含む。

20

【0061】

図12を参照して、コンピュータ340はさらに、CPU(中央処理装置)356と、CPU356、FDドライブ352およびCD-ROMドライブ350に接続されたバス366と、ハードディスク354と、ブートアッププログラム等を記憶する読出専用メモリ(ROM)358と、CPU356に接続され、アプリケーションプログラム命令、システムプログラム、及びデータ等を記憶するランダムアクセスメモリ(RAM)360とを含む。

【0062】

ここでは示さないが、コンピュータ340はさらにローカルエリアネットワーク(LAN)への接続を提供するネットワークアダプタボードを含んでもよい。

30

【0063】

コンピュータシステム330に上述の音声合成システムを実現させるためのコンピュータプログラムは、CD-ROMドライブ350またはFDドライブ352に挿入されるCD-ROM362またはFD364に記憶され、さらにハードディスク354に転送される。または、プログラムは図示しないネットワークを通じてコンピュータ340に送信されハードディスク354に記憶されてもよい。プログラムは実行の際にRAM360にロードされる。CD-ROM362から、FD364から、またはネットワークを介して、直接にRAM360にプログラムをロードしてもよい。

【0064】

図8から図10を参照して説明したこのプログラムは、コンピュータ340にこの実施例の音声合成システム40の機能ブロックを実現させるための複数の命令を含む。この方法を行なわせるのに必要な基本的機能のいくつかはコンピュータ340上で動作するオペレーティングシステム(OS)またはコンピュータ340にインストールされるサードパーティのプログラムにより提供される。従って、このプログラムはこの実施の形態のシステムおよび方法を実現するのに必要な機能全てを必ずしも含まなくてよい。このプログラムは、命令のうち、所望の結果が得られるように制御されたやり方で適切な関数または「ツール」を呼出すことにより、上述の処理を行う命令のみを含んでいてもよい。コンピュータシステム330の動作は周知であるので、ここでは繰返さない。

40

【0065】

B.2 動作

50

この実施例の、上述の音声合成システム40(図7を参照)は以下のように動作する。音声合成システム40の動作は3段階である。すなわち、基準発話からのF0ターゲットの抽出と、基準発話からの f_0 の計算と、F0ターゲット及び音声合成とである。これらの段階における音声合成システム40の動作を以下で説明する。

【0066】

B.2.1 基準発話からのF0ターゲットの抽出

図7を参照して、所定の話者の音声データを、声調1~声調4の全てについて録音し、基準発話として記憶装置50に記憶する。声調1~声調4の各々について、第1のターゲット抽出モジュール80により、基準発話からF0ターゲットが抽出される。声調1~声調4の各々について平均のF0ターゲットが記憶装置56に記憶される。

10

【0067】

B.2.2 基準発話からの f_0 の計算

基準発話と同じ話者のサンプル発話を録音し、記憶装置52に記憶する。サンプル発話の各々の各音節について、第2のターゲット抽出モジュール82がF0ターゲットを抽出する。その後、モジュール82から出力されたF0ターゲットの各々について、 f_0 計算モジュール84が f_0 を計算し、サンプル発話の各々について f_0 のシーケンスを生成する。

【0068】

B.2.3 F0ターゲット及び音声合成

ユーザは、入力テキスト62と、入力テキストをそのイントネーションで合成したいと考えているイントネーションを特定する関連のイントネーション情報60とを準備する。ユーザは、記憶装置58に記憶されている f_0 のシーケンスを調べることにより、イントネーション情報を準備することができる。

20

【0069】

イントネーション情報60と入力テキスト62とが準備されると、入力テキスト62の各音節について、 f_0 計算モジュール90が f_0 を計算し、これをF0計算モジュール92に出力する。例えば、 i 番目の音節に対し、 f_0 計算モジュール90は式(11)の逆関数に従って f_{0i} からこの音節の f_0 を計算する。

【0070】

F0計算モジュール92は、音節の各々に対し、このようにして計算された f_{0i} と、記憶装置56に記憶された f_0 と、定数 $f_0 = 0.156$ とに以下の関数を適用してF0ターゲット f_0 を計算する。

30

【0071】

$$f_0 = T_{f_0} (T (f_{0i}, f_0), f_{0i}) \quad (13)$$

この結果、入力テキスト62内の音節について、F0計算モジュール92により、 f_{0i} のシーケンスが出力される。このシーケンスが音声シンセサイザ66に与えられる。

【0072】

F0計算モジュール92から f_{0i} のシーケンスが与えられると、音声シンセサイザ66は、イントネーション情報60で指定されたイントネーションを備えた入力テキスト62の音声信号68を合成することができる。

40

【0073】

C. 実験結果

ここで提案した方法が、測定されたF0輪郭内の、語の声調よりも高いレベルのイントネーションの変化を明らかにすることが可能であると示すために、2つの実験結果を報告する。音声サンプルは中国語音声コーパスから選択され、専門のナレータに朗読してもらった。ナレータの声域 $[f_{0b}, f_{0t}]$ は $[100\text{Hz}, 500\text{Hz}]$ と一致し、ナレータによる語の声調の参考値は表1に示されるとおりである。太字は主要ターゲットを示す。これらの参考値に対応する声調パターンを図2(a)に見ることができる。

【0074】

【表1】

表1. ナレータの語の声調の参考値

	不完全指定	参考値 (Hz)
声調1	<P11, P12>	<370, 363>
声調2	<V21, P22>	<250, 355>
声調3	<P31, V31, P32>	<235, 200, 280>
声調4	<P41, V41>	<375, 225>

【0075】

【表2】

表2. 図13(a)のイントネーション変化のパラメータ

i	音節	f_{0i}	\hat{f}_{0i}	ζ_i	ζ_{ni}
1	ni3	243	250 (V21)	0.149	-0.045
2	(声調2)*	367	355 (P22)	0.169	0.024
3	hao3	197	235 (P31)	0.115	-0.263
4	(声調3)*	152	200 (V31)	0.090	-0.423
5		204	280 (P32)	0.091	-0.417

*声調3の後に声調3が続くと、前者は声調2に変化する

図13~図16に示される結果は、4つの慣用の挨拶を含むイントネーション変化の分析から得られた。4つの挨拶の実際のイントネーションは音韻論的には同じであるが、語の声調のためにF0輪郭は大きく起伏する。計算の例として、表2は、図13(a)に示されたサンプルからの観測値 f_{0i} 、 $i=1, \dots, 5$ 、対応の参考値 \hat{f}_{0i} 、及び結果として得られるパラメータ ζ_i 及び ζ_{ni} を列挙している。これらの結果は図13(b)に示される。

【0076】

この例では、文のアクセントは、声調2の主要ターゲット（最初の声調3の表面声調）である0.024から第2の声調3の-0.423まで ζ_{ni} が下降したことで示される。他の文の文アクセントもまた、基となる声調の種類に関わりなく一貫して下降するように思われる。この4つの挨拶で示される基本的な特徴は、(1)文のアクセントは発話の最後に位置し、もう1つの音節にかかること、(2)最後の声調（声調1~4）はその参考声調パターンを維持する（すなわち ζ_{ni} が変化しない）ことである。声調0は最後の非声調0である声調の連続したものであるとみなされる。この結果は上述の仮定と一致する。イントネーション変化の現象は、例えば非特許文献11で例示されているように、非声調言語でイントネーションを説明するのに通常用いられるいわゆる「ハットパターン」に非常に類似している。

【0077】

図17は声調及びイントネーションを合成する例を示す。図17(a)は基となる語の声調の参考値を示す。図17(b)は $\zeta_{ni}(t)$ によりイントネーションパターンをプロットする。図17(c)はこれらのF0ターゲット（丸）とこれらのターゲットによりモデルによって与えられる輪郭（連続線）とを示す。「+」のシーケンスはサンプル発話の測定されたF0輪郭を示す。

【0078】

図17から明らかのように、モデルによって与えられるF0輪郭は元のF0輪郭に非常に近い。

【0079】

図18は同じ話者にいくつかの数字列を読んでもらうことで得られたさらなる結果を示す。朗読した数字列は、言語学的意味がないため、中立である。明瞭な結果を求めるため、主要な声調ターゲットの ζ_{ni} 値のみを図にプロットする。加えて、これらの発話では休

10

20

30

40

50

止（ポーズ）がない。イントネーション変化には2つの形状が現れる。1つは最初から最後まで下がる線である（左側）。他方は、下降部とそれに続く平坦部とからなる線である。この下降は最初の2個の音節間で起こる。明らかになったイントネーション変化は、語の声調を越えた高いレベルで体系的である。

【0080】

3人の話者による約200個の中国語サンプルを分析した。これらのサンプルでは実際のイントネーションは多少変化するものの、分析した結果は、この方法により、上で示したとおりイントネーションの変化をはっきりと明らかにできることを示した。

【0081】

D. 結論

この発明の実施の形態は、測定されたF0輪郭から語の声調を除外したイントネーション変化を測定する方法に関する。イントネーション変化は語の声調パターンを構成する選択されたF0ターゲットを用いてサンプリングされ、時間軸上の1点のパラメータで特徴づけられる。実験結果から、この提案した方法が、F0輪郭に埋もれ、語の声調と混じりあった、実際のマンダリン語のイントネーションを分析するのに非常に有望であることがわかった。明らかにされた実際のイントネーションは、非声調言語で報告されたイントネーションとの類似性を示した。提案された方法は基となる語の声調をともなったF0輪郭の自動的な分析を試みるものであり、これは音声合成、認識、さらには理解において決定的に重要である。

【0082】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味および範囲内でのすべての変更を含む。

【図面の簡単な説明】

【0083】

【図1】 n_1 、 n_2 及び n 間の条件を示す図である。

【図2】 声調変形をマンダリン語の声調に適用した例を示す図である。

【図3】 声調変形をマンダリン語の声調に適用した別の例を示す図である。

【図4】 声調変形をマンダリン語の声調に適用した別の例を示す図である。

【図5】 F0ターゲットをベースとしてマンダリン語の声調を表す図である。

【図6】 n （丸）でのイントネーション変化の推定に用いられるF0ターゲット対と、元のF0輪郭のための、参考値（三角）と観測値（丸）とをプロットした図である。

【図7】 この発明の一実施の形態に従った音声合成システム40のブロック図である。

【図8】 第1のF0ターゲット抽出モジュール80を実現するコンピュータプログラムの制御構造をフローチャートで示す図である。

【図9】 第2のターゲット抽出モジュール82と n 計算モジュール84とを実現するコンピュータプログラムの制御構造をフローチャートで示す図である。

【図10】 F0シンセサイザ64を実現するコンピュータプログラムの制御構造をフローチャートで示す図である。

【図11】 一実施の形態に係るコンピュータプログラムを実行するコンピュータシステム330の斜視図である。

【図12】 システム330のブロック図である。

【図13】 慣用の挨拶「 $n i 3 h a o 3$ 」（こんにちは）のF0輪郭を示す図である。

【図14】 慣用の挨拶「 $z e n 3 m e 0 y a n g 4 a 0 ?$ 」（いかがお過ごしですか）のF0輪郭を示す図である。

【図15】 慣用の挨拶「 $n i 3 m a n g 2 m a 0 ?$ 」（お忙しいですか）のF0輪郭を示す図である。

【図16】 慣用の挨拶「 $n i 3 s h e n 1 t i 3 h a o 3 m a 0 ?$ 」（ごきげんいかがですか）のF0輪郭を示す図である。

10

20

30

40

50

【図17】語による韻律の特徴と、語によらない韻律の特徴とを合成する例を示す図である。

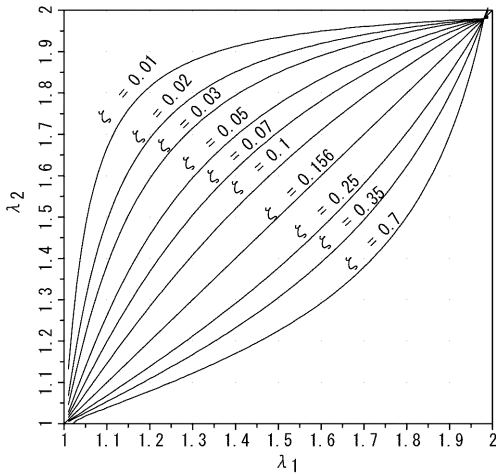
【図18】朗読された数字列での中立イントネーションの変化を示す図である。

【符号の説明】

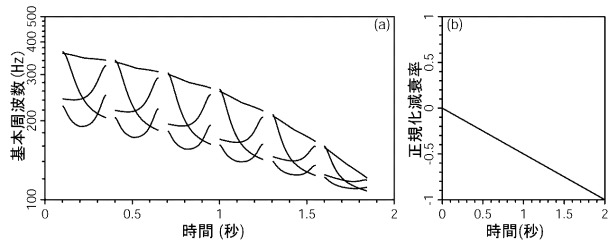
【0084】

- 40 音声合成システム
- 50、52、56、58 記憶装置
- 54 イントネーション抽出モジュール
- 60 イントネーション情報
- 62 入力テキスト
- 64 F0シンセサイザ
- 66 音声シンセサイザ
- 68 イントネーションのある音声信号
- 80 第1のF0ターゲット抽出モジュール
- 82 第2のF0ターゲット抽出モジュール
- 84 ζ_n 計算モジュール
- 90 計算モジュール
- 92 F0計算モジュール

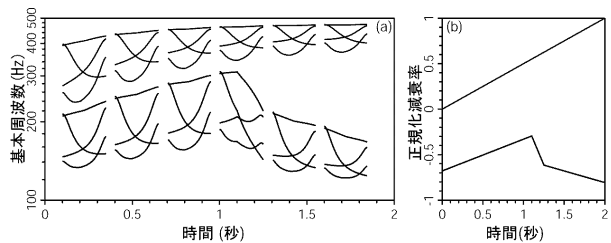
【図1】



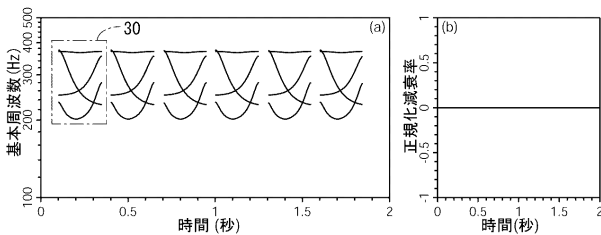
【図3】



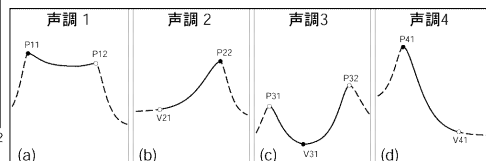
【図4】



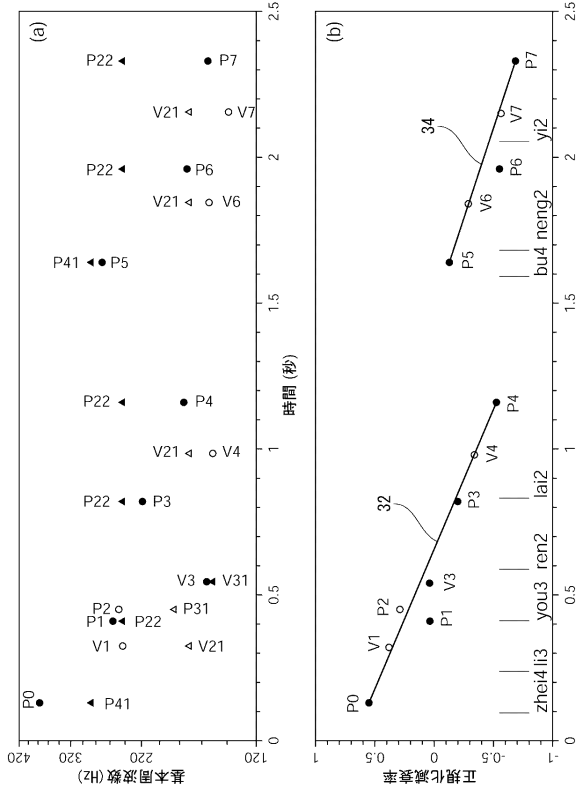
【図2】



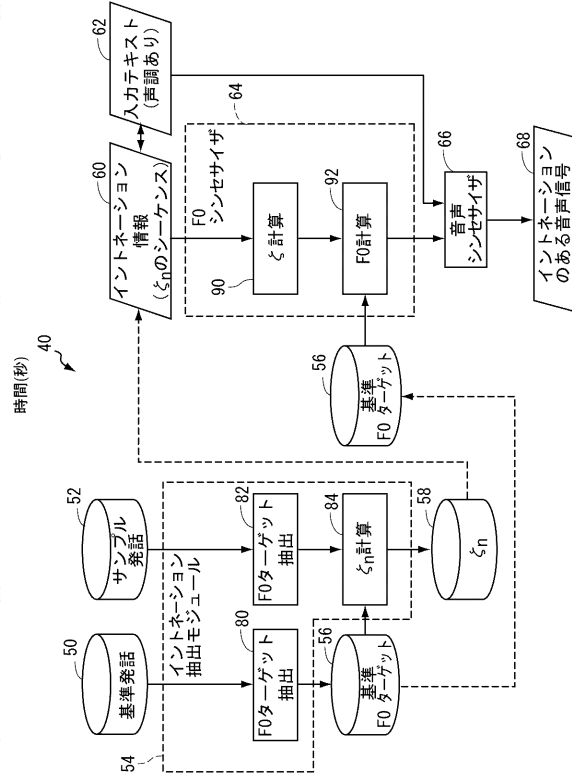
【図5】



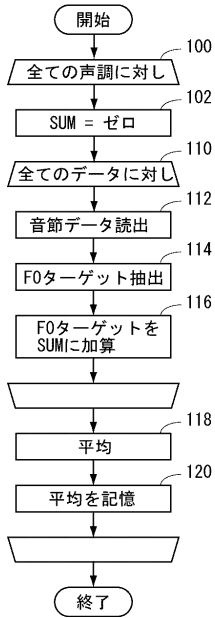
【図6】



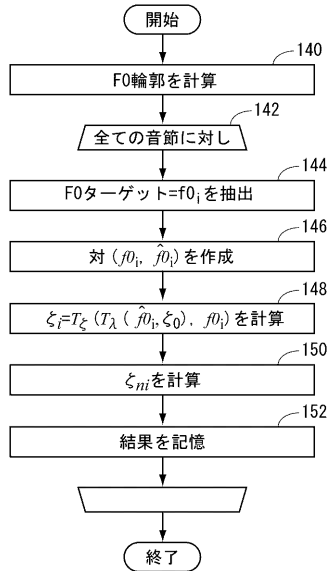
【図7】



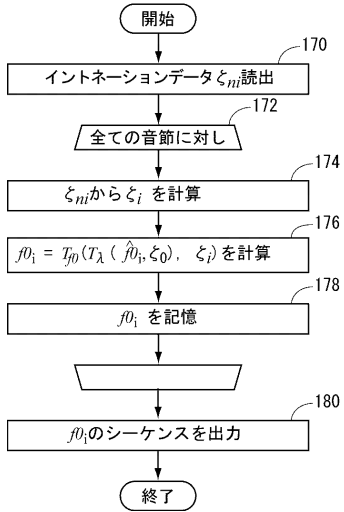
【図8】



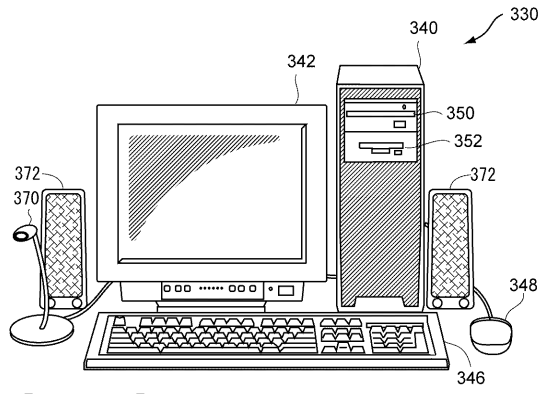
【図9】



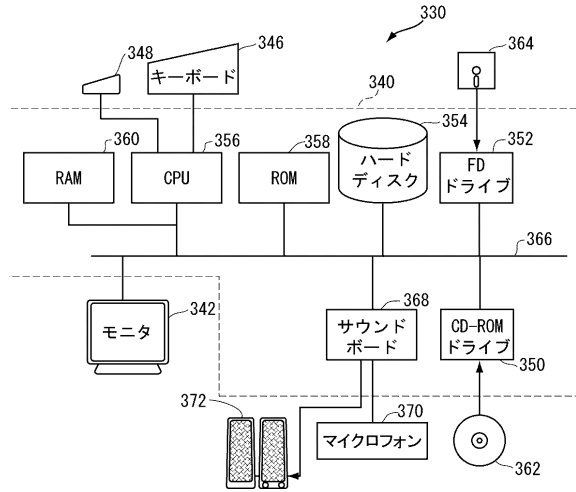
【図10】



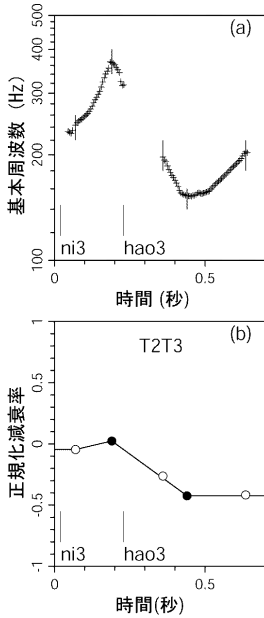
【図11】



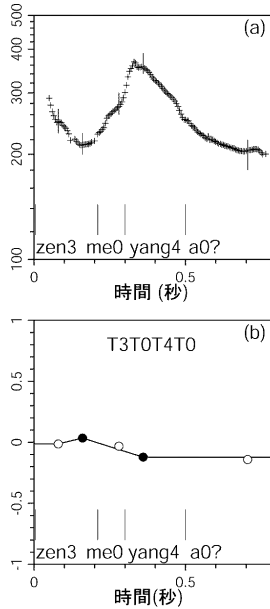
【図12】



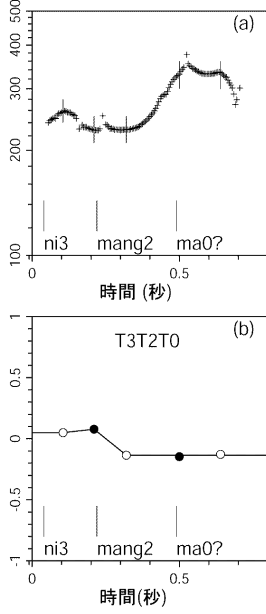
【図13】



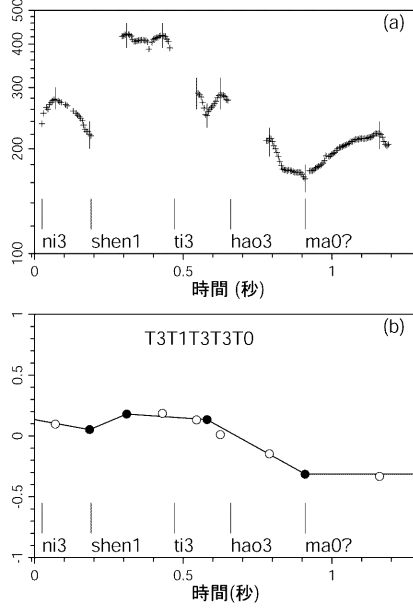
【図14】



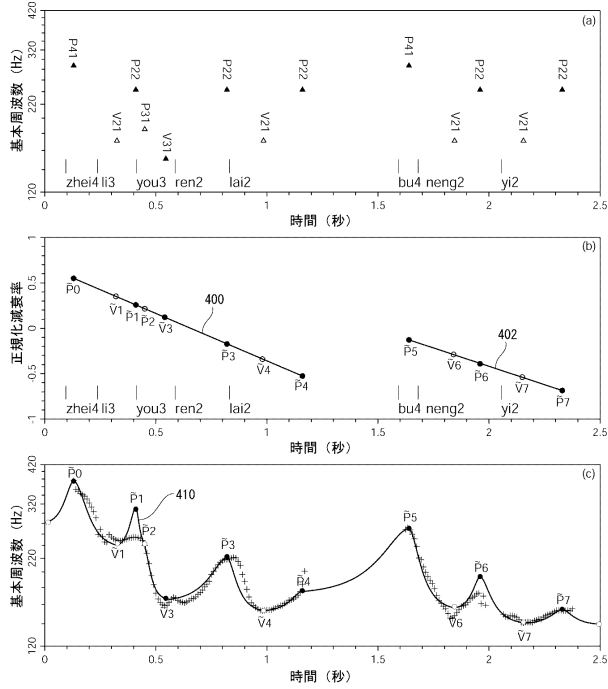
【 図 15 】



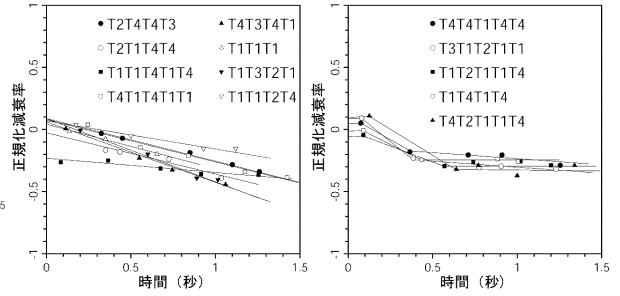
【 図 16 】



【 図 17 】



【 図 18 】



フロントページの続き

- (56)参考文献 特開2003-330482(JP,A)
特開2002-229590(JP,A)
特開2005-265955(JP,A)
特開2005-250264(JP,A)
広瀬啓吉他, "基本周波数パターン生成過程のモデルによる標準中国語音声の韻律的特徴の分析と定式化", 日本音響学会誌, Vol.50, No.3(1994-03), pp.177-187
森大毅他, "単語レベルのF0レンジを考慮した中国語音声の韻律制御", 日本音響学会1999年秋季研究発表会講演論文集-I-, 2-Q-20(1999-09), pp.319-320
徐大威他, "中国語単語のF0レンジの変化に対する許容度に関する検討", 日本音響学会2001年秋季研究発表会講演論文集-I-, 1-2-14(2001-10), pp.233-234
徐大威他, "中国語二音節単語の相対F0変化域の不変性", 電子情報通信学会技術研究報告, Vol.101, No.86, SP2001-13(2001-05), pp.29-34

(58)調査した分野(Int.Cl., DB名)

G10L 11/00 - 13/08