

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4223416号
(P4223416)

(45) 発行日 平成21年2月12日(2009.2.12)

(24) 登録日 平成20年11月28日(2008.11.28)

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

請求項の数 11 外国語出願 (全 20 頁)

<p>(21) 出願番号 特願2004-45855 (P2004-45855) (22) 出願日 平成16年2月23日(2004.2.23) (65) 公開番号 特開2005-234418 (P2005-234418A) (43) 公開日 平成17年9月2日(2005.9.2) 審査請求日 平成17年7月29日(2005.7.29)</p> <p>(出願人による申告)平成15年度通信・放送機構、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 菊池 智紀</p>
--	---

最終頁に続く

(54) 【発明の名称】 F0輪郭を合成する方法及びコンピュータプログラム

(57) 【特許請求の範囲】

【請求項1】

所定の声調言語の言語学的情報から基本周波数 (F_0) 輪郭を合成する、コンピュータにより実現される方法であって、

コンピュータが、声調基本周波数 (F_0) の山及び谷の時刻と周波数とを示す山及び谷パラメータの、前記声調言語における言語学的特徴に対する隠された依存性を、言語学的情報とそれに関連する発話データの F_0 輪郭とを含むトレーニングデータセットから得られる統計的な情報を用いて予測するための機械学習モデルを用いて、入力された言語学的情報に対応する F_0 の山及び谷パラメータを予測するステップと、

コンピュータが、前記予測された F_0 の山及び谷パラメータに、予め定められる関数モデルを適用することにより、前記関数モデルによって前記 F_0 の山及び谷パラメータに対応付けられた F_0 輪郭を推定するステップとを含む、 F_0 輪郭を合成する方法。

【請求項2】

前記予測するステップが、

コンピュータが、前記トレーニングデータセットの入力を受けるステップと、

コンピュータが、前記トレーニングデータセット内の前記発話データの前記 F_0 輪郭から F_0 の山パラメータを抽出するステップと、

コンピュータが、前記 F_0 輪郭と前記抽出するステップで抽出された前記 F_0 の山パラメータとから F_0 の谷パラメータを推定するステップと、

コンピュータが、前記抽出するステップ及び前記 F_0 の谷パラメータを推定するステッ

プでそれぞれ抽出及び推定された前記 F_0 の山及び谷パラメータと、前記トレーニングデータセット内の前記言語学的情報とを用いて、前記機械学習モデルが F_0 の山及び谷パラメータの言語学的情報に対する前記隠された依存性を予測できるように、前記機械学習モデルのパラメータを前記トレーニングデータセットを用いて統計的に算出するステップとを含む、請求項 1 に記載の F_0 輪郭を合成する方法。

【請求項 3】

前記 F_0 の山パラメータを抽出するステップが、
コンピュータが、前記トレーニングデータセット内の各発話の F_0 輪郭を、時間軸に沿って直列に並んだ連続した山型パターンで表される $RONDO - F_0$ 輪郭に変換するステップと、

10

コンピュータが、前記変換するステップで得られた前記 $RONDO - F_0$ 輪郭内における F_0 の山の位置を特定するステップとを含み、

前記 F_0 の谷パラメータを推定するステップが、

コンピュータが、前記変換するステップで得られた前記 $RONDO - F_0$ 輪郭内で、隣接する全ての F_0 の山の間、先行して隣接する F_0 の山からの減衰割合が予め定められた定数となる時点、 F_0 の谷に定めるステップを含む、請求項 2 に記載の F_0 輪郭を合成する方法。

【請求項 4】

前記 F_0 の谷に定めるステップが、

コンピュータが、前記 $RONDO - F_0$ 輪郭内の i 番目の F_0 の山と次の山との間に F_0 の谷 (t_{vfi}, v_{fi}) の初期候補を見出すステップと、

20

コンピュータが、前記初期候補から始めて、前記 F_0 の谷 (t_{vfi}, v_{fi}) が、先行して隣接する F_0 の山からの減衰割合が予め定められた定数となるまで t_{vfi} を所定の時間間隔で減じることにより、前記 $RONDO - F_0$ 輪郭上で F_0 の谷を探索するステップとを含む、請求項 3 に記載の F_0 輪郭を合成する方法。

【請求項 5】

前記初期候補を見出すステップは、コンピュータが、 $(/t_{vi}, /v_i)$ (以下本文中の「/」は上付きバーを示す) で表される最も低い窪みに、初期候補の F_0 の谷 (t_{vfi}, v_{fi}) を設定するステップを含む、請求項 4 に記載の F_0 輪郭を合成する方法。

【請求項 6】

30

前記探索するステップは、コンピュータが、前記初期候補 $(/t_{vi}, /v_i)$ から始めて

$t_{vfi} = t_{pi} - (/v_i - /v_i) \times C$ 、 C は所定の定数、または
 $t_{vfi} = t_{pi}$ となるまで、所定の時間間隔で t_{vfi} を減じることにより、前記 $RONDO - F_0$ 輪郭上で F_0 の谷を探索するステップを含む、請求項 5 に記載の F_0 輪郭を合成する方法。

【請求項 7】

前記定数 C が 0.95 に選ばれる、請求項 6 に記載の F_0 輪郭を合成する方法。

【請求項 8】

前記 F_0 輪郭 $F_0(t)$ 及び前記対応する $RONDO - F_0$ 輪郭 (t) が、以下の、
時間 t の関数

40

【数 1】

$$\frac{\ln F_0(t) - \ln f_{0b}}{\ln f_{0t} - \ln f_{0b}} = \frac{A(\Lambda(t)) - A(\lambda_b)}{A(\lambda_t) - A(\lambda_b)}, \text{ for } t \geq 0,$$

ただし

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

かつ

$$\Lambda(t) = A_{r_1}(t) + \sum_{i=1}^{n-1} \text{Min}(\Lambda_{f_i}(t), \Lambda_{r_{i+1}}(t)) + \Lambda_{f_n}(t)$$

ただし Λ は RONDOSケールでの F_0 周波数、として定義される、請求項 1 ~ 請求項 7 のいずれかに記載の F_0 輪郭を合成する方法。

【請求項 9】

コンピュータが、前記入力された言語学的情報と、前記生成するステップで推定された前記 F_0 輪郭とに基づいて、音声を合成するステップをさらに含む、請求項 1 ~ 請求項 8 のいずれかに記載の F_0 輪郭を合成する方法。

【請求項 10】

所定の声調言語が中国語である、請求項 1 ~ 請求項 9 のいずれかに記載の F_0 輪郭を合成する方法。

【請求項 11】

コンピュータ上で実行されると、コンピュータに請求項 1 ~ 請求項 10 のいずれかに記載のすべてのステップを行なわせる、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は音声合成方法に関し、特に、声調言語の音声合成方法及びシステムに関する。

【背景技術】

【0002】

高品質の合成音声を達成するために、テキストを分析して得られた言語学的情報を信頼性をもって伝えるよう、韻律生成の性能を改善するための努力がなされてきた。韻律生成の最も困難な点は、いかにして音声を適切な声調とイントネーションで合成するか、ということである。この韻律要素は、声調言語では特に重要である。

【0003】

中国語は典型的な声調言語の一つである。中国語の韻律的構成体系 (complex) において、特に基本周波数 (F_0) 輪郭に焦点をあてると、最も小さな特徴的な構成要素は声調である。声調が基本となるのは、韻律的構成体系それ自体を含むより大きな構造が、一つまたは二つ以上の声調を特定の形で変形したものだからである。この観点から、中国語の F_0 輪郭の合成では声調と、文脈によるその変動とに焦点を当てることができる。

【0004】

【非特許文献 1】チャオ、Y. R. 1968、「中国語話し言葉の文法」、バークレー、CA、カリフォルニア大学出版局 (Chao, Y. R., 1968. A Grammar of Spoken Chinese. Berkeley, CA. University of California Press.)

【非特許文献 2】シェン、X. S. 1990、「標準中国語の韻律」、カリフォルニア大学出版局。(Shen, X. S., 1990. The Prosody of Mandarin Chinese. University of California Publications.)

10

20

30

40

50

【非特許文献3】シェン、「漢語語調と語調類型」、中国語文、1994、3、221-228

【0005】

【数1】

(*Hànyǔyǔdiào hé yǔdiào lèixíng. Zhōngguó yǔwén*)

【非特許文献4】シュイ、Y. 1999、「F0輪郭の形成及びアライメントに対する声調と焦点の効果」、音声学ジャーナル、27、55-105。(Xu, Y., 1999. Effects of Tone and Focus on the Formation and Alignment of F0 Contours. Journal of Phonetics, 27, 55-105.)

10

【非特許文献5】ニ、J及びヒロセ、K. 2000、「標準中国語文の基本周波数輪郭の機能的モデリングに対する実験的評価」、ISCSLP2000、北京、319-322。(Ni, J. and Hirose, K., 2000. Experimental Evaluation of a Functional Modeling of Fundamental Frequency Contours of Standard Chinese Sentences. ISCSLP2000. Beijing, 319-322.)

【非特許文献6】ニ、J及びカワイ、H.、「パラメータ的モデリングによる声調特徴の抽出と合成ベース分析によるパターンマッチング」、ICASSP2003、pp.72-75、2003。(Ni, J. and Kawai, H., "Tone Feature Extraction through Parametric Modeling and Analysis-by-Synthesis-based Pattern Matching," ICASSP2003, pp. 72-75, 2003.)

20

【非特許文献7】カワハラ、H.、イクヨ、M. K.、チェイニー、A. 1999、「ピッチ適応時間周波数平滑化及び瞬時周波数ベースのF0抽出を用いた音声表現の再構築：音声における反復構造の果たし得る役割」、音声コミュニケーション、27、187-207。(Kawahara, H., Ikuyo, M. K., Cheneigne, A., 1999. Restructuring Speech Representations Using a Pitch-Adaptive Time-Frequency Smoothing and an Instantaneous-Frequency-Based F0 Extraction: Possible Role of a Repetitive Structure in Sounds. Speech Communication, 27, 187-207.)

【非特許文献8】コロベール、R.、ベンジオ、S.及びマリソン、J.、「トーチ：モジュラー機械学習ソフトウェアライブラリ」、技術報告IDIAP、pp.1-9、2002。(Collobert, R., Bengio, S., and Mariethon, J., "Torch: a Modular Machine Learning Software Library," Technical Report IDIAP, pp. 1-9, 2002.)

30

【発明の開示】

【発明が解決しようとする課題】

【0006】

ピッチターゲットは、中国語の声調及びイントネーションの表出において重要な役割を果たしていると考えられる。ピッチターゲットは基本的には、高低を含むが、これは英語や日本語のようなアクセント言語のイントネーションを表すのに通常用いられるものである。中国語には、第一声から第四声と呼ばれる四声があり、さらに、第0声と呼ばれる中間的な声調がある。もし話者の音声を、1.低、2.半低、3.中間、4.半高、5.高という点数で表した4個の等しい間隔に分割するとすれば、第一声から第四声はそれぞれ、55、35、214、及び51と表される。実際の間隔と絶対ピッチとは共に個々の音声と話すときの気分(mod)とに対し相対的なものであるから、この明細書で用いる「ピッチターゲット」という用語は、時間変化に対するF₀(基本周波数)の山と谷を意味する。

40

【0007】

他方で、声調とイントネーションとのパターンと、F₀輪郭との間には密接な関連が存在する。声調パターンの時間範囲は音節のサイズに限定されるのに対し、イントネーションパターンの時間範囲は音節以上のものをカバーし、発話全体に及ぶ場合もある。

【0008】

50

声調及びイントネーションについては、多くの研究があり、例えば非特許文献1から非特許文献4等の文献がそうである。過去の知覚試験及び器具を用いた分析から、発話の F_0 輪郭は声調とイントネーションを複合的に表し得る、という一致した見解が得られている。

【0009】

しかし、いくつかの基本的な問題について、少なくとも実務的には、明確な解答は得られていない。例えば、中国語の声調とイントネーションとを表すのにピッチターゲットで十分であるのか、またはテキスト音声変換において自然な音を達成するために、声調とイントネーションとの合成に必要な必須の特徴は何か、といったことは明確でない。このため、自然な中国語音声合成する信頼のおける方法はなかった。

10

【0010】

従って、この発明の目的の一つは、高い信頼性をもって自然な音声を合成する方法とコンピュータプログラムとを提供することである。

【課題を解決するための手段】

【0011】

この発明の一つの局面に従えば、所定の声調言語の言語学的情報から基本周波数(F_0)輪郭を合成する方法は、声調基本周波数の山及び谷パラメータの、声調言語の言語学的特徴に対する内部依存性をモデリングするための機械学習モデルを準備するステップと、機械学習モデルを用いて、入力された言語学的情報に対応する F_0 の山及び谷パラメータを予測するステップと、予測された F_0 の山及び谷パラメータに基づいて F_0 輪郭を生成するステップとを含む。

20

【0012】

好ましくは、準備するステップは、言語学的情報とそれに関連する発話データとを含むトレーニングデータセットを準備するステップと、トレーニングデータセット内の発話データから F_0 輪郭モデルパラメータを抽出するステップと、抽出するステップで抽出された F_0 輪郭モデルパラメータから F_0 の山及び谷パラメータを推定するステップと、推定するステップで推定された F_0 の山及び谷パラメータと、トレーニングデータセット内の言語学的情報とを用いて、機械学習モデルが F_0 の山及び谷パラメータの言語学的情報に対する内部依存性を学習するように、機械学習モデルをトレーニングするステップとを含む。

30

【0013】

より好ましくは、 F_0 の山及び谷パラメータを推定するステップが、トレーニングデータセット内の各発話の F_0 輪郭を、時間軸に沿って直列に並んだ連続した山型パターンで表されるRONDO- F_0 輪郭に変換するステップと、変換するステップで得られたRONDO- F_0 輪郭内における F_0 の山を特定するステップと、変換するステップで得られたRONDO- F_0 輪郭内で、隣接する全ての F_0 の山の間に、隣接する F_0 の山と F_0 の谷とが予め定められた条件を満たすように、 F_0 の谷を見出すステップとを含む。

【0014】

見出すステップは、RONDO- F_0 輪郭内の i 番目の F_0 の山と次の山との間に F_0 の谷(t_{vfi} , v_{fi})の初期候補を見出すステップと、初期候補から始めて、 F_0 の谷(t_{vfi} , v_{fi})が予め定められた条件を満たすまで t_{vfi} を所定の時間間隔で減じることにより、RONDO- F_0 輪郭上で F_0 の谷を探索するステップとを含んでも良い。

40

【0015】

好ましくは、初期候補を見出すステップが、(t_{vi} , v_i)で表される最も低い窪みに、初期候補の F_0 の谷(t_{vfi} , v_{fi})を設定するステップを含む。

【0016】

より好ましくは、探索するステップが、初期候補(t_{vi} , v_i)から始めて $v_{fi} - p_i$ ($v_i - p_i$) $\times C$ 、 C は所定の定数、または $t_{vfi} = t_{pi}$ となるまで、所定の時間間隔で t_{vfi} を減じることにより、RONDO- F_0 輪郭上で F_0 の谷を探

50

索するステップを含む。

【 0 0 1 7 】

定数 C がほぼ 0 . 9 5 となるように選択されても良い。

【 0 0 1 8 】

さらに好ましくは、 F_0 輪郭 $F_0(t)$ が、以下の、時間 t の関数

【 0 0 1 9 】

【数 2】

$$\frac{\ln F_0(t) - \ln f_{0b}}{\ln f_{0t} - \ln f_{0b}} = \frac{A(\Lambda(t)) - A(\lambda_b)}{A(\lambda_t) - A(\lambda_b)}, \text{ for } t \geq 0,$$

10

ただし

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

かつ

$$\Lambda(t) = A_{r_1}(t) + \sum_{i=1}^{n-1} \text{Min}(\Lambda_{f_i}(t), \Lambda_{r_{i+1}}(t)) + \Lambda_{f_n}(t)$$

として定義される。

20

【 0 0 2 0 】

好ましくは、この方法は入力された言語学的情報と、生成するステップで生成された F_0 輪郭とに基づいて、音声を合成するステップをさらに含む。

【 0 0 2 1 】

所定の声調言語は中国語であっても良い。

【 0 0 2 2 】

この発明の別の局面は、コンピュータ上で実行されると、コンピュータに上述の方法のいずれかのすべてのステップを行なわせる、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】

【 0 0 2 3 】

1 . はじめに

この発明は中国語の声調及びイントネーションの表示における、ピッチターゲットの役割に焦点をあてたものである。中国語の声調及びイントネーションを表すのにピッチターゲットで十分であるか否かを調べるために、ピッチターゲットを特に時間変化に関する F_0 の山及び谷として測定した。

【 0 0 2 4 】

各々がほぼ同一の声調マッピングで女性の母語話者により平叙文と疑問文として 2 回発話された 7 2 の文に対し、分析及び知覚実験が行なわれた。 F_0 輪郭から観察された声調及びイントネーションのパターンが、関数モデルを用いて定量的に分析され、その後測定されたピッチターゲットから予測されたモデルパラメータを用いて再合成された。二つの認知実験が行なわれた。一方では、予測された声調及びイントネーションのパターンとピッチターゲット及び原文との類似性を評価した。他方では、2 つの平叙文での最終的な声調（第 2 声及び第 4 声）のピッチターゲットを体系的に変化させた場合の人間による声調及びイントネーションの知覚を試験した。

40

【 0 0 2 5 】

実験結果は一貫して、ピッチターゲットが中国語の声調及びイントネーションパターンの規定に重要な役割を果たすことを示した。ピッチターゲットが与えられれば、 F_0 輪郭の正確な形状が予測可能である。この結果に基づき、中国語音声合成方法を構築できる。まず始めに実験について説明し、明細書の後半でこの発明の実施例を説明する。

【 0 0 2 6 】

50

2. 音声試料及び分析方法

2.1. 音声試料

ここで用いられた音声データは、72の中国語文を含み、そのほとんどすべてが非特許文献2から採用されたものである。これらの文を6個のグループに分けた。各々は12の基本文を含み、これをさらに3つのタイプに細分した。各タイプは4つの文を含み、それらは音節数が等しくさらに全文に対し同一の声調のマッピングで特徴付けられる同じ文法構造となっており、これは表1に示すとおりである。表においてT1、T2、T3、T4はそれぞれ、第1声、第2声、第3声及び第4声を示す。

<表1>

【0027】

【表1】

タイプ1	タイプ2	タイプ3
T1 T1 T1 T1	T1 T1 T1 T1 T1	T1 T1 T1 T1 T1 T1 T1 T1 T1
T3 T2 T2 T2	T2 T2 T2 T2 T2	T3 T2 T2 T2 T2 T2 T2 T2 T2
T3 T3 T3 T3	T3 T3 T3 T3 T3	T3 T3 T3 T3 T3 T3 T3 T3 T3
T4 T4 T4 T4	T4 T4 T4 T4 T4	T4 T4 T4 T4 T4 T4 T4 T4 T4

タイプ1は主語 動詞(SV)構造で4つの音節を含む。タイプ2は主語 動詞 目的語(SVO)構造で5個の音節を含む。タイプ3はタイプ1及び2の組合せである。すなわち、タイプ2がタイプ1にその文の目的語として付加され、9音節のSVO構造となっている。これらの文を以下の範疇にグループ分けした。

【0028】

【数3】

S: 平叙文におけるタイプ1、2、3

U: 語彙的にも文法的にもマークされていないイエス/ノーの質問文(以下マーク無し質問文)でのタイプ1、2、3

P: 文の最後の位置に疑問小辞 $ma0$ (吗)のあるイエス/ノーの質問文(マーク付疑問文)におけるタイプ1、2、3

N0: $shi4-bu2-shi4$ (是-不-是)構造のイエス/ノーの質問文におけるタイプ2

N1: $X-mei2-X$ (X-没-X)構造のイエス/ノーの質問文におけるタイプ2

N2: $X-le0mei2-X$ (X-了没-X)構造のイエス/ノーの質問文におけるタイプ2

Q0: $X-hai2shi4-Y$ (X-还是-X)構造の択一的質問文におけるタイプ2

Q1: $shi4-X-hai2shi4-Y$ (是-X-还是-Y)構造の質問文におけるタイプ2

Q2: $hai2shi4-X-hai2shi4-Y$ (还是-X-还是-Y)構造の質問文におけるタイプ2

W0: なぜ($wei4she2me0$) (为什么)の質問文におけるタイプ2

W1: いつ($she2me0shi2hou4$) (什么时候)の質問文におけるタイプ2

W2: 何($she2me0$) (什么)の質問文におけるタイプ2

これら72の文を、女性話者によって感情表現なしで防音室で2回録音した。

【0029】

2.2. F_0 輪郭の関数モデル

この応用では、関数モデルを用いて(非特許文献5を参照)、 F_0 輪郭をパラメータの形で表す。このモデルによれば、話者の声区(発話の周波数区)はまず、いわゆるRONDOスケール(対数スケールと同様)に変換される。その後RONDO- F_0 輪郭を時間軸に直列に並んだ連続した山形状のパターンとして表す。 F_0 輪郭 $F_0(t)$ は以下で与えられる。

【0030】

10

20

30

40

【数4】

$$\frac{\ln F_0(t) - \ln f_{0b}}{\ln f_{0t} - \ln f_{0b}} = \frac{A(\Lambda(t)) - A(\lambda_b)}{A(\lambda_t) - A(\lambda_b)}, \text{ for } t \geq 0, \quad (1)$$

where

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

and

$$\Lambda(t) = \Lambda_{r_1}(t) + \sum_{i=1}^{n-1} \text{Min}(\Lambda_{f_i}(t), \Lambda_{r_{i+1}}(t)) + \Lambda_{f_n}(t). \quad (3)$$

Min(z1, z2) は z1 及び z2 のうち、小さい方を意味する。式(1)及び(2)は合わせて声区の変換を示す。式(3)は R O N D O - F₀ 輪郭 (t) を表し、ここで $\Lambda_{r_i}(t)$ 及び $\Lambda_{f_i}(t)$ はそれぞれ i 番目の山形状パターンの上昇及び下降成分を示す。すなわち

【0031】

【数5】

$$\Lambda_{r_i}(t) = \begin{cases} \lambda_{p_i} + \Delta\lambda_{r_i}(1 - D_{r_i}(t_{p_i} - t)), & \text{for } t \leq t_{p_i}, \\ 0, & \text{otherwise,} \end{cases} \quad (4)$$

$$\Lambda_{f_i}(t) = \begin{cases} \lambda_{p_i} + \Delta\lambda_{f_i}(1 - D_{f_i}(t - t_{p_i})), & \text{for } t \geq t_{p_i}, \\ 0, & \text{otherwise,} \end{cases} \quad (5)$$

$$\text{where } D_{x_i}(t) = (1 + \frac{4.8t}{\Delta t_{x_i}}) e^{-\frac{4.8t}{\Delta t_{x_i}}}, \text{ for } t \geq 0. \quad (6)$$

パラメータ λ_{p_i} 、 $\Delta\lambda_{r_i}$ 及び $\Delta\lambda_{f_i}$ はそれぞれともに 0.237、1 及び 2 に固定され得る。(非特許文献5を参照。)これにより、周波数ドメインに、話者に依存するが発話には依存しない次の2個のパラメータ、

[f_{0b}, f_{0t}]: 声区の最高及び最低周波数、

が得られ、さらに R O N D O - 時間空間内に、発話に依存するが話者に依存しない5個のパラメータ、

n: 山形状パターンの数

t_{x_i}: i 番目の上昇/下降成分の応答時間Δt_{x_i}: i 番目の上昇/下降成分の振幅、x ∈ {r, f}

(t_{p_i}, λ_{p_i}): i 番目の山形状パターンの山(ピーク)、i = 1, ..., n
ができる。

【0032】

2.3.方法

観察された144個のF₀輪郭は最初に、非特許文献6の方法を用いて自動的に分析された。その後、F₀の山と谷とを、もとの声調を考慮しながらF₀輪郭を目で見て調べながらマニュアルで判断した。ある声調に対するF₀の山の数は声調モデリング(非特許文献6)に従って定められた。その後、隣接する山の間の輪郭を用いて、F₀の谷を決定した。モデルにより生成されたF₀輪郭により、これらの発話を再合成し、S T R A I G H T (非特許文献7)と呼ばれるツールを用いて知覚実験を行なった。3つの分析及び知覚

10

20

30

40

50

実験を行なった。実験 1 では、 F_0 の山及び谷に基づき F_0 輪郭の再合成の有効性を分析した。実験 2 では、 F_0 の山及び谷の変化と声調及びイントネーションとの相互作用の相関を調査した。実験 3 はピッチターゲットの変化により声調及びイントネーションが変化し得ることを示す。これらの実験結果に基づき、ピッチターゲットが声調及びイントネーションを規定することについて論じる。

【 0 0 3 3 】

3 . 結果

3 . 1 . 声調及びイントネーションパターンの再合成

実験 1 は、 F_0 の山及び谷に基づき F_0 輪郭の再合成の有効性を調べるために行なわれた。 (t_{v_i}, v_i) が i 番目と $i + 1$ 番目の山の間にある谷を示すこととする。山が与えられると、 F_0 輪郭の生成のために必要な他のモデルパラメータが上述のように計算される。

< 表 2 >

【 0 0 3 4 】

【 表 2 】

	個数	μ_c	σ_c	μ_p	σ_p	μ_e	σ_e
Δt_r	366	0.140	0.047	0.122	0.043	0.022	0.022
$\Delta \lambda_r$	366	0.224	0.147	0.215	0.143	0.013	0.035
Δt_f	382	0.139	0.047	0.134	0.055	0.019	0.027
$\Delta \lambda_f$	382	0.196	0.129	0.188	0.122	0.007	0.015

表 2 は音声試料から測定された声調に関するサンプルの統計的結果を示すものであり、ここで μ_c 及び σ_c はそれぞれマニュアルでチェックされたこれらのモデルパラメータの平均及び分散を示し (チェック済パラメータ)、 μ_p 及び σ_p は F_0 の山及び谷により予測されたもの (予測パラメータ) を示す。 μ_e 及び σ_e の欄はチェック済パラメータと予測パラメータとの間の誤差の平均及び分散を示す。

【 0 0 3 5 】

再合成された声調及びイントネーションパターンと原文との類似性を試験するため、チェック済パラメータを伴う 144 個の再合成発話と、予測パラメータを伴う 144 個の発話とを含む 288 個の刺激対で知覚実験を行なった。刺激は、無音室でヘッドフォンを用いて二人の母語話者に提示された。刺激対を聴いた後、聴者はそれらの声調及びイントネーションの類似度を 3 点スケール、すなわち 0 (非常に異なる)、1 (似ている)、2 (相違なし) で評価した。聴者は判断に先だって刺激を何回も聴くことが許された。チェック済パラメータと予測パラメータとの平均スコアはそれぞれ 1.93 と 1.89 であり、「非常に異なる」サンプルは生じなかった。この実験結果から、ピッチターゲット、すなわち時間変化に対する F_0 の山及び谷は、声調及びイントネーションパターンの特徴を捕捉するのに十分であることが示された。

【 0 0 3 6 】

3 . 2 . 声調とイントネーションとの相互作用

実験 2 では、12 個のカテゴリの各々について、 F_0 の山及び谷の分析により声調とイントネーションとの相互作用を検証した。主な結果を以下に説明する。まず、疑問文における発話の F_0 輪郭は、平叙文におけるそれに比べて、多少とも全体に上向きに動いた。この結果は非特許文献 2 及び非特許文献 3 の知見と一致する。同一の第 1 声及び第 4 声マッピングの発話では、その F_0 の山及び谷は同一の第 2 声及び第 3 声のマッピングのものより高い声区に上昇した。図 1 は平叙文とマーク無しの疑問文で発話された 2 つの文

【 0 0 3 7 】

10

20

30

40

【数2】

“hong2bi2tou2 mei2 quan2” (紅鼻头没权) 及び

“guo4lu4ke4 zhao4 xiang4” (过路客照相)

のF₀輪郭の例を示しており、その意味は、それぞれ、(a)「赤い鼻は権力を持っていない」及び(b)「通りがかりの人が写真を撮る」と、(c)「赤い鼻は権力を持っているか」及び(d)「通りがかりの人が写真を撮るか」である。

【0038】

第2に、マーク無しの疑問文とマーク付きの疑問文の両方でイントネーションを表す、文の最後の声調に依存するやり方がある。最後の声調を2つの組にグループ分けしてみる。第2声と第3声、及び第1声と第4声である。前者では上昇部分のF₀の山が高い声区に上げられ、そのため声調の範囲がかなり広がる。しかしながら後者では、F₀の山と谷(もしあれば)が共に高い声区まで上げられるため、声調範囲は狭くなり、F₀の谷のスケールが上に移動する。この現象は図1(c)と(d)に示される例で明らかに観察される。

10

【0039】

第3に、X - not - X構造のイエス/ノー疑問文、すなわちカテゴリN0、N1、N2と、whの疑問文、すなわちカテゴリW0、W1、W2とでは、上昇 下降パターンが存在する。上昇 下降パターンは基本的には、shi4-bu2-shi4等のように機能語の声調構造にF₀の山と谷とを配することによって表される。図2は4つの例を示す。(a)“ba01shen1gong1 shi4-bu2-shi4 ca1 che1?”(年季奉公の労働者は車を掃除するか?) ; (b)“lao3shou3zhang3 shi4-bu2-shi4 mai3 jiu3?”(年取った高官は酒を買うか?) ; (c)“bao1shen1gong1 ca1-mei2-ca1 che1?”(年季奉公の労働者は車を掃除したか?) ; (d)“lao3shou3zhang3 mai3-mei2-mai3 jiu3?”(年取った高官は酒を買ったか?)である。この図から、機能語に含まれる声調が上昇 下降パターンに適合するように調整されることが明らかである。もし声調がこれと衝突する場合は、声調はその基本的形状を失い、例えば図2(c)の音節ca1及びmei2等のように、上昇 下降パターンの軌跡に従う。

20

【0040】

第4に、択一的疑問文(カテゴリQ0、Q1、Q2)では、機能語 shi4 及び hais hi4 に対し「遷移パターン」が用いられる。遷移パターンという用語は、機能語中の声調がかなり狭いF₀範囲をとり、中間の声区に位置付けられることを意味する。これに対し、機能語周辺の句では通常、焦点現象(focus phenomena)(非特許文献4)が観察される。これらの観察から、声調とイントネーションとの相互作用がピッチターゲットにより良好に捕捉されることが明らかに示される。

30

【0041】

3.3.人工の声調及びイントネーションパターンの知覚

実験3ではピッチターゲットを体系的に変化させながら、声調及びイントネーションの知覚を調査した。図1(a)と1(b)とで示された平叙文の2つの発話をキャリア発話として用い、最終的な声調のF₀の山及び谷を2つ/3つの態様で変化させた。

【0042】

第1の様態は、原文に対しF₀の山(すなわち、モデルパラメータ p_i)をステップサイズ0.1で変化させ、一方F₀の谷は変化無しで固定するというものであった。図3は観察されたF₀輪郭(“+”シーケンス)とこれらのモデル生成輪郭(実線、人工声調及びイントネーションパターンと称する)とを表す。

40

【0043】

第2の様態は、F₀の山及び谷の双方を単純にその声調を同じステップサイズ0.1で上昇させるか下降させることによって変化させるというものであった。モデル生成F₀輪郭は図4で記号B1からB7とD1からD6で示される。

【0044】

50

第3の様態は、特に第4声について、 F_0 の山は固定しながら F_0 の谷をステップサイズ0.1で上に移動させ、谷を上昇させるというものであった。モデル生成 F_0 輪郭は図4で記号E1からE3で示される。2つのキャリア発話はすべてのモデル生成 F_0 輪郭で再合成された。

【0045】

3人の母語話者で、これら発話の知覚試験を行なった。刺激は聴者に2回、ランダムな順序で、無音室でヘッドフォンを通して提示された。刺激を聴いた後、聴者は3つの質問に答えた。

- (1) 発話は平叙文か疑問文か？
- (2) 最後の音節は強調されていたか、普通か、ニュートラルか？
- (3) 最後の音節で聴いたのはどの声調か？

実験結果を表3にまとめた。ここで、“Que”及び“Sta”はそれぞれ「疑問文」と「平叙文」を示し、“Emp”、“Nor”及び“Wea”はそれぞれ「強調」「普通」「弱いストレス」を示す。

<表3>

【0046】

【表 3】

要因	Que(%)	Sta(%)	Emp(%)	Nor(%)	Wea(%)	T1(%)	T2(%)	T3(%)	T4(%)
A1	66.7	33.3	100				100		
A2	66.7	33.3	100				100		
A3	66.7	33.3	66.7	33.3			100		
A4	33.3	66.7	66.7	33.3			100		
A5	16.7	83.3	66.7	33.3			100		
A6		100	33.3	66.7			100		
A7		100		66.7	33.3		66.7	33.3	
B1	100		100			100			
B2	100		83.3	16.7		50	50		
B3	83.3	16.7	100				100		
B4	16.7	83.3	100				100		
B5		100	83.3	16.7			100		
B6		100	16.7	83.3			100		
B7		100	16.7	50	33.3		100		
C1		100	100						100
C2		100	100						100
C3		100	66.7	33.3					100
C4		100	33.3	66.7					100
C5		100	16.7	50	33.3				100
C6		100		66.7	33.3			33.3	66.7
D1	66.7	33.3	100						100
D2		100	100						100
D3		100	66.7	33.3					100
D4		100	33.3	66.7					100
D5		100	16.7	66.7	16.7				100
D6		100		50	50			16.6	83.3
E1	66.7	33.3	83.3	16.7		50			50
E2	66.7	33.3	100						100
E3	66.7	33.3	100						100

この実験から3つの知見を得ることができる。第1に輪郭は、声調、ストレス及びイントネーションを表すことができる[非特許文献1及び2を参照]。F₀の山を上昇させると、音節は一貫して強調されたと知覚された。F₀の山を下げると、音節は弱いストレスで知覚された。第2に、これは疑問イントネーションが最終声調に依存して表わされることを証明した。第2声の場合、F₀の山が高いほど、その発話は容易に疑問文と判断された。第4声の場合、声調が高い声区にある場合のみ、発話は疑問文であると認識された。しかしながら、声調の谷が低い声区にあるときには、聴者は全て、発話を平叙文であると知

10

20

30

40

50

覚した。実験結果はまた、最終的な声調の特徴は疑問文と平叙文を区別するのに十分でないことを示した。テンポ等の他の特徴もまた、知覚の鍵となる。最後に、声調はその F_0 の山、谷及びそれらの音節との整列によって決定される[非特許文献1を参照]。第1声と第3声とはこの実験である条件下で知覚された。加えて、B1、B2及びE1と印をつけた F_0 輪郭が第1声と知覚されたという結果から、第1声は図2(c)に示されるように高い声区で上昇する輪郭を示して、イントネーションを表す必要性を満たすが、その知覚は失われない、という現象を説明する。

【0047】

4. 実施例

中国語の声調及びイントネーションパターンを研究するため、良好に設計された音声試料に対しいくつかの分析と知覚実験とを行なった。実験結果は、声調及びイントネーションパターンの規定においてピッチターゲットが重要な役割を果たすことを示した。例えば関数モデルを用いて、 F_0 の山と谷とから正確な F_0 輪郭を予測することができる。この結果に基づき、観察された F_0 輪郭を、それが伝える主たる言語学的及びパラ言語学的情報を失うことなく、 F_0 の山と谷のシーケンスとして骨格化できると仮定した。以下で説明する実施例はこの思想に基づくものである。

【0048】

4.1. 構造

図5はこの発明の一実施例に従った音声合成システムのブロック図である。図5を参照して、システム20は、トレーニングデータ30から F_0 の山及び谷のデータを抽出する F_0 パラメータ抽出モジュール34を含み、このデータは韻律的特徴と基になる言語学的情報との間を関連付けるために用いられる。システム20はさらに、関連付けされたパラメータの基になる言語学的情報に対する内部依存性を学習するのに用いられる機械学習モジュール36と、言語学的情報32から F_0 輪郭を推定し、適切な声調で中国語音声40を合成するための合成モジュール38とを含む。

【0049】

トレーニングデータ30は、言語学的情報と、付随する F_0 輪郭データとを備えたテキストを含む。

【0050】

F_0 パラメータ抽出モジュール34は、トレーニングデータ30の観察された F_0 輪郭からモデルパラメータの最適な推定を達成するための分析合成 (analysis-by-synthesis: ABS) ベースの分析モジュール50と、 F_0 の谷 ($t_{v_{x_i}}, v_{x_i}$) を推定するためのターゲット探索モジュール52とを含み、ここで F_0 の山は ABS - ベースの分析モジュール50の出力で得られる。

【0051】

ABS - ベースの分析モジュール50は観察された F_0 輪郭からモデルパラメータの最適な推定を達成しようとするものである。非特許文献6は、関数モデルに基づき、 F_0 輪郭から声調の山と下り勾配の特徴とを抽出するためのアルゴリズムを提案している。 F_0 輪郭を、その基礎になる F_0 の山と谷とに信頼性を持って骨格化するために、声調の下り勾配の特徴に関するモデルパラメータ、すなわち、 t_{x_i} 、 x_i 、を再推定するステップでこのアルゴリズムにいくつかの制約を新たに組み入れた。 $(\hat{v}_i, \hat{t}_{v_i})$ (式中では $\hat{\quad}$ は各文字の上に付す) が、 i 番目の山 (p_i, t_{p_i}) と、もしあれば次の山との間に観察された F_0 の谷を示すものとする。谷の決定における F_0 抽出誤差の影響を抑制するため、音声フレームの尤度を考慮する。 t_{x_i} 及び x_i を再推定するための制約を以下に列挙する。

【0052】

10

20

30

40

【数7】

$$\lambda_{p_i} + \Delta\lambda_{f_i} \leq \hat{\lambda}_{v_i} \times 1.1, \quad (7)$$

$$\lambda_{p_i} + \Delta\lambda_{f_i} = \lambda_{p_{i+1}} + \Delta\lambda_{r_{i+1}}, \quad (8)$$

$$\Delta\lambda_{f_i} \geq 0.02, \quad (9)$$

$$\Delta\lambda_{r_{i+1}} \geq 0.02, \quad (10)$$

$$0.05 \leq \Delta t_{f_i} \leq (\hat{t}_{v_i} - t_{p_i}) \times 1.1, \quad (11)$$

$$0.05 \leq \Delta t_{r_{i+1}} \leq (t_{p_{i+1}} - \hat{t}_{v_i}) \times 1.1. \quad (12)$$

ターゲット探索モジュール52はF₀の谷を推定する。2種類のピッチターゲット、F₀の山と谷とが、韻律的特徴と基になる言語学的情報との間の関連付けを行なうための関連付けパラメータとして用いられる。上述の通り、音響分析と知覚試験とから得られた実験結果は一貫して、中国語の声調及びイントネーションパターンを表すのにピッチターゲットで十分であることを示している。ピッチターゲットが与えられれば、関数モデルからF₀輪郭の正確な形状を予測可能である。

【0053】

F₀の山は推定されたモデルパラメータの組から入手可能なので、ターゲットの探索は特に、i番目の上昇または下降成分のいずれかについてF₀の谷(t_{v_xi}, t_{v_xi})に焦点をあて、i番目の山付近のRONDO輪郭に対して行なわれる。谷の初期候補、例えば(t_{v_fi}, t_{v_fi})は、i番目の山と次のものとの間のRONDO輪郭で(t_{v_i}, t_{v_i})で示される最も低い窪みの最初の組である。その後、

【0054】

【数8】

$$\lambda_{v_{f_i}} - \lambda_{p_i} \leq (\bar{\lambda}_{v_i} - \lambda_{p_i}) \times 0.95, \quad (13)$$

またはt_{v_fi} = t_{p_i}となるまで、ごく短いステップ間隔(たとえば0.005秒)t_{v_fi}を減少させることで、RONDO輪郭に沿ってi番目の山に向かって谷の候補を探索する。式(13)において、定数0.95は1から0.05への単位減衰に必要とされる応答時間としてのt_xの定義を考慮して決定された。すなわち

【0055】

【数9】

$$D_x(\Delta t_x) = (1 + \alpha \Delta t_x) e^{-\alpha \Delta t_x} = 1 - 0.95, \quad (14)$$

これは式(6)で用いられる $\alpha = 4.8 / t_x$ という表現に対応する。もしt_{v_fi}とt_{v_ri+1}(ここで「t_{v_ri+1}」の「i+1」は「r」の添え字である。)との差がしきい値より小さい場合、2つの谷は常に2つの谷候補の平均に固定される。

【0056】

機械学習モジュール36は関連付けパラメータの基になる言語学的情報に対する内部依存性を学習するのに用いられる。非特許文献8に記載の通り、いくつかの有効な機械学習方法が利用可能である。

【0057】

合成モジュール38は、機械学習モジュール36で用いられる機械学習方法に関連する回帰アルゴリズムを用いることにより、入力された言語学的情報からピッチターゲットを予測するための機械予測モジュール70と、機械予測モジュール70から与えられるF₀の山及び谷から声調の下り勾配の特徴に関連するモデルパラメータを計算するためのパラメータ変換モジュール72と、モデルパラメータと特定の声区[f_{0b}, f_{0t}]に関数モデルを適用することにより、F₀輪郭を合成するためのモーダルベースの合成モジュール74と、入力された言語学的情報32とモーダルベースの合成モジュール74から与えられるF₀輪郭に基づく適切な声調とに従って音声を作成するための音声合成モジュール76とを含む。これらのモジュールの各々を以下で説明する。

10

20

30

40

50

【 0 0 5 8 】

パラメータ変換モジュール 7 2 は所与の F_0 の山及び谷から声調の下り勾配の特徴に関連するモデルパラメータを計算する。具体的には、以下の通りである。

【 0 0 5 9 】

【 数 1 0 】

$$\Delta t_{r_i} = \max(0.05, t_{p_i} - t_{v_{r_i}}), \quad (15)$$

$$\Delta \lambda_{r_i} = \max(0.02, (\lambda_{v_{r_i}} - \lambda_{p_i}) \times 1.05), \quad (16)$$

$$\Delta t_{f_i} = \max(0.05, t_{v_{f_i}} - t_{p_i}), \quad (17)$$

$$\Delta \lambda_{f_i} = \max(0.02, (\lambda_{v_{f_i}} - \lambda_{p_i}) \times 1.05). \quad (18)$$

10

4 . 2 . 動作

システム 2 0 は以下のように動作する。動作には 3 局面がある。トレーニングデータ 3 0 から F_0 パラメータ抽出モジュール 3 4 により F_0 パラメータを抽出する。 F_0 パラメータ抽出モジュール 3 4 により抽出されたパラメータで機械学習を行なう。その後 トレーニングデータ 3 0 について F_0 輪郭を推定し、推定された F_0 輪郭に従った声調で、言語学的情報 3 2 に基づき中国語音声合成する。

【 0 0 6 0 】

第 1 の局面では、ABS - ベースの分析モジュール 5 0 が F_0 輪郭と言語学的情報とを含む入力されたトレーニングデータ 3 0 を分析し、入力データの声調の山と下り勾配の特徴パラメータとを出力する。

20

【 0 0 6 1 】

ターゲット探索モジュール 5 2 は式 (1 3) 及び (1 4) を用いて F_0 の谷を推定する。推定された F_0 の谷は、第 2 の局面での機械学習のために、 F_0 の山とともに集められる。

【 0 0 6 2 】

第 2 の局面では、集められた F_0 の山と谷とがトレーニングデータ 3 0 内の対応する言語学的情報とともに機械学習モジュール 3 6 に与えられ、機械学習モジュール 3 6 は関連付けられたパラメータの基になる言語学的情報への内部依存性を学習する。この局面が終わると、機械学習モジュール 3 6 を用い、入力された言語学的情報に起こる可能性の高い F_0 の山と谷とを推定することが可能になる。

30

【 0 0 6 3 】

第 3 の局面で、合成モジュール 3 8 は機械学習モジュール 3 6 を用いて言語学的情報から中国語音声 4 0 を合成する。

【 0 0 6 4 】

具体的には、機械予測モジュール 7 0 は、言語学的情報 3 2 が与えられると、機械学習モジュール 3 6 を用いて、入力された言語学的情報に起こる可能性の高いピッチターゲットを予測する。

【 0 0 6 5 】

40

パラメータ変換モジュール 7 2 は式 (1 5) から式 (1 8) により推定された F_0 の山及び谷から声調の下り勾配の特徴に関連するモデルパラメータを計算する。

【 0 0 6 6 】

モーダルベースの合成モジュール 7 4 は、モデルパラメータ及び声区 $[f_{0b}, f_{0t}]$ が与えられると、関数モデルを用いて F_0 輪郭を合成する。

【 0 0 6 7 】

音声合成モジュール 7 6 は言語学的情報 3 2 及びモーダルベースの合成モジュール 7 4 から与えられる F_0 輪郭を分析し、適切な声調で言語学的情報を伝える中国語音声合成する。

【 0 0 6 8 】

50

結果例

図6及び図7は提案されたシステム及び方法を一部示す例であって、F₀輪郭を基となるF₀の山及び谷に骨格化する局面と、F₀の山及び谷から輪郭を復元する局面とを含んでいる。

【0069】

図6(a)に示されるような測定されたF₀輪郭(“+”シーケンス)が与えられるとすると、これらはまず関数モデルに基づき図6(b)に示されるようなパラメータの形で表される。山のパラメータはモデルパラメータの組から利用可能であるが、谷は、ROND O時間空間におけるこれらの山の周囲の上昇/下降成分から探索される。図7の上部と図6(c)は山(黒)と谷(白)とをプロットしており、これによりF₀輪郭の骨格が与えられる。図6(d)は、それによって輪郭が合成されるモデルパラメータにF₀の山と谷とを変換することによって復元されたF₀輪郭を示し、図7下部は観察されたF₀輪郭と回復された輪郭とを実線で示す。

【0070】

言語学的情報から中国語のF₀輪郭と音声とを合成するこの実施例のシステム及び方法は関数モデルを使用する。先行技術と比較して、この実施例ではF₀の山及び谷を、韻律的及び言語学的特徴の間を関連付ける主たるパラメータとして導入し、それらの間の内部依存性を学習するのに機械学習技術を導入する。上述の通り、予備的実験によりこの実施例の有効性が確認された。

【0071】

この発明は中国語に関する実施例について説明されてきたが、この発明は、特徴を表す最も小さい構成要素が声調であるような言語であればどのようなものにも適用可能なことは容易に理解できるであろう。

【0072】

上述のシステムはコンピュータハードウェア及びオペレーティングシステム(OS)上で実行されるソフトウェア(コンピュータプログラム)で実現され得る。音声合成する場合、スピーカ等の音声生成装置を用いることになる。それ以外に特別なハードウェアの必要はない。

【0073】

このプログラムはFD(フレキシブルディスク)、CD-ROM(コンパクトディスク読み専用メモリ)、MO(光磁気ディスク)、またはDVD(デジタル多用途ディスク)等の記憶媒体に記憶されても良く、またはインターネット等の何らかのデータ通信ネットワーク上を送信しても良い。コンピュータ上で実行されると上述のシステムを実現するようなプログラムであれば、どのようなものでもこの発明の範囲内に含まれる。

【0074】

上述の実施の形態は単なる例示であって制限的なものと解してはならない。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及び範囲内でのすべての変更を含む。

【図面の簡単な説明】

【0075】

【図1】平叙文で発話された2つの中国語文“hong2bi2tou2 mei2 quan2”及び“guo4lu4ke4 zhao4 xiang4”のF₀輪郭の例を示す図である。

【図2】4つの例のイエス/ノー疑問文での上昇/下降パターンを示す図である。

【図3】観察されたF₀輪郭(“+”シーケンス)と、ある実験において最後のF₀の山を上向きに移動させた場合にモデルから生成された輪郭とを示す図である。

【図4】ある実験において声調を上向きまたは下向きに移動させて生成したF₀輪郭を示す図である。

【図5】この発明の一実施例の音声合成システムのブロック図である。

【図6】F₀輪郭骨格化及び復元の例を示す図である。

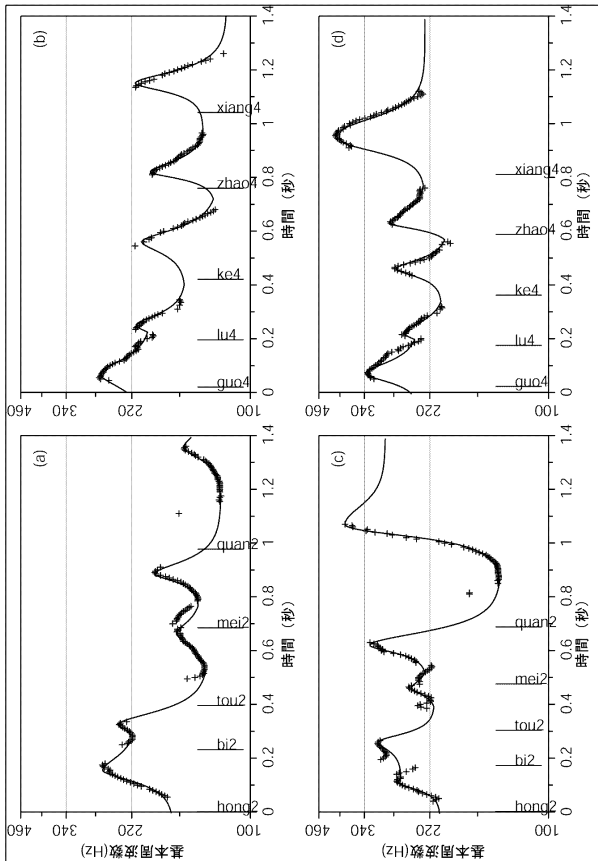
【図7】この発明の一実施例のシステム及び方法を例示する図である。

【符号の説明】

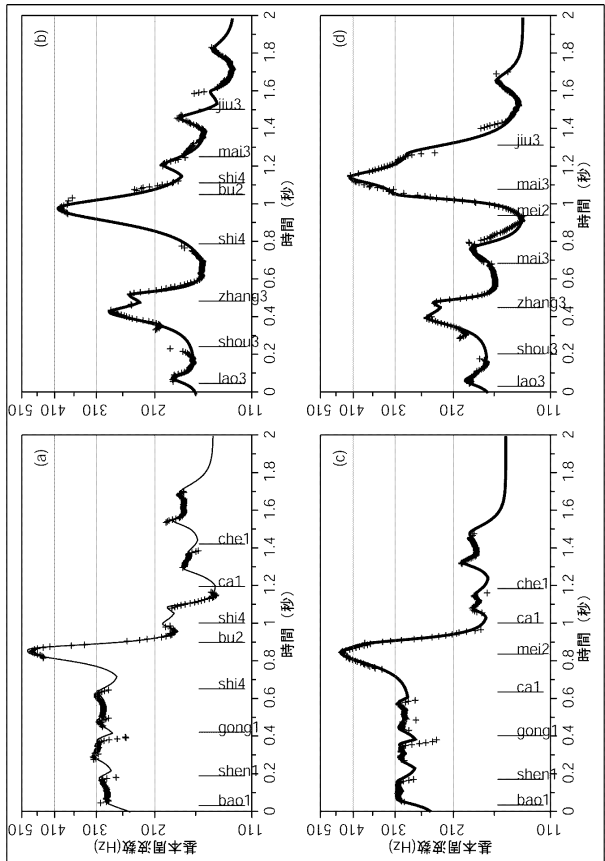
【0076】

20 中国語音声合成システム、30 トレーニングデータ、32 言語学的情報、34 F₀ パラメータ抽出モジュール、36 機械学習モジュール、38 合成モジュール、50 ABS-ベース分析モジュール、52 ターゲット探索モジュール、70 機械予測モジュール、72 パラメータ変換モジュール、74 モーダルベース合成モジュール、76 音声合成モジュール

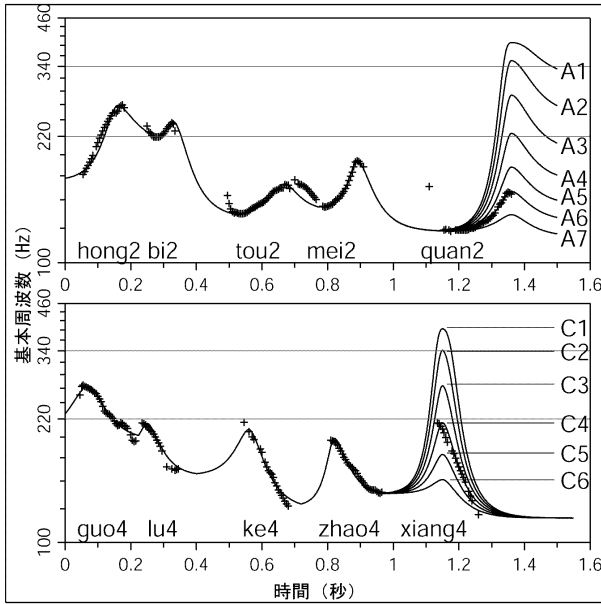
【図1】



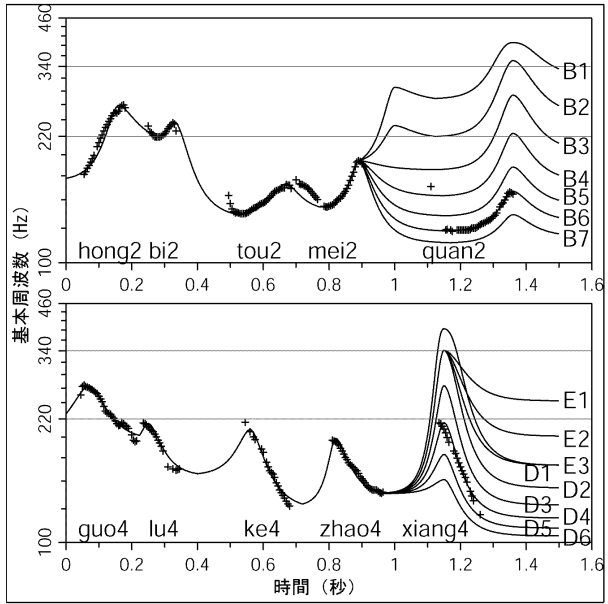
【図2】



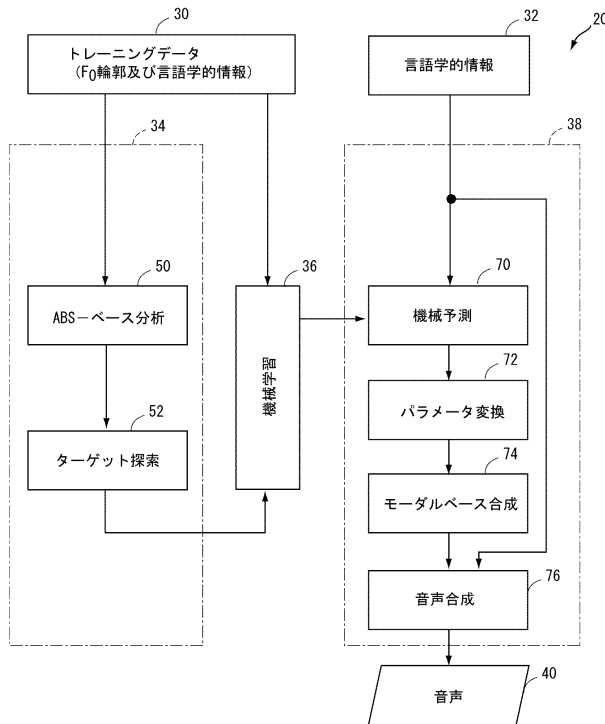
【図3】



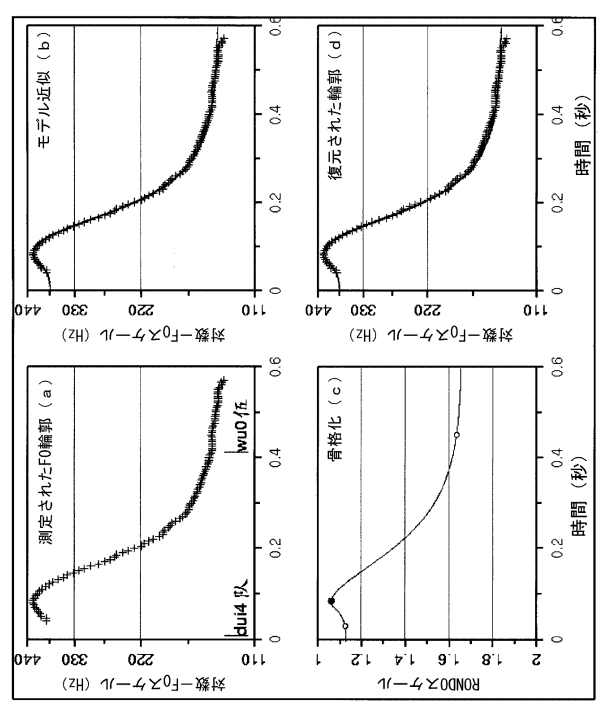
【図4】



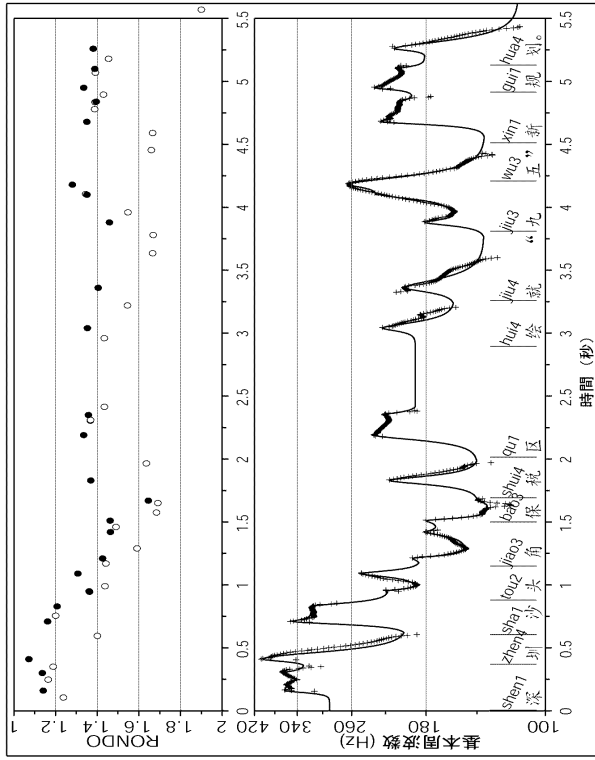
【図5】



【図6】



【 図 7 】



フロントページの続き

(56)参考文献 特開2003-330482(JP,A)

特開昭61-027597(JP,A)

呉亜棟 他, 正接関数を用いた中国語声調認識のための近似モデルの提案, 電子情報通信学会論文誌D-II, 1991年12月25日, Vol.74-D-II, No.12, p.1631-1638

陸金林 他, "中国語音声合成におけるHMMに基づく韻律生成について", 日本音響学会2002年秋季研究発表会講演論文集 - I -, 2002年 9月26日, 2-10-20, p.325-326

陸金林 他, "コーパスベース中国語音声合成システムの韻律制御について", 日本音響学会2001年春季研究発表会講演論文集 - I -, 2001年 3月14日, 2-6-15, p.311-312

豊田悟史 他, "話者の発話特徴を反映したDMPパターンによる音声合成器", 電子情報通信学会技術研究報告, 2002年 1月17日, Vol.101, No.603, p.1-8

(58)調査した分野(Int.Cl., DB名)

G10L 11/00 - 13/08

JSTPlus(JDreamII)