

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4680429号
(P4680429)

(45) 発行日 平成23年5月11日(2011.5.11)

(24) 登録日 平成23年2月10日(2011.2.10)

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

請求項の数 9 (全 36 頁)

<p>(21) 出願番号 特願2001-192778 (P2001-192778) (22) 出願日 平成13年6月26日 (2001.6.26) (65) 公開番号 特開2003-5775 (P2003-5775A) (43) 公開日 平成15年1月8日 (2003.1.8) 審査請求日 平成20年3月3日 (2008.3.3)</p>	<p>(73) 特許権者 308033711 OKIセミコンダクタ株式会社 東京都八王子市東浅川町550番地1 (74) 代理人 100079049 弁理士 中島 淳 (74) 代理人 100084995 弁理士 加藤 和詳 (74) 代理人 100099025 弁理士 福田 浩志 (72) 発明者 茅原 桂一 東京都港区虎ノ門1丁目7番12号 沖電 気工業株式会社内 審査官 毛利 太郎</p>
--	---

最終頁に続く

(54) 【発明の名称】 テキスト音声変換装置における高速読上げ制御方法

(57) 【特許請求の範囲】

【請求項1】

入力されたテキストから音韻・韻律記号列を生成するテキスト解析手段と、前記音韻・韻律記号列に対して少なくとも音声素片・音韻継続時間・基本周波数の合成パラメータを生成するパラメータ生成手段と、音声の基本単位となる音声素片が登録された素片辞書と前記パラメータ生成手段から生成される合成パラメータに基づいて前記素片辞書を参照しながら波形重畳を行って合成波形を生成する波形生成手段とを備えたテキスト音声変換装置における高速読み上げ制御方法であって、前記パラメータ生成手段は、音韻継続時間を予め経験的に求めた継続時間規則テーブルと、音韻継続時間を統計的手法を用いて予測した継続時間予測テーブルとを併せ持ち、ユーザから指定される発声速度が閾値を超えた時には前記継続時間規則テーブルを用い、閾値を超えていない時には前記継続時間予測テーブルを用いて音韻継続時間の決定を行う音韻継続時間決定手段を有することを特徴とするテキスト音声変換装置における高速読み上げ制御方法。

10

【請求項2】

前記閾値は、所定の最大発声速度であることを特徴とする請求項1記載のテキスト音声変換装置における高速読み上げ制御方法。

【請求項3】

入力されたテキストから音韻・韻律記号列を生成するテキスト解析手段と、前記音韻・韻律記号列に対して少なくとも音声素片・音韻継続時間・基本周波数の合成パラメータを生成するパラメータ生成手段と、音声の基本単位となる音声素片が登録された素片辞書と

20

前記パラメータ生成手段から生成される合成パラメータに基づいて前記素片辞書を参照しながら波形重畳を行って合成波形を生成する波形生成手段とを備えたテキスト音声変換装置における高速読み上げ制御方法であって、前記パラメータ生成手段は、アクセント成分及びフレーズ成分を決定するために必要となるデータを、予め経験的に求めた規則テーブルと、統計的手法を用いて予測した予測テーブルとを併せ持ち、ユーザから指定される発声速度が閾値を超えた時には前記規則テーブルを用い、閾値を超えていない時には前記予測テーブルを用いてアクセント成分及びフレーズ成分を決定することによりピッチパターンを決定するピッチパターン決定手段を有することを特徴とするテキスト音声変換装置における高速読み上げ制御方法。

【請求項 4】

前記閾値は、所定の最大発声速度であることを特徴とする請求項 3 記載のテキスト音声変換装置における高速読み上げ制御方法。

【請求項 5】

入力されたテキストから音韻・韻律記号列を生成するテキスト解析手段と、前記音韻・韻律記号列に対して少なくとも音声素片・音韻継続時間・基本周波数の合成パラメータを生成するパラメータ生成手段と、音声の基本単位となる音声素片が登録された素片辞書と前記パラメータ生成手段から生成される合成パラメータに基づいて前記素片辞書を参照しながら波形重畳を行って合成波形を生成する波形生成手段とを備えたテキスト音声変換装置における高速読み上げ制御方法であって、前記パラメータ生成手段は、前記音声素片を変形させて声質を切り換えるための声質変換係数テーブルを備え、ユーザから指定される発声速度が閾値を超えたときには、声質が変化しないような係数を前記声質変換係数テーブルから選択する声質係数決定手段を有することを特徴とするテキスト音声変換装置における高速読み上げ制御方法。

【請求項 6】

前記閾値は、所定の最大発声速度であることを特徴とする請求項 5 記載のテキスト音声変換装置における高速読み上げ制御方法。

【請求項 7】

入力されたテキストから音韻・韻律記号列を生成するテキスト解析手段と、前記音韻・韻律記号列に対して少なくとも音声素片・音韻継続時間・基本周波数の合成パラメータを生成するパラメータ生成手段と、音声の基本単位となる音声素片が登録された素片辞書と前記パラメータ生成手段から生成される合成パラメータに基づいて前記素片辞書を参照しながら波形重畳を行って合成波形を生成する波形生成手段とを備えたテキスト音声変換装置における高速読み上げ制御方法であって、前記パラメータ生成手段は、ユーザが指定した抑揚レベルに応じて修正したピッチパターンを出力するピッチパターン修正手段と、ユーザが指定した発声速度に応じて前記修正したピッチパターンを基底ピッチに加算するか否かを選択する切り換え手段とを有し、前記発声速度が所定の閾値を超えた場合には前記基底ピッチを変更しないように前記切り換え手段を制御することを特徴とするテキスト音声変換装置における高速読み上げ制御方法。

【請求項 8】

前記閾値は、所定の最大発声速度であることを特徴とする請求項 7 記載のテキスト音声変換装置における高速読み上げ制御方法。

【請求項 9】

前記ピッチパターン修正手段は、ユーザが指定した前記発声速度に応じて統計的手法によりフレーズ成分を算出するか或いは当該フレーズ成分を零とする処理を入力文章中に含まれる全フレーズについて行うフレーズ成分算出処理と、ユーザが指定した前記発声速度に応じて統計的手法によりアクセント成分を算出すると共にユーザが指定した前記抑揚レベルに応じて前記算出したアクセント成分を修正するか或いは当該アクセント成分を零とする処理を入力文章中の全ての単語について行う処理とを含むピッチパターン生成処理を行うことを特徴とする請求項 7 記載のテキスト音声変換装置における高速読み上げ制御方法。

【発明の詳細な説明】

10

20

30

40

50

【 0 0 0 1 】

【 発明の属する技術分野 】

本発明は、日常読み書きしている漢字・仮名混じり文を音声として出力するテキスト音声変換技術に係わり、特に高速読上げ時の韻律制御に関するものである。

【 0 0 0 2 】

【 従来技術 】

テキスト音声変換技術は、我々が日常読み書きしている漢字かな混じり文を入力し、それを音声に変換して出力するもので、出力語彙の制限がないことから録音・再生型の音声合成に代わる技術として種々の利用分野での応用が期待できる。

従来、この種の音声合成装置としては、図 1 5 に示すような処理形態となっているものが代表的である。

10

【 0 0 0 3 】

日常読み書きしている漢字仮名混じり文（以下テキストと呼ぶ）を入力すると、テキスト解析部 1 0 1 は、文字情報から音韻・韻律記号列を生成する。ここで、音韻・韻律記号列とは、入力文の読みに加えて、アクセント、イントネーション等の韻律情報を文字列として記述したものの（以下中間言語と呼ぶ）である。単語辞書 1 0 4 は個々の単語の読みやアクセント等が登録された発音辞書で、テキスト解析部 1 0 1 はこの発音辞書を参照しながら、形態素解析ならびに構文解析等の言語処理を施して中間言語を生成する。

【 0 0 0 4 】

テキスト解析部 1 0 1 で生成された中間言語に基づいて、パラメータ生成部 1 0 2 で、音声素片（音の種類）、声質変換係数（声色の種別）、音韻継続時間（音の長さ）、音韻パワー（音の強さ）、基本周波数（声の高さ、以下ピッチと呼ぶ）等の各パターンから成る合成パラメータが決定され、波形生成部 1 0 3 に送られる。

20

【 0 0 0 5 】

ここで音声素片とは、接続して合成波形を作るための音声の基本単位で、音の種類等に応じて様々なものが用意されている。一般的に、C V、V V、V C V、C V C（C：子音、V：母音）といった音韻連鎖で構成されている場合が多い。

【 0 0 0 6 】

パラメータ生成部 1 0 2 で生成された各種パラメータに基づいて、波形生成部 1 0 3 において音声素片等を蓄積する ROM 等から構成された素片辞書 1 0 5 を参照しながら、合成波形が生成され、スピーカを通して合成音声出力される。音声合成方法としては、予め音声波形にピッチマーク（基準点）を付けておき、その位置を中心に切り出して、合成時には合成ピッチ周期に合わせて、ピッチマーク位置をずらしながら重ね合わせる方法が知られている。以上がテキスト音声変換処理の簡単な流れである。

30

【 0 0 0 7 】

次に、パラメータ生成部 1 0 2 における従来処理を図 1 6 を参照して詳細に説明する。

【 0 0 0 8 】

パラメータ生成部 1 0 2 に入力される中間言語は、アクセント位置・ポーズ位置などの韻律情報を含んだ音韻文字列であり、これより、ピッチの時間的な変化（以下ピッチパターン）、音声パワー、それぞれの音韻継続時間、素片辞書内に格納されている音声素片アドレス等の波形を生成する上でのパラメータ（以下、総称して合成パラメータと呼ぶ）を決定する。またこの時、ユーザの好みに合わせた発声様式（発声速度、声の高さ、抑揚の大きさ、声の大きさ、発声話者、声質など）を指定するための制御パラメータも入力される場合がある。

40

【 0 0 0 9 】

入力された中間言語に対して、中間言語解析部 2 0 1 で文字列の解析が行われ、中間言語上に記された呼吸段落記号・単語区切り記号から単語境界を判定し、アクセント記号からアクセント核のモーラ（音節）位置を得る。呼吸段落とは、一息で発声する区間の区切り単位である。アクセント核とは、アクセントが下降する位置のことで、1 モーラ目にアクセント核が存在する単語を 1 型アクセント、n モーラ目にアクセント核が存在する単語を

50

n型アクセントと呼び、総称して起伏型アクセント単語と呼ぶ。逆に、アクセント核の存在しない単語（例えば「新聞」や「パソコン」）を0型アクセントまたは平板型アクセント単語と呼ぶ。これらの韻律に関わる情報は、ピッチパタン決定部202、音韻継続時間決定部203、音韻パワー決定部204、音声素片決定部205、声質係数決定部206に送られる。

【0010】

ピッチパタン決定部202は、中間言語上の韻律情報などからアクセント句あるいはフレーズ単位でのピッチ周波数の時間的変化パタンの算出を行う。従来では「藤崎モデル」と呼ばれる、臨界制動2次線形系で記述されるピッチ制御機構モデルが用いられてきた。声の高さの情報を与える基本周波数は、次のような過程で生成されると考えるのがピッチ制御機構モデルである。声帯振動の周波数、すなわち基本周波数は、フレーズの切り替わりごとに発せられるインパルス指令と、アクセントの上げ下げごとに発せられるステップ指令によって制御される。そのとき、生理機構の遅れ特性により、フレーズのインパルス指令は文頭から文末に向かう緩やかな下降曲線（フレーズ成分）となり、アクセントのステップ指令は局所的な起伏の激しい曲線（アクセント成分）となる。これらの二つの成分は、各指令の臨界制動2次線形系の応答としてモデル化され、対数基本周波数の時間変化パターンは、これら両成分の和（以降、抑揚成分と呼ぶ）として表現される。

【0011】

図18はピッチ制御機構モデルを示す。対数基本周波数 $\ln F_0(t)$ (t は時刻)は、次式のように定式化される。

$$\ln F_0(t) = \ln F_{\min} + \sum_{i=1}^I A_{p_i} G_{p_i}(t-T_{0_i}) + \sum_{j=1}^J A_{a_j} \{G_{a_j}(t-T_{1_j}) - G_{a_j}(t-T_{2_j})\} \dots (1)$$

ここで、 F_{\min} は最低周波数（以下、基底ピッチと呼ぶ）、 I は文中のフレーズ指令の数、 A_{p_i} は文中*i*番目のフレーズ指令の大きさ、 T_{0_i} は文中*i*番目のフレーズ指令の開始時点、 J は文内のアクセント指令の数、 A_{a_j} は文内*j*番目のアクセント指令の大きさ、 T_{1_j} 、 T_{2_j} はそれぞれ*j*番目のアクセント指令の開始時点と終了時点である。

【0012】

また、 $G_{p_i}(t)$ 、 $G_{a_j}(t)$ はそれぞれ、フレーズ制御機構のインパルス応答関数、アクセント制御機構のステップ応答関数であり、次式で与えられる。

$$G_{p_i}(t) = \frac{1}{\omega_i^2} t \exp(-\omega_i t) \dots (2)$$

$$G_{a_j}(t) = \min[1 - (1 + \omega_j t) \exp(-\omega_j t), 1] \dots (3)$$

上式は、 $t \geq 0$ の範囲での応答関数であり、 $t < 0$ では $G_{p_i}(t) = G_{a_j}(t) = 0$ である。式(3)の記号 $\min[x, y]$ は、 x, y のうち小さい方をとることを意味しており、実際の音声でアクセント成分が有限の時間で上限に達することに対応している。ここで、 ω_i は*i*番目のフレーズ指令に対するフレーズ制御機構の固有角周波数であり、例えば3.0などに選ばれる。 ω_j は*j*番目のアクセント指令に対するアクセント制御機構の固有角周波数であり、例えば20.0などに選ばれる。また、 β_j はアクセント成分の上限値であり、例えば0.9などに選ばれる。

【0013】

なおここで、基本周波数およびピッチ制御パラメータ($A_{p_i}, A_{a_j}, T_{0_i}, T_{1_j}, T_{2_j}, \omega_i, \omega_j, F_{\min}$)の値の単位は次のように定義される。すなわち、 $F_0(t)$ および F_{\min} の単位は[Hz]、 T_{0_i}, T_{1_j} および T_{2_j} の単位は[sec]、 ω_i および ω_j の単位は[rad/sec]とする。また A_{p_i} および A_{a_j} の値は、基本周波数およびピッチ制御パラメータの値の単位を上記のように定めたときの値を用

10

20

30

40

50

いる。

【 0 0 1 4 】

以上で述べた生成過程に基づき、ピッチパターン決定部 2 0 2 では、中間言語からピッチ制御パラメータの決定を行う。例えば、フレーズ指令の生起時点 T_{0_i} は中間言語上での句読点が存在する位置に設定し、アクセント指令の開始時点 T_{1_j} は単語境界記号直後に設定し、アクセント指令の終了時点 T_{2_j} はアクセント記号が存在する位置、あるいはアクセント記号がない平板型アクセント単語の場合は、次単語との単語境界記号直前に設定する。フレーズ指令の大きさを表わす A_{p_i} とアクセント指令の大きさを表わす A_{a_j} は、数量化 I 類などの統計的手法を用いて決定する場合が多い。数量化 I 類については公知であるのでここでは特に説明はしない。

10

【 0 0 1 5 】

図 1 9 にピッチパターン生成に関する機能ブロック図を示す。中間言語解析部 2 0 1 からの解析結果が制御要因設定部 5 0 1 に入力される。制御要因設定部 5 0 1 では、フレーズ成分、アクセント成分の大きさを予測するために必要な制御要因の設定を行う。フレーズ成分予測には、例えば、該当するフレーズを構成しているモーラ総数、文内位置、先頭単語のアクセント型といった情報が用いられ、フレーズ成分推定部 5 0 3 に送られる。一方、アクセント成分予測には、例えば、該当するアクセント句のアクセント型、構成しているモーラ総数、品詞、フレーズ内位置といった情報が用いられ、アクセント成分推定部 5 0 2 に送られる。それぞれの成分値予測には、自然発声データを基に数量化 I 類などの統計的手法を用いて予め学習した予測テーブル 5 0 6 を用いて行われる。

20

【 0 0 1 6 】

予測された結果は、ピッチパターン修正部 5 0 4 に送られ、ユーザから抑揚指定があった場合は、推定された値 A_{p_i} 、 A_{a_j} に対しての修正を行う。この機能は、文中のある単語を特に強調あるいは抑制したい時に用いることを想定した制御機構である。通常、抑揚指定は 3 ~ 5 段階に制御され、それぞれのレベルに対してあらかじめ割り当てられた定数を乗ずることにより行われる。抑揚指定がない場合は修正は行われぬ。

【 0 0 1 7 】

フレーズ・アクセント両成分値の修正が施された後、基底ピッチ加算部 5 0 5 に送られ、式 (1) に従ってピッチパタンの時系列データが生成される。この時、ユーザからの声の高さ指定レベルに従って、基底ピッチテーブル 5 0 7 から指定レベルに応じたデータが基底ピッチとして呼び出され加算される。ユーザから特に指定がない場合は、予め定められたデフォルト値が呼び出され加算される。対数化基底ピッチ $\ln F_{m i n}$ は合成音声の最低ピッチを表わしており、このパラメータが声の高さの制御に用いられている。通常 $\ln F_{m i n}$ は、5 ~ 1 0 段階に量子化されてテーブルとして保持されておりユーザの好みによって、全体的に声を高くしたい場合は $\ln F_{m i n}$ を大きくし、逆に声を低くしたい場合は $\ln F_{m i n}$ を小さくするといった処理を行う。

30

【 0 0 1 8 】

基底ピッチテーブル 5 0 7 は、男声音用と女声音用とに分けられており、ユーザから入力される話者指定によって読み出す基底ピッチを選択する。通常男性音の場合は 3 . 0 ~ 4 . 0 の範囲内、女性音の場合は 4 . 0 ~ 5 . 0 の範囲内で声の高さ指定の段階数に応じて量子化されている。以上がピッチパターン生成過程である。

40

【 0 0 1 9 】

次に音韻継続時間制御について述べる。音韻継続時間決定部 2 0 3 は、音韻文字列・韻律記号などからそれぞれの音韻の長さ、休止区間長を決定する。休止区間とは、フレーズ間、あるいは文章間でのポーズの長さである (以後ポーズ長と呼ぶ) 。音韻長は通常、音韻を構成している子音・母音の長さの他、破裂性を有する音韻 (p , t , k など) の直前に現れる無音長 (閉鎖区間長) を、それぞれ決定する。音韻継続時間長、ポーズ長を総称して継続時間長と呼ぶことにする。音韻継続時間の決定方法は通常、目標となる音韻の前後近傍の音韻の種別あるいは、単語内・呼気段落内の音韻位置などにより、数量化 I 類などの統計的手法が用いられる場合が多い。一方、ポーズ長は、前後隣接するフレーズのモー

50

ラ総数などにより同じく、数量化Ⅰ類などの統計的手法が用いられる。またこの時、ユーザから発声速度を指定された場合は、それに応じて音韻継続時間の伸縮を行う。通常、発声速度指定は、5～10段階程度に制御され、それぞれのレベルに対してあらかじめ割り当てられた定数を乗ずることにより行われる。発声速度を遅くしたい場合は音韻継続時間を長くし、発声速度を速くしたい場合は音韻継続時間を短くする。音韻継続時間制御に関しては、本発明の主題であるので後述する。

【0020】

音韻パワー決定部204は、音韻文字列からそれぞれの音韻の波形振幅値の算出を行う。波形振幅値は、/a, i, u, e, o/などの音韻の種類・呼気段落内での音節位置などから経験的に決められる。また、音節内においても、立ち上がりの徐々に振幅値が大きくなる区間と、定常状態にある区間と、立ち下りの徐々に振幅値が小さくなる区間のパワー遷移も同時に決定している。これらパワー制御は通常、テーブル化された係数値を用いることにより実行される。またこの時、ユーザからの声の大きさ指定があった場合は、それに応じて振幅値を増減する。通常、声の大きさ指定は、10段階程度に制御され、それぞれのレベルに対してあらかじめ割り当てられた定数を乗ずることにより行われる。

10

【0021】

音声素片決定部205は、音韻文字列を表現するために必要な音声素片の、素片辞書105内アドレスの決定を行う。素片辞書105は、例えば男声音と女性音といった具合に複数話者の音声素片が格納されており、ユーザからの話者指定により素片アドレスの決定を行う。素片辞書105に格納されている音声素片データは、CV、VCVなど前後の音韻環境に応じた形で様々な単位で構築されているため、入力テキストの音韻文字列の並びから最適な合成単位を選択する。

20

【0022】

声質係数決定部206は、ユーザから声質変換指定があった場合に、変換パラメータの決定を行う。声質変換とは、素片辞書105に登録されている素片データに、信号処理等の加工を施すことにより、聴感上、別話者として取り扱えるようにした機能である。一般に、素片データを線形に伸縮する処理を施して実現する場合が多い。伸長処理は、素片データのオーバーサンプリング処理で実現され、太い声となる。逆に縮小処理は、素片データのダウンサンプリング処理で実現され、細い声となる。通常、声質変換指定は、5～10段階程度に制御され、それぞれのレベルに対してあらかじめ割り当てられたリサンプリング・レートにより変換を行う。

30

【0023】

以上の処理により生成されたピッチパターン・音韻パワー・音韻継続時間・音声素片アドレス・伸縮パラメータは合成パラメータ生成部207に送られ、合成パラメータが生成される。合成パラメータは、フレーム(通常8ms程度の長さ)を一つの単位とした波形生成用のパラメータであり、波形生成部103に送られる。

【0024】

図17に波形生成部の機能ブロック図を示す。素片復号部301では、合成パラメータのうち、素片アドレスを参照ポイントとして素片辞書105から素片データをロードし、必要に応じて復号処理を行う。素片辞書105には、音声を合成するための元となる音声素片データが格納されており、何らかの圧縮処理が施されている場合は、復号処理を施す。復号された音声素片データは、振幅制御部302で振幅係数が乗じられてパワー制御が行われる。素片加工部303では、声質変換のための素片伸縮処理が施される。声質を太くする場合は素片全体を伸長し、声質を細くする場合は素片全体を縮小するといった処理が施される。重畳制御部304では、合成パラメータのうち、ピッチパターンや音韻継続時間といった情報から、素片データの重畳を制御し、合成波形を生成する。波形重畳が完了したデータから逐次DAリングバッファ305に書き込み、出力サンプリング周期でDAコンバータに転送し、スピーカから出力する。

40

【0025】

次に音韻継続時間制御について詳細に説明する。図20に従来技術による音韻継続時間決

50

定部の機能ブロック図を示す。中間言語解析部 201 から解析結果が制御要因設定部 601 に入力される。制御要因設定部 601 では、例えば、音韻個々の継続時間長あるいは、単語全体での継続時間長などを予測するために必要な制御要因の設定を行う。予測には、例えば、対象となる音韻、前後の音韻の種類、構成しているフレーズのモーラ総数、文内位置といった情報が用いられ、継続時間推定部 602 に送られる。アクセント成分、フレーズ成分の各成分値予測には、自然発声データを基に数量化 I 類などの統計的手法を用いて予め学習した継続時間予測テーブル 604 が用いられる。予測された結果は継続時間修正部 603 に送られ、ユーザから発声速度指定があった場合は予測値の修正が施される。通常、発声速度指定は、5 ~ 10 段階程度に制御され、それぞれのレベルに対してあらかじめ割り当てられた定数を乗ずることにより行われる。発声速度を遅くしたい場合は音韻継続時間を長くし、発声速度を速くしたい場合は音韻継続時間を短くする。例えば、発声速度レベルが 5 段階に制御され、レベル 0 からレベル 4 まで指定可能だとする。それぞれのレベル n に対応した定数 T_n を次のように定める。すなわち、
 $T_0 = 2.0$ 、 $T_1 = 1.5$ 、 $T_2 = 1.0$ 、 $T_3 = 0.75$ 、 $T_4 = 0.5$ とする。

【0026】

先に予測された音韻継続時間のうち、母音長とポーズ長に対して、ユーザから指定されたレベル n に対応した定数 T_n が乗じられる。レベル 0 の場合は 2.0 が乗じられるので生成される波形は長くなり発声速度は遅くなる。レベル 4 の場合は 0.5 が乗じられるので生成される波形は短くなり発声速度は速くなる。上記の例では、レベル 2 が通常発声速度（デフォルト）となっている。

【0027】

発声速度制御が施された合成波形の例を図 21 に示す。図示したように、音韻継続時間の発声速度制御は通常、母音のみで行う。閉鎖区間長あるいは子音長は、発声速度に依らずほぼ一定と考えられるからである。発声速度を速くした (a) 図では母音長だけが 0.5 倍されており、重畳される音声素片数を減じて実現している。逆に発声速度を遅くした (c) 図では母音長だけが 1.5 倍されており、重畳される音声素片数を繰り返し使うなどして実現している。また、ポーズ長に対しては母音長制御と同様に、指定レベルに応じた定数が乗じられるため、発声速度が遅くなるほどポーズ長も長くなり、発声速度が速くなるほどポーズ長も短くなる。

【0028】

ここで発声速度が速い場合を考える。前述の例ではレベル 4 に当たる。テキスト音声変換システムの利用特性上、最大発声速度レベルは「早聞き機能」という意味合いが大きい。読上げ対象となるテキストの中でも、ユーザにとって、重要な部分とそうでない部分が存在するため、重要でない部分は発声速度を速くして読み飛ばし、重要な部分は通常発声速度で合成する。このような利用方法が一般的である。最近のテキスト音声変換装置では、早聞き機能用のボタンがあり、このボタンを押下すると発声速度レベルが最大に設定され最高速度で合成され、ボタンを離すと発声速度レベルが以前の設定値に復帰するといったものがある。

【0029】

【発明が解決しようとする課題】

しかしながら上記の従来技術では、以下に述べる問題があった。

(1) 早聞き機能を有効にすると、単純に音韻の継続時間長を短くする、言い換えると、生成する波形の長さを短くする処理を施しているため、波形生成部に負荷がかかるといった問題があった。波形生成部では、波形重畳が完了し、生成された波形データから逐次 D A リングバッファに書き込むという処理を行っているため、生成される波形長が短い場合はその分、波形生成処理に費やすことのできる時間が短くなることになる。波形データ長が半分になると、処理時間も半分で終了させなければならない。例えば、音韻継続時間長が半分になったからといって、必ずしも演算量が半分になるわけではないため、D A コンバータへの転送処理に、波形生成処理が追いつかない場合は、合成音が途中で止まる「音切れ」現象が発生する場合がある。

10

20

30

40

50

【 0 0 3 0 】

(2) 早聞き機能を有効にすると、単純に音韻の継続時間長を短くする処理が施されるため、ピッチパターンも基本的に線形に縮小される。つまり抑揚も時間的に速い周期で変動することになり、これは、不自然なイントネーションで非常に聞き取りにくい合成音となっていた。早聞き機能は、読上げ対象となるテキストを完全にスキップするのではなく、聞き流すという用途で用いられるため、抑揚の激しい合成音は不向きであった。従来技術において早聞き機能有効時の合成音声は、抑揚変化が激しすぎるため聞き取りにくく理解しづらいものとなっていた。

【 0 0 3 1 】

(3) 早聞き機能を有効にすると、音韻継続時間と共に、文章間のポーズも同一比率で縮小される。そのため、文章と文章の境界がほとんどなくなり、切れ目が分かり難くなっていた。1文の合成音声を出力した直後に、さらに次の1文の合成音声出力されるため、従来技術において早聞き機能有効時の合成音声は、テキスト内容を理解しつつ読み飛ばす用途においては不向きであった。

10

【 0 0 3 2 】

(4) 早聞き機能を有効にすると、テキスト全体に渡って、発声速度が速くなるため、早聞き解除のタイミングを取ることが難しかった。通常早聞き機能使用方法は、ある文章の中から所望の部分までを読み飛ばし、以降を通常速度で合成するというものである。従来技術によると、ユーザが欲した部分の読上げが行われ、早聞き機能解除をした時点では、所望の部分を大きく通り越してしまうといった問題があった。この場合、早聞き機能を解除した後一旦、読上げ対象区間を前にさかのぼって設定した後に通常発声速度で合成開始するといった面倒な操作をしなければいけなかった。またユーザは、必要な部分と必要でない部分とを聞き分けながら、早聞き機能の有効化・無効化の動作を行わなければならない、非常に労力を必要としていた。

20

【 0 0 3 3 】

本発明は、(A) 発声速度を速くした時に高負荷になって音切れが発生するという問題点と、(B) 発声速度を速くした時にピッチ変動周期も速くなり、不自然なイントネーションになってしまうという問題点を解決したテキスト音声変換における高速読み上げ制御方法を提供することを目的とする。

【 0 0 3 4 】

【課題を解決するための手段】

この発明は、上記課題(A)を解決するために、ユーザの指定する発声速度が最高速に設定された場合、すなわち早聞き機能が有効となった場合に、パラメータ生成手段における音韻継続時間決定手段において、統計的手法を用いて予測した継続時間予測テーブルに替えて、予め経験的に求めた継続時間規則テーブルを用いて音韻継続時間を決定し、また、ピッチパターン決定手段において、統計的手法により算出した予測テーブルを用いる代わりに、予め経験的に求めた規則テーブルを使用してピッチパターンを決定し、更に、声質決定手段においては声質が変化しないような声質変換係数を選択する。

30

【 0 0 3 5 】

また、この発明は、上記課題(B)を解決するために、ユーザの指定する発声速度が最高速に設定された場合に、アクセント成分及びフレーズ成分の計算を行わないようにすると共に基底ピッチを変更しないようにしている。

40

【 0 0 3 8 】

【発明の実施の形態】

第1の実施の形態

〔構成〕

以下、第1の実施の形態における構成を図面を参照しながら詳細に説明する。従来技術と異なる点は、発声速度が最高速に設定された場合、すなわち、早聞き機能が有効となった場合に内部演算処理の一部を簡略化、省略を行うことによって負荷軽減させた点である。

【 0 0 3 9 】

50

図1は、第1の実施の形態におけるパラメータ生成部102の機能ブロック図である。パラメータ生成部102への入力は従来と同じく、テキスト解析部101から出力される中間言語および、ユーザが個別に指定する韻律制御パラメータである。中間言語解析部801には一文毎の中間言語が入力され、以降の韻律生成処理で必要となる音韻系列・フレーズ情報・アクセント情報などといった中間言語解析結果が、それぞれピッチパターン決定部802、音韻継続時間決定部803、音韻パワー決定部804、音声素片決定部805、声質係数決定部806に出力される。

【0040】

ピッチパターン決定部802には、前述の中間言語解析結果に加えてユーザからの抑揚指定・声の高さ指定・発声速度指定・話者指定の各パラメータが入力され、ピッチパターンが合成パラメータ生成部807に出力される。ピッチパターンとは基本周波数の時間的遷移のことである。

10

【0041】

音韻継続時間決定部803には、前述の中間言語解析結果に加えてユーザからの発声速度指定のパラメータが入力され、それぞれの音韻の音韻継続時間・ポーズ長といったデータが合成パラメータ生成部807に出力される。

【0042】

音韻パワー決定部804には、前述の中間言語解析結果に加えてユーザからの声の大きさ指定パラメータが入力され、それぞれの音韻の音韻振幅係数が合成パラメータ生成部807に出力される。

20

【0043】

音声素片決定部805には、前述の中間言語解析結果に加えてユーザからの話者指定パラメータが入力され、波形重畳するための必要な音声素片アドレスが合成パラメータ生成部807に出力される。

【0044】

声質係数決定部806には、前述の中間言語解析結果に加えてユーザからの声質指定・発声速度指定の各パラメータが入力され、声質変換パラメータが合成パラメータ生成部807に出力される。

【0045】

合成パラメータ生成部807は、入力された各韻律パラメータ(前述したピッチパターン、音韻継続時間、ポーズ長、音韻振幅係数、音声素片アドレス、声質変換係数)から、フレーム(通常8ms程度の長さ)を一つの単位とした波形生成用のパラメータを生成し、波形生成部103に出力する。

30

【0046】

パラメータ生成部102において、従来技術と比較して異なる点は、発声速度指定パラメータが音韻継続時間決定部803のほかに、ピッチパターン決定部802、声質係数決定部806のそれぞれに入力されている点と、ピッチパターン決定部802、音韻継続時間決定部803、声質係数決定部806のそれぞれの内部処理である。テキスト解析部101および波形生成部103においては、従来と同様であるため、その構成に関する説明は省略する。

40

【0047】

ピッチパターン決定部802の構成について図2を用いて説明する。第1の実施の形態においては、アクセント成分およびフレーズ成分の決定に、数量化I類等の統計的手法を用いる場合と規則による場合との2通りの構成を有する。規則による制御の場合は、予め経験的に求められた規則テーブル910を用い、統計的手法による制御の場合は、自然発声データを基に数量化I類などの統計的手法を用いて予め学習した予測テーブル909を用いる。予測テーブル909のデータ出力はスイッチ907のa端子に接続され、規則テーブル910のデータ出力はスイッチ907のb端子に接続される。いずれの端子が選択されるかは、セレクタ906の出力によって決定される。

【0048】

50

セレクタ 906 には、ユーザから指定される発声速度レベルが入力され、スイッチ 907 を制御するための信号がスイッチ 907 に接続される。発声速度が最高レベルの場合はスイッチ 907 を b 端子側に接続し、それ以外の場合はスイッチ 907 を a 端子側に接続する。スイッチ 907 の出力は、アクセント成分決定部 902 とフレーズ成分決定部 903 に接続される。

【0049】

中間言語解析部 801 からの出力は制御要因設定部 901 に入力され、アクセント・フレーズ両成分の決定のための要因パラメータの解析が行われ、その出力がアクセント成分決定部 902 とフレーズ成分決定部 903 に接続される。

【0050】

アクセント成分決定部 902 とフレーズ成分決定部 903 には、スイッチ 907 からの出力が接続されており、予測テーブル 909 もしくは規則テーブル 910 を用いてそれぞれの成分値を決定しピッチパターン修正部 904 に出力する。

【0051】

ピッチパターン修正部 904 には、ユーザから指定される抑揚指定レベルが入力され、該レベルに応じて予め定められた定数が乗じられ、その結果が基底ピッチ加算部 905 に接続される。

【0052】

基底ピッチ加算部 905 にはさらに、ユーザから指定される声の高さレベル・話者指定および、基底ピッチテーブル 908 が接続されている。基底ピッチテーブル 908 には、ユーザ指定された声の高さレベルと性別とに応じて予め定められた定数値が格納されており、ピッチパターン修正部 904 からの入力に加算してピッチパターン時系列データとして合成パラメータ生成部 807 に出力する。

【0053】

音韻継続時間決定部 803 の構成について図 3 を用いて説明する。第 1 の実施の形態においては、音韻継続時間の決定に、数量化 I 類等の統計的手法を用いる場合と規則による場合との 2 通りの構成を有する。規則による制御の場合は、予め経験的に求められた継続時間規則テーブル 1007 を用い、統計的手法による制御の場合は、自然発声データを基に数量化 I 類などの統計的手法を用いて予め学習した継続時間予測テーブル 1006 を用いる。継続時間予測テーブル 1006 のデータ出力はスイッチ 1005 の a 端子に接続され、継続時間規則テーブル 1007 のデータ出力はスイッチ 1005 の b 端子に接続される。いずれの端子が選択されるかは、セレクタ 1004 の出力によって決定される。

【0054】

セレクタ 1004 には、ユーザから指定される発声速度レベルが入力され、スイッチ 1005 を制御するための信号がスイッチ 1005 に接続される。発声速度が最高レベルの場合はスイッチ 1005 を b 端子側に接続し、それ以外の場合はスイッチ 1005 を a 端子側に接続する。スイッチ 1005 の出力は、継続時間決定部 1002 に接続される。

【0055】

中間言語解析部 801 からの出力は制御要因設定部 1001 に入力され、音韻継続時間決定のための要因パラメータの解析が行われ、その出力が継続時間決定部 1002 に接続される。

【0056】

継続時間決定部 1002 には、スイッチ 1005 からの出力が接続されており、継続時間予測テーブル 1006 もしくは継続時間規則テーブル 1007 を用いて音韻継続時間長を決定し継続時間修正部 1003 に出力する。継続時間修正部 1003 には、ユーザから指定される発声速度レベルが入力され、該レベルに応じて予め定められた定数が乗じられて修正が施され、その結果が合成パラメータ生成部 807 に出力される。

【0057】

声質係数決定部 806 の構成について図 4 を用いて説明する。この例では声質変換指定レベルは 5 段階となっている。ユーザから指定される発声速度レベルおよび声質指定レベル

10

20

30

40

50

がセクタ 1102 に入力され、スイッチ 1103 を制御するための信号がスイッチ 1103 に接続される。この時のスイッチ制御信号は、発声速度が最高レベルの場合は無条件で c 端子有効にし、それ以外の場合は、声質指定レベルに応じた端子が有効となる。すなわち、声質レベルが 0 の時は a 端子、レベル 1 の時は b 端子、以下同様にレベル 4 の時 e 端子がそれぞれ有効となる。スイッチ 1103 の a ~ e の各端子は、声質変換係数テーブル 1104 に接続され、それぞれに対応した声質変換係数データが呼び出され、スイッチ 1103 の出力として声質係数選択部 1101 に接続される。声質係数選択部 1101 は入力された声質変換係数を合成パラメータ生成部 807 に出力する。

【0058】

[動作]

以上のように構成された第 1 の実施の形態における動作について詳細に説明する。従来技術と異なる点は、パラメータ生成に関わる処理であるので、それ以外の処理については説明を省略する。

【0059】

テキスト解析部 101 で生成された中間言語は、パラメータ生成部 102 内部の中間言語解析部 801 に送られる。中間言語解析部 801 では、中間言語上に記述されているフレーズ区切り記号、単語区切り記号、アクセント核を示すアクセント記号、そして音韻記号列から、韻律生成に必要なデータを抽出して、ピッチパタン決定部 802、音韻継続時間決定部 803、音韻パワー決定部 804、音声素片決定部 805、声質係数決定部 806 のそれぞれの機能ブロックへ送る。

【0060】

ピッチパタン決定部 802 では、声の高さの遷移であるイントネーションが生成され、音韻継続時間決定部 803 では、音韻個々の継続時間のほか、フレーズとフレーズの切れ目あるいは、文と文との切れ目に挿入するポーズ長を決定する。また、音韻パワー決定部 804 では、音声波形の振幅値の遷移である音韻パワーが生成され、音声素片決定部 805 では合成波形を生成するために必要となる音声素片の、素片辞書 105 におけるアドレスを決定する。声質係数決定部 806 では、素片データを信号処理で加工するためのパラメータの決定が行われる。ユーザから指定される韻律制御指定のうち、抑揚指定および声の高さ指定はピッチパタン決定部 802 に、発声速度指定はピッチパタン決定部 802 と音韻継続時間決定部 803 と声質係数決定部 806 に、声の大きさ指定は音韻パワー決定部 804 に、話者指定はピッチパタン決定部 802 と音声素片決定部 805 に、声質指定は声質係数決定部 806 にそれぞれ送られている。

【0061】

以下に、それぞれの機能ブロックごとに動作の説明を行う。

まず、図 2 を用いて、ピッチパタン決定部 802 の動作を詳細に説明する。中間言語解析部 201 から解析結果が制御要因設定部 901 に入力される。制御要因設定部 901 では、フレーズ成分、アクセント成分の大きさを決定するために必要な制御要因の設定を行う。フレーズ成分の大きさの決定に必要なデータとは、例えば、該当するフレーズを構成しているモーラ総数、文内での相対位置、先頭単語のアクセント型といった情報である。一方、アクセント成分の大きさの決定に必要なデータとは、例えば、該当するアクセント句のアクセント型、構成しているモーラ総数、品詞、フレーズ内での相対位置といった情報である。これらの成分値を決定するために予測テーブル 909 あるいは、規則テーブル 910 が使用される。前者は、自然発声データを基に数量化 I 類などの統計的手法を用いて予め学習したテーブルであり、後者は、予備実験等の実施により経験的に導き出された成分値が格納されたテーブルである。数量化 I 類に関しては公知であるのでここでは説明を省略する。どちらが選択されるかはスイッチ 907 により制御され、スイッチ 907 が a 端子に接続された場合は予測テーブル 909 が、b 端子に接続された場合は規則テーブル 910 が選択されることになる。

【0062】

ピッチパタン決定部 802 には、ユーザから指定される発声速度レベルが入力されており

10

20

30

40

50

、これによりセレクタ906を介してスイッチ907が駆動されている。セレクタ906は、入力された発声速度レベルが最高速度であった時、スイッチ907をb端子側に接続するような制御信号を送信する。逆に、入力された発声速度レベルが最高速度ではない時、スイッチ907をa端子側に接続するような制御信号を送信する。例えば、発声速度が5段階、レベル0からレベル4まで設定でき、数値が大きくなる程発声速度が速くなる仕様の場合、セレクタ906は、入力された発声速度レベルが4の時だけスイッチ907をb端子に接続するような制御信号を送信し、それ以外の時はa端子に接続するような制御信号を送信する。すなわち、発声速度が最高速度の場合は規則テーブル910が選択され、そうでない場合は予測テーブル909が選択されることになる。

【0063】

アクセント成分決定部902とフレーズ成分決定部903は、選択されたテーブルを用いてそれぞれの成分値の算出を行う。予測テーブル909が選択された場合は、統計的手法を用いてアクセント・フレーズ両成分の大きさを決定する。規則テーブル910が選択された場合は、あらかじめ決められた規則に従ってアクセント・フレーズ両成分の大きさを決定する。例えばフレーズ成分の大きさの規則化の例としては、文内の位置で決定し、文先頭フレーズは一律に0.3、文終端フレーズは一律に0.1、それ以外の文中フレーズは0.2などが考えられる。アクセント成分の大きさに関しても、アクセント型が1型の時とそれ以外の時、フレーズ内での単語位置が先頭の場合とそうでない場合といった具合に場合分けして、それぞれの条件に対して成分値を割り当てておく。このような構成にすることで、フレーズ・アクセント両成分値の決定はテーブル参照を行うだけで行える。本発明におけるピッチパタン決定部の主題は、統計的手法を用いてフレーズ・アクセント成分の大きさを決定する場合と比較して、演算量が少なく済み、処理時間の短縮が図れるモードを有する構成にすることである。したがって、規則化手順は上記に限られるものではない。

【0064】

以上のような処理が施され決定したアクセント成分、フレーズ成分は、ピッチパタン修正部904で抑揚制御が行われ、基底ピッチ加算部905で声の高さ制御が施される。

【0065】

ピッチパタン修正部904はユーザから指定される抑揚制御レベルに応じた係数を乗ずる操作が行われる。ユーザからの抑揚制御指定は例えば、3段階で与えられ、レベル1が抑揚を1.5倍に、レベル2が抑揚を1.0倍に、レベル3が抑揚を0.5倍にといった具合に定められている。

【0066】

基底ピッチ加算部905では、抑揚修正されたアクセント成分、フレーズ成分に対して、ユーザから指定される声の高さレベルあるいは、話者指定(性別)に応じた定数を加算する操作が行われ、ピッチパタン時系列データとして合成パラメータ生成部807に送られる。例えば、声の高さレベルが5段階、レベル0からレベル4まで設定できるシステムの場合、基底ピッチテーブル908に格納されているデータは男声音の場合、3.0、3.2、3.4、3.6、3.8といった数値、女性音の場合は、4.0、4.2、4.4、4.6、4.8といった数値が良く用いられる。

【0067】

次に音韻継続時間制御について図3を用いてその動作について詳細に説明する。中間言語解析部201から解析結果が制御要因設定部1001に入力される。制御要因設定部1001では、音韻継続時間(子音長・母音長・閉鎖区間長)、ポーズ長を決定するために必要な制御要因の設定を行う。音韻継続時間の決定に必要なデータとは、例えば、目標となる音韻の種別、対象音節の前後近傍の音韻の種別あるいは、単語内・呼気段落内の音節位置といった情報である。一方、ポーズ長決定に必要なデータとは、前後隣接するフレーズのモーラ総数といった情報である。これらの継続時間長を決定するために継続時間予測テーブル1006あるいは、継続時間規則テーブル1007が使用される。前者は、自然発声データを基に数量化I類などの統計的手法を用いて予め学習したテーブルであり、後者

10

20

30

40

50

は、予備実験等の実施により経験的に導き出された成分値が格納されたテーブルである。どちらが選択されるかはスイッチ1005により制御され、スイッチ1005がa端子に接続された場合は継続時間予測テーブル1006が、b端子に接続された場合は継続時間規則テーブル1007が選択されることになる。

【0068】

音韻継続時間決定部803には、ユーザから指定される発声速度レベルが入力されており、これによりセクタ1004を介してスイッチ1005が駆動されている。セクタ1004は、入力された発声速度レベルが最高速度であった時、スイッチ1005をb端子側に接続するような制御信号を送信する。逆に、入力された発声速度レベルが最高速度ではない時は、スイッチ1005をa端子側に接続するような制御信号を送信する。例えば、発声速度が5段階、レベル0からレベル4まで設定でき、数値が大きくなる程発声速度が速くなる仕様の場合、セクタ1004は、入力された発声速度レベルが4の時だけスイッチ1005をb端子に接続するような制御信号を送信し、それ以外の時はa端子に接続するような制御信号を送信する。すなわち、発声速度が最高速度の場合は継続時間規則テーブル1007が選択され、そうでない場合は継続時間予測テーブル1006が選択されることになる。

10

【0069】

継続時間決定部1002は、選択されたテーブルを用いて音韻継続時間、ポーズ長の算出を行う。継続時間予測テーブル1006が選択された場合は、統計的手法を用いて決定する。継続時間規則テーブル1007が選択された場合は、あらかじめ決められた規則に従って決定する。例えば音韻継続時間の規則化の例としては、その音韻の種類、文内の位置などに応じて基本長を割り当てておく。大量の自然発声データから音韻毎に平均を算出し、これを基本長としてもよい。ポーズ長に関しては、一律に300msを割り当てるか、あるいは、テーブル参照を行うだけで決定できるような構成が望ましい。本実施の形態における音韻継続時間決定部の主題は、統計的手法を用いて継続時間を決定する場合と比較して、演算量が少なく済み、処理時間の短縮が図れるモードを有する構成にすることである。したがって、規則化手順は上記に限られるものではない。

20

【0070】

以上のような処理が施され決定した継続時間は、継続時間修正部1003に送られる。継続時間修正部1003には、ユーザから指定される発声速度レベルも同時に入力されており、このレベルに応じて音韻継続時間の伸縮を行う。通常、発声速度指定は、5~10段階程度に制御され、それぞれのレベルに対してあらかじめ割り当てられた定数を母音の継続時間長あるいは、ポーズ長に対して乗ずることにより行われる。発声速度を遅くしたい場合は音韻継続時間を長くし、発声速度を速くしたい場合は音韻継続時間を短くする。

30

【0071】

次に声質係数決定について図4を用いてその動作について詳細に説明する。声質係数決定部806には、ユーザから指定される声質変換レベルと、発声速度レベルが入力される。これらの韻律制御パラメータは、セクタ1102を介してスイッチ1103を制御するために用いられる。セクタ1102はまず、発声速度レベルの判定を行う。発声速度レベルが最高速度の場合は、スイッチ1103をc端子に接続し、最高速度以外の場合は、声質変換レベルの判定を行う。この時は、声質変換レベルに応じた端子に接続するようにスイッチ1103を制御する。声質指定レベルが0の時はa端子、レベル1の時はb端子、以下同様にレベル4の時はe端子に接続する。スイッチ1103のa~eの各端子は、声質変換係数テーブル1104に接続され、それぞれに対応した声質変換係数データが呼び出される機能になっている。

40

【0072】

声質変換係数テーブル1104には、音声素片の伸縮係数が格納されており、例えば声質変換レベルnに対応する伸縮係数を K_n を次のように定める。すなわち、 $K_0 = 2.0$ 、 $K_1 = 1.5$ 、 $K_2 = 1.0$ 、 $K_3 = 0.8$ 、 $K_4 = 0.5$ のように設定する。これらの数値は、元となる音声素片の長さを K_n 倍に伸縮した後に波

50

形重畳して合成音声を生成するという意味である。レベル2の時は、係数値が1.0なので声質変換のための処理は一切行われなくなる。スイッチ1103のa端子に接続されている場合は、係数 K_0 が選択されて声質係数選択部1101に送られる。スイッチ1103のb端子に接続されている場合は、係数 K_1 が選択されて声質係数選択部1101に送られるといった具合である。

【0073】

ここで、図5を参照しながら素片の線形伸縮の方法の一例について述べる。声質変換レベル n における音声素片のデータの第 m サンプル目を X_{nm} とする。このように定義すると、声質変換後のデータ系列は、変換前のデータ系列 X_{2n} を用いて以下のようにして算出することができる。即ち、

レベル0では、

$$X_{00} = X_{20}$$

$$X_{01} = X_{20} \times 1/2 + X_{21} \times 1/2$$

$$X_{02} = X_{21}$$

レベル1では、

$$X_{10} = X_{20}$$

$$X_{11} = X_{20} \times 1/3 + X_{21} \times 2/3$$

$$X_{12} = X_{21} \times 2/3 + X_{22} \times 1/3$$

$$X_{13} = X_{22}$$

レベル3では、

$$X_{30} = X_{20}$$

$$X_{31} = X_{21} \times 3/4 + X_{22} \times 1/4$$

$$X_{32} = X_{22} \times 1/2 + X_{23} \times 1/2$$

$$X_{33} = X_{23} \times 1/4 + X_{24} \times 3/4$$

$$X_{34} = X_{25}$$

レベル4では、

$$X_{40} = X_{20}$$

$$X_{41} = X_{22}$$

のようになる。上記は、声質変換のための一例であって、これに限られるものではない。本実施の形態における声質係数決定部の主題は、発声速度レベルが最高速の時に声質変換指定を無効とする機能を有することにより、処理時間の短縮を図ることである。

【0074】

以上詳細に説明したように、第1の実施の形態によれば、発声速度が既定値最大に設定された場合に、テキスト音声変換処理の中で演算負荷が大きい機能ブロックを簡略化あるいは、無効にする処理を施しているため、高負荷による音切れが発生する機会を減少させ、聞き易い合成音声を生成することが可能となる。

【0075】

この場合、発声速度が最高レベル以外に設定された時の合成音と比較して、ピッチや継続時間などの韻律性能の若干の違い、声質変換機能が有効とならない、といったことが起きるが、最高速度での合成音出力は通常、読み飛ばしという意味合いで利用される場合がほとんどある。したがって、音声出力されるテキストの内容を把握・理解できれば良い、という程度の使用方法なので声質変換機能の有無、あるいは韻律性能低下といった点は音切れ現象と比較すると許容できるものと考えられる。

【0076】

第2の実施の形態

[構成]

第2の実施の形態における構成を図面を参照しながら詳細に説明する。本実施の形態が従来技術と異なる点は、発声速度が最高速に設定された場合、すなわち、早聞き機能が有効となった時にピッチパタン生成処理を変更する点である。したがって、従来と異なるパラメータ生成部、ピッチパタン決定部についてのみ説明する。

【 0 0 7 7 】

図 6 は第 2 の実施の形態におけるパラメータ生成部の機能ブロック図を示しており、このブロック図を用いて説明する。パラメータ生成部 1 0 2 への入力は従来と同じく、テキスト解析部 1 0 1 から出力される中間言語および、ユーザが個別に指定する韻律制御パラメータである。中間言語解析部 1 3 0 1 には一文毎の中間言語が入力され、以降の韻律生成処理で必要となる音韻系列・フレーズ情報・アクセント情報などといった中間言語解析結果が、それぞれピッチパタン決定部 1 3 0 2、音韻継続時間決定部 1 3 0 3、音韻パワー決定部 1 3 0 4、音声素片決定部 1 3 0 5、声質係数決定部 1 3 0 6 に出力される。

【 0 0 7 8 】

ピッチパタン決定部 1 3 0 2 には、前述の中間言語解析結果に加えてユーザからの抑揚指定・声の高さ指定・発声速度指定・話者指定の各パラメータが入力され、ピッチパタンが合成パラメータ生成部 1 3 0 7 に出力される。

10

【 0 0 7 9 】

音韻継続時間決定部 1 3 0 3 には、前述の中間言語解析結果に加えてユーザからの発声速度指定のパラメータが入力され、それぞれの音韻継続時間・ポーズ長といったデータが合成パラメータ生成部 1 3 0 7 に出力される。

【 0 0 8 0 】

音韻パワー決定部 1 3 0 4 には、前述の中間言語解析結果に加えてユーザからの声の大きさ指定パラメータが入力され、それぞれの音韻振幅係数が合成パラメータ生成部 1 3 0 7 に出力される。

20

【 0 0 8 1 】

音声素片決定部 1 3 0 5 には、前述の中間言語解析結果に加えてユーザからの話者指定パラメータが入力され、波形重畳するための必要な音声素片アドレスが合成パラメータ生成部 1 3 0 7 に出力される。

【 0 0 8 2 】

声質係数決定部 1 3 0 6 には、前述の中間言語解析結果に加えてユーザからの声質指定・発声速度指定の各パラメータが入力され、声質変換パラメータが合成パラメータ生成部 1 3 0 7 に出力される。

【 0 0 8 3 】

合成パラメータ生成部 1 3 0 7 は、入力された各韻律パラメータ（前述したピッチパタン、音韻継続時間、ポーズ長、音韻振幅係数、音声素片アドレス、声質変換係数）を、フレーム（通常 8 m s 程度の長さ）を一つの単位とした波形生成用のパラメータに変換し、波形生成部 1 0 3 に出力する。

30

【 0 0 8 4 】

パラメータ生成部 1 0 2 において、従来技術と比較して異なる点は、発声速度指定パラメータが音韻継続時間決定部 1 3 0 3 のほかに、ピッチパタン決定部 1 3 0 2 に入力されている点と、ピッチパタン決定部 1 3 0 2 の内部処理である。テキスト解析部 1 0 1 および波形生成部 1 0 3 においては、従来と同様であるため、その構成に関する説明は省略する。また、パラメータ生成部 1 0 2 の内部機能ブロックにおいても、ピッチパタン決定部 1 3 0 2 以外は従来と同様であるため、その構成に関する説明は省略する。

40

【 0 0 8 5 】

ピッチパタン決定部 1 3 0 2 の構成について図 7 を用いて説明する。中間言語解析部 1 3 0 1 からの出力は制御要因設定部 1 4 0 1 に入力され、アクセント・フレーズ両成分の決定のための要因パラメータの解析が行われ、その出力がアクセント成分決定部 1 4 0 2 とフレーズ成分決定部 1 4 0 3 に接続される。

【 0 0 8 6 】

アクセント成分決定部 1 4 0 2 とフレーズ成分決定部 1 4 0 3 には、予測テーブル 1 4 0 8 が接続され、数量化 I 類等の統計的手法を用いてそれぞれの成分の大きさを予測する。予測されたアクセント成分値、フレーズ成分値はピッチパタン修正部 1 4 0 4 に接続される。

50

【0087】

ピッチパターン修正部1404にはユーザから指定される抑揚指定レベルが入力され、該レベルに応じて予め定められた定数が前述のアクセント成分、フレーズ成分に乘じられ、その結果がスイッチ1405のa端子に接続される。スイッチ1405にはさらにb端子が存在し、セレクト1406から出力される制御信号により、端子a、端子bのいずれかに接続されるように構成されている。

【0088】

セレクト1406には、ユーザから指定される発声速度レベルが入力され、発声速度が最高レベルの場合はスイッチ1405をb端子に接続し、それ以外の場合はスイッチ1405をa端子に接続する制御信号を出力する。スイッチ1405のb端子は常にグランドに接続されており、スイッチ1405は、a端子が有効の時はピッチパターン修正部1404からの出力を、b端子が有効の時は0を基底ピッチ加算部1407に出力する機能を有している。

10

【0089】

基底ピッチ加算部1407にはさらに、ユーザから指定される声の高さレベル・話者指定および、基底ピッチテーブル1409が接続されている。基底ピッチテーブル1409には、ユーザ指定された声の高さレベルと話者の性別に応じて予め定められた定数値が格納されており、スイッチ1405からの入力に加算してピッチパターン時系列データとして合成パラメータ生成部1307に出力する。

【0090】

[動作]

以上のように構成された本発明の第2の実施の形態における動作について詳細に説明する。

20

【0091】

まず、テキスト解析部101で生成された中間言語は、パラメータ生成部102内部の中間言語解析部1301に送られる。中間言語解析部1301では、中間言語上に記述されているフレーズ区切り記号、単語区切り記号、アクセント核を示すアクセント記号、そして音韻記号列から、韻律生成に必要なデータを抽出して、ピッチパターン決定部1302、音韻継続時間決定部1303、音韻パワー決定部1304、音声素片決定部1305、声質係数決定部1306のそれぞれの機能ブロックへ送る。

30

【0092】

ピッチパターン決定部1302では、声の高さの遷移であるイントネーションが生成され、音韻継続時間決定1303では、音韻個々の継続時間のほか、フレーズとフレーズの切れ目あるいは、文と文との切れ目に挿入するポーズ長を決定する。また、音韻パワー決定部1304では、音声波形の振幅値の遷移である音韻パワーが生成され、音声素片決定部1305では合成波形を生成するために必要となる音声素片の、素片辞書105におけるアドレスを決定する。声質係数決定部1306では、素片データを信号処理で加工するためのパラメータの決定が行われる。

【0093】

ユーザから指定される種々の韻律制御指定のうち、抑揚指定および声の高さ指定はピッチパターン決定部1302に、発声速度指定はピッチパターン決定部1302と音韻継続時間決定部1303に、声の大きさ指定は音韻パワー決定部1304に、話者指定はピッチパターン決定部1302と音声素片決定部1305に、声質指定は声質係数決定部1306にそれぞれ送られている。

40

【0094】

以下に図7を用いてピッチパターン決定部1302の動作に関して説明する。従来技術と異なる点は、ピッチパターン生成に関わる処理であるので、それ以外の処理については省略する。

【0095】

中間言語解析部201から解析結果が制御要因設定部1401に入力される。制御要因設

50

定部 1401 では、フレーズ成分、アクセント成分の大きさを予測するために必要な制御要因の設定を行う。フレーズ成分の大きさの予測に必要なデータとは、例えば、該当するフレーズを構成しているモーラ総数、文内での相対位置、先頭単語のアクセント型といった情報である。一方、アクセント成分の大きさの予測に必要なデータとは、例えば、該当するアクセント句のアクセント型、構成しているモーラ総数、品詞、フレーズ内での相対位置といった情報である。これらの成分値を決定するために予測テーブル 1408 が使用される。予測テーブル 1408 は、自然発声データを基に数量化 I 類などの統計的手法を用いて予め学習したテーブルである。数量化 I 類に関しては公知であるのでここでは説明を省略する。

【0096】

制御要因設定部 1401 で解析された予測制御要因は、アクセント成分決定部 1402 とフレーズ成分決定部 1403 に送られ、それぞれにおいてアクセント成分の大きさ、フレーズ成分の大きさが予測テーブル 1408 を用いて予測される。第 1 の実施の形態でも示したように、予測モデルを使わずに規則でそれぞれの成分値を決定しても構わない。算出されたアクセント成分、フレーズ成分は、ピッチパタン修正部 1404 に送られ、ユーザから指定される抑揚指定レベルに応じた係数を乗ずる操作が行われる。

【0097】

ユーザからの抑揚制御指定は例えば、3 段階で与えられ、レベル 1 が抑揚を 1.5 倍に、レベル 2 が抑揚を 1.0 倍に、レベル 3 が抑揚を 0.5 倍にといった具合に定められている。

【0098】

修正されたアクセント、フレーズ両成分はスイッチ 1405 の a 端子に送られる。スイッチ 1405 は、a、b、2 つの端子を有しており、セレクトア 1406 からの制御信号によりどちらかの端子に接続するような機能になっている。一方の b 端子は常に 0 が入力されるようになっている。

【0099】

セレクトア 1406 にはユーザからの発声速度レベルが入力されており、これにより出力制御が行われている。セレクトア 1406 は、入力された発声速度レベルが最高速度であった時、スイッチ 1405 を b 端子側に接続するような制御信号を送信する。逆に、入力された発声速度レベルが最高速度ではない時、スイッチ 1405 を a 端子側に接続するような制御信号を送信する。例えば、発声速度が 5 段階、レベル 0 からレベル 4 まで設定でき、数値が大きくなる程発声速度が速くなる仕様の場合、セレクトア 1406 は、入力された発声速度レベルが 4 の時だけスイッチ 1405 を b 端子に接続するような制御信号を送信し、それ以外の時は a 端子に接続するような制御信号を送信する。すなわち、発声速度が最高速度の場合は 0 が選択され、そうでない場合は、ピッチパタン修正部 1404 の出力である修正されたアクセント成分値とフレーズ成分値が選択されることになる。

【0100】

選択されたデータは基底ピッチ加算部 1407 に送られる。基底ピッチ加算部 1407 にはユーザからの声の高さ指定レベルが入力されており、基底ピッチテーブル 1409 から該レベルに対応する基底ピッチデータが読み出され、前述のスイッチ 1405 からの出力値との加算処理が施され、ピッチパタンの時系列データとして合成パラメータ生成部 1307 に出力される。

【0101】

例えば、声の高さレベルが 5 段階、レベル 0 からレベル 4 まで設定できるシステムの場合、基底ピッチテーブル 1409 に格納されているデータは男声音の場合、3.0、3.2、3.4、3.6、3.8 といった数値、女性音の場合は、4.0、4.2、4.4、4.6、4.8 といった数値が良く用いられる。

【0102】

上記の例では、ピッチパタン修正部 1404 の出力と数値 0 とをスイッチ 1405 で切り替える処理を行っているが、無論、発声速度指定が最高レベルの時は、制御要因設定部 1

10

20

30

40

50

401からピッチパタン修正部1404までの処理は不要になる。

【0103】

図8に第2の実施の形態におけるピッチパタン生成処理のフローチャートを示す。ここで図中の記号は以下の通りとする。すなわち、入力文章中に含まれるフレーズ総数をI、単語総数をJ、第i番目のフレーズ成分の大きさを A_{p_i} 、第j番目のアクセント成分の大きさを A_{a_j} 、第j番目のアクセント句に対して指定される抑揚制御係数 E_j 、とする。

【0104】

ステップST101からステップST106にかけては、フレーズ成分の大きさ A_{p_i} の算出を行う。まずステップST101で、フレーズカウンタ i を0に初期化する。次いでステップST102で発声速度レベルの判定を行い、発声速度が最高速度である場合はステップST104に進み、そうでない場合はステップST103に進む。ステップST104では、第i番目のフレーズ成分の大きさ A_{p_i} を0に設定してステップST105に進む。一方ステップST103では数量化I類などの統計的手法を用いて第i番目のフレーズ成分の大きさ A_{p_i} が予測され、ステップST105に進む。ステップST105においては、フレーズカウンタ i を1インクリメントする。次いでステップST106で入力文章中のフレーズ総数Iとの比較を行い、フレーズカウンタ i が文内フレーズ総数Iを超えた場合、すなわち全てのフレーズに対する処理が終了した場合にフレーズ成分生成処理を終え、ステップST107に進む。そうでない場合は、ステップST102に戻り次のフレーズに対する処理を前述と同様に繰り返す。

【0105】

ステップST107からステップST113にかけては、アクセント成分の大きさ A_{a_j} の算出を行う。まずステップST107で、単語カウンタ j を0に初期化する。次いでステップST108で発声速度レベルの判定を行い、発声速度が最高速度である場合はステップST111に進み、そうでない場合はステップST109に進む。ステップST111では、第j番目のアクセント成分の大きさ A_{a_j} を0に設定してステップST112に進む。一方ステップST109では数量化I類などの統計的手法を用いて第j番目のアクセント成分の大きさ A_{a_j} が予測され、ステップST110に進む。ステップST110では、第j番目のアクセント句に対して抑揚修正処理が下式により行われる。

$$A_{a_j} = A_{a_j} \times E_j \quad \dots (4)$$

【0106】

ここで E_j は、ユーザが指定する抑揚制御レベルに応じてあらかじめ定められている抑揚制御係数であり、先にも説明したように例えば抑揚制御レベルが3段階で与えられ、レベル0が抑揚を1.5倍に、レベル1が抑揚を1.0倍に、レベル2が抑揚を0.5倍にといった場合は以下のようなになる。

レベル0 (抑揚を1.5倍)	$E_j = 1.5$
レベル1 (抑揚を1.0倍)	$E_j = 1.0$
レベル2 (抑揚を0.5倍)	$E_j = 0.5$

【0107】

抑揚修正終了後ステップST112に進む。ステップST112においては、単語カウンタ j を1インクリメントする。次いでステップST113で入力文章中の単語総数Jとの比較を行い、単語カウンタ j が文内単語総数Jを超えた場合、すなわち全て単語に対する処理が終了した場合にアクセント成分生成処理を終え、ステップST114に進む。そうでない場合は、ステップST108に戻り次のアクセント句に対する処理を前述と同様に繰り返す。

【0108】

ステップST114では、上記の処理で決定されたフレーズ成分値 A_{p_i} とアクセント成分値 A_{a_j} 、基底ピッチテーブル1409を参照して得られる基底ピッチ $l_n F_{min}$ とから式(1)によりピッチパタンを生成する。

【0109】

以上詳細に説明したように本発明の第2の実施の形態によれば、発声速度が既定値最大に

10

20

30

40

50

設定された場合に、ピッチパタンの抑揚成分を0にしてピッチパタン生成を行うため、時間的に速い周期で抑揚が変動することがなくなり、非常に聞き取りにくい合成音となることが解消される。

【0110】

図9は従来技術における発声速度によるピッチパタンの違いの説明図である。上段(a)が通常発声速度の場合であり、下段(b)が最高速度の場合である。横軸が時間であり、図中点線で示す曲線がフレーズ成分を表わし、実線で示す曲線がアクセント成分に対応している。最高速度が通常速度の2倍だとすると、生成される波形は通常時の約1/2となる。 $(T_2 = T_1 / 2)$ ピッチパタンの遷移も発声速度に比例して速くなるため、合成音声の抑揚は非常に速い周期での変動となることが図を見ても分かる。しかし実際の発声においては発声速度に応じて、フレーズの結合によるフレーズ境界の消失、アクセント結合によるアクセント句境界の消失といった現象が見られるため図(b)のようにはならない。発声速度が速くなるにつれて、ピッチパタンの変化も相対的に緩やかになることが多い。

10

【0111】

例えば図9の例で言えば2つのフレーズで構成されているが、これが1つのフレーズとして結合するといった現象が確認されている。従来技術においては、この点を考慮に入れておらず、非常に聞きづらい合成音声となっていたが、第2の実施の形態によれば、抑揚成分を0にすることで聞き取り易い合成音声を生成することが可能となる。

【0112】

抑揚成分を0にすることで抑揚の全くない、平坦なロボット音声のようになってしまうが、最高速度での合成音出力は通常、読み飛ばしという意味合いで利用される場合がほとんどある。したがって、音声出力されるテキストの内容を把握・理解できれば良い、という程度の使用方法なので、抑揚のない合成音声は使用に耐え得るものである。

20

【0113】

第3の実施の形態

[構成]

発明の第3の実施の形態における構成を図面を参照しながら詳細に説明する。本実施の形態が従来技術と異なる点は、文章間に合図音を入れることで文と文との境界を明示する点である。

30

【0114】

図10は、第3の実施の形態におけるパラメータ生成部102の機能ブロック図であり、この図を用いて説明する。パラメータ生成部102への入力は従来と同じく、テキスト解析部101から出力される中間言語および、ユーザが個別に指定する韻律制御パラメータである。ユーザからの韻律制御指定には、従来技術あるいは第1、第2の実施の形態にはないパラメータとして、合図音指定入力がある。これは後述する、文章間に挿入する合図音の種類を指定するための入力である。

【0115】

中間言語解析部1701には一文毎の中間言語が入力され、以降の韻律生成処理で必要となる音韻系列・フレーズ情報・アクセント情報などといった中間言語解析結果が、それぞれピッチパタン決定部1702、音韻継続時間決定部1703、音韻パワー決定部1704、音声素片決定部1705、声質係数決定部1706に出力される。

40

【0116】

ピッチパタン決定部1702には、前述の中間言語解析結果に加えてユーザからの抑揚指定・声の高さ指定・発声速度指定・話者指定の各パラメータが入力され、ピッチパタンが合成パラメータ生成部1708に出力される。

【0117】

音韻継続時間決定部1703には、前述の中間言語解析結果に加えてユーザからの発声速度指定のパラメータが入力され、それぞれの音韻継続時間・ポーズ長といったデータが合成パラメータ生成部1708に出力される。

50

【 0 1 1 8 】

音韻パワー決定部 1 7 0 4 には、前述の中間言語解析結果に加えてユーザからの声の大きさ指定パラメータが入力され、それぞれの音韻振幅係数が合成パラメータ生成部 1 7 0 8 に出力される。

【 0 1 1 9 】

音声素片決定部 1 7 0 5 には、前述の中間言語解析結果に加えてユーザからの話者指定パラメータが入力され、波形重畳するための必要な音声素片アドレスが合成パラメータ生成部 1 7 0 8 に出力される。

【 0 1 2 0 】

声質係数決定部 1 7 0 6 には、前述の中間言語解析結果に加えてユーザからの声質指定パラメータが入力され、声質変換パラメータが合成パラメータ生成部 1 7 0 8 に出力される。

10

【 0 1 2 1 】

合図音決定部 1 7 0 7 には、ユーザからの発声速度指定・合図音指定パラメータが入力され、合図音の種類および制御用のための合図音制御信号が波形生成部 1 0 3 に出力される。

【 0 1 2 2 】

合成パラメータ生成部 1 7 0 8 は、入力された各韻律パラメータ（前述したピッチパターン、音韻継続時間、ポーズ長、音韻振幅係数、音声素片アドレス、声質変換係数）から、フレーム（通常 8 m s 程度の長さ）を一つの単位とした波形生成用のパラメータに変換し、波形生成部 1 0 3 に出力する。

20

【 0 1 2 3 】

パラメータ生成部 1 0 2 において、従来技術と比較して異なる点は、合図音決定部 1 7 0 7 が新たな機能ブロックとして存在していることと、その入力パラメータとしてユーザから合図音指定がある点および、波形生成部 1 0 3 の内部構成である。テキスト解析部 1 0 1 においては、従来と同様であるため、その構成に関する説明は省略する。

【 0 1 2 4 】

はじめに合図音決定部 1 7 0 7 の構成について図 1 1 を用いて説明する。図に示すように、合図音決定部 1 7 0 7 は単にスイッチの役割を果たす機能ブロックである。ユーザから指定される発声速度レベルはスイッチ 1 8 0 1 の制御用端子に接続され、同じくユーザから指定される合図音コードがスイッチ 1 8 0 1 の a 端子に接続される。スイッチ 1 8 0 1 の b 端子は常にグランドに接続されている。スイッチ 1 8 0 1 は、発声速度レベルによって、端子 a、端子 b のいずれかに接続されるように構成されている。発声速度が最高レベルの場合はスイッチ 1 8 0 1 を a 端子に接続し、それ以外の場合はスイッチ 1 8 0 1 を b 端子に接続する。すなわちスイッチ 1 8 0 1 は、発声速度が最高レベルの時には合図音コードを、それ以外の時には 0 を出力する構成となっている。スイッチ 1 8 0 1 の出力は、合図音制御信号として波形生成部 1 0 3 に出力される。

30

【 0 1 2 5 】

次に波形生成部 1 0 3 の構成について図 1 2 を用いて説明する。第 3 の実施の形態においては、波形生成部 1 0 3 は、素片復号部 1 9 0 1 と振幅制御部 1 9 0 2 と素片加工部 1 9 0 3 と重畳制御部 1 9 0 4 と合図音制御部 1 9 0 5 と D A リングバッファ 1 9 0 6 の各機能ブロック、および合図音辞書 1 9 0 7 とから構成されている。

40

【 0 1 2 6 】

前述したパラメータ生成部 1 0 2 からの出力は、合成パラメータとして素片復号部 1 9 0 1 に入力される。素片復号部 1 9 0 1 には素片辞書 1 0 5 が接続されており、入力された合成パラメータのうち、素片アドレスを参照ポイントとして素片辞書 1 0 5 から素片データをロードし、必要に応じて復号処理を行い、復号素片データを振幅制御部 1 9 0 2 に出力する。素片辞書 1 0 5 には、音声を合成するための元となる音声素片データが格納されており、記憶容量の節約のために何らかの圧縮処理が施されている場合がある。この時は復号処理を施し、その必要がない非圧縮素片の場合は、単に読み込んでくるだけの処理と

50

なる。

【 0 1 2 7 】

振幅制御部 1 9 0 2 には、前述の復号後の音声素片データと合成パラメータとが入力されており、合成パラメータのうち音韻振幅係数によって素片データのパワー制御が行われ、素片加工部 1 9 0 3 に出力される。

【 0 1 2 8 】

素片加工部 1 9 0 3 には、前述の振幅制御された素片データと合成パラメータとが入力されており、合成パラメータのうち声質変換係数によって素片データの伸縮処理が施され、重畳制御部 1 9 0 4 に出力される。

【 0 1 2 9 】

重畳制御部 1 9 0 4 には、前述の伸縮処理が施された素片データと合成パラメータとが入力されており、合成パラメータのうちピッチパターン、音韻継続時間、ポーズ長といったパラメータを用いて素片データの波形重畳処理を施す。重畳制御部 1 9 0 4 で生成される波形は、逐次 D A リングバッファ 1 9 0 6 に出力され書き込まれる。D A リングバッファ 1 9 0 6 に書き込まれたデータは、当該テキスト音声変換システムで設定されている出力サンプリング周期で、図示していない D A コンバータに送られ、合成音がスピーカなどから出力される。

【 0 1 3 0 】

波形生成部 1 0 3 には、前述したパラメータ生成部 1 0 2 からの出力として合図音制御信号が合図音制御部 1 9 0 5 に入力される。合図音制御部 1 9 0 5 にはさらに合図音辞書 1 9 0 7 が接続されており、これに格納されているデータを必要に応じて加工して D A リングバッファ 1 9 0 6 に出力する。ただし書き込むタイミングは、重畳制御部 1 9 0 4 が 1 文章分の合成波形を出力し終えた後あるいは、合成波形を書き込む前とする。

【 0 1 3 1 】

合図音辞書 1 9 0 7 には例えば、各種効果音データの P C M (P u l s e C o d e M o d u l a t i o n) データで構築されている構成でも、基準正弦波データが格納された構成でも、どの形態でも構わない。この場合、合図音制御部 1 9 0 5 は、前者の辞書構成においては合図音辞書 1 9 0 7 からデータを読み出してきて、そのまま D A リングバッファ 1 9 0 6 に出力し、後者の辞書構成においては合図音辞書 1 9 0 7 からデータを読み出し、それを繰り返しつなぎ合わせるなどして出力する。合図音制御部 1 9 0 5 に接続されている合図音制御信号が 0 の場合は、D A リングバッファ 1 9 0 6 に出力する処理は行わない。

【 0 1 3 2 】

[動作]

以上のように構成された第 3 の実施の形態における動作について図 1 0 ~ 図 1 2 を用いて詳細に説明する。従来技術と異なる点は、ピッチパターン生成と波形生成に関わる処理であるので、それ以外の処理については省略する。

【 0 1 3 3 】

まず、テキスト解析部 1 0 1 で生成された中間言語は、パラメータ生成部 1 0 2 内部の中間言語解析部 1 7 0 1 に送られる。中間言語解析部 1 7 0 1 では、中間言語上に記述されているフレーズ区切り記号、単語区切り記号、アクセント核を示すアクセント記号、そして音韻記号列から、韻律生成に必要なデータを抽出して、ピッチパターン決定部 1 7 0 2、音韻継続時間決定部 1 7 0 3、音韻パワー決定部 1 7 0 4、音声素片決定部 1 7 0 5、声質係数決定部 1 7 0 6 のそれぞれの機能ブロックへ送る。

【 0 1 3 4 】

ピッチパターン決定部 1 7 0 2 では、声の高さの遷移であるイントネーションが生成され、音韻継続時間決定部 1 7 0 3 では、音韻個々の継続時間のほか、フレーズとフレーズの切れ目あるいは、文と文との切れ目に挿入するポーズ長を決定する。また、音韻パワー決定部 1 7 0 4 では、音声波形の振幅値の遷移である音韻パワーが生成され、音声素片決定部 1 7 0 5 では合成波形を生成するために必要となる音声素片の、素片辞書 1 0 5 におけるア

10

20

30

40

50

ドレスを決定する。声質係数決定部1706では、素片データを信号処理で加工するためのパラメータの決定が行われる。ユーザから指定される韻律制御指定のうち、抑揚指定および声の高さ指定はピッチパターン決定部1702に、発声速度指定は音韻継続時間決定部1703と合図音決定部1707に、声の大きさ指定は音韻パワー決定部1704に、話者指定はピッチパターン決定部1702と音声素片決定部1705に、声質指定は声質係数決定部1706に、合図音指定は合図音決定部1707に、それぞれ送られている。

【0135】

各機能ブロックのうち、ピッチパターン決定部1702、音韻継続時間決定部1703、音韻パワー決定部1704、音声素片決定部1705、声質係数決定部1706については、従来技術と同様であるのでここでは説明を省略する。

10

【0136】

第3の実施の形態におけるパラメータ生成部102が従来技術と異なる点は、合図音決定部1707が新たに追加されたことであるので、合図音決定部1707の動作について図11を用いて説明する。図に示すように、合図音決定部1707は単にスイッチの役割を果たす機能ブロックである。スイッチ1801は、ユーザから指定される発声速度レベルによって制御されるような構成を有しており、これにより端子a、端子bのいずれかに接続されるようになっている。制御信号である発声速度レベルが最高速度の時は、スイッチ1801をa端子に接続し、それ以外の場合はスイッチ1801をb端子に接続する。a端子には、ユーザから指定される合図音コードが入力されており、b端子にはグラウンド・レベルすなわち0が入力されている。すなわちスイッチ1801は、発声速度が最高レベルの時には合図音コードを、それ以外時には0を出力する構成となっている。スイッチ1801の出力は、合図音制御信号として波形生成部103に送られる。

20

【0137】

次に波形生成部103の動作について図12を用いて説明する。パラメータ生成部102内の合成パラメータ生成部1708で生成された合成パラメータは、波形生成部103内の素片復号部1901と振幅制御部1902と素片加工部1903と重畳制御部1904に送られる。

【0138】

素片復号部1901では、合成パラメータのうち、素片アドレスを参照ポインタとして素片辞書105から素片データをロードし、必要に応じて復号処理を行い、復号素片データを振幅制御部1902に送る。素片辞書105には合成波形を生成するための元となる音声素片が格納されており、これをピッチパターンで示される周期で重ね合わせていくことにより音声波形を生成するしくみとなっている。

30

【0139】

ここで音声素片とは、接続して合成波形を作るための音声の基本単位で、音の種類等に応じて様々なものが用意されている。一般的に、CV、VV、VCV、CVC(C:子音、V:母音)といった音韻連鎖で構成されている場合が多い。上記のように、同じ音韻の素片であっても、前後の音韻環境によって様々な単位で構築されているためデータ容量は膨大となる。そのため通常は、ADPCM(Adaptive Differential PCM)符号化や、周波数パラメータと駆動音源データの対で構成するといった、圧縮技術を施す場合が多い。無論、圧縮を行わずPCMデータとして構築されている場合もある。素片復号部1901によって復元された音声素片データは、振幅制御部1902に送られパワー制御が施される。

40

【0140】

振幅制御部1902には、合成パラメータのうち振幅係数が入力されており、先の音声素片データに乗じられて振幅制御が施される。振幅係数は、ユーザから指定される声の大きさレベル、音韻の種類、呼気段落内での音節位置、該音韻内での位置(立ち上がり区間・定常区間・立ち下がり区間)など、様々な情報から経験的に決定されている。振幅制御された音声素片は、素片加工部1903に送られる。

【0141】

50

素片加工部 1903 では、ユーザから指定された声質変換レベルに応じて素片データの伸縮処理（リサンプリング）が施される。声質変換とは、素片辞書 105 に登録されている素片データに、信号処理等の加工を施すことにより、聴感上、別話者として取り扱えるようにした機能である。一般に、素片データを線形に伸縮する処理を施して実現する場合が多い。伸長処理は、素片データのオーバーサンプリング処理で実現され、太い声となる。逆に縮小処理は、素片データのダウンサンプリング処理で実現され、細い声となる。同一データで別話者を実現するための機能であるため、声質変換処理は上記の手法に限るものではない。また、ユーザからの声質変換指定がない場合は当然のことながら、素片加工部 1903 での処理は一切行われない。

【0142】

以上の処理によって生成された音声素片は、重畳制御部 1904 で波形重畳処理が施される。一般的に、ピッチパターンで示されたピッチ周期で素片データをずらしながら重ね合わせて加算するという手法が用いられる。

【0143】

このようにして生成された合成波形は、逐次 DA リングバッファ 1906 に書き込まれ、当該テキスト音声変換システムで設定されている出力サンプリング周期で、図示していない DA コンバータに送られ、合成音がスピーカなどから出力される。

【0144】

波形生成部 103 にはさらに、パラメータ生成部 102 内の合図音決定部 1707 から送られる合図音制御信号が入力されている。合図音制御信号は、合図音制御部 1905 を介して合図音辞書 1907 に登録されているデータを DA リングバッファ 1906 に書き込むための信号である。合図音制御信号が 0 の場合、すなわち前述したように、ユーザから指定される発声速度が最高速度レベルではない時は、合図音制御部 1905 は一切の処理を行わない。0 以外の場合、すなわち前述したように、ユーザから指定される発声速度が最高速度レベルの時は、合図音制御信号を合図音の種類とみなして合図音辞書 1907 からのデータロードを行う。

【0145】

例えば、合図音の種類を 3 種類設ける。合図音辞書 1907 には、例えば、500 Hz の正弦波データ、1 KHz の正弦波データ、2 KHz の正弦波データがそれぞれ 1 周期分格納されており、それらを複数回繰り返し接続することにより「ピッ」という合図音を生成することとする。合図音制御信号の取り得る値は、0、1、2、3 の 4 種類となり、0 の時は一切の処理を行わず、1 の時は合図音辞書 1907 から 500 Hz の正弦波データを読み出してきて、それらを既定回繰り返し接続して DA リングバッファ 1906 に書き込む。1 の時は合図音辞書 1907 から 1 KHz の正弦波データを読み出してきて、それらを既定回繰り返し接続して DA リングバッファ 1906 に書き込む。2 の時は合図音辞書 1907 から 2 KHz の正弦波データを読み出してきて、それらを既定回繰り返し接続して DA リングバッファ 1906 に書き込む。ただし書き込むタイミングは、重畳制御部 1904 が 1 文章分の合成波形を出力し終えた後あるいは、合成波形を書き込む前である。したがって、合図音が出力されるのは文章間ということになる。出力される正弦波データは、100 ms ~ 200 ms 程度が適当と思われる。

【0146】

また、正弦波データではなく、出力されるべき合図音を直接 PCM データとして合図音辞書 1907 に格納しておくという構成でも構わない。この場合、合図音辞書 1907 からデータを読み出してきて、そのまま DA リングバッファ 1906 に出力する処理が施されることになる。

【0147】

以上詳細に説明したように、第 3 の実施の形態によれば、発声速度が既定値最大に設定された場合に、文章と文章の間に合図音を挿入する機能を有しているため、早聞き機能有効時での従来技術での問題点である、文境界が把握しにくく、読上げテキストの内容理解が困難であるといったことが解消される。

10

20

30

40

50

【0148】

例えば、以下の文言をテキスト合成する場合を考える。

「出席予定者：開発部 山田部長。企画室 斉藤室長。営業1部 渡辺部長。」処理単位、すなわち1文章の区切り記号は句点「。」とすると、上記の文言は以下の3文章からなる。

(1) 「出席予定者：開発部 山田部長。」

(2) 「企画室 斉藤室長。」

(3) 「営業1部 渡辺部長。」

従来技術によれば、発声速度が速くなるとそれぞれの文終端におけるポーズ長も短くなるため、文章(1)の最後の「山田部長」という合成音声と、文章(2)の先頭の「企画室」という合成音声がほぼ連続して出力されるため、「山田部長」=「企画室」というような誤った認識を受ける場合も発生する。

10

【0149】

しかしながら、第3の実施の形態によれば、「山田部長」という合成音声と、「企画室」という合成音声の間に、例えば「ピッ」という合図音が挿入されるため、上記のような誤認識は発生しない。

【0150】

第4の実施の形態

[構成]

本発明の第4の実施の形態における構成を図13を参照しながら詳細に説明する。この実施の形態が従来技術と異なる点は、早聞き機能有効時の音韻継続時間の伸縮率決定の際に、現在処理中のテキストが文内における先頭単語あるいは先頭フレーズであるかを判定して、その結果により伸縮係数を決定する点である。したがって、従来と異なる音韻継続時間決定部についてのみ説明し、それ以外の機能ブロックすなわち、テキスト解析部、波形生成部、音韻継続時間決定部以外のパラメータ生成部内部モジュールについては説明を省略する。

20

【0151】

音韻継続時間決定部203への入力とは従来と同じく、中間言語解析部201からの音韻・韻律情報を含んだ解析結果および、ユーザからの指定される発声速度レベルである。1文章に対する中間言語解析結果は制御要因設定部2001と単語カウンタ2005とに接続されている。制御要因設定部2001では、音韻継続時間決定のために必要な制御要因パラメータの解析が行われ、その出力が継続時間推定部2002に接続される。継続時間の決定には数量化I類等の統計的手法を用いており、例えば、音韻長は通常、目標となる音韻の前後近傍の音韻の種別あるいは、単語内・呼気段落内の音節位置などにより予測され、ポーズ長は、前後隣接するフレーズのモーラ総数などといった情報から予測が行われる場合が多い。制御要因設定部2001はこれら予測に必要な情報の抽出を行っている。

30

【0152】

継続時間推定部2002には、継続時間予測テーブル2004が接続されており、これを用いて継続時間の予測が行われ、継続時間修正部2003に出力される。継続時間予測テーブル2004は、大量の自然発声データを基に数量化I類などの統計的手法を用いて予め学習されたデータである。

40

【0153】

一方、単語カウンタ2005では、現在解析中の音韻が、文章内における先頭単語あるいは先頭フレーズに含まれているのか、そうでないのかの判定を行い、その結果を伸縮係数決定部2006に出力する。

【0154】

伸縮係数決定部2006にはさらに、ユーザから指定される発声速度レベルが入力されており、現在処理中の音韻に対する音韻継続時間長の修正係数を決定する機能を有しており、これを継続時間修正部2003に接続している。

【0155】

50

継続時間修正部 2003 では、継続時間推定部 2002 で予測された音韻継続時間に対して、伸縮係数決定部 2006 で決定された伸縮係数を乗じることにより、音韻継続時間の修正を行い合成パラメータ生成部へ出力する。

【0156】

[動作]

以上のように構成された本発明の第4の実施の形態における動作について図13～図14を用いて詳細に説明する。従来技術と異なる点は、音韻継続時間決定に関わる処理であるので、それ以外の処理については省略する。

【0157】

中間言語解析部 2001 から1文章に対応する解析結果が制御要因設定部 2001 と単語カウンタ 2005 へ入力される。制御要因設定部 2001 では、音韻継続時間(子音長・母音長・閉鎖区間長)、ポーズ長を決定するために必要な制御要因の設定を行う。音韻継続時間の決定に必要なデータとは、例えば、目標となる音韻の種別、対象音節の前後近傍の音韻の種別あるいは、単語内・呼気段落内の音節位置といった情報である。一方、ポーズ長決定に必要なデータとは、前後隣接するフレーズのモーラ総数といった情報である。これらの継続時間長を決定するために継続時間予測テーブル 2004 が使用される。

【0158】

継続時間予測テーブル 2004 は、自然発声データを基に数量化I類などの統計的手法を用いて予め学習したテーブルである。継続時間推定部 2002 は、このテーブルを参照しながら音韻継続時間、ポーズ長の予測を行う。継続時間推定部 2002 で算出される個々の音韻継続時間長は、通常発声速度の場合のものである。これらは、継続時間修正部 2003 において、ユーザから指定された発声速度に応じて修正が施される構成となっている。通常、発声速度指定は、5～10段階程度に制御され、それぞれのレベルに対してあらかじめ割り当てられた定数を乗ずることにより行われる。発声速度を遅くしたい場合は音韻継続時間を長くし、発声速度を速くしたい場合は音韻継続時間を短くする。

【0159】

一方、単語カウンタ 2005 にも、中間言語解析部 2001 から1文章に対応する解析結果が入力されており、現在解析中の音韻が、文章内における先頭単語あるいは先頭フレーズに含まれているのか、そうでないのかの判定が行われる。本実施の形態では、文章内における先頭単語であるか否かの判定を行う機能として説明を行う。単語カウンタ 2005 から送られる判定結果は、該音韻が文内先頭単語に含まれている場合に TRUE、そうでない場合に FALSE を出力することとする。単語カウンタ 2005 での判定結果は伸縮係数決定部 2006 に送られる。

【0160】

伸縮係数決定部 2006 には前述の単語カウンタ 2005 からの判定結果に加えて、ユーザから指定される発声速度レベルが入力されており、これら2つのパラメータから該音韻の伸縮係数の算出を行う。例えば、発声速度レベルが5段階に制御され、発声速度が遅い方からレベル0、レベル1、レベル2、レベル3、レベル4まで指定可能だとする。それぞれのレベルnに対応した定数 T_n を次のように定める。すなわち、

$T_0 = 2.0$ 、 $T_1 = 1.5$ 、 $T_2 = 1.0$ 、 $T_3 = 0.75$ 、 $T_4 = 0.5$ とする。通常発声速度はレベル2となり、早聞き機能が有効とされると発声速度はレベル4に設定されることになる。単語カウンタ 2005 からの信号が TRUE の場合、発声速度レベルが0～3まで範囲であれば上記 T_n をそのまま継続時間修正部 2003 へ出力する。発声速度レベルが4であれば、通常発声時の T_2 の数値を出力する。単語カウンタ 2005 からの信号が FALSE の場合は、発声速度レベルに関わらず上記 T_n をそのまま継続時間修正部 2003 へ出力する。

【0161】

継続時間修正部 2003 では、継続時間推定部 2002 から送られる音韻継続時間長に対して、伸縮係数決定部 2006 からの伸縮係数を乗じて修正を施す。ただし修正を行うのは通常、母音長のみである。発声速度レベルに応じた修正が施された音韻継続時間は合成

10

20

30

40

50

パラメータ生成部へ送られる。

【0162】

さらに詳細に説明するために図14に継続時間決定処理のフローチャートを示す。ここで図中の記号は以下の通りとする。すなわち、入力文章中に含まれる単語総数を I 、第 i 番目の単語を構成する音韻に対する継続時間修正係数を TC_i 、ユーザから指定される発声速度レベルを lev （ただし範囲は0～4までの5段階とし、数値が多いほど速度が速いこととする）、発声速度がレベル n の時の伸縮係数を $T(n)$ 、第 i 番目の単語の第 j 番目の母音長を T_{ij} 、単語を構成する音節数はそれぞれの単語によって変わるがここでは簡単化のために一律 J とする。

【0163】

まずステップST201で単語数カウンタ i を0に初期化する。次いでステップST202で単語数と発声速度レベルの判定が行われる。現在処理中の単語数カウンタが0でかつ、発声速度レベルが4の時、これはすなわち、現在処理している音節が文内先頭単語に属しており、かつ発声速度が最高レベルの時であるが、この時はステップST204に進み、そうでないときはステップST203に進む。ステップST204では発声速度レベル2の値が修正係数として選択され、ステップST205に進む。すなわち、

$$TC_i = T(2) \quad \dots (5)$$

となる。

【0164】

ステップST203では、ユーザから指定されたレベル通りの修正係数が選択され、ステップST205に進む。すなわち、

$$TC_i = T(lev) \quad \dots (6)$$

となる。

【0165】

ステップST205では、音節カウンタ j が0に初期化されステップST206に進む。ステップST206では第 i 番目の単語の第 j 番目の母音の継続時間 T_{ij} が、先に求められた修正係数 TC_i によって下式を用いて行われる。

$$T_{ij} = T_{ij} \times TC_i \quad \dots (7)$$

【0166】

次いでステップST207で音節カウンタ j が1インクリメントされステップST208に進む。ステップST208では、音節カウンタ j と該単語の音節総数 J との比較を行い、音節カウンタ j が音節総数 J を超えた場合、すなわち該単語の全ての音節に対する処理が終了した場合にステップST209に進む。そうでない場合は、ステップST206に戻り次の音節に対する処理を前述と同様に繰り返す。

【0167】

ステップST209では単語数カウンタ i が1インクリメントされ、次のステップST210に進む。

【0168】

ステップST210では、単語数カウンタ i と単語総数 I との比較を行い、単語数カウンタ i が単語総数 I を超えた場合、すなわち入力文章中の全て単語に対する処理が終了した場合は処理を終了し、そうでない場合は、ステップST202に戻り次の単語に対する処理を前述と同様に繰り返す。

【0169】

上記の処理により、ユーザから指定される発声速度レベルが最高速度となっても、文章先頭単語だけは通常の発声速度での合成音が生成されることになる。

【0170】

以上詳細に説明したように、第4の実施の形態によれば、発声速度が既定値最大に設定された場合に、文先頭の単語に対して音韻継続時間制御を通常の発声速度として処理するため、ユーザが早聞き機能解除のタイミングを計りやすいという効果がある。例えば、ソフトウェア仕様書などのマニュアル類には、「第3章」あるいは「4.1.3」などの項目

10

20

30

40

50

番号が付与されている場合がほとんどある。こういったマニュアル類をテキスト音声変換で読上げを行う際に、第3章から聞きたい、あるいは4.1.3節から聞きたいといった場合に、従来技術においては、早聞き機能を有効にした後ユーザが、高速で出力される合成音声の中から「ダイサンショー」あるいは「ヨンテンイッテンサン」といったキーワードを聞き分け、早聞き機能を解除するといった面倒な操作が必要であった。第4の実施の形態によれば、ユーザに負担をかけずに早聞き機能の有効化・無効化を実現することが可能となる。

【0171】

尚、本発明は前述の実施の形態に限定されるものではなく、本発明の趣旨に基づいて種々変形させることが可能である。例えば、第1の実施の形態において、発声速度が既定値最大に設定された場合に、テキスト音声変換処理の中で演算負荷が大きい機能ブロックを簡略化あるいは、無効にする処理を施しているが、この処理は最大発声速度に限らない。つまり、ある閾値を設けて、その閾値を超えたときに前述の処理を施す構成でも構わない。また、高負荷処理として数量化I類による韻律パラメータの予測処理、声質変換のための素片データ加工処理を挙げているが、これに限るものではない。他に高負荷処理機能（例えばエコーや高域強調などの音響処理など）を有している場合は当然のことながら、これを無効化あるいは簡略化といった処理形態にすることが望ましい。また、声質変換処理として波形そのものを線形伸縮しているが、非線形伸縮でも、あるいは周波数パラメータに対して規定の変換関数に通して変形するといった方法でも構わない。また、音韻継続時間決定規則、ピッチパタン決定規則を挙げているが、本発明では演算量が少なく済み、処理時間の短縮が図れるモードを有する構成にすること目的としているため、規則化手順は上記に限られるものではない。逆に、通常発声速度の時には、統計的手法を用いた韻律パラメータの予測を行っているが、規則化手順よりも演算負荷がかかる処理であればこれに限るものではない。また、その予測に用いる制御要因を幾つか挙げているがこれはあくまでも一例である。

【0172】

第2の実施の形態において、発声速度が既定値最大に設定された場合に、ピッチパタンの抑揚成分を0にしてピッチパタン生成を行っているが、この処理は最大発声速度に限らない。即ち、ある閾値を設けて、その閾値を超えたときに前述の処理を施す構成でも構わない。また、抑揚成分を完全に0にしているが、通常時に比べて抑揚成分を弱めるといった方法でも構わない。例えば、発声速度が既定値最大に設定された時は、抑揚指定レベルを強制的に最低レベルに設定し、ピッチパタン修正部において抑揚成分を縮小するといった構成でも構わない。ただこの時の抑揚指定レベルは、高速合成時においても聞き易いイントネーションとなる必要がある。また、ピッチパタンのアクセント成分、フレーズ成分を数量化I類によって決定しているが規則によって決定しても無論構わない。また、予測を行う際にその制御要因を幾つか挙げているがこれはあくまでも一例である。

【0173】

第3の実施の形態において、発声速度が既定値最大に設定された場合に、文章と文章の間に合図音を挿入しているが、この処理は最大発声速度に限らない。即ち、ある閾値を設けて、その閾値を超えたときに前述の処理を施す構成でも構わない。また、実施例では基準正弦波の繰り返しにより合図音を生成しているが、ユーザの注意を引けるものであればこれに限らない。録音された効果音をそのまま出力する構成でも構わない。無論、実施例で示したような合図音辞書を持たずに、内部回路あるいはプログラムでその都度生成するような構成でも構わない。またこの実施の形態では1文の合成波形直後に合図音を挿入する構成となっているが、逆に合成波形直前でも構わない。発声速度が既定値最大に設定された時に、ユーザに対して文章境界が明示できればそれでよい。また、この実施の形態ではパラメータ生成部に合図音の種類を指定するための入力が存在するが、ハードウェア規模、ソフトウェア規模の制限などから、これを省略してもよい。しかしながら、ユーザの好みによって合図音を変えることのできる構成の方が好ましい。

【0174】

第4の実施の形態において、発声速度が既定値最大に設定された場合に、文先頭の単語に対して音韻継続時間制御を通常（デフォルト）の発声速度として処理しているが、この処理は最大発声速度に限らない。即ち、ある閾値を設けて、その閾値を超えたときに前述の処理を施す構成でも構わない。また、通常発声速度で処理する単位を文先頭の1単語としているが、先頭2単語あるいは先頭フレーズという構成でも構わない。また、通常発声速度ではなく、レベルを1段階落とすといった方法も十分考えられる。

【0175】

【発明の効果】

以上詳細に説明したように、請求項1に係る発明によれば、入力されたテキストから音韻・韻律記号列を生成するテキスト解析手段と、前記音韻・韻律記号列に対して少なくとも音声素片・音韻継続時間・基本周波数の合成パラメータを生成するパラメータ生成手段と、音声の基本単位となる音声素片が登録された素片辞書と前記パラメータ生成手段から生成される合成パラメータに基づいて前記素片辞書を参照しながら波形重畳を行って合成波形を生成する波形生成手段とを備えたテキスト音声変換装置における高速読み上げ制御方法であって、前記パラメータ生成手段は、音韻継続時間を予め経験的に求めた継続時間規則テーブルと、音韻継続時間を統計的手法を用いて予測した継続時間予測テーブルとを併せ持ち、ユーザから指定される発声速度が閾値を超えた時には前記継続時間規則テーブルを用い、閾値を超えていない時には前記継続時間予測テーブルを用いて音韻継続時間の決定を行う音韻継続時間決定手段を有する構成としたことにより、また、請求項3に係る発明によれば、前記パラメータ生成手段は、アクセント成分及びフレーズ成分を決定するために必要となるデータを、予め経験的に求めた規則テーブルと、統計的手法を用いて予測した予測テーブルとを併せ持ち、ユーザから指定される発声速度が閾値を超えた時には前記規則テーブルを用い、閾値を超えていない時には前記予測テーブルを用いてアクセント成分及びフレーズ成分を決定することによりピッチパタンを決定するピッチパタン決定手段を有する構成としたことにより、更に、請求項5に係る発明によれば、前記パラメータ生成手段は、前記音声素片を変形させて声質を切り換えるための声質変換係数テーブルを備え、ユーザから指定される発声速度が閾値を超えたときには、声質が変化しないような係数を前記声質変換係数テーブルから選択する声質係数決定手段を有する構成としたので、発声速度が既定値最大に設定された場合に、テキスト音声変換処理の中で演算負荷が大きい機能ブロックを簡略化あるいは、無効にする処理を施しているため、高負荷による音切れが発生する機会を減少させ、聞き易い合成音声を生成することが可能となる。

【0176】

また、請求項7に係る発明によれば、前記パラメータ生成手段は、ユーザが指定した抑揚レベルに応じて修正したピッチパタンを出力するピッチパタン修正手段と、ユーザが指定した発声速度に応じて前記修正したピッチパタンを基底ピッチに加算するか否かを選択する切り換え手段とを有し、前記発声速度が所定の閾値を超えた場合には前記基底ピッチを変更しないように前記切り換え手段を制御する構成としたので、発声速度が既定値最大に設定された場合に、ピッチパタンの抑揚成分を0にしてピッチパタン生成を行うため、時間的に速い周期で抑揚が変動することがなくなり、非常に聞き取りにくい合成音となることが解消される。

【図面の簡単な説明】

【図1】本発明の第1の実施の形態におけるパラメータ生成部の機能ブロック図である。

【図2】本発明の第1の実施の形態におけるピッチパタン決定部の機能ブロック図である。

【図3】本発明の第1の実施の形態における音韻継続時間決定部の機能ブロック図である。

【図4】本発明の第1の実施の形態における声質係数決定部の機能ブロック図である。

【図5】声質変換のためのデータのリサンプリング周期の説明図である。

【図6】本発明の第2の実施の形態におけるパラメータ生成部の機能ブロック図である。

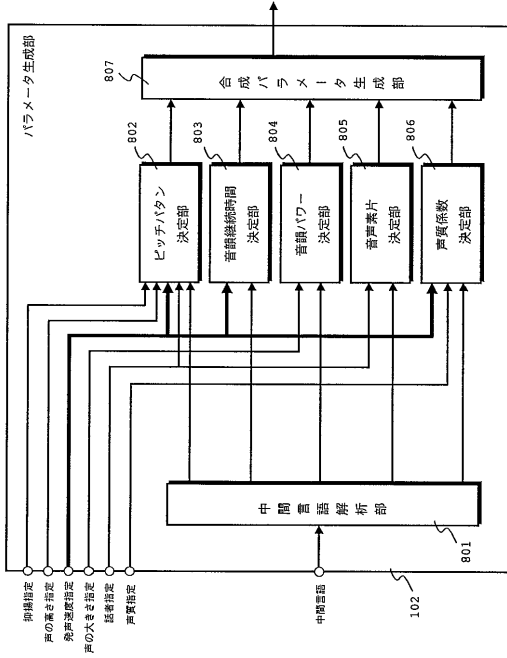
【図7】本発明の第2の実施の形態におけるピッチパタン決定部の機能ブロック図である。

- 。【図 8】本発明の第 2 の実施の形態におけるピッチパタン生成フローチャートである。
 【図 9】発声速度によるピッチパタンの違いの説明図である。
 【図 10】本発明の第 3 の実施の形態におけるパラメータ生成部の機能ブロック図である。
 。【図 11】本発明の第 3 の実施の形態における合図音決定部の機能ブロック図である。
 【図 12】本発明の第 3 の実施の形態における波形生成部の機能ブロック図である。
 【図 13】本発明の第 4 の実施の形態における音韻継続時間決定部の機能ブロック図である。
 【図 14】本発明の第 4 の実施の形態における継続時間決定フローチャートである。 10
 【図 15】一般的なテキスト音声変換処理の機能ブロック図である。
 【図 16】従来技術によるパラメータ生成部の機能ブロック図である。
 【図 17】従来技術による波形生成部の機能ブロック図である。
 【図 18】ピッチパタン生成過程モデルの説明図である。
 【図 19】従来技術によるピッチパタン決定部の機能ブロック図である。
 【図 20】従来技術による音韻継続時間決定部の機能ブロック図である。
 【図 21】発声速度の違いによる波形伸縮の説明図である。

【符号の説明】

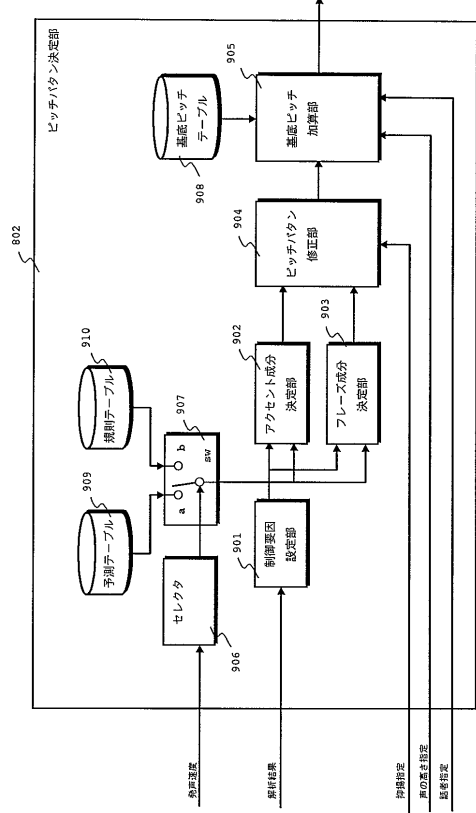
1 0 1	テキスト解析部	
1 0 2	パラメータ生成部	20
1 0 3	波形生成部	
1 0 4	単語辞書	
1 0 5	素片辞書	
8 0 1 , 1 3 0 1 , 1 7 0 1 ,	中間言語解析部	
8 0 2 , 1 3 0 2 , 1 7 0 2 ,	ピッチパタン決定部	
8 0 3 , 1 3 0 3 , 1 7 0 3	音韻継続時間決定部	
8 0 4 , 1 3 0 4 , 1 7 0 4	音韻パワー決定部	
8 0 5 , 1 3 0 5 , 1 7 0 5	音声素片決定部	
8 0 6 , 1 3 0 6 , 1 7 0 6	声質係数決定部	
1 7 0 7	合図音決定部	30
8 0 7 , 1 3 0 7 , 1 7 0 8	合成パラメータ生成部	

【図 1】



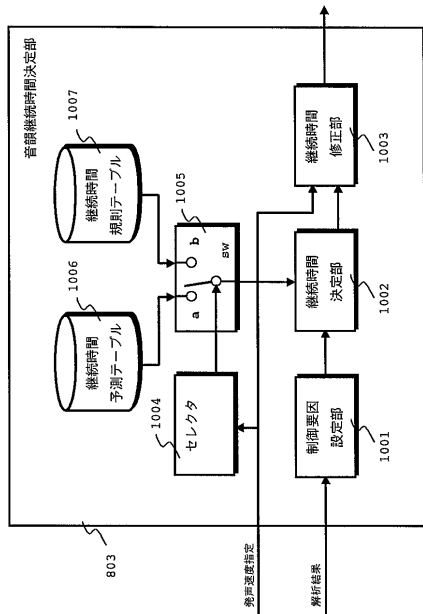
本発明の第 1 の実施の形態におけるパラメータ生成部の機能ブロック図

【図 2】



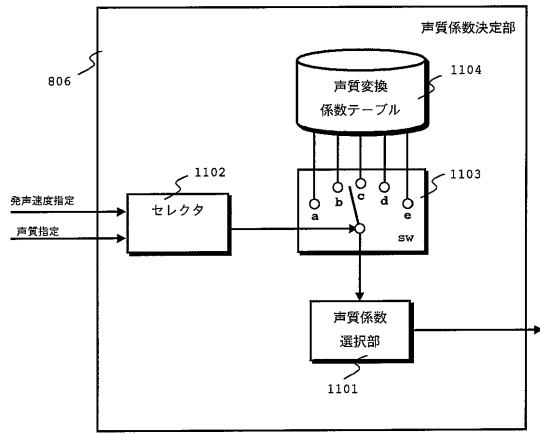
本発明の第 1 の実施の形態におけるピッチパタン決定部の機能ブロック図

【図 3】



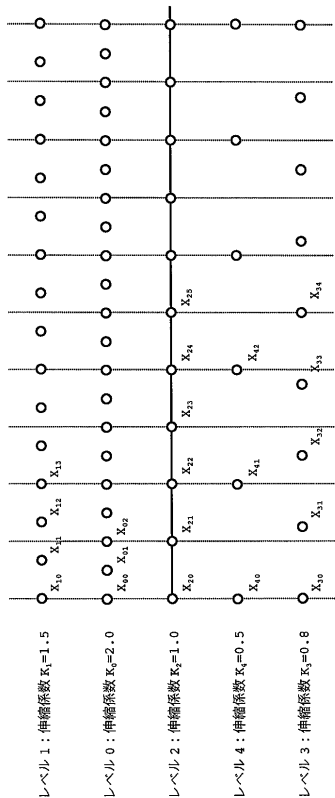
本発明の第 1 の実施の形態における音韻継続時間決定部の機能ブロック図

【図 4】



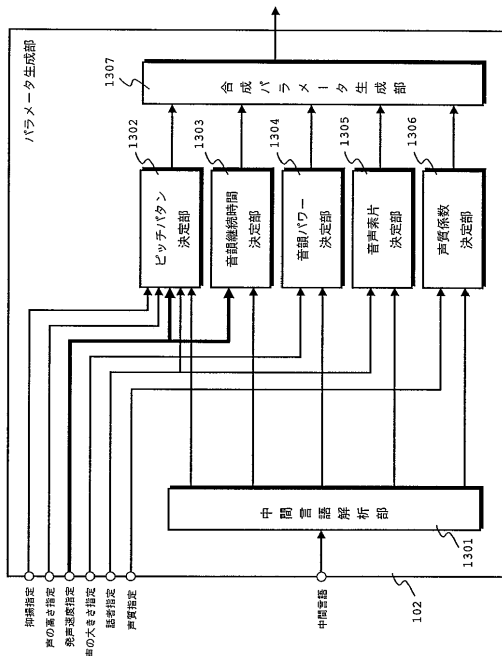
本発明の第 1 の実施の形態における声質係数決定部の機能ブロック図

【図 5】



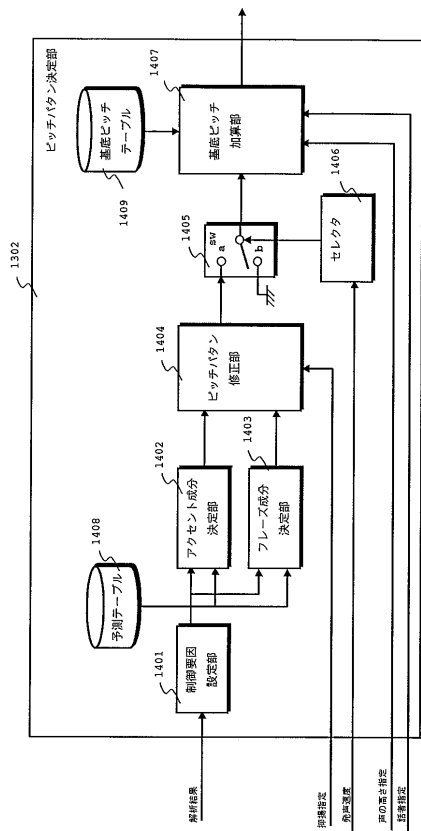
声質変換のためのデータのリスサンプリング周期の説明図

【図 6】



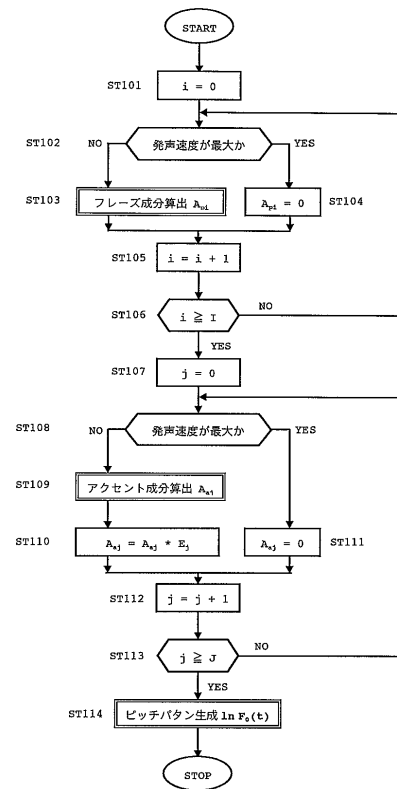
本発明の第 2 の変換の形態におけるパラメータ生成部の機能ブロック図

【図 7】



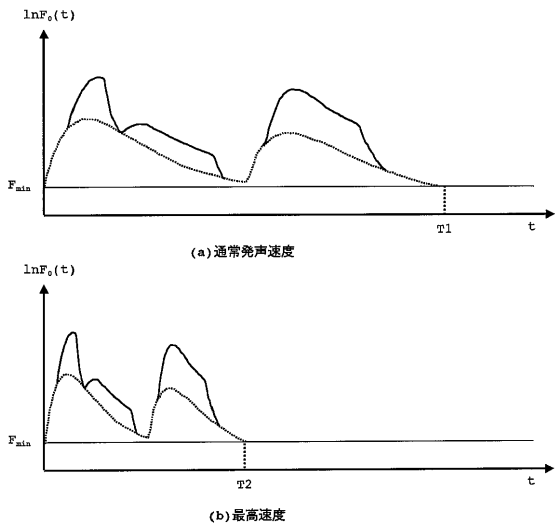
本発明の第 2 の変換の形態におけるヒッチボタン決定部の機能ブロック図

【図 8】



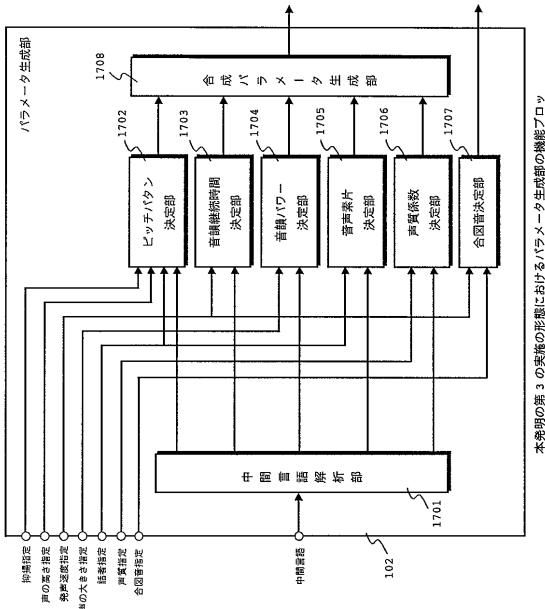
本発明の第 2 の変換の形態におけるピッチボタン生成フローチャ

【図9】



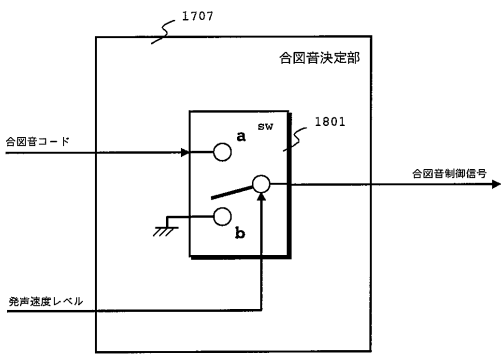
発声速度によるピッチパタンの違いの説明図

【図10】



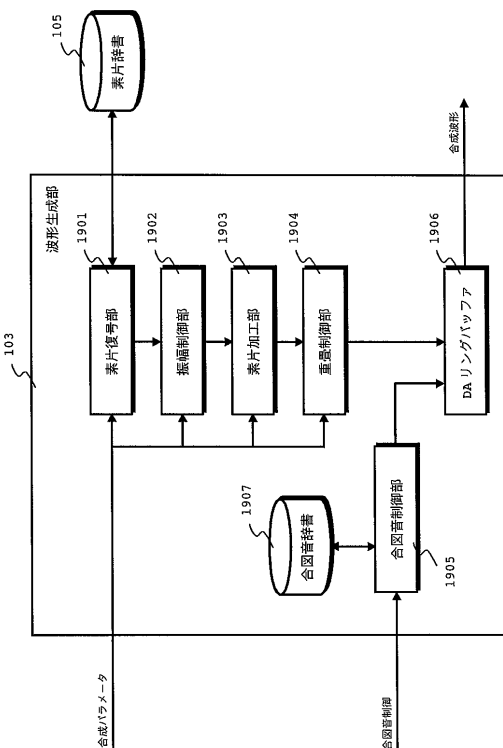
本発明の第3の実施の形態におけるパラメータ生成部の機能ブロック

【図11】



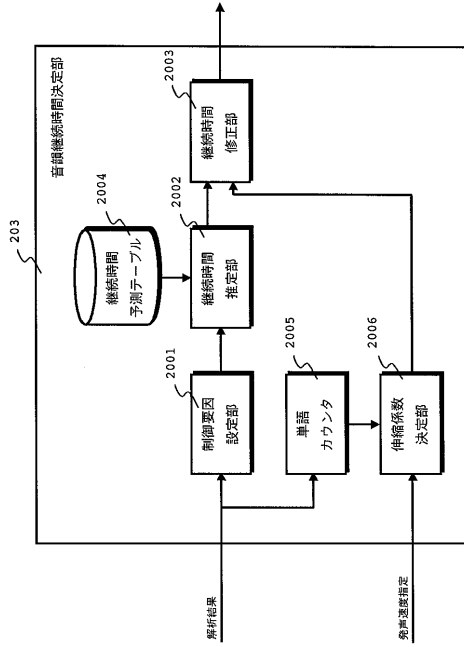
本発明の第3の実施の形態における合図音決定部の機能ブロック図

【図12】



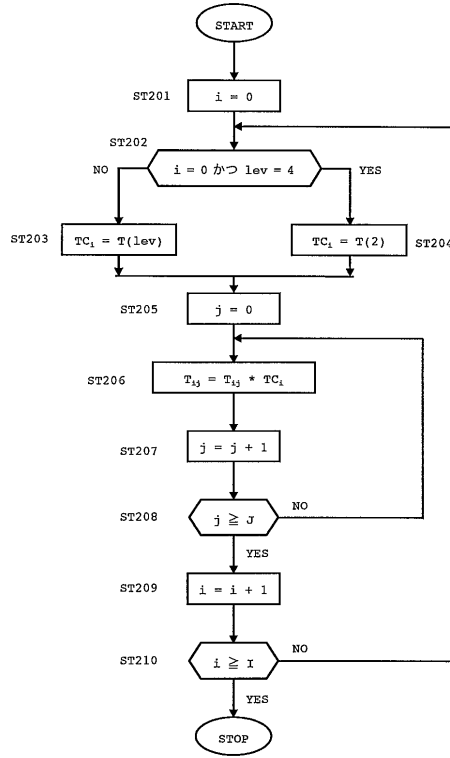
本発明の第3の実施の形態における波形生成部の機能ブロック

【図13】



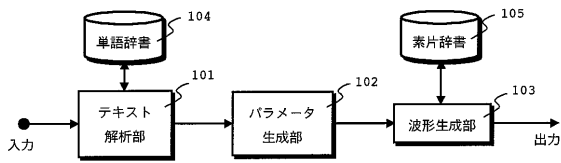
本発明の第4の実施の形態における音韻継続時間決定部の機能ブロック

【図14】



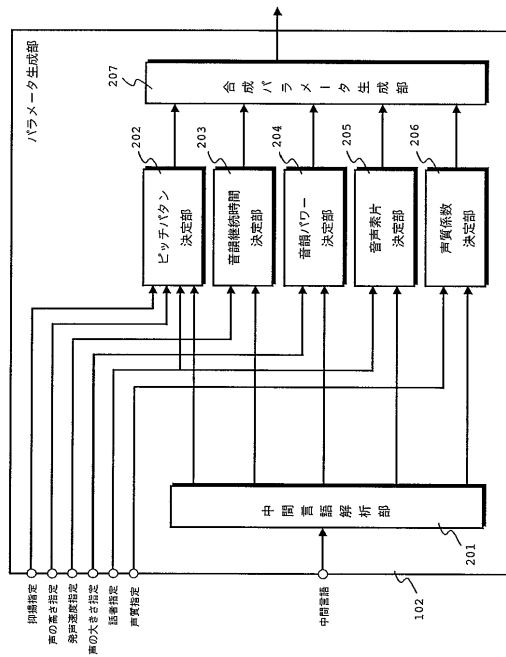
本発明の第4の実施の形態における継続時間決定フローチャート

【図15】



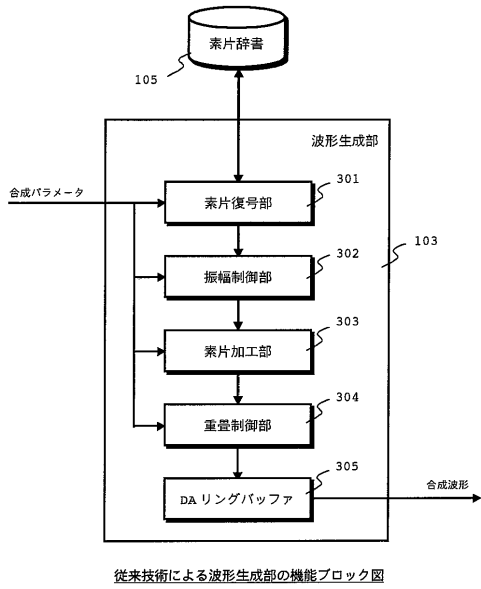
代表的なテキスト音声変換処理の機能ブロック図

【図16】

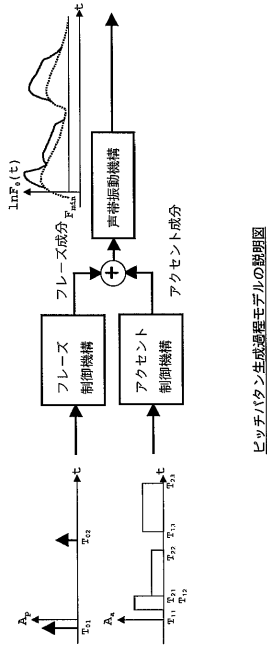


従来技術によるパラメータ生成部の機能ブロック図

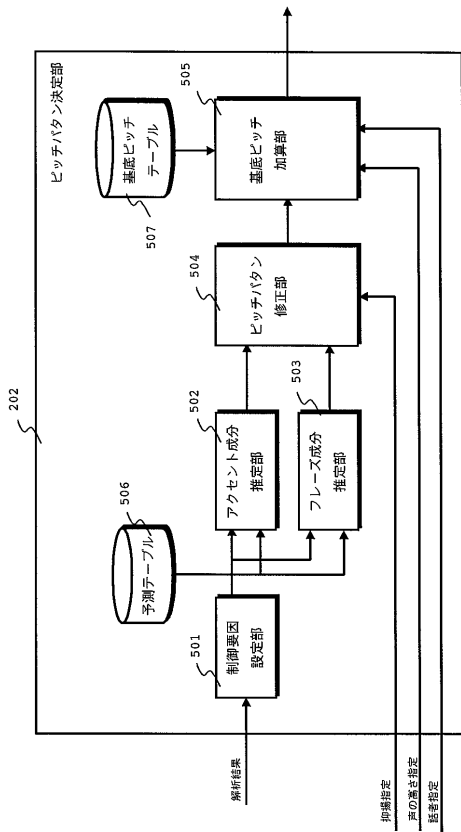
【図17】



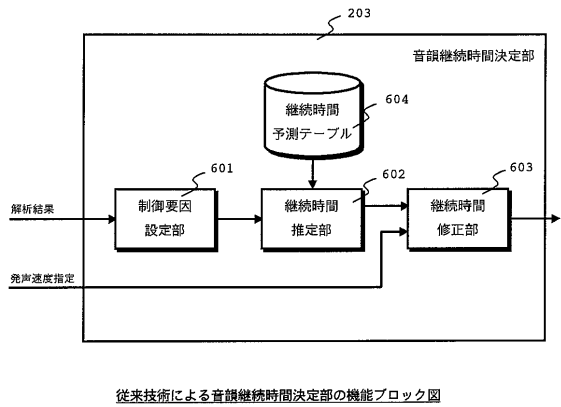
【図18】



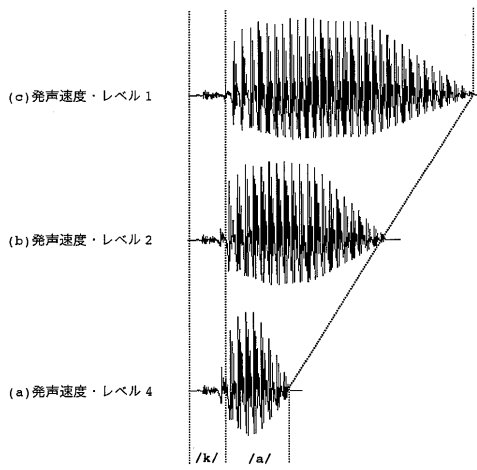
【図19】



【図20】



【 図 2 1 】



発声速度の違いによる波形伸縮の説明図

フロントページの続き

- (56)参考文献 特開平 1 1 - 1 6 7 3 9 8 (J P , A)
特開平 1 1 - 0 7 3 2 9 8 (J P , A)
特開平 0 2 - 1 9 5 3 9 7 (J P , A)
特開平 0 6 - 1 4 9 2 8 4 (J P , A)
特開 2 0 0 0 - 3 0 5 5 8 2 (J P , A)
特開平 0 9 - 1 7 9 5 7 7 (J P , A)
特開平 0 8 - 3 3 5 0 9 6 (J P , A)
実開昭 5 9 - 1 6 0 3 4 8 (J P , U)
特開 2 0 0 0 - 3 0 5 5 8 5 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)

G10L 13/00-13/08