

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3854713号

(P3854713)

(45) 発行日 平成18年12月6日(2006.12.6)

(24) 登録日 平成18年9月15日(2006.9.15)

(51) Int. Cl.

G 1 0 L 13/08 (2006.01)

F I

G 1 0 L 13/08 1 2 7 F

G 1 0 L 13/08 1 3 1 C

請求項の数 9 (全 16 頁)

(21) 出願番号	特願平10-57900	(73) 特許権者	000001007
(22) 出願日	平成10年3月10日(1998.3.10)		キヤノン株式会社
(65) 公開番号	特開平11-259095		東京都大田区下丸子3丁目30番2号
(43) 公開日	平成11年9月24日(1999.9.24)	(74) 代理人	100076428
審査請求日	平成16年5月27日(2004.5.27)		弁理士 大塚 康德
		(74) 代理人	100112508
			弁理士 高柳 司郎
		(74) 代理人	100115071
			弁理士 大塚 康弘
		(74) 代理人	100116894
			弁理士 木村 秀二
		(74) 代理人	100093908
			弁理士 松本 研一
		(74) 代理人	100101306
			弁理士 丸山 幸雄

最終頁に続く

(54) 【発明の名称】 音声合成方法および装置および記憶媒体

(57) 【特許請求の範囲】

【請求項1】

音韻系列に従って音声合成する音声合成装置であって、
 音韻の種類ごとに音韻時間長の標準偏差を格納する格納手段と、
 前記音韻系列の発声時間を示す発声時間情報を取得する発声時間取得手段と、
 前記音韻系列の各音韻に対応する第1の音韻時間長を取得する取得手段と、
 前記取得手段で取得した第1の音韻時間長の和を、前記発声時間から減じた値を、各音韻に対応する標準偏差の二乗和で割った値を係数とし、各音韻について、該係数と当該音韻の標準偏差の二乗との積を当該音韻の第1の音韻時間長に加えた値を第2の音韻時間長として設定する設定手段とを備えることを特徴とする音声合成装置。

10

【請求項2】

前記格納手段は更に音韻の種類ごとに音韻時間長の平均値を格納し、
 前記取得手段は、前記音韻系列の各音韻の前記平均値又は重回帰分析による音韻時間長推定値のいずれかをを用いて前記第1の音韻時間長を取得することを特徴とする請求項1記載の音声合成装置。

【請求項3】

前記取得手段は、各音韻の第1の音韻時間長として、当該音韻の平均値を中心に標準偏差の定数倍の範囲内に収まる音韻時間長を設定することを特徴とする請求項1記載の音声合成装置。

【請求項4】

20

前記格納手段は更に音韻の種類ごとに音韻時間長の最小値を格納し、

前記取得手段は、各音韻の第1の音韻時間長が当該音韻の前記最小値より小さい場合は該第1の音韻時間長を該最小値に設定することを特徴とする請求項1記載の音声合成装置。

【請求項5】

前記格納手段は各音韻の平均値、標準偏差、最小値を発声速度に基づいた分類毎に格納し、

前記取得手段は、前記音韻系列の発声時間から算出した発声速度に対応する各音韻の平均値、標準偏差、最小値を利用して各音韻の第1の音韻時間長を算出することを特徴とする請求項1乃至請求項4のいずれかに記載の音声合成装置。

10

【請求項6】

音声合成対象の文字系列を取得する文字系列取得手段と、

前記文字系列取得手段で取得した文字系列を音韻系列に変換する変換手段とを更に備え、

前記発声時間取得手段は、前記文字系列に含まれる発声速度を示す制御シーケンスに基づいて、前記発声時間情報を取得することを特徴とする請求項1記載の音声合成装置。

【請求項7】

音声合成対象の文字系列を取得する文字系列取得手段と、

前記文字系列取得手段で取得した文字系列を音韻系列に変換する変換手段とを更に備え、

前記発声時間取得手段は、ユーザによって設定された発声速度に基づいて、前記発声時間情報を取得することを特徴とする請求項1記載の音声合成装置。

20

【請求項8】

音韻系列に従って音声合成する音声合成方法であって、

前記音韻系列の発声時間を示す発声時間情報を取得する発声時間取得工程と、

前記音韻系列の各音韻に対応する第1の音韻時間長を取得する取得工程と、

音韻の種類ごとに音韻時間長の標準偏差を格納する格納手段から前記音韻系列の各音韻に対応する標準偏差を取得し、前記取得工程で取得した第1の音韻時間長の和を、前記発声時間から減じた値を、各音韻に対応する標準偏差の二乗和で割った値を係数とし、各音韻について、該係数と当該音韻の標準偏差の二乗との積を当該音韻の第1の音韻時間長に加えた値を第2の音韻時間長として設定する設定工程とを備えることを特徴とする音声合成方法。

30

【請求項9】

音韻系列に従って音声合成する音声合成方法をコンピュータに実行させるための制御プログラムを格納した記憶媒体であって、

前記音声合成方法が、

前記音韻系列の発声時間を示す発声時間情報を取得する発声時間取得工程と、

前記音韻系列の各音韻に対応する第1の音韻時間長を取得する取得工程と、

音韻の種類ごとに音韻時間長の標準偏差を格納する格納手段から前記音韻系列の各音韻に対応する標準偏差を取得し、前記取得工程で取得した第1の音韻時間長の和を、前記発声時間から減じた値を、各音韻に対応する標準偏差の二乗和で割った値を係数とし、各音韻について、該係数と当該音韻の標準偏差の二乗との積を当該音韻の第1の音韻時間長に加えた値を第2の音韻時間長として設定する設定工程とを備えることを特徴とする記憶媒体。

40

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、規則合成方式による音声合成方法および音声合成装置、および、音声合成方法を実装した、コンピュータが読むことができるプログラムを格納した記憶媒体に関する。

【0002】

50

【従来の技術】

従来の音声規則合成装置では、音韻時間長を制御する方法として、音韻時間長に関する統計量から導出した制御規則による方法（勾坂芳典、東倉洋一：“規則による音声合成のための音韻時間長制御”、電子通信学会論文誌、Vol.J67-A, No.7(1984)pp.629-636）、重回帰分析の一手法である数量化Ⅰ類を用いる方法（酒寄哲也、佐々木昭一、北川博雄：“規則合成のための数量化Ⅰ類を用いた韻律制御”、音響学会講演論文集、3-4-17(1986-10)）がある。

【0003】

【発明が解決しようとする課題】

しかしながら、上述した従来技術においては、音韻系列の発声時間を指定することが難しいという問題がある。たとえば、制御規則による方法では、指定された発声時間に対応した制御規則の導出が難しい。また、制御規則による方法で例外的な入力がある場合や数量化Ⅰ類を用いる方法で良い推定値が得られない場合に自然性を感じる音韻時間長に対する誤差が大きくなる、という問題がある。

10

【0004】

制御規則を用いて音韻時間長を制御する場合、統計量（平均値や標準偏差など）に対して前後の音韻の組み合わせを考慮した重み付けや、伸縮係数の設定などが必要になってくる。音韻の組み合わせの場合分けや、重み付けや伸縮係数などのパラメータなど操作する項目が多く、しかも、操作方法（制御規則）を経験則で決めていかなければならない。音韻系列の発声時間が指定されたときに、たとえ音韻の個数が同じでも、音韻の組み合わせは膨大になる。どのような音韻の組み合わせでも、音韻時間長の和が指定された発声時間に近くなるような、制御規則の導出は困難である。

20

【0005】

本発明は上記の問題点に鑑みてなされたものであり、指定した発声時間になるように音韻系列の音韻時間長を設定することを可能とし、発声時間の長短によらず自然な音韻時間長を与える音声合成方法および装置および記憶媒体を提供することを目的とする。

【0006】

【課題を解決するための手段】

上記の目的を達成するための本発明の一態様による音声合成装置は例えば以下の構成を備える。すなわち、

30

音韻系列に従って音声合成する音声合成装置であって、

音韻の種類ごとに音韻時間長の標準偏差を格納する格納手段と、

前記音韻系列の発声時間を示す発声時間情報を取得する発声時間取得手段と、

前記音韻系列の各音韻に対応する第1の音韻時間長を取得する取得手段と、

前記取得手段で取得した第1の音韻時間長の和を、前記発声時間から減じた値を、各音韻に対応する標準偏差の二乗和で割った値を係数とし、各音韻について、該係数と当該音韻の標準偏差の二乗との積を当該音韻の第1の音韻時間長に加えた値を第2の音韻時間長として設定する設定手段とを備える。

【0007】

また、本発明によれば、上記音声合成装置で実行される音声合成方法が提供される。更に、本発明によれば、上記音声合成方法をコンピュータに実現させるための制御プログラムを格納する記憶媒体が提供される。

40

【0008】

【発明の実施の形態】

以下、添付の図面を参照して本発明の好適な実施形態を説明する。

【0009】

〔第1の実施形態〕

図1は、第1の実施形態の音声合成装置の構成を示すブロック図である。101はCPUであり、本音声規則合成装置における各種制御を行なう。102はROMであり、各種パラメータやCPU101が実行する制御プログラムを格納する。103はRAMであり、

50

CPU101が実行する制御プログラムを格納するとともに、CPU101の作業領域を提供する。104はハードディスク、フロッピーディスク、CD-ROM等の外部記憶装置である。105は入力部であり、キーボード、マウス等から構成される。106はディスプレイであり、CPU101の制御により各種表示を行なう。6は音声合成部であり、合成音声を生成する。107はスピーカであり、音声合成部6より出力される音声信号(電気信号)を音声に変換して出力する。

【0010】

図2は、第1の実施形態による音声合成装置の機能構成を示すブロック図である。以下に示される各機能は、ROM102に格納された制御プログラムあるいは外部記憶装置104からRAM103にロードされた制御プログラムをCPU101が実行することによって実現される。

10

【0011】

1は文字系列入力部であり、入力部105より入力された合成すべき音声の文字系列、すなわち表音テキストの入力処理を行なう。例えば合成すべき音声「音声」であるときには、「おんせい」というような文字系列を入力する。また、この文字系列中には、発声速度や声の高さなどを設定するための制御シーケンスなどが含まれることもある。2は制御データ格納部であり、文字系列入力部1で制御シーケンスと判断された情報や、ユーザインタフェースより入力される発声速度や声の高さなどの制御データを内部レジスタに格納する。3は音韻系列生成部であり、文字系列入力部1より入力された文字系列を音韻系列へ変換する。例えば、「おんせい」という文字系列は、「o, X, s, e, i」という音韻系列へ変換される。4は音韻系列格納部であり、音韻系列生成部3で生成された音韻系列を内部レジスタに格納する。なお、上述の各内部レジスタとしてはRAM103を用いることが可能である。

20

【0012】

5は音韻時間長設定部であり、制御データ格納部2に格納された制御データの発声速度と音韻系列格納部4に格納された音韻の種類より、音韻時間長を設定する。6は音声合成部であり、音韻時間長設定部5で音韻時間長の設定された音韻系列と制御データ格納部2に格納された制御データの声の高さから、合成音声を生成する。

【0013】

次に、音韻時間長設定部5で行なわれる音韻時間長の設定について説明する。以下の説明において、音韻集合を Φ とする。 の例としては、

30

$$\Phi = \{ a, e, i, o, u, X(\text{撥音}), b, d, g, m, n, r, w, y, z, ch, f, h, k, p, s, sh, t, ts, Q(\text{促音}) \}$$

などを使用することができる。

【0014】

また、音韻時間長設定区間を呼気段落(ポーズとポーズの間の区間)とする。さて、音韻時間長設定区間の音韻系列 $i(1 \leq i \leq N)$ を、制御データ格納部2に格納された制御データの発声速度によって決定される発声時間Tで発声するように、当該音韻系列の各音韻 i の音韻時間長 d_i を決定する。すなわち、音韻系列の各 i (式(1a))の音韻時間長 d_i (式(1b))を、式(1c)を満足するように決定する。

【0015】

40

【数1】

$$\alpha_i \in \Omega \quad (1 \leq i \leq N) : \quad (1 a)$$

$$d_i \quad (1 \leq i \leq N) \quad (1 b)$$

$$T = \sum_{i=1}^N d_i \quad (1 c)$$

10

【 0 0 1 6 】

ここで、音韻 i の音韻時間長初期値を d_{i0} とする。また、音韻 i に関して、音韻時間長の平均、標準偏差、最小値をそれぞれ μ_i , σ_i , $d_{i \min}$ とする。そして、これらの値を用いて、以下に示す式 (2) に従って d_i を決定し、これを新たな音韻時間長初期値とする。すなわち、音韻時間長の平均値、標準偏差値、最小値を音韻の種類毎 (i 毎) に求め、これをメモリに格納しておき、これらの値を用いて音韻時間長の初期値を決定しなおす。

【 0 0 1 7 】

20

【 数 2 】

$$d_{\alpha_i} = \begin{cases} \max(\mu_{\alpha_i} - 3\sigma_{\alpha_i}, d_{\alpha_i \min}) & (d_{\alpha_i 0} < \max(\mu_{\alpha_i} - 3\sigma_{\alpha_i}, d_{\alpha_i \min})) \\ d_{\alpha_i 0} & (\max(\mu_{\alpha_i} - 3\sigma_{\alpha_i}, d_{\alpha_i \min}) \leq d_{\alpha_i 0} \leq \mu_{\alpha_i} + 3\sigma_{\alpha_i}) \\ \mu_{\alpha_i} + 3\sigma_{\alpha_i} & (\mu_{\alpha_i} + 3\sigma_{\alpha_i} < d_{\alpha_i 0}) \end{cases} \quad (2)$$

【 0 0 1 8 】

こうして得られた音韻時間長初期値 d_i を用いて、音韻時間長 d_i を式 (3 a) に従って設定する。なお、得られた d_i が閾値 θ_i (> 0) に対して $d_i < \theta_i$ となるときは、式 (3 b) に従って設定される。

30

【 0 0 1 9 】

【 数 3 】

$$d_i = d_{\alpha_i} + \rho(\sigma_{\alpha_i})^2$$

$$\rho = \frac{(T - \sum_{i=1}^N d_{\alpha_i})}{\sum_{i=1}^N (\sigma_{\alpha_i})^2} \quad (3 a)$$

40

$$d_i = \theta_i \quad (3 b)$$

【 0 0 2 0 】

すなわち、更新された音韻時間長の初期値の和を設定された発声時間 T から差引き、これを音韻時間長の標準偏差 σ_i の二乗和で割った値を係数 ρ とし、この係数 ρ と標準偏差 σ_i の二乗との積を当該音韻時間長の初期値 d_i に加えた値を、音韻時間長 d_i とする。

【 0 0 2 1 】

50

以上の動作を、図 3 のフローチャートを参照して説明する。

【 0 0 2 2 】

まず、ステップ S 1 で、文字系列入力部 1 より表音テキストが入力される。ステップ S 2 では、外部入力された制御データ（発声速度、声の高さ）と入力された表音テキスト中の制御データが制御データ格納部 2 に格納される。ステップ S 3 で、文字系列入力部 1 より入力された表音テキストから音韻系列生成部 3 において音韻系列が生成される。

【 0 0 2 3 】

次に、ステップ S 4 で、次の時間長設定区間の音韻系列が音韻系列格納部 4 に取り込まれる。ステップ S 5 で、音韻時間長設定部 5 において、音韻 i の種類に応じて音韻時間長初期値 d_i が設定される（式（2））。ステップ S 6 では、まず、制御データ格納部 2 に格納された制御データの発声速度から音韻時間長設定区間の発声時間 T を設定する。そして、音韻時間長設定区間の音韻系列の音韻時間長の和が音韻時間長設定区間の発声時間 T に等しくなるように、上記式（3a）、（3b）を用いて、音韻時間長設定区間の音韻系列の各音韻時間長を設定する。

10

【 0 0 2 4 】

ステップ S 7 で、音韻時間長設定部 5 で音韻時間長の設定された音韻系列と制御データ格納部 2 に格納された制御データの声の高さから、合成音声が生産される。そして、ステップ S 8 で、入力された文字列に対する最後の音韻時間長設定区間であるか否かが判別され、最後の音韻時間長設定区間でない場合はステップ S 10 で外部入力された制御データが制御データ格納部 2 に格納されてステップ S 4 に戻り、処理が続けられる。

20

【 0 0 2 5 】

一方、ステップ S 8 で最後の音韻時間長設定区間であると判定された場合はステップ S 9 に進み、入力が終了したか否かが判別される。入力が終了していない場合はステップ S 1 に戻り、上記処理が繰り返される。

【 0 0 2 6 】

なお、式（2）は、音韻時間長初期値が現実にはあり得ないような値や出現確率の低い値に設定されるのを防ぐためのものである。音韻時間長の確率密度が正規分布であると仮定したときに、平均値から標準偏差の ± 3 倍以内に入る確率は 0.996 となる。更に、音韻時間長が短くなりすぎるのを防ぐために、標本集団の最小値未満にはならないようにしている。

30

【 0 0 2 7 】

式（3a）は、式（2）で設定された音韻時間長初期値を平均値とする正規分布が各音韻時間長の確率密度関数であると仮定して、式（1c）の制約条件のもとで最尤推定（maximum likelihood estimation）を行った結果である。本例の最尤推定について説明すると次のとおりである。

【 0 0 2 8 】

音韻 i の音韻時間長の標準偏差を σ_i とする。音韻時間長の確率密度分布が正規分布であると仮定する（式（4a））。このとき、音韻時間長の対数尤度は式（4b）のようになる。ここで、対数尤度を最大にするのは、式（4c）の K を最小にするのと同値である。そこで、音韻時間長の対数尤度が最大になるように上述の式（1c）を満たす d_i を決定する。

40

【 0 0 2 9 】

【 数 4 】

$$P_{\alpha_i}(d_i) = \left(\sqrt{2\pi}\sigma_{\alpha_i}\right)^{-1} \exp\left(-\frac{(d_i - d_{\alpha_i})^2}{2(\sigma_{\alpha_i})^2}\right) \quad (4 a)$$

$$\begin{aligned} \log(L(d_i)) &= \log\left(\prod_{i=1}^N P_{\alpha_i}(d_i)\right) \\ &= -\sum_{i=1}^N \log\left(\sqrt{2\pi}\sigma_{\alpha_i}\right) - \frac{1}{2} \sum_{i=1}^N \frac{(d_i - d_{\alpha_i})^2}{(\sigma_{\alpha_i})^2} \end{aligned} \quad (4 b)$$

10

$$K = \sum_{i=1}^N \frac{(d_i - d_{\alpha_i})^2}{(\sigma_{\alpha_i})^2} \quad (4 c)$$

【 0 0 3 0 】

今、式 (5 a) のように変数変換を行うと、式 (4 c) 及び式 (1 c) は式 (5 b) 及び (5 c) のようになる。K が最小となるのは、球 (式 5 b) が平面 (式 (5 c)) に接するときであり、式 (5 d) の場合である。この結果、式 (3 a) が導かれる。

20

【 0 0 3 1 】

【 数 5 】

$$\rho_i = \frac{d_i - d_{\alpha_i}}{\sigma_{\alpha_i}} \quad (5 a)$$

$$K = \sum_{i=1}^N \rho_i^2 \quad (5 b)$$

30

$$\sum_{i=1}^N \rho_i \sigma_{\alpha_i} = T - \sum_{i=1}^N d_{\alpha_i} \quad (5 c)$$

$$\begin{aligned} \rho_i &= \rho \sigma_{\alpha_i} \\ \rho &= \frac{\left(T - \sum_{i=1}^N d_{\alpha_i}\right)}{\sum_{i=1}^N (\sigma_{\alpha_i})^2} \end{aligned} \quad (5 d)$$

40

【 0 0 3 2 】

式 (2) と式 (3 a)、(3 b) を総合して、自然発声の標本集団から求めた統計量 (平均値、標準偏差、最小値) を用いて、所望の発声時間 ((1 c) 式) を満たす最も確からしい (尤度が最大になる) 値に音韻時間長が設定される。したがって、所望の発声時間 ((1 c) 式) を満たすように自然発声したときに得られる音韻時間長に対する誤差が小さい、という意味で自然な音韻時間長が得られる。

【 0 0 3 3 】

[第 2 の実施形態]

第 1 の実施形態では、発声速度 (発声時間) や音韻のカテゴリにかかわらず、各音韻 i の音韻時間長 d_i を同一の規則で決定した。第 2 の実施形態では、発声速度や音韻のカテ

50

ゴリに応じて音韻時間長 d_i の決定規則を変化させ、より自然な音声合成を可能とする。
 なお、第 2 の実施形態によるハードウェア構成、機能構成は第 1 の実施形態（図 1、図 2）と同様である。

【 0 0 3 4 】

音韻 i に関して、発声速度でカテゴリーを分けて音韻時間長の平均値、標準偏差、最小値を求める。例えば、発声速度のカテゴリーを呼気段落の平均モーラ時間長で表すとして、

- 1 : 1 2 0 ミリ秒未満、
- 2 : 1 2 0 ミリ秒以上 1 4 0 ミリ秒未満、
- 3 : 1 4 0 ミリ秒以上 1 6 0 ミリ秒未満、
- 4 : 1 6 0 ミリ秒以上 1 8 0 ミリ秒未満、
- 5 : 1 8 0 ミリ秒以上

10

とする。なお、上述した項目の先頭の数字を発声速度に対応するカテゴリーのインデックスとする。発声速度に対応するカテゴリーのインデックスを n として音韻時間長の平均値、標準偏差、最小値を求め、それぞれ $\mu_i(n)$, $\sigma_i(n)$, $d_{imin}(n)$ とする。

【 0 0 3 5 】

音韻 i の音韻時間長初期値を d_{i0} とする。音韻時間長初期値 d_{i0} を平均値によって決定する音韻の集合を a 、重回帰分析の一手法である数量化 I 類（質的なデータから量的に測定される外的基準を予測したり、説明したりするための手法）によって決定する音韻の集合を r とする。ここで、 a の要素で、 r のどちらにも含まれない要素や、

20

【 0 0 3 6 】

【数 6】

$$\Omega_a \cup \Omega_r = \Omega \quad (6a)$$

$$\Omega_a \cap \Omega_r = \emptyset \quad (6b)$$

30

【 0 0 3 7 】

$i \in a$ のとき、すなわち i が a に属するときは、平均値によって音韻時間長初期値を決定する。すなわち、音声速度に対応するカテゴリーのインデックス n を求めて、以下の式（7）によって音韻時間長初期値を決定する。

【 0 0 3 8 】

【数 7】

$$d_{\alpha_i,0} = \mu_{\alpha_i}(n) \quad (7)$$

40

【 0 0 3 9 】

一方、 $i \in r$ のとき、すなわち i が r に属するときは、数量化 I 類によって音韻時間長初期値を決定する。ここで、要因のインデックスを j ($1 \leq j \leq J$)、各要因に対応するカテゴリーのインデックスを k ($1 \leq k \leq K(j)$) として、 (j, k) に対応する数量化 I 類の係数を、

a_{jk}

とする。

【 0 0 4 0 】

50

要因の一例として、

- 1：当該音韻の2つ前の先行音韻、
- 2：当該音韻の1つ前の先行音韻、
- 3：当該音韻、
- 4：当該音韻の1つ後の後続音韻、
- 5：当該音韻の2つ後の後続音韻、
- 6：呼気段落の平均モーラ時間長、
- 7：呼気段落内モーラ位置、
- 8：当該音韻を含む単語の品詞

などを使用することができる。上述した項目の先頭の数字が要因のインデックス j に対応する。 10

【0041】

さらに、各要因に対応するカテゴリーの例を述べる。音韻のカテゴリーは、

- 1：a、2：e、3：i、4：o、5：u、6：X、7：b、8：d、9：g、10：m、
- 11：n、12：r、13：w、14：y、15：z、16：+、17：c、18：f、19：h、20：k、
- 21：p、22：s、23：sh、24：t、25：ts、26：Q、27：ポーズ、とし、当該音韻のみ“ポーズ”をはずす。実施形態において、呼気段落を音韻時間長設定区間としているが、呼気段落はポーズを含まないので、当該音韻からポーズをはずす。なお、呼気段落という用語は、ポーズ（または文頭）とポーズ（または文末）の間の区間で、途中にポーズを含まないものという意味で使用している。 20

【0042】

また、呼気段落内の平均モーラ時間長のカテゴリは、

- 1：120ミリ秒未満
- 2：120ミリ秒以上140ミリ秒未満
- 3：140ミリ秒以上160ミリ秒未満
- 4：160ミリ秒以上180ミリ秒未満
- 5：180ミリ秒以上

とする。

【0043】

また、呼気段落内モーラ位置に関しては、 30

- 1：第1モーラ
- 2：第2モーラ
- 3：第3モーラ以降最後から第3番目のモーラまで
- 4：最後から2番目のモーラ
- 5：最後のモーラ

とする。

【0044】

更に、品詞のカテゴリーを

- 、1：名詞、2：副詞的名詞、3：代名詞、4：固有名詞、5：数、6：動詞、7：形容詞、8：形容動詞、9：副詞、10：連体詞、11：接続詞、12：感動詞、13：助動詞、14：格助詞、15：副助詞、16：並立助詞、17：準体助詞、18：接続助詞、19：終助詞、20：接頭辞、21：接尾辞、22：形動接尾、23：サ変接尾、24：形容詞接尾、25：動詞接尾、26：助数詞 40

とする。

【0045】

なお、要因（アイテムともいう）とは、数量化Ⅰ類での予測に使用する質的なデータの種類を意味する。カテゴリーは、各要因毎に取りうる選択肢を意味する。したがって、上記の例に即して説明すると、次のようになる。

【0046】

要因のインデックス j = 1：当該音韻の2つ前の先行音韻 50

インデックス $k = 1$ に対応するカテゴリー : a
 インデックス $k = 2$ に対応するカテゴリー : e
 インデックス $k = 3$ に対応するカテゴリー : i
 インデックス $k = 4$ に対応するカテゴリー : o

(中略)

インデックス $k = 26$ に対応するカテゴリー : Q
 インデックス $k = 27$ に対応するカテゴリー : ポーズ。

【0047】

要因のインデックス $j = 2$: 当該音韻の1つ前の先行音韻

インデックス $k = 1$ に対応するカテゴリー : a
 インデックス $k = 2$ に対応するカテゴリー : e
 インデックス $k = 3$ に対応するカテゴリー : i
 インデックス $k = 4$ に対応するカテゴリー : o

(中略)

インデックス $k = 26$ に対応するカテゴリー : Q
 インデックス $k = 27$ に対応するカテゴリー : ポーズ。

【0048】

要因のインデックス $j = 3$: 当該音韻

インデックス $k = 1$ に対応するカテゴリー : a
 インデックス $k = 2$ に対応するカテゴリー : e
 インデックス $k = 3$ に対応するカテゴリー : i
 インデックス $k = 4$ に対応するカテゴリー : o

(中略)

インデックス $k = 26$ に対応するカテゴリー : Q
 インデックス $k = 27$ に対応するカテゴリー : ポーズ。

【0049】

要因のインデックス $j = 4$: 当該音韻の1つ後の後続音韻

インデックス $k = 1$ に対応するカテゴリー : a
 インデックス $k = 2$ に対応するカテゴリー : e
 インデックス $k = 3$ に対応するカテゴリー : i
 インデックス $k = 4$ に対応するカテゴリー : o

(中略)

インデックス $k = 26$ に対応するカテゴリー : Q
 インデックス $k = 27$ に対応するカテゴリー : ポーズ。

【0050】

要因のインデックス $j = 5$: 当該音韻の2つ後の後続音韻

インデックス $k = 1$ に対応するカテゴリー : a
 インデックス $k = 2$ に対応するカテゴリー : e
 インデックス $k = 3$ に対応するカテゴリー : i
 インデックス $k = 4$ に対応するカテゴリー : o

(中略)

インデックス $k = 26$ に対応するカテゴリー : Q
 インデックス $k = 27$ に対応するカテゴリー : ポーズ。

【0051】

要因のインデックス $j = 6$: 呼気段落内の平均モーラ時間長

インデックス $k = 1$ に対応するカテゴリー : 120ミリ秒未満
 インデックス $k = 2$ に対応するカテゴリー : 120ミリ秒以上140ミリ秒未満
 インデックス $k = 3$ に対応するカテゴリー : 140ミリ秒以上160ミリ秒未満
 インデックス $k = 4$ に対応するカテゴリー : 160ミリ秒以上180ミリ秒未満
 インデックス $k = 5$ に対応するカテゴリー : 180ミリ秒以上。

10

20

30

40

50

【 0 0 5 2 】

要因のインデックス $j = 7$: 呼気段落内モーラ位置

インデックス $k = 1$ に対応するカテゴリー : 第 1 モーラ

インデックス $k = 2$ に対応するカテゴリー : 第 2 モーラ

(中略)

インデックス $k = 5$ に対応するカテゴリー : 最後のモーラ。

【 0 0 5 3 】

要因のインデックス $j = 8$: 当該音韻を含む単語の品詞

インデックス $k = 1$ に対応するカテゴリー : 名詞

インデックス $k = 2$ に対応するカテゴリー : 副詞的名詞

10

(中略)

インデックス $k = 26$ に対応するカテゴリー : 助数詞

となる。

【 0 0 5 4 】

上述した項目の先頭の数字がカテゴリーのインデックス k に対応する。

【 0 0 5 5 】

そして、各要因ごとに係数 a_{jk} の平均が 0 になるようにする。すなわち、式 (8) を満足するようにする。

【 0 0 5 6 】

【 数 8 】

20

$$\sum_{k=1}^{K(j)} a_{jk} = 0 \quad (1 \leq j \leq J) \quad (8)$$

【 0 0 5 7 】

また、音韻 i のダミー変数を、以下のように設定する。

【 0 0 5 8 】

【 数 9 】

30

$$\delta_i(j, k) = \begin{cases} 1 & (\text{音韻 } \alpha_i \text{ が要因 } j \text{ のカテゴリー } k \text{ に反応するとき}) \\ 0 & (\text{上記以外}) \end{cases} \quad (9)$$

【 0 0 5 9 】

係数とダミー変数の積和に加える定数を c_0 とする。このとき、音韻 i の音韻時間長の数量化 I 類による推定値は、式 (10) となる。

【 0 0 6 0 】

【 数 10 】

40

$$\hat{d}_{\alpha_i} = \sum_{j=1}^J \sum_{k=1}^{K(j)} a_{jk} \delta_i(j, k) + c_0 \quad (10)$$

【 0 0 6 1 】

そして、この推定値を用いて音韻 i の音韻時間長初期値を以下のように決定する。

【 0 0 6 2 】

【 数 11 】

$$d_{\alpha,0} = \hat{d}_{\alpha,i} \quad (11)$$

【0063】

さらに、発声速度と対応するカテゴリーのインデックス n を求めて、当該カテゴリーの音韻時間長の平均値、標準偏差、最小値を得て、これらを用いて音韻時間長初期値 d_{i0} を以下の式で更新する。こうして得られた d_{i0} を改めて音韻時間長初期値として設定する。

【0064】

10

【数12】

$$d_{\alpha,i} = \begin{cases} \max(\mu_{\alpha,i}(n) - r_{\sigma}\sigma_{\alpha,i}(n), d_{\alpha,i,\min}(n)) & (d_{\alpha,0} < \max(\mu_{\alpha,i}(n) - r_{\sigma}\sigma_{\alpha,i}(n), d_{\alpha,i,\min}(n))) \\ d_{\alpha,0} & (\max(\mu_{\alpha,i}(n) - r_{\sigma}\sigma_{\alpha,i}(n), d_{\alpha,i,\min}(n)) \leq d_{\alpha,0} \leq \mu_{\alpha,i}(n) + r_{\sigma}\sigma_{\alpha,i}(n)) \\ \mu_{\alpha,i}(n) + r_{\sigma}\sigma_{\alpha,i}(n) & (\mu_{\alpha,i}(n) + r_{\sigma}\sigma_{\alpha,i}(n) < d_{\alpha,0}) \end{cases} \quad (12)$$

【0065】

ここで、式中の標準偏差に掛ける係数の r は、例えば、 $r = 3$ とする。以上のようにして得られた音韻時間長初期値を用いて、第1の実施形態と類似の方法で音韻時間長を決定する。すなわち、以下の式(13a)を用いて音韻時間長 d_i を決定し、閾値 θ_i (> 0) に対して $d_i < \theta_i$ となるときは、式(13b)により音韻時間長 d_i を決定する。

20

【0066】

【数13】

$$d_i = d_{\alpha,i} + \rho(\sigma_{\alpha,i}(n))^2$$

$$\rho = \frac{(T - \sum_{i=1}^N d_{\alpha,i})}{\sum_{i=1}^N (\sigma_{\alpha,i}(n))^2} \quad (13a)$$

30

$$d_i = \theta_i \quad (13b)$$

【0067】

以上の動作を、図3のフローチャートを流用して説明する。ステップS1で、文字系列入力部1より表音テキストが入力される。ステップS2で、外部入力された制御データ(発声速度、音の高さ)と入力された表音テキスト中の制御データが制御データ格納部2に格納される。ステップS3で、文字系列入力部1より入力された表音テキストから音韻系列生成部3において音韻系列が生成される。ステップS4で、次の音韻時間長設定区間の音韻系列が音韻系列格納部4に取り込まれる。

40

【0068】

ステップS5では、音韻時間長設定部5において、制御データ格納部2に格納された制御データの発声速度、音韻時間長の平均値と標準偏差と最小値、および、数量化I類による音韻時間長推定値を用いて、上述した方法により、音韻の種類(カテゴリ)に応じて音韻時間長初期値が設定される。

【0069】

ステップS6では、音韻時間長設定部5において、制御データ格納部2に格納された制御データの発声速度から音韻時間長設定区間の発声時間を設定し、音韻時間長設定区間の音韻系列の音韻時間長の和が音韻時間長設定区間の発声時間に等しくなるように、音韻時間長設定区間の音韻系列の音韻時間長を上述した方法により設定する。

50

【 0 0 7 0 】

ステップ S 7 で、音韻時間長設定部 5 で音韻時間長の設定された音韻系列と制御データ格納部 2 に格納された制御データの声の高さから、合成音声が生産される。ステップ S 8 で、入力された文字列に対する最後の音韻時間長設定区間であるか否かが判別される。最後の音韻時間長設定区間でない場合はステップ S 1 0 へ進む。ステップ S 1 0 では、外部入力された制御データが制御データ格納部 2 に格納されてステップ S 4 に戻り、処理が続けられる。一方、最後の音韻時間長設定区間である場合はステップ S 9 に進み、入力が終了したか否かが判別され、終了していない場合はステップ S 1 に戻り、処理が続けられる。

【 0 0 7 1 】

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

【 0 0 7 2 】

(1) 上述した各実施形態において音韻集合 は一例であり、それ以外の集合も使用でき、言語や音韻の種類に応じて音韻集合の要素を決めることができる。また、本発明は日本語以外の言語にも適用可能である。

【 0 0 7 3 】

(2) 上述した実施形態において、呼気段落は音韻時間長設定区間の一例であり、他にも、単語、形態素、文節、文などを音韻時間長設定区間とすることができる。なお、文を音韻時間長設定区間とするときは、当該音韻のポーズを考慮する必要がある。

【 0 0 7 4 】

(3) 上述した実施形態において、音韻時間長の初期値として設定する値として、自然発声した音声の音韻時間長を使用することができる。また、他の音韻時間長制御規則によって決定した値や数量化 I 類を用いて推定した値を使用することもできる。

【 0 0 7 5 】

(4) 上述した第 2 の実施形態において、音韻時間長の平均値を求めるのに使用する発声速度のカテゴリーは一例を示すものであり、他のカテゴリーを用いても良い。

【 0 0 7 6 】

(5) 上述した第 2 の実施形態において、数量化 I 類の要因とカテゴリーは一例を示すものであり、他の要因やカテゴリーを用いても良い。

【 0 0 7 7 】

(6) 上述した実施形態において、音韻時間長初期値の設定に使用する標準偏差に掛ける係数 $r = 3$ は、一例を示すものであり、他の値を用いてもよい。

【 0 0 7 8 】

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

【 0 0 7 9 】

この場合、記憶媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記憶した記憶媒体は本発明を構成することになる。

【 0 0 8 0 】

プログラムコードを供給するための記憶媒体としては、例えば、フロッピーディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモリカード、ROMなどを用いることができる。

【 0 0 8 1 】

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

10

20

30

40

50

うまでもない。

【 0 0 8 2 】

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

【 0 0 8 3 】

【 発明の効果 】

以上説明したように、本発明によれば、指定した発声時間になるように音韻系列の音韻時間長を設定することが可能となり、発声時間の長短によらず自然な音韻時間長を与えることが可能である。

【 0 0 8 4 】

【 図面の簡単な説明 】

【 図 1 】 本発明の実施形態に係る音声合成装置の構成を示すブロック図である。

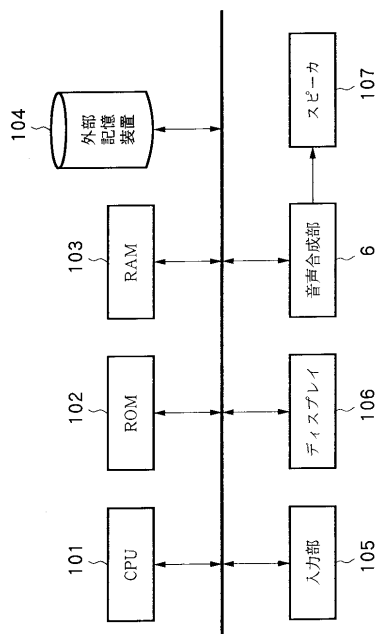
【 図 2 】 本発明の実施形態に係る音声合成装置の機能構成を示すブロック図である。

【 図 3 】 本発明の実施形態に係る音声合成手段を示すフローチャートである。

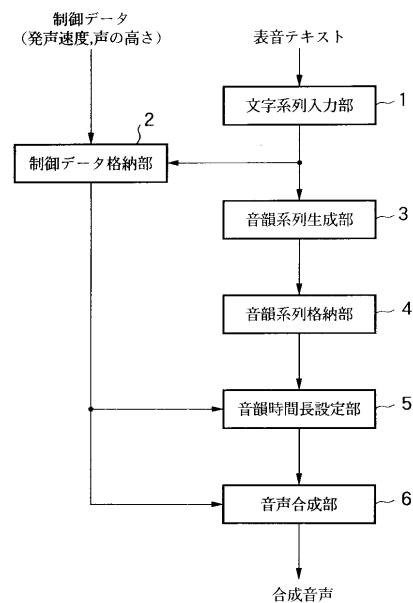
【 符号の説明 】

- 1 文字系列入力部
- 2 制御データ格納部
- 3 音韻系列生成部
- 4 音韻系列格納部
- 5 音韻時間長設定部
- 6 音声合成部

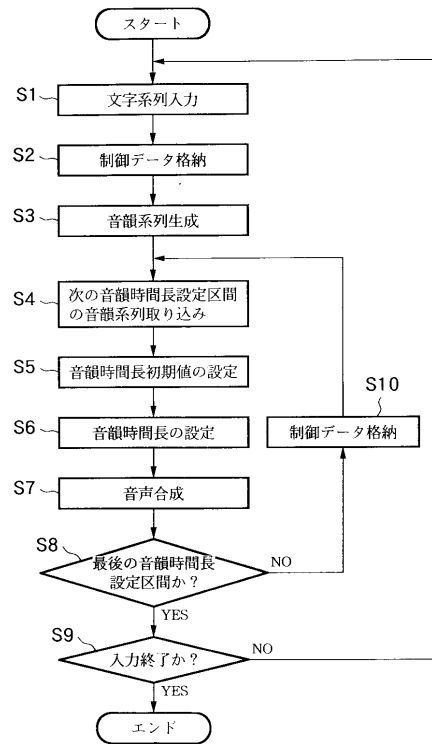
【 図 1 】



【 図 2 】



【図 3】



フロントページの続き

(72)発明者 大塚 充
東京都大田区下丸子3丁目30番2号 キヤノン株式会社内

審査官 荏原 雄一

(56)参考文献 特開平10-039896(JP,A)
特開平01-255899(JP,A)
特開平04-170600(JP,A)
規則合成のための数量化I類を用いた韻律制御,音講論,昭和61年秋,3-4-17

(58)調査した分野(Int.Cl.,DB名)
G10L 13/08
JSTPlus(JDream2)