

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号
特許第4489371号
(P4489371)

(45) 発行日 平成22年6月23日(2010. 6. 23)

(24) 登録日 平成22年4月9日(2010. 4. 9)

(51) Int.Cl.

F I

G 1 O L 13/02 (2006. 01)

G 1 O L 13/02 1 2 3

G 1 O L 19/06 (2006. 01)

G 1 O L 19/06 B

請求項の数 15 (全 14 頁)

(21) 出願番号	特願2003-125758 (P2003-125758)	(73) 特許権者	392026693
(22) 出願日	平成15年4月30日 (2003. 4. 30)		株式会社エヌ・ティ・ティ・ドコモ
(65) 公開番号	特開2003-323200 (P2003-323200A)		東京都千代田区永田町二丁目11番1号
(43) 公開日	平成15年11月14日 (2003. 11. 14)	(74) 代理人	100098084
審査請求日	平成18年4月28日 (2006. 4. 28)		弁理士 川▲崎▼ 研二
(31) 優先権主張番号	10/134281	(72) 発明者	ワイ チュウ
(32) 優先日	平成14年4月29日 (2002. 4. 29)		アメリカ合衆国、カリフォルニア州 9 5
(33) 優先権主張国	米国 (US)		1 1 2、サンノゼ、1 7 0 0、ノース・フ
			ァースト・ストリート 1 3 0
		(72) 発明者	コズロウ ラシュキャリ
			アメリカ合衆国、カリフォルニア州 9 4
			5 3 9、フリーモント、1 5 2 5、サラマ
			ンカ コート
		審査官	山下 剛史
			最終頁に続く

(54) 【発明の名称】 合成音声を最適化する方法、音声合成フィルタを生成する方法、音声最適化方法及び音声最適化装置

(57) 【特許請求の範囲】

【請求項 1】

合成音声を最適化する方法において、
原音サンプルを使用して音声合成多項式の第 1 の係数を計算して、第 1 合成音声サンプルを生成する過程と、
前記原音サンプルと前記第 1 合成音声サンプルとの間の第 1 誤差を計算する過程と、
前記第 1 誤差の勾配を計算する過程と、
前記勾配を使用して前記音声合成多項式の第 2 の係数を計算し、第 2 合成音声サンプルを生成する過程と、
前記原音サンプルと前記第 2 合成音声サンプルとの間の第 2 誤差を計算する過程と、
前記第 2 誤差が前記第 1 誤差より小さい場合は前記の第 2 の係数を選択し、前記第 2 誤差が前記第 1 誤差より大きい場合は前記第 1 の係数を選択する過程と、
システムの不安定さに関して前記第 2 の誤差をテストし、前記テスト結果が不安定であるとなった場合は前記第 1 係数を選択するテスト過程と
を有することを特徴とする方法。

【請求項 2】

請求項 1 に記載の方法において、
前記勾配に適用されたステップサイズを使用して、前記第 2 の係数が繰り返し計算され、
各繰り返しにおいて誤差が計算され、各繰り返しでの前記誤差が減少しなくなったら繰

り返しを終える

ことを特徴とする方法。

【請求項 3】

請求項 2 に記載の方法において、

前記第 2 の誤差の勾配を計算する過程と、

前記第 2 の誤差の勾配を使用して前記音声合成多項式の第 3 の係数を計算して、第 3 の合成音声サンプルを生成する過程と

を更に有し、

前記第 3 の係数は、前記第 2 の誤差の勾配に適用されるステップサイズを使用して繰り返し計算され、

各繰り返しにおいて 1 つの誤差が計算され、

各繰り返しにおける誤差が減少しなくなると繰り返しを終了し、

前記原音サンプルと前記第 3 の合成音声サンプルとの間の第 3 の誤差を計算する過程と

、

システムの不安定さに関して前記第 3 の誤差をテストする過程と、

前記テストにより不安定であるとされた場合、前記第 2 の係数を選択する過程と

を更に有することを特徴とする方法。

【請求項 4】

請求項 1 に記載の方法において、

前記第 2 の係数は前記勾配に適用される適応ステップサイズを用いて計算され、

前記適応ステップサイズは前記勾配と前記第 1 の係数との関数である

ことを特徴とする方法。

【請求項 5】

請求項 1 に記載の方法において、

前記テスト過程は更に、

前記第 2 誤差と前記第 1 誤差との間の差分を計算する過程と、

前記差分が終了閾値より大きいかどうかを判定する過程と

を有することを特徴とする方法。

【請求項 6】

請求項 1 に記載の方法において、

前記勾配に聴覚重み付けを与える過程を

更に有することを特徴とする方法。

【請求項 7】

請求項 1 に記載の方法において、

前記勾配に適用されたステップサイズを使用して前記第 2 の係数が繰り返し計算され、

各繰り返しにおいて 1 つの誤差が計算され、

前記誤差が各繰り返しで減少しなくなったら前記繰り返しを終了し、

前記勾配に聴覚重み付けを行う過程と

を更に有することを特徴とする方法。

【請求項 8】

音声合成フィルタを生成する方法において、

第 1 合成音声を生成する過程と、

前記第 1 合成音声に基づいて第 1 誤差エネルギーを計算する過程と、

前記第 1 誤差エネルギーに基づいて誤差エネルギー勾配を計算する過程と、

前記誤差エネルギー勾配を使用して第 2 合成音声を生成する過程と、

システムの不安定性さに関して前記第 2 の合成音声をテストする過程と、

前記テストにより不安定であるとされた場合、前記第 1 の合成音声を選択する過程と

を有することを特徴とする方法。

【請求項 9】

請求項 8 に記載の方法において、

10

20

30

40

50

前記誤差エネルギー勾配は、システム差分方程式から直接計算されることを特徴とする方法。

【請求項 1 0】

請求項 9 に記載の方法において、

前記第 2 合成音声から最小第 2 誤差エネルギーを繰り返し検索することにより、前記第 2 合成音声が生産される

ことを特徴とする方法。

【請求項 1 1】

請求項 1 0 に記載の方法において、

前記第 2 合成音声は、前記誤差エネルギー勾配に適用される適応ステップサイズを使用して計算され、

前記適応ステップサイズは、前記誤差エネルギー勾配と前記第 1 合成音声の関数であることを特徴とする方法。

【請求項 1 2】

第 1 の線形予測係数を計算する過程と、

システム差分方程式から直接誤差エネルギー勾配を計算する過程と、

前記誤差エネルギー勾配から第 2 の線形予測係数を計算する過程と、

前記第 1 の線形予測係数の前記誤差エネルギーと、前記第 2 の線形予測係数の前記誤差エネルギーとを比較する過程と、

前記第 1 の線形予測係数および前記第 2 の線形予測係数のうち誤差エネルギーの少ないほうを選択する過程と、

システム不安定性さに関して前記第 2 の線形予測係数をテストし、前記テストにより不安定であるとされた場合、前記第 1 の線形予測係数を選択する過程と

を有することを特徴とする音声最適化方法。

【請求項 1 3】

請求項 1 2 に記載の音声最適化方法において、

前記第 2 の線形予測係数を適応的に計算する過程を更に有する

ことを特徴とする音声最適化方法。

【請求項 1 4】

請求項 1 3 に記載の音声最適化方法において、

前記誤差エネルギー勾配に聴覚重み付けを行う過程を更に有する

ことを特徴とする音声最適化方法。

【請求項 1 5】

原音サンプルを使用して音声合成多項式の第 1 の線形予測係数を計算する手段と、

システム差分方程式から直接誤差エネルギー勾配を計算する手段と、

前記誤差エネルギー勾配から第 2 の線形予測係数を計算する手段と、

前記第 1 の線形予測係数の前記誤差エネルギーと、前記第 2 の線形予測係数の前記誤差エネルギーとを比較し、前記誤差エネルギーの少ない方を選択する手段と、

システム不安定性さに関して前記第 2 の線形予測係数をテストし、前記テストにより不安定であるとされた場合、前記第 1 の線形予測係数を選択する手段と

を有することを特徴とする音声最適化装置。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は音声符号化に関し、特に線形予測係数の最適化に関する。

【0 0 0 2】

【従来の技術】

音声符号化（あるいは圧縮）と言う技術は、データ送信のために音声をデジタルデータに符号化する広く知られた技術である。符号化されたデータを受信した受信者側は、その音声を再現する。デジタル化された音声データは符号化後、後に音声に復号されるまでの間

、様々なデジタル記録メディアに保存することが出来る。

【 0 0 0 3 】

音声符号化システムは、他のアナログ符号化システムおよびデジタル符号化システムとは異なっている。アナログ符号化システムおよびデジタル符号化システムでは、音声を高いビットレートでダイレクトサンプリングして、サンプリングされた生データを受信機に送信する。ダイレクトサンプリングシステムは、原音を高品質に再現し、再現音質が重要な場合に好まれる。ダイレクトサンプリングシステムが使われている一般的な例としては、アナログの音楽レコードやカセットテープ、またデジタルの音楽CDやDVDがある。ダイレクトサンプリングシステムの欠点は、データを送信するのに広い帯域幅が必要であり、またデータの保存にも大きな記憶容量が必要なことである。これゆえ、原音からサンプリングされた生の音声データを送信する典型的な符号化システムでは、毎秒128000ビットものデータレートが必要となることがある。

10

【 0 0 0 4 】

これに対して、音声符号化システムは、人間の音声発生の数学的モデルを使っている。発声モデル化の基本的技術は当技術分野で知られており、アメリカ音響協会の機関紙、1971年第50巻で、B・S・アタル(B・S・Atal)とスザンヌ・L・ハナー(Suzanne・L・Hanauer)による「音声分析と音声波の線形予測による合成(Speech Analysis and Synthesis by Linear Prediction of the Speech Wave)」に説明されている。音声符号化システムで使われる人間の音声発生のモデルは、普通ソース・フィルター・モデルと呼ばれている。一般に、このモデルには、肺と声帯によって発生する空気の流れを表している励起信号と、声道(すなわち声門、口、舌、鼻腔と唇)を表している合成フィルタが含まれている。よって、肺と声帯が声道にパルス状の空気の流れを発生させるように、励起信号は合成フィルタへの入力信号として働く。そして、声道が肺と声帯からの空気の流れに変更を加えるように、合成フィルタが励起信号に変更を加える。この結果、出来上がった合成音声は、ほぼ原音を表すようになる。

20

【 0 0 0 5 】

音声符号化システムの長所は、ダイレクトサンプリングシステムと比べて、原音をデジタル化した形で送信するのに必要な帯域幅が、ずっと小さくなり得ることである。比べてみると、ダイレクトサンプリングシステムでは原音を表す生の音響データが送られているのに対し、音声符号化システムでは、数学的な音声モデルを再現するのに必要なわずかな量の制御データが送られているだけである。結果として、典型的な音声符号化システムを使えば、音声を送るのに必要な帯域幅をほぼ毎秒2400~8000ビットまで減らすことができる。

30

【 0 0 0 6 】

【非特許文献1】

アメリカ音響協会、機関誌、1971年第50巻、B・S・アタル(B・S・Atal)とスザンヌ・L・ハナー(Suzanne・L・Hanauer)による「音声分析と音声波の線形予測による合成(Speech Analysis and Synthesis by Linear Prediction of the Speech Wave)」

【 0 0 0 7 】

【発明が解決しようとする課題】

40

音声符号化システムの欠点の1つは、ダイレクトサンプリングシステムに比べて、再現された音声の質がかなり低いことがあるということである。多くの音声符号化システムでは、受信者が正確に元の音声の内容を認知するのに十分な質は提供されている。しかし、いくつかの音声符号化システムでは、再現された音声は聞きやすい。つまり、受信者はもともと話された言葉を理解する事は出来るが、音声の質が低かったり、不快であったりする。従って、より正確な音声生成モデルを提供する音声符号化システムが望まれている。

【 0 0 0 8 】

音声符号化システムの質を改善する1つの方法として認知されているものが、ラシュキャリ(Lashkari)等によるアメリカ特許出願09/800071号に説明されている。簡単

50

に説明すると、この解決法は、原音サンプルと合成音声サンプルとの間の合成化誤差を最小化する方法である。この音声符号化システムで分かった難しい問題の1つは、合成化誤差がかなり非線形であるということである。このことで、この問題が数学的に難しくなっている。この問題を解くこの難しさは、合成フィルター多項式の解を、多項式の係数の代わりに使うことで克服されている。よって、合成フィルター多項式の解を探すための解最適化アルゴリズムが、そこで説明されている。

【0009】

上述の解法および当業者に知られている他の最適化解法に関して解決されない問題は、原音をエンコードするのに必要なコンピュータの処理能力が大きいことである。当業者なら分かるであろうが、原音を符号化するのに使用される様々な計算式を計算するのに、音声符号化システムでは、CPU (central processing unit) やDSP (digital signal processor) を使う必要がある。しばしば、携帯電話などの携帯ユニットで音声符号化が行われる場合、CPUやDSPは内臓のバッテリーから電気を得ている。よって、通常、音声符号化のために利用可能な計算能力は、CPUやDSPまたはバッテリー容量によって制限をうけることとなる。この問題は、どの音声符号化システムにおいても共通したものであるが、最適化アルゴリズムを使用するシステムにおいて、特に重大である。概して、最適化アルゴリズムは、通常の符号化アルゴリズムに加えて、追加の数式計算を含めることで、より質の高い音声を提供することができる。しかしながら、非効率な最適化アルゴリズムでは、CPUやDSPとして、値段が高く、計算能力があり、大きなものが必要になってしまう。非効率な最適化アルゴリズムでは、バッテリーも余計に使用し、バッテリーライフが短くなってしまう。よって、音声符号化システムでは、効率的な最適化アルゴリズムが望まれている。

【0010】

【課題を解決するための手段】

本発明は、音声生成の数学的モデルを最適化するためのアルゴリズムを提供する。この最適化アルゴリズムは、畳み込みを使用せず、また係数を解領域に変換することなしに、合成フィルタ多項式の係数を計算する。このように、係数の最適化に関する計算効率が改善される。原音と合成音声との間の誤差の勾配を使用して、係数がシステム差分方程式から直接計算される。システムの不安定さをテストするために終結閾値が使用され、不安定であるとされた場合、最適化アルゴリズムは止められ、最後の係数が返される。

【0011】

【発明の実施の形態】

図を参照して説明する。図1に、より正確に原音をモデル化するために合成誤差を最小化する音声符号化システムを示す。図1には音声の合成による分析 (A b S (analysis-by-synthesis)) システムが示されている。このシステムは、普通ソース・フィルター・モデルと呼ばれている。当技術分野でよく知られているように、ソース・フィルター・モデルは、人間の音声発生を数学的にモデル化する。このモデルで普通使われる仮定では、音声を生成している人間の音声発生メカニズムは、短い期間またはフレームで (例えば20 ~ 30 ms の分析フレーム) で、変化しないものとしている。更にまたこのモデルは、隣り合う期間の間に人間の音声発生メカニズムは変化するものと仮定している。このシステムでモデル化された物理的メカニズムには、声帯、声門、口、舌、鼻腔そして唇によって起こされる気圧変化が含まれている。よって、音声復号化器は、そのモデルを再現し、各期間用に制御データを少数セット使用するだけで、原音を再生成する。このように、通常の音声送信システムと異なり、原音の生サンプルデータは符号化器から復号化器へは送られない。この結果、送信されたり、記憶されたりするデジタル符号化されたデータ (つまり、バンド幅やビット数) は、典型的なダイレクトサンプリングシステムが必要とするのよりずっと少なくなる。

【0012】

図1において、デジタル化された元の音声10が励起モジュール12に送られている。励起モジュール12は、この原音の各サンプル $s(n)$ を解析して、励起関数 $u(n)$ を生

10

20

30

40

50

成する。励起関数 $u(n)$ は通常、連続したパルス信号であり、この連続したパルス信号は、声帯によって声道に突発的に放出される肺からの空気の流れを表している。原音サンプル $s(n)$ の性質によるが、励起関数 $u(n)$ は、有声音 13 か 14 か無声音 15 かである。

【0013】

音声符号化システムでの再生音質を改善する方法として、有声音の励起関数 $u(n)$ をより正確にする方法がある。今まで、励起関数 $u(n)$ は、決まったパルス間隔 P と大きさ G を持つパルス列 13 であった。当業者に知られている様に、この大きさ G と間隔 P を隣り合う期間で変えるようにしてもよい。大きさ G と間隔 P が固定されている今までのものに比べて、励起パルス 14 のパルスの大きさと間隔を変化させることで励起関数 $u(n)$ を最適化すると、よりよい音声合成がなされることが明らかにされている。この改良は、米国電気電子技術者協会 (IEEE) の音響、音声、信号処理に関する国際会議 (1982 年、614 頁 ~ 617 頁) の、ビシュヌ・S・アタル (Bishnu・S・Atal) とジョエル・R・レムデ (Joel R. Remde) による、「低ビットレートにおける自然な音声を生成するための LPC 励起の新しいモデル (A New Model of LPC Excitation For Producing Natural-Sounding Speech At Low Bit Rates)」に説明されている。

10

【0014】

この最適化技術では、原音 $s(n)$ を符号化するための計算量が増えるが、最近のコンピュータは励起関数 $u(n)$ の最適化に十分な計算能力があるので、重大な欠点ではない。この改良でのもっと重大な問題は、可変励起パルス 14 のデータを送信するのに必要な帯域幅が余計に必要であるということだ。この問題を解決する方法として、米国電気電子技術者協会 (IEEE) の音響、音声、信号処理に関する国際会議 (1985 年、937 頁 ~ 940 頁) の、マンフレッド・R・シュレッダー (Manfred R. Schroeder) とビシュヌ・S・アタル (Bishnu・S・Atal) による「符号励起線形予測化 (CELP): 低ビットレートにおける高品質音声 (Code-Excited Linear Prediction (CELP): High-Quality Speech At Very Low Bit Rates)」に説明されている符号化システムがある。

20

【0015】

この解決法では、多くの最適化された関数を分類して、関数ライブラリすなわちコードブックを作成する。そして、符号化励起モジュール 12 は、原音 $s(n)$ に最も近い合成音声生成する最適化された励起関数をコードブックから選択する。そして、コードブック内の最適な項目を特定するコードが復号化器に送られる。復号化器は送られてきたコードを受信し、対応するコードブックにアクセスし、選択された最適な励起関数 $u(n)$ を再生成する。

30

【0016】

励起モジュール 12 は、無声音 15 の励起関数 $u(n)$ も生成することが出来る。無声音 15 の励起関数 $u(n)$ は、話者の声帯が開いて、突発的な空気の流れが声道に起こされた時に使われる。多くの励起モジュール 12 は、この状態をモデル化するのに、パルスでなく白色ノイズ 15 (すなわちランダム信号) を有する励起関数 $u(n)$ を生成する。

【0017】

次に、合成フィルター 16 は、声道のモデル化と、声帯からの空気の流れに対する、声道の効果のモデル化を行う。普通は、合成フィルター 16 には、声道の様々な形を表す多項式を使う。多項式のパラメータつまり係数は通常、入力音声信号を使用して見積もりがなされる、また線形予測係数と呼ばれる。

40

【0018】

上述のアタル (Atal) とレムデ (Remde) によると、合成フィルター 16 は次の数式で表すことができる。

【数 1】

$$H(z) = G/A(z)$$

50

ここで、 G は音声の大きさを表している利得項である（簡潔のため、以降の式では利得項 G は省略する）。 $A(z)$ は M 次の多項式であり次の式で表される。

【数 2】

$$A(z) = 1 + \sum_{i=1}^M a_i z^{-i}$$

【0019】

多項式 $A(z)$ の次数は用途によって変わる。サンプリングレート8kHzの場合、10次の多項式が通常使用される。合成フィルタ16で決定される合成音声 $ss(n)$ と励起関数 $u(n)$ との関係は次の式で定義される。

10

【数 3】

$$ss(n) = u(n) - \sum_{i=1}^M a_i ss(n-i)$$

【0020】

ここで、表記「 ss 」は、本願の優先権主張の基礎である米国特許出願においては「 e 」の上に「 s 」を載せた表記となっていたものであるが、本願明細書においてそのような表記を用いることが困難であったため、その代わりに採用されたものである。従って、これ以降における表記「 ss 」は、実際には「 e 」の上に「 s 」を載せたものを表していると

20

【0021】

数3はシステム差分方程式とも呼ばれる。通常、この多項式の係数 $a_1 \dots a_M$ は、この分野で線形予測符号化(LPC)として知られる技術を使って計算される。LPCに基づく技術では、トータル予測誤差 E_p を最小にすることにより、多項式の係数 $a_1 \dots a_M$ を計算する。これにより、サンプル予測誤差 $e_p(n)$ が次の式により定義される。

【数 4】

$$e_p(n) = s(n) + \sum_{i=1}^M a_i s(n-i)$$

30

【0022】

トータル予測誤差 E_p は、次の式によって定義される。

【数 5】

$$E_p = \sum_{i=0}^{N-1} e_p^2(i)$$

ここで、 N はサンプルの数で表される分析フレームの長さである。多項式の係数 $a_1 \dots a_M$ は、よく知られた数学的方法を用いて、トータル予測誤差 E_p を最小化することにより解

40

【0023】

多項式の係数 $a_1 \dots a_M$ を計算するLPC技術に関する問題の1つは、トータル予測誤差だけが最小化されることである。このように、LPC技術では、原音 $s(n)$ と合成音声 $ss(n)$ との間の誤差を最小化していない。従って、サンプル合成誤差 $e_s(n)$ は次の式で定義できる。

【数 6】

$$e_s(n) = s(n) - ss(n)$$

50

【 0 0 2 4 】

トータル合成誤差 E_s は、合成誤差エネルギー J とも呼ばれるが、次の式で定義される。

【 数 7 】

$$J = E_s = \sum_{n=0}^{N-1} e_s^2(n) = \sum_{n=0}^{N-1} (s(n) - ss(n))^2$$

N はサンプルの数で表される分析フレームの長さである。上で述べたトータル予測誤差 E_p のように、合成誤差エネルギー J は、最適フィルタの係数 $a_1 \dots a_M$ を計算するために最小化されなければならない。しかしながら、この技術の問題は、数 3 で表される合成音声 $ss(n)$ によって、合成誤差エネルギー J がかなり非線形な関数になり、数学的に扱うのが難しいことである。

10

【 0 0 2 5 】

最適化アルゴリズムの中には、この難しさを、係数 $a_1 \dots a_M$ の代わりに、多項式 $A(z)$ の解を使用することで避けているものもある。合成フィルタ 16 の安定性に関して制御することができるが、この解決法では、多くの計算能力を要求する。更に、解領域最適化は、畳み込みに関連したインパルス応答に基づいた勾配を通常計算する。畳み込みでは、先行する音声サンプルの履歴を考慮することなしに、システムのゼロ状態応答がわかるだけである。

【 0 0 2 6 】

他の最適化アルゴリズムと比較して、フィルタ係数 $a_1 \dots a_M$ を解領域に変換することなしに、フィルタ係数 $a_1 \dots a_M$ を最適化するため、合成誤差エネルギー J の勾配が直接計算されるようにできる。よって、数 3、数 6、数 7 を使用して、合成誤差エネルギー J の勾配は次式で与えられる。

20

【 数 8 】

$$\frac{\partial J}{\partial a_i} = \sum_{n=0}^{N-1} \frac{\partial}{\partial a_i} (s[n] - ss[n])^2 = \sum_{n=0}^{N-1} 2(s[n] - ss[n]) \left(-\frac{\partial ss[n]}{\partial a_i} \right)$$

ここで $i = 1 \sim M$ である。

30

数 3 を使用して、係数 $a_1 \dots a_M$ に関する合成音声の勾配は、次式で表される。

【 数 9 】

$$\frac{\partial ss[n]}{\partial a_i} = -\sum_{j=1}^M a_j \frac{\partial ss[n-j]}{\partial a_i} - ss[n-i]$$

ここで、 $n = 0 \sim N-1$ 、かつ $i = 1 \sim M$ である。

係数 $a_1 \dots a_M$ は、長さ N のフレームだけで有効であると仮定される。このように、係数は、 $n = 0 \sim N-1$ のみで存在し、この間隔を外れたところでは、音声は、係数 $a_1 \dots a_M$ から独立している。

40

【 0 0 2 7 】

この最適化アルゴリズムの利点の 1 つは、数 9 に示される合成音声の偏導関数が、回帰的な方法により効率よく計算できることである。従って、数 9 の偏導関数は、以下の 2 次元配列で示される。

【 数 10 】

$$D[n, i] = \frac{\partial ss[n]}{\partial a_i}$$

ここで、 $n = 0 \sim N-1$ 、および $i = 1 \sim M$ である。

50

配列 $D[n, i]$ は、以下の繰り返しコード A を使用して計算される。

【 0 0 2 8 】

< 繰り返しコード A >

For $n=1$ to $N-1$

$D[n, i] = -ss[n-i]$

For $j = 1$ to M

If $n-j = 0$

$D[n, i] = -a_j * D[n-j, i]$

Next j

Next n

10

【 0 0 2 9 】

繰り返しコード A を使用して計算された偏導関数を数 8 に代入すると、合成誤差エネルギーの勾配を得ることができる。

合成誤差エネルギーの勾配ベクトルは、次式で得ることができる。

【 数 1 1 】

$$\nabla J = \left[\frac{\partial J}{\partial a_1} \quad \frac{\partial J}{\partial a_2} \quad \cdots \quad \frac{\partial J}{\partial a_M} \right]^T$$

係数 $a_1 \dots a_M$ のベクトルは、次式でも定義される。

20

【 数 1 2 】

$$a = [a_1 \quad a_2 \cdots a_M]^T$$

最適化係数は、次式を使用して計算できる。

【 数 1 3 】

$$a_{\text{new}} = a - \mu \nabla J$$

ここで、 μ は正の値を持ちステップサイズとして知られる。よって、係数の新たなベクトルは、勾配に対して負の方向に移動することによって計算される。ステップサイズ μ の大きさは、最適化プロセスのスピードと安定性を変えるために、増やされたり減らされたりする。

30

【 0 0 3 0 】

最適化アルゴリズムで使用されるフローチャートを図 2 と図 3 に示す。このフローチャートは、図 1 の合成フィルタ 1 6 および合成フィルタ最適化部 1 8 によって実行されるものである。図 2 に示すように、最適化アルゴリズムへ入力されるのは、原音 s 、励起関数 u 、線形予測係数 $a_1 \dots a_M$ 、そして合成フィルタのメモリである（ステップ 2 2）。CalculateErrorEnergy 関数を呼び出して、数 7 を使用して、合成誤差エネルギー J が計算される（ステップ 2 4）。変数 OldErrorEnergy が、次に初期化されて、合成誤差エネルギー J が代入される（ステップ 2 4）。次に、CalculateGradient 関数を呼び出して、合成誤差エネルギーの勾配が、数 8 もしくは繰り返しコード A を使用して計算される。

40

【 0 0 3 1 】

図 3 に示す GradientDescent 関数を呼び出すことで、変数 ErrorEnergy が、計算される（ステップ 2 8）。変数 ErrorEnergy は、更新された係数 $a_1 \dots a_M$ を使用した合成誤差エネルギー J を表している。それに対し、変数 OldErrorEnergy は、以前の係数 $a_1 \dots a_M$ を使った場合の合成誤差エネルギーを表している。次に、ErrorEnergy は、OldErrorEnergy と比較される。これは、合成誤差エネルギー J の変化が、終了閾値 Termination_Threshold 以下であるかどうかを決定するものである（ステップ 3 0）。合成誤差エネルギー J の変化が、終了閾値 Termination_Threshold 以下でない場合は、OldErrorEnergy に ErrorEnergy

50

が代入され、合成誤差エネルギー J の新たな勾配、新たな係数 $a_1 \dots a_M$ 、そして新たな合成誤差エネルギー J が、ステップ 32、26、28 で計算される。合成誤差エネルギー J の変化が、終了閾値 Termination_Threshold 以下である場合、最適化アルゴリズムは終了し、係数の修正された線形予測値が返される（ステップ 34）。

【0032】

図3にGradientDescent関数を示す（ステップ36）。GradientDescent関数は、CalculateErrorEnergyを呼び出し、数7を使用して、合成誤差エネルギー J を計算することで開始する（ステップ38）。変数OldErrorEnergyに合成誤差エネルギー J が代入される（ステップ38）。次に、合成誤差エネルギー J の勾配ベクトルの大きさが計算され、変数Gノルムが割り当てられる。また、係数 $a_1 \dots a_M$ の大きさが計算され、変数Aノルムが割り当てられる（ステップ40）。変数Gノルム、Aノルム、および変数ステップサイズに割り当てられた所定の値を使用して、適応ステップサイズ μ が計算される（ステップ42）。新たな線形予測係数 a_{new} が、数13を使用して計算される（ステップ44）。次に合成誤差エネルギー J が、CalculateErrorEnergyを呼び出すことで、数7を使用して計算される（ステップ46）。次に変数ErrorEnergyに、合成誤差エネルギー J が代入される（ステップ46）。次に、合成誤差エネルギー J に増減があるかを調べるために、OldErrorEnergyがErrorEnergyと比較される（ステップ48）。合成誤差エネルギー J が減少している場合、変数OldErrorEnergyに変数ErrorEnergyが代入され、線形予測係数 $a_1 \dots a_M$ が新たな線形予測係数 a_{new} へと更新される（ステップ50）。それから、新たな線形予測係数 a_{new} と合成誤差エネルギー J が計算される（ステップ44、46）。一方、合成誤差エネルギー J が増大していた場合、GradientDescent関数は終了して、現在の線形予測係数 $a_1 \dots a_M$ とOldErrorEnergyが返される（ステップ52）。

【0033】

合成モデルが決定され、係数 $a_1 \dots a_M$ が最適化されたら、モデル化のための制御データが、送信又は保存のために量子化されてデジタルデータへと変換される。量子化には業界で標準となっている方法が多くある。ある例では、量子化された制御データは、10つの合成フィルタ係数 $a_1 \dots a_M$ 、励起パルスの大きさを表す利得値 G を1つ、励起パルスの周波数のためのピッチ間隔を1つ、有声13もしくは無声15の励起関数 $u(n)$ を示す指示子を1つ、含んでいる。よって、この例では、各音声フレームの最後で、13の異なる変数を送信する必要がある。しかしながら、他の制御データを送るようにしてもよい。例えば、CELPEncodeにおいては、使用する最適化励起関数 $u(n)$ を特定するコードブックインデックスも送信される。通常CELPEncodeでは、制御データは、合計80ビットへと量子化される。この例では、最適化を算入した合成音声 $ss(n)$ を、毎秒8000ビット（80ビット/フレーム ÷ 0.01秒/フレーム）のバンド幅で送ることができる。

【0034】

エンコードに使用できるコンピュータの能力によるが、合成モデルをより正確にするために、追加のエンコードシーケンスも可能である。これらのシーケンス例を、図1に点線で示してある。例えば、励起関数 $u(n)$ は、合成モデルにおけるエンコードの間、様々な段階で再最適化を行うことができる。更に、数式とアルゴリズムは、特定のアプリケーションのために、変更することもできる。

【0035】

合成音声の主観的品質を更に改善する方法として、聴覚重み付けを使用する方法がある。この場合、合成誤差エネルギー J は、次式によって、聴覚重み付けも使用して定義される。

【数14】

$$J = \sum_{n=0}^{N-1} \{e[n] * h_w[n]\}^2 = \sum_{n=0}^{N-1} \{(s[n] - ss[n]) * h_w[n]\}^2$$

10

20

30

40

ここで、 $h_w[n]$ は、聴覚重み付けフィルタのインパルス応答である。
数 1 4 に畳み込み処理を行うことで、合成誤差エネルギーは次式になる。

【数 1 5】

$$J = \sum_{n=0}^{N-1} \left\{ \sum_{k=0}^n h_w[k]s[n-k] - \sum_{l=0}^n h_w[l]ss[n-l] \right\}^2$$

$$= \sum_{n=0}^{N-1} \left\{ \left(\sum_{k=0}^n h_w[k]s[n-k] \right)^2 - 2 \left(\sum_{k=0}^n h_w[k]s[n-k] \right) \left(\sum_{l=0}^n h_w[l]ss[n-l] \right) + \left(\sum_{l=0}^n h_w[l]ss[n-l] \right)^2 \right\}$$

10

【0 0 3 6】

次に、数 1 5 を合成フィルタ係数 $a_1 \dots a_M$ に関して微分することにより、合成誤差エネルギーの偏導関数は、以下の式になる。

【数 1 6】

$$\frac{\partial J}{\partial a_i} = \sum_{n=0}^{N-1} \left\{ -2 \left(\sum_{k=0}^n h_w[k](s[n-k] - ss[n-k]) \right) \left(\sum_{l=0}^n h_w[l] \frac{\partial}{\partial a_i} ss[n-l] \right) \right\}$$

聴覚重み付け誤差も、次式で定義することができる。

20

【数 1 7】

$$e_w[n] = e[n] * h_w[n]$$

数 1 0、数 1 4、数 1 6、数 1 7、を使用して、合成誤差エネルギーの偏導関数は、以下の式のようにになる。

【数 1 8】

$$\frac{\partial J}{\partial a_i} = \sum_{n=0}^{N-1} \{ -2e_w[n](D[n,i] * h_w[n]) \}$$

30

したがって、最適化アルゴリズムは、数 1 4 を数 7 の代わりに使い、更に数 1 8 を数 8 の代わりに使って、更に改善される。

【0 0 3 7】

当業者に明らかなように、この最適化アルゴリズムは、合成フィルタ多項式 $A(z)$ を最適化するのに必要な計算を著しく減らす。よって、エンコーダの効率が著しく改善される。あるいはこの効率化を、合成音声 $ss(n)$ の質を改善するのに使用することもできる。従来の最適化アルゴリズムを使用する場合、各サンプルを合成音声にするのに必要な計算が多かった。しかし、改善された最適化アルゴリズムは、解領域を使用することなしに、合成誤差エネルギー J の勾配をシステム差分方程式から直接計算することで、合成音声 $ss(n)$ を計算するのに必要な計算量を減らしている。本発明は、様々な音声エンコーダに適用することができるが、ラシュキャリ等による米国特許出願 0 9 / 8 0 0 0 7 1 号に記載の音声エンコーダに適用した場合は、すでに性能の改善が認められている。

40

【0 0 3 8】

この最適化アルゴリズムは多くの利点を持っている。例えば、解領域解決法が通常使う畳み込みを、この最適化アルゴリズムは使用しないので、ゼロ状態応答とゼロ入力応答を含む、システムのトータル応答も考慮されることになる。この最適化アルゴリズムはまた、各繰り返しにおいて、合成誤差エネルギーを所定の終了閾値でテストすることで、不安定さの問題を避けている。よって、テスト結果が、システムが不安定であるとなった場合、最適化アルゴリズムは終了して、最後に最適化された線形予測係数が使われる。適応ス

50

テップサイズも、最適化のスピードを改善するために使用される。更に、最適化アルゴリズムの計算効率を改善する反復アルゴリズムを使用して、合成誤差エネルギーの勾配を計算することもできる。この最適化アルゴリズムの他の利点は、合成音声の質を更に改善するために、聴覚重み付けを使うことができることである。

【 0 0 3 9 】

図 4 は、最適化アルゴリズムによる結果例を示す図である。図 4 は、原音サンプル、G . 7 2 9 エンコーダを使用した場合の合成音声、勾配降下最適化アルゴリズムを使用した G . 7 2 9 エンコーダを使用した場合の合成音声、それぞれの音声波形を示している。当業者は理解するであろうが、この G . 7 2 9 エンコーダは、様々な音声エンコーダの質を比較するために、研究者などに使用される標準化された音声エンコーダである。図から分かるように、勾配降下最適化アルゴリズムを使った場合の合成音声は、G . 7 2 9 エンコーダだけで生成された合成音声よりも、原音に一致している。

10

【 0 0 4 0 】

他のテスト結果では、最適化アルゴリズムにより良好になった質および効率が示されている。例えば、あるテストにおいて、最適化アルゴリズムの付属したものと付属していない G . 7 2 9 を使用して、男性と女性の標準化された音声データを符号化した。10 ミリ秒毎に部分 S N 比 (S S N R : segmental signal to noise ratio measurements) を計測すると、G . 7 2 9 エンコーダのみの場合、7 . 0 6 d B S S N R であり、解領域最適化アルゴリズム付属の G . 7 2 9 エンコーダを使用した場合、7 . 3 3 d B S S N R となり、説明を行ってきた勾配降下最適化アルゴリズム付属の G . 7 2 9 エンコーダでは、7 . 3 4 d B S S N R となった。当業者には当然であるが、S S N R 計測値が高いということは通常、合成音声が良好な聴覚品質を有しているということである。更に、勾配降下最適化アルゴリズムの計算量は、解領域最適化アルゴリズムに比べて、およそ 20 分の 1 ~ 30 分の 1 である。このテストは聴覚重み付けを使用せずに行った。聴覚重み付けを行うと、勾配降下最適化アルゴリズムでは、部分信号対重み付けノイズ比 (S S W N R : segmental signal to weighted noise ratio) の計測値が 1 4 . 1 9 d B S S W N R となるが、普通の G . 7 2 9 エンコーダでは、1 4 . 0 4 d B S S W N R である。

20

【 0 0 4 1 】

本発明の好ましい実施形態をここに説明したが、本発明はこれに限定されず、本発明の趣旨から外れずに変形することが可能である。本発明の範囲は、特許請求の範囲によって決まるものであり、文言上でも均等上でも特許請求の範囲内にある装置および方法は、本発明に含まれるものである。

30

【 0 0 4 2 】

【発明の効果】

以上説明したように、本発明によれば、より正確な合成音声を提供する音声符号化システムが提供される。

【図面の簡単な説明】

【図 1】 音声の合成による分析システムのブロック図である。

【図 2】 フィルタ係数のための最適化アルゴリズムのフローチャートである。

【図 3】 誤差勾配を見つけるための最適化アルゴリズムで使用される、勾配降下関数のフローチャートである。

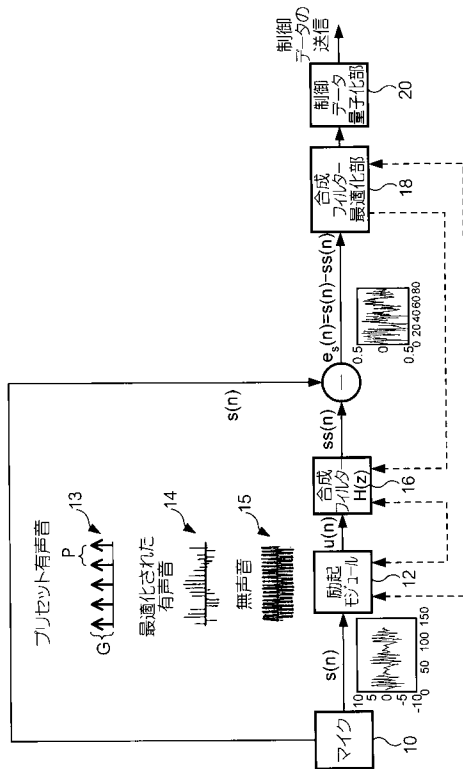
40

【図 4】 原音の波形と、G . 7 2 9 エンコーダを使用した場合の合成音声の波形と、勾配降下最適化を行う G . 7 2 9 エンコーダを使用した場合の合成音声の波形とを比較している図である。

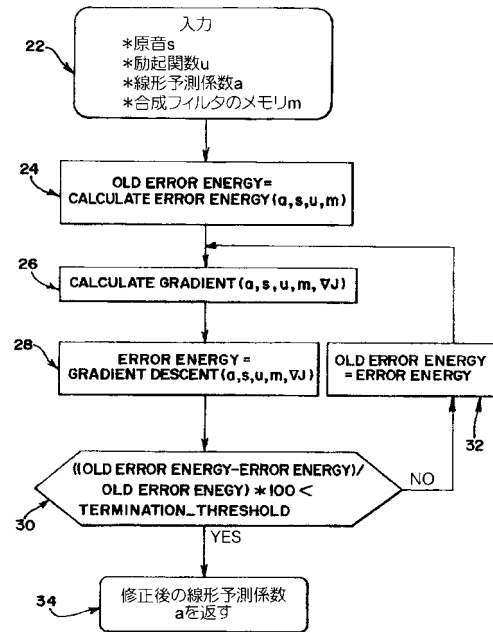
【符号の説明】

1 0マイク、1 2励起モジュール、1 6合成フィルタ、1 8合成フィルタ最適化部、2 0制御データ量子化部

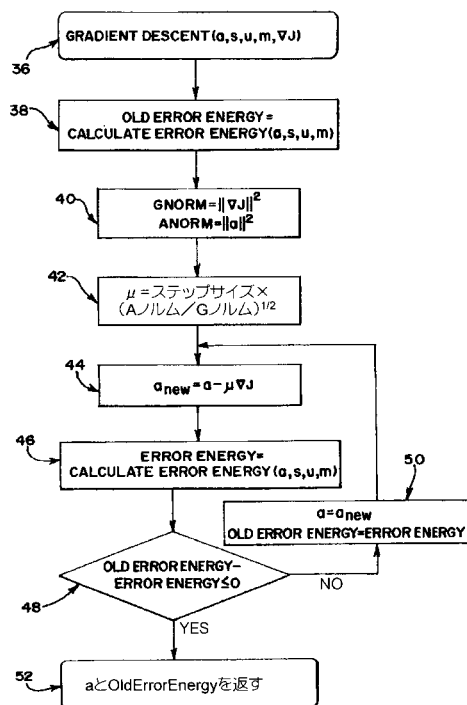
【 図 1 】



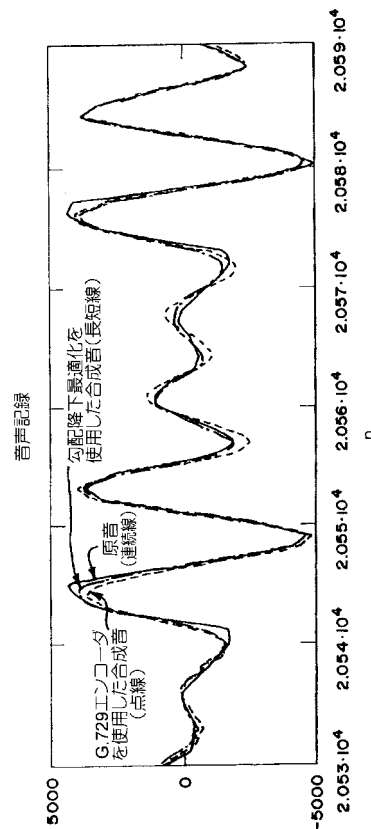
【 図 2 】



【圖 3】



【 図 4 】



フロントページの続き

- (56)参考文献 特開昭62-111299(JP,A)
特開平02-272500(JP,A)
特開2002-073097(JP,A)
特開2002-328692(JP,A)
小池伸一, "絶対値誤差の勾配に基づく適応制御ステップサイズサインアルゴリズムの解析",
1997年電子情報通信学会総合大会講演論文集, A-4-46(1997-03), p.177
小池伸一, "勾配法を用いた適応制御ステップサイズサイン - サインアルゴリズムの解析", 19
97年電子情報通信学会基礎・境界ソサイエティ大会講演論文集, A-4-14(1997-08), p.84

(58)調査した分野(Int.Cl., DB名)

G10L 13/00-13/08, 19/00-19/14

CiNii