

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2011-90311

(P2011-90311A)

(43) 公開日 平成23年5月6日(2011.5.6)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 O L 19/14 (2006.01)	G 1 O L 19/14 4 O O B	5 J O 6 4
G 1 O L 19/12 (2006.01)	G 1 O L 19/12 Z	
G 1 O L 19/02 (2006.01)	G 1 O L 19/02 1 5 O	
H O 3 M 7/30 (2006.01)	H O 3 M 7/30 Z	

審査請求 有 請求項の数 1 O L 外国語出願 (全 28 頁)

(21) 出願番号	特願2010-249991 (P2010-249991)	(71) 出願人	595020643
(22) 出願日	平成22年11月8日 (2010.11.8)		クアルコム・インコーポレイテッド
(62) 分割の表示	特願2001-564148 (P2001-564148) の分割		QUALCOMM INCORPORATED
原出願日	平成12年2月29日 (2000.2.29)		アメリカ合衆国、カリフォルニア州 92 121-1714、サン・ディエゴ、モア ハウス・ドライブ 5775
		(74) 代理人	100108855 弁理士 蔵田 昌俊
		(74) 代理人	100091351 弁理士 河野 哲
		(74) 代理人	100088683 弁理士 中村 誠
		(74) 代理人	100109830 弁理士 福原 淑弘

最終頁に続く

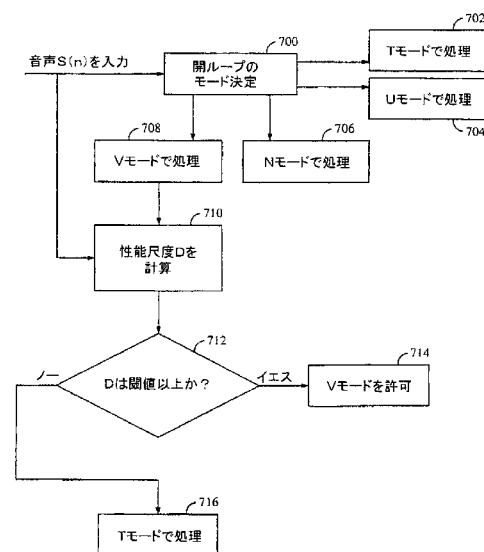
(54) 【発明の名称】 閉ループのマルチモードの混合領域の線形予測音声コーデ

(57) 【要約】 (修正有)

【課題】 音声を閉ループのマルチモードの混合領域でコード化する方法の提供。

【解決手段】 閉ループのマルチモードの混合領域の線形予測 (MDLP) の音声コーデは、高レートの時間領域コード化モードと、低レートの周波数領域コード化モードと、入力されたフレームの音声内容に基づいてモードを選択するモード選択機構とを含む。遷移音声のフレームは、高レート of CELP モードでコード化される。有声音声のフレームは、低レートの高調波モードでコード化される。位相パラメータは、準位相モデルによりモデル化される。初期位相値は、周波数領域モードでコード化された直前の音声フレームの初期位相値になり、直前の時間領域でコード化された音声フレーム情報から計算される。周波数領域モードでコード化される各音声フレームを、対応する入力音声フレームと比較して、性能尺度が所定の閾値よりも低いときは、時間領域モードでコード化される。

【選択図】 図 9



【特許請求の範囲】**【請求項 1】**

少なくとも 1 つの時間領域コード化モードおよび少なくとも 1 つの周波数領域コード化モードをもつコードと、

コードに接続され、かつ音声プロセッサによって処理されるフレーム内容に基づいてコードのコード化モードを選択するように構成されている閉ループのモード選択デバイスとを含むマルチモードの混合領域の音声プロセッサ。

【請求項 2】

コードが、音声フレームをコード化する請求項 1 記載の音声プロセッサ。

【請求項 3】

コードが、音声フレームの線形予測残余をコード化する請求項 1 記載の音声プロセッサ。

【請求項 4】

少なくとも 1 つの時間領域コード化モードが、第 1 のコード化レートでフレームをコード化するコード化モードを含み、少なくとも 1 つの周波数領域コード化モードが、第 2 のコード化レートでフレームをコード化するコード化モードを含み、第 2 のコード化レートが第 1 のコード化レートよりも低い請求項 1 記載の音声プロセッサ。

【請求項 5】

少なくとも 1 つの周波数領域コード化モードが、高調波のコード化モードを含む請求項 1 記載の音声プロセッサ。

【請求項 6】

コードに接続された比較回路であって、コード化されていないフレームを、少なくとも 1 つの周波数領域コード化モードでコード化されたフレームと比較して、比較に基づいて性能尺度を生成する比較回路をさらに含み、コードが、性能尺度が所定の閾値よりも低いときだけ、少なくとも 1 つの時間領域コード化モードを適用し、さもなければコードは、少なくとも 1 つの周波数領域コード化モードを適用する請求項 1 記載の音声プロセッサ。

【請求項 7】

コードが、少なくとも 1 つの時間領域コード化モードを、少なくとも 1 つの周波数領域コード化モードでコード化された所定数の連続的に処理されるフレームの直ぐ後の各フレームに適用する請求項 1 記載の音声プロセッサ。

【請求項 8】

少なくとも 1 つの周波数領域コード化モードが、周波数、位相、および振幅を含む 1 組のパラメータをもつ複数のシヌソイドで各フレームの短期間のスペクトルを表わし、位相は多項式表現および初期位相値でモデル化されていて、初期位相値が、(1) 前フレームが、少なくとも 1 つの周波数領域コード化モードでコード化されたときは、前フレームの推定された最終位相値であるか、または (2) 前フレームが、少なくとも 1 つの時間領域コード化モードでコード化されたときは、前フレームの短期間のスペクトルから求められる位相値である請求項 1 記載の音声プロセッサ。

【請求項 9】

各フレームにおけるシヌソイドの周波数が、フレームのピッチ周波数の整数倍である請求項 8 記載の音声プロセッサ。

【請求項 10】

各フレームにおけるシヌソイドの周波数が、0 ないし 2 の 1 組の実数から得られる請求項 8 記載の音声プロセッサ。

【請求項 11】

フレームを処理する方法であって、

閉ループのコード化モード選択プロセスを各連続する入力フレームへ適用して、入力フレームの音声内容に基づいて、時間領域コード化モードか、または周波数領域コード化モードの何れかを選択するステップと、

入力フレームの音声内容が定常状態の有声音の音声を示すときは、入力フレームを周波数領域でコード化するステップと、

10

20

30

40

50

入力フレームの音声内容が定常状態の有声音の音声以外のものを示すときは、入力フレームを時間領域でコード化するステップと、

周波数領域でコード化されたフレームと入力フレームとを比較して、性能尺度を求めるステップと、

性能尺度が所定の閾値よりも低いときは、入力フレームを時間領域でコード化するステップとを含むフレームを処理する方法。

【請求項 1 2】

フレームが、線形予測残余フレームである請求項 1 1 記載の方法。

【請求項 1 3】

フレームが音声フレームである請求項 1 1 記載の方法。

10

【請求項 1 4】

時間領域でコード化するステップが、第1のコード化レートでフレームをコード化することを含み、周波数領域でコード化するステップが、第2のコード化レートでフレームをコード化することを含み、第2のコード化レートが第1のコード化レートよりも低い請求項 1 1 記載の方法。

【請求項 1 5】

周波数領域でコード化するステップが、高調波でコード化することを含む請求項 1 1 記載の方法。

【請求項 1 6】

周波数領域でコード化するステップが、周波数、位相、および振幅を含む1組のパラメータをもつ複数のシヌソイドで各フレームの短期間のスペクトルを表わし、位相は多項式表現および初期位相値でモデル化されていて、初期位相値が、(1)前フレームが周波数領域でコード化されたときは、前フレームの推定された最終位相値であるか、または(2)前フレームが時間領域でコード化されたときは、前フレームの短期間のスペクトルから求められる位相値である請求項 1 1 記載の方法。

20

【請求項 1 7】

各フレームのシヌソイド周波数が、フレームのピッチ周波数の整数倍である請求項 1 6 記載の方法。

【請求項 1 8】

各フレームのシヌソイド周波数が、0 ないし 2 の1組の実数から得られる請求項 1 6 記載の方法。

30

【請求項 1 9】

マルチモードの混合領域の音声プロセッサであって、

開ループのコード化モード選択プロセスを入力フレームへ適用して、入力フレームの音声内容に基づいて、時間領域コード化モードか、または周波数領域コード化モードの何れかを選択する手段と、

入力フレームの音声内容が定常状態の有声音の音声を示すときは、入力フレームを周波数領域でコード化する手段と、

入力フレームの音声内容が定常状態の有声音の音声以外のものを示すときは、入力フレームを時間領域でコード化する手段と、

40

周波数領域でコード化されたフレームと入力フレームとを比較して、性能尺度を求める手段と、

性能尺度が所定の閾値よりも低いときは、入力フレームを時間領域でコード化する手段とを含むマルチモードの混合領域の音声プロセッサ。

【請求項 2 0】

フレームが線形予測残余フレームである請求項 1 9 記載の音声プロセッサ。

【請求項 2 1】

入力フレームが音声フレームである請求項 1 9 記載の音声プロセッサ。

【請求項 2 2】

時間領域でコード化する手段が、第1のコード化レートでフレームをコード化する手段を

50

含み、周波数領域でコード化する手段が、第 2 のコード化レートでフレームをコード化する手段を含み、第 2 のコード化レートが第 1 のコード化レートよりも低い請求項 19 記載の音声プロセッサ。

【請求項 23】

周波数領域でコード化する手段が、高調波コードを含む請求項 19 記載の音声プロセッサ。

【請求項 24】

周波数領域でコード化する手段が、周波数、位相、および振幅を含む 1 組のパラメータをもつ複数のシヌソイドで各フレームの短期間のスペクトルを表わす手段を含み、位相は多項式表現および初期位相値でモデル化されていて、初期位相値が、(1) 直前のフレームが周波数領域でコード化されたときは、直前のフレームの推定された最終位相値であるか、または (2) 直前のフレームが時間領域でコード化されたときは、直前のフレームの短期間のスペクトルから求められる位相値である請求項 19 記載の音声プロセッサ。

10

【請求項 25】

各フレームのシヌソイド周波数が、フレームのピッチ周波数の整数倍である請求項 24 記載の音声プロセッサ。

【請求項 26】

各フレームのシヌソイド周波数が、0 ないし 2 の 1 組の実数から得られる請求項 24 記載の音声プロセッサ。

【発明の詳細な説明】

20

【技術分野】

【0001】

本発明は、概ね音声処理の分野、とくに音声を閉ループのマルチモードの混合領域でコード化するための方法および装置に関する。

【背景技術】

【0002】

ディジタル技術による音声 (voice) の伝送は、とくに長距離のディジタル無線電話の応用において普及してきた。これにより、チャンネル上で送ることができる最少情報量を判断し、一方で再構成された音声の知覚品質を維持することに関心が生まれた。音声を単にサンプリングして、ディジタル形式にすることによって送るとき、従来のアナログ電話の音声品質を実現するには、毎秒 64 キロビット秒 (kbps) のオーダのデータレートが必要である。しかしながら、音声解析を使用し、その後で適切にコード化し、伝送し、受信機において再合成をすることによって、データレートを相当に低減することができる。

30

【0003】

人間の音声の生成モデルに係るパラメータを抽出することによって音声を圧縮する技術を採用したデバイスは、音声コードと呼ばれている。音声コードは、入力音声信号を時間のブロック、すなわち解析フレームに分割する。一般的に音声コードはエンコーダとデコーダとを含む。エンコーダは、入力音声フレームを解析して、一定の関連するパラメータを抽出して、パラメータを二値表現、すなわち 1 組のビットまたは二値データパケットに量子化する。データパケットは通信チャンネル上で受信機およびデコーダへ送られる。デコーダはデータパケットを処理し、非量子化して (unquantize) パラメータを生成し、非量子化したパラメータを使用して音声フレームを再合成する。

40

【0004】

音声コードの機能は、音声が本質的にもっている固有の冗長の全てを取去ることによって、ディジタル化された音声信号を低ビットレートの信号へ圧縮することである。ディジタル圧縮は、入力音声フレームを 1 組のパラメータで表わし、量子化を採用して、このパラメータを 1 組のビットで表わすことによって実現される。入力音声フレームが多数のビット N_i をもち、音声コードによって生成されるデータパケットが多数のビット N_0 をもつとき、音声コードによって得られる圧縮係数は、 $C_r = N_i / N_0$ である。デコードされた音声 (speech) の高い音声品質 (voice quality) を維持し、一方で目標の圧縮係数を

50

得ることが課題とされている。音声コードの性能は、(1) 音声モデル、すなわち上述の解析および合成プロセスの組合せがどのくらい適切に行われるか、および(2) パラメータ量子化プロセスが1フレーム当り N_0 ビットの目標ビットレートでどのくらい適切に実行されるかに依存する。したがって音声モデルは、各フレームごとの小さい組のパラメータを使用して、音声信号の本質(essence)、すなわち目標の音声品質を得ることを目的としている。

【0005】

音声コードは時間領域のコード、すなわち音声の小さいセグメント(一般的に5ミリ秒(milliseconds, ms)のサブフレーム)を一度にコード化する高度な時間分解処理(time-resolution processing)を採用することによって時間領域の音声波形を得ることを試みる時間領域のコードとして構成することができる。各サブフレームごとに、この技術において知られている種々のサーチアルゴリズムによって、コードブック空間から高精度の見本(representative)を見付ける。その代わりに、音声コードは周波数領域のコードとして構成されていてもよく、1組のパラメータを使用して入力音声フレームの短期間の音声スペクトルを捕らえて(解析)、対応する合成プロセスを採用して、スペクトルパラメータから音声波形を再現することを試みる。パラメータ量子化器は、文献(A Gersho & R.M. Gray, Vector Quantization and Signal Compression (1992))に記載されている既知の量子化技術にしたがって、コードベクトルの記憶されている表現を使用してパラメータを表わすことによってそれらのパラメータを保存する。

【0006】

よく知られている時間領域の音声コードは、C E L P (Code Excited Linear Predictive) コードであり、これはL.B. Rabiner & R.W. Schaferによる文献(Digital Processing of Speech Signals 396-453 (1978))に記載されており、ここでは参考文献として全体的にこれを取り上げている。C E L P コードでは、線形予測(linear prediction, LP)解析によって、短期間のフォルマントフィルタの係数を見付け、音声信号における短期間の相関関係、すなわち冗長を取去る。短期間の予測フィルタを入力音声フレームに適用して、L Pの残余信号(residue signal)を生成し、このL Pの残余信号をさらに長期間の予測フィルタパラメータおよび次の確率コードブックでモデル化して、量子化する。したがってC E L Pのコード化は、時間領域の音声波形をコード化するタスクを、L Pの短期間のフィルタ係数をコード化するタスクおよびL Pの残余をコード化するタスクの別々のタスクへ分ける。時間領域のコード化は、固定レート(すなわち、各フレームごとに、同数のビット N_0 を使用するレート)で、または可変レート(すなわち、異なるビットレートが異なるタイプのフレームの内容に対して使用されるレート)で実行することができる。可変レートのコードは、目標の品質を得るのに適したレベルまでコーデックパラメータをコード化するのに必要なビット量のみを使用することを試みる。例示的な可変レートのC E L Pのコードは米国特許第5,414,796号に記載されており、これは本発明の譲受人に譲渡され、ここでは参考文献として全体的に取り上げている。

【0007】

C E L P コードのような時間領域のコードは、通常は、フレームごとに多数のビット N_0 に依存して、時間領域の音声波形の精度を保持する。このようなコードは、通常はフレーム当りのビット数 N_0 が比較的が多いとき(例えば、8キロビット秒以上)、優れた音声品質を伝える。しかしながら低ビットレート(4キロビット秒以下)では、時間領域のコードは、使用可能なビット数が制限されているために、高品質で丈夫な性能を維持しない。低ビットレートではコードブック空間が制限されているので、従来の時間領域のコードには備えられている波形を整合する能力を取去って、より高レートの市販のアプリケーションにおいてこのようなコードを実行するのに成功した。

【0008】

現在、研究に対する関心および活発な商業上の要求が急激に高まり、中程度から低いビットレート(すなわち、2.4ないし4キロビット秒の範囲およびそれ以下)で動作する高品質の音声コードを発展させた。応用分野には、無線電話通信、衛星通信、インターネ

10

20

30

40

50

ット電話通信、種々のマルチメディアおよび音声ストリーミングアプリケーション、音声メール、並びに他の音声保存システムを含む。駆動力については、大きい容量が必要とされ、かつパケットが失われた状況下での丈夫な性能が要求されている。種々の最近の音声のコード化を標準化する努力は、低レートの音声コード化アルゴリズムの研究および発展を推進する別の直接的な駆動力に当てられている。低レートの音声コードは、許容可能な適用バンド幅ごとに、より多くのチャンネル、すなわちユーザを生成し、低レートの音声コードを適切なチャンネルコーディングの追加の層と結合して、コードの全体的なビット予定値 (bit budget) の仕様に適合させ、チャンネルを誤った状況のもとでも丈夫な性能を発揮させることができる。

【 0 0 0 9 】

より低いビットレートでコード化するために、音声のスペクトル、すなわち周波数領域でコード化する種々の方法が開発され、この方法では音声信号は、時間にしたがって漸進的に変化するスペクトル (time-varying evolution of spectra) として解析される。例えば、R.J. McAulay & T.F. Quatieriによる文献 (Sinusoidal Coding, in Speech Coding and Synthesis ch. 4 (W.B. Kleijin & K.K. Paliwal eds., 1995) 参照。スペクトルコードは、時間にしたがって変化する音声波形を精密にまねるのではなく、1組のスペクトルパラメータを使用して、音声の各入力フレームの短期間の音声スペクトルをモデル化、すなわち予測することを目的とする。スペクトルパラメータはコード化され、音声の出力フレームはデコードされたパラメータを使用して生成される。生成された合成された音声は、元の入力音声波形と整合しないが、同様の知覚品質を与える。この技術においてよく知られている周波数領域コードの例には、マルチバンド励起コード (multiband excitation coder, MBE)、シヌソイド変換コード (sinusoidal transform coder, STC)、高調波コード (harmonic coder, HC) を含む。このような周波数領域のコードは、低ビットレートで使用可能な少数のビットで正確に量子化できるコンパクトな組のパラメータをもつ高品質のパラメータモデルを与える。

【 0 0 1 0 】

それにも関わらず、低ビットレートのコード化は、制限されたコード化分解能、すなわち制限されたコードブック空間に重大な制約を加えて、単一のコード化機構の効果を制限し、コードが、等しい精度の種々の背景条件のもとで、種々のタイプの音声セグメントを表わすことができないようにしている。例えば、従来の低ビットレートの周波数領域のコードは、音声フレームの位相情報を送らない。その代わりに、位相情報は、ランダムな人工的に生成された初期位相値および線形補間技術 (linear interpolation technique) を使用することによって再構成される。例えば、H. Yang、他による文献 (Quadratic Phase Interpolation for Voiced Speech Synthesis in the MBE Model, in 29 Electronic Letters 856-57 (May 1993)) 参照。位相情報は人工的に生成されるので、シヌソイドの振幅は量子化 - 非量子化プロセスによって完全に保持されるときでも、周波数領域のコードによって生成される出力音声は元の入力音声と整合しない (例えば、大半のパルスは同期しない)。したがって、周波数領域のコードでは、例えば信号対雑音比 (signal-to-noise ratio, SNR) または知覚の SNR のような、閉ループの性能尺度 (performance measure) を採用することが難しいことが分かった。

【 0 0 1 1 】

閉ループのモード決定プロセスに関連して低レートの音声のコード化を行なうために、マルチモードコード化技術が採用された。1つのこのようなマルチモードコード化技術は、Amitava Das、他による文献 (Multimode and Variable-Rate Coding of Speech, in Speech Coding and Synthesis ch. 7 (W.B. Kleijin & K.K. Paliwal eds., 1995)) に記載されている。従来のマルチモードコードは異なるモード、すなわちコード化 - デコード化アルゴリズムを、異なるタイプの入力音声フレームへ適用する。各モード、すなわちコード化 - デコード化プロセスは、最も効率的なやり方で、例えば、有声音音声、無声音音声、または背景ノイズ (非音声 (nonspeech)) のような一定のタイプの音声セグメントを表わすために特化される。外部の開ループのモード決定機構は、入力音声フレームを検査し

10

20

30

40

50

て、何れのモードをフレームに適用するかに関して判断する。通常は、閉ループのモード決定は、入力フレームから多数のパラメータを抽出して、一定の時間およびスペクトルの特性に関するパラメータを評価して、この評価に対するモード決定に基づくことによって行われる。したがってモード決定は、出力音声の抽出状態、すなわち出力音声の音声品質または他の性能尺度に関して入力音声にどのくらい近くなるかを前もって知ることなく行われる。

【 0 0 1 2 】

上述に基づいて、位相情報をより精密に推定する低ビットレートの周波数領域のコードを用意することが望ましい。マルチモードの混合領域のコードを用意して、フレームの音声内容に基づいて、一定の音声フレームを時間領域でコード化し、他の音声フレームを周波数領域でコード化することがさらに好都合である。閉ループのコード化モード決定機構にしたがって、一定の音声フレームを時間領域でコード化して、他の音声フレームを周波数領域でコード化することができる混合領域のコードを用意することが、なおいっそう望ましい。したがって、コードによって生成される出力音声と、コードへ入力される元の音声との時間の同期性を保証する、閉ループのマルチモードの混合領域の音声コードが必要とされている。

10

【 発明の概要 】

【 0 0 1 3 】

本発明は、コードによって生成される出力音声と、コードへ入力される元の音声との時間の同期性を保証する、閉ループのマルチモードの混合領域の音声コードに関する。したがって、本発明の1つの態様では、マルチモードの混合領域の音声プロセッサが、少なくとも1つの時間領域コード化モードおよび少なくとも1つの周波数領域コード化モードをもつコードと、コードに接続され、かつ音声プロセッサによって処理されるフレーム内容に基づいてコードのコード化モードを選択するように構成されている閉ループのモード選択デバイスとを含むことが好都合である。

20

【 0 0 1 4 】

本発明の別の態様では、フレームを処理する方法は、各連続する入力フレームへ閉ループのコード化モード選択プロセスを適用して、入力フレームの音声内容に基づいて時間領域コード化モードか、または周波数領域コード化モードの何れかを選択するステップと、入力フレームの音声内容が定常状態の有声音の音声を示すときは、入力フレームを周波数領域でコード化するステップと、入力フレームの音声内容が定常状態の有声音の音声以外のものを示すときは、入力フレームを時間領域でコード化するステップと、周波数領域でコード化されたフレームと入力フレームとを比較して、性能尺度を求めるステップと、性能尺度が所定の閾値より低いときは入力フレームを時間領域でコード化するステップとを含むことが好都合である。

30

【 0 0 1 5 】

本発明の別の態様では、マルチモードの混合領域の音声プロセッサは、閉ループのコード化モード選択プロセスを入力フレームへ適用して、入力フレームの音声内容に基づいて、時間領域コード化モードか、または周波数領域コード化モードの何れかを選択する手段と、入力フレームの音声内容が定常状態の有声音の音声を示すときは、入力フレームを周波数領域でコード化する手段と、入力フレームの音声内容が定常状態の有声音の音声以外のものを示すときは、入力フレームを時間領域でコード化する手段と、周波数領域でコード化されたフレームと入力フレームとを比較して、性能尺度を求める手段と、性能尺度が所定の閾値より低いときは、入力フレームを時間領域でコード化する手段とを含むことが好都合である。

40

【 図面の簡単な説明 】

【 0 0 1 6 】

【 図 1 】 音声コードによって各端部で終端している通信チャンネルのブロック図。

【 図 2 】 マルチモードの混合領域の線形予測 (mixed-domain linear prediction, MDLP) の音声コードにおいて使用できるエンコーダのブロック図。

50

【図 3】マルチモードの M D L P の音声コードにおいて使用できるデコーダのブロック図。

【図 4】図 2 のエンコーダにおいて使用できる M D L P エンコーダによって実行される M D L P のコード化ステップを示すフローチャート。

【図 5】音声コード化決定プロセスを示すフローチャート。

【図 6】閉ループのマルチモードの M D L P の音声コードのブロック図。

【図 7】図 6 のコードまたは図 2 のエンコーダにおいて使用できるスペクトルコードのブロック図。

【図 8】高調波コードのシヌソイドの振幅を示す振幅対周波数のグラフ。

【図 9】マルチモードの M D L P の音声コードにおけるモード決定プロセスを示すフローチャート。

【図 10】音声信号の振幅対時間のグラフ（図 10 a）および線形予測（linear prediction, LP）の残余振幅対時間のグラフ（図 10 b）。

【図 11】閉ループのコード化決定のもとでのレート / モード対フレーム指標のグラフ（図 11 a）、閉ループの決定のもとでの知覚の信号対雑音比（perceptual signal-to-noise ratio, PSNR）対フレーム指標のグラフ（図 11 b）、閉ループのコード化決定がないときのレート / モードおよび P S N R の両者対フレーム指標のグラフ（図 11 c）。

【発明を実施するための形態】

【0017】

図 1 では、第 1 のエンコーダ10は、ディジタル形式の音声サンプル $s(n)$ を受信し、サンプル $s(n)$ をコード化して、伝送媒体12、すなわち通信チャンネル12上で第 1 のデコーダ14へ送る。デコーダ14はコード化された音声サンプルをデコードし、出力された音声信号 $S_{SYNTH}(n)$ を合成する。反対方向で伝送するには、第 2 のエンコーダ16がディジタル形式の音声サンプル $s(n)$ をコード化し、それを通信チャンネル18上で送る。第 2 のデコーダ20はコード化された音声サンプルを受信し、デコードし、合成された出力音声信号 $S_{SYNTH}(n)$ を生成する。

【0018】

音声サンプル $s(n)$ は、この技術において知られている種々の方法、例えばパルスコード変調（pulse code modulation, PMC）、コンパンドされた μ 法、すなわち A 法（companded μ -law, or A-law）を含む方法にしたがって、ディジタル形式にされて量子化された音声信号を表わしている。この技術において知られているように、音声サンプル $s(n)$ は、各々が所定数のディジタル形式の音声サンプル $s(n)$ を含む入力データのフレームへ編成される。例示的な実施形態では、8 キロヘルツのサンプリングレートが採用され、各 20 ミリ秒のフレームは 160 サンプルを含んでいる。別途記載する実施形態では、データ伝送レートはフレームごとに 8 キロビット秒（フルレート）から 4 キロビット秒（2 分の 1 レート）、2 キロビット秒（4 分の 1 レート）、1 キロビット秒（8 分の 1 レート）へ変化することが好都合である。その代わりに、他のデータレートを使用してもよい。ここで使用されているように、“フルレート（full rate）”または“高レート（high rate）”という用語は、通常は、8 キロビット秒以上のデータレートを指し、“2 分の 1 レート”または“低レート”という用語は、通常は、4 キロビット秒以下のデータレートを指す。比較的少ない音声情報を含むフレームに対して、より低いビットレートを選択的に採用できるので、データ伝送レートを変化させることが好都合である。当業者によって理解されるように、他のサンプリングレート、フレームサイズ、およびデータ伝送レートを使用してもよい。

【0019】

第 1 のエンコーダ10および第 2 のデコーダ20は共に第 1 の音声コード、すなわち音声コーデックを含む。同様に、第 2 のエンコーダ16および第 1 のデコーダ14は共に第 2 の音声コードを含む。音声コードはディジタル信号プロセッサ（digital signal processor, DSP）、特定用途向け集積回路（application-specific integrated circuit, ASIC）、離散的ゲート論理（discrete gate logic）、ファームウェア、または従来のプログラマブル

10

20

30

40

50

ソフトウェアモジュールおよびマイクロプロセッサで構成されていてもよいことが分かるであろう。ソフトウェアモジュールは、RAMメモリ、フラッシュメモリ、レジスタ、またはこの技術において知られている他の形態の書き込み可能な記憶媒体内にある。その代わりに、従来のプロセッサ、制御装置、または状態機械をマイクロプロセッサと置換してもよい。音声のコード化のために特別に設計されたASICの例は、本発明の譲受人に譲渡され、かつここでは参考文献として全面的に取り上げている米国特許第5,727,123号、および1994年2月16日に出願され、本発明の譲受人に譲渡され、かつここでは参考文献として全面的に取り上げている米国特許出願第08/197,417号（発明の名称：VOCODER ASIC）に記載されている。

【0020】

1つの実施形態にしたがって、図2に示されているように、音声コード内で使用できるマルチモードの混合領域の線形予測（mixed-domain linear prediction, MDLP）エンコーダ100は、モード決定モジュール102、ピッチ推定モジュール104、線形予測（linear prediction, LP）解析モジュール106、LP解析フィルタ108、LP量子化モジュール110、およびMDLP残余エンコーダ112を含む。入力音声フレーム $s(n)$ は、モード決定モジュール102、ピッチ推定モジュール104、LP解析モジュール106、およびLP解析フィルタ108へ供給される。モード決定モジュール102は、各入力音声フレーム $s(n)$ の周期性および他の抽出パラメータ、例えばエネルギー、スペクトルチルト、ゼロ交差レート、などに基づいて、モード指標 I_M およびモードMを生成する。周期性にしたがって音声フレームを分類する種々の方法は、米国特許出願第08/815,354号（発明の名称：METHOD AND APPARATUS FOR PERFORMING REDUCED RATE VARIABLE RATE VOCODING）に記載されており、これは1997年3月11日に出願され、本発明の譲受人に譲渡され、ここでは参考文献として全面的に取り上げている。このような方法は、米国電気通信工業会の業界暫定標準（Telecommunication Industry Association Industry Interim Standards）のTIA/EIA IS-127およびTIA/EIA IS-733にも採用されている。

10

20

【数 1】

ピッチ推定モジュール 104 は、各入力音声フレーム $s(n)$ に基づいて、ピッチ指標 I_P およびラグ値 P_0 を生成する。LP 解析モジュール 106 は、各入力音声フレーム $s(n)$ に対して線形予測解析を行って、LP パラメータ a を生成する。LP パラメータ a は LP 量子化モジュール 110 へ供給される。LP 量子化モジュール 110 はさらにモード M を受信し、モードに依存するやり方で量子化プロセスを行なう。LP 量子化モジュール 110 は、LP 指標 I_{LP} および量子化された LP パラメータ \hat{a} を生成する。LP 解析フィルタ 108 は、入力音声フレーム $s(n)$ に加えて、量子化された LP パラメータ \hat{a} を受信する。LP 解析フィルタ 108 は LP 残余信号 $R[n]$ を生成し、LP 残余信号 $R[n]$ は、入力音声フレーム $s(n)$ と、量子化された線形予測パラメータ \hat{a} に基づいて再構成された音声との誤差を表わしている。LP 残余信号 $R[n]$ 、モード M 、および量子化された LP パラメータ \hat{a} は MD LP 残余エンコーダ 112 へ供給される。これらの値に基づいて、MD LP 残余エンコーダ 112 は、図 4 のフローチャートを参照して別途記載するステップにしたがって、残余指標 I_R および量子化された残余信号 $\hat{R}[n]$ を生成する。

10

20

【数 2】

図 3 では、音声コードにおいて使用されるデコーダ 200 は、LP パラメータデコーディングモジュール 202、残余デコーディングモジュール 204、モードデコーディングモジュール 206、および LP 解析フィルタ 208 を含む。モードデコーディングモジュール 206 は、モード指標 I_M を受信してデコードし、モード M を生成する。LP パラメータデコーディングモジュール 202 はモード M および LP 指標 I_{LP} を受信する。LP パラメータデコーディングモジュール 202 は、受信した値をデコードして、量子化された LP パラメータ \hat{a} を生成する。残余デコーディングモジュール 204 は、残余指標 I_R 、ピッチ指標 I_P 、およびモード指標 I_M を受信する。残余デコーディングモジュール 204 は受信した値をデコードして、量子化された残余信号 $\hat{R}[n]$ を生成する。量子化された残余信号 $\hat{R}[n]$ および量子化された LP パラメータ \hat{a} は LP 合成フィルタ 208 へ供給され、LP 合成フィルタ 208 はデコードされた出力音声信号 $\hat{s}[n]$ を合成する。

10

20

【0021】

MDLP 残余エンコーダ 112 を除いて、図 2 のエンコーダ 100 および図 3 のデコーダ 200 の種々のモジュールの動作および構成はこの技術において知られており、上述の米国特許第 5,414,796 号および LB. Rabiner & R.W. Schafer による文献 (Digital Processing of Speech Signals 396-453 (1978)) に記載されている。

【0022】

1 つの実施形態にしたがって、MDLP エンコーダ (図示されていない) は、図 4 のフローチャートに示したステップを実行する。MDLP エンコーダは、図 2 の MDLP 残余エンコーダ 112 であってもよい。ステップ 300 では、MDLP エンコーダは、モード M がフルレート (full rate, FR) であるか、4 分の 1 レート (quarter rate, QR) であるか、または 8 分の 1 レート (eighth rate, ER) であるかを検査する。モード M が FR、QR、または ER であるときは、MDLP エンコーダはステップ 302 へ進む。ステップ 302 では、MDLP エンコーダは対応するレート (M の値に依存して - FR, QR、または ER) を残余指標 I_R へ適用する。時間領域のコード化は、FR モードでは高精度で高レートのコード化であり、かつ CELP のコード化であることが好都合であるが、この時間領域のコード化は、LP の残余フレーム、またはその代わりに音声フレームへ適用される。次にフレームは (デジタル対アナログ変換および変調を含む別の信号処理の後で) 送られる。1 つの実施形態では、フレームは、予測誤差を表わす LP 残余フレームである。代替の実施形態では、フレームは、音声サンプルを表わす音声フレームである。

30

40

【0023】

他方で、ステップ 300 では、モード M が FR、QR、または ER でなかったとき (すなわち、モード M が 2 分の 1 レート (half rate, HR) であるとき)、MDLP エンコーダはステップ 304 へ進む。ステップ 304 では、スペクトルのコード化、好ましくは高調波のコード化を 2 分の 1 のレートで LP 残余、またはその代わりに音声信号へ適用する。次に MDLP エンコーダはステップ 306 へ進む。ステップ 306 では、コード化された音声をデコードして、それを元の入力フレームと比較することによって、ひずみ尺度 D を得る。次に MDLP エンコーダは、ステップ 308 へ進む。ステップ 308 では、ひずみ尺度 D は所定の閾値

50

Tと比較される。ひずみ尺度Dが閾値Tよりも大きいときは、2分の1レートのスペクトル的にコード化されたフレームについて、対応する量子化されたパラメータが変調されて、送られる。他方で、ひずみ尺度Dが閾値T以下であるときは、MDLPエンコーダはステップ310へ進む。ステップ310では、デコードされたフレームは、この時間領域においてフルレートで再びコード化される。従来の高レートで高精度のコード化アルゴリズム、例えば好ましくはCELPのコード化を使用してもよい。次に、フレームと関係するFRモードの量子化されたパラメータが変調されて、送られる。

【0024】

図5のフローチャートに示したように、次に1つの実施形態にしたがって閉ループのマルチモードのMDLPの音声コードは、音声サンプルを処理して送る1組のステップにしたがう。ステップ400では、音声コードは、連続するフレーム内の音声信号のデジタルサンプルを受信する。所与のフレームを受信すると、音声コードはステップ402へ進む。ステップ402では、音声コードはフレームのエネルギーを検出する。エネルギーはフレームの音声活動(speech activity)の尺度である。音声検出は、デジタル形式の音声サンプルの振幅の平方を加算して、生成されたエネルギーを閾値と比較することによって行なわれる。1つの実施形態では、背景ノイズの変化レベルに基づいて閾値を採用する。例示的な可変閾値の音声活動検出器は、上述の米国特許第5,414,796号に記載されている。若干の無声音の音声は非常に低いエネルギーのサンプルであり、誤って背景ノイズとしてコード化されてしまうことがある。このようなことが発生するのを防ぐために、上述の米国特許第5,414,796号に記載されているように、低エネルギーサンプルのスペクトルのチ

10

20

【0025】

フレームのエネルギーを検出した後で、音声コードはステップ404へ進む。ステップ404では、音声コードは、音声情報を含んでいるかについてフレームを分類するのに、検出されたフレームエネルギーが十分であるかどうかを判断する。検出されたフレームエネルギーが所定の閾値レベルよりも低いときは、音声コードはステップ406へ進む。ステップ406では、音声コードは背景ノイズ(すなわち、非音声、または黙音)としてフレームをコード化する。1つの実施形態では、背景ノイズのフレームは、8分の1レート、すなわち1キロビット秒でコード化される時間領域である。ステップ404では、検出されたフレームのエネルギーが所定の閾値レベル以上であるとき、フレームは音声として分類され、音声コードはステップ408へ進む。

30

【0026】

ステップ408では、音声コードは、フレームが周期的であるかどうかを判断する。周期性を判断する種々の既知の方法には、例えばゼロ交差の使用および正規化された自動相関関数(normalized autocorrelation function, NACF)の使用を含む。とくに、ゼロ交差およびNACFを使用して、周期性を検出することは、米国出願第08/815,354号(発明の名称:METHOD AND APPARATUS FOR PERFORMING REDUCED RATE VARIABLE RATE VOCODING)に記載されており、これは1997年3月11日に出願され、本発明の譲受人に譲渡され、ここでは参考文献として全面的に取り上げている。さらに加えて、無声音の音声から有声音の音声を区別するのに使用される上述の方法は、米国電気通信工業会の業界暫定標準(Telecommunication Industry Association Industry Interim Standards)のTIA/EIA IS-127およびTIA/EIA IS-733に採用されている。ステップ408においてフレームが周期的でない

40

【0027】

ステップ412では、音声コードは、例えば上述の米国特許出願第08/815,354号に記載されているように、この技術において知られている周期性検出方法を使用して、フレームが

50

十分に周期的であるかどうかを判断する。フレームが十分に周期性でないと判断されるときは、音声コードはステップ414へ進む。ステップ414では、フレームは遷移音声 (transition speech) (すなわち、無声音の音声から有声音の音声への遷移) として時間領域でコード化される。1つの実施形態では、遷移音声フレームはフルレート、すなわち 8 キロビット秒で時間領域でコード化される。

【0028】

音声コードは、ステップ412においてフレームが十分に周期的であると判断すると、ステップ416へ進む。ステップ416では、音声コードは有声音の音声としてフレームをコード化する。1つの実施形態では、有声音の音声フレームは、とくに2分の1レート、すなわち 4 キロビット秒でスペクトル的にコード化される。図7を参照して別途記載するように、有声音の音声フレームは、高調波のコードでスペクトル的にコード化されることが好都合である。その代わりに、他のスペクトルコードは、この技術において知られているように、例えばシヌソイド変換コード (sinusoidal transmission coder) またはマルチバンド励起コード (multiband excitation coder) として使用されることが好都合である。次に音声コードはステップ418へ進む。ステップ418では、音声コードはコード化された有声音の音声フレームをデコードする。次に音声コードはステップ420へ進む。ステップ420では、デコードされた有声音の音声フレームを、このフレームの対応する入力音声サンプルと比較して、合成された音声のひずみ尺度を得て、2分の1レートの有声音音声のスペクトルコード化モデルが許容限度内で動作しているかどうかを判断する。次に音声コードはステップ422へ進む。

10

20

【0029】

ステップ422では、音声コードは、デコードされた有声音の音声フレームと、このフレームに対応する入力音声フレームとの誤差が所定の閾値より小さいかどうかを判断する。1つの実施形態では、この判断は、図6を参照して別途記載するやり方で行われる。コード化のひずみが所定の閾値よりも低いときは、音声コードはステップ426へ進む。ステップ426では、音声コードは、ステップ416のパラメータを使用して、フレームを有声音の音声として送る。ステップ422では、コード化のひずみが所定の閾値以上であるときは、音声コードはステップ414へ進み、ステップ400において受信したデジタル形式の音声サンプルのフレームを遷移音声としてフルレートで時間領域でコード化する。

30

【0030】

ステップ400ないし410は開ループのコード化決定モードを含むことに注目すべきである。他方で、ステップ412ないし426は閉ループのコード化決定モードを含む。

【0031】

1つの実施形態では、図6に示したように、閉ループのマルチモードのMDLPの音声コードはアナログ対デジタルコンバータ (analog-to-digital converter, A/D) 500を含み、A/D 500はフレームバッファ502に接続され、フレームバッファ502は制御プロセッサ504に接続される。エネルギー計算器506、有声音音声の検出器508、背景ノイズエンコード510、高レートの時間領域エンコード512、および低レートのスペクトルエンコード514は制御プロセッサ504へ接続される。スペクトルデコード516はスペクトルエンコード514に接続され、誤差計算器518はスペクトルデコード516および制御プロセッサ504へ接続される。閾値比較器520は、誤差計算器518および制御プロセッサ504へ接続される。バッファ522はスペクトルエンコード514、スペクトルデコード516、および閾値比較器520へ接続される。

40

【0032】

図6の実施形態では、音声コードの構成要素は、音声コード内にファームウェアまたは他のソフトウェア駆動モジュールとして構成されていることが好都合であり、音声コード自身はDSPまたはASIC内にあることが好都合である。当業者には、音声コードの構成要素は、多数の他の既知のやり方で同様に適切に構成されることが分かるであろう。制御プロセッサ504はマイクロプロセッサであることが好都合であるが、制御装置、状態機械、または離散的論理と共に構成されていてもよい。

50

【 0 0 3 3 】

図 6 のマルチモードのコードでは、音声信号は A / D 500 へ供給される。A / D 500 はアナログ信号をデジタル形式の音声サンプル $S(n)$ へ変換する。デジタル形式の音声サンプルは、フレームバッファ 502 へ供給される。制御プロセッサ 504 は、フレームバッファ 502 からデジタル形式の音声サンプルを得て、それらをエネルギー計算器 506 へ供給する。エネルギー計算器 506 は、次の式にしたがって音声サンプルのエネルギー E を計算する：

【 数 3 】

$$E = \sum_{n=0}^{159} S^2(n)$$

10

【 0 0 3 4 】

なお、フレームは 20 ミリ秒長であり、サンプリングレートは 8 キロヘルツである。計算されたエネルギー E は制御プロセッサ 504 へ送られる。

【 0 0 3 5 】

制御プロセッサ 504 は、計算された音声エネルギーを音声活動 (speech activity) の閾値と比較する。計算されたエネルギーが音声活動の閾値よりも小さいときは、制御プロセッサ 504 はデジタル形式の音声サンプルをフレームバッファ 502 から背景ノイズエンコーダ 510 へ送る。背景ノイズエンコーダ 510 は、背景ノイズの推定値を保持するために必要な最少数のビットを使用して、フレームをコード化する。

20

【 0 0 3 6 】

計算されたエネルギーが音声活動の閾値以上であるときは、制御プロセッサ 504 はデジタル形式の音声サンプルをフレームバッファ 502 から有声音音声の検出器 508 へ方向付ける。有声音音声の検出器 508 は、音声フレームの周期性が、低ビットレートのスペクトルのコード化を使用して効率的なコード化を可能にするかどうかを判断する。音声フレーム内の周期性のレベルを判断する方法は、この技術においてよく知られており、例えば正規化された自動相関関数 (normalized autocorrelation function, NACF) およびゼロ交差の使用を含む。これらの方法および他の方法は、上述の米国特許出願第 08/815,354 号に記載されている。

30

【 0 0 3 7 】

有声音音声の検出器 508 は、スペクトルエンコーダ 514 が効率的にコード化するのに十分な周期性をもつ音声を含むかどうかを示す信号を制御プロセッサ 504 へ供給する。有声音音声の検出器 508 が、音声フレームが十分な周期性を欠いていると判断するとき、制御プロセッサ 504 はデジタル形式の音声サンプルを高レートのエンコーダ 512 へ方向付け、エンコーダ 512 は所定の最大データレートで音声を時間領域でコード化する。1 つの実施形態では、所定の最大データレートは 8 キロビット秒であり、高レートのエンコーダ 512 は CELP のコードである。

40

【 0 0 3 8 】

有声音音声の検出器 508 が最初に、音声信号が、スペクトルエンコーダ 514 が効率的にコード化するのに十分な周期性をもつと判断するとき、制御プロセッサ 504 は、フレームバッファ 502 からスペクトルエンコーダ 514 へデジタル形式の音声サンプルを方向付ける。例示的なスペクトルエンコーダは、図 7 を参照して別途詳しく記載する。

【数 4】

スペクトルエンコーダ 514 は、推定されたピッチ周波数 F_0 、ピッチ周波数の高調波の振幅 A_1 、および音声情報 (voicing information) V_c を抽出する。スペクトルエンコーダ 514 はこれらのパラメータをバッファ 522 およびスペクトルデコーダ 516 へ送る。スペクトルデコーダ 516 は、従来の CELP エンコーダ内のエンコーダのデコーダに似ていることが好都合である。スペクトルデコーダ 516 は、(図 7 を参照して別途記載される) スペクトルデコーディングフォーマットのしたがって、合成された音声サンプル $\hat{s}[n]$ を生成し、合成された音声サンプルを誤差計算器 518 へ供給する。制御プロセッサ 504 は音声サンプル $S(n)$ を誤差計算器 518 へ送る。

10

【数 5】

20

誤差計算器 518 は、次の式にしたがって、各音声サンプル $S(n)$ と、各対応する合成された音声サンプル $\hat{s}[n]$ との間の平均平方誤差 (mean square error, MSE) を計算する：

$$MSE = \sum_{n=0}^{159} (S(n) - \hat{S}(n))^2$$

計算された MSE は閾値比較器 520 へ供給され、閾値比較器 520 は、ひずみレベルが許容範囲内であるかどうか、すなわちひずみレベルが所定の閾値よりも小さいかどうかを判断する。

30

【0039】

計算された MSE が許容範囲内であるときは、閾値比較器 520 は信号をバッファ 522 へ供給し、スペクトル的にコード化されたデータは音声コードから出力される。他方で、MSE が許容限界内でないときは、閾値の比較器 520 は信号を制御プロセッサ 504 へ送り、制御プロセッサ 504 はデジタル形式のサンプルをフレームバッファ 502 から高レートの時間領域のエンコーダ 512 へ方向付ける。時間領域のエンコーダ 512 は、所定の最大レートでフレームをコード化し、バッファ 522 の内容は捨てられる。

40

【0040】

図 6 の実施形態では、採用されたスペクトルのコード化のタイプは高調波のコード化であり、これについては図 7 を参照して別途記載するが、代わりの実施形態では、シヌソイド変換のコード化またはマルチバンド励起のコード化のような、スペクトルのコード化のタイプであってもよい。マルチバンド励起のコード化の使用は、米国特許第 5,195,166 号に記載されており、シヌソイド変換のコード化の使用は、例えば米国特許第 4,865,068 号に記載されている。

【0041】

遷移フレーム、および位相ひずみ閾値が周期性パラメータ以下である有声音フレームで

50

は、図 6 のマルチモードコードはフルレート、すなわち 8 キロビット秒で、高レートの時間領域のコード 512 によって、C E L P のコード化を採用することが好都合である。その代わりに、このようなフレームに対して、他の既知の形態の高レートの時間領域のコード化を使用してもよい。したがって、遷移フレーム（および十分に周期的でない有声音フレーム）は高い精度でコード化され、入力および出力における波形は適切に整合し、位相情報は適切に保持される。1 つの実施形態では、マルチモードコードは、閾値比較器 520 の判断と無関係に、閾値が周期性の尺度を越えている所定数の連続する有声音フレームを処理した後で、各フレームごとに 2 分の 1 レートのスペクトルのコード化からフルレートの C E L P のコード化へスイッチする。

【 0 0 4 2 】

制御プロセッサ 504 に関連して、エネルギー計算器 506 および有声音音声の検出器 508 は閉ループのコード化決定を含むことに注意すべきである。対照的に、制御プロセッサ 504 に関連して、スペクトルエンコード 514、スペクトルデコード 516、誤差計算器 518、閾値比較器 520、およびバッファ 522 は閉ループのコード化決定を含む。

【 0 0 4 3 】

図 7 を参照して記載した 1 つの実施形態では、スペクトルのコード化、好ましくは高調波のコード化を使用して、低ビットレートで十分に周期的な有声音フレームをコード化する。スペクトルコードは、一般的に、周波数領域内の各音声フレームをモデル化してコード化することによって知覚的に重要なやり方で音声スペクトル特性の時間にしたがう漸進的变化（time-evolution）を保持することを試みるアルゴリズムとして規定される。このようなアルゴリズムの本質的な部分では、（ 1 ）スペクトルの解析またはパラメータの推定、（ 2 ）パラメータの量子化、（ 3 ）出力された音声波形とデコードされたパラメータとの合成を行う。したがって、1 組のスペクトルパラメータをもつ短期間の音声スペクトルの重要な特性を保持し、デコードされたスペクトルパラメータを使用して、出力音声を合成することを目的とする。通常は、出力音声は、シヌソイドの重み付けされた和として合成される。シヌソイドの振幅、周波数、および位相は、解析中に推定されるスペクトルパラメータである。

【 0 0 4 4 】

“ 合成による解析 ” は C E L P のコード化においてよく知られた技術であるが、この技術はスペクトルのコード化には利用されていない。合成による解析がスペクトルコードに適用されない主な理由は、初期位相の情報の損失によって、音声モデルが知覚の観点から適切に機能していても、合成された音声の平均二乗エネルギー（mean square energy, MSE）が高いからである。したがって、初期位相を正確に生成すると、音声サンプルと再構成された音声とを直接に比較して、音声モデルが音声フレームを正確にコード化しているかどうかを判断できるといった別の長所がある。

【 0 0 4 5 】

スペクトルのコード化では、出力された音声フレームは次に示すように合成することができる：

$$S[n] = S_v[n] + S_{uv}[n], \quad n = 1, 2, \dots, N,$$
 なお、 N は 1 フレーム当りのサンプル数であり、 S_v および S_{uv} は、それぞれ有声音成分および無声音成分である。シヌソイド和合成プロセス（sum-of-sinusoid synthesis process）は次の式に示すように有声音成分を生成する：

10

20

30

40

【数 6】

$$s[n] = \sum_{k=1}^L A(k, n) \cdot \cos(2\pi n f_k + \theta(k, n))$$

なお、 L はシヌソイドの合計数であり、

f_k は短期間のスペクトルにおける目的の周波数であり、

$A(k, n)$ はシヌソイドの振幅であり、

$\theta(k, n)$ はシヌソイドの位相である。

10

【0046】

振幅、周波数、および位相パラメータは、スペクトル解析プロセスによって入力フレームの短期間のスペクトルから推定される。無声音成分は、単一のシヌソイド和合成において有声音部分と一緒に生成されるか、または専用の無声音合成プロセスによって別々に計算され、 S_v へ再び加えられる。

【0047】

図7の実施形態では、高調波コードと呼ばれる特定のタイプのスペクトルコードを使用して、低ビットレートで十分に周期的な有声音フレームをスペクトル的にコード化する。高調波のコードは、シヌソイド和としてフレームを特徴付け、フレームの小さいセグメントを解析する。シヌソイド和の中の各シヌソイドは、フレームのピッチ F_0 の整数倍の周波数をもつ。代替の実施形態では、高調波のコード以外の特定のタイプのスペクトルコードを使用し、各フレームに対するシヌソイド周波数は、0ないし2の1組の実数から得られる。図7の実施形態では、和の中の各シヌソイドの振幅および位相が選択されることが好都合であり、その結果、図8のグラフによって示したように、和は1期間において信号と最良に整合する。高調波のコードは一般的に外部の分類を採用し、各入力音声フレームは有声音または無声音として表示する。有声音フレームでは、シヌソイドの周波数は推定されたピッチ(F_0)の高調波に制限され、すなわち $f_k = k F_0$ である。無声音の音声では、短期間のスペクトルのピークを使用して、シヌソイドを判断する。次の式に示すように、振幅および位相が補間されて、フレームにおいて漸進的変化をまねる：

20

30

【数 7】

$$A(k, n) = C_1(k) \cdot n + C_2(k)$$

$$\theta(k, n) = B_1(k) \cdot n^2 + B_2(k) \cdot n + B_3(k)$$

係数 $[C_i(k), B_i(k)]$ は、ウィンドウ処理された入力音声フレームの短期間のフーリエ変換 (short-term Fourier Transform (STFT)) から、特定された周波数の位置 $f_k (= k f_0)$ における振幅、周波数、および位相の瞬間値によって推定される。

40

【0048】

シヌソイドごとに送られるパラメータは振幅および周波数である。位相は送られないが、その代わりに、例えば準位相モデル (quadratic phase model)、または位相の従来の多項式表現を含むいくつかの既知の技術にしたがってモデル化される。

【0049】

図7に示されているように、高調波コードはピッチ抽出器600を含み、ピッチ抽出器600はウィンドウ処理論理602へ接続され、ウィンドウ処理論理602は離散フーリエ変換 (Discrete Fourier Transform, DFT)、および高調波解析論理604へ接続される。入力として音

50

声サンプル $S(n)$ を受信するピッチ抽出器600は、DFTおよび高調波解析論理604へも接続される。DFTおよび高調波解析論理604は、残余エンコーダ606へ接続される。ピッチ抽出器600、DFTおよび高調波解析論理604、並びに残余エンコーダ606は、パラメータ量子化器608へそれぞれ接続される。パラメータ量子化器608はチャンネルエンコーダ610へ接続され、チャンネルエンコーダ610は送信機612へ接続される。送信機612は、例えば、符号分割多重アクセス (code division multiple access, CDMA) のような標準の無線周波数 (radio-frequency, RF) のインターフェイスによって空中インターフェイス (over-the-air interface) 上で、受信機614へ接続される。受信機614はチャンネルデコーダ616へ接続され、チャンネルデコーダ616は非量子化器618へ接続される。非量子化器618はシヌソイド和音声合成器620へ接続される。シヌソイド和音声合成器620へさらに接続されるのは位相推定器622であり、位相推定器622は入力として前フレーム情報を受信する。シヌソイド和音声合成器620は合成された音声出力 $S_{SYNTH}(n)$ を生成するように構成されている。

10

【0050】

ピッチ抽出器600、ウインドウ処理論理602、DFTおよび高調波解析論理604、残余エンコーダ606、パラメータ量子化器608、チャンネルエンコーダ610、チャンネルデコーダ616、非量子化器618、シヌソイド和音声合成器620、並びに位相推定器622は、例えばファームウェアまたはソフトウェアモジュールを含む、当業者によく知られている種々の異なるやり方で構成することができる。送信機612および受信機614は、当業者には知られている対応する標準のRFの構成要素で実行されていてもよい。

20

【0051】

図7の高調波コーダでは、入力サンプル $S(n)$ はピッチ抽出器600によって受信され、ピッチ抽出器600はピッチ周波数情報 F_0 を抽出する。次にサンプルは、ウインドウ処理論理602によって適切なウインドウ処理関数によって乗算され、音声フレームの小さいセグメントの解析を可能にしている。ピッチ抽出器600によって供給されるピッチ情報を使用して、DFTおよび高調波解析論理604はサンプルのDFTを計算して、複合のスペクトル点を生成し、この複合のスペクトル点から、図8のグラフによって示されているように、高調波の振幅 A_I を抽出し、なお図8において、 L は高調波の合計数を示している。DFTは残余エンコーダ606へ供給され、残余エンコーダ606は音声情報 (voicing information) V_c を抽出する。

30

【0052】

V_c パラメータは、図8に示されているように、周波数軸上の点を示し、 V_c がより高くなると、スペクトルは無声音の音声信号の特性を示し、最早高調波ではなくなることに注意すべきである。対照的に、点 V_c より低くなると、スペクトルは高調波であり、有声音の音声の特性を示す。

【0053】

A_I , F_0 , および V_c の成分は、パラメータ量子化器608へ供給され、パラメータ量子化器608では情報を量子化する。量子化された情報はパケットの形態でチャンネルエンコーダ610へ供給され、チャンネルエンコーダ610では、例えばハーフレート、すなわち4キロビット秒のような低ビットレートでパケットを量子化する。パケットは送信機612へ供給され、送信機612はパケットを変調して、生成された信号を受信機614へ空中で (over the air) 送る。受信機614は信号を受信して、復調して、コード化されたパケットをチャンネルデコーダ616へ送る。チャンネルデコーダ616はパケットをデコードして、デコードされたパケットを非量子化器618へ供給する。非量子化器618は情報を非量子化する。情報はシヌソイド和音声合成器620へ供給される。

40

【0054】

シヌソイド和音声合成器620は、 $S[n]$ についての上述の式にしたがって短期間の音声スペクトルをモデル化する複数のシヌソイドのモデリングを合成するように構成されている。シヌソイド f_k の周波数は、基本周波数 F_0 の倍数または高調波であり、準周期的な (すなわち、遷移の) 有声音の音声セグメントに対するピッチの周期性をもつ周波数であ

50

る。

【 0 0 5 5 】

さらに加えて、シヌソイド和の音声合成器620は位相推定器622から位相情報を受信する。位相推定器622は前フレームの情報、すなわち直前フレームについての A_I , F_0 , および V_c のパラメータを受信する。位相推定器622は、前フレームの再構成された N のサンプルも受信し、なお N はフレーム長（すなわち、 N は 1 フレーム当りのサンプル数）である。位相推定器622は、前フレームの情報に基づいて、フレームの初期位相を判断する。初期位相の判断は、シヌソイド和の音声合成器620へ供給される。現在のフレームに関する情報と、過去のフレーム情報に基いて位相推定器622によって行なわれた初期位相の計算とを基にして、シヌソイド和音声合成器620は上述のように音声フレームを生成する。

10

【 0 0 5 6 】

既に記載したように、高調波のコーダは、前フレームの情報を使用して、位相がフレームからフレームへ線形に変化することを予測することによって、音声フレームを合成、すなわち再構成する。上述の合成モデルは、一般的に準位相モデルと呼ばれており、このような合成モデルでは、係数 $B_3(k)$ は、現在の有声音フレームの初期位相が合成されていることを表わしている。位相を判断するとき、従来の高調波のコーダは初期位相をゼロに設定するか、または初期位相値をランダムに、あるいは疑似ランダム生成方法を使用して生成する。位相をより正確に予測するために、位相推定器622は、直前のフレームが有声音の音声フレーム（すなわち、十分に周期的なフレーム）であるか、または遷移音声フレームであるかに依存して、初期位相を判断するための2つの可能な方法の一方を使用する。前フレームが有声音の音声フレームであったときは、このフレームの推定された最終位相値は、現在のフレームの初期位相値として使用される。他方で、前フレームが遷移フレームとして分類されたときは、現在のフレームの初期位相値は、前フレームのスペクトルから得られ、これは前フレームのデコーダ出力の DFT を行なうことによって得られる。したがって位相推定器622は、（遷移フレームである前フレームがフルレートで処理されたので）既に使用可能である正確な位相情報を使用できる。

20

【 0 0 5 7 】

1つの実施形態では、閉ループのマルチモードのMDLPの音声コーダは、図9のフローチャート内に示されている音声処理ステップにしたがう。音声コーダは、最も適切なコード化モードを選択することによって、各入力音声フレームのLPの残余をコード化する。一定のモードは時間領域内でLPの残余、すなわち音声の残余をコード化し、一方で他のモードは周波数領域内でLPの残余、すなわち音声の残余を表わす。モードの組には、遷移フレームに対するフルレートの時間領域（Tモード）；有声音フレームに対する2分の1レートの周波数領域（Vモード）；無声音フレームに対する4分の1レートの時間領域（Uモード）；およびノイズフレームに対する8分の1レートの時間領域（Nモード）がある。

30

【 0 0 5 8 】

当業者には、図9に示したステップにしたがうことによって、音声信号または対応するLPの残余がコード化されることが分かるであろう。ノイズ、無声音、遷移、および有声音の音声の波形特性は、図10aのグラフにおいて時間関数として参照することができる。ノイズ、無声音、遷移、および有声音のLPの残余の波形特性は、図10bのグラフにおいて時間関数として参照することができる。

40

【 0 0 5 9 】

ステップ700では、4つのモード（T、V、U、またはN）の何れか1つに関して、閉ループのモード決定を行って、入力音声の残余 $S(n)$ へ適用する。Tモードが適用されるときは、ステップ702では、時間領域においてTモード、すなわちフルレートで音声の残余が処理される。Uモードが適用されるときは、ステップ704で、時間領域においてUモード、すなわち4分の1レートで音声の残余が処理される。Nモードが適用されるときは、ステップ706では、時間領域においてNモード、すなわち8分の1レートで音声の残

50

余が処理される。Vモードが適用されるときは、ステップ708では、周波数領域においてVモードで、すなわち2分の1レートで音声の残余が処理される。

【0060】

ステップ710では、ステップ708でコード化された音声デコードされ、入力音声の残余 $S(n)$ と比較され、性能尺度 D が計算される。ステップ712では、性能尺度 D が所定の閾値 T と比較される。性能尺度 D が閾値 T 以上であるときは、ステップ714では、ステップ708においてスペクトル的にコード化された音声の残余は送信を許可される。他方では、性能尺度 D が閾値 T よりも小さいときは、ステップ716では、入力音声の残余 $S(n)$ はTモードで処理される。別の実施形態では、性能尺度は計算されず、閾値は規定されない。その代わりに、所定数の音声残余フレームがVモードで処理された後で、次のフレームはTモードで処理される。

10

【0061】

図9に示した決定のステップでは、高ビットレートのTモードを必要なときだけ使用して、より低いビットレートのVモードで有声音の音声セグメントの周期性を活用することができ、一方でVモードが適切に実行されないときは、フルレートにスイッチすることによって品質の低下を防ぐことが好都合である。したがって、フルレートの音声品質に近く非常に高い音声品質を、フルレートよりも相当に低い平均レートで生成することができる。さらに、選択された性能尺度および選ばれた閾値によって、目標の音声品質を制御することができる。

【0062】

20

Tモードへの“更新”は、モデル位相追跡を入力音声の位相追跡の近くに維持することによって、後でVモードを適用する動作を向上することができる。Vモードの性能が不適切であるときは、ステップ710および712の閉ループの性能検査はTモードへスイッチし、初期位相値を“リフレッシュ”して、モデルの位相追跡を元の入力音声位相追跡に再び近付けることによって、次のVモードの処理の性能を向上することができる。例えば、図11aないしcのグラフに示したように、開始から5番目のフレームは、使用されているPSNRのひずみ尺度によって証明されているように、Vモードで適切に働かない。その結果、閉ループの決定および更新がないときは、モデル化された位相追跡は元の入力音声位相追跡から相当に外れ、図11cに示したように、PSNRを相当に劣化する。さらに、Vモードで処理される次のフレームの性能は劣化する。しかしながら、閉ループの決定のもとでは、5番目のフレームは、図11aに示したように、Tモードの処理へスイッチされる。5番目のフレームの性能は、図11bに示したように、PSNRにおける向上によって証明されているように、更新によって相当に向上する。さらに加えて、Vモードのもとで処理される次のフレームの性能も向上する。

30

【0063】

図9に示した決定のステップでは、非常に正確な初期位相推定値を与えることによって、Vモードの表現品質を向上し、生成されたVモードの合成された音声の残余信号は元の入力音声の残余 $S(n)$ と正確に時間的に整合することを保証する。最初のVモードで処理された音声の残余セグメントにおける初期位相は、次に示すやり方で直前のデコードされたフレームから求められる。各高調波では、前フレームがVモードで処理されたときは、初期位相は前フレームの推定された最終位相に等しく設定される。各高調波では、前フレームがTモードで処理されたときは、初期位相は前フレームの実際の高調波の位相に等しく設定される。前フレームの実際の高調波の位相は、全ての前フレームを使用して過去のデコードされた残余のDFTをとることによって求められる。その代わりに、前フレームの実際の高調波の位相は、前フレームの種々のピッチ期間を処理することによって、ピッチが同期するやり方で、過去のデコードされたフレームのDFTをとることによって求められる。

40

【0064】

本明細書では、斬新な閉ループのマルチモードの混合領域の線形予測(mixed-domain linear prediction, MDLP)の音声コードを記載した。当業者には、ここに開示した実施形

50

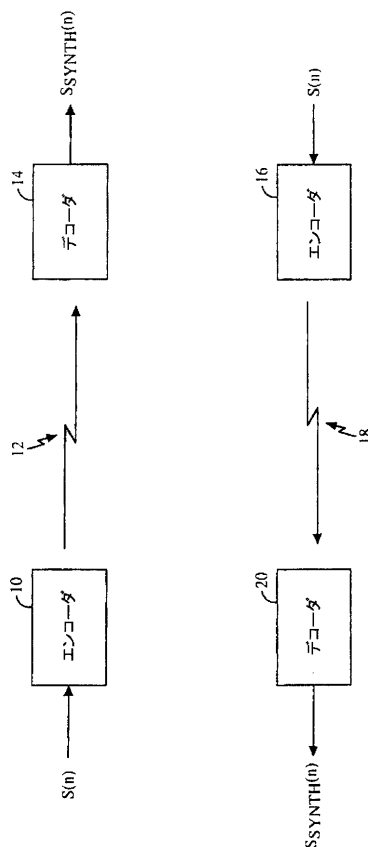
態に関して記載した種々の例示的な論理ブロックおよびアルゴリズムのステップが、デジタル信号プロセッサ (digital signal processor, DSP)、特定用途向け集積回路 (application specific integrated circuit, ASIC)、離散的ゲートまたはトランジスタ論理、例えばレジスタおよび FIFO のような離散的ハードウェア構成要素、1組のファームウェア命令を実行するプロセッサ、または従来のプログラマブルソフトウェアモジュールおよびプロセッサで構成または実行できることが分かるであろう。プロセッサは、マイクロプロセッサであることが好都合であるが、その代わりに従来のプロセッサ、制御装置、マイクロプロセッサ、または状態機械であってもよい。ソフトウェアモジュールは、RAMメモリ、フラッシュメモリ、レジスタ、またはこの技術において知られている他の形態の書き込み可能な記憶媒体内にあってもよい。当業者にはさらに、上述の記述全体で参照したデータ、命令、コマンド、情報、信号、ビット、符号、およびチップが、電圧、電流、電磁波、磁界または磁粒、光の範囲または粒子 (optical field or particles)、あるいはその組み合わせによって都合よく表わされることが分かるであろう。

10

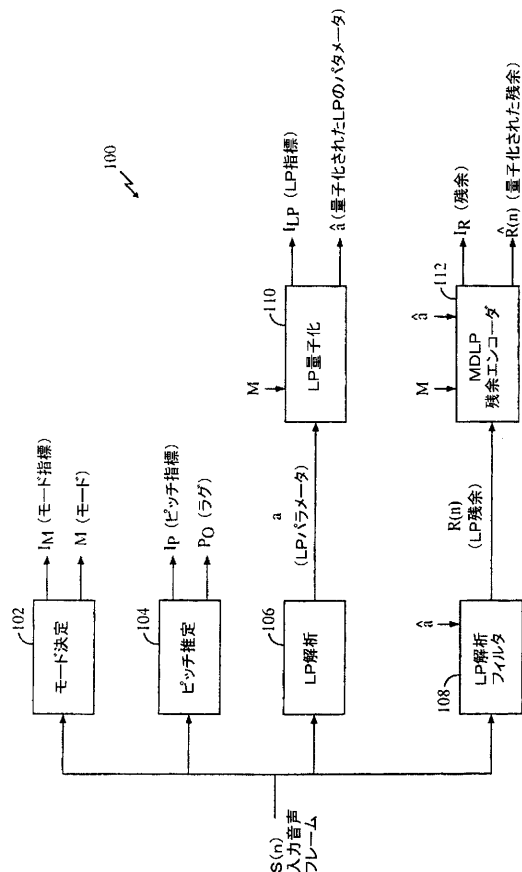
【0065】

本明細書では、本発明の好ましい実施形態を示し、記載した。しかしながら、当業者の一人には、ここに記載した実施形態に対して、本発明の意図または技術的範囲から逸脱せずに多数の変更を加えられることが分かるであろう。したがって、本発明は、特許請求項にしたがうことを除いて制限されない。

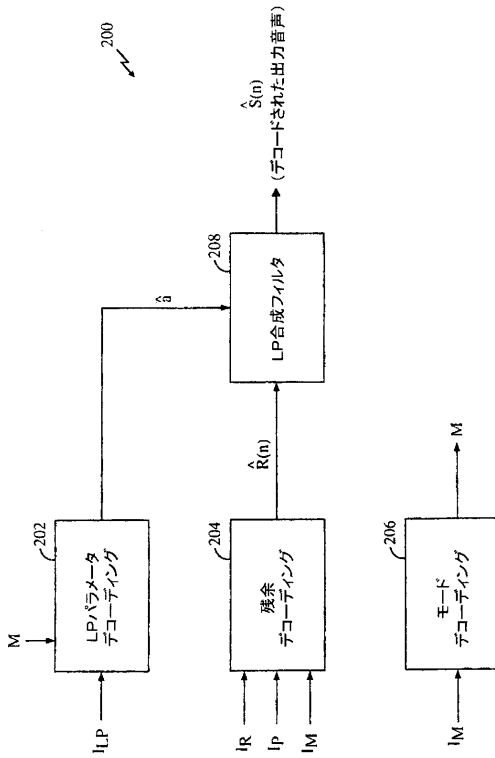
【図1】



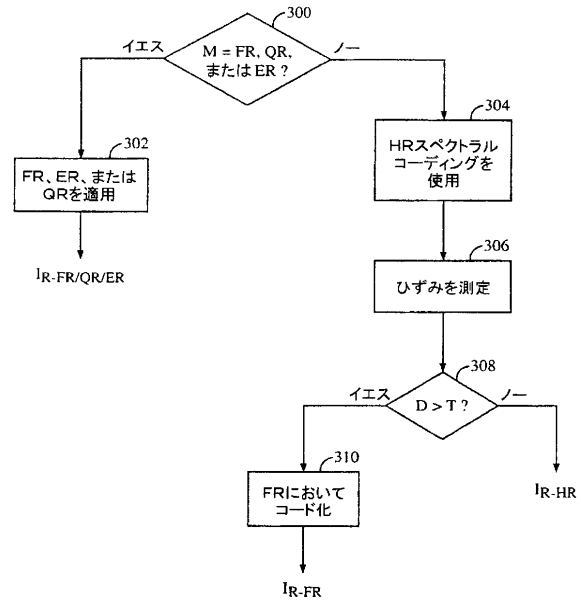
【図2】



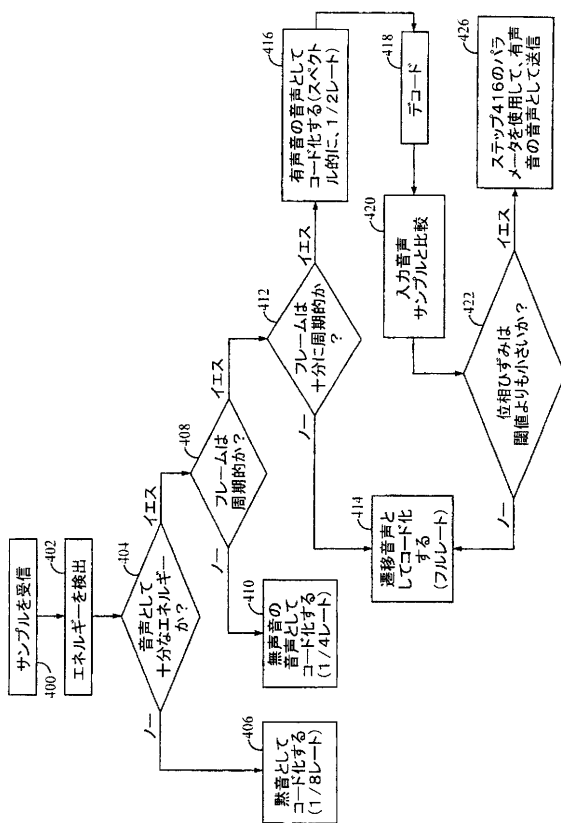
【図 3】



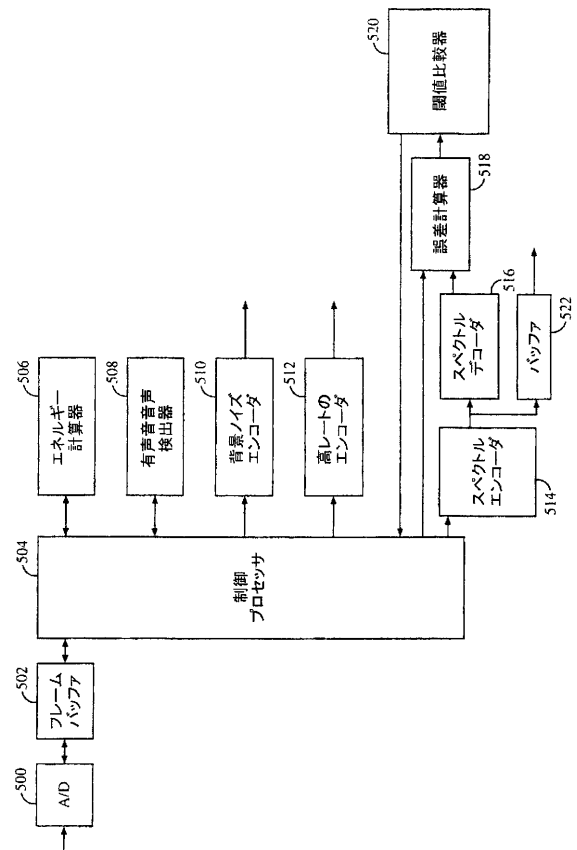
【図 4】



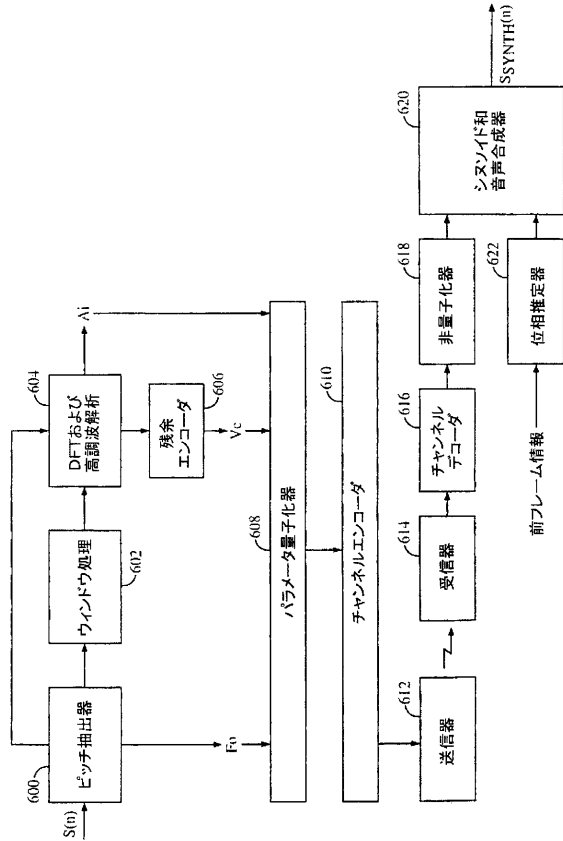
【図 5】



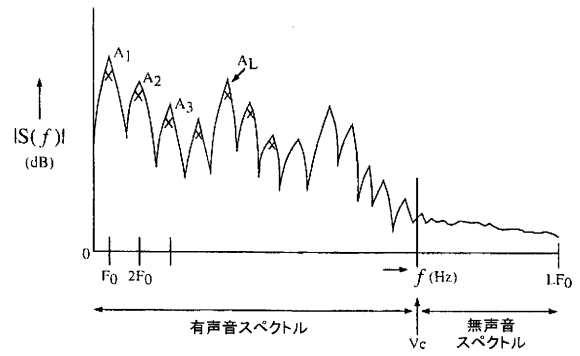
【図 6】



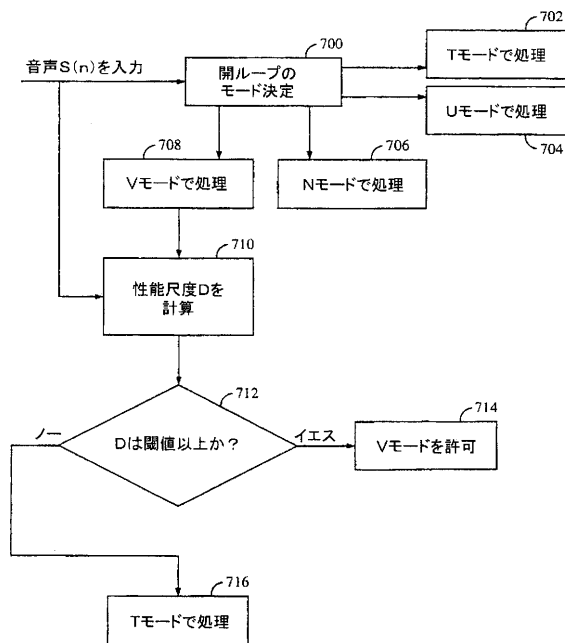
【図 7】



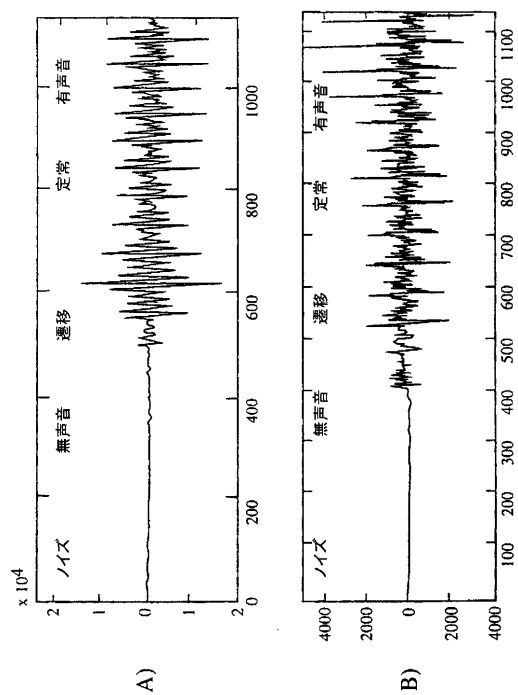
【図 8】



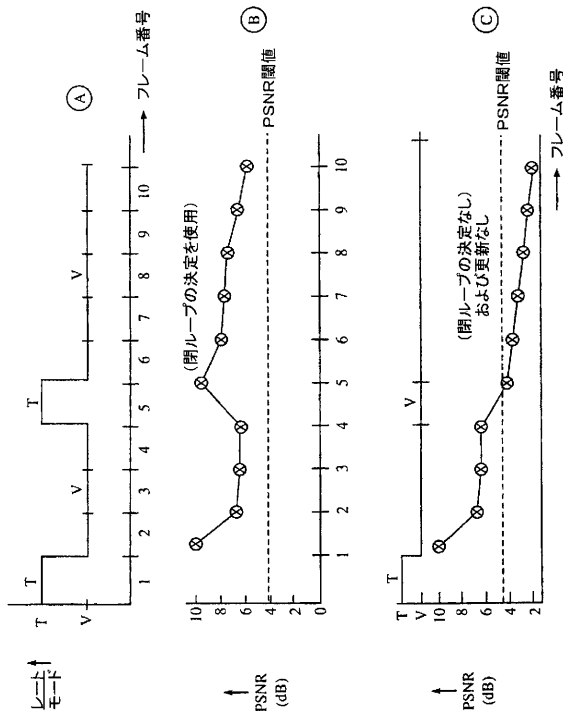
【図 9】



【図 10】



【図 1 1】



【手続補正書】

【提出日】平成22年12月8日(2010.12.8)

【手続補正 1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項 1】

少なくとも 1 つの時間領域コード化モードおよび少なくとも 1 つの周波数領域コード化モードをもつコードと、

コードに接続され、かつ音声プロセッサによって処理されるフレーム内容に基づいてコードのコード化モードを選択するように構成されている閉ループのモード選択デバイスとを含むマルチモードの混合領域の音声プロセッサ。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】0 0 6 5

【補正方法】変更

【補正の内容】

【0 0 6 5】

本明細書では、本発明の好ましい実施形態を示し、記載した。しかしながら、当業者の一人には、ここに記載した実施形態に対して、本発明の意図または技術的範囲から逸脱せずに多数の変更を加えられることが分かるであろう。したがって、本発明は、特許請求項にしたがうことを除いて制限されない。

付記

(1) 少なくとも1つの時間領域コード化モードおよび少なくとも1つの周波数領域コード化モードをもつコードと、

コードに接続され、かつ音声プロセッサによって処理されるフレーム内容に基づいてコードのコード化モードを選択するように構成されている閉ループのモード選択デバイスとを含むマルチモードの混合領域の音声プロセッサ。

(2) コードが、音声フレームをコード化する (2) 記載の音声プロセッサ。

(3) コードが、音声フレームの線形予測残余をコード化する (3) 記載の音声プロセッサ。

(4) 少なくとも1つの時間領域コード化モードが、第1のコード化レートでフレームをコード化するコード化モードを含み、少なくとも1つの周波数領域コード化モードが、第2のコード化レートでフレームをコード化するコード化モードを含み、第2のコード化レートが第1のコード化レートよりも低い (1) 記載の音声プロセッサ。

(5) 少なくとも1つの周波数領域コード化モードが、高調波のコード化モードを含む (1) 記載の音声プロセッサ。

(6) コードに接続された比較回路であって、コード化されていないフレームを、少なくとも1つの周波数領域コード化モードでコード化されたフレームと比較して、比較に基づいて性能尺度を生成する比較回路をさらに含み、コードが、性能尺度が所定の閾値よりも低いときだけ、少なくとも1つの時間領域コード化モードを適用し、さもなければコードは、少なくとも1つの周波数領域コード化モードを適用する請求項1記載の音声プロセッサ。

(7) コードが、少なくとも1つの時間領域コード化モードを、少なくとも1つの周波数領域コード化モードでコード化された所定数の連続的に処理されるフレームの直ぐ後の各フレームに適用する (1) 記載の音声プロセッサ。

(8) 少なくとも1つの周波数領域コード化モードが、周波数、位相、および振幅を含む1組のパラメータをもつ複数のシヌソイドで各フレームの短期間のスペクトルを表わし、位相は多項式表現および初期位相値でモデル化されていて、初期位相値が、(a) 前フレームが、少なくとも1つの周波数領域コード化モードでコード化されたときは、前フレームの推定された最終位相値であるか、または (b) 前フレームが、少なくとも1つの時間領域コード化モードでコード化されたときは、前フレームの短期間のスペクトルから求められる位相値である (1) 記載の音声プロセッサ。

(9) 各フレームにおけるシヌソイドの周波数が、フレームのピッチ周波数の整数倍である (8) 記載の音声プロセッサ。

(10) 各フレームにおけるシヌソイドの周波数が、0 ないし 2 の1組の実数から得られる (8) 記載の音声プロセッサ。

(11) フレームを処理する方法であって、

閉ループのコード化モード選択プロセスを各連続する入力フレームへ適用して、入力フレームの音声内容に基づいて、時間領域コード化モードか、または周波数領域コード化モードの何れかを選択するステップと、

入力フレームの音声内容が定常状態の有声音の音声を示すときは、入力フレームを周波数領域でコード化するステップと、

入力フレームの音声内容が定常状態の有声音の音声以外のものを示すときは、入力フレームを時間領域でコード化するステップと、

周波数領域でコード化されたフレームと入力フレームとを比較して、性能尺度を求めるステップと、

性能尺度が所定の閾値よりも低いときは、入力フレームを時間領域でコード化するステップとを含むフレームを処理する方法。

(12) フレームが、線形予測残余フレームである (11) 記載の方法。

(13) フレームが音声フレームである (11) 記載の方法。

(14) 時間領域でコード化するステップが、第1のコード化レートでフレームをコード化することを含み、周波数領域でコード化するステップが、第2のコード化レートで

フレームをコード化することを含み、第2のコード化レートが第1のコード化レートよりも低い(11)記載の方法。

(15) 周波数領域でコード化するステップが、高調波でコード化することを含む(11)記載の方法。

(16) 周波数領域でコード化するステップが、周波数、位相、および振幅を含む1組のパラメータをもつ複数のシヌソイドで各フレームの短期間のスペクトルを表わし、位相は多項式表現および初期位相値でモデル化されていて、初期位相値が、(a)前フレームが周波数領域でコード化されたときは、前フレームの推定された最終位相値であるか、または(b)前フレームが時間領域でコード化されたときは、前フレームの短期間のスペクトルから求められる位相値である(11)記載の方法。

(17) 各フレームのシヌソイド周波数が、フレームのピッチ周波数の整数倍である(16)記載の方法。

(18) 各フレームのシヌソイド周波数が、0ないし2の1組の実数から得られる(16)記載の方法。

(19) マルチモードの混合領域の音声プロセッサであって、

開ループのコード化モード選択プロセスを入力フレームへ適用して、入力フレームの音声内容に基づいて、時間領域コード化モードか、または周波数領域コード化モードの何れかを選択する手段と、

入力フレームの音声内容が定常状態の有声音の音声を示すときは、入力フレームを周波数領域でコード化する手段と、

入力フレームの音声内容が定常状態の有声音の音声以外のものを示すときは、入力フレームを時間領域でコード化する手段と、

周波数領域でコード化されたフレームと入力フレームとを比較して、性能尺度を求める手段と、

性能尺度が所定の閾値よりも低いときは、入力フレームを時間領域でコード化する手段とを含むマルチモードの混合領域の音声プロセッサ。

(20) フレームが線形予測残余フレームである(19)記載の音声プロセッサ。

(21) 入力フレームが音声フレームである(19)記載の音声プロセッサ。

(22) 時間領域でコード化する手段が、第1のコード化レートでフレームをコード化する手段を含み、周波数領域でコード化する手段が、第2のコード化レートでフレームをコード化する手段を含み、第2のコード化レートが第1のコード化レートよりも低い(19)記載の音声プロセッサ。

(23) 周波数領域でコード化する手段が、高調波コーダを含む(19)記載の音声プロセッサ。

(24) 周波数領域でコード化する手段が、周波数、位相、および振幅を含む1組のパラメータをもつ複数のシヌソイドで各フレームの短期間のスペクトルを表わす手段を含み、位相は多項式表現および初期位相値でモデル化されていて、初期位相値が、(a)直前のフレームが周波数領域でコード化されたときは、直前のフレームの推定された最終位相値であるか、または(b)直前のフレームが時間領域でコード化されたときは、直前のフレームの短期間のスペクトルから求められる位相値である(19)記載の音声プロセッサ。

(25) 各フレームのシヌソイド周波数が、フレームのピッチ周波数の整数倍である(24)記載の音声プロセッサ。

(26) 各フレームのシヌソイド周波数が、0ないし2の1組の実数から得られる(24)記載の音声プロセッサ。

フロントページの続き

(74)代理人 100075672
弁理士 峰 隆司

(74)代理人 100095441
弁理士 白根 俊郎

(74)代理人 100084618
弁理士 村松 貞男

(74)代理人 100103034
弁理士 野河 信久

(74)代理人 100119976
弁理士 幸長 保次郎

(74)代理人 100153051
弁理士 河野 直樹

(74)代理人 100140176
弁理士 砂川 克

(74)代理人 100101812
弁理士 勝村 紘

(74)代理人 100124394
弁理士 佐藤 立志

(74)代理人 100112807
弁理士 岡田 貴志

(74)代理人 100111073
弁理士 堀内 美保子

(74)代理人 100134290
弁理士 竹内 将訓

(74)代理人 100127144
弁理士 市原 卓三

(74)代理人 100141933
弁理士 山下 元

(72)発明者 アミタバ・ダス
アメリカ合衆国 カリフォルニア州 9 2 1 3 1 サン・ディエゴ、ラムスデル・コート 1 1 7
3 8

F ターム(参考) 5J064 AA02 BC11 BC14 BC16 BC22 BD02 BD03

【外国語明細書】
2011090311000001.pdf