

(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2013-160937

(P2013-160937A)

(43) 公開日 平成25年8月19日(2013.8.19)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 0 L 15/04 (2013.01)	G 1 0 L 15/04 3 0 0 C	5 D 0 1 5
G 1 0 L 25/78 (2013.01)	G 1 0 L 11/00 4 0 2 L	

審査請求 未請求 請求項の数 7 O L (全 11 頁)

(21) 出願番号	特願2012-23132 (P2012-23132)	(71) 出願人	000006013 三菱電機株式会社 東京都千代田区丸の内二丁目7番3号
(22) 出願日	平成24年2月6日(2012.2.6)	(74) 代理人	100123434 弁理士 田澤 英昭
		(74) 代理人	100101133 弁理士 濱田 初音
		(74) 代理人	100173934 弁理士 久米 輝代
		(74) 代理人	100156351 弁理士 河村 秀央
		(72) 発明者	太刀岡 勇気 東京都千代田区丸の内二丁目7番3号 三菱電機株式会社内
		Fターム(参考)	5D015 DD03

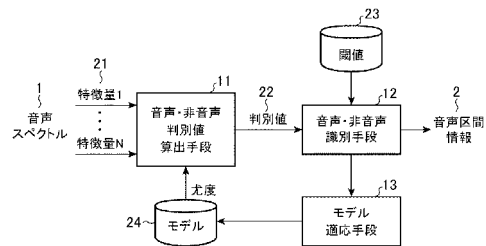
(54) 【発明の名称】 音声区間検出装置

(57) 【要約】

【課題】環境が変化した場合でも音声の検出を正しく行うことのできる音声区間検出装置を得る。

【解決手段】音声・非音声判別値算出手段11は、入力された音声スペクトル1から特徴量21を利用して、モデル24により音声・非音声区間それぞれにおいて異なる判別値22を算出する。音声・非音声識別手段12は、判別値22と予め設定した閾値23とを比較することで音声・非音声を識別する。モデル適応手段13は、音声・非音声識別手段12の識別結果に基づいてモデル24を適応する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

入力された音声のスペクトルからある特徴量もしくは複数の異なる特徴量の組み合わせを利用して、モデルにより音声・非音声区間それぞれにおいて異なる判別値を算出する音声・非音声判別値算出手段と、

前記判別値と予め設定した閾値とを比較することで音声・非音声を識別する音声・非音声識別手段と、

前記音声・非音声識別手段の識別結果に基づいて前記モデルを適応するモデル適応手段を備えたことを特徴とする音声区間検出装置。

【請求項 2】

モデルは音声モデルとノイズモデルであり、音声・非音声判別値算出手段は、これらの尤度の比または差を算出し、判別値とすることを特徴とする請求項 1 記載の音声区間検出装置。

【請求項 3】

モデルは、予め音声とノイズのデータから学習した密度比モデルであり、音声・非音声判別値算出手段は、前記密度比モデルに基づいて尤度比を算出し、判別値とすることを特徴とする請求項 1 記載の音声区間検出装置。

【請求項 4】

学習環境が異なる複数のモデルを用意すると共に、実際の環境に対応して前記複数のモデルから特定のモデルを選択するモデル選択手段を備えたことを特徴とする請求項 1 から請求項 3 のうちのいずれか 1 項記載の音声区間検出装置。

【請求項 5】

モデル適応手段は、音声・非音声識別手段で非音声または音声と判定されたフレームの特徴量を用いてモデルの適応を行うことを特徴とする請求項 4 記載の音声区間検出装置。

【請求項 6】

音声・非音声判別値算出手段および音声・非音声識別手段とは別の特徴量に基づいて音声と非音声とを識別する音声・非音声判別値算出および識別手段を備え、

モデル適応手段は、前記音声・非音声判別値算出および識別手段で非音声または音声と判定されたフレームの特徴量に基づいてモデルの適応を行うことを特徴とする請求項 1 から請求項 5 のうちのいずれか 1 項記載の音声区間検出装置。

【請求項 7】

モデル適応手段は、予め作成された音声モデルと適応したノイズモデルとを比較し、当該比較結果に基づいて適応の可否を判断することを特徴とする請求項 2 記載の音声区間検出装置。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、音声モデルやノイズモデルまたは尤度比モデルを用いて音声区間を検出する音声区間検出装置に関する。

【背景技術】**【0002】**

音声区間検出法としては音声のパワーがノイズのそれよりも大きなことを利用するパワーによるものがよく用いられている（例えば、非特許文献 1 参照）。また、音声区間、非音声区間の尤度比により音声区間検出を行う手法が提案されている（例えば、非特許文献 2 参照）。さらに、別途学習した音声モデルを用いたもの（例えば、特許文献 1 参照）が提案されている。

【先行技術文献】**【特許文献】****【0003】**

【特許文献 1】特開 2009 - 210647 号公報

10

20

30

40

50

【非特許文献】

【0004】

【非特許文献1】L.R. Rabiner and M.R. Sambur, "An algorithm for determining the endpoints of isolated utterances," Bell Syst. Tech., vol.54(2), pp.297-315, 1975.

【非特許文献2】J. Sohn, N.S. Kim, and W. Sung, "Statistical Model-Based Voice Activity Detection," IEEE Signal Processing Letters, vol.6(1), pp. 1-3, 1999.

【発明の概要】

【発明が解決しようとする課題】

【0005】

10

しかしながら、上記非特許文献1に記載されたような手法（以下、手法1とする）は音声のパワーがノイズに埋もれて検出できないため、低S/N下で有効でないという問題点があった。また、非特許文献2に記載されているようなノイズモデルを用いる手法（以下、手法2とする）は、ノイズらしさ音声らしさを判定できるため、手法1のような問題は起こりにくく、低S/N下でも音声を検出できるが、ノイズが非定常な場合には学習したモデルと実際の環境とのミスマッチにより、有効でないという問題があった。さらに、特許文献1に記載されたような手法（手法3とする）はオンラインでモデルを推定するため手法2のような問題は起こりにくいが、毎フレームモデルを合成するため、計算量が多いという問題があった。

【0006】

20

この発明は上記のような課題を解決するためになされたもので、環境が変化した場合でも音声の検出を正しく行うことのできる音声区間検出装置を得ることを目的とする。

【課題を解決するための手段】

【0007】

この発明に係る音声区間検出装置は、入力された音声のスペクトルからある特徴量もしくは複数の異なる特徴量の組み合わせを利用して、モデルにより音声・非音声区間それぞれにおいて異なる判別値を算出する音声・非音声判別値算出手段と、判別値と予め設定した閾値とを比較することで音声・非音声を識別する音声・非音声識別手段と、音声・非音声識別手段の識別結果に基づいてモデルを適応するモデル適応手段とを備えたものである。

30

【発明の効果】

【0008】

この発明の音声区間検出装置は、音声・非音声識別手段の識別結果に基づいてモデルを適応するモデル適応手段を備えたので、環境が変化した場合でも音声の検出を正しく行うことができる。

【図面の簡単な説明】

【0009】

【図1】この発明の実施の形態1による音声区間検出装置を示す構成図である。

【図2】この発明の実施の形態1による音声区間検出装置の具体例を示す構成図である。

【図3】この発明の実施の形態2による音声区間検出装置の構成図である。

40

【図4】この発明の実施の形態3による音声区間検出装置の構成図である。

【図5】この発明の実施の形態4による音声区間検出装置の構成図である。

【図6】この発明の実施の形態5による音声区間検出装置に係わるモデル更新の検証を示す説明図である。

【図7】この発明の実施の形態5による音声区間検出装置の構成図である。

【発明を実施するための形態】

【0010】

本発明では手法2の性能を改善することを目的とする。具体的には、学習したモデルと実際の環境とのミスマッチという問題に対応するために、モデルの適応を行う方法を提案する。また、モデルの頑健性を増し、頻繁にモデル適応をしなくてもよいように、尤度比

50

モデルの使用により頑健性を高める方法に関しても述べる。

【 0 0 1 1 】

実施の形態 1 .

図 1 は、この発明の実施の形態 1 による音声区間検出装置の構成図である。

図 1 に示す音声区間検出装置は、音声・非音声判別値算出手段 1 1、音声・非音声識別手段 1 2、モデル適応手段 1 3 を備えている。音声・非音声判別値算出手段 1 1 は、音声スペクトル 1 を入力し、その特徴量 (1 ~ N) 2 1 とモデル 2 4 との尤度の比または差に基づいて、判別値 2 2 を出力する手段である。音声・非音声識別手段 1 2 は、音声・非音声判別値算出手段 1 1 から送出された判別値 2 2 と、予め求められた閾値 2 3 とを比較し、音声と非音声との識別結果をモデル適応手段 1 3 に出力すると共に、音声と判定した場合は、その音声区間情報 2 を出力する手段である。モデル適応手段 1 3 は、音声・非音声識別手段 1 2 の音声・非音声の識別結果に基づいて、モデル 2 4 を適応する手段である。

10

【 0 0 1 2 】

このように構成された音声区間検出装置は、入力の音声スペクトル 1 から音声区間情報 2 を出力する。まず、一つもしくは複数の特徴量 2 1 を音声・非音声判別値算出手段 1 1 に入力する。特徴量としては、帯域別のスペクトル、パワー、MFCC (メル周波数ケプストラム係数)、ゼロ交差点数などを用いることができる。音声・非音声判別値算出手段 1 1 は、モデル 2 4 により尤度を得て、これらを統合して判別値 2 2 を出力する。音声・非音声識別手段 1 2 では、受け取った判別値 2 2 を、事前に設定しておいた閾値 2 3 と比較し、音声区間であれば音声区間情報 2 として出力する。また、その結果をモデル適応手段 1 3 にフィードバックし、モデル適応手段 1 3 はモデル 2 4 を適応する。

20

【 0 0 1 3 】

この処理を具体的に説明するために、図 1 のモデル 2 4 として音声モデル 2 4 a とノイズモデル 2 4 b を用いた構成を図 2 に示す。ここでは、特徴量の例として複素スペクトルを用いた手法 2 に即してモデル適応の方法を説明する。非音声区間 H_0 と音声区間 H_1 において、音声とノイズの離散フーリエ変換 (D F T) の係数の L 次元のベクトル (複素スペクトル) をそれぞれ S , N として観測音 X は下式のように表される。それぞれは $S = (S_0, S_1, \dots, S_k, \dots, S_{L-1})$, $N = (N_0, N_1, \dots, N_k, \dots, N_{L-1})$, $X = (X_0, X_1, \dots, X_k, \dots, X_{L-1})$ である。

$$H_0 : X = N, H_1 : X = N + S$$

30

【 0 0 1 4 】

手法 2 ではここでそれぞれの状態の D F T 係数の確率密度関数が、下式のように独立なガウス分布で表せると仮定する。

$$p(X | H_0) = \prod_{k=0}^{L-1} \frac{1}{\pi \lambda_N(k)} \exp \left\{ -\frac{|X_k|^2}{\lambda_N(k)} \right\}$$

$$p(X | H_1) = \prod_{k=0}^{L-1} \frac{1}{\pi \{ \lambda_N(k) + \lambda_S(k) \}} \exp \left\{ -\frac{|X_k|^2}{\lambda_N(k) + \lambda_S(k)} \right\}$$

40

ここで $\lambda_N(k)$, $\lambda_S(k)$ は N_k , S_k の分散を表す。

すると k 番目の周波数帯域の音声・非音声の尤度比は下式で表される。

$$\Lambda_k = \frac{p(X|H_1)}{p(X|H_0)} = \frac{1}{1+\xi_k} \exp\left(\frac{\gamma_k \xi_k}{1+\xi_k}\right)$$

【0015】

モデルとしてノイズパワーの分散 $\sigma_N(k)$ と音声パワーの分散 $\sigma_S(k)$ を予め学習しておく。このモデルを用いて、音声・非音声判別値算出手段 11 で上式（確率の比（音声の確率 / ノイズの確率））を計算する。これが判別値（確率密度比）22 となる。音声・非音声識別手段 12 において、判別値 22 が予め設定した閾値 23 より大きかった場合には音声であると判断し、小さかった場合にはノイズであると判断する。

10

【0016】

この方法によると学習時の $\{\sigma_S, \sigma_N\}(k)$ と実際の環境でのそれが異なった場合には推定精度が大きく低下する。この場合にモデル適応手段 13 では、判別値 22 が閾値 23 より非常に大きかった場合には $\sigma_S(k)$ の学習データに加え、再学習する。一方、判別値 22 が閾値 23 より非常に小さかった場合にはノイズのそれを再学習する。これにより $\{\sigma_S, \sigma_N\}(k)$ が徐々に変化した場合には追従が可能となり、 $\{\sigma_S, \sigma_N\}(k)$ が大きく変化した場合にも時間の経過と共に徐々に追従することができる。このようにモデル適応手段 13 を備えたことにより、音声と騒音のモデルが常に最新のものに更新されるので、環境が変化した場合にも追従することができるという効果が得られる。

20

【0017】

以上説明したように、実施の形態 1 の音声区間検出装置によれば、入力された音声のスペクトルからある特徴量もしくは複数の異なる特徴量の組み合わせを利用して、モデルにより音声・非音声区間それぞれにおいて異なる判別値を算出する音声・非音声判別値算出手段と、判別値と予め設定した閾値とを比較することで音声・非音声を識別する音声・非音声識別手段と、音声・非音声識別手段の識別結果に基づいてモデルを適応するモデル適応手段とを備えたので、環境が変化した場合でも音声の検出を正しく行うことができる。

【0018】

また、実施の形態 1 の音声区間検出装置によれば、モデルは音声モデルとノイズモデルであり、音声・非音声判別値算出手段は、これらの尤度の比または差を算出し、それを判別値とするようにしたので、音声と騒音のモデルが常に最新のものに更新されるので、環境が変化した場合にも追従することができるという効果が得られる。

30

【0019】

実施の形態 2 .

実施の形態 1 の図 2 に示した構成では、音声とノイズの二つのモデルを用いていた。ところで、尤度比を求めるためには、ある特徴量の音声区間の確率密度関数 $p(X|H_1)$ およびノイズ区間の確率密度関数 $p(X|H_0)$ がそれぞれ求まる必要はなく、それらの比の確率密度関数が直接推定できればよい。このように確率密度の比を直接推定する枠組みが密度比推定であり、これについては、例えば、杉山将，“密度比に基づく機械学習の新たなアプローチ，”統計数理 vol. 58, no. 2, pp. 141 - 155, 2010. とした文献（以下、この文献を参考文献 1 とする）に示されている手法を用いることができる。この文献は、二つの異なる分布の出力値の尤度比を算出する一般的な枠組みであるが、ここではそれを音声区間検出に応用することで新たな効果が生まれることを説明する。

40

【0020】

図 3 は、実施の形態 2 の音声区間検出装置を示す構成図である。

実施の形態 2 の音声区間検出装置では、音声・非音声判別値算出手段 11a が用いるモデルとして密度比モデル 25 となっている点が実施の形態 1 と異なる部分であり、音声・非音声判別値算出手段 11a は密度比モデル 25 に基づいて尤度比を算出し、これを判別値 22 として出力するよう構成されている。他の構成は図 1 に示した実施の形態 1 と同様

50

であるため対応する部分に同一符号を付してその説明を省略する。

【 0 0 2 1 】

図示のように、実施の形態 1 では二つあったモデルが一つに減らせている。密度比モデル 2 5 は事前の学習データにより学習しておく。この際に学習環境とのミスマッチを低減するために、密度比モデル 2 5 を適応することで、推定精度を向上させることができる。

【 0 0 2 2 】

次に、密度比の推定法に関して説明する。密度比の推定にはいくつかの手法が提案されているが、ここでは最も基本的なカルバックライブラー重要度推定法 (Kullback - Leibler importance estimation procedure ; K L I E P) を例に述べる。

10

K L I E P では、2 組のデータ列 $\{x_i\}_{i=1}^n$ $\{x'_j\}_{j=1}^{n'}$ の確率密度関数 q , q' の比 $r(x) = \frac{q'(x)}{q(x)}$ を下式 (1) の線形モデルでモデル化する。

$$\hat{r}(x) = \sum_{l=1}^b \alpha_l \phi_l(x) \quad (1)$$

20

【 0 0 2 3 】

α_l は非負のパラメータであり、 ϕ_l は非負の値をとる基底関数である。 x が音声に対応する特徴量のデータ列であり、 x' がノイズに対応する特徴量のデータ列であると考え、参考文献 1 の方法に従って基底関数 ϕ_l を事前に求める。参考文献 1 では下式のガウシアンカーネル K を用いその幅 σ の検定にはクロスバリデーションを使っている。

$$K_{\sigma}(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

30

上式において、 $\|\cdot\|$ はユークリッドノルムである。

【 0 0 2 4 】

ここで特徴量としてはパワー、FFT 後の帯域別パワー (log パワー)、MFCC、ゼロ交差数など音声とノイズを分けるのに有効な特徴量であればいずれを用いてもよい。またそれらを組み合わせて用いることもできる。ここでは log パワーを用いた場合を説明する。ノイズに比べて音声のパワーは狭い帯域に偏る傾向があるため、この特徴量を用いることでトータルのエネルギーとしての S / N を見るよりも判別性能を向上させることができる。

40

log パワーの特徴量のベクトルを $LP = (LP_0, LP_1, \dots, LP_{L-1})$ とする。

学習の際は音声・非音声でラベル付されたデータ $\{x_i\}_{i=1}^n$ $\{x'_j\}_{j=1}^{n'}$ より α , ϕ を学習する。

【 0 0 2 5 】

この方法で学習された密度比モデルは実際の環境で動作させる際にはミスマッチがあることが多いので、モデル適応手段13でモデルを適応する必要がある。最も単純には学習

時のノイズに対応する特徴量の平均 $\frac{1}{N} \sum LP$ と、始めの数フレームから算出された実際の環境の平均を取りその差を引き去ることで適応できる。また μ としてガウシアンカーネルを用いた場合は、分散 σ を正規化することでも適応できる。

【0026】

10

実施の形態1ではどちらかのモデルの推定に誤りが生じた場合に推定精度が大きく低下する。それに対して実施の形態2では、密度比モデル25の一つだけから尤度比を直接計算できるので、判別値算出が頑健に行えるようになる。そのため、モデルの学習環境と実環境がマッチしている場合には従来法よりも性能が良くなる効果が得られる。加えてモデル適応手段13を備えたことにより、密度比モデル25を適応することができるようになる。上述の通り、マッチしていれば従来法よりも性能が良くなる効果が得られるので、モデル適応によりモデルの学習環境と実環境を近づけることで、従来法よりも性能が良くなる効果が得られる。

【0027】

20

以上説明したように、実施の形態2の音声区間検出装置によれば、モデルは、予め音声とノイズのデータから学習した密度比モデルであり、音声・非音声判別値算出手段は、密度比モデルに基づいて尤度比を算出し、判別値とするようにしたので、判別値算出が頑健に行えるようになり、頻りにモデル適応をしなくてもよいといった効果を得ることができる。

【0028】

実施の形態3

音声とノイズの組み合わせによって、尤度比モデルや音声・ノイズモデルが異なるものになるのは明らかであるから、これらを複数用意し、学習環境が実際の環境に最も近いものを選び出すようにしてもよく、このような例を実施の形態3として以下に説明する。

【0029】

30

図4は、実施の形態3の音声区間検出装置を示す構成図である。

図示の音声区間検出装置は、それぞれ学習環境が異なる複数の密度比モデル25-1~25-nを選択するモデル選択手段14を備えている。それ以外の構成は図3に示した実施の形態2と同様であるため、対応する部分に同一符号を付してその説明を省略する。なお、実施の形態1の構成に追加するようにしてもよく、この場合は、図示のように、モデル選択手段14は音声モデル24a-1とノイズモデル24b-1~音声モデル24a-nとノイズモデル24b-nのいずれかを選択する。

【0030】

モデル選択手段14は、モデルの学習環境が実際の環境に最も近いものを選び出す手段であるが、学習環境が実際の環境に近いかどうかを測る基準には、例えば背景ノイズの類似性、話者の類似性を用いることができる。まず、学習時のノイズと音声のモデルを作成しておく(いずれか一方でもよい)。例えば、下式のGMM(Gaussian Mixture Model)により、EMアルゴリズムなどを用いてモデル化することができる。

40

$$p(x) = \sum_{k=1}^K \pi_k N(\mu_k, \Sigma_k)$$

50

上式において、 N は平均 μ_k , 共分散 σ_k^2 , 混合率 π_k の正規分布である。

【0031】

モデル選択手段14は、最も尤度 $p(x)$ の高いモデルに対応する密度比モデル25-1~25-nを選択することで、学習時と実際の環境のミスマッチの少ないモデルを選び出すことができる。また、ノイズモデル24b-1~24b-nはノイズが観測されれば一通りに選べるが、音声モデル24a-1~24a-nは音声が入力するまで選べないので、あるノイズに対する複数の音声モデル24a-1~24a-nを同時に用いて音声・非音声識別手段12で音声区間検出を行い、検出された部分の特徴量から音声モデル24a-1~24a-nを選び出すこともできる。

【0032】

以上説明したように、実施の形態3の音声区間検出装置によれば、学習環境が異なる複数のモデルを用意すると共に、実際の環境に対応して複数のモデルから特定のモデルを選択するモデル選択手段を備えたので、環境にマッチしたモデルを選択することができ、モデル適応の効果がより出やすいという効果が得られる。

【0033】

また、実施の形態3の音声区間検出装置によれば、モデル適応手段は、音声・非音声識別手段で非音声または音声と判定されたフレームの特徴量を用いてモデルの適応を行うようにしたので、種々の環境でモデルの適応を行うことができる。

【0034】

実施の形態4

音声・非音声判別値算出手段11aと音声・非音声識別手段12による音声・非音声の判別基準と、モデル適応手段13における判別基準が同じであると、音声・非音声判別値算出手段11aや音声・非音声識別手段12での誤りがモデル適応手段13の適応によって訂正される効果を期待できなくなる。そこで、別のM個の特徴量により音声区間検出を行い、その結果に従ってモデルの適応を行うことで音声・非音声判別値算出手段11aと音声・非音声識別手段12の誤りを正す効果が期待できる。別のM個の特徴量としては、例えば音声・非音声判別値算出手段11aで帯域別のパワーを用いた場合には、これとは異なるゼロ交差数を用いればよい。このような例を実施の形態4として図5に示す。

【0035】

図5に示す音声区間検出装置では、実施の形態2の構成に加えて、音声・非音声判別値算出および識別手段15を備えたものである。音声・非音声判別値算出および識別手段15は、音声・非音声判別値算出手段11aと音声・非音声識別手段12の機能を備えたものであり、これら音声・非音声判別値算出手段11aおよび音声・非音声識別手段12とは別の特徴量 ($N+1 \sim N+M$) 21aに基づいて音声と非音声とを識別する手段である。また、モデル適応手段13aは、音声・非音声判別値算出および識別手段15から出力されるモデル適応可否フラグが真である場合にのみモデル適応を行うよう構成されている。なお、図5では実施の形態2に対して音声・非音声判別値算出および識別手段15を追加した例を示しているが、実施の形態1に対して追加するよう構成してもよい。

【0036】

このように構成された音声区間検出装置では、音声・非音声判別値算出および識別手段15において、モデルを更新することができると判断された場合(確実に非音声または音声と判定された場合)のみモデル適応可否フラグを真とする。モデル適応手段13aはモデル適応可否フラグが真であった場合にのみモデル適応を行う。

【0037】

以上説明したように、実施の形態4の音声区間検出装置によれば、音声・非音声判別値算出手段および音声・非音声識別手段とは別の特徴量に基づいて音声と非音声とを識別する音声・非音声判別値算出および識別手段を備え、モデル適応手段は、音声・非音声判別値算出および識別手段で非音声または音声と判定されたフレームの特徴量に基づいてモデルの適応を行うようにしたので、判別誤りを訂正することができ、モデル適応がより効果的に行われるという効果が得られる。

10

20

30

40

50

【0038】

実施の形態5 .

実施の形態1においては音声モデルがあるのでその情報を積極的に使うことで、モデル適応の可否を判断することができる。図6のように信号が観測されたとする。図6の上図は信号波形、下図は適当な手段を用いて算出された尤度比Rである。閾値を超えたものを音声区間とみなすことにする。実際の音声は上図Voiceの部分であるので、音声区間の終端がノイズと誤って識別されている。このままモデル適応を行うとノイズモデルの適応データに音声データが混ざり、検出性能が著しく低下する。そこで、音声区間中の音声モデル(Speech model)とノイズモデル(Noise model)を比較することによって類似している場合には、モデル適応を行わないことにする。

10

【0039】

$$\chi = \sum_k \frac{\lambda_N + \lambda_S}{\lambda_N}$$

例えば、実施の形態1では χ をモデル間の類似性を測る尺度として用いることができる。

図6の例ではNoise model 1においては音声の学習データが混入しているため λ_N が大きくなって χ が小さくなり、Noise model 1とSpeech model 1が類似していると判断される。このため適応が不可となる。Noise model 2は正しく適応されているため、類似性が小さく適応が可となる。

20

【0040】

もしくは予め音声と分かっているテストデータを用意しておき、それに対する尤度 $p(x | H_0)$ と $p(x | H_1)$ を測り、それぞれの尤度差(比)による方法も考えられる。この方法を実施する音声区間検出装置を図7に示す。図7において、音声モデル26a、ノイズモデル26bは、予め音声と分かっているテストデータに基づいて適応したモデルである。モデル識別手段16は、これら音声モデル26aとノイズモデル26bとが類似しているか否かを識別し、類似していないと判断した場合にモデル適応可否フラグを真とする手段である。また、モデル更新手段17a, 17bは、モデル適応可否フラグが真であった場合にそれぞれ音声モデル24aとノイズモデル24bとを更新する手段である。なお、モデル適応手段13、モデル識別手段16およびモデル更新手段17a, 17bは、予め作成された音声モデルと適応したノイズモデルとを比較し、その比較結果に基づいて適応の可否を判断するモデル適応手段を構成している。これ以外の構成は図2に示した実施の形態1の構成と同様であるため、対応する部分に同一符号を付してその説明を省略する。

30

【0041】

このように、ノイズモデル26bと音声モデル26aの双方を識別するモデル識別手段16により、それらのモデルが類似していないと判断された場合には音声モデル24aとノイズモデル24bを更新する。従って、ノイズモデルの適応データに音声データが混入し検出性能が著しく低下するのを避けることができる効果が得られる。

40

【0042】

以上説明したように、実施の形態5の音声区間検出装置によれば、モデル適応手段は、予め作成された音声モデルと適応したノイズモデルとを比較し、比較結果に基づいて適応の可否を判断するようにしたので、ノイズモデルの適応データに音声データが混入することで検出性能が著しく低下するといったことを避けることができる。

【0043】

なお、本願発明はその発明の範囲内において、各実施の形態の自由な組み合わせ、あるいは各実施の形態の任意の構成要素の変形、もしくは各実施の形態において任意の構成要

50

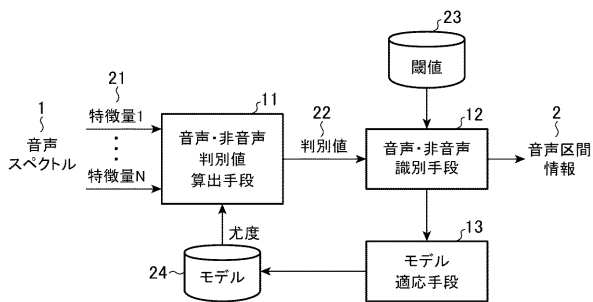
素の省略が可能である。

【符号の説明】

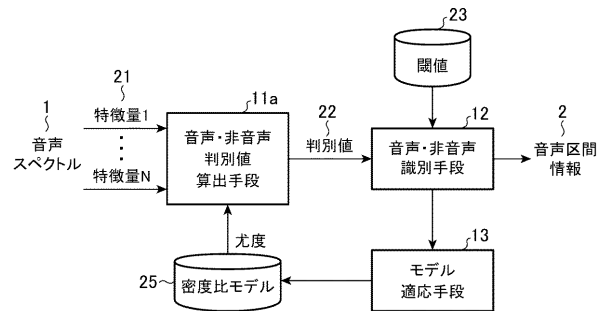
【0044】

- 1 音声スペクトル、2 音声区間情報、11, 11a 音声・非音声判別値算出手段、12 音声・非音声識別手段、13, 13a モデル適応手段、14 モデル選択手段、15 音声・非音声判別値算出および識別手段、16 モデル識別手段、17a, 17b モデル更新手段、21, 21a 特徴量、22 判別値、23 閾値、24 モデル、24a, 26a 音声モデル、24b, 26b ノイズモデル、25, 25-1~25-n 密度比モデル。

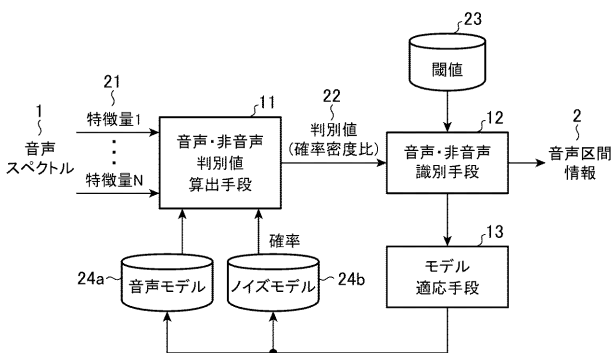
【図1】



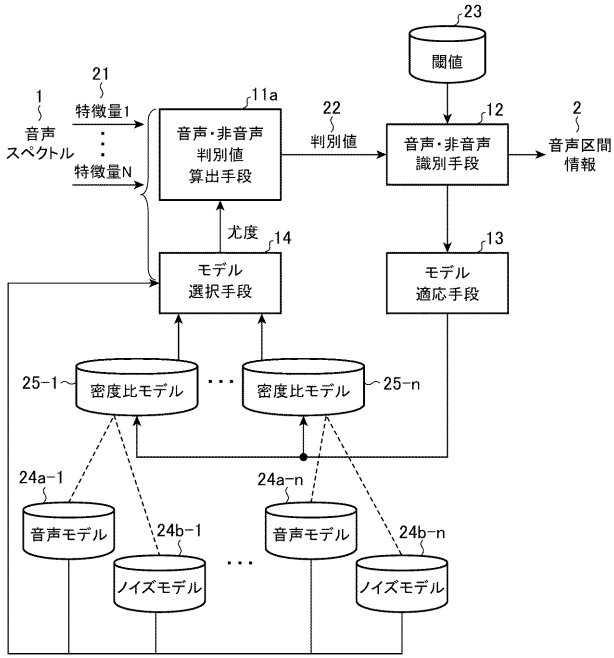
【図3】



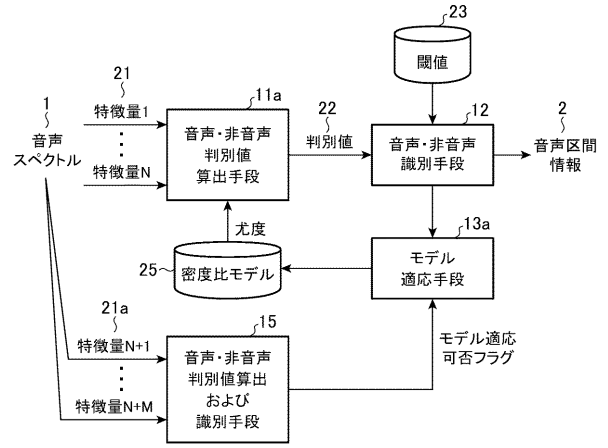
【図2】



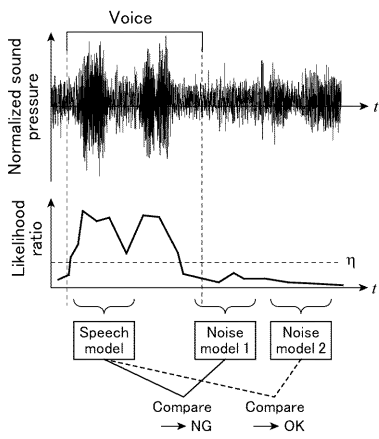
【 図 4 】



【 図 5 】



【 図 6 】



【 図 7 】

