

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4521673号
(P4521673)

(45) 発行日 平成22年8月11日(2010.8.11)

(24) 登録日 平成22年6月4日(2010.6.4)

(51) Int.Cl. F I
G 1 O L 15/04 (2006.01) G 1 O L 15/04 3 0 0 A
G 1 O L 11/02 (2006.01) G 1 O L 11/02

請求項の数 13 (全 23 頁)

(21) 出願番号	特願2004-101094 (P2004-101094)	(73) 特許権者	393031586 株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2
(22) 出願日	平成16年3月30日(2004.3.30)	(74) 代理人	100099933 弁理士 清水 敏
(65) 公開番号	特開2005-31632 (P2005-31632A)	(72) 発明者	スーン フランク ガービン 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
(43) 公開日	平成17年2月3日(2005.2.3)	(72) 発明者	中村 哲 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
審査請求日	平成19年3月28日(2007.3.28)	(72) 発明者	葦苺 豊 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
(31) 優先権主張番号	特願2003-174416 (P2003-174416)		
(32) 優先日	平成15年6月19日(2003.6.19)		
(33) 優先権主張国	日本国(JP)		

最終頁に続く

(54) 【発明の名称】 発話区間検出装置、コンピュータプログラム及びコンピュータ

(57) 【特許請求の範囲】

【請求項1】

音声データを逐次フレーム化するためのフレーム化手段と、

前記フレーム化手段によりフレーム化された音声のエネルギー値をフレームごとに算出し、FIFO(First-In First-Out)形式で第1の個数のフレームのエネルギー値を記憶するためのフレームエネルギー算出及び記憶手段と、

前記フレームエネルギー算出及び記憶手段に、第2の個数のフレームのエネルギー値が格納されたことに応答して、前記第2の個数のフレームのエネルギー値を所定の統計的手法に従って処理することにより、前記音声データに含まれる環境雑音の推定値の初期値を算出するための初期値算出手段と、

前記推定値の初期値と、フレームエネルギー算出及び記憶手段に逐次記憶される音声のエネルギー値とに基づいて、前記音声データに含まれる環境雑音の変化に追従して変化するように、前記発話区間を検出するためのエネルギー値のしきい値をフレームごとに逐次算出するための手段と、

前記しきい値に基づいて、前記第2の個数のフレーム以降のフレームの中で、前記音声データの発話区間の開始位置又は終了位置に対応するフレームを推定するための発話区間推定手段とを含む、発話区間検出装置であって、

前記初期値算出手段は、

前記第2の個数のフレームを、各フレームのエネルギー値の大きさによって、第1のエネルギー値を中心とする第1のクラスと、前記第1のエネルギーよりも大きな第2のエネルギー

値を中心とする第2のクラスタとにクラスタ化するための手段と、

前記第1のエネルギー値を前記環境雑音の推定値の初期値として出力するための手段とを含む、発話区間検出装置。

【請求項2】

前記クラスタ化するための手段は、

前記第2の個数のフレームを前記第1及び第2のクラスタにクラスタ化するための境界値を決定するための手段と、

前記境界値よりも小さなエネルギー値を持つフレームを前記第1のクラスタに、それ以外のフレームを前記第2のクラスタに、それぞれ分類するための手段とを含む、請求項1に記載の発話区間検出装置。

10

【請求項3】

前記境界値を決定するための手段は、

前記第2の個数のフレームのうち、エネルギー値をキーとしてソートしたときに予め定める第1のソート順位及び第2のソート順位となる二つのフレームを選択するための手段と、

、

前記選択された二つのフレームのエネルギー値の平均値を算出するための第1の平均値算出手段と、

前記第1の平均値算出手段により算出された平均値より小さいエネルギー値を持つか否かを基準として、前記第2の個数のフレームを第1及び第2のグループに分類するための手段と、

20

前記第1及び第2のグループに属するフレームのエネルギー値の平均値をそれぞれ算出するための第2の平均値算出手段と、

前記第2の平均値算出手段により算出された二つの平均値の平均値をさらに算出し、前記境界値として出力するための第3の平均値算出手段とを含む、請求項2に記載の発話区間検出装置。

【請求項4】

前記しきい値をフレームごとに逐次算出するための手段は、

前記フレームエネルギー算出及び記憶手段に格納されているフレームのエネルギー値と、前記環境雑音の推定値の初期値とに基づいて、前記フレームエネルギー算出及び記憶手段に格納されているフレームの環境雑音のエネルギー値をフレームごとに推定するための手段と、

30

前記フレームエネルギー算出及び記憶手段に格納されているフレームのエネルギー値のうち、定常的な背景雑音及び発話音声の合計のエネルギー値の最大値をフレームごとに逐次推定するための手段と、

前記推定された環境雑音のエネルギー値と、前記推定された背景雑音及び発話音声の合計のエネルギー値とに基づいて、前記発話区間を検出するためのエネルギーのしきい値をフレームごとに算出するための手段とを含む、請求項1に記載の発話区間検出装置。

【請求項5】

前記発話区間推定手段は、前記しきい値に基づいて、前記第2の個数のフレーム以降のフレームの状態を判定するための手段を含み、

前記状態は、非発話状態を含み、

40

前記環境雑音のエネルギー値をフレームごとに逐次推定するための手段は、

1フレーム前の時点において推定された前記環境雑音のエネルギー値を記憶するための手段と、

前記環境雑音の推定値の初期値が算出された時点で前記記憶するための手段に前記環境雑音の推定値の初期値を記憶させるための手段と、

前記記憶するための手段に記憶された値、前記フレームエネルギー算出及び記憶手段に含まれるフレームのエネルギー値、及び前記フレームの状態を判定する手段による判定結果に基づいて、以下の式

$$b(t) = b(t-1) \times \quad + E(t) \times (1 - \quad) \quad (\text{状態が非発話状態の場合})$$

$$b(t) = b(t-1) \quad (\text{状態が非発話状態以外の場合})$$

50

ただし α は所定の忘却係数、 $E(t)$ は時刻 t におけるフレームのエネルギー値、に從つて時刻 t における背景雑音 $b(t)$ を算出するための手段とを含み、

前記記憶するための手段は、算出された前記背景雑音 $b(t)$ を記憶する、請求項 4 に記載の発話区間検出装置。

【請求項 6】

前記合計のエネルギー値の最大値をフレームごとに推定するための手段は、

前記フレームエネルギー算出及び記憶手段に格納されているフレームを、エネルギー値をキーとしてソートするための手段と、

前記ソートするための手段によりソートされた結果所定の順位となるフレームのエネルギー値を前記合計のエネルギー値の最大値 $E_{\max}(t)$ として選択するための手段を含む、請求項 5 に記載の発話区間検出装置。

10

【請求項 7】

前記しきい値をフレームごとに逐次算出するための手段は、

時刻 t における発話開始位置検出のためのしきい値 $E_{th1}(t)$ を、

$$E_{th1}(t) = b(t) + \max(\alpha, E_{\max}(t) - b(t)) \times \text{第 1 の定数}$$

ただし α は音声データ信号の最低ダイナミックレンジとして予め定められた定数、に從つて算出するための手段を含む、請求項 6 に記載の発話区間検出装置。

【請求項 8】

前記しきい値をフレームごとに逐次算出するための手段は、さらに、

時刻 t における発話終了位置検出のためのしきい値 $E_{th2}(t)$ を、

$$E_{th2}(t) = b(t) + \max(\alpha, E_{\max}(t) - b(t)) \times \text{第 2 の定数}$$

ただし第 2 の定数 < 第 1 の定数、

α は音声データ信号の最低ダイナミックレンジとして予め定められた前記定数、に從つて算出するための手段を含む、請求項 7 に記載の発話区間検出装置。

20

【請求項 9】

さらに、発話の先頭からの各フレームの音声データの最大エネルギー値又は所定のデフォルト基準値のいずれか大きい方を用いて各フレームの音声データを正規化し、各フレームの音声特徴パラメータとして出力するための音声エネルギー正規化手段を含む、請求項 1 ~ 請求項 8 のいずれかに記載の発話区間検出装置。

【請求項 10】

前記音声エネルギー正規化手段は、

正規化の基準値を記憶するための基準値記憶手段と、

前記フレームエネルギー算出及び記憶手段により算出された音声エネルギーが、前記基準値記憶手段に記憶された基準値を超えていることを検出し、検出信号を出力するための検出手段と、

前記検出手段により出力される前記検出信号にตอบสนองして、前記基準値記憶手段に記憶された基準値を、前記フレームエネルギー算出及び記憶手段により算出された値で置換するための手段と、

前記フレームエネルギー算出及び記憶手段により算出された音声エネルギー値を、前記基準値記憶手段に記憶された基準値で除算することにより、当該フレームの音声エネルギーを正規化するための除算手段とを含む、請求項 9 に記載の発話区間検出装置。

30

40

【請求項 11】

前記発話区間推定手段により、発話区間の終了位置に対応するフレームが推定されたことにตอบสนองして、前記基準値記憶手段の記憶内容を、所定のデフォルト値で置換するための手段をさらに含む、請求項 10 に記載の発話区間検出装置。

【請求項 12】

前記所定のデフォルト値を、前記発話区間検出装置の起動時に与えられたオプション値に基づいて設定するための手段をさらに含む、請求項 10 又は請求項 11 に記載の発話区間検出装置。

【請求項 13】

50

コンピュータにより実行されると、当該コンピュータを請求項 1 から請求項 1_2 のいずれかに記載の発話区間検出装置として動作させる、発話区間検出のためのコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は音声認識処理などの前処理として発話区間を検出するための装置に関し、特に、実時間での音声認識処理において、環境雑音による発話区間の誤検出を避けることができる発話区間検出装置、ならびにフレームごとの特徴量として正規化した音声エネルギーを算出するための音声エネルギー正規化装置に関する。

10

【背景技術】

【0002】

音声認識などの処理においては、音声認識に先立って入力信号中の発話区間とそれ以外の区間（無音区間と呼ぶ。）との区別をすることが必要である。さもなければ、発話のない区間を音声認識することにより無意味な結果がもたらされるためである。

【0003】

従来、このような発話区間（又は無音区間）の検出は、入力される音声信号のパワー（エネルギー）を算出し、その値が予め定められたしきい値以上になれば発話区間、しきい値未満であれば無音区間とする、という手法により行なわれている。このとき、そうした条件の成立が持続した時間をも考慮に入れて発話区間又は無音区間の検出がされるのが通常である。

20

【0004】

そのような技術が特許文献 1 に開示されている。特許文献 1 は、音声付の映像情報から要約を自動的に作成するために、要約の対象となる個所を抽出するための技術を開示している。音声付の映像では、その内容（ジャンル）により、環境雑音の大きさが異なることが知られている。例えばニュース番組では環境雑音が小さく、スポーツ中継等の番組では環境雑音が大きいの、などである。そのため、同じしきい値を用いて発話区間を検出しようとすると、映像情報のジャンルによりその結果が異なってしまうという問題がある。そのために特許文献 1 に開示の技術では、映像情報に、そのジャンルを示す付帯情報をもたせておき、付帯情報に従って各ジャンルに予め割当てられたしきい値を選択している。

30

【0005】

【特許文献 1】特開 2003 - 101939（段落 209、210、図 1 及び図 7）

【発明の開示】

【発明が解決しようとする課題】

【0006】

しかし、上記した特許文献 1 に記載の技術では、一つの映像情報には一種類のしきい値しか使用できない。そのため、番組の中で環境雑音に変化した場合には、発話区間の検出に問題が生じるという問題がある。

【0007】

特に、実時間の音声認識を行なう場合には、上記したような付帯情報が利用可能となるとは考えられない。また、電話による自動応答などに音声認識を用いる場合、音声信号の背景に存在する環境雑音がどのようなものになるかは予想できない。たとえば突発的な環境雑音が生じた場合、発話区間の検出を誤る可能性が高い。

40

【0008】

また、音声認識においては発話中の音声エネルギーの最大値で各フレームの音声エネルギーを正規化した特徴量を用いると有効であることが知られている。しかしそのためには、発話の終了まで待って発話中での最大パワーを算出した後、算出された最大パワーを用いて当該発話中の各フレームの音声エネルギーを正規化する必要がある。しかし、発話の終了まで待っていると実時間の音声認識を行なうことができないという問題がある。

【0009】

50

従って、本発明の目的は、環境雑音にかかわらず発話区間の検出を適切に行なうことができる発話区間検出装置を提供することである。

【0010】

本発明の他の目的は、環境雑音が変化しても発話区間の検出を適切に行なうことができる発話区間検出装置を提供することである。

【0011】

本発明のさらに他の目的は、環境雑音が変化しても発話区間の検出を実時間で適切に行なうことができる発話区間検出装置を提供することである。

【0012】

本発明のさらに他の目的は、突発的な環境雑音の変化があっても発話区間の検出を実時間で適切に行なうことができる発話区間検出装置を提供することである。

10

【0013】

本発明の他の目的は、実時間で各フレームの音声エネルギーを正規化することができる音声エネルギー正規化装置を提供することである。

【課題を解決するための手段】

【0014】

本発明の第1の局面に係る発話区間検出装置は、音声データを逐次フレーム化するためのフレーム化手段と、フレーム化手段によりフレーム化された音声のエネルギー値をフレームごとに算出し、FIFO(First-In First-Out)形式で第1の個数のフレームのエネルギー値を記憶するフレームエネルギー算出及び記憶手段と、フレームエネルギー算出及び記憶手段に、第2の個数のフレームのエネルギー値が格納されたことに応答して、第2の個数のフレームのエネルギー値を所定の統計的手法に従って処理することにより、音声データに含まれる環境雑音の推定値の初期値を算出するための初期値算出手段と、推定値の初期値と、フレームエネルギー算出及び記憶手段に逐次記憶される音声のエネルギー値とに基づいて、音声データに含まれる環境雑音の変化に追従して変化する様に、発話区間を検出するためのエネルギー値のしきい値をフレームごとに逐次算出するための手段と、しきい値に基づいて、第2の個数のフレーム以降のフレームの中で、音声データの発話区間の開始位置又は終了位置に対応するフレームを推定するための発話区間推定手段とを含む。

20

【0015】

環境雑音の推定値の初期値が、第2の個数のフレームのエネルギー値を統計的に処理することにより算出される。以後は、この推定値の初期値と、フレームエネルギー算出及び記憶手段に逐次記憶される音声のエネルギー値とに基づいて、音声データに含まれる環境雑音の変化に追従して変化する様に、発話区間を検出するためのエネルギー値のしきい値をフレームごとに逐次算出する。そのしきい値を用いて音声データの発話区間の開始位置又は終了位置に対応するフレームを推定する。しきい値が、環境雑音の変化に追従して変化するので、正確に発話区間の開始位置又は終了位置を推定できる。

30

【0016】

好ましくは、初期値算出手段は、第2の個数のフレームを、各フレームのエネルギー値の大きさによって、第1のエネルギー値を中心とする第1のクラスと、第1のエネルギーよりも大きな第2のエネルギー値を中心とする第2のクラスとにクラス化するための手段と、第1のエネルギー値を環境雑音の推定値の初期値として出力するための手段とを含む。

40

【0017】

音声信号には、環境雑音と発話音声とが含まれる。各フレームをクラス化すると、環境雑音のみのフレームと、環境雑音と発話音声とを含むフレームとの二つのグループに分類されると思われる。フレームをエネルギーの大きさに従って二つのクラスにクラス化すると、エネルギーの小さな第1のフレームからなるクラスにおいて、環境雑音のみからなるフレームの占める割合が高くなる。そこで、この第1のクラスのフレームのエネルギー値の平均を環境雑音の推定値の初期値とすれば、環境雑音の初期値を信頼性高く推定することができる。

50

【 0 0 1 8 】

より好ましくは、クラスタ化するための手段は、第2の個数のフレームを第1及び第2のクラスタにクラスタ化するための境界値を決定するための手段と、境界値よりも小さなエネルギー値を持つフレームを第1のクラスタに、それ以外のフレームを第2のクラスタに、それぞれ分類するための手段とを含む。

【 0 0 1 9 】

境界値を決定するための手段は、第2の個数のフレームのうち、エネルギー値をキーとしてソートしたときに予め定める第1のソート順位及び第2のソート順位となる二つのフレームを選択するための手段と、選択された二つのフレームのエネルギー値の平均値を算出するための第1の平均値算出手段と、第1の平均値算出手段により算出された平均値より小さいエネルギー値を持つか否かを基準として、第2の個数のフレームを第1及び第2のグループに分類するための手段と、第1及び第2のグループに属するフレームのエネルギー値の平均値をそれぞれ算出するための第2の平均値算出手段と、第2の平均値算出手段により算出された二つの平均値の平均値をさらに算出し、境界値として出力するための第3の平均値算出手段とを含んでもよい。

10

【 0 0 2 0 】

好ましくは、しきい値をフレームごとに逐次算出するための手段は、フレームエネルギー算出及び記憶手段に格納されているフレームのエネルギー値と、環境雑音の推定値の初期値とに基づいて、フレームエネルギー算出及び記憶手段に格納されているフレームの環境雑音のエネルギー値をフレームごとに推定するための手段と、フレームエネルギー算出及び記憶手段に格納されているフレームのエネルギー値のうち、定常的な背景雑音及び発話音声の合計のエネルギー値の最大値をフレームごとに逐次推定するための手段と、推定された環境雑音のエネルギー値と、推定された背景雑音及び発話音声の合計のエネルギー値とに基づいて、発話区間を検出するためのエネルギーのしきい値をフレームごとに算出するための手段とを含む。

20

【 0 0 2 1 】

より好ましくは、発話区間推定手段は、しきい値に基づいて、第2の個数のフレーム以降のフレームの状態を判定するための手段を含み、状態は、非発話状態を含み、環境雑音のエネルギー値をフレームごとに逐次推定するための手段は、1フレーム前の時点において推定された環境雑音のエネルギー値を記憶するための手段と、環境雑音の推定値の初期値が算出された時点で記憶するための手段に環境雑音の推定値の初期値を記憶させるための手段と、記憶するための手段に記憶された値、フレームエネルギー算出及び記憶手段に含まれるフレームのエネルギー値、及びフレームの状態を判定する手段による判定結果に基づいて、以下の式

30

$$b(t) = b(t-1) \times \quad + E(t) \times (1 - \quad) \quad (\text{状態が非発話状態の場合})$$

$$b(t) = b(t-1) \quad (\text{状態が非発話状態以外の場合})$$

ただし \quad は所定の忘却係数、 $E(t)$ は時刻 t におけるフレームのエネルギー値、 \quad に従って時刻 t における背景雑音 $b(t)$ を算出するための手段とを含み、記憶するための手段は、算出された背景雑音 $b(t)$ を記憶する。

【 0 0 2 2 】

40

合計のエネルギー値の最大値をフレームごとに推定するための手段は、フレームエネルギー算出及び記憶手段に格納されているフレームを、エネルギー値をキーとしてソートするための手段と、ソートするための手段によりソートされた結果所定の順位となるフレームのエネルギー値を合計のエネルギー値の最大値 $E_{\max}(t)$ として選択するための手段を含んでもよい。

【 0 0 2 3 】

好ましくは、しきい値をフレームごとに逐次算出するための手段は、時刻 t における発話開始位置検出のためのしきい値 $E_{th_1}(t)$ を、

$E_{th_1}(t) = \quad b(t) + \max(\quad, E_{\max}(t) - b(t)) \times \text{第1の定数}$
に従って算出するための手段を含む。

50

【 0 0 2 4 】

さらに好ましくは、しきい値をフレームごとに逐次算出するための手段は、さらに、時刻 t における発話終了位置検出のためのしきい値 $E_{th_2}(t)$ を、

$$E_{th_2}(t) = b(t) + \max(\quad , E_{\max}(t) - b(t)) \times \text{第2の定数}$$

ただし第2の定数 < 第1の定数、に従って算出するための手段を含む。

【 0 0 2 5 】

発話区間検出装置はさらに、発話の先頭からの各フレームの音声データの最大エネルギー値又は所定のデフォルト基準値のいずれか大きい方を用いて各フレームの音声データを正規化し、各フレームの音声特徴パラメータとして出力するための音声エネルギー正規化手段を含んでもよい。

10

【 0 0 2 6 】

発話の先頭からの各フレームの音声データの最大エネルギー値又は所定のデフォルト基準値のいずれか大きい方を用いて正規化するので、発話の終了まで待たずに、擬似的にはあるが実時間で正規化することが可能になる。したがって、音声特徴パラメータの一つとして音声エネルギーを実時間で得ることができる。

【 0 0 2 7 】

好ましくは、音声エネルギー正規化手段は、正規化の基準値を記憶するための基準値記憶手段と、フレームエネルギー算出及び記憶手段により算出された音声エネルギーが、基準値記憶手段に記憶された基準値を超えていることを検出し、検出信号を出力するための検出手段と、検出手段により出力される検出信号に応答して、基準値記憶手段に記憶された基準値を、フレームエネルギー算出及び記憶手段により算出された値で置換するための手段と、フレームエネルギー算出及び記憶手段により算出された音声エネルギー値を、基準値記憶手段に記憶された基準値で除算することにより、当該フレームの音声エネルギーを正規化するための除算手段とを含む。

20

【 0 0 2 8 】

さらに好ましくは、発話区間検出装置は、発話区間推定手段により、発話区間の終了位置に対応するフレームが推定されたことに応答して、基準値記憶手段の記憶内容を、所定のデフォルト値で置換するための手段をさらに含む。

【 0 0 2 9 】

発話区間検出装置は、所定のデフォルト値を、発話区間検出装置の起動時に与えられたオプション値に基づいて設定するための手段をさらに含んでもよい。

30

【 0 0 3 0 】

本発明の第2の局面に係るコンピュータプログラムは、上記したいずれかの発話区間検出装置としてコンピュータを動作させるためのものである。

【 0 0 3 1 】

本発明の第3の局面にかかる音声エネルギー正規化装置は、フレーム化された音声データの正規化音声エネルギーを実時間で算出するための音声エネルギー正規化装置であって、正規化の基準値を記憶するための基準値記憶手段と、フレームごとの音声データの音声エネルギーを算出するための手段と、音声エネルギー算出手段により算出された音声エネルギーが、基準値記憶手段に記憶された基準値を超えていることを検出し、検出信号を出力するための手段と、検出手段により出力される検出信号に応答して、基準値記憶手段に記憶された基準値を、音声エネルギー算出手段により算出された値で置換するための手段と、音声エネルギー算出手段により算出された音声エネルギーを、基準値記憶手段に記憶された基準値で除算することにより、当該フレームの音声エネルギーを正規化するための除算手段とを含む。

40

【 0 0 3 2 】

発話区間の最初においては、デフォルトの値を基準値として音声エネルギーを正規化する。発話区間の途中でフレームの音声エネルギーが基準値を超えると、フレームの音声エネルギーを新たな基準値として音声エネルギーを正規化する。発話区間の終了まで到達しなくても擬似的にはあるが音声エネルギーの実時間での正規化が可能になる。発話区間の最初では誤差が生ずるが、実際に音声エネルギーが発話区間中での最大値まで到達すると、後は正確

50

な正規化が行なえる。またデフォルトの値を適切に選ぶことにより、発話区間の最初に生ずる誤差も小さく抑えることができる。

【 0 0 3 3 】

好ましくは、音声エネルギー正規化装置は、発話区間の終了を検出して発話終了検出信号を出力するための手段と、発話終了検出信号に応答して、基準値記憶手段の記憶内容を、所定のデフォルト値で置換するための手段とをさらに含む。

【 0 0 3 4 】

発話区間が終了すると、基準値を再びデフォルトの値に再設定できる。音声エネルギーを、フレームごとに適切な基準値を使用して正規化できる。

【 0 0 3 5 】

さらに好ましくは、音声エネルギー正規化装置は、所定のデフォルト値を、音声エネルギー正規化装置の起動時に与えられたオプション値に基づいて設定するための手段をさらに含む。

【 0 0 3 6 】

起動時のオプション値によってデフォルト値を設定できるので、様々なオプション値をデフォルト値として音声エネルギー正規化装置を動作させることができる。その結果、音声エネルギーの正規化処理をより適切に実現することが容易になる。

【 0 0 3 7 】

本発明の第 4 の局面に係るコンピュータプログラムは、上記したいずれかの音声エネルギー正規化装置としてコンピュータを動作させるためのものである。

【 0 0 3 8 】

本発明の第 5 の局面に係るコンピュータは、上記した第 2 の局面に係るコンピュータプログラム、又は第 4 の局面に係るコンピュータプログラムによりプログラムされ、発話区間検出装置又は音声エネルギー正規化装置として動作する。

【 発明を実施するための最良の形態 】

【 0 0 3 9 】

本実施の形態に係る発話区間検出装置は、フレーム化して入力される音声信号に基づき、統計的手法によって発話区間検出の際のしきい値を変化させる。その際、装置の立上がり時の遅延をできるだけ少なくするとともに、突発的な雑音があっても安定して発話区間の検出を行なうことができるよう、統計的手法を工夫している。また、音声認識のための特徴量パラメータとしてフレームの正規化した音声エネルギーを算出する際、実時間処理によって、擬似的な正規化ができるような工夫をしている。

【 0 0 4 0 】

[発話区間の検出原理]

図 1 に、音声信号と、本実施の形態において発話区間の検出に使用される手法で使用される様々なパラメータとを示す。図 1 を参照して、音声信号 2 0 に対し、発話開始しきい値 2 2 と発話終了しきい値 2 4 という二つのしきい値を用いて発話の開始位置 2 6 及び終了位置 2 8 を判定する。これら発話開始しきい値 2 2 及び発話終了しきい値 2 4 は、入力波形データからフレーム単位で算出されるエネルギーから統計的手法により定められる。これらを定める手法については後述する。

【 0 0 4 1 】

図 1 において、発話区間の判定の際に使用される時間的パラメータ T 1 から T 6 は以下の意味を持つ。

【 0 0 4 2 】

T 1 : プリロール時間 あるフレームが発話の開始位置であると判定されたとき、そのフレームからさらにこのプリロール時間だけさかのぼった位置 (図 1 の参照符号 2 6) のフレームに、発話開始フレームとしてのマークが付される。

【 0 0 4 3 】

T 2 : 発話開始判定時間 発話が開始したと判定されるための第 1 の条件として、フレーム単位のエネルギー値が連続して発話開始しきい値を超えなければならない時間。

【 0 0 4 4 】

T 3 : 最短発話時間 発話開始と判定されるために、フレーム単位のエネルギー値が連続して超えなければならない最小時間。エネルギー値が発話開始しきい値を T 2 時間連続して超え、かつ T 3 時間連続して超えてはじめて発話開始と判定される。

【 0 0 4 5 】

T 4 : 最長無音時間 発話状態でフレーム単位のエネルギー値が発話終了しきい値を下回っても、発話終了と判定されない最長の時間。

【 0 0 4 6 】

T 5 : 発話終了判定時間 発話が終了したと判定されるための第 1 の条件として、フレーム単位のエネルギー値が連続して発話終了しきい値を下回らなければならない時間。エネルギー値が発話終了しきい値を T 5 時間連続して下回り、かつ T 4 時間連続して下回った場合、発話終了と判定される。

10

【 0 0 4 7 】

T 6 : アフタロール時間 あるフレームで発話終了と判定されたとき、そのフレームからさらにこのアフタロール時間だけ下った位置のフレーム (図 1 の参照符号 2 8) に、発話終了フレームとしてのマークが付される。

【 0 0 4 8 】

図 1 の水平軸付近に記載されている S 1 から S 4 の符号は、後述する手法により決定される、各フレームの状態を示す。図 2 に、フレームの状態の遷移を示す。

【 0 0 4 9 】

20

図 2 を参照して、フレームは 4 つの状態 (非発話状態 (S 1) 3 0、発話開始状態 (S 2) 3 2、発話状態 (S 3) 3 4、及び発話終了状態 (S 4) 3 6) の間を遷移する。状態間の遷移は以下の様にして行なわれる。

【 0 0 5 0 】

(1) 非発話状態 (S 1) 3 0 で、フレームのエネルギー値が発話開始しきい値 2 2 を上回ると状態は発話開始状態 (S 2) 3 2 に遷移する (アーク 4 2) 。

【 0 0 5 1 】

(2) 発話開始状態 (S 2) 3 2 が、一定時間 T 3 だけ継続すると状態は発話状態 (S 3) 3 4 となる (アーク 4 8) 。

【 0 0 5 2 】

30

(3) 発話開始状態 (S 2) 3 2 で、フレームのエネルギー値が発話開始しきい値 2 2 を下回ると状態は非発話状態 (S 1) 3 0 に遷移する (アーク 4 6) 。

【 0 0 5 3 】

(4) 発話状態 (S 3) 3 4 で、フレームのエネルギー値が発話終了しきい値 2 4 を下回ると状態は発話終了状態 (S 4) 3 6 に遷移する (アーク 5 2) 。

【 0 0 5 4 】

(5) 発話終了状態 (S 4) 3 6 が、一定時間 T 4 だけ継続すると状態は非発話状態 (S 1) 3 0 に遷移する (アーク 5 8) 。

【 0 0 5 5 】

(6) 発話終了状態 (S 4) 3 6 で、フレームのエネルギー値が発話終了しきい値 2 4 を上回ると状態は発話状態 (S 3) 3 4 に戻る (アーク 5 4) 。

40

【 0 0 5 6 】

(7) それ以外の場合、状態は現在の状態を維持する (アーク 4 0、4 4、5 0 及び 5 6) 。

【 0 0 5 7 】

上記した種々のパラメータは、本実施の形態の装置では、装置の起動時に手操作により設定される。設定のないものはデフォルト値が用いられる。パラメータ設定の部分は本発明と直接関係をもたないため、以下の説明では詳細には説明しない。

【 0 0 5 8 】

[フレームの構成]

50

後述する様に、本実施の形態に係る装置は、音声入力信号をフレーム単位で処理する。図3にフレーム及びフレームシフトの概念を説明するための模式図を示す。

【0059】

図3を参照して、各フレーム70、72、74、...はフレーム長 $T_w = 30$ ミリ秒の長さの音声信号である。本実施の形態では、このフレームを10ミリ秒単位で時間軸上を移動させながら順次音声信号をフレーム化する。この移動量をフレームシフト量と呼ぶ。従って、本実施の形態の装置の処理対象となる音声データは、フレーム長30ミリ秒、フレームシフト量10ミリ秒である。

【0060】

また、各フレームのエネルギーは、当該フレーム中のデータに窓関数80（ハミング窓）で示される値を乗算して総和を計算することにより得られる。フレームごとのエネルギーの算出方法については後述する。

【0061】

本実施の形態の装置では、通常は100フレームのデータを統計的に処理することにより発話開始しきい値22及び発話終了しきい値24を動的に計算する。この様に動的な処理を行なう場合、ある程度のデータが集積されないと処理を開始することができない。他方で、あまり多くのデータを使用して統計的処理を行なおうとすると、装置が適切に動作するまでの時間的遅延が長くなり、発話の最初を正しく検出できなくなるおそれがある。

【0062】

そこで、本実施の形態の装置では、処理の開始後、最初の400ミリ秒までは無音状態であると仮定し、この間に40フレーム分のデータをフレームバッファに収集する。この40フレーム分のデータを用いて環境雑音の初期値を求め、その値を用いてさらにしきい値の初期値を決める。以後、100フレーム分のデータが収集されるまで、フレームデータをフレームバッファに蓄積しながら、収集したデータを用いてしきい値を動的に計算する。100フレームに達したら、以後、FIFO（First-In First-Out）形式でフレームデータを100個に維持しながらしきい値の計算を行なう。なお、この最大のフレーム数（フレームバッファ内に記憶され使用される最大のフレーム数）をフレームバッファサイズと呼ぶことにする。また、環境雑音の初期値を求めるために使用するフレームの数を初期バッファサイズと呼ぶ。すなわち、本実施の形態の装置ではフレームバッファサイズは100、初期バッファサイズは40である。

【0063】

なお、これらのフレームバッファサイズ及び初期バッファサイズは一例であって、これ以外の値を用いることも考えられる。

【0064】

以下の説明では、入力されるフレームの番号を t （0 t ）で表す。フレームは10ミリ秒ごとに入力されるので、 t はまた時刻も表す。従って、以下の説明では単に「 t 番目のフレーム」を「時刻 t におけるフレーム」という表現で表すこともある。

【0065】

こうした処理を行なうことで、処理開始時の遅延は400ミリ秒となり、実用上の問題は見られない。通常は100個のフレームデータを用いてしきい値を計算するので、信頼性高く発話区間の検出を行なうことができる。

【0066】

[装置の構成]

図4は、本実施の形態に係る発話区間検出装置の構成を示す機能的ブロック図である。図4を参照して、この発話区間検出装置100は、マイク102から与えられる音声信号の中で発話区間を検出するためのものである。発話区間検出装置100は、マイク102から与えられる音声信号を標本化し、量子化することによりデジタル化し、さらに上記した形式のフレームデータとして10ミリ秒ごとに出力するとともに、フレームデータを出力したことを示すフレーム出力信号124を出力するための音声入力部104と、音声入力部104から与えられる複数個のフレームデータを記憶するための入力バッファ106

10

20

30

40

50

とを含む。

【 0 0 6 7 】

発話区間検出装置 1 0 0 はさらに、入力バッファ 1 0 6 からフレームデータを読み出してエネルギー値などのフレーム情報を算出するためのフレーム情報算出部 1 0 8 と、フレーム情報算出部 1 0 8 の出力するフレーム情報を記憶するためのフレームバッファ 1 1 0 とを含む。フレームバッファ 1 1 0 のバッファサイズは、前述した通り 1 0 0 フレーム分である。フレームバッファ 1 1 0 は、入力されたフレーム情報を F I F O 形式で 1 0 0 個保持することができる。

【 0 0 6 8 】

本実施の形態では、フレーム情報算出部 1 0 8 は、次の式に従って時刻 t におけるフレームの音声エネルギー $E(t)$ を算出する。

【 0 0 6 9 】

【 数 1 】

$$E(t) = \log \frac{\sum_{i=1}^N |S_i| \times H_i}{N} \times 20$$

ただし、 N は 1 フレーム中のデータサンプル数、 S_i ($i = 1 \sim N$) はデータの値、 H_i ($i = 1 \sim N$) はハミング窓関数の値を、それぞれ示す。

【 0 0 7 0 】

発話区間検出装置 1 0 0 はさらに、フレーム情報算出部 1 0 8 が算出したフレームの音声エネルギーを、発話中の最大パワーを基準として正規化し、フレームの特徴ベクトルの一要素として入力バッファ 1 0 6 中に書込むためのフレーム音声エネルギー正規化処理部 1 2 6 を含む。フレームの音声エネルギーの大きさを一発話中の最大エネルギーで正規化し、特徴量の一つとして音声認識に利用すると効果があることが認められている。しかし、そのためには発話の終了まで待ってフレームエネルギーの最大値を算出する必要がある。しかしそれでは実時間処理を行なうことができない。

【 0 0 7 1 】

そこでフレーム音声エネルギー正規化処理部 1 2 6 は、音声エネルギーのダイナミックレンジを実時間に更新することにより、擬似的にはあるが音声エネルギーを実時間で正規化する機能を持つ。フレーム音声エネルギー正規化処理部 1 2 6 はそのため、図 5 に示すような構成を持つ。

【 0 0 7 2 】

図 5 を参照して、フレーム音声エネルギー正規化処理部 1 2 6 は、発話の先頭部分でまた音声エネルギーの十分大きなフレームがないときに、最大音声エネルギーのデフォルト値として使用されるデフォルト最大値を記憶するためのデフォルト最大値記憶部 1 3 2 と、発話の最初の部分ではデフォルト最大値記憶部 1 3 2 から与えられたデフォルト最大値を記憶し、発話途中でデフォルト最大値より大きな音声エネルギーを持つフレームが検出された場合に、当該音声エネルギーの値を記憶するための最大値記憶部 1 3 4 と、フレーム情報算出部 1 0 8 からの音声エネルギー 1 2 8 を最大値記憶部 1 3 4 に記憶されている最大値で除算し、結果を入力バッファ 1 0 6 の該当フレームの特徴量の一つとして書込むための除算部 1 3 6 と、最大値記憶部 1 3 4 の出力とフレーム情報算出部 1 0 8 からの音声エネルギー 1 2 8 とを受けて両者の値を比較し、比較結果信号 1 3 9 を最大値記憶部 1 3 4 に与えるための比較部 1 3 8 とを含む。比較結果信号 1 3 9 は、音声エネルギー 1 2 8 により示される値が最大値記憶部 1 3 4 に記憶された最大値を上回ると H (ハイ) レベルとなり、それ以外の場合は L (ロー) レベルとなる。なお、デフォルトの値は、オプションとしてこの装置 (プログラム) 起動時に与えられた値があれば、その値で書換えられる。

【 0 0 7 3 】

最大値記憶部 1 3 4 は、状態判定部 1 1 8 から与えられる信号 2 0 0 によって発話が終了したことが示されると、デフォルト最大値記憶部 1 3 2 の値を新たな最大値として記憶

10

20

30

40

50

し、比較部 138 からの比較結果信号 139 が H レベルとなると、音声エネルギー 128 により示される値を新たな最大値として記憶する。したがって、最大値記憶部 134 に記憶される値は、発話開始時にはデフォルト最大値記憶部 132 に記憶されていたデフォルト値となり、発話の進行とともに音声エネルギーがデフォルト値を上回るものが出現するとその音声エネルギーとなる。以下、発話の進行中には同様の処理が繰返される。この値を発話中の音声エネルギーの最大値として使用して各フレームの音声エネルギーを正規化することにより、擬似的にはあるが、実時間で音声エネルギーの正規化を行なうことができる。

【0074】

なお、デフォルトの値は予め実験により適切な値を決めておくことが望ましい。

【0075】

発話区間検出装置 100 はさらに、音声入力部 104 からのフレーム出力信号 124 を受け、入力バッファ 106、フレーム情報算出部 108 及びフレームバッファ 110 の読出ポイント及び書込ポイント、並びにそれらへの書込み・読出しのタイミングを管理するための入出力・アドレス管理部 114 と、発話区間検出装置 100 の処理開始後 400 ミリ秒までの間にフレームバッファ 110 に格納されるフレームデータ 160 を読出し、初期環境雑音を算出するための初期環境雑音算出部 112 と、フレームバッファ 110 からのフレームデータ 192、初期環境雑音算出部 112 からの初期環境雑音の推定値 194、及び現在の状態が非発話状態 (S1) 30 (図 2 参照) か否かを示す信号 190 を受け、それらから発話開始しきい値 22 及び発話終了しきい値 24 を動的に算出し、しきい値の値を示す信号 198 とし出力するための動的しきい値算出部 116 とを含む。

【0076】

入力バッファ 106、フレームバッファ 110 などは半導体記憶装置などにより実現される。入出力・アドレス管理部 114 はタイマを装備しており、音声入力部 104 による音声データのデジタル化に同期して、入力バッファ 106、フレームバッファ 110 などへの書込みのポイント、それらからの読出しポイントを管理する。入出力・アドレス管理部 114 はまた、起動後 400 ミリ秒までのフレームを処理する際には H レベル、それ以後は L レベルの値をとる初回フラグ 196 を動的しきい値算出部 116 に与える機能も持つ。動的しきい値算出部 116 の処理は、初回フラグ 196 及び信号 190 の値によって制御される。

【0077】

発話区間検出装置 100 はさらに、動的しきい値算出部 116 から出力されたしきい値の値を示す信号 198 及びフレームバッファ 110 からのフレームデータ 192 とから、後述する方法に従ってフレームの状態を判定し、状態を表す信号 200 を出力するための状態判定部 118 と、状態判定部 118 の出力する状態を表す信号 200 を受け、入力バッファ 106 から状態の確定したフレームに対応する入力データを読出して予め定められた算出方法によってこのフレームの音声の特徴ベクトルを算出し、さらに発話区間の開始又は終了フレームの場合には、それらを示すマークを当該特徴ベクトル 122 に付して出力するための特徴ベクトル出力部 120 とを含む。状態判定部 118 はまた、現在の状態が非発話状態 (S1) 30 か否かを示す信号 190 を生成し、動的しきい値算出部 116 に与える機能も持つ。

【0078】

図 6 は初期環境雑音算出部 112 のブロック図であって、初期環境雑音算出部 112 は、フレームバッファ 110 から与えられるフレーム情報のうち、フレームごとのエネルギー値を昇順にソートしてソート後フレームエネルギー記憶部 142 に格納させるためのソート処理部 140 と、ソート処理部 140 に格納されたフレームごとのエネルギー値のうち、下位から 25% 及び 75% の大きさにあたる位置のフレームのエネルギーを算出し、それぞれ後述するクラスタリング処理のシードとなる値 e_{m1} 及び e_{m2} として出力するためのサイズ算出部 144 と、この値 e_{m1} 及び e_{m2} を記憶するための記憶部 146 とを含む。

【0079】

初期環境雑音算出部 112 はさらに、記憶部 146 から値 e_{m1} 及び e_{m2} を読出し、

10

20

30

40

50

その平均値 e_{average} を算出するための第 1 の平均値算出部 148 と、第 1 の平均値算出部 148 が出力する平均値を境界値としてそれより大きいエネルギー値を持つか否かを基準として、ソート後フレームエネルギー記憶部 142 中の各フレームを二つのクラスタ C1 及び C2 に分類するためのフレーム分類部 150 と、フレーム分類部 150 により得られた二つのクラスタ C1 及び C2 の各々に属するフレームのエネルギー値の平均値 E_{m1} 及び E_{m2} を次の式に従って算出するための第 2 の平均値算出部 152 とを含む。

【0080】

【数 2】

$$Em1 = \frac{1}{I1} \sum_{i_1=1}^N E(i_1) \quad Em2 = \frac{1}{I2} \sum_{i_2=I1+1}^N E(i_2)$$

10

ただし、N はフレームバッファ 110 内のフレーム数、I1 は e_{average} より小さいエネルギー値を持ち、クラスタ C1 に属するフレームの数、I2 は e_{average} より大きいエネルギー値を持ち、クラスタ C2 に属するフレームの数を、それぞれ表す。

【0081】

初期環境雑音算出部 112 はさらに、第 2 の平均値算出部 152 によって算出された二つの平均値 E_{m1} 及び E_{m2} をそれぞれ新たな値 e_{m1} 及び e_{m2} として記憶部 146 に記憶させ、さらに第 1 の平均値算出部 148、フレーム分類部 150、及び第 2 の平均値算出部 152 に先ほどの処理を繰返し実行させ、その結果得られた平均値 E_{m1} を初期環境雑音の推定値 (e_{m1}) 194 として図 4 に示す動的しきい値算出部 116 に与えるための判定部 154 とを含む。

20

【0082】

以下に、第 1 の平均値算出部 148、フレーム分類部 150 及び第 2 の平均値算出部 152 により行なわれる処理について、図 4、及び図 6 から図 9 を参照して説明する。一般に、図 4 に示すフレームバッファ 110 に記憶されている各フレームのエネルギー値は、図 7 に示される様に、入力音声信号のエネルギーの大きさに従って変動する。これをエネルギーの大きさに従って昇順にソートすると図 8 の様になると想定される。ソート処理部 140 が行なうソート処理はこうした処理であり、ソート後フレームエネルギー記憶部 142 に記憶されているフレーム情報は図 8 に示すものに対応している。

【0083】

30

図 8 の様にソートすることで、エネルギー値のヒストグラムを容易にとることができる。図 9 にその例を示す。音声信号に環境雑音と発話成分とが含まれているとすれば、環境雑音のみのフレームのエネルギー値と、発話成分を含むフレームのエネルギー値とは、それぞれ別々の値を中心として分布することになると思われる。そして、それらは図 9 に示されるようなヒストグラムにおいて、エネルギーの比較的低い部分のピークと、エネルギーの比較的高い部分のピークとの二つのピークを形成することになるであろう。

【0084】

図 6 に示す第 1 の平均値算出部 148、フレーム分類部 150、及び第 2 の平均値算出部 152 が行なっているのは、最初にエネルギー値の 25% と 75% の部分とをピークの初期位置として、上記した二つのピークをその後の演算により求め、ソート後フレームエネルギー記憶部 142 に格納されている各フレームを、環境雑音側のピークに近いフレームと、発話部分側のピークに近いフレームとの二つのクラスタにクラスタ化する処理である。

40

【0085】

図 10 は、図 4 に示す動的しきい値算出部 116 の機能的ブロック図である。図 10 を参照して、動的しきい値算出部 116 は、フレームデータ 192 を受け、フレームバッファ 110 に格納されているソート後のフレーム情報のうち、小さい方から 90% の位置にあるフレームのエネルギーを、t 番目までのフレームバッファサイズ分の数のフレームにおける最大エネルギー $e_{\text{max}}(t)$ (最大エネルギー信号 182) として出力するための最大エネルギー算出部 176 と、フレームデータ 192 を受け、後述する式に従って環境雑音の推定値を算出するための環境雑音算出部 170 と、1 フレーム分だけ前の処理で算出された

50

環境雑音の推定値 $b(t-1)$ を記憶するための記憶部 174 とを含む。

【0086】

動的しきい値算出部 116 はさらに、記憶部 174 に記憶されている 1 フレーム分だけ前の推定値 $b(t-1)$ と、環境雑音算出部 170 から与えられる環境雑音の推定値と、初期環境雑音の推定値 ($em1$) 194 とを受けて、初回フラグ 196 が H レベルであれば初期環境雑音の推定値 ($em1$) 194 を、初回フラグ 196 が L レベルでかつ状態を示す信号 190 が非発話状態を示す値であれば環境雑音算出部 170 の出力を、初回フラグ 196 が L レベルでかつ状態を示す信号 190 が非発話状態を示す値でなければ記憶部 174 の出力を、それぞれ選択して t 番目のフレームに対する環境雑音 $b(t)$ として出力するための選択部 172 とを含む。選択部 172 の出力は記憶部 174 に与えられ記憶される。

10

【0087】

動的しきい値算出部 116 はさらに、最大エネルギー算出部 176 からの最大エネルギーと、選択部 172 からの t 番目のフレームにおける環境雑音 $b(t)$ とに基づいて発話開始しきい値 22 及び発話終了しきい値 24 を動的に算出するためのしきい値算出部 178 を含む。しきい値算出部 178 の出力する、しきい値を表す信号 198 は状態判定部 118 に与えられ、状態判定に用いられる。

【0088】

環境雑音算出部 170 は、フレームバッファ 110 に記憶されたフレームデータ 192 の中で t 番目のフレームのエネルギー $E(t)$ 、及び記憶部 174 に記憶された $t-1$ 番目のフレームに対する環境雑音 $b(t-1)$ とから次の式 1 に従って環境雑音の推定値 $b'(t)$ を算出する。

20

[式 1]

$$b'(t) = b(t-1) \times \alpha + E(t) \times (1 - \alpha)$$

ここで、 α は予め定められた忘却係数、 $E(t)$ は t 番目のフレームのエネルギーを表す。忘却係数は 0 以上 1 以下の値であるが、本実施の形態では 0.8 を用いる。

【0089】

選択部 172 は、状態が非発話状態以外であれば記憶部 174 から出力される $t-1$ 番目のフレームに対する環境雑音 $b(t-1)$ を選択する。従ってこの場合には環境雑音は変化しない。状態が非発話状態であれば、選択部 172 は、環境雑音算出部 170 から出力される環境雑音の推定値 $b'(t)$ を選択する。

30

【0090】

従って、環境雑音算出部 170 から出力される時刻 t における環境雑音 $b(t)$ は以下の通りの式で表される。ただし $E(t)$ は時刻 t におけるフレームのエネルギー値、 α は前述の忘却係数である。

[式 2]

$$b(t) = b(t-1) \times \alpha + E(t) \times (1 - \alpha) \quad (\text{状態が非発話状態の場合})$$

$$b(t) = b(t-1) \quad (\text{状態が非発話状態以外の場合})$$

しきい値算出部 178 は以下の式に従って発話開始しきい値 E_{th1} 及び発話終了しきい値 c を動的に算出する。

40

[式 3]

0 $t < 400$ ミリ秒では

$$E_{th1}(t) = b(t) + \alpha_1$$

$$E_{th2}(t) = b(t) + \alpha_2$$

400 ミリ秒 t では

$$E_{th1}(t) = b(t) + \max(\alpha_1, E_{max}(t) - b(t)) \times \alpha_1$$

$$E_{th2}(t) = b(t) + \max(\alpha_2, E_{max}(t) - b(t)) \times \alpha_2$$

ただし、 α_1 は発話の最低ダイナミックレンジで、本実施の形態では 20 dB である。また α_1 及び α_2 はそれぞれ発話開始しきい値比率及び発話終了しきい値比率であり、それぞれ実験的に定められる、0 以上で 1 以下の定数である。本実施の形態では $\alpha_1 = 0.25$ 、

50

$2 = 0.20$ を用いる。

【0091】

こうして算出された発話開始しきい値 E_{th1} 及び発話終了しきい値 E_{th2} が、図1を参照して説明した発話区間の検出時の発話開始しきい値 2_2 及び発話終了しきい値 2_4 として用いられる。

【0092】

[装置の動作]

以上構成を述べた装置は以下のように動作する。

【0093】

-起動時-

起動時には、処理に必要となるバッファ及びオプションの値を格納するためのエリアを記憶装置に確保する。起動時に与えられるオプションの値を調べ、オプションの値に誤りがなければ当該オプションに、与えられた値を設定する。オプションの値が与えられなかったものにはデフォルト値を設定する。与えられたオプションの値に誤りがあれば、その旨のメッセージを表示して処理を終了する。図5に示すフレーム音声エネルギー正規化処理部126のデフォルト最大値記憶部132については、起動時にオプションの値が与えられれば、その値をデフォルトの値として記憶し、さらに最大値記憶部134に記憶する。オプションの値が与えられなければ、プログラム上のデフォルト値をデフォルト最大値記憶部132に記憶し、さらに最大値記憶部134に記憶する。

【0094】

各バッファの書き込みポイント及び読出しポイントをそれぞれ初期値に設定する。

【0095】

なお、起動後、実際の処理を開始する時刻(フレーム番号)を $t = 0$ とする。このときのフレームの状態は非発話状態に設定される。以後、図4に示す音声入力部104は、マイク102からの電気信号を10ミリ秒ごとに、30ミリ秒のフレーム長でデジタル化する。

【0096】

-0ミリ秒から400ミリ秒まで-

入出力・アドレス管理部114からの初回フラグ196はHレベルである。音声入力部104は、発話判定に必要なデータ数が集まると、1回の処理で引き渡す数として予め定められた数のデータを入力バッファ106の、バッファ書き込みポイントにより指定されるアドレスに書込む。

【0097】

フレーム情報算出部108は、入力バッファ106の、読出しポイントにより指定されるアドレスから1フレーム分のデータを読出し、フレームエネルギーを算出してフレームバッファ110の当該フレームに対応するエリアに書込む。フレーム情報算出部108はまた、算出されたフレームエネルギーをこのフレームの音声エネルギー128として図5に示す除算部136、比較部138及び最大値記憶部134に与える。比較部138は、最大値記憶部134に記憶された値と音声エネルギー128により示される値とを比較し、比較結果信号139を最大値記憶部134に与える。音声エネルギー128により示される値が最大値記憶部134に記憶されている値を上回ったことが検出されると、比較結果信号139はHレベルとなり、最大値記憶部134は比較結果信号139がHレベルとなったことに応答して、これまで記憶していた値に代えて音声エネルギー128により表される値を記憶する。

【0098】

除算部136は、音声エネルギー128により表される値を最大値記憶部134に記憶された値で除算して正規化された音声エネルギーを算出する。正規化された音声エネルギー130は、入力バッファ106中で該当するフレームの、正規化音声エネルギーのフィールドに書込まれる。以後、フレーム情報算出部108とフレーム音声エネルギー正規化処理部126とは、これと同様の動作を各フレームに対して繰返す。

10

20

30

40

50

【 0 0 9 9 】

初期環境雑音算出部 1 1 2 は、フレーム情報算出部 1 0 8 によりフレームバッファ 1 1 0 に書込まれたフレームエネルギーを読み出し、初期環境雑音の算出を行なう。時刻 0 ミリ秒から 4 0 0 ミリ秒の間は、状態の判定は行なわない。

【 0 1 0 0 】

次に、図 6 を参照して、初期環境雑音算出部 1 1 2 の動作について説明する。ソート処理部 1 4 0 は、フレームバッファ 1 1 0 から読み出したフレームエネルギーの値 1 6 0 をソートし、ソート後フレームエネルギー記憶部 1 4 2 に格納する。t = 0 では読み出されるフレームエネルギーの値は一つ (E (0)) だけなので、その値をソート後フレームエネルギー記憶部 1 4 2 の 1 番目の領域に書込む。2 回目以後は、ソート後フレームエネルギー記憶部 1 4 2 に前のソートの結果が既に書込まれており、そこに新たに一つのフレームエネルギーをその大きさに従った位置に追加するだけでよい (ヒープソート)。従って、ソート処理は少ない計算量で実行できる。

10

【 0 1 0 1 】

起動後、0 ミリ秒から 4 0 0 ミリ秒の間は、シーズ算出部 1 4 4 以後の処理部は動作しない。

【 0 1 0 2 】

- 4 0 0 ミリ秒経過時 -

起動後 4 0 0 ミリ秒が経過すると、フレームバッファ 1 1 0 には 4 0 個のフレームデータ (E (0) ~ E (39)) のエネルギー値が格納されている。この状態が図 7 に対応する。ソート後フレームエネルギー記憶部 1 4 2 には、これら 4 0 個のフレームのエネルギー値を昇順にソートしたものが格納されている。この状態が図 8 に対応する。

20

【 0 1 0 3 】

フレーム情報算出部 1 0 8 及びフレーム音声エネルギー正規化処理部 1 2 6 は、4 0 0 ミリ秒経過までと同様に動作する。

【 0 1 0 4 】

除算部 1 3 6 は、音声エネルギー 1 2 8 により表される値を最大値記憶部 1 3 4 に記憶された値で除算して正規化された音声エネルギーを算出する。正規化された音声エネルギー 1 3 0 は、入力バッファ 1 0 6 中で該当するフレームの、正規化音声エネルギーのフィールドに書込まれる。

30

【 0 1 0 5 】

シーズ算出部 1 4 4 は、ソート後フレームエネルギー記憶部 1 4 2 に格納されている 4 0 個のフレームエネルギーのうち、小さい方から 2 5 % 及び 7 5 % に相当する値を算出する。この値は記憶部 1 4 6 に記憶され、第 1 の平均値算出部 1 4 8、フレーム分類部 1 5 0 及び第 2 の平均値算出部 1 5 2 により行なわれるクラスタリングのシードとなる。

【 0 1 0 6 】

第 1 の平均値算出部 1 4 8 は、記憶部 1 4 6 からこのシード e m 1 及び e m 2 の平均値を算出しフレーム分類部 1 5 0 に与える。フレーム分類部 1 5 0 は、全てのフレームについて、そのエネルギー値がシード e m 1 及び e m 2 のいずれに近いかを基準として、4 0 個のフレームを二つのクラスタに分類し、分類した結果を第 2 の平均値算出部 1 5 2 に与える。

40

【 0 1 0 7 】

第 2 の平均値算出部 1 5 2 は、それら二つのクラスタの各々について、そのクラスタに属するフレームのエネルギー値の平均値 E m 1 及び E m 2 を算出し判定部 1 5 4 に与える。

【 0 1 0 8 】

判定部 1 5 4 は、第 2 の平均値算出部 1 5 2 から与えられた E m 1 及び E m 2 を記憶部 1 4 6 に新たな e m 1 及び e m 2 として記憶させ、先ほどと同じ処理を第 1 の平均値算出部 1 4 8、フレーム分類部 1 5 0、及び第 2 の平均値算出部 1 5 2 に実行させる。こうして再び得られた E m 1 及び E m 2 のうち、E m 1 を初期環境雑音の推定値 1 9 4 (e m 1) として動的しきい値算出部 1 1 6 に与える。

50

【 0 1 0 9 】

図 10 を参照して、動的しきい値算出部 1 1 6 の動作について説明する。動的しきい値算出部 1 1 6 の選択部 1 7 2 は、 $b(t)$ の初期値として初期環境雑音の推定値 1 9 4 である e_{m1} を選択し、記憶部 1 7 4 及びしきい値算出部 1 7 8 に与える。記憶部 1 7 4 はこの値を記憶する。

【 0 1 1 0 】

一方、最大エネルギー算出部 1 7 6 は、ソート後フレームエネルギー記憶部 1 4 2 に記憶されているソートされているフレームエネルギーの値のうち、小さい方から 90% に相当するエネルギー値を計算し、最大エネルギー値 (E_{max}) 1 8 2 としてしきい値算出部 1 7 8 に与える。

10

【 0 1 1 1 】

しきい値算出部 1 7 8 は、選択部 1 7 2 から与えられる環境雑音の推定値 e_{m1} と、最大エネルギー算出部 1 7 6 からの最大エネルギー値 (E_{max}) 1 8 2 とに基づき、前述の式 3 によって発話開始しきい値 2 2 及び発話終了しきい値 2 4 を算出し (1 9 8)、図 4 に示す状態判定部 1 1 8 に与える。

【 0 1 1 2 】

状態判定部 1 1 8 は、動的しきい値算出部 1 1 6 から与えられる発話開始しきい値 2 2 及び発話終了しきい値 2 4 に基づき、図 1 及び図 2 に示す判定方法に従ってフレームの状態を判定し、その結果を表す信号 2 0 0 を特徴ベクトル出力部 1 2 0 及びフレーム音声エネルギー正規化処理部 1 2 6 に与える。状態判定部 1 1 8 はまた、フレームの状態が非発話状態か否かを示す信号 1 9 0 を動的しきい値算出部 1 1 6 に与える。

20

【 0 1 1 3 】

フレーム音声エネルギー正規化処理部 1 2 6 の最大値記憶部 1 3 4 (図 5 参照) は、状態を表す信号 2 0 0 により発話区間が終了したことが示されると、それまで記憶していた値に代えてデフォルト最大値記憶部 1 3 2 の値を記憶する。この処理により、次の発話に対する音声エネルギーの正規化処理の開始時には、最大パワーとしてデフォルトの値 (又はオプションとして与えられた値) が再び使用される。

【 0 1 1 4 】

特徴ベクトル出力部 1 2 0 は、状態判定部 1 1 8 の処理によって状態が確定したフレームのデータを入力バッファ 1 0 6 から読み出し、そのフレームの特徴ベクトルを算出し、出力 (1 2 2) する。特徴ベクトル出力部 1 2 0 はこのとき、当該フレームが発話開始フレーム又は発話終了フレームであれば、それを示すマークを当該特徴ベクトルに付して出力する。

30

【 0 1 1 5 】

- 4 0 0 ミリ秒から 1 秒まで -

入出力・アドレス管理部 1 1 4 からの初回フラグ 1 9 6 はオフとなる。40 番目のフレームの後、100 番目までのフレームについては、40 番目のフレームに対する処理とほぼ同様である。この間の処理では、フレームバッファ 1 1 0 には 10 ミリ秒ごとに 1 フレーム分のデータが追加されていく。そして、その結果フレームバッファ 1 1 0 に格納されている全てのフレーム情報を用いて状態判定が実行される。

40

【 0 1 1 6 】

また、図 10 に示す動的しきい値算出部 1 1 6 においては、既に記憶部 1 7 4 に一つ前のフレームに対する処理で計算された環境雑音の推定値 $b(t-1)$ が記憶されている。環境雑音算出部 1 7 0 は、記憶部 1 7 4 に記憶された環境雑音の推定値 $b(t-1)$ と、フレームデータ 1 9 2 から得られる t 番目のフレームのエネルギー $E(t)$ とから、式 1 に従って環境雑音の推定値 $b'(t)$ を算出し選択部 1 7 2 に与える。

【 0 1 1 7 】

選択部 1 7 2 は、初回フラグ 1 9 6 の値がオフなので、記憶部 1 7 4 の出力と、環境雑音算出部 1 7 0 の出力とのいずれかを状態を示す信号 1 9 0 の値に従って選択する。すなわち、信号 1 9 0 の表す状態が非発話状態であれば選択部 1 7 2 は環境雑音算出部 1 7 0

50

の出力を選択し、それ以外であれば記憶部 174 の出力を選択する。選択部 172 は、選択した値を示す信号を、記憶部 174 及びしきい値算出部 178 に与える。

【0118】

他の点では、動的しきい値算出部 116 は、40 番目のフレームに対する処理と同様の処理を実行する。状態判定部 118、特徴ベクトル出力部 120、及びフレーム音声エネルギー正規化処理部 126 の動作も同様である。

【0119】

- 1 秒以後 -

101 番目のフレーム以降の処理も、400 ミリ秒から 1 秒までの処理とほぼ同様である。ただしこの処理では、フレームバッファ 110 に記憶されているフレーム情報に新たなフレーム情報を追加する際には、最も古いフレーム情報が削除される。すなわちフレームバッファ 110 は F I F O 形式でデータを格納する。その結果、フレームバッファ 110 には常に 100 フレーム分のフレーム情報が維持される。ソート処理部 140 によるソート処理も同様である。ソート後フレームエネルギー記憶部 142 のうち、最も古いフレームのエネルギー値が削除され、新たなフレームのエネルギー値が、その大きさに従って決まる位置に書込まれる。

10

【0120】

初期環境雑音算出部 112、動的しきい値算出部 116、状態判定部 118 及び特徴ベクトル出力部 120 は、いずれもフレームバッファ 110 に記憶された 100 フレーム分のデータに基づいて、背景雑音の推定、しきい値の算出、状態の判定、及び特徴ベクトル

20

【0121】

こうして、特徴ベクトル出力部 120 から出力されるフレームごとの特徴ベクトル 122 には、そのフレームが発話開始位置であれば発話開始マーカが、発話終了位置であれば発話終了マーカが、それぞれ付されている。このマーカにより、最初の音声データの発話区間（発話開始位置から発話終了位置まで）を検出する事ができる。

【0122】

また、特徴ベクトル 122 にはフレームごとの音声エネルギーを正規化した値が含まれ、これを特徴量として音声認識で利用することができる。またこの音声エネルギーは、発話全体にわたって調べることで算出された最大値ではなく、発話の最初からの最大値によって

30

【0123】

図 11 を参照して、この正規化処理により定められる音声エネルギーの最大値の推移について説明する。図 11 を参照して、従前の方法によれば、発話の終了まで完了した時点で発話の音声エネルギーの最大値を調べ、その値によって音声エネルギーを正規化する。図 11 において、この音声エネルギーの最大値は点線 212 とその後続く太い実線 218 により表される。

【0124】

これに対し上記した実施の形態では、発話の開始時点では一定のデフォルト値（又はオプション値）214 で、点線 212 で示される音声エネルギーの最大値を近似する。さらに音声エネルギーの値がこのデフォルト値より大きくなると（図 11 における太い実線の曲線 216 の部分）、その値で音声エネルギーの最大値の近似値を置換する。発話中で実際の音声エネルギーの最大値位置に到達した後は、この近似値は実際の最大値と等しくなる（太い実線 218 の部分）。

40

【0125】

この正規化処理によって、実時間で音声エネルギーの正規化を行なうことができる。各発話の先頭部分ではデフォルトの値が最大値として使用されるため、多少の誤差は生じるが、デフォルトの値を適当な大きさに定めておくことにより、擬似的な正規化ではあっても十分な効果を得ることができる。

50

【 0 1 2 6 】

-実施の形態の効果-

以上説明した本実施の形態の装置によれば、発話の開始及び終了のための発話開始しきい値及び発話終了しきい値が、実際の音声データを統計的に処理する事により、実際の音声データに従って動的に変化される。環境雑音の変化に追従して変化するしきい値を用いて発話区間の検出ができる。その結果、環境雑音の影響を最小限に抑えて、正しく発話区間を検出する事ができる。

【 0 1 2 7 】

上記した実施の形態の装置では、しきい値を算出する際に用いられるフレームの最大エネルギー値として、実際の最大値の90%のものを用いている。そのため、環境雑音の突発的な変化によるしきい値の大きな変化を抑制する事ができる。また、フレームバッファサイズだけの量のフレームに対する統計的処理によりしきい値を算出しているため、一部のフレームで突出したエネルギー値の変化があっても、しきい値にその変化が与える影響は比較的少なく済む。その結果、安定してしきい値を算出できる。

【 0 1 2 8 】

本実施の形態の装置ではさらに、フレームデータが40個となった時点で状態の判定を開始している。統計処理にはある程度数が必要なので、あまり少ない数のフレームデータを用いたしきい値算出では、状態判定結果の信頼性が低くなる。従って、最低で300ミリ秒程度、望ましくは本実施の形態の装置のように400ミリ秒程度の音声データに基づいてしきい値算出を開始する事がよい。また、処理対象のフレーム数が40個となった時点で状態判定を開始するので、起動後、状態判定の開始までの遅延時間はほぼ400ミリ秒程度となる。この程度の遅延の大きさは実用上で問題とならない程度である。あまり大きな遅延をとるようにすると、発話区間の開始の検出に失敗するおそれがある。また上記実施の形態では、遅延は400ミリ秒であるが、しきい値判定には1000ミリ秒分のデータが使用されるので、少ない遅延で信頼性の高いしきい値算出を行なう事ができる。

【 0 1 2 9 】

[変形例]

上記した実施の形態では、フレームのエネルギー算出の際の窓関数としてハミング窓を用いた。しかし使用可能な窓関数はこれに限らない。ハミング窓、ブラックマン、カイザー、ブラックマン-ハリスなど種々の窓関数のうち、適切と思われるものを用いればよい。

【 0 1 3 0 】

上記実施の形態では、フレームバッファサイズを100、初期バッファサイズを40とした。これらの値は一例であって、これ以外の組合せをとる事もできる。ただし、フレームバッファサイズをあまり大きくとると環境雑音の変化に追従してしきい値を変化させる事が困難になる。またフレームバッファサイズをあまり小さくとると、環境雑音のちょっとした変化に対応してしきい値が変化し、発話区間の検出が安定してできなくなる。また、初期バッファサイズをあまり大きくとると環境雑音の推定までの遅延が大きくなり、発話区間の先頭の検出に失敗する可能性が高くなる。また、当然の事ながら初期バッファサイズはフレームバッファサイズ以下でなければならない。従って、フレームバッファサイズとしては300~2000ミリ秒程度、初期バッファサイズとしては200から500ミリ秒程度がよい。特に、フレームバッファサイズが600~1000ミリ秒程度、初期バッファサイズとして300から450ミリ秒程度が適している。

【 0 1 3 1 】

また、上記した実施の形態では、音声エネルギーの正規化のため、予め算出された固定的な値をデフォルト値として使用している。しかし本発明はそのような実施の形態には限定されない。例えば、このデフォルト値を発話の終了時に直前の発話での最大パワーによって更新することもできる。このとき、最大エネルギーに対して所定の係数 a ($0 < a < 1$ 、好ましくは $0.7 < a < 0.9$ 、さらに好ましくは $a = 0.8$ 程度)を乗算しておくともよい。また、直前の発話だけでなく、過去の複数個の発話での最大エネルギーの関数としてこのデフォルトの値を更新するようにしてもよい。

【 0 1 3 2 】

また、上記した実施の形態では、フレーム内の各音声データの絶対値に窓関数の値を乗じた値の平均値の対数を取り、さらに係数20を掛けることにより求めた対数音声エネルギーを正規化したものを音声エネルギーの特徴パラメータとしている。しかし本発明はそのような実施の形態には限定されず、例えば各音声データの二乗に窓関数の値を乗じた値の平均値の対数を取り、さらに係数10を掛けることで対数音声エネルギーを算出するようにした場合にも本発明は同様に適用できる。

【 0 1 3 3 】

上記した実施の形態の装置は、DSP (Digital Signal Processor) などのプロセッサと、そうしたプロセッサ上で実行されるプログラム (マイクロプログラムを含む。) とにより実現される事が想定される。上記した説明により、そうしたプログラムを作成する事は、当業者には容易であろう。

10

【 0 1 3 4 】

今回開示された実施の形態は単に例示であって、本発明が上記した実施の形態のみに制限されるわけではない。本発明の範囲は、発明の詳細な説明の記載を参酌した上で、特許請求の範囲の各請求項によって示され、そこに記載された文言と均等の意味及び範囲内でのすべての変更を含む。

【 図面の簡単な説明 】

【 0 1 3 5 】

【 図 1 】 本発明での発話区間判定の方式及びそのためのパラメータとを説明するための図である。

20

【 図 2 】 本発明での発話区間処理における状態遷移図である。

【 図 3 】 フレーム長及びフレームシフト量を説明するための図である。

【 図 4 】 本発明の一実施の形態に係る発話区間検出装置の機能的ブロック図である。

【 図 5 】 図 4 に示す装置の、音声エネルギー正規化処理部の機能的ブロック図である。

【 図 6 】 図 4 に示す装置の、初期環境雑音算出部の機能的ブロック図である。

【 図 7 】 フレームエネルギーの変化の例を示す図である。

【 図 8 】 フレームエネルギーを昇順にソートした結果を示す図である。

【 図 9 】 フレームエネルギーのヒストグラムである。

【 図 10 】 図 4 に示す装置の動的しきい値算出部の機能的ブロック図である。

30

【 図 11 】 本発明の一実施の形態における音声エネルギー正規化処理を説明するための図である。

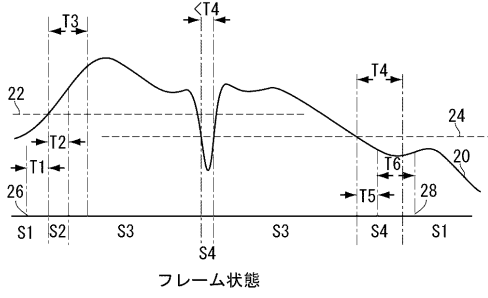
【 符号の説明 】

【 0 1 3 6 】

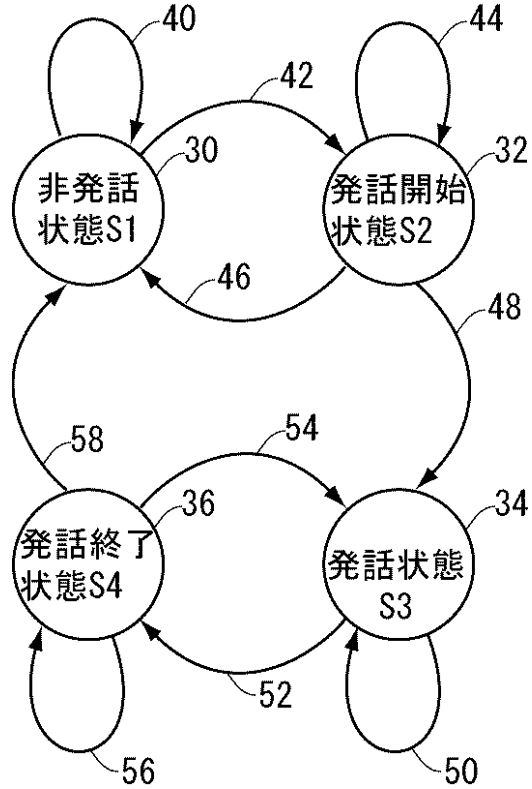
20 音声信号、22 発話開始しきい値、24 発話終了しきい値、30 非発話状態 (S1)、32 発話開始状態 (S2)、34 発話状態 (S3)、36 発話終了状態 (S4)、100 発話区間検出装置、102 マイク、104 音声入力部、106 入力バッファ、108 フレーム情報算出部、110 フレームバッファ、112 初期環境雑音算出部、114 入出力・アドレス管理部、116 動的しきい値算出部、118 状態判定部、120 特徴ベクトル出力部、122 特徴ベクトル、124 フレーム出力信号、126 フレーム音声エネルギー正規化処理部、140 ソート処理部、142 ソート後フレームエネルギー記憶部、144 シーズ算出部、146、174 記憶部、148 第1の平均値算出部、150 フレーム分類部、152 第2の平均値算出部、154 判定部、160 フレームデータ、170 環境雑音算出部、172 選択部、176 最大エネルギー算出部、178 しきい値算出部

40

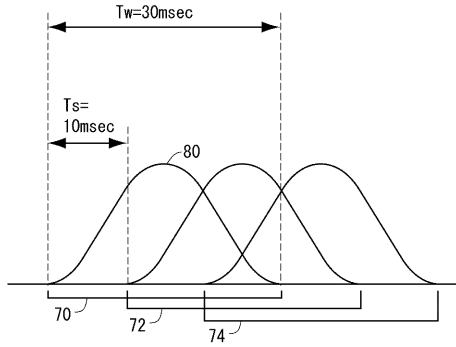
【図1】



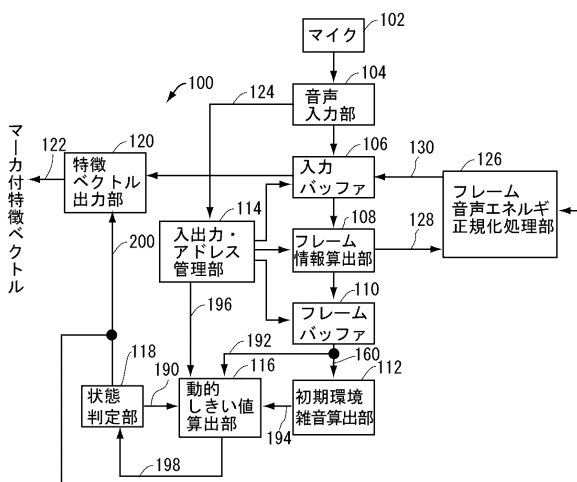
【図2】



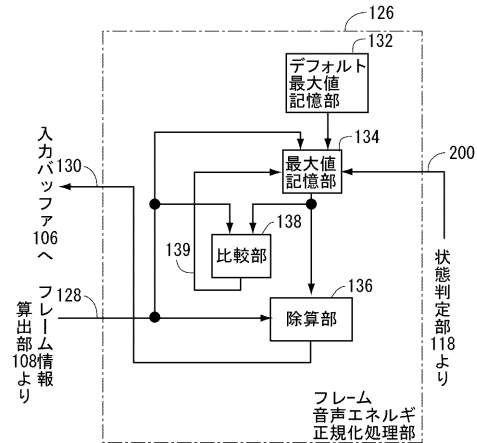
【図3】



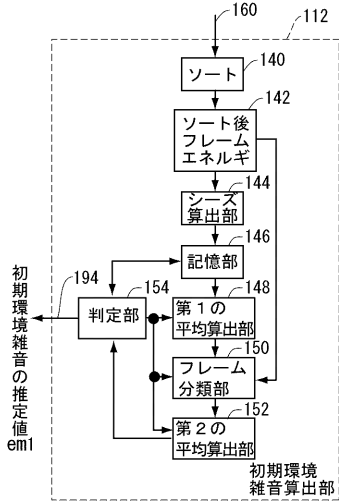
【図4】



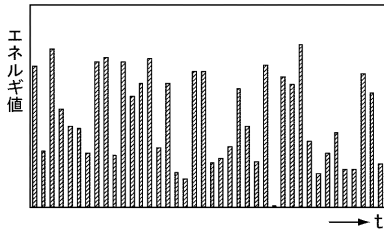
【図5】



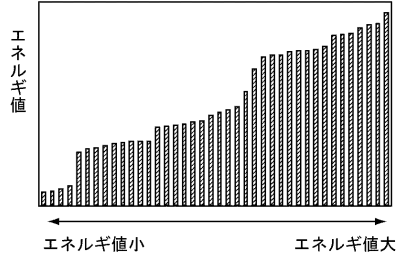
【図 6】



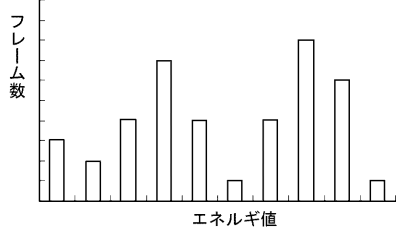
【図 7】



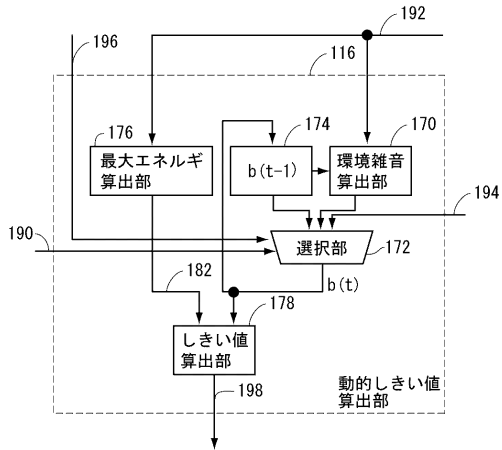
【図 8】



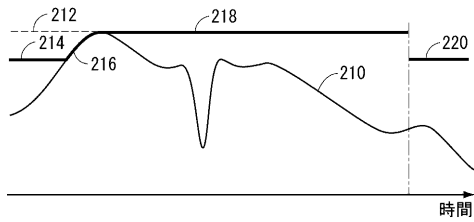
【図 9】



【図 10】



【図 11】



フロントページの続き

(72)発明者 伊藤 玄

京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

審査官 涌井 智則

(56)参考文献 特公平08-023756(JP, B2)

特開2002-258882(JP, A)

特開平08-187368(JP, A)

特開平08-032526(JP, A)

特開平08-314500(JP, A)

特開平10-301593(JP, A)

特開昭61-273596(JP, A)

特開昭58-076899(JP, A)

(58)調査した分野(Int.Cl., DB名)

G10L 11/00-21/06

CiNii