

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5229234号  
(P5229234)

(45) 発行日 平成25年7月3日 (2013.7.3)

(24) 登録日 平成25年3月29日 (2013.3.29)

(51) Int. Cl.

F I

G 1 O L 15/04 (2013.01)

G 1 O L 15/04 3 O O B

G 1 O L 25/78 (2013.01)

G 1 O L 11/02

請求項の数 7 (全 32 頁)

(21) 出願番号	特願2009-546107 (P2009-546107)	(73) 特許権者	000005223
(86) (22) 出願日	平成19年12月18日 (2007.12.18)		富士通株式会社
(86) 国際出願番号	PCT/JP2007/074274		神奈川県川崎市中原区上小田中4丁目1番1号
(87) 国際公開番号	W02009/078093	(74) 代理人	100078868
(87) 国際公開日	平成21年6月25日 (2009.6.25)		弁理士 河野 登夫
審査請求日	平成22年2月16日 (2010.2.16)	(72) 発明者	鷲尾 信之
前置審査			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	早川 昭二
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		審査官	菊池 智紀
			最終頁に続く

(54) 【発明の名称】 非音声区間検出方法及び非音声区間検出装置

(57) 【特許請求の範囲】

【請求項 1】

音を標準化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出方法において、

各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関関数の比の絶対値を導出し、

導出した絶対値が、所定の閾値以上であるか否かを判定し、

前記閾値以上であると判定したフレームが連なる数を計数し、

計数した数が前記閾値に応じて定める所定数以上であるか否かを判定し、

所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する

ことを特徴とする非音声区間検出方法。

【請求項 2】

音を標準化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出方法において、

各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関の比を導出し、

導出した比について、前フレームとの変化量の絶対値を導出し、

10

20

導出した変化量の絶対値が、所定の閾値以下であるか否かを判定し、  
前記閾値以下であると判定したフレームが連なる数を計数し、  
計数した数が前記閾値に応じて定める所定数以上であるか否かを判定し、  
所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する

ことを特徴とする非音声区間検出方法。

【請求項 3】

音を標準化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出装置において、

各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関関数の比の絶対値を導出する導出手段と、

導出した絶対値が、所定の閾値以上であるか否かを判定する判定手段と、

前記閾値以上であると判定したフレームが連なる数を計数する手段と、

計数した数が前記閾値に応じて定める所定数以上であるか否かを判定する手段と、

所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する検出手段と

を備えることを特徴とする非音声区間検出装置。

【請求項 4】

音を標準化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出装置において、

各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関の比を導出する導出手段と、

導出した比について、前フレームとの変化量の絶対値を導出する第 2 の導出手段と、

導出した変化量の絶対値が所定の閾値以下であるか否かを判定する判定手段と、

前記閾値以下であると判定したフレームが連なる数を計数する手段と、

計数した数が前記閾値に応じて定める所定数以上であるか否かを判定する手段と、

所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する検出手段と

を備えることを特徴とする非音声区間検出装置。

【請求項 5】

前記第 2 の導出手段が導出した変化量が、前記閾値より大きい第 2 の閾値を超えるか否かを判定する第 2 の判定手段を備え、

前記検出手段は、前記第 2 の判定手段が第 2 の閾値を超えると判定した場合、該判定が成立するフレームを含めて第 2 の所定数だけ連なるフレームからなる区間を、非音声区間の検出対象から除外するように構成してあることを特徴とする請求項 4 に記載の非音声区間検出装置。

【請求項 6】

前記第 2 の判定手段の判定が成立するフレームが連なる数を計数する手段と、

計数した数が所定数以下であるか否かを判定する手段と、

所定数以下であると判定した場合、該判定が成立するフレーム及び前記第 2 の所定数未満のフレームが連なる区間が、非音声区間に挟まれているときに、前記非音声区間に挟まれた区間を非音声区間として検出する第 2 の検出手段と

を備えることを特徴とする請求項 5 に記載の非音声区間検出装置。

【請求項 7】

前記尺度は、音データの N 次（N は 0 以上の整数）の自己相関関数に対する M 次（M は N と異なる 0 以上の整数）の自己相関関数の比であることを特徴とする請求項 3 乃至 6 の何れかに記載の非音声区間検出装置。

【発明の詳細な説明】

10

20

30

40

50

## 【技術分野】

## 【0001】

本発明は、音を標本化した音データから所定の時間長のフレームを生成し、非音声区間を検出する非音声区間検出方法、該非音声区間検出方法を適用した非音声区間検出装置に関し、特に非音声の特徴を有する物理量と所定の閾値との比較に基づいて、非音声区間を検出する非音声区間検出方法及び非音声区間検出装置に関する。

## 【背景技術】

## 【0002】

カーナビゲーション装置に代表される車載装置に多く用いられる音声認識装置では、一般的には音声区間を検出し、検出した音声区間について算出した音声の特徴量に基づいて、単語列を認識する。特に音声区間の検出を誤った場合、当該区間における音声の認識率が低下するため、音声区間を的確に検出すること、又は非音声区間を検出して音声認識の対象から除外することが重要である。

10

## 【0003】

音声区間の基本的な検出方式として、入力音声のパワーが、その時の推定背景雑音レベルに閾値を加えた基準値を超えた区間を、音声区間として扱うものがある。この場合は、ブザー音のようにパワー変動が大きい雑音、ワイパーの摺動音、及び音声プロンプトのエコー等、何れも非正常性が強い雑音を含む区間を、音声区間として誤検出する可能性が高い。そこで、直近の発声中の最大音声パワー及びその時の音声認識結果より補正係数を導出し、推定背景雑音レベルと併せて、以後の基準値を補正する技術が、特許文献1に開示されている。

20

【特許文献1】特開平7-92989号公報

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0004】

しかしながら、特許文献1に開示されている技術では、発声前後の非音声区間は除外できても、発声がない場合に基準値を補正することができず、雑音のみの区間を音声区間として誤検出することがある問題は解消されない。

## 【0005】

本発明は斯かる事情に鑑みてなされたものであり、音データの周波数スペクトルに偏りを有するフレームが、音声らしからぬ程度に連なる区間、又は周波数スペクトルの偏り、パワー若しくはピッチについての变化に乏しい音データを有するフレームが音声らしからぬ程度に連なる区間を、非音声区間として検出することにより、パワーの大きい雑音若しくは非正常性の強い雑音、又はパワー変動の大きい雑音が発生する環境下においても、発声前後か否かに拘わらず、高精度に非音声区間を検出することが可能な非音声区間検出方法、及び該非音声区間検出方法を適用した非音声区間検出装置を提供することを目的とする。

30

## 【課題を解決するための手段】

## 【0006】

第1の非音声区間検出方法は、音を標本化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出方法において、各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0次の自己相関関数に対する1次の自己相関関数の比の絶対値を導出し、導出した絶対値が、所定の閾値以上であるか否かを判定し、前記閾値以上であると判定したフレームが連なる数を計数し、計数した数が前記閾値に応じて定める所定数以上であるか否かを判定し、所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出することを要件とする。

40

## 【0007】

第2の非音声区間検出方法は、音を標本化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声

50

区間を検出する非音声区間検出方法において、各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関の比を導出し、導出した比について、前フレームとの変化量の絶対値を導出し、導出した変化量の絶対値が、所定の閾値以下であるか否かを判定し、前記閾値以下であると判定したフレームが連なる数を計数し、計数した数が前記閾値に依りて定める所定数以上であるか否かを判定し、所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出することを要件とする。

【0008】

第3の非音声区間検出装置は、音を標本化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出装置において、各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関関数の比の絶対値を導出する導出手段と、導出した絶対値が、所定の閾値以上であるか否かを判定する判定手段と、前記閾値以上であると判定したフレームが連なる数を計数する手段と、計数した数が前記閾値に依りて定める所定数以上であるか否かを判定する手段と、所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する検出手段とを備えることを要件とする。

10

【0009】

第4の非音声区間検出装置は、音を標本化した音データから所定の時間長の複数のフレームを生成し、人が発声した音声に基づく音声データを含まないフレームを有する非音声区間を検出する非音声区間検出装置において、各フレームの音データを周波数軸上の成分に変換したスペクトルについて、0 次の自己相関関数に対する 1 次の自己相関の比を導出する導出手段と、導出した比について、前フレームとの変化量の絶対値を導出する第2の導出手段と、導出した変化量の絶対値が所定の閾値以下であるか否かを判定する判定手段と、前記閾値以下であると判定したフレームが連なる数を計数する手段と、計数した数が前記閾値に依りて定める所定数以上であるか否かを判定する手段と、所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する検出手段とを備えることを要件とする。

20

【0010】

第5の非音声区間検出装置は、第4の装置において、前記第2の導出手段が導出した変化量が、前記閾値より大きい第2の閾値を超えるか否かを判定する第2の判定手段を備え、前記検出手段は、前記第2の判定手段が第2の閾値を超えると判定した場合、該判定が成立するフレームを含めて第2の所定数だけ連なるフレームからなる区間を、非音声区間の検出対象から除外するように構成してあることを要件とする。

30

【0011】

第6の非音声区間検出装置は、第5の装置において、前記第2の判定手段の判定が成立するフレームが連なる数を計数する手段と、計数した数が所定数以下であるか否かを判定する手段と、所定数以下であると判定した場合、該判定が成立するフレーム及び前記第2の所定数未満のフレームが連なる区間が、非音声区間に挟まれているときに、前記非音声区間に挟まれた区間を非音声区間として検出する第2の検出手段とを備えることを要件とする。

40

【0012】

本願の非音声区間検出装置は、前記第2の導出手段による変化量の導出の対象となったフレームを含めて、所定数だけ連なるフレームについて、変化量の最大値を導出する第3の導出手段を備え、前記判定手段は、前記第3の導出手段が導出した最大値を、前記第2の導出手段が導出した変化量として扱うように構成してあることを要件とする。

【0013】

第7の非音声区間検出装置は、第3の装置乃至第6の装置の何れかにおいて、前記尺度は、音データのN次（Nは0以上の整数）の自己相関関数に対するM次（MはNと異なる0以上の整数）の自己相関関数の比であることを要件とする。

50

## 【 0 0 1 4 】

本願の非音声区間検出装置は、前記導出手段が、各フレームについてスペクトルの偏倚を導出した場合、前記各フレームに夫々時系列に前後する複数のフレームについて、スペクトルの偏倚の最大値、最小値、平均値及び中央値の少なくとも一を導出して、導出した値を前記各フレーム夫々についてのスペクトルの偏倚として扱うように構成してあることを要件とする。

## 【 0 0 1 5 】

本願の非音声区間検出装置は、前記判定手段が判定の対象とした全フレームの数に対する、前記判定が成立するフレームの数の割合を算出する手段と、算出した割合が、所定の割合以上であるか否かを判定する手段と、該判定が成立するフレームが連なる数を計数する手段と、計数した数が所定数以上であるか否かを判定する手段と、所定数以上であると判定したときに、前記フレームが連なる区間を非音声区間として検出する第3の検出手段とを備えることを要件とする。

10

## 【 0 0 1 6 】

本願の非音声区間検出装置は、非音声区間として検出されたフレームの音データ、及び前記非音声区間以外のフレームの音データに基づいて、信号対雑音比を導出する手段と、導出した信号対雑音比に基づいて、前記閾値を変更する手段とを備えることを要件とする。

## 【 0 0 1 7 】

本願の非音声区間検出装置は、各フレームの音データについて、ピッチの各周波数成分の強度の最大値を導出する手段と、導出した強度の最大値に基づいて、前記閾値を変更する手段とを備えることを要件とする。

20

## 【 0 0 1 8 】

本願の非音声区間検出装置は、人が発声した音データについて、予め準備された複数の候補閾値に対し、前記判定手段の判定が成立するフレームが連なる個数を夫々集計する手段と、集計した結果に基づいて、複数の候補閾値の中から前記閾値を決定する手段とを備えることを要件とする。

## 【 0 0 1 9 】

本願の非音声区間検出装置は、各フレームの音データのパワーを導出する第4の導出手段と、各フレームの1又は複数の前フレームの音データのパワーに基づいて、夫々のフレームの背景雑音パワーを推定する推定手段と、各フレームについて前記第4の導出手段が導出したパワーが、夫々のフレームについて前記推定手段が推定した背景雑音パワーより、所定の閾値以上大きいかな否かを判定する手段と、前記背景雑音パワーより前記閾値以上大きいと判定したフレームからなる区間を音声区間として検出する第4の検出手段とを備え、前記推定手段は、前記第4の検出手段が検出した音声区間のフレームについて、前フレームの背景雑音パワーを維持するように構成してあり、更に、前記第4の検出手段が検出した音声区間のうち、前記検出手段によって非音声区間として検出されたフレームについて、背景雑音パワーを推定するように構成してあることを要件とする。

30

## 【 0 0 2 0 】

本願の非音声区間検出装置は、各フレームの音データのパワーを導出する第4の導出手段と、各フレームの1又は複数の前フレームの音データのパワーに基づいて、夫々のフレームの背景雑音パワーを推定する推定手段と、各フレームについて前記第4の導出手段が導出したパワーが、夫々のフレームについて前記推定手段が推定した背景雑音パワーより、所定の閾値以上大きいかな否かを判定する手段と、前記背景雑音パワーより前記閾値以上大きいと判定したフレームからなる区間を音声区間として検出する第4の検出手段とを備え、前記推定手段は、前記第4の検出手段が検出した音声区間のフレームについて、前フレームの背景雑音パワーを維持するように構成してあり、更に、前記第4の検出手段が検出した音声区間の全部又は一部が、前記検出手段によって非音声区間として検出された回数を計数する手段と、計数した回数が所定回数以上であるかな否かを判定する手段と、所定回数以上であると判定した場合、該判定が成立した際のフレームの音データのパワーを、

40

50

背景雑音パワーとして更新する手段とを備えることを要件とする。

【 0 0 2 1 】

第 1 の方法及び第 3 の装置では、音データを周波数軸上の成分に変換したスペクトルにおける高周波側又は低周波側への偏りの大きさを示す尺度が所定の閾値以上となるフレームが所定数以上連なる区間を、非音声区間として検出することにより、音データの周波数スペクトルに偏りを有するフレームが音声らしからぬ程度に連なる区間を非音声区間として検出するので、パワーの大きい雑音又は非定常性の強い雑音が発生する環境下においても、高精度に非音声区間を検出することが可能である。

【 0 0 2 2 】

第 2 の方法及び第 4 の装置では、音データの周波数スペクトルにおける高周波側又は低周波側への偏りの大きさを示す尺度、パワー及びピッチの少なくとも一つについて前フレームとの変化量が所定の閾値以下となるフレームが、所定数以上連なる区間を非音声区間として検出することにより、周波数スペクトルにおける前記尺度、パワー若しくはピッチについての变化に乏しい音データを有するフレームが音声らしからぬ程度に連なる区間を非音声区間として検出するので、パワー変動の大きい雑音が発生する環境下においても、高精度に非音声区間を検出することが可能である。

【 0 0 2 3 】

第 5 の装置では、導出した指標の前フレームとの変化量が前記閾値より大きい第 2 の閾値を超えるフレームを含めて第 2 の所定数だけ連なるフレームからなる区間を、非音声区間として検出することがないので、音声データを含む可能性のあるフレームからなる区間を、非音声区間として誤検出することを防止することが可能である。

【 0 0 2 4 】

第 6 の装置では、導出した指標の前フレームとの変化量が第 2 の閾値を超えて所定数以下だけ連なるフレーム及び第 2 の所定数以下のフレームからなる区間が、非音声区間に挟まれている場合に、その挟まれた区間を非音声区間として検出することにより、音データの単発的な変化が発生した場合であっても、高精度に非音声区間を検出することが可能である。

【 0 0 2 5 】

本願の装置では、連なる所定数のフレームについて、夫々導出した指標の前フレームとの変化量の最大値を、一のフレームについての前フレームとの変化量として扱うことにより、各フレームの指標について当初導出した前フレームとの変化量が近傍のフレームについての当該変化量の最大値と置き換わるので、音声データを含む可能性のあるフレームからなる区間を、非音声区間として誤検出することを抑止することが可能である。

【 0 0 2 6 】

第 7 の装置では、音データの自己相関関数の N 次の値に対する M 次の値の比が、音データのスペクトルの包絡を近似する指標であるので、これをスペクトルにおける高周波側又は低周波側への偏りの大きさを示す尺度とすることにより、音データの周波数スペクトルの偏りが的確に把握されて、高精度に非音声区間を検出することが可能である。

【 0 0 2 7 】

本願の装置では、前後する所定数のフレームについて、夫々導出したスペクトルの偏倚の最大値、最小値、平均値及び中央値の少なくとも一を、一のフレームについてのスペクトルの偏倚として扱うことにより、スペクトルの偏倚が短時間に変化した場合であっても、高精度に非音声区間を検出することが可能である。

【 0 0 2 8 】

本願の装置では、音データの周波数スペクトルの偏倚が正の値（又は負の値）の場合、所定の閾値以上（又は所定の閾値以下）となるフレーム、又は導出した指標の前フレームとの変化量が前記閾値と異なる他の閾値以下となるフレームが、所定の割合以上で所定数以上連なる区間を、非音声区間として検出することにより、音データの周波数スペクトルの偏倚、又は導出した指標の前フレームとの変化量が、短時間に変動する場合にも、高精度に非音声区間を検出することが可能である。

10

20

30

40

50

## 【 0 0 2 9 】

本願の装置では、検出した非音声区間の音データ及び非音声区間以外の音データより導出した信号対雑音比に基づいて、前記閾値を変更することにより、例えば信号対雑音比が低下して、スペクトルの偏倚又は導出した指標の前フレームとの変化量が変動した場合に、前記閾値を適切に調整して、非音声区間の誤検出を抑止することができ、高精度に非音声区間を検出することが可能である。

## 【 0 0 3 0 】

本願の装置では、ピッチの各周波数成分の強度についての最大値に基づいて、前記閾値を調整することにより、ピッチが明瞭に現れる度合いに応じて前記閾値を適切に調整することができるので、高精度に非音声区間を検出することが可能である。

10

## 【 0 0 3 1 】

本願の装置では、予め準備した複数の候補閾値を所定の音声データに適用し、夫々の閾値以上（又は閾値以下）となるフレームが連なる個数を集計した結果に基づいて、前記閾値を決定することにより、事前の学習に基づいて前記閾値を決定することができるので、高精度に非音声区間を検出することが可能である。

## 【 0 0 3 2 】

本願の装置では、非音声区間のフレームの音データのパワーに基づいて推定した背景雑音パワーより、所定の閾値以上大きいパワーを有するフレームからなる区間を音声区間として検出し、検出した音声区間のうち、非音声区間として検出されたフレームについて、背景雑音パワーを推定するので、音データのパワーに基づいて音声検出した結果を適正に修正することが可能である。

20

## 【 0 0 3 3 】

本願の装置では、非音声区間のフレームの音データのパワーに基づいて推定した背景雑音パワーより、所定の閾値以上大きいパワーを有するフレームからなる区間を音声区間として検出し、検出した音声区間の全部又は一部が、所定回数だけ非音声区間として検出された際のフレームの音データのパワーを、背景雑音パワーとして更新するので、背景雑音パワーの推定値が上がり過ぎて、音声区間が検出できなくなることを抑止することができる。

## 【 発明の効果 】

## 【 0 0 3 4 】

開示の非音声区間検出方法、及び非音声区間検出装置は、各フレームの音データを周波数軸上の成分に変換したスペクトルにおける高周波側又は低周波側への偏りの大きさを示す尺度が所定の閾値以上であるかを判定し、前記閾値以上と判定したフレームが連なる数が所定数以上かを判定し、そして所定数以上と判定したフレームが連なる区間を非音声区間として検出する。

30

## 【 0 0 3 5 】

この構成により、開示の方法及び装置では、スペクトルの偏りに係る閾値とフレームが連なる数に係る閾値とを組み合わせ、非音声の特徴を有するフレームが音声らしからぬ程度に連なる区間を非音声区間として検出し、人の発声による基準値の補正を要しない。従って、パワーの大きい雑音、又は非定常性の強い雑音が発生する環境下においても、発声前後か否かに拘わらず、高精度に非音声区間を検出することが可能である等、優れた効果を奏する。

40

## 【 0 0 3 6 】

また、開示の非音声区間検出方法、及び非音声区間検出装置は、各フレームの音データを周波数軸上の成分に変換したスペクトルのにおける高周波側又は低周波側への偏りの大きさを示す尺度を少なくとも用いて、前フレームとの変化量が所定の閾値以下であるかを判定し、前記閾値以下と判定したフレームが連なる数が所定数以上かを判定し、そして所定数以上と判定したフレームが連なる区間を非音声区間として検出する。

## 【 0 0 3 7 】

この構成により、開示の方法及び装置では、周波数スペクトルの偏り、パワー若しくは

50

ピッチについての変化に係る閾値とフレームが連なる数に係る閾値とを組み合わせ、非音声の特徴を有するフレームが音声らしからぬ程度に連なる区間を非音声区間として検出し、人の発声による基準値の補正を要しない。従って、パワー変動の大きい雑音が発生する環境下においても、発声前後か否かに拘わらず、高精度に非音声区間を検出することが可能である等、優れた効果を奏する。

【図面の簡単な説明】

【0038】

【図1】本発明の実施の形態1に係る非音声区間検出装置の一実施例である音声認識装置の構成例を示すブロック図である。

【図2】制御手段の音声認識に係る処理構成例を示すブロック図である。

10

【図3】制御手段の音声認識処理の一例を示すフローチャートである。

【図4】非音声区間検出のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図5】鼻をすする音について、パワー及び高域・低域強度等のデータを示す図である。

【図6】踏切の警報音について、パワー及び高域・低域強度等のデータを示す図である。

【図7】発声音（「えーテスト中です」）について、パワー及び高域・低域強度等のデータを示す図である。

【図8】発声音（「経営（けーえー）」）について、パワー及び高域・低域強度等のデータを示す図である。

【図9】本発明の実施の形態2に係る非音声区間検出装置の一実施例である音声認識装置について、制御手段の音声認識に係る処理構成例を示すブロック図である。

20

【図10】本発明の実施の形態3に係る非音声区間検出装置の一実施例である音声認識装置について、制御手段の音声認識に係る処理構成例を示すブロック図である。

【図11】制御手段の音声認識処理の一例を示すフローチャートである。

【図12】非音声区間検出のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図13】非音声区間検出除外のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図14】非音声区間検出除外のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

30

【図15】非音声区間検出確定のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図16】非音声区間検出確定のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図17】本発明の実施の形態4に係る非音声検出装置の一実施例である音声認識装置について、非音声区間検出のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図18】本発明の実施の形態4に係る非音声検出装置の一実施例である音声認識装置について、非音声区間検出のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

40

【図19】本発明の実施の形態5に係る非音声検出装置の一実施例である音声認識装置について、制御手段の音声認識処理の一例を示すフローチャートである。

【図20】本発明の実施の形態6に係る非音声検出装置の一実施例である音声認識装置について、非音声区間検出のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図21】本発明の実施の形態6に係る非音声検出装置の一実施例である音声認識装置について、非音声区間検出のサブルーチンに係る制御手段の処理手順を示すフローチャートである。

【図22】本発明の実施の形態7に係る非音声検出装置の一実施例である音声認識装置について、制御手段の音声認識処理の一例を示すフローチャートである。

50



## 【符号の説明】

## 【 0 0 3 9 】

- 1 音声認識装置
- 2 制御手段（第 3 の導出手段、第 3 の検出手段）
- 3 記録手段
- 4 記憶手段
- 5 音取得手段
- 2 0 フレーム生成部
- 2 1 スペクトルの偏倚導出部（導出手段）
- 2 1 a スペクトルの偏倚 / パワー / ピッチ導出部（導出手段） 10
- 2 1 b 変化量導出部（第 2 の導出手段）
- 2 2 非音声区間検出部（判定手段、検出手段）
- 2 2 a 非音声区間検出部（判定手段、検出手段）
- 2 2 b 非音声区間検出部（判定手段、検出手段、第 2 の判定手段、第 2 の検出手段）

## 【発明を実施するための最良の形態】

## 【 0 0 4 0 】

以下、本発明をその実施の形態を示す図面に基づいて詳述する。

## 実施の形態 1

図 1 は、本発明の実施の形態 1 に係る非音声区間検出装置の一実施例である音声認識装置の構成例を示すブロック図である。図中 1 は、例えば車両に搭載されるナビゲーション装置のようなコンピュータを用いた音声認識装置であり、音声認識装置 1 は、装置全体を制御する CPU（Central Processing Unit）及び DSP（Digital Signal Processor）等の制御手段 2 と、プログラム及びデータ等の各種情報を記録するハードディスク及び ROM 等の記録手段 3 と、一時的に発生するデータを記録する RAM からなる記憶手段 4 と、外部から音を取得するマイクロホンからなる音取得手段 5 と、音を出力するスピーカからなる音出力手段 6 と、液晶モニタからなる表示手段 7 と、目的地までの経路指示のようなナビゲーションに係る処理を実行するナビゲーション手段 8 とを備えている。

## 【 0 0 4 1 】

記録手段 3 には、本発明に係る非音声区間検出方法を実行するコンピュータプログラム 3 0 が記録されており、記録されているコンピュータプログラム 3 0 に含まれる各種手順を記録手段 3 に記憶して制御手段 2 の制御にて実行することにより、コンピュータは、本発明の非音声区間検出装置としても動作する。

## 【 0 0 4 2 】

また、記録手段 3 の記録領域の一部は、音声認識用の音響モデルを記録している音響モデルデータベース（音響モデル DB）3 1、音響モデルに対応する音素又は音節定義で表記された認識語彙及び文法を記録している認識辞書 3 2 等の各種データベースとして用いられている。

## 【 0 0 4 3 】

記憶手段 4 の記憶領域の一部は、音取得手段 5 が取得したアナログ信号である音を所定の周期で標本化（サンプリング）してデジタル化した音データを記録する音データバッファ 4 1、及び音データを所定の時間長に区分したフレームから抽出した特徴量をはじめとするデータを記憶するフレームバッファ 4 2、及び一時的に発生した情報を記憶するワークメモリ 4 3 として用いられる。

## 【 0 0 4 4 】

ナビゲーション手段 8 は、GPS（Global Positioning System）のような位置検出機構と、地図情報を記録する DVD（Digital Versatile Disk）及びハードディスク等の記録媒体とを有し、現在地から目的地までの経路検索及び経路指示等のナビゲーション処理を実行し、地図及び経路を表示手段 7 に表示し、音声による案内を音出力手段 6 から出力する。

## 【 0 0 4 5 】

10

20

30

40

50

尚、図 1 に示した構成例はあくまでも一例であり、様々な形態に展開することが可能である。例えば、音声認識に係る機能を一又は複数の V L S I チップとして構成し、ナビゲーション装置に組み込むことも可能であり、音声認識用の専用装置をナビゲーション装置に外付けすることも可能である。また、制御手段 2 を音声認識及びナビゲーションの双方の処理で共用するようにしても、夫々専用の回路を設けるようにしてもよく、更には音声認識に関する特定の演算、例えば後述する F F T (Fast Fourier Transform)、D C T (Discrete Cosine Transform) 及び I D C T (Inverse Discrete Cosine Transform) 等の処理を実行するコプロセッサを制御手段 2 に組み込んでよい。また、音データバッファ 4 1 を音取得手段 5 の付属回路とし、フレームバッファ 4 2 及びワークメモリ 4 3 を制御手段 2 が備えるメモリ上に構成するようにしてもよい。更に、本発明の音声認識装置 1 は、ナビゲーション装置のような車載装置に限らず、音声認識を行う様々な用途の装置に用いることが可能である。

10

#### 【 0 0 4 6 】

次に本発明の実施の形態 1 に係る非音声区間検出装置の一実施例である音声認識装置 1 の処理について説明する。図 2 は、制御手段 2 の音声認識に係る処理構成例を示すブロック図である。また、図 3 は、制御手段 2 の音声認識処理の一例を示すフローチャートである。

制御手段 2 は、音データからフレームを生成するフレーム生成部 2 0、生成されたフレームについてスペクトルの偏倚を導出するスペクトルの偏倚導出部 2 1、導出されたスペクトルの偏倚に基づく判定基準を用いて非音声区間を検出する非音声区間検出部 2 2、検出された非音声区間をもとに音声区間の開始 / 終了を確定させる音声区間判定部 2 3、及び判定された音声区間について音声を認識する音声認識部 2 4 を備えている。

20

#### 【 0 0 4 7 】

制御手段 2 は、音取得手段 5 によって外部の音をアナログ信号として取得し (ステップ S 1 1)、取得した音を所定の周期で標本化してデジタル化した音データを、音データバッファ 4 1 に記録する (ステップ S 1 2)。ステップ S 1 1 にて取得する外部の音とは、人が発声する音声、定常雑音及び非定常雑音等の様々な音が重畳された音である。人が発声する音声は、音声認識装置 1 による認識の対象となる音声である。定常雑音は、ロードノイズ及びエンジン音等の雑音であり、既に提案及び確立されている様々な除去方法が適用される。非定常雑音としては、車両に配設されたハザード、ウインカーのようなリレー音、及びワイパーの摺動音のような機構による雑音を例示することができる。

30

#### 【 0 0 4 8 】

そして制御手段 2 のフレーム生成部 2 0 は、音データバッファ 4 1 に記憶した音データより、10 msec のフレーム長で 5 msec ずつオーバーラップさせたフレームを生成し (ステップ S 1 3)、生成したフレームをフレームバッファ 4 2 に記憶させる (ステップ S 1 4)。尚、フレーム生成部 2 0 は、音声認識の分野における一般的なフレーム処理として、フレーム分割前のデータに対して高域強調フィルタリング処理を施した後に、フレームに分割する。このようにして生成された各フレームに対し、以下の処理が行われる。

#### 【 0 0 4 9 】

スペクトルの偏倚導出部 2 1 は、フレーム生成部 2 0 からフレームバッファ 4 2 を介して与えられたフレームについて、後述するスペクトルの偏倚を導出し (ステップ S 1 5)、導出したスペクトルの偏倚をフレームバッファ 4 2 に書き込む。この場合、書き込まれたフレーム及びスペクトルの偏倚を夫々参照するのに用いられるフレームバッファ 4 2 へのポインタ (アドレス) が、ワークメモリ 4 3 上に設けてあり、前記ポインタを介して、フレームバッファ 4 2 に記憶したスペクトルの偏倚にアクセスする。

40

尚、スペクトルの偏倚を導出する前に、ノイズキャンセル処理及びスペクトルサブトラクション処理を行って、雑音の影響を除外してもよい。

#### 【 0 0 5 0 】

非音声区間検出部 2 2 は、フレームバッファ 4 2 を介してスペクトルの偏倚導出部 2 1 より与えられたフレームについて、スペクトルの偏倚に基づく判定基準により非音声区間

50

を検出するサブルーチンを呼び出す（ステップS 1 6）。非音声区間検出部 2 2 が判定基準を用いて検出した非音声区間のフレームは、フレームバッファ 4 2 を介して順次音声区間判定部 2 3 に与えられる。判定結果が未確定のフレーム、即ち後続するフレームによっては非音声区間になり得るフレームは、判定基準が用い尽くされるまで、非音声区間検出部 2 2 によって保留される。

【 0 0 5 1 】

音声区間判定部 2 3 は、非音声区間検出部 2 2 が非音声区間として検出できなかった区間を音声区間とみなし、音声区間長が既定の最短音声区間長  $L_1$  を超えた場合に音声区間開始と判定して、音声区間開始フレームを確定させる。そして音声区間が途切れたフレームを、音声区間終了点候補とする。その後、既定の最大ポーズ長  $L_2$  を超えるまでに次の音声区間が始まった場合は、前述の音声区間終了点候補を棄却して、再び音声区間が途切れるのを待つ。

既定の最大ポーズ長  $L_2$  を超えても次の音声区間が始まらなかった場合、音声区間判定部 2 3 は、音声区間終了候補を音声区間終了フレームとして確定させる。音声区間の開始 / 終了フレームを確定したことにより、音声区間判定部 2 3 は、一つの音声区間の判定を終える（ステップS 1 7）。このようにして検出された音声区間は、フレームバッファ 4 2 を介して音声認識部 2 4 に与えられる。

尚、音声区間の検出誤りを回避するため、音声区間判定部 2 3 が判定した音声区間よりも、例えば前後に 1 0 0 msec だけ広い区間を、確定させた音声区間としてもよい。

【 0 0 5 2 】

音声認識部 2 4 は、音声認識の分野で一般的な技術を用いて、音声区間のフレームのデジタル信号から特徴ベクトルを抽出し、抽出した特徴ベクトルに基づいて、音響モデルデータベース 3 1 に記録している音響モデル並びに認識辞書 3 2 に記憶している音響語彙及び文法を参照し、入力されたフレームバッファ 4 2 の最後（音声区間の最後）まで、音声認識処理を実行する（ステップS 1 8）。

【 0 0 5 3 】

図 3 は、一音声区間が確定した場合に、音声認識処理を実行して終了する構成であるが、音声区間を検出した場合に、計算可能なフレームから音声認識処理を実行してレスポンスタイムを短縮する構成、又は一定時間について、音声区間が検出できない場合に、処理を終了する構成としてもよい。

【 0 0 5 4 】

ここで、図 3 を用いて説明したステップS 1 5 におけるスペクトルの偏倚について、更に詳述する。

本実施の例では、音データの各フレームにおけるスペクトルの傾き、即ち、スペクトルの高域 / 低域での偏りを示す尺度として高域・低域強度を定義する。高域・低域強度は、そのままスペクトルの偏倚として用いることができるが、本実施の例では、スペクトルの偏倚を、高域・低域強度の絶対値で表すものとする。高域・低域強度は、スペクトル包絡を近似する指標であって、音データのパワーを示す 0 次の自己相関関数に対する、遅れ時間が 1 サンプルの 1 次の自己相関関数の比で表すことができる。

自己相関関数は、音データを分析単位である 1 フレーム毎（例えば、フレーム幅： $N = 256$  サンプル）に抽出し、ハミング窓をかけた音データの波形  $\{x(n)\}$  から、短時間自己相関関数  $\{c(\ )\}$  として、下記の式 1 より算出することができる。

【 0 0 5 5 】

10

20

30

40

【数 1】

$$c(\tau) = \frac{1}{N-1} \sum_{n=0}^{N-2} x(n)x(n+\tau) \quad , \quad 0 \leq \tau \leq 1 \quad \cdots \text{式 1}$$

【0056】

また、0 次及び 1 次の自己相関関数の比を用いるので、夫々について共通の係数である  $1/(N-1)$  を除いて、下記の式 2 としてもよい。 10

【0057】

【数 2】

$$c(\tau) = \sum_{n=0}^{N-2} x(n)x(n+\tau) \quad , \quad 0 \leq \tau \leq 1 \quad \cdots \text{式 2}$$

【0058】

20

また、自己相関関数  $c(\quad)$  は、Wiener-Khinchine の定理により、短時間スペクトル  $S(\quad)$  を逆フーリエ変換 ( I D F T : Inverse Discrete Fourier Transform ) して算出することもできる。短時間スペクトル  $S(\quad)$  は、音データを分析単位である 1 フレーム毎 ( 例えば、フレーム幅 :  $N = 256$  サンプル ) に抽出し、各フレームに対してハミング窓をかけ、窓かけ後のフレームのデータに対して D F T ( Discrete Fourier Transform ) を行うことで算出できる。

尚、算出に伴う処理量を削減するため、I D F T / D F T に替えて I D C T / D C T を用いることができる。

【0059】

上述のようにして求めた自己相関関数  $c(\quad)$  について、0 次及び 1 次の比を用いて、高域・低域強度  $A$  を下記の式 3 及び式 4 のとおり定義する。 30

【0060】

$$A = c(1) / c(0) \quad ( c(0) \neq 0 ) \quad \cdots \text{式 3}$$

$$A = 0 \quad ( c(0) = 0 ) \quad \cdots \text{式 4}$$

【0061】

この場合、 $A$  は、 $-1 \leq A \leq 1$  の範囲の値をとり、1 ( 又は -1 ) に近い値であるほどスペクトルの低域 ( 又は高域 ) の強度が大きいことを示す。

尚、高域・低域強度としては、上述した  $A$  に限定されるものではなく、0 次及び 1 次以外の異なる次数についての自己相関関数の比、所定周波数帯域のパワー、所定の異なる周波数帯域についてのパワーの比、M F C C、対数スペクトラムを逆フーリエ変換したケプストラム、又は推定したフォルマントのうち所定の異なるフォルマントについての周波数の比若しくはパワーの比の少なくとも一であってもよい。複数の高域・低域強度を導出した場合は、夫々導出した値に基づいて、非音声区間の判定を並列的に実行することができる。 40

【0062】

図 5 乃至 8 は、夫々鼻をすする音、踏切の警報音及び 2 種類の発声音 ( 「えーテスト中です」、「経営 ( けーえー ) 」 ) について、パワー及び高域・低域強度等のデータを示す図である。図 5 乃至 8 の各図において、横軸は時間であり、縦軸は、上から音データの波形、音データのパワー ( 鎖線、左軸 )、高域・低域強度  $A$  ( 実線、右軸 ) 及びスペクトログラム ( 左軸 ) である。 50

## 【 0 0 6 3 】

図 5 では、スペクトログラムにおいて、黒の濃い領域が高域である上方に偏っているため、当該区間で A の値は - 1 に近づいている。

図 6 では、警報のトーン信号により、スペクトログラムの下半分に黒の濃い線が出現して、低域に偏っているため。A の値は 1 に近づいている。

## 【 0 0 6 4 】

図 7 では、発声されている音素によって、高域 / 低域が強い、又はどちらでもない、という区間が出現しており、A の値は概ね  $-0.7 < A < 0.7$  の範囲で大きく変動している。即ち、発声中の区間では、A の値は長時間特定の値に留まることがなく、ある程度の範囲で変動するといえる。発声中であっても A の値が安定するのは、図 7 の発声末尾の「す」のように、同じ音素が継続している場合である。この場合、「す」が無声化して、高域が強い摩擦音 /s/ が継続しているため、A の値は - 1 に近い - 0.7 近辺で約 0.3 秒間に渡り安定している。また、同じように 1 音素が継続する区間であっても、発声される音素によって A の値は変動する。例えば、図 7 では、「テスト中」末尾の「う」近辺で、母音 /u/ が継続しているが、A の値はプラス方向に振れ、0.6 前後の値をとっている。

## 【 0 0 6 5 】

一方、日本語の語彙においては、特定の母音 / 子音が無意味に連なることはないため、一般的な音声認識処理では、一つの音素が長時間発声されることは考慮する必要がない。このため、一般の単語又は文の発声において各音素が継続され得る時間長と、各音素の発声において A の値が取り得る範囲とを想定することにより、音素が想定外に継続した場合、又は A の値が想定外となった場合は、当該単語又は文は音声でないと見做すことができる。例えば、図 8 では、「経営」を「けーえー」と発声する場合があります、最初の /k/ 以外は、/e/ が約 4 モーラ長だけ継続する。この場合は、日本語において同一の音素が最も長時間継続する場合と想定され、その継続時間は、ゆっくりと発声された場合であっても高々 1.2 秒程度である。

## 【 0 0 6 6 】

上述した内容及び図 5 乃至 8 に示された事項より、スペクトルの偏倚 |A| について、例えば音声区間では、|A| 0.7 とはならないこと、また、音素は高々 1.2 秒しか継続せず、当該区間で |A| 0.5 とはならないことがいえるため、非音声区間について、例えば下記のような判定を行うことが可能である。

(a) : |A| 0.7 が 0.1 秒以上継続する場合、当該区間は非音声とする。

(b) : |A| 0.5 が 1.2 秒以上継続する場合、当該区間は非音声とする。

また、上記の判定を更に細分化して、以下のような判定を行うことも可能である。

(c) : |A| 0.6 が 0.5 秒以上継続する場合、当該区間は非音声とする。

尚、フレームが継続する時間に係る閾値は、フレーム長が一定であるため、フレームが継続する数に係る閾値に置き換えることができる。また、音取得手段 5 のマイクロホンの特性を含む音入力系の伝達特性によっては、高域・低域のバランスが変動してスペクトルの偏倚 |A| も変化することが想定されるため、入力系の伝達特性に応じて上述した判定の閾値を調整することが望ましい。

## 【 0 0 6 7 】

上述した内容を踏まえて、非音声区間検出のサブルーチンについて説明する。図 4 は、非音声区間検出のサブルーチンに係る制御手段 2 の処理手順を示すフローチャートである。非音声区間検出のサブルーチンが呼び出された場合、制御手段 2 は、そのときのポイントが示すフレームのスペクトルの偏倚が、所定の閾値（例えば上述した 0.7）以上であるか否かを判定する（ステップ S 2 1）。所定の閾値未満であると判定した場合（ステップ S 2 1 : NO）、制御手段 2 は、ワークメモリ 4 3 に記憶されたフレームバッファ 4 2 へのポイントを 1 フレーム後方に更新して（ステップ S 2 2）、リターンする。

これにより、制御手段 2 は、非音声区間を検出することなくリターンする。

## 【 0 0 6 8 】

所定の閾値以上であると判定した場合（ステップ S 2 1 : YES）、制御手段 2 は、そ

のときのポインタが示すフレームのフレーム番号を「開始フレーム番号」としてワークメモリ43上に記憶する(ステップS23)。そして、制御手段2は、ワークメモリ43上に設けた「フレームカウント」の記憶値を「1」に初期化する(ステップS24)。ここで、「フレームカウント」は、スペクトルの偏倚と所定の閾値との比較判定を行ったフレーム数を計数するものである。

【0069】

その後、制御手段2は、「フレームカウント」の記憶内容が所定数(例えば上述した0.1秒間に含まれるフレームの数である10)以上であるか否かを判定し(ステップS25)、所定数未満であると判定した場合(ステップS25:NO)、制御手段2は、「フレームカウント」の記憶内容に「1」を加算すると共に(ステップS26)、フレームバッファへのポインタを1フレーム後方に更新する(ステップS27)。そして、制御手段2は、そのときのポインタが示すフレームのスペクトルの偏倚が、所定の閾値以上であるか否かを判定する(ステップS28)。

10

【0070】

スペクトルの偏倚が所定の閾値以上であると判定した場合(ステップS28:YES)、制御手段2は、処理をステップS25に戻す。

スペクトルの偏倚が所定の閾値未満であると判定した場合(ステップS28:NO)、制御手段2は、「開始フレーム番号」の内容を消去して(ステップS29)、リターンする。

これにより、制御手段2は、非音声区間を検出することなくリターンする。

20

【0071】

ステップS25で「フレームカウント」の記憶内容が所定数以上であると判定した場合(ステップS25:YES)、制御手段2は、非音声区間の終了フレームを検出する処理に移り、フレームバッファへのポインタを1フレーム後方に更新する(ステップS30)。そして、制御手段2は、そのときのポインタが示すフレームのスペクトルの偏倚が、所定の閾値以上であるか否かを判定する(ステップS31)。

【0072】

スペクトルの偏倚が所定の閾値以上であると判定した場合(ステップS31:YES)、制御手段2は、処理をステップS30に戻す。スペクトルの偏倚が所定の閾値未満であると判定した場合(ステップS31:NO)、制御手段2は、そのときのポインタが示すフレームの1つ前のフレーム番号を「終了フレーム番号」としてワークメモリ43上に記憶し(ステップS32)、リターンする。

30

これにより、「開始フレーム番号」及び「終了フレーム番号」で区切られた区間が、検出された非音声区間となる。

【0073】

このように、本発明の実施の形態1では、各フレームの音データより導出したスペクトルの偏倚 $|A|$ が、例えば0.7以上となるフレームが、継続時間にして0.1秒に相当する数以上連なる場合、スペクトルの偏倚が最初に0.7以上となったフレームから、最後に0.7以上となったフレームまでを非音声区間として検出する。

これにより、本実施の形態1では、スペクトルの偏倚が大きくて非音声の特徴を有するフレームが、音声らしからぬ程度まで連なる区間を非音声区間として検出し、人の発声による基準値の補正を要しない。従って、パワーの大きい雑音、又は非定常性の強い雑音が発生する環境下においても、発声前後か否かに拘わらず、高精度に非音声区間を検出することが可能である。

40

【0074】

実施の形態2

実施の形態2は、推定背景雑音パワーを基本とした音声区間検出装置と、実施の形態1に係る非音声区間検出装置とを併用した形態である。

図9は、本発明の実施の形態2に係る非音声区間検出装置の一実施例である音声認識装置1について、制御手段2の音声認識に係る処理構成例を示すブロック図である。

50

## 【 0 0 7 5 】

制御手段 2 は、フレーム生成部 2 0、スペクトルの偏倚導出部 2 1、導出されたスペクトルの偏倚に基づく判定基準を用いて非音声区間を検出する非音声区間検出部 2 2 a、検出された非音声区間をもとに音声区間の開始 / 終了を確定させる音声区間判定部 2 3 a、確定された音声区間について音声認識の照合に用いる特徴量を算出する特徴量算出部 2 8、及び算出された特徴量を用いて音声認識のための照合処理を行う照合部 2 9 を備えている。

制御手段 2 は、更に、フレーム生成部 2 0 で生成されたフレームについて、音データのパワーを導出するパワー導出部 2 6、導出したパワーに基づいて背景雑音パワーを推定する背景雑音パワー推定部 2 7、及び音声区間判定部 2 3 a に修正すべきフレーム番号を通知する音声区間修正部 2 5 を備える。

10

## 【 0 0 7 6 】

非音声区間検出部 2 2 a は、検出した非音声区間のフレーム番号を音声区間判定部 2 3 a 及び音声区間修正部 2 5 に与える。

音声区間修正部 2 5 は、非音声区間検出部 2 2 a が非音声区間として検出したフレームが、音声区間判定部 2 3 a では音声区間と判定されていた場合に、音声区間判定部 2 3 a に対して、所定の修正信号及び修正すべきフレーム番号を与える。

## 【 0 0 7 7 】

パワー導出部 2 6 は、フレーム生成部 2 0 から与えられた各フレームについて音データのパワーを導出し、導出したパワーを背景雑音パワー推定部 2 7 に与える。

20

尚、パワーを算出する前に、ノイズキャンセル処理及びスペクトルサブトラクション処理を行って、雑音の影響を除外してもよい。

## 【 0 0 7 8 】

背景雑音パワー推定部 2 7 は、音データの先頭フレームを無条件に雑音とみなし、当該フレームの音データのパワーを推定背景雑音パワーの初期値とする。その後、背景雑音パワー推定部 2 7 は、音声区間判定部 2 3 a から通知された音声区間のフレームを除いて、音データの 2 フレーム目以降について、直近の 2 フレームのパワーの単純移動平均をとり、導出した移動平均値によって推定背景雑音パワーをフレーム毎に更新する。尚、推定背景雑音パワーの更新値を、パワーの単純移動平均から導出するのではなく、IIR ( Infinite Impulse Response ) フィルタによって導出するようにしてもよい。

30

また、背景雑音パワー推定部 2 7 は、音声区間判定部 2 3 a より後述する推定背景雑音パワーの修正を通知された場合、非音声区間に修正されたフレームのうち、その時の最新のフレームの音データから導出されたパワーにより、推定背景雑音パワーを上書きして修正する。

## 【 0 0 7 9 】

尚、背景雑音パワー推定部 2 7 は、音声区間判定部 2 3 a より推定背景雑音パワーの修正を通知された場合、非音声区間に修正されたフレームの音データについて、推定背景雑音パワーを導出するようにしてもよい。また、所定の N 回目 ( N は 2 以上の自然数 ) の修正を通知された場合に初めて、その時の最新のフレームの音データから導出されたパワーにより、推定背景雑音パワーを上書きするようにしてもよい。これにより、背景雑音レベルが上下に変動した場合に、推定背景雑音レベルが上がり過ぎて音声区間が検出できなくなるのを防止することができる。

40

## 【 0 0 8 0 】

音声区間判定部 2 3 a は、各フレームの音データのパワーが、「推定背景雑音パワー + 所定の閾値」以上となった場合、当該フレームを音声区間と判定する。また、音声区間判定部 2 3 a は、音声区間修正部 2 5 より上述した所定の修正信号を与えられた場合、修正すべきフレーム番号に基づいて、音声区間の判定結果を修正する。そして、音声区間判定部 2 3 a は、判定した音声区間が最短入力時間長以上、且つ最長入力時間長以下だけ継続した場合、その時の音声区間を確定させ、確定させた音声区間を特徴量算出部 2 8、照合部 2 9 及び背景雑音パワー推定部 2 7 に通知する。

50

更に、音声区間判定部 23 a は、背景雑音パワー推定部 27 に対し、非音声区間に修正されたフレームの音データにより、推定背景雑音パワーを修正するように通知する。

【0081】

特徴量算出部 28 は、音声区間判定部 23 a が最終的に音声区間と確定させた区間について、音声認識の照合に用いる特徴量を算出する。ここでの特徴量とは、例えば音響モデルデータベース 31 に記録している音響モデルとの類似度計算が可能な特徴ベクトルであり、フレーム処理されたデジタル信号を変換することにより導出される。本実施の形態における特徴量は M F C C (Mel Frequency Cepstrum Coefficient) であるが、L P C (Linear Predictive Coding) ケプストラム又は L P C 係数であってもよい。M F C C は、フレーム処理されたデジタル信号を F F T にて変換し、振幅スペクトルを求め、中心周波数がメル周波数領域で一定間隔であるメルフィルタバンクにて処理し、処理の結果の対数を D C T にて変換し、1 次乃至 14 次等の低次の係数を M F C C と呼ばれる特徴ベクトルとして用いる。尚、次数については、標準化周波数及びアプリケーション等の要因により決定され、数値は限定されない。

10

【0082】

照合部 29 は、音声区間判定部 23 a が音声と判定し確定させた音声区間について、特徴量算出部 28 が導出した特徴量である特徴ベクトルに基づいて、音響モデルデータベース 31 に記録している音響モデル並びに認識辞書 32 に記録している認識語彙及び文法を参照し、音声認識処理を実行する。また、認識結果に基づいて、音出力手段 6 及び表示手段 7 等の他の入出力手段に対して出力を制御する。

20

【0083】

その他、実施の形態 1 に対応する部分には同一符号を付して、それらの説明を省略する。

【0084】

このように、本発明の実施の形態 2 では、音データのパワーを基本とした音声区間検出装置の検出結果を、本発明に係る非音声区間検出装置により修正することが可能となり、全体として音声区間検出の精度を向上させることができる。

【0085】

実施の形態 3

実施の形態 3 は、実施の形態 1 及び 2 でスペクトルの偏倚に基づいて非音声区間を検出するのに対し、スペクトルの偏倚、音データのパワー又は音データのピッチについての前フレームとの変化量に基づいて、非音声区間を検出する形態である。また、非音声区間の検出対象から除外する区間を検出し、更に検出対象から除外された区間を復活させる処理をも含む形態である。図 10 は、本発明の実施の形態 3 に係る非音声区間検出装置の一実施例である音声認識装置 1 について、制御手段 2 の音声認識に係る処理構成例を示すブロック図である。また、図 11 は、制御手段 2 の音声認識処理の一例を示すフローチャートである。

30

【0086】

制御手段 2 は、音データからフレームを生成するフレーム生成部 20、生成されたフレームについて、音データのスペクトルの偏倚 / パワー / ピッチを導出するスペクトルの偏倚 / パワー / ピッチ導出部 21 a、導出されたスペクトルの偏倚 / パワー / ピッチについて前フレームとの変化量を導出する変化量導出部 21 b、導出された変化量に基づく判定基準を用いて非音声区間を検出する非音声区間検出部 22 b、検出された非音声区間をもとに音声区間の開始 / 終了を確定させる音声区間判定部 23 b、及び判定された音声区間について音声を認識する音声認識部 24 を備えている。

40

【0087】

ステップ S 41 乃至 S 44 の処理は、夫々図 3 のステップ S 11 乃至 S 14 と同様であるので、説明を省略する。ステップ S 41 乃至 S 44 の処理で生成された各フレームに対し、以下の処理が行われる。

【0088】

50



スペクトルの偏倚／パワー／ピッチ導出部 2 1 a は、フレーム生成部 2 0 からフレームバッファ 4 2 を介して与えられたフレームについて、音データのスペクトルの偏倚、音データのパワー及び音データのピッチの少なくとも一を導出し（ステップ S 4 5 ）、導出したスペクトルの偏倚、パワー及びピッチの少なくとも一をフレームバッファ 4 2 に書き込む。

尚、ここで導出する値は、スカラー量であるスペクトルの偏倚／パワー／ピッチに限定されるものではなく、音響的な特性を表すベクトルであるパワースペクトル、振幅スペクトル、M F C C、L P C ケプストラム、L P C 係数、P L P 係数又は L S P パラメータであってもよい。

【 0 0 8 9 】

10

変化量導出部 2 1 b は、フレームバッファ 4 2 に書き込まれたスペクトルの偏倚、音データのパワー及び音データのピッチの少なくとも一について、前フレームとの変化量を導出してフレームバッファ 4 2 に書き込む（ステップ S 4 6 ）。この場合、書き込まれたフレーム及び変化量を夫々参照するのに用いられるフレームバッファ 4 2 へのポインタ（アドレス）が、ワークメモリ 4 3 上に設けられ、初期化される。

【 0 0 9 0 】

非音声区間検出部 2 2 b は、フレームバッファ 4 2 を介して変化量導出部 2 1 b より与えられたフレームについて、変化量に基づく判定基準により非音声区間を検出するサブルーチン呼び出す（ステップ S 4 7 ）。非音声区間検出部 2 2 b が判定基準を用いて検出した非音声区間のフレームは、フレームバッファ 4 2 を介して順次音声区間判定部 2 3 b に与えられる。その後、音声区間判定部 2 3 b は、音声区間の開始／終了フレームを確定して音声区間の判定を行う（ステップ S 4 8 ）。そして、音声認識部 2 4 は、入力されたフレームバッファ 4 2 の最後（音声区間の最後）まで、音声認識処理を実行する（ステップ S 4 9 ）。

20

【 0 0 9 1 】

ここで、図 1 1 を用いて説明したステップ S 4 6 における変化量について、更に詳述する。

人が発声した場合の音データは、スペクトルの偏倚、パワー及びピッチの何れについても、時間と共にある程度の変動が生じるのを避けられない。逆に音データの上記指標に変動が観測されない場合は、非音声であると見做するのが適当である。

30

例えば、t 番目のフレーム（以下、フレーム t という。t = 1、2、・・・）における高域・低域強度 A を A (t) とするとき、フレーム t での変化量を下記の式 5 及び式 6 のとおり定義する。

【 0 0 9 2 】

$$C(t) = |A(t) - A(t-1)|, \quad t > 1 \quad \dots \text{式 5}$$

$$C(t) = 0, \quad t = 1 \quad \dots \text{式 6}$$

【 0 0 9 3 】

この場合、非音声区間について、例えば下記のような判定を行うことが可能である。

(d) : C(t) 0 . 0 5 のフレームが 0 . 5 秒以上継続する場合は、非音声とする。

(e) : C(t) 0 . 1 のフレームが 1 . 2 秒以上継続する場合は、非音声とする。

40

【 0 0 9 4 】

尚、C(t) による判定は、上記 (d)、(e) に限定されるものではなく、変化量に係る閾値と継続時間に係る閾値との組み合わせにより、異なる条件を設定することが可能である。また、フレームが継続する時間に係る閾値は、フレーム長が一定であるため、フレームの継続する数に係る閾値に置き換えることができる。

更に、スペクトルの偏倚、音データのパワー及び音データのピッチ夫々について変化量を別々に導出し、夫々の変化量について、図 1 1 のステップ S 4 7 を実行して、非音声区間を別々に検出することも可能である。

【 0 0 9 5 】

一方、上述の (d)、(e) の判定基準とは逆に、変化量が大きいフレームは非音声で

50

ない可能性があるため、例えば下記 (f) の判定を加えることが有効である。

(f) :  $C(t) > 0.5$  の場合、 $t - w + 1$  (例えば  $w = 3$ ) から  $t + w - 1$  のフレームを非音声区間の検出対象から除外する。即ちそのときのフレームを含めて前後に  $w$  だけ連なるフレームからなる区間を、非音声区間の検出対象から除外する。

【0096】

また、上記 (f) の判定に拘わらず、変化量が大きいフレームが連なる区間が所定数より短い場合は、単発的に変化量が増大した非音声区間である可能性があるため、例えば下記 (g) の判定を更に加えることが望ましい。

(g) : (f) により、変化量が大きいと判定されるフレームが連なる数が所定数以下であって、(f) により非音声区間の検出対象から除外されている区間が、非音声区間に挟まれている場合は、(f) の判定を覆して非音声区間として検出する。

10

【0097】

上述した内容を踏まえて、非音声区間検出のサブルーチンについて説明する。図12は、非音声区間検出のサブルーチンに係る制御手段2の処理手順を示すフローチャートである。非音声区間検出のサブルーチンが呼び出された場合、制御手段2は、そのときのポインタが示すフレームの変化量が、所定の閾値 (例えば上述した  $0.05$ ) 以下であるか否かを判定する (ステップ S51)。所定の閾値以下であると判定した場合 (ステップ S51 : YES)、制御手段2は、非音声区間検出確定のサブルーチンを読み出し (ステップ S52)、その後リターンする。

【0098】

20

変化量が所定の閾値を超えると判定した場合 (ステップ S51 : NO)、制御手段2は、変化量が第2の閾値 (例えば上述した  $0.5$ ) を超えるか否かを判定する (ステップ S53)。第2の閾値を超えないと判定した場合 (ステップ S53 : NO)、制御手段2はそのままリターンする。

変化量が第2の閾値を超えると判定した場合 (ステップ S53 : YES)、制御手段2は、非音声区間検出除外のサブルーチンを読み出し (ステップ S54)、その後リターンする。

【0099】

図13及び図14は、非音声区間検出除外のサブルーチンに係る制御手段2の処理手順を示すフローチャートであり、図15及び図16は、非音声区間検出確定のサブルーチンに係る制御手段2の処理手順を示すフローチャートである。図13及び図14について、非音声区間検出除外のサブルーチンが呼び出された場合、制御手段2は、そのときのポインタが示すフレームのフレーム番号を「開始フレーム番号」としてワークメモリ43上に記憶する (ステップ S61)。そして、制御手段2は、ワークメモリ43上に設けた「フレームカウンタ」の記憶値を「1」に初期化する (ステップ S62)。ここで、「フレームカウンタ」は、変化量と第2の閾値との比較判定を行ったフレーム数を計数するものである。

30

【0100】

その後、制御手段2は、「フレームカウンタ」の記憶内容が所定数 (例えば  $30$  msecの間に含まれるフレームの数である  $3$ ) 以下であるか否かを判定し (ステップ S63)、所定数以下であると判定した場合 (ステップ S63 : YES)、制御手段2は、「フレームカウンタ」の記憶内容に「1」を加算すると共に (ステップ S64)、フレームバッファへのポインタを1フレーム後方に更新する (ステップ S65)。そして、制御手段2は、そのときのポインタが示すフレームの変化量が、上述した所定の閾値より大きい第2の閾値を超えないか否かを判定する (ステップ S66)。

40

【0101】

変化量が第2の閾値を超えると判定した場合 (ステップ S66 : YES)、制御手段2は、処理をステップ S63に戻す。変化量が第2の閾値以下であると判定した場合 (ステップ S66 : NO)、即ち単発的に変化量が増大した区間が終了した場合、制御手段2は、「開始フレーム番号」に記憶しているフレームに対して「第2の所定数」フレーム前 (

50

ここでは、上述のwフレーム前)が、非音声区間であるか否かを判定する(ステップS67)。「第2の所定数」フレーム前が非音声区間であると判定した場合(ステップS67: YES)、制御手段2は、単発的に変化量が増大した区間が、後に非音声区間と判定される可能性があるものとして、当該区間に「非音声候補区間」のマークを付与する(ステップS68)。

【0102】

ステップS63で「フレームカウント」の記憶内容が所定数を超えると判定した場合(ステップS63: NO)、即ち、変化量の大きい区間が単発的とは言えない程度に継続した場合、制御手段2は、当該区間の終了フレームを検出する処理に移り、フレームバッファへのポインタを1フレーム後方に更新する(ステップS69)。そして、制御手段2は、そのときのポインタが示すフレームの変化量が、第2の閾値を超えるか否かを判定する(ステップS70)。変化量が第2の閾値を超えると判定した場合(ステップS70: YES)、制御手段2は、処理をステップS69に戻す。

【0103】

変化量が第2の閾値以下であると判定した場合(ステップS70: NO)、即ち変化量が第2の閾値より増大した区間が終了した場合、又はステップS67で「第2の所定数」フレーム前が非音声区間でないと判定した場合(ステップS67: NO)、制御手段2は、変化量が増大した区間を非音声区間の検出対象から除外するために、当該区間に「非音声除外区間」のマークを付与する(ステップS71)。

【0104】

ステップS71の処理を終えた場合、又はステップS68の処理を終えた場合、制御手段2は、「開始フレーム番号」の内容から「第2の所定数(ここでは上述のw)-1」を減じる処理を行う(ステップS72)。更に、制御手段2は、そのときのポインタが示すフレームの1つ前のフレーム番号に「第2の所定数(ここでは上述のw)-1」を加えた数を「終了フレーム番号」としてワークメモリ43上に記憶し(ステップS73)、リターンする。

これにより、変化量が第2の閾値を超えた区間を、前後に「w-1」だけ拡張した区間が、「非音声候補区間」又は「非音声除外区間」の扱いとなる。

【0105】

次に、図15及び図16について、非音声区間検出確定のサブルーチンが呼び出された場合、制御手段2は、そのときのポインタが示すフレームのフレーム番号を「開始フレーム番号」としてワークメモリ43上に記憶する(ステップS81)。そして、制御手段2は、ワークメモリ43上に設けた「フレームカウント」の記憶値を「1」に初期化する(ステップS82)。ここで、「フレームカウント」は、変化量と所定の閾値との比較判定を行ったフレーム数を計数するものである。

【0106】

その後、制御手段2は、「フレームカウント」の記憶内容が、ステップS63での所定数とは異なる所定数(例えば上述の0.5秒の間に含まれるフレームの数)以上であるか否かを判定し(ステップS83)、所定数未満であると判定した場合(ステップS83: NO)、制御手段2は、「フレームカウント」の記憶内容に「1」を加算すると共に(ステップS84)、フレームバッファへのポインタを1フレーム後方に更新する(ステップS85)。そして、制御手段2は、そのときのポインタが示すフレームの変化量が、所定の閾値以下であるか否かを判定する(ステップS86)。

【0107】

変化量が所定の閾値以下であると判定した場合(ステップS86: YES)、制御手段2は、処理をステップS83に戻す。変化量が所定の閾値を超えると判定した場合(ステップS86: NO)、即ち変化量が所定の閾値以下であるフレームが所定数未満しか継続しなかった場合、制御手段2は、非音声区間を検出しなかったものとし、「開始フレーム番号」に記憶したフレームの直前のフレームが、非音声候補区間に含まれるか否かを判定する(ステップS87)。

## 【0108】

直前のフレームが非音声候補区間に含まれていると判定した場合（ステップS87：YES）、制御手段2は、当該非音声候補区間を非音声除外区間に変更する（ステップS88）。直前のフレームが非音声候補区間に含まれていないと判定した場合（ステップS87：NO）、又はステップS88の処理を終えた場合、制御手段2は、「開始フレーム番号」の記憶内容を消去して（ステップS89）、リターンする。

## 【0109】

ステップS83で「フレームカウント」の記憶内容が所定数以上であると判定した場合（ステップS83：YES）、制御手段2は、非音声区間の終了フレームを検出する処理に移り、フレームバッファへのポインタを1フレーム後方に更新する（ステップS90）。そして、制御手段2は、そのときのポインタが示すフレームの変化量が、所定の閾値以下であるか否かを判定する（ステップS91）。変化量が所定の閾値以下であると判定した場合（ステップS91：YES）、制御手段2は、処理をステップS90に戻す。

10

## 【0110】

変化量が所定の閾値を超えると判定した場合（ステップS91：NO）、即ち検出した非音声区間が終了した場合、制御手段2は、「開始フレーム番号」に記憶したフレームの直前のフレームが、非音声候補区間に含まれるか否かを判定する（ステップS92）。直前のフレームが非音声候補区間に含まれていると判定した場合（ステップS92：YES）、制御手段2は、当該非音声候補区間のマークを消去して、非音声区間に確定させる（ステップS93）。

20

## 【0111】

直前のフレームが非音声候補区間に含まれていないと判定した場合（ステップS92：NO）、又はステップS93の処理を終えた場合、制御手段2は、そのときのポインタが示すフレームの1つ前のフレーム番号を「終了フレーム番号」としてワークメモリ43上に記憶し（ステップS94）、リターンする。

これにより、「開始フレーム番号」及び「終了フレーム番号」で区切られた区間が、新たに検出された非音声区間となる。

## 【0112】

その他、実施の形態1又は2に対応する部分には同一符号を付して、それらの説明を省略する。

30

## 【0113】

このように、本発明の実施の形態3では、各フレームの音データより導出したスペクトルの偏倚、パワー及びピッチの少なくとも一つについて、前フレームとの変化量 $C(t)$ が、例えば0.05以下となるフレームが、継続時間にして0.5秒に相当する数以上連なる場合、変化量が最初に0.05以下となったフレームから、最後に0.05以下となったフレームまでを非音声区間として検出する。また、単発的に変化量の大きい区間は非音声区間の検出対象から除外し、更に当該区間が非音声区間に挟まれている場合は、判定を覆して非音声区間として検出する。

これにより、本実施の形態3では、変化量が小さくて非音声の特徴を有するフレームが、音声らしからぬ程度まで連なる区間を非音声区間として検出し、人の発声による基準値の補正を要しない。従って、パワー変動の大きい雑音が発生する環境下においても、発声前後か否かに拘わらず、高精度に非音声区間を検出することが可能である。また、単発的に変化量が大きい区間（例えば、エアコンの風量が変動して、定量的な雑音が変化した瞬間）についても、適切に非音声区間の検出を行うことが可能となる。

40

## 【0114】

尚、実施の形態3にあっては、変化量導出部21bがフレーム $t$ において導出する変化量 $C(t)$ は、上述の式5及び式6に限定されるものではなく、フレーム $t$ の前後 $v$ （例えば $v=2$ ）フレームの区間、即ちフレーム $t-v$ からフレーム $t+v$ の区間において、下記の式7又は式8で定義される最大値であってもよい。

## 【0115】

50

【数 3】

$$D(t) = \max_{j \leq i \leq t+v} A(i) - \min_{j \leq i \leq t+v} A(i) \quad , \quad j = \max(0, t-v) \quad \cdots \text{式 7}$$

【数 4】

$$E(t) = \max_{j \leq i \leq t+v} C(i) \quad , \quad j = \max(0, t-v) \quad \cdots \text{式 8}$$

10

【0116】

これにより、変化量は $C(t)$ 近傍のフレームにおける変化量の最大値と置き換わるため、非音声区間が検出され難くなって、非音声区間を誤検出することを抑止することができる。

【0117】

また、実施の形態1（又は実施の形態3）にあつては、スペクトルの偏倚導出部21（又はスペクトルの偏倚／パワー／ピッチ導出部21a）は、フレーム $t$ の前後 $z$ （例えば $z=3$ ）フレームの区間、即ちフレーム $t-z$ からフレーム $t+z$ の区間におけるスペクトルの偏倚の最大値、最小値、平均値及び中央値の少なくとも一を導出して、導出した値を夫々フレーム $t$ についてのスペクトルの偏倚としてもよい。これらの統計的な集計値を用いることにより、短時間で急激な信号変化があった場合に、スペクトルの偏倚の誤認識を防止することができる。この場合、新たに導出した夫々のスペクトルの偏倚について、非音声区間を別々に検出することが可能である。

20

【0118】

実施の形態4

実施の形態4は、実施の形態1において、スペクトルの偏倚が所定の閾値以上となるフレームが、所定数以上連なる区間を非音声区間として検出するのに対し、スペクトルの偏倚が所定の閾値以上となるフレームが、所定の割合を超える区間について、当該区間が所定数以上のフレームに亘って連なる場合、当該区間を非音声区間として検出する形態である。

30

図17及び図18は、本発明の実施の形態4に係る非音声検出装置の一実施例である音声認識装置1について、非音声区間検出のサブルーチンに係る制御手段2の処理手順を示すフローチャートである。

【0119】

非音声区間検出のサブルーチンが呼び出された場合、制御手段2は、そのときのポイントが示すフレームのスペクトルの偏倚が、所定の閾値以上であるか否かを判定する（ステップS111）。所定の閾値未満であると判定した場合（ステップS111：NO）、制御手段2は、ワークメモリ43に記憶されたフレームバッファ42へのポイントを1フレーム後方に更新して（ステップS112）、リターンする。

40

これにより、制御手段2は、非音声区間を検出することなくリターンする。

【0120】

所定の閾値以上であると判定した場合（ステップS111：YES）、制御手段2は、そのときのポイントが示すフレームのフレーム番号を「開始フレーム番号」としてワークメモリ43上に記憶する（ステップS113）。そして、制御手段2は、ワークメモリ43上に設けた「フレームカウント1」の記憶値を「1」に初期化し（ステップS114）、更に「フレームカウント2」の記憶値を「1」に初期化する（ステップS115）。ここで、「フレームカウント1」は、スペクトルの偏倚と所定の閾値との比較判定を行ったフレーム数を計数するものである。また、「フレームカウント2」は、スペクトルの偏倚

50

が所定の閾値以上となったフレーム数を計数するものである。

【0121】

その後、制御手段2は、「フレームカウント1」の記憶内容が所定数以上であるか否かを判定し(ステップS116)、所定数未満であると判定した場合(ステップS116:NO)、制御手段2は、「フレームカウント1」の記憶内容に「1」を加算すると共に(ステップS117)、フレームバッファへのポインタを1フレーム後方に更新する(ステップS118)。そして、制御手段2は、そのときのポインタが示すフレームのスペクトルの偏倚が、所定の閾値以上であるか否かを判定する(ステップS119)。

【0122】

スペクトルの偏倚が所定の閾値以上であると判定した場合(ステップS119:YES)、制御手段2は、「フレームカウント2」の記憶内容に「1」を加算して(ステップS120)、処理をステップS116に戻す。スペクトルの偏倚が所定の閾値未満であると判定した場合(ステップS119:NO)、制御手段2は、「フレームカウント1」の記憶内容に対する「フレームカウント2」の記憶内容の比、即ちスペクトルの偏倚を判定した全フレームに対する、スペクトルの偏倚が所定の閾値以上となったフレームの割合が、所定の割合(例えば0.8)以上であるか否かを判定する(ステップS121)。

【0123】

所定の割合以上であると判定した場合(ステップS121:YES)、制御手段2は、処理をステップS116に戻す。所定の割合未満であると判定した場合(ステップS121:NO)、制御手段2は、「開始フレーム番号」の内容を消去して(ステップS122)、リターンする。

これにより、制御手段2は、非音声区間を検出することなくリターンする。

【0124】

ステップS116で「フレームカウント1」の記憶内容が所定数以上であると判定した場合(ステップS116:YES)、制御手段2は、非音声区間の終了フレームを検出する処理に移り、「フレームカウント」の記憶内容に「1」を加算すると共に(ステップS123)、フレームバッファへのポインタを1フレーム後方に更新する(ステップS124)。そして、制御手段2は、そのときのポインタが示すフレームのスペクトルの偏倚が、所定の閾値以上であるか否かを判定する(ステップS125)。

【0125】

スペクトルの偏倚が所定の閾値以上であると判定した場合(ステップS125:YES)、制御手段2は、「フレームカウント2」の記憶内容に「1」を加算する(ステップS126)。ステップS126の処理を終えた場合、又はスペクトルの偏倚が所定の閾値未満であると判定した場合(ステップS125:NO)、制御手段2は、「フレームカウント1」の記憶内容に対する「フレームカウント2」の記憶内容の比が、所定の割合以上であるか否かを判定する(ステップS127)。

【0126】

所定の割合以上であると判定した場合(ステップS127:YES)、制御手段2は、処理をステップS123に戻す。所定の割合未満であると判定した場合(ステップS127:NO)、制御手段2は、そのときのポインタが示すフレームの1つ前のフレーム番号を「終了フレーム番号」としてワークメモリ43上に記憶し(ステップS128)、リターンする。

これにより、「開始フレーム番号」及び「終了フレーム番号」で区切られた区間が、検出された非音声区間となる。

【0127】

その他、実施の形態1に対応する部分には同一符号を付して、それらの説明を省略する。

【0128】

このように、本発明の実施の形態4では、各フレームの音データより導出したスペクトルの偏倚が所定の閾値以上となるフレームが、所定の割合を超える区間について、当該区

10

20

30

40

50

間が所定数以上のフレームに亘って連なる場合、スペクトルの偏倚が最初に所定の閾値以上となったフレームから、スペクトルの偏倚が所定の閾値以上となるフレームの割合が所定の割合未満となる直前のフレームまでを非音声区間として検出する。

これにより、スペクトルの偏倚が、短時間に変動する場合であっても、高精度に非音声区間を検出することができる。

【0129】

尚、検出する非音声区間の先頭フレームは、最初に所定の閾値以上となったフレームに限定されず、スペクトルの偏倚が所定の閾値以上となるフレームの割合が所定の割合以上である範囲において、前方のフレームまで遡ったフレームを先頭フレームとしてもよい。

【0130】

実施の形態5

実施の形態5は、実施の形態1に対し、信号対雑音比を導出し、導出した信号対雑音比に応じて、スペクトルの偏倚に係る所定の閾値を変更する形態である。

図19は、本発明の実施の形態5に係る非音声検出装置の一実施例である音声認識装置1について、制御手段2の音声認識処理の一例を示すフローチャートである。

【0131】

ステップS131乃至S135の処理は、夫々図3のステップS11乃至S15と同様であるので、説明を省略する。ステップS131乃至S135の処理で生成されてフレームバッファ42に書き込まれたスペクトルの偏倚に対し、以下の処理が行われる。

【0132】

非音声区間検出部22は、フレームバッファ42を介してスペクトルの偏倚導出部21より与えられたフレームについて、非音声区間を検出するサブルーチン呼び出す(ステップS136)。その後、制御手段2は、非音声区間として検出されたフレームの音データ、及び非音声区間以外のフレームの音データに基づいて信号対雑音比を導出し(ステップS137)、導出した信号対雑音比の高/低に応じて、所定の閾値を下降/上昇させるように変更する(ステップS138)。

【0133】

音声区間判定部23は、非音声区間検出部22が非音声区間として検出できなかった区間を音声区間とみなし、そして、音声区間開始フレーム及び音声区間終了フレームを確定させて、一つの音声区間の判定を終える(ステップS139)。このようにして検出された音声区間は、フレームバッファを介して音声認識部24に与えられる。

音声認識部24は、音声認識の分野で一般的な技術を用いて、入力されたフレームバッファ42の最後まで、音声認識処理を実行する(ステップS140)。

【0134】

その他、実施の形態1に対応する部分には同一符号を付して、それらの説明を省略する。

【0135】

このように、本発明の実施の形態5では、非音声区間として検出されたフレームの音データ、及び非音声区間以外のフレームの音データに基づいて信号対雑音比を導出し、導出した信号対雑音比の高/低に応じて、スペクトルの偏倚に係る所定の閾値を下降/上昇させるように変更する。

これにより、信号対雑音比が低下した場合に、雑音の影響により、スペクトルの偏倚が変動して、非音声区間を誤検出することを防止できる。

【0136】

実施の形態6

実施の形態6は、実施の形態1に対し、ピッチの各周波数成分の強度の最大値(以下、ピッチ強度という)を導出し、導出したピッチ強度に応じて、スペクトルの偏倚に係る所定の閾値を変更する形態である。

図20及び図21は、本発明の実施の形態6に係る非音声検出装置の一実施例である音声認識装置1について、非音声区間検出のサブルーチンに係る制御手段2の処理手順を示

10

20

30

40

50

すフローチャートである。

【0137】

非音声区間検出のサブルーチンが呼び出された場合、制御手段2は、そのときのポインタが示すフレームのピッチ強度を導出し(ステップS151)、導出したピッチ強度の大/小に応じて、所定の閾値を下降/上昇させるように変更する(ステップS152)。その後、制御手段2は、当該フレームのスペクトルの偏倚が、所定の閾値以上であるか否かを判定する(ステップS153)。所定の閾値未満であると判定した場合(ステップS153:NO)、制御手段2は、ワークメモリ43に記憶されたフレームバッファ42へのポインタを1フレーム後方に更新して(ステップS154)、リターンする。

これにより、制御手段2は、非音声区間を検出することなくリターンする。

10

【0138】

所定の閾値以上であると判定した場合(ステップS153:YES)、制御手段2は、そのときのポインタが示すフレームのフレーム番号を「開始フレーム番号」としてワークメモリ43上に記憶する(ステップS155)。そして、制御手段2は、ワークメモリ43上に設けた「フレームカウント」の記憶値を「1」に初期化する(ステップS156)。ここで、「フレームカウント」は、スペクトルの偏倚と所定の閾値との比較判定を行ったフレーム数を計数するものである。

【0139】

その後、制御手段2は、「フレームカウント」の記憶内容が所定数以上であるか否かを判定し(ステップS157)、所定数未満であると判定した場合(ステップS157:NO)、制御手段2は、「フレームカウント」の記憶内容に「1」を加算すると共に(ステップS158)、フレームバッファ42へのポインタを1フレーム後方に更新する(ステップS159)。その後、制御手段2は、そのときのポインタが示すフレームのピッチ強度を導出し(ステップS160)、導出したピッチ強度に基づいて所定の閾値を変更する(ステップS161)。

20

【0140】

次いで、制御手段2は、スペクトルの偏倚が所定の閾値以上であるか否かを判定する(ステップS162)。所定の閾値以上であると判定した場合(ステップS162:YES)、制御手段2は、処理をステップS157に戻す。所定の閾値未満であると判定した場合(ステップS162:NO)、制御手段2は、「開始フレーム番号」の内容を消去して(ステップS163)、リターンする。

30

これにより、制御手段2は、非音声区間を検出することなくリターンする。

【0141】

ステップS157で「フレームカウント」の記憶内容が所定数以上と判定した場合(ステップS157:YES)、制御手段2は、非音声区間の終了フレームを検出する処理に移り、フレームバッファ42へのポインタを1フレーム後方に更新する(ステップS164)。その後、制御手段2は、そのときのポインタが示すフレームのピッチ強度を導出し(ステップS165)、導出したピッチ強度に基づいて所定の閾値を変更する(ステップS166)。

【0142】

40

次いで、制御手段2は、当該フレームのスペクトルの偏倚が所定の閾値以上であるか否かを判定する(ステップS167)。所定の閾値以上であると判定した場合(ステップS167:YES)、制御手段2は、処理をステップS164に戻す。所定の閾値未満であると判定した場合(ステップS167:NO)、制御手段2は、そのときのポインタが示すフレームの1つ前のフレーム番号を「終了フレーム番号」としてワークメモリ43上に記憶し(ステップS168)、リターンする。

これにより、「開始フレーム番号」及び「終了フレーム番号」で区切られた区間が、検出された非音声区間となる。

【0143】

ここで、図20図21を用いて説明したステップS151、S160及びS165にお

50



けるピッチ強度について詳述する。

ピッチ強度  $B$  は、短時間スペクトル  $S()$  の自己相関関数  $( )$  を用いて、以下の式 9 を用いて導出することができる。

【 0 1 4 4 】

$B = \text{argmax} ( ), 1 \quad \text{max}, \quad \dots \dots \text{式 9}$

但し、 $\text{max}$  は、想定される最高ピッチ周波数に対応する値。

【 0 1 4 5 】

例えば、8 0 0 0 Hz サンプリングで、1 フレーム長が 2 5 6 サンプルの場合、短時間スペクトルは、0 ~ 4 0 0 0 Hz を 1 2 9 次元ベクトルで表現できる。この場合、最高ピッチ周波数を 5 0 0 Hz としたとき、短時間スペクトル上では、 $5 0 0 / 4 0 0 0 \times 1 2 8 = 1$  6 より、 $\text{max} = 1 6$  となる。

10

【 0 1 4 6 】

その他、実施の形態 1 に対応する部分には同一符号を付して、それらの説明を省略する。

【 0 1 4 7 】

このように、本発明の実施の形態 6 では、各フレームの音データについて、ピッチ強度を導出し、導出したピッチ強度の大 / 小に応じて、スペクトルの偏倚に係る所定の閾値を下降 / 上昇させる。例えば、ピッチ強度が大きい場合、即ち、ピッチが明確に現れている場合は、音データが音声の母音又は半母音であることが想定される。この場合、スペクトルの偏倚が取り得る値は制限される。従って所定の閾値を下げて非音声区間を検出する判定条件を緩めても、誤検出を抑止して高精度に非音声区間を検出することができる。

20

【 0 1 4 8 】

尚、導出したピッチ強度に応じて所定の閾値を変更するのではなく、例えば下記 (h) の判定を加えてもよい。

(h) : ピッチ強度  $B$  所定の強度、且つ、 $|A| \quad 0.5$  が 0.5 秒以上継続する場合、当該区間は非音声とする。(上述した (b) 又は (c) の判定とピッチ強度とを組合せて改良したもの)

【 0 1 4 9 】

実施の形態 7

実施の形態 7 は、実施の形態 1 において、スペクトルの偏倚に係る所定の閾値を、事前

30

の学習によって決定する形態である。

図 2 2 は、本発明の実施の形態 7 に係る非音声検出装置の一実施例である音声認識装置 1 について、制御手段 2 の音声認識処理の一例を示すフローチャートである。

【 0 1 5 0 】

ステップ S 1 7 1 乃至 S 1 7 4 の処理は、夫々図 3 のステップ S 1 1 乃至 S 1 4 と同様であるので、説明を省略する。ステップ S 1 7 1 乃至 S 1 7 4 の処理で生成された各フレームに対し、以下の処理が行われる。

【 0 1 5 1 】

制御手段 2 は、フレームバッファ 4 2 を介して与えられたフレームについて、音データにおける発声区間をマーキングする (ステップ S 1 7 5)。この場合、学習用の音声データには、音素ラベリングがされているため、容易に発声区間をマーキングすることが可能である。更に、制御手段 2 は、スペクトルの偏倚  $|A|$  が取り得る値の範囲  $[-1, -1]$  内に  $N$  個の閾値を設定する (ステップ S 1 7 6)。そして、制御手段 2 は、 $N$  個の閾値のうち 1 つの閾値について、当該閾値以上となるフレームが継続する最大数を集計する (ステップ S 1 7 7)。

40

【 0 1 5 2 】

次いで、制御手段 2 は、 $N$  個の閾値全てについての集計を終了したか否かを判定する (ステップ S 1 7 8)。未集計の閾値があると判定した場合 (ステップ S 1 7 8 : NO)、制御手段 2 は、処理をステップ S 1 7 7 に戻す。 $N$  個の閾値全てについての集計を終了したと判定した場合 (ステップ S 1 7 8 : YES)、制御手段 2 は、集計した結果に基づい

50

て、スペクトルの偏倚に係る所定の閾値を決定する（ステップS179）。

この場合、所定の閾値を大きめに（又は小さめに）決定して、非音声区間の誤検出を抑止することが好ましい。

【0153】

このように、本発明の実施の形態7では、既存の音声データのマーキングされた発声区間について、予め複数の閾値候補を準備し、所定の閾値以上となるフレームが継続する最大数を集計した結果に基づいて、複数の閾値候補の中から、スペクトルの偏倚に係る所定の閾値の最適値を決定する。

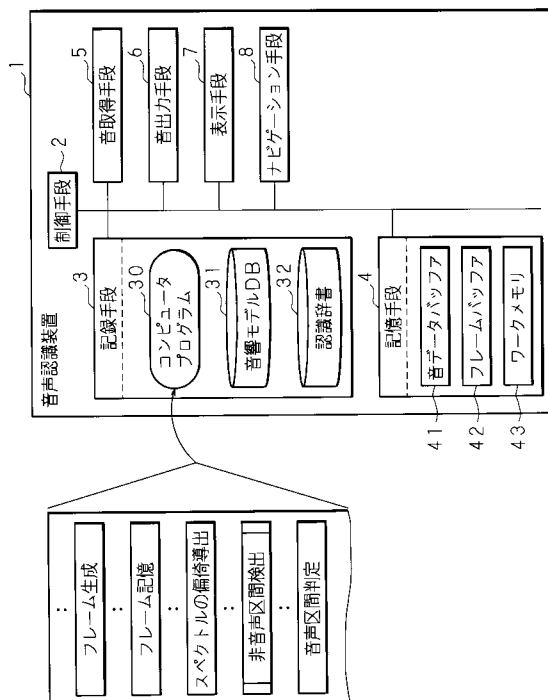
これにより、高精度に非音声区間を検出することができる。

【0154】

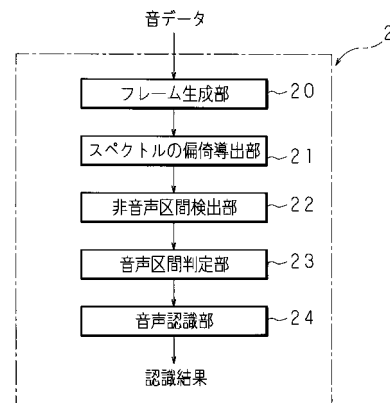
実施の形態1乃至7にあっては、高域・低域強度の絶対値 $|A|$ をスペクトルの偏倚とし、スペクトルの偏倚が所定の正の閾値以上であるか否かを判定する場合について説明したが、高域・低域強度 $A$ をスペクトルの偏倚とし、スペクトルの偏倚が正の値（又は負の値）の場合、所定の正の閾値以上（又は所定の負の閾値以下）であるか否かを判定するようにしてもよい。

10

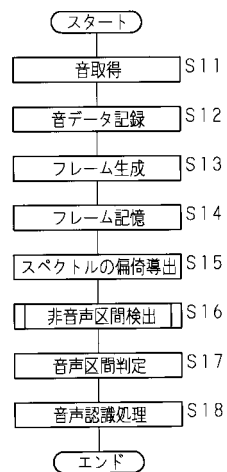
【図1】



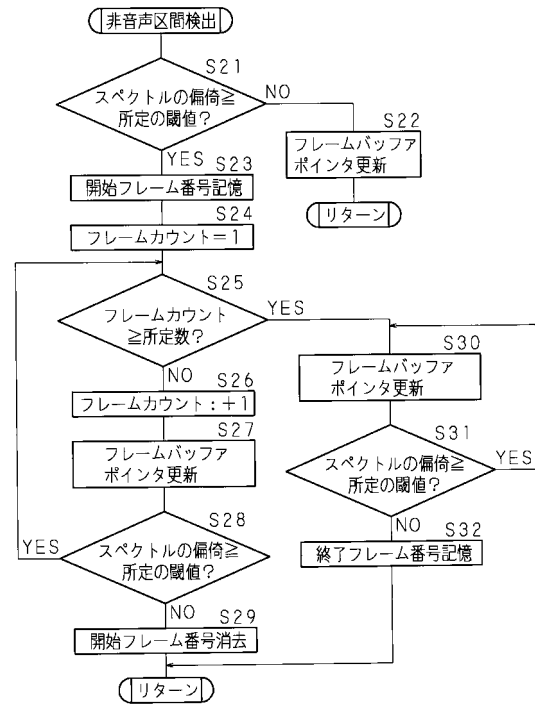
【図2】



【図 3】



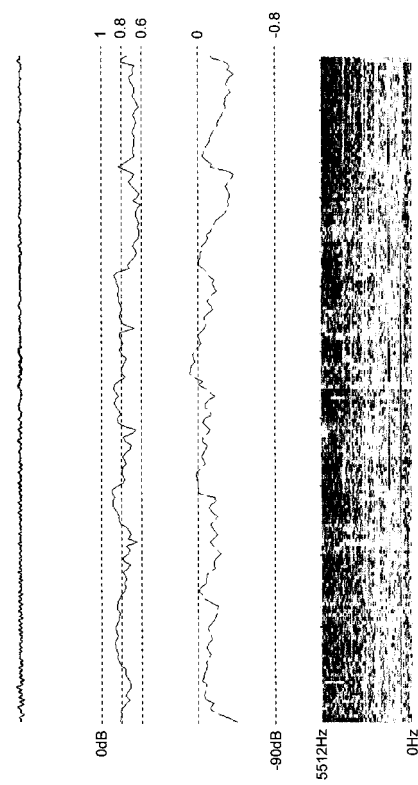
【図 4】



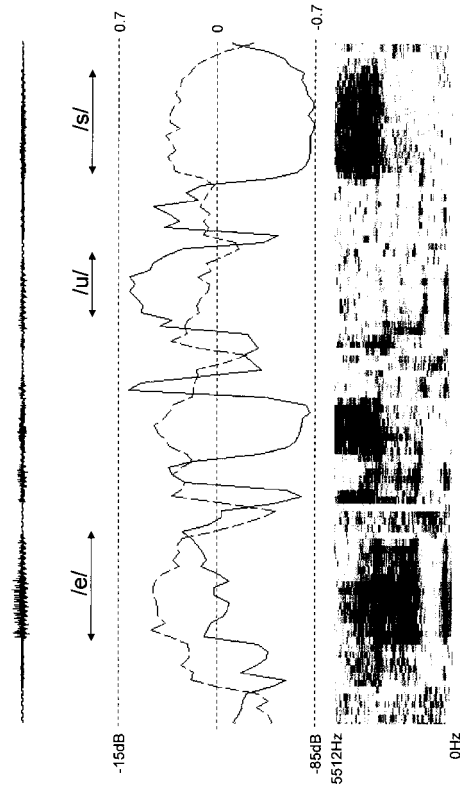
【図 5】



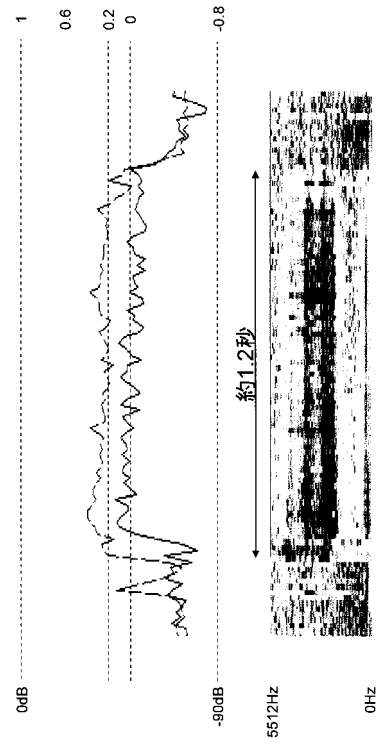
【図 6】



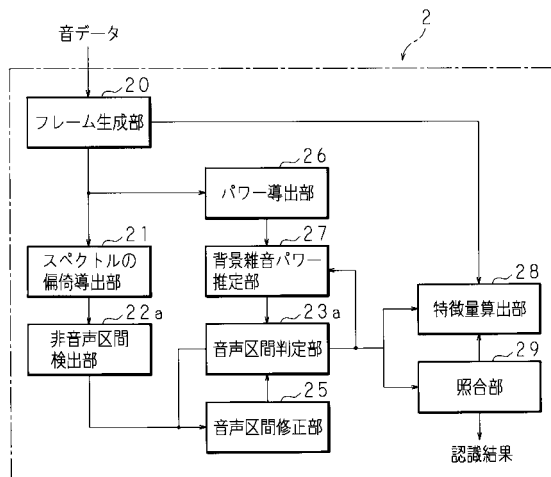
【図 7】



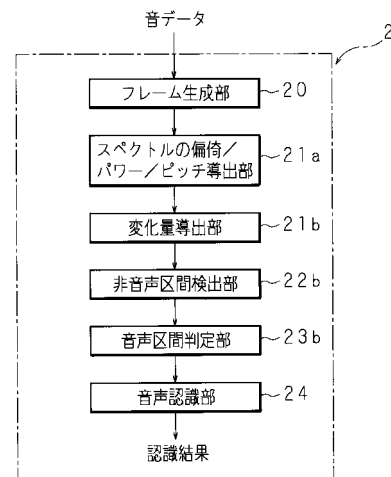
【図 8】



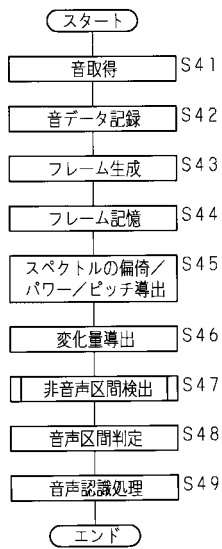
【図 9】



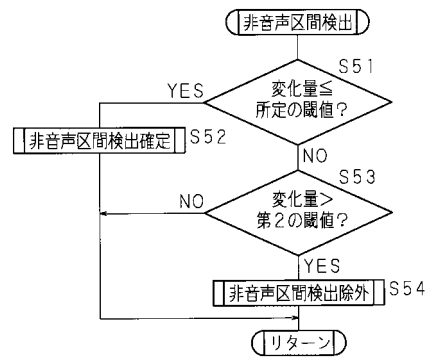
【図 10】



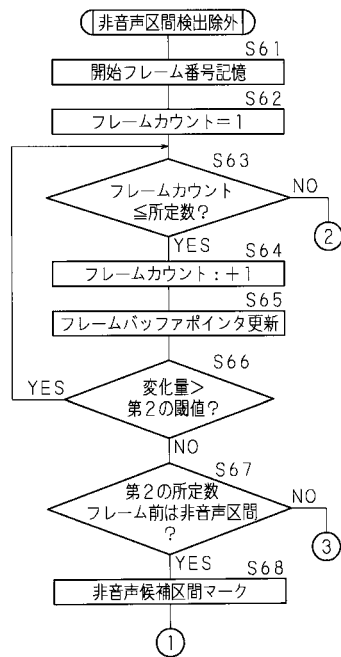
【図 1 1】



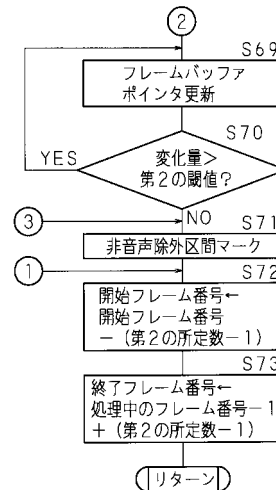
【図 1 2】



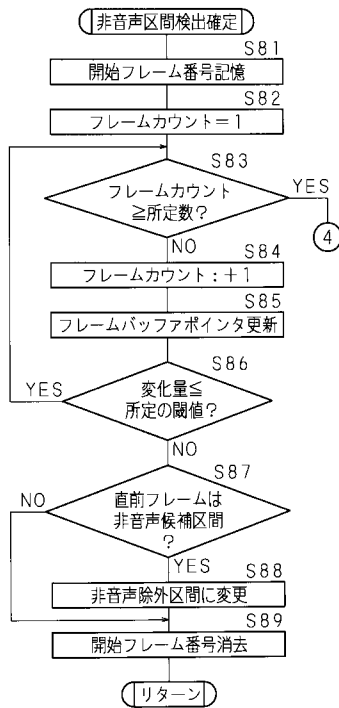
【図 1 3】



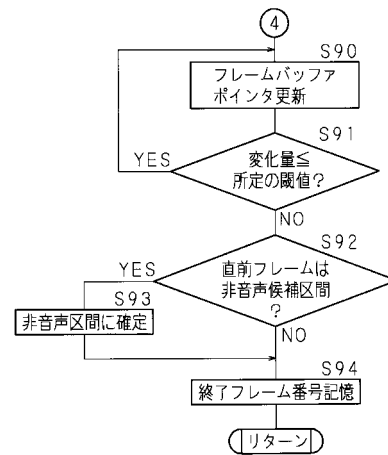
【図 1 4】



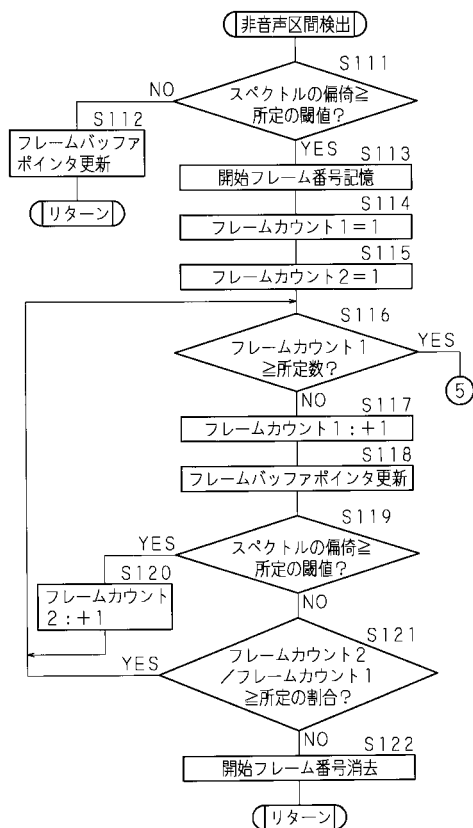
【図 15】



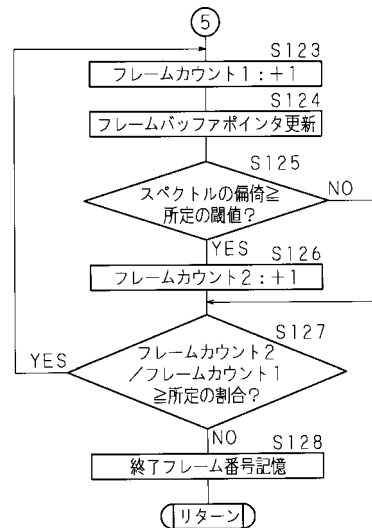
【図 16】



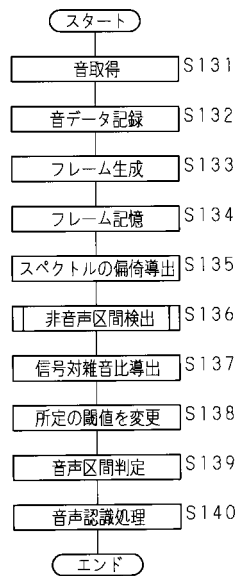
【図 17】



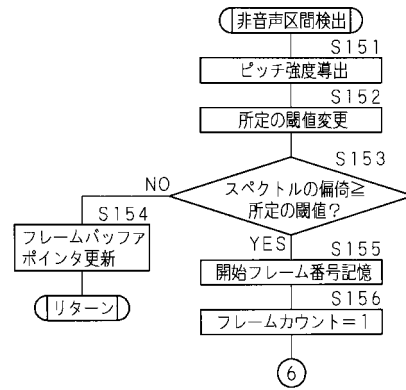
【図 18】



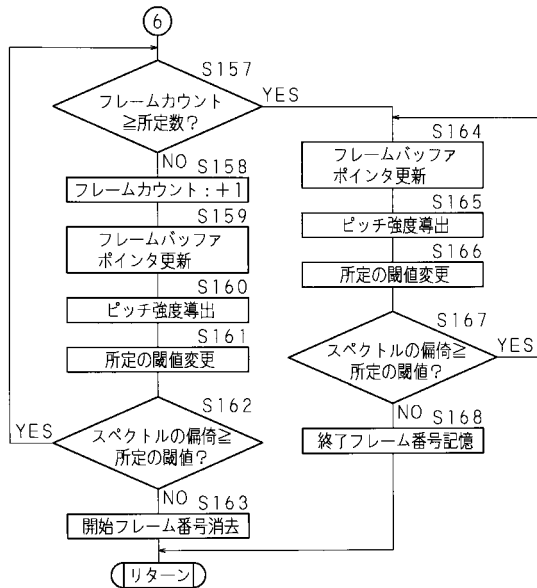
【図 19】



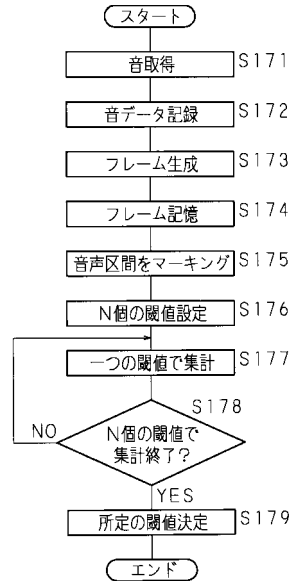
【図 20】



【図 21】



【図 22】



---

フロントページの続き

(56)参考文献 特開平09-152894(JP,A)  
特開2001-350488(JP,A)  
特開平06-083391(JP,A)  
特開2005-156887(JP,A)  
特開2006-209069(JP,A)  
特開2007-233267(JP,A)  
特開平07-212296(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G10L 15/04, 25/78  
JSTPlus(JDreamII)