

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6859499号  
(P6859499)

(45) 発行日 令和3年4月14日(2021.4.14)

(24) 登録日 令和3年3月30日(2021.3.30)

(51) Int.Cl. F I  
G 1 O L 25/84 (2013.01) G 1 O L 25/84

請求項の数 17 (全 19 頁)

|                    |                               |           |                     |
|--------------------|-------------------------------|-----------|---------------------|
| (21) 出願番号          | 特願2019-520035 (P2019-520035)  | (73) 特許権者 | 520015461           |
| (86) (22) 出願日      | 平成29年9月26日 (2017.9.26)        |           | アドバンスド ニュー テクノロジーズ  |
| (65) 公表番号          | 特表2019-535039 (P2019-535039A) |           | カンパニー リミテッド         |
| (43) 公表日           | 令和1年12月5日 (2019.12.5)         |           | 英国領ケイマン諸島 グランド ケイマン |
| (86) 国際出願番号        | PCT/CN2017/103489             |           | ケーワイ1-9008 ジョージ タウ  |
| (87) 国際公開番号        | W02018/068636                 |           | ン ホスピタル ロード 27 ケイマン |
| (87) 国際公開日         | 平成30年4月19日 (2018.4.19)        |           | コーポレート センター         |
| 審査請求日              | 令和1年6月12日 (2019.6.12)         | (74) 代理人  | 100188558           |
| (31) 優先権主張番号       | 201610890946.9                |           | 弁理士 飯田 雅人           |
| (32) 優先日           | 平成28年10月12日 (2016.10.12)      | (74) 代理人  | 100205785           |
| (33) 優先権主張国・地域又は機関 | 中国 (CN)                       |           | 弁理士 ▲高▼橋 史生         |
| 早期審査対象出願           |                               |           |                     |
|                    |                               |           | 最終頁に続く              |

(54) 【発明の名称】 音声信号検出方法及び装置

(57) 【特許請求の範囲】

【請求項 1】

コンピュータにより実施される方法であって、  
ユーザ端末により、オーディオ信号を取得するステップと；  
所定の音声信号のサンプリングレートと前記所定の音声信号の周波数との比率を特定するステップと；

前記ユーザ端末により、前記オーディオ信号を、前記比率で示される数のサンプルを含む、最大数量の短時間エネルギーフレームに分割するステップと；

前記ユーザ端末により、各短時間エネルギーフレームのエネルギーを特定するステップと；

前記ユーザ端末により、各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号が音声信号を含むかどうかを特定するステップと；を備える、

コンピュータにより実施される方法。

【請求項 2】

前記オーディオ信号は、前記サンプリングレートで収集され、パルス符号変調 (PCM) 方式である、

請求項 1 に記載の方法。

【請求項 3】

前記取得したオーディオ信号が、非 PCM 方式であり、

前記オーディオ信号を分割する前に、

10

20

前記オーディオ信号をパルス符号変調（PCM）方式に変換するステップと；  
前記オーディオ信号の前記サンプリングレートを識別するステップと；を更に備える

請求項 1 に記載の方法。

【請求項 4】

前記各短時間エネルギーフレームのエネルギーは、各短時間エネルギーフレームの各サンプリングポイントに関連付けられたエネルギーの合計であり、前記各サンプリングポイントに関連付けられたエネルギーは、前記短時間エネルギーフレームのサンプリングポイントに対応する前記オーディオ信号の振幅に基づいて特定される、

請求項 1 に記載の方法。

10

【請求項 5】

前記オーディオ信号が音声信号を含むかどうかを特定するステップは、

複数の高エネルギーフレームを特定するステップであって、前記複数の高エネルギーフレームの各高エネルギーフレームは、エネルギーが所定の閾値よりも大きい短時間エネルギーフレームである、ステップと；

前記オーディオ信号に含まれる前記短時間エネルギーフレームの量に対する前記複数の高エネルギーフレームの量の比によって表される高エネルギーフレーム比率を特定するステップと；

前記高エネルギーフレーム比率が所定の値より大きいかどうかを特定するステップと；

前記高エネルギーフレーム比率が前記所定の値より大きいと特定された場合に、

20

前記オーディオ信号には音声信号が含まれていると特定するステップ；又は、

前記高エネルギーフレーム比率が前記所定の値より大きくないと特定された場合に、

前記オーディオ信号には音声信号が含まれていないと特定するステップ；を備える、

請求項 1 に記載の方法。

【請求項 6】

前記高エネルギーフレーム比率が所定の値よりも大きいと特定され、更に、

前記オーディオ信号に含まれる前記短時間エネルギーフレームから、所定数の連続する短時間エネルギーフレームがあるかどうかを特定するステップであって、前記所定数の連続する短時間エネルギーフレームのそれぞれは、前記所定の閾値よりも大きいエネルギーを有する、ステップと；

30

肯定の場合に、前記オーディオ信号には音声信号が含まれていると特定するステップ；又は、

肯定でない場合に、前記オーディオ信号には音声信号が含まれていないと特定するステップ；を備える、

請求項 5 に記載の方法。

【請求項 7】

所定の操作を実行するためにコンピュータシステムによって実行可能な 1 又は複数の命令を格納する、非一時的なコンピュータ読取可能媒体であって、前記操作は、

ユーザ端末により、オーディオ信号を取得するステップと；

所定の音声信号のサンプリングレートと前記所定の音声信号の周波数との比率を特定するステップと；

40

前記ユーザ端末により、前記オーディオ信号を、前記比率で示される数のサンプルを含む、最大数量の短時間エネルギーフレームに分割するステップと；

前記ユーザ端末により、各短時間エネルギーフレームのエネルギーを特定するステップと；

前記ユーザ端末により、各短時間エネルギーフレームのエネルギーに基づいて、前記オーディオ信号が音声信号を含むかどうかを特定するステップと；を備える、

非一時的なコンピュータ読取可能媒体。

【請求項 8】

前記オーディオ信号は、前記サンプリングレートで収集され、パルス符号変調（PCM）

50

方式である、

請求項 7 に記載の非一時的なコンピュータ読取可能媒体。

【請求項 9】

前記取得したオーディオ信号が、非 PCM 方式であり、

前記オーディオ信号を分割する前に、

前記オーディオ信号をパルス符号変調 (PCM) 方式に変換するステップと；

前記オーディオ信号の前記サンプリングレートを識別するステップと；を更に備える

、

請求項 7 に記載の非一時的なコンピュータ読取可能媒体。

【請求項 10】

前記各短時間エネルギーフレームのエネルギーは、各短時間エネルギーフレームの各サンプリングポイントに関連付けられたエネルギーの合計であり、前記各サンプリングポイントに関連付けられたエネルギーは、前記短時間エネルギーフレームのサンプリングポイントに対応する前記オーディオ信号の振幅に基づいて特定される、

請求項 7 に記載の非一時的なコンピュータ読取可能媒体。

【請求項 11】

前記オーディオ信号が音声信号を含むかどうかを特定するステップは、

複数の高エネルギーフレームを特定するステップであって、前記複数の高エネルギーフレームの各高エネルギーフレームは、エネルギーが所定の閾値よりも大きい短時間エネルギーフレームである、ステップと；

前記オーディオ信号に含まれる前記短時間エネルギーフレームの量に対する前記複数の高エネルギーフレームの量の比によって表される高エネルギーフレーム比率を特定するステップと；

前記高エネルギーフレーム比率が所定の値よりも大きいかどうかを特定するステップと；

前記高エネルギーフレーム比率が前記所定の値よりも大きいと特定された場合に、

前記オーディオ信号には音声信号が含まれていると特定するステップ；又は、

前記高エネルギーフレーム比率が前記所定の値よりも大きくないと特定された場合に、

前記オーディオ信号には音声信号が含まれていないと特定するステップ；を備える、

請求項 7 に記載の非一時的なコンピュータ読取可能媒体。

【請求項 12】

前記高エネルギーフレーム比率が所定の値よりも大きいと特定され、更に、

前記オーディオ信号に含まれる前記短時間エネルギーフレームから、所定数の連続する短時間エネルギーフレームがあるかどうかを特定するステップであって、前記所定数の連続する短時間エネルギーフレームのそれぞれは、前記所定の閾値よりも大きいエネルギーを有する、ステップと；

肯定の場合に、前記オーディオ信号には音声信号が含まれていると特定するステップ；又は、

肯定でない場合に、前記オーディオ信号には音声信号が含まれていないと特定するステップ；を備える、

請求項 11 に記載の非一時的なコンピュータ読取可能媒体。

【請求項 13】

コンピュータにより実施されるシステムであって、

1 又は複数のコンピュータと；

前記 1 又は複数のコンピュータと相互運用可能に接続され、前記 1 又は複数のコンピュータによって実行されると 1 又は複数の操作を実行する 1 又は複数の命令を格納する有形の非一時的な機械読取可能媒体を備えた、1 又は複数のコンピュータメモリデバイスであって、前記 1 又は複数の操作は、

ユーザ端末により、オーディオ信号を取得するステップと；

所定の音声信号のサンプリングレートと前記所定の音声信号の周波数との比率を特定するステップと；

10

20

30

40

50

前記ユーザ端末により、前記オーディオ信号を、前記比率で示される数のサンプルを含む、最大数量の短時間エネルギーフレームに分割するステップと；

前記ユーザ端末により、各短時間エネルギーフレームのエネルギーを特定するステップと；

前記ユーザ端末により、各短時間エネルギーフレームのエネルギーに基づいて、前記オーディオ信号が音声信号を含むかどうかを特定するステップと；を備える、前記 1 又は複数のコンピューターメモリデバイスと；を備える、

コンピュータにより実施されるシステム。

【請求項 1 4】

前記オーディオ信号は、前記サンプリングレートで収集され、パルス符号変調 ( P C M ) 方式である、

請求項 1 3 に記載のコンピュータにより実施されるシステム。

【請求項 1 5】

前記取得したオーディオ信号が、非 P C M 方式であり、

前記オーディオ信号を分割する前に、

前記オーディオ信号をパルス符号変調 ( P C M ) 方式に変換するステップと；

前記オーディオ信号の前記サンプリングレートを識別するステップと；を更に備える、

請求項 1 3 に記載のコンピュータにより実施されるシステム。

【請求項 1 6】

前記各短時間エネルギーフレームのエネルギーは、各短時間エネルギーフレームの各サンプリングポイントに関連付けられたエネルギーの合計であり、前記各サンプリングポイントに関連付けられたエネルギーは、前記短時間エネルギーフレームのサンプリングポイントに対応する前記オーディオ信号の振幅に基づいて特定される、

請求項 1 3 に記載のコンピュータにより実施されるシステム。

【請求項 1 7】

前記オーディオ信号が音声信号を含むかどうかを特定するステップは、

複数の高エネルギーフレームを特定するステップであって、前記複数の高エネルギーフレームの各高エネルギーフレームは、エネルギーが所定の閾値よりも大きい短時間エネルギーフレームである、ステップと；

前記オーディオ信号に含まれる前記短時間エネルギーフレームの量に対する前記複数の高エネルギーフレームの量の比によって表される高エネルギーフレーム比率を特定するステップと；

前記高エネルギーフレーム比率が所定の値より大きいかどうかを特定するステップと；

前記高エネルギーフレーム比率が前記所定の値より大きいと特定された場合に、

前記オーディオ信号には音声信号が含まれていると特定するステップ；又は、

前記高エネルギーフレーム比率が前記所定の値より大きくないと特定された場合に、

前記オーディオ信号には音声信号が含まれていないと特定するステップ；を備える、

請求項 1 3 に記載のコンピュータにより実施されるシステム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本願はコンピュータ技術の分野に関し、特に、音声信号検出方法及び装置に関する。

【背景技術】

【0 0 0 2】

人々は実生活の中でスマートデバイス（例えば、スマートフォンやタブレットコンピュータ）を使って音声メッセージを送信することが多い。しかし、スマートデバイスを使って音声メッセージを送信する場合、通常は、音声メッセージを送信する前にスマートデバイスのスクリーン上の開始ボタン又は終了ボタンをタップする必要がある、これらのタップ操作はユーザにとって非常に不便である。

## 【 0 0 0 3 】

ユーザがボタンをタップすることなく音声メッセージの送信を終えるには、スマートデバイスが連続的に、又は、所定の周期に基づいて録音を実行し、取得されたオーディオ信号 ( a u d i o   s i g n a l ) が音声信号 ( v o i c e   s i g n a l ) を含むかどうか特定する必要がある。取得されたオーディオ信号が音声信号を含む場合、スマートデバイスは音声信号を抽出してから、音声信号を処理して送信する。そのようにして、スマートデバイスは音声メッセージの送信を終える。

## 【 0 0 0 4 】

既存の技術では、取得されたオーディオ信号が音声信号を含むかどうかを検出するために、通常は、二重閾値法、自己相関最大値に基づく検出法、及びウェーブレット変換に基づく検出法などの音声信号検出法が用いられる。しかし、これらの方法では、通常、フーリエ変換のような複雑な計算を用いてオーディオ情報の周波数特性を求め、更にその周波数特性に基づいてオーディオ情報が音声信号を含むかどうか特定する。したがって、より多くのバッファデータを計算する必要があり、メモリ使用量が比較的多くなり、比較的多くの計算が必要であり、処理速度は比較的遅く、消費電力も比較的大きくなる。

## 【 発明の概要 】

## 【 0 0 0 5 】

本願の実施は音声信号検出方法及び装置を提供し、既存の技術における音声信号検出方法では処理速度が比較的低く、リソース消費が比較的高いという問題を軽減する。

## 【 0 0 0 6 】

以下の技術的解決策が本願の実施で用いられる。

## 【 0 0 0 7 】

音声信号検出方法が提供され、この方法は：オーディオ信号を取得するステップと；所定の音声信号の周波数に基づいて、前記オーディオ信号を複数の短時間エネルギーフレームに分割するステップと；各短時間エネルギーフレームのエネルギーを特定するステップと；各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号が音声信号を含んでいるかどうかを検出するステップと；を含む。

## 【 0 0 0 8 】

音声信号検出装置が提供され、この装置は：オーディオ信号を取得するよう構成された取得モジュールと；所定の音声信号の周波数に基づいて、前記オーディオ信号を複数の短時間エネルギーフレームに分割するよう構成された分割モジュールと；各短時間エネルギーフレームのエネルギーを特定するよう構成された特定モジュールと；各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号は音声信号を含んでいるかどうかを検出するよう構成された検出モジュールと；を含む。

## 【 0 0 0 9 】

本願の実施において用いられる先に述べた技術的解決策の少なくとも1つは、以下の有益な効果を奏する。

## 【 0 0 1 0 】

既存の技術では、フーリエ変換のような複雑な計算を通して、オーディオ信号が音声信号を含むかどうか特定される。対照的に、本願の実施で用いられる音声信号検出方法では、フーリエ変換のような複雑な計算を行う必要はない。取得されたオーディオ信号は、所定の音声信号の周波数に基づいて複数の短時間エネルギーフレームに分割され、各短時間エネルギーフレームのエネルギーが更に特定され、そして、各短時間エネルギーフレームのエネルギーに基づいて、取得されたオーディオ信号が音声信号を含むかどうかを検出できる。したがって、本願の実施で提供される音声信号検出方法においては、既存の技術における音声信号検出方法では処理速度が比較的遅くリソース消費が比較的高い、という問題を軽減できる。

## 【 図面の簡単な説明 】

## 【 0 0 1 1 】

本明細書で述べる添付図面は本願の更なる理解を提供し、本願の一部を構成するもので

10

20

30

40

50

ある。本願の例示の実施とその記述は本願を説明するものであり、本願に制限を設けるものではない。添付図面について以下のとおり説明する。

【 0 0 1 2 】

【 図 1 】 図 1 は、本願の実施に係る音声信号検出方法を示すフローチャートである。

【 0 0 1 3 】

【 図 2 】 図 2 は、本願の実施に係る別の音声信号検出方法を示すフローチャートである。

【 0 0 1 4 】

【 図 3 】 図 3 は、本願の実施に係る所定の持続時間の音声信号を示す表示図である。

【 0 0 1 5 】

【 図 4 】 図 4 は、本願の実施に係る音声信号検出装置の構造を示す概略図である。

10

【 発明を実施するための形態 】

【 0 0 1 6 】

本願の目的、技術的解決策及び利点を明瞭にするために、以下では、本願の具体的な実施及び添付図面を参照しながら本願の技術的解決策を明確且つ包括的に記述する。記述するこれらの実施は本願の実施の全てではなく、むしろそのいくつかに過ぎないことは言うまでもない。創造的な努力なく本願の実施に基づいて当業者により得られるその他の全ての実施は、本願の保護範囲に含まれる。

【 0 0 1 7 】

本願の実施で提供される技術的解決策を、添付の図面を参照して、以下詳細に説明する。

20

【 0 0 1 8 】

既存の技術の音声信号検出方法における比較的低い処理速度及び比較的高いリソース消費という問題を軽減するために、本願の実施は音声信号検出方法を提供する。

【 0 0 1 9 】

本方法を実行する主体は、携帯電話、タブレットコンピュータ、又はパーソナルコンピュータ ( Personal Computer、PC ) などのユーザ端末であってもよいが、これらに限定されず、これらユーザ端末上で作動するアプリケーション ( APP : 以後「アプリ」とする ) であっても、サーバなどのデバイスであってもよい。

【 0 0 2 0 】

説明を容易にするために、本方法を実行する主体がアプリである実施例を用いて、本方法の実施を、以下説明する。言うまでもなく本方法はアプリによって実行されるが、これは説明のための例にすぎず、本方法に対する限定として解釈されるべきではない。

30

【 0 0 2 1 】

図 1 は、本方法の手順の概略図である。本方法は以下のステップを含む。

【 0 0 2 2 】

ステップ 1 0 1 : オーディオ信号を取得する。

【 0 0 2 3 】

オーディオ信号は、オーディオ収集デバイスを用いてアプリにより収集されたオーディオ信号であっても、アプリにより受信されたオーディオ信号であってもよく、例えば、別のアプリ又はデバイスによって送信されたオーディオ信号であってもよい。実施については本願で限定されない。オーディオ信号を得た後、アプリはオーディオ信号をローカルに格納できる。

40

【 0 0 2 4 】

本願は、オーディオ信号に対応するサンプリングレート、持続時間、方式 ( フォーマット )、サウンドチャンネルなどに対して制限しない。

【 0 0 2 5 】

本願のこの実施において提供される音声信号検出方法では、アプリがオーディオ信号を取得することができ、取得されたオーディオ信号に対して音声信号検出を実行できるのであれば、アプリは、チャットアプリや決済アプリなどの任意のタイプのアプリであってもよい。

50

## 【 0 0 2 6 】

ステップ 1 0 2 : 所定の音声信号の周波数に基づいて、オーディオ信号を複数の短時間エネルギーフレームに分割する。

## 【 0 0 2 7 】

短時間エネルギーフレームは、実際には、ステップ 1 0 1 で取得されたオーディオ信号の一部である。

## 【 0 0 2 8 】

具体的には、所定の音声信号の周波数に基づいて所定の音声信号の周期を特定でき、この特定された周期に基づいて、ステップ 1 0 1 で取得されたオーディオ信号が、対応する持続時間が周期である複数の短時間エネルギーフレームに分割される。例えば、ステップ 1 0 1 で取得されたオーディオ信号の持続時間に基づいて、所定の音声信号の周期が 0 . 0 1 秒であると仮定すると、オーディオ信号を、持続時間が 0 . 0 1 秒であるいくつかの短時間エネルギーフレームに分割できる。注記すると、ステップ 1 0 1 で取得されたオーディオ信号を分割する場合、代替として、オーディオ信号を、実際の状態と所定の音声信号の周波数とに基づいて、少なくとも 2 つの短時間エネルギーフレームに分割してもよい。後に続く説明を分かり易くするために、オーディオ信号が複数の短時間エネルギーフレームに分割される例を本願のこの実施で用いて、以下説明する。

## 【 0 0 2 9 】

更に、ステップ 1 0 1 でアプリがオーディオ収集デバイスを用いてオーディオ信号を収集する場合、一般に、オーディオ信号を収集することは、ある特定のサンプリングレートで、実際にはデジタル信号を形成するためのアナログ信号であるオーディオ信号、すなわちパルスコード変調 ( P u l s e C o d e M o d u l a t i o n 、 P C M ) 方式のオーディオ信号を収集することであるため、オーディオ信号は、オーディオ信号のサンプリングレートと所定の音声信号の周波数とに基づいて、更に複数の短時間エネルギーフレームに分割できる。

## 【 0 0 3 0 】

具体的には、所定の音声信号の周波数に対するオーディオ信号のサンプリングレートの比率  $m$  を特定でき、次いで、収集されたデジタルオーディオ信号内の各  $m$  個のサンプリング点は、比率  $m$  に基づいて 1 つの短時間エネルギーフレームにグループ化される。比率  $m$  が正の整数である場合、オーディオ信号を、 $m$  に基づいて最大数量の短時間エネルギーフレームに分割でき、 $m$  が正の整数ではない場合、オーディオ信号を、正の整数に丸められる ( 端数処理する )  $m$  に基づいて最大数量の短時間エネルギーフレームに分割できる。注記すると、ステップ 1 0 1 で取得されたオーディオ信号に含まれるサンプリング点の数量が  $m$  の整数倍でない場合、オーディオ信号が最大数量の短時間エネルギーフレームに分割された後に、残りのサンプリング点を破棄してもよい、又は、その代わりに、残りのサンプリング点を後続の処理のための短時間エネルギーフレームとして用いてもよい。所定の音声信号の周期における、ステップ 1 0 1 で取得されたオーディオ信号に含まれるサンプリング点の数量を表すために  $M$  を用いる。

## 【 0 0 3 1 】

例えば、所定の音声信号の周波数が 8 2 H z の場合、ステップ 1 0 1 で取得されたオーディオ信号の持続時間は 1 秒であり、サンプリングレートは 1 6 0 0 0 H z であり、比率  $m = 1 6 0 0 0 / 8 2 = 1 9 5 . 1$  である。ここで、 $m$  は正の整数ではないので、1 9 5 . 1 は正の整数 1 9 5 に丸められる。オーディオ信号の持続時間とサンプリングレートとに基づき、オーディオ信号に含まれるサンプリング点の数量は 1 6 0 0 0 であると特定できる。オーディオ信号に含まれるサンプリング点の数量は 1 9 5 の整数倍ではないので、オーディオ信号が 8 2 の短時間エネルギーフレームに分割された後、残りの 1 0 のサンプリング点は破棄してもよい。各短時間エネルギーフレームに含まれるサンプリング点の数量は 1 9 5 である。

## 【 0 0 3 2 】

ステップ 1 0 1 で取得されたオーディオ信号が別のアプリ又はデバイスによって送信さ

10

20

30

40

50

れた受信オーディオ信号である場合、オーディオ信号は、前述の方法のうちのいずれか 1 つを用いて複数の短時間エネルギーフレームに分割できる。注記すると、オーディオ信号の方式が P C M 方式ではない場合がある。前述の方法でオーディオ信号のサンプリングレートと所定の音声信号の周波数とに基づいて分割することにより短時間エネルギーフレームが得られる場合、受信オーディオ信号を P C M 方式のオーディオ信号に変換する必要がある。更に、オーディオ信号を受信したときには、オーディオ信号のサンプリングレートを特定する必要がある。オーディオ信号のサンプリングレートを識別する方法は、既存の技術における識別方法であってよい。ここでは説明を簡単にするために詳細は省略する。

【 0 0 3 3 】

ステップ 1 0 3 : 各短時間エネルギーフレームのエネルギーを特定する。

10

【 0 0 3 4 】

本願のこの実施では、P C M 方式のオーディオ信号が、前述の方法で、同じく P C M 方式のいくつかの短時間エネルギーフレームに分割されるとき、短時間エネルギーフレームのエネルギーは、短時間エネルギーフレーム内の各サンプリング点に対応するオーディオ信号の振幅に基づいて特定できる。具体的には、短時間エネルギーフレーム内の各サンプリング点に対応するオーディオ信号の振幅に基づいて各サンプリング点のエネルギーを特定し、次いで、サンプリング点のエネルギーを合計する。最終的に取得されたエネルギーの合計は、短時間エネルギーフレームのエネルギーとして用いられる。

【 0 0 3 5 】

例えば、短時間エネルギーフレームのエネルギーは以下の式を用いて特定できる。

20

【 数 1 】

$$\text{エネルギー} = \sum_i^{i+n} (A_i[t])^2$$

式中、 $i$  はオーディオ信号の  $i$  番目のサンプリング点を表し、 $n$  は短時間エネルギーフレームに含まれるサンプリング点の数量であり、 $A_i[t]$  は  $i$  番目のサンプリング点に対応するオーディオ信号の振幅であり、短時間エネルギーフレームの振幅の値の範囲は、 $-32768$  から  $32767$  である。

【 0 0 3 6 】

30

更に、本願のこの実施においては、計算を簡素化し、リソースを節約するために、振幅を  $32768$  で除した値を更に短時間エネルギーフレームの正規化振幅として使用できる。振幅は、オーディオ信号が収集されたときに得られる。短時間エネルギーフレームの正規化振幅の値の範囲は、 $-1$  から  $1$  である。

【 0 0 3 7 】

短時間エネルギーフレームが P C M 方式ではない場合、振幅計算関数を各瞬間における短時間エネルギーフレームの振幅に基づいて特定でき、積分はその関数の 2 乗に対して実行される。そして最終的に得られる積分結果は短時間エネルギーフレームのエネルギーである。

【 0 0 3 8 】

40

ステップ 1 0 4 : 各短時間エネルギーフレームのエネルギーに基づいて、オーディオ信号に音声信号が含まれているかどうかを検出する。

【 0 0 3 9 】

具体的には、オーディオ信号に音声信号が含まれているかどうかを特定するために、次の 2 つの方法を用いることができる。

【 0 0 4 0 】

方法 1 : 全ての短時間エネルギーフレームの総量に対する、エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の比率 (以下、高エネルギーフレーム比率と呼ぶ) が特定され、特定された高エネルギーフレーム比率は所定の比率より大きいかどうか特定される。それが肯定であれば、オーディオ信号は音声信号を含むと特定され、そうで

50



なければ、オーディオ信号は音声信号を含まないと特定される。

【 0 0 4 1 】

所定の閾値の値及び所定の比率の値は、実際の要求に基づいて設定できる。本願のこの実施において、所定の閾値は2に設定でき、所定の比率は20%に設定できる。高エネルギーフレーム比率が20%より大きい場合、オーディオ信号は音声信号を含むと特定され、そうでなければ、オーディオ信号は音声信号を含まないと特定される。

【 0 0 4 2 】

本願のこの実施では、人が話すとき、実生活の中では外部環境にいくらかのノイズがあり、このノイズのエネルギーは、一般に、人の声よりも低いので、方法1を用いてオーディオ信号が音声信号を含むかどうか特定できる。この場合、エネルギーが所定の閾値より10

【 0 0 4 3 】

方法2：最終的な検出結果をより正確にするために、方法1を用いて、高エネルギーフレーム比率を特定し、特定された高エネルギーフレーム比率が所定の比率より大きいかどうかを特定できる。否定であれば、オーディオ信号は音声信号を含まないと特定される。肯定であれば、エネルギーが所定の閾値より大きい短時間エネルギーフレーム内に少なくともN個の連続する短時間エネルギーフレームがある場合、オーディオ信号は音声信号を含むと特定され、エネルギーが所定の閾値より大きい短時間エネルギーフレーム内に少なく20

【 0 0 4 4 】

具体的には、方法1に基づいて、方法2では、オーディオ信号が音声信号を含むかどうか特定するために以下の要件が追加される。すなわち、エネルギーが所定の閾値より大きい短時間エネルギーフレーム内に、少なくともN個の連続する短時間エネルギーフレームがあるかどうか特定される。そのようにして、ノイズを効果的に減らすことができる。実生活では、ノイズは人の声よりもエネルギーが低く、オーディオ信号はランダムである。方法2では、オーディオ信号が過度のノイズを含む場合を効果的に排除でき、外部環境におけるノイズの影響が低減され、ノイズリダクション機能を果たす。30

【 0 0 4 5 】

注記すると、本願のこの実施において提供される音声信号検出方法は、モノラルオーディオ信号、バイノーラルオーディオ信号、マルチチャンネルオーディオ信号等の検出に適用できる。1つのサウンドチャンネルを用いて収集されたオーディオ信号はモノラルオーディオ信号であり、2つのサウンドチャンネルを用いて収集されたオーディオ信号はバイノーラルオーディオ信号であり、複数のサウンドチャンネルを用いて収集されたオーディオ信号はマルチチャンネルオーディオ信号である。

【 0 0 4 6 】

図1に示す方法でバイノーラルオーディオ信号及びマルチチャンネルオーディオ信号を検出する場合、ステップ101乃至ステップ104で説明した操作を実行することにより、各チャンネルの取得されたオーディオ信号を検出でき、最後に、各チャンネルのオーディオ信号の検出結果に基づいて、取得されたオーディオ信号が音声信号を含むかどうかを特定する。40

【 0 0 4 7 】

具体的には、ステップ101で取得されたオーディオ信号がモノラルオーディオ信号である場合、そのオーディオ信号に対してステップ101乃至ステップ104で説明した操作を、直接、実行でき、検出結果が最終的な検出結果として用いられる。

【 0 0 4 8 】

ステップ101で取得されたオーディオ信号がモノラルオーディオ信号ではなくバイノ 50

ーラルオーディオ信号又はマルチチャンネルオーディオ信号である場合、ステップ101乃至ステップ104で説明した操作を実行することによって各チャンネルの音声信号を処理できる。各チャンネルのオーディオ信号が音声信号を含まないことが検出された場合、ステップ101で取得されたオーディオ信号は音声信号を含まないと特定される。少なくとも1つのチャンネルのオーディオ信号が音声信号を含むことが検出された場合、ステップ101で取得されたオーディオ信号は音声信号を含むと特定される。

#### 【0049】

更に、ステップ102で説明した所定の音声信号の周波数は、任意の音声の周波数とすることができる。実施は本願において限定されない。実際には、現実のケースに基づいて、ステップ101で取得された異なるオーディオ信号に対して異なる周波数の所定の音声信号を設定できる。注記すると、所定の音声信号の周波数は、分割を通して最終的に得られる短時間エネルギーフレームが以下の要求、すなわち短時間エネルギーフレームに対応する持続時間は、ステップ101で取得されたオーディオ信号に対応する周期以上であるとの要求、を満たすという条件で、最高音（ソプラノ）の音声周波数又は最低音（バス）の音声周波数などの任意の音声信号の周波数であってよい。より良好な検出効果を確保して、できるだけ多くのリソースを節約し、処理速度を向上させるために、本願のこの実施では、所定の音声信号の周波数を、人の最低音声周波数、すなわち82Hz、に設定できる。周期は周波数の逆数であるので、所定の音声信号の周波数が人の最低音声周波数である場合、所定の音声信号の周期は人の最高音声周期である。したがって、ステップ101で取得されたオーディオ信号の周期にかかわらず、短時間エネルギーフレームに対応する持続時間は、先に取得されたオーディオ信号の周期以上である。

#### 【0050】

注記すると、本願のこの実施では、人の音声の特徴に基づいてオーディオ信号が音声信号を含むかどうか特定するためにここで論じた検出方法が用いられるので、短時間エネルギーフレームに対応する持続時間は、ステップ101で取得されたオーディオ信号の周期以上であることが要求される。ノイズと比較して、人の音声はより高いエネルギーを持ち、より安定しており、そして連続的である。短時間エネルギーフレームに対応する持続時間がステップ101で取得されたオーディオ信号の周期より短い場合、短時間エネルギーフレームに対応する波形は全周期（completion period）の波形を含まず、短時間エネルギーフレームの期間は比較的短い。この場合、高エネルギーフレーム比率が所定比率よりも大きく、エネルギーが所定の閾値よりも大きい短時間エネルギーフレーム内に少なくともN個の連続する短時間エネルギーフレームがある場合でも、それはオーディオ信号が音響信号（sound signal）を含むことを単に示すだけであり、この音響信号が音声信号であることを示すものではない。したがって、本願のこの実施では、ステップ101で取得されたオーディオ信号の持続時間は、人の最高音声周期よりも長くなければならない。

#### 【0051】

更に、本願のこの実施において提供される音声信号検出方法は、特に、ユーザのタップ操作なくチャットアプリを用いることによって音声メッセージの送信を終えることができるアプリケーションシナリオに適用可能である。シナリオに基づいて、本願のこの実施において提供される音声信号検出方法を、以下、詳細に説明する。このシナリオでは、図2は、本方法の手順の概略図である。本方法は以下のステップを含む。

#### 【0052】

ステップ201：リアルタイムでオーディオ信号を収集する。

#### 【0053】

ユーザは、アプリを起動した後に、タップ操作をせずにチャットアプリが音声メッセージの送信を終えることを期待する場合がある。この場合、アプリは、外部環境を連続的に録音してリアルタイムでオーディオ信号を収集し、ユーザの音声の抜けを減らす。更に、オーディオ信号を収集した後、アプリはオーディオ信号をリアルタイムでローカルに格納できる。ユーザがアプリを停止した後、アプリは録音を停止する。

## 【 0 0 5 4 】

ステップ 2 0 2 : リアルタイムで収集したオーディオ信号から所定の持続時間を持つオーディオ信号を切り取る。

## 【 0 0 5 5 】

アプリがオーディオ信号をリアルタイムで検出する代わりに録音を続けると、音声メッセージはリアルタイムで送信されない。したがって、アプリは、ステップ 2 0 1 で収集されたオーディオ信号から、所定の持続時間を持つオーディオ信号をリアルタイムに切り取り、所定の持続時間を持つオーディオ信号に対して後続の検出を実行できる。

## 【 0 0 5 6 】

所定の持続時間を持つ現在切り取られたオーディオ信号は、現在のオーディオ信号 ( c u r r e n t a u d i o s i g n a l ) と呼ぶことができ、所定の持続時間を持つ最後に切り取られたオーディオ信号は、最後に取得されたオーディオ信号 ( l a s t o b t a i n e d a u d i o s i g n a l ) と呼ぶことができる。

10

## 【 0 0 5 7 】

ステップ 2 0 3 : 所定の音声信号の周波数に基づいて、所定の持続時間内のオーディオ信号を複数の短時間エネルギーフレームに分割する。

## 【 0 0 5 8 】

ステップ 2 0 4 : 各短時間エネルギーフレームのエネルギーを特定する。

## 【 0 0 5 9 】

ステップ 2 0 5 : 各短時間エネルギーフレームのエネルギーに基づいて、所定の持続時間内のオーディオ信号が音声信号を含むかどうかを検出する。

20

## 【 0 0 6 0 】

現在のオーディオ信号が音声信号を含むことが検出された場合、最後に取得されたオーディオ信号が音声信号を含むかどうか特定される。最後に取得されたオーディオ信号が音声信号を含まないと特定されると、現在のオーディオ信号の開始点を音声信号の開始点として特定でき、最後に取得されたオーディオ信号が音声信号を含むと特定されると、現在のオーディオ信号の開始点は音声信号の開始点ではない。

## 【 0 0 6 1 】

現在のオーディオ信号が音声信号を含まないことが検出されると、最後に取得されたオーディオ信号が音声信号を含むかどうか特定される。最後に取得されたオーディオ信号が音声信号を含むと特定されると、最後に取得されたオーディオ信号の終了点は音声信号の終了点として特定でき、最後に取得されたオーディオ信号が音声信号を含まないと特定されると、現在のオーディオ信号の終了点も、最後に取得されたオーディオ信号の終了点も音声信号の終了点ではない。

30

## 【 0 0 6 2 】

例えば、図 3 に示すように、A、B、C、D は、所定の持続時間を持つ 4 つの隣接するオーディオ信号である。オーディオ信号 A と D とは音声信号を含まず、オーディオ信号 B と C とは音声信号を含む。この場合、オーディオ信号 B の開始点を音声信号の開始点と特定し、オーディオ信号 C の終了点を音声信号の終了点と特定できる。

## 【 0 0 6 3 】

40

時に、現在のオーディオ信号がユーザの文言の開始部分又は終了部分であり、そのオーディオ信号には少しの音声信号が含まれていることがある。この場合、アプリは、オーディオ信号が音声信号を含まない、と間違えて特定する可能性がある。現在のオーディオ信号は音声信号を含むことが検出された後、間違った特定によるユーザの音声の抜けを減らすために、最後に取得されたオーディオ信号が音声信号を含むかどうか特定でき、最後に取得されたオーディオ信号は音声信号を含まないと特定された場合、最後に取得されたオーディオ信号の開始点を音声信号の開始点として特定できる。更に、現在のオーディオ信号が音声信号を含まないことが検出された後、最後に取得されたオーディオ信号が音声信号を含むかどうか特定でき、最後に取得されたオーディオ信号が音声信号を含むと特定されると、現在のオーディオ信号の終了点を音声信号の終了点として特定できる。前述の例

50

においては、オーディオ信号 A の開始点を音声信号の開始点と特定し、オーディオ信号 D の終了点を音声信号の終了点として特定できる。

【0064】

現在のオーディオ信号が音声信号を含むことを検出した後、アプリはオーディオ信号を音声識別装置へ送信でき、その結果、音声識別装置はオーディオ信号に対して音声処理を実行して音声結果を取得することができる。その後、音声識別装置はオーディオ信号を後続の処理装置へ送信し、最後に音声メッセージの形式でオーディオ信号を送信する。送信された音声メッセージ内のユーザの音声が必要な文章であることを保証するために、音声信号の特定された開始点と特定された終了点との間の全てのオーディオ信号を音声識別装置へ送信した後、アプリはオーディオ停止信号を音声識別装置へ送信してユーザが現在述べているこの文章が完了した旨を音声識別装置に通知でき、それにより、音声識別装置は全てのオーディオ信号を後続の処理装置へ送信する。最終的に、オーディオ信号は音声メッセージの形式で送信される。

10

【0065】

更に、正確な特定を確実にするために、現在のオーディオ信号を得た後、所定の時間周期を持つ副信号を、最後に取得されたオーディオ信号から更に切り取ることが可能である。現在のオーディオ信号と切り取られた副信号とが連結されて、取得されたオーディオ信号（以下、連結オーディオ信号（concatenated audio signal）と呼ぶ）として機能する。更に、後続の音声信号検出は、連結オーディオ信号に対して実行される。

20

【0066】

副信号は現在のオーディオ信号の前に連結できる。所定の時間周期は、最後に取得されたオーディオ信号のテール時間周期であってよく、時間周期に対応する持続時間は任意の持続時間であってよい。最終的な検出結果がより正確であることを保証するために、本願のこの実施では、所定の時間周期に対応する持続時間は、所定の比率と連結オーディオ信号に対応する持続時間との積以下である値に設定できる。

【0067】

連結オーディオ信号が音声信号を含むことが検出されると、最後に取得された連結オーディオ信号が音声信号を含むかどうかを特定できる。最後に取得された連結オーディオ信号が音声信号を含まないと特定されると、連結オーディオ信号の開始点を音声信号の開始点として用いることができる。連結オーディオ信号が音声信号を含まないことが検出されると、最後に取得された連結オーディオ信号が音声信号を含むかどうかを特定できる。最後に取得された連結オーディオ信号が音声信号を含むと特定されると、連結オーディオ信号の終了点を音声信号の終了点として用いることができる。

30

【0068】

本願のこの実施において、連続的な録音に加えて、アプリは周期的に録音を実行できる。実施は本願のこの実施において限定されない。

【0069】

本願のこの実施で提供される音声信号検出方法は、音声信号検出装置を用いて更に実施できる。図 4 に、この装置の概略構造図を示す。音声信号検出装置は、主に以下のモジュール、すなわち、オーディオ信号を取得するよう構成された取得モジュール 41 と；所定の音声信号の周波数に基づいてオーディオ信号を複数の短時間エネルギーフレームに分割するよう構成された分割モジュール 42 と；各短時間エネルギーフレームのエネルギーを特定するよう構成された特定モジュール 43 と；各短時間エネルギーフレームのエネルギーに基づいて、オーディオ信号が音声信号を含むかどうかを検出するよう構成された検出モジュール 44 と；を含む。

40

【0070】

実施において、取得モジュール 41 は：現在のオーディオ信号を取得し；最後に取得されたオーディオ信号から所定の周期を持つ副信号を切り取り；そして、取得されたオーディオ信号として機能するように、現在のオーディオ信号と切り取られた副信号とを連結す

50

るよう構成される。

【0071】

実施において、分割モジュール42は：所定の音声信号の周波数に基づいて所定の音声信号の周期を特定し；そして、特定された周期に基づいて、オーディオ信号を、対応する持続時間がその周期である複数の短時間エネルギーフレームに分割するよう構成される。

【0072】

実施において、検出モジュール44は：エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の、全ての短時間エネルギーフレームの総量に対する比率を特定し；比率が所定の比率より大きいかどうか特定し；肯定であればオーディオ信号は音声信号を含むと特定し；否定であればオーディオ信号は音声信号を含まないと特定する；よう構成される。

10

【0073】

実施において、検出モジュール44は、エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の、全ての短時間エネルギーフレームの総量に対する比率を特定し；比率が所定の比率より大きいかどうか特定し；否定であればオーディオ信号は音声信号を含まない、と特定し；肯定であればエネルギーが所定の閾値より大きい短時間エネルギーフレーム内に少なくともN個の連続した短時間エネルギーフレームがあるとき、オーディオ信号は音声信号を含む、と特定し；エネルギーが所定の閾値よりも大きい短時間エネルギーフレーム内に少なくともN個の連続する短時間エネルギーフレームがないとき、オーディオ信号は音声信号を含まない、と特定するよう構成される。

20

【0074】

既存の技術では、フーリエ変換のような複雑な計算を通して、オーディオ信号が音声信号を含むかどうか特定される。対照的に、本願の実施で用いられる音声信号検出方法では、フーリエ変換のような複雑な計算を実行する必要はない。取得されたオーディオ信号は、所定の音声信号の周波数に基づいて複数の短時間エネルギーフレームに分割され、各短時間エネルギーフレームのエネルギーが更に特定され、そして、各短時間エネルギーフレームのエネルギーに基づいて、取得されたオーディオ信号が音声信号を含むかどうかを検出できる。したがって、本願の実施において提供される音声信号検出方法では、既存の技術における音声信号検出方法における、処理速度が比較的低く、リソース消費が比較的高いという問題を軽減できる。

30

【0075】

本開示は、本願の実施に係る方法、デバイス（システム）、コンピュータプログラム製品のフローチャート及び／又はブロック図を参照して説明されている。フローチャート及び／又はブロック図内の各プロセス及び／又は各ブロック、並びにフローチャート及び／又はブロック図内のプロセス及び／又はブロックの組み合わせを実施するために、コンピュータプログラム命令を使用することができることを理解されたい。これらのコンピュータプログラム命令は、汎用コンピュータ、専用コンピュータ、組み込みプロセッサ、又はあらゆるその他のプログラマブルデータ処理デバイスに、マシンを生成するために提供されることができ、これにより、コンピュータ、又はあらゆるその他のプログラマブルデータ処理デバイスのプロセッサが、フローチャートの1つ以上のプロセスにおける、及び／又は、ブロック図の1つ以上のブロックにおける、特定の機能を実施するデバイスを生成できるようになる。

40

【0076】

このコンピュータプログラム命令を、コンピュータ又はあらゆるその他のプログラマブルデータ処理デバイスにある方法で機能するように命令することができるコンピュータ読取可能なメモリに記憶して、これらのコンピュータ読取可能なメモリに記憶された命令が、命令装置を含むアーチファクトを作り出すようにすることができる。この命令装置は、フローチャート内の1つ以上のプロセスにおける、及び／又はブロック図内の1つ以上のブロックにおける特定の機能を実施する。

【0077】

50

これらのコンピュータプログラム命令をコンピュータ又はその他のプログラマブルデータ処理デバイスにロードして、コンピュータ又はその他のプログラマブルデバイス上で一連の操作及びステップが実行されるようにし、コンピュータで実施される処理を生成することができる。これにより、コンピュータ又はその他のプログラマブルデバイス上で実行される命令が、フローチャート内の1つ以上のプロセス及び/又はブロック図内の1つ以上のブロックにおける特定の機能を実施するデバイスを提供することを可能とする。

【0078】

典型的な構成では、計算デバイスは1つ以上の中央処理演算装置(CPU)、1つ以上の入出力インターフェース、1つ以上のネットワークインターフェース、及び1つ以上のメモリを含む。

【0079】

メモリは、揮発性メモリ、ランダムアクセスメモリ(RAM)、不揮発性メモリ、及び/又はリードオンリーメモリ(ROM)やフラッシュメモリ(flash RAM)のようなコンピュータ読取可能な媒体を含んでよい。メモリはコンピュータ読取可能な媒体の一例である。

【0080】

コンピュータ読取可能な媒体には、任意の方法又は技術を用いて情報を記憶できる、永続的、非永続的、移動可能な、及び移動不能な媒体が含まれる。この情報はコンピュータ読取可能な命令、データ構造、プログラムモジュール、又はその他のデータであってよい。コンピュータの記憶媒体の例として、相変化ランダムアクセスメモリ(PRAM)、スタティックランダムアクセスメモリ(SRAM)、ダイナミックランダムアクセスメモリ(DRAM)、別タイプのランダムアクセスメモリ、リードオンリーメモリ(ROM)、電氣的に消去可能でプログラム可能なROM(EEPROM)、フラッシュメモリ、又は別のメモリ技術、コンパクトディスクROM(CD-ROM)、デジタル多用途ディスク(DVD)、又は別の光学記憶装置、カセット磁気テープ、磁気テープ/磁気ディスクストレージ、他の磁氣的記憶装置、又は他の任意の非伝送媒体があるが、これに限定されない。このコンピュータの記憶媒体は、計算デバイスによってアクセスできる情報を記憶するように構成することができる。本願の定義に基づき、コンピュータ読取可能な媒体は、変調されたデータ信号及び搬送波のような一時的な媒体(transitory media)を含まない。

【0081】

さらに、用語「含む」、「備える」、又はこれらのその他任意の応用形は、非限定的な包含を網羅するものであるため、一連の要素を含んだ工程、方法、商品、デバイスはこれらの要素を含むだけでなく、ここで明確に挙げていないその他の要素をも含む、あるいは、このような工程、方法、商品、デバイスに固有の要素をさらに含むことができる点に留意することが重要である。「(1つの)~を含む」との用語を付けて示された要素は、それ以上の制約がなければ、その要素を含んだ工程、方法、商品、デバイス内に別の同一の要素をさらに含むことを排除しない。

【0082】

当業者は、本願の実施が方法、システム、又はコンピュータプログラム製品として提供できることを理解するはずである。そのため、本発明は、ハードウェアのみの実施、ソフトウェアのみの実施、又は、ソフトウェアとハードウェアとの組み合わせによる実施を用いることができる。さらに、本発明は、コンピュータで使用可能なプログラムコードを含んだ1台以上のコンピュータで使用可能な記憶媒体(ディスクメモリ、CD-ROM、光学メモリ等を含むがこれに限定されない)上で実施されるコンピュータプログラム製品を使用できる。

【0083】

上述のものは本願の一実施の形態であり、本願を限定することを意図するものではない。当業者は、本願に対して様々な修正及び変更を加えることができる。本願の主旨及び原理から逸脱せずに為されるあらゆる修正、均等物による置換、改善は、本願の特許請求の

10

20

30

40

50

範囲に含まれるものである。

以下、本発明の実施の態様の例を列挙する。

[ 第 1 の局面 ]

音声信号検出方法であって：

オーディオ信号を取得するステップと；

所定の音声信号の周波数に基づいて、前記オーディオ信号を複数の短時間エネルギーフレームに分割するステップと；

各短時間エネルギーフレームのエネルギーを特定するステップと；

各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号が音声信号を備えているかどうかを検出するステップと；を備える、

音声信号検出方法。

[ 第 2 の局面 ]

オーディオ信号を取得する前記ステップは：

現在のオーディオ信号を取得するステップと；

最後に取得されたオーディオ信号から、所定の時間周期を持つ副信号を切り取るステップと；

前記取得されたオーディオ信号として機能するよう、前記現在のオーディオ信号と前記切り取られた副信号とを連結するステップと；を備える、

第 1 の局面に記載の方法。

[ 第 3 の局面 ]

所定の音声信号の周波数に基づいて、前記オーディオ信号を複数の短時間エネルギーフレームに分割する前記ステップは：

前記所定の音声信号の周波数に基づいて前記所定の音声信号の周期を特定するステップと；

前記特定された周期に基づいて、前記オーディオ信号を、対応する持続時間が前記周期である複数の短時間エネルギーフレームに分割するステップと；を備える、

第 1 の局面に記載の方法。

[ 第 4 の局面 ]

各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号が音声信号を備えているかどうかを検出する前記ステップは：

エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の、全短時間エネルギーフレームの総量に対する比率を特定するステップと；

前記比率が所定の比率より大きいかどうかを特定するステップと；

肯定であれば、前記オーディオ信号は音声信号を備える、と特定し、否定であれば、前記オーディオ信号は音声信号を備えないと特定するステップと；を備える、

第 1 の局面に記載の方法。

[ 第 5 の局面 ]

各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号が音声信号を備えているかどうかを検出する前記ステップは：

エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の、全短時間エネルギーフレームの総量に対する比率を特定するステップと；

前記比率が所定の比率より大きいかどうかを特定するステップと；

否定であれば、前記オーディオ信号は音声信号を備えない、と特定し、

肯定であれば、

エネルギーが前記所定の閾値より大きい前記短時間エネルギーフレーム内に少なくとも N 個の連続する短時間エネルギーフレームがあるときには前記オーディオ信号は音声信号を備える、と特定し、

エネルギーが前記所定の閾値より大きい前記短時間エネルギーフレーム内に少なくとも N 個の連続する短時間エネルギーフレームがないときには前記オーディオ信号は音声信号を備えない、と特定するステップと；を備える、

10

20

30

40

50

第 1 の局面に記載の方法。

[ 第 6 の局面 ]

音声信号検出装置であって：

オーディオ信号を取得するよう構成された取得モジュールと；

所定の音声信号の周波数に基づいて、前記オーディオ信号を複数の短時間エネルギーフレームに分割するよう構成された分割モジュールと；

各短時間エネルギーフレームのエネルギーを特定するよう構成された特定モジュールと；

各短時間エネルギーフレームの前記エネルギーに基づいて、前記オーディオ信号は音声信号を備えているかどうかを検出するよう構成された検出モジュールと；を備える、  
音声信号検出装置。

10

[ 第 7 の局面 ]

前記取得モジュールは、

現在のオーディオ信号を取得し、

最後に取得されたオーディオ信号から、所定の時間周期を持つ副信号を切り取り、

前記取得されたオーディオ信号として機能するよう、前記現在のオーディオ信号と前記切り取られた副信号とを連結するよう構成される、

第 6 の局面に記載の装置。

[ 第 8 の局面 ]

前記分割モジュールは、

前記所定の音声信号の周波数に基づいて前記所定の音声信号の周期を特定し、

前記特定された周期に基づいて、前記オーディオ信号を、対応する持続時間が前記周期である複数の短時間エネルギーフレームに分割するよう構成される、

第 6 の局面に記載の装置。

20

[ 第 9 の局面 ]

前記検出モジュールは、

エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の、全短時間エネルギーフレームの総量に対する比率を特定し、

前記比率が所定の比率より大きいかどうかを特定し、

肯定であれば、前記オーディオ信号は音声信号を備える、と特定し、

否定であれば、前記オーディオ信号は音声信号を備えない、と特定するよう構成される、

30

第 6 の局面に記載の装置。

[ 第 10 の局面 ]

前記検出モジュールは、

エネルギーが所定の閾値よりも大きい短時間エネルギーフレームの量の、全短時間エネルギーフレームの総量に対する比率を特定し、

前記比率が所定の比率より大きいかどうかを特定し、

否定であれば、前記オーディオ信号は音声信号を備えない、と特定し、

肯定であれば、

エネルギーが前記所定の閾値より大きい前記短時間エネルギーフレーム内に少なくとも N 個の連続する短時間エネルギーフレームがあるときには前記オーディオ信号は音声信号を備える、と特定し、

エネルギーが前記所定の閾値より大きい前記短時間エネルギーフレーム内に少なくとも N 個の連続する短時間エネルギーフレームがないときには前記オーディオ信号は音声信号を備えない、と特定するよう構成される、

第 6 の局面に記載の装置。

40

【 符号の説明 】

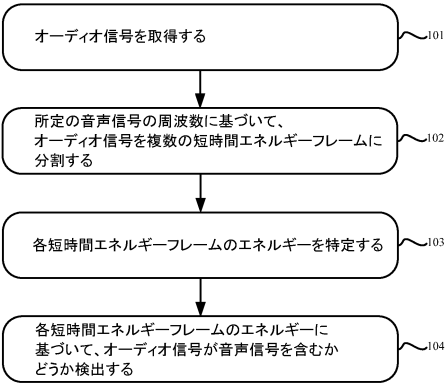
【 0 0 8 4 】

50

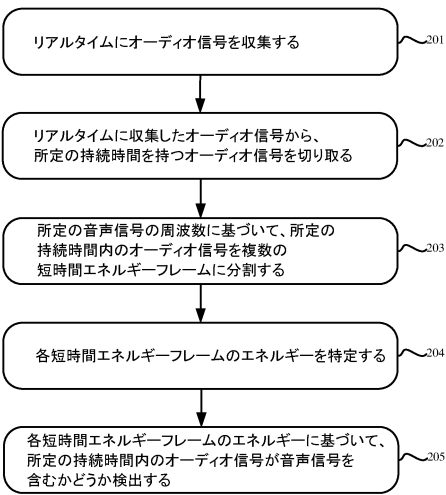


- 4 1 取得モジュール
- 4 2 分割モジュール
- 4 3 特定モジュール
- 4 4 検出モジュール

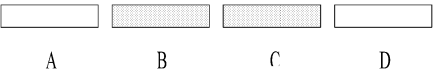
【図 1】



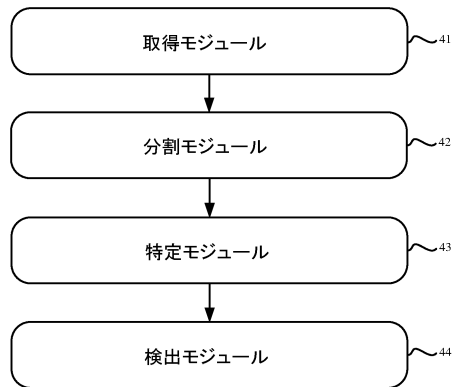
【図 2】



【図 3】



【図 4】



## フロントページの続き

(72)発明者 ジャオ, レイ

中華人民共和国, ディアンシアン 3 1 1 1 2 1, ハンヂョウ, ユ ハンディストリクト, ウェスト ウェン イー ロード ナンバー 9 6 9, ビルディング 3, 5 / エフ, アリババ グループ リーガル デパートメント

(72)発明者 グァン, イェンチュ

中華人民共和国, ディアンシアン 3 1 1 1 2 1, ハンヂョウ, ユ ハンディストリクト, ウェスト ウェン イー ロード ナンバー 9 6 9, ビルディング 3, 5 / エフ, アリババ グループ リーガル デパートメント

(72)発明者 ツァン, シャオドン

中華人民共和国, ディアンシアン 3 1 1 1 2 1, ハンヂョウ, ユ ハンディストリクト, ウェスト ウェン イー ロード ナンバー 9 6 9, ビルディング 3, 5 / エフ, アリババ グループ リーガル デパートメント

(72)発明者 リン, ファン

中華人民共和国, ディアンシアン 3 1 1 1 2 1, ハンヂョウ, ユ ハンディストリクト, ウェスト ウェン イー ロード ナンバー 9 6 9, ビルディング 3, 5 / エフ, アリババ グループ リーガル デパートメント

審査官 菊池 智紀

(56)参考文献 特開平10 - 301600 (JP, A)

特開2000 - 200100 (JP, A)

特表2013 - 508744 (JP, A)

(58)調査した分野(Int.Cl., DB名)

G10L 25/78 - 25/87, 15/04