

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7636088号  
(P7636088)

(45)発行日 令和7年2月26日(2025.2.26)

(24)登録日 令和7年2月17日(2025.2.17)

(51)国際特許分類 F I  
G 1 0 L 21/0264(2013.01) G 1 0 L 21/0264 C  
G 1 0 L 21/0232(2013.01) G 1 0 L 21/0232

請求項の数 15 (全37頁)

(21)出願番号	特願2023-527431(P2023-527431)	(73)特許権者	517392436
(86)(22)出願日	令和4年1月26日(2022.1.26)		騰 訊 科 技 ( 深 セ ン ) 有 限 公 司
(65)公表番号	特表2023-548707(P2023-548707 A)		TENCENT TECHNOLOGY (SHENZHEN) COMPANY LIMITED
(43)公表日	令和5年11月20日(2023.11.20)		中華人民共和國 5 1 8 0 5 7 広 東 省 深 セ ン 市 南 山 区 高 新 区 科 技 中 一 路 騰 訊 大 廈 3 5 層
(86)国際出願番号	PCT/CN2022/074003		3 5 / F , T e n c e n t B u i l d i n g , K e j i z h o n g y i R o a d , M i d w e s t D i s t r i c t o f H i - t e c h P a r k , N a n s h a n D i s t r i c t , S h e n z h e n , G u a n g d o n g 5 1 8
(87)国際公開番号	WO2022/166710		最終頁に続く
(87)国際公開日	令和4年8月11日(2022.8.11)		
審査請求日	令和5年5月8日(2023.5.8)		
(31)優先権主張番号	202110181389.4		
(32)優先日	令和3年2月8日(2021.2.8)		
(33)優先権主張国・地域又は機関	中国(CN)		

(54)【発明の名称】 音声強調方法、装置、機器及びコンピュータプログラム

(57)【特許請求の範囲】

【請求項1】

コンピュータ機器によって実行される、音声強調方法であって、  
 目標音声フレームの対応する複素スペクトルに基づいて前記目標音声フレームに対して  
 プリエンファシス処理を行い、第1複素スペクトルを得るステップと、  
 前記第1複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前  
 記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得るステップと、  
 前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目  
 標音声フレームの対応する強調音声信号を得るステップと  
 を含む、音声強調方法。

【請求項2】

目標音声フレームの対応する複素スペクトルに基づいて前記目標音声フレームに対して  
 プリエンファシス処理を行い、第1複素スペクトルを得る前記ステップは、  
 前記目標音声フレームの対応する複素スペクトルを第1ニューラルネットワークに入力  
 するステップであって、前記第1ニューラルネットワークはサンプル音声フレームの対応  
 する複素スペクトルと前記サンプル音声フレームにおける元の音声信号の対応する複素ス  
 ペクトルとに基づいてトレーニングを行って得られ、前記サンプル音声フレームは、前記  
 元の音声信号とノイズ信号とを組み合わせることにより得られる、ステップと、  
 前記第1ニューラルネットワークによって、前記目標音声フレームの対応する複素ス  
 ペクトルに基づいて前記第1複素スペクトルを出力するステップと

を含む、請求項 1 に記載の方法。

【請求項 3】

前記第 1 ニューラルネットワークは複素畳み込み層、ゲート付き回帰型ユニット層及び全結合層を含み、

前記第 1 ニューラルネットワークによって、前記目標音声フレームの対応する複素スペクトルに基づいて前記第 1 複素スペクトルを出力する前記ステップは、

前記複素畳み込み層によって前記目標音声フレームに対応する複素スペクトルにおける実部及び虚部に基づいて複素畳み込み処理を行うステップと、

前記ゲート付き回帰型ユニット層によって前記複素畳み込み層の出力に対して変換処理を行うステップと、

前記全結合層によって前記ゲート付き回帰型ユニットの出力に対して全結合処理を行い、前記第 1 複素スペクトルを出力するステップと

を含む、請求項 2 に記載の方法。

【請求項 4】

前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得る前記ステップは、

前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得るステップと、

前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して励起信号予測を行い、前記目標音声フレームの対応する励起信号を得るステップと、

前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームに対してゲイン予測を行い、前記目標音声フレームの対応するゲインを得るステップと

を含む、請求項 1 に記載の方法。

【請求項 5】

前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得る前記ステップは、

前記第 1 複素スペクトルを第 2 ニューラルネットワークに入力するステップであって、前記第 2 ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームの対応する声門パラメータとに基づいてトレーニングを行って得られるものである、ステップと、

前記第 2 ニューラルネットワークによって、前記第 1 複素スペクトルに基づいて前記目標音声フレームの対応する声門パラメータを出力するステップと

を含む、請求項 4 に記載の方法。

【請求項 6】

前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得る前記ステップは、

前記第 1 複素スペクトルと前記目標音声フレームの前の履歴音声フレームの対応する声門パラメータとを第 2 ニューラルネットワークに入力するステップであって、前記第 2 ニューラルネットワークはサンプル音声フレームの対応する複素スペクトル、サンプル音声フレームの前の履歴音声フレームの対応する声門パラメータ及びサンプル音声フレームの対応する声門パラメータに基づいてトレーニングを行って得られるものである、ステップと、

前記第 2 ニューラルネットワークによって、前記第 1 複素スペクトルと前記目標音声フレームの前の履歴音声フレームの対応する声門パラメータとに基づいて前記目標音声フレームの対応する声門パラメータを出力するステップと

を含む、請求項 4 に記載の方法。

【請求項 7】

前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音

10

20

30

40

50

声フレームに対してゲイン予測を行い、前記目標音声フレームの対応するゲインを得る前記ステップは、

前記目標音声フレームの前の履歴音声フレームの対応するゲインを第3ニューラルネットワークに入力するステップであって、前記第3ニューラルネットワークはサンプル音声フレームの前の履歴音声フレームの対応するゲインと前記サンプル音声フレームの対応するゲインとに基づいてトレーニングを行って得られるものである、ステップと、

前記第3ニューラルネットワークによって、前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームの対応するゲインを出力するステップと

を含む、請求項4に記載の方法。

10

【請求項8】

前記第1複素スペクトルに基づいて前記目標音声フレームに対して励起信号予測を行い、前記目標音声フレームの対応する励起信号を得る前記ステップは、

前記第1複素スペクトルを第4ニューラルネットワークに入力するステップであって、前記第4ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームに対応する励起信号の周波数領域表現とに基づいてトレーニングを行って得られるものである、ステップと、

前記第4ニューラルネットワークによって、前記第1複素スペクトルに基づいて前記目標音声フレームに対応する励起信号の周波数領域表現を出力するステップと

を含む、請求項4に記載の方法。

20

【請求項9】

前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得る前記ステップは、

声門フィルターにより前記目標音声フレームの対応する励起信号に対してフィルタリングを行い、フィルタリング出力信号を得るステップであって、前記声門フィルターは前記目標音声フレームの対応する声門パラメータに基づいて構築されるものである、ステップと、

前記目標音声フレームの対応するゲインに応じて前記フィルタリング出力信号に対して増幅処理を行い、前記目標音声フレームの対応する強調音声信号を得るステップと

を含む、請求項4に記載の方法。

30

【請求項10】

前記第1複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得る前記ステップは、

前記第1複素スペクトルに基づいてパワースペクトルを計算して取得するステップと、

前記パワースペクトルに基づいて自己相関係数を計算して取得するステップと、

前記自己相関係数に基づいて前記声門パラメータを計算して取得するステップと、

前記声門パラメータと前記自己相関係数とに基づいて前記ゲインを計算して取得するステップと、

前記ゲインと声門フィルターのパワースペクトルとに基づいて前記励起信号のパワースペクトルを計算して取得するステップであって、前記声門フィルターは前記声門パラメータに基づいて構築されるフィルターである、ステップと

を含む、請求項1に記載の方法。

40

【請求項11】

前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得る前記ステップは、

前記声門フィルターのパワースペクトルと前記励起信号のパワースペクトルとに基づいて第1振幅スペクトルを生成するステップと、

前記ゲインに応じて前記第1振幅スペクトルに対して増幅処理を行い、第2振幅スペクトルを得るステップと、

50

前記第 2 振幅スペクトルと前記第 1 複素スペクトル中から抽出された位相スペクトルとに基づいて、前記目標音声フレームの対応する強調音声信号を決定するステップとを含む、請求項 10 に記載の方法。

【請求項 12】

前記第 2 振幅スペクトルと前記第 1 複素スペクトル中から抽出された位相スペクトルとに基づいて、前記目標音声フレームの対応する強調音声信号を決定する前記ステップは、

前記第 2 振幅スペクトルと前記第 1 複素スペクトル中から抽出された位相スペクトルとを組み合わせ、第 2 複素スペクトルを得るステップと、

前記第 2 複素スペクトルを時間領域に変換し、前記目標音声フレームに対応する強調音声信号の時間領域信号を得るステップと

を含む、請求項 11 に記載の方法。

【請求項 13】

音声強調装置であって、

目標音声フレームの複素スペクトルに基づいて前記目標音声フレームに対してプリエンファシス処理を行い、第 1 複素スペクトルを得ることに用いられるプリエンファシスモジュールと、

前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得ることに用いられる音声分解モジュールと、

前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得ることに用いられる合成処理モジュールと

を含む、音声強調装置。

【請求項 14】

電子機器であって、

プロセッサと、

メモリであって、前記メモリ上にコンピュータ可読指令が記憶され、前記コンピュータ可読指令が前記プロセッサによって実行されるときに、請求項 1 ~ 12 のいずれか一項に記載の方法を実現するメモリと

を含む、電子機器。

【請求項 15】

コンピュータプログラムであって、プロセッサによって実行されるときに、請求項 1 ~ 12 のいずれか一項に記載の方法を実現する、コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本願は音声処理の技術分野に関し、具体的に言えば、音声強調方法、装置、機器及び記憶媒体に関する。

【0002】

本願は 2021 年 2 月 8 日に中国特許庁に提出された、出願番号が第 202110181389.4 号、発明の名称が「音声強調方法、装置、機器及び記憶媒体」である中国特許出願の優先権を主張し、その全内容は引用により本願に組み込まれている。

【背景技術】

【0003】

音声通信の利便性及び適時性により、音声通信の応用はますます幅広くなっており、たとえば、クラウド会議の会議参加者の間で音声信号が伝送される。ただし、音声通信においては、音声信号中にはノイズが混入される可能性があり、音声信号中に混入されるノイズが通信品質の劣化を招き、ユーザーの聴覚的体験に極めて大きな影響を与えることがある。従って、如何に音声に対して強調処理を行うことでノイズを除去するかは従来技術において早急に解決する技術的課題である。

【発明の概要】

10

20

30

40

50

**【発明が解決しようとする課題】****【0004】**

本願の実施例は音声強調方法、装置、機器及び記憶媒体を提供することで、音声強調を実現し、音声信号の品質を向上させる。

**【0005】**

本願のその他特性及び利点は以下の詳細な記述により明らかになるか、又は部分的に本願の実践により把握されて得られる。

**【課題を解決するための手段】****【0006】**

本願の実施例の一態様によれば、音声強調方法を提供し、目標音声フレームの対応する複素スペクトルに基づいて前記目標音声フレームに対してプリエンファシス処理を行い、第1複素スペクトルを得るステップと、前記第1複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得るステップと、前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得るステップとを含む。

10

**【0007】**

本願の実施例の別の一態様によれば、音声強調装置を提供し、目標音声フレームの複素スペクトルに基づいて前記目標音声フレームに対してプリエンファシス処理を行い、第1複素スペクトルを得ることに用いられるプリエンファシスモジュールと、前記第1複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得ることに用いられる音声分解モジュールと、前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得ることに用いられる合成処理モジュールとを含む。

20

**【0008】**

本願の実施例の別の一態様によれば、電子機器を提供し、プロセッサと、メモリであって、前記メモリ上にコンピュータ可読指令が記憶され、前記コンピュータ可読指令が前記プロセッサによって実行されるときに、上記に記載の音声強調方法を実現するメモリとを含む。

30

**【0009】**

本願の実施例の別の一態様によれば、コンピュータ可読記憶媒体を提供し、その上にコンピュータ可読指令が記憶され、前記コンピュータ可読指令がプロセッサによって実行されるときに、上記に記載の音声強調方法を実現する。

**【発明の効果】****【0010】**

本願の解決手段においては、まず目標音声フレームに対してプリエンファシスを行って第1複素スペクトルを得て、次に第1複素スペクトルを基礎として目標音声フレームに対して音声分解と合成を行い、2段階に分けて目標音声フレームに対して強調を行うことを実現するため、音声強調効果を効果的に保証することができる。そして、目標音声フレームに対してプリエンファシスを行って得られた第1複素スペクトルを基礎として、目標音声フレームに対して音声分解を行い、プリエンファシス前の目標音声フレームに比べて、第1複素スペクトルにおけるノイズの情報がより少なくなる。一方、音声分解過程において、ノイズが音声分解の正確性に影響を与えることがあり、従って、第1複素スペクトルを音声分解の基礎とすることで、音声分解の難度を低減させ、音声分解で得られた声門パラメータ、励起信号及びゲインの正確性を向上させ、さらに後続で取得された強調音声信号の正確性を保証することができる。そして、プリエンファシスで得られた第1複素スペクトル中には位相情報と振幅情報とが含まれ、該第1複素スペクトルにおける位相情報と振幅情報とを基礎として音声分解及び音声合成を行うことで、得られた目標音声フレームに対応する強調音声信号の振幅と位相の精度が保証されている。

40

50

## 【 0 0 1 1 】

理解すべきことは、以上の一般的な記述と後述の細部の記述は例示的で解釈的なものに過ぎず、本願を限定し得るものではないことである。

## 【 0 0 1 2 】

ここでの図面は、明細書に組み込まれ、且つ本明細書の一部を構成しており、本願にマッチングする実施例を示し、且つ明細書とともに本願の原理を解釈することに用いられる。明らかなように、以下の記述における図面は本願のいくつかの実施例に過ぎず、当業者にとって、創造的な労働を必要としない前提において、これらの図面に基づいてその他の図面を取得することもできる。図面において以下のとおりである。

## 【 図面の簡単な説明 】

## 【 0 0 1 3 】

【 図 1 】 1つの具体的な実施例に基づいて示される V o I P システムにおける音声通信リンクの模式図である。

【 図 2 】 音声信号が生じているデジタルモデルの模式図を示す。

【 図 3 】 1つの元の音声信号に基づいて励起信号と声門フィルタを分解する周波数応答の模式図を示す。

【 図 4 】 本願の一実施例に基づいて示される音声強調方法のフローチャートである。

【 図 5 】 1つの具体的な実施例に基づいて示される複素畳み込み層が複素数に対して畳み込み処理を行う模式図である。

【 図 6 】 1つの具体的な実施例に基づいて示される第 1 ニューラルネットワークの構造模式図である。

【 図 7 】 1つの具体的な実施例に基づいて示される第 2 ニューラルネットワークの模式図である。

【 図 8 】 別の一実施例に基づいて示される第 2 ニューラルネットワークの入力と出力の模式図である。

【 図 9 】 1つの具体的な実施例に基づいて示される第 3 ニューラルネットワークの模式図である。

【 図 1 0 】 1つの具体的な実施例に基づいて示される第 4 ニューラルネットワークの模式図である。

【 図 1 1 】 一実施例に基づいて示されるステップ 4 3 0 のフローチャートである。

【 図 1 2 】 1つの具体的な実施例に基づいて示される音声強調方法のフローチャートである。

【 図 1 3 】 一実施例に基づいて示されるステップ 4 2 0 のフローチャートである。

【 図 1 4 】 別の一実施例に基づいて示されるステップ 4 3 0 のフローチャートである。

【 図 1 5 】 別の 1 つの具体的な実施例に基づいて示される音声強調方法のフローチャートである。

【 図 1 6 】 1つの具体的な実施例に基づいて示される短時間フーリエ変換における窓掛け・オーバーラップの模式図である。

【 図 1 7 】 一実施例に基づいて示される音声強調装置のブロック図である。

【 図 1 8 】 本願の実施例を実現するための電子機器に適するコンピュータシステムの構造模式図を示す。

## 【 発明を実施するための形態 】

## 【 0 0 1 4 】

これより、図面を参照しながら例示的な実施形態をより全面的に記述する。しかしながら、例示的な実施形態は複数種の形式で実施でき、且つここで述べられた例に限定されると理解すべきでない。逆に、これらの実施形態の提供により、本願はより全面的で完全になり、且つ例示的な実施形態の発想は当業者に全面的に伝達される。

## 【 0 0 1 5 】

この他、記述される特徴、構造又は特性は、任意の適切な方式で 1 つ又はより多くの実施例に組み込まれてもよい。以下の記述において、多くの具体的な細部を提供することで

10

20

30

40

50

本願の実施例に対する十分な理解を与える。しかしながら、当業者は、特定の細部のうちの1つ又はより多くがなかったとしても、又はその他の方法、エレメント、装置、ステップ等を採用したとしても本願の技術的手段を實踐できることを認識することができる。その他の状況においては、公知の方法、装置、實現又は操作を詳細に示さない、又は記述しないことによって、本願の各態様を不明瞭にすることを回避する。

【0016】

図面において示されるブロック図は、単なる機能エンティティであり、必ずしも物理的に独立したエンティティに対応するわけではない。すなわち、ソフトウェアの形式を採用することでこれらの機能エンティティを實現する、又は1つ又は複数のハードウェアモジュール又は集積回路においてこれらの機能エンティティを實現する、又は異なるネットワーク及び/又はプロセッサ装置及び/又はマイクロ制御器装置においてこれらの機能エンティティを實現することができる。

10

【0017】

図面において示されるフローチャートは例示的な説明に過ぎず、必ずしもあらゆる内容と操作/ステップを含むわけではなく、必ずしも記述された順序で実行されるわけでもない。たとえば、ある操作/ステップはさらに分解でき、一方、ある操作/ステップは併せることができ、又は部分的に併せることができ、従って、実際に実行される順序は実際の状況に応じて変化する可能性がある。

【0018】

説明する必要がある点として、本明細書中に言及される「複数」は2つ又は2つ以上を指す。「及び/又は」は関連対象の関連関係を記述し、3種の関係が存在できることを表し、たとえば、A及び/又はBは、Aが単独で存在すること、AとBが同時に存在すること、Bが単独で存在することの3種の状況を表すことができる。文字「/」は一般的に前後の関連対象が「又は」の関係であることを表す。

20

【0019】

音声信号におけるノイズが、音声品質を極めて大きく低減させ、ユーザーの聴覚的体験に影響を与えることがあり、従って、音声信号の品質を向上させるために、音声信号に対して強調処理を行うことで、ノイズを最大限に除去し、信号における元の音声信号(すなわち、ノイズを含まない純粋な信号)を保留する必要がある。音声に対して強調処理を行うことを實現するために、本願の解決手段が提案されている。

30

【0020】

本願の解決手段は、音声通話の応用シーンにおいて適用でき、たとえば、インスタントメッセージングアプリケーションを介して行われる音声通信、ゲームアプリケーションにおける音声通話である。具体的には、音声の送信端、音声の受信端、又は音声通信サービスを提供するサーバ端末で本願の解決手段に従って音声強調を行うことができる。

【0021】

クラウド会議はオンライン業務実行における1つの重要な過程であり、クラウド会議において、クラウド会議の参加者の音収集装置が発言者の音声信号を収集した後に、収集された音声信号をその他の会議参加者に送信する必要がある。該過程に関わる音声信号は複数の参加者の間で伝送されて再生され、音声信号中に混入されたノイズ信号に対して処理を行われなければ、会議参加者の聴覚的体験に極めて大きな影響を与えることがある。このようなシーンにおいて、本願の解決手段を応用してクラウド会議中の音声信号に対して強調を行うことができ、これにより、会議参加者が聞き取る音声信号は強調された後の音声信号とすることができ、音声信号の品質を向上させることができる。

40

【0022】

クラウド会議は、クラウドコンピューティング技術に基づく高効率で、便利な、低コストの会議形式である。ユーザーはインターネットインターフェースを介して、簡単で使いやすい操作を行うだけで、迅速且つ高効率に世界的なチーム及び顧客と音声、データファイル及びビデオを同期して共有することができ、一方、会議中のデータの伝送、処理等の複雑な技術はクラウド会議サービス提供者がユーザーを補助することにより操作され得る。

50

## 【 0 0 2 3 】

現在、中国国内のクラウド会議は主に SaaS (Software as a Service、ソフトウェア・アズ・ア・サービス) モードを主体とするサービス内容に焦点を当てて、電話、ネットワーク、ビデオ等のサービス形式を含み、クラウドコンピューティングに基づくビデオ会議はクラウド会議と呼ばれる。クラウド会議の時代においては、データの伝送、処理、記憶はすべてビデオ会議提供者のコンピュータリソースにより処理され、ユーザーはさらに高価なハードウェアを購入したり煩雑なソフトウェアをインストールしたりする必要が全くなく、クライアント端末を開いて対応するインターフェースにアクセスするだけで、高効率な遠隔会議を行うことができる。

## 【 0 0 2 4 】

クラウド会議システムは、マルチサーバの動的クラスター配置をサポートし、且つ複数台の高性能サーバを提供し、会議の安定性、安全性、可用性を大幅に高める。近年、ビデオ会議はコミュニケーション効率を大幅に向上させ、コミュニケーションコストを連続的に低減させ、内部管理レベルのアップグレードをもたらしうることができるため、多くのユーザーに人気があり、すでに政府、軍隊、交通、輸送、金融、オペレータ、教育、企業等の各分野に幅広く応用されている。

## 【 0 0 2 5 】

図 1 は、1 つの具体的な実施例に基づいて示される VoIP (Voice over Internet Protocol、ネットワーク電話) システムにおける音声通信リンクの模式図である。図 1 に示すように、送信端 110 と受信端 120 のネットワーク接続に基づき、送信端 110 と受信端 120 は音声伝送を行うことができる。

## 【 0 0 2 6 】

図 1 に示すように、送信端 110 は収集モジュール 111、前強調処理モジュール 112 及び符号化モジュール 113 を含み、ここで、収集モジュール 111 は、音声信号を収集することに用いられ、それは収集した音響信号をデジタル信号に変換することができ、前強調処理モジュール 112 は、収集された音声信号に対して強調を行うことで、収集された音声信号中のノイズを除去し、音声信号の品質を向上させることに用いられる。符号化モジュール 113 は、強調された後の音声信号に対して符号化を行うことで、音声信号の伝送過程中の干渉抵抗性を向上させることに用いられる。前強調処理モジュール 112 は、本願の方法に従って音声強調を行い、音声に対して強調を行った後、さらに符号化圧縮及び伝送を行うことができ、このように、受信端が受信した信号がノイズに影響されなくなることを保証できる。

## 【 0 0 2 7 】

受信端 120 は復号モジュール 121、後強調モジュール 122 及び再生モジュール 123 を含む。復号モジュール 121 は受信した符号化音声信号に対して復号を行い、復号後の音声信号を得ることに用いられ、後強調モジュール 122 は復号後の音声信号に対して強調処理を行うことに用いられ、再生モジュール 123 は強調処理後の音声信号を再生することに用いられる。後強調モジュール 122 は本願の方法に従って音声強調を行うこともできる。いくつかの実施例では、受信端 120 はさらに音響効果調節モジュールを含んでもよく、該音響効果調節モジュールは強調された後の音声信号に対して音響効果調節を行うことに用いられる。

## 【 0 0 2 8 】

具体的な実施例において、受信端 120 のみ、又は送信端 110 のみで本願の方法に従って音声強調を行うことができ、もちろん、さらに送信端 110 と受信端 120 の両方で本願の方法に従って音声強調を行うこともできる。

## 【 0 0 2 9 】

いくつかの応用シーンにおいて、VoIP システムにおける端末機器は VoIP 通信をサポートできる以外に、さらにその他のサードパーティプロトコル、たとえば従来の PSTN (Public Switched Telephone Network、公共交換電話網) 回路ドメイン電話をサポートすることもできる。一方、従来の PSTN サービス

10

20

30

40

50

は音声強調を行うことができず、このようなシーンにおいては、受信端としての端末において本願の方法に従って音声強調を行うことができる。

【0030】

本願の解決手段に対して具体的な説明を行う前に、音声信号が生じるということについて説明を行う必要がある。音声信号は、人体の発音器官の脳制御における生理的運動によって生じるものであり、すなわち、気管のところで一定のエネルギーのノイズのような衝撃信号（励起信号に相当）が生じ、衝撃信号が人間の声帯（声帯が声門フィルターに相当）に衝撃を与え、略周期的な開閉が生じ、口腔を通じて増幅した後に、音を発する（音声信号を出力）。

【0031】

図2は、音声信号が生じているデジタルモデルの模式図を示しており、該デジタルモデルにより音声信号が生じる過程を記述することができる。図2に示すように、励起信号は声門フィルターに衝撃を与えた後、さらにゲイン制御を行って、その後音声信号を出力し、ここで、声門フィルターは声門パラメータにより限定される。該過程は下式で表すことができる。

$$x(n) = G \cdot r(n) \cdot a_r(n) \quad (\text{式1})$$

ここで、 $x(n)$ は入力された音声信号を表し、 $G$ はゲインを表し、線形予測ゲインと呼ばれることもでき、 $r(n)$ は励起信号を表し、 $a_r(n)$ は声門フィルターを表す。

【0032】

図3は、1つの元の音声信号に基づいて励起信号と声門フィルターを分解する周波数応答の模式図を示す。図3aは該元の音声信号の周波数応答の模式図を示し、図3bは該元の音声信号に基づいて分解された声門フィルターの周波数応答の模式図を示し、図3cは該元の音声信号に基づいて分解された励起信号の周波数応答の模式図を示す。図3に示すように、該元の音声信号の周波数応答の模式図における波形部分は声門フィルターの周波数応答の模式図におけるピーク位置に対応し、励起信号は該元の音声信号に対してLP（Linear Prediction、線形予測）分析を行った後の残差信号に相当し、従って、その対応する周波数応答が比較的緩やかである。

【0033】

上記からわかるように、1つの元の音声信号（すなわち、ノイズを含まない音声信号）に基づいて励起信号、声門フィルター及びゲインを分解することができ、分解された励起信号、声門フィルター及びゲインは該元の音声信号を表現することに用いられてもよく、ここで、声門フィルターは声門パラメータにより表現できる。逆に、1つの元の音声信号の対応する励起信号、声門フィルターを決定することに用いられる声門パラメータ及びゲインが知られていれば、対応する励起信号、声門フィルター及びゲインに基づいて該元の音声信号を再構成することができる。

【0034】

本願の解決手段は、該原理に基づき、音声フレームの対応する声門パラメータ、励起信号及びゲインに基づいて該音声フレームにおける元の音声信号を再構成し、音声強調を実現することである。

【0035】

以下、本願の実施例の技術的手段を詳細に述べる。

【0036】

図4は、本願の一実施例に基づいて示される音声強調方法のフローチャートであり、該方法は処理能力を備えるコンピュータ機器により実行されてもよく、たとえば、端末、サーバ等であり、ここで具体的な限定を行わない。図4に示されるものを参照すると、該方法は少なくともステップ410～430を含み、以下のように詳細に説明される。

【0037】

ステップ410：目標音声フレームの対応する複素スペクトルに基づいて前記目標音声フレームに対してプリエンファシス処理を行い、第1複素スペクトルを得る。

【0038】

10

20

30

40

50

音声信号は緩やかでランダムに変化するのではなく経時的に変化するものであるが、短時間で音声信号が強い相関を有する、すなわち、音声信号が短時間相関性を有する。従って、本願の解決手段において、音声フレームを単位として音声強調を行う。目標音声フレームとは現在の強調処理対象の音声フレームを指す。

【0039】

目標音声フレームの対応する複素スペクトルは該目標音声フレームの時間領域信号に対して時間周波数変換を行うことにより取得することができ、時間周波数変換はたとえば短時間フーリエ変換 (Short-term Fourier transform、STFT) であってもよい。目標音声フレームの対応する複素スペクトルにおける実部の係数は該目標音声フレームの振幅情報を指示することに用いられ、虚部の係数は目標音声フレームの位相情報を指示することに用いられる。

10

【0040】

目標音声フレームに対してプリアンファシスを行うことにより、目標音声フレームにおける一部のノイズを除去することができ、従って、目標音声フレームの対応する複素スペクトルに比べて、プリアンファシスで得られた第1複素スペクトルにおけるノイズ含有量がより少ない。

【0041】

本願のいくつかの実施例では、深層学習の方式を採用して目標音声フレームに対してプリアンファシスを行うことができる。1つのニューラルネットワークモデルをトレーニングすることにより、音声フレームの対応する複素スペクトルに基づいて音声フレームにおけるノイズの複素スペクトルを予測し、次に音声フレームの複素スペクトルと予測されたノイズの複素スペクトルとを減算し、第1複素スペクトルを得る。記述の便宜のために、音声フレームにおけるノイズの複素スペクトルを予測することに用いられる該ニューラルネットワークモデルをノイズ予測モデルと呼ぶ。トレーニング終了後に、該ノイズ予測モデルは入力された音声フレームの複素スペクトルに基づいて予測されたノイズの複素スペクトルを出力することができ、次に音声フレームの複素スペクトルとノイズの複素スペクトルとを減算すると、第1複素スペクトルを得られる。

20

【0042】

本願のいくつかの実施例では、さらに1つのニューラルネットワークモデルをトレーニングすることで、音声フレームの複素スペクトルに基づいて強調された後の該音声フレームの第1複素スペクトルを予測することができる。記述の便宜のために、強調された後の複素スペクトルを予測することに用いられる該ニューラルネットワークモデルを強調複素スペクトル予測モデルと呼ぶ。トレーニング過程において、サンプル音声フレームの複素スペクトルを該強調複素スペクトル予測モデル中に入力し、該強調複素スペクトル予測モデルによって強調された後の複素スペクトルを予測し、且つ予測された強調された後の複素スペクトルと該サンプル音声フレームのラベル情報とに基づいて強調複素スペクトル予測モデルのパラメータを調整し、予測された強調された後の複素スペクトルとラベル情報が指示した複素スペクトルとの間の差異が所定の要件を満たすまで続ける。サンプル音声フレームのラベル情報はサンプル音声フレームにおける元の音声信号の複素スペクトルを指示することに用いられる。トレーニング終了後に、該強調複素スペクトル予測モデルは目標音声フレームの複素スペクトルに基づいて第1複素スペクトルを出力することができる。

30

【0043】

ステップ420：前記第1複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得る。

【0044】

音声分解で得られた目標音声フレームの対応する声門パラメータ、対応するゲイン及び対応する励起信号は、図2に示される過程に従って目標音声フレームにおける元の音声信号を再構成することに用いられる。

40

50

## 【 0 0 4 5 】

上記の記述のように、1つの元の音声信号は、励起信号が声門フィルタに衝撃を与えてからゲイン制御を行うことにより得られるものである。該第1複素スペクトル中には目標音声フレームの元の音声信号の情報が含まれており、従って、該第1複素スペクトルに基づき線形予測分析を行い、目標音声フレームにおける元の音声信号を再構成することに用いられる声門パラメータ、励起信号及びゲインを逆方向に決定する。

## 【 0 0 4 6 】

声門パラメータとは、声門フィルタを構築することに用いられるパラメータを指し、声門パラメータが決定されると、声門フィルタが対応して決定され、声門フィルタはデジタルフィルタである。声門パラメータは線形予測符号化 (Linear Prediction Coefficients、LPC) 係数であってもよく、さらに線スペクトル周波数 (Line Spectral Frequency、LSF) パラメータであってもよい。目標音声フレームに対応する声門パラメータの数量は声門フィルタの次数に関連しており、前記声門フィルタがK次フィルタである場合、前記声門パラメータはK次LSFパラメータ又はK次LPC係数を含み、ここで、LSFパラメータとLPC係数との間が相互に転換することができる。

10

## 【 0 0 4 7 】

1つのp次の声門フィルタは、

$A_p(z) = 1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_p z^{-p}$  (式2)として表されてもよい。

20

ここで、 $a_1$ 、 $a_2$ 、 $\dots$ 、 $a_p$ はLPC係数であり、pは声門フィルタの次数であり、zは声門フィルタの入力信号である。

## 【 0 0 4 8 】

式2を基礎として、

$$P(z) = A_p(z) - z^{-(p+1)} A_p(z^{-1}) \quad (\text{式3})$$

$Q(z) = A_p(z) + z^{-(p+1)} A_p(z^{-1})$  (式4)のように設定する場合、以下[数1] (式5)を得ることができる。

## 【 0 0 4 9 】

## 【数1】

$$A_p(z) = \frac{P(z) + Q(z)}{2}$$

30

## 【 0 0 5 0 】

物理的には、 $P(z)$ と $Q(z)$ は、それぞれ声門開放と声門閉鎖の周期的な変化規律を代表する。多項式 $P(z)$ と $Q(z)$ の根は複素平面上で交互に出現し、それは複素平面単位円上に分布する一連の角周波数であり、LSFパラメータはすなわち $P(z)$ と $Q(z)$ の根の複素平面単位円上の対応する角周波数であり、第nフレームの音声フレームの対応するLSFパラメータ $LSF(n)$ はnとして表されてもよい。もちろん、第nフレームの音声フレームの対応するLSFパラメータ $LSF(n)$ はさらに該第nフレームの音声フレームに対応する $P(z)$ の根と対応する $Q(z)$ 根で直接的に示されることができる。

40

## 【 0 0 5 1 】

第nフレームの音声フレームに対応する $P(z)$ と $Q(z)$ の複素平面での根を $n$ として定義すると、第nフレームの音声フレームの対応するLSFパラメータは、

以下[数2] (式6)として表される。

## 【 0 0 5 2 】

## 【数2】

50

$$\omega_n = \tan^{-1} \left( \frac{\text{Re}\{\theta_n\}}{\text{Im}\{\theta_n\}} \right)$$

## 【 0 0 5 3 】

ここで、 $\text{Re}\{\theta_n\}$  は複素数  $\theta_n$  の実部を表し、 $\text{Im}\{\theta_n\}$  は複素数  $\theta_n$  の虚部を表す。

## 【 0 0 5 4 】

本願のいくつかの実施例では、深層学習の方式を採用して音声分解を行うことができる。まず、それぞれ声門パラメータ予測を行うこと、励起信号予測を行うこと、及びゲイン予測を行うことに用いられるニューラルネットワークモデルをトレーニングすることができ、該3つのニューラルネットワークモデルが第1複素スペクトルに基づき目標音声フレームの対応する声門パラメータ、励起信号及びゲインをそれぞれ予測できるようにする。

10

## 【 0 0 5 5 】

本願のいくつかの実施例では、さらに線形予測分析の原理に従って、第1複素スペクトルに基づいて信号処理を行い、且つ目標音声フレームの対応する声門パラメータ、励起信号及びゲインを計算することができ、具体的な過程は下記の記述を参照する。

## 【 0 0 5 6 】

ステップ430：前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得る。

20

## 【 0 0 5 7 】

目標音声フレームの対応する声門パラメータが決定される場合に、その対応する声門フィルターは対応して決定される。それを基に、図2に示される元の音声信号の生成過程に基づいて、目標音声フレームの対応する励起信号が決定される声門フィルターに衝撃を与え、且つ目標音声フレームの対応するゲインに応じてフィルタリングで得られた信号に対してゲイン制御を行うことにより、元の音声信号の再構成を実現することができ、再構成で取得された信号はすなわち目標音声フレームの対応する強調音声信号である。

## 【 0 0 5 8 】

本願の解決手段において、まず、目標音声フレームに対してプリアンファシスを行って第1複素スペクトルを得て、次に第1複素スペクトルを基礎として目標音声フレームに対して音声分解と合成を行い、2段階に分けて目標音声フレームに対して強調を行うことを実現し、音声強調効果を効果的に保証することができる。そして、目標音声フレームに対してプリアンファシスを行って得られた第1複素スペクトルを基礎として、目標音声フレームに対して音声分解を行い、目標音声フレームがプリアンファシスされる前のスペクトルに比べて、第1複素スペクトルにおけるノイズの情報がより少なくなる。音声分解過程においては、ノイズが音声分解の正確性に影響を与えることがあり、従って、第1複素スペクトルを音声分解の基礎とすることで、音声分解の難度を低減させ、音声分解で得られた声門パラメータ、励起信号及びゲインの正確性を向上させ、さらに後続で取得された強調音声信号の正確性を保証することができる。プリアンファシスで得られた第1複素スペクトル中には位相情報と振幅情報が含まれ、該第1複素スペクトルにおける位相情報と振幅情報を基礎として音声分解及び音声合成を行うことで、得られた目標音声フレームに対応する強調音声信号の振幅と位相の精度が保証されている。

30

40

## 【 0 0 5 9 】

本願のいくつかの実施例では、ステップ410は、前記目標音声フレームの対応する複素スペクトルを第1ニューラルネットワークに入力するステップであって、前記第1ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームにおける元の音声信号の対応する複素スペクトルとに基づいてトレーニングを行って得られるものである、ステップと、前記第1ニューラルネットワークによって、前記目標音声フレームの対応する複素スペクトルに基づいて前記第1複素スペクトルを出力するステップとを含む。

50

## 【 0 0 6 0 】

第1ニューラルネットワークは、長・短期記憶ニューラルネットワーク、畳み込みニューラルネットワーク、回帰型ニューラルネットワーク、全結合ニューラルネットワーク、ゲート付き回帰型ユニット等により構築されたモデルであってもよく、ここで具体的な限定を行わない。

## 【 0 0 6 1 】

本願のいくつかの実施例では、サンプル音声信号に対してフレーム分割を行うことにより、複数のサンプル音声フレームを得ることができる。ここで、サンプル音声信号は、知られている元の音声信号と知られているノイズ信号とを組み合わせることにより得ることができ、このように、元の音声信号が知られている場合に、対応してサンプル音声フレームにおける元の音声信号に対して時間周波数変換を行って、サンプル音声フレームにおける元の音声信号の対応する複素スペクトルを得ることができる。サンプル音声フレームの対応する複素スペクトルは、該サンプル音声フレームの時間領域信号に対して時間周波数変換を行うことにより得ることができる。

10

## 【 0 0 6 2 】

トレーニング過程において、サンプル音声フレームの対応する複素スペクトルを第1ニューラルネットワークに入力し、第1ニューラルネットワークによって、サンプル音声フレームの対応する複素スペクトルに基づいて予測を行い、予測された第1複素スペクトルを出力し、次に予測された第1複素スペクトルと該サンプル音声フレームにおける元の音声信号の対応する複素スペクトルとを比較し、両方の間の類似度が所定の要件を満たさなければ、第1ニューラルネットワークのパラメータを調整し、第1ニューラルネットワークが出力した予測された第1複素スペクトルと該サンプル音声フレームにおける元の音声信号の対応する複素スペクトルとの間の類似度が所定の要件を満たすまで続ける。ここで、該所定の要件は、予測された第1複素スペクトルと該サンプル音声フレームにおける元の音声信号の対応する複素スペクトルとの間の類似度が類似度閾値以上であることであってもよく、該類似度閾値はニーズに応じて設定を行うことができ、たとえば、100%、98%等である。上記のようなトレーニング過程により、該第1ニューラルネットワークは入力された複素スペクトルに基づいて第1複素スペクトルを予測する能力を学習することができる。

20

## 【 0 0 6 3 】

本願のいくつかの実施例では、前記第1ニューラルネットワークは複素畳み込み層、ゲート付き回帰型ユニット層及び全結合層を含む。上記した前記第1ニューラルネットワークによって、前記目標音声フレームの複素スペクトルに基づいて前記第1複素スペクトルを出力するステップは、さらに、前記複素畳み込み層によって前記目標音声フレームに対応する複素スペクトルにおける実部及び虚部に基づいて複素畳み込み処理を行うステップと、前記ゲート付き回帰型ユニット層によって前記複素畳み込み層の出力に対して変換処理を行うステップと、前記全結合層によって前記ゲート付き回帰型ユニットの出力に対して全結合処理を行い、前記第1複素スペクトルを出力するステップとを含む。

30

## 【 0 0 6 4 】

具体的な実施例において、第1ニューラルネットワークは1層又は複数層の複素畳み込み層を含んでもよく、同様に、ゲート付き回帰型ユニット層と全結合層も1層又は複数層であってもよく、具体的には、複素畳み込み層、ゲート付き回帰型ユニット層及び全結合層の数量は実際のニーズに応じて設定を行うことができる。

40

## 【 0 0 6 5 】

図5は、1つの具体的な実施例に基づいて示される複素畳み込み層が複素数に対して畳み込み処理を行う模式図であり、複素畳み込み層の入力複素数が  $E + jF$  であり、複素畳み込み層の加重が  $A + jB$  であると仮定する。図5に示すように、複素畳み込み層は2次元畳み込み層 (Real\_conv、Imag\_conv)、結合層 (Concat) 及び活性化層 (Leaky\_ReLU) を含む。入力複素数中の実部  $E$  と虚部  $F$  とを2次元畳み込み層に入力した後に、該2次元畳み込み層は複素畳み込み層の加重に応じて畳み込

50

みを行い、それが畳み込み演算を行う過程は下式で示される。

$$(E + jF) * (A + jB) = (E * A - F * B) + j(E * B + F * A) \quad (\text{式7})$$

$C = E * A - F * B$ 、 $D = E * B + F * A$ に設定する場合、上式7はさらに、

$$(E + jF) * (A + jB) = C + jD \quad (\text{式8}) \text{に転換する。}$$

【0066】

図5に示すように、2次元畳み込み層が畳み込まれた後の実部と虚部を出力した後に、結合層によって実部と虚部とを結合し、結合結果を得て、次に、活性化層によって結合結果に対して活性化を行う。図5において、活性化層に使用された活性化関数が `Leaky_ReLU` 活性化関数である。`Leaky_ReLU` 活性化関数の表現式は、 $f(x) = \max(ax, x)$  ( $a$ が定数である) (式9)である。

10

【0067】

その他の実施例において、活性化層に使用された活性化関数はさらにその他の関数、たとえば `zReLU` 関数等であってもよく、ここで具体的な限定を行わない。

【0068】

図6は、1つの具体的な実施例に基づいて示される第1ニューラルネットワークの構造模式図であり、図6に示すように、該第1ニューラルネットワークは、順にカスケード接続された6層の複素畳み込み層 (`Conv`)、1層のゲート付き回帰型ユニット (`Gated Recurrent Unit`、`GRU`) 層及び2層の全結合 (`Full Connected`、`FC`) 層を含む。目標音声フレームに対応する複素スペクトル  $S(n)$  を該第1ニューラルネットワークに入力した後に、まず6層の複素畳み込み層によって順に複素畳み込み処理を行い、次に `GRU` 層によって変換を行い、さらに2層の `FC` 層によって順次に全結合を行い、且つ最後の1層の `FC` 層によって第1複素スペクトルを出力する。ここで、各層の括弧内の数字は該層が出力した変数の次元を表す。図6に示される第1ニューラルネットワークにおいて、最後の1層の `FC` 層が出力した次元は322次元であり、161個の `STFT` 係数中の実部と虚部を示すことに用いられる。

20

【0069】

本願のいくつかの実施例では、ステップ420は、前記第1複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得るステップと、前記第1複素スペクトルに基づいて前記目標音声フレームに対して励起信号予測を行い、前記目標音声フレームの対応する励起信号を得るステップと、前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームに対してゲイン予測を行い、前記目標音声フレームの対応するゲインを得るステップとを含む。

30

【0070】

本願のいくつかの実施例では、声門パラメータ予測を行うことに用いられるニューラルネットワークモデル(第2ニューラルネットワークとして仮定)、ゲイン予測を行うニューラルネットワークモデル(第3ニューラルネットワークとして仮定)、及び励起信号予測を行うニューラルネットワークモデル(第4ニューラルネットワークとして仮定)をそれぞれトレーニングすることができる。ここで、該3種のニューラルネットワークモデルは長・短期記憶ニューラルネットワーク、畳み込みニューラルネットワーク、回帰型ニューラルネットワーク、全結合ニューラルネットワーク等により構築されたモデルであってもよく、ここで具体的な限定を行わない。

40

【0071】

本願のいくつかの実施例では、上記した前記第1複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得るステップは、さらに、前記第1複素スペクトルを第2ニューラルネットワークに入力するステップであって、前記第2ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームの対応する声門パラメータとに基づいてトレーニングを行って得られるものである、ステップと、前記第2ニューラルネットワークによって、前記第1複素スペクトルに基づいて前記目標音声フレームの対応

50

する声門パラメータを出力するステップとを含む。

【0072】

サンプル音声フレームの対応する複素スペクトルは、サンプル音声フレームの時間領域信号に対して時間周波数変換を行うことにより得られるものである。本願のいくつかの実施例では、サンプル音声信号に対してフレーム分割を行い、複数のサンプル音声フレームを得ることができる。サンプル音声信号は知られている元の音声信号と知られているノイズ信号とを組み合わせることにより得ることができる。このように、元の音声信号が知られている場合に、元の音声信号に対して線形予測分析を行うことによりサンプル音声フレームの対応する声門パラメータを得ることができ、換言すれば、サンプル音声フレームの対応する声門パラメータとはサンプル音声フレームにおける元の音声信号を再構成することにより得られる声門パラメータを指す。

10

【0073】

トレーニング過程においては、サンプル音声フレームの複素スペクトルを第2ニューラルネットワークに入力した後に、第2ニューラルネットワークによって、サンプル音声フレームの複素スペクトルに基づいて声門パラメータ予測を行い、予測声門パラメータを出力し、次に、予測声門パラメータと該サンプル音声フレームの対応する声門パラメータとを比較し、両方が一致しなければ、第2ニューラルネットワークのパラメータを調整し、第2ニューラルネットワークがサンプル音声フレームの複素スペクトルに基づいて出力した予測声門パラメータが該サンプル音声フレームの対応する声門パラメータと一致するまで続ける。トレーニング終了後に、該第2ニューラルネットワークは、入力された音声フレームの複素スペクトルに基づいて該音声フレームにおける元の音声信号を再構成することにより得られる声門パラメータを正確に予測する能力を学習している。

20

【0074】

図7は、1つの具体的な実施例に基づいて示される第2ニューラルネットワークの模式図である。図7に示すように、該第2ニューラルネットワークは、1層のLSTM (Long-Short Term Memory、長・短期記憶ネットワーク)層と3層のカスケード接続されたFC (Full Connected、全結合)層とを含む。ここで、LSTM層は1つの隠れ層であり、それは256個のユニットを含み、LSTM層の入力は第nフレームの音声フレームの対応する第1複素スペクトル $S'(n)$ である。本実施例において、LSTM層の入力は321次元である。3層のカスケード接続されたFC層において、前の2層のFC層中には活性化関数( )が設定され、設定された活性化関数は第2ニューラルネットワークの非線形発現能力を増加することに用いられ、最後の1層のFC層中には活性化関数が設定されず、該最後の1層のFC層は分類器として分類出力を行う。図7に示すように、入力から出力への方向に沿って、3層のFC層中にはそれぞれ512、512、16個のユニットが含まれ、最後の1層のFC層の出力は該第nフレームの音声フレームに対応する16次元の線スペクトル周波数係数LSF(n)、すなわち16次元線スペクトル周波数パラメータである。

30

【0075】

本願のいくつかの実施例では、音声フレームの間に相関性があり、隣接する2つの音声フレームの間の周波数領域特徴の類似性が比較的高く、従って、目標音声フレームの前の履歴音声フレームの対応する声門パラメータと組み合わせて目標音声フレームの対応する声門パラメータを予測することができる。一実施例において、上記した前記第1複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得るステップは、さらに、前記第1複素スペクトルと前記目標音声フレームの前の履歴音声フレームの対応する声門パラメータとを第2ニューラルネットワークに入力するステップであって、前記第2ニューラルネットワークはサンプル音声フレームの対応する複素スペクトル、サンプル音声フレームの前の履歴音声フレームの対応する声門パラメータ及びサンプル音声フレームの対応する声門パラメータに基づいてトレーニングを行って得られるものである、ステップと、前記第1ニューラルネットワークによって、前記第1複素スペクトルと前記目標音声フレームの前の履歴音声

40

50

フレームの対応する声門パラメータとに基づいて前記目標音声フレームの対応する声門パラメータを出力するステップとを含む。

【0076】

履歴音声フレームと目標音声フレームとの間に相関性があり、目標音声フレームの履歴音声フレームに対応する声門パラメータと目標音声フレームの対応する声門パラメータとの間に類似性があるため、目標音声フレームの履歴音声フレームの対応する声門パラメータを参照として、目標音声フレームの声門パラメータの予測過程に対して監視を行うことで、声門パラメータ予測の正確率を向上させることができる。

【0077】

本願のいくつかの実施例では、音声フレームが時間的により近いほど声門パラメータの類似性がより高いため、目標音声フレームに比較的近い履歴音声フレームの対応する声門パラメータを参照とすることで、予測正確率をさらに保証することができ、たとえば、目標音声フレームの直前音声フレームの対応する声門パラメータを参照とすることができる。具体的な実施例において、参照としての履歴音声フレームの数量は1フレームであってもよく、又はマルチフレームであってもよく、具体的には、実際のニーズに応じて選択して用いることができる。

10

【0078】

目標音声フレームの履歴音声フレームに対応する声門パラメータは該履歴音声フレームに対して声門パラメータ予測を行うことにより得られた声門パラメータであってもよい。換言すれば、声門パラメータの予測過程において、履歴音声フレームについて予測された声門パラメータを現在の音声フレームの声門パラメータ予測過程の参照として多重化する。

20

【0079】

本実施例における第2ニューラルネットワークのトレーニング過程は、前の一実施例における第2ニューラルネットワークのトレーニング過程に類似しており、ここではトレーニングの過程を繰り返し説明しない。

【0080】

図8は、別の一実施例に基づいて示される第2ニューラルネットワークの入力と出力の模式図である。ここで、図8における第2ニューラルネットワークの構造は図7におけるものと同じであり、図7と比べて、図8における第2ニューラルネットワークの入力は、さらに該第nフレームの音声フレームの直前音声フレーム(すなわち第n-1フレーム)の線スペクトル周波数パラメータLSF(n-1)を含む。図8に示すように、第2層のFC層中に第nフレームの音声フレームの直前音声フレームの線スペクトル周波数パラメータLSF(n-1)を埋め込んで参照情報とする。隣接する2つの音声フレームのLSFパラメータの類似性が非常に高く、従って、第nフレームの音声フレームの履歴音声フレームの対応するLSFパラメータを参照情報とすれば、LSFパラメータの予測正確率を高めることができる。

30

【0081】

本願のいくつかの実施例では、上記した前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームに対してゲイン予測を行い、前記目標音声フレームの対応するゲインを得るステップは、さらに、前記目標音声フレームの前の履歴音声フレームの対応するゲインを第3ニューラルネットワークに入力するステップであって、前記第3ニューラルネットワークはサンプル音声フレームの前の履歴音声フレームの対応するゲインと前記サンプル音声フレームの対応するゲインとに基づいてトレーニングを行って得られるものである、ステップと、前記第3ニューラルネットワークによって、前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームの対応するゲインを出力するステップとを含むことができる。

40

【0082】

目標音声フレームの履歴音声フレームの対応するゲインは、該第3ニューラルネットワークが該履歴音声フレームのゲイン予測を行うことにより得られるのもであってもよく、換言すれば、履歴音声フレームについて予測されたゲインを目標音声フレームに対してゲ

50

イン予測を行う過程における第3ニューラルネットワークモデルの入力として多重化する。

【0083】

サンプル音声フレームはサンプル音声信号に対してフレーム分割を行うことにより得られてもよく、サンプル音声信号は知られている元の音声信号と知られているノイズ信号とを組み合わせることにより得ることができる。このようにして、サンプル音声中の元の音声信号が知られている場合に、該元の音声信号に対して線形予測分析を行って、該元の音声信号を再構成することに用いられる声門パラメータ、すなわちサンプル音声フレームの対応する声門パラメータを得ることができる。

【0084】

図9は、1つの具体的な実施例に基づいて示される第3ニューラルネットワークの模式図である。図9に示すように、第3ニューラルネットワークは1層のLSTM層と1層のFC層とを含み、ここで、LSTM層は1つの隠れ層であり、それは128個のユニットを含み、FC層の入力の次元が512であり、出力が1次元のゲインである。1つの具体的な実施例において、第nフレームの音声フレームの履歴音声フレームの対応するゲイン  $G\_pre(n)$  は第nフレームの音声フレームの最初の4つ音声フレームに対応するゲインとして定義することができ、すなわち、  
 $G\_pre(n) = \{G(n-1), G(n-2), G(n-3), G(n-4)\}$  である。

10

【0085】

もちろん、ゲイン予測に用いられるものとして選択された履歴音声フレームの数量は上記のような例に限定されず、具体的には、実際のニーズに応じて選択して用いることができる。

20

【0086】

上記のように示される第2ニューラルネットワークと第3ニューラルネットワークは全体的に  $M - to - N$  のマッピング関係 ( $N < M$ ) を呈し、すなわち、ニューラルネットワークモデルの入力情報の次元が  $M$  であり、出力情報の次元が  $N$  であり、ニューラルネットワークモデルの構造を極めて大きく簡略化して、ニューラルネットワークモデルの複雑さを低減させている。

【0087】

本願のいくつかの実施例では、上記した前記第1複素スペクトルに基づいて前記目標音声フレームに対して励起信号予測を行い、前記目標音声フレームの対応する励起信号を得るステップは、さらに、前記第1複素スペクトルを第4ニューラルネットワークに入力するステップであって、前記第4ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームに対応する励起信号の周波数領域表現とに基づいてトレーニングを行って得られるものである、ステップと、前記第4ニューラルネットワークによって、前記第1複素スペクトルに基づいて前記目標音声フレームに対応する励起信号の周波数領域表現を出力するステップとを含むことができる。

30

【0088】

サンプル音声フレームの対応する励起信号は、サンプル音声フレームにおける知られている元の音声信号に対して線形予測分析を行うことにより得られるものであってもよい。周波数領域表現は振幅スペクトルであってもよく、又は複素スペクトルであってもよく、ここで具体的な限定を行わない。

40

【0089】

第4ニューラルネットワークをトレーニングする過程において、サンプル音声フレームの複素スペクトルを第4ニューラルネットワークモデル中に入力し、次に第4ニューラルネットワークによって、入力されたサンプル音声フレームの複素スペクトルに基づいて励起信号予測を行い、予測励起信号の周波数領域表現を出力し、次に予測励起信号の周波数領域表現と該サンプル音声フレームに対応する励起信号の周波数領域表現とに基づいて第4ニューラルネットワークのパラメータを調整する。すなわち、予測励起信号の周波数領域表現と該サンプル音声フレームに対応する励起信号の周波数領域表現との類似度が所定

50

の要件を満たさなければ、第4ニューラルネットワークのパラメータを調整し、第4ニューラルネットワークがサンプル音声フレームについて出力された予測励起信号の周波数領域表現と該サンプル音声フレームに対応する励起信号の周波数領域表現との間の類似度が所定の要件を満たすまで続ける。上記のようなトレーニング過程により、第4ニューラルネットワークに、音声フレームの振幅スペクトルに基づいて該音声フレームの対応する励起信号の周波数領域表現を予測する能力を学習させることができ、それにより励起信号の予測を正確に行う。

【0090】

図10は、1つの具体的な実施例に基づいて示される第4ニューラルネットワークの模式図である。図10に示すように、該第4ニューラルネットワークは、1層のLSTM層と3層のFC層を含み、ここで、LSTM層は1つの隠れ層であり、256個のユニットを含み、LSTMの入力は第nフレームの音声フレームの対応する第1複素スペクトル $S(n)$ であり、その次元が321次元であってもよい。3層のFC層中に含まれるユニットの数量はそれぞれ512、512及び321であり、最後の1層のFC層は321次元の第nフレームの音声フレームに対応する励起信号の周波数領域表現 $R(n)$ を出力する。入力から出力へ方向に沿って、3層のFC層のうちの最初の2層のFC層中に活性化関数が設定され、モデルの非線形発現能力を高めることに用いられ、最後の1層のFC層中に活性化関数がなく、分類出力を行うことに用いられる。

10

【0091】

上記に示される第1ニューラルネットワーク、第2ニューラルネットワーク、第3ニューラルネットワーク及び第4ニューラルネットワークの構造は単に例示的なものであり、その他の実施例において、深層学習のオープンソースプラットフォーム中に相応な構造のニューラルネットワークモデルを設置し、且つ対応してトレーニングを行うこともできる。

20

【0092】

本願のいくつかの実施例では、図11に示すように、ステップ430は、ステップ1110とステップ1120を含み、

【0093】

ステップ1110：声門フィルターにより前記目標音声フレームの対応する励起信号に対してフィルタリングを行い、フィルタリング出力信号を得る。前記声門フィルターは前記目標音声フレームの対応する声門パラメータに基づいて構築されるものである。

30

【0094】

ステップ1120：前記目標音声フレームの対応するゲインに応じて前記フィルタリング出力信号に対して増幅処理を行い、前記目標音声フレームの対応する強調音声信号を得る。

【0095】

声門パラメータがLPC係数であれば、直接的に上式(2)にしたがって声門フィルターの構築を行うことができる。声門フィルターがp次フィルターであれば、目標音声フレームの対応する声門パラメータはp次LPC係数、すなわち上式(2)における $a_1$ 、 $a_2$ 、...、 $a_p$ を含み、その他の実施例において、上式(2)における定数1はLPC係数としてもよい。

40

【0096】

声門パラメータがLSFパラメータであれば、LSFパラメータをLPC係数に変換し、次に対応して上式(2)にしたがって声門フィルターを構築することができる。

【0097】

フィルタリング処理は、すなわち時間領域上の畳み込みであり、従って、上記のように声門フィルターにより励起信号に対してフィルタリングを行う過程は時間領域に変換して行うことができる。目標音声フレームに対応する励起信号の周波数領域表現を予測して得ることに加えて、励起信号の周波数領域表現を時間領域に変換し、目標音声フレームに対応する励起信号の時間領域信号を得る。

【0098】

50

本願の解決手段において、目標音声フレーム中には複数のサンプル点を含む。声門フィルターにより励起信号に対してフィルタリングを行い、すなわち1つのサンプル点の前の履歴サンプル点と該声門フィルターにより畳み込みを行い、該サンプル点の対応する目標信号値を得る。

【0099】

本願のいくつかの実施例では、前記目標音声フレームは複数のサンプル点を含み、前記声門フィルターは $p$ 次フィルターであり、 $p$ が正の整数であり、前記励起信号は前記目標音声フレームにおける複数のサンプル点のそれぞれの対応する励起信号値を含む。上記のようなフィルタリング過程に従って、ステップ1120は、さらに、前記目標音声フレームにおける各サンプル点の前の $p$ 個のサンプル点に対応する励起信号値と前記 $p$ 次フィルターを畳み込み、前記目標音声フレームにおける各サンプル点の目標信号値を得るステップと、時間順序に応じて前記目標音声フレームにおける全部サンプル点の対応する目標信号値を組み合わせ、前記第1音声信号を得るステップとを含む。ここで、 $p$ 次フィルターの表現式は上式(1)を参照することができる。つまり、目標音声フレームにおける各サンプル点に対しては、その前の $p$ 個のサンプル点に対応する励起信号値を利用して $p$ 次フィルターと畳み込みを行い、各サンプル点の対応する目標信号値を得る。

【0100】

理解できることとして、目標音声フレームにおける最初のサンプル点に対しては、該目標音声フレームの直前音声フレームにおける最後の $p$ 個のサンプル点の励起信号値を借りて該最初のサンプル点の対応する目標信号値を計算する必要があり、同様に、該目標音声フレームにおける2番目のサンプル点は、目標音声フレームの直前音声フレームにおける最後の $(p-1)$ 個のサンプル点の励起信号値及び目標音声フレームにおける最初のサンプル点の励起信号値と $p$ 次フィルターを借りて畳み込みを行って、目標音声フレームにおける2番目のサンプル点に対応する目標信号値を得る必要がある。

【0101】

要約すると、ステップ1120はさらに目標音声フレームの履歴音声フレームに対応する励起信号値の参加を必要とする。所要の履歴音声フレームにおけるサンプル点の数量は声門フィルターの次数に関連し、すなわち、声門フィルターが $p$ 次であれば、目標音声フレームの直前音声フレームにおける最後の $p$ 個のサンプル点に対応する励起信号値の参加を必要とする。

【0102】

関連する技術において、スペクトル推定とスペクトル回帰予測の方式で音声強調を行うことが存在する。スペクトル推定の音声強調方式は一段の混合音声に音声部分とノイズ部分が含まれると考えるため、統計モデル等によりノイズを推定することができるものであり、混合音声の対応するスペクトルからノイズの対応するスペクトルを減算すれば、残るのは音声スペクトルであり、これにより、混合音声の対応するスペクトルに基づいてノイズの対応するスペクトルを減算して得られたスペクトルはクリーンな音声信号を復元することになる。スペクトル回帰予測の音声強調方式は、ニューラルネットワークにより音声フレームの対応するマスキング閾値を予測し、該マスキング閾値は該音声フレームにおける各々の周波数点における音声成分とノイズ成分の割合を反映し、次に該マスキング閾値に基づいて混合信号スペクトルに対してゲイン制御を行い、強調された後のスペクトルを取得するということである。

【0103】

上記のスペクトル推定とスペクトル回帰予測による音声強調方式は、ノイズスペクトル事後確率に基づく推定であり、推定されるノイズが不正確である。たとえば、キーボード叩き等の過渡ノイズが存在する可能性があり、瞬時に発生するため、推定されるノイズスペクトルは非常に不正確であり、ノイズ抑制の効果が良くないことを引き起こす。ノイズスペクトル予測が不正確である場合に、推定されるノイズスペクトルに応じて元の混合音声信号に対して処理を行えば、混合音声信号における音声の歪みを引き起こす、又はノイズ抑制効果の劣化を引き起こす可能性があり、従って、このような状況においては、音声

10

20

30

40

50

忠実度とノイズ抑制との間で妥協を行う必要がある。

【0104】

声門パラメータ、励起信号及びゲイン予測に基づき音声強調を実現する上記実施例において、声門パラメータが音声生成の物理的過程における声門特徴と強い相関を有するため、予測された声門パラメータが目標音声フレームにおける元の音声信号の音声構造を効果的に保証し、従って、音声分解で得られた声門パラメータ、励起信号及びゲインに対して合成を行うことにより目標音声フレームの強調音声信号を得ることは、元の音声が増減されることを効果的に回避することができ、音声構造を効果的に保護し、且つ、目標音声フレームの対応する声門パラメータ、励起信号及びゲインを得た後、元のノイズ付きの音声に対して処理を行うことがなくなるため、音声忠実度とノイズ抑制との両方の間に妥協を行う必要がなくなる。

10

【0105】

図12は、別の1つの具体的な実施例に基づいて示される音声強調方法のフローチャートである。図12に示される実施例においては、上記第2ニューラルネットワーク、第3ニューラルネットワーク及び第4ニューラルネットワークを結合して音声分解を行う。第nフレームの音声フレームを目標音声フレームとすると仮定すると、該第nフレームの音声フレームの時間領域信号は $s(n)$ である。図12に示すように、該音声強調方法はステップ1210~1270を含む。

【0106】

ステップ1210：時間周波数変換であって、第nフレームの音声フレームの時間領域信号 $s(n)$ を第nフレームの音声フレームの対応する複素スペクトル $S(n)$ に変換する。

20

【0107】

ステップ1220：プリエンファシスであって、複素スペクトル $S(n)$ に基づいて第nフレームの音声フレームに対してプリエンファシスを行い、第1複素スペクトル $S'(n)$ を得る。

【0108】

ステップ1230：第2ニューラルネットワークにより声門パラメータを予測する。該ステップにおいて、第2ニューラルネットワークの入力は第1複素スペクトル $S'(n)$ のみを有してもよく、第1複素スペクトル $S'(n)$ と該第nフレームの音声フレームの履歴音声フレームの対応する声門パラメータ $P\_pre(n)$ とを含んでもよく、該第2ニューラルネットワークは該第nフレームの音声フレームの対応する声門パラメータ $a_r(n)$ を出力し、該声門パラメータはLPC係数であってもよく、LSFパラメータであってもよい。

30

【0109】

ステップ1240：第3ニューラルネットワークにより励起信号を予測する。第3ニューラルネットワークの入力は第1複素スペクトル $S'(n)$ であり、出力は該第nフレームの音声フレームに対応する励起信号の周波数領域表現 $R(n)$ である。次にステップ1250によって $R(n)$ に対して周波数時間変換を行い、第nフレームの音声フレームに対応する励起信号の時間領域信号 $r(n)$ を得ることができる。

40

【0110】

ステップ1260：第4ニューラルネットワークによりゲインを予測する。第4ニューラルネットワークの入力は第nフレームの音声フレームの履歴音声フレームに対応するゲイン $G\_pre(n)$ であり、出力は第nフレームの音声フレームの対応するゲイン $G(n)$ である。

【0111】

第nフレームの音声フレームの対応する声門パラメータ $a_r(n)$ 、対応する励起信号 $r(n)$ 及び対応するゲイン $G(n)$ を取得した後に、該3種のパラメータに基づきステップ1270で合成フィルタリングを行い、該第nフレームの音声フレームに対応する強調音声信号の時間領域信号 $s\_e(n)$ を得る。ステップ1270の合成フィルタリン

50

グの過程は、図 11 に示される過程を参照して行うことができる。

【0112】

本願の別のいくつかの実施例において、図 13 に示すように、ステップ 420 は、ステップ 1310 ~ ステップ 1350 を含む。

【0113】

ステップ 1310：前記第 1 複素スペクトルに基づいてパワースペクトルを計算して取得する。

【0114】

第 1 複素スペクトルが  $S'(n)$  であれば、ステップ 1310 において得られたパワースペクトル  $Pa(n)$  は、

$Pa(n) = \text{Real}(S'(n))^2 + \text{Imag}(S'(n))^2$  (式 10) である。

【0115】

ここで、 $\text{Real}(S'(n))$  は第 1 複素スペクトル  $S'(n)$  の実部を表し、 $\text{Imag}(S'(n))$  は第 1 複素スペクトル  $S'(n)$  の虚部を表す。ステップ 1310 において計算されて取得されたパワースペクトルは、すなわち目標音声フレームに対してプリエンファシスを行った後の信号のパワースペクトルである。

【0116】

ステップ 1320：前記パワースペクトルに基づいて自己相関係数を計算して取得する。

【0117】

ウィナーヒンチンの定理に従う：定常なランダム過程のパワースペクトルとその自己相関関数とは一對のフーリエ変換関係である。本解決方法において、1 フレームの音声フレームは定常なランダム信号と見なされる。従って、目標音声フレームに対応するプリエンファシスされた後のパワースペクトルを得たことに加えて、目標音声フレームに対応するプリエンファシスされた後のパワースペクトルに対して逆フーリエ変換を行い、該プリエンファシスされた後のパワースペクトルの対応する自己相関係数を得ることができる。

【0118】

具体的には、ステップ 1320 は、前記パワースペクトルに対して逆フーリエ変換を行い、逆変換結果を得て、前記逆変換結果中の実部を抽出し、前記自己相関係数を得ることを含む。すなわち、

$$AC(n) = \text{Real}(i\text{FFT}(Pa(n))) \quad (\text{式 11})$$

$AC(n)$  は第  $n$  フレームの音声フレームの対応する自己相関係数を表し、 $i\text{FFT}$  (Inverse Fast Fourier Transform、逆高速フーリエ変換) とは  $\text{FFT}$  (Fast Fourier Transform、高速フーリエ変換) の逆変換を指し、 $\text{Real}$  は逆高速フーリエ変換で得られた結果の実部を表す。 $AC(n)$  は  $p$  個のパラメータを含み、 $p$  が声門フィルターの次数であり、 $AC(n)$  中の係数はさらに  $AC_j(n)$  として表されてもよく、 $1 \leq j \leq p$  である。

【0119】

ステップ 1330：前記自己相関係数に基づいて前記声門パラメータを計算して取得する。

【0120】

Yule-Walker (ユール - ウォーカー) 方程式にしたがって、第  $n$  フレームの音声フレームに対して、その対応する自己相関係数と対応する声門パラメータとの間に以下の関係が存在する

$$k - KA = 0 \quad (\text{式 12})$$

ここで、 $k$  は自己相関ベクトルであり、 $K$  は自己相関行列であり、 $A$  は LPC 係数行列である。具体的には、[数 3] である。

【0121】

【数 3】

10

20

30

40

50

$$r = \begin{bmatrix} AC_0(n) \\ AC_1(n) \\ \vdots \\ AC_p(n) \end{bmatrix}, R = \begin{bmatrix} AC_0(n) & AC_1(n) & \cdots & AC_{p-1}(n) \\ AC_1(n) & AC_0(n) & \cdots & AC_{p-2}(n) \\ \vdots & \vdots & \ddots & \vdots \\ AC_{p-1}(n) & AC_{p-1}(n) & \cdots & AC_0(n) \end{bmatrix},$$

$$A = \begin{bmatrix} a_1(n) \\ a_2(n) \\ \vdots \\ a_p(n) \end{bmatrix}$$

10

【 0 1 2 2 】

ここで、 $AC_j(n) = E[s(n)s(n-j)]$ ,  $0 \leq j \leq p$  (式 1 3)

【 0 1 2 3 】

$p$  は声門フィルターの次数であり、 $a_1(n)$ 、 $a_2(n)$ 、...、 $a_p(n)$  はいずれも第  $n$  フレームの音声フレームに対応する LPC 係数であり、それぞれ上式 2 における  $a_1$ 、 $a_2$ 、...、 $a_p$  であり、 $a_0(n)$  が定数 1 であるため、 $a_0(n)$  を第  $n$  フレームの音声フレームに対応する 1 つの LPC 係数として見なすこともできる。

【 0 1 2 4 】

自己相関係数を得たことに加えて、自己相関ベクトルと自己相関行列は対応して決定することができ、次に式 1 2 を求めることにより、LPC 係数を得ることができる。具体的な実施例において、Levinson-Durbin アルゴリズムを採用して式 1 2 を求めることができ、Levinson-Durbin アルゴリズムは自己相関行列の対称性を利用し、反復の方式を利用して、自己相関係数を計算して取得する。

20

【 0 1 2 5 】

LSF パラメータと LPC 係数との間は相互に変換することができ、従って、LPC 係数を計算して取得する時に、LSF パラメータに対応して決定することができる。換言すれば、声門パラメータが LPC 係数であるか LSF パラメータであるかにかかわらず、いずれも上記のような過程によって決定することができる。

30

【 0 1 2 6 】

ステップ 1 3 4 0 : 前記声門パラメータと前記自己相関パラメータ集合とに基づいて前記ゲインを計算して取得する。

【 0 1 2 7 】

以下の式 [ 数 4 ] にしたがって第  $n$  フレームの音声フレームの対応するゲインを計算することができる。

$$[ \text{数 4} ] \quad ( \text{式 1 4} )$$

【 0 1 2 8 】

【 数 4 】

$$G(n) = \sum_{j=0}^p AC_j(n) * a_j(n)$$

40

【 0 1 2 9 】

式 1 4 にしたがって計算して取得した  $G(n)$  は時間領域表示上の目標音声フレームに対応するゲインの二乗である。

【 0 1 3 0 】

ステップ 1 3 5 0 : 前記ゲインと声門フィルターのパワースペクトルとに基づいて前記励起信号のパワースペクトルを計算して取得する。前記声門フィルターは前記声門パラメータに基づいて構築されるフィルターである。

【 0 1 3 1 】

50

目標音声フレームの対応する複素スペクトルが  $m$  ( $m$  が正の整数) 個のサンプル点に対してはフーリエ変換を行って得られるものと仮定すると、声門フィルターのパワースペクトルを計算するためには、まず第  $n$  フレームの音声フレームのために次元が  $m$  の全 0 の配列  $s\_AR(n)$  を構造し、次に、 $(p+1)$  次元の  $a_j(n)$  を該全 0 の配列の最初の  $(p+1)$  次元に代入し、ここで  $j = 0, 1, 2, \dots, p$  であり、 $m$  個のサンプル点の高速フーリエ変換 (Fast Fourier Transform、FFT) を呼び出すことにより、FFT 係数を取得する。

$$S\_AR(n) = FFT(s\_AR(n)) \quad (\text{式 15})$$

FFT 係数  $S\_AR(n)$  を得たことに加えて、下式 16 にしたがって 1 つずつのサンプルについて第  $n$  フレームの音声フレームに対応する声門フィルターのパワースペクトルを取得することができる、

$$AR\_LPS(n, k) = (\text{Real}(S\_AR(n, k)))^2 + (\text{Imag}(S\_AR(n, k)))^2 \quad (\text{式 16})$$

ここで、 $\text{Real}(S\_AR(n, k))$  は  $S\_AR(n, k)$  の実部を表し、 $\text{Imag}(S\_AR(n, k))$  は  $S\_AR(n, k)$  の虚部を表し、 $k$  は FFT 係数の数列を表し、 $0 \leq k < m$ 、 $k$  は正の整数である。

#### 【0132】

第  $n$  フレームの音声フレームに対応する声門フィルターの周波数応答  $AR\_LPS(n)$  を得た後に、計算を便利にするために、式 17 にしたがって声門フィルターのパワースペクトル  $AR\_LPS(n)$  を自然数領域から対数領域に変換し、

$$AR\_LPS_1(n) = \log_{10}(AR\_LPS(n)) \quad (\text{式 17})$$

上記  $AR\_LPS_1(n)$  を下式 18 にしたがって反転し、すなわち、声門フィルターの逆対応するパワースペクトル  $AR\_LPS_2(n)$  を得て、

$$AR\_LPS_2(n) = -1 * AR\_LPS_1(n) \quad (\text{式 18})$$

次に下式 19 にしたがって目標音声フレームに対応する励起信号のパワースペクトル  $R(n)$  を計算して取得することができる。

$$R(n) = Pa(n) * (G_1(n))^2 * AR\_LPS_3(n) \quad (\text{式 19})$$

ここで、[数 5] (式 20)

[数 6] (式 21)

#### 【0133】

##### 【数 5】

$$G_1(n) = \frac{1}{\sqrt{G(n)}}$$

##### 【数 6】

$$AR\_LPS_3(n) = 10^{\frac{AR\_LPS_2(n)}{10}}$$

#### 【0134】

上記のような過程により、目標音声フレームに対応する声門パラメータ、ゲイン及び励起信号の周波数応答、及び声門パラメータにより限定される声門フィルターの周波数応答を計算して取得する。

#### 【0135】

目標音声フレームに対応するゲイン、対応する励起信号のパワースペクトル、及び声門パラメータに限定される声門フィルターのパワースペクトルを得た後に、図 14 に示される過程に基づいて合成処理を行うことができる。図 14 に示すように、ステップ 430 は、ステップ 1410 ~ ステップ 1430 を含む。

#### 【0136】

10

20

30

40

50

ステップ 1 4 1 0 : 前記声門フィルターのパワースペクトルと前記励起信号のパワースペクトルとに基づいて第 1 振幅スペクトルを生成する。

【 0 1 3 7 】

以下の式 2 2 にしたがって第 1 振幅スペクトル  $S\_filt(n)$  を計算して取得することができる。

[ 数 7 ] ( 式 2 2 )

【 0 1 3 8 】

【 数 7 】

$$S\_filt(n) = \sqrt{10^{R_1(n) + AR_1\_LPS(n)}}.$$

10

【 0 1 3 9 】

ここで、 $R_1(n) = 10 * \log_{10}(R(n))$  ( 式 2 3 )

【 0 1 4 0 】

ステップ 1 4 2 0 : 前記ゲインに応じて前記第 1 振幅スペクトルに対して増幅処理を行い、第 2 振幅スペクトルを得る。

【 0 1 4 1 】

下式にしたがって第 2 振幅スペクトル  $S\_e(n)$  を得ることができる。

$$S\_e(n) = G_2(n) * S\_filt(n) \quad ( 式 2 4 )$$

ここで、[ 数 8 ] ( 式 2 5 )

20

【 0 1 4 2 】

【 数 8 】

$$G_2(n) = \sqrt{G(n)}$$

【 0 1 4 3 】

ステップ 1 4 3 0 : 前記第 2 振幅スペクトルと前記第 1 複素スペクトル中から抽出された位相スペクトルとに基づいて、前記目標音声フレームの対応する強調音声信号を決定する。

【 0 1 4 4 】

本願のいくつかの実施例では、ステップ 1 4 3 0 は、さらに、前記第 2 振幅スペクトルと前記第 1 複素スペクトル中から抽出された位相スペクトルとを組み合わせ、第 2 複素スペクトルを得るステップ、換言すれば、第 2 振幅スペクトルを第 2 複素スペクトルの実部とし、第 1 複素スペクトル中から抽出された位相スペクトルを第 2 複素スペクトルの虚部とし、前記第 2 複素スペクトルを時間領域に変換し、前記目標音声フレームに対応する強調音声信号の時間領域信号を得るステップを含む。

【 0 1 4 5 】

図 1 5 は、1 つの具体的な実施例に基づいて示される音声強調方法のフローチャートであり、第 n フレームの音声フレームを目標音声フレームとし、第 n フレームの音声フレームの時間領域信号が  $s(n)$  である。図 1 5 に示すように、具体的には、ステップ 1 5 1 0 ~ 1 5 6 0 を含む。

【 0 1 4 6 】

ステップ 1 5 1 0 : 時間周波数変換であって、ステップ 1 5 1 0 により第 n フレームの音声フレームの時間領域信号  $s(n)$  を変換して第 n フレームの音声フレームの対応する複素スペクトル  $S(n)$  を得る。

【 0 1 4 7 】

ステップ 1 5 2 0 : プリエンファシスであって、第 n フレームの音声フレームの対応する複素スペクトル  $S(n)$  に基づき該第 n フレームの音声フレームに対してプリエンファシス処理を行い、第 n フレームの音声フレームの第 1 複素スペクトル  $S(n)$  を得る。

【 0 1 4 8 】

50

ステップ1530：スペクトル分解であって、第1複素スペクトル $S(n)$ に対してスペクトル分解を行うことにより、第1複素スペクトル $S(n)$ の対応するパワースペクトル $P_a(n)$ と対応する位相スペクトル $P_h(n)$ とを得る。

【0149】

ステップ1540：音声分解であって、第 $n$ フレームの音声フレームのパワースペクトル $P_a(n)$ に基づき音声分解を行い、第 $n$ フレームの音声フレームの対応する声門パラメータ集合 $P(n)$ と第 $n$ フレームの音声フレームに対応する励起信号の周波数領域表現 $R(n)$ とを決定する。声門パラメータ集合 $P(n)$ は声門パラメータ $a_r(n)$ とゲイン $G(n)$ を含む。具体的な音声分解の過程は図13に示されてもよく、声門パラメータを取得し、且つ声門フィルターのパワースペクトル $AR\_LPS(n)$ 、励起信号のパワースペクトル $R(n)$ 、及びゲイン $G(n)$ を対応して取得する。

10

【0150】

ステップ1550：音声合成する。具体的な音声合成の過程は図14に示されてもよく、第 $n$ フレームの音声フレームに対応する声門フィルターの周波数応答 $AR\_LPS(n)$ 、励起信号の周波数応答 $R(n)$ 、及びゲイン $G(n)$ に対して合成を行って第2振幅スペクトル $S_e(n)$ を得る。

【0151】

ステップ1560：周波数時間変換する。第1複素スペクトル $S(n)$ から抽出された位相スペクトル $P_h(n)$ を多重化し、位相スペクトル $P_h(n)$ と第2振幅スペクトル $S_e(n)$ を組み合わせることで第 $n$ フレームの音声フレームに対応する強調された後の複素スペクトルを得る。得られた強調された後の複素スペクトルを時間領域に変換すると、第 $n$ フレームの音声フレームに対応する強調音声信号の時間領域信号 $s_e(n)$ を得る。

20

【0152】

本実施例の解決手段において、目標音声フレームに対してプリアンファシスを行うことにより得られた第1複素スペクトルに基づいて音声分解を行い、プリアンファシスする過程において、一部のノイズの情報が除外され、従って、第1複素スペクトルにおけるノイズ情報がより少なくなる。従って、第1複素スペクトルに基づいて音声分解を行うことで、ノイズによる音声分解への影響を減少し、音声分解の難度を低減させ、音声分解で得られた声門パラメータ、励起信号及びゲインの正確性を向上させ、さらに後続で取得された強調音声信号の正確性を保証することができる。また、本解決方法において、音声合成過程において、振幅スペクトルのみ注目することができ、位相情報に注目する必要がなく、第1複素スペクトル中から抽出された位相スペクトルを直接的に多重化することにより、音声合成過程における計算量を減少させる。第1複素スペクトルはプリアンファシスを行って得られるものであり、そのノイズ含有量がより少なく、従って、ある程度で位相情報の精度を保証する。

30

【0153】

図15に示される実施例においては、ステップ1510において、第1ニューラルネットワークによってプリアンファシスを実現することができる。ステップ1540は図13に示される過程にしたがって実現でき、ステップ1550は図14に示される過程にしたがって実現でき、それにより、従来信号処理と深層学習とを深く組み合わせ、且つ目標音声フレームに対して二次強調を行うことが実現される。従って、本願の実施例は目標音声フレームに対して複数段階の強調を行うことを実現する。すなわち、第1段階では、深層学習の方式を採用して目標音声フレームの振幅スペクトルに基づいてプリアンファシスを行い、第2段階における音声分解して声門パラメータ、励起信号及びゲインを取得する難しさを低減させることができ、第2段階では、信号処理の方式により元の音声信号を再構成することに用いられる声門パラメータ、励起信号及びゲインを取得する。そして、第2段階において、音声が生じているデジタルモデルにしたがって音声合成を行い、目標音声フレームの信号に対して処理を直接的に行わず、従って、第2段階における音声削減状況の出現を回避することができる。

40

【0154】

50

本願のいくつかの実施例では、ステップ410の前に、該方法は、さらに、前記目標音声フレームの時間領域信号を取得するステップと、前記目標音声フレームの時間領域信号に対して時間周波数変換を行い、前記目標音声フレームの複素スペクトルを得るステップとを含む。

【0155】

時間周波数変換は短時間フーリエ変換 (short-term Fourier transform、STFT) であってもよい。短時間フーリエ変換において窓掛け・オーバーラップの操作を採用してフレームの間の不平滑化を解消する。図16は1つの具体的な実施例に基づいて示される短時間フーリエ変換における窓掛け・オーバーラップの模式図であり、図16において、50%窓掛け・オーバーラップの操作を採用し、短時間フーリエ変換が640個のサンプル点に対するものであれば、該窓関数の重なったサンプル数 (hop-size) は320である。窓掛けに使用される窓関数はハニング (Hanning) 窓、ハミング窓等であってもよく、もちろん、その他の窓関数を採用してもよく、ここで具体的な限定を行わない。

10

【0156】

その他の実施例において、50%ではない窓掛け・オーバーラップの操作を採用してもよい。たとえば、短時間フーリエ変換が512個のサンプル点に対するものであれば、この場合には、1つの音声フレーム中に320個のサンプル点が含まれれば、直前音声フレームの192個のサンプル点をオーバーラップするだけでよい。

【0157】

本願のいくつかの実施例では、目標音声フレームの時間領域信号を取得するステップは、さらに、処理対象の音声信号を取得するステップであって、前記処理対象の音声信号は収集された音声信号又は符号化音声に対して復号を行って得られた音声信号である、ステップと、前記処理対象の音声信号に対してフレーム分割を行い、前記目標音声フレームの時間領域信号を得るステップとを含む。

20

【0158】

いくつかの実例において、設定されたフレーム長さに応じて処理対象の音声信号に対してフレーム分割を行うことができ、該フレーム長さは実際のニーズに応じて設定を行うことができ、たとえば、フレーム長さが20msに設定される。フレーム分割を行うことにより、複数の音声フレームを得ることができ、各音声フレームはいずれも本願における目標音声フレームとすることができる。

30

【0159】

上記の記述のように、本願の解決手段は送信端に適用され音声強調を行うことができ、受信端に適用され音声強調を行うこともできる。本願の解決手段が送信端に適用される場合に、該処理対象の音声信号は送信端が収集した音声信号であり、その場合、処理対象の音声信号に対してフレーム分割を行い、複数の音声フレームを得る。フレーム分割の後、処理対象の音声信号は複数の音声フレームに分割され、次に各音声フレームを目標音声フレームとし且つ上記ステップ410~430の過程にしたがって目標音声フレームに対して強調を行うことができる。さらには、目標音声フレームの対応する強調音声信号を得た後に、さらに該強調音声信号に対して符号化を行うこともでき、それにより、得られた符号化に基づき音声伝送を行う。

40

【0160】

一実施例において、直接収集された音声信号はアナログ信号であるため、信号処理を便利に行うために、フレーム分割を行う前に、音声信号をさらにデジタル化し、時間的に連続する音声信号を時間的に離散する音声信号に変換する必要もある。デジタル化を行う過程において、設定されたサンプリングレートに応じて収集された音声信号に対してサンプリングを行うことができ、設定されたサンプリングレートは16000Hz、8000Hz、32000Hz、48000Hz等であってもよく、具体的には、実際のニーズに応じて設定を行うことができる。

【0161】

50

本願の解決手段が受信端に適用される場合に、該処理対象の音声信号は受信された符号化音声に対して復号を行って得られた音声信号である。このような場合に、送信端が、伝送する必要がある音声信号に対して強調を行っていない可能性があり、従って、信号品質を向上させるためには、受信端で音声信号に対して強調を行う必要がある。処理対象の音声信号に対してフレーム分割を行って複数の音声フレームを得た後に、それを目標音声フレームとし、且つ上記のようなステップ 4 1 0 ~ 4 3 0 の過程にしたがって目標音声フレームに対して強調を行い、目標音声フレームの強調音声信号を得る。さらに、目標音声フレームの対応する強調音声信号に対して再生を行うこともでき、得られた強調音声信号は目標音声フレームの強調前の信号に比べて、ノイズが既に除去されているため、音声信号の品質がより高く、従って、ユーザーにとって、聴覚的体験がより高い。

10

**【 0 1 6 2 】**

以下、本願の上記実施例における方法を実行することに用いることができる本願の装置の実施例を説明する。本願の装置実施例において披露されない細部に対しては、本願の上記方法実施例を参照されたい。

**【 0 1 6 3 】**

図 1 7 は、一実施例に基づいて示される音声強調装置のブロック図である。図 1 7 に示すように、該音声強調装置は、目標音声フレームの複素スペクトルに基づいて前記目標音声フレームに対してプリエンファシス処理を行い、第 1 複素スペクトルを得ることに用いられるプリエンファシスモジュール 1 7 1 0 と、前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して音声分解を行い、前記目標音声フレームの対応する声門パラメータ、ゲイン及び励起信号を得ることに用いられる音声分解モジュール 1 7 2 0 と、前記声門パラメータ、前記ゲイン及び前記励起信号に基づいて合成処理を行い、前記目標音声フレームの対応する強調音声信号を得ることに用いられる合成処理モジュール 1 7 3 0 とを含む。

20

**【 0 1 6 4 】**

本願のいくつかの実施例では、プリエンファシスモジュール 1 7 1 0 は、前記目標音声フレームの対応する複素スペクトルを第 1 ニューラルネットワークに入力することに用いられる第 1 入力ユニットであって、前記第 1 ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームにおける元の音声信号の対応する複素スペクトルとに基づいてトレーニングを行って得られるものである、第 1 入力ユニットと、前記第 1 ニューラルネットワークによって、前記目標音声フレームの対応する複素スペクトルに基づいて前記第 1 複素スペクトルを出力することに用いられる第 1 出力ユニットとを含む。

30

**【 0 1 6 5 】**

本願のいくつかの実施例では、前記第 1 ニューラルネットワークは複素畳み込み層、ゲート付き回帰型ユニット層及び全結合層を含み、第 1 出力ユニットは、前記複素畳み込み層によって前記目標音声フレームに対応する複素スペクトルにおける実部及び虚部に基づいて複素畳み込み処理を行うことに用いられる複素畳み込みユニットと、前記ゲート付き回帰型ユニット層によって前記複素畳み込み層の出力に対して変換処理を行うことに用いられる変換ユニットと、前記全結合層によって前記ゲート付き回帰型ユニットの出力に対して全結合処理を行い、前記第 1 複素スペクトルを出力することに用いられる全結合ユニットとを含む。

40

**【 0 1 6 6 】**

本願のいくつかの実施例では、音声分解モジュール 1 7 2 0 は、前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して声門パラメータ予測を行い、前記目標音声フレームの対応する声門パラメータを得ることに用いられる声門パラメータ予測ユニットに用いられる第 1 振幅スペクトル取得ユニットと、前記第 1 複素スペクトルに基づいて前記目標音声フレームに対して励起信号予測を行い、前記目標音声フレームの対応する励起信号を得ることに用いられる励起信号予測ユニットと、前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームに対してゲイン予測を行い

50

、前記目標音声フレームの対応するゲインを得ることに用いられるゲイン予測ユニットとを含む。

【0167】

本願のいくつかの実施例では、声門パラメータ予測ユニットは、前記第1複素スペクトルを第2ニューラルネットワークに入力することに用いられる第2入力ユニットであって、前記第2ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームの対応する声門パラメータとに基づいてトレーニングを行って得られるものである、第2入力ユニットと、前記第2ニューラルネットワークによって、前記第1複素スペクトルに基づいて前記目標音声フレームの対応する声門パラメータを出力することに用いられる第2出力ユニットとを含む。

10

【0168】

本願の別のいくつかの実施例において、声門パラメータ予測ユニットは、前記第1複素スペクトルと前記目標音声フレームの前の履歴音声フレームの対応する声門パラメータとを第2ニューラルネットワークに入力することに用いられる第3入力ユニットであって、前記第2ニューラルネットワークはサンプル音声フレームの対応する複素スペクトル、サンプル音声フレームの前の履歴音声フレームの対応する声門パラメータ及びサンプル音声フレームの対応する声門パラメータに基づいてトレーニングを行って得られるものである、第3入力ユニットと、前記第1ニューラルネットワークによって、前記第1複素スペクトルと前記目標音声フレームの前の履歴音声フレームの対応する声門パラメータとに基づいて前記目標音声フレームの対応する声門パラメータを出力することに用いられる第3出力ユニットとを含む。

20

【0169】

本願のいくつかの実施例では、ゲイン予測ユニットは、前記目標音声フレームの前の履歴音声フレームの対応するゲインを第3ニューラルネットワークに入力することに用いられる第4入力ユニットであって、前記第3ニューラルネットワークはサンプル音声フレームの前の履歴音声フレームの対応するゲインと前記サンプル音声フレームの対応するゲインとに基づいてトレーニングを行って得られるものである、第4入力ユニットと、前記第3ニューラルネットワークによって、前記目標音声フレームの前の履歴音声フレームの対応するゲインに基づいて前記目標音声フレームの対応するゲインを出力することに用いられる第4出力ユニットとを含む。

30

【0170】

本願のいくつかの実施例では、励起信号予測ユニットは、前記第1複素スペクトルを第4ニューラルネットワークに入力することに用いられる第5入力ユニットであって、前記第4ニューラルネットワークはサンプル音声フレームの対応する複素スペクトルと前記サンプル音声フレームに対応する励起信号の周波数領域表現とに基づいてトレーニングを行って得られるものである、第5入力ユニットと、前記第4ニューラルネットワークによって、前記第1複素スペクトルに基づいて前記目標音声フレームに対応する励起信号の周波数領域表現を出力することに用いられる第5出力ユニットとを含む。

【0171】

本願のいくつかの実施例では、合成処理モジュール1730は、声門フィルターにより前記目標音声フレームの対応する励起信号に対してフィルタリングを行い、フィルタリング出力信号を得ることに用いられるフィルタリングユニットであって、前記声門フィルターは前記目標音声フレームの対応する声門パラメータに基づいて構築されるものである、フィルタリングユニットと、前記目標音声フレームの対応するゲインに応じて前記フィルタリング出力信号に対して増幅処理を行い、前記目標音声フレームの対応する強調音声信号を得ることに用いられる増幅処理ユニットとを含む。

40

【0172】

本願のいくつかの実施例では、音声分解モジュール1720は、前記第1複素スペクトルに基づいてパワースペクトルを計算して取得することに用いられるパワースペクトル計算ユニットと、前記パワースペクトルに基づいて自己相関係数を計算して取得することに

50

用いられる自己相関係数計算ユニットと、前記自己相関係数に基づいて前記声門パラメータを計算して取得することに用いられる声門パラメータ計算ユニットと、前記声門パラメータと前記自己相関パラメータ集合とに基づいて前記ゲインを計算して取得することに用いられるゲイン計算ユニットと、前記ゲインと声門フィルターのパワースペクトルとに基づいて前記励起信号のパワースペクトルを計算して取得することに用いられる励起信号決定ユニットであって、前記声門フィルターは前記声門パラメータに基づいて構築されるフィルターである、励起信号決定ユニットとを含む。

【0173】

本願のいくつかの実施例では、合成処理モジュール1730は、前記声門フィルターのパワースペクトルと前記励起信号のパワースペクトルとに基づいて第1振幅スペクトルを生成することに用いられる第2振幅スペクトル生成ユニットと、前記ゲインに応じて前記第1振幅スペクトルに対して増幅処理を行い、第2振幅スペクトルを得ることに用いられる第3振幅スペクトル決定ユニットと、前記第2振幅スペクトルと前記第1複素スペクトル中から抽出された位相スペクトルとに基づいて、前記目標音声フレームの対応する強調音声信号を決定することに用いられる強調音声信号決定ユニットとを含む。

10

【0174】

本願のいくつかの実施例では、強調音声信号決定ユニットは、前記第2振幅スペクトルと前記第1複素スペクトル中から抽出された位相スペクトルとを組み合わせ、第2複素スペクトルを得ることに用いられる第2複素スペクトル計算ユニットと、前記第2複素スペクトルを時間領域に変換し、前記目標音声フレームに対応する強調音声信号の時間領域信号を得ることに用いられる時間領域変換ユニットとを含む。

20

【0175】

図18は、本願の実施例を実現するための電子機器に適するコンピュータシステムの構造模式図を示す。

【0176】

説明する必要があることとして、図18に示される電子機器のコンピュータシステム1800は一例に過ぎず、本願の実施例の機能及び使用範囲に対して何ら制限をもたらすべきではない。

【0177】

図18に示すように、コンピュータシステム1800は中央処理ユニット(Central Processing Unit、CPU)1801を含み、それは読み出し専用メモリ(Read-Only Memory、ROM)1802において記憶されたプログラム又は記憶部分1808からランダムアクセスメモリ(Random Access Memory、RAM)1803中にアップロードされたプログラムに基づいて各種の適当な動作と処理を実行することができ、たとえば、上記実施例における方法を実行する。RAM 1803において、システム操作に必要な各種のプログラムとデータも記憶されている。CPU1801、ROM1802及びRAM 1803はバス1804を介して互いに連結される。入力/出力(Input/Output、I/O)インターフェース1805もバス1804に接続される。

30

【0178】

以下の部材がI/Oインターフェース1805に接続される。キーボード、マウス等を含む入力部分1806、陰極線管(Cathode Ray Tube、CRT)、液晶ディスプレイ(Liquid Crystal Display、LCD)等のようなもの及びスピーカ等を含む出力部分1807、ハードディスク等を含む記憶部分1808、及びLAN(Local Area Network、ローカルエリアネットワーク)カード、モデム等のようなネットワークインタフェースカードを含む通信部分1809である。通信部分1809は、インターネットのようなネットワークを介して通信処理を実行する。ドライバ1810もニーズに応じてI/Oインターフェース1805に接続される。着脱可能な媒体1811、例えば磁気ディスク、光ディスク、光磁気ディスク、半導体メモリ等は、ニーズに応じてドライバ1810上に装着され、それにより、その上から読み出さ

40

50

れたコンピュータプログラムがニーズに応じて記憶部分 1808 にインストールされる。

【0179】

特に、本願の実施例に基づき、上記のフローチャートを参照して記述される過程はコンピュータソフトウェアプログラムとして実現できる。たとえば、本願の実施例は、1種のコンピュータプログラム製品を含み、それはコンピュータ可読媒体上に担持されるコンピュータプログラムを含み、該コンピュータプログラムはフローチャートに示された方法を実行することに用いられるプログラムコードを含む。このような実施例において、該コンピュータプログラムは通信部分 1809 によりネットワーク上からダウンロードされインストールすることができ、且つ/又は着脱可能な媒体 1811 からインストールされる。該コンピュータプログラムが中央処理ユニット (CPU) 1801 によって実行されるときに、本願のシステム中に限定される各種の機能を実行する。

10

【0180】

説明する必要があることとして、本願の実施例に示されるコンピュータ可読媒体はコンピュータ可読信号媒体、又はコンピュータ可読記憶媒体又は上記両方の任意の組み合わせであってもよい。コンピュータ可読記憶媒体は、たとえば、電気、磁気、光、電磁、赤外線、又は半導体のシステム、装置又はデバイス、又は以上の任意の組み合わせであってもよいがこれらに限定されない。コンピュータ可読記憶媒体のより具体的な例は、1つ又は複数の導線を有する電氣的接続、ポータブルコンピュータ磁気ディスク、ハードディスク、ランダムアクセスメモリ (RAM)、読み出し専用メモリ (ROM)、消去可能なプログラマブル読み出し専用メモリ (Erasable Programmable Read Only Memory、EPROM)、フラッシュメモリ、光ファイバー、ポータブルコンパクト磁気ディスク読み出し専用メモリ (Compact Disc Read-Only Memory、CD-ROM)、光記憶デバイス、磁気記憶デバイス、又は上記の任意の適切な組み合わせを含んでもよいがこれらに限定されない。本願において、コンピュータ可読記憶媒体は、プログラムを含む又は記憶する任意の有形媒体であってもよく、該プログラムは指令実行システム、装置又はデバイスに使用され又はそれと組み合わせて使用することができる。本願において、コンピュータ可読の信号媒体は、ベースバンド中における又は搬送波の一部として伝播されるデータ信号を含んでもよく、その中でコンピュータ可読のプログラムコードが担持されている。このような伝播されるデータ信号は複数種の形式を採用することができ、電磁信号、光信号又は上記任意の適切な組み合わせを含むがこれらに限定されない。コンピュータ可読の信号媒体はさらにコンピュータ可読記憶媒体以外の任意のコンピュータ可読媒体であってもよく、該コンピュータ可読媒体は、指令実行システム、装置又はデバイスに使用され又はそれと組み合わせて使用されることに用いられるプログラムを送信、伝播又は伝送することができる。コンピュータ可読媒体上に含まれるプログラムコードは任意の適当な媒体で伝送でき、無線、有線等、又は上記の任意の適切な組み合わせを含むがこれらに限定されない。

20

30

【0181】

図面におけるフローチャートとブロック図は、本願の各種の実施例に係るシステム、方法及びコンピュータプログラム製品の実現可能な体系アーキテクチャ、機能及び操作を図示する。ここで、フローチャート又はブロック図における各ブロックは1つのモジュール、プログラムセグメント、又はコードの一部を代表することができ、上記モジュール、プログラムセグメント、又はコードの一部は規定されるロジック機能を実現することに用いられる1つ又は複数の実行可能な指示を含む。また、注意すべきことは、代替としてのいくつかの実現形式において、ブロック中にマークされる機能は図面中にマークされる順序と異なるものとして生じさせることができる点である。たとえば、連続的に示される2つのブロックは実際には基本的に並行して実行することができ、場合によって、それらは逆の順序で実行することもでき、これは関連する機能によって定められる。また注意する必要があるのは、ブロック図又はフローチャートにおける各ブロック、及びブロック図又はフローチャートにおけるブロックの組み合わせは、規定される機能又は操作を実行する専用のハードウェアに基づくシステムで実現することができ、又は専用ハードウェアとコン

40

50

コンピュータ指令の組み合わせで実現することもできる。

【0182】

本願の実施例においてに記述されて言及されるユニットはソフトウェアの方式で実現されても、又はハードウェアの方式で実現されてもよく、記述されるユニットはプロセッサ中に設置されてもよい。ここで、これらのユニットの名称がある場合には、該ユニット自体に対する限定を構成しない。

【0183】

別の態様として、本願はコンピュータ可読記憶媒体をさらに提供し、該コンピュータ可読媒体は上記実施例に記述される電子機器に含まれてもよく、単独で存在し、該電子機器中に組み立てられなくてもよい。上記コンピュータ可読記憶媒体はコンピュータ可読指令を担持し、該コンピュータ可読記憶指令がプロセッサによって実行されるときに、上記いずれかの実施例における方法を実現する。

10

【0184】

本願の一態様によれば、電子機器をさらに提供し、それは、プロセッサと、メモリであって、メモリ上にコンピュータ可読指令が記憶され、コンピュータ可読指令がプロセッサによって実行されるときに、上記いずれかの実施例における方法を実現するメモリを含む。

【0185】

本願の実施例の一態様によれば、コンピュータプログラム製品、又はコンピュータプログラムを提供し、該コンピュータプログラム製品、又はコンピュータプログラムはコンピュータ指令を含み、該コンピュータ指令がコンピュータ可読記憶媒体中に記憶される。コンピュータ機器のプロセッサはコンピュータ可読記憶媒体から該コンピュータ指令を読み取り、プロセッサは該コンピュータ指令を実行し、該コンピュータ機器に上記いずれかの実施例における方法を実行させる。

20

【0186】

注意すべきことは、上記詳細な記述において動作実行用の機器の複数のモジュール又はユニットが言及されているが、このような分割は強制的ではないことである。実際には、本願の実施形態によれば、上記で記述された2つ又はより多くのモジュール又はユニットの特徴と機能は1つのモジュール又はユニットにおいて具現化され得る。逆に、上記で記述された1つのモジュール又はユニットの特徴と機能はさらに複数のモジュール又はユニットにより具現化されるように分割されてもよい。

30

【0187】

以上の実施形態の記述により、当業者が容易に理解できることは、ここで記述される例示的な実施形態はソフトウェアで実現されてもよく、ソフトウェアと必要なハードウェアを組み合わせた方式で実現されてもよい。従って、本願の実施形態に係る技術的手段は、ソフトウェア製品の形式で体现されてもよく、該ソフトウェア製品は1つの不揮発性記憶媒体(CD-ROM、Uディスク、モバイルディスク等であってもよい)中に又はネットワーク上に記憶されてもよく、幾つかの指令を含むことで一台の計算機器(パソコンコンピュータ、サーバ、タッチ端末、又はネットワーク機器等であってもよい)に本願の実施形態に係る方法を実行させる。

40

【0188】

当業者は明細書を考慮し、且つここで開示される実施形態を实践した後に、本願のその他の実施形態を容易に想到することができる。本願は本願の任意の変形、用途又は適応的な変化をカバーすることを目的としており、これらの変形、用途又は適応的な変化は本願の一般原理に従い、且つ本願に開示されていない本技術分野における公知の知識又は一般的な技術手段を含む。

【0189】

理解すべきことは、本願は上記において記述され、且つ図面中に示される正確な構造には限定されず、且つその範囲を逸脱することなく、各種の修正や変更を行うことができる。本願の範囲は添付の請求項の記載のみによって制限される。

50

## 【符号の説明】

## 【0190】

110	送信端	
111	収集モジュール	
112	前強調処理モジュール	
113	符号化モジュール	
120	受信端	
121	復号モジュール	
122	後強調モジュール	
123	再生モジュール	10
1710	プリエンファシスモジュール	
1720	音声分解モジュール	
1730	合成処理モジュール	
1800	コンピュータシステム	
1801	中央処理ユニット(CPU)	
1804	バス	
1805	I/Oインターフェース	
1805	出力(Input/Output、I/O)インターフェース	
1806	入力部分	
1807	出力部分	20
1808	記憶部分	
1809	通信部分	
1810	ドライバ	
1811	媒体	

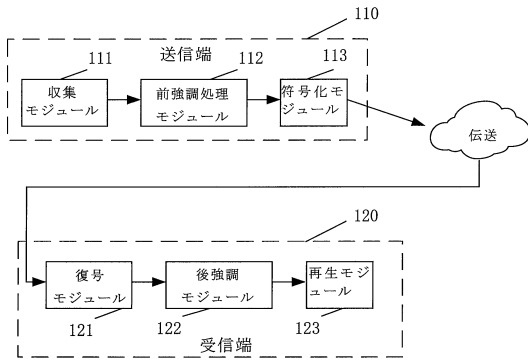
30

40

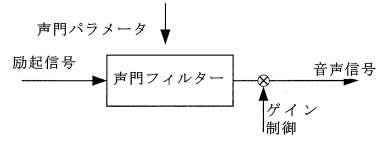
50

【図面】

【図 1】

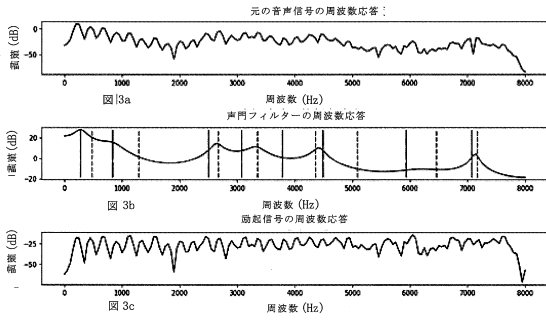


【図 2】

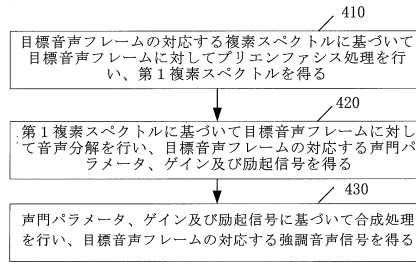


10

【図 3】



【図 4】



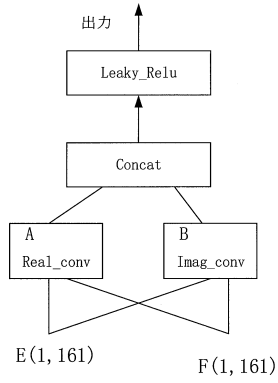
20

30

40

50

【 図 5 】



【 図 6 】

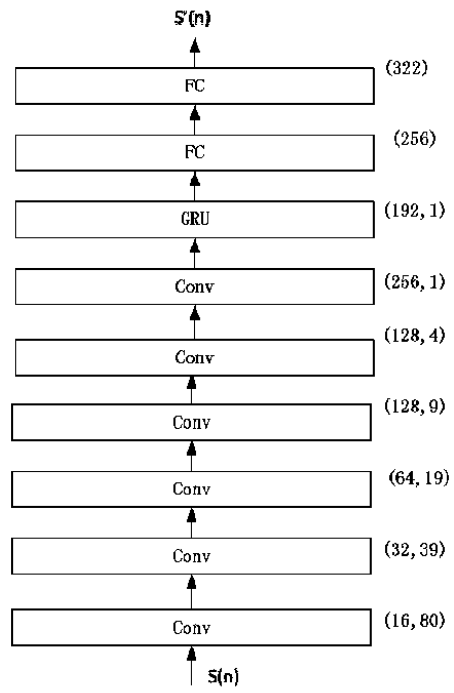


图 6

10

20

【 图 7 】

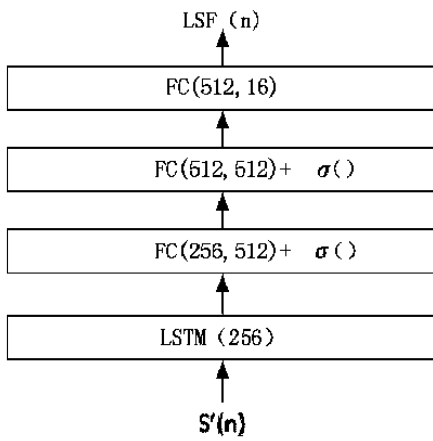


图 7

【 图 8 】

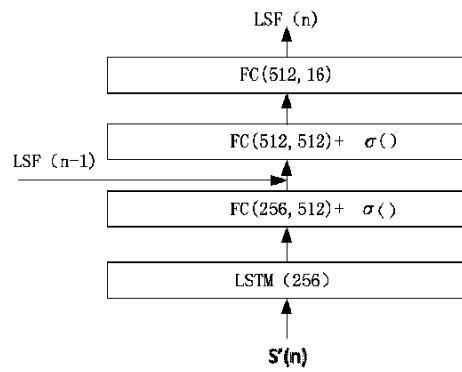


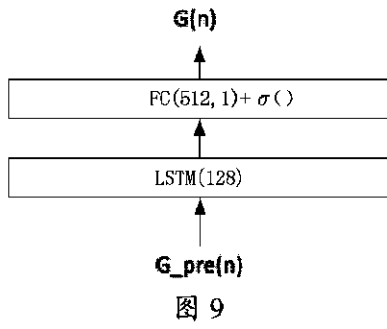
图 8

30

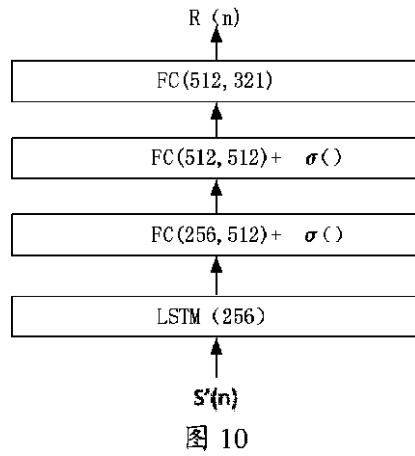
40

50

【 図 9 】

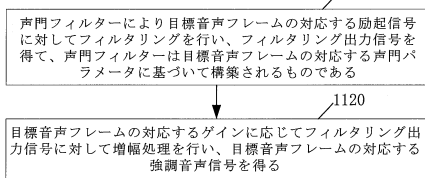


【 図 10 】

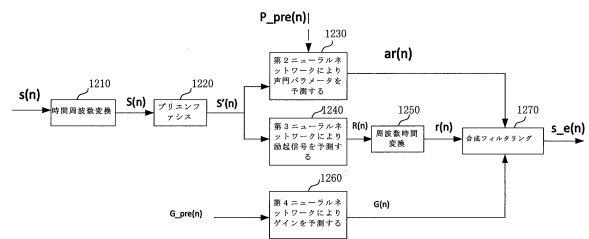


10

【 図 11 】



【 図 12 】



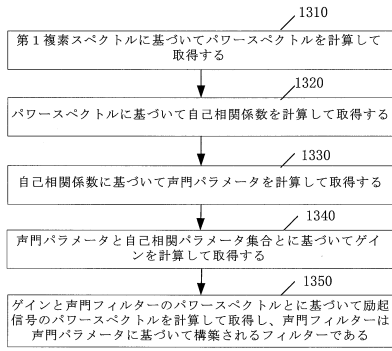
20

30

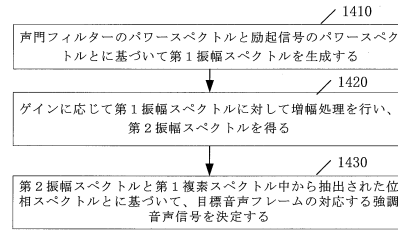
40

50

【図 13】

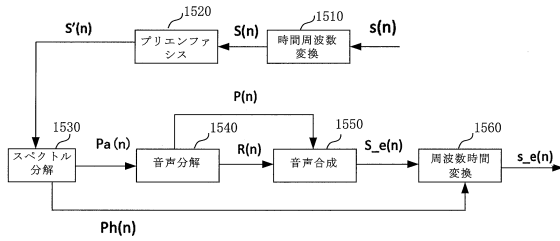


【図 14】



10

【図 15】



【図 16】

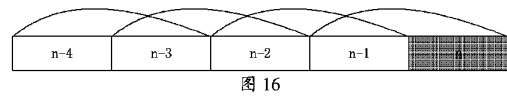
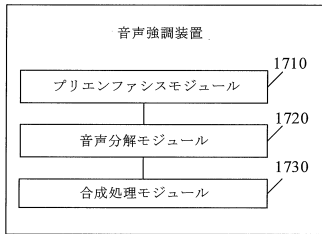


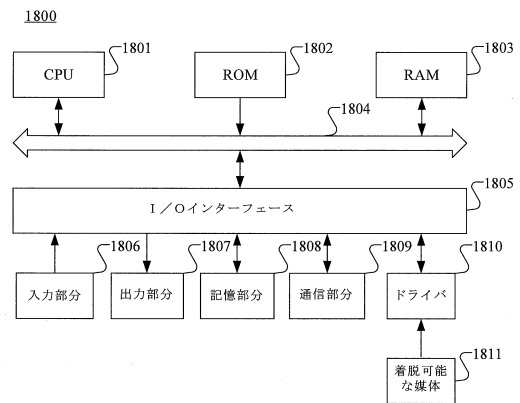
図 16

20

【図 17】



【図 18】



30

40

50

## フロントページの続き

- 057, CHINA
- (74)代理人 100110364  
弁理士 実広 信哉
- (74)代理人 100150197  
弁理士 松尾 直樹
- (72)発明者 肖 ウェイ  
中華人民共和国518057 広 東 省深 セン 市南山区高新区科技中一路 騰 訊  
大厦35 層
- (72)発明者 史 裕 鵬  
中華人民共和国518057 広 東 省深 セン 市南山区高新区科技中一路 騰 訊  
大厦35 層
- (72)発明者 王 蒙  
中華人民共和国518057 広 東 省深 セン 市南山区高新区科技中一路 騰 訊  
大厦35 層
- 審査官 中村 天真
- (56)参考文献 特開2000-347698(JP, A)  
特開2020-060612(JP, A)  
特開平02-137900(JP, A)  
特開2020-122896(JP, A)  
特開平10-190498(JP, A)  
特開2002-041085(JP, A)  
特開2002-366200(JP, A)  
米国特許第05148488(US, A)
- (58)調査した分野 (Int.Cl., DB名)  
G10L 13/00 - 13/10  
19/00 - 99/00