

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2014-228779

(P2014-228779A)

(43) 公開日 平成26年12月8日(2014.12.8)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 O L 19/005 (2013.01)	G 1 O L 19/00 3 3 O E	
G 1 O L 21/0388 (2013.01)	G 1 O L 21/04 1 3 O A	

審査請求 未請求 請求項の数 11 O L (全 14 頁)

(21) 出願番号	特願2013-109897 (P2013-109897)	(71) 出願人	000003078 株式会社東芝 東京都港区芝浦一丁目1番1号
(22) 出願日	平成25年5月24日 (2013.5.24)	(74) 代理人	100089118 弁理士 酒井 宏明
		(74) 代理人	100112656 弁理士 宮田 英毅
		(72) 発明者	大谷 大和 東京都港区芝浦一丁目1番1号 株式会社東芝内
		(72) 発明者	森田 真弘 東京都港区芝浦一丁目1番1号 株式会社東芝内

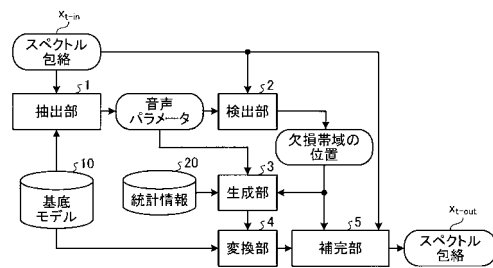
(54) 【発明の名称】 音声処理装置、方法およびプログラム

(57) 【要約】

【課題】 任意の周波数帯域で欠損した音声成分を適切に補完することができる音声処理装置、方法およびプログラムを提供する。

【解決手段】 実施形態の音声処理装置は、抽出部と、検出部と、生成部と、変換部と、補完部とを備える。抽出部は、入力音声のスペクトル包絡から、細分化された周波数帯域ごとの音声成分を表現する音声パラメータを抽出する。検出部は、入力音声のスペクトル包絡において音声成分が欠損している周波数帯域である欠損帯域を検出する。生成部は、欠損帯域の位置と、統計情報と、入力音声のスペクトル包絡から抽出された音声パラメータとに基づいて、欠損帯域に対応する音声パラメータを生成する。変換部は、欠損帯域に対応する音声パラメータを、欠損帯域のスペクトル包絡に変換する。補完部は、欠損帯域のスペクトル包絡と入力音声のスペクトル包絡とを合成して、欠損帯域が補完されたスペクトル包絡を生成する。

【選択図】 図1



**【特許請求の範囲】****【請求項 1】**

入力音声のスペクトル包絡から、細分化された周波数帯域ごとの音声成分を表現する音声パラメータを抽出する抽出部と、

前記入力音声のスペクトル包絡において音声成分が欠損している周波数帯域である欠損帯域を検出する検出部と、

検出された前記欠損帯域の位置と、音声成分が欠損していない音声のスペクトル包絡から抽出された前記音声パラメータを用いて事前に作成された統計情報と、前記入力音声のスペクトル包絡から抽出された前記音声パラメータとに基づいて、前記欠損帯域に対応する前記音声パラメータを生成する生成部と、

生成された前記欠損帯域に対応する前記音声パラメータを、前記欠損帯域のスペクトル包絡に変換する変換部と、

前記欠損帯域のスペクトル包絡と前記入力音声のスペクトル包絡とを合成して、前記欠損帯域が補完されたスペクトル包絡を生成する補完部と、を備える音声処理装置。

**【請求項 2】**

前記音声パラメータは、細分化された前記周波数帯域の各々に対応する複数の基底ベクトルを用いて算出される値であり、

前記基底ベクトルの数は、音声のスペクトル包絡の分析に用いた分析点数よりも少ないことを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 3】**

前記基底ベクトルに対応する前記周波数帯域の範囲は、周波数軸上で隣り合う範囲の一部が重複していることを特徴とする請求項 2 に記載の音声処理装置。

**【請求項 4】**

前記音声パラメータは、複数の前記基底ベクトルと各基底ベクトルに対応する重みベクトルとの線形結合と、音声のスペクトル包絡と、の誤差が最小になるように決定された前記重みベクトルであることを特徴とする請求項 2 または 3 に記載の音声処理装置。

**【請求項 5】**

前記検出部は、前記入力音声のスペクトル包絡または該スペクトル包絡から抽出された前記音声パラメータの包絡形状を解析して、前記欠損帯域を検出することを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 6】**

前記統計情報は、音声成分が欠損していない複数の話者の音声から抽出された前記音声パラメータを学習データとして構築された統計モデルであることを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 7】**

前記統計情報は、音声成分が欠損していない複数の話者の音声から抽出された前記音声パラメータの系列と、該音声パラメータの系列から抽出された時間変動成分と、を学習データとして構築された統計モデルであることを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 8】**

前記生成部は、前記欠損帯域の位置と前記統計情報とに基づいて、前記欠損帯域を除く周波数帯域である残存帯域に対応する前記音声パラメータから前記欠損帯域に対応する前記音声パラメータを生成する規則を構築し、該規則を用いて、前記入力音声の音声スペクトル包絡から抽出された前記音声パラメータから、前記欠損帯域に対応する前記音声パラメータを生成することを特徴とする請求項 1 に記載の音声処理装置。

**【請求項 9】**

前記変換部は、前記欠損帯域に対応する前記音声パラメータとして生成された前記重みベクトルと、前記欠損帯域に対応する前記基底ベクトルとを線形結合することにより、前記欠損帯域に対応する前記音声パラメータを前記欠損帯域のスペクトル包絡に変換することを特徴とする請求項 4 に記載の音声処理装置。

10

20

30

40

50

**【請求項 10】**

音声処理装置において実行される音声処理方法であって、  
前記音声処理装置が、入力音声のスペクトル包絡から、細分化された周波数帯域ごとの音声成分を表現する音声パラメータを抽出するステップと、  
前記音声処理装置が、前記入力音声のスペクトル包絡において音声成分が欠損している周波数帯域である欠損帯域を検出するステップと、  
前記音声処理装置が、検出された前記欠損帯域の位置と、音声成分が欠損していない音声のスペクトル包絡から抽出された前記音声パラメータを用いて事前に作成された統計情報と、前記入力音声のスペクトル包絡から抽出された前記音声パラメータとに基づいて、前記欠損帯域に対応する前記音声パラメータを生成するステップと、  
前記音声処理装置が、生成された前記欠損帯域に対応する前記音声パラメータを、前記欠損帯域のスペクトル包絡に変換するステップと、  
前記音声処理装置が、前記欠損帯域のスペクトル包絡と前記入力音声のスペクトル包絡とを合成して、前記欠損帯域が補完されたスペクトル包絡を生成するステップと、を含む音声処理方法。

10

**【請求項 11】**

コンピュータに、  
入力音声のスペクトル包絡から、細分化された周波数帯域ごとの音声成分を表現する音声パラメータを抽出する機能と、  
前記入力音声のスペクトル包絡において音声成分が欠損している周波数帯域である欠損帯域を検出する機能と、  
検出された前記欠損帯域の位置と、音声成分が欠損していない音声のスペクトル包絡から抽出された前記音声パラメータを用いて事前に作成された統計情報と、前記入力音声のスペクトル包絡から抽出された前記音声パラメータとに基づいて、前記欠損帯域に対応する前記音声パラメータを生成する機能と、  
生成された前記欠損帯域に対応する前記音声パラメータを、前記欠損帯域のスペクトル包絡に変換する機能と、  
前記欠損帯域のスペクトル包絡と前記入力音声のスペクトル包絡とを合成して、前記欠損帯域が補完されたスペクトル包絡を生成する機能と、を実現させるためのプログラム。

20

**【発明の詳細な説明】**

30

**【技術分野】****【0001】**

本発明の実施の形態は、音声処理装置、方法およびプログラムに関する。

**【背景技術】****【0002】**

従来、例えば携帯電話機や音声収録装置等の音声品質を向上させる技術として、帯域拡張が知られている。帯域拡張は、狭帯域音声から広帯域音声を構築する技術であり、例えば、入力音声において欠損している高周波帯域の音声成分を、欠損していない音声成分を用いて補完することができる。

**【0003】**

40

しかし、従来の帯域拡張では、入力音声において欠損している高周波帯域の音声成分や、予め定められた特定の周波数帯域の音声成分を補完することはできるが、任意の周波数帯域の音声成分が部分的に欠損した場合に対応できない。音声処理装置に入力される音声信号は、伝送路の静的特性等の何らかの影響によって、任意の周波数帯域の音声成分が部分的に欠損することがあり、任意の周波数帯域の音声成分を適切に補完できるようにすることが求められる。

**【先行技術文献】****【特許文献】****【0004】**

【特許文献 1】特開 2012 - 83790 号公報

50

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0005】

本発明が解決しようとする課題は、任意の周波数帯域で欠損した音声成分を適切に補完することができる音声処理装置、方法およびプログラムを提供することである。

## 【課題を解決するための手段】

## 【0006】

実施形態の音声処理装置は、抽出部と、検出部と、生成部と、変換部と、補完部と、を備える。抽出部は、入力音声のスペクトル包絡から、細分化された周波数帯域ごとの音声成分を表現する音声パラメータを抽出する。検出部は、前記入力音声のスペクトル包絡において音声成分が欠損している周波数帯域である欠損帯域を検出する。生成部は、検出された前記欠損帯域の位置と、音声成分が欠損していない音声のスペクトル包絡から抽出された前記音声パラメータを用いて事前に作成された統計情報と、前記入力音声のスペクトル包絡から抽出された前記音声パラメータとに基づいて、前記欠損帯域に対応する前記音声パラメータを生成する。変換部は、生成された前記欠損帯域に対応する前記音声パラメータを、前記欠損帯域のスペクトル包絡に変換する。補完部は、前記欠損帯域のスペクトル包絡と前記入力音声のスペクトル包絡とを合成して、前記欠損帯域が補完されたスペクトル包絡を生成する。

10

## 【図面の簡単な説明】

## 【0007】

【図1】実施形態の音声処理装置の構成を示すブロック図。

【図2】実施形態の音声処理装置が実行する処理の流れを示すフローチャート。

【図3】検出部による欠損帯域の検出方法の一例を示す図。

【図4】補完部による処理の一例を示す図。

【図5】補完部による処理の他の例を示す図。

## 【発明を実施するための形態】

## 【0008】

本実施形態の音声処理装置は、任意の周波数帯域の音声成分が欠損している入力音声のスペクトル包絡から、欠損している成分を補完したスペクトル包絡を生成する。入力音声は、主に、人の発話音声を想定している。図1は、実施形態の音声処理装置の構成を示すブロック図である。図2は、実施形態の音声処理装置が実行する処理の流れを示すフローチャートである。

30

## 【0009】

本実施形態の音声処理装置は、図1に示すように、抽出部1と、検出部2と、生成部3と、変換部4と、補完部5と、を備える。

## 【0010】

抽出部1は、入力音声のスペクトル包絡  $t\_in$  から、基底モデル10を用いて、細分化された周波数帯域ごとの音声成分を表現する音声パラメータを抽出する(図2のステップS101)。なお、入力音声からスペクトル包絡  $t\_in$  を生成する処理は、音声処理装置の内部で行ってもよいし、外部で行ってもよい。

40

## 【0011】

基底モデル10は、音声のスペクトル包絡  $t$  によって形成される空間の部分空間の基底を表す基底ベクトルのセットである。本実施形態では、基底モデル10として、下記の参考文献1に記載されたサブバンド基底スペクトルモデル(以下、SBMという。)を用いる。基底モデル10は、音声処理装置内の図示しない記憶部に予め格納されてもよいし、音声処理装置の動作時に外部から取得されて保持されてもよい。

参考文献1: M Tamura, T Kagoshima, and M Akamine, "Sub-band basis spectrum model for pitch-synchronous log-spectrum and phase based on approximation of sparse coding," in Proceeding Interspeech 2010, pp. 2046 - 2049

50

, Sept . 2010 .

【 0 0 1 2 】

参考文献 1 によれば、S B M の基底は、以下の ( 1 ) ~ ( 3 ) に示す特徴を持つ。

( 1 ) 周波数軸上で単一の最大値を与えるピーク周波数を含む所定の周波数帯域に値が存在し、その周波数帯域の外側は値を零とし、フーリエ変換やコサイン変換で用いられるような周期的な基底のように同じ最大値を複数持たない。

( 2 ) 基底の数は、スペクトル包絡がもつ分析点数よりも少なく、その数は分析点数の半分未満の数となる。

( 3 ) ピーク周波数位置が隣りあう 2 つの基底間に重なりを持つ、すなわちピーク周波数が隣り合う基底は、値の存在する周波数の範囲の一部が重なる。

10

【 0 0 1 3 】

また、参考文献 1 によれば、S B M の基底を表す基底ベクトルは、下記式 ( 1 ) により定義される。

【 数 1 】

$$\phi_n(k) = \begin{cases} 0.5 - 0.5 \cos\left(\frac{k - \tilde{\Omega}(n-1)}{\tilde{\Omega}(n) - \tilde{\Omega}(n-1)} \pi\right) & (\tilde{\Omega}(n-1) \leq k < \tilde{\Omega}(n)) \\ 0.5 - 0.5 \cos\left(\frac{k - \tilde{\Omega}(n)}{\tilde{\Omega}(n+1) - \tilde{\Omega}(n)} \pi + \frac{\pi}{2}\right) & (\tilde{\Omega}(n) \leq k < \tilde{\Omega}(n+1)) \quad \dots (1) \\ 0 & (\text{otherwise}) \end{cases} \quad 20$$

ここで、 $\phi_n(k)$  は  $n$  番目の基底ベクトルの  $k$  番目の成分である。また、 $(n)$  [rad] は  $n$  番目の基底ベクトルのピーク周波数であり、下記式 ( 2 ) のように定義される。

【 数 2 】

$$\tilde{\Omega}(n) = \begin{cases} \Omega + 2 \tan^{-1} \frac{\alpha \sin \Omega}{1 - \alpha \cos \Omega} & (0 \leq n < N_w) \\ \frac{n - N_w}{N - N_w} \pi + \frac{\pi}{2} & (N_w \leq n < N) \end{cases} \quad \dots (2) \quad 30$$

ここで、 $\alpha$  は伸縮係数、 $\Omega$  は周波数 [rad]、 $N_w$  は  $(N_w) = \lfloor N/2 \rfloor$  を満たす値である。

【 0 0 1 4 】

また、S B M は、上記のような特徴を持つ基底の重み付け線形結合により、 $t$  フレーム目のスペクトル包絡  $\mathbf{t} = [t(1), t(1), \dots, t(k), \dots, t(K)]^T$  を、下記式 ( 3 ) のように表現する。

【 数 3 】

$$\mathbf{X}_t = \exp\left(\frac{1}{2} \phi \mathbf{c}_t\right) \quad \dots (3)$$

ここで、 $\mathbf{c}_t = [c_t(0), c_t(2), \dots, c_t(n), \dots, c_t(N-1)]^T$  は、S B M の基底ベクトルに対する  $t$  フレーム目の重みベクトルであり、 $\phi = [0, 1, \dots, n, \dots, N-1]$  は基底ベクトルを行列化したものである。

【 0 0 1 5 】

本実施形態では、S B M の各基底ベクトルに対応する重みベクトル  $\mathbf{c}_t$  を、音声パラメータとして扱う。この音声パラメータは、参考文献 1 に記載されている非負最小二乗誤差

50

法を用いて、スペクトル包絡  $t$  から抽出することができる。すなわち、音声パラメータとしての重みベクトル  $c$   $t$  は、音声パラメータの値が必ず零以上になるとの制約のもとで、各基底ベクトルと重みベクトル  $c$   $t$  との線形結合と、スペクトル包絡  $t$  と、の誤差が最小となるように最適化を行うことで求められる。

#### 【0016】

本実施形態では、スペクトル包絡  $t$  の分析に用いた分析点数が160以上であることを想定し、SBMの基底の数を80とする。これらの基底のうち、周波数軸上で0ラジアンから  $\pi/2$ ラジアンまでの低い周波数帯域を表現する1番目の基底から55番目の基底までは、メルケプストラム分析で用いられるオールパスフィルタの伸縮係数値（ここでは0.35）に基づいたメル尺度で作成する。また、周波数軸上で  $\pi/2$ ラジアン以上の高い周波数帯域を表現する56番目から80番目の基底は、線形尺度に基づいて作成されたものを用いる。なお、上述した低い周波数帯域の基底は、メル尺度以外の尺度、例えば線形尺度やバーク尺度、ERB尺度などを用いて作成されたものを用いてもよい。

10

#### 【0017】

なお、本実施形態では、スペクトル包絡  $t$  から音声パラメータを抽出するための基底モデル10としてSBMを用いている。しかし、スペクトル包絡  $t$  から、細分化された局所的な周波数帯域ごとの音声成分を表現した音声パラメータを抽出でき、かつ、抽出した音声パラメータから元のスペクトル包絡  $t$  を再現できるものであれば、どのような基底モデル10を用いてもよい。例えば、スパースコーディング法により求めた基底モデルや、非負値行列分解によって求めた基底行列を、スペクトル包絡  $t$  から音声パラメータを抽出するための基底モデル10として用いることができる。また、スペクトル包絡  $t$  から、細分化された局所的な周波数帯域ごとの音声成分を表現した音声パラメータを抽出でき、かつ、抽出した音声パラメータから元のスペクトル包絡  $t$  を再現できるのであれば、サブバンド分割やフィルタバンクによる表現を用いて、音声パラメータを抽出してもよい。

20

#### 【0018】

検出部2は、入力音声のスペクトル包絡  $t_{in}$ 、または、このスペクトル包絡  $t_{in}$  から抽出部1によって抽出された音声パラメータの包絡形状を解析し、入力音声のスペクトル包絡  $t_{in}$  において音声成分が欠損している周波数帯域である欠損帯域を検出する（図2のステップS102）。

30

#### 【0019】

検出部2は、例えば、入力音声のスペクトル包絡  $t_{in}$ 、または、このスペクトル包絡  $t_{in}$  から抽出された音声パラメータに対して、周波数軸方向の1次の変化の割合および2次の変化の割合を用いて、欠損帯域を検出することができる。

#### 【0020】

図3は、検出部2による欠損帯域の検出方法の一例を示す図である。図3に示す例は、入力音声が低域通過特性を持つ伝送路を通過することで高周波側の成分が欠損した場合の例であり、スペクトル包絡  $t_{in}$  から抽出された音声パラメータの包絡形状を解析して欠損帯域を検出する例である。図の横軸は周波数軸であり、数値は基底の番号を表している。図3(a)は、入力音声のスペクトル包絡  $t_{in}$  から抽出部1により抽出された音声パラメータの周波数軸方向の変化を表すグラフ図であり、縦軸は音声パラメータの値を示している。また、図3(b)は、図3(a)に示した音声パラメータの周波数軸方向の1次変化の割合を表すグラフ図であり、縦軸は音声パラメータを1次微分した値を示している。また、図3(b)は、図3(a)に示した音声パラメータの周波数軸方向の2次変化の割合を表すグラフ図であり、縦軸は音声パラメータを2階微分した値を示している。

40

#### 【0021】

検出部2は、まず、図3(b)に示す音声パラメータの1次の変化の割合から、値が最小となる次元（以下、第1の基準位置という。）を、次元が大きい方から探索して決定する。次に、検出部2は、第1の基準位置とこの位置から数次元小さい次元との間の範囲を

50

探索範囲として、図3(c)に示す音声パラメータの2次の変化の割合から、探索範囲内で値が最小となる次元(以下、第2の基準位置という。)を求める。そして、検出部2は、第2の基準点より1つ小さい次元の位置を、欠損帯域の低周波側の端部である開始位置とする。また、図3に示す例では、高周波側の成分が欠損している場合を想定しているため、欠損帯域の高周波側の端部である終了位置は、最大の次元の位置とする。検出部2は、上記のように決定された開始位置と終了位置との間の周波数帯域を、欠損帯域として検出することができる。

#### 【0022】

入力音声が高域通過特性を持つ伝送路を通過することで低周波側の成分が欠損している場合には、次元の小さい方から上記と同様の処理を行うことで、欠損帯域を検出することができる。すなわち、検出部2は、まず、音声パラメータの1次の変化の割合を次元が小さいほうから探索して、第1の基準位置を決定する。次に、検出部2は、第1の基準位置とこの位置から数次元大きい次元との間の範囲を探索範囲として、音声パラメータの2次の変化の割合から、第2の基準位置を求める。そして、検出部2は、第2の基準位置より1つ大きい次元の位置を、欠損帯域の高周波側の端部である終了位置とする。また、この場合は、欠損帯域の低周波側の端部である開始位置は、最小の次元の位置とする。検出部2は、上記のように決定された開始位置と終了位置との間の周波数帯域を、欠損帯域として検出することができる。

#### 【0023】

また、入力音声帯域遮断特性を持つ伝送路を通過することで、低周波と高周波の間の任意の周波数帯域の成分が欠損している場合には、検出部2は、例えば以下の方法で欠損帯域を検出することができる。すなわち、検出部2は、まず、スペクトル傾斜情報を取り除いた音声パラメータに対して、低次元側からの1次の変化の割合および2次の変化の割合を求め、1次の変化の割合の最小値および最大値となる次元をそれぞれ求めて、これらを第1の基準位置とする。次に、検出部2は、最小値となる第1の基準位置より小さい次元において2次の変化の割合が最小となる点を求める。同様に、検出部2は、最大値となる第1の基準位置より大きな次元において変化の割合が最小となる点を求め、それぞれを第2の基準位置とする。そして、検出部2は、これら2つの第2の基準位置に基づいて、低次元側を開始位置、高次元側を終了位置として定める。検出部2は、上記のように定めた開始位置と終了位置との間の周波数帯域を、欠損帯域として検出することができる。

#### 【0024】

入力音声の伝送路の特性によって欠損帯域が生じる場合、欠損帯域は入力音声ごとに一定であることが想定される。したがって、検出部2は、入力音声の少なくとも1つのフレームに対して上述した処理を行うことで、欠損帯域の検出が可能である。ただし、検出部2は、入力音声の複数のフレームを対象として対して上述した処理を行うようにすれば、欠損帯域の検出をより精度よく行うことができる。この場合、検出部2は、例えば、複数フレームの音声パラメータの平均値を次元ごとに求め、求めた平均値の1次の変化の割合および2次の変化の割合を用いて、欠損位置を精度よく検出することができる。また、検出部2は、複数フレームの音声パラメータに対してそれぞれ上述した処理をそれぞれ行って、得られた結果をマージすることで、最終的な欠損帯域を検出するようにしてもよい。

#### 【0025】

また、検出部2は、入力音声の各フレームに対して上述した処理を繰り返し行うようにすれば、突発的な要因によって入力音声における欠損帯域がフレーム間で異なる場合であっても、フレーム間で異なる欠損位置をそれぞれ検出することができる。

#### 【0026】

なお、上述した処理は、入力音声のスペクトル包絡  $t\_in$  から抽出された音声パラメータを処理対象としたが、入力音声のスペクトル包絡  $t\_in$  そのものを処理対象としても、同様の処理によって欠損帯域を検出することができる。すなわち、入力音声のスペクトル包絡  $t\_in$  に対して、周波数軸方向の1次の変化の割合および2次の変化の割合を用いて上記と同様の処理を行うようにしても、欠損帯域を検出することができる。

## 【 0 0 2 7 】

生成部 3 は、検出部 2 により検出された欠損帯域の位置と、統計情報 2 0 と、入力音声のスペクトル包絡  $t\_in$  から抽出部 1 によって抽出された音声パラメータとに基づいて、欠損帯域に対応する音声パラメータを生成する（図 2 のステップ S 1 0 3）。

## 【 0 0 2 8 】

統計情報 2 0 は、音声成分が欠損していない音声のスペクトル包絡から抽出された音声パラメータ（抽出部 1 が入力音声のスペクトル包絡  $t\_in$  から抽出する音声パラメータと同様の音声パラメータ）を用いて事前に作成されている。ここで、統計情報とは、音声パラメータベクトルの平均、分散やヒストグラムなどにより、音声パラメータをモデル化したものであり、例えばコードブック、混合分布モデル、隠れマルコフモデルなどである。本実施形態では、統計情報 2 0 として混合正規分布モデル（以下、GMM という。）を用いる。統計情報 2 0 は、音声処理装置内の図示しない記憶部に予め格納されてもよいし、音声処理装置の動作時に外部から取得されて保持されてもよい。

10

## 【 0 0 2 9 】

GMM では、重みベクトル  $c_t$  の確率密度関数は、下記式（4）のように表される。

## 【数 4】

$$P(c_t|\lambda) = \sum_{m=1}^M P(c_t|m, \lambda) = \sum_{m=1}^M \alpha_m N(c_t; \mu_m^{(c)}, \Sigma_m^{(cc)}) \quad \dots (4)$$

ここで、 $\lambda$  は GMM のパラメータセットを表し、 $N(c_t; \mu_m^{(c)}, \Sigma_m^{(cc)})$  は平均ベクトル  $\mu_m^{(c)}$ 、全共分散行列  $\Sigma_m^{(cc)}$  をもつ GMM の  $m$  番目の正規分布であり、 $\alpha_m$  は  $m$  番目の正規分布の重みである。

20

## 【 0 0 3 0 】

なお、本実施形態において、残存帯域（欠損帯域以外の帯域）に対応するパラメータ成分（以下、残存帯域成分という。）の数と、欠損帯域に対応するパラメータ成分（以下、欠損帯域成分という。）の数が異なることを想定している。このため全共分散行列、すなわち、行列のすべての成分にある値を有するものを用いている。しかし、実施形態において残存帯域成分の数と欠損帯域成分の数が常に同数である場合には、全共分散行列の代わりに、行列の対角成分と事前に決定した残存帯域成分とそれに対応する欠損帯域成分とに値を有し、それ以外の成分は零であるような分散行列を用いてもよい。

30

## 【 0 0 3 1 】

本実施形態では、音声成分が欠損していない（欠損帯域のない）複数の話者の発話音声から抽出された音声パラメータを学習データとして用いて事前に構築された統計モデルである不特定話者 GMM を、統計情報 2 0 として用いる。統計情報 2 0 の構築には、例えば、LGB アルゴリズムや EM アルゴリズムなどを用いることができる。

## 【 0 0 3 2 】

生成部 3 は、統計情報 2 0 としての GMM を用いて、残存帯域成分から欠損帯域成分を生成するための規則を、次のような手順で求める。

40

## 【 0 0 3 3 】

生成部 3 は、まず、統計情報 2 0 としての GMM を、検出部 2 により検出された欠損帯域の位置、すなわち、上述した開始位置および終了位置に基づいて、音声パラメータベクトル、平均ベクトル  $\mu_m(c)$ 、および共分散行列  $\Sigma_m(cc)$  を分割して、下記式（5）のように変形する。



【数 5】

$$P(c_t|\lambda) = \sum_{m=1}^M \alpha_m N \left( \begin{bmatrix} c_t^{(r)} \\ c_t^{(l)} \end{bmatrix}; \begin{bmatrix} \mu_m^{(r)} \\ \mu_m^{(l)} \end{bmatrix}, \begin{bmatrix} \sum_m^{(r)} \sum_m^{(rl)} \\ \sum_m^{(lr)} \sum_m^{(l)} \end{bmatrix} \right) \quad \dots (5)$$

ここで、 $c_t^{(r)}$ は残存帯域に関する音声パラメータベクトル、 $c_t^{(l)}$ は欠損帯域に関する音声パラメータベクトル、 $\mu_m^{(r)}$ は残存帯域に関する平均ベクトル、 $\mu_m^{(l)}$ は欠損帯域に関する平均ベクトル、 $\sum_m^{(r)}$ は残存帯域に関する自己共分散行列、 $\sum_m^{(l)}$ は欠損帯域に関する自己共分散行列、 $\sum_m^{(lr)}$ は欠損帯域と残存帯域に関する相互共分散行列である。

10

【0034】

次に、生成部3は、この変形したGMMを、下記式(6)に示すように、残存帯域の音声パラメータベクトルに対する欠損帯域の音声パラメータベクトルの条件付き確率分布へと変形する。そして、生成部3は、式(6)に示す条件付き確率分布を規則として用いて、残存帯域成分(入力音声のスペクトル包絡  $t\_in$  から抽出された音声パラメータ)から、欠損帯域成分(欠損帯域に対応する音声パラメータ)を生成する。

20

【数 6】

$$P(c_t^{(l)}|c_t^{(r)}, \lambda) = \sum_{m=1}^M P(m|c_t^{(r)}, \lambda) P(c_t^{(l)}|c_t^{(r)}, m, \lambda) \quad \dots (6)$$

ここで、

$$P(m|c_t^{(r)}, \lambda) = \frac{\alpha_m N(c_t^{(r)}; \mu_m^{(r)}, \sum_m^{(r)})}{\sum_{m=1}^M \alpha_m N(c_t^{(r)}; \mu_m^{(r)}, \sum_m^{(r)})} \quad \dots (7)$$

30

$$P(c_t^{(l)}|c_t^{(r)}, m, \lambda) = N(c_t^{(l)}; E_{m,t}^{(l)}, D_m^{(l)}) \quad \dots (8)$$

$$E_{m,t}^{(l)} = \sum_m^{(lr)} \sum_m^{(rr)^{-1}} (c_t^{(r)} - \mu_m^{(r)}) + \mu_m^{(l)} \quad \dots (9)$$

$$D_m^{(l)} = \sum_m^{(ll)} - \sum_m^{(lr)} \sum_m^{(rr)^{-1}} \sum_m^{(rl)} \quad \dots (10)$$

である。

これより欠損帯域の音声パラメータ $\tilde{c}_t^{(l)}$ は最小二乗誤差基準により、下記式(11)のように求まる。

$$\tilde{c}_t^{(l)} = \sum_{m=1}^M P(m|c_t^{(r)}, \lambda) E_{m,t}^{(l)} \quad \dots (11)$$

40

【0035】

本実施形態においては、上述したように、1つの入力音声における欠損帯域がフレーム間で一定であることを想定している。この場合、上述したように、フレームごとに欠損帯域に対応する音声パラメータを生成すると、フレーム間で不連続が生じることが考えられる。そこで、この不連続を緩和させるために、生成部3は、当該フレームと前後数フレー

50

ムを用いて移動平均フィルタ、中央値フィルタ、加重平均フィルタ、ガウスフィルタなどにより平滑化処理を行うことで、欠損帯域に対応する音声パラメータのフレーム間における不連続性を緩和させてもよい。

【0036】

また、生成部3により生成された欠損帯域に対応する音声パラメータは、汎化されたGMMの影響により平滑化されている。そのため、生成部3は、欠損帯域に対応する音声パラメータを生成した後に、下記の参考文献2で示される系列内変動（以下、GVという）の統計情報や音声パラメータのヒストグラムを用いたパラメータ強調を行ってもよい。

参考文献2：藤敦渉、他4名、「GMMに基づく最尤変換法による携帯電話音声の帯域拡張」，社団法人 情報処理学会 研究報告（IPSS SIG Technical Report），2007年7月21日，p.63-68

【0037】

さらに、生成部3は、フレーム間の不連続性や音声パラメータの平滑化を防ぐために、参考文献2で示されている、動的特徴量を用いた尤度最大化基準によるGMM変換手法を用いて、欠損帯域に対応する音声パラメータを生成してもよい。この場合、GMMの学習においては、音声パラメータである重みベクトル $c_t$ と、この重みベクトル $c_t$ の時間変化成分 $\dot{c}_t$ とを結合した下記式(12)で示す特徴量 $C_t$ を用意し、下記式(13)に示すGMMを構築して、これを統計情報20として保持する。

【数7】

$$C_t = [c_t^\Gamma, \Delta c_t^\Gamma]^\Gamma \quad \dots (12)$$

$$P(C_t|\lambda) = \sum_{m=1}^M P(C_t|m, \lambda) = \sum_{m=1}^M \alpha_m N(C_t; \mu_m^{(CC)}, \Sigma_m^{(CC)}) \quad \dots (13)$$

ここで、 $\mu_m^{(C)}$ は $m$ 番目の分布が保持する結合特徴量の平均ベクトルであり、 $\Sigma_m^{(CC)}$ は $m$ 番目の分布が保持する結合特徴量の全共分散行列である。

【0038】

式(13)に示すGMMを統計情報20として用いる場合においても、生成部3は、まず、検出部2により検出された欠損帯域の位置（開始位置および終了位置）に基づいてGMMを残存帯域成分と欠損帯域成分とに分割し、式(13)を下記式(14)のように変形する。

【数8】

$$P(C_t|\lambda) = \sum_{m=1}^M \alpha_m N \left( \begin{bmatrix} C_t^{(R)} \\ C_t^{(L)} \end{bmatrix}; \begin{bmatrix} \mu_m^{(R)} \\ \mu_m^{(L)} \end{bmatrix}, \begin{bmatrix} \Sigma_m^{(RR)} & \Sigma_m^{(RL)} \\ \Sigma_m^{(LR)} & \Sigma_m^{(LL)} \end{bmatrix} \right) \quad \dots (14)$$

【0039】

次に、生成部3は、式(14)に示すGMMを、下記式(15)に示すように、残存帯域の音声パラメータベクトルに対する欠損帯域の音声パラメータベクトルの条件付き確率分布へと変形する。

【数9】

$$P(C_t^{(L)}|C_t^{(R)}, \lambda) = \sum_{m=1}^M P(m|C_t^{(R)}, \lambda) P(C_t^{(L)}|C_t^{(R)}, m, \lambda) \quad \dots (15)$$

【0040】

そして、生成部3は、尤度最大化基準で、下記式(16)および下記式(17)に示す

10

20

30

40

50

ように、欠損帯域の音声パラメータを生成する。

【数 1 0】

$$\left[ \tilde{c}_1^{(l)\Gamma}, \tilde{c}_2^{(l)\Gamma}, \dots, \tilde{c}_T^{(l)\Gamma} \right]^\Gamma = \arg \max_{c^{(l)}} \prod_{t=1}^T \sum_{m=1}^M P(m | C_t^{(R)}, \lambda) P(C_t^{(L)}, C_t^{(R)}, m, \lambda) \quad \dots (16)$$

$$\text{Subject to } \left[ C_1^{(L)\Gamma}, C_2^{(L)\Gamma}, \dots, C_T^{(L)\Gamma} \right]^\Gamma = W \left[ c_1^{(l)\Gamma}, c_2^{(l)\Gamma}, \dots, c_T^{(l)\Gamma} \right]^\Gamma \quad \dots (17)$$

ここで、Wは音声パラメータ系列から音声パラメータと時間変化量成分との結合特徴量系列へと変換するための行列を表す。 10

【0041】

また、生成部3は、式(16)の代わりに、参考文献2で示される準最尤分布系列からのパラメータ生成やGVを用いたパラメータ生成法を用いて、欠損帯域に対応する音声パラメータを生成してもよいし、式(16)による音声パラメータの生成後に、GVやヒストグラムを用いたパラメータ強調を行ってもよい。

【0042】

なお、本実施形態では、統計情報20として不特定話者GMMを使用することを想定している。しかし、不特定話者GMMのほかに、複数の特定話者GMMを統計情報20として用いてもよい。この場合、生成部3は、入力音声のスペクトル包絡  $t\_in$  から抽出された音声パラメータに最も適合した特定話者GMM、または適合度に合わせて複数の特定話者GMMを線形結合したものをを用いて、欠損帯域に対応する音声パラメータの生成を行う。これにより、欠損帯域の音声パラメータを、入力音声のスペクトル包絡  $t\_in$  から抽出された音声パラメータに適合するように生成することができる。 20

【0043】

さらに、入力音声のスペクトル包絡  $t\_in$  から抽出された音声パラメータとの適合性を向上させるために、不特定話者GMMないしは特定話者GMMに対して、線形回帰や最大事後確率推定などの統計的な音声認識や音声合成で用いられている話者適応手法を適用し、入力音声のスペクトル包絡  $t\_in$  から抽出された音声パラメータと適合したGMMを用いて、欠損帯域に対応する音声パラメータを生成してもよい。 30

【0044】

変換部4は、生成部3が生成した欠損帯域に対応する音声パラメータを、基底モデル10を用いて、欠損帯域のスペクトル包絡に変換する(図2のステップS104)。

【0045】

本実施形態では、基底モデル10としてSBMを用いるため、上記式(3)に示したような処理を行うことで、欠損帯域に対応する音声パラメータとして生成された重みベクトル  $c_t$  を、欠損帯域の音声スペクトル包絡  $\sim t$  に変換することができる。すなわち、変換部4は、欠損帯域に対応する音声パラメータである重みベクトル  $c_t$  と、この欠損帯域に対応する基底ベクトルとを線形結合することにより、欠損帯域のスペクトル包絡  $\sim t$  を求めることができる。 40

【0046】

補完部5は、変換部4により得られた欠損帯域のスペクトル包絡  $\sim t$  と、入力音声のスペクトル包絡  $t\_in$  とを合成して、欠損帯域が補完されたスペクトル包絡  $t\_out$  を生成する(図2のステップS105)。

【0047】

補完部5は、例えば、入力音声のスペクトル包絡  $t\_in$  のうち、検出部2により検出された欠損帯域の位置(開始位置と終了位置との間の帯域)に、変換部4により得られた欠損帯域のスペクトル包絡  $\sim t$  を当てはめるとともに、不連続性を緩和させる処理を行ってこれらを合成することで、欠損帯域が補完されたスペクトル包絡  $t\_out$  を生成することができる。 50

## 【0048】

図4は、補完部5による処理の一例を示す図である。図4に示す例は、低域通過特性を持つ伝送路により高周波側の成分が欠損した入力音声のスペクトル包絡  $t\_in$  から、欠損帯域が補完されたスペクトル包絡  $t\_out$  を生成する例である。

## 【0049】

入力音声のスペクトル包絡  $t\_in$  の欠損帯域の位置に、変換部4により得られた欠損帯域のスペクトル包絡  $\sim t$  をそのまま当てはめると、欠損帯域の境界位置にて2つのスペクトル包絡の値が大きくなり、不連続性が発生する場合がある。そこで、補完部5は、まず、欠損帯域の境界位置における2つのスペクトル包絡の差分  $d$  を計測する(図4(a))。そして、補完部5は、計測した差分  $d$  に基づき、変換部4により得られた欠損帯域のスペクトル包絡  $\sim t$  の全体にバイアス補正を行う(図4(b))。

10

## 【0050】

次に、補完部5は、入力音声のスペクトル包絡  $t\_in$  と欠損帯域のスペクトル包絡  $\sim t$  とが滑らかに接続されるように、それぞれのスペクトル包絡の境界位置周辺の成分に対して片側ハン窓をかけ(図4(c))、該当する箇所のスペクトル包絡の成分を加算することで、入力音声のスペクトル包絡  $t\_in$  と欠損帯域のスペクトル包絡  $\sim t$  とを合成する(図4(d))。これにより、欠損帯域が補完されたスペクトル包絡  $t\_out$  が生成される。

## 【0051】

なお、高域通過特性を持つ伝送路により低周波側の成分が欠損した入力音声のスペクトル包絡  $t\_in$  から、欠損帯域が補完されたスペクトル包絡  $t\_out$  を生成する場合も、上記と同様の手順で、欠損帯域が補完されたスペクトル包絡  $t\_out$  を適切に生成することができる。

20

## 【0052】

図5は、補完部5による処理の他の例を示す図である。図5に示す例は、帯域遮断特性を持つ伝送路により低周波と高周波の間の任意の周波数帯域の成分が欠損した入力音声のスペクトル包絡  $t\_in$  から、欠損帯域が補完されたスペクトル包絡  $t\_out$  を生成する例である。

## 【0053】

図5の例の場合、補完部5は、欠損帯域の開始位置における2つのスペクトル包絡の差分  $d_s$  を計測するとともに、欠損帯域の終了位置における2つのスペクトル包絡の差分  $d_e$  を計測する(図5(a))。そして、補完部5は、欠損帯域の開始位置で計測された差分  $d_s$  と、欠損帯域の終了位置で計測された差分  $d_e$  とに基づき、変換部4により得られた欠損帯域のスペクトル包絡  $\sim t$  に対して傾斜補正をかける(図5(b))。

30

## 【0054】

次に、補完部5は、欠損帯域の開始位置と終了位置の双方において、入力音声のスペクトル包絡  $t\_in$  と欠損帯域のスペクトル包絡  $\sim t$  とが滑らかに接続されるように、これら開始位置および終了位置の周辺におけるそれぞれのスペクトル包絡の成分に対して片側ハン窓をかけ(図5(c))、該当する箇所のスペクトル包絡の成分を加算することで、入力音声のスペクトル包絡  $t\_in$  と欠損帯域のスペクトル包絡  $\sim t$  とを合成する(図5(d))。これにより、欠損帯域が補完されたスペクトル包絡  $t\_out$  が生成される。

40

## 【0055】

本実施形態の音声処理装置は、補完部5により生成された、欠損帯域が補完されたスペクトル包絡  $t\_out$  を外部に出力することができる。また、本実施形態の音声処理装置は、欠損帯域が補完されたスペクトル包絡  $t\_out$  から音声を復元し、復元した音声を出力するようにしてもよい。

## 【0056】

以上、具体的な例を挙げながら詳細に説明したように、本実施形態の音声処理装置によれば、任意の周波数帯域で欠損した音声成分を適切に補完することができる。

50

## 【 0 0 5 7 】

なお、本実施形態の音声処理装置は、例えば、汎用のコンピュータ装置を基本ハードウェアとして用いて実現することが可能である。すなわち、本実施形態の音声処理装置は、汎用のコンピュータ装置に搭載されたプロセッサにプログラムを実行させることにより実現することができる。このとき、音声処理装置は、上記のプログラムをコンピュータ装置にあらかじめインストールすることで実現してもよいし、CD-ROMなどの記憶媒体に記憶して、あるいはネットワークを介して上記のプログラムを配布して、このプログラムをコンピュータ装置に適宜インストールすることで実現してもよい。また、上記のプログラムをサーバーコンピュータ装置上で実行させ、ネットワークを介してその結果をクライアントコンピュータ装置で受け取ることにより実現してもよい。

10

## 【 0 0 5 8 】

また、本実施形態の音声処理装置で使用する各種情報は、上記のコンピュータ装置に内蔵あるいは外付けされたメモリ、ハードディスクもしくはCD-R、CD-RW、DVD-RAM、DVD-Rなどの記録媒体を適宜利用して格納しておくことができる。例えば、本実施形態の音声処理装置が使用する基底モデル10や統計情報20は、これら記録媒体を適宜利用して格納しておくことができる。

## 【 0 0 5 9 】

本実施形態の音声処理装置で実行されるプログラムは、音声処理装置を構成する各処理部（抽出部1、検出部2、生成部3、変換部4および補完部5）を含むモジュール構成となっており、実際のハードウェアとしては、例えば、プロセッサが上記記憶媒体からプログラムを読み出して実行することにより、上記各処理部が主記憶装置上にロードされ、主記憶装置上に生成されるようになっている。

20

## 【 0 0 6 0 】

以上、本発明の実施形態を説明したが、ここで説明した実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。ここで説明した新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。ここで説明した実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

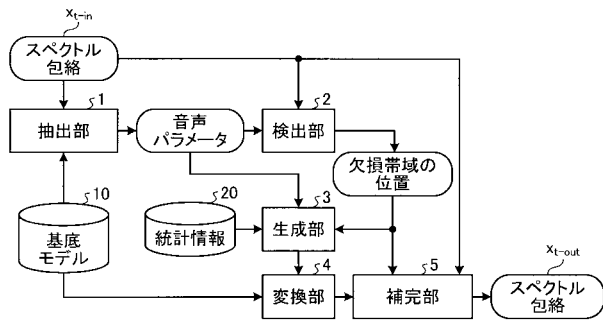
## 【 符号の説明 】

30

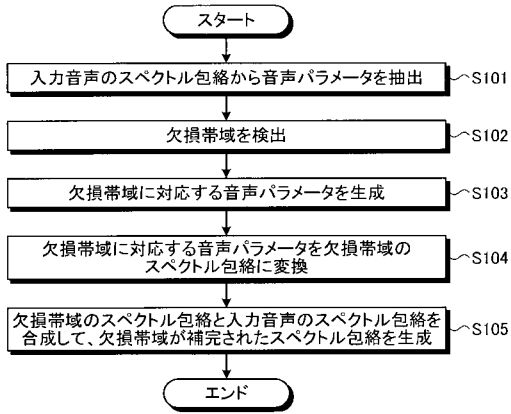
## 【 0 0 6 1 】

- 1 抽出部
- 2 検出部
- 3 生成部
- 4 変換部
- 5 補完部
- 10 基底モデル
- 20 統計情報

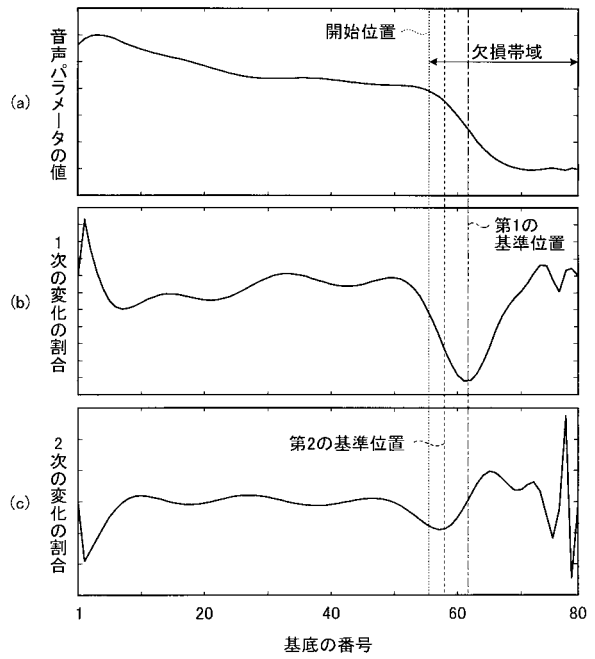
【 図 1 】



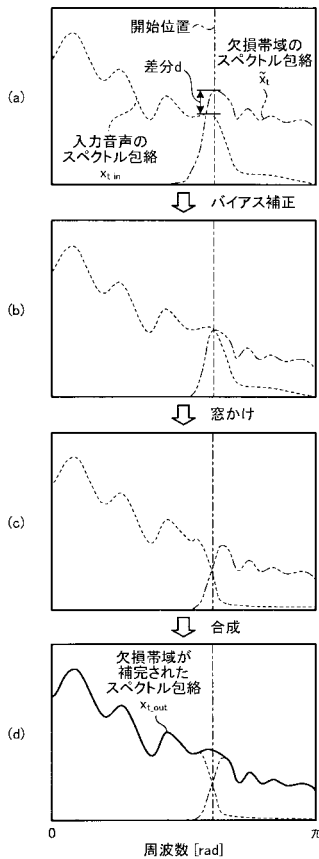
【 図 2 】



【 図 3 】



【 図 4 】



【 図 5 】

