

(19)日本国特許庁(JP)

## (12)特許公報(B2)

(11)特許番号

特許第7023934号

(P7023934)

(45)発行日 令和4年2月22日(2022.2.22)

(24)登録日 令和4年2月14日(2022.2.14)

(51)国際特許分類	F I
G 1 0 L 15/16 (2006.01)	G 1 0 L 15/16
G 1 0 L 17/18 (2013.01)	G 1 0 L 17/18

請求項の数 19 (全25頁)

(21)出願番号	特願2019-510589(P2019-510589)	(73)特許権者	511050697
(86)(22)出願日	平成29年8月24日(2017.8.24)		アリババ グループ ホウルディング リ
(65)公表番号	特表2019-528476(P2019-528476		ミテッド
	A)		英国領ケイマン諸島 グランド ケイマン
(43)公表日	令和1年10月10日(2019.10.10)		ジョージ タウン ピーオーボックス 8
(86)国際出願番号	PCT/US2017/048499		47 ワン キャピタル プレイス フォー
(87)国際公開番号	WO2018/039500		ス フロア
(87)国際公開日	平成30年3月1日(2018.3.1)	(74)代理人	100079108
審査請求日	令和2年7月6日(2020.7.6)		弁理士 稲葉 良幸
(31)優先権主張番号	201610741622.9	(74)代理人	100109346
(32)優先日	平成28年8月26日(2016.8.26)		弁理士 大貫 敏史
(33)優先権主張国・地域又は機関	中国(CN)	(74)代理人	100117189
			弁理士 江口 昭彦
		(74)代理人	100134120
			弁理士 内藤 和彦

最終頁に続く

(54)【発明の名称】 音声認識方法及び装置

## (57)【特許請求の範囲】

## 【請求項1】

音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出することと、  
前記話者認識特徴を含む前記ベクトルに従って第2のニューラルネットワーク内のバイアスを補償することと、  
前記第2のニューラルネットワークに基づく音響モデルを介し前記音声データ内の音声を認識することと、を含む音声認識方法。

## 【請求項2】

前記話者認識特徴を含む前記ベクトルに従って前記第2のニューラルネットワーク内のバイアスを補償することは、前記話者認識特徴を含む前記ベクトルに前記第2のニューラルネットワークのバイアス項となるべき重み行列を掛けることを含む、請求項1に記載の音声認識方法。

## 【請求項3】

前記第1のニューラルネットワーク、前記第2のニューラルネットワーク、及び前記重み行列は、前記第1のニューラルネットワーク及び前記第2のニューラルネットワークをそれぞれトレーニングし、次に前記トレーニングされた第1のニューラルネットワーク、前記重み行列、及び前記トレーニングされた第2のニューラルネットワークを一括してトレーニングすることによりトレーニングされる、請求項2に記載の音声認識方法。

## 【請求項4】

前記第 1 のニューラルネットワーク、前記第 2 のニューラルネットワーク、及び前記重み行列を初期化することと、

所定客観的判定基準に従って逆伝搬アルゴリズムを使用することにより前記重み行列を更新することと、

所定客観的判定基準に従って前記逆伝搬アルゴリズムを使用することにより前記第 2 のニューラルネットワーク及び接続行列を更新することと、をさらに含む請求項 3 に記載の音声認識方法。

【請求項 5】

前記話者認識特徴は少なくとも話者声紋情報を含む、請求項 1 乃至 4 のいずれか一項に記載の音声認識方法。

10

【請求項 6】

前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内のバイアスを補償することは、前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することを含み、

前記話者認識特徴を含む前記ベクトルは前記第 1 のニューラルネットワーク内の最後の隠れ層の出力ベクトルである、請求項 1 に記載の音声認識方法。

【請求項 7】

前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することは、前記第 1 のニューラルネットワークの前記最後の隠れ層においてノードにより出力された前記話者認識特徴を含む前記ベクトルを、前記第 2 のニューラルネットワーク内の前記入力層を除く層の前記すべて又は一部に対応するバイアスノードへ送信することを含む、請求項 6 に記載の音声認識方法。

20

【請求項 8】

前記音声データは、収集された元音声データ又は前記収集された元音声データから抽出された音声特徴である、請求項 1 に記載の音声認識方法。

【請求項 9】

前記話者認識特徴は、様々なユーザ又は様々なユーザのクラスタに対応する、請求項 1 に記載の音声認識方法。

30

【請求項 10】

一組の命令を格納する非一時的コンピュータ可読媒体であって、前記一組の命令は、装置の 1 つ又は複数のプロセッサによって、前記装置に音声認識の方法を行わせるように、実行可能であり、前記方法は、

音声データから話者認識特徴を含むベクトルを第 1 のニューラルネットワークを介し抽出することと、

前記話者認識特徴を含む前記ベクトルに従って第 2 のニューラルネットワーク内のバイアスを補償することと、

前記第 2 のニューラルネットワークに基づく音響モデルを介し前記音声データ内の音声を認識することと、を含む、非一時的コンピュータ可読媒体。

40

【請求項 11】

前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内のバイアスを補償することは、前記話者認識特徴を含む前記ベクトルに前記第 2 のニューラルネットワークのバイアス項となるべき重み行列を掛けることを含む、請求項 10 に記載の非一時的コンピュータ可読媒体。

【請求項 12】

前記第 1 のニューラルネットワーク、前記第 2 のニューラルネットワーク、及び前記重み行列は、前記第 1 のニューラルネットワーク及び前記第 2 のニューラルネットワークをそれぞれトレーニングし、次に前記トレーニングされた第 1 のニューラルネットワーク、前記重み行列、及び前記トレーニングされた第 2 のニューラルネットワークを一括してトレ

50

ーニングすることによりトレーニングされる、請求項 1 1 に記載の非一時的コンピュータ可読媒体。

【請求項 1 3】

前記一組の命令は、前記装置の前記 1 つ又は複数のプロセッサにより、前記装置に、前記第 1 のニューラルネットワーク、前記第 2 のニューラルネットワーク、及び前記重み行列を初期化することと、

所定客観的判定基準に従って逆伝搬アルゴリズムを使用することにより前記重み行列を更新することと、

所定客観的判定基準に従って前記逆伝搬アルゴリズムを使用することにより前記第 2 のニューラルネットワーク及び接続行列を更新することと、をさらに行わせるように、実行可能である、請求項 1 2 に記載の非一時的コンピュータ可読媒体。

10

【請求項 1 4】

前記話者認識特徴は少なくとも話者声紋情報を含む、請求項 1 0 に記載の非一時的コンピュータ可読媒体。

【請求項 1 5】

前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内のバイアスを補償することは、前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することを含み、

前記話者認識特徴を含む前記ベクトルは前記第 1 のニューラルネットワーク内の最後の隠れ層の出力ベクトルである、請求項 1 0 に記載の非一時的コンピュータ可読媒体。

20

【請求項 1 6】

前記話者認識特徴を含む前記ベクトルに従って前記第 2 のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することは、前記第 1 のニューラルネットワークの前記最後の隠れ層においてノードにより出力された前記話者認識特徴を含む前記ベクトルを、前記第 2 のニューラルネットワーク内の前記入力層を除く層の前記すべて又は一部に対応するバイアスノードへ送信することを含む、請求項 1 5 に記載の非一時的コンピュータ可読媒体。

【請求項 1 7】

前記音声データは、収集された元音声データ又は前記収集された元音声データから抽出された音声特徴である、請求項 1 0 に記載の非一時的コンピュータ可読媒体。

30

【請求項 1 8】

前記話者認識特徴は、様々なユーザ又は様々なユーザのクラスタに対応する、請求項 1 0 に記載の非一時的コンピュータ可読媒体。

【請求項 1 9】

音声データから話者認識特徴を含むベクトルを第 1 のニューラルネットワークを介し抽出するように構成された抽出ユニットと、

前記話者認識特徴を含む前記ベクトルに従って第 2 のニューラルネットワーク内のバイアスを補償し、前記第 2 のニューラルネットワークに基づく音響モデルを介し前記音声データ内の音声を認識するように構成された認識ユニットと、を含む音声認識装置。

40

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

[001] 本出願は、参照のためその全体を本明細書に援用する 2016 年 8 月 26 日申請の中国特許出願第 201610741622.9 号への優先権の恩恵を主張する。

【0002】

技術分野

[002] 本出願は、音声認識に関し、より具体的には音声認識方法及び装置に関する。

【背景技術】

50

## 【 0 0 0 3 】

## 背景

[003] 現在のところ、大きな進歩が話者非依存 ( S I : speaker independent ) 音声認識システムに対しなされてきた。しかし、様々なユーザ間の差異が特定ユーザの音声認識システムの性能劣化を生じさせ得る。

## 【 0 0 0 4 】

[004] 話者依存 ( S D : speaker dependent ) 音声認識システムは S I 音声認識システムの性能劣化の問題を解決し得る。しかし、 S D 音声認識システムはトレーニングのための大量のユーザ音声データの入力を必要とし、これはユーザの大きな不都合と高コストとを生じる。

## 【 0 0 0 5 】

[005] 話者適応化 ( speaker adaptation ) 技術は S I 及び S D 音声認識システムの欠点のある程度まで補い得る。話者適応化技術により、 S D 音声特徴は S I 音声特徴へ変換され得、 S I 音声特徴は次に認識のために S I 音響モデルへ提供される。代替的に、 S I 音響システムは S D 音響システムへ変換され得る。次に、 S D 音声特徴が認識される。

## 【 0 0 0 6 】

[006] S I 音声認識システムと比較して、話者適応化技術は、ユーザ個人差を有する音声特徴を考慮し、したがってより良好な認識性能を有し得る。 S D 認識システムと比較して、話者適応化技術は、 S I システムの事前情報を導入し、したがって必要とされるユーザ音声データの量は著しく低減される。

## 【 0 0 0 7 】

[007] 話者適応化技術は、ユーザ音声データが予め取得されるかどうかによって依存してオフライン話者適応化技術とオンライン話者適応化技術とに分割され得る。オンライン話者適応化技術により、音声認識システムのパラメータは、現在のユーザ音声入力に従って等間隔 (例えば 6 0 0 m s ) で調整され得、これにより話者適応化を実現する。

## 【 0 0 0 8 】

[008] 現時点で、オンライン話者適応化方法の解決策が図 1 に示される。この解決策は、ユーザの音声特徴とユーザに関して抽出された  $i$  ベクトル (すなわち識別可能ベクトル) とを繋ぐことを含み得る。この解決策はまた、繋がれた特徴を音声認識のためにディープニューラルネットワーク ( D N N : deep neural network ) 内へ入力することを含み得る。  $i$  ベクトルの抽出プロセスは、平均スーパーベクトルを取得するために音声の音響特性をガウス混合モデルに入力することと、  $i$  ベクトルを取得するために平均スーパーベクトルに T 行列を掛けることを含み得る。ユーザが話している時、この解決策によると、  $i$  ベクトルがユーザの音声の始めの部分から抽出され得る。抽出された  $i$  ベクトルは、ユーザの音声の残りの音声認識のために使用され、こうしてオンライン話者適応化を実現する。

## 【 0 0 0 9 】

[009] この解決策は主として以下の問題を有する。オンライン話者適応化技術では、  $i$  ベクトル抽出プロセスは、複雑であり、一定時間長の音声データを必要とするので、  $i$  ベクトルを抽出するための音声データと、音声認識のための音声データとは、互いに異なる。音声認識では、  $i$  ベクトルを抽出するための音声データは、認識されるべきそれらの音声データの予備的 ( preliminary ) 音声データである。このため、  $i$  ベクトルは、認識される必要がある音声データと整合しなく、したがって音声認識の性能に影響を与える。

## 【 発明の概要 】

## 【 課題を解決するための手段 】

## 【 0 0 1 0 】

## 概要

[010] 本開示の実施形態は、余りに大きな計算複雑性を導入することなくオンライン話者適応化における音声認識の性能を効果的に改善し得る音声認識方法及び装置を提供する。

## 【 0 0 1 1 】

[011] これらの実施形態は音声認識方法を含む。本方法は、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出することを含み得る。本方法はまた、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することを含み得る。本方法はさらに、第2のニューラルネットワークに基づく音響モデルを介し音声データ内の音声を認識することを含み得る。話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルに第2のニューラルネットワークのバイアス項となるべき重み行列を掛けることを含み得る。

【0012】

[012] 第1のニューラルネットワーク、第2のニューラルネットワーク、及び重み行列は、第1のニューラルネットワーク及び第2のニューラルネットワークをそれぞれトレーニングし、次にトレーニングされた第1のニューラルネットワーク、重み行列、及びトレーニングされた第2のニューラルネットワークを一括してトレーニングすることによりトレーニングされ得る。

10

【0013】

[013] 加えて、本方法は、第1のニューラルネットワーク、第2のニューラルネットワーク、及び重み行列を初期化することを含み得る。本方法はまた、所定客観的判定基準に従って逆伝搬 (back propagation) アルゴリズムを使用することにより重み行列を更新することを含み得る。本方法はさらに、所定客観的判定基準に従って誤差逆伝搬アルゴリズムを使用することにより第2のニューラルネットワーク及び接続行列を更新することを含み得る。話者認識特徴は少なくとも話者声紋情報を含み得る。

20

【0014】

[014] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することを含み得る。話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層の出力ベクトルであり得る。

【0015】

[015] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信することを含み得る。第1のニューラルネットワークは再帰型 (recursive) ニューラルネットワークであり得る。音声データは、収集された元音声データ又は収集された元音声データから抽出される音声特徴であり得る。話者認識特徴は様々なユーザに対応してもよいし、様々なユーザのクラスタに対応してもよい。

30

【0016】

[016] これらの実施形態はまた、音声認識方法を含む。本方法は音声データを収集することを含み得る。本方法はまた、収集された音声データを第1のニューラルネットワークに入力することにより、話者認識特徴を含むベクトルを抽出することを含み得る。本方法はさらに、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することを含み得る。加えて、本方法は収集された音声データを第2のニューラルネットワークに入力することにより音声を認識することを含み得る。

40

【0017】

[017] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルに第2のニューラルネットワークのバイアス項となるべき重み行列を掛けることを含み得る。話者認識特徴は少なくとも話者声紋情報を含み得る。第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

【0018】

50

[018] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信することを含み得る。  
【0019】

[019] さらに、これらの実施形態は音声認識装置を含む。音声認識装置は、音声認識のプログラムを格納するように構成されたメモリを含み得る。音声認識装置はまた、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出するために音声認識のプログラムを実行するように構成されたプロセッサを含み得る。プロセッサはまた、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成され得る。プロセッサはさらに、音声データ内の音声データを第2のニューラルネットワークに基づく音響モデルを介し認識するように構成され得る。

10

【0020】

[020] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、話者認識特徴を含むベクトルに、第2のニューラルネットワークのバイアス項となるべき重み行列を掛けるように構成され得ることを含み得る。話者認識特徴は少なくとも話者声紋情報を含み得る。第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

【0021】

[021] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信するように構成されることを含み得る。

20

【0022】

[022] さらに、これらの実施形態は音声認識装置を含む。音声認識装置は、音声認識のプログラムを格納するように構成されたメモリを含み得る。音声認識装置はまた、音声データを収集するために音声認識のプログラムを実行するように構成されたプロセッサを含み得る。プロセッサはまた、収集された音声データを第1のニューラルネットワークに入力することにより、話者認識特徴を含むベクトルを抽出するように構成され得る。プロセッサはさらに、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成され得る。加えて、プロセッサは、収集された音声データを第2のニューラルネットワークに入力することにより音声を認識するように構成され得る。

30

【0023】

[023] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、話者認識特徴を含むベクトルに、第2のニューラルネットワークのバイアス項となるべき重み行列を掛けるように構成されることを含む。話者認識特徴は少なくとも話者声紋情報を含み得る。第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

【0024】

[024] 話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信するように構成されることを含み得る。

40

【0025】

[025] これらの実施形態はまた、音声認識装置を含む。音声認識装置は、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出するように構成された抽出ユニットを含み得る。音声認識装置はまた、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償し、第2のニューラルネットワー

50

クに基づく音響モデルを介し音声データ内の音声を認識するように構成された認識ユニットを含み得る。

【0026】

[026] これらの実施形態はさらに音声認識装置を含む。音声認識装置は音声データを収集するように構成された収集ユニットを含み得る。音声認識装置はまた、収集された音声データを第1のニューラルネットワークに入力することにより話者認識特徴を含むベクトルを抽出し、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成された抽出及び補償ユニットを含み得る。音声認識装置はさらに、収集された音声データを第2のニューラルネットワークに入力することにより音声を認識するように構成された認識ユニットを含み得る。

10

【0027】

図面の簡単な説明

[027] 本明細書の一部を構成する添付図面は、いくつかの実施形態を示し、開示された原理について本明細書と共に説明する役目を果たす。

【図面の簡単な説明】

【0028】

【図1】 [028] 例示的 i ベクトルベースオンライン話者適応化解決策の概要図である。

【図2】 [029] 本開示のいくつかの実施形態による例示的音声認識方法のフローチャートである。

【図3】 [030] 本開示のいくつかの実施形態による音声認識のための例示的システムアーキテクチャの概要図である。

20

【図4】 [031] 本開示のいくつかの実施形態による例示的ニューラルネットワークの概要図を示す。

【図5】 [032] 本開示のいくつかの実施形態による例示的システムアーキテクチャの概要図である。

【図6】 [033] 本開示のいくつかの実施形態による例示的音声認識装置の概要図である。

【図7】 [034] 本開示のいくつかの実施形態による例示的音声認識方法のフローチャートである。

【図8】 [035] 本開示のいくつかの実施形態による音声認識方法の例示的实施形態プロセスの概要図である。

30

【図9】 [036] 本開示のいくつかの実施形態による例示的音声認識装置の概要図である。

【発明を実施するための形態】

【0029】

詳細な説明

[037] 多くの詳細が、本開示の包括的理解を容易にするために以下の明細書に示される。本開示における方法及び装置は本明細書で説明されるものとは異なる多くの他のやり方で実現され得る。当業者は、本開示の暗示するものから逸脱することなく同様な拡張をなし得る。したがって、本開示は以下に開示される特定実施形態に限定されない。

【0030】

[038] 本出願の技術的解決策は添付図面及び実施形態を参照して詳細に説明される。本出願の保護範囲内に入るすべての本出願の実施形態及び実施形態における様々な特徴は相反しない限り互いに組み合わせられ得るということに注意すべきである。加えて、論理的順番がフローチャート内に示されるが、いくつかのケースでは、示される又は説明される工程は本明細書のものとは異なる順番で行われ得る。

40

【0031】

[039] いくつかの実施形態では、音声認識方法を実行するコンピュータデバイスは1つ又は複数のプロセッサ(CPU)、入出力インターフェース、ネットワークインターフェース、及びメモリを含み得る。

【0032】

[040] メモリは、非恒久的メモリ、ランダムアクセスメモリ(RAM: random access

50

memory)、及び/又はコンピュータ可読媒体内の読み取り専用メモリ(ROM: read-only memory)又はフラッシュメモリ(フラッシュRAM)などの非揮発性メモリを含み得る。メモリはコンピュータ可読媒体の一例である。メモリはモジュール1、モジュール2、...、モジュールNを含み得る、ここでNは2より大きい整数である。

【0033】

[041] コンピュータ可読媒体は、恒久的及び非恒久的ストレージ媒体、着脱可能及び着脱不能ストレージ媒体を含む。ストレージ媒体は、任意の方法又は技術により情報格納を実現し得る。情報はコンピュータ可読命令、データ構造、プログラムモジュール、又は他のデータであり得る。コンピュータ記憶媒体の例は、限定しないが、相転移メモリ(PRAM: phase change memory)、スタティックランダムアクセスメモリ(SRAM: static random access memory)、ダイナミックランダムアクセスメモリ(DRAM: dynamic random access memory)、他のタイプのランダムアクセスメモリ(RAM)、読み取り専用メモリ(ROM)、電氣的消去可能PROM(EEPROM: electrically erasable programmable read-only memory)、フラッシュメモリ又は他のメモリ技術、コンパクトディスク読み出し専用メモリ(CD-ROM: compact disk read-only memory)、デジタルバーサタイルディスク(DVD: digital versatile disc)又は他の光学的ストレージ、磁気カセットテープ、磁気ディスクストレージ又は他の磁気ストレージデバイス、又はコンピュータデバイスによりアクセスされ得る情報を格納するために使用し得る任意の他の非伝送媒体を含む。本明細書で規定されるように、コンピュータ可読媒体は変調データ信号及び搬送波などの過渡的媒体を含まない。

【0034】

[042] 本開示の実施形態は多くの利点を提供する。これらの実施形態うちのいくつかは、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出し、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することを含む。この結果、音声認識のためのニューラルネットワークをSI音響システムからSD音響システムへ変換し、これにより認識性能を改善する。話者認識特徴を含むベクトルが抽出される音声データは、音声認識のための音声データと同じであるので、認識性能は著しく改善され得る。さらに、話者認識特徴を含むベクトルが話者認識のためにニューラルネットワークを介し抽出される場合、抽出はニューラルネットワークの順方向プロセスを介し実現し得る。

【0035】

[043] 図2は本開示のいくつかの実施形態による例示的音声認識方法のフローチャートである。音声認識方法は工程S110、S120を含み得る。

【0036】

[042] 工程S110では、音声データからの話者認識特徴を含むベクトルが第1のニューラルネットワークを介し抽出される。いくつかの実施形態では、工程S110後、本方法はさらに、話者認識特徴を含む抽出されたベクトルの長さを正規化することを含み得る。いくつかの実施形態では、話者認識特徴を含む抽出されたベクトルは、長さ正規化無しに直接使用され得る。

【0037】

[045] 工程S120では、バイアスは、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内で補償され、音声データ内の音声は、第2のニューラルネットワークに基づく音響モデルを介し認識される。

【0038】

[046] 用語「第1」及び「第2」は、例えば様々なニューラルネットワークを識別するためだけに使用される。例示的実施形態の範囲から逸脱することなく、第1のニューラルネットワークは第2のニューラルネットワークと呼ばれ得る。同様に、第2のニューラルネットワークは第1のニューラルネットワークと呼ばれ得る。

【0039】

[047] 第1のニューラルネットワークは話者を類別するためのニューラルネットワークで

あり得る。第1のニューラルネットワークは、限定しないが話者声紋情報などの入力音声データに従って話者認識特徴を抽出し得る。第2のニューラルネットワークは音声認識のためのニューラルネットワークであり得る。第2のニューラルネットワークは入力音声データに従ってテキスト情報を認識し得る。

【0040】

[048] 本方法は、第1のニューラルネットワーク及び第2のニューラルネットワークを含むシステムへ適用され得る。音声データは、認識のために第1のニューラルネットワーク及び第2のニューラルネットワーク内に入力され得る。本方法は、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出することを含み得る。本方法はまた、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することを含み得る。本方法はまた、音声データからテキスト情報を取得するために、第2のニューラルネットワークに基づく音響モデルを介し音声認識することを含み得る。

10

【0041】

[049] 話者認識特徴は、話者の個人差を効果的に特徴付け得る特徴を指す。話者の個人差は声道差により引き起こされ得る。話者の個人差はまた、環境又はチャンネルにより引き起こされ得る。話者適応化技術は、声道補償と、雑音環境及びオフィス環境などの様々な話者環境の補償とのために使用され得る補償技術である。話者適応化技術はまた、電話チャンネル及びマイクロホンチャンネルなどの様々なチャンネルのための補償に使用され得る。異なる環境内では異なるチャンネルを介して同じユーザから収集された音声データは、異なる話者認識特徴のために、異なる話者の音声データと見なされ得る。

20

【0042】

[050] 本明細書の実施形態では、ニューラルネットワークは、互いに接続された複数のニューロンノードを含み得る。1つのニューロンノードの出力は別のニューロンノードの入力であり得る。ニューラルネットワークは複数のニューラル層を含み得る。それらの機能及び性質に従って、ニューラルネットワーク内のニューラル層は入力層、隠れ層、及び出力層に分割され得る。隠れ層はユーザに見えない層を指す。入力層は入力を受信し隠れ層へ分配する責任を負う。1つ又は複数の隠れ層が存在し得る。最後の隠れ層の出力結果が出力層へ提供される。ユーザは出力層から出力される最終結果を見ることができる。ニューラルネットワーク及びバイアス補償については以下に詳細に説明される。

30

【0043】

[051] 本方法では、第2のニューラルネットワークに基づく音響モデルは、話者非依存であり得、バイアス補償のための話者認識特徴を含むベクトルを導入することによりSD音響モデルに変換され得る。したがって、音声認識の性能が改善され得る。

【0044】

[052] ニューラルネットワークは、多次元の音響特性を並列に計算し得る一種のグローバルモデルである。別の態様では、 $i$ ベクトルを抽出するために使用されるガウスモデルがローカルモデルであり、次元のそれぞれを別々に計算する必要がある。したがって、第1のニューラルネットワークが本方法における話者認識特徴を含むベクトルを抽出するために使用される場合、抽出がより良好なリアルタイム性能を実現するとともに実製品において実現可能となるように短い音声データが使用され得る。

40

【0045】

[053] 加えて、ガウス混合モデルを使用するシステムはニューラルネットワークベースシステムとは異なるので、これらの2つのシステムの同時最適化は容易ではないかもしれない。それにもかかわらず、本出願の実施形態では、第1のニューラルネットワーク及び第2のニューラルネットワークが全体として最適化され得る。さらに、話者認識特徴を含むベクトルが第1のニューラルネットワークを介し抽出されるプロセスは単純であり、少ない演算量を含む。同プロセスはオンライン話者適応化認識のリアルタイム要件を満たし得る。その上、短時間データが抽出のために使用され得る。

【0046】

50

[054] 短期的データが抽出のために使用され得る場合、話者認識特徴を含むベクトルが抽出される音声データは、認識されるべき音声データであり得る。換言すれば、話者認識特徴を含む抽出されたベクトルは、認識されるべき音声データに良く整合し得る。したがって、音声認識の性能は著しく改善され得る。

【0047】

[055] 図3は、本開示のいくつかの実施形態による音声認識のための例示的システムアーキテクチャの概要図である。本システムは、音声収集デバイス11、音声認識デバイス12、及び出力デバイス13を含む。

【0048】

[056] 音声認識デバイス12は、上記工程S110を実行するように構成された話者認識ユニット121と、上記工程S120を実行するように構成された音声認識ユニット122とを含む。換言すれば、話者認識ユニット121は、話者認識特徴を含む抽出されたベクトルを音声認識ユニット122へ送信するように構成され得る。代替的に、音声認識ユニット122は、話者認識特徴を含むベクトルを話者認識ユニット121から取得するように構成され得る。

10

【0049】

[057] 音声収集デバイス11は、元音声データを収集し、元音声データ又は元音声データから抽出された音声特徴を、話者認識ユニット121及び音声認識ユニット122それぞれへ出力するように構成され得る。

【0050】

20

[058] 出力デバイス13は、音声認識ユニット122の認識結果を出力するように構成される。出力デバイス13の出力方式は、限定しないが、認識結果をデータベース内に格納すること、認識結果を所定デバイスへ送信すること、又は認識結果を所定デバイス上に表示することのうちの1つ又は複数を含み得る。

【0051】

[059] いくつかの実施形態では、音声収集デバイス11及び音声認識デバイス12は1つのデバイスに一体化され得る。代替的に、音声収集デバイス11は、元音声データ又は抽出された音声特徴を接続線、無線接続などを介し音声認識デバイス12へ送信し得る。いくつかの実施形態では、音声認識デバイス12がネットワーク側に配置される場合、音声収集デバイス11は、元音声データ又は抽出された音声特徴をインターネットを介し音声認識デバイス12へ送信し得る。

30

【0052】

[060] 出力デバイス13及び音声認識デバイス12は1つのデバイスに一体化され得る。代替的に、出力デバイス13は、認識結果を接続線、無線接続などを介し音声認識デバイス12から受信又は取得するように構成され得る。いくつかの実施形態では、音声認識デバイス12がネットワーク側に配置される場合、出力デバイス13は、認識結果をインターネットを介し音声認識デバイス12から受信又は取得するように構成され得る。

【0053】

[061] 音声認識デバイス12はさらに、話者認識特徴を含むベクトルに重み行列を掛けるための計算ユニットを含み得る。話者認識特徴を含むベクトルは話者認識ユニット121により抽出される。音声認識デバイス12は乗算の積を音声認識ユニット122へ提供するように構成され得る。代替的に、話者認識ユニット121又は音声認識ユニット122は、話者認識特徴を含むベクトルに重み行列を掛けるように構成され得る。

40

【0054】

[062] 音声認識デバイス12は、単独のデバイスでなくてもよい。例えば、話者認識ユニット121及び音声認識ユニット122は2つのデバイスに分散され得る。話者認識ユニット121又は音声認識ユニット122はまた1つ又は複数の分散デバイスにより実現され得る。

【0055】

[063] 図4は、本開示のいくつかの実施形態による例示的ニューラルネットワークの概要

50

図を示す。図 4 に示すように、ニューラルネットワークは、入力層 L 1、隠れ層 L 2、及び出力層 L 3 を含む。入力層 L 1 は 3 つのニューロンノード X 1、X 2 及び X 3 を含む。隠れ層 L 2 は 3 つのニューロンノード Y 1、Y 2 及び Y 3 を含む。出力層 L 3 は 1 つのニューロンノード Z を含む。図 4 に示すニューラルネットワークは、ニューラルネットワークの原理を単に示すために使用されており、上述の第 1 のニューラルネットワーク及び第 2 のニューラルネットワークを規定するようには意図されていない。

【 0 0 5 6 】

[064] 図 4 において、バイアスノード B 1 は、隠れ層 L 2 に対応しており、隠れ層 L 2 におけるバイアス補償のためのバイアス項を格納するために使用される。バイアスノード B 1 におけるバイアス項及び入力層 L 1 内の各ニューロンノードの出力は、隠れ層 L 2 内の各ニューロンノードの入力を提供する。バイアスノード B 2 は、出力層 L 3 に対応しており、出力層 L 3 におけるバイアス補償のためのバイアス項を格納するために使用される。バイアスノード B 2 におけるバイアス項と隠れ層 L 2 における各ニューロンノードの出力は、出力層 L 3 内の各ニューロンノードの入力を提供する。バイアス項は、事前設定され得るか、又は外部デバイスからニューラルネットワーク内へ入力され得るかのいずれかである。

10

【 0 0 5 7 】

[065] バイアス項は、バイアス補償のために使用されるベクトルを指す。ある層におけるバイアス補償は、同層のニューロンノード毎に、同層に対応するバイアスノードにより提供されるバイアス項内の、ニューロンノードに対応する値を加えた、前の層のすべてのニューロンノードの出力値の加重和の結果に基づく計算を指す。

20

【 0 0 5 8 】

[066] 例えば、入力層 L 1 内のニューロンノード X 1、X 2 及び X 3 の出力値がそれぞれ  $x_1$ 、 $x_2$  及び  $x_3$  であると仮定する。隠れ層 L 2 内のニューロンノード Y 1 に関しては、出力値は次のようになる：

【数 1】

$$f(W_{11}^{L1} x_1 + W_{12}^{L1} x_2 + W_{13}^{L1} x_3 + b_1^{B1})$$

ここで、 $f$  は、括弧内の内容に対しニューロンノード Y 1 によりなされた計算を表し、括弧内の内容は、ニューロンノード Y 1 により受信された入力値を表す。

30

【数 2】

$$W_{ij}^{L1}$$

は、例えば Y 1、 $i = 1, 2, 3$  に関しては、層 L 1 内の第  $j$  ニューロンノードと次層（すなわち層 L 2）内の第  $i$  ニューロンノード間の重み付けを指し、

【数 3】

$$b_s^{B1}$$

40

は、隠れ層 L 2 内の  $s$  番目ニューロンノード（ $s = 1, 2, 3$ ）に対応するバイアスノード B 1 内のバイアス項の値を指し、例えば、ニューロンノード Y 1 に対応するバイアスノード B 1 内のバイアス項の値は

【数 4】

$$b_s^{B1}$$

である。

【 0 0 5 9 】

50

[067] 図2に戻って参照すると、工程S120において、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルを線型変換し、これを第2のニューラルネットワークにおける入力層以外の層又はいくつかの層のバイアス項として採用することを指し得る。線形変換は、限定しないが重み行列による乗算のやり方で行われ得る。

【0060】

[068] 第1のニューラルネットワークは3つの隠れ層を含み得る。いくつかの実施形態では、第1のニューラルネットワークは1つ又は2つの隠れ層を含んでもよいし、4以上の隠れ層を含んでもよい。いくつかの実施形態では、第2のニューラルネットワークは3つの隠れ層を含み得る。いくつかの実施形態では、第2のニューラルネットワークは1つ又は2つの隠れ層を含んでもよいし、4以上の隠れ層を含んでもよい。

10

【0061】

[069] いくつかの実施形態では、話者認識特徴は少なくとも話者声紋情報を含み得る。話者声紋情報は異なるユーザの音声データを識別するために使用され得る。換言すれば、異なるユーザの音声データから抽出される話者声紋情報は異なる。いくつかの実施形態では、話者認識特徴は、話者声紋情報、環境情報、及びチャンネル情報のうちの1つ又は複数を含み得る。環境情報は、音声データが収集される環境の特徴を特徴付けるために使用され得る。チャンネル情報は、音声データが収集されるチャンネルの特徴を特徴付けるために使用され得る。

【0062】

20

[070] いくつかの実施形態では、第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。再帰型ニューラルネットワークは、1つ又は複数のフィードバックループを有するニューラルネットワークを指し、非線形システムのリアルな動的モデリングを実現し得る。再帰型ニューラルネットワークが、話者認識特徴を含むベクトルを抽出するために使用される場合、抽出は短期データに対し行われ得る。再帰型ニューラルネットワークは、限定しないが、LSTM (long-short term memory) 再帰型ニューラルネットワークであり得る。

【0063】

[071] いくつかの実施形態では、工程S120における話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルに、第2のニューラルネットワークのバイアス項となるべき重み行列を掛けることを含み得る。

30

【0064】

[072] いくつかの実施形態では、重み行列が単位行列である場合、重み行列が掛けられた後、話者認識特徴を含むベクトルは変化しなくてもよい。話者認識特徴を含むベクトルは、第2のニューラルネットワークのバイアス項として直接採用され得る。

【0065】

[073] いくつかの実施形態では、第1のニューラルネットワーク、第2のニューラルネットワーク、及び重み行列は、第1のニューラルネットワーク及び第2のニューラルネットワークをそれぞれトレーニングし、次にトレーニングされた第1のニューラルネットワーク、重み行列、及びトレーニングされた第2のニューラルネットワークを一括してトレーニングすることによりトレーニングされ得る。一括的にトレーニングすることは、トレーニングするための音声データを、第1のニューラルネットワーク及び第2のニューラルネットワーク内にそれぞれ入力することと、第1のニューラルネットワークにより抽出された話者認識特徴を含むベクトルに重み行列を掛けた後に、第2のニューラルネットワーク上のバイアスを補償することとを指し得る。トレーニングは、限定しないがグラフィック処理ユニット (GPU: graphics processing unit) により行われ得る。

40

【0066】

[074] いくつかの実施形態では、トレーニングされた第1のニューラルネットワーク、重み行列、及びトレーニングされた第2のニューラルネットワークを一括してトレーニング

50

した後、本方法はさらに、第1のニューラルネットワーク、第2のニューラルネットワーク、及び重み行列を初期化することを含み得る。本方法はまた、所定客観的判定基準に従って逆伝搬アルゴリズムを使用することにより重み行列を更新することを含み得る。加えて、本方法は、所定客観的判定基準に従って誤差逆伝搬アルゴリズムを使用することにより第2のニューラルネットワーク及び接続行列を更新することを含み得る。重み行列に関する初期化は、ガウス分布による乱数的初期化であり得る。上記所定客観的判定基準は、限定しないが、標的最小二乗平均誤差(LMS)、再帰型最小二乗(RLS: recursive least square)、及び正規化最小二乗平均誤差(NLMS: normalized least mean square error)を含み得る。

【0067】

[075] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することを含む。話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層の出力ベクトルであり得る。例えば、第2のニューラルネットワークが入力層、3つの隠れ層、及び1つの出力層を含むと仮定すると、入力層を除くすべての層は、出力層及び3つの隠れ層を指し得る。入力層を除くいくつかの層は、4つの層(すなわち、出力層及び3つの隠れ層)のうちの1つ又は複数を指し得る。

【0068】

[076] 話者認識特徴を含むベクトルに基づく第2のニューラルネットワーク内のある層上のバイアス補償は、話者認識特徴を含むベクトルに重み行列を掛けることにより取得されるベクトルを、同層のバイアス項として採用することを指し得る。例えば、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層以外のすべての層上のバイアス補償は、話者認識特徴を含むベクトルに重み行列を掛けることにより取得されたベクトルを、第2のニューラルネットワーク内の出力層及び3つの隠れ層のそれぞれのバイアス項として採用することを指し得る。

【0069】

[077] いくつかの実施形態では、話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層の出力ベクトルであり得る。最後の隠れ層の出力ベクトルは出力層の出力ベクトルより少ない次元を有し、これにより過剰フィッティング(overfitting)を回避する。

【0070】

[078] いくつかの実施形態では、話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層以外の隠れ層の出力ベクトルであってもよいし、出力層の出力ベクトルであってもよい。

【0071】

[079] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償することは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信することを含み得る。ある層に対応するバイアスノードは、同層の上のバイアス補償のために使用されるバイアス項を格納し得る。話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層内の複数のニューロンノードのそれぞれの出力値からなるベクトルであり得る。

【0072】

[080] 話者認識特徴を含むベクトルをバイアスノードへ送信することは、話者認識特徴を含むベクトルをバイアスノードへ直接送信することを指してもよいし、話者認識特徴を含むベクトルを線型変換し、次にこれをバイアスノードへ送信することを指してもよい。

【0073】

[081] 第2のニューラルネットワーク内の複数層が同じベクトルによりバイアス補償を受

10

20

30

40

50

ける場合、複数層は個別のバイアスノードに対応してもよいし、同じバイアスノードに対応してもよい。例えば、第1のニューラルネットワークにより抽出された話者認識特徴を含むベクトルは、複数のバイアスノードそれぞれへ送信され得る。複数のバイアスノードは、第2のニューラルネットワーク内のバイアス補償を必要とする複数の層に1対1ベースで対応する。別の例として、第1のニューラルネットワークにより抽出された話者認識特徴を含むベクトルはまた、1つのバイアスノードへ送信され得る。当該バイアスノードは、第2のニューラルネットワーク内のバイアス補償を必要とする複数の層に対応する。

【0074】

[082] いくつかの実施形態では、音声データは、収集された元音声データであるか又は収集された元音声データから抽出された音声特徴である。音声特徴は、限定しないが、メル周波数ケプストラム係数(MFCC: Mel frequency cepstral coefficient)、知覚線形予測係数(PLP: perceived linear prediction coefficient)、フィルタバンク特徴、又はそれらの任意の組み合わせを含み得る。

10

【0075】

[083] いくつかの実施形態では、話者認識特徴は、様々なユーザに1対1ベースで対応してもよいし、様々なユーザのクラスタに1対1ベースで対応してもよい。様々なユーザに1対1ベースで対応する話者認識特徴は、第1のニューラルネットワークの出力層がユーザの識別子を出力することを意味する。様々なユーザのクラスタに1対1ベースで対応する話者認識特徴は、ユーザがクラスタ化された後、第1のニューラルネットワークの出力層がカテゴリ識別子を出力することを意味する。

20

【0076】

[084] クラスタは1つ又は複数のパターンを含み得、パターンは測定のベクトルを指してもよいし、多次元空間内の点であってもよい。クラスタ化操作は類似性に基づいており、同じクラスタ内のパターンは、異なるクラスタ内のパターンより高い類似性を有する。クラスタ化のためのアルゴリズムは分割方法、階層的な方法、密度アルゴリズム、グラフ理論クラスタリング法、グリッドアルゴリズム、及びモデルアルゴリズムに分割され得る。例えば、これらのアルゴリズムはK平均法(K-MEANS)、K-MEDOIDS、Clara又はClaransであり得る。

【0077】

[085] ユーザをクラスタ化することは、トレーニング中に様々なユーザの話者認識特徴間の類似性に従って複数のユーザの話者認識特徴を複数のクラスタに類別することと、クラスタに対応する話者認識特徴を含むベクトルを取得するために、1つのクラスタに分類された複数の話者認識特徴を計算(例えば重み付け平均化)することとを指し得る。カテゴリ識別子は、1つのクラスタを表すために使用される識別子であり得る。カテゴリ識別子はクラスタに1対1ベースで対応する。

30

【0078】

[086] 非常に多くのユーザの音声認識が必要とされる場合、クラスタ化操作が行われれば、一組の出力結果は、話者認識特徴を含む限定数のベクトルであり得る。例えば、何百万のユーザが存在する場合、ユーザが何千ものクラスタに分類されれば、話者認識特徴を含む何千ものベクトルだけが存在し、これにより実装の複雑性を著しく低減する。

40

【0079】

[087] 話者認識特徴が、話者認識特徴間の類似性に従って複数のクラスタに分類される場合、様々な次元の類似性(例えば、声紋情報、環境情報、チャンネル情報などの様々なタイプの話者認識特徴)に従って、様々なクラスタ化結果が取得され得る。例えば、同様な声紋を有する話者認識特徴は1つのクラスタと見なされ得る。別の例として、同じ又は同様な環境に対応する話者認識特徴が1つのクラスタと見なされ得る。代替的に、同様なチャンネルに対応する話者認識特徴が1つのクラスタと見なされ得る。

【0080】

[088] 図5は、本開示のいくつかの実施形態による音声認識のための例示的システムアーキテクチャの概要図である。図5に示すように、本システムは話者分類器21と、音声認

50

識システム 2 3 とを含み得る。本システムにおける話者認識特徴は話者声紋情報である。話者分類器 2 1 は上記工程 S 1 1 0 を実行するように構成される。音声認識システム 2 3 は上記工程 S 1 2 0 を実行するように構成される。

【 0 0 8 1 】

[089] 話者声紋情報を含むベクトルが接続行列 2 2 により線形に変換され得る。接続行列は限定しないが重み行列であり得る。

【 0 0 8 2 】

[090] 話者認識特徴を含むベクトルを抽出するための第 1 のニューラルネットワークを利用する話者分類器 2 1 は、入力層 2 1 1、1 つ又は複数の隠れ層 2 1 2、及び出力層 2 1 3 を含み得る。いくつかの実施形態では、隠れ層 2 1 2 の数は 3 であり得る。代替的に、1 つ又は複数の隠れ層 2 1 2 が存在し得る。

10

【 0 0 8 3 】

[091] 音声を認識するための第 2 のニューラルネットワークを利用する音声認識システム 2 3 は、入力層 2 3 1、1 つ又は複数の隠れ層 2 3 2、及び出力層 2 3 3 を含み得る。いくつかの実施形態では、隠れ層 2 1 2 の数は 3 であり得る。いくつかの実施形態では、1 つ又は複数の隠れ層 2 1 2 が存在し得る。

【 0 0 8 4 】

[092] 話者分類器 2 1 内の第 1 のニューラルネットワークの入力層 2 1 1 により受信される音声データは、音声認識システム 2 3 内の第 2 のニューラルネットワークの入力層 2 3 1 により受信されるものと同じであり得る。音声データは、収集された元音声データであり得る。代替的に、音声データは、元音声データから抽出された音声特徴であり得る。

20

【 0 0 8 5 】

[093] したがって、話者分類器 2 1 内の第 1 のニューラルネットワークは、音声認識システム 2 3 内の第 2 のニューラルネットワークと同じ入力を有し得る。すなわち、話者声紋情報を含むベクトルが取得される音声データは、音声認識のための音声データと同じであり得る。したがって、話者声紋情報を含むベクトルによる第 2 のニューラルネットワーク上のバイアス補償は、認識されるべき音声データと完全に整合し得る。その結果、音声認識の性能は効果的に改善され得る。第 1 のニューラルネットワーク及び第 2 のニューラルネットワークはそれぞれ、全結合ニューラルネットワーク ( D N N : fully connected neural network )、畳み込みニューラルネットワーク ( C N N : convolution neural network )、及び再帰型ニューラルネットワーク ( R N N : recurrent neural network ) の任意の 1 つ、又はそのいくつかの組み合わせを含み得る。

30

【 0 0 8 6 】

[094] 話者声紋情報を含むベクトル表現は、話者分類器 2 1 内の最後の隠れ層の出力ベクトルであり得る。

【 0 0 8 7 】

[095] 音声認識システム 2 3 において、出力層 2 3 3 及び 1 つ又は複数の隠れ層 2 3 2 のそれぞれは、話者声紋情報を含む線形変換されたベクトル表現をバイアス項として採用し得る。いくつかの実施形態では、出力層 2 3 3 及び 1 つ又は複数の隠れ層 2 3 2 において、少なくとも 1 つ又は複数の層は、話者声紋情報を含む線型変換されたベクトル表現をバイアス項として採用し得る。

40

【 0 0 8 8 】

[096] 接続行列 2 2 はまた、話者声紋情報を含むベクトルに対し長さ正規化を行うように構成され得る。いくつかの実施形態では、話者分類器により出力された話者声紋情報を含むベクトルは、長さ正規化を受けることなく、重み付けを掛けられた後、音声認識システム 2 3 へ直接提供され得る。

【 0 0 8 9 】

[097] 話者分類器 2 1 の出力層 2 1 3 によるデータ出力は、様々なユーザのタグ ID であってもよいし、ユーザがクラスタ化された後のクラスタのタグ ID であってもよい。出力層の出力データはトレーニングのためにだけ使用され得る。音声認識システム 2 3 の出力

50

層 2 3 3 から出力される認識結果は、状態レベル、音素レベル、又は単語レベルタグ I D であり得る。

【 0 0 9 0 】

[098] 図 5 に示す例示的システムアーキテクチャはさらに、以下の機能を実行し得る。

【 0 0 9 1 】

[099] トレーニングデータを使用することにより、話者分類器の第 1 のニューラルネットワークと、第 2 のニューラルネットワークに基づく音響モデル（例えば図 2 において参照される音響モデル）とをトレーニングすること。第 1 及び第 2 のニューラルネットワークは所望の音声認識性能又は話者認識性能をそれぞれ実現し得る。さらに、このトレーニングは、第 1 のニューラルネットワーク、接続行列、及び第 2 のニューラルネットワークを一括してトレーニングすることを含み得る。GPU がこれらのトレーニングを加速するために使用され得る。

10

【 0 0 9 2 】

[0100] システムアーキテクチャは、ネットワーク初期化のための音響モデル及び話者分類器として、トレーニングされた音響モデル及び話者分類器を使用し得る。いくつかの実施形態では、ネットワーク初期化はまた図 5 の接続行列をランダムに初期化することを含み得る。

【 0 0 9 3 】

[0101] 所定客観的判定基準に従って、システムアーキテクチャは、収束状態に到達するように接続行列を更新するために、逆伝播（BP）アルゴリズムを使用し得る。

20

【 0 0 9 4 】

[0102] 所定客観的判定基準に従って、システムアーキテクチャは、収束状態に到達するように音響モデル及び接続行列を更新するために BP アルゴリズムを使用し得る。所定客観的判定基準は、実アプリケーションにおけるニーズに従って設定され得る。

【 0 0 9 5 】

[0103] さらに、システムアーキテクチャは、収集された元音声データから音声特徴を抽出し得る。抽出された音声特徴は話者分類器により処理され、その結果、音声特徴に対応する話者声紋情報を含むベクトルが取得される。このベクトルは接続行列により線形変換され、音声認識システムへ送信される。抽出された音声特徴は、音声認識システムにおける第 2 内のニューラルネットワークに基づき音響モデルにより復号化される。最終的に、音声認識結果が取得され得る。音声認識システムでは、第 2 のニューラルネットワークの出力層及び 3 つの隠れ層のバイアス項は、話者声紋情報を含む形型変換されたベクトルであり得る。

30

【 0 0 9 6 】

[0104] 本出願はまた、音声認識装置に関する。音声認識装置は、音声認識のプログラムを格納するように構成されたメモリを含む。音声認識装置はまた、音声認識のプログラムを実行するように構成されたプロセッサを含む。プロセッサは、音声認識のプログラムを実行する際、第 1 のニューラルネットワークを介し、音声データから話者認識特徴を含むベクトルを抽出するように構成され得る。

【 0 0 9 7 】

40

[0105] プロセッサはまた、音声認識のプログラムを実行する際、話者認識特徴を含むベクトルに従って第 2 のニューラルネットワーク内のバイアスを補償するように構成され得る。プロセッサはさらに、音声認識のプログラムを実行する際、第 2 のニューラルネットワークに基づく音響モデルを介し音声データ内の音声を認識するように構成される。

【 0 0 9 8 】

[0106] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第 2 のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、話者認識特徴を含むベクトルに、第 2 のニューラルネットワークのバイアス項となるべき重み行列を掛けるように構成されることを含み得る。

【 0 0 9 9 】

50

[0107] いくつかの実施形態では、話者認識特徴は少なくとも話者声紋情報を含み得る。

【0100】

[0108] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償するように構成されることを含み得る。話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層の出力ベクトルであり得る。

【0101】

[0109] いくつかの実施形態では、話者認識特徴は、様々なユーザに1対1ベースで対応してもよいし、様々なユーザのクラスタに1対1ベースで対応してもよい。様々なユーザに1対1ベースで対応する話者認識特徴は、第1のニューラルネットワークの出力層がユーザの識別子を出力することを意味する。様々なユーザのクラスタに1対1ベースで対応する話者認識特徴は、ユーザがクラスタ化された後に、第1のニューラルネットワークの出力層がカテゴリ識別子を出力することを意味する。

【0102】

[0110] いくつかの実施形態では、第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

【0103】

[0111] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成されたプロセッサは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信するように構成されることを含む。

【0104】

[0112] さらに、プロセッサは、音声認識のプログラムを実行する際、上記工程S110、S120を実行するように構成され得る。音声認識のプログラムを実行する際にプロセッサにより実行される動作のさらなる詳細は上に見出され得る。

【0105】

[0113] 本出願はさらに音声認識装置に関する。図6は、本開示のいくつかの実施形態による例示的音声認識装置の概要図である。音声認識装置は、抽出ユニット31及び認識ユニット32を含む。

【0106】

[0114] 一般的に、これらのユニット（そして任意の副ユニット）は、他の部品（例えば集積回路の一部）と共に及び/又は関連機能の特定機能を実行するプログラム（コンピュータ可読媒体上に格納された）の一部と共に使用するために設計されたパッケージ化機能ハードウェアユニットであり得る。このユニットは入口点及び出口点を有し得、例えばJava（登録商標）、Lua、C、又はC++などのプログラミング言語で書かれ得る。ソフトウェアユニットは、コンパイルされ、実行可能プログラム内へリンクされ、動的リンクライブラリ内にインストールされてもよいし、例えばBASIC、Perl、又はPythonなどのインタープリット型プログラミング言語で書かれてもよい。ソフトウェアユニットは他のユニット又は自身から呼出し可能であり得る及び/又は検出された事象又は割り込みに応答して呼び出され得るということが理解される。コンピュータデバイス上で実行するように構成されたソフトウェアユニットは、コンパクトディスク、デジタルビデオディスク、フラッシュドライブ、磁気ディスク、又は任意の他の非一時的媒体などのコンピュータ可読媒体上に提供されてもよいし、デジタルダウンロードとして提供されてもよい（そして実行に先立って、インストール、圧縮解除、又は解読を必要とする圧縮された又はインストール可能なフォーマットで元々格納され得る）。このようなソフトウェアコードは、コンピュータデバイスによる実行のために実行コンピュータデバイスのメモリデバイス上に部分的又は完全に格納され得る。ソフトウェア命令はEPROMなどのファームウェアで埋め込まれ得る。ハードウェアユニットはゲート及びフリップフロップなどの接続された論

10

20

30

40

50

理ユニットで構成され得る及び/又はプログラマブルゲートアレイ又はプロセッサなどのプログラム可能ユニットで構成され得るといことがさらに理解される。本明細書で説明されたユニット又はコンピュータデバイス機能は好適にはソフトウェアユニットとして実現されるが、ハードウェア又はファームウェアで表され得る。一般的に、本明細書で説明されたユニットは、他のユニットと組み合わせられ得る又は物理的編成又はストレージにもかかわらず副ユニットに分割され得る論理ユニットを指す。

【0107】

[0115] 抽出ユニット31は、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出するように構成され得る。抽出ユニット31は、上記装置内の話者認識特徴を含むベクトルを抽出するためのものと同様な動作を実行するように構成され得る。

10

【0108】

[0116] 認識ユニット32は、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償し、第2のニューラルネットワークに基づく音響モデルを介し音声データ内の音声を認識するように構成され得る。認識ユニット32は上記装置内の音声を認識するためのものと同様な動作を実行するように構成され得る。

【0109】

[0117] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成された認識ユニット32は、話者認識特徴を含むベクトルに、第2のニューラルネットワークのバイアス項となるべき重み行列を掛けるように構成されることを含み得る。

20

【0110】

[0118] いくつかの実施形態では、話者認識特徴は少なくとも話者声紋情報を含み得る。

【0111】

[0119] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成された認識ユニット32は、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内の入力層を除く層のすべて又は一部においてバイアスを補償するように構成されることを含み得る。代替的に、話者認識特徴を含むベクトルは、第1のニューラルネットワーク内の最後の隠れ層の出力ベクトルであり得る。

30

【0112】

[0120] いくつかの実施形態では、第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

【0113】

[0121] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成された認識ユニット32は、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信するように構成されることを含み得る。

【0114】

[0122] いくつかの実施形態では、抽出ユニット31は、図3に示すシステムアーキテクチャ内の話者認識ユニット121として構成され得る。認識ユニット32は、図3に示すシステムアーキテクチャ内の音声認識ユニット122として構成され得る。図6の装置は、図3に示すシステムアーキテクチャの音声認識装置として構成され得る。図6の装置のさらなる詳細動作は音声認識図3に示す装置に関して上に説明したものが参照され得る。

40

【0115】

[0123] さらに、抽出ユニット31及び認識ユニット32により実行される動作は、上記音声認識方法における工程S110、S120と同様であり得る。抽出ユニット31及び認識ユニット32により実行される動作のさらなる詳細も上に見出され得る。

【0116】

50

[0124] 本出願はまた音声認識方法に向けられる。図7は、本開示のいくつかの実施形態による例示的音声認識方法のフローチャートである。この方法は、図3の音声認識デバイス及び/又は図5のシステムアーキテクチャにより行われ得る。図7に示すように、音声認識方法は、以下の工程S410、S420、及び工程S430を含む。

【0117】

[0125] 工程S410では、システムアーキテクチャは音声データを収集する。

【0118】

[0126] 工程S420では、システムアーキテクチャは、収集された音声データを第1のニューラルネットワークに入力することにより、話者認識特徴を含むベクトルを抽出し、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償する。

10

【0119】

[0127] 工程S430では、システムアーキテクチャは、収集された音声データを第2のニューラルネットワークに入力することにより音声を認識する。

【0120】

[0128] 工程S410、S420及び工程S430は、収集プロセス中に連続的に行われ得る。一団の音声データが収集されるたびに、工程S420、S430が、その一団の音声データの音声認識の結果を取得するために、当該一団の音声データに対し行われ得る。一団の音声データのサイズは限定しないが1つ又は複数のフレームであり得る。

【0121】

20

[0129] 図8は、本開示のいくつかの実施形態による図7の音声認識方法の例示的实施形態プロセスの概要図である。

【0122】

[0130] この実施形態プロセスはユーザの音声を収集することを含む。この実施形態プロセスはまた、収集された音声データ又はそれから抽出された音声特徴を、第1のニューラルネットワーク及び第2のニューラルネットワークへ直接入力することを含む。この実施形態プロセスはさらに、話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出し、このベクトルを第2のニューラルネットワークへバイアス項として送信することを含む。この実施形態プロセスはまた、第2のニューラルネットワークから音声データの認識結果を出力することを含む。

30

【0123】

[0131] 収集された元音声データは、第1のニューラルネットワーク及び第2のニューラルネットワークへ直接提供され得る。代替的に、音声特徴は、収集された元音声データから抽出され得、抽出された音声特徴は次に第1のニューラルネットワーク及び第2のニューラルネットワークへ提供される。

【0124】

[0132] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、話者認識特徴を含むベクトルに、第2のニューラルネットワークのバイアス項となるべきに重み行列を掛けることを含み得る。

【0125】

40

[0133] いくつかの実施形態では、話者認識特徴は少なくとも話者声紋情報を含み得る。

【0126】

[0134] いくつかの実施形態では、第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

【0127】

[0135] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償することは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信することを含む。

50

## 【0128】

[0136] 第1のニューラルネットワーク、第2のニューラルネットワーク、話者認識特徴を含むベクトルの抽出、話者認識特徴を含むベクトルに従って第2のニューラルネットワークにおけるバイアス補償、及び第2のニューラルネットワークに基づく音声認識のさらなる詳細は、音声認識方法について上で説明したものと同様である。

## 【0129】

[0137] 本出願に開示された実施形態はさらに音声認識装置に関する。音声認識装置は音声認識のプログラムを格納するように構成されたメモリを含む。音声認識装置はまた、音声データから話者認識特徴を含むベクトルを第1のニューラルネットワークを介し抽出するために、音声認識のプログラムを実行するように構成されたプロセッサを含む。プロセッサはまた、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するために、音声認識のプログラムを実行するように構成される。プロセッサはさらに、音声データ内の音声を第2のニューラルネットワークに基づく音響モデルを介し認識するために、音声認識のプログラムを実行するように構成される。

10

## 【0130】

[0138] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するために音声認識のプログラムを実行するように構成されたプロセッサは、話者認識特徴を含むベクトルに、第2のニューラルネットワークのバイアス項となるべき重み行列を掛けるように構成されることを含み得る。

## 【0131】

[0139] いくつかの実施形態では、話者認識特徴は少なくとも話者声紋情報を含み得る。

20

## 【0132】

[0140] いくつかの実施形態では、第1のニューラルネットワークは再帰型ニューラルネットワークであり得る。

## 【0133】

[0141] いくつかの実施形態では、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するために音声認識のプログラムを実行するように構成されたプロセッサは、第1のニューラルネットワークの最後の隠れ層においてニューロンノードにより出力された話者認識特徴を含むベクトルを、第2のニューラルネットワーク内の入力層を除く層のすべて又は一部に対応するバイアスノードへ送信するように構成されることを含み得る。

30

## 【0134】

[0142] プロセッサが音声認識のプログラムを読み実行するように構成される場合、音声データを収集するように構成されたプロセッサは、図3の音声収集装置の動作を参照し得る。動作のさらなる詳細は図3の説明に見出され得る。話者認識特徴を含むベクトルを抽出すること、話者認識特徴を含むベクトルに従って第2のニューラルネットワークにおいてバイアスを補償すること、及び第2のニューラルネットワークに基づき音声を認識することのさらなる動作詳細もまた、本音声認識方法における上記それらの説明を参照することができる。

## 【0135】

[0143] 本出願はさらに音声認識装置に関する。図9は、本開示のいくつかの実施形態による例示的音声認識装置の概要図である。図9に示すように、音声認識装置は、音声データを収集するように構成された収集ユニット61を含む。音声認識装置はまた、収集された音声データを第1のニューラルネットワークに入力することにより、話者認識特徴を含むベクトルを抽出し、話者認識特徴を含むベクトルに従って第2のニューラルネットワーク内のバイアスを補償するように構成された抽出及び補償ユニット62を含む。音声認識装置はさらに、収集された音声データを第2のニューラルネットワークに入力することにより音声を認識するように構成された認識ユニット63を含む。これらのユニット（そして任意の副ユニット）は、他の部品（例えば集積回路の一部）と共に及び/又は関連機能の特定機能を実行するプログラム（コンピュータ可読媒体上に格納された）の一部と共に

40

50

使用するために設計されたパッケージ化機能ハードウェアユニットであり得る。

【0136】

[0144] 収集ユニット61は、上記装置内の音声データを収集するためのものと同様な動作を実行するように構成され得る。

【0137】

[0145] 抽出及び補償ユニット62は、上記装置内の話者認識特徴を含むベクトルを抽出するためのもの、及び第2のニューラルネットワークにおいてバイアスを補償するためのものと同様な動作を実行するように構成され得る。

【0138】

[0146] 認識ユニット63は、音声を認識するためのものと同様な動作を実行するように構成され得る。

10

【0139】

[0147] 収集ユニット61は、独立したデバイス内に装備され得る。代替的に、収集ユニット61は、抽出及び補償ユニット62、認識ユニット63と共に同じデバイス内に装備され得る。

【0140】

[0148] 収集ユニット61は、図3に示す音声収集装置を参照して実現され得る。抽出及び補償ユニット62並びに認識ユニット63による、第1のニューラルネットワークによる話者認識特徴を含むベクトルを抽出すること、話者認識特徴を含むベクトルに従って第2のニューラルネットワークにおいてバイアスを補償すること、及び音声を認識することのさらなる実施詳細は、上記音声認識方法におけるそれらの説明として参照され得る。

20

【0141】

[0149] 上に示すように、上記方法のすべて又はいくつかの工程は、プログラムを介した命令下で当該ハードウェアにより完了され得るということが理解される。プログラムは、読み取り専用メモリ、磁気ディスク又はコンパクトディスクなどのコンピュータ可読記憶媒体内に格納され得る。任意選択的に、上述の実施形態のすべて又はいくつかの工程はまた、1つ又は複数の集積回路を使用して実現され得る。したがって、上記実施形態における様々なモジュール/ユニットはハードウェアの形式で実装されてもよいし、ソフトウェア機能モジュールの形式で実装されてもよい。本出願はハードウェア及びソフトウェアの組み合わせのいかなる特定形式にも限定されない。

30

【0142】

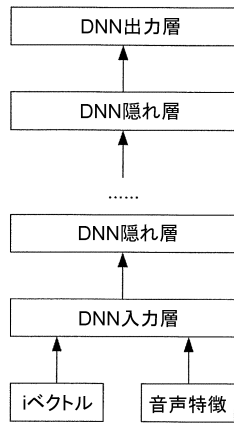
[0150] 確かに、本出願の様々な他の実施形態が存在し得る。当業者は本出願の精神及び本質から逸脱することなく本出願に従って様々な変更及び変形をなすことができるだろう。すべてのこれらの対応する変更及び変形はすべて本出願の特許請求の範囲に入るべきである。

40

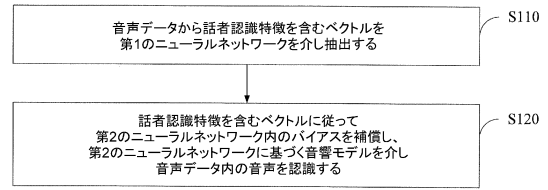
50

【図面】

【図 1】



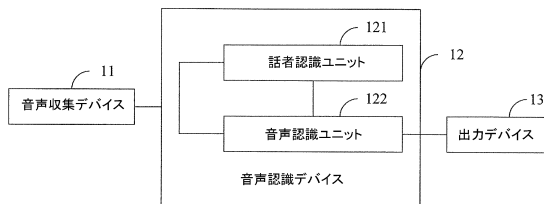
【図 2】



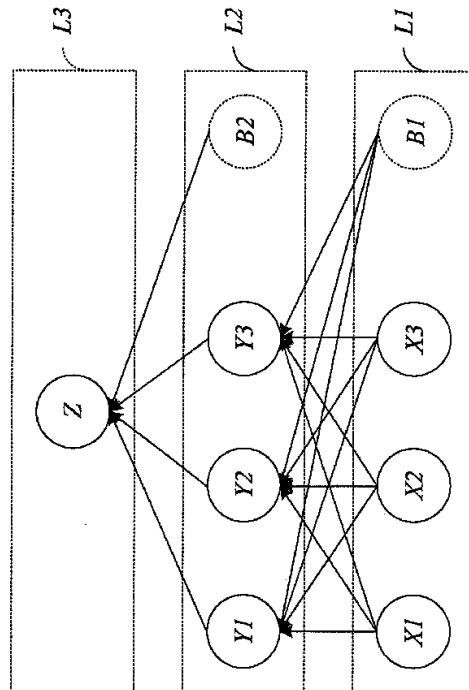
10

20

【図 3】



【図 4】



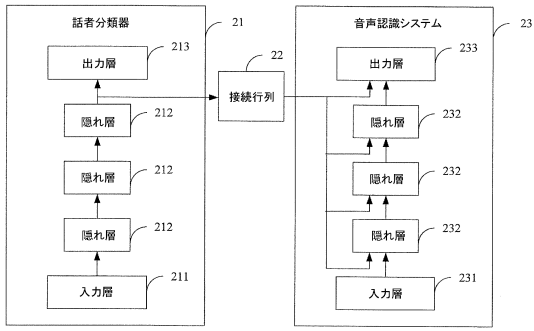
30

40

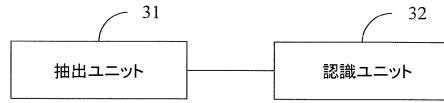
Fig. 4

50

【図5】

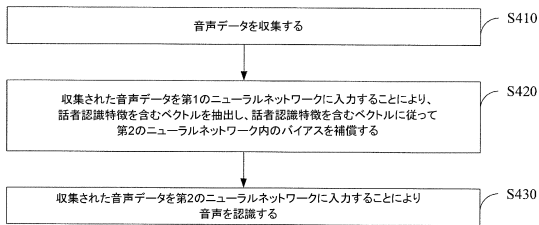


【図6】

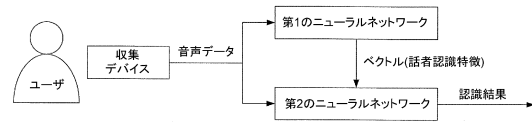


10

【図7】



【図8】



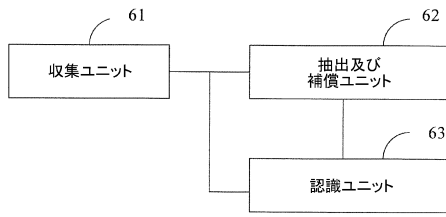
20

30

40

50

【図 9】



10

20

30

40

50

## フロントページの続き

- (72)発明者 ファン, ジーイン  
中華人民共和国, 3 1 1 1 2 1, ハンチョウ, ユ ハン ディストリクト, ウェスト ウェン イ 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 シュエ, シャオフェイ  
中華人民共和国, 3 1 1 1 2 1, ハンチョウ, ユ ハン ディストリクト, ウェスト ウェン イ 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 ヤン, ジージエ  
中華人民共和国, 3 1 1 1 2 1, ハンチョウ, ユ ハン ディストリクト, ウェスト ウェン イ 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 山下 剛史
- (56)参考文献 特開 2 0 1 5 - 1 0 2 8 0 6 ( J P , A )  
柏木陽佑他, 話者コードに基づく話者正規化学習を利用したニューラルネット音響モデルの適応, 電子情報通信学会技術研究報告, 2014年12月, Vol.114, No.365, pp.105-110  
柏木陽佑他, 制約付き話者コードの同時推定によるニューラルネット音響モデルの話者正規化学習, 日本音響学会 2 0 1 4 年秋季研究発表会講演論文集, 2014年09月, pp.7-10  
Shaofei XUE et al., Fast Adaptation of Deep Neural Network Based on Discriminant Codes for Speech Recognition, IEEE/ACM Transactions on Audio, Speech and Language Processing, 2014年12月, Vol.22, No.12, pp.1713-1725  
Zhiying HUANG et al., Speaker adaptation of RNN-BLSTM for speech recognition based on speaker code, 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2016), 2016年03月, pp.5305-5309
- (58)調査した分野 (Int.Cl., D B 名)  
G 1 0 L 1 5 / 0 0 - 1 7 / 2 6  
I E E E X p l o r e