



特許証
(CERTIFICATE OF PATENT)

特許第6954680号
(PATENT NUMBER)

発明の名称
(TITLE OF THE INVENTION)

話者の確認方法及び話者の確認装置

特許権者
(PATENTEE)

中華人民共和国北京市海淀区清華大学 郵編 1
00084
国籍・地域 中華人民共和国
清華大学

発明者
(INVENTOR)

王 李 鄭
東 藍 方 天

出願番号
(APPLICATION NUMBER)

特願2019-553913

出願日
(FILING DATE)

平成29年12月 1日 (December 1, 2017)

登録日
(REGISTRATION DATE)

令和 3年10月 4日 (October 4, 2021)

この発明は、特許するものと確定し、特許原簿に登録されたことを証する。
(THIS IS TO CERTIFY THAT THE PATENT IS REGISTERED ON THE REGISTER OF THE JAPAN PATENT OFFICE.)

特許庁長官
(COMMISSIONER, JAPAN PATENT OFFICE)

令和 3年10月 4日 (October 4, 2021)

森



(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第6954680号
(P6954680)

(45) 発行日 令和3年10月27日(2021.10.27)

(24) 登録日 令和3年10月4日(2021.10.4)

(51) Int.Cl.

F 1

G 1 O L 17/18

(2013.01)

G 1 O L 17/18

G 1 O L 15/06

(2013.01)

G 1 O L 15/06

G 0 6 N 3/04

(2006.01)

G 0 6 N 3/04

5 0 0 P

1 4 5

請求項の数 10 (全 12 頁)

(21) 出願番号	特願2019-553913 (P2019-553913)	(73) 特許権者	502192546
(86) (22) 出願日	平成29年12月1日 (2017.12.1)	清華大学	Tsinghua University
(65) 公表番号	特表2020-515905 (P2020-515905A)	中華人民共和国北京市海淀区清華大学 郵 編100084	中華人民共和国北京市海淀区清華大学 郵 編100084
(43) 公表日	令和2年5月28日 (2020.5.28)	Tsinghua University	, Haidian District, Beijing 100084, P. R. China
(86) 國際出願番号	PCT/CN2017/114293	(74) 代理人	100108833
(87) 國際公開番号	W02018/176894	弁理士	早川 裕司
(87) 國際公開日	平成30年10月4日 (2018.10.4)	(74) 代理人	100162156
審査請求日	令和1年11月29日 (2019.11.29)	弁理士	村雨 圭介
(31) 優先権主張番号	201710214666.0		
(32) 優先日	平成29年4月1日 (2017.4.1)		
(33) 優先権主張国・地域又は機関	中国(CN)		

最終頁に続く

(54) 【発明の名称】 話者の確認方法及び話者の確認装置

【特許請求の範囲】

【請求項 1】

第2の音声を取得することと、

予め取得した第1の音声と前記第2の音声を、対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換することと、

畳み込みニューラルネットワークを使用して、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得することと、

時間遅延ニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得することと、

前記第3の特徴と前記第4の特徴により、話者を確認することと、を含み、

前記の、時間遅延ニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得することは、

前記第1の音声スペクトログラムの前後のフレームに対応する前記第1の特徴に対してつなぎ合わせを行い、前記第2の音声スペクトログラムの前後のフレームに対応する前記第2の特徴に対してつなぎ合わせを行うことと、

つなぎ合わせた後の前記第1の特徴とつなぎ合わせた後の前記第2の特徴のそれぞれに対して線形変換と次元削減を行い、対応する前記第3の特徴と前記第4の特徴を取得することと、を含むことを特徴とする話者の確認方法。

【請求項 2】

前記の、予め取得した第1の音声と前記第2の音声を対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換することは、具体的に、

前記第1の音声と前記第2の音声のそれにおけるフレームをつなぎ合わせることと、

前記第1の音声におけるフレームをつなぎ合わせた後の音声を、対応する第1の音声スペクトログラムに変換し、前記第2の音声におけるフレームをつなぎ合わせた後の音声を、対応する第2の音声スペクトログラムに変換することと、を含むことを特徴とする請求項1に記載の話者の確認方法。

【請求項3】

前記の、畳み込みニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得することは、

前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して畠み込み処理を行うことと、

畠み込んだ後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対してブーリング処理を行うことと、

ブーリング後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して次元削減を行い、前記対応する第1の特徴と第2の特徴を得ることと、を含むことを特徴とする請求項1または2に記載の話者の確認方法。

【請求項4】

畠み込みニューラルネットワークを使用して前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行う前に、更に

前記畠み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングすることを含むことを特徴とする請求項1から3のいずれか1項に記載の話者の確認方法。

【請求項5】

前記畠み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングすることは、

交差エントロピー関数を目的関数として使用して、前記畠み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングすることを含むことを特徴とする請求項4に記載の話者の確認方法。

【請求項6】

第2の音声を取得するための取得ユニットと、

予め取得した第1の音声と前記第2の音声を、対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換するための変換ユニットと、

畠み込みニューラルネットワークを使用して、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得するための第1の抽出ユニットと、

時間遅延ニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得するための第2の抽出ユニットと、

前記第3の特徴と前記第4の特徴により、話者を確認するための確認ユニットと、を含み、

前記第2の抽出ユニットは、具体的に、

前記第1の音声スペクトログラムの前後のフレームに対応する前記第1の特徴に対してつなぎ合わせを行い、前記第2の音声スペクトログラムの前後のフレームと対応する前記第2の特徴に対してつなぎ合わせを行うための第2のつなぎ合わせサブユニットと、

つなぎ合わせた後の前記第1の特徴とつなぎ合わせた後の前記第2の特徴のそれに対して線形変換と次元削減を行い、対応する前記第3の特徴と前記第4の特徴を取得するための第2の変換サブユニットを含むことを特徴とする話者の確認装置。

【請求項7】

前記変換ユニットは、

具体的に、それぞれ前記第1の音声と前記第2の音声におけるフレームをつなぎ合わせるためのつなぎ合わせサブユニットと、

前記第1の音声におけるフレームをつなぎ合わせた後の音声を、対応する第1の音声スペクトログラムに変換し、前記第2の音声におけるフレームをつなぎ合わせた後の音声を、対応する第2の音声スペクトログラムに変換するための変換サブユニットと、を含むことを特徴とする請求項6に記載の話者の確認装置。

【請求項8】

前記第1の抽出ユニットは、

前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して置み込み処理を行うための置み込みサブユニットと、

置み込んだ後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対してブーリング処理を行うためのブーリングサブユニットと、

ブーリング後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して次元削減を行い、前記対応する第1の特徴と第2の特徴を得るために次元削減サブユニットと、を含むことを特徴とする請求項6または7に記載の話者の確認装置。

【請求項9】

更に、前記置み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングするためのトレーニングユニットを含むことを特徴とする請求項6から8のいずれか1項に記載の話者の確認装置。

【請求項10】

前記トレーニングユニットは、具体的に、交差エントロピー関数を目的関数として使用して、前記置み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングするために使われることを特徴とする請求項9に記載の話者の確認装置。

【発明の詳細な説明】

【参照関係】

【0001】

本願は、2017年4月1日に中国に出願された、特許名称が「話者の確認方法及び装置」である中国特許出願2017102146660号を引用し、その全体が本出願において参照により引用されている。

【技術分野】

【0002】

本発明は音声情報処理分野に関し、より具体的には、話者の確認方法及び話者の確認装置に関する。

【背景技術】

【0003】

話者の確認方法とは、音声に含まれる声紋特徴により、話者に対して身分検証を行う方法である。話者の確認を行う時に、ユーザは、音声を予め保存しておき、そして検証音声を入力する。検証音声とシステムに予め保存した音声を比較すると、検証音声がそのユーザの発した音声であるかを判断することができ、ユーザの身分認証が実現される。

【0004】

現在、話者の確認方法は、統計モデルを主にしており、性能が良い話者の確認方法は、一般的に*i-vector*モデルとPLDAモデルに基づくものである。*i-vector*モデルは、音声信号に対して、下記の線形モデルを構築する。

【0005】

$$X = Tw + v$$

【0006】

ただし、*X*は音声信号のMFCC特徴であり、*T*は低次元行列であり、*w*はセンテンスベクトル、すなわち*i-vector*であり、*v*はガウス雑音である。当該モデルは、本

質的に確率的PCAモデルである。実際の応用では、一般的に、音声空間は複数の領域に分けられ、それぞれの領域に対して上記線形モデリングが行われ、全ての領域はセンテンスペクトル w を共有する。 w は低次元ベクトルであり、その中に話者、音声内容、チャンネル等の情報を含む。話者情報の区別性を高めるために、次のようなPLDAモデルを取り入れる。

【0007】

$$w = H u + K c + n$$

【0008】

ただし、 u は話者ベクトルであり、 c は表現ベクトルであり、発音方式、チャンネル等を含み、 n はガウス雑音である。PLDAモデルは、話者特徴と表現特徴を区別する。

【0009】

上記モデルは、一般に公知のMFCC特徴に基づくものであり、モデルにより話者情報を分離する。当該方法は、音声信号の分布状態に基づいてモデリングを行うものであるので、より良い結果を得るにはより多くのデータが必要であり、その計算量が多く、かつチャンネル、ノイズ及び時間的変化等の要因の影響を受けやすい。

【発明の概要】

【発明が解決しようとする課題】

【0010】

本発明は、上記の必要なデータが多く、計算量が多くかつロバスト性が悪い等の課題又は少なくともその一部の課題を解決するために、話者の確認方法及び話者の確認装置を提供することを目的とする。

【課題を解決するための手段】

【0011】

本発明のある局面に係る話者の確認方法は、第2の音声を取得することと、予め取得した第1の音声と前記第2の音声を、対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換することと、畳み込みニューラルネットワークを使用して、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得することと、時間遅延ニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得することと、前記第3の特徴と前記第4の特徴により、話者を確認することと、を含む。

【0012】

具体的に、前記の、予め取得した第1の音声と前記第2の音声を対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換することは、それぞれ前記第1の音声と前記第2の音声におけるフレームをつなぎ合わせることと、それぞれ前記第1の音声におけるフレームをつなぎ合わせた後の音声を、対応する第1の音声スペクトログラムに変換し、前記第2の音声におけるフレームをつなぎ合わせた後の音声を、対応する第2の音声スペクトログラムに変換することと、を含む。

【0013】

具体的に、前記の、畳み込みニューラルネットワークを使用して、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得することは、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して畳み込み処理を行うことと、畳み込んだ後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対してブーリング処理を行うことと、ブーリング後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して次元削減を行い、前記対応する第1の特徴と第2の特徴を得ることと、を含む。

【0014】

具体的に、前記の、時間遅延ニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得すること

は、前記第1の音声スペクトログラムの前後のフレームに対応する前記第1の特徴に対してつなぎ合わせを行い、前記第2の音声スペクトログラムの前後のフレームに対応する前記第2の特徴に対してつなぎ合わせを行うことと、つなぎ合わせた後の前記第1の特徴とつなぎ合わせた後の前記第2の特徴のそれぞれに対して線形変換と次元削減を行い、対応する前記第3の特徴と前記第4の特徴とを取得することと、を含む。

【0015】

具体的に、畳み込みニューラルネットワークを使用して前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行う前に、更に、前記畳み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングすることを含む。

【0016】

具体的に、前記畳み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングすることは、交差エントロピー関数を目的関数として使用して、前記畳み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングすることを含む。

【0017】

本発明の他の局面に係る話者の確認装置は、第2の音声を取得するための取得ユニットと、予め取得した第1の音声と前記第2の音声を、対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換するための変換ユニットと、畳み込みニューラルネットワークを使用して、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得するための第1の抽出ユニットと、時間遅延ニューラルネットワークを使用して、前記第1の特徴と前記第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得するための第2の抽出ユニットと、前記第3の特徴と前記第4の特徴により、話者を確認するための確認ユニットと、を含む。

【0018】

具体的に、更に、前記畳み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングするためのトレーニングユニットを含む。

【0019】

具体的に、前記変換ユニットは、更に、前記第1の音声と前記第2の音声のそれぞれにおけるフレームをつなぎ合わせるためのつなぎ合わせサブユニットと、前記第1の音声におけるフレームをつなぎ合わせた後の音声を、対応する第1の音声スペクトログラムに変換し、前記第2の音声におけるフレームをつなぎ合わせた後の音声を、対応する第2の音声スペクトログラムに変換するための変換サブユニットと、を含む。

【0020】

具体的に、前記第1の抽出ユニットは、前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して畳み込み処理を行うための畳み込みサブユニットと、畳み込んだ後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対してブーリング処理を行うためのブーリングサブユニットと、ブーリング後の前記第1の音声スペクトログラムと前記第2の音声スペクトログラムのそれぞれに対して次元削減を行うための次元削減サブユニットと、を含む。

【0021】

具体的に、前記第2の抽出ユニットは、前記第1の音声スペクトログラムの前後のフレームに対応する前記第1の特徴に対してつなぎ合わせを行い、前記第2の音声スペクトログラムの前後のフレームに対応する前記第2の特徴に対してつなぎ合わせを行うための第2のつなぎ合わせサブユニットと、前記第1の特徴をつなぎ合わせた後の特徴と前記第2の特徴をつなぎ合わせた後の特徴のそれぞれに対して線形変換と次元削減を行い、対応する第3の特徴と第4の特徴を取得するための第2の変換サブユニットと、を含む。

【0022】

具体的に、前記トレーニングユニットは、交差エントロピー関数を目的関数として使用

して、前記置み込みニューラルネットワークと前記時間遅延ニューラルネットワークとをトレーニングするために使われる。

【発明の効果】

【0023】

本発明は、話者の確認方法及び話者の確認装置を提供しており、置み込みニューラルネットワークと時間遅延ニューラルネットワークとを組み合わせることで、第1の音声と第2の音声に対して二回の特徴抽出を行い、最終的に抽出した第3の特徴と第4の特徴とを比較することにより、話者の確認が実現される。本発明は、計算が簡単で、口バスト性が強く、良好な認識効果を達成することができる。

【図面の簡単な説明】

【0024】

【図1】本発明の実施例に係る話者の確認方法のフローチャートである。

【図2】置み込みニューラルネットワークと時間遅延ニューラルネットワークのモデルの構成図である。

【図3】本発明の実施例に係る話者の確認装置の構成図である。

【図4】本発明の他の実施例に係る話者の確認装置の構成図である。

【発明を実施するための形態】

【0025】

以下、添付の図面と実施例を参照して、本発明を実施するための形態について、より詳細に説明する。以下の実施例は、本発明を説明するためのものに過ぎず、本発明の範囲を限定するものではない。

【0026】

図1は、本発明の実施例に係る話者の確認方法のフローチャートであり、第2の音声を取得するステップS1と、予め取得した第1の音声と第2の音声を、対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換するステップS2と、置み込みニューラルネットワークを使用して、第1の音声スペクトログラムと第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得するステップS3と、時間遅延ニューラルネットワークを使用して、第1の特徴と第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得するステップS4と、第3の特徴と第4の特徴により、話者を確認するステップS5、を含む。

【0027】

具体的に、ステップS1において、第2の音声を取得する。第2の音声は、話者が新たに入力した音声であり、検証する必要がある音声である。

【0028】

ステップS2において、第1の音声は、話者が予め入力した音声であり、それぞれの第1の音声は、一つの話者ラベルと対応し、話者ラベルにより、話者を一意に確認することができる。第1の音声は、複数の話者の音声を含むことができ、それぞれの話者は、一つ又は複数の第1の音声と対応することができる。第1の音声におけるフレームをつなぎ合わせて、第1の音声スペクトログラムに変換し、第2の音声におけるフレームをつなぎ合わせて、第2の音声スペクトログラムに変換する。第1の音声スペクトログラムと第2の音声スペクトログラムの横軸は時間を表し、縦軸は周波数を表し、色又は輝度は振幅を表す。

【0029】

ステップS3において、置み込みニューラルネットワークを使用して、第1の音声スペクトログラムと第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得することは、第1の音声スペクトログラムと第2の音声スペクトログラムのそれぞれに対して置み込み処理を行うことと、置み込んだ後の第1の音声スペクトログラムと第2の音スペクトルのそれぞれに対してブーリング処理を行うことと、ブーリング後の第1の音声スペクトログラムと第2の音声スペクトログラムのそれぞれに対して次元削減を行い、対応する第1の特徴と第2の特徴を得ることと、を含む。

【0030】

具体的に、畳み込みネットワークは、複数の畳み込み層とプーリング層を含んでもよい。それぞれの畳み込み層の畳み込みカーネルの数とサイズは必要に応じて調整してもよい。それぞれの畳み込みカーネルを使用して、第1の音声スペクトログラムと第2の音声スペクトログラムに対して畳み込みを行う時に、何れも一つの特徴プレーンを生成することができる。それぞれの畳み込み層の後ろには、一つのプーリング層を接続してもよい。プーリング層は、最大プーリング層又は平均プーリング層であってもよい。プーリング層のウインドウが重なり合っていてもよく、重なり合っていないてもよい。プーリング層のウインドウのサイズは必要に応じて調整してもよい。最後のプーリング層から得られた特徴プレーンに対して線形変換を行い、次元削減することで、対応する第1の特徴と第2の特徴を得る。もちろん、その他の次元削減方法を使用してもよく、本発明は、これを限定しない。

【0031】

ステップS4において、時間遅延ネットワークは、複数の時間遅延層を含んでもよく、それぞれの時間遅延層の入力特徴は、直前の時間遅延層の出力特徴であり、第1の音声と対応する第1の時間遅延層の入力は、第1の特徴であり、第2の音声と対応する第1の時間遅延層の入力は、第2の特徴である。第1の特徴と第2の特徴は、それぞれ第1の音声スペクトログラムと第2の音声スペクトログラムにより得られるので、第1の特徴と第2の特徴は、音声スペクトログラムにおけるフレームと1対1に対応する。第1の音声スペクトログラムについて、第1の時間遅延層を通過する時に、まず、第1の音声スペクトログラムにおける現在フレームの前後のいくつかのフレームに対応する第1の特徴をつなぎ合わせて、つなぎ合わせた後の特徴に対して、線形変換と次元削減を実行して、第1の時間遅延層の出力特徴を得る。第1の時間遅延層の出力特徴を第2の時間遅延層の入力特徴とする。

【0032】

第2の時間遅延層を通過する時に、まず、第1の音声スペクトログラムにおける現在フレームの前後のいくつかのフレームに対応する第1の時間遅延層の出力特徴をつなぎ合わせて、つなぎ合わせた後の特徴に対して、線形変換と次元削減を実行して、第2の時間遅延層の出力特徴を得る。ただし、本発明は具体的な次元削減の方法を限定しない。

【0033】

第2の音声スペクトログラムと第1の音声スペクトログラムとは、時間遅延ネットワークにおけるプロセスが同じである。第1の音声スペクトログラムについて、最後の時間遅延層の出力特徴は、第3の特徴である。第2の音声スペクトログラムについて、最後の時間遅延層の出力特徴は、第4の特徴である。

【0034】

ステップS5において、第1の音声と第2の音声に対して、ニューラルネットワークによりフォワード演算を行い、それぞれ第1の音声と第2の音声におけるそれぞれのフレームに対応する第3の特徴と第4の特徴を抽出する。それぞれのフレームの特徴によって、任意の統計モデルを使用して、話者を確認することができる。例えば、第1の音声と第2の音声のそれぞれにおけるそれぞれのフレームの特徴の平均値を算出し、第1の音声におけるそれぞれのフレームの第3の特徴の平均値と、第2の音声におけるそれぞれのフレームの第4の特徴の平均値と、の間の距離を算出する。距離は、コサイン類似度であってもよいが、この距離に限定されない。算出したコサイン類似度がブリセット閾値よりも大きい時に、第2の音声に対応する話者ラベルにより、現在の話者を確認する。

【0035】

本実施例は、畳み込みニューラルネットワークと時間遅延ニューラルネットワークを組み合わせることで、第1の音声と第2の音声に対して二回の特徴抽出を行い、最終的に抽出した第3の特徴と第4の特徴とを比較することにより、話者の確認が実現される。本発明は、計算が簡単で、ロバスト性が強く、良好な認識効果を達成することができる。

【0036】

図2は、畳み込みニューラルネットワークと畳み込みニューラルネットワークの構成図であり、図2に示すように、示された畳み込みニューラルネットワークの入力は、スペクトログラムである。畳み込みニューラルネットワークは、二つの畳み込み層があり、一つ目の畳み込み層の畳み込みカーネルは、128個であり、それぞれの畳み込みカーネルのサイズは、 6×33 であり、一つ目のブーリング層のブーリングウィンドウのサイズは、 3×11 である。二つ目の畳み込み層の畳み込みカーネルは、256個であり、それぞれの畳み込みカーネルのサイズは、 2×8 である。二つ目のブーリング層のブーリングウィンドウのサイズは、 2×4 である。二つ目のブーリング層から得られた256個の特徴ブレーンに対して、次元削減を行い、512個の特徴まで次元削減され、512個のニューロンに対応する。時間遅延ネットワークは、二つの時間遅延層があり、それぞれの時間遅延層は、タイミングつなぎ合わせにより前後のコンテキスト情報を拡張する。一つ目の時間遅延層は、前後の各2フレームの音声スペクトログラムに対応する第1の特徴をつなぎ合わせ、二つ目の時間遅延層は、前後の各4フレームの音声スペクトログラムに対応する一つ目の時間遅延層の出力特徴をつなぎ合わせる。それぞれの時間遅延層は、まず、一回線形変換され、そして、更に一つの次元削減層をつなぎ合わせる。それぞれの次元削減層は、時間遅延層に対して次元削減を行い、400個の特徴を出力する。最後に、二つ目の時間遅延層の次元削減層から出力された400個の特徴に対して、線形変換を行う。

【0037】

上記実施例に基づいて、本実施において畳み込みニューラルネットワークを使用して第1の音声スペクトログラムと第2の音声スペクトログラムに対して特徴抽出を行う前に、更に畳み込みニューラルネットワークと時間遅延ニューラルネットワークをトレーニングすることを含む。

【0038】

具体的に、トレーニングする前に、確認する必要がある話者の音声を取得し、確認する必要がある話者の音声を、トレーニングセットとする。トレーニングを行う時に、音声におけるそれぞれのフレームを学習サンプルとして、畳み込みニューラルネットワークと時間遅延ニューラルネットワークにより二回の特徴抽出を行ってから、最後に算出された特徴に対して、線形変換を行い、その後、当該フレームによって予測された話者ラベルは当該フレームと対応する実際の話者ラベルであるかを確認し、話者が確認した誤差情報を使用して、畳み込みニューラルネットワークと時間遅延ニューラルネットワークにおけるパラメータを逆調整する。目的関数は、交差エントロピー関数である。トレーニングする時に使った逆伝播アルゴリズムは、NステップSGD（Natural Step SGD）アルゴリズム又は任意のディープニューラルネットワートトレーニング方法であってもよい。

【0039】

本実施例において、音声におけるそれぞれのフレームを学習サンプルとして使用して、畳み込みニューラルネットワークと時間遅延ニューラルネットワークに対してトレーニングを行い、畳み込みニューラルネットワークと時間遅延ニューラルネットワークにおけるパラメータを調整する。当該トレーニング方法に必要なデータ量は少なく、最適化した後のパラメータを使用するので、話者の確認の精度を向上させることができる。

【0040】

図3は、本発明の実施例に係る話者の確認装置の構成図であり、図3に示すように、取得ユニット1、変換ユニット2、第1の抽出ユニット3、第2の抽出ユニット4、及び確認ユニット5を含む。ただし、取得ユニット1は、第2の音声を取得するためのものであり、変換ユニット2は、予め取得した第1の音声と第2の音声を対応する第1の音声スペクトログラムと第2の音声スペクトログラムに変換するためのものであり、第1の抽出ユニット3は、畳み込みニューラルネットワークを使用して、第1の音声スペクトログラムと第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得するためのものであり、第2の抽出ユニット4は、時間遅延ニューラルネットワ

ークを使用して、第1の特徴と第2の特徴に対して特徴抽出を行い、対応する第3の特徴と第4の特徴を取得するためのものであり、確認ユニット5は、第3の特徴と第4の特徴により、話者を確認するためのものである。

【0041】

具体的に、取得ユニット1は、第2の音声を取得する。第2の音声は、話者が新たに入力した音声であり、検証する必要な音声である。

【0042】

変換ユニット2は、具体的に、第1の音声と第2の音声のそれそれにおけるフレームをつなぎ合わせるためのつなぎ合わせサブユニットと、第1の音声におけるフレームをつなぎ合わせた後の音声を、対応する第1の音声スペクトログラムに変換し、第2の音声におけるフレームをつなぎ合わせた後の音声を、対応する第2の音声スペクトログラムに変換するための変換サブユニットと、を含む。

【0043】

具体的に、変換ユニット2は、第1の音声におけるフレームをつなぎ合わせてから、第1の音声スペクトログラムに変換し、第2の音声におけるフレームをつなぎ合わせてから、第2の音声スペクトログラムに変換する。第1の音声は、話者が予め入力した音声であり、それぞれの第1の音声は、一つの話者ラベルと対応し、話者ラベルにより、話者を一意に確認することができる。第1の音声は、複数の話者の音声を含むことができ、それぞれの話者は、一つ又は複数の第1の音声と対応することができる。第1の音声スペクトログラムと第2の音声スペクトログラムの横軸は時間を表し、縦軸は周波数を表し、色又は輝度は振幅を表す。

【0044】

第1の抽出ユニット3は、第1の音声スペクトログラムと第2の音声スペクトログラムのそれぞれに対して置み込み処理を行うための置み込みサブユニットと、置み込んだ後の第1の音声スペクトログラムと第2の音声スペクトログラムのそれぞれに対してブーリング処理を行うためのブーリングサブユニットと、ブーリング後の第1の音声スペクトログラムと第2の音声スペクトログラムのそれぞれに対して次元削減を行うための次元削減サブユニットと、を含む。

【0045】

具体的に、第1の抽出ユニット3は、置み込みニューラルネットワークを使用して、第1の音声スペクトログラムと第2の音声スペクトログラムに対して特徴抽出を行い、対応する第1の特徴と第2の特徴を取得する。置み込みネットワークは、複数の置み込み層とブーリング層を含んでもよい。それぞれの置み込み層の置み込みカーネルの数とサイズは必要に応じて調整してもよい。それぞれの置み込み層の後ろには、一つのブーリング層を接続してもよい。ブーリング層は、最大ブーリング層又は平均ブーリング層であってもよい。ブーリング層のウィンドウが重なり合っていてもよく、重なり合っていないなくてもよい。ブーリング層のウィンドウのサイズは必要に応じて調整してもよい。最後のブーリング層から得られた特徴プレーンに対して線形変換を行い、次元削減することで、対応する第1の特徴と第2の特徴を得る。もちろん、その他の次元削減方法を使用してもよく、本発明は、これを限定しない。

【0046】

ただし、置み込みニューラルネットワークを使用して、第1の音声スペクトログラムと第2の音声スペクトログラムに対して特徴抽出を行う時に、それぞれの置み込みカーネルは、一枚の特徴プレーンを生成する。置み込みカーネルの数が多いと、複数枚の特徴プレーンを生成することができ、それぞれのプレーンは、多数の特徴があり、それぞれの置み込み層の後ろに一つのブーリング層をつなぎ合わせるが、特徴の数が依然として多く、算出スピードを大幅に低減させる。そこで、置み込みニューラルネットワークにおける最後のブーリング層から得られた特徴プレーンに対して次元削減を行う必要がある。低次元行列を使用して次元削減を行うことができるが、本実施例は、次元削減の方法に限定されない。

【0047】

本実施例は、最後のプーリング層から得られた特徴プレーンに対して次元削減を行うことで、第1の音声に対応する第1の特徴と、第2の音声に対応する第2の特徴が得られ、算出スピードを大幅に向上させる。

【0048】

更に、第2の抽出ユニット4が使用した時間遅延ネットワークは、例えば複数の全接続の時間遅延層のような複数の時間遅延層を含むことができ、それぞれの時間遅延層における第2の抽出ユニット4は、前後の各フレームをつなぎ合わせることで、前後のコンテキスト情報に対して拡張を行う。前後の各フレームをつなぎ合わせる数は、必要に応じて設置することができる。

【0049】

具体的に、第2の抽出ユニット4は、第1の音声スペクトログラムの前後フレームに対応する第1の特徴に対してつなぎ合わせを行い、第2の音声スペクトログラムの前後フレームに対応する第2の特徴に対してつなぎ合わせを行うための第2のつなぎ合わせサブユニットと、第1の特徴をつなぎ合わせた後の特徴と、第2の特徴をつなぎ合わせた後の特徴のそれぞれに対して線形変換と次元削減を行い、対応する第3の特徴と第4の特徴を取得するための第2の変換サブユニットと、を含む。

【0050】

具体的に、それぞれの時間遅延層において、第2のつなぎ合わせサブユニットがつなぎ合わせるフレームの数は同じであってもよいし、異なる時間遅延層において、第2のつなぎ合わせサブユニットがつなぎ合わせるフレームの数は異なってもよい。つなぎ合わせウインドウは重なり合っていてもよい。

【0051】

第1の音声スペクトログラムと第2の音声スペクトログラムとは、時間遅延ネットワークにおけるプロセスが同じである。第1の音声スペクトログラムについて、最後の時間遅延層の出力特徴は、第3の特徴である。第2の音声スペクトログラムについて、最後の時間遅延層の出力特徴は、第4の特徴である。

【0052】

本発明の実施例は、時間遅延ニューラルネットワークを使用して、第1の特徴と第2の特徴に対して、特徴抽出を行い、対応する第3の特徴と第4の特徴を取得する。時間遅延ニューラルネットワークは、特徴に対して強い抽出能力を有し、話者の正確な確認に対して基礎を築く。

【0053】

確認ユニット5は、第1の音声と第2の音声に対して、畳み込みニューラルネットワークと時間遅延ニューラルネットワークによりフォワード演算を行い、第1の音声と第2の音声のそれぞれにおけるそれぞれのフレームに対応する第3の特徴と第4の特徴を抽出する。それぞれのフレームの特徴によって、任意の統計モデルを使用して、話者を確認することができる。例えば、第1の音声と第2の音声のそれぞれにおけるそれぞれのフレームの特徴の平均値を算出し、第1の音声におけるそれぞれのフレームの第3の特徴の平均値と、第2の音声におけるそれぞれのフレームの第4の特徴の平均値と、の間の距離を算出する。距離は、コサイン類似度であってもよいが、この距離に限定されない。算出したコサイン類似度がプリセット閾値よりも大きい時に、第2の音声に対応する話者ラベルにより、現在の話者を確認する。

【0054】

本実施例は、畳み込みニューラルネットワークと時間遅延ニューラルネットワークを組み合わせることで、第1の音声と第2の音声に対して二回の特徴抽出を行い、最終的に抽出した第3の特徴と第4の特徴を比較することにより、話者の確認が実現される。本発明は、計算が簡単で、ロバスト性が強く、良好な認識効果を達成することができる。

【0055】

図4は、本発明の実施例に係る話者の確認装置の構成図であり、図4に示すように、上

記各実施例に基づいて、装置は、更に畳み込みニューラルネットワークと時間遅延ニューラルネットワークとをトレーニングするためのトレーニングユニット6を含む。

【0056】

具体的に、トレーニングする前に、確認する必要がある話者の音声を取得し、確認する必要がある話者の音声を、トレーニングセットとする。トレーニングを行う時に、トレーニングユニット6は、音声中におけるそれぞれのフレームを学習サンプルとして、畳み込みニューラルネットワークと時間遅延ニューラルネットワークにより二回の特徴抽出を行ってから、最後に算出された特徴に対して、線形変換を行い、その後、当該フレームによって予測された話者ラベルは当該フレームと対応する実際の話者ラベルであるかを確認し、話者が確認した誤差情報を使用して、畳み込みニューラルネットワークと時間遅延ニューラルネットワークにおけるパラメータを逆調整する。目的関数は、交差エントロピー関数である。トレーニングする時に使った逆伝播アルゴリズムは、NステップSGD（Natural Steps Stochastic Gradient Descent Steps Scent、確率的勾配降下法）アルゴリズム又は任意のディープニューラルネットワークトレーニング方法であってもよい。

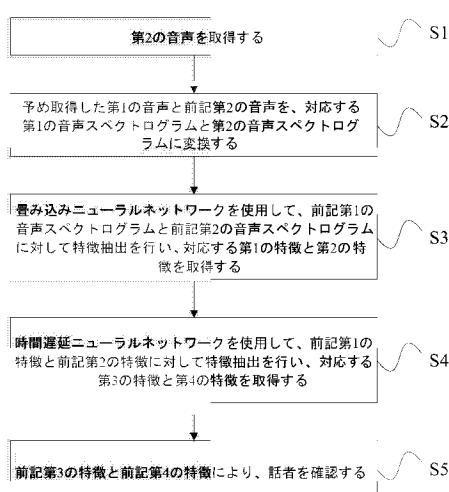
【0057】

本実施例において、音声におけるそれぞれのフレームを学習サンプルとして使用して、畳み込みニューラルネットワークと時間遅延ニューラルネットワークに対してトレーニングを行い、畳み込みニューラルネットワークと時間遅延ニューラルネットワークにおけるパラメータを調整する。当該トレーニング方法に必要なデータ量が少なく、最適化した後のパラメータを使用するので、話者の確認の精度を向上させることができる。

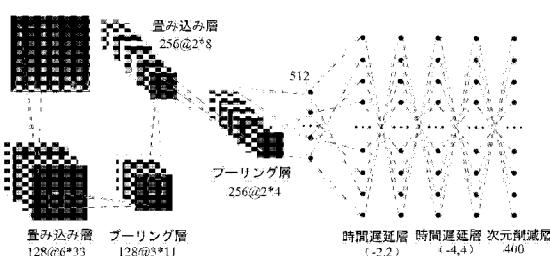
【0058】

最後に、以上説明した方法は、一つの好適な実施案に過ぎず、本発明の保護範囲を限定するものではない。本発明の要旨と原則を逸脱しない範囲においてなされる様々な修正、等価交換、改善等は、何れも本開示の保護範囲に含まれる。

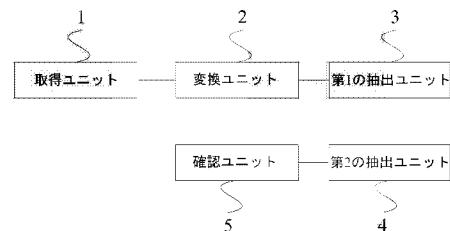
【図1】



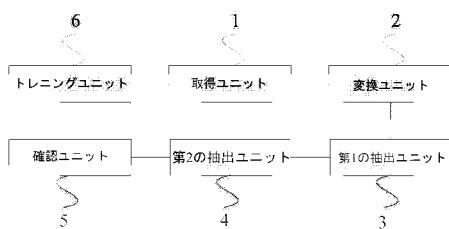
【図2】



【図3】



【図4】



フロントページの続き

発明者 王 東
中華人民共和国100084北京市海淀区双清路30号清华大学
発明者 李 藍天
中華人民共和国100084北京市海淀区双清路30号清华大学
発明者 鄭 方
中華人民共和国100084北京市海淀区双清路30号清华大学

審査官 渡部 幸和

参考文献 中国特許出願公開第102034472(CN, A)
中国特許出願公開第106128465(CN, A)
中川 聖一, "再訪: ニューラルネットワークによる音声処理", 電子情報通信学会技術研究報告 Vol. 113 No. 161, 年月日, ~

調査した分野 , DB名
G10L 15/00-17/26
G06N 3/04