

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2018-517927  
(P2018-517927A)

(43) 公表日 平成30年7月5日(2018.7.5)

(51) Int.Cl.		F I		テーマコード (参考)
<b>G 1 0 L 17/18</b>	<b>(2013.01)</b>	G 1 0 L	17/18	
<b>G 1 0 L 17/00</b>	<b>(2013.01)</b>	G 1 0 L	17/00	2 0 0 C

審査請求 有 予備審査請求 未請求 (全 36 頁)

(21) 出願番号 特願2017-556908 (P2017-556908)  
 (86) (22) 出願日 平成28年7月27日 (2016.7.27)  
 (85) 翻訳文提出日 平成29年12月20日 (2017.12.20)  
 (86) 国際出願番号 PCT/US2016/044181  
 (87) 国際公開番号 WO2017/039884  
 (87) 国際公開日 平成29年3月9日 (2017.3.9)  
 (31) 優先権主張番号 14/846, 187  
 (32) 優先日 平成27年9月4日 (2015.9.4)  
 (33) 優先権主張国 米国 (US)

(71) 出願人 502208397  
 グーグル エルエルシー  
 アメリカ合衆国 カリフォルニア州 94  
 043 マウンテン ビュー アンフィシ  
 アター パークウェイ 1600  
 (74) 代理人 100108453  
 弁理士 村山 靖彦  
 (74) 代理人 100110364  
 弁理士 実広 信哉  
 (74) 代理人 100133400  
 弁理士 阿部 達彦  
 (72) 発明者 ゲオルク・ハイゴルト  
 アメリカ合衆国・カリフォルニア・940  
 43・マウンテン・ビュー・アンフィシ  
 アター・パークウェイ・1600

最終頁に続く

(54) 【発明の名称】 話者検証のためのニューラルネットワーク

(57) 【要約】

本明細書は一般に、(i) 話者検証モデルに対するニューラルネットワークをトレーニングするステップと、(ii) クライアントデバイスでユーザを加入させるステップと、(iii) ユーザのアイデンティティを前記ユーザの音声の特性に基づいて検証するステップとを含む、話者検証に関連するシステム、方法、デバイス、および他の技術を説明する。幾つかの実装はコンピュータ実行型の方法を含む。前記方法は、コンピューティングデバイスで、前記コンピューティングデバイスのユーザの発声を特徴付けるデータを受信するステップを含むことができる。話者表現を、前記コンピューティングデバイスで、前記コンピューティングデバイス上のニューラルネットワークを用いて前記発声に対して生成することができる。前記ニューラルネットワークは、それぞれが、(i) 第1の発声を特徴付けるデータおよび1つまたは複数の第2の発声を特徴付けるデータを含み、(ii) マッチング話者サンプルまたは非マッチング話者サンプルとしてラベル付けされる、複数のトレーニング・サンプルに基づいてトレーニングされることことができる。

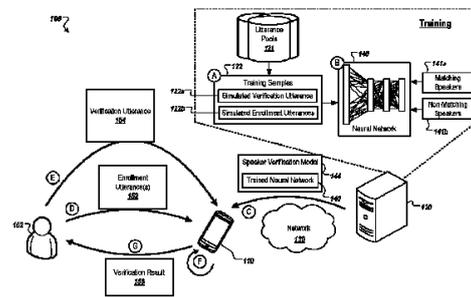


FIG. 1

## 【特許請求の範囲】

## 【請求項 1】

コンピューティングデバイスで、前記コンピューティングデバイスのユーザの発声を特徴付けるデータを受信するステップと、

前記コンピューティングデバイスで、前記コンピューティングデバイス上のニューラルネットワークを用いて前記発声に対する話者表現を生成するステップであって、前記ニューラルネットワークが、それぞれ、

( i ) 第 1 の発声を特徴付けるデータおよび 1 つまたは複数の第 2 の発声を特徴付けるデータを含み、

( i i ) 前記第 1 の発声の話者が前記 1 つまたは複数の第 2 の発声の話者と同一であるかどうかに従って、マッチングサンプルまたは非マッチングサンプルとしてラベル付けされる、

複数のトレーニング・サンプルに基づいてトレーニングされている、ステップと、

前記コンピューティングデバイスで、前記コンピューティングデバイスの認証されたユーザに対する話者モデルにアクセスするステップと、

前記コンピューティングデバイスで、前記話者モデルに関する前記発声に対する前記話者表現を評価して、前記発声の前記コンピューティングデバイスの前記認証されたユーザにより話された可能性があったかどうかを判定するステップと、

を含む、コンピュータ実行型の方法。

## 【請求項 2】

前記複数のトレーニング・サンプルの各々は、各発声グループが前記各発声グループに対する前記対応する話者の発声のみから構成されるように、前記第 1 の発声および前記 1 つまたは複数の第 2 の発声を異なる話者に対応する発声グループから選択することにより生成された、請求項 1 に記載のコンピュータ実行型の方法。

## 【請求項 3】

前記コンピューティングデバイスの前記認証されたユーザの 1 組の発声を取得するステップと、

前記発声に対する各話者表現を生成するために、各発声を前記 1 組の発声から前記ニューラルネットワークに入力するステップと、

前記認証されたユーザの前記 1 組の発声における前記発声に対する前記各話者表現の平均に基づいて前記コンピューティングデバイスの前記認証されたユーザに対する前記話者モデルを生成するステップと、

をさらに含む、請求項 1 または 2 に記載のコンピュータ実行型の方法。

## 【請求項 4】

前記ニューラルネットワークがトレーニングされている前記複数のトレーニング・サンプルの何れも、前記コンピューティングデバイスの前記ユーザの前記発声を特徴付けるデータを含まない、請求項 1 乃至 3 の何れか 1 項に記載のコンピュータ実行型の方法。

## 【請求項 5】

前記コンピューティングデバイスで、前記発声に対する前記話者表現を生成するステップは、前記発声の全体を特徴付けるデータを前記ニューラルネットワークで単一のパスで前記ニューラルネットワークを通じて処理するステップを含む、請求項 1 乃至 4 の何れか 1 項に記載のコンピュータ実行型の方法。

## 【請求項 6】

前記ユーザの前記発声の前記コンピューティングデバイスの前記認証されたユーザにより話された可能性があったと判定したことに応答して、機能を前記コンピューティングデバイス上で実施するステップをさらに含む、請求項 1 乃至 5 の何れか 1 項に記載のコンピュータ実行型の方法。

## 【請求項 7】

第 1 の 1 組の発声に対するニューラルネットワークの出力に基づいて特定の話者に対する話者モデルを決定するステップであって、前記第 1 の 1 組の発声は前記特定の話者の複

10

20

30

40

50

数の異なる発声を含む、ステップと、

前記第 1 の 1 組の発声内にはない特定の発声に対する前記ニューラルネットワークの出力に基づいて話者表現を決定するステップと、

前記話者表現を前記特定の話者に対する前記話者モデルと比較して、前記特定の発声を前記特定の話者の発声としてまたは前記特定の話者と異なる話者の発声として分類するステップと、

前記発声の前記特定の話者の発声としてまたは前記特定の話者と異なる話者の発声としての前記分類が正しかったかどうかに基づいて前記ニューラルネットワークを調節するステップと、

を含む、コンピュータ実行型の方法。

10

【請求項 8】

コンピューティングシステムで、複数の異なる発声のセットを、前記ニューラルネットワークをトレーニングするためのトレーニングデータとして選択するステップであって、各 1 組の発声は、

( i ) 前記各 1 組の発声に対する第 1 の第 1 の話者の複数の異なる発声と、

( i i ) 前記各 1 組の発声に対する前記第 1 の話者または前記第 1 の話者と異なる前記各 1 組の発声に対する第 2 の話者の何れかの第 2 の発声と、

を含む、ステップと、

前記複数の異なる発声のセットを用いて、前記ニューラルネットワークのトレーニングを反復するステップであって、前記複数の異なる発声のセットからの各 1 組の発声は前記ニューラルネットワークをトレーニングするために異なるトレーニングの反復で使用される、ステップと、

20

を含み、

前記第 1 の 1 組の発声は前記複数の異なる発声のセットから選択され、前記特定の話者は前記第 1 の 1 組の発声に対する前記第 1 の話者である、

請求項 7 に記載のコンピュータ実行型の方法。

【請求項 9】

コンピューティングシステムで、ニューラルネットワークをトレーニングするためのトレーニングデータの複数の異なるサブセットを選択するステップであって、トレーニングデータの各サブセットは、第 1 の話者の各発声を分類する複数の第 1 のコンポーネントおよび前記第 1 の話者または第 2 の話者の発声を特徴付ける第 2 のコンポーネントを含む、ステップと、

30

トレーニングデータの前記選択されたサブセットの各々に対して、

前記第 1 のコンポーネントの各々に対応する各第 1 の話者表現を生成するために、前記第 1 のコンポーネントの各々を前記ニューラルネットワークに入力するステップと、

前記第 2 のコンポーネントに対応する第 2 の話者表現を生成するために、前記第 2 のコンポーネントを前記ニューラルネットワークに入力するステップと、

前記複数の第 1 のコンポーネントに対する前記各第 1 の話者表現の平均に基づいて前記第 1 の話者に対するシミュレートされた話者モデルを決定するステップと、

前記第 2 の話者表現を前記シミュレートされた話者モデルと比較して、前記第 2 のコンポーネントにより特徴付けられた前記発声を前記第 1 の話者の発声としてまたは前記第 1 の話者と異なる話者の発声として分類するステップと、

40

前記第 2 のコンポーネントにより特徴付けられた前記発声の前記第 1 の話者の発声としてまたは前記第 1 の話者と異なる話者の発声として正しく分類されたかどうかに基づいて、前記ニューラルネットワークを調節するステップと、

を含む、コンピュータ実行型の方法。

【請求項 10】

前記第 2 のコンポーネントを前記ニューラルネットワークに入力するステップにตอบสนองして、単一のパスで前記ニューラルネットワークを通じて、前記第 2 のコンポーネントにより特徴付けられた前記発声の全体に対するデータを処理することで、前記第 2 の話者表現

50

を前記ニューラルネットワークで生成するステップをさらに含む、請求項 9 に記載のコンピュータ実行型の方法。

【請求項 11】

トレーニングデータの前記複数の異なるサブセットからトレーニングデータの第 1 のサブセットを選択するステップは、

各発声グループが前記各発声グループに対する前記対応する話者の発声のみを特徴づけるデータから構成されるように、それぞれ異なる話者に対応する複数の発声グループからの前記第 1 の話者に対応する第 1 の発声グループを選択するステップと、

前記第 1 の話者または前記第 2 の話者に対応する第 2 の発声グループを前記複数の発声グループから選択するステップと、

前記複数の第 1 のコンポーネントを前記第 1 の発声グループ内の発声の特徴付ける前記データから決定するステップと、

前記第 2 のコンポーネントを前記第 2 の発声グループ内の発声の特徴付ける前記データから決定するステップと、

を含む、請求項 9 または 10 に記載のコンピュータ実行型の方法。

【請求項 12】

前記第 1 の発声グループおよび前記第 2 の発声グループの少なくとも 1 つは前記複数の発声グループからランダムに選択される、請求項 11 に記載のコンピュータ実行型の方法。

【請求項 13】

前記第 1 の話者はトレーニングデータの前記複数の異なるサブセットの少なくとも幾つかの間で互いに異なり、

前記第 2 の話者はトレーニングデータの前記複数の異なるサブセットの少なくとも幾つかの間で互いに異なる、

請求項 9 乃至 12 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 14】

トレーニングデータの前記複数の異なるサブセットのうちトレーニングデータの第 1 のサブセット内の第 1 のコンポーネントの総数は、トレーニングデータの前記複数の異なるサブセットのうちトレーニングデータの第 2 のサブセット内の第 1 のコンポーネントの総数と異なる、請求項 9 乃至 13 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 15】

記第 2 の話者表現を前記シミュレートされた話者モデルと比較して、前記第 2 のコンポーネントにより特徴付けられた前記発声を前記第 1 の話者の発声としてまたは前記第 1 の話者と異なる話者の発声として分類するステップは、前記第 2 の話者表現からの値および前記シミュレートされた話者モデルからの値の間の距離を決定するステップと、ロジスティック回帰機能を前記距離に適用するステップとを含む、請求項 9 乃至 14 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 16】

前記ニューラルネットワークは複数の隠蔽されたレイヤを含み、前記ニューラルネットワークはソフトマックス出力レイヤを有さない、請求項 9 乃至 15 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 17】

前記ニューラルネットワークは、ローカルに接続された隠蔽されたレイヤと、前記レイヤに続く複数の完全に接続された隠蔽されたレイヤを有するディープ・ニューラルネットワークを含む、請求項 9 乃至 16 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 18】

前記複数の第 1 のコンポーネントにより特徴付けられた前記発声と前記第 2 のコンポーネントにより特徴付けられた前記発声とは全て固定長を有する、請求項 17 に記載のコンピュータ実行型の方法。

【請求項 19】

10

20

30

40

50

前記ニューラルネットワークは、可変長を有する発声を特徴付けるデータでトレーニングされるように構成された長短期メモリ再帰型ニューラルネットワークを含む、請求項 9 乃至 18 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 20】

トレーニングデータの前記複数の異なるサブセットにわたって前記各複数の第 1 のコンポーネントにより特徴付けられた前記発声と、トレーニングデータの前記複数の異なるサブセットにわたって前記各第 2 のコンポーネントにより特徴付けられた前記発声とは全て、同一の単語またはフレーズの発声である、請求項 9 乃至 19 の何れか 1 項に記載のコンピュータ実行型の方法。

【請求項 21】

トレーニングデータの前記複数の異なるサブセットにわたって前記各複数の第 1 のコンポーネントにより特徴付けられた前記発声の少なくとも幾つかと、トレーニングデータの前記複数の異なるサブセットにわたって前記各第 2 のコンポーネントにより特徴付けられた前記発声の少なくとも幾つかとは、異なる単語または異なるフレーズの発声である、請求項 9 乃至 19 の何れか 1 項に記載のコンピュータ実行型の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本明細書の主題は一般に、話者検証タスクで使用されるニューラルネットワークおよび他のモデルに関する。

【背景技術】

【0002】

話者検証は一般に人のアイデンティティを、当該人の音声の特性に基づいて検証することに関する。幾つかのコンピューティングデバイスは、ユーザに、ユーザにより話された会話の 1 つまたは複数のサンプルを当該デバイスに提供することで当該デバイスに「加入」させることができる。当該サンプルから、ユーザの音声を表す話者モデルが決定される。当該デバイスで受信された後続の会話サンプルについて、当該話者モデルに関して処理し評価して、ユーザのアイデンティティを検証してもよい。

【発明の概要】

【課題を解決するための手段】

【0003】

本明細書では一般に、話者検証のためのニューラルネットワーク、または他のタイプのモデルをトレーニングし使用するためのシステム、方法、デバイス、および他の技術を説明する。幾つかの実装では、ニューラルネットワークは、話者検証を実施するコンピューティングデバイスによりアクセス可能な話者検証モデルのコンポーネントであってもよい。一般に、ニューラルネットワークは、それぞれ話者の加入および発声の検証をシミュレートする反復でトレーニングされてもよい。例えば、各トレーニングの反復において、所与の発声に対してニューラルネットワークにより生成された話者表現を話者モデルに関して評価してもよい。シミュレートされた検証発声に対する話者表現の、1 つまたは複数のシミュレートされた加入発声に対する話者表現の組合せ（例えば、平均）との比較に基づいて、ニューラルネットワークのパラメータを、所与の発声を同一の人または加入した人とは異なる人により話されているとして分類するように当該話者検証モデルの能力を最適化するように更新してもよい。これが、当該システムの信頼性を増大する点で利点を有することは理解される。ニューラルネットワークはさらに、当該発声のフレームを独立にまたは逐次的に処理するのではなく、単一のパスでニューラルネットワークを通じて、発声全体を特徴づけるデータを処理するように構成されてもよい。これらのおよび他の実装は下記でより完全に説明され、図面で示される。

【0004】

本明細書で説明する主題の幾つかの実装はコンピュータ実行型の方法を含む。当該方法は、コンピューティングシステムで、ニューラルネットワークをトレーニングするための

10

20

30

40

50

トレーニングデータの複数の異なるサブセットを選択するステップを含むことができる。トレーニングデータの各サブセットは、第1の話者の各発声を分類する複数の第1のコンポーネントと、第1の話者または第2の話者の発声を特徴付ける第2のコンポーネントとを含むことができる。トレーニングデータの当該選択されたサブセットごとに、当該方法は、第1のコンポーネントの各々に対応する各第1の話者表現を生成するために、第1のコンポーネントの各々をニューラルネットワークに入力するステップと、第2のコンポーネントに対応する第2の話者表現を生成するために、第2のコンポーネントをニューラルネットワークに入力するステップと、当該複数の第1のコンポーネントに対する当該各第1の話者表現の平均に基づいて第1の話者に対するシミュレートされた話者モデルを決定するステップと、第2の話者表現を当該シミュレートされた話者モデルと比較して、第2のコンポーネントにより特徴付けられた発声を第1の話者の発声としてまたは第1の話者と異なる話者の発声として分類するステップと、第2のコンポーネントにより特徴付けられた発声が第1の話者の発声としてまたは第1の話者と異なる話者の発声として正しく分類されたかどうかに基づいて、ニューラルネットワークを調節するステップとを含むことができる。

10

【0005】

これらのおよび他の実装は以下の特徴のうち1つまたは複数を含むことができる。

【0006】

第2のコンポーネントをニューラルネットワークに入力するステップにตอบสนองして、第2の話者表現を、単一のパスでニューラルネットワークを通じて、第2のコンポーネントにより特徴付けられた当該発声の全体に対するデータを処理することで、ニューラルネットワークで生成することができる。

20

【0007】

トレーニングデータの当該複数の異なるサブセットからトレーニングデータの第1のサブセットを選択するステップは、各発声グループが当該各発声グループに対する当該対応する話者の発声のみを特徴づけるデータから構成されるように、それぞれ異なる話者に対応する複数の発声グループからの第1の話者に対応する第1の発声グループを選択するステップと、第1の話者または第2の話者に対応する第2の発声グループを発声の当該複数のグループから選択するステップと、当該複数の第1のコンポーネントを第1の発声グループ内の発声を特徴付ける当該データから決定するステップと、第2のコンポーネントを第2の発声グループ内の発声を特徴付ける当該データから決定するステップとを含むことができる。

30

【0008】

第1の発声グループおよび第2の発声グループの少なくとも1つを発声の当該複数のグループからランダムに選択することができる。

【0009】

第1の話者は、互いにトレーニングデータの当該複数の異なるサブセットのうち少なくとも幾つかと異なることができる。第2の話者は、互いにトレーニングデータの当該複数の異なるサブセットのうち少なくとも幾つかと異なることができる。

【0010】

トレーニングデータの当該複数の異なるサブセットのうちトレーニングデータの第1のサブセット内の第1のコンポーネントの総数は、トレーニングデータの当該複数の異なるサブセットのうちトレーニングデータの第2のサブセット内の第1のコンポーネントの総数と異なることができる。

40

【0011】

第2の話者表現を当該シミュレートされた話者モデルと比較して、第2のコンポーネントにより特徴付けられた当該発声を第1の話者の発声としてまたは第1の話者と異なる話者の発声として分類するステップは、第2の話者表現からの値および当該シミュレートされた話者モデルからの値の間の距離を決定するステップと、ロジスティック回帰機能を当該距離に適用するステップとを含むことができる。

50

## 【 0 0 1 2 】

ニューラルネットワークは複数の隠蔽されたレイヤを含むことができる。ニューラルネットワークはソフトマックス出力レイヤを有さなくてもよい。

## 【 0 0 1 3 】

ニューラルネットワークは、ローカルに接続された隠蔽されたレイヤと、当該レイヤに続く複数の完全に接続された隠蔽されたレイヤを有するディープ・ニューラルネットワークを含むことができる。

## 【 0 0 1 4 】

当該複数の第1のコンポーネントにより特徴付けられた当該発声、および第2のコンポーネントにより特徴付けられた当該発声は全て、固定長を有することができる。

10

## 【 0 0 1 5 】

ニューラルネットワークは、可変長を有する発声を特徴付けるデータでトレーニングされるように構成される長短期メモリ再帰型ニューラルネットワークであることができる。

## 【 0 0 1 6 】

トレーニングデータの当該複数の異なるサブセットにわたって各複数の第1のコンポーネントにより特徴付けられた発声、およびトレーニングデータの当該複数の異なるサブセットにわたって各第2のコンポーネントにより特徴付けられた発声は全て、同一の単語またはフレーズの発声であることができる。

## 【 0 0 1 7 】

トレーニングデータの当該複数の異なるサブセットにわたって各複数の第1のコンポーネントにより特徴付けられた発声の少なくとも幾つか、およびトレーニングデータの当該複数の異なるサブセットにわたって各第2のコンポーネントにより特徴付けられた当該発声の少なくとも幾つかは異なる単語または異なるフレーズから成る発声であることができる。

20

## 【 0 0 1 8 】

当該トレーニングされたニューラルネットワークを、話者検証を当該1つまたは複数のコンピューティングデバイスで実施する際に使用するための当該コンピューティングシステムとは別の1つまたは複数のコンピューティングデバイスに提供することができる。

## 【 0 0 1 9 】

第1の話者表現の各々を、第1の話者表現に対応するニューラルネットワークに入力された当該各第1のコンポーネントに基づいてニューラルネットワーク生成することができる。第2の話者表現を、第2の話者表現に対応するニューラルネットワークに入力された第2のコンポーネントに基づいてニューラルネットワークにより生成することができる。

30

## 【 0 0 2 0 】

本明細書で説明する主題の幾つかの実装はコンピューティングデバイスを含むことができる。コンピューティングデバイスは、1つまたは複数のコンピュータプロセッサと、当該1つまたは複数のプロセッサにより実行されたとき、動作を実施させる命令を格納した1つまたは複数のコンピュータ可読媒体とを備えることができる。当該動作は、コンピューティングデバイスで、コンピューティングデバイスのユーザの発声を特徴付けるデータを受信するステップと、コンピューティングデバイスで、コンピューティングデバイス上のニューラルネットワークを用いて当該発声に対する話者表現を生成するステップであって、ニューラルネットワークが、それぞれ(i)第1の発声を特徴付けるデータおよび1つまたは複数の第2の発声を特徴付けるデータを含み、(ii)第1の発声の話者が当該1つまたは複数の第2の発声の話者と同一であるかどうかに従って、マッチングサンプルまたは非マッチングサンプルとしてラベル付けされる、複数のトレーニング・サンプルに基づいてトレーニングされている、ステップと、コンピューティングデバイスで、コンピューティングデバイスの認証されたユーザに対する話者モデルにアクセスするステップと、コンピューティングデバイスで、当該話者モデルに関する当該発声に対する話者表現を評価して、当該発声がコンピューティングデバイスの認証されたユーザにより話された可能性があったかどうかを判定するステップとを含むことができる。

40

50

## 【0021】

各発声グループが当該各発声グループに対する対応する話者の発声のみから構成されるように、第1の発声および当該1つまたは複数の第2の発声を異なる話者に対応する発声グループから選択することによって、当該複数のトレーニング・サンプルの各々を生成することができる。

## 【0022】

当該動作はさらに、コンピューティングデバイスの認証されたユーザの1組の発声を取得するステップと、当該発声に対する各話者表現を生成するために、各発声を当該1組の発声からニューラルネットワークに入力するステップと、当該認証されたユーザの当該1組の発声における当該発声に対する各話者表現の平均に基づいてコンピューティングデバイスの認証されたユーザに対する話者モデルを生成するステップとを含むことができる。

10

## 【0023】

ニューラルネットワークがトレーニングされている当該複数のトレーニング・サンプルの何れも、コンピューティングデバイスのユーザの発声を特徴付けるデータを含まなくてもよい。

## 【0024】

コンピューティングデバイスで、当該発声に対する話者表現を生成するステップは、当該発声の全体をニューラルネットワークで単一のパスでニューラルネットワークを通じて特徴付けるデータを処理するステップを含むことができる。

## 【0025】

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ることができる。ユーザの発声は第1の時間長を有することができる。当該発声に対する話者表現は、当該発声の第1の時間長の全体に対して当該発声を特徴付けるデータを処理するステップを含むことができる。当該動作はさらに、コンピューティングデバイスのユーザの別の発声を受信するステップであって、当該他の発声は、第1の時間長と異なる第2の時間長を有する、ステップと、当該他の発声の第2の時間長の全体に対して当該他の発声を特徴付けるデータを処理することでユーザの当該他の発声に対する第2の話者表現を生成するステップとを含むことができる。

20

## 【0026】

当該動作はさらに、ユーザの発声がコンピューティングデバイスの認証されたユーザにより話された可能性があったと判定したことに応答して、機能をコンピューティングデバイスで実施するステップを含むことができる。当該機能は、コンピューティングデバイスの状態をロックされた状態からロックされていない状態に変更するステップであって、コンピューティングデバイスは、当該ロックされた状態のコンピューティングデバイスの1つまたは複数の能力へのアクセスをブロックするように構成され、コンピューティングデバイスは当該ロックされていない状態のコンピューティングデバイスの当該1つまたは複数の能力へのアクセスを許可するように構成される、ステップを含むことができる。当該話者表現は、決定された発声に基づいて決定されるユーザの音声の区別的な特徴を示すニューラルネットワークの出力を含むことができる。

30

## 【0027】

本明細書で説明する主題の幾つかの実装はコンピュータ実行型の方法を含むことができる。当該方法は、コンピューティングデバイスで、コンピューティングデバイスのユーザの発声を特徴付けるデータを受信するステップを含むことができる。話者表現を、コンピューティングデバイスで、コンピューティングデバイス上のニューラルネットワークを用いて当該発声に対して生成することができる。ニューラルネットワークを、それぞれが(i)第1の発声を特徴付けるデータおよび1つまたは複数の第2の発声を特徴付けるデータを含み、(ii)第1の発声の話者が当該1つまたは複数の第2の発声の話者と同一であるかどうかに従って、マッチングサンプルまたは非マッチングサンプルとしてラベル付けされる、複数のトレーニング・サンプルに基づいてトレーニングすることができる。話者モデルに、コンピューティングデバイスで、コンピューティングデバイスの認証された

40

50

ユーザに対してアクセスすることができる。当該発声に対する話者表現を、コンピューティングデバイスで、当該話者モデルに関して評価して、当該発声がコンピューティングデバイスの認証されたユーザにより話された可能性があったかどうかを判定することができる。

【0028】

本明細書で説明する主題の幾つかの実装はコンピュータ実行型の方法を含むことができる。当該コンピュータ実行型の方法は、第1の1組の発声に対するニューラルネットワークの出力に基づいて特定の話者に対する話者モデルを決定するステップであって、第1のセットは複数の異なる特定の話者の発声を含む、ステップと、第1の1組の発声内にはない特定の話者に対するニューラルネットワークの出力に基づいて話者表現を決定するステップと、当該話者表現を当該特定の話者に対する当該話者モデルと比較して、当該特定の発声を当該特定の話者の発声としてまたは当該特定の話者と異なる話者の発声として分類するステップと、当該発声の当該特定の話者の発声としてまたは当該特定の話者と異なる話者の発声としての当該分類が正しかったかどうかに基づいてニューラルネットワークを調節するステップとを含むことができる。

10

【0029】

これらのおよび他の実装は以下の特徴のうち1つまたは複数を含むことができる。複数の異なる発声のセットを、ニューラルネットワークをトレーニングするためのトレーニングデータとして選択することができる。各1組の発声は、(i)当該各1組の発声に対する第1の第1の話者の複数の異なる発声、および(ii)当該各1組の発声に対する第1の話者、または、第1の話者と異なる当該各1組の発声に対する第2の話者の何れかの第2の発声を含むことができる。当該複数の異なる発声のセットを、ニューラルネットワークを繰り返しトレーニングするために使用することができる。当該複数の異なる発声のセットからの各1組の発声を、ニューラルネットワークをトレーニングするために異なるトレーニングの反復において使用することができる。第1の1組の発声を当該複数の異なる発声のセットから選択することができる。当該特定の話者は第1の1組の発声に対する第1の話者である。

20

【0030】

本明細書で説明する主題のさらなる実装は、当該方法のアクションを実施するように構成され、コンピュータ記憶デバイスで符号化された対応するシステム、装置、およびコンピュータプログラムを含むことができる。1つまたは複数のコンピュータのシステムを、動作中に当該システムに当該アクションを実施させるソフトウェア、ファームウェア、ハードウェア、または当該システムにインストールされたそれらの組合せにより構成することができる。1つまたは複数のコンピュータプログラムを、1つまたは複数のデータ処理装置により実行されたとき、当該装置に当該アクションを実施させる命令を有することによりそのように構成することができる。

30

【0031】

本明細書で説明する主題の幾つかの実装は以下の利点の1つまたは複数を実現しうる。ニューラルネットワークを、ユーザの音声の特性に基づいて話者のアイデンティティをより正確に検証できる話者検証モデルで使用するための話者表現を生成するために、トレーニングすることができる。ターゲット性能レベルを、有限のコンピューティングリソースを有するモバイルコンピューティングデバイスに格納しそこで使用しうるコンパクトなニューラルネットワークで実現することができる。さらに、ニューラルネットワークを、話者検証プロセスの検証および加入段階をシミュレートする方式でトレーニングしてもよい。したがって、ニューラルネットワークは、当該話者検証プロセスのトレーニング段階と検証および加入段階との間の対称性のため良好な性能を実現しうる。本明細書で説明するアプローチに従うニューラルネットワークをトレーニングする利益は、事前選択された数の話者の間の特定の話者に属するとして入力を分類するために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他のアプローチと対照的に、より多くの数の様々な話者を、当該ネットワークをトレーニングするために使用しうるということである。さらに、高信

40

50

頼のトレーニングを保証するためのトレーニング話者ごとに要求された最小数のトレーニング発声がなくともよい。さらに、ニューラルネットワークは、独立なパス内のニューラルネットワークを通じて当該発声のフレームを処理する必要なく、発声全体を特徴づけるデータを単一のパスでニューラルネットワークを通じて処理するように構成されてもよい。上で参照した態様は、発声がコンピューティングデバイスの認証されたユーザにより話された可能性があったかどうかを評価することに寄与しうる。当該態様はかかる評価を特に高信頼としうる。これを、少なくとも、コンピューティングデバイスの認証されたユーザのような、発声が特定の人により話されたかどうかを評価する際のエラーを減らすことで、より有効な音声認識システムとすることができる。当該高信頼性は、上で参照した態様に関連付けられた広範囲のシステムのセキュリティを増大させうる。例えば、当該態様が認証されたユーザを認識するために使用され、応答して、コンピューティングデバイスの状態をロックされた状態からロックされていない状態に変更する場合、上で参照したエラーの減少は、コンピューティングデバイスのアンロックをより安全にする。実際、当該システムの高信頼性は、コンピューティングデバイスを、非認証されたユーザからの詐欺的なアンロック試行に対してあまり脆弱でなくしうる。当該エラーの減少はまた、エラー訂正の必要性を削減でき、これはコンピューティングリソースをエラー訂正に割り当てる必要性を削減できる。これは、コンピュータリソースがより制限されうるモバイルコンピューティングデバイスにおいて特に有利である。

【図面の簡単な説明】

【0032】

【図1】(i)ニューラルネットワークをトレーニングし、(ii)ユーザをコンピューティングデバイスで加入させ、(iii)ユーザの音声の区別的な特徴に基づいてコンピューティングデバイスのユーザの発声を検証する動作を実行する、例示的なクライアントデバイスおよびコンピューティングシステムの略図である。

【図2】話者検証タスクで使用するための話者表現を決定するためにニューラルネットワークをトレーニングするためのブロック図である。

【図3】音声の区別的な特性を示す話者表現を決定するためにニューラルネットワークをトレーニングするための例示的なプロセスの流れ図である。

【図4A】発声プール内の異なる話者に対する発声のグループからのトレーニング発声のサブセットの例示的な選択を示す概念図である。

【図4B】ニューラルネットワークをトレーニングするための発声プールからのトレーニングデータのバッチの例示的な選択を示す概念図である。

【図5A】発声の少なくとも一部を特徴づけるデータを処理し、当該発声の当該少なくとも一部を特徴づけるデータに基づいて話者表現を生成するように構成された例示的なディープ・ニューラルネットワークのブロック図である。

【図5B】話者検証モデルで使用するために構成される長短期メモリレイヤを有する例示的な再帰型ニューラルネットワークのブロック図である。

【図6】コンピューティングデバイス上のニューラルネットワークを用いて、ユーザの発声から決定されたユーザの音声の特性に基づいてユーザのアイデンティティを検証するための例示的なプロセスの流れ図である。

【図7】本明細書で説明するコンピュータ実行型の方法および他の技術を実行する際に使用できるコンピューティングデバイスおよびモバイルコンピューティングデバイスの1例を示す図である。

【発明を実施するための形態】

【0033】

図1は、話者検証モデルに対するニューラルネットワークをトレーニングし、当該モデルを用いて話者検証のプロセスを実行するための例示的なシステム100の略図である。一般に、話者検証とは、当該話者の1つまたは複数の発声から決定される当該話者の音声の特性に基づいて話者のアイデンティティ要求を受理または拒否するタスクである。図1に示すように、話者検証は一般に3つの段階、即ち(i)当該話者検証モデルに対する二

10

20

30

40

50

ユーラルネットワークのトレーニング、( i i ) 新たな話者の加入、および ( i i i ) 当該加入した話者の検証を含むことができる。

【 0 0 3 4 】

システム 1 0 0 はクライアントデバイス 1 1 0、コンピューティングシステム 1 2 0、およびネットワーク 1 3 0を含む。幾つかの実装では、コンピューティングシステム 1 2 0 は、トレーニングされたニューラルネットワーク 1 4 0に基づいて話者検証モデル 1 4 4 をクライアントデバイス 1 1 0 に提供してもよい。幾つかの実装では、話者検証モデル 1 4 4 は、例えば、オペレーティング・システムまたはアプリケーションのコンポーネントとしてクライアントデバイス 1 1 0 に事前インストールされてもよい。他の実装では、話者検証モデル 1 4 4 がネットワーク 1 3 0 上で受信されてもよい。クライアントデバイス 1 1 0 は、話者検証モデル 1 4 4 を使用してユーザ 1 0 2 を当該話者検証プロセスに加入させてもよい。後の時点でユーザ 1 0 2 のアイデンティティを検証する必要があるとき、クライアントデバイス 1 1 0 は、話者検証モデル 1 4 4 を用いてユーザ 1 0 2 のアイデンティティを検証するためにユーザ 1 0 2 の会話発声を受信してもよい。話者検証モデル 1 4 4 をクライアントデバイス 1 1 0 にローカルに格納してもよいので、クライアントデバイス 1 1 0 はネットワーク 1 3 0 上の通信なしに話者検証判定を行うことができる。

10

【 0 0 3 5 】

図 1 には示していないが、幾つかの実装では、コンピューティングシステム 1 2 0 は、クライアントデバイス 1 1 0 に格納されるニューラルネットワーク 1 4 0 ではなくまたはそれに加えて、当該トレーニングされたニューラルネットワーク 1 4 0 に基づいて話者検証モデル 1 4 4 を格納してもよい。これらの実装では、クライアントデバイス 1 1 0 は、ネットワーク 1 3 0 を介してコンピューティングシステム 1 2 0 と通信して話者検証モデル 1 4 4 にリモートにアクセスし、それをユーザ 1 0 2 の加入のために使用してもよい。後の時点でユーザ 1 0 2 のアイデンティティを検証する必要があるとき、クライアントデバイス 1 1 0 はユーザ 1 0 2 の会話発声を受信してもよく、ネットワーク 1 3 0 を介してコンピューティングシステム 1 2 0 と通信して、当該リモートに配置された話者検証モデル 1 4 4 を用いてユーザ 1 0 2 のアイデンティティを検証してもよい。コンピューティングシステム 1 2 0 およびコンピューティング・デバイス 1 1 0 は、互いと異なってもよく、物理的に別々であってもよい。

20

【 0 0 3 6 】

システム 1 0 0 では、クライアントデバイス 1 1 0 は、例えば、デスクトップコンピュータ、ラップトップコンピュータ、タブレットコンピュータ、時計、ウェアブルコンピュータ、セルラ電話、スマートフォン、音楽プレイヤー、eブックリーダー、ナビゲーションシステム、またはユーザが対話しうる任意の他の適切なコンピューティングデバイスであることができる。幾つかの実装では、クライアントデバイス 1 1 0 はモバイルコンピューティングデバイスであってもよい。コンピューティングシステム 1 2 0 は 1 つまたは複数のコンピュータを含むことができ、当該コンピュータのうち独立なものに対して機能を実施してもよく、または当該機能は、複数のコンピュータにわたって実施するために分散されてもよい。ネットワーク 1 3 0 は、有線またはワイヤレスまたはその両方の組合せであることができ、インターネットを含むことができる。

30

40

【 0 0 3 7 】

幾つかの実装では、ユーザの電話のようなクライアントデバイス 1 1 0 は、話者検証モデル 1 4 4 をクライアントデバイス 1 1 0 にローカルに格納して、クライアントデバイス 1 1 0 が、加入または検証処理のためにリモートサーバ(例えば、コンピューティングシステム 1 2 0)でのモデルに頼ることなくユーザのアイデンティティを検証できるようにしてもよく、したがって通信帯域幅および時間を節約することができる。さらに、幾つかの実装では、1 つまたは複数の新たなユーザが加入するとき、本明細書で説明した話者検証モデル 1 4 4 は当該新たなユーザを用いた話者検証モデル 1 4 4 の任意の再トレーニングを必要とせず、これはまた、計算的に効率的でありうる。他の実装では、ニューラルネットワーク(およびしたがって当該話者検証モデル)を新たに収集されたトレーニングデ

50

ータの使用に基づいて定常的に更新できるように、加入、検証、またはその両方のために提供された所与のユーザの発声をコンピューティングシステム 120 に提供し、当該トレーニングデータに追加してもよい。

#### 【0038】

当該トレーニングされたニューラルネットワーク 140 を含む話者検証モデル 144 のサイズは、クライアントデバイス 110 の記憶空間およびメモリ空間が制限されるので、コンパクトでありうるのが望ましい。後述のように、話者検証モデル 144 はトレーニングされたニューラルネットワーク 140 に基づく。話者検証モデル 144 は、発声を特徴付けるデータに基づいて、当該発声の話者の音声の区分的な特徴を示す話者表現を生成するためのニューラルネットワーク 140 を含んでもよい。話者検証モデル 144 は、当該発声の話者のアイデンティティ要求を検証できるように、当該話者表現を処理し、当該発声の話者の音声が入力したユーザの音声と十分に同様であるかどうかを判定するためのさらなるコンポーネントを含んでもよい。

10

#### 【0039】

幾つかの実装では、ニューラルネットワークは大規模な 1 組のトレーニングデータを用いてトレーニングされてもよい。様々な技術を当該トレーニングデータの前処理の間に、トレーニング自体の間に、または後トレーニング段階の間に適用して、ニューラルネットワークのサイズを強制および/または削減し、コンパクトなモデルサイズを実現してもよい。例えば、話者検証モデル 144 を、ニューラルネットワーク 140 の特定のレイヤのみを選択することで構築してもよい。これは、クライアントデバイス 110 に格納するのに適したコンパクトな話者検証モデルをもたらす。当該話者モデルに対する話者表現を生成する際にソフトマックスまたは他の分類レイヤなしに、加入を実施してもよい。

20

#### 【0040】

図 1 はまた、段階 (A) 乃至 (F) で示す、データの例示的なフローを示す。段階 (A) 乃至 (F) を、当該示されたシーケンスで行ってもよく、または、当該示されたシーケンスと異なるシーケンスで行ってもよい。幾つかの実装では、段階 (A) 乃至 (F) の 1 つまたは複数をオフラインで行ってもよい。コンピューティングシステム 120 は、クライアントデバイス 110 がネットワーク 130 に接続されないときに計算を実施してもよい。段階 (A) および (B) は一般に上で参照した当該トレーニング段階の間に発生する。段階 (D) は一般に加入段階の間に発生する。段階 (E) 乃至 (G) は一般に検証段階の間に発生する。

30

#### 【0041】

段階 (A) で、コンピューティングシステム 120 は、ニューラルネットワーク 140 の監視されたトレーニングをニューラルネットワーク 140 に提供するためのトレーニング発声のサンプルを選択する。幾つかの実装では、トレーニング・サンプル 122 における発声がそれぞれ、多数の異なるトレーニング話者により話された 1 つまたは複数の所定の単語から構成されてもよい。当該発声は以前に記録されておりコンピューティングシステム 120 により使用するためにアクセス可能とされている。各トレーニング話者は所定の発声をコンピューティングデバイスに話してもよく、コンピューティングデバイスは当該発声を含むオーディオ信号を記録してもよい。例えば、各トレーニング話者を、当該トレーニングフレーズ「Hello Phone」を話すように促してもよい。幾つかの実装では、各トレーニング話者を、同一のトレーニングフレーズ複数回を話すように促してもよい。各トレーニング話者の当該記録されたオーディオ信号がコンピューティングシステム 120 に送信されてもよく、コンピューティングシステム 120 は当該記録されたオーディオ信号を多数の異なるコンピューティングデバイスおよび多数の異なるトレーニング話者から収集してもよい。幾つかの実装では、ユーザのアイデンティティが当該予め定義されたトレーニングフレーズの発声から決定されたユーザの音声の特性に基づいて検証されてもよいという点で、ニューラルネットワーク 140 をテキスト依存の話者検証に対して最適化してもよい。かかる実装では、ニューラルネットワーク 140 は、全ての、または実質的に全ての、当該予め定義されたトレーニングフレーズを含む発声でトレーニングされても

40

50

よい。他の実装では、ニューラルネットワーク 140 は、ユーザのアイデンティティが多種多様な単語またはフレーズの発声から決定されたユーザの音声の特性に基づいて検証されてもよいという点で、テキスト独立な話者検証を可能とするためにトレーニングされてもよく、当該多種多様な単語またはフレーズの発声は予め定義されていなくてもよい。例えば、ユーザは、自分のアイデンティティを検証するために、どの単語またはフレーズを彼または彼女が話したいかを独立に判定でき、当該トレーニングされたニューラルネットワーク 140 に基づく話者検証モデルは次いで、当該話された単語またはフレーズが与えられた場合にユーザを認証することができる。テキスト独立な話者検証を可能とするために、ニューラルネットワーク 140 は、多数の異なるトレーニング話者により話された多種多様な単語またはフレーズの発声でトレーニングされてもよい。

10

**【0042】**

段階 (B) で、ニューラルネットワーク 140 は、クライアントデバイスでユーザの加入および検証と平行する方式でトレーニングされてもよい。したがって、コンピューティングシステム 120 は、各トレーニング・サンプル 122 において、1組のシミュレートされた加入発声 122b およびシミュレートされた検証発声 122a を選択することができる。シミュレートされた加入発声 122b は、シミュレートされた話者モデルをトレーニング・サンプル 122 ごとに決定できるように、同一のトレーニング話者の全ての発声であってもよい。シミュレートされた検証発声 122a は、シミュレートされた加入発声 122b の話者と同一の話者の発声であってもよく、または異なる話者の発声であってもよい。トレーニング・サンプル 122 を次いでニューラルネットワーク 140 に提供でき、シミュレートされた検証発声 122a がシミュレートされた加入発声 122b の話者と同一の話者により話されたか、またはシミュレートされた加入発声 122b の話者と異なる話者により話されたかどうかに関して、ニューラルネットワーク 140 の出力に基づいて分類を行ってもよい。ニューラルネットワーク 140 は次いで、当該話者判定が正しかったかどうかに基づいて更新することができる。幾つかの実装では、各トレーニング・サンプル 122 を、2つのクラス、即ち、(当該シミュレートされた検証発声および加入発声の話者が同一であるサンプルに関する) マッチング話者クラス 141a および (当該シミュレートされた検証発声および加入発声の話者が異なるサンプルに関する) 非マッチング話者クラス 141b のうち1つに属するとしてラベル付けしてもよい。これらのラベルは、同一の話者が発声 122a および発声 122b を話したかどうかのグランド・トゥースを示す。トレーニング・サンプル 122 の分類の正確性を当該サンプルのラベルに基づいて決定することができる。幾つかの実装では、ニューラルネットワークに対する調節は、入力サンプルの分類の正確性に厳密に基づかなくてもよいが、一般にシミュレートされた検証発声 122a およびシミュレートされた加入発声 122b に対するニューラルネットワークにより生成された話者表現の比較から決定された1つまたは複数のメトリックに基づいてもよい。幾つかの実装では、トレーニング・サンプル 122 をトレーニングデータのリポジトリから選択してもよい。当該トレーニングデータのリポジトリを、発声プール 121 に編成してもよい。発声プール 121 の各々は当該発声のトレーニング話者によりグループ化されたトレーニング発声を含んでもよい。

20

30

**【0043】**

ニューラルネットワーク 140 は、トレーニング・サンプル 122 における発声に関する情報を入力するための入力レイヤ、およびサンプル 122 を処理するための幾つかの隠蔽されたレイヤを含んでもよい。当該トレーニングされたニューラルネットワーク 140 が、サンプル 122 のシミュレートされた検証発声および加入発声のうちマッチングまたは非マッチング話者の何れかを有するとのトレーニング・サンプル 122 の所望の分類を生成させる出力を話者検証モデル 144 が生成するように、1つまたは複数の隠蔽されたレイヤの重みまたは他のパラメータを調節してもよい。幾つかの実装では、ニューラルネットワーク 140 のパラメータをコンピューティングシステム 120 により自動的に調節してもよい。幾つかの他の実装では、ニューラルネットワーク 140 のパラメータをコンピューティングシステム 120 のオペレータにより手動で調節してもよい。ニューラルネ

40

50

ットワークのトレーニング段階を、例えば図 2、3、4 A - B、および 5 A - B の説明において以下でより詳細に説明する。

【0044】

段階(C)で、ニューラルネットワーク140がトレーニングされると、当該トレーニングされたニューラルネットワーク140に基づく話者検証モデル144が、例えば、ネットワーク130を通じてコンピューティングシステム120からクライアントデバイス110に送信される。幾つかの実装では、当該トレーニングされたニューラルネットワーク140、またはその一部は、話者検証モデル144のコンポーネントであってもよい。話者検証モデル144を、1つまたは複数のユーザ102の発声から決定されたユーザの音声の特性に基づいてユーザ102のアイデンティティを検証するように構成することができる。モデル144を、ユーザの音声の区別的な特徴を示すユーザ102に対する話者表現を生成するために、ユーザ102の発声を特徴づけるデータを当該トレーニングされたニューラルネットワーク140への入力として提供するように構成してもよい。当該話者表現を次いで、以前に決定されたユーザの音声のモデルと比較することができる。当該話者表現がユーザの話者モデルと十分に同様である場合、話者検証モデル144は、ユーザ102のアイデンティティが正当であるという指示を出力することができる。対照的に、当該話者表現がユーザの話者モデルと十分に同様でない場合、話者検証モデル144は、ユーザ102のアイデンティティが無効である(検証されない)という指示を出力することができる。

10

【0045】

段階(D)で、クライアントデバイス110で自分の音声を加入させたいと望むユーザ102が、加入段階で1つまたは複数の加入発声152をクライアントデバイス110に提供する。一般に、ユーザ102は、音声ニューラルネットワーク140のトレーニングで使用されなかったトレーニング話者の1つではない。幾つかの実装では、クライアントデバイス110は、トレーニング・サンプル122の発声において当該1組のトレーニング話者により話された同一のフレーズである加入フレーズを話すことをユーザ102に促してもよい。幾つかの実装では、クライアントデバイス110は、加入フレーズを何回か話すようにユーザに促してもよく、当該話された加入発声に対するオーディオ信号を加入発声152として記録してもよい。

20

【0046】

クライアントデバイス110は、加入発声152を使用して、クライアントデバイス110の話者検証システムにおいてユーザ102を加入させる。一般に、ユーザ102の加入はニューラルネットワーク140の再トレーニングなしに行われる。同一の話者検証モデル144の夫々のインスタンスは、多数の異なる話者を加入させるために、ニューラルネットワーク140における重み値または他のパラメータを変更するのを必要とせず、多数の異なるクライアントデバイスで使用されてもよい。話者検証モデル144を、再トレーニングニューラルネットワーク140なしに任意のユーザを加入させるために使用できるので、制限された処理要件でクライアントデバイス110で加入を実施してもよい。

30

【0047】

幾つかの実装では、加入発声152に関する情報は話者検証モデル144へ入力され、話者検証モデル144は、ユーザ102に対応する基準ベクトルまたは他の1組の値を出力してもよい。当該基準ベクトルまたは他の1組の値は、ユーザの音声の区別的な特徴を特徴付ける話者モデルを構成してもよい。当該話者モデルをクライアントデバイス110に、またはクライアントデバイス110から離れたコンピューティングシステムに格納してもよく、その結果、後にクライアントデバイス110により受信された発声に基づいて生成された話者表現を当該話者モデルと比較して、後に受信された発声の各話者がユーザ102であるかまたは他の話者であるかどうかを検証してもよい。

40

【0048】

段階(E)で、ユーザ102が、音声認証を用いてクライアントデバイス110へのアクセスを得ようと試みる。ユーザ102は、検証段階で検証発声154をクライアントデ

50

バイス 1 1 0 に提供する。幾つかの実装では、検証発声 1 5 4 は、加入発声 1 5 2 と同一の、話されたフレーズの発声である。検証発声 1 5 4 は話者検証モデル 1 4 4 への入力として使用される。

【 0 0 4 9 】

段階 ( F ) で、クライアントデバイス 1 1 0 は、ユーザの音声が入力したユーザの当該音声にマッチするかどうかを判定する。幾つかの実装では、ニューラルネットワーク 1 4 0 は検証発声 1 5 4 を特徴付けるデータを処理してもよく、検証発声 1 5 4 に基づいてユーザ 1 0 2 に対する話者表現を出力してもよい。幾つかの実装では、クライアントデバイス 1 1 0 は、ユーザ 1 0 2 に対する話者表現を入力したユーザに対する当該話者モデルと比較して、検証発声 1 5 4 が入力したユーザにより話されたかどうかを判定してもよい。ニューラルネットワークの検証段階を、例えば図 6 に関して以下でより詳細に説明する。

10

【 0 0 5 0 】

段階 ( G ) で、クライアントデバイス 1 1 0 が、検証結果 1 5 6 を表す指示をユーザ 1 0 2 に提供する。幾つかの実装では、クライアントデバイス 1 1 0 がユーザ 1 0 2 のアイデンティティを受理した場合、クライアントデバイス 1 1 0 は、検証が成功したという視覚的指示またはオーディオ指示をユーザ 1 0 2 に送信してもよい。幾つかの他の実装では、クライアントデバイス 1 1 0 がユーザ 1 0 2 のアイデンティティを受理した場合、クライアントデバイス 1 1 0 はユーザ 1 0 2 に次の入力を促してもよい。例えば、クライアントデバイス 1 1 0 は、メッセージ「Device enabled. Please enter your search」を当該ディスプレイに出力してもよい。幾つかの他の実装では、クライアントデバイス 1 1 0 がユーザ 1 0 2 のアイデンティティを受理した場合、クライアントデバイス 1 1 0 は、さらなる入力をユーザ 1 0 2 から待つことなく後続のアクションを実施してもよい。例えば、ユーザ 1 0 2 は検証段階の間に、クライアントデバイス 1 1 0 に対して「Hello Phone , search the nearest coffee shop」を話してもよい。クライアントデバイス 1 1 0 は、検証フレーズ「Hello Phone」を用いてユーザ 1 0 2 のアイデンティティを検証してもよい。ユーザ 1 0 2 のアイデンティティが受理された場合、クライアントデバイス 1 1 0 は、ユーザ 1 0 2 にさらなる入力を求めることなく、最も近いコーヒー店の検索を実施してもよい。一般に、幾つかの実装では、クライアントデバイス 1 1 0 がユーザ 1 0 2 のアイデンティティを受理した場合、クライアントデバイス 1 1 0 はロックされた状態からロックされていない状態に遷移することで応答してもよい。当該ロックされた状態において、クライアントデバイス 1 1 0 の 1 つまたは複数の能力は無効化またはブロックであり、当該ロックされていない状態では、当該能力は有効であり、または、ユーザ 1 0 2 がアクセスするのに利用可能とされる。同様に、クライアントデバイス 1 1 0 は、成功した検証に応答してより完全に特徴付けられた状態に「活性化」または遷移してもよい。

20

30

【 0 0 5 1 】

幾つかの実装では、クライアントデバイス 1 1 0 がユーザ 1 0 2 のアイデンティティを拒絶した場合、クライアントデバイス 1 1 0 は、検証が拒絶されたという視覚的指示またはオーディオ指示をユーザ 1 0 2 に送信してもよい。幾つかの実装では、クライアントデバイス 1 1 0 がユーザ 1 0 2 のアイデンティティを拒絶した場合、クライアントデバイス 1 1 0 はユーザ 1 0 2 に別の発声試行を促してもよい。幾つかの実装では、試行の数が閾値を超過した場合、クライアントデバイス 1 1 0 は、ユーザ 1 0 2 が自分のアイデンティティをさらに検証するのを試みるのをブロックしてもよい。

40

【 0 0 5 2 】

図 2 を参照すると、ニューラルネットワーク 2 0 6 をトレーニングするための例示的なシステム 2 0 0 のブロック図が示されている。図 2 により示されたトレーニング段階が完了すると、当該トレーニングされたニューラルネットワーク 2 0 6 は、話者の発声を特徴付けるデータを処理し、当該話者の音声の区別的な特徴を示す当該話者に対する話者表現を生成することができる。当該話者表現はついで、加入段階の間に当該話者に対する話者モデルを生成するか、または検証段階の間に当該話者のアイデンティティを検証する際に話者検証モデルにより使用されてもよい。

50

## 【 0 0 5 3 】

一般に、図 2 は、ニューラルネットワーク 2 0 6 が、後に話者検証タスクを実施するクライアントデバイスで発生する加入および検証段階と平行する方式でトレーニングされてもよいことを示す。トレーニング発声を有限数の話者から当該話者の各々に対する対応するクラスに分類するためにニューラルネットワーク 2 0 6 をトレーニングする幾つかのアプローチとは異なり、図 2 のニューラルネットワーク 2 0 6 は、所与の発声の特定の話者を決定するためにトレーニングされない。その代わりに、ニューラルネットワーク 2 0 6 は、当該発声の何れかを特定の話者アイデンティティと必ずしもマッチすることなく、区別的かつ所与の発声の話者が別の 1 組の発声の話者と同一であるか否かを判定するために使用可能である話者表現を生成するためにトレーニングされる。このように、トレーニングの間に最適化された損失関数は、検証段階の間に当該話者検証モデルにより利用される同一の関数である。換言すれば、検証の間に、検証発声に基づく話者表現は加入したユーザに対する話者モデルと比較される。当該話者表現が当該話者モデルと十分に同様である場合、検証発声の話者のアイデンティティが検証される。図 2 に示す当該アプローチはトレーニングの間に同様な技術を使用する。即ち、シミュレートされた話者モデル 2 1 4 が 1 つまたは複数の加入発声に対する話者表現に基づいて生成され、話者表現 2 0 8 はまた、シミュレートされた検証発声 2 0 2 に対して生成される。ニューラルネットワーク 2 0 6 の重み値および他のパラメータは、シミュレートされた検証発声 2 0 2 をシミュレートされた加入発声 2 0 4 a 乃至 n と同一または異なる話者により話されているとして分類する際のエラーを最小化するために、トレーニングの間に調節される。

10

20

## 【 0 0 5 4 】

図 2 は、シミュレートされた検証発声 2 0 2 を特徴づけるデータおよび 1 つまたは複数のシミュレートされた加入発声 2 0 4 a 乃至 n を特徴づけるデータを含むトレーニングデータのサンプルに基づく単一のトレーニングの反復の前方パスを示す。実際に、ニューラルネットワーク 2 0 6 は、多数の反復およびトレーニングデータの多数の異なるサンプルにわたってトレーニングされる。各反復により、ニューラルネットワーク 2 0 6 を、当該各反復に対するトレーニングデータの対応するサンプルを処理した結果に基づいて調節してもよい。図 4 A および 4 B は、さらに以下で説明するように、シミュレートされた検証発声 2 0 2 およびシミュレートされた加入発声 2 0 4 a 乃至 n が選択されうる例示的な技術を示す。特定のサンプルに対するシミュレートされた加入発声 2 0 4 a 乃至 n は一般に、同一のトレーニング話者により話された全ての発声である。シミュレートされた加入発声 2 0 4 a 乃至 n の話者は異なるトレーニングの反復に対するトレーニングデータの異なるサンプルの間で異なってもよいが、所与のトレーニング反復に対する所与のトレーニング・サンプルにおいて、加入発声 2 0 4 a 乃至 n の全ては一般に同一のトレーニング話者により話される。シミュレートされた検証発声 2 0 2 はシミュレートされた加入発声 2 0 4 a 乃至 n の話者と同一のトレーニング話者により話されているかもしれず、または、シミュレートされた加入発声 2 0 4 a 乃至 n の話者と異なるトレーニング話者により話されているかもしれない。当該話者がシミュレートされた検証発声 2 0 2 とシミュレートされた加入発声 2 0 4 a 乃至 n の両方の間で同一であるトレーニングデータのサンプルに対して、当該サンプルを「マッチング」サンプルとしてラベル付けしてもよい。当該話者がシミュレートされた検証発声 2 0 2 およびシミュレートされた加入発声 2 0 4 a 乃至 n の間で異なるトレーニングデータのサンプルに対して、当該サンプルを「非マッチング」サンプルとしてラベル付けしてもよい。当該ラベルは当該トレーニング・サンプルの真の分類を表してもよく、トレーニングの前の前処理段階で決定してもよい。幾つかの実装では、ニューラルネットワーク 2 0 6 の出力に基づくトレーニング・サンプルの当該推定された分類は、当該トレーニング・サンプルに対するラベルにより示される真の分類と比較して、ニューラルネットワーク 2 0 6 を調節するかどうかを判定してもよい。

30

40

## 【 0 0 5 5 】

幾つかの実装では、当該トレーニング・サンプルにおけるデータは、当該シミュレートされた検証発声および加入発声 2 0 2 、 2 0 4 a 乃至 n に対する生のオーディオ信号はで

50

なくてもよい。その代わりに、発声 202、204 a 乃至 n が、ニューラルネットワーク 206 により処理するための適切なフォーマットに処理および変換されていてもよい。例えば、当該トレーニング・サンプルにおけるデータは、生のオーディオ信号自体ではなく、当該シミュレートされた検証発声および加入発声 202、204 a 乃至 n の各特徴を特徴付けてもよい。幾つかの実装では、当該トレーニング・サンプル内のシミュレートされた発声 202、204 a 乃至 n の各々を表すデータは各発声に対する 1 つまたは複数のログフィルタバンクを含んでもよい。幾つかの実装では、各発声を当該発声に対する複数のフレームに分割してもよく、別々のログフィルタバンクを当該発声のフレームごとに生成することができる。例えば、当該発声の各フレームは例えば 40 個のログフィルタバンクにより表されてもよい。

10

#### 【0056】

幾つかの実装では、シミュレートされた検証発声 202 を特徴づけるデータおよびシミュレートされた加入発声 204 a 乃至 n の各々を特徴づけるデータをニューラルネットワーク 206 を通じて一度（即ち、単一のパスで）処理することができる。したがって、所与の発声に対するトレーニングデータが各々各 1 組のログフィルタバンクにより表された複数のフレームに分割されているが、当該発声の全体に対するフレームの全てを特徴づけるデータを、ニューラルネットワークを通じて単一のパスで処理するために、ニューラルネットワーク 206 に（例えば、それぞれ 40 個のログフィルタバンクを有する 80 個のフレームに対する  $80 \times 40$  個の特徴ベクトルとして）入力することができる。これは、当該フレームを別々に処理するために、当該発声のフレームごとにデータをニューラルネットワーク 206 に独立して入力するのは対照的である。他の実装では、各発声 202、204 a 乃至 n の全体を特徴づけるデータを単一のパスでニューラルネットワーク 206 を通じて処理するためにニューラルネットワーク 206 をトレーニングするのではなく、発声 202、204 a 乃至 n の独立なフレームを特徴づけるデータを、ニューラルネットワーク 206 への入力として提供することができる。

20

#### 【0057】

幾つかの実装では、シミュレートされた検証発声および加入発声 202、204 a 乃至 n を 1 つまたは複数の追加の技術に従って事前に処理してもよい。例えば、ニューラルネットワーク 206 の構造は、トレーニング発声全体が固定された長さ（例えば、0.8 秒のオーディオ）を有するのを要求してもよい。少なくとも幾つかの発声 202、204 a 乃至 n はしたがって、長い発声を均一な長さに刈り込み、かつ/または幾つかの短い発声をパディングして長いクリップを作成した結果であってもよい。他の実装では、しかし、ニューラルネットワーク 206 は可変長さ発声を処理できてもよく、この場合、トレーニングデータ内の発声 202、204 a 乃至 n を固定された長さに刈り込むかまたはパディングしてもよい。発声 202、204 a 乃至 n に対するオーディオはまた均一化されていてもよく、雑音の存在において堅牢に実施されるようにニューラルネットワークがトレーニングされることを保証するために、トレーニング発声 202、204 a 乃至 n に雑音が追加されているかまたはそこから雑音が除去されていてもよい。

30

#### 【0058】

点線ボックス 215 内のシステム 200 の部分は、複数のシミュレートされた加入発声 204 a 乃至 n を特徴づけるデータがシミュレートされた加入発声 204 a 乃至 n の特定のトレーニング話者に対するシミュレートされた話者モデル 214 を生成するために使用されるという点で、話者検証プロセスの加入段階をシミュレートする。シミュレートされた加入発声 204 a 乃至 n の各々を特徴づけるそれぞれのデータはニューラルネットワーク 206 ニューラルネットワーク 206 の入力レイヤに別々に入力される。ニューラルネットワーク 206 は 1 つまたは複数の隠蔽されたレイヤを通じて当該データを処理し、シミュレートされた加入発声 204 a 乃至 n の各々に対する各話者表現 210 a 乃至 n を生成する。例えば、図 2 に示すように、話者表現 1 (210 a) はシミュレートされた加入発声 1 に基づいてニューラルネットワーク 206 により生成される (204 a)。同様に、話者表現 2 (210 b) はシミュレートされた加入発声 2 に基づいてニューラルネット

40

50

ワーク 206 により生成される (204b)。話者表現はしたがって、シミュレートされた加入発声 204a 乃至 n の各々に対するニューラルネットワーク 206 により生成されてもよい。幾つかの実装では、話者表現 210a 乃至 n を、ニューラルネットワーク 206 を通じてシミュレートされた加入発声 204a 乃至 n の各々を逐次的に処理することで生成してもよい。幾つかの実装では、話者表現 210a 乃至 n を、発声 204a 乃至 n を特徴付けるデータを、シミュレートされた加入発声 204a 乃至 n の各々に対するニューラルネットワーク 206 の各インスタンスと並列に処理することで並列に生成することができる。話者表現 210a 乃至 n は一般にそれぞれ、シミュレートされた加入発声 204a 乃至 n の対応する 1 つに基づいてニューラルネットワーク 206 により決定された、当該シミュレートされた加入トレーニング話者の音声の区別的な特性を表す値の集合を含む。幾つかの実装では、話者表現 210a 乃至 n はニューラルネットワーク 206 の最後の隠蔽されたレイヤの重み値または他のパラメータを示してもよい。幾つかの実装では、話者表現 210a 乃至 n は、ニューラルネットワーク 206 がソフトマックス出力レイヤなしに構成されたときの、ニューラルネットワーク 206 の最終的な出力であってもよい。

10

20

30

40

50

**【0059】**

シミュレートされた話者モデル 214 を生成するために、話者表現 210a 乃至 n を図 2 のボックス 212 に示すように平均化することができる。したがって、シミュレートされた話者モデル 214 は、シミュレートされた加入発声 204a 乃至 n のトレーニング話者の音声の区別的な特性を表す値の集合を定義してもよい。シミュレートされた話者モデル 214 を決定するために複数の話者表現 210a 乃至 n を平均化することで、当該異なるシミュレートされた加入発声 204a 乃至 n の間の話者の音声の変形を平滑化することができる。シミュレートされた話者モデル 214 は、独立な話者表現 210a 乃至 n の何れかよりも高信頼な話者の音声の表現であってもよい。これは、所与のシミュレートされた加入発声 204a 乃至 n の特質を独立に反映してもよい。

**【0060】**

幾つかの実装では、トレーニングデータの各サンプル内のシミュレートされた加入発声 204a 乃至 n の総数はトレーニングの反復ごとに変化してもよい。例えば、第 1 のトレーニングの反復に対する第 1 のトレーニング・サンプルは 9 個のシミュレートされた加入発声 204a 乃至 n を含んでもよい。第 2 のトレーニングの反復に対する第 2 のトレーニング・サンプルは、しかし、4 個のシミュレートされた加入発声 204a 乃至 n のみを含んでもよい。他の実装では、トレーニングデータの各サンプル内のシミュレートされた加入発声 204a 乃至 n の総数はトレーニングの反復ごとに固定されていてもよい。例えば、ニューラルネットワーク 206 は、当該 1 組のトレーニングデータが反復ごとに全体で 5 個のシミュレートされた加入発声 204a 乃至 n を含む一連の反復でトレーニングされてもよい。幾つかの実装では、当該トレーニングの反復の 1 つ、一部または全部を、単一のシミュレートされた加入発声 204a 乃至 n のみを含むトレーニング・サンプルで実施してもよい。

**【0061】**

話者表現 210a 乃至 n が、シミュレートされた加入発声 204a 乃至 n を特徴付けるデータから生成されたのと同じ方式で、話者表現 208 を、シミュレートされた検証発声 202 を特徴付けるデータから生成することができる。当該シミュレートされた検証発声 202 を特徴付けるデータ (例えば、検証発声 202 の特徴を特徴づけるログフィルタバンク値) をニューラルネットワーク 206 の入力レイヤに提供することができる。ニューラルネットワーク 206 次いで、当該ネットワークの 1 つまたは複数の隠蔽されたレイヤを通じて入力を処理する。ニューラルネットワーク 206 の出力は、シミュレートされた検証発声 202 を話した話者の音声の区別的な特性を示す値の集合を定義する話者表現 208 である。

**【0062】**

ニューラルネットワーク 206 のトレーニングの間の検証段階とさらに並行するために、シミュレートされた検証発声 202 に基づく話者表現 208 は、例えば、検証段階の間

に話者検証モデルによりクライアントデバイスで行われるのと同じの方式でシミュレートされた話者モデル 214 と比較されることができる。幾つかの実装では、当該比較を、(1)シミュレートされた話者表現 208 に対して値の集合を定義する第 1 のベクトルおよび(2)シミュレートされた話者モデル 214 に対する値の集合を定義する第 2 のベクトルの余弦距離を(ブロック 216 に示すように)取得することにより実施することができる。ロジスティック回帰 218 を次いで当該距離に適用して、シミュレートされた検証発声 202 を話したトレーニング話者が、シミュレートされた加入発声 204 a 乃至 n を話したトレーニング話者と同じであるかまたは異なるかどうかを推定することができる。これは、マッチング話者クラスに対する第 1 のブロック 220 a、および非マッチング話者クラスに対する第 2 のブロック 220 b により図 2 で表されている。幾つかの実装では、ロジスティック回帰 218 と異なる分類技術を適用して、シミュレートされた検証発声 202 を話したトレーニング話者が、シミュレートされた加入発声 204 a 乃至 n を話したトレーニング話者と同じであるかまたは異なるかどうかに関する判定を行ってもよい。例えば、ヒンジレイヤまたはソフトマックスレイヤを幾つかの代替的な分類に対して使用してもよい。図 2 に示すような 2 つのクラスモデルでは、当該ソフトマックスおよびロジスティック回帰技術は同一のまたは同様な最適化機能を使用してもよい。

10

#### 【0063】

ニューラルネットワーク 206 の重み値または他のパラメータを次いで、ブロック 222 により表されるように、シミュレートされた検証発声 202 に対する話者表現 208 のシミュレートされた話者モデル 214 との比較の結果に基づいて、調節することができる。例えば、当該トレーニング・サンプルが真に非マッチング話者を有するとしてラベル付けされ、不正確に分類されたトレーニング・サンプルがマッチング話者を有するとして分類された場合、ニューラルネットワーク 206 はエラーを補正するように自動的に調節されてもよい。より一般に、ニューラルネットワーク 206 を最適化して、話者サンプルをマッチングするための類似性スコアを最大化するか、または、ロジスティック回帰によるスコア出力を最適化してもよく、ニューラルネットワーク 206 をまた最適化して、非マッチング話者サンプルに対する類似性スコアを最小化するか、またはロジスティック回帰によるスコア出力を最適化してもよい。幾つかの実装では、ニューラルネットワーク 206 に対する調節を、トレーニングの反復ごとに各トレーニング・サンプルの結果に回答して行うことができ、またはニューラルネットワーク 206 をトレーニングの反復の幾つかのみの結果に基づいて調節してもよい。幾つかの実装では、ニューラルネットワーク 206 を、話者表現 208 および非マッチング話者に対するシミュレートされた話者モデル 214 の間の距離(即ち、差異の最大化)を最大化し、話者表現 208 およびマッチング話者に対するシミュレートされた話者モデル 214 の間の距離を最小化(即ち、差異を最小化)するように調節してもよく。幾つかの実装では、トレーニング・サンプルをマッチング話者クラス 220 a または非マッチング話者クラス 220 b の何れかに属するとして分類するための硬判定をトレーニング段階の間に行わなくてもよいことに留意されたい。寧ろ、ニューラルネットワーク 206 を、ロジスティック回帰レイヤ 218 により出力される当該スコアを最適化するか、または、1 つまたは複数の他のメトリックを最適化する方式で調節してもよい。

20

30

40

#### 【0064】

次に図 3 を参照すると、話者検証モデルで使用されうるニューラルネットワークをトレーニングするための例示的なプロセス 300 の流れ図が示されている。幾つかの実装では、プロセス 300 を、図 1 からのコンピューティングシステム 120 および図 2 からのコンピューティングシステム 200 のような本明細書で説明するコンピューティングシステムにより実行してもよい。

#### 【0065】

プロセス 300 は段階 302 で開始し、第 1 の 1 組のトレーニングデータが選択される(即ち、第 1 のトレーニング・サンプル)。第 1 の 1 組のトレーニングデータは、シミュレートされた検証発声の特徴づけるデータおよび 1 つまたは複数のシミュレートされた加

50

入発声の特徴づけるデータを含むことができる。当該トレーニングセットにおける発声は、それらがトレーニング段階の間に話者検証の加入段階および検証段階を並行させるかまたは「シミュレート」する方式でトレーニングプロセスで使用されるという点で、「シミュレートされる」。しかし、当該発声自体は一般に、トレーニング話者により話された記録された会話の実際のスニペットである。当該トレーニング話者は一般に、当該話者検証プロセスの実際的な加入および検証段階の間に発声を提供した同一の話者ではない。下記でさらに説明する図 4 A および 4 B は、当該シミュレートされた検証発声および加入発声を選択するための例示的な技術を示す。

**【 0 0 6 6 】**

当該選択された 1 組のトレーニングデータ（即ち、選択されたサンプル）を、それがマッチング話者の会話または非マッチング話者に対するサンプルを表すかどうかに従ってラベル付けしてもよい。当該シミュレートされた検証発声の話者が当該シミュレートされた加入発声の話者と同じである場合、当該 1 組のトレーニングデータはマッチング話者サンプルとしてラベル付けされる。当該シミュレートされた検証発声の話者が当該シミュレートされた加入発声の話者と異なる場合、当該 1 組のトレーニングデータが非マッチング話者サンプルとしてラベル付けされる。幾つかの実装では、当該ラベルを、マッチングまたは非マッチングサンプルの何れかであるとして当該 1 組のトレーニングデータの推定された分類が正確であるか否かを判定するために、後にトレーニングプロセス 3 0 0 において使用することができる。

10

**【 0 0 6 7 】**

幾つかの実装では、選択された 1 組のトレーニングデータは、シミュレートされた検証発声および加入発声に対する生のオーディオ信号ではなく、その代わり当該発声の特徴を特徴づけるデータを含んでもよい。例えば、当該 1 組のトレーニングデータで表された各発声は、当該発声の固定長のフレームに対して決定された 1 組のログフィルタバンクにより特徴づけられることができる。当該発声のフレームごとの当該ログフィルタバンクをついで、ニューラルネットワークへの入力として提供され当該発声の全体を分類する単一の 1 組の入力値に連結してもよい。

20

**【 0 0 6 8 】**

プロセス 3 0 0 の段階 3 0 4 および 3 0 6 で、話者表現は、第 1 の 1 組のトレーニングデータで分類される発声の各々に対して決定される。当該話者表現はそれぞれ、当該各話者表現に対する対応する発声を話した当該トレーニング話者の音声の区別的な特徴を示す値の集合であることができる。例えば、第 1 の話者表現が当該シミュレートされた検証発声に基づいて生成されてもよく、各第 2 の話者表現が当該シミュレートされた加入発声の各々に基づいて生成されてもよい。当該話者表現を生成するために、発声の特徴づけるデータが、トレーニングされているニューラルネットワークの入力レイヤに提供される。ニューラルネットワークは次いで、当該ネットワークの 1 つまたは複数の隠蔽されたレイヤを通じて当該入力データを処理する。当該話者表現は次いでニューラルネットワークの出力である。幾つかの実装では、当該出力は、ソフトマックスレイヤではない出力レイヤで出力される。当該出力を提供する最終的なレイヤは完全に接続された線形レイヤであってもよい。幾つかの実装では、当該話者表現は、ソフトマックス出力レイヤの出力ではなく、ニューラルネットワークの最後の隠蔽されたレイヤで生成された値またはその活性化を含んでもよい。幾つかの実装では、ニューラルネットワークをソフトマックス出力レイヤなしで構成してもよい。

30

40

**【 0 0 6 9 】**

段階 3 0 8 で、当該シミュレートされた加入発声に対応する話者表現が、シミュレートされた話者モデルを生成するために結合される。当該シミュレートされた話者モデルは、当該シミュレートされた加入発声に対する話者表現の平均であることができる。当該話者表現を平均化することで、当該トレーニング話者の音声を特徴づける高信頼のモデルを決定することができる。例えば、当該話者が当該シミュレートされた加入発声の各々を話した方式の変形を、当該シミュレートされた検証発声に対する話者表現が比較される堅牢な

50

ベースラインで当該話者モデルを使用できるように、平滑化してもよい。幾つかの実装では、プロセス300は、当該シミュレートされた加入発声に対する話者表現のサブセットのみを選択して、当該シミュレートされた話者モデルを生成する際に結合してもよい。例えば、当該シミュレートされた加入発声の各々または当該対応するシミュレートされた加入発声の品質の測定値を決定してもよい。プロセス300はついで、当該シミュレートされた話者モデルを生成するために使用される1組の表現に含めるために、閾値品質スコアを満たすこれらの話者表現のみ、または対応する発声が閾値品質スコアを満たすこれらの話者表現を選択してもよい。

#### 【0070】

段階310で、シミュレートされた検証発声に対する話者表現がシミュレートされた話者モデルと比較される。幾つかの実装では、二進クラシファイアが、マッチング話者を表すかまたは表さないとしてデータサンプルを分類するために使用される。幾つかの実装では、当該比較は、当該シミュレートされた検証発声に対する話者表現および当該シミュレートされた話者モデルの間の類似性の測定値を決定するステップを含むことができる。例えば、類似性の測定値は、当該話者表現に対する値のベクトルと当該シミュレートされた話者モデルに対する値のベクトルの間の余弦距離であってもよい。類似性の測定値はついで、マッチング話者サンプルまたは非マッチング話者サンプルの何れかとして第1の1組のトレーニングデータの分類を推定するために使用されてもよい。例えば、類似性の測定値が十分に高い（例えば、閾値類似性スコアを満たす）場合、ロジスティック回帰を、当該1組のトレーニングデータをマッチング話者のクラスにマップするために使用してもよい。他方、類似性の測定値が低すぎる（例えば、閾値類似性スコアを満たさない）場合、ロジスティック回帰を、当該1組のトレーニングデータを非マッチング話者のクラスにマップするために使用してもよい。

#### 【0071】

次に、段階312で、ニューラルネットワークの1つまたは複数のパラメータを、シミュレートされた検証発声に対する話者表現およびシミュレートされた話者モデルの間の段階310での比較の結果に基づいて調節してもよい。例えば、当該トレーニングデータが非マッチング話者サンプルとしてラベル付けされた場合に、隠蔽されたレイヤにおける様々なノードの重み、またはニューラルネットワークの他のパラメータを調節して、当該話者表現および当該シミュレートされた話者モデルの間の距離を増大（類似性スコアを減少）させてもよい。さらに、当該トレーニングデータがマッチング話者サンプルとしてラベル化された場合に、ニューラルネットワークの重みまたは他のパラメータを調節して、当該話者表現および当該シミュレートされた話者モデルの間の距離を削減（当該類似性スコアを増大）させてもよい。一般に、トレーニングプロセス300の各反復は各加入段階および各検証段階をシミュレートすることを意図しているので、ニューラルネットワークを調節して、話者検証の間に実際の加入および検証段階に適用されるたものと同一の損失関数を最適化してもよい。このアプローチの1つの利益は、ニューラルネットワークが、より正確な検証話者のアイデンティティのための話者検証モデルにおいて使用できる話者表現をより良く生成するためにトレーニングされるということである。例えば、幾つかの実装では、ニューラルネットワークをトレーニングするときに考慮されない発声の実際の検証の間に追加の後処理ステップが行われず。これらの技術を、ニューラルネットワークをトレーニングするための「エンド・ツー・エンド」のアプローチと考えるもよい。

#### 【0072】

最後に、段階314で、次に1組のトレーニングデータが、トレーニングニューラルネットワークの別の反復に対して選択される。再度、この段階で選択された当該1組のトレーニングデータはシミュレートされた検証発声の特徴付けるデータおよび1つまたは複数のシミュレートされた加入発声の特徴付けるデータを含んでもよい。プロセス300はついで、段階304乃至312を反復し、追加のトレーニングの反復に対するトレーニングデータの追加のセットを限界に達するまで選択し続けてもよい。幾つかの実装では、当該限界が、利用可能なトレーニングデータの全てが期限切れになることから生じてもよい。

幾つかの実装では、プロセス300をターゲット性能レベルに到達するまで続けてもよい。例えば、何回かのトレーニングの反復の後、ニューラルネットワークを、トレーニングプロセス300の間に使用されなかった差し出された1組のデータに対してテストしてもよい。トレーニングを、当該差し出されたセット上の試験が、ニューラルネットワークが少なくともターゲット性能レベルを達成したことを示すまで、続けてもよい。

#### 【0073】

次に図4Aおよび4Bを参照すると、話者検証モデルに対するニューラルネットワークをトレーニングする際に使用するためのトレーニングデータのセットを選択するための例示的な技術を示す略図が示されている。幾つかの実装では、図4Aおよび4Bに関して説明された技術は、多数のトレーニングの反復にわたって選択されたトレーニング発声の多様性を保証することができる。これは、所与の数のトレーニング発声に対してニューラルネットワークをより良く実施することをもたらす。

10

#### 【0074】

幾つかの実装では、当該利用可能なトレーニング発声の全部のまたは一部を複数のグループ410a乃至nにクラスタ化してもよい。グループ410a乃至nを、トレーニング発声のグループの集合を含む発声プール408にさらに構成してもよい。当該トレーニング発声を幾つかの実装では話者によりグループ化してもよい。例えば、グループ410aは、全て第1の話者により話された複数の発声を含み、グループ410nは全て別の話者により話された複数の発声を含む。したがって、グループ410a乃至nの各々は異なる話者に対応してもよい。グループ410a乃至nは全て同一のトレーニング発声の数を含んでもよく、またはトレーニング発声の数はグループ410a乃至nの異なるもの間で変化してもよい。

20

#### 【0075】

トレーニングの反復ごとに、発声プール408にアクセスしてもよく、特定の発声を、各トレーニングの反復における入力として適用されるトレーニングデータのサンプルに対して選択してもよい。例えば、図4Aは、入力サンプル402としてトレーニングの反復のために発声プール408からランダムに選択された1つの1組のトレーニングデータを示す。第1の話者に対応する第1の発声グループは、当該シミュレートされた話者モデルを生成するために使用するための発声プール408内のグループ410a乃至nから選択することができる。当該グループをランダムにまたは別の方式で選択してもよい。当該選択されたグループ、例えば、図4Aのグループ410aから、当該第1の話者の発声のサブセットが、入力サンプル402内のシミュレートされた加入発声406として選択される。このサブセットは一般に複数の発声を含み、或るトレーニングの反復から別の反復へと、同一数または異なる数の発声を含んでもよい。当該選択されたグループ、例えば、グループ410aからの発声をランダムに選択してもよく、その結果、当該発声の異なる組合せが、異なるトレーニングの反復において第1の話者に対する異なるシミュレートされた話者モデルを生成するために使用される。

30

#### 【0076】

発声404もシミュレートされた検証発声として選択される。発声404は、当該トレーニングの反復が加入発声406とのマッチまたは非マッチの1例であるかどうか依存して、第1の話者または異なる話者の発声であってもよい。マッチングおよび非マッチングの例の両方がトレーニングで使用される。結果として、幾つかのトレーニングの反復に対して、発声404は、第1の話者の発声、例えば、グループ410aからの発声である。他のトレーニングの反復に対して、発声404は、図4Aに示すように、第1の話者と異なる第2の話者の発声であり、その結果入力サンプル402はシミュレートされた検証発声404とシミュレートされた加入発声406の間のマッチを表さない。

40

#### 【0077】

図4Aの例では、特定の発声が、シミュレートされた検証発声404として第2の発声グループ410nから選択される(例えば、ランダムに選択される)。幾つかの実装では、(発声404がそこから選択される)第2の発声グループを、発声プール408内のグ

50

グループ410 a乃至nからランダムに、または、グループ410 a乃至nの変化する選択のパターンに従って選択してもよい。他の実装では、当該シミュレートされた加入発声の話者と同一の話者からの別の発声が当該シミュレートされた検証発声として適用されるべきかどうかに関してランダムな選択を行ってもよい。したがって、おそらく当該ランダムな選択がバイアスされ、その結果、シミュレートされた検証発声404がシミュレートされた加入発声の話者406と同一の話者の発声であるという50パーセントの可能性が存在する。当該ランダムな選択の結果が、入力サンプル402がマッチング話者サンプルであるというものである場合、シミュレートされた検証発声404を、シミュレートされた加入発声406が選択された発声のグループと同一の発声のグループ410から選択することができる。しかし、当該ランダムな選択の結果が、入力サンプル402が非マッチング話者サンプルであるというものである場合、シミュレートされた検証発声404を、シミュレートされた加入発声406がそこから選択された発声のグループと異なる話者に対応する発声410の異なるグループから選択することができる。

#### 【0078】

一般に、図4Aにより示される選択技術は、話者の異なる組合せからの発声を異なるトレーニングの反復で適用させることができる。例えば、第1のトレーニングの反復において、当該シミュレートされた加入発声は第1の話者により話されているかもしれず、当該シミュレートされた検証発声はまた第1の話者により話されているかもしれない。第2のトレーニングの反復において、当該シミュレートされた加入発声は第2の話者により話されているかもしれず、当該シミュレートされた検証発声は第3の話者により話されているかもしれない。次に第3のトレーニングの反復において、当該シミュレートされた加入発声は第1の話者により話されているかもしれず、当該シミュレートされた検証発声は第2の話者により話されているかもしれない。幾つかの実装では、異なる順列を生成するかまたは当該シミュレートされた検証発声と加入発声の話者の間の入力サンプル402における順列数を最大化する方式で発声410 a乃至nのグループをランダムに選択せずその代わりに決定的に発声410 a乃至nのグループを選択する、選択アルゴリズムを使用してもよい。単純な例として、3つの異なるトレーニング話者からの発声の3つのグループA、B、およびCが発声プール408で利用可能であった場合、9個の異なる入力サンプル402が9個のトレーニングの反復、即ち、(A、A)、(A、B)、(A、C)、(B、A)、(B、B)、(B、C)、(C、A)、(C、B)、および(C、C)に関して生成されてもよい。トレーニングの反復はまた、これらの同一のグループのペアとともに生じうるが、当該グループ内では異なる発声を選択される。

#### 【0079】

本明細書で説明する当該トレーニングアプローチの1つの履歴は、事前に選択された数の話者のうち特定の話者に属するとして入力を分類するために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他のアプローチと対照的に、より多くの様々な話者を当該ネットワークをトレーニングするために使用してもよい。さらに、高信頼のトレーニングを保証するためにトレーニング話者ごとに要求される(トレーニング話者ごとに実際に使用される当該1つまたは複数の発声以外の)最小数のトレーニング発声はない。なぜならば、当該ネットワークは特定の話者に対してトレーニングされず、その代わりに所与の入力サンプル402が当該シミュレートされた検証発声および加入発声の中にマッチング話者または非マッチング話者があるかどうかに基づいてトレーニングされるからである。

#### 【0080】

図4Bは、ニューラルネットワークのトレーニングの間に入力サンプルに関する発声を選択するためのシャッフル技術の略図400bを示す。本図に示すように、トレーニング・サンプルのバッチにおけるサンプルは全て、当該バッチ内のトレーニング・サンプルの間の発声の良好なシャッフルおよび多様性を取得するために異なるプールから来ることができる。当該シャッフル技術は、より堅牢なおよび高信頼のニューラルネットワークのトレーニングをもたらす。

#### 【0081】

10

20

30

40

50

図 5 A および 5 B を参照すると、話者検証モデルで使用されうる例示的なニューラルネットワーク 5 0 2、5 1 2 のブロック図が示されている。幾つかの実装ではニューラルネットワーク 5 0 2、5 1 2 の何れかを、図 1 乃至 4 B および 6 に関して説明した技術を実装するために使用してもよく、図 2 乃至 4 B に関して説明したトレーニング技術を含む。

【 0 0 8 2 】

図 5 A のディープ・ニューラルネットワーク 5 0 2 のアーキテクチャは、ローカルに接続されたレイヤ 5 0 4 を含み、それに続いて 1 つまたは複数の完全に接続された隠蔽されたレイヤ 5 0 6 a 乃至 n を含む。ローカルに接続されたレイヤ 5 0 4 および完全に接続されたレイヤ 5 0 6 a 乃至 n は整流線形ユニット ( R e L U ) を有してもよい。ネットワーク 5 0 2 の最後のレイヤは完全に接続された線形レイヤ 5 0 8 であり、これは、入力発声 ( または発声のフレーム ) 5 0 3 a に基づいて話者表現 5 1 0 a を出力する。表現 5 1 0 a の前の最後のレイヤ 5 0 2 は、幾つかの実装では、非負の活性化を全空間にマップし、投影を決定するための線形レイヤである。当該全空間は、R e L u 活性化が  $y = \max(x, 0)$  のような関数でありうるという概念を指す。したがって、当該話者表現を形成する活性化 ( y ) は常に正のベクトルであってもよい。かかる活性化関数が線形活性化関数  $y = x$  により変更される場合、当該話者表現を潜在的に正負の値を有するベクトルとして生成することができる。後者は、例えば、それに余弦距離比較関数が続くとき、当該話者のより適切な表現であることができる。

【 0 0 8 3 】

ニューラルネットワーク 5 0 2 の構成は一般に、固定された長さのトレーニング発声、または固定された数の発声のフレームを処理することができる。ニューラルネットワーク 5 0 2 がトレーニングされ、後に加入および検証段階における実行時の間に使用されるとき、発声を適切に刈り込みまたはパディングして、当該発声がニューラルネットワーク 5 0 2 により処理するために要求される固定長を有することを保証してもよい。結果として、ニューラルネットワーク 5 0 2 は、単一のパス、例えば、ディープ・ニューラルネットワーク 5 0 2 を通じた単一の前方伝播で話者表現を計算することができる。これにより、当該話者表現を、発声の異なる部分の逐次的処理を含む技術より低いレイテンシで生成することができる。

【 0 0 8 4 】

次に、図 5 B に示すニューラルネットワーク 5 1 2 は、再帰型ニューラルネットワークである。ニューラルネットワーク 5 0 2 のアーキテクチャと異なり、ニューラルネットワーク 5 1 2 は可変長の入力発声を処理することができる。例えば、発声 5 0 3 b は、ニューラルネットワーク 5 1 2 が使用されているコンテキストに依存して、トレーニング発声、加入発声、または検証発声であってもよい。発声 5 0 3 b を複数のフレームに分割してもよく、当該複数のフレームは固定長を有してもよい。ニューラルネットワーク 5 1 2 に入力されたフレームの数は、発声 5 0 3 b の全体の長さの関数であってもよい。換言すれば、長い発声はより多くのフレームを有してもよく、短い発声はより少ないフレームを有してもよい。発声 5 0 3 b のフレームは長短期メモリ ( L S T M ) レイヤ 5 1 6 に入力される。1 つまたは複数の追加の隠蔽されたレイヤが L S T M レイヤ 5 1 6 に続いてもよい。ネットワーク 5 1 2 の最後のレイヤは再度、完全に接続された線形レイヤ 5 1 8 である。幾つかのケースでは、完全に接続された線形レイヤ 5 1 8 は、非負の活性化を当該全空間にマッピングし、投影を決定することで話者表現 5 1 0 b を出力してもよい。ニューラルネットワーク 5 1 2 は可変長さ発声を扱うことができるので、これは発声の単語またはフレーズが予め定義されず異なる発声の間で変化するテキスト独立な話者検証に良く適しうる。

【 0 0 8 5 】

図 5 A および 5 B に示すニューラルネットワーク 5 0 2 および 5 1 2 は特定の構成を有するとして示されているが、本明細書で説明する当該技術で使用されうるニューラルネットワークはこれらの例により限定されない。例えば、ニューラルネットワークの隠蔽されたトポロジは異なる数および配置のレイヤを有してもよく、当該レイヤは、完全に接続さ

10

20

30

40

50

れたレイヤ、ローカルに接続されたレイヤ、または長短期メモリレイヤのような任意の回帰レイヤを含んでも含まなくてもよい。幾つかの実装では、ニューラルネットワークは従来型のニューラルネットワークであってもよい。

【0086】

図6は、本明細書で説明する技術に従ってトレーニングされている話者検証モデルおよびニューラルネットワークを用いて発声を検証するための例示的なプロセス600の流れ図である。プロセス600は一般に図1に示す検証段階(段階E乃至G)に対応する。図6で参照されるニューラルネットワークは、幾つかの実装では、図2乃至4Bに関して説明した技術に従ってトレーニングされてもよく、図5Aまたは5Bに示すような構造を有してもよい。

10

【0087】

段階602で、発声はコンピューティングデバイスのユーザから受信されることができ。例えば、ユーザは、自分のスマートフォンをアンロックするか、または、幾つかの他の機能をコンピューティングデバイスで実施したいかもしれない。しかし、スマートフォンは、それがロックされない前に、または所望の機能が実施される前にユーザを認証するようにユーザに要求してもよい。幾つかの実装では、当該認証を、スマートフォンの話者検証モデルを用いてユーザの音声の特性に基づいて実施してもよい。スマートフォンは、検証発声を話すようにユーザに促してもよく、検証発声を段階602でスマートフォンにより受信し記録してもよい。

【0088】

段階604で、スマートフォンは、当該受信された発声に基づいて話者表現を生成するためにニューラルネットワークにアクセスする。ニューラルネットワークは、ローカルにスマートフォンに格納されてもよく、または例えば、アプリケーションプログラミングインタフェース(API)を介してリモートコンピューティングシステム上でアクセスされてもよい。ニューラルネットワークは本明細書で説明する技術に従ってトレーニングされてもよく、それぞれシミュレートされた検証発声および複数のシミュレートされた加入発声を含むデータのサンプルに基づいてトレーニングされているかもしれない。ニューラルネットワークは、単一のパス内でニューラルネットワークを通じて、発声の全体を特徴付けるデータを処理するように構成されてもよい。段階606で、当該受信された発声を特徴付けるデータがニューラルネットワークへの入力として提供される。ニューラルネットワークは、当該入力を処理し、ユーザの音声の区別的な特性を示す話者表現を生成する。

20

30

【0089】

段階608で、話者モデルがスマートフォンでアクセスされる。当該話者モデルは加入したユーザの音声の区別的な特徴を示してもよい。幾つかの実装では、当該話者モデルが、加入したユーザの各発声を特徴付けるデータからニューラルネットワークにより生成された複数の話者表現の平均に基づいてもよい。段階610で、検証発声に基づいて段階606で生成された話者表現が、当該話者モデルと比較され、または、そうでない場合は当該話者モデルに関して評価される。幾つかの実装では、当該比較または他の評価がユーザのスマートフォン上で話者検証モデルにより実施される。当該話者検証モデルは、検証発声に対する話者モデルと話者表現の間の類似性の距離または他の測定値を決定してもよい。類似性の距離または他の測定値に基づいて、当該話者検証モデルは、ユーザの音声が入加入したユーザの音声と十分に同様である場合に、ユーザを認証してもよい。そうでない場合、ユーザの音声の類似性が加入したユーザの音声に関して少なくとも閾値類似性スコアを満たさない場合に、当該話者検証モデルはユーザが認証されないという指示を生成してもよい。

40

【0090】

幾つかの実装では、当該話者検証モデルが、検証発声が入加入した話者により話されたことを十分な確信度で判定した場合、加入したユーザに対する話者モデルをついで、検証発声に基づいて更新してもよい。例えば、当該デバイスが以下の3つの検証発声にどのように応答しうるかを考える。当該話者検証モデルが第1の検証発声を話したユーザのアイデ

50

ンティティを拒否する（例えば、デバイスは第1の検証発声に回答してアンロックを拒否してもよい）ように、3つの検証発声のうち最初のものに対する類似性スコアは第1の閾値より小さい。当該3つの検証発声のうち2番目のものに対する類似性スコアは、第2の検証発声を話したユーザのアイデンティティが受理されるように、第1の閾値を満たしてもよい。しかし、第2の検証発声に対する類似性スコアは、加入したユーザの話者モデルを第2の検証発声に基づいて更新するために十分に高い。最後に、検証発声のうち第3のものに対する類似性スコアが、第3の検証発声を話したユーザのアイデンティティが受理される（例えば、デバイスのアンロックのような第1の1組のアクションを実施してもよい）ように、第1の閾値を満たし、また、加入したユーザに対する当該話者モデルが第3の検証発声に基づいて更新されるように、より高い第2の閾値を満たす。当該話者モデルが、第3の検証発声に対してニューラルネットワークにより生成された話者表現を第1のインスタンス内の話者モデルを生成するために使用されたユーザの加入発声からの他の話者表現と結合（例えば、平均化）することで更新されてもよい。

10

20

30

40

50

#### 【0091】

段階612で、スマートフォンは次いで、ユーザが認証されたかどうかに基づいて動作を行うことができる。例えば、当該発声を提供したユーザが加入したユーザであるという判定に回答して、スマートフォンを起動またはアンロックしてもよい。しかし、当該発声を提供したユーザが加入したユーザでないか、または、複数の加入したユーザの1つでないと判定された場合、スマートフォンはロックされたままであってもよく、または、そうでない場合はユーザが実施するために選択された1つまたは複数の機能の実施をブロックしてもよい。別のアプリケーションでは、本明細書で説明する話者検証技術をユーザデバイス（例えば、スマートフォン、ノートブックコンピュータ、ウェアブルデバイス）で使用して、当該デバイスにより非認証されたユーザ（例えば、音声当該デバイスで加入されていないユーザ）から検出された会話入力を拒否してもよい。例えば、当該デバイスがロックされていない状態であるとき、当該デバイスは、ユーザが当該デバイスに実施してほしいアクション（例えば「Navigate to the football game」または「Play my music collection」）を示すデバイスの認証されたユーザにより話された音声コマンドをリッスンしてもよい。幾つかの実装では、当該音声コマンドが当該認証されたユーザにより話されたことと判定できる場合には、当該デバイスは当該音声コマンドにより示される当該要求されたアクションをのみを実施してもよい。このように、他の非認証されたユーザからの雑談を、例えば、拒絶してもよい。

#### 【0092】

図7は、本明細書で説明する技術を実装するために使用できるコンピューティングデバイス700およびモバイルコンピューティングデバイスの1例を示す。コンピューティングデバイス700は、ラップトップ、デスクトップ、ワークステーション、携帯情報端末、サーバ、ブレードサーバ、メインフレーム、および他の適切なコンピュータのような様々な形態のデジタルコンピュータを表すことを意図している。当該モバイルコンピューティングデバイスは、携帯情報端末、セルラスマートフォン、スマートフォン、および他の同様なコンピューティングデバイスのような様々な形態のモバイルデバイスを表すことを意図している。ここで示したコンポーネント、それらの接続および関係、およびそれらの機能は例示的なものにすぎず、本明細書で説明および/またはクレームした発明の実装を限定しようとするものではない。

#### 【0093】

コンピューティングデバイス700は、プロセッサ702、メモリ704、記憶デバイス706、メモリ704および複数の高速拡張ポート710に接続する高速インタフェース708、および低速拡張ポート714および記憶デバイス706に接続する低速インタフェース712を含む。プロセッサ702、メモリ704、記憶デバイス706、高速インタフェース708、高速拡張ポート710、および低速インタフェース712の各々は様々なバスを用いて相互接続され、必要に応じて共通のマザーボード上でまたは他の方式でマウントされてもよい。プロセッサ702は、高速インタフェース708に接続される

ディスプレイ 716 のような外部入力/出力デバイスに GUI に関するグラフィカル情報を表示するためのメモリ 704 または記憶デバイス 706 に格納された命令を含む、コンピューティングデバイス 700 内で実行するための命令を処理することができる。他の実装では、複数のプロセッサおよび/または複数のバスは、必要に応じて、複数のメモリおよびメモリのタイプに沿って使用されてもよい。また、複数のコンピューティングデバイスは、(例えば、サーババンク、ブレードサーバのグループ、またはマルチプロセッサシステムとして)必要な動作の部分を提供する各デバイスと接続されてもよい。

#### 【0094】

メモリ 704 は情報をコンピューティングデバイス 700 内に格納する。幾つかの実装では、メモリ 704 は揮発性メモリユニットまたはユニットである。幾つかの実装では、メモリ 704 は非揮発性メモリユニットまたはユニットである。メモリ 704 はまた、磁気または光ディスクのような別の形態のコンピュータ可読媒体であってもよい。

10

#### 【0095】

記憶デバイス 706 は大容量記憶をコンピューティングデバイス 700 に提供することができる。幾つかの実装では、記憶デバイス 706 は、記憶領域ネットワークまたは他の構成でのデバイスを含む、フロッピーディスクデバイス、ハードディスクデバイス、光ディスクデバイス、またはテープデバイス、フラッシュ・メモリまたは他の同様な固体状態メモリデバイス、またはデバイスのアレイのようなコンピュータ可読媒体であってもよい。当該コンピュータプログラム製品はまた、実行されたとき上述したような 1 つまたは複数の方法を実施する命令を含んでもよい。当該コンピュータプログラム製品をまた、プロセッサ 702 上のメモリ 704、記憶デバイス 706、またはメモリのようなコンピュータまたはマシン可読媒体で有形に具体化することができる。

20

#### 【0096】

高速インタフェース 708 はコンピューティングデバイス 700 に対する帯域幅集約的な動作を管理し、低速インタフェース 712 はより低い帯域幅集約的な動作を管理する。かかる機能の割当ては例示的なものにすぎない。幾つかの実装では、高速インタフェース 708 は、(例えば、グラフィックスプロセッサまたはアクセラレータを通じて)メモリ 704、ディスプレイ 716 に接続され、高速拡張ポート 710 に接続される。高速拡張ポート 710 は様々な拡張カード(図示せず)を受け付けてもよい。当該実装では、低速インタフェース 712 は記憶デバイス 706 および低速拡張ポート 714 に接続される。低速拡張ポート 714 は、様々な通信ポート(例えば、USB、Bluetooth(登録商標)、イーサネット(登録商標)、ワイヤレスイーサネット(登録商標))を含んでもよく、キーボード、ポインティングデバイス、スキャナのような 1 つまたは複数の入力/出力デバイス、またはスイッチまたはルータのようなネットワークデバイスに、例えば、ネットワークアダプタを通じて接続してもよい。

30

#### 【0097】

コンピューティングデバイス 700 を本図に示すように幾つかの異なる形態で実装してもよい。例えば、それを標準サーバ 720 として、またはかかるサーバのグループにおいて複数回、実装してもよい。さらに、ラップトップコンピュータ 722 のようなパーソナルコンピュータで実装してもよい。また、ラックサーバシステム 724 の一部として実装してもよい。代替的に、コンピューティングデバイス 700 からのコンポーネントは、モバイルコンピューティングデバイス 750 のようなモバイルデバイス(図示せず)内の他のコンポーネントと結合されてもよい。かかるデバイスの各々はコンピューティングデバイス 700 およびモバイルコンピューティングデバイス 750 の 1 つまたは複数を含んでもよく、システム全体は互いと通信する複数のコンピューティングデバイスで構成されてもよい。

40

#### 【0098】

モバイルコンピューティングデバイス 750 は、プロセッサ 752、メモリ 764、他のコンポーネントのうちディスプレイ 754、通信インタフェース 766、および送受信機 768 のような入力/出力デバイスを含む。モバイルコンピューティングデバイス 75

50

0にまた、追加の記憶を提供するための、マイクロドライブまたは他のデバイスのような記憶デバイスが提供されてもよい。プロセッサ752、メモリ764、ディスプレイ754、通信インタフェース766、および送受信機768の各々は様々なバスを用いて相互接続され、当該コンポーネントの幾つかは必要に応じて共通のマザーボード上でまたは他の方式でマウントされてもよい。

【0099】

プロセッサ752は、メモリ764に格納された命令を含む、モバイルコンピューティングデバイス750内の命令を実行することができる。プロセッサ752を、別々のおよび複数のアナログおよびデジタルプロセッサを含むチップから成るチップ・セットとして実装してもよい。プロセッサ752は、例えば、モバイルコンピューティングデバイス750により実行されるユーザインタフェース、アプリケーション、およびモバイルコンピューティングデバイス750によるワイヤレス通信の制御のような、モバイルコンピューティングデバイス750の他のコンポーネントの協調を提供してもよい。

10

【0100】

プロセッサ752は、ディスプレイ754に接続される制御インタフェース758およびディスプレイインタフェース756を通じてユーザと通信してもよい。ディスプレイ754は、例えば、TFT（薄膜トランジスタ液晶ディスプレイ）ディスプレイまたはOLED（有機発光ダイオード）ディスプレイ、または他の適切なディスプレイ技術であってもよい。ディスプレイインタフェース756は、ディスプレイ754を駆動してグラフィカルおよび他の情報をユーザに提供するための適切な回路を備えてもよい。制御インタフェース758は、ユーザからコマンドを受信し、プロセッサ752に送信するために当該コマンドを変換してもよい。さらに、外部インタフェース762は、他のデバイスとのモバイルコンピューティングデバイス750の近領域通信を可能するために、プロセッサ752との通信を提供してもよい。外部インタフェース762は、例えば、幾つかの実装では有線通信を提供し、または他の実装ではワイヤレス通信を提供してもよく、複数のインタフェースをまた使用してもよい。

20

【0101】

メモリ764はモバイルコンピューティングデバイス750内に格納する。メモリ764を、コンピュータ可読媒体または媒体、揮発性メモリユニットまたはユニット、または非揮発性メモリユニットまたはユニットの1つまたは複数として実装することができる。拡張メモリ774はまた、拡張インタフェース772を通じてモバイルコンピューティングデバイス750に提供され接続されてもよい。拡張インタフェース772は、例えば、SIMM（Single In Line Memory Module）カードインタフェースを含んでもよい。拡張メモリ774はモバイルコンピューティングデバイス750に対する追加の記憶空間を提供してもよく、または、モバイルコンピューティングデバイス750に対するアプリケーションまたは他の情報を格納してもよい。特に、拡張メモリ774は、上述したプロセスを実行または補完する命令を含んでもよく、セキュア情報を含んでもよい。したがって、例えば、拡張メモリ774を、モバイルコンピューティングデバイス750に対するセキュリティモジュールとして提供してもよく、モバイルコンピューティングデバイス750の安全な使用を許可する命令でプログラムされてもよい。さらに、セキュアアプリケーションを、ハック不能な方式でSIMMカードに識別情報を配置するといった、当該SIMMカードを介して追加の情報とともに提供してもよい。

30

40

【0102】

当該メモリは、以下で説明するように、例えば、フラッシュ・メモリおよび/またはNVRAMメモリ（非揮発性ランダム・アクセスメモリ）を含んでもよい。当該コンピュータプログラム製品は、実行されたとき、上述したものののような1つまたは複数の方法を実施する命令を含む。当該コンピュータプログラム製品は、メモリ764、拡張メモリ774、またはプロセッサ752上のメモリのようなコンピュータまたはマシン可読媒体であることができる。幾つかの実装では、当該コンピュータプログラム製品を、伝播信号で、例えば、送受信機768または外部インタフェース762上で受信することができる。

50

## 【 0 1 0 3 】

モバイルコンピューティングデバイス 7 5 0 は通信インタフェース 7 6 6 を通じて無線で通信してもよい。通信インタフェース 7 6 6 は、必要な場合はデジタル信号処理回路を含んでもよい。通信インタフェース 7 6 6 は、とりわけ G S M (登録商標) 通話 (Global System for Mobile communications)、S M S (Short Message Service)、E M S (Enhanced Messaging Service)、または M M S メッセージング (Multimedia Messaging Service)、C D M A (code division multiple access)、T D M A (time division multiple access)、P D C (Personal Digital Cellular)、W C D M A (登録商標) (Wideband Code Division Multiple Access)、C D M A 2 0 0 0、または G P R S (General Packet Radio Service) のような、様々なモードまたはプロトコルの下での通信を提供してもよい。かかる通信を、例えば、送受信機 7 6 8 を通じて無線周波数を用いて行ってもよい。さらに、短波通信を、例えば B l u e t o o t h (登録商標)、W i F i、または他のかかる送受信機 (図示せず) を用いて行ってもよい。さらに、G P S (Global Positioning System) 受信器モジュール 7 7 0 は追加のナビゲーション位置関連のワイヤレスデータをモバイルコンピューティングデバイス 7 5 0 に提供してもよい。モバイルコンピューティングデバイス 7 5 0 は、モバイルコンピューティングデバイス 7 5 0 で実行されているアプリケーションにより必要に応じて使用されてもよい。

10

## 【 0 1 0 4 】

モバイルコンピューティングデバイス 7 5 0 はまた、オーディオコーデック 7 6 0 を用いて可聴的に通信してもよく、オーディオコーデック 7 6 0 は話された情報をユーザから受信し、それを使用可能なデジタル情報に変換してもよい。オーディオコーデック 7 6 0 は、話者を通じて、例えば、モバイルコンピューティングデバイス 7 5 0 のハンドセットでユーザに対する可聴音を同様に生成してもよい。かかる音は音声通話からの音を含んでもよく、記録された音 (例えば、音声メッセージ、音楽ファイル等) を含んでもよく、また、モバイルコンピューティングデバイス 7 5 0 で動作しているアプリケーションにより生成された音を含んでもよい。

20

## 【 0 1 0 5 】

モバイルコンピューティングデバイス 7 5 0 を、本図で示すように幾つかの異なる形態で実装してもよい。例えば、セルラスマートフォン 7 8 0 として実装してもよい。また、スマートフォン 7 8 2、携帯情報端末、または他の同様なモバイルデバイスの一部として実装してもよい。

30

## 【 0 1 0 6 】

本明細書で説明したシステムおよび技術の様々な実装を、デジタル電子回路、集積回路、特別に設計された A S I C (特殊用途向け集積回路)、コンピュータハードウェア、ファームウェア、ソフトウェア、および/またはその組合せで実現することができる。これらの様々な実装は、少なくとも 1 つのプログラム可能プロセッサを含むプログラム可能システムで実行可能および/または解釈可能である 1 つまたは複数のコンピュータプログラムでの実装を含むことができる。当該少なくとも 1 つのプログラム可能プロセッサは、特殊目的または一般的な目的であってもよく、記憶システム、少なくとも 1 つの入力デバイス、および少なくとも 1 つの出力デバイスとデータおよび命令を送受信するために接続されてもよい。

40

## 【 0 1 0 7 】

これらのコンピュータプログラム (プログラム、ソフトウェア、ソフトウェアアプリケーションまたはコードとしても知られる) はプログラム可能プロセッサに対するマシン命令を含み、高レベル手続き型および/またはオブジェクト指向プログラミング言語で、および/またはアセンブリ/マシン言語で実装することができる。本明細書で使用する際、当該用語マシン可読媒体およびコンピュータ可読媒体は、マシン命令をマシン可読信号として受信するマシン可読媒体を含む、マシン命令および/またはデータをプログラム可能プロセッサに提供するために使用される、任意のコンピュータプログラム製品、装置および/またはデバイス (例えば、磁気ディスク、光ディスク、メモリ、プログラム可能論理デ

50

バイス ( P L D ) ) を指す。マシン可読信号という用語は、マシン命令および / またはデータをプログラム可能プロセッサに提供するために使用される任意の信号を指す。

【 0 1 0 8 】

ユーザとの対話を提供するために、本明細書で説明したシステムおよび技術を、ユーザが当該コンピュータへの入力を提供できるユーザおよびキーボードおよびポインティングデバイス ( 例えば、マウスまたはトラックボール ) に情報を表示するための、ディスプレイデバイス ( 例えば、C R T ( カソード・レイ・チューブ ) または L C D ( 液晶ディスプレイ ) モニタ ) を有するコンピュータで実装することができる。他の種類のデバイスを、ユーザとの対話を提供するために使用することができる。例えば、ユーザに提供されるフィードバックは任意の形態のセンサフィードバック ( 例えば、視覚フィードバック、可聴フィードバック、または触覚フィードバック ) であることができ、ユーザからの入力を音響、会話、または触覚入力を含む任意の形態で受信することができる。

10

【 0 1 0 9 】

本明細書で説明したシステムおよび技術を、( 例えば、データサーバとして ) バックエンドコンポーネントを含む、またはミドルウェアコンポーネント ( 例えば、アプリケーションサーバ ) を含む、またはフロントエンドコンポーネント ( 例えば、ユーザがそれを通じて本明細書で説明したシステムおよび技術の実装と対話できるグラフィカルユーザインタフェースまたはウェブ・ブラウザを有するクライアントコンピュータ ) を含む、またはかかるバックエンド、ミドルウェア、またはフロントエンドコンポーネント任意の組合せを含む、コンピューティングシステムで実装することができる。当該システムの当該コンポーネントは、デジタルデータ通信 ( 例えば、通信ネットワーク ) の任意の形態または媒体により相互接続されることができる。通信ネットワークの例はローカル・エリア・ネットワーク ( L A N ) 、広帯域ネットワーク ( W A N ) 、およびインターネットを含む。

20

【 0 1 1 0 】

当該コンピューティングシステムはクライアントおよびサーバを含むことができる。クライアントおよびサーバは一般に互いから離れており、一般に通信ネットワークを通じて対話する。クライアントおよびサーバの関係は、当該各コンピュータで実行され互いにクライアントサーバ関係を有するコンピュータプログラムにより生ずる。

【 0 1 1 1 】

様々な実装を上で詳細に説明したが、他の修正が可能である。さらに、本図で示した論理フローは、所望の結果を実現するために、示した特定の順序、または逐次的な順序を必要としない。さらに、他のステップを提供してもよく、またはステップを当該説明したフローから削除してもよく、他のコンポーネントを当該説明したシステムに追加し、または、そこから削除してもよい。したがって、他の実装は添付の特許請求の範囲内にある。

30

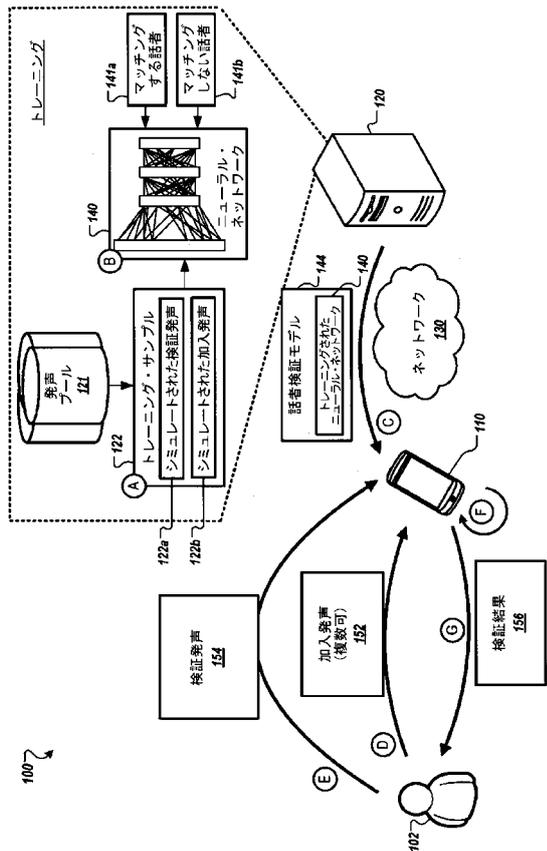
【 符号の説明 】

【 0 1 1 2 】

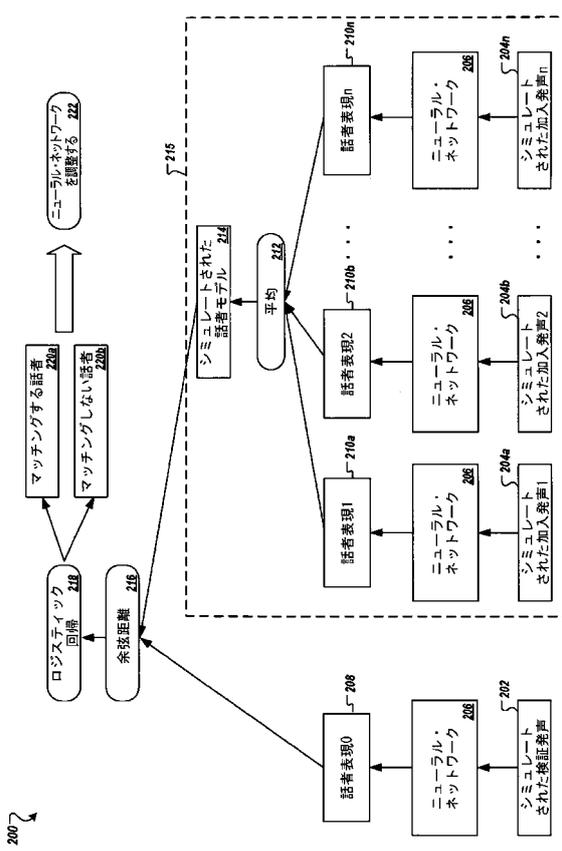
- 1 2 1 発声プール
- 1 2 2 トレーニング・サンプル
- 1 2 2 a シミュレートされた検証発声
- 1 2 2 b シミュレートされた加入発声
- 1 3 0 ネットワーク
- 1 4 0 トレーニングされたニューラル・ネットワーク
- 1 4 1 a マッチングする話者
- 1 4 1 b マッチングしない話者
- 1 4 4 話者検証モデル
- 1 5 2 加入発声 ( 複数可 )
- 1 5 4 検証発声
- 1 5 6 検証結果

40

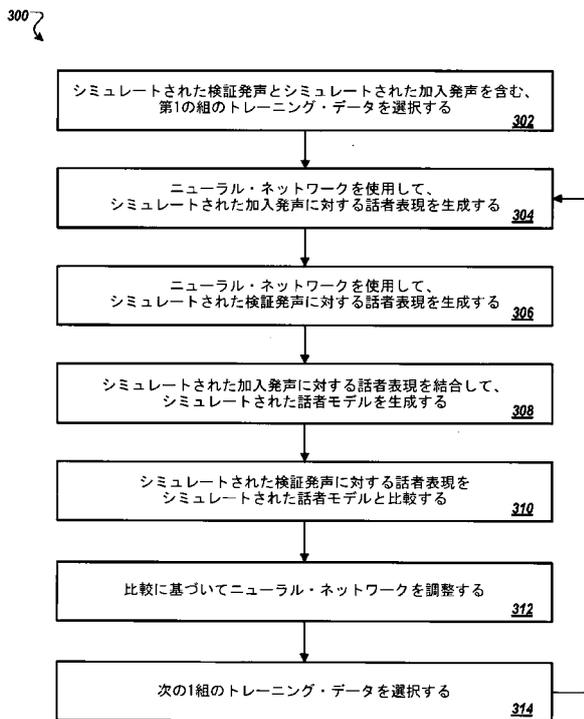
【図1】



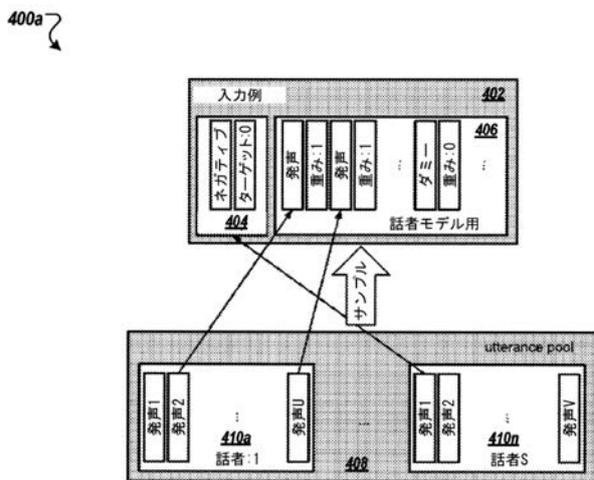
【図2】



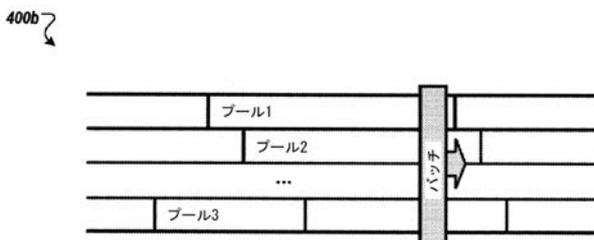
【図3】



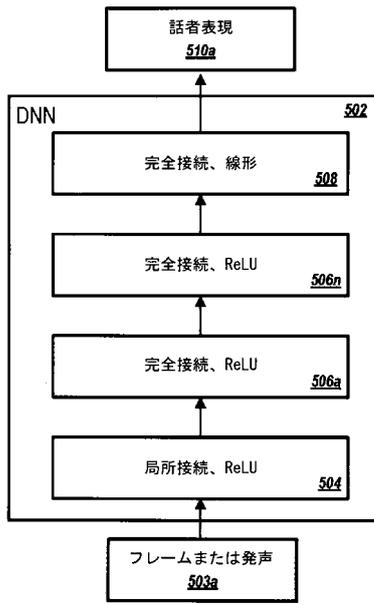
【図4A】



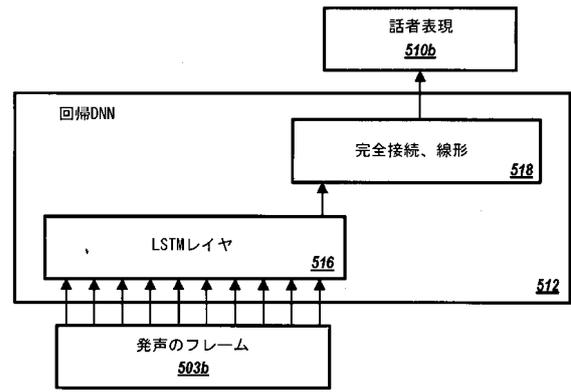
【図4B】



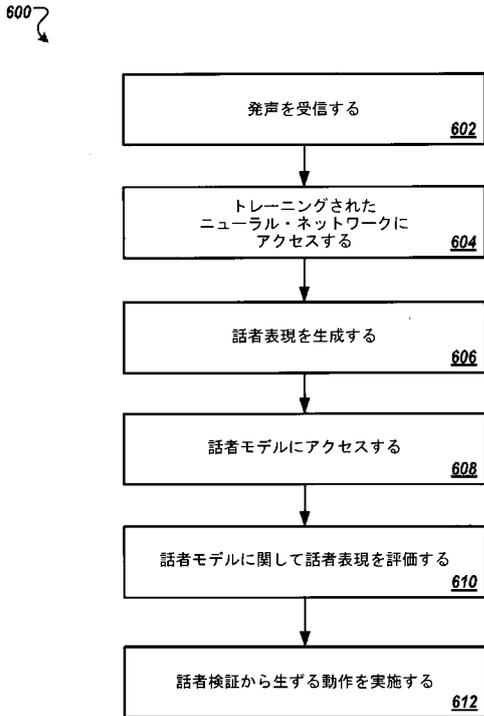
【図5A】



【図5B】



【図6】



【図7】

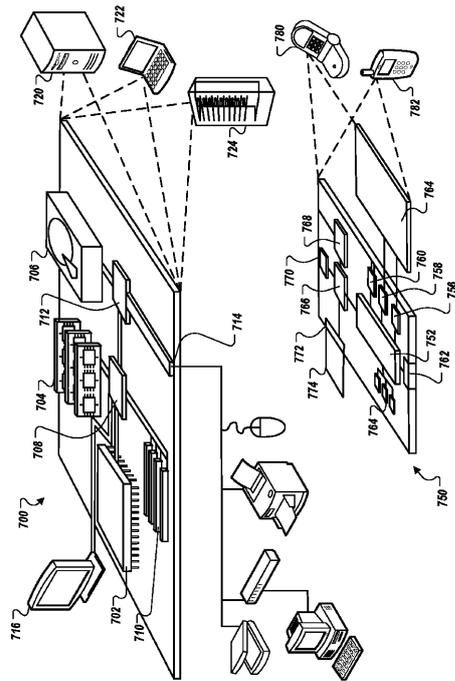


FIG. 7

## 【 国際調査報告 】

## INTERNATIONAL SEARCH REPORT

International application No  
PCT/US2016/044181

A. CLASSIFICATION OF SUBJECT MATTER INV. G10L17/04 G10L17/18 G10L17/02 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G10L		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal, WPI Data		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X A	US 2015/127336 A1 (LEI XIN [US] ET AL) 7 May 2015 (2015-05-07) paragraph [0030] paragraphs [0035] - [0040] paragraphs [0042], [0043] paragraph [0045] paragraphs [0060], [0061] paragraph [0068] paragraphs [0076] - [0079] ----- -/--	1-6 7,9
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents :		
"A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed		"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family
Date of the actual completion of the international search  6 October 2016		Date of mailing of the international search report  19/10/2016
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer  Ramos Sánchez, U

1

INTERNATIONAL SEARCH REPORT

International application No  
PCT/US2016/044181

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X,P	<p>HEIGOLD GEORG ET AL: "End-to-end text-dependent speaker verification", 2016 IEEE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH AND SIGNAL PROCESSING (ICASSP), IEEE, 20 March 2016 (2016-03-20), pages 5115-5119, XP032901578, DOI: 10.1109/ICASSP.2016.7472652 [retrieved on 2016-05-18] section '6. Summary &amp; Conclusion', first 4 lines; page 5118 section '3. D-Vector Baseline Approach'; page 5116, left-hand column section '4. End-To-End Speaker Verification'; page 5116, right-hand column, paragraph 3 - page 5117, left-hand column, paragraph 4 page 5117, right-hand column, paragraphs 2,3</p> <p style="text-align: center;">-----</p>	1-21

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No

PCT/US2016/044181

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2015127336	A1	NONE	
-----			

## フロントページの続き

(81)指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, TM), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(72)発明者 サミー・ベンジオ

アメリカ合衆国・カリフォルニア・94043・マウンテン・ビュー・アンフィシアター・パーク  
ウェイ・1600

(72)発明者 イグナシオ・ロペス・モレーノ

アメリカ合衆国・カリフォルニア・94043・マウンテン・ビュー・アンフィシアター・パーク  
ウェイ・1600