

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2024-50983

(P2024-50983A)

(43)公開日 令和6年4月10日(2024.4.10)

(51)国際特許分類

F I

G 1 0 L 15/32 (2013.01)

G 1 0 L 15/32 2 0 0 Z

G 1 0 L 15/16 (2006.01)

G 1 0 L 15/16

審査請求 有 請求項の数 20 O L (全22頁)

(21)出願番号 特願2024-24188(P2024-24188)  
 (22)出願日 令和6年2月21日(2024.2.21)  
 (62)分割の表示 特願2023-558803(P2023-558803)  
 )の分割  
 原出願日 令和4年3月22日(2022.3.22)  
 (31)優先権主張番号 63/166,916  
 (32)優先日 令和3年3月26日(2021.3.26)  
 (33)優先権主張国・地域又は機関  
 米国(US)

(71)出願人 502208397  
 グーグル エルエルシー  
 Google LLC  
 アメリカ合衆国 カリフォルニア州 9 4  
 0 4 3 マウンテン ビュー アンフィシ  
 アター パークウェイ 1 6 0 0  
 1 6 0 0 Amphitheatre P  
 arkway 9 4 0 4 3 Mounta  
 in View, CA U.S.A.  
 (74)代理人 100142907  
 弁理士 本田 淳  
 (72)発明者  
 ガウル、 ニーラジ  
 アメリカ合衆国 9 4 0 4 3 カリフォル  
 ニア州 マウンテン ビュー アンフィシ  
 アター パークウェイ 1 6 0 0

最終頁に続く

(54)【発明の名称】 自動音声認識のための多言語再スコアリングモデル

(57)【要約】

【課題】 N個の候補仮説の中から、最も高い各々の総合スコアを有している候補仮説を、発話の最終トランスクリプション(120)として選択する。

【解決手段】 方法(400)は、発話(106)に対応する音声データから抽出済みの音響フレーム(110)のシーケンスを受信する。第1パス(301)中、音響フレームのシーケンスを処理して、発話に対するN個の候補仮説(204)を生成する。第2パス(302)中、および各候補仮説に対して、本方法は：各々の非正規化尤度スコア(325)を生成する工程と、各々の外部言語モデルスコア(315)を生成する工程と、対応する候補仮説の事前統計をモデル化する単体スコア(205)を生成する工程と、非正規化尤度スコア、外部言語モデルスコア、および単体スコアに基づき、候補仮説に対する各々の総合スコア(355)を生成する工程と、を備える。

【選択図】 図1

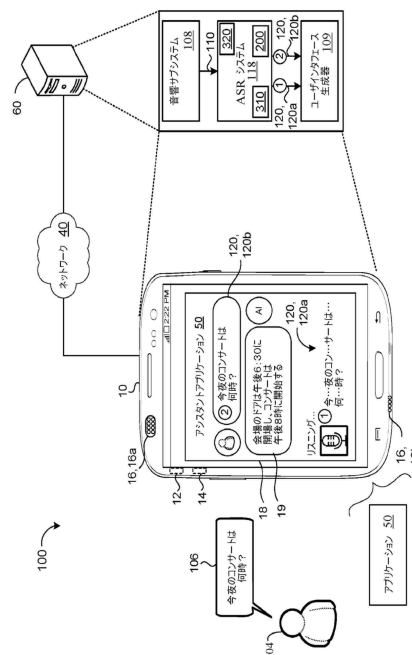


FIG. 1

## 【特許請求の範囲】

## 【請求項 1】

データ処理ハードウェア上で実行されたときに前記データ処理ハードウェアに動作を実行させるコンピュータ実装方法としての方法であって、前記動作は、

発話のグラウンドトゥールーストランスクリプションとでペアにされた前記発話に対応するトレーニング音声データを備えている、書き起こされた音声トレーニングデータを受信する工程と、

第 1 パス中に、音声認識モデルを使用することで前記トレーニング音声データを処理することによって、前記発話のための N 個の候補仮説を生成する工程であって、N 個の前記候補仮説のうちの対応する前記各候補仮説は各々の第 1 パススコアを有している、N 個の前記候補仮説を生成する工程と、

10

第 2 パス中に、N 個の前記候補仮説のうちの対応する前記候補仮説ごとに

ニューラルネットワーク再スコアリングモデルを使用することで、対応する前記候補仮説についての各々の前記第 1 パススコアに基づき、各々の第 2 パススコアを生成するとともに、

対応する前記候補仮説と前記グラウンドトゥールーストランスクリプションとの間の各々の負の編集距離に、ソフトマックス関数を適用する工程と、

前記グラウンドトゥールーストランスクリプションと N 個の前記候補仮説のうちの対応する前記各候補仮説との間の各々の前記負の編集距離に適用される前記ソフトマックス関数に基づき、前記ニューラルネットワーク再スコアリングモデルのモデルパラメータを最適化する工程と、

20

を備えている、方法。

## 【請求項 2】

前記音声認識モデルは、リカレントニューラルネットワーク - トランスデューサ (RNN - T) アーキテクチャを備えている、

請求項 1 に記載の方法。

## 【請求項 3】

N 個の前記候補仮説のうちの対応する前記各候補仮説は、単語ラベルの各々のシーケンスを備えており、

前記各単語ラベルは、各々の埋込ベクトルによって表わされる、

30

請求項 1 に記載の方法。

## 【請求項 4】

N 個の前記候補仮説のうちの対応する前記各候補仮説は、サブワードラベルの各々のシーケンスを備えており、

前記各サブワードラベルは、各々の埋込ベクトルによって表わされる、

請求項 1 に記載の方法。

## 【請求項 5】

前記第 1 パス中に前記音声認識モデルを使用することで生成済みの N 個の前記候補仮説は、候補仮説の N ベストリストを備えている、

請求項 1 に記載の方法。

40

## 【請求項 6】

前記音声認識モデルは、複数のコンフォーマ層を有しているコンフォーマエンコーダを備えているエンコーダデコーダアーキテクチャを備えている、

請求項 1 に記載の方法。

## 【請求項 7】

前記音声認識モデルは、複数のトランスフォーマ層を有しているトランスフォーマエンコーダを備えているエンコーダデコーダアーキテクチャを備えている、

請求項 1 に記載の方法。

## 【請求項 8】

前記動作はさらに、N 個の前記候補仮説について生成済みの前記各第 2 パススコアに基

50

づき、前記発話の最終トランスクリプションとしてN個の前記候補仮説のうちの1つを選択する工程を備えている、

請求項1に記載の方法。

【請求項9】

前記ニューラルネットワーク再スコアリングモデルは、言語固有のニューラルネットワーク再スコアリングモデルを備えている、

請求項1に記載の方法。

【請求項10】

前記ニューラルネットワーク再スコアリングモデルは、多言語ニューラルネットワーク再スコアリングモデルを備えている、

請求項1に記載の方法。

【請求項11】

システムであって、前記システムは、

データ処理ハードウェアと、

前記データ処理ハードウェアに通信するメモリハードウェアであって、前記データ処理ハードウェア上で実行されると前記データ処理ハードウェアに動作を実行させる命令を記憶している前記メモリハードウェアと、

を備えており、前記動作は、

発話のグラウンドトゥーストランスクリプションとペアにされた前記発話に対応するトレーニング音声データを備えている、書き起こされた音声トレーニングデータを受信する工程と、

第1パス中に、音声認識モデルを使用することで前記トレーニング音声データを処理することによって、前記発話のためのN個の候補仮説を生成する工程であって、N個の前記候補仮説のうちの対応する前記各候補仮説は各々の第1パススコアを有している、N個の前記候補仮説を生成する工程と、

第2パス中に、N個の前記候補仮説のうちの対応する前記候補仮説ごとに、

ニューラルネットワーク再スコアリングモデルを使用することで、対応する前記候補仮説についての各々の前記第1パススコアに基づき、各々の第2パススコアを生成するとともに、

対応する前記候補仮説と前記グラウンドトゥーストランスクリプションとの間の各々の負の編集距離に、ソフトマックス関数を適用する工程と、

前記グラウンドトゥーストランスクリプションとN個の前記候補仮説のうちの対応する前記各候補仮説との間の各々の前記負の編集距離に適用される前記ソフトマックス関数に基づき、前記ニューラルネットワーク再スコアリングモデルのモデルパラメータを最適化する工程と、

を備えている、システム。

【請求項12】

前記音声認識モデルは、リカレントニューラルネットワーク-トランスデューサ(RNN-T)アーキテクチャを備えている、

請求項11に記載のシステム。

【請求項13】

N個の前記候補仮説のうちの対応する前記各候補仮説は、単語ラベルの各々のシーケンスを備えており、

前記各単語ラベルは、各々の埋込ベクトルによって表わされる、

請求項11に記載のシステム。

【請求項14】

N個の前記候補仮説のうちの対応する前記各候補仮説は、サブワードラベルの各々のシーケンスを備えており、

前記各サブワードラベルは、各々の埋込ベクトルによって表わされる、

請求項11に記載のシステム。

10

20

30

40

50

## 【請求項 15】

前記第1パス中に前記音声認識モデルを使用することで生成済みのN個の前記候補仮説は、前記候補仮説のNベストリストを備えている、

請求項11に記載のシステム。

## 【請求項 16】

前記音声認識モデルは、複数のコンフォーマ層を有しているコンフォーマエンコーダを備えているエンコーダデコーダアーキテクチャを備えている、

請求項11に記載のシステム。

## 【請求項 17】

前記音声認識モデルは、複数のトランスフォーマ層を有しているトランスフォーマエンコーダを備えているエンコーダデコーダアーキテクチャを備えている、

請求項11に記載のシステム。

10

## 【請求項 18】

前記動作はさらに、N個の前記候補仮説について生成済みの前記各第2パススコアに基づき、前記発話の最終トランスクリプションとしてN個の前記候補仮説のうちの1つを選択する工程を備えている、

請求項11に記載のシステム。

## 【請求項 19】

前記ニューラルネットワーク再スコアリングモデルは、言語固有のニューラルネットワーク再スコアリングモデルを備えている、

請求項11に記載のシステム。

20

## 【請求項 20】

前記ニューラルネットワーク再スコアリングモデルは、多言語ニューラルネットワーク再スコアリングモデルを備えている、

請求項11に記載のシステム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本開示は、自動音声認識のための多言語再スコアリングモデルに関する。

## 【背景技術】

30

## 【0002】

自動音声認識（ASR）システムは、一般的にモバイル機器（移動装置）やその他の機器で使用される技術を提供する。一般に、自動音声認識ASRシステムは、ユーザがモバイル機器に話した内容の正確なトランスクリプション（転写、書き起こし、採録）を提供しようとする。より具体的には、自動音声認識ASRシステムは複数のトランスクリプション候補を生成するとともに、音声入力に一致する可能性が最も高いトランスクリプション候補を出力する。場合によっては、自動音声認識ASRシステムは、ユーザが実際に話した内容には一致しない不正確なトランスクリプションを出力する。このような場合、自動音声認識ASRシステムは複数のトランスクリプション候補を再スコアリングするとともに、音声入力に一致する正確なトランスクリプションを出力する。

40

## 【先行技術文献】

## 【非特許文献】

## 【0003】

【非特許文献1】OGAWA ATSUNORI ET AL, "Rescoring N-Best Speech Recognition List Based on One-on-One Hypothesis Comparison Using Encoder-Classifer Model", 2018 IEEE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH AND SIGNAL PROCESSING (ICASSP), IEEE, 15 April 2018 (2018-04-15), page 6099-6103, XP033403971, DOI: 10.1109/ICASSP.2018.8461405

## 【発明の概要】

50

## 【発明が解決しようとする課題】

## 【0004】

しかし、再スコアリングの1つの課題は、自動音声認識ASRシステムが複数のトランスクリプション候補を正確に再スコアリングするべく、音声入力の言語情報に依存することである。そのため、自動音声認識ASRシステムが多言語の音声環境で再スコアリングを行なうのは、しばしば面倒な作業となる。

## 【課題を解決するための手段】

## 【0005】

本開示の一態様は、データ処理ハードウェア上で実行されるとデータ処理ハードウェアに、自動音声認識のための多言語再スコアリングモデルを使用する動作（操作、オペレーション）を実行させる、コンピュータ実装方法を提供する。動作は、発話に対応する音声データから抽出済みの音響フレームのシーケンスを受信する工程を備えている。第1パス中に、動作は、発話に対するN個の候補仮説（ヒポセシス）を生成するべく、多言語音声認識モデルを使用することで、音響フレームのシーケンスを処理する工程を備えている。第2パス中に、N個の候補仮説の各候補仮説について、本方法は以下を備えている：すなわち、ニューラルオラクルサーチ（NOS）モデルを使用することで、音響フレームのシーケンスと、対応する候補仮説と、に基づき各々の非正規化（正規化されていない）尤度スコアを生成する工程と、言語モデルを使用することで、各々の外部言語モデルスコアを生成する工程と、第1パス中に生成済みの対応する候補仮説の事前統計（プライアスタティスティクス）をモデル化する単体（スタンドアロン、単独）スコアを生成する工程と、非正規化スコア、外部言語モデルスコア、および単体スコア、に基づき候補仮説の各々の総合（オーバーオール、全体）スコアを生成する工程と、を備えている。動作はまた、N個の候補仮説の中から、最も高い各々の総合（全体）スコアを有している候補仮説を、発話の最終トランスクリプションとして選択する工程も備えている。

10

20

## 【0006】

本開示の実施形態は、以下の任意の特徴の1つまたは複数を備えていることができる。いくつかの実装では、N個の候補仮説の各候補仮説は、単語またはサブワードラベル（ワードラベルまたはサブワードラベル）の各々のシーケンスを備えている。ここで、各単語またはサブワードラベルは、各々の埋込ベクトルによって表わされる。外部言語モデルは、テキストのみのデータでトレーニング（訓練、学習）されてもよい。いくつかの例では、ニューラルオラクルサーチNOSモデルは言語固有（言語特異的）のニューラルオラクルサーチNOSモデルを備えている。これらの例では、動作はさらに、発話の言語を示す言語識別子を受信する工程と、異なる各々の言語について各々が訓練（トレーニング）された複数の言語固有ニューラルオラクルサーチNOSモデルの中から、言語固有ニューラルオラクルサーチNOSモデルを選択する工程と、を備えている。

30

## 【0007】

オプションとして、ニューラルオラクルサーチNOSモデルは、多言語ニューラルオラクルサーチNOSモデルを備えていることができる。いくつかの実装では、外部言語モデルは、言語固有の外部言語モデルを備えている。これらの実装では、動作はさらに、発話の言語を示す言語識別子を受信する工程と、異なる各々の言語で各々が訓練された複数の言語固有の外部言語モデルの中から言語固有の外部言語モデルを選択する工程と、を備えている。ニューラルオラクルサーチNOSモデルは、2つの単方向長短期記憶（一方向LSTM）層を備えていることができる。いくつかの例では、音声認識モデルは、複数のコンフォーマ層を有しているコンフォーマエンコーダと、2つのLSTM層を有しているLSTMデコーダと、を備えているエンコーダデコーダアーキテクチャを備えている。

40

## 【0008】

本開示の別の態様は、データ処理ハードウェアと、データ処理ハードウェア上で実行されるとデータ処理ハードウェアに動作（オペレーション、操作）を実行させる命令を記憶するメモリハードウェアと、を備えているシステムを提供する。動作は、発話に対応する音声データから抽出済みの音響フレームのシーケンスを受信する工程を備えている。第1

50

パス中に、動作は、発話に対するN個の候補仮説を生成するべく、多言語音声認識モデルを使用することで、音響フレームのシーケンスを処理する工程を備えている。第2パス中に、N個の候補仮説の各候補仮説について、本方法は以下を備えている：すなわち、ニューラルオラクルサーチ(NOS)モデルを使用することで、音響フレームのシーケンスと、対応する候補仮説と、に基づき各々の非正規化尤度スコアを生成する工程と、言語モデルを使用することで、各々の外部言語モデルスコアを生成する工程と、第1パス中に生成済みの対応する候補仮説の事前統計をモデル化する単体スコアを生成する工程と、非正規化スコア、外部言語モデルスコア、および単体スコア、に基づき候補仮説の各々の総合スコアを生成する工程と、を備えている。動作はまた、N個の候補仮説の中から最も高い各々の総合スコアを有している候補仮説を、発話の最終トランスクリプションとして選択する工程も備えている。

10

【0009】

本開示の実施形態は、以下の任意の特徴の1つまたは複数を備えていることができる。いくつかの実装では、N個の候補仮説の各候補仮説は、単語またはサブワードラベルの各々のシーケンスを備えている。ここで、各単語またはサブワードラベルは、各々の埋込ベクトルによって表わされる。外部言語モデルは、テキストのみのデータでトレーニングされてもよい。いくつかの例では、ニューラルオラクルサーチNOSモデルは言語固有のニューラルオラクルサーチNOSモデルを備えている。これらの例では、動作はさらに、発話の言語を示す言語識別子を受信する工程と、異なる各々の言語について各々が訓練された複数の言語固有ニューラルオラクルサーチNOSモデルの中から、言語固有ニューラルオラクルサーチNOSモデルを選択する工程と、を備えている。

20

【0010】

オプションとして、ニューラルオラクルサーチNOSモデルは多言語ニューラルオラクルサーチNOSモデルを備えていることができる。いくつかの実装では、外部言語モデルは、言語固有の外部言語モデルを備えている。これらの実装では、動作はさらに、発話の言語を示す言語識別子を受信する工程と、異なる各々の言語で各々が訓練された複数の言語固有の外部言語モデルの中から言語固有の外部言語モデルを選択する工程と、を備えている。ニューラルオラクルサーチNOSモデルは、2つの単方向長短期記憶(一方向LSTM)層を備えていることができる。いくつかの例では、音声認識モデルは、複数のコンフォーマ層を有しているコンフォーマエンコーダと、2つのLSTM層を有しているLSTMデコーダと、を備えているエンコーダデコーダアーキテクチャを備えている。

30

【0011】

本開示の1つまたは複数の実施態様の詳細は、添付の図面および以下の説明に記載されている。他の態様、特徴、および利点、は説明および図面、ならびに特許請求の範囲から明らかになるであろう。

【図面の簡単な説明】

【0012】

【図1】一例の音声認識モデルを実行する、音声環境の概略図。

【図2】図1の例示的な音声認識モデルの概略図。

【図3A】複数の言語固有のニューラルオラクル探索(NOS)モデルを使用する、例示的な再スコアリング処理の概略図。

40

【図3B】多言語ニューラルオラクルサーチNOSモデルを用いた、再スコアリング処理の一例を示す概略図。

【図4】自動音声認識のために多言語の再スコアリングモデルを使用する、方法の動作の配置例のフローチャート。

【図5】本明細書に記載のシステムおよび方法を実施するべく使用され得る、例示的なコンピューティング装置の概略図。

【発明を実施するための形態】

【0013】

様々な図面における同様の参照符号は、同様の要素を示す。

50

自動音声認識（ＡＳＲ）システムは、ユーザが話した内容のより正確なトランスクリプション（転写、書き起こし、音声記録）を提供するべく、ユーザ装置（ユーザデバイス、ユーザ機器）にますます普及している。しかし、自動音声認識ＡＳＲシステムは、ユーザが実際に話した内容を誤認識した不正確なトランスクリプションを生成する場合もありうる。いくつかの構成では、自動音声認識ＡＳＲシステムは、発声済みの発話に対してＮ個の最良候補仮説を生成するとともに、最良の候補仮説を最終トランスクリプションとして出力する。しかし、Ｎ個の最良候補仮説の構成は、１つの最良仮説の構成とで比較して、単語誤り率（ＷＥＲ）がほぼ５０％低くなる。したがって、いくつかの実装では、自動音声認識ＡＳＲシステムは、単語誤り率ＷＥＲを高めるべく追加情報を統合することによって、Ｎ個の最良候補仮説を再スコアリングする。このような再スコアリングの実装は、多言語音声環境における言語情報（すなわち、ユーザが話した言語識別子）に依存しており、わずかな単語誤り率ＷＥＲの改善しか提供しない。上述した課題によって、Ｎ個の最良候補仮説構成を使用する自動音声認識ＡＳＲシステムと、１個の最良候補仮説構成を使用する自動音声認識ＡＳＲシステムと、の間に単語誤り率ＷＥＲ性能のギャップがあることが明らかになった。

10

**【 0 0 1 4 】**

したがって、本明細書の実装は、対応する発話に対してＮ個の候補仮説を生成する再スコアリング処理を実行するとともに、最も可能性の高い候補仮説を選択して最終トランスクリプションとして出力するような、方法およびシステムに向けられている。特に、第１パス中、再スコアリング処理は、多言語音声認識モデルを使用することで、Ｎ個の候補仮説を生成する。その後、第２パス中、再スコアリング処理は、各候補仮説について、ニューラルオラクルサーチ（ＮＯＳ）モデルを使用することで各々の非正規化（正規化されていない）尤度スコアを生成したり、外部言語モデルのスコアを生成したり、候補仮説の事前統計をモデル化した単体（スタンドアロン、単独）スコアを生成したりする。以下で明らかになるように、ニューラルオラクルサーチＮＯＳモデルは、言語固有のニューラルオラクルサーチＮＯＳモデルまたは多言語ニューラルオラクルサーチＮＯＳモデルであってもよい。さらに、第２パス中、再スコアリング処理は、非正規化尤度スコア、外部言語モデルスコア、および単体スコア、に基づき各候補仮説の総合スコア（オーバーオールスコア）を生成する。再スコアリング処理は、総合スコアが最も高い候補仮説を、発話に対する最終トランスクリプションとして選択する。

20

30

**【 0 0 1 5 】**

図１は、音声（スピーチ、発話）環境１００の一例である。音声環境１００において、ユーザ１０４がユーザ装置１０などのコンピューティング装置とで対話する方法は、音声入力であってもよい。ユーザ装置１０は、音声環境１００内の１人または複数のユーザ１０４からの音（例えば、ストリーミング音響データ）を取り込む（キャプチャする、捕捉する）ように構成される。ここで、ストリーミング音響（オーディオ）データは、可聴クエリ、ユーザ装置１０に対するコマンド、またはユーザ装置１０によって捕捉された可聴コミュニケーション、として機能するユーザ１０４による発話（話し言葉、音声発話）１０６を指す場合がある。ユーザ装置１０の音声対応システムは、クエリに応答することによって、および／またはコマンドを１つまたは複数の下流アプリケーションによって実行／履行させることによって、クエリまたはコマンドをフィールド化することができる。

40

**【 0 0 1 6 】**

ユーザ装置１０は、ユーザ１０４に関連付けられているとともに音響データを受信することができる任意のコンピューティング装置に対応することができる。ユーザ装置１０のいくつかの例には、モバイル機器（例えば、携帯電話、タブレット、ラップトップ、など）、コンピュータ、ウェアラブル機器（例えば、スマートウォッチ）、スマート家電、モノのインターネット（ＩｏＴ）機器、車両エン터테인먼트システム、スマートディスプレイ、スマートスピーカ、などが含まれるが、これらに限定されない。ユーザ装置１０は、データ処理ハードウェア１２と、データ処理ハードウェア１２に通信するメモリハードウェア１４と、を備えている。メモリハードウェア１４は、データ処理ハードウェア

50

12によって実行されるとデータ処理ハードウェア12に1つまたは複数の動作を実行させる命令を記憶する。ユーザ装置10はさらに音声システム16を備えている。音声システム16は、音声環境100内の発話106を捕捉およびカバーして電気信号に変換するための音声キャプチャ装置（例えば、マイクロフォン）16、16aと、可聴音声信号を（例えば、ユーザ装置10からの出力音声データとして）伝達するための音声出力装置（例えば、スピーカ）16、16bと、を有している。ユーザ装置10は、図示の例では単一の音声（音響）キャプチャ装置16aを実装しているが、ユーザ装置10は、本開示の範囲から逸脱することなく音声キャプチャ装置16aのアレイを実装してもよい。それによって、アレイ内の1つまたは複数の音声キャプチャ装置16aは、ユーザ装置10上に物理的に存在しなくてもよいが、音声（音響）システム16に通信していてもよい。

10

#### 【0017】

音声環境100において、音声認識モデル（すなわち、自動音声認識ASRモデル）200を実装する自動音声認識（ASR）システム118は、ユーザ104のユーザ装置10上に、および/または、ネットワーク40を介してユーザ装置10に通信するリモートコンピューティング装置60（例えば、クラウドコンピューティング環境で実行される分散システムの1つまたは複数のリモートサーバ）上に、存在（常駐）する。自動音声認識ASRシステム118はまた、1つまたは複数の外部言語モデル310およびニューラルオラクルサーチ（NOS）モデル320を実装することができる。ユーザ装置10および/またはリモートコンピューティング装置（すなわち、リモートサーバ）60は、音声サブシステム108も備えている。音声サブシステム108は、ユーザ104によって発声済みの発話106であって、音声キャプチャ装置16aによってキャプチャ済みの発話106を、受信するとともに、当該発話106を、自動音声認識ASRシステム118によって処理可能な入力音響フレーム（音響フレーム110）に関連付けられた対応するデジタルフォーマットに変換するように構成されている。図示の例では、ユーザが各々の発話106を話している一方で、音声サブシステム108は当該発話106を、自動音声認識ASRシステム118に入力するための対応する音声データ（例えば、音響フレーム110）に変換する。その後、音声認識モデル200は、入力として、発話106に対応する音声データ（110）を受信するとともに、出力として、発話106の対応するトランスクリプション120（例えば、音声認識結果/仮説）を生成/予測する。以下にさらに詳細に説明するように、音声認識モデル200は、発話106によって指定済みのクエリが待ち時間に対してどの程度敏感であるか、および/またはユーザ104が待ち時間に対してどの程度寛容であるか、に応じて音声認識を実行する際に、音声認識モデル200が推論中に、先読み音声コンテキストの異なる継続時間を設定できるようにするべく、可変の先読み音声コンテキストでトレーニング済みのエンドツーエンドの音声認識モデル200を備えていることができる。例えば、ユーザ装置10上で実行されるデジタルアシスタントアプリケーション50は、発話106によって指定済みのクエリが待ち時間（レイテンシ）に対してどの程度敏感であるか、および/またはユーザ104が待ち時間に対してどの程度の許容範囲を持っているか、に応じて音声認識を要求することができる。

20

30

#### 【0018】

いくつかの実装では、音声認識モデル200は、N個の候補仮説204（図3A、図3B）を生成するための第1パス中、音声データ110に対してストリーミング音声認識を実行している。ニューラルオラクルサーチNOSモデル320および言語モデル310は、最終トランスクリプション（転写、書き起こし）120を生成するための第2パス中、N個の候補仮説204を再スコアリングする。例えば、図示の例では、音声認識モデル200は、（N個の候補仮説204に基づき）部分音声認識結果（すなわち、部分トランスクリプション）120、120aを生成するべく、音声データ110に対してストリーミング音声認識を実行している。言語モデル310およびニューラルオラクルサーチNOSモデル320は、最終音声認識結果（すなわち、最終トランスクリプション）120、120bを生成するべく、N個の候補仮説204を再スコアリングする。特に、音声認識モデル200は、部分音声認識結果120aを生成するべく、ゼロ（または約240ミリ秒

40

50



)に設定され得る先読み音声コンテキストを使用してもよい。したがって、入力発話(発話106)に対する最終音声認識結果120bは、入力発話に対する部分音声認識結果120aから遅れてもよい。

#### 【0019】

ユーザ装置10および/またはリモートコンピューティング装置60はまた、発話106のトランスクリプション120の表現を、ユーザ装置10のユーザ104に提示するように構成されたユーザインタフェース生成部109を実行する。以下にさらに詳細に説明するように、ユーザインタフェース生成部109は、第1時間1中に部分音声認識結果120aをストリーミング方式で表示することができる。その後、ユーザインタフェース生成部109は、第2時間2中に最終音声認識結果120bを表示することができる。いくつかの構成では、自動音声認識ASRシステム118から出力されたトランスクリプション120は、例えば、ユーザ装置10またはリモートコンピューティング装置60上で実行される自然言語理解(NLU)モジュールによって処理されることで、発話106によって指定済みのユーザコマンド/クエリを実行する。追加的または代替的に、テキスト音声合成システム(図示せず)(例えば、ユーザ装置10またはリモートコンピューティング装置60の任意の組み合わせ上で実行される)は、ユーザ装置10および/または別の機器による可聴出力用に、トランスクリプションを合成音声に変換してもよい。

10

#### 【0020】

図示の例では、ユーザ104がデジタルアシスタントアプリケーション50に通信している。デジタルアシスタントアプリケーション50は、ユーザ104とデジタルアシスタントアプリケーション50との間の会話を描写するべく、ユーザ装置10の画面上にデジタルアシスタントインタフェース18を表示する。この例では、ユーザ104はデジタルアシスタントアプリケーション50に、「今夜のコンサートは何時?」と質問する。ユーザ104からのこの質問は、音声キャプチャ装置16aによってキャプチャされるとともに、ユーザ装置10の音声システム16によって処理される発話106である。この例では、音声システム16は発話106を受信するとともに、当該発話106を、自動音声認識ASRシステム118に入力するための音響フレーム110に変換する。

20

#### 【0021】

この例を続けると、音声認識モデル200は、ユーザ104が話す発話106に対応する音響フレーム(すなわち、音声データ)110を受信しながら、音響フレーム110を符号化(エンコード)するだけでなく、さらに符号化(エンコード)済みの音響フレーム110を部分音声認識結果120aに復号化(デコード)する。第1時間1中、ユーザインタフェース生成部109は、デジタルアシスタントインタフェース18を介して、発話106の部分音声認識結果120aの表現を、単語、単語片、および/または個々の文字、が発声されるとすぐに画面上に現れるストリーミング方式で、ユーザ装置10のユーザ104に提示する。

30

#### 【0022】

第2パス中、および発話106に対応する全ての音響フレーム110が受信された後、自動音声認識ASRシステム118は、言語モデル310およびニューラルオラクルサーチNOSモデル320を使用することで、N個の候補仮説204のうちの各候補仮説204を再スコアリングするとともに、N個の候補仮説204の中から、発話106の正確なトランスクリプション(転写)120である可能性(尤度)が最も高い候補仮説204を選択する。第2時間2中、ユーザインタフェース生成部109は、デジタルアシスタントインタフェース18を介して、発話106の最終音声認識結果120bの表現を、ユーザ装置10のユーザ105に提示する。いくつかの実装では、ユーザインタフェース生成部109は、部分音声認識結果120aの表現を、最終音声認識結果120bの表現によって置き換える。例えば、最終音声認識結果120bは、先読み音声コンテキストを活用せずに生成済みの部分音声認識結果120aよりも、正確であると推定される。よって、最終的にトランスクリプション120として表示される最終音声認識結果120bは、部分音声認識結果120aにおいて誤認識されたかもしれない用語を、修正することができる

40

50

。この例では、ストリーミングの部分音声認識結果 120 a は、音声認識モデル 200 によって出力されているとともに、第 1 時間 1 にユーザ装置 10 の画面上に表示されており、低レイテンシ（低い待ち時間）に関連付けられている。よって、自分のクエリが処理されているという応答性をユーザ 104 に提供する。一方、第 2 時間 2 に画面上に表示される最終音声認識結果 120 b は、精度の点で音声認識品質を向上させるが、待ち時間（レイテンシ）が増大する。しかし、部分音声認識結果 120 a はユーザが発話 106 を話すときに表示されるので、最終認識結果を生成するとともに最終的に表示することに関連する高い待ち時間は、ユーザ 104 には気づかれない。

#### 【0023】

図 1 に示す例では、デジタルアシスタントアプリケーション 50 は、自然言語処理を使用することで、ユーザ 104 によって提起された質問に回答することができる。自然言語処理は、一般に、書かれた言語（例えば、部分音声認識結果 120 a および / または最終音声認識結果 120 b）を解釈するとともに、書かれた言語が何らかのアクションを促すかどうかを決定する、といった処理を指す。この例では、デジタルアシスタントアプリケーション 50 は自然言語処理を使用することで、ユーザ 104 からの質問が、ユーザのスケジュールに関するものであることを、より詳細にはユーザのスケジュール上のコンサートに関するものであることを、認識する。自然言語処理でこれらの詳細を認識することによって、自動アシスタントは、ユーザの問い合わせ（クエリ）に対して、応答 19 を返している。応答 19 は、「会場のドアは午後 6 時 30 分に開き、コンサートは午後 8 時に始まります」と述べる。いくつかの構成では、自然言語処理は、ユーザ装置 10 のデータ処理ハードウェア 12 に通信しているリモートサーバ 60 上で、行なわれる。

#### 【0024】

図 2 を参照すると、例示的なフレームアライメントベースのトランスデューサモデル 200 a は、対話型アプリケーションに関連する待ち時間制約を遵守するリカレントニューラルネットワーク - トランスデューサ（RNN - T）モデルアーキテクチャを備えている。リカレントニューラルネットワーク - トランスデューサ RNN - T モデルアーキテクチャの使用は例示的なものであり、フレームアライメントベースのトランスデューサモデル 200 は、特に、トランスフォーマ - トランスデューサおよびコンフォーマ - トランスデューサモデルアーキテクチャなどの、他のアーキテクチャを含み得る。リカレントニューラルネットワーク - トランスデューサ RNN - T モデル 200 は、小さな計算フットプリントを提供しているとともに、従来の自動音声認識 ASR アーキテクチャよりも少ないメモリ要件を利用するので、リカレントニューラルネットワーク - トランスデューサ RNN - T モデルアーキテクチャは、ユーザ装置 102 上で完全に音声認識を実行するのに適している（たとえば、リモートサーバとの通信は必要無い）。リカレントニューラルネットワーク - トランスデューサ RNN - T モデル 200 は、エンコーダネットワーク 210、予測ネットワーク 220、および結合（ジョイント）ネットワーク 230、を備えている。エンコーダネットワーク 210 は、従来の自動音声認識 ASR システムにおける音響モデル（AM：アコースティックモデル）にほぼ類似しているので、スタックされたロングショートターム（LSTM）層のリカレントネットワークを備えていることができる。例えば、エンコーダは、 $d$  次元特徴ベクトルのシーケンス（例えば、音響フレーム 110（図 1） $x = (X_1, X_2, \dots, X_T)$ 、ここで  $X_t \in \mathbb{R}^d$ 、 $R$  は白抜き文字）を読み取ることによって、各出力ステップで高次特徴表現を生成する。この高次特徴表現は、 $h_1^{enc}, \dots, h_T^{enc}$  と表記される。

#### 【0025】

同様に、予測ネットワーク 220 も LSTM ネットワーク（すなわち、LSTM デコーダ）であり、言語モデル（LM）のように、これまでの最終ソフトマックス層 240 によって出力された非空白記号シーケンス（すなわち、ラベル履歴）245（ $y_0, \dots, y_{u_i-1}$ ）を、密な表現  $p_{u_i}$  に変換する。最後に、リカレントニューラルネットワーク - トランスデューサ RNN - T モデルアーキテクチャでは、エンコーダおよび予測 / デコーダネットワーク 210、220 によって生成済みの表現同士は、結合ネットワーク 2

10

20

30

40

50

30によって結合される。予測ネットワーク220は、密な表現同士を処理する代わりに、ルックアップ（先読み）されたスパス（疎）な埋め込みを出力することによってレイテンシ（待ち時間）を改善するべく、埋込ルックアップテーブルによって置き換えられてもよい。次に、結合ネットワークは、次の出力記号に対する分布である、 $P(y_i | X_{t_i}, y_0, \dots, y_{u_i-1})$ を予測する。別の言い方をすれば、結合ネットワーク230は、各出力ステップ（例えば、時間ステップ）において、可能性のある音声認識仮説に対する確率分布を生成する。ここで、「可能性のある（ポシブル）音声認識仮説」は、指定済みの自然言語における記号（シンボル）/文字を各々表わす、出力ラベルのセットに対応する。例えば、自然言語が英語である場合、出力ラベルのセットは、27個のシンボル（記号）を備えていることができ、例えば、英語のアルファベットにおける26個の文字の各々に対する1個のラベルと、スペースを指定する1個のラベルと、を備えていることができる。従って、結合ネットワーク230は、所定の出力ラベルの集合（セット）の各々の発生の尤度を示す値の集合を出力することができる。この値のセットはベクトルとすることができるとともに、出力ラベルのセットにわたる確率分布を示すことができる。場合によっては、出力ラベルは書記素（例えば、個々の文字であり、潜在的には句読点や他の記号）であるが、出力ラベルのセットはそれほど限定されない。例えば、出力ラベルのセットは、書記素（グラフェムズ）に加えて、または書記素の代わりに、単語片および/または単語全体を備えていることができる。結合ネットワーク230の出力分布は、異なる出力ラベル同士の各々に対する事後確率値を備えていることができる。したがって、異なる書記素または他の記号を表わす100個の異なる出力ラベルが存在する場合、結合ネットワーク230の出力 $y_i$ は、出力ラベルごとに1つずつの確率値であるように、100個の異なる確率値を備えていることができる。次に、確率分布は、トランスクリプション120を決定するためのビーム探索処理（例えば、ソフトマックス層240による）において、候補となる正書法（オルソグラフィック）要素（例えば、書記素（グラフェムズ）、単語片（ワードピース）、および/または単語（ワード））を選択するとともにスコアを割り当てるべく、使用することができる。

10

20

#### 【0026】

ソフトマックス層240は、対応する出力ステップにおいてリカレントニューラルネットワーク-トランスデューサRNN-Tモデル200によって予測される次の出力記号（アウトプットシンボル）として、分布において最も高い確率を有している出力ラベル/記号を選択する任意の技術を採用することができる。このように、リカレントニューラルネットワーク-トランスデューサRNN-Tモデル200は、条件付き独立性の仮定を行わず、むしろ、各記号の予測は、音響だけでなく、これまでに出力されたラベルのシーケンスにも条件付けられる。リカレントニューラルネットワーク-トランスデューサRNN-Tモデル200は、出力記号が将来の音響フレーム110から独立していると仮定している。これによって、リカレントニューラルネットワーク-トランスデューサRNN-Tモデルをストリーミング方式で採用することができる。

30

#### 【0027】

いくつかの例では、リカレントニューラルネットワーク-トランスデューサRNN-Tモデル200のエンコーダネットワーク（例えば、音響エンコーダ）210は、コンフォーマ層のスタックを備えているコンフォーマベースのエンコーダを有しているエンコーダデコーダアーキテクチャである。ここで、各コンフォーマ層は、一連（シリーズ）の多頭自己アテンション層、深度ワイズ畳み込み層、およびフィードフォワード層、を備えている。いくつかの例では、コンフォーマベースのエンコーダは、17個のコンフォーマ層のスタックを備えていることができる。エンコーダネットワーク210は、多頭自己アテンション機構を有している、他のタイプのエンコーダを含んでもよい。たとえば、エンコーダネットワーク210は、トランスフォーマベースのエンコーダ、または軽量畳み込み（LConv）ベースのエンコーダであってもよい。また、エンコーダネットワーク210は、LSTM層のシーケンスを備えているRNNベースであってもよい。予測ネットワーク220は、2つの2048次元LSTM層を有しているLSTMデコーダであってもよ

40

50

く、各 L S T M 層には 6 4 0 次元の投影層も続く。あるいは、予測ネットワーク 2 2 0 は、L S T M 層の代わりに、変換器（トランスフォーマ）またはコンフォーマブロックのスタック、または埋込ルックアップテーブルを含んでもよい。最後に、結合ネットワーク 2 3 0 も 6 4 0 個の隠れユニットを持つことがある。ソフトマックス層 2 4 0 は、複数のトレーニングデータセットに含まれる全ての固有の単語片（ワードピース）または書記素を用いて生成済みの、統一された単語片または書記素セットで構成されてもよい。

#### 【 0 0 2 8 】

ここで図 3 A および図 3 B を参照すると、いくつかの実装では、リモートサーバ 6 0（図 1）は、第 1 パス 3 0 1 中に自動音声認識 A S R モデル 2 0 0 によって生成済みの N 個の候補仮説 2 0 4 を再スコアリングするための例示的な再スコアリング処理 3 0 0 を実行する。あるいは、ユーザ装置 1 0（図 1）は、リモートサーバ 6 0（図 1）に加えてまたはリモートサーバ 6 0（図 1）の代わりに、例示的な再スコアリング処理 3 0 0 を実行してもよい。再スコアリング処理 3 0 0 は、発話 1 0 6 に対応する音響フレーム 1 1 0 のシーケンス（ $X_1, X_2, \dots, X_T$ ）に対して、N 個の候補仮説 2 0 4、 $2 0 4 a \sim 2 0 4 n$ （ $H_1, H_2, \dots, H_N$ ）を生成する第 1 パス 3 0 1 を備えている。さらに、再スコアリング処理 3 0 0 は、N 個の候補仮説 2 0 4 のうちの各候補仮説 2 0 4 を、以下でさらに詳細に説明する追加情報源（インフォメーションリソース）を統合することによって再スコアリングする第 2 パス 3 0 2 を備えている。このように、第 2 パス 3 0 2 は、N 個の候補仮説 2 0 4 の中から、発話 1 0 6 の正確なトランスクリプションである可能性（尤度）が最も高い候補仮説 2 0 4 を選択するように構成されたシーケンス分類オブジェクトを備えている。

#### 【 0 0 2 9 】

特に、自動音声認識 A S R モデル 2 0 0 は、発話 1 0 6 に対応する音声データから抽出済みの音響フレーム 1 1 0 のシーケンスを受信する。第 1 パス 3 0 1 中、自動音声認識 A S R モデル 2 0 0 は、音響フレーム 1 1 0 のシーケンスを処理することで、発話 1 0 6 に対する N 個の候補仮説 2 0 4 を生成する。ここで、各候補仮説 2 0 4 は、発話 1 0 6 の候補トランスクリプション 1 2 0 に対応しており、各々の埋込ベクトルによって表わされる単語、サブワード、および/または書記素ラベルの各々のシーケンスによって表わされる。さらに各候補仮説 2 0 4 は、対応する候補仮説 2 0 4 の事前統計をモデル化する単体スコア 2 0 5 を備えている。すなわち、単体スコア 2 0 5 は、対応する候補仮説 2 0 4 が発話 1 0 6 の正確なトランスクリプションであるという信頼度（確信度、コンフィデンス）を示すことができる。単体スコア 2 0 5 の信頼度は、以前に実現された発話 1 0 6 の頻度（例えば、候補仮説 2 0 4 が以前に発声された回数）を示すこともある。

#### 【 0 0 3 0 】

自動音声認識 A S R モデル 2 0 0 は、任意の数の候補仮説 2 0 4 を生成してもよい（例えば、N は任意の整数値であってもよい）。いくつかの例では、自動音声認識 A S R モデル 2 0 0 は、予め定義されたパラメータに基づき、指定済みの数の候補仮説 2 0 4 を出力する。例えば、自動音声認識 A S R モデル 2 0 0 は、全ての発話 1 0 6 に対して 5 つの候補仮説 2 0 4（すなわち、 $N = 5$ ）を出力する。例えば、N 個の候補仮説 2 0 4 は、最も高い単体スコア 2 0 5 を有している N 個の候補仮説に関連する候補仮説の N ベストリストに対応することができる。他の例では、自動音声認識 A S R モデル 2 0 0 は、閾値を満たす単体スコア 2 0 5 を有している全ての候補仮説 2 0 4 を出力する。

#### 【 0 0 3 1 】

図示の例では、自動音声認識 A S R モデル 2 0 0 は、ユーザ 1 0 4 によって発声済みの発話 1 0 6 「プレイ\_\_ネクスト\_\_ソング」（次の歌を再生して）に対応する音響フレーム 1 1 0 のシーケンスを処理する。そして自動音声認識 A S R モデル 2 0 0 は、3 つの候補仮説 2 0 4（すなわち、 $N = 3$ ）を生成する。すなわち、複数の候補仮説 2 0 4 は、0.6 の単体スコア 2 0 5 を有している「プレイ\_\_ネクスト\_\_ソング」（次の歌を再生して）と、0.3 の単体スコア 2 0 5 を有している「ハイ\_\_ネクスト\_\_ロング」（やあ次の長い）と、および 0.8 の単体スコア 2 0 5 を有している「プレイ\_\_ネクスト\_\_ポン」（次の

10

20

30

40

50

ポン ( p o n g ) を再生して) と、を備えている。ここで、再スコアリング処理 3 0 0 は、候補仮説 2 0 4 「プレイ\_\_ネクスト\_\_ポン」( 次のポン ( p o n g ) を再生して) が最も高い単体スコア 2 0 5 を有しているため、当該候補仮説 2 0 4 「プレイ\_\_ネクスト\_\_ポン」を部分(部分的な、パーシャルな)トランスクリプション 1 2 0 a ( 図 1 ) として出力することができる。あるいは、再スコアリング処理 3 0 0 は、再スコアリング処理が最終トランスクリプションを生成するまで、部分トランスクリプションの出力を控えてもよい。特にこの例では、最も高い単体スコア 2 0 5 を有している候補仮説 2 0 4 は、ユーザ 1 0 4 によって発声済みの発話 1 0 6 の、不正確なトランスクリプションである。

#### 【 0 0 3 2 】

自動音声認識 A S R モデル 2 0 0 は、多言語(複数の言語)で話された発話 1 0 6 を認識するように構成された、多言語自動音声認識 A S R モデルであってもよい。すなわち、単一の自動音声認識 A S R モデル 2 0 0 は、第 1 言語で発話 1 0 6 を受信することによって第 1 言語で N 個の候補仮説 2 0 4 を生成したり、異なる第 2 言語で別の発話 1 0 6 を受信することによって第 2 言語で N 個の候補仮説 2 0 4 を生成したり、することができる。さらに、単一の自動音声認識 A S R モデルは、第 1 言語および第 2 言語の両方の用語を備えているコード混合音声を用意している発話 1 0 6 を受信することができる。このように、再スコアリング処理 3 0 0 は、単一の多言語自動音声認識 A S R モデル 2 0 0 を、多言語音声環境で実装することができる。

#### 【 0 0 3 3 】

いくつかの実装では、第 2 パス 3 0 2 は、第 1 パス 3 0 1 から N 個の候補仮説 2 0 4 を受け取るとともに、各候補仮説 2 0 4 の追加情報を統合することによって、対応する総合スコア 3 5 5 を生成することができる。総合スコア 3 5 5 は、各候補仮説 2 0 4 が正確なトランスクリプションであるかどうかについての、第 1 パス 3 0 1 からの単体スコア 2 0 5 よりも一層正確な信頼度を示すことがある。その後、第 2 パス 3 0 2 は、最も高い総合スコア 3 5 5 を有している候補仮説 2 0 4 を、トランスクリプション 1 2 0 ( すなわち、最終トランスクリプション 1 2 0 b ( 図 1 ) ) として選択することができる。

#### 【 0 0 3 4 】

より具体的には、第 2 パス 3 0 2 中に、外部言語モデル ( L M ) 3 1 0 は、N 個の候補仮説 2 0 4 を受信することで、各候補仮説 2 0 4 について各々の外部言語モデルスコア 3 1 5 を生成する。いくつかの実装では、外部言語モデル L M 3 1 0 は、リカレントニューラルネットワーク言語モデル R N N \_ L M を備えている。ここで、外部言語モデル L M 3 1 0 は、各々が特定の言語のテキストのみのデータ(すなわち、ペア(対)になっていないデータ)上でトレーニング済みの、複数の言語固有の外部言語モデル L M 3 1 0 、 3 1 0 a ~ 3 1 0 n を備えていることができる。このように、「外部言語モデル L M 3 1 0 」および「言語固有の外部言語モデル L M 3 1 0 」は、本明細書において互換的に使用される場合がある。したがって、各言語固有の外部言語モデル L M 3 1 0 は、各々の言語の発話 1 0 6 の外部言語モデルスコア(すなわち、言語モデルスコア) 3 1 5 を生成するように構成される。例えば、英語のテキストのみのデータでトレーニング済みの第 1 言語固有の外部言語モデル L M 3 1 0 、 3 1 0 a は、英語で発声済みの発話 1 0 6 の言語モデルスコア 3 1 5 を生成する。スペイン語のテキストのみのデータでトレーニング済みの第 2 言語固有の外部言語モデル L M 3 1 0 、 3 1 0 b は、スペイン語で発声済みの発話 1 0 6 の言語モデルスコア 3 1 5 を生成する。複数の外部言語モデル L M 3 1 0 は、任意の数の言語でトレーニングすることができる。ここで各外部言語モデル L M 3 1 0 は、異なる各々の言語のテキストのみのデータでトレーニングされる。

#### 【 0 0 3 5 】

したがって、外部言語モデル L M 3 1 0 は、発話 1 0 6 の言語を示す言語識別子 1 0 7 を受信することで、複数の言語固有の外部言語モデル L M 3 1 0 の中から、発話 1 0 6 の言語に対応する言語固有の外部言語モデル L M 3 1 0 を選択することができる。別の言い方をすれば、再スコアリング処理 3 0 0 は、言語識別子 1 0 7 に基づき、言語固有の外部言語モデル L M 3 1 0 を選択することができる。いくつかの例では、自動音声認識 A S R

10

20

30

40

50

モデル 200 は、発話 106 の音響フレーム 110 のシーケンスの処理に基づき、言語識別子 107 を決定する。他の例では、自動音声認識 ASR モデル 200 は、外部ソースから言語識別子 107 を取得する。例えば、ユーザは、特定の言語用に自動音声認識 ASR モデルを設定（構成）することができる。他の例では、自動音声認識 ASR モデル 200 は、発話 106 を発声したユーザ 104 のアイデンティティを決定するとともに、識別済みのユーザ 104 に関連付けられた言語に基づき言語識別子 107 を識別することができる。

#### 【0036】

したがって、第 2 パス 302 中、再スコアリング処理 300 は、言語識別子 107 に基づき発話 106 の言語に対応する外部言語モデル LM 310 を選択するとともに、各候補仮説（仮説候補）204 の言語モデルスコア 315 を生成する。言語モデルスコア 315 は、候補仮説 204 のうちの一連（シーケンス）の仮説用語が、ユーザ 104 によって発声される尤度（可能性）を示す。例えば、外部言語モデル LM 310 は、候補仮説 204 「今日の天気は何？」（ワット\_\_イズ\_\_ザ\_\_ウェザー\_\_トゥデイ）に対して、候補仮説 204 「天気は何？ホーレイ」（ワット\_\_イズ\_\_ザ\_\_ウェザー\_\_ホーレイ（h o o r a y、万歳））とは対照的に、より高い言語モデルスコア 315 を生成する。特に、外部言語モデル LM 310 は、「今日の天気は何？」（ワット\_\_イズ\_\_ザ\_\_ウェザー\_\_トゥデイ）に対して、より高い言語モデルスコア 315 を生成する。なぜなら、この仮説用語のシーケンスは、「天気は何？ホーレイ」（ワット\_\_イズ\_\_ザ\_\_ウェザー\_\_ホーレイ（h o o r a y、万歳））よりも頻繁に、テキストのみのトレーニング（学習）データに含まれている可能性があるのである。

#### 【0037】

例示的な再スコアリング処理 300 はまた、ニューラルオラクルサーチ（NOS）モデル 320 を備えている。ニューラルオラクルサーチ NOS モデル 320 は、N 個の候補仮説 204、音響フレーム 110 のシーケンス、およびラベル履歴 245（例えば、以前に出力された単語、単語片、および/または書記素）、を受信する。ラベル履歴 245（y<sub>0:i-1</sub>）は、自動音声認識 ASR モデル 200、（例えば、再スコアラ 350 を介した）再スコアリング処理 300 の第 2 パス 302、またはそれらの組み合わせ、によって出力され得る。いくつかの例では、ラベル履歴 245 は、ユーザ 104 によって発声済みの以前の発話 106 のトランスクリプション（文字起こし）を備えている。例えば、ユーザ 104 は、「明日はどうする？」（ワット\_\_アバウト\_\_トゥデイ？）という現在の発話 106 に対するラベル履歴 245 を表わす、「今日は私に何か会議がある？」（ドゥ\_\_アイ\_\_ハブ\_\_エニ\_\_ミーティング\_\_トゥデイ？）という以前の発話 106 を以前に発声したことがあるかもしれない。他の例では、ラベル履歴 245 は、発話の現在のラベルに先行する、全ての用語を備えている。例えば、発話 106 の「プレイ\_\_マイ\_\_プレイリスト」（私の再生リストを再生して）について、ラベル履歴 245 は、発話 106 の現在の用語（例えば、次の仮説用語）が「プレイリスト」（再生リスト）である、用語「プレイ\_\_マイ」（私の再生して）に対応することができる。任意選択で、ニューラルオラクルサーチ NOS モデル 320 は、ユーザ 104 によって発声済みの発話 106 の言語を示す言語識別子 107 を受信することができる。

#### 【0038】

図 3A は、複数の言語固有ニューラルオラクルサーチ NOS モデル 320 S、320 S<sub>a</sub> ~ 320 S<sub>n</sub> を備えている再スコアリング処理 300、300 a の一例を示す。ここで、各言語固有ニューラルオラクルサーチ NOS モデル 320 S は、特定の言語のペアワイズデータ（すなわち、書き起こされた音響トレーニングデータ）上でトレーニングされる。したがって、第 2 パス 302 中、再スコアリング処理 300 は、言語識別子 107 に基づき、複数の言語固有ニューラルオラクルサーチ NOS モデル 320 S の中から、発話 106 の言語に対応する言語固有ニューラルオラクルサーチ NOS モデル 320 S を選択する。このように、例示的な再スコアリング処理 300 a は、正しい言語固有ニューラルオラクルサーチ NOS モデル 320 S を選択するべく、言語識別子 107 が利用可能にされ

10

20

30

40

50

ていると仮定する。

【0039】

代替的に、図3Bは、多言語ニューラルオラクルサーチNOSモデル320、320Mを備えている例示的な再スコアリング処理300、300bを示す。この例では、多言語ニューラルオラクルサーチNOSモデル320Mは、任意の数の言語のペアワイズデータ（すなわち、書き起こされた音声トレーニングデータ）上でトレーニングされる。したがって、例示的な再スコアリング処理300bは、多言語音声環境において、単一の多言語ニューラルオラクルサーチNOSモデル320Mを実装することができる。注目すべきことに、発話106の言語に関連付けられた言語固有のニューラルオラクルサーチNOSモデル320S（図3Aを参照して説明したような）の選択が要求されないので、例示的な再スコアリング処理300bは、任意の言語識別子107の使用を必要にしない。したがって、発話106は、2つ以上の言語にまたがる音声のコードミキシングを備えている多言語発話を備えていることができる。本明細書で使用されるように、ニューラルオラクルサーチNOSモデル320は、再スコアリング処理300aが言語識別子107に基づき選択する言語固有ニューラルオラクルサーチNOSモデル320S（図3A）か、または多言語ニューラルオラクルサーチNOSモデル（図3B）か、のいずれかを備えていることができる。

10

【0040】

図3Aおよび図3Bを引き続き参照すると、ニューラルオラクルサーチNOSモデル320は、ラベル履歴245を与えられた次のラベル $Y_i$ を予測する、事前モデルを備えている。すなわち、事前モデルは、以前に認識された単語、単語片、および/または書記素、に基づき次のラベルの事前スコアを予測する。ニューラルオラクルサーチNOSモデル320の事前モデルは、1層あたり512ユニットである、2層の一方向LSTMを備えていることができる。事前モデルは、ラベル付けされた音声トレーニングデータとクロスエントロピー損失とを使用することで、トレーニング（学習）する。さらに、ニューラルオラクルサーチNOSモデル320は、ラベル履歴245と、第1パス301からの音響フレーム110のシーケンスと、をラベル同期方式で組み合わせることによって、事後スコアを予測する事後モデルを備えている。ニューラルオラクルサーチNOSモデル320の事後モデルは、1層あたり512ユニットである2層の一方向LSTMと、1層あたり128ユニットである2層のラベル同期アテンション（注意）メカニズムと、を備えていることができる。事後モデルは、ラベル履歴245と、音響フレーム110のシーケンスと、が与えられると、次のラベル $Y_i$ を予測するべく、ラベル付き音響トレーニングデータとクロスエントロピー損失とでトレーニング（学習）する。ニューラルオラクルサーチNOSモデル320は、トークンレベルの事前スコアと、トークンレベルの事後スコアと、を合計することで非正規化尤度スコア325を生成する。このように、非正規化尤度スコア325は、以下のような和によって表わされるシーケンスレベルのスコアである。

20

30

【0041】

【数1】

40

50

$$\begin{aligned}
S_{\theta_1}(X|Y = y_{0:U}) &= \sum_{i=0}^U \phi(X, y_{0:i-1}|Y_i = y) \\
&\propto \sum_{i=0}^U \log P(Y_i = y|X, y_{0:i-1}) \\
&\quad - \sum_{i=0}^U \log P(Y_i = y|y_{0:i-1}) \\
&= \log P(Y|X) - \log P(Y) \\
&\propto \log P(X|Y) \tag{1}
\end{aligned}$$

10

【0042】

式1において、 $S_{\theta_1}$ は非正規化尤度スコア325を表わす。

20

再スコアラ-350は、N個の候補仮説204のうちの各候補仮説204について、単体(スタンドアロン、単独)スコア205、言語モデルスコア315、および非正規化尤度スコア325、を受け取るとともに、各々の総合スコア355を生成する。特に、再スコアラ-350は、単体スコア205、言語モデルスコア315、および非正規化尤度スコア325、の任意の組み合わせに基づき、各候補仮説204の総合スコア355を生成する。いくつかの例では、再スコアラ-350は、単体(スタンドアロン、単独)スコア205、言語モデルスコア315、および非正規化尤度スコア325、を線形に合計することで、以下の式で表わされるシーケンスレベルの総合スコア355を決定する。

【0043】

【数2】

30

$$\tilde{P}(\text{Oracle} = i|X, H_{1:N}) = \frac{\exp(\text{Score}(H_i, X))}{\sum_j \exp(\text{Score}(H_j, X))} \tag{2}$$

$$\text{Score}(H_i, X) = \lambda_1 S_{\theta_1}(X|H_i) + \lambda_2 S_{\theta_2}(H_i) + S_{\theta_3}(i) \tag{3}$$

【0044】

40

式(3)において、 $S_{\theta_1}$ は、非正規化尤度スコア325を表している。 $S_{\theta_2}$ は外部言語モデルスコア315を表している。 $S_{\theta_3}$ は単体スコア205を表わす。トレーニング中に再スコアラ-350のモデルパラメータを最適化するべく、再スコアリング処理300は、事後スコアとシーケンスレベルのグラウンドトゥルス分布との間の、クロスエントロピーオブジェクトを使用する。いくつかの例では、トレーニング処理は、全(トータル)グラウンドトゥルス分布をグラウンドトゥルーストランスクリプションに割り当てる一方で、他の全ての候補仮説をゼロに割り当てる。他の例では、トレーニング処理は、最良の候補仮説(すなわち、グラウンドトゥルーストランスクリプション)未満の単語誤り率(WER)を有している全ての候補仮説にわたって一様に、全グラウンドトゥルス分布を割り当てる。さらに他の例では、トレーニング処理は、各候補仮説とグラウンド

50



トゥルーストランスクリプションとの間の負の編集距離に、ソフトマックス関数を適用する。

#### 【 0 0 4 5 】

その後、再スコアラ-350は、N個の候補仮説204の中から、最も高い総合スコア355を有している候補仮説204を、発話106の最終トランスクリプション120として選択する。図示の例では、候補仮説204には、0.9の総合スコア355を有している「プレイ\_\_ネクスト\_\_ソング」（次の歌を再生して）、0.3の総合スコア355を有している「ヘイ\_\_ネクスト\_\_ロング」（やあ次の長い）、および0.5の総合スコア355を有している「プレイ\_\_ネクスト\_\_ポン」（次のポン（pong）を再生して）、が含まれる。この例を続けると、再スコアラ-350は、0.9という最も高い総合スコア355を有している「プレイ\_\_ネクスト\_\_ソング」（次の歌を再生して）（実線のボックスで示される）の候補仮説204を、トランスクリプション120（例えば、最終トランスクリプション120b（図1））として選択する。注目すべきは、最も高い単体スコア205（すなわち、正しいトランスクリプションである可能性）を持つ候補仮説204は、正しい候補仮説204ではないということ。そして、最も高い総合スコア355を持つ候補仮説が、第2パス302からの正しいトランスクリプションであるということである。

10

#### 【 0 0 4 6 】

図4は、自動音声認識のために多言語の再採点（再スコアリング）モデルを使用するコンピュータ実装方法400の動作（操作、オペレーション）の、例示的な配置（アレンジメント、構成）のフローチャートである。動作402において、方法400は、発話106に対応する音声データから抽出済みの音響フレーム110のシーケンスを受信する工程を備えている。動作404において、第1パス301中、方法400は、音響フレーム110のシーケンスを処理して、多言語音声認識モデル（すなわち、自動音声認識ASRモデル）200を使用することで、発話106に対するN個の候補仮説204、204a~204nを生成する工程を備えている。第2パス302中、N個の候補仮説204のうちの各候補仮説204について、方法400は動作406~412を実行する。動作406において、方法400は、ニューラルオラクルサーチNOSモデル320を用いて、各々の非正規化尤度スコア325を生成する工程を備えている。ここで、ニューラルオラクルサーチNOSモデル320は、音響フレーム110のシーケンスと、対応する候補仮説204と、に基づき非正規化尤度スコア325を生成する。動作408において、方法400は、言語モデル310を用いて各々の外部言語モデルスコア315を生成する工程を備えている。動作410において、方法400は、第1パス301中に生成済みの対応する候補仮説204の事前統計をモデル化する、単体スコア205を生成する工程を備えている。動作412において、方法400は、非正規化尤度スコア325、外部言語モデルスコア315、および単体スコア205、に基づき候補仮説255の各々の総合スコア355を生成する工程を備えている。動作414において、方法400は、N個の候補仮説204の中から最も高い（最高の）各々の総合スコア355を有している候補仮説204を、発話106の最終トランスクリプション120として選択する工程を備えている。

20

30

#### 【 0 0 4 7 】

図5は、本書に記載されるシステムおよび方法を実施するべく使用され得る例示的なコンピューティング装置500の概略図である。コンピューティング装置500は、ラップトップ、デスクトップ、ワークステーション、パーソナルデジタルアシスタント、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、などの様々な形態のデジタルコンピュータを表わすことを意図している。ここに示された構成要素、それらの接続および関係、ならびにそれらの機能は、例示的なものであることのみを意図しているのであり、本書で説明および/または特許請求される発明の実施を制限することを意図していない。

40

#### 【 0 0 4 8 】

コンピューティング装置500は、プロセッサ510と、メモリ520と、ストレージ

50

デバイス 530 と、メモリ 520 および高速拡張ポート 550 に接続する高速インタフェース/コントローラ 540 と、および低速バス 570 およびストレージデバイス 530 に接続する低速インタフェース/コントローラ 560 と、を備えている。各構成要素 510、520、530、540、550、および 560、はさまざまなバスを使用して相互接続されており、共通のマザーボード上に、または適切な他の方法で、実装することができる。プロセッサ 510 は、高速インタフェース 540 に結合されたディスプレイ 580 などの外部入出力デバイスにグラフィカルユーザインタフェース (GUI) のためのグラフィカル情報を表示するべく、メモリ 520 または記憶デバイス 530 に記憶された命令を備えている、コンピューティング装置 500 内で実行するための命令を処理することができる。他の実施態様では、複数のプロセッサおよび/または複数のバスは、複数のメモリおよびメモリの種類とともに、適宜、使用されてもよい。また、複数のコンピューティング装置 500 が接続されることで、各装置 (デバイス) が必要な動作 (操作) の一部を提供してもよい (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)。

10

#### 【0049】

メモリ 520 は、コンピューティング装置 500 内で情報を非遷移的 (非一時的、非一過性) に記憶する。メモリ 520 は、コンピュータ可読媒体、揮発性メモリユニット (複数可)、または不揮発性メモリユニット (複数可)、であってもよい。不揮発性メモリ 520 は、コンピューティング装置 500 によって使用されるプログラム (例えば、命令のシーケンス) またはデータ (例えば、プログラム状態情報) を一時的または永続的に記憶するべく使用される物理的デバイスであってもよい。不揮発性メモリの例としては、フラッシュメモリ、読み出し専用メモリ (ROM) / プログラマブル読み出し専用メモリ (PROM) / 消去可能プログラマブル読み出し専用メモリ (EPROM) / 電子消去可能プログラマブル読み出し専用メモリ (EEPROM) (例えば、ブートプログラムなどのファームウェアに通常使用される) が挙げられるが、これらに限定されない。揮発性メモリの例としては、ランダムアクセスメモリ (RAM)、ダイナミックランダムアクセスメモリ (DRAM)、スタティックランダムアクセスメモリ (SRAM)、相変化メモリ (PCM)、およびディスクやテープ、などがあるが、これらに限定されるものではない。

20

#### 【0050】

記憶装置 530 は、コンピューティング装置 500 に大容量記憶装置を提供することができる。いくつかの実施態様において、記憶装置 530 は、コンピュータ読み取り可能な媒体である。様々な異なる実装において、記憶装置 530 は、フロッピー (登録商標) ディスク装置、ハードディスク装置、光ディスク装置、またはテープ装置、フラッシュメモリまたは他の同様のソリッドステートメモリデバイス、またはストレージエリアネットワークまたは他の構成のデバイスを備えているデバイスのアレイ、であってもよい。追加の実施態様において、コンピュータプログラム製品は、情報キャリアに具体化される。コンピュータプログラム製品は、実行されると上述したような 1 つまたは複数の方法を実行する命令を備えている。情報キャリアは、メモリ 520、記憶装置 530、またはプロセッサ 510 上のメモリ、などのコンピュータ可読媒体または機械可読媒体である。

30

#### 【0051】

高速コントローラ 540 は、コンピューティング装置 500 の帯域幅集約的な動作を管理しており、低速コントローラ 560 は、帯域幅集約的ではない (低帯域幅集約的な) 動作を管理する。このような任務の割り当ては、例示的なものに過ぎない。一部の实装では、高速コントローラ 540 は、メモリ 520、ディスプレイ 580 (例えば、グラフィックプロセッサまたはアクセラレータを介して)、および高速拡張ポート 550 に結合されているとともに、様々な拡張カード (図示せず) を受け入れることができる。いくつかの実装では、低速コントローラ 560 は、ストレージデバイス 530 および低速拡張ポート 590 に結合される。低速拡張ポート 590 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)) を備えているとともに、キーボード、ポインティングデバイス、スキャナ、など

40

50

の1つまたは複数の入出力デバイスに、またはネットワークアダプタを介してスイッチャー ルータなどのネットワークデバイスに、結合される。

【0052】

コンピューティング装置500は、図示のように、多数の異なる形態で実装されてもよい。例えば、標準サーバ500aとして、またはそのようなサーバ500aのグループ内の複数倍(回)、ラップトップコンピュータ500bとして、またはラックサーバシステム500cの一部として、実装することができる。

【0053】

本明細書で説明するシステムおよび技術の様々な実装は、デジタル電子回路および/または光回路、集積回路、特別に設計されたASIC(特定用途向け集積回路)、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/またはそれらの組み合わせ、で実現することができる。これらの様々な実装は、特殊目的であっても汎用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、からデータおよび命令を受信したり、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、にデータおよび命令を送信したり、するように結合された少なくとも1つのプログラマブルプロセッサを備えているプログラマブルシステム上で実行可能および/または解釈可能な1つまたは複数のコンピュータプログラムにおける実装を備えていることができる。

【0054】

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、プログラマブルプロセッサ用の機械命令を備えており、高レベルの手続き型および/またはオブジェクト指向プログラミング言語、および/またはアセンブリ/機械言語で実装することができる。本明細書で使用される場合、「機械可読媒体」および「コンピュータ可読媒体」という用語は、機械命令を機械可読信号として受信する機械可読媒体を含めて、機械命令および/またはデータをプログラマブルプロセッサに提供すべく使用される、任意のコンピュータプログラム製品、非一過性コンピュータ可読媒体、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、プログラマブルロジックデバイス(PLD))を指す。「機械可読信号」という用語は、機械命令および/またはデータをプログラマブルプロセッサに提供すべく使用される、あらゆる信号を指す。

【0055】

本明細書で説明する処理および論理フローは、データ処理ハードウェアとも呼ばれる1つまたは複数のプログラマブルプロセッサが、1つまたは複数のコンピュータプログラムを実行することで、入力データに対して動作するとともに出力を生成することによって機能を実行することで実行することができる。処理および論理フローは、特殊用途の論理回路、例えばFPGA(フィールドプログラマブルゲートアレイ)またはASIC(特定用途向け集積回路)によっても実行できる。コンピュータプログラムの実行に適したプロセッサには、一例として、汎用および特殊用途のマイクロプロセッサが、およびあらゆる種類のデジタルコンピュータの任意の1つまたは複数のプロセッサが、含まれる。一般に、プロセッサは、読み取り専用メモリまたはランダムアクセスメモリ、あるいはその両方から命令とデータを受け取る。コンピュータの本質的な要素は、命令を実行するためのプロセッサと、命令やデータを格納するための1つまたは複数のメモリ装置と、である。一般に、コンピュータは、データを格納するための1つまたは複数の大容量記憶装置、例えば磁気ディスク、光磁気ディスク、光ディスク、などからもデータを受け取るか、それら記憶装置にデータを転送するか、あるいはその両方、を行なうようにそれら記憶装置に動作可能に結合されている。しかし、コンピュータがそのような装置を備えている必要はない。コンピュータプログラム命令およびデータを記憶するのに適したコンピュータ可読媒体には、あらゆる形態の不揮発性メモリ、媒体およびメモリデバイスが含まれ、例えば、半導体メモリデバイス、例えば、EPROM、EEPROMおよびフラッシュメモリデバイス;磁気ディスク、例えば、内蔵ハードディスクまたはリムーバブルディスク;光磁気デ

10

20

30

40

50

ディスク；およびCDROMおよびDVD-ROMディスクが含まれる。プロセッサとメモリは、特殊用途の論理回路によって補足されるか、または特殊用途の論理回路に組み込まれる。

【0056】

ユーザとの相互作用を提供するべく、本開示の1つまたは複数の態様は、ユーザに情報を表示するためのディスプレイデバイス、例えばCRT（陰極線管）、LCD（液晶ディスプレイ）モニタ、またはタッチスクリーンと、任意選択で、ユーザがコンピュータに入力を提供することができる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と、を有しているコンピュータ上で実施することができる。他の種類のデバイスも同様に、ユーザとの対話を提供するべく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、視覚フィードバック、聴覚フィードバック、または触覚フィードバック、などの任意の形式の感覚フィードバックとすることができる。ユーザからの入力は、音響入力、音声入力、または触覚入力、を備えている任意の形式で受信することができる。さらに、コンピュータは、例えば、ウェブブラウザから受信した要求に応じて、ユーザのクライアントデバイス上のウェブブラウザにウェブページを送信することによって、ユーザが使用するデバイスにドキュメントを送信したり、デバイスからドキュメントを受信したり、することによってユーザと対話することができる。

10

【0057】

多くの実施態様を説明してきた。それにもかかわらず、本開示の精神および範囲から逸脱することなく、様々な変更がなされ得ることが理解されるであろう。従って、他の実施態様も以下の特許請求の範囲に含まれる。

20

【図面】

【図1】

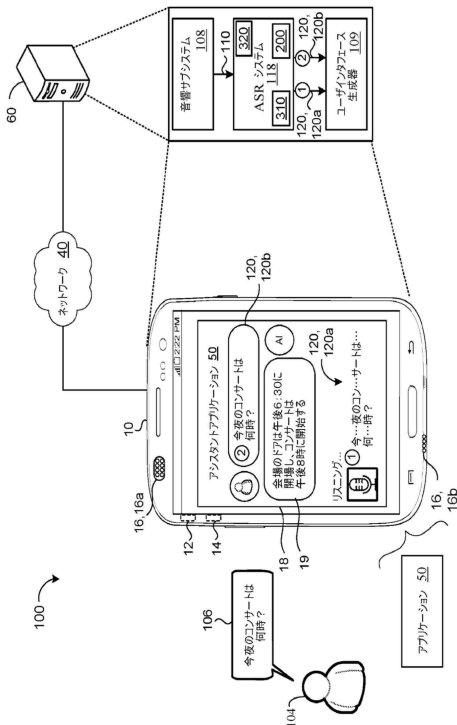


FIG. 1

【図2】

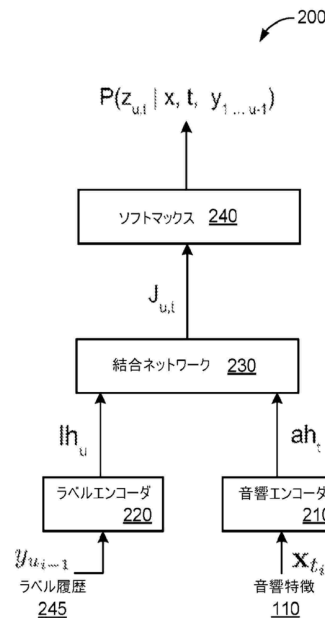


FIG. 2

30

40

50

【 図 3 A 】

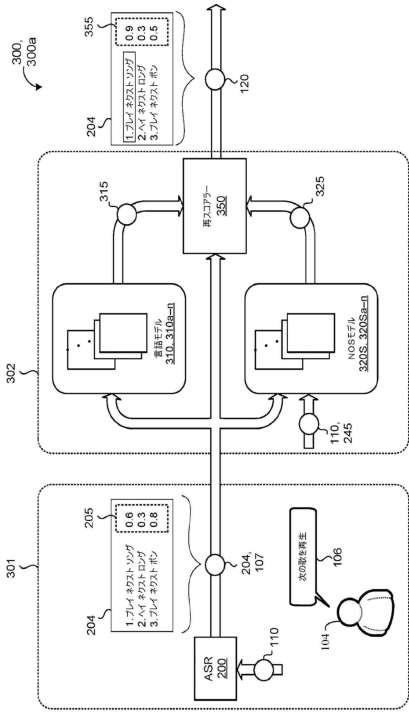


FIG. 3A

【 図 3 B 】

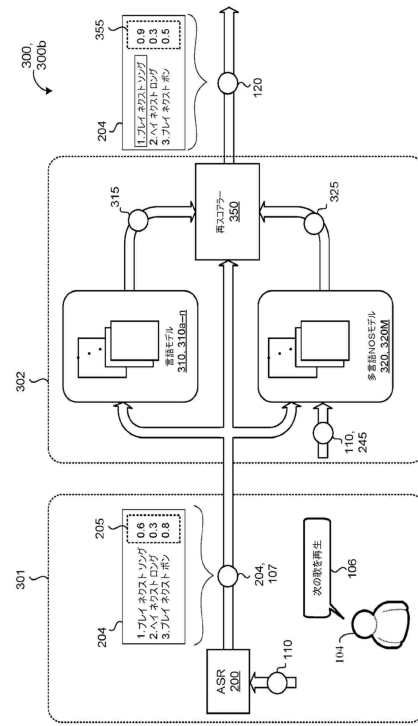


FIG. 3A

10

20

【 図 4 】

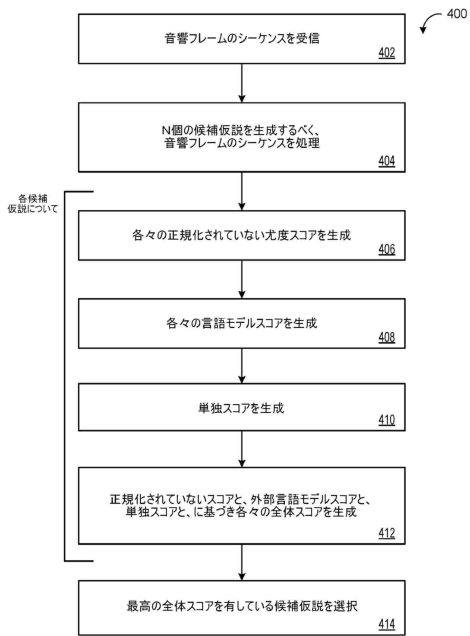


FIG. 4

【 図 5 】

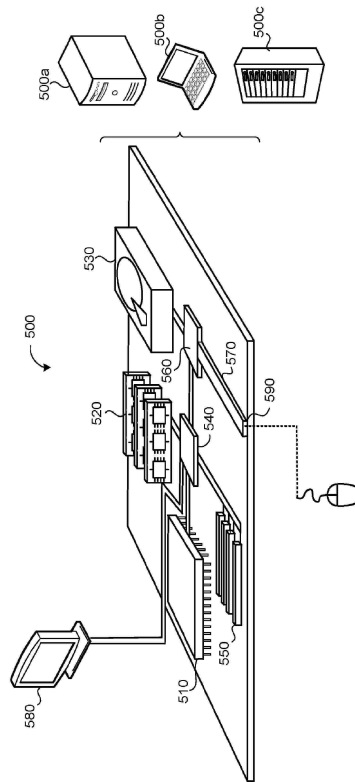


FIG. 5

30

40

50

---

フロントページの続き

- (72)発明者 チェン、トンジョウ  
アメリカ合衆国 9 4 0 4 3 カリフォルニア州 マウンテン ビュー アンフィシアター パークウ  
エイ 1 6 0 0
- (72)発明者 ヴァリアニ、エフサン  
アメリカ合衆国 9 4 0 4 3 カリフォルニア州 マウンテン ビュー アンフィシアター パークウ  
エイ 1 6 0 0
- (72)発明者 ラマバドラン、ブバナ  
アメリカ合衆国 1 0 5 4 9 ニューヨーク州 マウント キスコ ビクトリア ドライブ 4 0 0 2
- (72)発明者 ハガニ、パリサ  
アメリカ合衆国 9 4 0 4 3 カリフォルニア州 マウンテン ビュー アンフィシアター パークウ  
エイ 1 6 0 0
- (72)発明者 モレノ メンヒバル、ペドロ ジェイ .  
アメリカ合衆国 9 4 0 4 3 カリフォルニア州 マウンテン ビュー アンフィシアター パークウ  
エイ 1 6 0 0