

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7092953号
(P7092953)

(45)発行日 令和4年6月28日(2022.6.28)

(24)登録日 令和4年6月20日(2022.6.20)

(51)国際特許分類	F I
G 1 0 L 15/32 (2013.01)	G 1 0 L 15/32 2 1 0 E
G 1 0 L 15/16 (2006.01)	G 1 0 L 15/16
G 1 0 L 15/08 (2006.01)	G 1 0 L 15/08 2 0 0 Z

請求項の数 20 (全20頁)

(21)出願番号	特願2021-564950(P2021-564950)	(73)特許権者	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 94043 マウンテン ビュー アンフィシ アター パークウェイ 1600 1600 Amphitheatre P arkway 94043 Mounta in View, CA U.S.A.
(86)(22)出願日	令和2年4月28日(2020.4.28)	(74)代理人	100142907 弁理士 本田 淳
(65)公表番号	特表2022-523883(P2022-523883 A)	(72)発明者	フー、キー アメリカ合衆国 94043 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 1600 最終頁に続く
(43)公表日	令和4年4月26日(2022.4.26)		
(86)国際出願番号	PCT/US2020/030321		
(87)国際公開番号	WO2020/226948		
(87)国際公開日	令和2年11月12日(2020.11.12)		
審査請求日	令和4年1月24日(2022.1.24)		
(31)優先権主張番号	62/842,571		
(32)優先日	令和1年5月3日(2019.5.3)		
(33)優先権主張国・地域又は機関	米国(US)		
特許法第30条第2項適用	令和1年6月21日にウェブサイトのアドレス https://arxiv.org 最終頁に続く		

(54)【発明の名称】 エンドツーエンドモデルによる多言語音声認識のための音素に基づく文脈解析

(57)【特許請求の範囲】

【請求項1】

データ処理ハードウェア(610)において、第1言語のネイティブスピーカ(110)によって話される発話(106)を符号化する音声データを受け取る工程と、前記データ処理ハードウェア(610)において、前記第1言語とは異なる第2言語の1つまたは複数の用語を備えているバイアス用語リスト(105)を受け取る工程と、前記データ処理ハードウェア(610)において、音声認識モデル(200)を用いて、前記音声データから得られる音響特徴(104)を処理して、前記第1言語の語句と、対応する音素シーケンスとの両方に対する音声認識スコアを生成する工程と、前記データ処理ハードウェア(610)によって、前記バイアス用語リスト(105)内の前記1つまたは複数の用語に基づき、前記音素シーケンスに対する前記音声認識スコアを再スコアリングする工程と、前記データ処理ハードウェア(610)によって、前記語句に対する前記音声認識スコアと、前記音素シーケンスに対する再スコアリングされた音声認識スコアとを用いて、復号グラフ(400)を実行して、前記発話(106)に対する転写(116)を生成する工程と、を備えている方法(500)。

【請求項2】

前記音素シーケンスに対する前記音声認識スコアを再スコアリングする工程は、バイアス有限状態変換器(FST)を使用して、前記音素シーケンスに対する前記音声認識スコア

を再スコアリングする工程を備えている、
請求項 1 に記載の方法 (5 0 0)。

【請求項 3】

前記方法はさらに、

前記データ処理ハードウェア (6 1 0) によって、前記バイアス用語リスト (1 0 5) の各用語を、前記第 2 言語の対応する音素シーケンスにトークン化する工程と、

前記データ処理ハードウェア (6 1 0) によって、前記第 2 言語における各対応する音素シーケンスを、前記第 1 言語における対応する音素シーケンスに写像する工程と、

前記データ処理ハードウェア (6 1 0) によって、前記第 1 言語における各対応する音素シーケンスに基づき、前記バイアス有限状態変換器 (3 0 0) を生成する工程と、

を備えている、請求項 2 に記載の方法 (5 0 0)。

10

【請求項 4】

前記音声認識モデル (2 0 0) は、エンドツーエンド語句 - 音素モデル (2 0 0) を備えている、

請求項 1 ~ 3 のいずれか一項に記載の方法 (5 0 0)。

【請求項 5】

前記エンドツーエンド語句 - 音素モデル (2 0 0) は、リカレントニューラルネットワーク - 変換器 (R N N - T) を備えている、

請求項 4 に記載の方法 (5 0 0)。

【請求項 6】

前記復号グラフ (4 0 0) の実行中に、前記復号グラフ (4 0 0) は、前記バイアス用語リスト (1 0 5) 内の前記 1 つまたは複数の用語のいずれかを有利にするように、前記転写 (1 1 6) をバイアスする、

請求項 1 ~ 5 のいずれか一項に記載の方法 (5 0 0)。

20

【請求項 7】

前記音声認識モデル (2 0 0) は、前記第 1 言語のみの学習発話で学習される、

請求項 1 ~ 6 のいずれか一項に記載の方法 (5 0 0)。

【請求項 8】

前記バイアス用語リスト (1 0 5) 内の用語のいずれも、前記音声認識モデル (2 0 0) を学習するために使用されなかった、

請求項 1 ~ 7 のいずれか一項に記載の方法 (5 0 0)。

30

【請求項 9】

前記データ処理ハードウェア (6 1 0) および前記音声認識モデル (2 0 0) は、ユーザ装置 (1 0 2) 上に存在する、

請求項 1 ~ 8 のいずれか一項に記載の方法 (5 0 0)。

【請求項 10】

前記データ処理ハードウェア (6 1 0) および前記音声認識モデル (2 0 0) は、リモート計算装置 (2 0 1) 上に存在し、

前記発話 (1 0 6) を符号化する前記音声データを受け取る工程は、前記リモート計算装置 (2 0 1) に通信しているユーザ装置 (1 0 2) から、前記発話 (1 0 6) を符号化する前記音声データを受け取る工程を備えている、

請求項 1 ~ 9 のいずれか一項に記載の方法 (5 0 0)。

40

【請求項 11】

データ処理ハードウェア (6 1 0) と、

前記データ処理ハードウェア (6 1 0) に通信するメモリハードウェア (6 2 0) であって、前記メモリハードウェア (6 2 0) は、前記データ処理ハードウェア (6 1 0) 上で実行されると前記データ処理ハードウェア (6 1 0) に、以下を備えている動作を実行させる命令を格納する、前記メモリハードウェア (6 2 0) と、

を備えているシステム (1 0 0) であって、前記動作は、

第 1 言語のネイティブスピーカ (1 1 0) によって話される発話 (1 0 6) を符号化する

50

音声データを受け取る工程と、
 前記第 1 言語とは異なる第 2 言語による 1 つまたは複数の用語を備えているバイアス用語リスト (1 0 5) を受け取る工程と、
 音声認識モデル (2 0 0) を用いて、前記音声データから得られる音響特徴 (1 0 4) を処理して、前記第 1 言語の語句と、対応する音素シーケンスとの両方に対する音声認識スコアを生成する工程と、
 前記バイアス用語リスト (1 0 5) の前記 1 つまたは複数の用語に基づき、前記音素シーケンスに対する前記音声認識スコアを再スコアリングする工程と、
 前記語句に対する前記音声認識スコアと、前記音素シーケンスに対する再スコアリングされた音声認識スコアとを用いて、復号グラフ (4 0 0) を実行して、前記発話 (1 0 6) に対する転写 (1 1 6) を生成する工程と、
 を備えている、システム (1 0 0) 。

10

【請求項 1 2】

前記音素シーケンスに対する前記音声認識スコアを再スコアリングする工程は、バイアス有限状態変換器 (F S T) を使用して、前記音素シーケンスに対する前記音声認識スコアを再スコアリングする工程を備えている、
 請求項 1 1 に記載のシステム (1 0 0) 。

【請求項 1 3】

前記動作はさらに、
 前記バイアス用語リスト (1 0 5) の各用語を、前記第 2 言語の対応する音素シーケンスにトークン化する工程と、
 前記第 2 言語における各対応する音素シーケンスを、前記第 1 言語における対応する音素シーケンスに写像する工程と、
 前記第 1 言語における各対応する音素シーケンスに基づき、前記バイアス有限状態変換器 (3 0 0) を生成する工程と、
 を備えている、請求項 1 2 に記載のシステム (1 0 0) 。

20

【請求項 1 4】

前記音声認識モデル (2 0 0) は、エンドツーエンド語句 - 音素モデル (2 0 0) を備えている、
 請求項 1 1 ~ 1 3 のいずれか一項に記載のシステム (1 0 0) 。

30

【請求項 1 5】

前記エンドツーエンド語句 - 音素モデル (2 0 0) は、リカレントニューラルネットワーク - 変換器 (R N N - T) を備えている、
 請求項 1 4 に記載のシステム (1 0 0) 。

【請求項 1 6】

前記復号グラフ (4 0 0) の実行中に、前記復号グラフ (4 0 0) は、前記バイアス用語リスト (1 0 5) 内の前記 1 つまたは複数の用語のいずれかを有利にするように、前記転写 (1 1 6) をバイアスする、
 請求項 1 1 ~ 1 5 のいずれか一項に記載のシステム (1 0 0) 。

【請求項 1 7】

前記音声認識モデル (2 0 0) は、前記第 1 言語のみの学習発話で学習される、
 請求項 1 1 ~ 1 6 のいずれか一項に記載のシステム (1 0 0) 。

40

【請求項 1 8】

前記バイアス用語リスト (1 0 5) 内の用語のいずれも、前記音声認識モデル (2 0 0) を学習するために使用されなかった、
 請求項 1 1 ~ 1 7 のいずれか一項に記載のシステム (1 0 0) 。

【請求項 1 9】

前記データ処理ハードウェア (6 1 0) および前記音声認識モデル (2 0 0) は、ユーザ装置 (1 0 2) 上に存在する、
 請求項 1 1 ~ 1 8 のいずれか一項に記載のシステム (1 0 0) 。

50

【請求項 20】

前記データ処理ハードウェア(610)および前記音声認識モデル(200)は、リモート計算装置(201)上に存在し、
前記発話(106)を符号化する前記音声データを受け取る工程は、前記リモート計算装置(201)に通信しているユーザ装置(102)から、前記発話(106)を符号化する前記音声データを受け取る工程を備えている、
請求項11~19のいずれか一項に記載のシステム(100)。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、エンドツーエンドモデルにおける多言語(クロスリンガル)音声認識のための音素(phoneme)ベースのコンテキスト化(文脈解析)に関する。

【背景技術】

【0002】

音声の文脈(コンテキスト)を認識することは、自動音声認識(ASR)システムの目標である。しかし、人が話す言葉は多種多様であり、アクセントや発音にも違いがあるので、音声の文脈を認識することは困難である。多くの場合、人が話す単語やフレーズの種類(タイプ)は、その人が置かれている文脈に応じて変化する。

【0003】

文脈的(コンテクスチュアル)自動音声認識ASRは音声認識を、ユーザ自身のプレイリスト、連絡先、地理的な地名など、与えられた文脈(コンテキスト)に偏らせる(バイアスする)。文脈情報には、通常、認識すべき関連フレーズのリストが含まれており、このリストには、珍しいフレーズや、学習(トレーニング)ではあまり見られない外国語が含まれていることが多い。文脈バイアスを行うべく、従来の自動音声認識ASRシステムでは、文脈情報をn-gram重み付き有限状態変換器(WFST: weighted Finite State Transducer)を用いて、独立した文脈言語モデル(LM: Language Model)でモデル化し、その独立した文脈言語モデルLMをベースライン言語モデルLMと合成して、オンザフライ(OTF)再スコアリングを行うことがある。

【先行技術文献】

【特許文献】

【0004】

【文献】米国特許出願公開第2016/104482号明細書

【発明の概要】

【発明が解決しようとする課題】

【0005】

近年、エンドツーエンド(E2E)モデルが自動音声認識ASRに大きな期待を寄せており、従来のオンデバイスモデルと比較して、ワードエラーレート(WER)やレイテンシの指標(メトリックス)が改善されている。これらのE2Eモデルは、音響モデル(AM)、発音モデル(PM: Pronunciation Model)、および言語モデルLMを単一のネットワークに折り畳んで、音声とテキストの写像(スピーチツーテキストマッピング)を直接学習するものであり、音響モデルAM、発音モデルPM、および言語モデルLMを個別に有する従来のASRシステムと比較して、競争力のある結果を示している。代表的なE2Eモデルには、単語ベースのCTC(Connectionist temporal Classification)モデルと、RNN-T(リカレントニューラルネットワークトランスデューサ)モデルと、LAS(Listen, Attend, およびSpell)などの注意ベースモデル(アテンションベースモデル)とがある。E2Eモデルは、ビーム検索復号時(ビームサーチデコーディング時)に限られた数の認識候補を保持しているため、文脈的自動音声認識ASRはE2Eモデルにとって困難である。

【課題を解決するための手段】

10

20

30

40

50

【0006】

本開示の一態様は、バイアス用語 (biasing term) リストに存在する用語に音声認識結果をバイアスする (偏らせる) 方法を提供する。この方法は、データ処理ハードウェアにおいて、第1言語のネイティブスピーカによって話される発話を符号化 (エンコーディング) する音声データ (オーディオデータ) を受け取る工程と、データ処理ハードウェアにおいて、第1言語とは異なる第2言語の1つまたは複数の用語 (terms) を備えているバイアス用語リストを受け取る工程とを備えている。本方法は、データ処理ハードウェアによって、音声認識モデルを用いて、音声データから得られた音響特徴を処理して、第1言語における語句 (ワードピース (Word piece)) と、対応する音素シーケンス (音素列) との両方に対する音声認識スコアを生成する工程も備えている。また、本方法は、データ処理ハードウェアによって、バイアス用語リスト内の1つまたは複数の用語に基づき、音素シーケンスに対する音声認識スコアを再スコアリングする工程を備えている。本方法はまた、データ処理ハードウェアによって、語句に対する音声認識スコアと、音素シーケンスに対する再スコアリングされた音声認識スコアとを用いて、発話に対する転写 (トランスクリプション) を生成するための復号グラフ (デコーディンググラフ) を実行する工程を備えている。

10

【0007】

本開示の実装は、以下のオプション機能のうちの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、音素シーケンスに対する音声認識スコアを再スコアリングする工程は、バイアスのかかった有限状態変換器 (biasing Finite State Transducer. バイアスFST) を使用して、音素シーケンスに対する音声認識スコアを再スコアリングする工程を備えている。これらの実装では、本方法は、データ処理ハードウェアによって、バイアス用語リスト内の各用語を、第2言語における対応する音素シーケンスにトークン化する工程と、データ処理ハードウェアによって、第2言語における各対応する音素シーケンスを、第1言語における対応する音素シーケンスに写像する工程と、データ処理ハードウェアによって、第1言語における各対応する音素シーケンスに基づき、バイアス有限状態変換器FSTを生成する工程と、を備えていることもできる。

20

【0008】

いくつかの例では、音声認識モデルは、エンドツーエンドの語句 - 音素モデルを備えている。特定の例では、エンドツーエンドの語句 - 音素モデルは、リカレントニューラルネットワーク - 変換器 (RNN - T) を備えている。

30

【0009】

いくつかの実装では、復号グラフの実行中に、復号グラフは、バイアス用語リスト内の1つまたは複数の用語のいずれかを有利にするように転写 (トランスクリプション) をバイアスする。音声認識モデルは、第1言語のみの学習発話で学習されてもよい。さらに、バイアス用語リスト内のいずれの用語も、音声認識モデルの学習に使用されなくてもよい。

【0010】

データ処理ハードウェアおよび音声認識モデルは、ユーザ装置上、またはユーザ装置に通信するリモート計算装置上に存在してもよい。データ処理ハードウェアおよび音声認識モデルがリモート計算装置上に存在する場合、発話を符号化する音声データを受け取る工程は、ユーザ装置から、発話を符号化する音声データを受け取る工程を備えてもよい。

40

【0011】

本開示の別の態様は、バイアス用語リストに存在する用語に音声認識結果を偏らせるシステムを提供する。このシステムは、データ処理ハードウェアと、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、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を実行させる命令を格納するメモリハードウェアとを備えている。この動作は、第1言語のネイティブスピーカによって話された発話を符号化する音声データを受け取る工程と、第1言語とは異なる第2言語の1つまたは複数の用語を備えているバイアス用語リストを受け取る工程と、音声認識モデルを使用して、音声データから

50

得られた音響特徴を処理して、第1言語の語句と、対応する音素シーケンスとの両方に対する音声認識スコアを生成する工程とを備えている。また、動作は、バイアス用語リスト内の1つまたは複数の用語に基づき、音素シーケンスに対する音声認識スコアを再スコアリングする工程と、語句に対する音声認識スコアと、音素シーケンスに対する再スコアリングされた音声認識スコアとを使用して、復号グラフを実行して、発話に対する転写を生成する工程とを備えている。

【0012】

この態様は、以下のオプション機能の1つまたは複数を用意することができる。いくつかの実装において、音素シーケンスに対する音声認識スコアを再スコアリングする工程は、バイアスのかかった有限状態変換器(FST)を使用して、音素シーケンスに対する音声認識スコアを再スコアリングする工程を備えている。これらの実装では、動作は、バイアス用語リスト内の各用語を第2言語の対応する音素シーケンスにトークン化する工程と、第2言語の各対応する音素シーケンスを第1言語の対応する音素シーケンスに写像する工程と、第1言語の各対応する音素シーケンスに基づきバイアス有限状態変換器FSTを生成する工程と、を備えていることもできる。

10

【0013】

いくつかの例では、音声認識モデルは、エンドツーエンドの語句-音素モデルを備えている。特定の例では、エンドツーエンドの語句-音素モデルは、リカレントニューラルネットワーク-変換器(RNN-T)を備えている。

【0014】

いくつかの実装では、復号グラフの実行中に、復号グラフは、バイアス用語リスト内の1つまたは複数の用語のいずれかに有利になるように転写(トランスクリプション)をバイアスする。音声認識モデルは、第1言語のみの学習発話で学習されてもよい。さらに、バイアス用語リスト内のいずれの用語も、音声認識モデルの学習に使用されなくてもよい。

20

【0015】

データ処理ハードウェアおよび音声認識モデルは、ユーザ装置上、またはユーザ装置に通信するリモート計算装置上に存在してもよい。データ処理ハードウェアおよび音声認識モデルがリモート計算装置上に存在する場合、発話を符号化する音声データを受け取る工程は、ユーザ装置から、発話を符号化する音声データを受け取る工程を備えてもよい。

【0016】

本開示の1つまたは複数の実装の詳細は、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。

30

【図面の簡単な説明】

【0017】

【図1】バイアス用語リストに存在する用語に向けて音声認識結果をバイアスする音声認識モデルを備えている、自動音声認識システムの例を示す概略図。

【図2】図1の音声認識モデルのアーキテクチャの一例を示す概略図。

【図3】例示的なバイアス有限状態変換器の模式図。

【図4】語句と、対応する音素シーケンスとに基づく、復号グラフの例の概略図。

40

【図5】バイアス用語リストに存在する用語に向けて音声認識結果をバイアスする方法のための動作の、例示的な配置のフローチャート。

【図6】本明細書に記載されたシステムおよび方法を実施するべく使用することができる、例示的な計算装置の概略図。

【発明を実施するための形態】

【0018】

様々な図面における同様の参照記号は、同様の要素を示す。

本明細書の実装は、他の動作の中でも、外国語音素セットを自動音声認識ASRモデルの言語(例えば、アメリカ英語)の音素セットに写像(マッピング)して、音素レベルでバイアスかける有限状態変換器(FST)において外国語のモデリングを可能にすること

50

で、外国語を認識する文脈的（コンテクチュアル）自動音声認識（ASR）モデルを強化することに向けられている。さらなる実装は、自動音声認識ASRモデルが、モデリング空間における自動音声認識ASRモデルの言語（例えば、アメリカ英語）のための語句（ワードピース）および音素を備えている語句（ワードピース）-音素モデルを組み込むことに向けられている。例として、文脈的自動音声認識ASRモデルは、語句-音素モデルおよび文脈的バイアス有限状態変換器FSTを使用して音声発話（スポークン発話）を復号（デコード）し、発話の転写を文脈的に1つまたは複数の外国語に偏らせるように構成される。たとえば、アメリカ英語を話す人が、クレテイユ（Creteil。Creのeの上にアクセントが付けられている）という単語がフランス語である、「クレテイユまでの道順」（Directions to Creteil）という発話をする、文脈的自動音声認識ASRモデルは、アメリカ英語以外の言語の単語で学習されていないにもかかわらず、語句-音素モデルと文脈的バイアス有限状態変換器FSTを利用して、外国語であるクレテイユ（Creteil）を認識するように転写を偏らせることができる。この例では、外国語のクレテイユ（Creteil）は、現在の文脈に基づきバイアスをかけた単語リストに含まれる複数のフランス語のうちの一つである可能性がある。例えば、ユーザが現在フランスにいて車を運転している場合、現在の文脈（コンテキスト）は、フランスの都市名/地域名に関連していることを示している可能性があり、したがって、文脈的（コンテキストに基づく）自動音声認識ASRモデルは、これらのフランスの都市名/地域名に偏っている（バイアスしている）可能性がある。

10

【0019】

20

図1を参照すると、いくつかの実装では、強化（エンハンス）された自動音声認識ASRシステム100は、外国語の単語（ワード）を認識するように強化されている。示された例では、自動音声認識ASRシステム100は、ユーザ110のユーザ装置102上、および/または、ユーザ装置に通信するリモート計算装置（リモート計算装置）201（例えば、クラウド計算環境で実行される分散システムの1つまたは複数のサーバ）上に存在する。ユーザ装置102は、モバイル計算装置（例えば、スマートフォン）として描かれているが、ユーザ装置102は、限定されないが、タブレットデバイス、ラップトップ/デスクトップコンピュータ、ウェアラブルデバイス、デジタルアシスタントデバイス、スマートスピーカ/ディスプレイ、スマートアプライアンス、自動車インフォテイメントシステム、またはIoT（インターネットオブシングス）デバイスなどの任意のタイプの計算装置に対応してもよい。

30

【0020】

ユーザ装置102は音声サブシステム103を備えており、音声サブシステム103は、ユーザ110によって話された発話106を受け取り（例えば、ユーザ装置102は、話された発話106を記録するための1つまたは複数のマイクロフォンを備えてもよい）、発話106を、自動音声認識ASRシステム100によって処理可能なパラメータ化された入力音響フレーム104に関連する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。示されている例では、ユーザは「クレテイユまでの道順」（Directions to Creteil）というフレーズに対するそれぞれの発話106を話し、音声サブシステム103は発話106を、自動音声認識ASRシステム100に入力するための対応する音響フレーム104に変換する。例えば、音響フレーム104は、短い、例えば25msのウィンドウで計算され、数ミリ秒、例えば10ミリ秒ごとにシフトされた、それぞれが80次元のLog-Mel特徴を備えている一連（シリーズ）のパラメータ化された入力音響フレームであってもよい。

40

【0021】

その後、自動音声認識ASRシステム100は、入力として、発話106に対応する音響フレーム104を受け取り、出力として、発話106に対応する転写（トランスクリプション。例えば、認識結果/認識仮説）116を生成/予測する。図示の例では、ユーザ装置102および/またはリモート計算装置201は、ユーザ装置102のユーザインタフェース136において、発話106の転写116の表現（レプレゼンテーション）をユー

50

ザ 1 1 0 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システム 1 0 7 も実行する。いくつかの例では、ユーザインタフェース 1 3 6 は、ユーザ装置 1 0 2 に通信しているスクリーン上に表示されてもよい。

【 0 0 2 2 】

いくつかの構成では、自動音声認識 A S R システム 1 0 0 から出力された転写 1 1 6 は、例えば、ユーザ装置 1 0 2 またはリモート計算装置 2 0 1 上で実行される自然言語理解 (N L U) モジュールによって、ユーザコマンドを実行するべく処理される。さらに、または代替として、音声合成 (テキストツースピーチ) システム (例えば、ユーザ装置 1 0 2 またはリモート計算装置 2 0 1 の任意の組み合わせ上で実行される) は、別のデバイスによる可聴出力のために、転写を合成音声に変換してもよい。例えば、元の発話 1 0 6 は、ユーザ 1 1 0 が友人に送信しているメッセージに対応していてもよく、その場合、転写 1 1 6 は、元の発話 1 0 6 で伝えられたメッセージを聞くべく、友人への可聴出力のために合成音声に変換される。

10

【 0 0 2 3 】

強化された自動音声認識 A S R システム 1 0 0 は、バイアス構成要素 1 1 5 と、語句 - 音素モデル 2 0 0 およびバイアス有限状態変換器 F S T 3 0 0 を有する音声認識装置 1 5 0 と、学習構成要素 1 1 4 とを備えている。バイアス構成要素 1 1 5 は、バイアス有限状態変換器 F S T 3 0 0 を生成するように構成され、学習構成要素 1 1 4 は、音素レベルで外国語を再スコアリングすることで文脈的バイアスを実行するように、語句 - 音素モデル 2 0 0 およびバイアス有限状態変換器 F S T 3 0 0 を学習するように構成される。明らかになるように、音声認識装置 1 5 0 は、学習された語句 - 音素モデル 2 0 0 およびバイアス有限状態変換器 F S T 3 0 0 を使用して、外国語の単語に向かってバイアスをかけることで、文脈的な音声認識を実行する。

20

【 0 0 2 4 】

学習構成要素 1 1 4 は、単一の言語、例えば、アメリカ英語のテキストのコーパスを有する辞書 (レキシコン、語彙集) 1 1 7 と、頻度チェッカ 1 1 8 と、モデル学習器 1 2 0 とを備えている。頻度 (フリーケンシ) チェッカ 1 1 8 は、コーパスのテキストの中での単一言語の用語の相対的な頻度を決定するように構成され、モデル学習器 1 2 0 は、テキストコーパスの用語の語句と音素の両方に基づき語句 - 音素モデル 2 0 0 を学習し、モデリング空間に語句と音素の両方を含めるように構成される。いくつかの例では、語句 - 音素モデル 2 0 0 は、単一の言語のみ、例えば、アメリカ英語のみからの語句 - 音素セットを含む一方で、他の言語からの語句 - 音素セットを除外した学習データを用いて、モデル学習器 1 2 0 によってエンドツーエンドで学習される。モデル学習器 1 2 0 は、単語頻度ベースのサンプリング戦略を採用して、辞書 1 1 7 を用いて、稀な単語をターゲットシーケンスの音素にランダムにトークン化してもよい。段階 (ステージ) A において、学習構成要素 1 1 4 は、辞書 1 1 7 からのテキストを使用して、語句 - 音素モデル 2 0 0 を学習する。

30

【 0 0 2 5 】

いくつかの例では、辞書 1 1 7 は約 5 0 万個の単語を含み、その頻度は音素シーケンスを使用するタイミングを決定するべく使用される。辞書 1 1 7 は、学習データからの単語とその頻度を含み、同音異義語 (ホモフォン。例えば、「 f l o w e r 」 (花) と「 f l o u r 」 (小麦粉))、同形異義語 (ホモグラフ。例えば、動詞または形容詞としての「 l i v e 」 (生きる、生の))、および発音変種 (プロナンシエイションバリエント。例えば、「 e i t h e r 」 (イーザーまたはアイザー)) を除去してトリミングされる。このように、辞書 1 1 7 には、綴りから発音へまたはその逆の場合に、曖昧さが無い項目のみが含まれている。

40

【 0 0 2 6 】

いくつかの実装では、モデル学習器 (トレーナ) 1 2 0 は、学習入力発話を 2 5 m s のフレームに分割し、1 0 m s のレート (速度) で窓を開けシフトする。各フレームで 8 0 次元の l o g - M e l 特徴が抽出され、現在のフレームと左隣の 2 つのフレームが連結され

50

て240次元のlog-Mel特徴が生成される。これらの特徴は、その後、30msのレートでダウンサンプリングされる。

【0027】

いくつかの実装では、語句-音素モデル200は、シーケンスツーマルチヘッドモデルを備えている。いくつかの例では、語句-音素モデル200は、RNN-T（リカレントニューラルネットワーク-トランスデューサ）シーケンスツーマルチヘッドモデルアーキテクチャを備えている。他の例では、語句-音素モデル200は、リッスン、アテンド、スペルのシーケンスツーマルチヘッドモデルアーキテクチャを備えている。

【0028】

語句-音素モデル200は、学習において少数の語句を選択的に音素に分解することができる点で、語句のみのモデルとは異なる。このモデルの出力は、記号セット（シンボルセット）が語句記号と音素記号との組み合わせである、単一のソフトマックスである。単語の音素シーケンスを得るためには、発音辞書（レキシコン）が用いられる。音素は、希少な単語の認識に強みを発揮するので、これらの単語はより頻繁に音素として提示される。ターゲット文では、 i 番目の単語が確率 $p(i) = p_0 \cdot \min(T / (c(i)), 1.0)$ でランダムに音素として提示される。ここで p_0 と T は定数であり、 $c(i)$ は、学習コーパス全体での単語の出現回数を表す整数である。出現回数が T 回以下の単語は、確率 p_0 で音素として提示される。 T 回よりも多く出現する単語については、頻度が高いほど音素として提示されないことになる。いくつかの例では、 T は10に等しく、 p_0 は0.5に等しくなっているが、他の例では異なる値を選択することができる。なお、単語と音素のどちらを使用するかは、勾配のイテレーションごとにランダムに行われるので、ある文は、異なるエポックで異なるターゲットシーケンスを持つ可能性がある。いくつかの実装では、音素は文脈に依存しない音素である。

【0029】

図2を参照すると、語句-音素モデル200は、インタラクティブアプリケーションに関連付けられたレイテンシ制約に準拠したエンドツーエンド（E2E）のRNN-Tモデル200を備えていることができる。RNN-Tモデル200は、小さな計算フットプリントを提供し、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、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（例えば、リモートサーバとの通信は必要とされない）。RNN-Tモデル200は、符号化器ネットワーク（エンコーダネットワーク）210と、予測ネットワーク220と、結合ネットワーク（ジョイントネットワーク）230とを備えている。符号化器ネットワーク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）にほぼ類似しており、積層されたLSTM（Long Short-Term Memory）層のリカレントネットワークを備えている。例えば符号化器は、 x_t Rd（Rは白抜き文字）である d 次元特徴ベクトル（例えば、音響フレーム104（図1））のシーケンス $x = (x_1, x_2, \dots, x_T)$ を読み込み、各時間ステップで高次の特徴表現を生成する。この高次の特徴表現は、 $h_1^{enc}, \dots, h_T^{enc}$ のように示される。

【0030】

同様に、予測ネットワーク220もLSTMネットワークであり、言語モデル（LM）のように、これまでに最終ソフトマックス層240が出力した非空白記号のシーケンス y_0, \dots, y_{u_i-1} を処理して、高密度の表現 P_{u_i} にする。最後に、RNN-Tモデルのアーキテクチャでは、符号化器ネットワーク210および予測ネットワーク220によって生成された表現同士は、結合ネットワーク230によって結合される。結合ネットワーク230は、次の出力記号に対する分布である予測 $P(y_i | x_1, \dots, x_{t_i}, y_0, \dots, y_{u_i-1})$ を行う。別の言い方をすると、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、「可能性のある音声認識仮説」（ポシブルスピーチレコグニションヒポセシス）は、指定された自然言語の記号/文字（キャラクタ）をそれぞれ

10

20

30

40

50

が表す出力ラベルの第1セットと、指定された自然言語の音素をそれぞれが表す出力ラベルの第2セットとに対応する。したがって、結合ネットワーク230は、所定の出力ラベルのセットのそれぞれの発生の可能性（ライクリフッドオブオカレンス）を示す一連の値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、出力ラベルのセットに対する確率分布を示すことができる。いくつかのケースでは、出力ラベルは、第1セットでは書記素（graphemes。例えば、個々の文字、および潜在的には句読点および他の記号）であり、第2セットでは音素であるが、出力ラベルのセットはそのように限定されない。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それぞれに対する事後確率値（ポステリアプロバビリティバリュー）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す100個の異なる出力ラベルがある場合、結合ネットワーク230の出力 y_i は、各出力ラベルに対して1つずつになるように、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、転写116を決定するためのビーム探索プロセス（例えば、ソフトマックス層240による）において、正書法の候補要素（candidate orthographic element）（例えば、書記素、語句、単語、音素）を選択し、スコアを割り当てるべく使用することができる。

10

【0031】

ソフトマックス層240は、対応する出力ステップでモデル200によって予測される次の出力記号として、分布内で最も高い確率を持つ出力ラベル/記号を選択するべく、任意の技術を採用することができる。このようにして、RNN-Tモデル200は条件付き独立性仮定を行わず、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ている。RNN-Tモデル200は、出力記号が将来の音響フレーム104から独立していると仮定しており、これによって、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。

20

【0032】

いくつかの例では、RNN-Tモデル200の符号化器ネットワーク210は、8個の2048次元LSTM層で構成され、それぞれの後に640次元の投影（プロジェクション）層が続く。モデルのレイテンシを低減するべく、符号化器の第2LSTM層の後に、低減（リダクション）係数が2の時間低減層を挿入してもよい。また、予測ネットワーク220は、2個の2048次元LSTM層を有していてもよく、それぞれの後に640次元の投影層が続いている。最後に、結合ネットワーク230は、640個の隠れユニットと、その後に続く4096個のソフトマックス出力も有していてもよい。具体的には、出力ユニットは、41個の文脈非依存音素を含み、残りは語句（ワードピース）である。

30

【0033】

図1に戻って、自動音声認識ASRシステム100のバイアス構成要素115は、バイアスされるべき外国語のバイアス用語リスト105からの用語を外国語音素にトークン化するように構成されたトークン化器121と、トークン化された用語の外国語音素を単一言語、例えば、アメリカ英語に関連する類似の音素に写像（マッピング）するように構成された音素写像器（マッパー）123とを備えている。音素写像器123は、人間が生成したソース言語からターゲット言語への音素ペアを備えている辞書によって表されてもよく、X-SAMPA音素セットはすべての言語に使用される。注目すべきは、音素写像器123は、語句-音素モデル200が、単一の言語、例えば、アメリカ英語に関連する音素のみを備えている場合に有用である。

40

【0034】

例えば、ナビゲーションクエリ「クレティユまでの道順」（directions to Creteil）の発話106と、フランス語の単語「クレティユ」（Creteil）がバイアス用語リスト105内にあるという仮定が与えられた場合、「クレティユ」（Creteil）は、まずトークン化器121によって「k R e t E j」としてフランス語の音素にトークン化され、次に音素写像器123によって「k r ¥ E t E j」として英語の音素に写像されて、音素レベルのバイアス有限状態変換器FST300の生成

50

に使用される。語句 - 音素モデル 200 が単一の言語、例えば、アメリカ英語からの音素のみをモデリングユニットとして備えているので、音素写像は使用される。

【0035】

本開示は、どのような用語がバイアス用語リスト 105 に含まれるか、または用語がバイアス用語リスト 105 に含まれるようにどのように選択されるかに限定されない。バイアス用語リスト 105 は、関連する文脈（コンテキスト）に基づき、動的に更新されてもよい。例えば、文脈情報は、ユーザ装置 102 上でどのようなアプリケーションが開いていて使用中であるか、ユーザの連絡先リストからの連絡先名、ユーザ 110 のメディアライブラリ内のアーティスト名 / アルバム名、ユーザ 110 の位置などを示してもよい。例えば、ユーザ 110 はアメリカ英語を話すことができ、ナビゲーション / 地図アプリケーションがユーザ装置 102 上で開かれていることと、ユーザ 110 の場所がフランスであることとを示す文脈情報に基づき、バイアス用語リスト 105 は、フランスの都市名および / または地域名に関連する用語を備えていることができる。

10

【0036】

また、バイアス構成要素 115 は、音素レベルのバイアス有限状態変換器 FST 生成器 125 を備えており、音素レベルのバイアス有限状態変換器 FST 生成器 125 は、バイアス用語リスト 105 内の外国語（例えば、フランス語）用語のそれぞれを表す母語（例えば、アメリカ英語）の音素シーケンスに基づき、バイアス有限状態変換器 FST 300 を生成するように構成されている。いくつかの例では、バイアス有限状態変換器 FST 生成器 125 は、音素レベルで重みを割り当てるべく押す重み（weight pushing）を使用し、過剰バイアスを避けるべく失敗アーク（failure arcs）を追加する。いくつかの実装では、復号化において、すべてのバイアス語を使用して、各アーク（arc）が同じ重みを持つ文脈的有限状態変換器 FST を構築する。これらの重みは、異なるモデルに対して独立して、調整することができる。

20

【0037】

音声認識装置 150 は、バイアス構成要素 115 によって生成されたバイアス有限状態変換器 FST 300 を使用して、語句 - 音素モデル 200 によって出力された音素を再スコアリングし、一方、復号グラフ 400 は、バイアス有限状態変換器 FST 300 からの再スコアリングされた音素と、語句 - 音素モデル 200 によって出力された語句とを消費して、転写 116 に含めるための語句を生成する。復号グラフ 400 は、発話 106 に対する 1 つまたは複数の転写候補を決定するビーム探索復号処理に対応してもよい。

30

【0038】

いくつかの例では、語句 - 音素モデル 200 による復号中に、バイアス有限状態変換器 FST 300 は、語句 - 音素モデル 200 によって出力された英語音素記号を消費し、外国語辞書および音素写像を使用して語句を生成してもよく、すなわち「k r ¥ E t E j」クレテイユ（Creteil）である。復号グラフ 400 によって出力された語句は、連結器（コンカチネータ）134 によって、ユーザ装置 102 の他の構成要素に出力される転写 116 の単語（ワード）に連結され、ここでユーザ装置 102 の他の構成要素は、例えば、ユーザインタフェース生成システム 107 や、他の自然言語処理構成要素である。

40

【0039】

図 3 は、音素レベルでの単語「クレテイユ」（Creteil）に対する、例示的なバイアス有限状態変換器 FST 300 を示す。そして、このバイアス有限状態変換器 FST は、以下の式（1）を用いて、語句 - 音素モデルの音素出力をオンザフライで再スコアリングするべく使用される。

【0040】

【数 1】

$$y^* = \arg \max_y \log P(y|x) + \lambda \log P_C(y) \quad (1)$$

50

【 0 0 4 1 】

式 (1) において、 x は音響観測値であり、 y はサブ単語 (サブワード) ユニットシーケンスであり、 P は E 2 E モデルからの確率推定であり、 P_c はバイアス再スコアリング確率である。 λ は、再スコアリングにおける文脈言語モデル L M の重みを制御する。

【 0 0 4 2 】

図 1 に戻ると、語句 - 音素モデル 2 0 0 は、モデリングユニットとして語句だけでなく音素も組み込み、バイアス用語リスト 1 0 5 内の外国語用語に向けた文脈上のバイアスのためにバイアス有限状態変換器 F S T 3 0 0 を使用する。すべての音素モデルとは対照的に、音素と語句の両方をモデル化する語句 - 音素モデル 2 0 0 は、通常の単語 (レギュラーワード) を認識する際の回帰を緩和する。

10

【 0 0 4 3 】

語句 - 音素モデル 2 0 0 が段階 A で学習 (トレーニング) された後、段階 B で、ユーザ 1 1 0 は、発話 1 0 6 「クレティユへの道順」 (*directions to Creteil*) をユーザ装置 1 0 2 に話す。段階 C において、音声サブシステム 1 0 3 は、例えばマイクロフォンを使用して、発話を受け取り、受け取った発話を、一連のパラメータ化された入力音響フレーム 1 0 4 に変換する。例えば、パラメータ化された入力音響フレーム 1 0 4 はそれぞれ、8 0 次元の \log - M e l 特徴を備えてもよく、ここで 8 0 次元の \log - M e l 特徴は、短い、例えば 2 5 m s のウィンドウで計算されるとともに、数ミリ秒ごと、例えば 1 0 ミリ秒ごとにシフトされる。

【 0 0 4 4 】

段階 D において、自動音声認識 A S R システム 1 0 0 は、上述したようにパラメータ化された入力音響フレームを処理し、文脈的にバイアスされた転写 1 1 6、すなわちテキスト「クレティユ (*Creteil*) への道順」を出力する。段階 E において、ユーザインタフェース生成システム 1 0 7 は、転写の表現を備えているグラフィカルユーザインタフェース 1 3 6 のためのコンピュータコードを生成し、段階 F において、ユーザインタフェース 1 3 6 に表示するべく、そのコンピュータコードをモバイル装置 (1 0 2) に送信する。

20

【 0 0 4 5 】

自動音声認識 A S R システム 1 0 0 によって実行される追加の詳細は、段階 D の期間内に発生する可能性がある。例えば、段階 D の期間内に、バイアス構成要素 1 1 5 は、用語「クレティユ」 (*Creteil*) を備えているバイアス用語リスト 1 0 5 の受け取りに基づき、バイアス有限状態変換器 F S T 3 0 0 を生成する。段階 D において、音声認識装置 1 5 0 の学習された語句 - 音素モデル 2 0 0 は、ユーザ 1 1 0 の発話 1 0 6 に基づき、語句と、対応する音素シーケンスとの両方に対する音声認識スコアを生成し、音素に対する音声認識スコアは、バイアス有限状態変換器 F S T 3 0 0 によって再スコアリングおよび再写像され、語句と、再スコアリング / 再写像された音素とは、段階 D において、転写 1 1 6 で出力するための語句を生成するべく、復号グラフ 4 0 0 によって消費される。復号グラフ 4 0 0 および連結器 1 3 4 は、文脈的にバイアスされた転写 1 1 6 を生成し、出力用の転写を、例えば、ユーザ装置 1 0 2 の G U I 1 3 6 に表示するべくユーザインタフェース生成システム 1 0 7 に提供する。注目すべきは、バイアス用語リスト 1 0 5 内の用語のいずれかに対応する音素シーケンスをバイアス有限状態変換器 F S T 3 0 0 が再スコアリングした後に、復号グラフ 4 0 0 は実行されることである。このように、バイアス用語リスト 1 0 5 内の外国語に対応する低い音声認識スコアを有する語句は、早々には剪定 (*prune*) されない。

30

40

【 0 0 4 6 】

テスト中、語句 - 音素モデル 2 0 0 およびバイアス有限状態変換器 F S T 3 0 0 を採用して認識結果をバイアス用語リスト 1 0 5 内の用語に向けて文脈的にバイアスする音声認識装置 1 5 0 は、書記素のみのバイアスモデルと語句のみのバイアスモデルとの両方よりも顕著に優れた W E R 率で、外国語単語の認識に成功することが示された。また、語句 - 音素モデル 2 0 0 は、モデルのスケラビリティの問題なく、他の外国語に直接適用してバイアスをかけることができるという利点がある。

50

【 0 0 4 7 】

図 4 は、音声認識装置 1 5 0 が音声認識結果を文脈的にバイアスするべく実行する、例示的な復号グラフ 4 0 0 を示す。具体的には、例示的な復号グラフ 4 0 0 は、英語のクロスリンガル発音「k r ¥ E S」を有する単語「クレイシュ」(creche。creのeの上にアクソングラフが付されている。英語では「デイケア」(daycare))と、発音「k r ¥ E t E j」を有する単語「クレテイユ」(Cretail。フランスの都市)とに対する復号を描いている。なお、わかりやすくするべく、「0」という状態の語句はほとんど省略している。

【 0 0 4 8 】

復号グラフ 4 0 0 は、語句 - 音素モデル 2 0 0 から出力された音素と語句の両方を入力として受け取るように構成されている。音声復号化処理は、復号グラフ(デコーディンググラフ) 4 0 0 を検索して、出力として単語を生成する。図示の例では、復号有限状態変換器 F S T は、状態 0 を中心とした語句ループを有するが、発音有限状態変換器 F S T、すなわち状態 1 ~ 1 4 を有し、それら状態は音素を入力とし、対応する語句を出力とする接頭辞(prefix) ツリーを備えている。発音有限状態変換器 F S T は、すべてのバイアス用語について、バイアス時に使用された発音と同じ発音を用いて構築される。常に語句である最終出力記号は、(例えば、図 1 の連結器 1 3 4 によって) 単語(ワード)に連結される。

10

【 0 0 4 9 】

図 4 の復号グラフ 4 0 0 は、全体的な復号戦略に 2 つの追加の改善をもたらす。第 1 に、復号グラフ 4 0 0 の性質を考慮すると、同じコストで同じ入力を消費するが、同じ出力を持たないいくつかの仮説が存在する可能性がある。例えば、状態 7 で終了する仮説は、状態 9 で終了する仮説と同じコストを持つことになる。このため、すべてが等価な多くの仮説によって、ビームが埋め尽くされてしまうという問題が生じる。本明細書に記載されている強化された自動音声認識 A S R 技術は、このように、状態 9 で終わる仮説のみを保持することで、ビームを刈り取る(pruneする)。

20

【 0 0 5 0 】

第 2 改善点は、結合(マージ)された経路(path)に関する。学習と復号との性質を考慮すると、与えられた単語は、直接語句で出力されるか、または、音素から語句に変換される。同等の仮説が追跡され、それらの確率を加算することで再結合され、最も可能性の高い仮説に合計確率を割り当て、他のものをビームから削除する。

30

【 0 0 5 1 】

語句 - 音素モデル 2 0 0 のバイアスの結果を、語句のみのモデルと、書記素のみのモデルとに対して比較するテストが行われた。後者の 2 つのモデルは語句 - 音素モデル 2 0 0 と同じ構造を有しており、違いは、書記素モデルが出力として 7 6 個の書記素を有している一方で、語句モデルが 4 0 9 6 個の語句(ワードピース)を有することである。この違いによって、書記素モデルと語句モデルのパラメータは、それぞれ約 1 1 7 M 個と 1 2 0 M 個になる。なお、この 2 つのモデルの出力記号は英語であり、全英語データを用いて学習されている。これらの 2 つのモデルでは、フランス語のバイアス単語(ワード)の英語音訳版を使用して、書記素レベルまたは語句レベルのみでバイアスが行われる。

40

【 0 0 5 2 】

一般的に、テストでは、3 つのモデルはバイアスをかけなくても同じように動作することが示された。これは、地名がフランス語であり、それらが学習では見られたことがないためであり、すなわち、ほぼ 1 0 0 % の単語 O O V 率である。さらに、すべてのモデルは、バイアスをかけることで大幅に性能が向上する。バイアスをかけない場合と比較して、W E R の減少が顕著である。

【 0 0 5 3 】

異なるバイアス戦略を比較すると、語句 - 音素モデル 2 0 0 が最も優れた性能を示し、書記素モデルおよび語句モデルの両方よりも有意に良好に動作した。語句 - 音素モデルの優れた性能は、O O V の単語に対する音素のロバスト性に起因する。語句 - 音素モデル 2 0

50

0 は、モデリングユニットとして語句と音素の両方を備えているので、音素有限状態変換器 F S T に加えて語句有限状態変換器 F S T を構築することで、音素ベースのバイアスに加えて語句バイアスを実行することができる。この語句単位の有限状態変換器 F S T を追加することで、W E R がさらに減少することが実証されており、語句単位のバイアスと音素単位のバイアスとは相互に補完し合う関係にあることがわかる。音素と語句のバイアスに使用する重みは、同じでもよいし、異なってもよい。観察によると、長いユニットをマッチングする際のスパース性の問題から、語句単位の方が、書記素 (g r a p h e m e) 単位よりも性能が高い場合がある。

【 0 0 5 4 】

テストの結果、バイアスは外国の地名を認識するのに役立つことがわかった。例えば、バイアスをかけると、正しいフランス語の単語 (w o r d) が生成され、逆にバイアスをかけないと、音韻的には似ているが間違った英語の単語が生成される。誤りは、フランス語の音韻的に類似した単語が原因であることが多い。

10

【 0 0 5 5 】

バイアスなしのシナリオでの回帰がないことをより確実にするべく、通常の英語の発話の復号 (デコーディング) で 3 つのモデルを比較した。復号では、バイアスフレーズの空リストを使用することで、バイアスメカニズムをオフにした。テストの結果、語句モデルは、書記素 (g r a p h e m e) モデルよりも優れた性能を示すことがわかった。語句 - 音素モデルは、書記素モデルよりもやや良好な結果となったが、これは学習時に語句の頻度が高かったことに起因していると考えられる。語句モデルと比較して、語句 - 音素モデルは非常にわずかに劣化している。これは、モデリングに電話を導入したことによる。回帰性を向上させるための潜在的なアプローチとしては、語句ベースの再スコアリングと同様に、再スコアリングに音素の英語外部言語モデルを組み込むことが考えられる。しかし、全音素 (a l l - p h o n e m e) モデルに比べて、回帰が著しく小さくなる。

20

【 0 0 5 6 】

図 5 は、バイアス用語リスト内の外国語用語に向かって転写を文脈的にバイアスする方法の動作の例示的な配置のフローチャートである。動作 5 0 2 において、方法 5 0 0 は、第 1 言語のネイティブスピーカ (1 1 0) によって話される発話 1 0 6 を符号化 (エンコーディング) する音声データを受け取る工程を備えている。発話 1 0 6 は、第 1 言語とは異なる第 2 言語の 1 つまたは複数の外国語を備えてもよい。動作 5 0 4 において、方法 5 0 0 は、第 2 言語の 1 つまたは複数の用語を備えているバイアス用語リスト 1 0 5 を受け取る工程を備えている。

30

【 0 0 5 7 】

動作 5 0 6 において、方法 5 0 0 は、音声認識モデル 2 0 0 を使用して、音声データから導出された音響特徴 (1 0 4) を処理して、第 1 言語における語句と、対応する音素シーケンスとの両方に対する音声認識スコアを生成する工程も備えている。動作 5 0 8 において、方法 5 0 0 は、バイアス用語リスト内の 1 つまたは複数の用語に基づき、音素シーケンスに対する音声認識スコアを再スコアリングする工程も備えている。動作 5 0 6 において、方法 5 0 0 は、語句に対する音声認識スコアと、音素シーケンスに対する再スコアリングされた音声認識スコアとを用いて、復号グラフ (デコーディンググラフ) 4 0 0 を実行して、発話 1 0 6 に対する転写 (トランスクリプション) 1 1 6 を生成する工程を備えている。

40

【 0 0 5 8 】

ソフトウェアアプリケーション (すなわち、ソフトウェアリソース) は、計算装置にタスクを実行させるコンピュータソフトウェアを指すことがある。いくつかの例では、ソフトウェアアプリケーションは、「アプリケーション」、「アプリ」、または「プログラム」と呼ばれることがある。アプリケーションの例としては、システム診断アプリケーション、システム管理アプリケーション、システムメンテナンスアプリケーション、ワープロアプリケーション、表計算アプリケーション、メッセージングアプリケーション、メディアストリーミングアプリケーション、ソーシャルネットワーキングアプリケーション、ゲー

50

ムアプリケーションなどがあるが、これらに限定されない。

【0059】

非一過性メモリは、計算装置が使用するためのプログラム（例えば、命令のシーケンス）またはデータ（例えば、プログラムの状態情報）を一時的または永久的に保存するべく使用される物理デバイスであってもよい。非一時的メモリは、揮発性および/または不揮発性のアドレス可能な半導体メモリであってもよい。不揮発性メモリの例としては、フラッシュメモリ、リードオンリーメモリ（ROM）/プログラマブルリードオンリーメモリ（PROM）/消去可能プログラマブルリードオンリーメモリ（EPROM）/電子的消去可能プログラマブルリードオンリーメモリ（EEPROM）（例えば、ブートプログラムなどのファームウェアに典型的に使用される）などがあるが、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）、フェイズチェンジメモリ（PCM）のほか、ディスクやテープなどが挙げられるが、これらに限定されるものではない。

10

【0060】

図6は、本書で説明したシステムおよび方法を実施するべく使用することができる例示的な計算装置600の概略図である。計算装置600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、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。ここに示されている構成要素、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を意図しており、本書に記載および/または請求されている発明の実施を制限することを意図していない。

20

【0061】

計算装置600は、プロセッサ610と、メモリ620と、記憶装置（ストレージデバイス）630と、メモリ620および高速拡張ポート650に接続する高速インタフェース/コントローラ640と、および低速バス670および記憶装置（ストレージデバイス）630に接続する低速インタフェース/コントローラ660とを備えている。構成要素610、620、630、640、650、660のそれぞれは、様々なバスを用いて相互に接続されており、共通のマザーボードに搭載されていてもよいし、適宜他の態様で搭載されていてもよい。プロセッサ610は、高速インタフェース640に結合されたディスプレイ680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620または記憶装置630に格納された命令を備えている、計算装置600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/または複数のバスが、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、複数の計算装置600が接続され、各デバイスが必要な動作の一部を提供してもよい（例えば、サーババンク、ブレードサーバ群、またはマルチプロセッサシステムとして）。

30

【0062】

メモリ620は、計算装置600内の情報を非一時的に格納する。メモリ620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620は、計算装置600による使用のために、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）を一時的または永久的に格納するべく使用される物理デバイスであってもよい。不揮発性メモリの例には、フラッシュメモリおよびリードオンリーメモリ（ROM）/プログラマブルリードオンリーメモリ（PROM）/消去可能プログラマブルリードオンリーメモリ（EPROM）/電子的消去可能プログラマブルリードオンリーメモリ（EEPROM）（例えば、ブートプログラムなどのファームウェアに典型的に使用される）が含まれるが、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）、フェイズチェンジメモリ（PCM）のほか、ディスクやテープな

40

50

どが挙げられるが、これらに限定されるものではない。

【0063】

記憶装置630は、計算装置600に大容量記憶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630は、コンピュータ可読媒体である。様々な異なる実装において、記憶装置(ストレージデバイス)630は、フロッピー(登録商標)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類似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。追加の実装では、コンピュータプログラム製品が、情報キャリアに有形的に具現化される。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
10

【0064】

高速コントローラ640は、計算装置600のための帯域幅集中型の動作を管理し、低速コントローラ660は、より低い帯域幅集中型の動作を管理する。このような職務の割り当ては、例示的なものに過ぎない。いくつかの実装では、高速コントローラ640は、メモリ620と、ディスプレイ680(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)と、および、様々な拡張カード(図示せず)を受け入れ得る高速拡張ポート650とに結合される。いくつかの実装では、低速コントローラ660は、記憶装置630および低速拡張ポート690に結合される。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))を備えてもよい低速拡張ポート690は、キーボード、ポインティングデバイス、スキャナなどの1つまたは複数の入出力デバイスに、またはスイッチやルータなどの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介して結合されてもよい。
20

【0065】

計算装置600は、図に示すように、いくつかの異なる形態で実装されてもよい。例えば、計算装置は、標準的なサーバ600aまたはそのようなサーバ600aのグループにおける複数倍として、ラップトップコンピュータ600bとして、またはラックサーバシステム600cの一部として、実装されてもよい。
30

【0066】

本明細書に記載されたシステムおよび技術の様々な実装は、デジタル電子および/または光学回路、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、プログラム可能な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での実装を備えていることができ、ここでプログラム可能なシステムは、データおよび命令を記憶装置から受け取り、データおよび命令を記憶装置に送信するように記憶装置に結合された、特殊目的または汎用の少なくとも1つのプログラム可能なプロセッサと、少なくとも1つの入力装置と、および少なくとも1つの出力装置とを備えている。
40

【0067】

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、またはコードとも呼ばれる)は、プログラマブルプロセッサのための機械命令を含み、高レベルの手続き型および/またはオブジェクト指向の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。本明細書において、「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、機械可読信号として機械命令を受け取る機械可読媒体を含んでいる、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るべく使用される任意のコンピュータプログラム製品、非一時的なコンピュータ可読媒体、装置(アパレイタス)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))を意味する。「
50

機械可読信号」とは、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
するべく使用されるあらゆる信号を指す。

【0068】

本明細書に記載されてい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
ラムを実行して、入力データを操作して出力を生成することで機能を実行することができ
る。また、FPGA(Field Programmable Gate Array)やASIC(Application Specific Integrated Circui
t)などの特殊な論理回路によっても処理や論理フローを実行することができる。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、汎用および特殊目的のマ
イクロプロセッサ、および任意の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
数のメモリデバイスである。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、または光ディスク
を備えているか、またはデータを受け取るか、またはデータを転送するか、もしくは両方
であるように動作可能に結合される。しかし、コンピュータはそのようなデバイスを持っ
ている必要はない。コンピュータプログラムの命令やデータを格納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、例として、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、光磁気ディスク、およびCD-ROMおよびDVD-ROMディスクが挙げられる
。プロセッサとメモリは、特別な目的の論理回路によって補完されるか、またはそれに組
み込まれることができる。

10

20

【0069】

ユーザとの対話(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(cathode ray
tube)、LCD(liquid crystal display)モニタ、またはタ
ッチスクリーンと、任意で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
はトラックボールを有し、それによってユーザがコンピュータに入力を提供することが
できるコンピュータ上で実装することができる。同様にユーザに相互作用を提供できる多
の種類装置が使用でき、例えば、ユーザに提供されるフィードバックは、視覚的なフィ
ードバック、聴覚的なフィードバック、触覚的なフィードバックなど、あらゆる形態の感
覚的なフィードバックであり、ユーザからの入力、音響的な入力、音声的な入力、触覚
的な入力など、あらゆる形態で受け取ることができる。さらに、コンピュータは、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け取
ったりすることで、ユーザと対話することができる。例えば、ウェブブラウザから受け
取った要求に応答して、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で、ユーザと対話することができる。

30

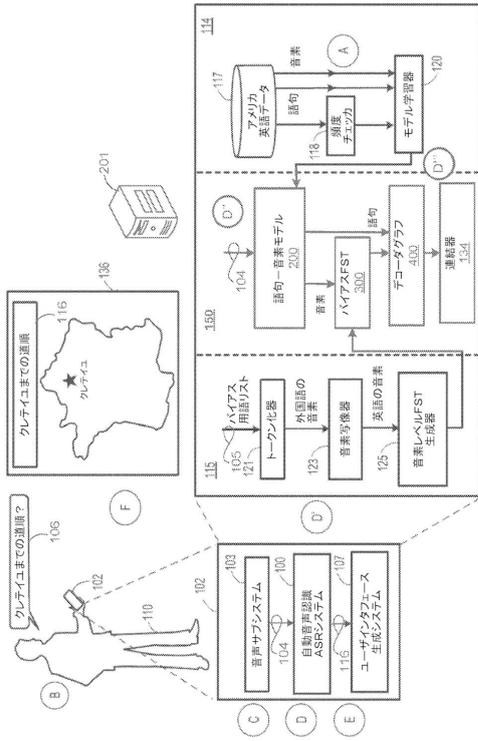
40

【0070】

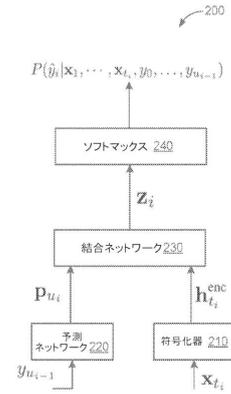
多数の実施例を説明してきた。それにもかかわらず、本開示の精神および範囲から逸脱
することなく、様々な変更を行うことができることが理解されるであろう。したがって、他
の実施態様は、以下の請求項の範囲内にある。

50

【図面】
【図 1】



【図 2】



10

20

【図 3】

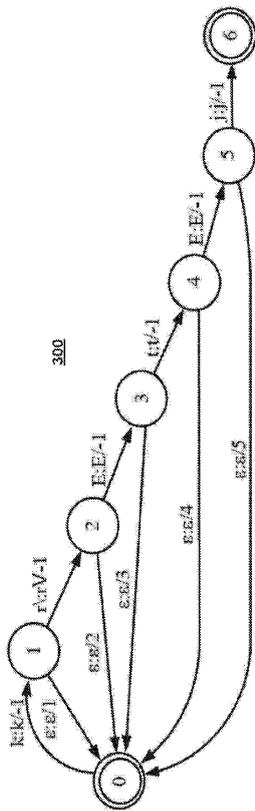


FIG. 3

【図 4】

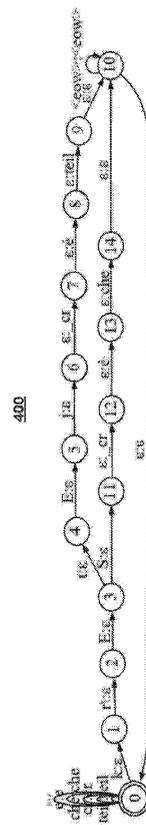


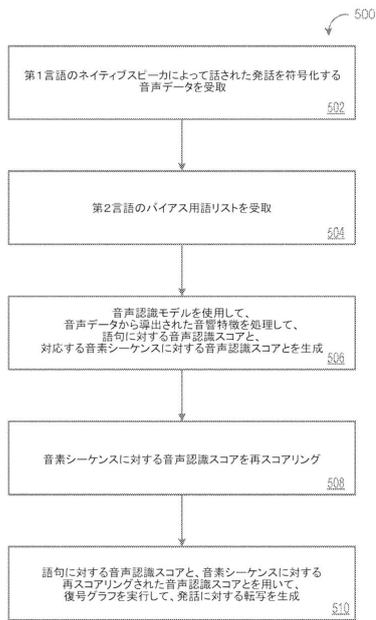
FIG. 4

30

40

50

【 図 5 】



【 図 6 】

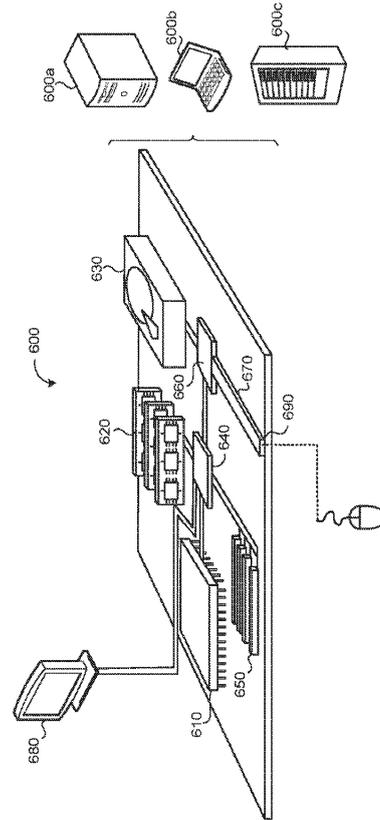


FIG. 6

10

20

30

40

50

フロントページの続き

g / a b s / 1 9 0 6 . 0 9 2 9 2 にて発表

早期審査対象出願

(72)発明者 ブルギエ、アントワーヌ ジャン

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 サイナス、ターラ エヌ .

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 プラバーバルカル、ロヒット プラカーシュ

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 ブンダック、ゴラン

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

審査官 菊池 智紀

(56)参考文献 特表 2 0 1 9 - 5 0 7 3 6 2 (J P , A)

PATEL, Ami et al. , "CROSS-LINGUAL PHONEME MAPPING FOR LANGUAGE ROBUST CONT
EXTUAL SPEECH RECOGNITION" , Proc. of the 2018 IEEE ICASSP , 2018年04月15日 , pp.
5924-5928

(58)調査した分野 (Int.Cl. , D B 名)

G 1 0 L 1 5 / 0 0 - 1 5 / 3 4

I E E E X p l o r e