

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7351018号
(P7351018)

(45)発行日 令和5年9月26日(2023.9.26)

(24)登録日 令和5年9月15日(2023.9.15)

(51)国際特許分類 F I
G 1 0 L 15/06 (2013.01) G 1 0 L 15/06 3 0 0 Z
G 1 0 L 15/16 (2006.01) G 1 0 L 15/16

請求項の数 16 (全22頁)

(21)出願番号	特願2022-545879(P2022-545879)	(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3年1月15日(2021.1.15)	(74)代理人	100142907 弁理士 本田 淳
(65)公表番号	特表2023-503717(P2023-503717 A)	(72)発明者	ベイザー、チャールズ ケイレブ 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5年1月31日(2023.1.31)		
(86)国際出願番号	PCT/US2021/013759		
(87)国際公開番号	WO2021/154520		
(87)国際公開日	令和3年8月5日(2021.8.5)		
審査請求日	令和4年10月26日(2022.10.26)		
(31)優先権主張番号	62/966,823		
(32)優先日	令和2年1月28日(2020.1.28)		
(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2年4月9日にウェブ サイトのアドレス https://ieeexplor 最終頁に続く		

(54)【発明の名称】 エンド・ツー・エンド音声認識における固有名詞認識

(57)【特許請求の範囲】

【請求項1】

データ処理ハードウェア(510)上での実行時に、前記データ処理ハードウェア(510)に動作を実行させるコンピュータが実施する方法(400)であって、前記動作は、最小単語誤り率損失関数を使用して音声認識モデル(200)をトレーニングすることであって、

固有名詞を含むトレーニング例(302)を受信すること、

前記トレーニング例(302)に対応する複数の仮説(222)を生成すること、前記複数の仮説(222)の各仮説(222)は、前記固有名詞を表し、かつ仮説(222)が前記固有名詞を表す可能性を示す対応する確率を含んでおり、

10

前記複数の仮説(222)のうちの1つに関連付けられた対応する確率がペナルティ基準を満たすことを決定すること、前記ペナルティ基準は、

前記対応する確率が確率しきい値を満たしていること、および

前記対応する確率を含む前記仮説(222)が固有名詞を誤って表していることを示しており、

前記最小単語誤り率損失関数にペナルティ(332)を適用すること

によって、前記音声認識モデル(200)をトレーニングすることを含んでおり、

前記音声認識モデル(200)が、

リカレントニューラルネットワークトランスデューサ(以下、RNN-Tとする)デコー

ダ(220)を含む第1のパスのネットワーク(206)と、

20

リッスン・アテンド・スペル（以下、LASとする）デコーダ（230）を含む第2のパスのネットワーク（208）とを備えており、
前記最小単語誤り率損失関数を使用したトレーニングは、前記LASデコーダ（230）
において行われる、コンピュータが実施する方法（400）。

【請求項2】

前記対応する確率が他の仮説（222）に関連付けられた対応する確率よりも大きい場合、前記対応する確率が前記確率しきい値を満たす、請求項1に記載のコンピュータが実施する方法（400）。

【請求項3】

前記音声認識モデルは、共有エンコーダをさらに備え、前記共有エンコーダは、前記第1のパスのネットワーク（206）および前記第2のパスのネットワーク（208）の各々に対して音響フレーム（212）をエンコードする、請求項1に記載の方法（400）。

10

【請求項4】

前記動作は、
前記RNN-Tデコーダ（220）をトレーニングすること、
前記最小単語誤り率損失関数を使用して前記LASデコーダ（230）をトレーニングする前に、トレーニングされた前記RNN-Tデコーダ（220）のパラメータが固定された状態で、前記LASデコーダ（230）をトレーニングすること、をさらに含む、請求項1または3に記載のコンピュータが実施する方法（400）。

【請求項5】

前記動作は、確率を前記複数の仮説（222）の各仮説（222）に割り当てることをさらに含む、請求項1乃至4のいずれか一項に記載のコンピュータが実施する方法（400）。

20

【請求項6】

前記動作は、
誤った仮説（222）を受信すること、
個別の確率を前記誤った仮説（222）に割り当てること、をさらに含み、
前記ペナルティ基準は、仮説（222）が生成された誤った仮説（222）を含むという表示をさらに含む、請求項1乃至5のいずれか一項に記載のコンピュータが実施する方法（400）。

30

【請求項7】

前記誤った仮説（222）は、固有名詞に対する音声学的類似性を含む、請求項6に記載のコンピュータが実施する方法（400）。

【請求項8】

前記動作が、前記複数の仮説（222）のうちの生成された仮説に対して前記誤った仮説を置換することをさらに含む、請求項6または7に記載のコンピュータが実施する方法（400）。

【請求項9】

システム（500）であって、
データ処理ハードウェア（510）と、
前記データ処理ハードウェア（510）と通信するメモリハードウェア（520）と、
を備え、前記メモリハードウェア（520）は、命令を格納しており、前記命令は、前記データ処理ハードウェア（510）上での実行時に、前記データ処理ハードウェア（510）に動作を実行させ、前記動作は、

40

最小単語誤り率損失関数を使用して音声認識モデル（200）をトレーニングすることであって、

固有名詞を含むトレーニング例（302）を受信すること、

前記トレーニング例（302）に対応する複数の仮説（222）を生成すること、

前記複数の仮説（222）の各仮説（222）は、前記固有名詞を表し、かつ仮説（222）が前記固有名詞を表す可能性を示す対応する確率を含んでおり、

50

前記複数の仮説(222)のうちの1つに関連付けられた対応する確率がペナルティ基準を満たすことを決定すること、前記ペナルティ基準は、

前記対応する確率が確率しきい値を満たしていること、および前記対応する確率を含む前記仮説(222)が固有名詞を誤って表していることを示しており、

前記最小単語誤り率損失関数にペナルティ(332)を適用することによって、前記音声認識モデル(200)をトレーニングすることを含んでおり、前記システムは、

リカレントニューラルネットワークトランスデューサ(以下、RNN-Tとする)デコーダ(222)を含む第1のパスのネットワーク(206)と、

リッスン・アテンド・スペル(以下、LASとする)デコーダ(230)を含む第2のパスのネットワーク(208)と、をさらに備えており、

前記音声認識モデル(200)は、前記第1のパスのネットワーク(206)および前記第2のパスのネットワーク(208)を含んでおり、

前記最小単語誤り率損失関数を使用したトレーニングは、前記LASデコーダ(230)において行われる、システム(500)。

【請求項10】

前記対応する確率が他の仮説(222)に関連付けられた対応する確率よりも大きい場合、前記対応する確率が前記確率しきい値を満たす、請求項9に記載のシステム(500)。

【請求項11】

前記第1のパスのネットワーク(206)および前記第2のパスのネットワーク(208)の各々に対して音響フレーム(212)をエンコードするように構成された共有エンコーダ(210)をさらに備える、請求項9に記載のシステム(500)。

【請求項12】

前記動作は、前記RNN-Tデコーダ(220)をトレーニングすること、前記最小単語誤り率損失関数を使用してLASデコーダ(230)をトレーニングする前に、トレーニングされた前記RNN-Tデコーダ(220)のパラメータが固定された状態で、前記LASデコーダ(230)をトレーニングすること、をさらに含む、請求項9または11に記載のシステム(500)。

【請求項13】

前記動作は、確率を前記複数の仮説(222)の各仮説(222)に割り当てることをさらに含む、請求項9乃至12のいずれか一項に記載のシステム(500)。

【請求項14】

前記動作は、誤った仮説(222)を受信すること、個別の確率を前記誤った仮説(222)に割り当てること、をさらに含み、前記ペナルティ基準は、仮説(222)が生成された誤った仮説(222)を含むという表示をさらに含む、請求項9乃至13のいずれか一項に記載のシステム(500)。

【請求項15】

前記誤った仮説は、前記固有名詞に対する音声学的類似性を含む、請求項14に記載のシステム(500)。

【請求項16】

前記動作は、前記複数の仮説(222)のうちの生成された仮説(222)に対して前記誤った仮説(222)を置換することをさらに含む、請求項14または15に記載のシステム(500)。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

本開示は、エンド・ツー・エンド音声認識における固有名詞認識に関する。

【背景技術】

【0002】

最新の自動音声認識 (ASR: automated speech recognition) システムは、高品質 (例えば、低い単語誤り率 (WER: word error rate)) のみならず、低遅延 (例えば、ユーザが話してから文字起こし (transcription) が表示されるまでの短い遅延) を提供することに重点を置いている。さらに、現在 ASR システムを使用する場合、ASR システムは、リアルタイムに対応するか、またはリアルタイムよりもさらに高速に対応するストリーミング方式で発話をデコードすることが要求されている。例えば、ユーザとの直接対話を行う携帯電話に ASR システムが搭載されている場合、ASR システムを使用する携帯電話上のアプリケーションは、単語が話されるとすぐに画面上に表示されるように音声認識がストリーミングされることを必要とする場合がある。ここで、携帯電話のユーザは、遅延に対する許容度が低い可能性もある。この低い許容度により、音声認識は、ユーザエクスペリエンスに悪影響を与える可能性のある遅延および不正確性による影響を最小限に抑えるようにモバイルデバイス上で動作することを目指している。

10

【発明の概要】

【0003】

本開示の一態様は、データ処理ハードウェア上での実行時に、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が実施する方法を提供し、動作は、最小単語誤り率損失関数を使用して音声認識モデルをトレーニングすることであって、固有名詞を含むトレーニング例を受信すること、トレーニング例に対応する複数の仮説を生成すること (複数の仮説の各仮説は、固有名詞を表し、かつ仮説が固有名詞を表す可能性を示す対応する確率を含んでいる)、複数の仮説のうちの1つに関連付けられた対応する確率がペナルティ基準を満たすと決定すること、最小単語誤り率損失関数にペナルティを適用することによって、音声認識モデルをトレーニングすることを含む。ペナルティ基準は、対応する確率が確率しきい値を満たしていること、および関連付けられた仮説が固有名詞を誤って表していることを示す。

20

【0004】

本開示の実施形態は、以下の任意の特徴のうちの1つまたは複数を含み得る。いくつかの実施形態では、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) デコーダを含む第1のパスのネットワークと、リッスン・アテンド・スペル (LAS) デコーダを含む第2のパスのネットワークとを含む2パスアーキテクチャを含む。これらの実施形態では、音声認識モデルは、第1のパスのネットワークおよび第2のパスのネットワークの各々に対して音響フレームをエンコードする共有エンコーダをさらに含み得る。これらの実施形態における最小単語誤り率損失関数を使用したトレーニングは、LASエンコーダにおいて行われ得る。動作は、RNN-Tデコーダをトレーニングすること、および最小単語誤り率損失関数を使用してLASデコーダをトレーニングする前に、トレーニングされたRNN-Tデコーダのパラメータが固定された状態で、LASデコーダをトレーニングすることをさらに含み得る。

30

【0005】

いくつかの例では、対応する確率が他の仮説に関連付けられた対応する確率よりも大きい場合、対応する確率が確率しきい値を満たす。動作は、確率を複数の仮説の各仮説に割り当てることをさらに含み得る。いくつかの実施形態では、動作は、誤った仮説を受信すること、および個別の確率を誤った仮説に割り当てることをさらに含み、ペナルティ基準は、仮説が生成された誤った仮説を含むという表示をさらに含む。これらの例では、誤った仮説は、固有名詞に対する音声学的 (phonetically) 類似性を含み得、かつ/または動作は、複数の仮説のうちの生成された仮説に対して誤った仮説に置換することをさらに含み得る。

40

【0006】

50

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と通信するメモリハードウェアとを含むシステムを提供し、メモリハードウェアは、命令を格納し、命令は、データ処理ハードウェアによる実行時に、データ処理ハードウェアに動作を実行させ、動作は、最小単語誤り率損失関数を使用して音声認識モデルをトレーニングすることであって、固有名詞を含むトレーニング例を受信すること、トレーニング例に対応する複数の仮説を生成すること（複数の仮説の各仮説は固有名詞を表し、かつ仮説が固有名詞を表す可能性を示す対応する確率を含んでいる）、複数の仮説のうちの1つに関連付けられた対応する確率がペナルティ基準を満たすと決定すること、最小単語誤り率損失関数にペナルティを適用することによって、音声認識モデルをトレーニングすることを含む。ペナルティ基準は、対応する確率が確率しきい値を満たしていること、および関連付けられた仮説が固有名詞を誤って表していることを示す。

10

【0007】

この態様は、以下の任意の特徴のうちの1つまたは複数を含み得る。いくつかの実施形態では、システムは、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）デコーダを備える第1のパスのネットワークと、リッスン・アテンド・スペル（LAS）デコーダを備える第2のパスのネットワークとをさらに含み、音声認識モデルは、第1のパスのネットワークと第2のパスのネットワークとを備える。これらの実施形態では、システムは、第1のパスのネットワークおよび第2のパスのネットワークの各々に対して音響フレームをエンコードするように構成された共有エンコーダをも含み得る。これらの実施形態における最小単語誤り率損失関数を使用したトレーニングは、LASデコーダで行われ得る。動作は、RNN-Tデコーダをトレーニングすること、および最小単語誤り率損失関数を使用してLASデコーダをトレーニングする前に、トレーニングされたRNN-Tデコーダのパラメータが固定された状態で、LASデコーダをトレーニングすることをさらに含み得る。

20

【0008】

いくつかの例では、対応する確率が他の仮説に関連付けられた対応する確率よりも大きい場合、対応する確率は確率しきい値を満たす。動作は、確率を複数の仮説の各仮説に割り当てることをさらに含み得る。いくつかの実施形態では、動作は、誤った仮説を受信すること、および個別の確率を誤った仮説に割り当てることをさらに含み、ペナルティ基準は、仮説が生成された誤った仮説を含むという表示をさらに含む。これらの例では、誤った仮説は、固有名詞に対する音声学的類似性を含み得、かつ/または動作は、複数の仮説のうちの生成された仮説に対して誤った仮説に置換することをさらに含み得る。

30

【0009】

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

【図面の簡単な説明】

【0010】

【図1A】ジョイント音響モデルおよびテキストモデルを備えた2パス音声認識アーキテクチャを使用する例示的な音声環境の概略図である。

40

【図1B】ジョイント音響モデルおよびテキストモデルを備えた2パス音声認識アーキテクチャを使用する例示的な音声環境の概略図である。

【図2】音声認識のための例示的な2パス音声認識アーキテクチャの概略図である。

【図3A】図2の2パス音声認識アーキテクチャをトレーニングするための例示的なトレーニング手順の概略図である。

【図3B】図2の2パス音声認識アーキテクチャをトレーニングするための例示的なトレーニング手順の概略図である。

【図3C】図2の2パス音声認識アーキテクチャをトレーニングするための例示的なトレーニング手順の概略図である。

【図4】図2の2パス音声認識アーキテクチャをトレーニングする方法のための動作の例

50

示的な構成のフローチャートである。

【図5】本明細書で説明されるシステムおよび方法を実施するために使用され得る例示的なコンピューティングデバイスの概略図である。

【発明を実施するための形態】

【0011】

様々な図面の同様の参照記号は、同様の構成要素を示す。

音声認識は、モバイル環境の非拘束性および機敏性の要求を満たすために進化し続けている。自動音声認識システム（ASR）の品質を向上させるために、新たな音声認識アーキテクチャまたは既存のアーキテクチャの改良が引き続き開発されている。例えば、音声認識は、当初、各モデルが専用の目的を持つ複数のモデルを採用していた。例えば、ASRシステムは、音響モデル（AM）、発音モデル（PM）、および言語モデル（LM）を含んでいた。音響モデルは、音声のセグメント（即ち、音声のフレーム）を音素（phonemes）にマッピングした。発音モデルは、これらの音素をつなぎ合わせて単語を形成し、言語モデルは、所与のフレーズの可能性（即ち、単語のシーケンスの確率）を表現するために使用された。これらの個々のモデルは連携して機能したが、各モデルは個別にトレーニングされ、多くの場合、異なるデータセットで手動で設計された。

【0012】

個別のモデルの手法により、特に所与のモデルに対するトレーニングコーパス（即ち、トレーニングデータの集合体）がモデルの有効性に対応している場合に、音声認識システムの精度をかなり向上させることが可能になったが、個別のモデルを個別にトレーニングする必要性により、それ自体が複雑になるため、統合モデルを備えたアーキテクチャが採用された。これらの統合モデルは、単一のニューラルネットワークを使用して、音波形（即ち、入力シーケンス）を出力センテンス（即ち、出力シーケンス）に直接マッピングしようとするものである。これにより、任意の音声特徴のシーケンスが与えられると、単語（または書記素（graphemes））のシーケンスが生成されるシーケンス・ツー・シーケンスの手法が実現された。シーケンス・ツー・シーケンスモデルの例には、「アテンションベース」モデルおよび「リッスン・アテンド・スペル」（LAS）モデルが含まれる。LASモデルは、リスナー（listener）コンポーネント、アテンダ（attender）コンポーネント、およびスペラー（speller）コンポーネントを使用して、音声の発話を文字に変換する。ここで、リスナーは、音声入力（例えば、音声入力の時間周波数表現）を受信し、音声入力をより高レベルの特徴表現にマッピングするリカレントニューラルネットワーク（RNN：recurrent neural network）エンコーダである。アテンダは、より高レベルの特徴をアテンションして、入力特徴と予測されるサブワード単位（例えば、書記素または単語ピース）との間のアラインメントを学習する。スペラーは、アテンションベースのRNNデコーダであり、仮定単語のセットに対して確率分布を生成することによって、入力から文字シーケンスを生成する。統合化された構造により、モデルの全てのコンポーネントを単一のエンド・ツー・エンド（E2E：end-to-end）ニューラルネットワークとして共同でトレーニングさせることができる。ここで、E2Eモデルとは、アーキテクチャが全てニューラルネットワークで構成されているモデルを指す。完全なニューラルネットワークは、外部コンポーネントおよび/または手動で設計したコンポーネント（例えば、有限状態トランسدューサ、辞書（lexicon）、またはテキスト正規化モジュール）なしで機能する。さらに、E2Eモデルをトレーニングする場合、これらのモデルは通常、決定木からのブートストラップ、または別のシステムからの時間調整を必要としない。

【0013】

初期のE2Eモデルは正確であり、個別にトレーニングされたモデルよりもトレーニングが改善されたが、LASモデルなどのこれらのE2Eモデルは、出力テキストを生成する前に入力シーケンス全体を確認することによって機能していたため、入力が受信されたときに出力をストリーミングすることはできなかった。ストリーミング機能がないと、LASモデルは、リアルタイムの音声文字起こし（voice transcription）

10

20

30

40

50

n) を実行することができない。この欠陥のため、遅延に敏感な、かつ/またはリアルタイムの音声文字起こしを必要とする音声アプリケーションに対してLASモデルを搭載すると、問題が発生する可能性がある。このため、リアルタイムアプリケーション(例えば、リアルタイム通信アプリケーション)に依存することが多いモバイル技術(例えば、携帯電話)にとって、LASモデルだけでは、理想的なモデルではない。

【0014】

さらに、音響モデル、発音モデル、および言語モデル、またはそれらが共に構成されているモデルを有する音声認識システムは、これらのモデルに関連する比較的大規模のサーチグラフをサーチする必要があるデコーダに依存し得る。大規模のサーチグラフでは、この種の音声認識システムを完全オンデバイスでホストするのに有利ではない。ここで、音声認識システムが「オンデバイス(on-device)」でホストされている場合、音声入力を受信するデバイスは、そのプロセッサ(単数または複数)を使用して音声認識システムの機能を実行する。例えば、音声認識システムが完全にオンデバイスでホストされている場合、デバイスのプロセッサは、音声認識システムの機能を実行するために、デバイス外のコンピューティングリソースと連携する必要はない。完全にオンデバイスではない音声認識を実行するデバイスは、音声認識システムの少なくとも一部の機能を実行するために、リモートコンピューティング(例えば、リモートコンピューティングシステムまたはクラウドコンピューティング)、従ってオンライン接続に依存している。例えば、音声認識システムは、サーバベースのモデルとのネットワーク接続を使用して、大規模なサーチグラフによりデコーディングを実行する。

【0015】

残念ながら、リモート接続に依存している状態では、音声認識システムは、遅延の問題および/または通信ネットワークに固有の信頼性の低さに対して脆弱になる。これらの問題を回避することによって音声認識の有用性を向上させるために、音声認識システムは、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)として知られるシーケンス・ツー・シーケンス(sequence-to-sequence)モデルの形態に再び進化した。RNN-Tはアテンション機構を採用しておらず、かつ出力(例えば、センテンス)を生成するためにシーケンス全体(例えば、音声波形)を処理する必要がある他のシーケンス・ツー・シーケンスモデルとは異なり、RNN-Tは、入力サンプルを連続的に処理して、出力シンボルをストリーミングするという、リアルタイム通信にとって特に魅力的な特徴を有している。例えば、RNN-Tを使用した音声認識では、話した通りに文字が1つずつ出力され得る。ここで、RNN-Tは、モデルによって予測されたシンボルを自身にフィードバックするフィードバックループを使用して、次のシンボルを予測する。RNN-Tのデコーディングは、大規模なデコーダグラフではなく、単一のニューラルネットワークを介したビームサーチを含むため、RNN-Tは、サーバベースの音声認識モデルの数分の1のサイズにスケールアップすることができる。サイズの縮小により、RNN-Tは完全にオンデバイスで搭載され、オフラインで(即ち、ネットワーク接続なしで)動作させることができるため、通信ネットワークの信頼性の問題を回避することができる。

【0016】

音声認識システムが低遅延で動作することに加えて、音声認識システムには、音声を正確に認識することが求められる。音声認識を実行するモデルの場合、モデルの精度を定義するメトリックとして、単語誤り率(WER)が用いられることが多い。WERは、実際に話された単語の数と比較して、どれだけ単語が変更されたかを示す尺度である。一般に、これらの単語の変更は、置換(即ち、単語が置き換えられる場合)、挿入(即ち、単語が追加される場合)、および/または削除(即ち、単語が省略される場合)を指す。例えば、話者は「カー(car)」と言っているが、ASRシステムは、「カー(car)」という単語を「バー(bar)」と文字起こしする。これは、音声の(phonetic)類似性による置換の例である。他のASRシステムと比較してASRシステムの能力を測定する場合、WERは、別のシステムまたはあるベースラインと比較して、改善または

10

20

30

40

50

品質保証能力の尺度を示すことができる。

【0017】

RNN-Tモデルは、オンデバイスの音声認識に関する有力な候補モデルとして有望であることを示したが、RNN-Tモデルのみでは、品質（例えば、音声認識精度）の観点で、大規模な最先端の従来モデル（例えば、別個のAM、PM、およびLMを備えたサーバベースのモデル）に遅れをとっている。しかし、非ストリーミングE2E、LASモデルは、大規模な最先端の従来モデルに匹敵する音声認識品質を備えている。非ストリーミングE2E LASモデルの品質を活用するために、RNN-Tネットワークの第1のパスのコンポーネントと、それに続くLASネットワークの第2のパスのコンポーネントを含む2パス音声認識システム（例えば、図2に示す）が開発された。この設計により、2パスモデルは、低遅延のRNN-Tモデルのストリーミング特性の恩恵を受け、LASネットワークを組み込んだ第2のパスを通じてRNN-Tモデルの精度を向上させている。LASネットワークは、RNN-Tモデルのみと比較して遅延を増加させるが、遅延の増加は、適度にわずかであり、かつオンデバイス動作に関する遅延制約に適合している。精度に関しては、2パスモデルは、RNN-T単独と比較した場合に17~22%のWER低減を達成し、大規模な従来モデルと比較した場合に同程度のWERを有している。

10

【0018】

RNN-Tネットワークの第1のパスとLASネットワークの第2のパスを備えた2パスモデルでも、特に稀少な単語または一般的でない単語の場合に、トレードオフがある。これらのタイプの単語は、テイル発話（tail utterances）と呼ばれることがあり、曖昧さ、トレーニングでの希少性、または特殊な言語化によって、音声システムが文字起こしするのが本質的に困難である。テイル発話の例には、アクセントのある話し言葉、異言語間の話し言葉、数字、および固有名詞が含まれる。例えば、固有名詞は、2パスモデルを使用してASRをストリーミングする際の課題を提示する。これは、特定の名前がトレーニング中にまれにしか登場にないか、またはまったく登場しない場合があるものの、より一般的な単語に似た発音を有する可能性があるためである。これまで、従来のモデルは、固有名詞の発音に関する知識を注入することによって、発音モデル（PM：pronunciation model）を最適化して、テイル性能を改善することができる。残念ながら、2パスアーキテクチャには、固有名詞の発音で特別にトレーニングすることができる明示的な発音モデル（PM）と、固有名詞に多く触れる大規模なコーパスでトレーニングすることができる言語モデル（LM：language model）を欠いている。ストリーミング2パスシステムにおいて、適切な既知の知識を注入するための特定の場所としてのPMがなければ、固有名詞の発音などの特定の要件をモデル化することはより困難となる。一部のモデルでは、追加のトレーニングデータまたはモデルを組み込むことによって、一般的でない単語/稀少な単語による問題を改善しようという試みがなされているが、これらの技法では、モデルのサイズ、トレーニング時間、および/または推論コストが増加する。

20

30

【0019】

固有名詞および/またはその他のテイル発話に対する2パスモデルの有効性を高めるために、2パスアーキテクチャでは、カスタマイズされた最小単語誤り率（MWER：minimum word error rate）損失基準が使用される。この損失基準は、特に固有名詞の認識を強調することを目的としている。損失基準を使用して固有名詞の認識を向上させることによって、音声認識システムは、トレーニング時の新たなデータまたは推論時の外部モデルを必要としない。ここで、損失基準の2つの異なる方法を固有名詞認識に関して使用することができる。第1の方法は、グラウンドトゥールの文字起こしにおいて固有名詞を識別して、トレーニング中に固有名詞を外している仮説の損失を増加させるエンティティタグ付けシステムを含む。第2の方法は、MWERビームに追加の仮説を注入することであり、追加の仮説は、音声学的に類似した代替語（alternatives）に置換された固有名詞に対応する。例えば、「ウォルマート（Walmart）」に音声学的に類似した代替語として、「ホールマーク（Hallmark）」という

40

50

追加の仮説が追加される。第2の手法では、トレーニングのプロセスによって、可能性のある間違いと潜在的な代替語をモデルに認識させる。様々な固有名詞テストセットにおいて、これらのカスタム損失基準方法は、カスタム損失基準のない従来の2パスアーキテクチャと比較して、WERを相対的に2～7%削減することができる。

【0020】

図1Aおよび図1Bは、発話環境100の例である。発話環境100において、ユーザデバイス110などのコンピューティングデバイスと対話するユーザの10の方法は、音声入力を介するものであり得る。ユーザデバイス110（一般にデバイス110とも呼ばれる）は、発話対応環境100内の1人または複数人のユーザ10からの音（例えば、ストリーミング音声データ）をキャプチャするように構成されている。ここで、ストリーミング音声データ12は、デバイス110によってキャプチャされる可聴の問い合わせ（クエリ）、デバイス110に対する命令（コマンド）、または可聴の会話（コミュニケーション）としての役割を持つ、ユーザ10によって話された発話を指すことができる。デバイス110の発話対応システムは、問い合わせに応答することによって、かつ/またはコマンドを実行させることによって、問い合わせまたは命令を処理し得る。

【0021】

ユーザデバイス110は、ユーザ10に関連付けられ、かつ音声データ12を受信することが可能な任意のコンピューティングデバイスに対応し得る。ユーザデバイス110のいくつかの例は、モバイルデバイス（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、ウェアラブルデバイス（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）デバイス、スマートスピーカなどを含むが、これらに限定されない。ユーザデバイス110は、データ処理ハードウェア112と、データ処理ハードウェア112と通信するメモリハードウェア114とを含み、メモリハードウェア114は、命令を格納し、命令は、データ処理ハードウェア112による実行時に、データ処理ハードウェア112に1つまたは複数の動作を実行させる。ユーザデバイス110は、発話対応システム100内で話された発話12をキャプチャして電気信号に変換するための音声キャプチャデバイス（例えば、マイクロフォン）116、116aと、可聴音声信号を（例えば、デバイス110からの出力音声データとして）伝達するための発話出力デバイス（例えばスピーカ）116、116bとを有する音声サブシステム116をさらに含む。図示される例では、ユーザデバイス110は単一の音声キャプチャデバイス116aを実装しているが、ユーザデバイス110は、本開示の範囲から逸脱することなく、音声キャプチャデバイス116aのアレイを実装してもよく、それにより、アレイ内の1つまたは複数のキャプチャデバイス116aは、ユーザデバイス110上に物理的に存在していないが、音声サブシステム116と通信状態になり得る。（例えば、ハードウェア112、114を使用する）ユーザデバイス110は、音声認識器200を使用して、ストリーミング音声データ12に対して音声認識処理を実行するようにさらに構成される。いくつかの例では、音声キャプチャデバイス116aを含むユーザデバイス110の音声サブシステム116は、音声データ12（例えば、話された発話）を受信し、音声データ12を音声認識器200と互換性のあるデジタル形式に変換するように構成される。デジタル形式は、メルフレーム（mel frames）などの音響フレーム（例えば、パラメータ化された音響フレーム）に対応し得る。例えば、パラメータ化された音響フレームは、ログメルフィルタバンク（log-mel filterbank）エネルギーに対応する。

【0022】

図1Aなどのいくつかの例では、ユーザ10は、音声認識器200を使用するユーザデバイス110のプログラムまたはアプリケーション118と対話する。例えば、図1Aは、ユーザ10が自動アシスタントアプリケーションと通信している状態を示している。この例では、ユーザ10が自動アシスタントに「今夜のコンサートは何時から？（What time is the concert tonight?）」と尋ねている。ユーザ10からのこの質問は、音声キャプチャデバイス116aによってキャプチャされ、ユーザ

10

20

30

40

50

デバイス 110 の音声サブシステム 116 によって処理される話された発話 12 である。この例では、ユーザデバイス 110 の音声認識器 200 は、「今夜のコンサートは何時から」という音声入力 202 を（例えば、音響フレームとして）受信し、音声入力 202 を文字起こし 204（例えば、「今夜のコンサートは何時から？」というテキスト表現）に転写する。ここで、アプリケーション 118 の自動アシスタントは、自然言語処理を使用して、ユーザ 10 によって投げかけられた質問に対して回答し得る。自然言語処理とは、一般に、書き言葉（例えば、文字起こし 204）を解釈し、書き言葉が何らかのアクションを促しているかどうかを判断するプロセスを指す。この例では、自動アシスタントは、自然言語処理を使用して、ユーザ 10 からの質問がユーザのスケジュール、より具体的にはユーザのスケジュールでのコンサートに関するものであることを認識する。自動アシスタントは、自然言語処理でこれらの詳細を認識することによって、ユーザの問い合わせに対して、「今夜のコンサートは午後 8 時 30 分に開場します (Doors open at 8:30 pm for the concert tonight)」という回答を返す。いくつかの構成では、自然言語処理は、ユーザデバイス 110 のデータ処理ハードウェア 112 と通信するリモートシステム上で行われ得る。

【0023】

図 1B は、音声認識器 200 を用いた音声認識の別の例である。この例では、ユーザデバイス 110 に関連付けられたユーザ 10 は、通信アプリケーション 118 を用いてジェーン・ドゥという名前の友人と通信している。ここで、テッドという名前のユーザ 10 は、音声認識器 200 に自身の音声入力を文字起こしさせることによって、ジェーンと通信する。音声キャプチャデバイス 116 は、これらの音声入力をキャプチャし、それらを音声認識器 200 にデジタル形式（例えば、音響フレーム）で伝達する。音声認識器 200 は、これらの音響フレームを、通信アプリケーション 118 を介してジェーンに送信されるテキストに文字起こしする。この種類のアプリケーション 118 はテキストを介して通信するため、音声認識器 200 からの文字起こし 204 は、さらなる処理（例えば、自然言語処理）なしでジェーンに送信され得る。

【0024】

図 2 などのいくつかの例では、音声認識器 200 は、2 パス音声認識アーキテクチャ（または単に「2 パスアーキテクチャ」）で構成される。一般的に、音声認識器 200 の 2 パスアーキテクチャは、少なくとも 1 つのエンコーダ 210、RNN-T デコーダ 220、および LAS デコーダ 230 を含む。2 パスデコーディングにおいて、第 2 のパス 208（例えば、LAS デコーダ 230 として示される）は、第 1 のパス 206（例えば、RNN-T デコーダ 220 として示される）からの初期出力を格子再スコアリング (lattice rescoring) または n ベスト再ランク付け (n-best re-ranking) などの技術を用いて改善し得る。言い換えれば、RNN-T デコーダ 220 はストリーミング予測（例えば、1 組の N ベスト仮説 (N-best hypotheses)）を生成し、LAS デコーダ 230 は、予測を確定する（例えば、1 つのベストの再スコアリングされた仮説を識別する）。ここで、具体的には、LAS デコーダ 230 は、RNN-T デコーダ 220 からストリーミングされた仮説 y_R を再スコアリングする。一般に、LAS デコーダ 230 は、RNN-T デコーダ 220 からストリーミングされた仮説 y_R を再スコアリングする再スコアリングモードで機能すると説明されているが、LAS デコーダ 230 は、設計またはその他の要因（例えば、発話の長さ）に応じてビームサーチモードなどの異なるモードで動作することも可能である。

【0025】

少なくとも 1 つのエンコーダ 210 は、ストリーミング音声データ 12 に対応する音響フレームを音声入力 202 として受信するように構成される。音響フレームは、音声サブシステム 116 によってパラメータ化された音響フレーム（例えば、メルフレームおよび/またはスペクトルフレーム）に事前に処理され得る。いくつかの実施形態では、パラメータ化された音響フレームは、ログメル特徴 (log-mel features) を有するログメルフィルタバンクエネルギーに対応する。例えば、音声サブシステム 116 に

10

20

30

40

50

よって出力され、エンコーダ 210 に入力されるパラメータ化された入力音響フレームは、 $x = (x_1, \dots, x_T)$ として表すことができ、ここで、

【0026】

【数1】

$$x_t \in \mathbb{R}^d$$

は、ログメルフィルタバンクエネルギーであり、 T は x のフレーム数を示し、 d はログメル特徴の数を表す。いくつかの例では、各パラメータ化された音響フレームは、短いシフトウィンドウ（例えば、32ミリ秒、10ミリ秒ごとにシフト）内で計算された128次元のログメル特徴を含む。各特徴は、前のフレーム（例えば、3つ前のフレーム）と重ねられて、より高次元のベクトル（例えば、3つ前のフレームを使用した512次元のベクトル）が形成されてもよい。次に、ベクトルを形成する特徴は、（例えば、30ミリ秒のフレームレートに）ダウンサンプリングされ得る。エンコーダ 210 は、音声入力 202 に基づいて、エンコーディング e を生成するように構成される。例えば、エンコーダ 210 は、エンコードされた音響フレーム（例えば、エンコードされたメルフレームまたは音響埋め込み（acoustic embeddings））を生成する。

【0027】

エンコーダ 210 の構造は、異なる方法で実施することができるが、いくつかの実施形態では、エンコーダ 210 は、長短期記憶（LSTM: long-short term memory）ニューラルネットワークである。例えば、エンコーダ 210 は、8個の LSTM 層を含む。ここで、各層は、2,048個の隠れユニットと、それに続く640次元の射影層（projection layer）とを含む。いくつかの例では、エンコーダ 210 の第2の LSTM 層の後に、短縮係数（reduction factor） $N=2$ を有する時間短縮層（time-reduction layer）が挿入される。

【0028】

いくつかの構成では、エンコーダ 210 は共有エンコーダネットワークである。言い換えれば、各パスネットワーク 206、208 がそれ自体の別個のエンコーダを有する代わりに、各パス 206、208 は単一のエンコーダ 210 を共有する。エンコーダを共有することによって、2パスアーキテクチャを使用する ASR 音声認識器 200 は、そのモデルサイズおよび/またはその計算コストを削減することができる。ここで、モデルサイズの縮小は、音声認識器 200 が完全にオンデバイス（on-device）で良好に機能することを可能にするのに役立つ。

【0029】

いくつかの例では、図2の音声認識器 200 は、LAS デコーダ 230 の第2のパス 208 に適したものとなるようにエンコーダ 210 の出力 212 を適応させるための音響エンコーダ 240 などの追加のエンコーダをも含む。音響エンコーダ 240 は、出力 212 をエンコードされた出力 252 にさらにエンコードするように構成される。いくつかの実施形態では、音響エンコーダ 240 は、エンコーダ 210 からの出力 212 をさらにエンコードする LSTM エンコーダ（例えば、2層 LSTM エンコーダ）である。追加のエンコーダを含むことによって、エンコーダ 210 は、パス 206、208 の間の共有エンコーダとして依然として保持され得る。

【0030】

第1のパス 206 を通じて、エンコーダ 210 は、音声入力 202 の各音響フレームを受信して、出力 212（例えば、音響フレームのエンコーディング e として示される）を生成する。RNN-T デコーダ 220 は、各フレームの出力 212 を受信して、仮説 y_R として示される出力 222 を各タイムステップにおいてストリーミング方式で生成する。言い換えれば、RNN-T デコーダ 220 は、フレーム毎の埋め込み e または出力 212 を消費して、単語ピース出力 222 を仮説として生成し得る。いくつかの例では、RNN-T デコーダ 220 は、受信したエンコードされた音響フレーム 212 に基づいてビーム

10

20

30

40

50

サーチを実行することによって、Nベスト仮説222を生成する。RNN-Tデコーダ220の構造に関して、RNN-Tデコーダ220は、予測ネットワークおよび結合ネットワークを含み得る。ここで、予測ネットワークは、2,048個の隠れユニットおよび640次元の射影(層ごと)の2つのLSTM層、並びに128ユニットの埋め込み層を有し得る。エンコーダ210および予測ネットワークの出力212は、ソフトマックス予測層を含む結合ネットワークに供給され得る。いくつかの例では、RNN-Tデコーダ220の結合ネットワークは、640個の隠れユニットと、それに続く4,096個の大文字と小文字が混在する単語ピースを予測するソフトマックス層とを含む。

【0031】

図2の2パスモデルにおいて、第2のパス208を通じて、LASデコーダ230は、各フレームに関してエンコーダ210からの出力212を受信し、仮説 y_L として指定された出力232を生成する。LASデコーダ230がビームサーチモードで動作する場合、LASデコーダ230は、出力212のみから出力232を生成し、RNN-Tデコーダ220の出力222を無視する。LASデコーダ230が再スコアリングモードで動作する場合、LASデコーダ230は、RNN-Tデコーダ220から(例えば、RNN-Tデコーダ220によって生成されたNベスト仮説に対応する)トップK仮説222、 y_R を取得し、次いで、LASデコーダ230は、出力212をアテンション(attention)しつつ、教師強制モードで各シーケンスに対して動作して、スコアを計算する。例えば、スコアは、シーケンスの対数確率とアテンションカバレッジペナルティ(attention coverage penalty)とを組み合わせたものである。LASデコーダ230は、最も高いスコアを有するシーケンスを出力232として選択する。言い換えれば、LASデコーダ230は、RNN-Tデコーダ220からの仮説222のNベストラストから、最大尤度(maximum likelihood)を有する単一の仮説 y_R を選択し得る。ここで、再スコアリングモードでは、LASデコーダ230は、出力212をアテンションするために、(例えば、4つのヘッドを有する)マルチヘッドアテンション(multi-headed attention)を含み得る。さらに、LASデコーダ230は、予測のためのソフトマックス層を備えた2層LASデコーダ230であり得る。例えば、LASデコーダ230の各層は、2,048個の隠れユニットと、それに続く640次元の射影とを有する。ソフトマックス層は、RNN-Tデコーダ220のソフトマックス層から同じ大文字と小文字が混在する単語ピースを予測するために、4,096次元を含み得る。

【0032】

ニューラルネットワークは、通常、損失関数(例えば、クロスエントロピー損失関数)を定義するバックプロパゲーションによってトレーニングされる。例えば、損失関数は、ネットワークの実際の出力とネットワークの所望の出力との差として定義される。クロスエントロピー(CE)損失関数を使用してモデルをトレーニングするために、モデルは、トレーニングデータの対数尤度(log-likelihood)を最大化することによって、CE損失関数を最適化するようにトレーニングする。図3A~図3Cを参照すると、トレーニング手順300は、音声認識器200の各コンポーネントを、対応する組のトレーニングデータ302、302a-dでトレーニングすることができる。トレーニング手順300は、データ処理ハードウェア510(図5)と、データ処理ハードウェアと通信するメモリハードウェア520(図5)を含むシステム500上で実行することができ、メモリハードウェア520は命令を格納し、命令は、データ処理ハードウェア510上での実行時に、処理ハードウェア510に動作を実行させる。例えば、図2の音声認識器200の2パスモデルアーキテクチャをトレーニングするためのトレーニング手順300は、3つの段階310、320、330で行われ得る。第1段階310の間、トレーニング手順300は、(例えば、CE損失関数を使用して)エンコーダ210およびRNN-Tデコーダ220をトレーニングする。いくつかの例では、トレーニング手順300は、 $P(y_R = y | x)$ を最大化するようにエンコーダ210およびRNN-Tデコーダ220をトレーニングする。第2段階320の間、トレーニング手順300は、エンコーダ

210またはRNN-Tデコーダ220のパラメータを更新することなく、LASデコーダ230をトレーニングする。いくつかの実施形態では、トレーニング手順300は、損失を強制することを教示するクロスエントロピーを使用してLASデコーダ230をトレーニングする。例えば、トレーニング手順300は、 $P(y_L = y | x)$ を最大化するようにLASデコーダ230をトレーニングする。第3段階330の間、トレーニング手順300はさらに、 n ベスト仮説を使用することによって期待される単語誤り率を最適化するために、最小WER(MWER)損失を用いてLASデコーダ230をトレーニングする。例えば、WER目的関数は、 N ベストビームの仮説222における単語誤りの加重平均として損失をモデル化する。この第3段階330の間、LASデコーダ230は、以下の式によって表されるWER目的関数に従って微調整され得る。

【0033】

【数2】

$$L_{MWER}(x, y^*) = \sum_{y \in B_{LAS}} P(y|x) \hat{W}(y|y^*) \quad (1)$$

ここで、 y^* は、グラウンドトゥルースであり、 B_{LAS} は、ビームサーチ中のLASデコーダ230からの仮説の N ベストリストであり、 $P(y|x)$ は、仮説 y の正規化された事後確率(normalized posterior)であり、

【0034】

【数3】

$$\hat{W}(y|y^*)$$

は、仮説 y における単語誤り数とビーム全体の単語誤りの平均数との間の差を表す。いくつかの実施形態では、LASデコーダ230が再スコアリング器として機能する場合、LASデコーダ230は、RNN-Tデコーダ220からの最良の仮説 y_R に高い可能性を割り当てることを最適化するようにトレーニングする。ここで、この損失最適化関数は次の式で表すことができる。

【0035】

【数4】

$$L_{MWER}(x, y^*) = \sum_{y \in B_{RNN-T}} P(y|x) \hat{W}(y|y^*) \quad (2)$$

ここで、 B_{RNN-T} は、RNN-Tデコーダ220でのビームサーチから取得される。ここで、これらの最適化モデルの各々は、損失基準が、音声認識器200またはその一部が確率質量を割り当てるように学習すべき分布を表していることを示している。

【0036】

図3Bを参照すると、いくつかの実施形態では、第3トレーニング段階330または微調整段階の間、トレーニング手順300は、MWER損失を使用してトレーニングを実行するが、修正損失関数 $MWER_{AUG}$ を使用する。ここで、修正損失関数 $MWER_{AUG}$ は、固有名詞の損失拡張(loss augmentation)の一形態である。このトレーニング手法では、損失は、トレーニングにおける固有名詞の性能を強調するように構成される。いくつかの例では、モデルが固有名詞を正確に外している仮説 y に高い確率を割り当てるときに、モデル(例えば、LASデコーダ230)に適用されるペナルティ332を増加させることによって、損失は、固有名詞の性能を強調する。例示すると、図3Bは、トレーニング手順300の第3段階330の間に、LASデコーダ230が、入力

10

20

30

40

50

302dを予測する1組の可能性のある仮説 y_L 、 y_{L1-3} を生成することを示している。ここで、入力302dは、固有名詞Pnを含んでいるが、LASデコーダ230は、固有名詞Pnを実際には含んでいないにも関わらず、入力302dに対する最高確率の仮説 y_L である第1の仮説 y_{L1} を識別する。この例では、LASデコーダ230が固有名詞Pnを誤って識別した仮説 y_L に最高確率の仮説 y_L を割り当てたため、修正損失関数MWERAugは、ペナルティ332を適用する。いくつかの構成では、トレーニング手順300は、モデル(例えば、LASデコーダ230)がペナルティ基準を満たす仮説yに確率を割り当てたと判断する。ペナルティ基準は、モデルが、確率しきい値を満たす(例えば、確率しきい値に割り当てられた値を超える)、固有名詞に対する誤った仮説に確率を割り当てたことを含み得る。ここで、確率しきい値は、誤った仮説に対する許容可能なレベルまたは値を示す事前構成された値であり得る。これらの例では、トレーニング手順300が、モデル(例えば、LASデコーダ230)がペナルティ基準を満たす仮説yに確率を割り当てたと判断した場合、トレーニング手順300は、修正損失関数にペナルティ332を適用する。いくつかの例では、固有名詞の損失拡張に対する修正損失関数は、次の式で表される。

10

【0037】

【数5】

$$L_{AUG}(x, y^*) = \sum_{y \in B_{RNN-T}} P(y|x) \hat{W}(y|y^*) \cdot C_\lambda(y|y^*) \quad (3)$$

20

ここで

$$C_\lambda(y, y^*) = \begin{cases} \lambda & y^* \text{が } y \text{ にない固有名詞を含む場合} \\ 1, & \text{その他の場合} \end{cases} \quad (4)$$

定数 $\lambda > 1$ である。ここで、 λ は、一般的な発話12に対する音声認識器200の性能に関して、固有名詞の認識の有効性をバランスさせるために選択されたハイパーパラメータを指す。例えば、ハイパーパラメータの設定は、他の誤りのタイプとのトレードオフで、固有名詞の誤りに起因する勾配の増加を回避しようとするものである。いくつかの構成では、各グラウンドトゥルース文字起こし(例えば、トレーニングデータ302d)に対する固有名詞Pnは、固有名詞識別システム340によるトレーニングの前に識別される。仮説yが固有名詞Pnを含むことを保証するために、仮説yが固有名詞Pnの単語列全体を適切な順序で含む場合に、仮説yが固有名詞Pnを含むものとして定義される。例えば、固有名詞Pn「シーダー・ラピッズ(Cedar Rapids)」は、仮説「シーダー・ラピッズの人口(Population of Cedar Rapids)」には含まれているが、仮説「シーダー・ツリーの高さ(Cedar tree height)」または「シーダー・ラピッドエスエスエス(Cedar Rapids s s s)」には含まれていない。

30

40

【0038】

図3Cは、ファズトレーニング(fuzz training)を適用して、音声認識器200が固有名詞を区別する能力を最適化するトレーニング手順300の別の例を示す。この手法では、ファズトレーニングは、固有名詞と、音声学的に類似した誤った代替語とを区別する方法を音声認識器200に教えることを目的としている。言い換えれば、ファズトレーニングにおいて、音声認識器200のようなモデルは、モデルが起こり得る間違いおよび同音異綴り(alternative spellings)の知識を得ることを可能にする。トレーニング中に、モデル(例えば、LASデコーダ230)が固有名詞の間違いに高い可能性を割り当てると、トレーニング手順300は、モデルにペナルティ332を課す。ペナルティ332を課すことによって、トレーニングは、将来同様の誤

50

りの可能性を減少させることを意図している。

【 0 0 3 9 】

これらの潜在的な間違いに関して音声認識器 2 0 0 (例えば、音声認識器 2 0 0 の L A S デコーダ 2 3 0) をトレーニングするために、ファストレーニングは、ビーム修正を実行し得る。一般に、ビームサーチは、最適なポテンシャル解 (例えば、仮説または候補) を評価する数を指定するビームサイズまたはビーム幅パラメータ B を含む。ファストレーニングは、ビームサーチからの仮説 y を置換するか、ビームサーチからの仮説 y の数を拡張するかのいずれかによって、ビームサイズ B を活用することができる。例示すると、図 3 C は、5 個の仮説 y_L、y_{L 1 - 5} に対応する 5 のビームサイズ、または 5 個の仮説 y_L、y_{L 1 - 5} に拡張された 3 のビームサイズを有するビームサーチを示す例である。この例では、ビームサイズが 5 である場合、ファジリングシステム 3 5 0 は、仮説のうちの 2 個を誤った固有名詞の代替語 3 5 2、3 5 2 a - b に置換し得る。同様に、ビームサイズが 3 である場合、ファジリングシステム 3 5 0 は、誤った固有名詞の代替語 3 5 2、3 5 2 a - b を含む追加の仮説 y を生成し得る。いくつかの実施形態では、ファジリングシステム 3 5 0 は、トレーニングデータ 3 0 2 に含まれる固有名詞 P_n に音声学的に類似する代替語 3 5 2 を生成する音声ファジリングと呼ばれる技法を使用して、固有名詞の代替語 3 5 2 を生成する。音声ファジリングにより、ファジリングシステム 3 5 0 は、トレーニングデータ 3 0 2 の従来型のコーパスでは強調されていないか、または含まれていない可能性がある新たな単語または同音異綴りを生成し得る。y_{BRNN-T} であり、かつグラウンドトゥルース y* に対応する仮説に関して、ファズ処理は次の式で表すことができる。

10

20

【 0 0 4 0 】

【数 6】

$$Fuzz(y, y^*) = \begin{cases} y^{fuzz} & y^* \text{ および } y \text{ が固有名詞を共有している場合} \\ y, & \text{その他の場合} \end{cases} \quad (5)$$

いくつかの構成では、ファズ仮説 y^{fuzz} は、y をコピーし、固有名詞 P_n の発生を音声学的に類似した代替語 3 5 2 に置換することによって形成される。ファストレーニングでは、損失関数は、RNN-T デコーダ 2 2 0 からの元のビームを代替語 3 5 2 (ファズまたはファズ仮説とも呼ばれる) と組み合わせることによって定義される。以下の式は、ファストレーニングによるトレーニング手順 3 0 0 間の損失関数を表し得る。

30

【 0 0 4 1 】

【数 7】

$$L_{Fuzz}(x, y^*) = \sum_{y \in BRNN-T \cup Fuzz(BRNN-T)} P(y|x) \hat{W}(y|y^*) \quad (6)$$

ここで、P (y | x) は、(例えば、付加項「F u z z (B R N N - T)」で表されるような) 修正されたビームサイズを考慮した再正規化された事後確率に対応する。いくつかの実施形態では、ファストレーニングの損失関数 L は、ハイパーパラメータ (0 1) をも含み、ハイパーパラメータは、ファストレーニング損失関数 L_{F u z z} を使用する確率を規定するようになっている。これらの実施形態において、トレーニング手順 3 0 0 がファストレーニング損失関数 L_{F u z z} を使用しない場合、トレーニング手順 3 0 0 は式 (2) によって表される損失関数を使用する。ハイパーパラメータ は任意の確率に設定され得るが、いくつかの構成では、ハイパーパラメータは、トレーニング手順 3 0 0 が常にファストレーニング損失関数 L_{F u z z} を組み込むように、1 に設定される。

40

【 0 0 4 2 】

いくつかの構成では、トレーニング手順 3 0 0 は、ファストレーニングの前に、トレー

50

ニングデータセット 302 に含まれる各固有名詞 P_n に対して一定数の代替語 352 (例えば、25 個の代替語 352) を決定する。ここで、ファズトレーニングの前に生成される代替語 352 の数は、計算コストを最小限に抑えつつ、代替語 352 の多様性を確保するように構成することができる。トレーニング手順 300 が、ファズトレーニングの前に一定数の代替語 352 を生成する場合、ファズトレーニング中に、トレーニング手順 300 は、必要に応じて、既に生成されたランダムな代替語 352 を選択し得る。

【0043】

引き続き図 3C を参照すると、第 3 段階 330 の間に、トレーニング手順 300 は、ファズトレーニングを使用して LAS デコーダ 230 をトレーニングする。ここで、LAS デコーダ 230 は、固有名詞 P_n を含むトレーニングデータ 302、302d を受信し、
 トレーニングデータ 302d の固有名詞 P_n に対応する 5 個の仮説 y_L 、 y_{L1-5} (例えば、ビーム幅 $B = 5$) を生成する。また、LAS デコーダ 230 は、各仮説 y_L に、特定の仮説が入力 (例えば、トレーニングデータ 302) を正しく識別すると LAS デコーダ 230 が考える可能性を示す確率 (たとえば、0.2、0.2、0.1、0.4、0.1 として示される) を割り当てる。この例では、ファジングシステム 350 は、(例えば、代替語 352 がファズトレーニングの前に生成された場合) 1 組の潜在的な仮説 y_L 、「ベルモント (Belmont)」および「ブームント (Boomundt)」を含むように 2 個のファズ仮説 352a-b を生成または選択する。この例に示されるように、LAS デコーダ 230 は、誤った代替語「ベルモント」352a に最も高い可能性 (例えば、0.4 として示される) を割り当てる。LAS デコーダ 230 は、誤った代替語 352 に最も高い可能性を割り当てているので、トレーニング手順 300 はファズトレーニング損失関数 L_{Fuzz} にペナルティ 332 を適用する。ここで、ペナルティ 332 などのペナルティは、ニューラルネットワークの重みまたはパラメータを調整するためにトレーニング中にフィードバックを提供する。一般に、ペナルティは、特定の入力に適用される重みを操作して、不所望の出力または不正確な出力ではなく、意図された出力に近づくか、または示すように機能する。言い換えれば、ペナルティ 332 は、LAS デコーダ 230 が将来、誤った代替語 352 が最良の仮説 y である可能性が高いことを示す可能性を低減するように機能する。

【0044】

図 4 は、音声認識モデル (例えば、音声認識器 200) をトレーニングする方法 400 の例示的な動作の構成のフローチャートである。方法 400 は、動作 402 - 408 によって、最小単語誤り率 (MWER) 損失関数を用いて音声認識モデルをトレーニングする。動作 402 において、方法 400 は、固有名詞 P_n を含むトレーニング例 302 を受信する。動作 404 において、方法 400 は、トレーニング例 302 に対応する複数の仮説 y を生成する。ここで、複数の仮説の各仮説 y は、固有名詞 P_n を表し、各仮説には、個々の仮説 y に対する可能性を示す確率が割り当てられる。動作 406 において、方法 400 は、仮説 y に関連付けられた確率がペナルティ基準を満たすことを決定する。ペナルティ基準は、(i) 確率が確率しきい値を満たしていること、および (ii) 仮説が固有名詞を誤って表していることを示す。動作 408 において、方法 400 は、ペナルティ 332 を最小単語誤り率損失関数に適用する。

【0045】

図 5 は、本明細書で説明されるシステム (例えば、音声認識器 200) および方法 (例えば、方法 400) を実施するために使用され得る例示的なコンピューティングデバイス 500 の概略図である。コンピューティングデバイス 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代表することが意図されている。本明細書に示された構成要素、それらの接続および関係、およびそれらの機能は、例示的なものに過ぎず、本明細書に記載および/または特許請求の範囲に記載される本発明の実施形態を限定するものではない。

【0046】

コンピューティングデバイス 500 は、プロセッサ 510（例えば、データ処理ハードウェア）、メモリ 520（例えば、メモリハードウェア）、ストレージデバイス 530、メモリ 520 および高速拡張ポート 540 に接続する高速インタフェース/コントローラ 540、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 を含む。構成要素 510、520、530、540、550、および 560 の各々は、様々なバスを使用して相互接続され、かつ共通のマザーボード上に、または適切な他の方法で搭載され得る。プロセッサ 510 は、メモリ 520 またはストレージデバイス 530 に格納された命令を含むコンピューティングデバイス 500 内での実行のための命令を処理して、高速インタフェース 540 に接続されたディスプレイ 580 などの外部入力/出力デバイス上にグラフィカルユーザインタフェース（GUI）用のグラフィカル情報を表示する。他の実施形態では、複数のメモリおよび複数のタイプのメモリと共に、複数のプロセッサおよび/または複数のバスが適宜使用されてもよい。また、複数のコンピューティングデバイス 500 が接続され、各デバイスが（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）必要な処理の一部を提供してもよい。

【0047】

メモリ 520 は、コンピューティングデバイス 500 内に非一時的に情報を記憶する。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、または不揮発性メモリユニットであってもよい。非一時的メモリ 520 は、コンピューティングデバイス 500 による使用のための一時的または永久的な基準でプログラム（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）を格納するために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には、これらに限定されないが、フラッシュメモリおよび読み出し専用メモリ（ROM）/プログラム可能読み出し専用メモリ（PROM）/消去可能プログラム可能読み出し専用メモリ（EPROM）/電子消去可能プログラム可能読み出し専用メモリ（EEPROM）（例えば、通常、ブートプログラムなどのファームウェアに使用される）が含まれる。揮発性メモリの例には、これらに限定されないが、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）、およびディスクまたはテープが含まれる。

【0048】

ストレージデバイス 530 は、コンピューティングデバイス 500 の大容量ストレージデバイスを提供することができる。いくつかの実施形態では、ストレージデバイス 530 は、コンピュータ可読媒体である。種々の異なる実施形態では、ストレージデバイス 530 は、フロッピーディスク（登録商標）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その他の構成におけるデバイスを含むデバイスのアレイであり得る。追加の実施形態では、コンピュータプログラム製品は、情報媒体に有形的に具体化される。コンピュータプログラム製品は、実行時に、上記したような 1 つまたは複数の方法を実行する命令を含む。情報媒体は、メモリ 520、ストレージデバイス 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。

【0049】

高速コントローラ 540 は、コンピューティングデバイス 500 の帯域幅を大量に使用する処理を管理し、低速コントローラ 560 は、より低い帯域幅を大量に使用する処理を管理する。このような役割の配分は、例示的なものに過ぎない。いくつかの実施形態では、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する）、および各種拡張カード（図示せず）を受け入れる高速拡張ポート 550 に接続される。いくつかの実施形態では、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 に接続される。様々な通信ポート（例えば、USB、ブルートゥース（登録商標）、イーサネット（登録商標

10

20

30

40

50

)、無線イーサネット(登録商標)を含む低速拡張ポート590は、キーボード、ポインティングデバイス、スキャナ、または例えばネットワークアダプターを介するスイッチまたはルータなどのネットワークデバイスなどの1つまたは複数の入力/出力デバイスに接続され得る。

【0050】

コンピューティングデバイス500は、図面に示されるように、いくつかの異なる形態で実施することができる。例えば、標準サーバ500aとして、またはそのようなサーバ500aのグループ内で複数回、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の一部として実施することができる。

【0051】

本明細書に記載のシステムおよび技術の様々な実施形態は、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において実現することができる。これらの様々な実施形態は、ストレージシステム、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、それらにデータおよび命令を送信するように接続された、特別または一般的な目的であってもよい、少なくとも1つのプログラム可能なプロセッサを含む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における実施形態を含むことができる。

【0052】

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ている)は、プログラマブルプロセッサ用の機械命令を含み、高水準の手続き型言語および/またはオブジェクト指向のプログラミング言語、および/またはアセンブリ言語/機械語で実施することができる。本明細書で使用する場合、「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、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、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む、プログラマブルプロセッサに機械命令および/またはデータを提供するために使用される装置および/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るために使用される任意の信号を指す。

【0053】

本明細書で説明するプロセスおよび論理フローは、入力データを処理して出力を生成することによって機能を実行する1つまたは複数のコンピュータプログラムを実行する1つまたは複数のプログラマブルプロセッサによって実行することができる。プロセスおよび論理フローは、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)などの特定用途論理回路によっても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は、一例として、汎用マイクロプロセッサおよび専用マイクロプロセッサの両方、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を含む。一般に、プロセッサは、読み出し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およびデータを受信する。コンピュータの必須要素は、命令を実行するプロセッサと、命令および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ストレージデバイス(例えば、磁気ディスク、光磁気ディスク、または光ディスク)からのデータを受信するか、またはデータを転送するか、あるいはその両方を行うように動作可能に結合される。しかしながら、コンピュータはそのようなデバイスを有する必要はない。コンピュータプログラム命令およびデータを格納するのに適したコンピュータ可読媒体には、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CDROMおよびDVD-

10

20

30

40

50

R O Mディスクを含む全ての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る。プロセッサおよびメモリは、特定用途論理回路によって補完または特定用途論理回路に組み込むことができる。

【 0 0 5 4 】

ユーザとのインタラクションを提供するために、本開示の1つまたは複数の態様は、例えば、C R T（陰極線管）、L D C（液晶ディスプレイ）モニタ、またはタッチスクリーンなどのユーザに情報を表示するためのディスプレイデバイスと、任意選択でユーザがコンピュータに入力を提供するキーボードおよびポインティングデバイス（例えば、マウスやトラックボール）とを有するコンピュータ上で実施することができる。他の種類の装置を使用して、例えば、任意の形態の感覚フィードバック（例えば、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）であり得るユーザに提供されるフィードバックとともにユーザとのインタラクションを提供することもでき、ユーザからの入力は、音響、音声、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。さらに、コンピュータは、ユーザによって使用されるデバイスとの間でドキュメントを送受信することによって（例えば、ウェブブラウザから受信した要求に回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によって）、ユーザとインタラクションすることができる。

10

【 0 0 5 5 】

いくつかの実施形態が説明されている。それにもかかわらず、本開示の技術思想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、他の実施形態も以下の特許請求の範囲内にある。

20

30

40

50

【図面】

【図 1 A】

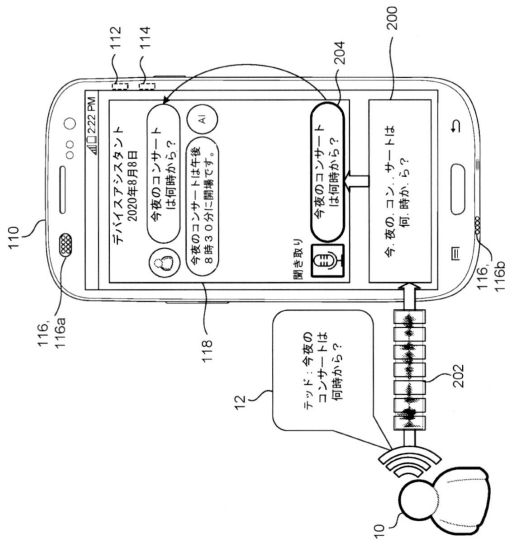


FIG. 1A

【図 1 B】

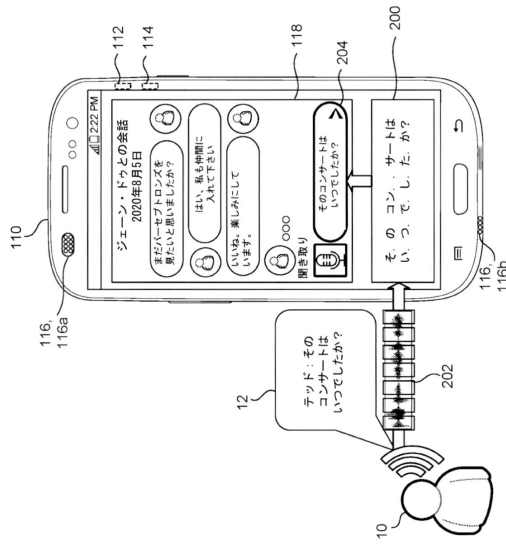


FIG. 1B

【図 2】

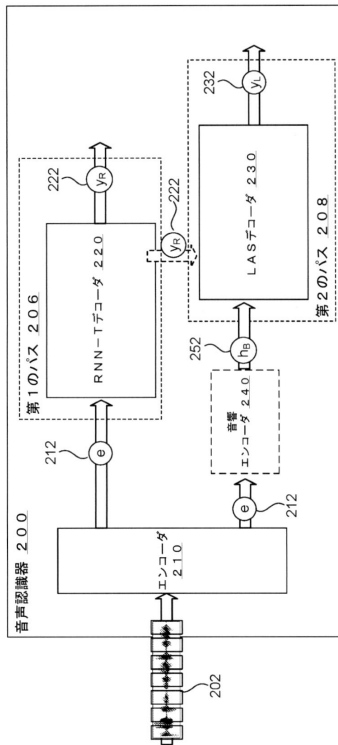


FIG. 2

【図 3 A】

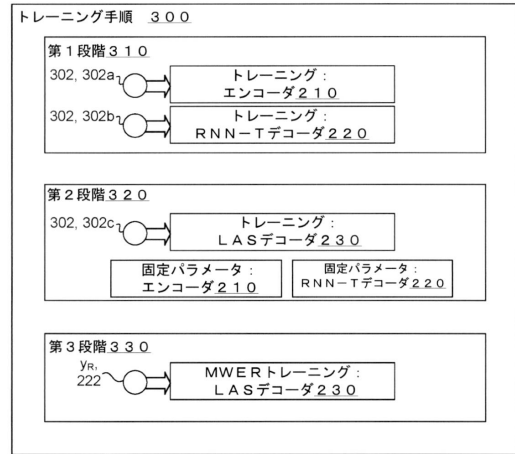


FIG. 3A

10

20

30

40

50

【図 3 B】

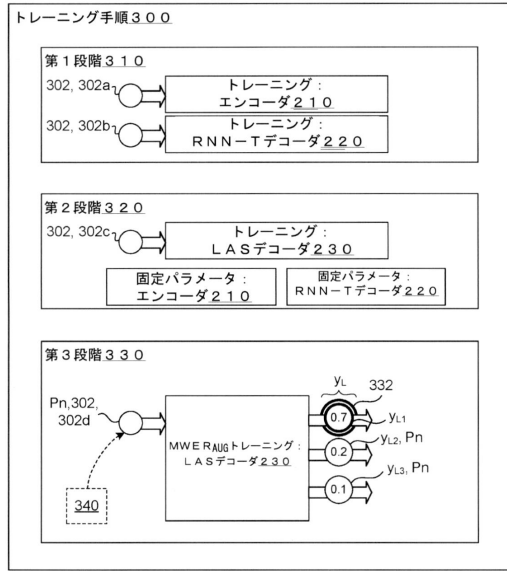


FIG. 3B

【図 3 C】

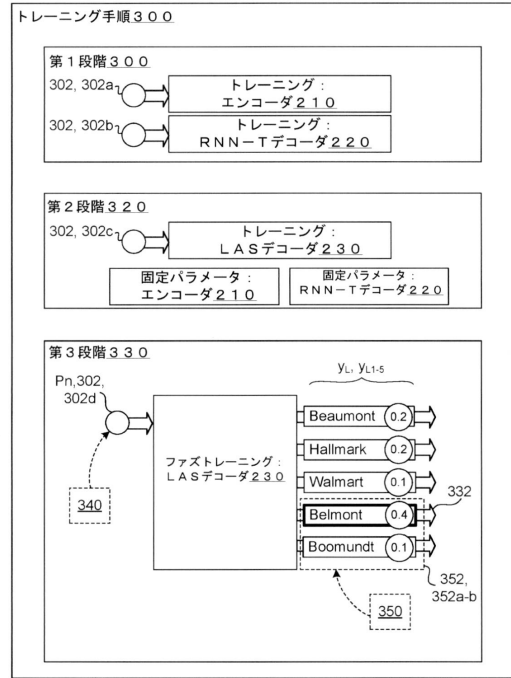


FIG. 3C

【図 4】

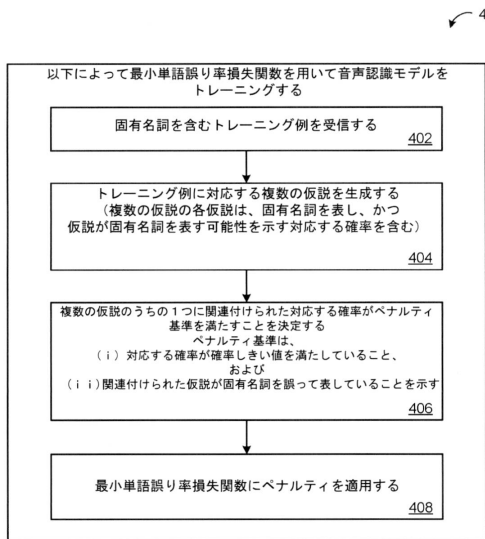


FIG. 4

【図 5】

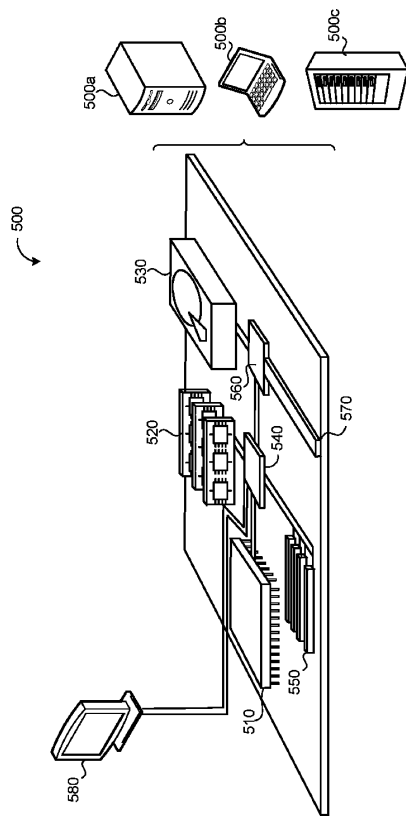


FIG. 5

10

20

30

40

50

フロントページの続き

e . i e e e . o r g / d o c u m e n t / 9 0 5 4 2 3 5 にて発表

早期審査対象出願

(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 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 特開 2 0 2 0 - 1 9 4 4 9 4 (J P , A)

国際公開第 2 0 1 2 / 1 6 5 5 2 9 (W O , A 1)

(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