

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7549061号
(P7549061)

(45)発行日 令和6年9月10日(2024.9.10)

(24)登録日 令和6年9月2日(2024.9.2)

(51)国際特許分類 F I
G 1 0 L 15/19 (2013.01) G 1 0 L 15/19
G 1 0 L 15/22 (2006.01) G 1 0 L 15/22 3 0 0 U

請求項の数 20 (全23頁)

(21)出願番号	特願2023-21323(P2023-21323)	(73)特許権者	502208397
(22)出願日	令和5年2月15日(2023.2.15)		グーグル エルエルシー
(62)分割の表示	特願2021-531511(P2021-531511)		Google LLC
)の分割		アメリカ合衆国 カリフォルニア州 9 4
原出願日	令和1年11月27日(2019.11.27)		0 4 3 マウンテン ビュー アンフィシ
(65)公開番号	特開2023-53331(P2023-53331A)		アター パークウェイ 1 6 0 0
(43)公開日	令和5年4月12日(2023.4.12)		1 6 0 0 Amphitheatre P
審査請求日	令和5年2月15日(2023.2.15)		arkway 9 4 0 4 3 Mounta
(31)優先権主張番号	62/774,507		in View, CA U.S.A.
(32)優先日	平成30年12月3日(2018.12.3)	(74)代理人	100142907
(33)優先権主張国・地域又は機関	米国(US)		弁理士 本田 淳
		(72)発明者	アレクシッチ、ペーター
			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】 音声入力処理

(57)【特許請求の範囲】

【請求項1】

コンピュータにより実装される方法であって、データ処理ハードウェアによって実行される

とき、ユーザに関連付けられているコンピューティングデバイスがデジタルアシスタントアプリケーションを実行していることを示す文脈を判定することと、

ユーザによって行われて前記コンピューティングデバイスによってキャプチャされた発話のオーディオデータを受信することと、

前記オーディオデータから、前記発話に対応する複数の音素を判定することと、

前記発話に対応する前記複数の音素を用いて、前記発話についての1つ以上の候補転写を生成することであって、前記1つ以上の候補転写の各々是对応する転写信頼スコアを有する、生成することと、

前記コンピューティングデバイスが前記デジタルアシスタントアプリケーションを実行していることを示す前記文脈に基づいて、複数の文法から文法を選択することと、

最高の前記転写信頼スコアを有する前記1つ以上の候補転写のうちの候補転写について、選択された前記文法を用いてパースを行って、前記コンピューティングデバイスが実行すべきアクションを識別することと、

を含む動作を前記データ処理ハードウェアに実行させ、

前記複数の文法は、互いに異なった、複数の用語からなる予め指定された構造を含む、方法。

【請求項 2】

前記動作は、識別された前記アクションを実行するよう前記コンピューティングデバイスに命令することをさらに含む、

請求項 1 に記載の方法。

【請求項 3】

前記複数の文法からアラームの設定に用いられるアラーム文法が選択される可能性は、前記コンピューティングデバイスがアラームアプリケーションを実行している場合、前記コンピューティングデバイスが前記デジタルアシスタントアプリケーションを実行している場合と比べて高い、

請求項 1 に記載の方法。

10

【請求項 4】

前記文脈が、前記コンピューティングデバイスが近くのデバイスと通信していることを含む、

請求項 1 に記載の方法。

【請求項 5】

前記コンピューティングデバイスは、携帯電話またはウェアラブルデバイスを含む、

請求項 1 に記載の方法。

【請求項 6】

前記文法は、デフォルト文法を含む、

請求項 1 に記載の方法。

20

【請求項 7】

前記 1 つ以上の候補転写の各々は複数の用語を含む、

請求項 1 に記載の方法。

【請求項 8】

前記複数の用語の各々は、対応する用語信頼スコアを含む、

請求項 7 に記載の方法。

【請求項 9】

前記転写信頼スコアは、前記候補転写が前記ユーザによって行われた前記発話と一致する可能性を示している、

請求項 1 に記載の方法。

30

【請求項 10】

前記データ処理ハードウェアは、前記コンピューティングデバイス上にある、

請求項 1 に記載の方法。

【請求項 11】

システムであって、

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と通信するメモリハードウェアと、を備え、前記メモリハードウェアは命令を記憶しており、前記命令は前記データ処理ハードウェア上において実行されるとき、

ユーザに関連付けられているコンピューティングデバイスがデジタルアシスタントアプリケーションを実行していることを示す文脈を判定することと、

ユーザによって行われて前記コンピューティングデバイスによってキャプチャされた発話のオーディオデータを受信することと、

前記オーディオデータから、前記発話に対応する複数の音素を判定することと、

前記発話に対応する前記複数の音素を用いて、前記発話についての 1 つ以上の候補転写を生成することであって、前記 1 つ以上の候補転写の各々は対応する転写信頼スコアを有する、生成することと、

前記コンピューティングデバイスが前記デジタルアシスタントアプリケーションを実行していることを示す前記文脈に基づいて複数の文法から文法を選択することと、

最高の前記転写信頼スコアを有する前記 1 つ以上の候補転写のうちの候補転写について

50

選択された前記文法を用いてパースを行って、前記コンピューティングデバイスが実行すべきアクションを識別することと、

を含む動作を前記データ処理ハードウェアに実行させ、

前記複数の文法は、互いに異なった、複数の用語からなる予め指定された構造を含む、システム。

【請求項 1 2】

前記動作は、識別された前記アクションを実行するよう前記コンピューティングデバイスに命令することをさらに含む、

請求項 1 1 に記載のシステム。

【請求項 1 3】

前記複数の文法からアラームの設定に用いられるアラーム文法が選択される可能性は、前記コンピューティングデバイスがアラームアプリケーションを実行している場合、前記コンピューティングデバイスが前記デジタルアシスタントアプリケーションを実行している場合と比べて高い、

請求項 1 1 に記載のシステム。

【請求項 1 4】

前記文脈が、前記コンピューティングデバイスが近くのデバイスと通信していることを含む、

請求項 1 1 に記載のシステム。

【請求項 1 5】

前記コンピューティングデバイスは、携帯電話またはウェアラブルデバイスを含む、

請求項 1 1 に記載のシステム。

【請求項 1 6】

前記文法は、デフォルト文法を含む、

請求項 1 1 に記載のシステム。

【請求項 1 7】

前記 1 つ以上の候補転写の各々は複数の用語を含む、

請求項 1 1 に記載のシステム。

【請求項 1 8】

前記複数の用語の各々は、対応する用語信頼スコアを含む、

請求項 1 7 に記載のシステム。

【請求項 1 9】

前記転写信頼スコアは、前記候補転写が前記ユーザによって行われた前記発話と一致する可能性を示している、

請求項 1 1 に記載のシステム。

【請求項 2 0】

前記データ処理ハードウェアは、前記コンピューティングデバイス上にある、

請求項 1 1 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本明細書は、一般に音声認識に関する。

【背景技術】

【0 0 0 2】

音声入力を使用してコンピュータとの対話を実行できるようにすることがますます求められている。これには、入力処理、特に自然言語データを処理および分析するようにコンピュータをプログラムする方法の開発が必要である。このような処理には、コンピュータによる話された言語の認識およびテキストへの変換を可能にする計算言語学の分野である音声認識が伴う場合がある。

【発明の概要】

10

20

30

40

50

【発明が解決しようとする課題】**【0003】**

本明細書は、一般に音声認識に関する。

【課題を解決するための手段】**【0004】**

ユーザが音声を通じてコンピューティングデバイスに入力を提供できるようにするため、音声入力処理システムは、文脈を使用して、自動音声認識部によって生成された複数の候補転写に適用する複数の文法を識別できる。各文法は、話者の異なる意図、あるいはシステムが同じ候補転写に対して実行するアクションを示している可能性がある。システムは、候補転写をパースする文法と、文法がユーザの意図と一致する可能性と、候補転写がユーザの発言と一致する可能性とに基づいて、文法および候補転写を選択し得る。次に、システムは、選択された候補転写に含まれる詳細を使用して、文法に対応するアクションを実行することができる。

10

【0005】

より詳細には、音声処理システムは、ユーザからの発話を受け取り、単語ラティスを生成する。単語ラティスは、発話における可能性の高い単語と、各単語の信頼スコアとを反映するデータ構造である。システムは、単語ラティスから、複数の候補転写と各候補転写の転写信頼スコアとを特定する。システムは、ユーザの特性、システムの位置、システムの特性、システム上で実行しているアプリケーション（例えば、現在アクティブなアプリケーションまたはフォアグラウンドで実行しているアプリケーション）、またはその他の類似の文脈データに基づいて、現在の文脈を特定する。文脈に基づいて、システムは、各候補転写をパースする複数の文法の文法信頼スコアを生成する。複数の文法が同じ候補転写に適用される可能性がある場合、システムは複数の文法信頼スコアのうちの一部を調整する場合がある。システムは、調整された複数の文法信頼スコアと複数の転写信頼スコアとの組み合わせに基づいて、文法と候補転写を選択する。

20

【0006】

本出願に記載される主題の革新的な態様によれば、音声入力処理方法は、コンピューティングデバイスによって、発話のオーディオデータを受信することと、コンピューティングデバイスによって、音響モデルおよび言語モデルを使用して、発話の複数の候補転写を含むとともに、複数の転写信頼スコアを含み、該複数の転写信頼スコアのそれぞれが、それぞれの候補転写が発話と一致する可能性を反映する単語ラティスを生成することと、コンピューティングデバイスによって、コンピューティングデバイスの文脈を判定することと、コンピューティングデバイスの文脈に基づいて、コンピューティングデバイスによって、複数の候補転写に対応する複数の文法を特定することと、現在の文脈に基づいて、コンピューティングデバイスによって、複数の候補転写のそれぞれについて、それぞれの文法がそれぞれの候補転写と一致する可能性を反映する複数の文法信頼スコアを判定することと、複数の転写信頼スコアおよび複数の文法信頼スコアに基づいて、コンピューティングデバイスによって、複数の候補転写の中から候補転写を選択することと、コンピューティングデバイスによる出力のために、選択された候補転写を発話の転写として提供することと、に係るアクションを含む。

30

40

【0007】

これらおよびその他の実装には、それぞれ任意選択により次のフィーチャの1つ以上を含めることができる。アクションは、複数の文法のうちの2つ以上が複数の候補転写のうちの1つに対応すると判定することと、複数の文法のうちの2つ以上が複数の候補転写のうちの1つに対応すると判定することに基づいて、2つ以上の文法について複数の文法信頼スコアを調整することと、を含む。コンピューティングデバイスは、複数の候補転写の中から、転写信頼スコアおよび調整された文法信頼スコアに基づいて、候補転写を選択する。2つ以上の文法について複数の文法信頼スコアを調整するアクションは、2つ以上の文法のそれぞれについて、複数の文法信頼スコアのそれぞれを係数で増加させることを含む。アクションは、複数の候補転写のそれぞれについて、それぞれの転写信頼スコアとそ

50

それぞれの文法信頼スコアとの積を判定することを含む。コンピューティングデバイスは、複数の候補転写の中から、転写信頼スコアとそれぞれの文法信頼スコアとの積に基づいて、候補転写を選択する。コンピューティングデバイスによって、コンピューティングデバイスの文脈を判定するアクションは、コンピューティングデバイスの位置と、コンピューティングデバイスのフォアグラウンドで実行しているアプリケーションと、時刻とに基づいている。言語モデルは、単語ラティスに含まれる語のシーケンスについて、確率を特定するように構成されている。音響モデルは、オーディオデータの一部に一致する音素を特定するように構成されている。アクションは、コンピューティングデバイスによって、選択された候補転写と、選択された候補転写に一致する文法とに基づくアクションを実行することを含む。

10

【0008】

この態様の他の実施形態は、対応するシステム、装置、およびコンピュータ記憶装置に記録されたコンピュータプログラムを含み、それぞれが方法の動作を実行するように構成されている。

【0009】

本明細書に記載されている主題の特定の実施形態は、以下の利点のうちの1つまたは複数を実現するように実施することができる。音声認識システムは、受信した音声入力と判定された文脈の両方を使用して、受信した音声入力をさらに処理してコンピューティングデバイスにアクションを実行させるために使用される文法を選択することができる。このように、音声認識システムは、限られた数の文法を複数の候補転写に適用することにより、ヒューマンマシンインタフェースの待ち時間を短縮することができる。音声認識システムは、システムが予期しない入力を受信したときに音声認識システムが転写を出力できるように、言語のすべてまたはほぼすべての単語を含む語彙を使用することができる。

20

【0010】

本明細書に記載主題の1つまたは複数の実施形態の詳細は、添付の図面および以下の説明に記載されている。主題の他の特徴、態様、および利点は、説明、図面、および特許請求の範囲から明らかとなる。

【図面の簡単な説明】

【0011】

【図1】文脈に基づいて発話に適用する文法を選択する例示的システムを示す図。

30

【図2】文脈に基づいて発話に適用する文法を選択する例示的システムの構成要素を示す図。

【図3】コンテキストに基づいて文法を選択する例示的プロセスを示すフローチャート。

【図4】コンピューティングデバイスおよびモバイルコンピューティングデバイスの例を示す図。

【発明を実施するための形態】

【0012】

さまざまな図面での同様の参照番号と指示は、同様の要素を示す。

図1は、文脈に基づいて発話に適用する文法を選択する例示的システム100を示す。以下で簡潔に、および、より詳細に説明するように、ユーザ102は発話104を行う。コンピューティングデバイス106は、発話104を検出し、発話に応じてアクション108を実行する。コンピューティングデバイス106は、コンピューティングデバイス106の文脈110を使用して、ユーザ102の可能性の高い意図を判定することができる。ユーザ102の可能性の高い意図に基づいて、コンピューティングデバイス106は、適切なアクションを選択する。

40

【0013】

より詳細には、ユーザ102は、コンピューティングデバイス106の近くで発話104を行う。コンピューティングデバイス106は、オーディオを検出するように構成された任意のタイプのコンピューティングデバイスであり得る。例えば、コンピューティングデバイス106は、携帯電話、ラップトップコンピュータ、ウェアラブルデバイス、スマ

50

ートアプライアンス、デスクトップコンピュータ、テレビ、またはオーディオデータを受信することができる任意の他のタイプのコンピューティングデバイスであってよい。コンピューティングデバイス106は、発話104のオーディオデータを処理し、単語ラティス112を生成する。単語ラティス112は、複数の単語からなる互いに異なる組み合わせを表し、それぞれについて、発話および信頼スコアに対応し得る。

【0014】

図1に示すように、ユーザ102は、明瞭に話しておらず、「フライツアウト (f l i g h t s o u t) 」のように聞こえる発話104を行う。コンピューティングデバイス106は、ユーザの音声のオーディオを検出し、自動音声認識を使用して、ユーザ102が何を言ったかを判定する。コンピューティングデバイス106は、音響モデルをオーディオデータに適用する。音響モデルは、オーディオデータに対応する可能性が高いさまざまな音素を判定するように構成され得る。例えば、音響モデルは、オーディオデータが、他の音素に加えて「f」音素の音素が含まれていると判定する場合がある。コンピューティングデバイス106は、言語モデルを音素に適用することができる。言語モデルは、音素に一致する、可能性の高い単語および一連の単語を特定するように構成されてよい。一部の実装形態では、コンピューティングデバイス106は、オーディオデータを別のコンピューティングデバイス、例えば、サーバに送信することができる。サーバは、オーディオデータに対して音声認識を実行する場合がある。

10

【0015】

言語モデルは、発話104に一致する複数の語からなる互いに異なる組み合わせを含む単語ラティス112を生成する。単語ラティスは、発話104の最初の単語であり得る2つの単語を含む。「ライツ (l i g h t s) 」という単語114が最初の単語である場合もあれば、「フライツ (f l i g h t s) 」という単語116が最初の単語である場合もある。言語モデルは、単語ラティス内の単語がユーザが話した単語である可能性を反映する信頼スコアを計算し得る。例えば、言語モデルは、単語114「ライツ (l i g h t s) 」の信頼スコア118が0.65であり、単語116「フライツ (f l i g h t s) 」の信頼スコア120が0.35であると判定できる。単語ラティス112は、2番目の単語としてあり得る単語を含む。この例では、言語モデルは、2番目の単語として1つのあり得る単語122「アウト (o u t) 」のみを特定した。

20

【0016】

一部の实装形態では、コンピューティングデバイス106は、音響モデルおよび言語モデルの代わりに、シーケンスツリーニューラルネットワーク、または他のタイプのニューラルネットワークを使用することができる。ニューラルネットワークは、1つまたは複数の隠れ層を有することができる。各サンプルの発話に対応するサンプルの発話および転写のオーディオデータを含む機械学習およびトレーニングデータを使用してトレーニングすることができる。この例では、シーケンスツリーニューラルネットワークは、単語ラティス112と信頼スコア118および120を生成することができる。シーケンスツリーニューラルネットワークは、音響モデルや言語モデルのように、音素および単語の組み合わせに対して個別の信頼スコアを生成しない場合がある。代わりに、シーケンスツリーニューラルネットワークは、音響モデルによって生成された音素の信頼スコアと、言語モデルによって生成された単語の組み合わせの信頼スコアと、の組み合わせである信頼スコアを生成する場合がある。

30

40

【0017】

単語ラティス112に基づいて、コンピューティングデバイス106は、発話104の2つの候補転写を特定する。第1の候補転写は「消灯 (l i g h t s o u t) 」であり、第2の候補転写は「フライツアウト (f l i g h t s o u t) 」である。第1の候補転写の信頼スコアは0.65であり、第2の候補転写の信頼スコアは0.35である。各候補転写の信頼スコアは、候補転写の各単語の信頼スコアの積であってもよい。フレーズ (p h a s e) 「ファイツアウト (f i g h t s o u t) 」は、発話104に最も近い音響一致であり得るが、言語モデルに基づいて、「ファイト (f i g h t) 」と「アウト

50

(out)」の組み合わせは生じる可能性が低い。したがって、「ファイツアウト (figh t s o u t)」は転写の候補とはならない。音響モデルと言語モデルの代わりにニューラルネットワークが使用される実装形態では、ニューラルネットワークは、音響モデルと言語モデルを使用して個別の信頼スコアを生成しない場合がある。

【0018】

コンピューティングデバイス106は、コンピューティングデバイス106の文脈110を判定することができる。文脈110は、コンピューティングデバイス106上またはその周囲に存在する要因の任意の組み合わせに基づいていてもよい。例えば、文脈110は、コンピューティングデバイス106がユーザ102の家に位置していることを含み得る。文脈110は、現在の時刻が午後10時であり、曜日が火曜日であることを含み得る。文脈110は、コンピューティングデバイス106が、コンピューティングデバイス106のフォアグラウンドでデジタルアシスタントアプリケーションを実行している携帯電話であることを含み得る。

10

【0019】

一部の实装形態では、文脈110は、追加情報を含み得る。例えば、文脈110は、テーブル上またはユーザの手の上で平らであるなど、コンピューティングデバイス106の配向に関するデータを含み得る。文脈110は、バックグラウンドで実行しているアプリケーションを含み得る。文脈は、発話104を受信する前にコンピューティングデバイス106によって出力されたオーディオまたはコンピューティングデバイスの画面に表示されたデータを含み得る。例えば、「こんにちは、いかがいたしましょう？」というプロンプトを含む表示によって、コンピューティングデバイス106がフォアグラウンドで仮想アシスタントアプリケーションを実行していることが示され得る。文脈110は、天気、ユーザ102の識別情報、ユーザ102の人口統計データ、および連絡先などのコンピューティングデバイス106に記憶されているかまたはコンピューティングデバイス106によってアクセス可能なデータを含み得る。

20

【0020】

コンピューティングデバイス106は、文脈を使用して、転写に適用する文法を選択することができる。文法は、バックス・ナウア記法 (Backus - Naur form) などの一般的な表記法を使用して記述できる複数の単語からなる任意の構造とすることができる。各文法は、特定のユーザの意図に対応している場合がある。例えば、ユーザの意図は、ホームオートメーションコマンドまたはメディア再生コマンドを発行することである可能性がある。文法の一例には、アラームの文法が含まれ得る。アラーム文法では、 $\$DIGIT = (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)$ という表記を使用して、数字を0、1、2、3、4、5、6、7、8、9、または0として定義できる。アラーム文法は、時刻に2桁、その次にコロン、その次に2桁、その次に「am」または「pm」とすることを示す $\$TIME = \$DIGIT \$DIGIT : \$DIGIT \$DIGIT (am | pm)$ という表記を使用して時刻を定義できる。アラーム文法は、アラームをアラームモードにするかタイマーモードにするかを示す $\$MODE = (alarm | timer)$ という表記を使用して、アラームのモードを定義できる。最後に、アラーム文法は、ユーザが「午前6時にアラームを設定 (set alarm for 6:00 am)」または「20:00にタイマーを設定 (set timer for 20:00)」とすることができることを示す、アラームシンタックスを $\$ALARM = set \$MODE for \$TIME$ として定義できる。コンピューティングデバイス106は、文法を使用して、話されたコマンドまたはタイプされたコマンドをパースし、コンピューティングデバイス106が実行すべきアクションを特定する。したがって、文法によって、コマンドをパースするためにコンピューティングデバイス106を特定の態様で動作させる機能データが提供される。

30

40

【0021】

各文法は、特定の文脈でアクティブになる場合がある。例えば、コンピューティングデバイス106のフォアグラウンドで実行しているアプリケーションがアラームアプリケー

50

ションである場合、アラーム文法はアクティブであり得る。さらには、コンピューティングデバイスのフォアグラウンドで実行しているアプリケーションがデジタルアシスタントアプリケーションである場合、アラーム文法がアクティブになる可能性がある。ユーザがアラームアプリケーションでアラームを設定しようとする可能性は、デジタルアシスタントアプリケーションの場合よりも高い可能性があるため、コンピューティングデバイス106は、コンピューティングデバイス106がアラームアプリケーションを実行している場合、コマンドおよびユーザの意図と一致するアラーム文法の可能性に0.95の確率を割り当ててもよく、コンピューティングデバイス106がデジタルアシスタントアプリケーションを実行している場合、0.03の確率を割り当ててもよい。

【0022】

図1に示す例では、コンピューティングデバイス106は、文法に適用され得る互いに異なる候補転写に基づいて、各文法の信頼スコアを判定する。換言すると、コンピューティングデバイス106は、候補転写のそれぞれをパースする文法を特定する。適用された文法124は、候補転写である「消灯 (lights out)」をパースする文法を示し、これにより、各文法がユーザの意図に対して正しい文法であるという文法信頼スコアが示されている。より高い文法信頼スコアは、その文法は所与の転写と文脈に基づいている可能性が高いことを示す。転写が「消灯 (lights out)」の場合、ホームオートメーション文法の文法信頼スコア126は0.35である。転写が「消灯 (lights out)」の場合、映画コマンド文法の文法信頼スコア128は0.25である。転写が「消灯 (lights out)」の場合、音楽コマンド文法の文法信頼スコア130は0.20である。各転写の文法信頼スコアの合計は1.0となるべきである。その場合、コンピューティングデバイス106は、残りの文法信頼スコアをデフォルト意図文法に割り当てることができる。この例では、転写が「消灯 (lights out)」の場合、デフォルト意図文法の文法信頼スコア132は0.20である。デフォルト意図文法は、特定のアクションまたは意図に限定されない場合があり、すべてまたはほぼすべての転写をパース可能である。

【0023】

一部の実装形態では、いずれの文法も特定の転写をパースできない場合がある。この場合、デフォルト意図文法が適用される。例えば、「フライトアウト (flight out)」をパースする文法がなくともよい。このため、転写が「フライトアウト (flight out)」の場合、デフォルト意図文法の文法信頼スコア132は1.0である。

【0024】

一部の実装形態では、コンピューティングデバイス106は、転写信頼スコアと文法信頼スコアとの積を計算することに基づいて、文法と転写を選択する。ボックス134は、転写信頼スコアと文法信頼スコアとの積に基づいて、転写および文法を選択した結果を示している。例えば、「消灯 (lights out)」とホームオートメーション文法とによる組み合わせ信頼スコア136は0.2275である。「消灯 (lights out)」と映画コマンド文法とによる組み合わせ信頼スコア138は0.1625である。「消灯 (lights out)」と音楽コマンド文法とによる組み合わせ信頼スコア140は0.130である。「消灯 (lights out)」とデフォルト意図文法とによる組み合わせ信頼スコア142は0.130である。「フライトアウト (flight out)」とデフォルト意図文法とによる組み合わせ信頼スコア144は0.350である。この実装形態では、組み合わせ信頼スコア144が最高スコアである。したがって、この手法を使用して、コンピューティングデバイス106は、コマンド「フライトアウト (flight out)」に対応するアクションを実行することができる。これは、ユーザ102にとっては妥当でない結果である可能性が高く、検索エンジンに転写「フライトアウト (flight out)」を提供する結果となる可能性がある。

【0025】

インターネットで「フライトアウト (flight out)」を検索することがユ

10

20

30

40

50

ーザ102の意図ではない場合、ユーザは発話104を繰り返す必要がある。そうすることによって、コンピューティングデバイス106が、追加の発話を処理するために追加のコンピューティングおよび電力リソースを消費することが必要となる。ユーザ102は、コンピューティングデバイス102の画面を追加の期間においてアクティブにすることによって、追加の処理および電力リソースを使用し得るコンピューティングデバイス102に手動で所望のコマンドを入力することになる。

【0026】

コンピューティングデバイス106がユーザの意図と一致しないアクションを実行する可能性を減らすために、コンピューティングデバイス106は、適用された文法124の文法信頼スコアを正規化することができる。適用された文法124の文法信頼スコアを正規化するために、コンピューティングデバイス106は、最高の文法信頼スコアを1.0に増加させるために必要な係数を計算することができる。換言すると、係数と最高の文法信頼スコアとの積は1.0となるべきである。次に、コンピューティングデバイス106は、他の文法信頼スコアを係数で乗算して、他の文法信頼スコアを正規化することができる。コンピューティングデバイス106によって実行される正規化は、確率の合計が1になるように調整される従来の確率正規化とは異なる場合がある。最高の文法信頼スコアを1.0に増加させるので、正規化された確率の合計は1にならない。増加させられた信頼スコアは、従来の意味での確率ではなく、疑似確率を表す場合がある。本願の他の場所で説明されている正規化処理でも、同様の疑似確率が生成される場合がある。

【0027】

正規化された文法信頼スコア146に示されているように、コンピューティングデバイス106は、ホームオートメーション文法の文法信頼スコアを最高文法信頼スコアとして特定し得る。文法信頼スコア126を1.0に増加させるために、コンピューティングデバイスは、文法信頼スコア126に $1.0 / 0.35 = 2.857$ を乗算する。コンピューティングデバイス106は、文法信頼スコア128に 2.857 を乗算することによって、正規化された文法信頼スコア150を計算する。コンピューティングデバイス106は、文法信頼スコア130に 2.857 を乗算することによって、正規化された文法信頼スコア152を計算する。コンピューティングデバイス106は、文法信頼スコア132に 2.857 を乗算することによって、正規化された文法信頼スコア154を計算する。転写が「フライツアウト (flights out)」の場合のデフォルト意図文法の文法信頼スコアは1.0であり、これは、転写「フライツアウト (flights out)」をパースする文法が他にないためである。したがって、転写が「フライツアウト (flights out)」である場合、スコアがすでに1.0であるので、デフォルト意図文法の文法信頼スコアを正規化する必要はない。

【0028】

ボックス134に示されているように、ボックス124の文法信頼スコアに単語ラティス112の転写信頼スコアを乗算する代わりに、コンピューティングデバイス106は、単語ラティス112からの転写信頼スコアと正規化された文法信頼スコアとを使用して、組み合わせられた信頼スコアを計算する。特に、コンピューティングデバイス106は、正規化された文法信頼スコアのそれぞれを、それぞれの転写の転写信頼スコアで乗算する。デフォルト意図文法のみが文法に適用される場合は、転写信頼スコアに1.0を乗算するので、対応する転写信頼スコアは変化しない。

【0029】

ボックス156に示すように、コンピューティングデバイス106は、「消灯 (lights out)」の転写信頼スコアに正規化された文法信頼スコア148を乗算して 0.650 の結果を得ることによって、組み合わせ信頼スコア158を計算する。コンピューティングデバイス106は、「消灯 (lights out)」の転写信頼スコアに正規化された文法信頼スコア150を乗算して 0.464 の結果を得ることによって、組み合わせ信頼スコア160を計算する。コンピューティングデバイス106は、「消灯 (lights out)」の転写信頼スコアに正規化された文法信頼スコア152を乗算し

10

20

30

40

50

て0.371の結果を得ることによって、組み合わせ信頼スコア162を計算する。コンピューティングデバイス106は、「消灯(light s o u t)」の転写信頼スコアに正規化された文法信頼スコア154を乗算して0.371の結果を得ることによって、組み合わせ信頼スコア164を計算する。コンピューティングデバイス106は、転写を「フライツアウト(f l i g h t s o u t)」をパースする文法ないので、「フライツアウト(f l i g h t s o u t)」の転写信頼スコアを0.350に維持する。

【0030】

一部の実装形態では、コンピューティングデバイス106は、コンピューティングデバイス106の現在の文脈の現在のユーザ文脈を考慮して、可能性が高い文法に対応できるよう、ボックス156内の信頼スコアを調整できる。例えば、ユーザ102は、コンピューティングデバイス106を使用して音楽を聴いている場合がある。コンピューティングデバイス106は、曲を再生するメディアデバイスであり得る。この場合、コンピューティングデバイス106は、信頼スコア162を調整し得る。コンピューティングデバイス106は、信頼スコアを係数で乗算することによって、プリセット値を信頼スコアに割り当てることによって、または別の技術によって、信頼スコアを増加させることができる。例えば、音楽コマンドの現在の確率は、メディアデバイスであるコンピューティングデバイス102を介して音楽を聴いているユーザ102の文脈に基づいて、0.9である。ボックス124内の音楽コマンドの確率は、0.20である信頼スコア130であり得る。この場合、コンピューティングデバイス106は、信頼スコア162に $0.9 / 0.2 = 4.5$ の比率を乗算してよい。結果として得られる信頼スコア162は、 $0.371 * 4.5 = 1.67$ となる。ボックス156の他の信頼スコアは、それぞれの信頼スコアについて同様の比率で調整することができる。例えば、ホームオートメーションコマンドの現在の確率は0.04である可能性がある。この場合、比率は $0.04 / 0.35 = 0.116$ となり、これは、0.04を信頼スコア126で割ったものである。コンピューティングデバイス106は、調整された、またはバイアスされた、0.10の信頼スコアを計算するために、信頼スコア158に0.16を乗算することができる。この場合、最高の信頼スコアは、音楽コマンドに対応するものとなる。

【0031】

一部の实装形態では、この追加の調整ステップは、コンピューティングデバイス106が選択する候補転写に影響を与える可能性がある。例えば、候補転写「フライツアウト(f l i g h t s o u t)」がビデオゲームの名前であり、ユーザがビデオゲームを起動することが期待される場合、コンピューティングデバイス106は、上記で計算されたものと同様の比率に基づいて、そしてユーザがビデオゲームを起動する確率を使用して、信頼スコア166を0.8と調整することができる。これは、コンピューティングデバイス106および/またはユーザ102の現在のコンテキストに基づいている。

【0032】

一部の实装形態では、複数の文法が同じ候補転写をパースする場合など、コンピューティングデバイス106は、この再スコアリングステップを使用せずに、ユーザの意図および音声認識の検出を改善することができる。この再スコアリングステップは、コンピューティングデバイス106が、音声認識信頼スコアに基づいて可能性高いとはされない、異なる候補転写を選択することを可能にし得る。

【0033】

一部の实装形態では、コンピューティングデバイス106は、ボックス124の信頼スコアを考慮せずに、この再スコアリングステップを使用することができる。例えば、コンピューティングデバイスは、単語ラティス112から特定された転写信頼スコアに再スコアリングステップを適用することができる。この場合、コンピューティングデバイス106は、ボックス146およびボックス156の両方に示されている調整を実行しないと考えられる。

【0034】

コンピューティングデバイス106は、最高の組み合わせ信頼スコアを有する文法およ

10

20

30

40

50

び転写を選択する。例えば、ホームオートメーション転写および転写「消灯 (l i g h t s o u t) 」は、0 . 6 5 0 の最高の組み合わせ信頼スコアを有し得る。この場合、コンピューティングデバイス 1 0 6 は、ホームオートメーション文法に基づいて「消灯 (l i g h t s o u t) 」を実行する。コンピューティングデバイス 1 0 6 は、コンピューティングデバイス 1 0 6 が配置されている家またはユーザ 1 0 2 の家の照明を消すことができる。コンピューティングデバイス 1 0 6 が映画コマンド文法を使用した場合、コンピューティングデバイス 1 0 6 は、映画「ライトアウト (L i g h t s O u t) 」を再生することができる。コンピューティングデバイス 1 0 6 が音楽コマンド文法を使用した場合、コンピューティングデバイス 1 0 6 は、曲「ライトアウト (L i g h t s O u t) 」を再生することができる。

10

【 0 0 3 5 】

図 1 に示すように、コンピューティングデバイス 1 0 6 のディスプレイは、コンピューティングデバイス 1 0 6 によって実行されるアクションを示し得る。最初に、コンピューティングデバイス 1 0 6 は、デジタルアシスタントのためのプロンプト 1 6 8 を表示することができる。コンピューティングデバイス 1 0 6 は、発話 1 0 4 を受信し、プロンプト 1 7 0 を表示することによって、コンピューティングデバイス 1 0 6 が照明をオフにしていることを示す。

【 0 0 3 6 】

図 2 は、文脈に基づいて発話に適用する文法を選択する例示的システム 2 0 0 の構成要素を示す。システム 2 0 0 は、音声オーディオを受信および処理するように構成された任意の種類コンピューティングデバイスであり得る。例えば、システム 2 0 0 は、図 1 のコンピューティングデバイス 1 0 6 と同様であり得る。システム 2 0 0 の構成要素は、単一のコンピューティングデバイスに実装することも、複数のコンピューティングデバイスに分散させることもできる。単一のコンピューティングデバイスに実装されているシステム 2 0 0 は、プライバシー上の理由から有益である可能性がある。

20

【 0 0 3 7 】

システム 2 0 0 は、オーディオサブシステム 2 0 5 を含む。オーディオサブシステム 2 0 5 は、マイクロフォン 2 1 0、アナログ - デジタル変換器 2 1 5、バッファ 2 2 0、および他の様々なオーディオフィルタを含み得る。マイクロフォン 2 1 0 は、音声などの周囲領域の音を検出するように構成することができる。アナログ - デジタル変換器 2 1 5 は、マイクロフォン 2 1 0 によって検出されたオーディオデータをサンプリングするように構成され得る。バッファ 2 2 0 は、システム 2 0 0 による処理のために、サンプリングされたオーディオデータを記憶することができる。一部の实装形態では、オーディオサブシステム 2 0 5 は、継続的にアクティブであり得る。この場合、マイクロフォン 2 1 0 は常に音を検出している可能性がある。アナログ - デジタル変換器 2 1 5 は、検出されたオーディオデータを常にサンプリングし得る。バッファ 2 2 0 は、音の最後の 1 0 秒などの最新のサンプリングされたオーディオデータを記憶することができる。システム 2 0 0 の他の構成要素がバッファ 2 2 0 におけるオーディオデータを処理しない場合、バッファ 2 2 0 は前のオーディオデータを上書きすることができる。

30

【 0 0 3 8 】

オーディオサブシステム 2 0 5 は、処理されたオーディオデータを音声認識部 2 2 5 に提供する。音声認識部は、音響モデル 2 3 0 への入力として音声データを提供する。音響モデル 2 3 0 は、オーディオデータにおける音に対応する可能性のある音素を特定するように訓練されてよい。例えば、ユーザが「セット (s e t) 」と言った場合、音響モデル 2 3 0 は、「s」音、「e」母音、および「t」音に対応する音素を特定し得る。音声認識部 2 2 5 は、特定された音素を言語モデル 2 3 5 への入力として提供する。言語モデル 2 3 5 は、単語ラティスを生成する。単語ラティスは、言語モデル 2 3 5 によって特定された候補用語のそれぞれについての用語信頼スコアを含む。例えば、単語ラティスは、最初の用語が「セット (s e t) 」の可能性が高いことを示し得る。言語モデル 2 3 5 は、最初の用語について他の可能性のある用語を特定しない場合がある。言語モデル 2 3 5 は

40

50

、第2の用語について2つのあり得る用語を特定し得る。例えば、単語ラティスは、あり得る第2の用語として「時間 (t i m e) 」および「チャイム (c h i m e) 」という用語を含み得る。言語モデル 2 3 5 は、各用語に用語信頼スコアを割り当てることができる。「時間」の用語信頼スコアは 0 . 7 5 であり、「チャイム」の用語信頼スコアは 0 . 2 5 である可能性がある。

【 0 0 3 9 】

音声認識部 2 2 5 は、単語ラティスに基づいて候補転写を生成することができる。各候補転写は、話者が転写における用語を話した可能性を反映する転写信頼スコアを有する可能性がある。例えば、候補転写は「時間をセット (s e t t i m e) 」であり、0 . 7 5 の転写信頼スコアを有する可能性がある。別の候補転写は「チャイムをセット (s e t c h i m e) 」であり、0 . 2 5 の転写信頼スコアを有する可能性がある。

10

【 0 0 4 0 】

音声認識部 2 2 5 が単語ラティス、候補転写、および転写信頼スコアを生成している間、文脈特定部 2 4 0 は、システム 2 0 0 の現在の文脈を示す文脈データを収集し得る。文脈特定部 2 4 0 は、システムの任意のセンサからセンサデータ 2 4 5 を収集することができる。センサは、位置センサ、温度計、加速度計、ジャイロスコープ、重力センサ、時間と曜日、および他の同様のセンサを含み得る。センサデータ 2 4 5 はまた、システム 2 0 0 のステータスに関連するデータを含み得る。例えば、ステータスは、システム 2 0 0 のバッテリーレベル、信号強度、システム 2 0 0 が通信している、または認識している可能性のある近くのデバイス、およびシステム 2 0 0 の他の同様の状態を含み得る。

20

【 0 0 4 1 】

文脈特定部 2 4 0 はまた、システム 2 0 0 が実行しているプロセスを示すシステムプロセスデータ 2 5 0 を収集することができる。プロセスデータ 2 5 0 は、各プロセスに割り当てられたメモリ、各プロセスに割り当てられたプロセスリソース、システム 2 0 0 が実行しているアプリケーション、システム 2 0 0 のフォアグラウンドまたはバックグラウンドで実行しているアプリケーション、システムのディスプレイ上のインターフェースの内容、および同様のプロセスデータを示し得る。

【 0 0 4 2 】

例として、文脈特定部 2 4 0 は、システム 2 0 0 がユーザの自宅であり、時刻が午後 6 時であり、曜日が月曜日であり、フォアグラウンドアプリケーションがデジタルアシスタントアプリケーション、デバイスがタブレットであることを示すシステムプロセスデータ 2 5 0 およびセンサデータ 2 4 5 を受信することができる。

30

【 0 0 4 3 】

文法スコア生成部 2 5 5 は、文脈特定部 2 4 0 からシステム 2 0 0 の文脈を受信し、音声認識部 2 2 5 から単語ラティスを受信する。文法スコア生成部 2 5 5 は、単語ラティスの候補転写のそれぞれをパースする文法 2 6 0 を特定する。場合によっては、文法 2 6 0 のいずれも候補転写をパースしない。この場合、文法スコア生成部 2 5 5 は、文法 2 6 0 のいずれによってもパース (p a r a b l e) できない候補転写にデフォルト意図文法の 1 . 0 の文法信頼スコアを与える。

【 0 0 4 4 】

一部の例では、単一の文法 2 6 0 が、候補転写をパースすることができる。この場合、文法スコア生成部 2 5 5 は、候補転写が実際の転写であると仮定して、単一の文法の文法信頼スコアを判定してもよい。文法信頼スコアは確率を表すので、文法信頼スコアは 1 . 0 未満である可能性がある。文法スコア生成部 2 5 5 は、1 . 0 と文法信頼スコアとの間の差を、候補転写のデフォルトの意図文法の文法信頼スコアに割り当てる。

40

【 0 0 4 5 】

一部の例では、複数の文法 2 6 0 が、候補転写をパースすることができる。この場合、文法スコア生成部 2 5 5 は、候補転写が実際の転写であると仮定して、複数の文法の各々の文法信頼スコアを判定してもよい。文法信頼スコアは確率の集合を表すので、文法信頼スコアの合計は 1 . 0 未満である可能性がある。文法スコア生成部 2 5 5 は、1 . 0 と文

50

法信頼スコアの合計との間の差を、候補転写のデフォルトの意図文法の文法信頼スコアに割り当てる。

【 0 0 4 6 】

文法スコア正規化部 2 6 5 は、候補転写のそれぞれについての文法信頼スコアを受信し、それらの文法信頼スコアを正規化する。一部の実装形態では、文法スコア正規化部 2 6 5 は、デフォルト意図文法以外の文法の文法信頼スコアのみを正規化する。文法スコア生成部 2 5 5 が、特定の候補転写について、1つの文法に対して1つの文法信頼スコアを生成する場合、文法スコア正規化部 2 6 5 は、当該文法信頼スコアを 1 . 0 に増やす。文法スコア生成部 2 5 5 が、特定の候補転写について、複数の文法に対して文法信頼スコアを生成しない場合、文法スコア正規化部 2 6 5 は、デフォルト意図文法の文法信頼スコアを 1 . 0 に維持する。

10

【 0 0 4 7 】

文法スコア生成部 2 5 5 が、特定の候補転写について、複数の文法の各々に対して複数の文法信頼スコアを生成する場合、文法スコア正規化部 2 6 5 は、特定の候補転写について最高の文法信頼スコアを特定する。文法スコア正規化部 2 6 5 は、係数と最高の文法信頼スコアとの積が 1 . 0 になるように、最高の文法信頼スコアを 1 . 0 に増加させるための係数を計算する。文法スコア正規化部 2 6 5 は、係数を他の文法信頼スコアのそれぞれに乗算することによって、同じ特定の転写についての他の文法信頼スコアを増加させる。

【 0 0 4 8 】

文法および転写選択部 2 7 0 は、正規化された文法信頼スコアおよび転写信頼スコアを受信する。文法信頼スコアと転写信頼スコアの両方を使用して、文法および転写選択部 2 7 0 は、組み合わせ信頼スコアを計算することによって、受信された発話および話者の意図にほぼ一致する可能性が高い文法および転写を特定する。文法および転写選択部 2 7 0 は、各正規化された文法信頼スコアおよび対応する候補転写の転写信頼スコアの積を計算することによって、それぞれの組み合わせ信頼スコアを判定する。デフォルト意図文法が特定の候補転写の唯一の文法である場合、文法および転写選択部 2 7 0 は、組み合わせ信頼スコアとして転写信頼スコアを維持する。文法および転写選択部 2 7 0 は、最高の組み合わせ信頼スコアを有する文法および候補転写を選択する。

20

【 0 0 4 9 】

アクション特定部 2 7 5 は、選択された文法および選択された転写を文法および転写選択部 2 7 0 から受信し、システム 2 0 0 が実行するアクションを特定する。選択された文法は、アラームの設定、メッセージの送信、曲の再生、人への電話、またはその他の同様のアクションなどのアクションのタイプを示している場合がある。選択された転写は、アラームを設定する時間、メッセージの送信先、再生する曲、呼び出し先、またはアクションのその他の同様の詳細など、アクションのタイプの詳細を示す場合がある。文法 2 6 0 は、特定の文法のアクションのタイプに関する情報を含み得る。文法 2 6 0 は、アクション特定部 2 7 5 が候補転写をどのようにパースしてアクションのタイプの詳細を判定すべきかについての情報を含み得る。例えば、文法は \$ アラーム (\$ A L A R M) 文法である可能性があり、アクション特定部 2 7 5 は選択された転写をパースして、タイマーを 2 0 分に設定することを判定する。アクション特定部 2 7 5 は、アクションを実行するか、またはシステム 2 0 0 の別の部分、例えば、プロセッサに命令を提供することができる。

30

40

【 0 0 5 0 】

一部の实装形態では、ユーザインタフェース生成部 2 8 0 は、アクションの標示またはアクションの実行の標示、あるいはその両方を表示する。例えば、ユーザインタフェース生成部 2 8 0 は、2 0 分からカウントダウンするタイマー、または家の中で照明をオフとすることの確認を表示することができる。一部の例では、ユーザインタフェース生成部 2 8 0 は、実行されたアクションの標示を提供しない場合がある。例えば、アクションはサーモスタットを調整している可能性がある。システム 2 0 0 は、システム 2 0 0 上に表示するためのユーザインタフェースを生成することなく、サーモスタットを調整してもよい。一部の实装形態では、ユーザインタフェース生成部 2 8 0 は、ユーザがアクションと対

50

話するか、またはアクションを確認するためのインタフェースを生成することができる。例えば、アクションは母親に電話することであってもよい。ユーザインタフェース生成部 280 は、システム 200 がアクションを実行する前に、母親を呼び出すアクションをユーザが確認するためのユーザインタフェースを生成することができる。

【0051】

図3は、文脈に基づいて文法を選択する例示的プロセス300を示すフローチャートである。一般に、プロセス300は、オーディオに対して音声認識を実行し、オーディオの転写およびオーディオをパースする文法に基づいて実行するアクションを特定する。プロセス300は、文法信頼スコアを正規化して、話者が意図した可能性が高いアクションを特定する。プロセス300は、1つまたは複数のコンピュータを含むコンピュータシステム、例えば、図1のコンピューティングデバイス106または図2のシステム200によって実行されるものとして説明される。

10

【0052】

システムは、発話のオーディオデータを受信する(310)。例えば、ユーザは、「消灯(lights out)」または「フライツアウト(flights out)」のように聞こえる発話を話すことがある。システムは、マイクを介して発話を検出するか、発話のオーディオデータを受信する。システムは、オーディオサブシステムを使用してオーディオデータを処理することができる。

【0053】

システムは、音響モデルおよび言語モデルを使用して、発話の複数の候補転写を含むとともに、複数の転写信頼スコアを含み、該複数の転写信頼スコアのそれぞれが、それぞれの候補転写が発話と一致する可能性を反映する単語ラティスを生成する(320)。システムは自動音声認識を使用して単語ラティスを生成する。自動音声認識プロセスは、オーディオデータの各部分に一致する様々な音素を特定する音響モデルへの入力としてオーディオデータを提供することを含み得る。自動音声認識プロセスは、発話中の各候補単語の信頼スコアを含む単語ラティスを生成する言語モデルへの入力として音素を提供することを含み得る。言語モデルは、語彙から単語ラティスの単語を選択する。語彙は、システムが認識するように構成されている言語の単語を含む場合がある。例えば、システムは英語のために構成されてよく、語彙には英語の単語を含んでいてよい。一部の実装形態では、プロセス300のアクションは、システムが認識できる語彙内の単語を制限することを含まない。換言すると、プロセス300は、システムが認識するように構成された言語の任意の単語を用いて転写を生成することができる。システムは言語における各単語を認識できるので、システムの音声認識プロセスは、ユーザが予期しないことを言ったとき、システムが認識するように構成されている言語内である限り機能できる。

20

30

【0054】

システムは、単語ラティスを使用して、様々な候補転写を生成できる。各候補転写は、ユーザが転写を話した可能性を反映する異なる転写信頼スコアを含む場合がある。例えば、候補転写「消灯(lights out)」の信頼スコアは0.65であり、候補転写「フライツアウト(flights out)」の信頼スコアは0.35である可能性がある。

40

【0055】

システムは、システムの文脈を判定する(330)。一部の实装形態では、文脈は、システムの種類、コンピューティングデバイスのフォアグラウンドで実行しているアプリケーション、ユーザの人口統計情報、連絡先データなど、システムに記憶されているかシステムからアクセスできるデータ、以前のユーザのクエリまたはコマンド、時刻、日付、曜日、天気、システムの向き、およびその他の同様のタイプの情報に基づいている。

【0056】

システムは、システムの文脈に基づいて、複数の候補転写に対応する複数の文法を特定する(340)。文法は、システムが実行できる様々なコマンドに対して様々な構造を含む場合がある。例えば、文法は、アラームの設定、映画の再生、インターネット検索の実

50

行、ユーザのカレンダーのチェック、またはその他の同様のアクションのためのコマンド構造を含む場合がある。

【 0 0 5 7 】

システムは、現在の文脈に基づいて、複数の候補転写のそれぞれについて、それぞれの文法がそれぞれの候補転写と一致する可能性を反映する複数の文法信頼スコアを判定する (3 5 0)。文法信頼スコアは、候補転写の1つが発話の転写であるとして、システムの文脈に基づいて文法が話者の意図と一致する可能性の条件付き確率と考えられ得る。例えば、転写が「消灯 (l i g h t s o u t)」の場合のホームオートメーション文法の文法信頼スコアは 0 . 3 5 である可能性がある。換言すると、転写が「消灯 (l i g h t s o u t)」であるとして、ホームオートメーション文法の条件付き確率は 0 . 3 5 である。システムは、候補転写をパースする文法ごとに文法信頼スコアを生成する場合がある。一部の実装形態では、文脈に起因して、システムは、候補転写をパースする一部の文法についてのみ、文法信頼スコアを生成する場合がある。各候補転写の文法信頼スコアの合計が 1 . 0 になるよう、システムは、残りの確率をデフォルト意図文法に割り当てる場合がある。

10

【 0 0 5 8 】

システムは、転写信頼スコアおよび文法信頼スコアに基づいて、複数の候補転写の中から1つの候補転写を選択する (3 6 0)。一部の实装形態では、システムは、複数の文法に一致する候補転写を増やすことにより、文法の信頼スコアを正規化する。候補転写ごとに、システムは、最高の文法信頼スコアに、文法信頼スコアを 1 . 0 に正規化するために必要な係数を乗算することができる。システムは、同じ候補転写に対する他の文法信頼スコアに対して、同じ係数で乗算することができる。システムは、正規化された文法信頼スコアに転写信頼スコアを掛け合わせて、組み合わせ信頼スコアを生成することができる。システムは、組み合わせられた信頼スコアが最も高い文法と転写を選択する。

20

【 0 0 5 9 】

システムは、出力のために、選択された候補転写を発話の転写として提供する (3 7 0)。一部の实装形態では、システムは、文法と候補転写に基づいてアクションを実行する。文法は、取るべき行動を示している場合がある。アクションには、映画の再生、連絡先への電話、照明の点灯、メッセージの送信、またはその他の同様のタイプのアクションが含まれる場合がある。選択された転写には、電話をかける連絡先、送信するメッセージ、メッセージの受信者、映画の名前、またはその他の同様の詳細が含まれる場合がある。

30

【 0 0 6 0 】

一部の实装形態では、システムは、有限状態トランスデューサー (F S T) を使用して、所与の文法セットの単語ラティスにタグを付けるプロセスを使用する。一部の实装形態では、システムは、現在の文脈に一致するすべての文法の和集合として文法を制約する場合がある。一部の实装形態では、システムは、事前にオフラインで、または実行時に動的に照合する場合がある。

【 0 0 6 1 】

一部の实装形態では、プロセス 3 0 0 は、単語ラティスのエッジに重みを付けることができる。システムは、文法を制約する重み付き有限状態トランスデューサへと文法をコンパイルすることができる。アークの重み、すなわち、この有限状態トランスデューサのエッジの重みは、所与の文法の単語の確率の量をエンコードする場合があり、これは、負の対数の重みである場合がある。一部の实装形態では、単語ラティスに関連するほぼすべての文法のすべてが、共に統合される場合がある。

40

【 0 0 6 2 】

プロセス 3 0 0 は、各文法の文脈依存確率を判定するシステムを続行することができ、これは、文脈を所与として、文法の確率であり得る。システムは、提供された構成要素と対応する重みによって文脈依存の確率を判定することができ、これは、オープニングデコーラアークなどの文法制約部において、アークまたはエッジにエンコードされてもよい。

【 0 0 6 3 】

50

プロセス300は、単語ラティスの重みをクリアするシステムを続行することができる。システムは、文法制約部を使用して単語ラティスを構成する場合があり、これにより、文法に一致するスパンがマークされる場合がある。例えば、文法に一致するスパンは、スパンについては<media_commands>曲を再生(play song)</media_commands>などのように開始および終了デコレータタグで囲まれてもよい(または「曲を再生(play song)」の転写)。

【0064】

プロセス300は、確率を正規化し続けてよい。システムは、デコレータタグが削除されたタグが付けられるラティスのコピーを生成する。システムは、トロピカルセミング(tropical semiring)のアーコストを判定および最小化し、その結果、それぞれの一意の単語経路が得られ、単語経路は、単語経路を通じた最も高い確率をエンコードする。確率は、文法を所与として単語経路の確率に文法の確率を乗算した確率であってもよい。確率は、負の対数の重みを用いてアーコストの重みの符号を反転することにより、反転させることができ、この反転したラティスは、タグ付きラティスで構成されてよい。反転格子で構成することは、その格子の重みで除算を実行することと実質的に同じである可能性がある(この場合、各単語シーケンスまたは経路ごとに、文法を所与として単語経路の確率の最大値に文法の確率を乗算したもの)。したがって、システムは最適なタグ付き経路の確率をそれ自体で分割し、それぞれが1.0になる可能性がある。一部の実装形態では、最適でない経路は、より低い疑似確率を受け取り、所与の単語経路の文法の疑似確率である所望の量を含むラティスを生成する。

【0065】

一部の实装形態では、プロセス300は、ビームプルーニング(beam pruning)を使用することによって、低い条件付き確率で文法を破棄することを含み得る。例えば、発話「消灯(lights out)」に一致する文法が多すぎる場合である(ここで、ホームオートメーションの確率は0.32、映画コマンド0.25、音楽コマンド0.20、一般検索0.10、オンラインビジネス0.03、社会集団0.03、およびオンライン百科事典検索0.02である)。システムは、閾値を下回る、または最も可能性の高い解釈すなわち文法の閾値パーセンテージを下回る文法を取り除くことによって、処理の負担を軽減することができる。例えば、システムは、最も可能性の高い文法の10分の1未満の可能性の文法を取り除いてよい。次に、システムは、所与の転写に対して0.032未満の文法を除去する場合がある。プルーニングの背後にある理論的根拠は、スコアの低いタグ付けの確率が非常に低いため、バイアシングがあっても、より可能性の高い解釈よりもその解釈が選択されるといった合理性がないことによる。プルーニングにおいては、システムが、所望のプルーニング重みで得られたラティスに動作する有限状態トランスデューサを追加する必要がある場合があり、これは、取り除くためにタグ付けが有する最良の仮定がどの遅れているかを指定する。

【0066】

プロセス300の利点によれば、システムが文法をラティス上で直接バイアスすることを可能となる。プロセス300は、他の単語ラティスタグ付けアルゴリズムよりも複雑さが低い。プロセス300は、最大正規化を実装することによって、タグ付けの曖昧さによって引き起こされる確率の断片化を解決する。プロセス300は、ビームプルーニングを使用することによって、比較的可能性の低い提案されたタグ付けを破棄することを可能にする。

【0067】

一部の实装形態では、プロセス300は、ラティスからn個の最良の仮説を抽出し、すべての仮説および可能性に個別にタグを付け、オプションでそれらを最後にラティスに再結合することができる。システムは、タグが他の点では同一であると想定してよい。これによって、レイテンシが長くなり、および/または、リコールが小さくなる同じ結果が得られる可能性がある。上位N(N=100など)の後のn個の最良の仮説が削除された場合、リコールは小さくなる可能性がある。これはレイテンシの管理に役立つ場合がある

10

20

30

40

50

。ラティス内の個別の文の数は、その中の単語の数に対して指数関数的になる可能性がある。前述の最良N制限（これにより、調べる選択肢の数が制限される場合がある）を使用しない限り、一部の回避策は最悪の場合に指数関数的に遅くなる可能性がある。

【0068】

図4は、本明細書に記載した技術を実装するために使用され得るコンピューティング装置400およびモバイルコンピューティング装置450の例を示す。コンピューティング装置400は、ラップトップ、デスクトップ、ワークステーション、携帯情報端末、サーバ、ブレードサーバ、メインフレーム、および他の適切なコンピュータ等、様々な形態のデジタルコンピュータを表すよう意図されている。モバイルコンピューティング装置450は、携帯情報端末、携帯電話、スマートフォン、および他の同様のコンピューティング装置等、様々な形態のモバイル装置を表すよう意図されている。本明細書に示された構成要素、構成要素の接続および関係、ならびに構成要素の機能は、例示的であることのみを意図されており、限定的であることは意図されていない。

10

【0069】

コンピューティング装置400は、プロセッサ402、メモリ404、記憶装置406、メモリ404および複数の高速拡張ポート410に接続する高速インタフェース408、ならびに低速拡張ポート414および記憶装置406に接続する低速インタフェース412を含む。プロセッサ402、メモリ404、記憶装置406、高速インタフェース408、高速拡張ポート410、および低速インタフェース412の各々は、様々なバスを介して相互接続され、共通のマザーボード上に、または必要に応じて他の方法で実装されてよい。プロセッサ402は、高速インタフェース408に接続されたディスプレイ416等の外部入力/出力装置にGUIのグラフィック情報を表示するために、メモリ404または記憶装置406に記憶された命令を含む、コンピューティング装置400内で実行するための命令を処理することが可能である。他の実装においては、複数のメモリおよび種類のメモリとともに、必要に応じて複数のプロセッサおよび/または複数のバスが使用されてよい。また、複数のコンピューティング装置が接続され、各装置が必要な動作の一部を提供してよい（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）。

20

【0070】

メモリ404は、コンピューティング装置400内の情報を記憶する。いくつかの実装においては、メモリ404は、単数または複数の揮発性メモリユニットである。いくつかの実装においては、メモリ404は、単数または複数の不揮発性メモリユニットである。メモリ404は、磁気ディスクまたは光学ディスク等の別の形態のコンピュータ可読媒体であってもよい。

30

【0071】

記憶装置406は、コンピューティング装置400に大容量記憶を提供することができる。いくつかの実装においては、記憶装置406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、光学ディスク装置、またはテープ装置、フラッシュメモリ、もしくは他の同様のソリッドステートメモリ装置等のコンピュータ可読媒体、またはストレージエリアネットワークもしくは他の構成における装置を含む装置のレイであるか、もしくはそれを含んでよい。命令は情報担体に記憶されてよい。命令は、1つまたは複数の処理装置（例えば、プロセッサ402）によって実行されると、上記のような1つまたは複数の方法を実行する。命令は、コンピュータ可読媒体または機械可読媒体（例えば、メモリ404、記憶装置406、またはプロセッサ402上のメモリ）等の1つまたは複数の記憶装置に記憶されてもよい。

40

【0072】

高速インタフェース408は、コンピューティング装置400の帯域幅集中型の動作を管理する一方、低速インタフェース412は、より低帯域幅集中型の動作を管理する。このような機能の割り当ては一例に過ぎない。いくつかの実装形態では、高速インタフェース408は、メモリ404、ディスプレイ416（例えば、グラフィックプロセッサまた

50

はアクセラレータを通じて)、および様々な拡張カード(図示せず)を受け得る高速拡張ポート410に接続される。実装において、低速インタフェース412は、記憶装置406および低速拡張ポート414に接続される。低速拡張ポート414は、種々の通信ポート(例えば、USB、Bluetooth(登録商標)、イーサネット(登録商標)、無線イーサネット)を含んでよく、例えばネットワークアダプタを通じて、キーボード、ポインティングデバイス、スキャナー、または、スイッチまたはルータ等のネットワーク装置等の、1つまたは複数の入力/出力装置に接続されてもよい。

【0073】

コンピューティング装置400は、図に示すように、多くの異なる形態で実装されてよい。例えば、標準のサーバ420として、またはそのようなサーバのグループに複数回実装されてよい。また、ラップトップコンピュータ422等のパーソナルコンピュータに実装されてよい。また、ラックサーバシステム424の一部として実装されてもよい。あるいは、コンピューティング装置400からの構成要素は、モバイルコンピューティング装置450等のモバイル装置(図示せず)内の他の構成要素と組み合わせられてよい。そのような装置のそれぞれは、コンピューティング装置400およびモバイルコンピューティング装置450のうちの1つまたは複数を含んでよく、システム全体は、互いに通信する複数のコンピューティング装置で構成されてよい。

10

【0074】

モバイルコンピューティング装置450は、構成要素の中でも特に、プロセッサ452、メモリ464、ディスプレイ454等の入力/出力装置、通信インタフェース46、送受信機468を含む。モバイルコンピューティング装置450は、追加のストレージを提供するために、マイクロドライブまたは他の装置等の記憶装置を備えてもよい。プロセッサ452、メモリ464、ディスプレイ454、通信インタフェース466、および送受信機468は、様々なバスを介して各々に相互接続され、複数の構成要素は、共通のマザーボード上に、または必要に応じて他の方法で実装されてよい。

20

【0075】

プロセッサ452は、メモリ464に記憶された命令を含む、モバイルコンピューティング装置450内の命令を実行することができる。プロセッサ452は、別個で複数のアナログおよびデジタルプロセッサを含むチップのチップセットとして実装されてよい。プロセッサ452は、例えば、ユーザインタフェース、モバイルコンピューティング装置450によって実行されるアプリケーション、およびモバイルコンピューティング装置450による無線通信の制御等、モバイルコンピューティング装置450の他の構成要素の調整を提供してよい。

30

【0076】

プロセッサ452は、ディスプレイ454に接続された制御インタフェース458およびディスプレイインタフェース456を通じてユーザと通信してよい。ディスプレイ454は、例えば、TFT(薄膜トランジスタ液晶ディスプレイ)ディスプレイ、OLED(有機発光ダイオード)ディスプレイ、または他の適切なディスプレイ技術であってよい。ディスプレイインタフェース456は、ディスプレイ454を駆動してグラフィックおよび他の情報をユーザに提示するための適切な回路を含んでよい。制御インタフェース458は、ユーザからコマンドを受け取り、プロセッサ452に供給するために、コマンドを変換してよい。さらに、外部インタフェース462は、プロセッサ452との通信を提供して、モバイルコンピューティング装置450と他の装置との近領域通信を可能にしてよい。外部インタフェース462は、例えば、いくつかの実装では有線通信を提供し、他の実装では無線通信を提供してよく、複数のインタフェースを使用してもよい。

40

【0077】

メモリ464は、コンピューティング装置450内の情報を記憶する。メモリ464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の1つまたは複数として実装することができる。拡張メモリ474が提供され、拡張インタフェース472を介してモバイルコンピューティング装置450に接続されてよい。拡張イン

50

タフェース 472 は、例えば、S I M M (シングルインラインメモリモジュール) カードインタフェースを含んでよい。拡張メモリ 474 は、モバイルコンピューティング装置 450 に余分な記憶領域を提供してよく、またはモバイルコンピューティング装置 450 のアプリケーションまたは他の情報を記憶してもよい。具体的には、拡張メモリ 474 は、上述のプロセスを実行または補完する命令を含んでよく、保安情報を含んでもよい。したがって、例えば、拡張メモリ 474 は、モバイルコンピューティング装置 450 のセキュリティモジュールとして提供されてよく、モバイルコンピューティング装置 450 の安全な使用を可能にする命令でプログラムされてよい。さらに、S I M M カードを介して、ハッキング不能な方法で S I M M カードに識別情報を配置する等の追加情報とともに、保安アプリケーションが提供されてよい。

10

【0078】

メモリは、以下で検討されるように、例えば、フラッシュメモリおよび/または N V R A M メモリ (不揮発性ランダムアクセスメモリ) を含んでよい。いくつかの実装においては、命令は、情報担体に記憶される。命令は、1 つまたは複数の処理装置 (例えば、プロセッサ 452) によって実行されると、上記のような 1 つまたは複数の方法を実行する。命令は、1 つまたは複数のコンピュータ可読媒体または 1 つまたは複数の機械可読媒体 (例えば、メモリ 464、拡張メモリ 474、またはプロセッサ 452 上のメモリ) 等の 1 つまたは複数の記憶装置に記憶されてもよい。いくつかの実装においては、命令は、例えば送受信機 468 または外部インタフェース 462 を通じて、伝搬信号で受信することができる。

20

【0079】

モバイルコンピューティング装置 450 は、必要に応じてデジタル信号処理回路を含み得る通信インタフェース 466 を通じて無線で通信してよい。通信インタフェース 466 は、特に、G S M (登録商標) 音声通話 (モバイル通信用グローバルシステム)、S M S (ショートメッセージサービス)、E M S (エンハンスドメッセージングサービス)、M M S メッセージング (マルチメディアメッセージングサービス)、C D M A (符号分割多元接続)、T D M A (時分割多元接続)、P D C (パーソナルデジタルセルラ)、W C D M A (登録商標) (広帯域符号分割多元接続)、C D M A 2 0 0 0、または G P R S (汎用パケット無線サービス) 等、様々な態様またはプロトコルの下で通信を提供してよい。このような通信は、例えば、無線周波数を使用する送受信機 468 を通じて行われてよい。

加えて、B l u e t o o t h、W i F i (登録商標)、または他のそのような送受信機 (図示せず) の使用等により、近距離通信が発生してよい。加えて、G P S (全地球測位システム) 受信機モジュール 470 は、モバイルコンピューティング装置 450 上で実行されるアプリケーションにより適宜使用され得る、追加的なナビゲーションおよび位置関連の無線データをモバイルコンピューティング装置 450 に提供してよい。

30

【0080】

モバイルコンピューティング装置 450 は、オーディオコーデック 460 を使用して可聴的に通信してよく、オーディオコーデック 460 は、ユーザから口頭の情報を受信し、それを使用可能なデジタル情報に変換してよい。オーディオコーデック 460 は、同様に、例えば、モバイルコンピューティング装置 450 のハンドセット内のスピーカ等を通じて、ユーザに可聴音を生成してよい。このような音は、音声通話の音を含んでよく、録音された音 (例えば、音声メッセージ、音楽ファイル等) を含んでよく、モバイルコンピューティング装置 450 上で動作するアプリケーションによって生成される音を含んでもよい。

40

【0081】

モバイルコンピューティング装置 450 は、図に示すように、多くの異なる形態で実装されてよい。例えば、モバイルコンピューティング装置 450 は、携帯電話 480 として実装されてよい。モバイルコンピューティング装置 450 は、スマートフォン 482、携帯情報端末、または他の同様のモバイル装置の一部として実装されてもよい。

【0082】

50

本明細書で説明するシステムおよび技法のさまざまな実装は、デジタル電子回路、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、記憶装置、1つ以上の入力装置、および1つ以上の出力装置に対してデータおよび命令を送信すると共にこれらからデータおよび命令を受信するよう接続された、特定目的または汎用目的の1つ以上のプログラマブルプロセッサを備えたプログラマブルシステム上で実行可能および/または翻訳可能な1つまたは複数のコンピュータプログラムでの実装を含んでよい。

【0083】

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のための機械語命令を含み、高水準手続き型および/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/機械語で実装することができる。本明細書で使用されるように、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、プログラマブルプロセッサに機械語命令及び/又はデータを供給するために使用される、機械可読信号として機械語命令を受け取る機械可読媒体を含む、任意のコンピュータプログラムプロダクト、装置及び/又はデバイス（例えば、磁気ディスク、光学ディスク、メモリ、PLD（プログラマブルロジックデバイス））を指す。機械可読信号という用語は、機械語命令および/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信号を指す。

【0084】

ユーザとの対話を提供するため、本明細書で説明したシステムおよび技術は、ユーザに情報を表示するための表示装置（例えば、CRT（陰極線管）またはLCD（液晶ディスプレイ）モニター）と、ユーザが入力をコンピュータに提供可能なキーボードおよびポインティング装置（例えば、マウスまたはトラックボール）と、を含むコンピュータ上で実装されてよい。他の種類の装置を使用して、ユーザとの対話を提供してもよい。例えば、ユーザに提供されるフィードバックは、任意の形式の感覚的なフィードバック（例えば、視覚的フィードバック、聴覚的フィードバック、または触覚的フィードバック）であってよく、ユーザからの入力、音響的入力、音声的入力、または触覚的入力を含む任意の形式で取り込まれてよい。

【0085】

本明細書で説明したシステムおよび技術は、バックエンドコンポーネント（例えば、データサーバ）を備えたコンピュータシステム、ミドルウェアコンポーネント（例えば、アプリケーションサーバ）を備えたコンピュータシステム、フロントエンドコンポーネント（例えば、ユーザが本明細書で説明されたシステムおよび技術の実装と対話することが可能なグラフィカルインタフェースまたはウェブブラウザを有するクライアントコンピュータ）を備えたコンピュータシステム、または、このようなバックエンド、ミドルウェア、またはフロントエンドコンポーネントの任意の組合せを備えたコンピュータシステムで実施されてよい。システムの構成要素は、デジタルデータ通信（例えば、通信ネットワーク）の任意の形式または媒体によって相互接続されてよい。通信ネットワークの例は、LAN（ローカルエリアネットワーク）、WAN（ワイドエリアネットワーク）、およびインターネットを含む。一部の实装形態では、本明細書において説明するシステムおよび技術は、音声認識および他の処理がデバイス上で直接実行される組み込みシステム上に実装することができる。

【0086】

コンピューティングシステムは、クライアントおよびサーバを含んでよい。クライアントとサーバは、一般には相互に離れており、典型的には通信ネットワークを通じて対話する。クライアントとサーバの関係は、各コンピュータ上で実行され、相互にクライアント・サーバ関係を有するコンピュータプログラムにより発生する。

【0087】

いくつかの実装が詳細に説明されたが、他の変更も可能である。例えば、クライアント

10

20

30

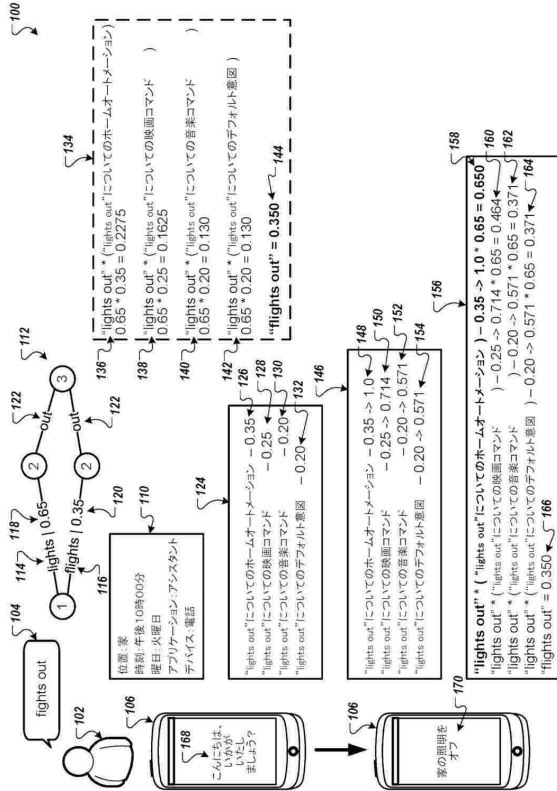
40

50

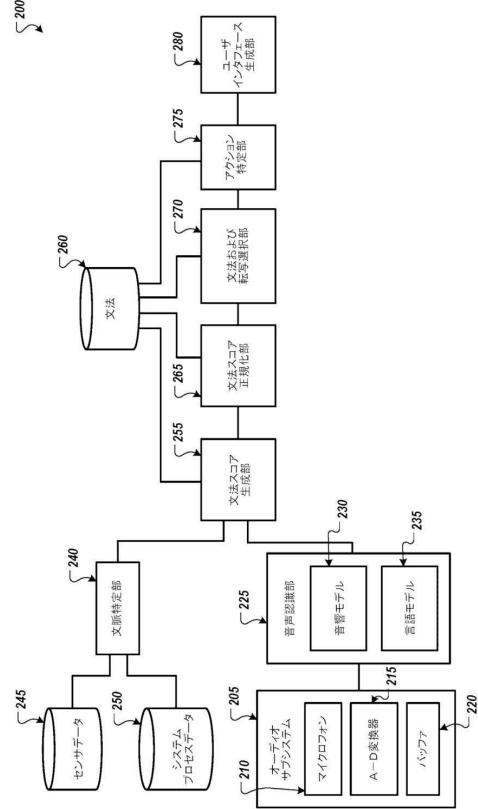
アプリケーションはデリゲートにアクセスするものとして説明されているが、他の実装では、デリゲートは、1つまたは複数のサーバで実行されるアプリケーション等、1つまたは複数のプロセッサによって実装された他のアプリケーションによって使用されてよい。さらに、図に示された論理の流れは、望ましい結果を得るために、示された特定の順序または連続した順序を必要とはしない。さらに、説明された流れに他の動作が提供されたり、または流れから除去されてよく、説明されたシステムに他の構成要素が追加されたり、またはシステムから除去されてよい。したがって、他の実装は、以下の特許請求の範囲内にある。

【図面】

【図 1】



【図 2】



10

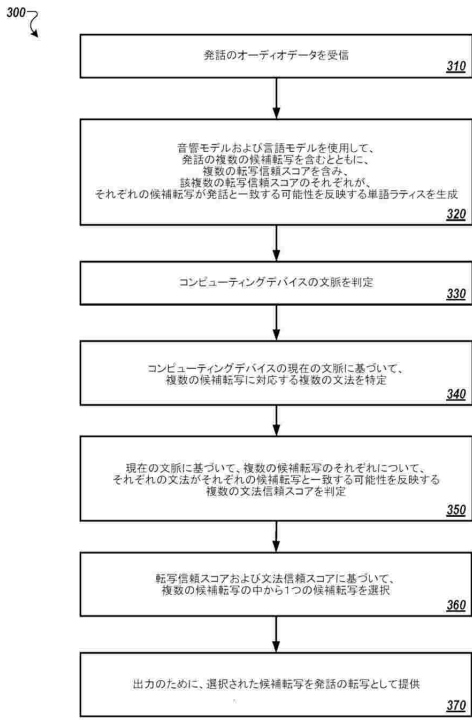
20

30

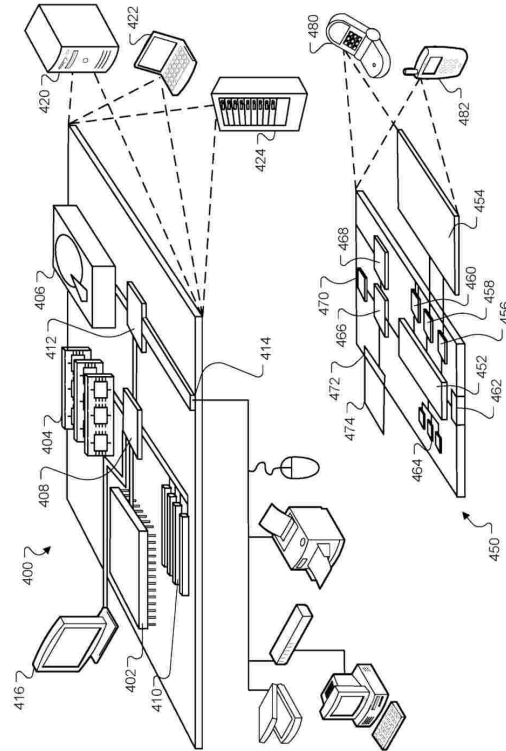
40

50

【 図 3 】



【 図 4 】



10

20

30

40

50

フロントページの続き

(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

(72)発明者 ベリコビッチ、レオニード
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 山下 剛史

(56)参考文献 特表 2 0 1 3 - 5 1 0 3 4 1 (J P , A)
特開 2 0 0 0 - 2 9 3 1 9 6 (J P , A)
特開 2 0 0 1 - 5 4 8 9 (J P , A)
米国特許出願公開第 2 0 1 8 / 0 0 5 3 5 0 2 (U S , A 1)
特開 2 0 0 1 - 1 5 7 1 3 7 (J P , A)
特開 2 0 0 3 - 9 1 2 9 8 (J P , A)
特開 2 0 1 8 - 1 8 2 6 9 2 (J P , A)
特表 2 0 1 8 - 5 3 3 0 3 6 (J P , A)

(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