

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7522177号
(P7522177)

(45)発行日 令和6年7月24日(2024.7.24)

(24)登録日 令和6年7月16日(2024.7.16)

(51)国際特許分類		F I			
G 0 6 F	3/01 (2006.01)	G 0 6 F	3/01	5 7 0	
G 0 6 F	3/16 (2006.01)	G 0 6 F	3/16	6 5 0	
		G 0 6 F	3/16	6 7 0	
		G 0 6 F	3/01	5 1 0	

請求項の数 15 (全36頁)

(21)出願番号	特願2022-500128(P2022-500128)	(73)特許権者	595020643
(86)(22)出願日	令和2年7月10日(2020.7.10)		クゥアルコム・インコーポレイテッド
(65)公表番号	特表2022-539794(P2022-539794 A)		QUALCOMM INCORPORATED
(43)公表日	令和4年9月13日(2022.9.13)		アメリカ合衆国、カリフォルニア州 9
(86)国際出願番号	PCT/US2020/041499		2 1 2 1 - 1 7 1 4、サン・ディエゴ、
(87)国際公開番号	WO2021/011331		モアハウス・ドライブ 5 7 7 5
(87)国際公開日	令和3年1月21日(2021.1.21)	(74)代理人	110003708
審査請求日	令和5年6月12日(2023.6.12)		弁理士法人鈴榮特許総合事務所
(31)優先権主張番号	62/873,775	(74)代理人	100108855
(32)優先日	令和1年7月12日(2019.7.12)		弁理士 蔵田 昌俊
(33)優先権主張国・地域又は機関	米国(US)	(74)代理人	100158805
(31)優先権主張番号	16/685,946		弁理士 井関 守三
(32)優先日	令和1年11月15日(2019.11.15)	(74)代理人	100112807
	最終頁に続く		弁理士 岡田 貴志
			最終頁に続く

(54)【発明の名称】 マルチモーダルユーザインターフェース

(57)【特許請求の範囲】

【請求項 1】

マルチモーダルユーザ入力のためのデバイスであって、

第 1 の入力デバイスから受信された第 1 のデータを処理することと、前記第 1 のデータは、第 1 の入力モードに基づくユーザからの第 1 の入力を示し、前記第 1 の入力は、コマンドに対応し、

前記第 1 のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることと、ここにおいて、前記フィードバックメッセージは、前記第 1 の入力モードとは異なる第 2 の入力モードに基づいて、前記第 1 の入力に関連するコマンドを識別する第 2 の入力を提供するように前記ユーザに命令する、

第 2 の入力デバイスから第 2 のデータを受信することと、前記第 2 のデータは、前記第 2 の入力を示し、

前記第 1 の入力のマッピングを、前記第 2 の入力に関連付けられた行為に対して更新することと、

を行うように構成された 1 つまたは複数のプロセッサを備える、デバイス。

【請求項 2】

前記第 1 の入力モードは、スピーチモード、ジェスチャーモード、またはビデオモードのうちの一つであり、前記第 2 の入力モードは、前記スピーチモード、前記ジェスチャーモード、または前記ビデオモードのうちの一つである、請求項 1 に記載のデバイス。

【請求項 3】

前記フィードバックメッセージは、前記第 1 の入力をディスアンビギュエートするために前記第 2 の入力を提供するように前記ユーザに命令し、

オプションで、前記 1 つまたは複数のプロセッサは、前記第 1 の入力の認識処理に関連する確信度レベルが確信度しきい値を満たすことに失敗したことに応答して、前記フィードバックメッセージを送るようさらに構成された、請求項 1 に記載のデバイス。

【請求項 4】

前記更新されたマッピングは、前記第 1 の入力と前記第 2 の入力との組合せを前記コマンドに関連付ける、請求項 1 に記載のデバイス。

【請求項 5】

前記 1 つまたは複数のプロセッサは、マルチモーダル認識エンジンを含み、前記マルチモーダル認識エンジンは、

組み合わせられた埋め込みベクトルを生成するために、前記第 1 の入力モードに関連する第 1 の埋め込みネットワークと、前記第 2 の入力モードに関連する第 2 の埋め込みネットワークとの出力を組み合わせるように構成された融合埋め込みネットワークと、

前記組み合わせられた埋め込みベクトルを特定のコマンドにマッピングするように構成された分類器と、

を含む、請求項 1 に記載のデバイス。

【請求項 6】

前記ユーザに対応する第 1 の埋め込みネットワークデータおよび第 1 の重みデータと、第 2 のユーザに対応する第 2 の埋め込みネットワークデータおよび第 2 の重みデータと、前記第 1 の埋め込みネットワークデータは、前記ユーザと前記第 2 のユーザとの間の入力コマンドの差に基づいて前記第 2 の埋め込みネットワークデータとは異なり、前記第 1 の重みデータは、前記ユーザと前記第 2 のユーザとの間の入力モード信頼性の差に基づいて前記第 2 の重みデータとは異なり、

を記憶するように構成されたメモリをさらに備える、請求項 5 に記載のデバイス。

【請求項 7】

前記第 1 の入力モードは、ビデオモードに対応し、前記 1 つまたは複数のプロセッサは、照明しきい値を下回る値を有する周辺光メトリックに応答して前記フィードバックメッセージを送るよう構成された、請求項 1 に記載のデバイス。

【請求項 8】

前記第 1 の入力モードは、スピーチモードに対応し、前記 1 つまたは複数のプロセッサは、雑音しきい値を超える値を有する雑音メトリックに応答して前記フィードバックメッセージを送るよう構成された、請求項 1 に記載のデバイス。

【請求項 9】

グラフィカルユーザインターフェースを表すように構成されたディスプレイをさらに備える、請求項 1 に記載のデバイス。

【請求項 10】

1 つまたは複数のキーワードまたは音声コマンドを含むオーディオ入力をキャプチャするように構成された 1 つまたは複数のマイクロフォンをさらに備える、請求項 1 に記載のデバイス。

【請求項 11】

1 つまたは複数のジェスチャーまたは視覚的コマンドを含むビデオ入力をキャプチャするように構成された 1 つまたは複数のカメラをさらに備える、請求項 1 に記載のデバイス。

【請求項 12】

ジェスチャー入力を示すデータを受信するように構成された 1 つまたは複数のアンテナをさらに備える、請求項 1 に記載のデバイス。

【請求項 13】

前記フィードバックメッセージをレンダリングするかまたは前記ユーザにダイレクトするように構成された 1 つまたは複数のラウドスピーカーをさらに備える、請求項 1 に記載のデバイス。

10

20

30

40

50

【請求項 1 4】

マルチモーダルユーザ入力のための方法であって、

デバイスの1つまたは複数のプロセッサにおいて、第1の入力デバイスから受信された第1のデータを処理することと、前記第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、前記第1の入力は、コマンドに対応し、

前記1つまたは複数のプロセッサから、前記第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることと、ここにおいて、前記フィードバックメッセージは、前記第1の入力モードとは異なる第2の入力モードに基づいて、前記第1の入力に関連するコマンドを識別する第2の入力を提供するように前記ユーザに命令する、

10

前記1つまたは複数のプロセッサにおいて、第2の入力デバイスから第2のデータを受信することと、前記第2のデータは、前記第2の入力を示し、

前記第1の入力のマッピングを、前記第2の入力に関連付けられた行為に対して更新することと、

を備える、方法。

【請求項 1 5】

デバイスの1つまたは複数のプロセッサによって実行されたとき、前記1つまたは複数のプロセッサに、

第1の入力デバイスから受信された第1のデータを処理することと、前記第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、前記第1の入力が、コマンドに対応し、

20

前記第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることと、ここにおいて、前記フィードバックメッセージは、前記第1の入力モードとは異なる第2の入力モードに基づいて、前記第1の入力に関連するコマンドを識別する第2の入力を提供するように前記ユーザに命令する、

第2の入力デバイスから第2のデータを受信することと、前記第2のデータは、前記第2の入力を示し、

前記第1の入力のマッピングを、前記第2の入力に関連付けられた前記コマンドに対して更新することと、

を行わせる命令を備える非一時的コンピュータ可読媒体。

30

【発明の詳細な説明】

【優先権の主張】

【0001】

[0001]本出願は、それらの各々の内容がそれらの全体として参照により本明細書に明確に組み込まれる、本願の譲受人が所有する2019年7月12日に出願された米国仮特許出願第62/873,775号と、2019年11月15日に出願された米国非仮特許出願第16/685,946号との優先権の利益を主張する。

【技術分野】

【0002】

[0002]本開示は、一般にユーザインターフェースに関係し、より詳細には、ユーザ入力の複数のモダリティ(modalities)をサポートするユーザインターフェースに関係する。

40

【背景技術】

【0003】

[0003]多くのユーザインターフェースは、自動音声認識(ASR)および自然言語処理(NLP)に基づき、大規模カスタマーベース上で有用であるように多くの異なるコマンド、アクセント、および言語にわたってトレーニングされる。様々なユーザの間の広い適用可能性のためにそのようなユーザインターフェースをトレーニングすることは、広範なリソースを必要とし、ユーザインターフェースを大規模カスタマーベースのために一般的に適用可能にするためのトレーニングの大部分は、各個々のユーザが典型的には単一の言語、アクセント、およびサポートされるコマンドのサブセットのみを使用するので、ユー

50

ザごとのベースでは「浪費」になる。

【発明の概要】

【0004】

[0004]本開示の一実装形態によれば、マルチモーダルユーザ入力のためのデバイスが、第1の入力デバイスから受信された第1のデータを処理するように構成された1つまたは複数のプロセッサを含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマンドに対応する。1つまたは複数のプロセッサは、第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送るように構成される。フィードバックメッセージは、第1の入力モードとは異なる第2の入力モードに基づいて、第1の入力に関連するコマンドを識別する第2の入力を提供するよう

10

【0005】

[0005]本開示の別の実装形態によれば、マルチモーダルユーザ入力のための方法が、デバイスの1つまたは複数のプロセッサにおいて、第1の入力デバイスから受信された第1のデータを処理することを含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマンドに対応する。本方法は、第1のデータを処理することに基づいて1つまたは複数のプロセッサから出力デバイスにフィードバックメ

20

【0006】

[0006]本開示の別の実装形態によれば、マルチモーダルユーザ入力のための装置が、第1の入力デバイスから受信された第1のデータを処理するための手段を含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマ

30

【0007】

[0007]本開示の別の実装形態によれば、非一時的コンピュータ可読媒体が、デバイスの1つまたは複数のプロセッサによって実行されたとき、1つまたは複数のプロセッサに、第1の入力デバイスから受信された第1のデータを処理させる命令を含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマ

40

50

マンドに第 1 の入力に関連付けるようにマッピングを更新させる。

【図面の簡単な説明】

【0008】

【図 1】[0008]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能なデバイスを含むシステムの特定の例示的な実装形態の図。

【図 2】[0009]本開示のいくつかの例による、図 1 のデバイスの構成要素の特定の实装形態の図。

【図 3】[0010]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能なデバイスを含むシステムの別の特定の实装形態の図。

【図 4】[0011]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能なデバイスを含むシステムの別の特定の实装形態の一例の図。

10

【図 5】[0012]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能なデバイスの別の実装形態の図。

【図 6】[0013]本開示のいくつかの例による、図 1 のデバイスによって実施され得るマルチモーダルユーザ入力を処理する方法の実装形態の図。

【図 7】[0014]本開示のいくつかの例による、図 1 のデバイスによって実施され得るマルチモーダルユーザ入力を処理する方法の別の実装形態の図。

【図 8】[0015]本開示のいくつかの例による、図 1 のデバイスによって実施され得るマルチモーダルユーザ入力を処理する方法の別の実装形態の図。

【図 9】[0016]本開示のいくつかの例による、図 1 のデバイスによって実施され得るマルチモーダルユーザ入力を処理する方法の別の実装形態の図。

20

【図 10】[0017]本開示のいくつかの例による、図 1 のデバイスによって実施され得るマルチモーダルユーザ入力を処理する方法の別の実装形態の図。

【図 11】[0018]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能な車両の図。

【図 12A】[0019]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能な仮想現実または拡張現実ヘッドセットの図。

【図 12B】[0020]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能なウェアラブル電子デバイスの図。

【図 13】[0021]本開示のいくつかの例による、マルチモーダルユーザ入力を処理するように動作可能であるデバイスの特定の例示的な例のブロック図。

30

【発明を実施するための形態】

【0009】

[0022]複数の入力モダリティを使用したユーザ対話を可能にするためのデバイスおよび方法について説明される。多くのユーザインターフェースは、自動音声認識 (ASR) および自然言語処理 (NLP) に基づき、大規模カスタマーベース上で有用であるように多くの異なるコマンド、アクセント、および言語にわたってトレーニングされる。様々なユーザの間の広い適用可能性のためにそのようなユーザインターフェースをトレーニングすることは、広範なリソースを必要とし、ユーザインターフェースを大規模カスタマーベースのために一般的に適用可能にするためのトレーニングの大部分は、各個々のユーザが典型的には単一の言語、アクセント、およびサポートされるコマンドのサブセットのみを使用するので、ユーザごとのベースでは「浪費」になる。

40

【0010】

[0023]ユーザコマンドの解釈を個人化する能力とともに、マルチモーダルユーザ対話を可能にすることによって、本明細書で説明される技法は、マルチモーダルユーザインターフェースが、特定のユーザによる使用のためにトレーニングされることを可能にし、それにより、従来のユーザインターフェースの広い適用可能性のための広範なトレーニングを低減するかまたはなくす。いくつかの実装形態では、異なる埋め込 (embedding) ネットワークは、異なる入力モダリティのために使用され (たとえば、スピーチ用の埋め込みネットワーク、視覚的入力用の埋め込みネットワーク、ジェスチャー入力用の埋め込みネッ

50

トワークなど)、それぞれのモダリティを使用して受信される異なるコマンド間で区別するように構成される。例示のために、「埋め込みネットワーク」は、埋め込みベクトルを生成するために、スピーチデータ(たとえば、時間領域スピーチデータまたは周波数領域スピーチデータ)などの入力データを処理するように構成された(たとえば、トレーニングされた)、1つまたは複数のニューラルネットワークレイヤを含むことができる。「埋め込みベクトル」は、入力データと比較して比較的次元であり、入力データを表し、入力データの異なるインスタンス間で区別するために使用され得る、ベクトル(たとえば、複数の値のセット)である。異なる埋め込みネットワーク出力は、共通の埋め込み空間に変換され、組み合わせられた埋め込みベクトルに融合される。たとえば、スピーチ入力の n 次元のスピーチ埋め込みベクトルは、 k 次元の第1の埋め込みベクトルに変換され得、ジェスチャー入力の m 次元のジェスチャー埋め込みベクトルは、 k 次元の第2の埋め込みベクトルに変換され得る(ここで、 m 、 n 、および k は、互いに等しいか、または異なり得る)。 k 次元のベクトル空間(たとえば、共通の埋め込み空間)中で、 k 次元の第1の埋め込みベクトルは、スピーチ入力を表し、 k 次元の第2の埋め込みベクトルは、ジェスチャー入力を表す。 k 次元の第1の埋め込みベクトルと、 k 次元の第2の埋め込みベクトルとは、組み合わせられた埋め込みベクトルを生成するために、ベクトル加算などによって組み合わせられ得る。分類器は、出力を生成するために、組み合わせられた埋め込みベクトルを解釈する。

10

【0011】

[0024]埋め込みネットワークと分類器との各々は、様々なモダリティを介して受信されたユーザコマンドの認識を改善するように個々のユーザによって更新(たとえば、トレーニング)され得る。たとえば、高い確信度で解釈され得ない、話されたユーザコマンドが受信された場合、ユーザインターフェースは、話されたコマンドの意味に関してユーザに問い合わせることができ、ユーザは、ユーザインターフェースによって認識されるジェスチャー入力を実施することなどによって、異なるモダリティを使用して意味を入力することができる。

20

【0012】

[0025]いくつかの実装形態では、ユーザインターフェースは、ユーザが入力モダリティを変更することを要求することができる。たとえば、再生ボリュームを上げるためのユーザの話されたコマンド「アップ」が、別のコマンド(たとえば、「オフ」)から確実に区別され得ない場合、ユーザインターフェースは、コマンドをより良く区別するためにユーザが別のモダリティを追加することを要求する(たとえば、話されたまた表示された)フィードバックメッセージを生成することができる。たとえば、ユーザは、「ボリュームを上げる」コマンドのために上方にポインティングすることなど、視覚的入力を追加することができる。ユーザインターフェースは、再生ボリュームを上げるためのマルチモーダルコマンドとして、話された入力「アップ」と、上方ポインティングの視覚的入力との組合せを認識するように更新され得る。したがって、コマンド認識精度を改善するために、(たとえば、シングルモーダルからマルチモーダルへの)個人化された更新が使用され得る。

30

【0013】

[0026]いくつかの実装形態では、ユーザインターフェースは、ユーザ入力をより容易にディスアンビグエートする(disambiguate)ために、ユーザが入力モダリティを変更することを要求する。たとえば、(たとえば、動いている車両中で)オーディオ雑音ユーザのスピーチの解釈を損なう実装形態では、ユーザインターフェースは、ユーザがモダリティを視覚的またはジェスチャーモダリティなどに変更することを要求するフィードバックメッセージを生成することができる。別の例として、低い光レベルがユーザの視覚的入力の解釈を損なう実装形態では、ユーザインターフェースは、ユーザがモダリティを、スピーチモダリティ、または手の移動および配向を検出するためにウェアラブル電子デバイス(たとえば、「スマートウォッチ」)の動き検出器を使用するジェスチャーモダリティなどに変更することを要求する、フィードバックメッセージを生成することができる。したがって、入力モダリティを変更するようにユーザに命令することは、コマンド認識精

40

50

度を改善するために使用され得る。

【0014】

[0027]いくつかの実装形態では、ユーザインターフェースは、多因子認証プロセスの一部としてユーザが入力モダリティを変更することを要求する。たとえば、音声認証を実施するために、話されたユーザ入力を受信した後に、ユーザインターフェースは、次に、ユーザが視覚的またはジェスチャー入力を提供することを要求し得る。別の入力モダリティを使用して追加のユーザ入力を提供するようにとの要求は、スピーチ入力ユーザの記録されたスピーチの再生を示す特性を有するという検出など、前のユーザ入力における異常によってトリガされ得る。代替または追加として、要求は、ランダムに、または多因子認証プロセスのための確立された一連の認証入力の一部として生成され得る。入力モダリティを変更するようにユーザに命令することは、したがって、より高い精度、よりロバストなユーザ認証のために使用され得る。本明細書で使用されるとき、多因子認証プロセスのための認証入力に対応するユーザ入力は、ユーザコマンドに対応するユーザ入力とは別個である。例示のために、コマンドに対応するユーザ入力は、コマンド（たとえば、「ライトをオンにする」）に関連する行為または「スキル」を実施するための命令としてユーザインターフェースによって解釈される一方で、認証入力に対応するユーザ入力は、（たとえば、生体データまたは他のユーザ識別データの比較を介して）ユーザ入力が、記憶されたユーザプロファイルに関連付けられた同じユーザから生起するという尤度を決定するために、記憶されたユーザプロファイルのデータと比較される。

10

【0015】

[0028]その文脈によって明確に限定されない限り、「発生すること」という用語は、計算すること、生成すること、および/または提供することなど、その通常の意味のいずれかを示すために使用される。その文脈によって明確に限定されない限り、「提供すること」という用語は、計算すること、生成すること、および/または発生することなど、その通常の意味のいずれかを示すために使用される。その文脈によって明確に限定されない限り、「結合」されるという用語は、直接的または間接的な電氣的接続または物理的接続を示すために使用される。接続が間接的である場合、「結合」されている構造の間に他のブロックまたは構成要素があり得る。たとえば、ラウドスピーカーは、ラウドスピーカーから壁への（またはその逆への）波（たとえば、音）の伝搬を可能にする介する媒体（たとえば、空気）を介して近くの壁に音響的に結合され得る。

20

30

【0016】

[0029]「構成」という用語は、その特定の文脈によって示されるように、方法、装置、デバイス、システム、またはそれらの任意の組合せに関して使用され得る。「備える」という用語は、本明細書および特許請求の範囲において使用される場合、他の要素または動作を除外しない。（「AはBに基づく」などにおけるような）「に基づく」という用語は、（i）「に少なくとも基づく」（たとえば、「AはBに少なくとも基づく」）、および特定の文脈において適切な場合、（ii）「に等しい」（たとえば、「AはBに等しい」）という場合を含む、その通常の意味のいずれかを示すために使用される。「AはBに基づく」が「に少なくとも基づく」を含む場合（i）、これは、AがBに結合される構成を含み得る。同様に、「に応答して」という用語は、「に少なくとも応答して」を含む、その通常の意味のいずれかを示すために使用される。「少なくとも1つ」という用語は、「1つまたは複数」を含む、その通常の意味のいずれかを示すために使用される。「少なくとも2つ」という用語は、「2つ以上」を含む、その通常の意味のいずれかを示すために使用される。

40

【0017】

[0030]「装置」および「デバイス」という用語は、特定の文脈によって別段に規定されていない限り、総称的および互換的に使用される。別段に規定されていない限り、特定の特徴を有する装置の動作のいかなる開示も、類似の特徴を有する方法を開示すること（その逆も同様）をも明確に意図され、特定の構成による装置の動作のいかなる開示も、類似の構成による方法を開示すること（その逆も同様）をも明確に意図される。「方法」、「

50

プロセス」、「手順」、および「技法」という用語は、特定の文脈によって別段に規定されていない限り、総称的および互換的に使用される。「要素」および「モジュール」という用語は、より大きい構成の一部を示すために使用され得る。「パッケージ」という用語は、ヘッダ部分とペイロード部分とを含むデータのユニットに対応し得る。文書の一部分の参照による任意の組込みはまた、その部分内で参照される用語または変数の定義が、文書内の他の場所、ならびに組み込まれた部分で参照される任意の図に現れる場合、そのような定義を組み込んでいると理解されたい。

【0018】

[0031]本明細書で使用されるとき、「通信デバイス」という用語は、ワイヤレス通信ネットワークを介した音声および/またはデータ通信のために使用され得る電子デバイスを指す。通信デバイスの例は、スマートスピーカー、スピーカーバー、セルラーフォン、携帯情報端末(PDA)、ハンドヘルドデバイス、ヘッドセット、ウェアラブルデバイス、ワイヤレスモデム、ラップトップコンピュータ、パーソナルコンピュータなどを含む。

10

【0019】

[0032]図1は、ユーザ102がマルチモーダルユーザ入力のデバイス110と対話するシステム100を示す。デバイス110は、第1の入力デバイス112と、第2の入力デバイス114と、場合によっては第3の入力デバイス116などの1つまたは複数の追加の入力デバイスと、出力デバイス120と、制御ユニット104とを含む。いくつかの実装形態では、デバイス110は、例示的および非限定的な例として、ポータブル通信デバイス(たとえば、「スマートフォン」)、ウェアラブルデバイス(たとえば、「スマートウォッチ」)、車両システム(たとえば、自動車エンターテインメントシステムとともに使用するための可動もしくはリムーバブルディスプレイ、ナビゲーションシステム、または自動運転制御システム)、あるいは仮想現実または拡張現実ヘッドセットを含むことができる。

20

【0020】

[0033]第1の入力デバイス112は、第1の入力モードに基づく第1のユーザ入力を検出するように構成される。一例では、第1の入力デバイス112は、マイクロフォンを含み、第1の入力モードは、(たとえば、ASR/NLPのための)スピーチモードを含む。例示のために、第1の入力デバイス112は、1つまたは複数のキーワードまたは音声コマンドを含むオーディオ入力をキャプチャするように構成された1つまたは複数のマイクロフォンを含むことができる。

30

【0021】

[0034]第2の入力デバイス114は、第2の入力モードに基づく第2のユーザ入力を検出するように構成される。一例では、第2の入力デバイス114は、カメラを含み、第2の入力モードは、(たとえば、サムズアップ(thumbs-up)またはサムズダウン(thumbs-down)の手の位置、顔の表情など、ユーザ102の視覚的態様を検出するための)ビデオモードを含む。例示のために、第2の入力デバイス114は、1つまたは複数のジェスチャーまたは視覚的コマンドを含むビデオ入力をキャプチャするように構成された1つまたは複数のカメラを含むことができる。

【0022】

40

[0035]第3の入力デバイス116は、第3の入力モードに基づく第3のユーザ入力を検出するように構成される。一例では、第3の入力デバイス116は、ジェスチャトラッカーを含み、第3の入力モードは、ジェスチャーモードを含む。第3の入力デバイス116は、ジェスチャー入力を示すデータ(たとえば、動きデータ)を受信するように構成された1つまたは複数のアンテナを含むことができる。例示のために、ユーザ102は、ユーザの手の移動を追跡する動きセンサー(たとえば、加速度計、ジャイロスコープなど)を含み、動きデータを第3の入力デバイス116に送信する、ブレスレットまたはウォッチを着用することができる。他の実装形態では、動き追跡電子デバイスは、人間ユーザ102中のサイバネティックインプラントなど、ユーザ102と一体化され得るか、またはユーザ102がロボットである実装形態では、ユーザ102の構成要素であり得る。

50

【 0 0 2 3 】

[0036]出力デバイス120は、ラウドスピーカーを使用した可聴出力、ディスプレイを使用した視覚的出力の生成を介して、1つまたは複数の他の出力モダリティ（たとえば、ハプティック）を介して、あるいはそれらの任意の組合せなどで、ユーザ102のために情報を出力するように構成される。たとえば、出力デバイス120は、以下でさらに説明されるように、制御ユニット104からメッセージデータ（たとえば、フィードバックメッセージ144）を受信することができ、ユーザ102への出力（たとえば、命令146）を生成することができる。特定の例では、出力デバイス120は、グラフィカルユーザインターフェースを表現するように構成されたディスプレイ、フィードバックメッセージ144をレンダリングするかまたはユーザ102にダイレクトするように構成された1つまたは複数のラウドスピーカー、あるいはそれらの組合せを含む。

10

【 0 0 2 4 】

[0037]制御ユニット104は、入力デバイス112～116からユーザ入力に対応するデータを受信し、出力デバイス120を介してユーザ102に提供されるべきフィードバックメッセージを生成するように構成される。制御ユニット104は、プロセッサ108と呼ばれる、1つまたは複数のプロセッサに結合されたメモリ106を含む。図2を参照しながらさらに説明されるように、メモリ106は、プロセッサ108による使用のためにアクセス可能な、1つまたは複数の埋め込みネットワークを表すデータと、組み合わされた埋め込み空間への埋め込みベクトルの1つまたは複数の変換を表すデータと、1つまたは複数の分類器を表すデータとを含むことができる。メモリ106はまた、マルチモーダル認識エンジン130、フィードバックメッセージ生成器132、またはそれらの両方を実装するためにプロセッサ108によって実行可能な命令を含むことができる。

20

【 0 0 2 5 】

[0038]プロセッサ108は、マルチモーダル認識エンジン130と、フィードバックメッセージ生成器132とを含む。いくつかの実装形態では、プロセッサ108は、マルチモーダル認識エンジン130とフィードバックメッセージ生成器132とを実装するための命令を実行するように構成された1つまたは複数の処理コアを含む。いくつかの実装形態では、プロセッサ108は、マルチモーダル認識エンジン130とフィードバックメッセージ生成器132との一方または両方を実装するように構成された専用回路を含む。一例では、プロセッサ108は、集積回路（IC）として実装される。

30

【 0 0 2 6 】

[0039]マルチモーダル認識エンジン130は、入力デバイス112～116のうちの1つまたは複数からデータを受信し、出力を生成するために受信データを処理するように構成される。たとえば、出力は、受信された入力に最も密接に一致するコマンドと、コマンドに関連する確信度（または尤度）インジケータとを含むことができる。いくつかの実装形態では、マルチモーダル認識エンジン130は、各入力モダリティについて、各入力モダリティの埋め込みベクトルを生成することなどによって、特定のトレーニングされたユーザ入力を他のトレーニングされたユーザ入力から区別するためのデータを生成するように構成される。マルチモーダル認識エンジン130は、ユニモーダルまたはマルチモーダルユーザ入力の一部として、入力デバイス112～116の各々を介して（もしあれば）どの認識されたユーザ入力が発見されたかを示す、組み合わされた埋め込みベクトルを生成するために、異なる入力モダリティに関連する埋め込みベクトルを組み合わせるように構成され得る。組み合わされた埋め込みベクトルは、組み合わされた埋め込みベクトルをコマンドにマッピングするようにトレーニングされた分類器を使用することなどによって、出力を決定するように処理される。マルチモーダル認識エンジン130において実装され得る構成要素の例示的な例については、図2に関して説明される。

40

【 0 0 2 7 】

[0040]フィードバックメッセージ生成器132は、出力デバイス120を介してユーザ102に出力されるべきフィードバックメッセージデータを生成するように構成される。たとえば、フィードバックメッセージ生成器132は、確信度レベルがしきい値を下回る

50

特定のコマンドであることが予測されるなど、適切に認識されなかったユーザ入力を繰り返すようにユーザ102に命令するために、フィードバックメッセージ144を出力デバイス120に送ることができる。他の例として、フィードバックメッセージ生成器132は、入力モダリティを変更するように、または1つの入力モダリティを使用して行われる入力を異なる入力モダリティを使用して行われる別の入力でオーグメントするようにユーザ102に命令するために、フィードバックメッセージ144を出力デバイス120に送ることができる。他の例は、ユーザ102がエミュレートするためのユーザ入力の記録サンプル、ユーザ102が識別するためのユーザの入力の記録サンプル、またはユーザ102がデバイス110を使用するのを支援するための他の情報を提供する、フィードバックメッセージデータを生成することを含む。例示的な例は、ユーザ102からの問合せを受信したことに応答して「アップ」に対応する動きを示すモーションビデオを表示すること、アップジェスチャーの動きに関連する最も類似している発話のオーディオ再生を生成すること、またはユーザ定義の動きにすでに密接に関連付けられている関係する発話のオーディオ再生を生成することなど、クロスモーダルサンプル取出しを含む。いくつかの例では、フィードバックメッセージ生成器132は、以下でより詳細に説明されるように、多因子認証プロセスに従って次の認証入力を提供するようにユーザ102に命令するために、フィードバックメッセージ144を生成するように構成される。

【0028】

[0041]動作中に、ユーザ102は、第1の入力デバイス112によって検出された第1の入力モード（たとえば、バーバルコマンド）に基づいて、第1の入力140を提供する。第1の入力デバイス112は、第1の入力140を示す第1のデータ142を生成し、第1のデータ142を制御ユニット104に提供する。

【0029】

[0042]プロセッサ108（たとえば、マルチモーダル認識エンジン130）は、第1の入力モード（たとえば、スピーチ）に基づくユーザ102からの第1の入力140を示す第1のデータ142を処理する。プロセッサ108（たとえば、フィードバックメッセージ生成器132）は、第1のデータ142の処理に基づいて出力デバイス120にフィードバックメッセージ144を送る。フィードバックメッセージ144は、たとえば、話された命令146のプレイアウトを介して、異なる入力モードを使用して第2の入力148を提供するようにユーザ102に命令する。第2の入力148は、第1の入力モードとは異なる第2の入力モード（たとえば、ビデオ）に基づき、マルチモーダル認識エンジン130が第1の入力140にどのように応答するかを更新するために使用され得る。本明細書で使用されるとき、異なる入力モードを使用することは、同じタイプの入力を使用するのではなく、異なるタイプの入力を使用することを意味する。各異なるタイプの入力は、様々な異なるセンサーを使用する。たとえば、スピーチ入力モードは、1つまたは複数のマイクロフォンを使用し得る。ジェスチャー入力モードは、動き検出を使用し得る。ビデオ入力モードは、カメラと、フレームのシーケンスとを使用し得る。概して、各入力モードは、その入力を提供するために使用され得る異なるタイプのセンサーを提供する。

【0030】

[0043]いくつかの実装形態では、第1の入力140は、コマンドであり、フィードバックメッセージ144は、第1の入力140をディスアンビギュエートするための第2の入力148を提供するようにユーザ102に命令する。マルチモーダル認識エンジン130は、出力の不確実性（たとえば、話された入力が「アップ」を示すのか「オフ」を示すかの不確実性）を示す、第1の入力140の認識処理に関連する確信度レベルが確信度しきい値を満たすのに失敗したことに応答して、フィードバックメッセージ144を送り得る。ユーザ102は、第2の入力148（たとえば、上方ポインティング）を提供し得、第2の入力148を示す第2のデータ150に基づいて、マルチモーダル認識エンジン130は、図2においてさらに詳細に説明されるように、第1の入力140（たとえば、スピーチ「アップ」）のマッピングを、第2の入力148に関連付けられた行為（たとえば、音楽ボリュームを上げる）に対して更新することができる。

10

20

30

40

50

【0031】

[0044]別の実装形態では、マルチモーダル認識エンジン130は、第2の入力148と組み合わせられた第1の入力140のマッピングを、第2の入力148に関連する行為に対して更新する。たとえば、雑音条件が、話された「アップ」コマンドの信頼できる認識を妨げるとき、マルチモーダル認識エンジン130は、ボリュームをアップするための単一のコマンドとして、ユーザの話された「アップ」コマンドと併せてユーザの「アップ」ビデオ入力（たとえば、上方ポインティング）を認識するように更新される。

【0032】

[0045]したがって、いくつかの実装形態では、ユーザ102は、フィードバックメッセージ144のフィードバック機構と第2の入力148とを介して特定の行為を実施するためのコマンドとして特定の入力を認識するようにデバイス110を個人化することができる。例示のために、ユーザ102は、マルチモーダル認識エンジン130によって現在認識されないコマンド（第1の入力140）を話すことができ、フィードバックメッセージ144に回答して、ユーザ102は、認識されたコマンド（第2の入力148）を入力することによって、この認識されないコマンドにマッピングされるべき行為を識別することができる。同様に、デバイス110は、ユーザの選定されたモードが信頼できなくなったとき、入力モードを変更するようにユーザ102に命令することができる。たとえば、デバイス110が車両（たとえば、カーナビゲーションおよび/またはエンターテインメントシステム）中に実装されたとき、夜間運転中に、ユーザ102は、（低い照明条件により）ビデオの代わりにスピーチ入力またはジェスチャー入力を使用するように命令され得、ウィンドウが開いた状態で運転しているとき、ユーザ102は、（高い風雑音により）スピーチの代わりにジェスチャー入力またはビデオ入力を使用するように命令され得る。デバイス110が、仮想現実または拡張現実ヘッドセットなどのヘッドセット中に実装されたとき、入力モードを変更するようにユーザ102に命令するための同様の動作が実施され得る。

【0033】

[0046]他の実装形態では、デバイス110は、多因子認証を実施するために使用される。たとえば、第1の入力140は、ユーザ102の第1の認証行為（たとえば、スピーカー検証のための話されたパスコード）に対応し得、フィードバックメッセージ144は、多因子認証手順の一部として、第2の認証行為として第2の入力148を提供する（たとえば、ユーザ102によって以前に選択された特定の手の構成を表示する）ようにユーザ102に命令する。デバイス110は、認証行為を実施するようにユーザ102に命令するための認証入力モードの数およびタイプをランダムにまたはアルゴリズム的に選択することができる。たとえば、デバイス110は、スピーチ入力（たとえば、第1の入力140）が、プレイアウトされている記録スピーチであり得るというインジケーションに回答して、命令146を生成することができ、カメラ（たとえば、第2の入力デバイス114）にウィンクするようにユーザ102に命令することなどによって、「ライブリネス」確認を要求し得る。

【0034】

[0047]上記の例では、第2の入力148が第1の入力140とは異なるモードを使用することについて説明しているが、他の実装形態では、第2の入力148は、第1の入力140と同じモードを使用することができる。たとえば、第1の入力140の話されたコマンドは、解釈するのが困難であり得るが（たとえば、周辺雑音の存在下での「アップ」対「オフ」）、別の話されたコマンド（たとえば、「より大きく」）は、正しい行為（たとえば、ボリュームを上げる）を選択するために、他のマッピングされたコマンドとは十分に異なり得る。別の例として、トレーニングプロセス中に、ユーザ102は、トレーニングされていないスピーチコマンドとして「より大きく」発話し得、デバイス110は、「より大きく」という発話に関連付けられるべき行為を識別するように、命令146を介してユーザ102に命令し得る。ユーザ102は、ボリュームを上げるためのコマンドとしてデバイス110によって認識される第2の話された発話「アップ」を提供し得、マルチ

10

20

30

40

50

モーダル認識エンジン 130 は、「より大きく」を「ボリュームを上げる」行為にマッピングするように、ユーザ入力のマッピングを更新し得る。

【0035】

[0048]図2は、特定の実装形態による、メモリ106と、マルチモーダル認識エンジン130と、プロセッサ108によって実行可能である1つまたは複数のアプリケーション240とを含む、制御ユニット104の構成要素の一例を示す。マルチモーダル認識エンジン130は、第1のユーザ入力(たとえば、スピーチ入力)を第1の埋め込みベクトル(たとえば、第1の埋め込みベクトル「E1」)にコンバートするように構成された第1の埋め込みネットワーク202を含む。第2の埋め込みネットワーク204は、第2のユーザ入力(たとえば、ジェスチャー入力)を第2の埋め込みベクトル(たとえば、第2の埋め込みベクトル「E2」)にコンバートするように構成される。マルチモーダル認識エンジン130は、第Nのユーザ入力(たとえば、ビデオ入力)を第Nの埋め込みベクトル(たとえば、第Nの埋め込みベクトル「En」)にコンバートするように構成された第Nの埋め込みネットワーク206を含む、1つまたは複数の追加の埋め込みネットワークを含み得る。マルチモーダル認識エンジン130は、本開示のいくつかの実施形態による任意の数の埋め込みネットワークを含み得る。

10

【0036】

[0049]融合埋め込みネットワーク(fusion embedding network)220は、埋め込みネットワーク202~206の出力を組み合わせて、組み合わせられた埋め込みベクトル「C」228など、組み合わせられた埋め込みベクトルを生成するように構成される。たとえば、第1の変換212は、第1の共通の埋め込みベクトル222を生成するために、スピーチ埋め込みベクトルを「共通」の埋め込み空間にコンバートすることができる。第2の変換214は、第2の共通の埋め込みベクトル224を生成するために、ジェスチャー埋め込みベクトルを共通の埋め込み空間にコンバートすることができ、第Nの変換216は、第Nの共通の埋め込みベクトル226を生成するために、ビデオ埋め込みベクトルを共通の埋め込み空間にコンバートすることができる。共通の埋め込みベクトル222~226の各々は、それぞれ、対応する重みW1、W2、およびW3で重み付けされ、融合埋め込みネットワーク220において組み合わせられ得る。マッピング230は、組み合わせられた埋め込みベクトル228に対応する出力232と確信度レベル234とを選択するように構成される。たとえば、マッピング230は、組み合わせられた埋め込みベクトルを特定の行為にマッピングするように構成された分類器231を含むことができる。例示のために、複数の埋め込みネットワーク202~206への組み合わせられた入力から生じる出力232を決定するために、各モダリティ入力について個々の分類器を使用するのではなく、単一の分類器231が使用される。

20

30

【0037】

[0050]マルチモーダル認識エンジン130によって使用される1つまたは複数のパラメータを示すデータは、メモリ106に記憶される。第1のユーザプロファイル250は、第1のユーザ(たとえば、ユーザ102)に関連付けられ、第1の埋め込みネットワークデータ252と、第1の重みデータ254と、第1の一時的調整データ256と、第1の履歴データ258とを含む。第1の埋め込みネットワークデータ252は、第1のユーザに対応すべき、第1の埋め込みネットワーク202と、第2の埋め込みネットワーク204と、第Nの埋め込みネットワーク206と、融合埋め込みネットワーク220とを含む、埋め込みネットワークを構成するためのデータ(たとえば、重みまたは他のパラメータもしくは値)を含む。第1の重みデータ254は、第1のユーザに対応すべき重み(たとえば、W1、W2、W3)を構成するための重み値を含む。第1の一時的調整データ256は、以下でさらに説明されるように、(たとえば、雑音の多い環境では重みW1を低減し、重みW2およびW3を増加させるための)一時的条件に基づいてマルチモーダル認識エンジン130の構成を調整するための値を含む。第1の履歴データ258は、第1のユーザに関連するヒストリカルデータを含み、マルチモーダル認識エンジン130によって処理される第1のユーザのマルチモーダル入力に対応する履歴傾向に基づいて、プロセッ

40

50

サ 1 0 8 が第 1 の埋め込みネットワークデータ 2 5 2、第 1 の重みデータ 2 5 4、またはそれらの両方を更新することを可能にする。

【 0 0 3 8 】

[0051]同様に、メモリ 1 0 6 は、第 2 のユーザのための第 2 の埋め込みネットワークデータ 2 6 2 と、第 2 の重みデータ 2 6 4 と、第 2 の一時的調整データ 2 6 6 と、第 2 の履歴データ 2 5 8 とを含む、第 2 のユーザに関連付けられた第 2 のユーザプロファイル 2 6 0 を含む。第 1 の埋め込みネットワークデータ 2 5 2 は、第 1 のユーザと第 2 のユーザとの間の入力コマンドの差に基づいて第 2 の埋め込みネットワークデータ 2 6 2 とは異なる。たとえば、第 1 のユーザと第 2 のユーザは、ビデオ入力を実施するときの異なるアクセント、異なるスタイルのジェスチャリング、異なる身体力学、またはそれらの任意の組合せを有し得る。第 1 の埋め込みネットワークデータ 2 5 2 は、第 1 のユーザのための埋め込みネットワークデータのデフォルトセットからのユーザ固有の変動を認識するように埋め込みネットワーク 2 0 2 ~ 2 0 6 および 2 2 0 をトレーニングした結果を表し得、第 2 の埋め込みネットワークデータ 2 6 2 は、第 2 のユーザのための埋め込みネットワークデータのデフォルトセットからのユーザ固有の変動を認識するように埋め込みネットワーク 2 0 2 ~ 2 0 6 および 2 2 0 をトレーニングした結果を表し得る。ただ 2 つのユーザプロファイル 2 5 0、2 6 0 が示されているが、デバイス 1 1 0 の複数のユーザのためのマルチモーダル認識エンジン 1 3 0 の動作をカスタマイズするために、任意の数のユーザプロファイルが含まれ得る。

10

【 0 0 3 9 】

[0052]異なるアクセント、ジェスチャースタイル、および身体力学など、個々のユーザ変動について調整することに加えて、第 1 の埋め込みネットワークデータ 2 5 2 はまた、第 1 のユーザによって決定されたユーザ入力の第 1 のカスタマイズされたセットを認識するように埋め込みネットワーク 2 0 2 ~ 2 0 6 および 2 2 0 をトレーニングした結果を表し得、第 2 の埋め込みネットワークデータ 2 6 2 はまた、第 2 のユーザによって決定されたユーザ入力の第 2 のカスタマイズされたセットを認識するように埋め込みネットワーク 2 0 2 ~ 2 0 6 および 2 2 0 をトレーニングした結果を表し得る。たとえば、第 1 のユーザは、オーディオ再生動作が進行中の間、スピーチコマンド「アップ」をボリュームを上げるためのコマンドとして認識するようにマルチモーダル認識エンジン 1 3 0 をカスタマイズ（たとえば、トレーニング）し得る。対照的に、第 2 のユーザは、オーディオ再生動作が進行中の間、スピーチコマンド「アップ」をプレイリスト上の前のオーディオトラックを選択するためのコマンドとして認識するようにマルチモーダル認識エンジン 1 3 0 をカスタマイズ（たとえば、トレーニング）し得る。

20

30

【 0 0 4 0 】

[0053]第 1 の重みデータ 2 5 4 は、第 1 のユーザと第 2 のユーザとの間の入力モード信頼性の差に基づいて第 2 の重みデータ 2 6 4 とは異なり得る。たとえば、プロセッサ 1 0 8 は、第 1 の履歴データ 2 5 8 などに基づいて、第 1 のユーザからのスピーチ入力、第 1 のユーザからのジェスチャー入力と比較してあまり確実に解釈されないと決定し得る。その結果、第 1 のユーザからのスピーチ入力への依拠を低減し、ジェスチャー入力への依拠を増加させるために、第 1 の重みデータ 2 5 4 において、重み W 1 は、デフォルト W 1 値から低減され得、重み W 2 は、デフォルト W 2 値から増加され得る。対照的に、プロセッサ 1 0 8 は、第 2 の履歴データ 2 6 8 などに基づいて、第 2 のユーザからのスピーチ入力、第 2 のユーザからのジェスチャー入力と比較してより確実であると決定し得る。その結果、第 2 のユーザからのジェスチャー入力への依拠を低減し、スピーチ入力への依拠を増加させるために、第 2 の重みデータ 2 6 4 において、重み W 1 は、デフォルト W 1 値から増加され得、重み W 2 は、デフォルト W 2 値から減少され得る。

40

【 0 0 4 1 】

[0054]アプリケーション 2 4 0 は、一時的調整器 2 9 0 と、データ調整器 2 9 2 とを含む。一時的調整器 2 9 0 は、一時的条件に基づいて、埋め込みネットワーク 2 0 2、2 0 4、2 0 6、または 2 2 0 のうちの 1 つまたは複数の調整、重み W 1 ~ W 3 のうちの 1 つ

50

または複数の調整、あるいはそれらの組合せを決定するように構成される。たとえば、一時的調整器 290 は、検出された条件に基づいて、1 つまたは複数の入力モダリティを強調するように、1 つまたは複数の入力モダリティを強調しないように、あるいはそれらの組合せを行うように、重み $W_1 \sim W_3$ のうちの 1 つまたは複数の調整を調整することができる。例示的および非限定的な例として、検出された条件は、以下でさらに詳細に説明されるように、周辺雑音データ 272、周辺光データ 274、ロケーションデータ 276、またはユーザ選好 278 のうちの 1 つまたは複数によって示され得る。

【0042】

[0055] データ調整器 292 は、一時的条件に基づかないと決定された変化を表すように埋め込みネットワークデータと重みデータとを更新するために、埋め込みネットワーク 202、204、206、または 220 のうちの 1 つまたは複数の調整、重み $W_1 \sim W_3$ のうちの 1 つまたは複数の調整、あるいはそれらの組合せを決定するように構成される。いくつかの実装形態では、データ調整器 292 は、たとえば、マルチモーダル認識エンジン 130 がユーザ入力をより正確に認識する（たとえば、話されたコマンドのユーザの発音とデフォルトスピーチ認識モデルとの間の差に適應する）のを助けるユーザからのディスアンビギュエーションフィードバックを受信したことに応答して、または特定のコマンドへの入力のカスタムマッピングを示すユーザ入力（たとえば、ユーザが、以前不明であったビデオ入力として両手で「サムズアップ」ジェスチャーを入力し、このビデオ入力により、デバイス 110 がアラームをオフすべきであることを示す）に応答して、特定のコマンドへのユーザ入力の更新されたマッピングを示すために、埋め込みネットワーク 202、204、206、または 220 のうちの 1 つまたは複数に対して更新トレーニングを実施するように構成される。

【0043】

[0056] 図 1 のシステム 100 中に実装されたマルチモーダル認識エンジン 130 の動作の例示的な例では、ユーザ 102 は、顔認識、音声認識、または何らかの他の形態のユーザ認識などを介して、デバイス 110 へのマルチモーダル入力のソースとして識別される。ユーザ 102 からの入力を認識するようにマルチモーダル認識エンジン 130 を構成（たとえば、カスタマイズ）するために、埋め込みネットワーク 202 ~ 206 は、第 1 の埋め込みネットワークデータ 252 に基づいて更新され、重み W_1 、 W_2 、および W_3 は、第 1 の重みデータ 254 に基づいて更新され、いずれかの一時的調整は、第 1 の一時的調整データ 256 に基づいて適用される。

【0044】

[0057] ユーザ 102 は、コマンドとして第 1 の入力 140 を提供する。第 1 の入力 140 は、十分な信頼性で何らかの特定のコマンドとして認識されず、フィードバックメッセージ 144 は、第 1 の入力 140 をディスアンビギュエートするための第 2 の入力 148 を提供するようにユーザ 102 に命令する。たとえば、フィードバックメッセージ 144 は、出力 232 の不確実性（たとえば、話された入力が「アップ」を示すのか「オフ」を示すのかの不確実性）を示す、第 1 の入力 140 の認識処理に関連する確信度レベル 234 が確信度しきい値 294 を満たすのに失敗したことに応答して送られ得る。他の実装形態では、フィードバックメッセージ 144 は、1 つまたは複数の環境条件が検出されたことに応答して送られる。

【0045】

[0058] たとえば、第 1 の入力 140 がビデオモードを介して受信される実装形態では、フィードバックメッセージ 144 は、照明しきい値 286 を下回る値を有する周辺光メトリック 284 に応答して送られる。たとえば、周辺光データ 274 は、デバイス 110 の 1 つまたは複数のセンサーを介して受信され、周辺光メトリック 284 を生成するために処理され得る。周辺光メトリック 284 は、周辺照明が、信頼できるビデオモード入力のためには薄暗すぎるかどうかを決定するために、照明しきい値 286 と比較され得る。フィードバックメッセージ 144 は、薄暗い照明がビデオ入力モードを信頼できないものに行っていることをユーザに通知し得、別のモダリティ（たとえば、スピーチ）を使用して入

10

20

30

40

50

力を繰り返すようにユーザに命令し得る。

【0046】

[0059]別の例として、第1の入力140がスピーチモードを介して受信される実装形態では、フィードバックメッセージ144は、雑音しきい値282を上回る値を有する雑音メトリック280（たとえば、信号対雑音比（SNR）または周辺雑音測定値）にตอบสนองして送られる。たとえば、周辺雑音データ272は、デバイス110の1つまたは複数のセンサーを介して受信され（あるいはマイクロフォン入力信号の音声アクティビティ検出処理中に測定され）、雑音メトリック280を生成するために処理され得る。雑音メトリック280は、周辺雑音が、信頼できるスピーチモード入力のためには大きすぎるかどうかを決定するために、雑音しきい値282と比較され得る。フィードバックメッセージ144は、雑音環境がスピーチ入力モードを信頼できないものになっていることをユーザに通知し得、別のモダリティ（たとえば、ビデオ）を使用して入力を繰り返すようにユーザに命令し得る。

10

【0047】

[0060]ユーザ102は、第2の入力148（たとえば、上方ポインティング）を提供し得、第2の入力148を示す第2のデータ150に基づいて、マルチモーダル認識エンジン130は、第1の入力140（たとえば、スピーチ「アップ」）のマッピングを、第2の入力148に関連付けられた行為（たとえば、音楽ボリュームを上げる）に対して更新することができる。例示のために、第1の埋め込みネットワーク202、第1の変換212、重みW1、融合埋め込みネットワーク220、またはマッピング230のうちの1つまたは複数は、マルチモーダル認識エンジン130が、音楽ボリュームを上げるためのコマンドとしてユーザの話された「アップ」をより正確に認識することを引き起こすように、データ調整器292によって調整され得る。

20

【0048】

[0061]動作の例示的な例では、1つの入力モダリティが低精度条件を有すると決定された場合、（たとえば、一時的調整器290によって生成された一時的調整データにตอบสนองして）マルチモーダル認識エンジン130は、組み合わされた埋め込みベクトル228の生成のためにそのモダリティを使用する入力の影響を低減するかまたはなくすように1つまたは複数の設定を調整する。（たとえば、雑音しきい値282を超える雑音メトリック280により）スピーチモダリティが信頼できないと決定される、予測される、または推定される一方で、ジェスチャーおよびビデオモダリティが入力認識のために十分に信頼できると決定される例では、一時的調整器290は、スピーチ入力に関連する共通の埋め込みベクトル222に適用される重みW1を「0」値に設定し得る。ジェスチャー入力に関連する共通の埋め込みベクトル224に適用される重みW2と、ビデオ入力に関連する共通の埋め込みベクトル226に適用される重みW3とは、非0値に設定される（たとえば、ジェスチャー入力とビデオ入力等しく信頼できるように扱われる実装形態では、 $W2 = W3 = 0.5$ ）。重みW1を「0」値に設定することにより、スピーチ入力に信頼できない状態である間、スピーチ入力に、得られた組み合わされた埋め込みベクトル228に影響を及ぼすのを防止する。

30

【0049】

[0062]上記の例によれば、重みの初期設定は、各モダリティが入力認識について等しい重要性または信頼性を有することを示す、 $W1 = W2 = W3 = 1/3$ を割り当て得る。スピーチモダリティは、（たとえば、雑音メトリック280が雑音しきい値282を超えるという検出、もしくは車両が動いている間に車両ウィンドウが開いているという検出を介した）大量の周辺雑音の検出により、またはスピーチ入力のしきい値数が所定の時間期間中に正確に認識されることに失敗することなどにより、信頼できないと後で決定または予測され得る。スピーチモダリティが信頼できないと決定または予測されたことにตอบสนองして、一時的調整器290は、入力認識に対するスピーチ入力の影響を除去するために、重みW1、W2、およびW3を、それぞれ0、1/2、および1/2に調整する。スピーチ入力モダリティがもはや信頼できなくはないという後続の決定（たとえば、風雑音が雑音し

40

50

きい値を下回るか、ウィンドウが閉じられるか、または車両が移動するのを止めた)に
答して、重みW 1、W 2、およびW 3は、1 / 3のそれらの初期値にそれぞれ戻され得る。

【0050】

[0063]別の例として、代わりに、ビデオモダリティが、周辺光の低い量の検出(たと
えば、周辺光メトリック284が照明しきい値286を下回る)により、またはビデオ入力
のしきい値数が所定の時間期間中に正確に認識されることに失敗することなどにより、信
頼できないと決定または予測され得る。ビデオモダリティが信頼できないと決定または予
測されたことに応答して、一時的調整器290は、入力認識に対するビデオ入力の影響を
除去するために、重みW 1、W 2、およびW 3を、それぞれ1 / 2、1 / 2、および0に
調整する。ビデオ入力モダリティがもはや信頼できなくはないという後続の決定(たと
えば、周辺光が照明しきい値を超えることを決定される)に
答して、重みW 1、W 2、お
よびW 3は、1 / 3のそれらの初期値にそれぞれ戻され得る。

10

【0051】

[0064]いくつかの実装形態では、複数の重みは、入力認識に対する複数の入力モダリ
ティのインパクトを低減または除去するように調整される。たとえば、スピーチモダリ
ティのみが使用されるべきであるという決定が行われる実装形態では、W 1は「1」に設定さ
れ、W 2とW 3とは「0」に設定される。例示のために、デバイス110は、低い周辺照
明条件を検出し得、また、アクティブなジェスチャー検出デバイスが検出されない(たと
えば、ユーザのスマートウォッチが存在しないかまたは動きデータを送信していない)と
決定し得る。別の例として、ユーザ102は、スピーチ入力のみを処理するように入力認
識を制限するように、ユーザ選好278を入力することなどを介して、デバイス110に
命令し得る。別の例として、1つまたは複数の入力モダリティを制限すべきかどうかを決
定するために、ロケーションデータ276が使用され得る。たとえば、ユーザが車両を操
作していることを示すロケーションデータ276に
応答して、一時的調整器290は、ユ
ーザの注意散漫を防ぐために、および車両の安全な操作を奨励するためになど、ジェス
チャー入力とビデオ入力との認識を防止するようにユーザ入力モードを制限し得る。ユ
ーザがもはや車両を操作しておらず、ユーザの自宅にいることを示すロケーションデータ27
6に
応答して、一時的調整器290は、ジェスチャー入力とビデオ入力との認識を有効に
するようにユーザ入力モードを復元し得る。

20

【0052】

[0065]上記の例では重み値の例について説明されているが、そのような例示的な重み値
は、例示的であり、限定的ではない。例示のために、重みを「0」に設定するのではなく
、重みは、全体的な入力認識に対する関連する入力モダリティの影響を減少させるが
なくしはしない、低減された値に設定され得る。別の例として、「信頼できる」入力モ
ダリティは、入力モダリティの相対的信頼性を示し得る、等しくない重みを有し得る。
例示のために、ジェスチャー入力が十分に信頼できると見なされ、ビデオ入力が
ジェスチャー入力よりも信頼できると見なされ、スピーチが信頼できないと決定され
た場合、重みは、W 1 = 0 . 1、W 2 = 0 . 4、およびW 3 = 0 . 5などの値に設定され得る。
上記の例では、重みW 1、W 2、およびW 3の和は1に等しいが、他の実装形態では、
重みW 1、W 2、およびW 3の和は、どんな特定の値にも制限されない。

30

40

【0053】

[0066]信頼できないと決定された入力モダリティの影響を低減するかまたはなくすよ
うに1つまたは複数の重みを調整することの追加または代替として、いくつかの実装形
態では、マルチモーダル認識エンジン130は、関連する埋め込みネットワークの出力を、
利用可能なスキルの中から「なし」出力に強制するか、変換の出力を、「0」値を有する
埋め込みベクトルへの「なし」カテゴリー入力のために共通の埋め込み空間に強制する
か、またはそれらの組合せを行い得る。

【0054】

[0067]いくつかの実装形態では、マルチモーダル認識エンジン130を含むデバイス1
10は、複数の入力モダリティの環境アウェア融合を実施する。たとえば、ユーザ102

50

が車を運転していると決定したことに応答して、ジェスチャー入力に関連する重みW2は、車を運転している間に安全でない手の動きを阻止するために、ユーザの手の動きがジェスチャー入力としてよりもむしろ車の操作に対応する可能性があることを示す、「0」に設定され得る。別の例として、ユーザ102が暗い部屋の中にいると決定したことに応答して、ビデオ入力に関連する重みW3は、「0」に設定され得る。別の例として、ユーザ102が雑音の多い環境の中にいると決定したことに応答して、スピーチ入力に関連する重みW1は、「0」に設定され得る。環境条件の決定は、デバイス110に組み込まれた1つまたは複数のセンサー（たとえば、周辺光センサー、周辺雑音センサー）、（たとえば、デバイス110と、ホームオートメーションシステム、モノのインターネットシステム、または別のシステムの1つまたは複数の構成要素との間の通信を介した）デバイス110の外部にある1つまたは複数のセンサー、あるいはそれらの任意の組合せに基づくことができる。

10

【0055】

[0068]図3は、ヘッドセット302を着用しているユーザが、スマートフォンなどの別のデバイス、車などの車両システム、またはワイヤレスデジタルアシスタントアプリケーションを組み込んでいるスピーカーシステム（たとえば、「スマートスピーカー」と通信している、マルチモーダルユーザ入力のためのシステム300の一例を示す。ヘッドセット302は、図1のデバイス110に対応することができ、拡張現実（「AR」）、仮想現実（「VR」）、または複合現実（「MR」）オーディオおよびビデオ出力を着用者に提供するために、ディスプレイと、イヤバッド308または他のウェアラブル雑音生成デバイスなどのトランスデューサとを含むことができる。

20

【0056】

[0069]ヘッドセット302は、ユーザ入力を検出するために、1つまたは複数のマイクロフォン、1つまたは複数のカメラなど、複数のセンサーを含むことができる。たとえば、1つまたは複数のマイクロフォンを介して受信されたオーディオ入力は、ヘッドセット302に組み込まれたかまたはそれに結合されたプロセッサにおいて1つまたは複数の動作310を実施するために使用され得る。たとえば、音環境分類を可能にするための機械学習、ヘッドセット302の着用者がいつ話しているかを決定するための自己音声の音声アクティビティ検出（VAD）、音響イベント検出、およびモード制御（たとえば、シーケンススペースのユーザインターフェース）を使用することなど、オーディオ入力に対応するオーディオ信号を処理することが実施され得る。

30

【0057】

[0070]1つまたは複数の動作310の結果は、1つまたは複数の行為312を生成するために使用され得る。たとえば、行為312は、アクティブ雑音除去（ANC）フィルタをチューニングすること、1つまたは複数の支援的リスニング特徴を実装すること、マルチマイクロフォン音キャプチャのフィールドを調整すること（たとえば、「AudioZoom」）、あるいは拡張現実レンダリング、仮想現実レンダリング、または複合現実レンダリング（まとめて「XR」レンダリングと呼ばれる）を実施することを含むことができる。たとえば、結果は、空間透過モードでヘッドセット302にレンダリングされ得る。

【0058】

40

[0071]ヘッドセット302において（たとえば、1つまたは複数のマイクロフォン、動き検出器、ジェスチャー検出器、カメラなどを介して）検出されたユーザ入力は、自動音声認識および自然言語処理、探索もしくは問合せ応答、またはそれらの両方など、1つまたは複数のスピーチベースの動作304の実施を開始するために使用され得る。1つまたは複数のスピーチベースの動作304は、ヘッドセット302と通信しているスマートフォンまたは他のポータブル通信デバイスなどにおいて、機械学習を使用して実施され得る。データ通信305（たとえば、ワイヤレスネットワーク通信、ワイヤライン通信、またはそれらの両方）は、外部処理リソース306（たとえば、機械学習を組み込んでいるクラウドベースASR/NLPおよび探索サーバ）にオーディオスピーチデータを送ることを含み得る。探索および問合せ結果は、ヘッドセット302を介してユーザに返信され得

50

る。

【 0 0 5 9 】

[0072]図4は、例示的および非限定的な例では図3のヘッドセット302などによって実施され得る、マルチマイクロフォン音キャプチャのフィールドを調整すること（たとえば、「AudioZoom」）の例400を示す。代表的なマイクロフォン412、414、および416など、複数のマイクロフォンが、ユーザの周りに配置される。ユーザは、極座標系の中心におり、0度角度方向を向くように配向されるものとして示されている。マイクロフォン412、414、および416は、指向性マイクロフォン、無指向性マイクロフォン、またはそれらの両方を含み、ユーザの周囲のオーディオ環境をキャプチャすることができる。第1の構成402では、マイクロフォン412～416からのオーディオの音処理は、ユーザ指示調整なしのオーディオ環境を表す、（たとえば、イヤホンまたはイヤパッドを介した）ユーザへの可聴出力を生じる。

10

【 0 0 6 0 】

[0073]第2の構成404では、マルチモーダルインターフェース（たとえば、例示的な例として、ユーザジェスチャー、発話、ビデオ入力、またはそれらの組合せ）を介したユーザ入力に応答して、マイクロフォン412～416からのオーディオの音処理は、特定の空間領域420（たとえば、90度角度方向の、またはユーザの左側の領域）から生起または到着する音を強調（たとえば、増幅）する一方で、空間領域420外のエリアから生起する音を減衰させるように調整される。第2の構成404に遷移することを生じるユーザ入力の例は、例示的および非限定的な例として、スピーチモダリティに基づく「左にズームする」スピーチシーケンス、ジェスチャーモダリティに基づく「手を左側にポインティングする」または「指を左側にポインティングする」ジェスチャーシーケンス、あるいはオーディオ（非スピーチ）モダリティに基づく「スナップ音を起こす」オーディオシーケンスを含むことができる。

20

【 0 0 6 1 】

[0074]いくつかの実装形態では、図1～図4を参照しながら上記で説明されたマルチモーダルインターフェースは、ユーザの近傍にあるロケーションまたはアクティビティ（たとえば、リビングルームでテレビジョンを見ること、またはキッチンで皿を洗うこと）などのコンテキストに応答する。たとえば、ウォッチまたはアームバンドベースの加速度計を使用してキャプチャされたジェスチャーは、検出されたコンテキストに基づいて解釈され得る。たとえば、手を振ることは、ターゲットコマンド「ライトをオンにする」として解釈され得、手を左側に反転させることは、「次の曲」または「次のチャンネル」として解釈され得、手を右側に反転させることは、「前の曲」、「前のチャンネル」、または「ドアオープン」として解釈され得る。例示的および非限定的な例として、閉じられた拳が形成される「クラブ」ジェスチャーは、「電話を取る」または「チャンネルを選択する」として解釈され得、長いクラブは、「曲を止める」、「アラームをキャンセルする」、または「ドアクローズ」として解釈され得、指を伸ばしている手の反時計回りの回転は、「ホームデバイスを発見する」として解釈され得る。コンテキストは、検出された音響イベント/環境シーケンスとの関連付けを介して決定され得る。たとえば、様々な音響イベントは、音響環境（たとえば、ユーザがどこにいるか）を推論するために、または適切なフィードバックタイミングを監視するために検出され得る。そのような検出可能な音響イベントの例は、ヘアドライヤー、掃除機、音楽、キッチンフード、料理、食事、皿の洗浄、屋内空調、電子レンジ、洗濯機、乾燥機、シャワー、およびテレビジョンを見ることを含む。

30

40

【 0 0 6 2 】

[0075]手のジェスチャー認識のためのデータセットは、手の移動を示す（たとえば、x、yおよびz軸に沿った）3次元（3D）加速度計およびジャイロスコープセンサーデータを含むことができる。（たとえば、加速度計とジャイロスコープとからの）センサー信号の各成分は、3秒のウィンドウ（たとえば、150の読取り/ウィンドウ）など、固定幅のウィンドウであり得る。例示的および非限定的な例として、次、前、アップ/増加、

50

ダウン/減少、オン、オフ、および不明など、複数のジェスチャークラスが実装され得る。置換、時間ワーピング、スケーリング、大きさワーピング、ジッタ、およびクロッピングなど、1つまたは複数のデータオーグメンテーション技法が実装され得る。

【0063】

[0076]手のジェスチャー認識のデータセットの統計的特徴などに基づく、特徴抽出が実施され得る。例示のために、抽出された特徴は、例示的および非限定的な例として、最小、最大、分散、平均、標準偏差、MSE（最小2乗誤差）、ACF（自己相関）、ACV（自己共分散）、ゆがみ、尖度、平均交差率、ジッタ、または3分位数に対応することができる。

【0064】

[0077]サポートベクターマシン（SVM）、勾配ブースティング、分類器、積層長短期記憶リカレントニューラルネットワーク（LSTM-RNN）、シーケンスツーシーケンスエンコーダデコーダモデルウィズアテンション、1つまたは複数の他のモデル、あるいはそれらの任意の組合せなど、1つまたは複数のモデルが手のジェスチャー認識のために使用され得る。

【0065】

[0078]いくつかの態様では、マルチモーダル認識エンジン130は、ターゲット行為に直接マッピングされたシーケンス埋め込みベクトルを生成することを学習またはトレーニングすることができる。入力シーケンスの例は、（たとえば、ジェスチャー入力のための）加速度計もしくはジャイロスコープ時系列、スピーチコマンド時系列、またはオーディオ時系列を含む。エンコーダデコーダLSTM-RNNウィズアテンションは、入力シーケンスに関連するターゲット行為クラスを示すためのソフトマックスレイヤへの出力を生成するためになど、可変長時系列信号を固定長および弁別ベクトルとして表す埋め込みベクトルを生成することを学習するために使用され得る。

【0066】

[0079]いくつかの態様では、マルチモーダル認識エンジン130は、異なる行為クラスの登録と設計とのために埋め込みベクトルを使用することができる。たとえば、いくつかの異なる入力シーケンスが登録され得、1つまたは複数の分類器は、各ターゲット行為にマッピングされた埋め込みベクトルを使用して設計され得る。たとえば、埋め込みをターゲット行為にマッピングするために、SVM、K平均、k近傍法（KNN）、コサイン（cos）距離、または他の設計が実装され得る。更新されたシステムの精度を検証するために、ユーザシーケンスのテストが実施され得る。

【0067】

[0080]いくつかの態様では、登録およびSVM/K平均/KNN設計の後に、分類器評価に関連するメトリックは、クラス間の分離があまりにあいまいであり、シーケンス整形が実施され得ることを示す。そのような場合、フィードバックメッセージ生成器132は、他のクラスとの混同を引き起こすいくつかの問題があるシーケンスをユーザに示すためのフィードバックを生成することができる。たとえば、混同されたクラスの動き、オーディオ、またはスピーチシーケンスは、出力デバイス120などを介して、ユーザに再生され得る。ユーザは、どのシーケンスが混同を引き起こすかを了解することができ、ターゲットクラス間の分離を改善しディスプレイを介して新しいシーケンスを発話する/ジェスチャーで示すことができる。代替的に、混同を招く入力シーケンスは、入力シーケンス間のあいまいさが未決定にレンダリングされるように、ユーザによって、マルチモーダルユーザインターフェースを介して、同じ行為/クラスに一致させられ得る。ユーザフィードバックを受信した後に、マルチモーダル認識エンジン130は、SVM/K平均/KNN設計を再登録および修正することができ、フィードバックメッセージ生成器132は、混同がある場合に、シーケンスマッピングが互いに十分に別個になるまで、入力シーケンスを繰り返すようにユーザに再プロンプトすることができる。たとえば、「混同行列」は、異なるシーケンス間のあいまいさの量を表すことができ、トレーニングは、混同行列が準対角になるまで繰り返され得る。

10

20

30

40

50

【 0 0 6 8 】

[0081]いくつかの態様では、他の行為クラスとの混同を引き起こす「問題がある」入力シーケンスを検出したことに応答して、マルチモーダル入力を用いたシーケンス整形が実施され得る。デバイス 1 1 0 は、ユーザが、問題がある入力シーケンスの各々のためにマルチモーダル入力を使用することを望むかどうかをユーザに要求することができる。たとえば、「オフ」および「ボリュームダウン」のためのユーザの特定のジェスチャーが、マルチモーダル認識エンジン 1 3 0 にとって区別するのが困難である場合、出力デバイス 1 2 0 は、「あなたは、『オフ』カテゴリーのためにバーバルコマンド『オフにする』を使用したいですか？」という問合せをユーザに出力し得る。別の例として、出力デバイス 1 2 0 は、「あなたは、『ボリュームダウン』カテゴリーのためにバーバルコマンド『ボリュームを下げる』を使用したいですか？」という問合せを出力し得る。ユーザが（たとえば、ジェスチャー混同によりバーバルコマンドを追加するために）マルチモーダル入力を使用することを選択したことに応答して、マルチモーダルキューがアクティブにされ得、デバイス 1 1 0 は、マルチモーダル入力シーケンスを使用した混同の確率を含めるように混同行列を調整することができる。

10

【 0 0 6 9 】

[0082]いくつかの態様では、入力シーケンスをディスアンビギュエートするために、対話型連続検証が使用され得る。たとえば、ユーザは、どのカテゴリーがどのマルチモーダル入力に登録されたかを忘れることがある。ユーザとデバイス 1 1 0 との間でダイアログベースの対話が行われ得る。たとえば、ジェスチャー入力が「オフ」カテゴリーとして検出された場合、出力デバイス 1 2 0 は、「あなたは、『オフ』カテゴリーまたは『次』カテゴリーを意図していますか？」をユーザに問い合わせ得る。ユーザは、「オフ」と答えることがあり、マルチモーダル認識エンジン 1 3 0 は、「オフ」コマンドをアクティブにし得る。

20

【 0 0 7 0 】

[0083]図 5 は、図 1 3 に関してさらに説明されるように、半導体チップまたはパッケージなどの個別構成要素に組み込まれたマルチモーダル認識エンジン 1 3 0 とフィードバックメッセージ生成器 1 3 2 とを含む、デバイス 5 0 2 の実装形態 5 0 0 を示す。例示のために、デバイス 5 0 2 は、マルチモーダル認識エンジン 1 3 0 とフィードバックメッセージ生成器 1 3 2 とに関して説明される動作を実施するために、記憶された命令を実行するように構成された 1 つまたは複数のプロセッサ（たとえば、プロセッサ 1 0 8 ）を含むことができる。デバイス 5 0 2 は、図 1 の入力デバイス 1 1 2 ~ 1 1 6 のうちの 1 つまたは複数からのデータなど、センサーデータ 5 0 4 がデバイス 5 0 2 の外部の 1 つまたは複数のセンサーから受信されることを可能にするために、第 1 のバスインターフェースなどのセンサーデータ入力 5 1 0 を含む。デバイス 5 0 2 はまた、（たとえば、出力デバイス 1 2 0 に）フィードバックメッセージ 1 4 4 を送ることを可能にするために、第 2 のバスインターフェースなどの出力 5 1 2 を含む。デバイス 5 0 2 は、図 1 1 に示されている車両、図 1 2 A に示されている仮想現実もしくは拡張現実ヘッドセット、図 1 2 B に示されているウェアラブル電子デバイス、または図 1 3 に示されているワイヤレス通信デバイスなどの中に、複数のセンサーと出力デバイスとを含むシステム中の構成要素として、マルチモーダルユーザインターフェース処理の実装を可能にする。

30

40

【 0 0 7 1 】

[0084]図 6 を参照すると、例示的および非限定的な例として、図 1 のデバイス 1 1 0 もしくは制御ユニット 1 0 4、図 5 のデバイス 5 0 2、またはそれらの両方によって実施され得る、マルチモーダルユーザ入力を処理する方法 6 0 0 の特定の实装形態が示されている。

【 0 0 7 2 】

[0085]方法 6 0 0 は、6 0 2 において、デバイスのプロセッサにおいて、第 1 の入力デバイスから受信された第 1 のデータを処理することを含む。第 1 のデータは、第 1 の入力モードに基づくユーザからの第 1 の入力を示す。たとえば、図 1 を参照すると、プロセッ

50

サ 1 0 8 は、第 1 の入力デバイス 1 1 2 から受信された第 1 のデータ 1 4 2 を処理する。第 1 のデータ 1 4 2 は、第 1 の入力モードに基づくユーザ 1 0 2 からの第 1 の入力 1 4 0 を示す。

【 0 0 7 3 】

[0086]方法 6 0 0 はまた、6 0 4 において、デバイスのプロセッサから、第 1 のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることを含む。フィードバックメッセージは、第 1 の入力モードとは異なる第 2 の入力モードに基づく第 2 の入力を提供するようにユーザに命令する。たとえば、図 1 を参照すると、制御ユニット 1 0 4 は、第 1 のデータ 1 4 2 を処理することに基づいて出力デバイス 1 2 0 にフィードバックメッセージ 1 4 4 を送る。フィードバックメッセージ 1 4 4 は、第 2 の入力モードに基づく第 2 の入力 1 4 8 を提供するようにユーザ 1 0 2 に命令する。

10

【 0 0 7 4 】

[0087]方法 6 0 0 は、フィールドプログラマブルゲートアレイ (F P G A) デバイス、特定用途向け集積回路 (A S I C)、中央処理ユニット (C P U) などの処理ユニット、デジタル信号プロセッサ (D S P)、コントローラ、別のハードウェアデバイス、ファームウェアデバイス、またはそれらの任意の組合せによって実装され得る。一例として、方法 6 0 0 は、本明細書で説明されるように、命令を実行するプロセッサによって実施され得る。

【 0 0 7 5 】

[0088]図 7 を参照すると、例示的および非限定的な例として、図 1 の制御ユニット 1 0 4、図 5 のデバイス 5 0 2、またはそれらの両方によって実施され得る、マルチモーダルユーザ入力を処理する方法 7 0 0 の特定の実装形態が示されている。

20

【 0 0 7 6 】

[0089]方法 7 0 0 は、7 0 2 において、第 1 の入力デバイスから受信された第 1 のデータを処理することを含む。第 1 のデータは、第 1 の入力モードに基づくユーザからのコマンドに対応する第 1 の入力を示す。たとえば、図 1 を参照すると、プロセッサ 1 0 8 は、第 1 の入力デバイス 1 1 2 から受信された第 1 のデータ 1 4 2 を処理する。第 1 のデータ 1 4 2 は、第 1 の入力モードに基づくユーザ 1 0 2 からのコマンドに対応する第 1 の入力 1 4 0 を示す。

【 0 0 7 7 】

[0090]方法 7 0 0 はまた、7 0 4 において、第 1 のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることを含む。フィードバックメッセージは、第 1 の入力をディスアンビグユエートする (disambiguate) ために、第 1 の入力モードとは異なる第 2 の入力モードに基づく第 2 の入力を提供するようにユーザに命令する。たとえば、図 1 を参照すると、制御ユニット 1 0 4 は、第 1 のデータ 1 4 2 を処理することに基づいて出力デバイス 1 2 0 にフィードバックメッセージ 1 4 4 を送る。フィードバックメッセージ 1 4 4 は、第 1 の入力 1 4 0 をディスアンビグユエートするために、第 1 の入力モードとは異なる第 2 の入力モードに基づく第 2 の入力 1 4 8 を提供するようにユーザ 1 0 2 に命令する。

30

【 0 0 7 8 】

[0091]方法 7 0 0 は、フィールドプログラマブルゲートアレイ (F P G A) デバイス、特定用途向け集積回路 (A S I C)、中央処理ユニット (C P U) などの処理ユニット、D S P、コントローラ、別のハードウェアデバイス、ファームウェアデバイス、またはそれらの任意の組合せによって実装され得る。一例として、方法 7 0 0 は、本明細書で説明されるように、命令を実行するプロセッサによって実施され得る。

40

【 0 0 7 9 】

[0092]図 8 を参照すると、例示的および非限定的な例として、図 1 の制御ユニット 1 0 4、図 5 のデバイス 5 0 2、またはそれらの両方によって実施され得る、マルチモーダルユーザ入力を処理する方法 8 0 0 の特定の実装形態が示されている。

【 0 0 8 0 】

50

[0093]方法 800 は、802 において、第 1 の入力デバイスから受信された第 1 のデータを処理することを含む。第 1 のデータは、第 1 の入力モードに基づくユーザからの第 1 の入力を示し、第 1 のデータは、ユーザの第 1 の認証行為 (authentication action) に対応する。たとえば、図 1 を参照すると、プロセッサ 108 は、第 1 の入力デバイス 112 から受信された第 1 のデータ 142 を処理する。第 1 のデータ 142 は、第 1 の入力モードに基づくユーザ 102 からの第 1 の入力 140 を示し、第 1 のデータ 142 は、ユーザ 102 の第 1 の認証行為に対応する。

【 0081 】

[0094]方法 800 はまた、804 において、第 1 のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることを含む。フィードバックメッセージは、多因子認証手順 (multi-factor authentication procedure) の一部として、第 2 の認証行為として、第 1 の入力モードとは異なる第 2 の入力モードに基づいて、第 2 の入力を提供するようにユーザに命令する。たとえば、図 1 を参照すると、制御ユニット 104 は、第 1 のデータ 142 を処理することに基づいて出力デバイス 120 にフィードバックメッセージ 144 を送る。フィードバックメッセージ 144 は、多因子認証手順の一部として、第 2 の認証行為として、異なるフォーム第 1 の入力モードである第 2 の入力モードに基づいて、第 2 の入力 148 を提供するようにユーザ 102 に命令する。

【 0082 】

[0095]方法 800 は、フィールドプログラマブルゲートアレイ (F P G A) デバイス、特定用途向け集積回路 (A S I C)、中央処理ユニット (C P U) などの処理ユニット、 D S P、コントローラ、別のハードウェアデバイス、ファームウェアデバイス、またはそれらの任意の組合せによって実装され得る。一例として、方法 800 は、本明細書で説明されるように、命令を実行するプロセッサによって実施され得る。

【 0083 】

[0096]図 9 を参照すると、例示的および非限定的な例として、図 1 の制御ユニット 104、図 5 のデバイス 502、またはそれらの両方によって実施され得る、マルチモーダルユーザ入力を処理する方法 900 の特定の実装形態が示されている。

【 0084 】

[0097]方法 900 は、902 において、第 1 の入力モードに基づく第 1 のユーザ入力を検出することを含む。たとえば、図 1 を参照すると、第 1 の入力デバイス 112 は、第 1 の入力モードに基づく第 1 のユーザ入力 140 を検出する。

【 0085 】

[0098]方法 900 はまた、904 において、第 2 の入力モードに基づく第 2 のユーザ入力を検出することを含む。たとえば、図 1 を参照すると、第 2 の入力デバイス 114 は、第 2 の入力モードに基づく第 2 のユーザ入力 148 を検出する。

【 0086 】

[0099]方法 900 はまた、906 において、第 1 のユーザ入力を第 1 の埋め込みベクトルにコンバートするように構成された第 1 の埋め込みネットワークを使用して、第 1 の埋め込みベクトルを生成することを含む。たとえば、図 2 を参照すると、第 1 の埋め込みネットワーク 202 は、第 1 のユーザ入力を第 1 の埋め込みベクトルにコンバートすることによって、第 1 の埋め込みベクトルを生成する。

【 0087 】

[0100]方法 900 はまた、908 において、第 2 のユーザ入力を第 2 の埋め込みベクトルにコンバートするように構成された第 2 の埋め込みネットワークを使用して、第 2 の埋め込みベクトルを生成することを含む。たとえば、図 2 を参照すると、第 2 の埋め込みネットワーク 204 は、第 2 のユーザ入力を第 2 の埋め込みベクトルにコンバートすることによって、第 2 の埋め込みベクトルを生成する。

【 0088 】

[0101]方法 900 はまた、910 において、組み合わされた埋め込みベクトルを生成するために、第 1 の埋め込みネットワークと第 2 の埋め込みネットワークとの出力を組み合

10

20

30

40

50

わせるように構成された融合埋め込みネットワークを使用して、組み合わせられた埋め込みベクトルを生成することを含む。たとえば、図2を参照すると、融合埋め込みネットワーク220は、組み合わせられた埋め込みベクトルを生成するために、第1の埋め込みネットワーク202と第2の埋め込みネットワーク204との出力を組み合わせる。

【0089】

[0102]方法900はまた、912において、分類器を使用して、組み合わせられた埋め込みベクトルを特定の行為にマッピングすることを含む。たとえば、図2を参照すると、マッピング230は、組み合わせられた埋め込みベクトルを特定の行為にマッピングする。

【0090】

[0103]方法900は、フィールドプログラマブルゲートアレイ(FPGA)デバイス、
特定用途向け集積回路(ASIC)、中央処理ユニット(CPU)などの処理ユニット、
DSP、コントローラ、別のハードウェアデバイス、ファームウェアデバイス、またはそれらの任意の組合せによって実装され得る。一例として、方法900は、本明細書で説明されるように、命令を実行するプロセッサによって実施され得る。

【0091】

[0104]図10を参照すると、例示的および非限定的な例として、図1の制御ユニット104、図5のデバイス502、またはそれらの両方によって実施され得る、マルチモジュールユーザ入力を処理する方法1000の特定の実装形態が示されている。

【0092】

[0105]方法1000は、1002において、第1の入力デバイスから受信された第1のデータを処理することを含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマンドに対応する。たとえば、図1を参照すると、プロセッサ108は、第1の入力デバイス112から受信された第1のデータ142を処理する。第1のデータ142は、第1の入力モードに基づくユーザ102からの第1の入力140を示す。

【0093】

[0106]方法1000はまた、1004において、第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることを含む。フィードバックメッセージは、第1の入力モードとは異なる第2の入力モードに基づいて、第1の入力に関連するコマンドを識別する第2の入力を提供するようにユーザに命令する。たとえば、図1を参照すると、制御ユニット104は、第1のデータ142を処理することに基づいて出力デバイス120にフィードバックメッセージ144を送る。フィードバックメッセージ144は、別の入力モードを使用して、第1の入力140に関連するコマンドを再び入力するようにユーザ102に命令する。一例では、第1の入力モードは、スピーチモード、ジェスチャーモード、またはビデオモードのうちの一つであり、第2の入力モードは、スピーチモード、ジェスチャーモード、またはビデオモードのうち異なる一つである。

【0094】

[0107]いくつかの実装形態では、フィードバックメッセージは、第1の入力をディスプレイギューエートするために第2の入力を提供するようにユーザに命令する。たとえば、フィードバックメッセージは、確信度レベル234が確信度しきい値294よりも小さいときなど、第1の入力の認識処理に関連する確信度レベルが確信度しきい値を満たすことに失敗したことに応答して送られ得る。いくつかの例では、第1の入力モードは、ビデオモードに対応し、フィードバックメッセージは、照明しきい値286よりも小さい値を有する周辺光メトリック284など、照明しきい値を下回る値を有する周辺光メトリックに応答して送られる。他の例では、第1の入力モードは、スピーチモードに対応し、フィードバックメッセージは、雑音しきい値282よりも大きい値を有する雑音メトリック280など、雑音しきい値を超える値を有する雑音メトリックに応答して送られる。

【0095】

[0108]方法1000はまた、1006において、第2の入力デバイスから第2のデータを受信することを含み、第2のデータは、第2の入力を示す。たとえば、図1を参照する

10

20

30

40

50

と、制御ユニット104は、第1の入力140に関連する特定のコマンドを識別する、第2の入力148に関連する第2のデータ150を受信する。

【0096】

[0109]方法1000はまた、1008において、第2の入力によって識別されるコマンドに第1の入力を関連付けるようにマッピングを更新することを含む。たとえば、図1を参照すると、制御ユニット104は、特定のコマンドに第1の入力140を関連付けるように、コマンドへのユーザ入力のマッピングを更新する。いくつかの実装形態では、更新されたマッピングは、コマンドが、第1の入力モードを介した第1の入力と第2の入力とモードを介した第2の入力とのコンカレントな（たとえば、少なくとも部分的に時間的に重複している）受信を介してより確実に認識されるように、ユーザのためにコマンドへの入力のマッピングをカスタマイズするためになど、第1の入力と第2の入力との組合せをコマンドに関連付ける。いくつかの実装形態では、マッピングを更新することは、ユーザに関連する埋め込みネットワークデータ（たとえば、第1の埋め込みネットワークデータ252）を更新すること、またはユーザに関連する重みデータ（たとえば、第1の重みデータ254）を更新することのうちの少なくとも1つを含む。

10

【0097】

[0110]方法1000は、フィールドプログラマブルゲートアレイ（FPGA）デバイス、特定用途向け集積回路（ASIC）、中央処理ユニット（CPU）などの処理ユニット、DSP、コントローラ、別のハードウェアデバイス、ファームウェアデバイス、またはそれらの任意の組合せによって実装され得る。一例として、方法1000は、本明細書で説明されるように、命令を実行するプロセッサによって実施され得る。

20

【0098】

[0111]図11は、車ダッシュボードデバイス1102などの車両ダッシュボードデバイスに組み込まれたデバイス110の実装形態1100の一例を示す。複数のセンサー1150は、1つまたは複数のマイクロフォン、カメラ、または他のセンサーを含むことができ、図1の入力デバイス112~116に対応することができる。単一のロケーションに示されているが、他の実装形態では、車両操作者からおよび各搭乗者からのマルチモーダル入力を検出するために車両中の各シートに近接して位置する1つまたは複数のマイクロフォンおよび1つまたは複数のカメラのアレイなど、センサー1150のうちの1つまたは複数は、車両のキャビン内の様々なロケーションに分散されるなど、車両の他のロケーションに配置され得る。

30

【0099】

[0112]ディスプレイ1120などの視覚的インターフェースデバイスは、出力デバイス120に対応することができ、車の運転者に見えるように車ダッシュボードデバイス1102内に取り付けられるかまたはその上に配置される（たとえば、車両ハンドセットマウントに着脱可能に固定される）。マルチモーダル認識エンジン130とフィードバックメッセージ生成器132とは、マルチモーダル認識エンジン130とフィードバックメッセージ生成器132とは、車両の乗員に見えないことを示すために、破線の境界で示されている。マルチモーダル認識エンジン130とフィードバックメッセージ生成器132とは、図1のデバイス110中のように、ディスプレイ1120およびセンサー1150をも含むデバイス中に実装され得るか、または図5のデバイス502中のように、ディスプレイ1120およびセンサー1150とは別個であり、それらに結合され得る。

40

【0100】

[0113]図12Aは、仮想現実、拡張現実、または複合現実ヘッドセットなど、ヘッドセット1202に組み込まれたマルチモーダル認識エンジン130とフィードバックメッセージ生成器132との一例を示す。ディスプレイ1220などの視覚的インターフェースデバイスは、出力デバイス120に対応することができ、ヘッドセット1202が着用されている間、ユーザへの拡張現実または仮想現実の画像またはシーンの表示を可能にするために、ユーザの目の前に配置される。センサー1250は、1つまたは複数のマイクロフォン、カメラ、または他のセンサーを含むことができ、図1の入力デバイス112~1

50

16に対応することができる。単一のロケーションに示されているが、他の実装形態では、マルチモーダル入力を検出するためにヘッドセット1202の周りに分散された1つまたは複数のマイクロフォンおよび1つまたは複数のカメラのアレイなど、センサー1250のうちの1つまたは複数は、ヘッドセット1202の他のロケーションに配置され得る。
【0101】

[0114]図12Bは、ディスプレイ1220とセンサー1250とを含む、「スマートウォッチ」として示されている、ウェアラブル電子デバイス1204に組み込まれたマルチモーダル認識エンジン130とフィードバックメッセージ生成器132との一例を示す。センサー1250は、たとえば、ビデオ、スピーチ、およびジェスチャーなどのモダリティに基づくユーザ入力の検出を可能にする。また、単一のロケーションに示されているが、他の実装形態では、センサー1250のうちの1つまたは複数は、ウェアラブル電子デバイス1204の他のロケーションに配置され得る。

10

【0102】

[0115]図13は、ワイヤレス通信デバイス実装形態（たとえば、スマートフォン）またはデジタルアシスタントデバイス実装形態などにおける、マルチモーダル認識エンジン130を含むデバイス1300の特定の例示的な実装形態のブロック図を示す。様々な実装形態では、デバイス1300は、図13に示されているものよりも多いまたは少ない構成要素を有し得る。例示的な実装形態では、デバイス1300は、デバイス110に対応し得る。例示的な実装形態では、デバイス1300は、図1～図12Bを参照しながら説明された1つまたは複数の動作を実施し得る。

20

【0103】

[0116]特定の実装形態では、デバイス1300は、マルチモーダル認識エンジン130を含むプロセッサ1306（たとえば、プロセッサ108に対応する中央処理ユニット（CPU））を含む。デバイス1300は、1つまたは複数の追加のプロセッサ1310（たとえば、1つまたは複数のDSP）を含み得る。プロセッサ1310は、スピーチおよび音楽コーデック（コーデック）1308を含み得る。スピーチおよび音楽コーデック1308は、音声コーデック（「ボコーデック」）エンコーダ1336、ボコーデックデコーダ1338、またはそれらの両方を含み得る。

【0104】

[0117]デバイス1300は、メモリ1386と、コーデック1334とを含み得る。メモリ1386は、メモリ106に対応し得、マルチモーダル認識エンジン130、フィードバックメッセージ生成器132、アプリケーション240のうちの1つまたは複数、あるいはそれらの任意の組合せに関して説明された機能を実装するためにプロセッサ1306（あるいは1つまたは複数の追加のプロセッサ1310）によって実行可能である命令1356を含み得る。デバイス1300は、トランシーバ1350を介して1つまたは複数のアンテナ1352に結合されたワイヤレスコントローラ1340を含み得る。いくつかの実装形態では、1つまたは複数のアンテナ1352は、ジェスチャー入力を示すデータを受信するように構成された1つまたは複数のアンテナを含む。

30

【0105】

[0118]デバイス1300は、ディスプレイコントローラ1326に結合されたディスプレイ1328（たとえば、出力デバイス120）を含み得る。ディスプレイ1328は、フィードバックメッセージ144（たとえば、命令146）を出力するグラフィカルユーザインターフェースを表現するように構成され得る。コーデック1334は、デジタルアナログコンバータ（DAC）1302と、アナログデジタルコンバータ（ADC）1304とを含み得る。特定の実装形態では、コーデック1334は、1つまたは複数のマイクロフォン1312（たとえば、1つまたは複数のキーワードまたは音声コマンドを含むオーディオ入力をキャプチャするように構成された第1の入力デバイス112）からアナログ信号を受信し、アナログデジタル変換器1304を使用してアナログ信号をデジタル信号にコンバートし、デジタル信号をスピーチおよび音楽コーデック1308に提供し得る。スピーチおよび音楽コーデック1308は、デジタル信号を処理し得る。

40

50

【0106】

[0119]特定の実装形態では、スピーチおよび音楽コーデック1308は、オーディオ再生信号を表すデジタル信号をコーデック1334に提供し得る。コーデック1334は、デジタルアナログコンバータ1302を使用してデジタル信号をアナログ信号にコンバートし得、可聴信号を生成するために、アナログ信号を1つまたは複数のラウドスピーカー1314に提供し得る。1つまたは複数のラウドスピーカー1314は、出力デバイス120に対応することができ、図1のフィードバックメッセージ144をレンダリングするか、またはフィードバックメッセージ144をユーザにダイレクトするように構成され得る。

【0107】

[0120]特定の実装形態では、デバイス1300は、1つまたは複数の入力デバイス1330を含む。入力デバイス1330は、図1の入力デバイス112~116のうちの1つまたは複数に対応することができる。たとえば、入力デバイス1330は、1つまたは複数のジェスチャーまたは視覚的コマンドを含むビデオ入力をキャプチャするように構成された1つまたは複数のカメラを含むことができる。

【0108】

[0121]特定の実装形態では、デバイス1300は、システムインパッケージまたはシステムオンチップデバイス1322中に含まれ得る。特定の実装形態では、メモリ1386と、プロセッサ1306と、プロセッサ1310と、ディスプレイコントローラ1326と、コーデック1334と、ワイヤレスコントローラ1340とは、システムインパッケージまたはシステムオンチップデバイス1322中に含まれる。特定の実装形態では、入力デバイス1330（たとえば、図1の入力デバイス112~116のうちの1つまたは複数の）と、電源1344とは、システムインパッケージまたはシステムオンチップデバイス1322に結合される。その上、特定の実装形態では、図13に示されているように、ディスプレイ1328と、入力デバイス1330と、マイクロフォン1312と、アンテナ1352と、電源1344とは、システムインパッケージまたはシステムオンチップデバイス1322の外部にある。特定の実装形態では、ディスプレイ1328と、入力デバイス1330と、マイクロフォン1312と、ラウドスピーカー1314と、アンテナ1352と、電源1344との各々は、インターフェースまたはコントローラなど、システムインパッケージまたはシステムオンチップデバイス1322の構成要素に結合され得る。

【0109】

[0122]デバイス1300は、例示的および非限定的な例として、モバイル通信デバイス、スマートフォン、セルラーフォン、ラップトップコンピュータ、コンピュータ、タブレット、携帯情報端末、ディスプレイデバイス、テレビジョン、ゲーミングコンソール、音楽プレーヤ、ラジオ、デジタルビデオプレーヤ、デジタルビデオディスク(DVD)またはBlu-ray(登録商標)ディスクプレーヤ、チューナー、カメラ、ナビゲーションデバイス、仮想現実または拡張現実ヘッドセット、ウェアラブル電子デバイス、車両コンソールデバイス、あるいはそれらの任意の組合せを含み得る。

【0110】

[0123]説明される実装形態に関連して、マルチモーダルユーザ入力のためのデバイスは、第1の入力デバイスから受信された第1のデータを処理するマルチモーダル認識エンジンを含む。第1のデータは、第1の入力モード（たとえば、スピーチモード、ジェスチャーモード、またはビデオモード）に基づくユーザからの第1の入力を示す。フィードバックメッセージ生成器は、第1のデータを処理することに基づいて、第1の入力モードとは異なる第2の入力モードに基づく第2の入力を提供するようにユーザに命令するフィードバックメッセージを出力デバイスに送る。

【0111】

[0124]説明される実装形態に関連して、マルチモーダルユーザ入力のためのデバイスは、第1の入力デバイスから受信された第1のデータを処理するマルチモーダル認識エンジ

10

20

30

40

50

ンを含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示す。フィードバックメッセージ生成器は、第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送る。フィードバックメッセージは、第1の入力に関連付けられるべき行為を識別するようにユーザに命令する。マルチモーダル認識エンジンは、第1の入力に関連付けられるべき特定の行為を識別する第2の入力を受信し、特定の行為に第1の入力を関連付けるように行為へのユーザ入力のマッピングを更新する。

【0112】

[0125]説明される実装形態に関連して、マルチモーダルユーザ入力のための装置は、第1の入力デバイスから受信された第1のデータを処理するための手段を含む。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマンドに対応する。たとえば、第1のデータを処理するための手段は、プロセッサ108、マルチモーダル認識エンジン130、プロセッサ1306、1310によって実行可能な命令1356、1つまたは複数の他のデバイス、モジュール、回路、構成要素、あるいはそれらの組合せを含むことができる。

10

【0113】

[0126]本装置は、第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送るための手段を含む。フィードバックメッセージは、第1の入力モードとは異なる第2の入力モードに基づいて、第1の入力に関連するコマンドを識別する第2の入力を提供するようにユーザに命令する。たとえば、送るための手段は、プロセッサ108、マルチモーダル認識エンジン130、フィードバックメッセージ生成器132、プロセッサ1306、1310によって実行可能な命令1356、1つまたは複数の他のデバイス、モジュール、回路、構成要素、あるいはそれらの組合せを含むことができる。

20

【0114】

[0127]本装置は、第2の入力デバイスから第2のデータを受信するための手段を含み、第2のデータは、第2の入力を示す。たとえば、第2のデータを受信するための手段は、プロセッサ108、マルチモーダル認識エンジン130、プロセッサ1306、1310によって実行可能な命令1356、1つまたは複数の他のデバイス、モジュール、回路、構成要素、あるいはそれらの組合せを含むことができる。

【0115】

[0128]本装置はまた、第2の入力によって識別されるコマンドに第1の入力を関連付けるようにマッピングを更新するための手段を含む。たとえば、更新するための手段は、プロセッサ108、マルチモーダル認識エンジン130、データ調整器292、プロセッサ1306、1310によって実行可能な命令1356、1つまたは複数の他のデバイス、モジュール、回路、構成要素、あるいはそれらの組合せを含むことができる。

30

【0116】

[0129]いくつかの実装形態では、非一時的コンピュータ可読媒体（たとえば、メモリ106、メモリ1386、またはそれらの任意の組合せ）は、デバイスの1つまたは複数のプロセッサ（たとえば、プロセッサ108、プロセッサ1306、プロセッサ1310、またはそれらの任意の組合せ）によって実行されたとき、図6～図10の方法のうちの1つまたは複数の全部または一部に対応する動作を実施することなどによって、マルチモーダルユーザ入力を処理するための動作を1つまたは複数のプロセッサに実施させる命令（たとえば、命令1356）を含む。一例では、命令は、1つまたは複数のプロセッサによって実行されたとき、第1の入力デバイスから受信された第1のデータを1つまたは複数のプロセッサに処理させる。第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、第1の入力は、コマンドに対応する。命令は、1つまたは複数のプロセッサによって実行されたとき、1つまたは複数のプロセッサに、第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送らせる。フィードバックメッセージは、第1の入力モードとは異なる第2の入力モードに基づいて、第1の入力に関連するコマンドを識別する第2の入力を提供するようにユーザに命令する。命令は、1つまたは複数のプロセッサによって実行されたとき、1つまたは複数のプロセッサに、第2

40

50

の入力デバイスから第2のデータを受信させ、第2のデータは、第2の入力を示す。命令はまた、1つまたは複数のプロセッサによって実行されたとき、1つまたは複数のプロセッサに、第2の入力によって識別されるコマンドに第1の入力を関連付けるようにマッピングを更新させる。

【0117】

[0130]さらに、本明細書で開示される実装形態に関して説明される様々な例示的な論理ブロック、構成、モジュール、回路、およびアルゴリズムステップは、電子ハードウェア、プロセッサによって実行されるコンピュータソフトウェア、または両方の組合せとして実装され得ることを当業者は諒解されよう。様々な例示的な構成要素、ブロック、構成、モジュール、回路、およびステップについて、上記では概して、それらの機能に関して説明された。そのような機能がハードウェアとして実装されるか、プロセッサ実行可能命令として実装されるかは、特定の適用例および全体的なシステムに課された設計制約に依存する。当業者は、説明された機能を、特定の適用例ごとに様々な方法で実装し得、そのような実装の決定は、本開示の範囲からの逸脱を引き起こすと解釈されるべきではない。

【0118】

[0131]本明細書で開示される実装形態に関して説明された方法またはアルゴリズムのステップは、ハードウェアで直接実施されるか、プロセッサによって実行されるソフトウェアモジュールで実施されるか、またはその2つの組合せで実施され得る。ソフトウェアモジュールは、ランダムアクセスメモリ(RAM)、フラッシュメモリ、読取り専用メモリ(ROM)、プログラマブル読取り専用メモリ(PROM)、消去可能プログラマブル読取り専用メモリ(EPROM)、電気的消去可能プログラマブル読取り専用メモリ(EEPROM(登録商標))、レジスタ、ハードディスク、リムーバブルディスク、コンパクトディスク読取り専用メモリ(CD-ROM)、または当技術分野で知られている任意の他の形態の非一時的記憶媒体中に常駐し得る。例示的な記憶媒体は、プロセッサが記憶媒体から情報を読み取り、記憶媒体に情報を書き込むことができるように、プロセッサに結合される。代替として、記憶媒体はプロセッサと一体であり得る。プロセッサと記憶媒体とは、特定用途向け集積回路(ASIC)中に存在し得る。ASICは、コンピューティングデバイスまたはユーザ端末中に存在し得る。代替として、プロセッサと記憶媒体とは、コンピューティングデバイスまたはユーザ端末中に個別構成要素として存在し得る。

【0119】

[0132]開示される実装形態の前の説明は、開示される実装形態を当業者が製作または使用することを可能にするために提供される。これらの実装形態への様々な変更は当業者には容易に明らかになり、本明細書で定義された原理は本開示の範囲から逸脱することなく他の実装形態に適用され得る。したがって、本開示は、本明細書に示された実装形態に限定されるものではなく、以下の特許請求の範囲によって定義されるような原理および新規の特徴に一致する可能な最も広い範囲を与えられるべきである。

以下に本願の出願当初の特許請求の範囲に記載された発明を付記する。

[C1]

マルチモーダルユーザ入力のためのデバイスであって、

第1の入力デバイスから受信された第1のデータを処理することと、前記第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、前記第1の入力は、コマンドに対応し、

前記第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることと、ここにおいて、前記フィードバックメッセージは、前記第1の入力モードとは異なる第2の入力モードに基づいて、前記第1の入力に関連するコマンドを識別する第2の入力を提供するように前記ユーザに命令する、

第2の入力デバイスから第2のデータを受信することと、前記第2のデータは、前記第2の入力を示し、

前記第2の入力によって識別される前記コマンドに前記第1の入力を関連付けるようにマッピングを更新することと、

10

20

30

40

50

を行うように構成された1つまたは複数のプロセッサを備える、デバイス。

[C 2]

前記第1の入力モードは、スピーチモード、ジェスチャーモード、またはビデオモードのうちの一つであり、前記第2の入力モードは、前記スピーチモード、前記ジェスチャーモード、または前記ビデオモードのうち異なる1つである、C1に記載のデバイス。

[C 3]

前記フィードバックメッセージは、前記第1の入力をディスアンビギュエートするために前記第2の入力を提供するように前記ユーザに命令する、C1に記載のデバイス。

[C 4]

前記1つまたは複数のプロセッサは、前記第1の入力の認識処理に関連する確信度レベルが確信度しきい値を満たすことに失敗したことに応答して、前記フィードバックメッセージを送るようにさらに構成された、C3に記載のデバイス。

10

[C 5]

前記更新されたマッピングは、前記第1の入力と前記第2の入力との組合せを前記コマンドに関連付ける、C1に記載のデバイス。

[C 6]

前記1つまたは複数のプロセッサは、マルチモーダル認識エンジンを含み、前記マルチモーダル認識エンジンは、

組み合わせられた埋め込みベクトルを生成するために、前記第1の入力モードに関連する第1の埋め込みネットワークと、前記第2の入力モードに関連する第2の埋め込みネットワークとの出力を組み合わせるように構成された融合埋め込みネットワークと、

20

前記組み合わせられた埋め込みベクトルを特定のコマンドにマッピングするように構成された分類器と、

を含む、C1に記載のデバイス。

[C 7]

前記ユーザに対応する第1の埋め込みネットワークデータおよび第1の重みデータと、第2のユーザに対応する第2の埋め込みネットワークデータおよび第2の重みデータと、前記第1の埋め込みネットワークデータは、前記ユーザと前記第2のユーザとの間の入力コマンドの差に基づいて前記第2の埋め込みネットワークデータとは異なり、前記第1の重みデータは、前記ユーザと前記第2のユーザとの間の入力モード信頼性の差に基づいて前記第2の重みデータとは異なり、

30

を記憶するように構成されたメモリをさらに備える、C6に記載のデバイス。

[C 8]

前記第1の入力モードは、ビデオモードに対応し、前記1つまたは複数のプロセッサは、照明しきい値を下回る値を有する周辺光メトリックに応答して前記フィードバックメッセージを送るように構成された、C1に記載のデバイス。

[C 9]

前記第1の入力モードは、スピーチモードに対応し、前記1つまたは複数のプロセッサは、雑音しきい値を超える値を有する雑音メトリックに応答して前記フィードバックメッセージを送るように構成された、C1に記載のデバイス。

40

[C 1 0]

グラフィカルユーザインターフェースを表すように構成されたディスプレイをさらに備える、C1に記載のデバイス。

[C 1 1]

1つまたは複数のキーワードまたは音声コマンドを含むオーディオ入力をキャプチャするように構成された1つまたは複数のマイクロフォンをさらに備える、C1に記載のデバイス。

[C 1 2]

1つまたは複数のジェスチャーまたは視覚的コマンドを含むビデオ入力をキャプチャするように構成された1つまたは複数のカメラをさらに備える、C1に記載のデバイス。

50

[C 1 3]

ジェスチャー入力を示すデータを受信するように構成された1つまたは複数のアンテナをさらに備える、C 1に記載のデバイス。

[C 1 4]

前記フィードバックメッセージをレンダリングするかまたは前記ユーザにダイレクトするように構成された1つまたは複数のラウドスピーカーをさらに備える、C 1に記載のデバイス。

[C 1 5]

前記ユーザは、ロボットまたは他の電子デバイスを含む、C 1に記載のデバイス。

[C 1 6]

前記第1の入力デバイスと前記出力デバイスとは、仮想現実ヘッドセットまたは拡張現実ヘッドセットに組み込まれる、C 1に記載のデバイス。

[C 1 7]

前記第1の入力デバイスと前記出力デバイスとは、車両に組み込まれる、C 1に記載のデバイス。

[C 1 8]

マルチモーダルユーザ入力のための方法であって、

デバイスの1つまたは複数のプロセッサにおいて、第1の入力デバイスから受信された第1のデータを処理することと、前記第1のデータは、第1の入力モードに基づくユーザからの第1の入力を示し、前記第1の入力は、コマンドに対応し、

前記1つまたは複数のプロセッサから、前記第1のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることと、ここにおいて、前記フィードバックメッセージは、前記第1の入力モードとは異なる第2の入力モードに基づいて、前記第1の入力に関連するコマンドを識別する第2の入力を提供するように前記ユーザに命令する、

前記1つまたは複数のプロセッサにおいて、第2の入力デバイスから第2のデータを受信することと、前記第2のデータは、前記第2の入力を示し、

前記1つまたは複数のプロセッサにおいて、前記第2の入力によって識別される前記コマンドに前記第1の入力を関連付けるようにマッピングを更新することと、

を備える、方法。

[C 1 9]

前記第1の入力モードは、スピーチモード、ジェスチャーモード、またはビデオモードのうちの1つであり、前記第2の入力モードは、前記スピーチモード、前記ジェスチャーモード、または前記ビデオモードのうちの異なる1つである、C 1 8に記載の方法。

[C 2 0]

前記フィードバックメッセージは、前記第1の入力をディスアンビギュエートするために前記第2の入力を提供するように前記ユーザに命令する、C 1 8に記載の方法。

[C 2 1]

前記フィードバックメッセージは、前記第1の入力の認識処理に関連する確信度レベルが確信度しきい値を満たすことに失敗したことに応答して送られる、C 2 0に記載の方法。

[C 2 2]

前記更新されたマッピングは、前記第1の入力と前記第2の入力との組合せを前記コマンドに関連付ける、C 1 8に記載の方法。

[C 2 3]

前記マッピングを更新することは、

前記ユーザに関連する埋め込みネットワークデータを更新すること、または

前記ユーザに関連する重みデータを更新すること、

のうちの少なくとも1つを含む、C 1 8に記載の方法。

[C 2 4]

前記第1の入力モードは、ビデオモードに対応し、前記フィードバックメッセージは、

10

20

30

40

50

照明しきい値を下回る値を有する周辺光メトリックに応答して送られる、C 1 8に記載の方法。

[C 2 5]

前記第 1 の入力モードは、スピーチモードに対応し、前記フィードバックメッセージは、雑音しきい値を超える値を有する雑音メトリックに応答して送られる、C 1 8に記載の方法。

[C 2 6]

マルチモーダルユーザ入力のための装置であって、

第 1 の入力デバイスから受信された第 1 のデータを処理するための手段と、前記第 1 のデータは、第 1 の入力モードに基づくユーザからの第 1 の入力を示し、前記第 1 の入力は、コマンドに対応し、

前記第 1 のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送るための手段と、ここにおいて、前記フィードバックメッセージは、前記第 1 の入力モードとは異なる第 2 の入力モードに基づいて、前記第 1 の入力に関連するコマンドを識別する第 2 の入力を提供するように前記ユーザに命令する、

第 2 の入力デバイスから第 2 のデータを受信するための手段と、前記第 2 のデータは、前記第 2 の入力を示し、

前記第 2 の入力によって識別される前記コマンドに前記第 1 の入力に関連付けるようにマッピングを更新するための手段と、

を備える、装置。

[C 2 7]

前記更新されたマッピングは、前記第 1 の入力と前記第 2 の入力との組合せを前記コマンドに関連付ける、C 2 6 に記載の装置。

[C 2 8]

デバイスの 1 つまたは複数のプロセッサによって実行されたとき、前記 1 つまたは複数のプロセッサに、

第 1 の入力デバイスから受信された第 1 のデータを処理することと、前記第 1 のデータは、第 1 の入力モードに基づくユーザからの第 1 の入力を示し、前記第 1 の入力が、コマンドに対応し、

前記第 1 のデータを処理することに基づいて出力デバイスにフィードバックメッセージを送ることと、ここにおいて、前記フィードバックメッセージは、前記第 1 の入力モードとは異なる第 2 の入力モードに基づいて、前記第 1 の入力に関連するコマンドを識別する第 2 の入力を提供するように前記ユーザに命令する、

第 2 の入力デバイスから第 2 のデータを受信することと、前記第 2 のデータは、前記第 2 の入力を示し、

前記第 2 の入力によって識別される前記コマンドに前記第 1 の入力に関連付けるようにマッピングを更新することと、

を行わせる命令を備える非一時的コンピュータ可読媒体。

[C 2 9]

前記第 1 の入力モードは、ビデオモードに対応し、前記フィードバックメッセージは、照明しきい値を下回る値を有する周辺光メトリックに応答して送られる、C 2 8 に記載の非一時的コンピュータ可読媒体。

[C 3 0]

前記第 1 の入力モードは、スピーチモードに対応し、前記フィードバックメッセージは、雑音しきい値を超える値を有する雑音メトリックに応答して送られる、C 2 8 に記載の非一時的コンピュータ可読媒体。

10

20

30

40

50

【図面】

【図 1】

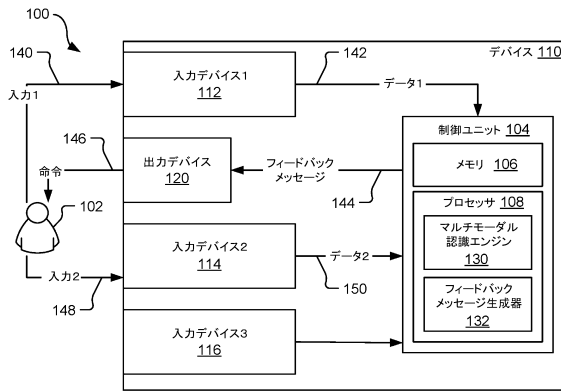


FIG. 1

【図 2】

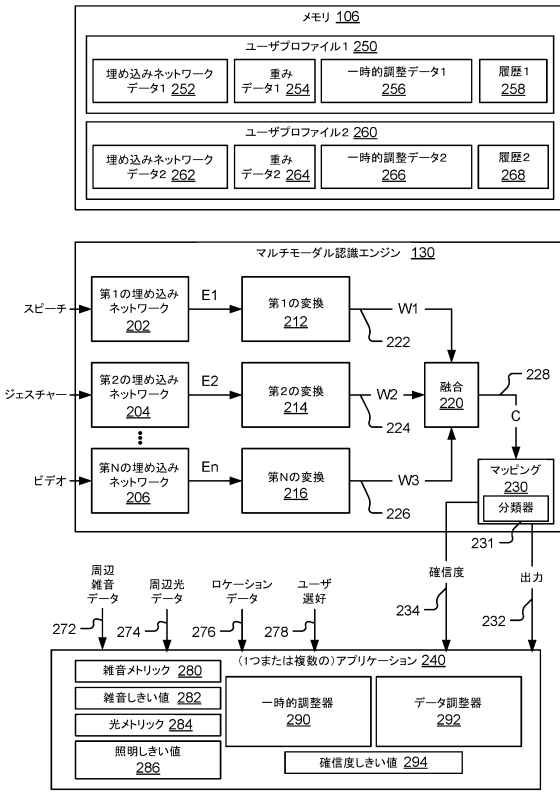


FIG. 2

【図 3】

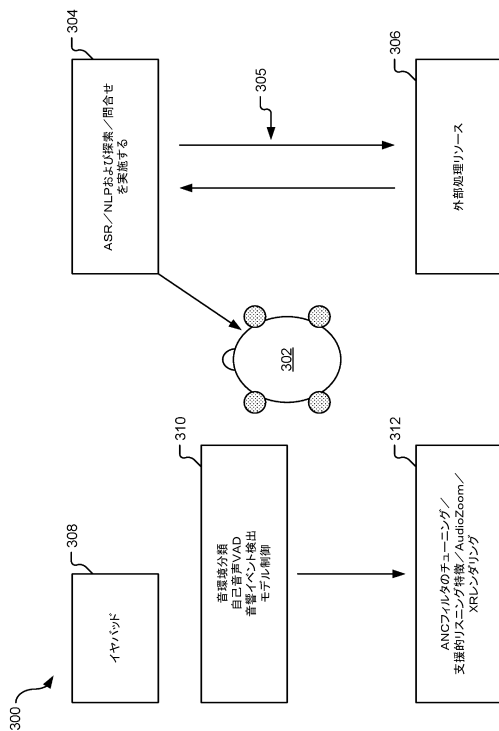


FIG. 3

【図 4】

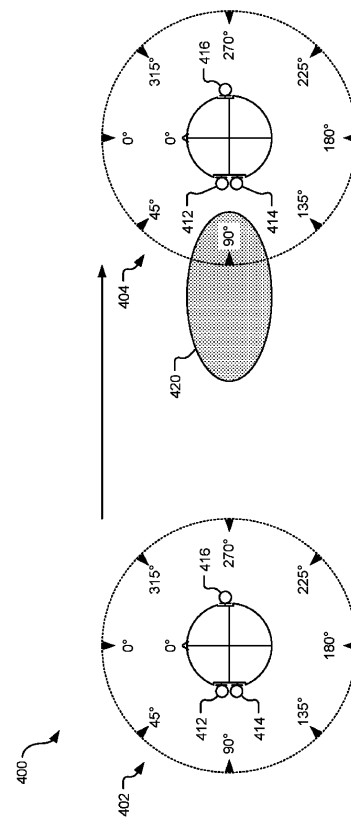


FIG. 4

10

20

30

40

50

【 図 5 】

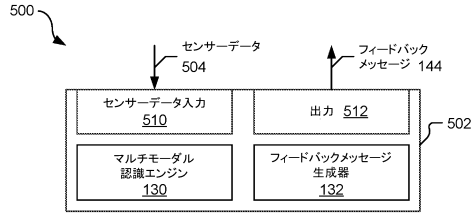


FIG. 5

【 図 6 】

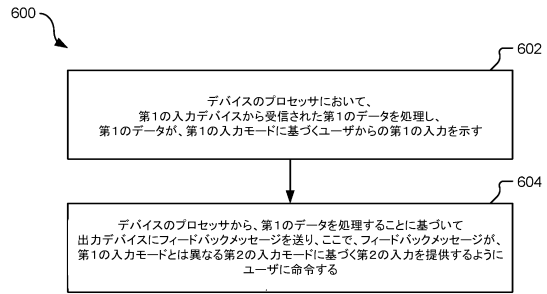


FIG. 6

10

【 図 7 】

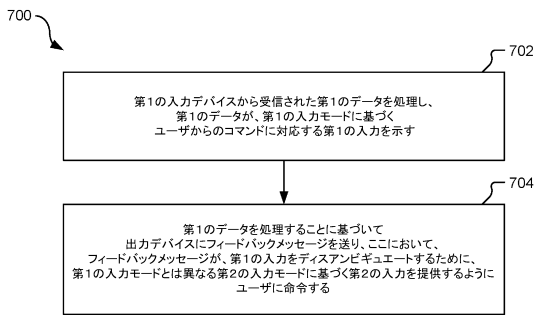


FIG. 7

【 図 8 】

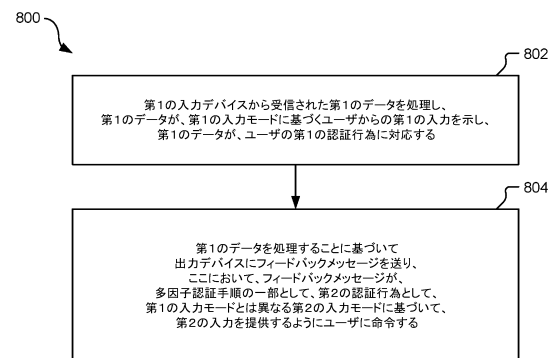


FIG. 8

20

30

40

50

【図 9】

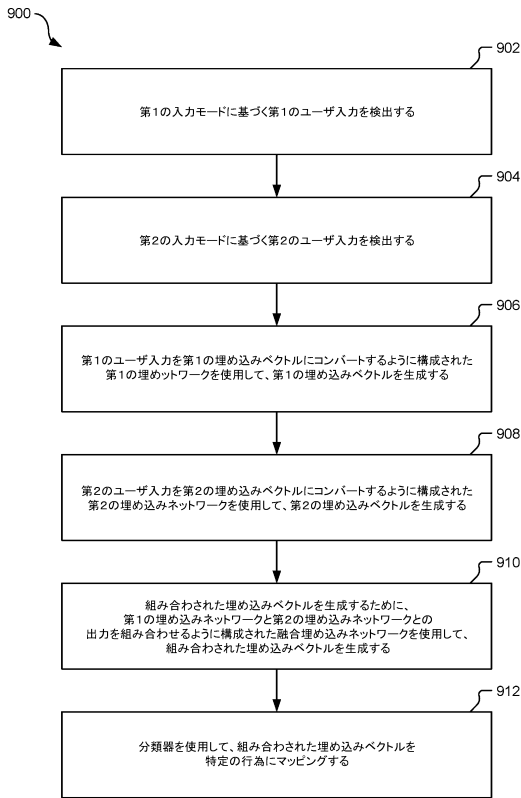


FIG. 9

【図 10】

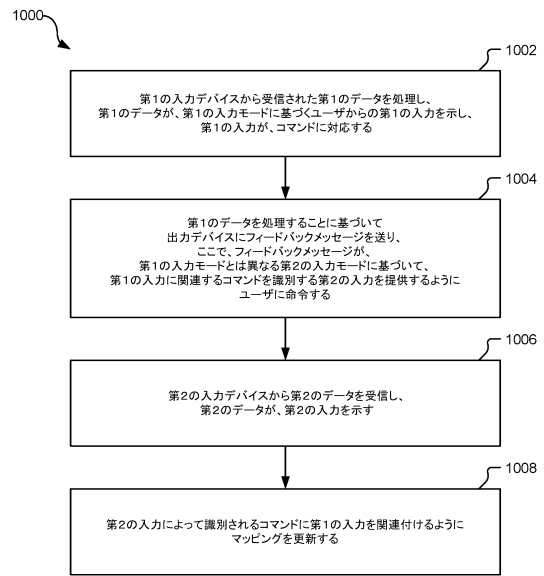


FIG. 10

【図 11】

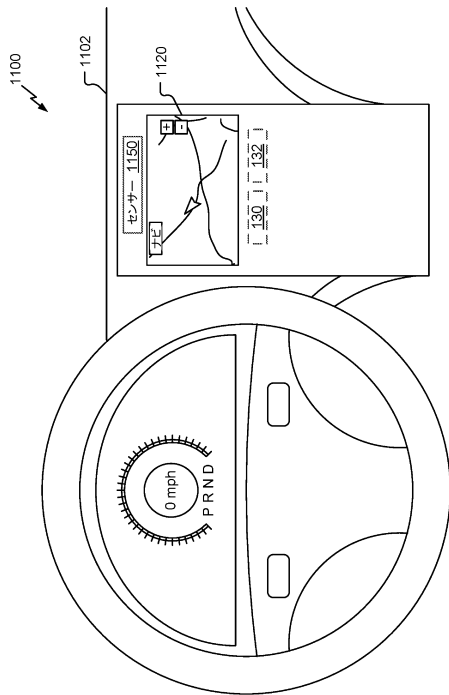


FIG. 11

【図 12 A】

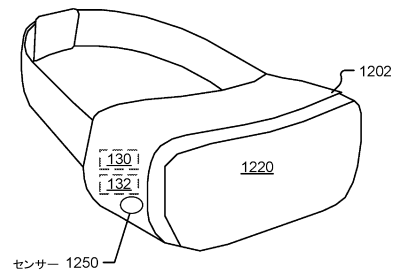


FIG. 12A

10

20

30

40

50

フロントページの続き

(33)優先権主張国・地域又は機関

米国(US)

(72)発明者 チョウドハリー、ラビ

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 キム、レ・フン

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ムン、ソクク

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 グオ、インイー

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 サキ、ファテメ

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ビッサー、エリック

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

審査官 富永 昌彦

(56)参考文献 米国特許出願公開第 2 0 1 8 / 0 3 2 9 6 7 7 (US , A 1)

国際公開第 2 0 1 4 / 0 7 0 8 7 2 (WO , A 2)

国際公開第 2 0 1 9 / 0 2 6 6 1 7 (WO , A 1)

特開 2 0 1 8 - 0 3 6 9 0 2 (JP , A)

米国特許出願公開第 2 0 2 0 / 0 1 5 2 1 9 1 (US , A 1)

中国特許出願公開第 1 1 0 9 9 8 7 1 8 (CN , A)

(58)調査した分野 (Int.Cl. , D B 名)

G 0 6 F 3 / 0 1

G 0 6 F 3 / 0 4 8 - 0 4 8 9 5

G 0 6 F 3 / 1 6