

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7432556号
(P7432556)

(45)発行日 令和6年2月16日(2024.2.16)

(24)登録日 令和6年2月7日(2024.2.7)

(51)国際特許分類	F I
G 0 6 T 13/40 (2011.01)	G 0 6 T 13/40
G 1 0 L 15/00 (2013.01)	G 1 0 L 15/00 2 0 0 Z
G 1 0 L 13/00 (2006.01)	G 1 0 L 13/00 1 0 0 M
G 0 6 F 3/16 (2006.01)	G 0 6 F 3/16 6 5 0
G 0 6 F 3/01 (2006.01)	G 0 6 F 3/16 6 9 0
請求項の数 17 外国語出願 (全22頁) 最終頁に続く	

(21)出願番号	特願2021-87333(P2021-87333)	(73)特許権者	514322098
(22)出願日	令和3年5月25日(2021.5.25)		ベイジン バイドゥ ネットコム サイエ ンス テクノロジー カンパニー リミテ ッド
(65)公開番号	特開2021-168139(P2021-168139 A)		Beijing Baidu Netco m Science Technolog y Co., Ltd.
(43)公開日	令和3年10月21日(2021.10.21)		中華人民共和国 ベキン 100085, ハイディアン ディストリクト, シャン ディ テンス ストリート, 10番, バ イドゥ キャンパス 2階
審査請求日	令和3年5月25日(2021.5.25)		2/F Baidu Campus, N o.10, Shangdi 10th Street, Haidian Dis trict, Beijing 1000
(31)優先権主張番号	202011598915.9		最終頁に続く
(32)優先日	令和2年12月30日(2020.12.30)		
(33)優先権主張国・地域又は機関	中国(CN)		

(54)【発明の名称】 マンマシンインタラクションのための方法、装置、機器および媒体

(57)【特許請求の範囲】

【請求項1】

受信した音声信号に基づいて、前記音声信号に対する回答の回答テキストを生成することと、

音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットを含む前記回答テキストに対応する回答音声信号を生成し、生成した前記回答音声信号は前記1セットのテキストユニットに対応する1セットの音声信号ユニットを含むことと、

前記回答テキストに基づいて、仮想オブジェクトによって表現される表情および/または動作の標識を確定することと、

前記回答音声信号、前記表情および/または動作の標識に基づいて、前記仮想オブジェクトを含む出力ビデオを生成し、前記出力ビデオは前記回答音声信号に基づいて確定された、前記仮想オブジェクトによって表現される唇形シーケンスを含むこととを含み、

前記回答音声信号を生成することは、

前記回答テキストを1セットのテキストユニットに分割することと、

音声信号ユニットとテキストユニットとのマッピング関係に基づいて、前記1セットのテキストユニットにおけるテキストユニットに対応する音声信号ユニットを取得することとであって、前記1セットのテキストユニットから前記テキストユニットを選択することと、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、音声ライブラリから前記テキストユニットに対応する前記音声信号ユニットを検索することとを含む、音

声信号ユニットを取得することと、

前記取得された音声信号ユニットに基づいて、前記回答音声信号を生成することとを含み、

前記音声ライブラリには音声信号ユニットとテキストユニットとの前記マッピング関係が記憶され、前記音声ライブラリにおける音声信号ユニットは、取得した、前記仮想オブジェクトに関する音声記録データを分割することで得られるものであり、前記音声ライブラリにおけるテキストユニットは、分割で得られた音声信号ユニットに基づいて確定されるものであり、

前記出力ビデオを生成することは、

前記回答音声信号を1セットの音声信号ユニットに分割することと、

10

前記1セットの音声信号ユニットに対応する前記仮想オブジェクトの唇形シーケンスを取得することと、

対応する前記表情および/または動作の標識に基づいて、前記仮想オブジェクトについての当該表情および/または動作のビデオセグメントを取得することと、

前記唇形シーケンスを前記ビデオセグメントに結合して前記出力ビデオを生成することとを含み、

前記仮想オブジェクトについての当該表情および/または動作の前記ビデオセグメントを取得することは、

表情および/または動作の標識とビデオセグメントとの間の事前に記憶されたマッピング関係を利用して、対応する前記表情および/または動作の標識に基づいて当該表情および/または動作の前記ビデオセグメントを取得することを含む、

20

マンマシンインタラクションのための方法。

【請求項2】

前記回答テキストを生成することは、

前記受信した音声信号を識別して入力テキストを生成することと、

前記入力テキストに基づいて、前記回答テキストを取得することとを含む、請求項1に記載の方法。

【請求項3】

前記入力テキストに基づいて、前記回答テキストを取得することは、

入力テキストと前記仮想オブジェクトの人格属性を用いて回答テキストを生成する機械学習モデルである対話モデルに、前記入力テキストと前記仮想オブジェクトの人格属性を入力して前記回答テキストを取得することを含む、請求項2に記載の方法。

30

【請求項4】

前記対話モデルは、前記仮想オブジェクトの人格属性と、入力テキストサンプルと回答テキストサンプルを含む対話サンプルトとを利用してトレーニングすることで得られるものである、請求項3に記載の方法。

【請求項5】

前記表情および/または動作の標識を確定することは、

テキストを用いて表情および/または動作の標識を確定する機械学習モデルである表情および動作識別モデルに、前記回答テキストを入力して、前記表情および/または動作の標識を取得することを含む、請求項1に記載の方法。

40

【請求項6】

前記唇形シーケンスを前記ビデオセグメントに結合して前記出力ビデオを生成することは、

前記ビデオセグメントにおける時間軸での所定の時間位置におけるビデオフレームを確定することと、

前記唇形シーケンスから前記所定の時間位置に対応する唇形を取得することと、

前記唇形を前記ビデオフレームに結合して前記出力ビデオを生成することとを含む、請求項1に記載の方法。

【請求項7】

50

前記回答音声信号と前記出力ビデオとを関連付けて出力することとをさらに含む、請求項 1 に記載の方法。

【請求項 8】

受信した音声信号に基づいて、前記音声信号に対する回答の回答テキストを生成するように構成される回答テキスト生成モジュールと、

音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1 セットのテキストユニットを含む前記回答テキストに対応する回答音声信号を生成し、生成された前記回答音声信号は前記 1 セットのテキストユニットに対応する 1 セットの音声ユニットを含むように構成される第 1 回答音声信号生成モジュールと、

前記回答テキストに基づいて、仮想オブジェクトによって表現される表情および/または動作の標識を確定するように構成される標識確定モジュールと、

前記回答音声信号、前記表情および/または動作の標識に基づいて、前記仮想オブジェクトを含む出力ビデオを生成し、前記出力ビデオは、前記回答音声信号に基づいて確定された、前記仮想オブジェクトによって表現される唇形シーケンスを含むように構成される第 1 出力ビデオ生成モジュールとを含み、

前記第 1 回答音声信号生成モジュールは、

前記回答テキストを 1 セットのテキストユニットに分割するように構成されるテキストユニット分割モジュールと、

音声信号ユニットとテキストユニットとのマッピング関係に基づいて、前記 1 セットのテキストユニットにおけるテキストユニットに対応する音声信号ユニットを取得する音声信号ユニット取得モジュールであって、1 セットのテキストユニットから前記テキストユニットを選択するように構成されるテキストユニット選択モジュールと、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、音声ライブラリから前記テキストユニットに対応する前記音声信号ユニットを検索するように構成される検索モジュールとを含む、音声信号ユニット取得モジュールと、

前記取得された音声信号ユニットに基づいて、前記回答音声信号を生成するように構成される第 2 回答音声信号生成モジュールとを含み、

前記音声ライブラリには音声信号ユニットとテキストユニットとの前記マッピング関係が記憶され、前記音声ライブラリにおける音声信号ユニットは、取得された、前記仮想オブジェクトに関する音声記録データを分割することで得られるものであり、前記音声ライブラリにおけるテキストユニットは、分割で得られた音声信号ユニットに基づいて確定されるものであり、

前記第 1 出力ビデオ生成モジュールは、

前記回答音声信号を 1 セットの音声信号ユニットに分割するように構成される音声信号分割モジュールと、

前記 1 セットの音声信号ユニットに対応する前記仮想オブジェクトの唇形シーケンスを取得するように構成される唇形シーケンス取得モジュールと、

対応する前記表情および/または動作の標識に基づいて、前記仮想オブジェクトについての当該表情および/または動作のビデオセグメントを取得するように構成されるビデオセグメント取得モジュールと、

前記唇形シーケンスを前記ビデオセグメントに結合して前記出力ビデオを生成するように構成される第 2 出力ビデオ生成モジュールとを含み、

前記仮想オブジェクトについての当該表情および/または動作の前記ビデオセグメントを取得することは、

表情および/または動作の標識とビデオセグメントとの間の事前に記憶されたマッピング関係を利用して、対応する前記表情および/または動作の標識に基づいて当該表情および/または動作の前記ビデオセグメントを取得することを含む、

マンマシンインタラクションのための装置。

【請求項 9】

前記回答テキスト生成モジュールは、

前記受信した音声信号を識別して入力テキストを生成するように構成される入力テキスト生成モジュールと、

前記入力テキストに基づいて、前記回答テキストを取得するように構成される回答テキスト取得モジュールとを含む、請求項 8 に記載の装置。

【請求項 10】

前記回答テキスト取得モジュールは、

入力テキストと前記仮想オブジェクトの人格属性を用いて回答テキストを生成する機械学習モデルである対話モデルに前記入力テキストと前記仮想オブジェクトの人格属性を入力して前記回答テキストを取得するように構成される、モデルに基づく回答テキスト取得モジュールを含む、請求項 9 に記載の装置。

10

【請求項 11】

前記対話モデルは、前記仮想オブジェクトの人格属性および入力テキストサンプルと回答テキストサンプルとを含む対話サンプルトを利用してトレーニングすることで得られるものである、請求項 10 に記載の装置。

【請求項 12】

前記標識確定モジュールは、

テキストを用いて表情および/または動作の標識を確定する機械学習モデルである表情および動作識別モデルに前記回答テキストを入力して、前記表情および/または動作の標識を取得するように構成される表情動作標識取得モジュールを含む、請求項 8 に記載の装置。

20

【請求項 13】

前記第 2 出力ビデオ生成モジュールは、

前記ビデオセグメントにおける時間軸での所定の時間位置におけるビデオフレームを確定するように構成されるビデオフレーム確定モジュールと、

前記唇形シーケンスから前記所定の時間位置に対応する唇形を取得するように構成される唇形取得モジュールと、

前記唇形を前記ビデオフレームに結合して前記出力ビデオを生成するように構成される結合モジュールとを含む、請求項 8 に記載の装置。

【請求項 14】

前記回答音声信号と前記出力ビデオとを関連付けて出力するように構成される出力モジュールをさらに含む、請求項 8 に記載の装置。

30

【請求項 15】

少なくとも 1 つのプロセッサ、および

前記少なくとも 1 つのプロセッサに通信接続されたメモリを含み、

前記メモリには、前記少なくとも 1 つのプロセッサによって実行可能なコマンドが記憶され、前記コマンドは前記少なくとも 1 つのプロセッサによって実行されることにより、前記少なくとも 1 つのプロセッサが請求項 1 ~ 7 のいずれか一項に記載の方法を実行する、電子機器。

【請求項 16】

コンピュータに請求項 1 ~ 7 のいずれか一項に記載の方法を実行させるためのコンピュータコマンドが記憶された非一時的コンピュータ可読記憶媒体。

40

【請求項 17】

プロセッサによって実行されると、請求項 1 ~ 7 のいずれか一項に記載の方法を実現するコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、人工知能の分野に関し、特にディープラーニング、音声技術およびコンピュータビジョン分野におけるマンマシンインタラクションのための方法、装置、機器および媒体に関する。

50

【背景技術】

【0002】

コンピュータ技術の急速な発展に伴って、人間と機械のインタラクションがますます多くなっている。ユーザの体験を向上させるために、マンマシンインタラクション技術が急速に発展している。ユーザが音声コマンドを出した後、計算機器は音声識別技術によってユーザの音声を識別する。識別を完了した後に、ユーザの音声コマンドに応じる操作を実行する。このような音声インタラクション方式はマンマシンインタラクションの体験を改善する。しかしながら、マンマシンインタラクションのプロセスにおいては、多くの解決する必要のある問題がまだ存在している。

【発明の概要】

【0003】

本開示は、マンマシンインタラクションのための方法、装置、機器および媒体を提供する。

本開示の第1態様によれば、マンマシンインタラクションのための方法が提供される。この方法は、受信した音声信号に基づいて、音声信号に対する回答の回答テキストを生成することを含む。この方法は、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットを含む回答テキストに対応する回答音声信号を生成することをさらに含む。この方法は、回答テキストに基づいて、仮想オブジェクトによって表現される表情および/または動作の標識を確定することをさらに含む。この方法は、回答音声信号、表情および/または動作の標識に基づいて、仮想オブジェクトを含む出力ビデオを生成することを含み、出力ビデオは、回答音声信号に基づいて確定された、仮想オブジェクトによって表現される唇形シーケンスを含む。

【0004】

本開示の第2態様によれば、マンマシンインタラクションのための装置が提供される。この装置は、受信した音声信号に基づいて、音声信号に対する回答の回答テキストを生成するように構成される回答テキスト生成モジュールと、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットを含む回答テキストに対応する回答音声信号を生成し、生成された回答音声信号は1セットのテキストユニットに対応する1セットの音声ユニットを含むように構成される第1回答音声信号生成モジュールと、回答テキストに基づいて、仮想オブジェクトによって表現される表情および/または動作の標識を確定する標識確定モジュールと、回答音声信号、表情および/または動作の標識に基づいて、仮想オブジェクトを含む出力ビデオを生成し、出力ビデオは回答音声信号に基づいて確定された、仮想オブジェクトによって表現される唇形シーケンスを含むように構成される第1出力ビデオ生成モジュールとを含む。

【0005】

本開示の第3態様によれば、電子機器が提供される。この電子機器は、少なくとも1つのプロセッサ、および少なくとも1つのプロセッサに通信接続されるメモリを含み、ここで、メモリには、少なくとも1つのプロセッサによって実行可能なコマンドが記憶され、コマンドは少なくとも1つのプロセッサによって実行されることにより、少なくとも1つのプロセッサが本開示の第1態様の方法を実行することができる。

【0006】

本開示の第4態様によれば、コンピュータに本開示の第1態様の方法を実行させるためのコンピュータコマンドが記憶された非一時的コンピュータ可読記憶媒体が提供される。

本開示の第5態様によれば、コンピュータプログラムを含むコンピュータプログラム製品が提供される。前記コンピュータプログラムはプロセッサによって実行されると、本開示の第1態様の方法を実現する。

【0007】

理解できるように、この部分に説明される内容は、本開示の実施形態の肝心または重要な特徴を示すことを目的とせず、本開示の保護範囲を限定するためのものではないことである。本開示の他の特徴は、以下の明細書によって理解されやすくなる。

10

20

30

40

50

【図面の簡単な説明】**【0008】**

図面は、本発明をより良く理解するためのものであり、本開示に対する限定を構成していない。

【図1】本開示の複数の実施形態を実現することができる環境100を示す概略図である。

【図2】本開示のいくつかの実施形態によるマンマシンインタラクションのためのプロセス200を示すフローチャートである。

【図3】本開示のいくつかの実施形態によるマンマシンインタラクションのための方法300を示すフローチャートである。

【図4】本開示のいくつかの実施形態による対話モデルをトレーニングするための方法400を示すフローチャートである。

10

【図5A】本開示のいくつかの実施形態による対話モデルネットワーク構造を示す例である。

【図5B】本開示のいくつかの実施形態によるマスクテーブルを示す例である。

【図6】本開示のいくつかの実施形態による回答音声信号を生成するための方法600を示すフローチャートである。

【図7】本開示のいくつかの実施形態による表情および/または動作の説明例700を示す概略図である。

【図8】本開示のいくつかの実施形態による表情および動作識別モデルを取得して使用するための方法800を示すフローチャートである。

20

【図9】本開示のいくつかの実施形態による出力ビデオを生成するための方法900を示すフローチャートである。

【図10】本開示のいくつかの実施形態による出力ビデオを生成するための方法1000を示すフローチャートである。

【図11】本開示の実施形態によるマンマシンインタラクションのための装置1100を示す概略的ブロック図である。

【図12】本開示の複数の実施形態を実施することができる機器1200を示すブロック図である。

【発明を実施するための形態】**【0009】**

30

以下、図面に合わせて本開示の例示的な実施形態を説明し、それに含まれる本開示の実施形態における様々な詳細が理解を助けるためのもので、それらは単なる例示的なものと考えられるべきである。したがって、当業者であれば、本開示の範囲および精神から逸脱することなく、本明細書で説明される実施形態に対して様々な変更および修正を行うことができることを認識すべきである。同様に、明瞭と簡潔のために、以下の説明では公知の機能および構造についての説明を省略する。

【0010】

本開示の実施形態の説明において、用語「含む」およびその類似用語はオープンな包含であり、すなわち「含むが、これらに限定されない」ことを理解されたい。用語「に基づいて」は、「少なくとも部分的に基づいて」ことを理解されたい。用語「一実施形態」または「該実施形態」は、「少なくとも1つの実施形態」ことを理解されたい。用語「第1」、「第2」などは異なるまたは同じオブジェクトを指すことができる。以下には他の明示的および暗示的な定義をさらに含む可能性もある。

40

【0011】

機械を人間のように人間と対話させることは人工知能の重要な目標である。現在、機械と人間のインタラクションの形式がインターフェースによるインタラクションから言語によるインタラクションへと進化している。しかしながら、従来の技術案では、ただ内容に限られたインタラクションだけであり、または音声の出力しか実行できない。例えばインタラクションの内容は主に、「天気を調べる」、「音楽を再生しろ」、「アラームを設定しろ」など、限られた分野でのコマンド型のインタラクションに限られる。また、イン

50

タラクションのモードも単一で、音声またはテキストによるインタラクションのみがある。また、マンマシンインタラクションには人格属性を欠けて、機械は対話する人よりも、ツールのようなものである。

【0012】

上述した問題を解決するために、本開示の実施形態によれば、改善案が提供される。この案において、計算機器は、受信した音声信号に基づいて、音声信号に対する回答の回答テキストを生成する。次に、計算機器は回答テキストに対応する回答音声信号を生成する。計算機器は、回答テキストに基づいて、仮想オブジェクトによって表現される表情および/または動作の標識を確定する。続いて、計算機器は、回答音声信号、表情および/または動作の標識に基づいて、仮想オブジェクトを含む出力ビデオを生成する。この方法により、インタラクションの内容の範囲を著しく増加させ、マンマシンインタラクションの品質とレベルを向上させ、ユーザ体験を向上させることができる。

10

【0013】

図1は、本開示の複数の実施形態を実現することができる環境100の概略図を示す。この例示的な環境は、マンマシンインタラクションを実現するために利用できる。この例示的な環境100は、計算機器108および端末機器104を含む。

【0014】

端末104における仮想人物などの仮想オブジェクト110は、ユーザ102と対話するために利用できる。インタラクションプロセスにおいて、ユーザ102は、端末104に問い合わせまたはチャット語句を送信することができる。端末104は、ユーザ102の音声信号を取得し、ユーザから入力された音声信号に対する回答を仮想オブジェクト110によって表現するために使用され、これによって人間と機械の対話を実現することができる。

20

【0015】

端末104は任意のタイプの計算機器として実現されることができ、携帯電話（例えばスマートフォン）、ラップトップコンピュータ、ポータブルデジタルアシスタント（PDA）、電子ブックリーダー、ポータブルゲームコンソール、ポータブルメディアプレーヤ、ゲームコンソール、セットトップボックス（STB）、スマートテレビ（TV）、パーソナルコンピュータ、車載コンピュータ（例えば、ナビゲーションユニット）、ロボットなどを含むがこれらに限定されない。

30

【0016】

端末104は、取得された音声信号をネットワーク106を介して計算機器108に送信する。計算機器108は、端末104から取得された音声信号に基づいて、対応する出力ビデオと出力音声信号を生成して、端末104上における仮想オブジェクト110によって表現することができる。

【0017】

図1は、計算機器108において、入力された音声信号に基づいて出力ビデオおよび出力音声信号を取得するプロセスを示しており、これは一例に過ぎず、本開示への具体的な限定ではない。このプロセスは、端末104上で実現されてもよく、または一部が計算機器108上で、他の一部が端末104上で実現されてもよい。いくつかの実施形態では、計算機器108と端末104は一体に統合されてもよい。図1は、計算機器108がネットワーク106を介して端末104に接続されていることを示す。これは一例に過ぎず、本開示への具体的な限定ではない。計算機器108は、他の方法で端末104と接続することもでき、例えば、ネットワークケーブルで直接的に接続される。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。

40

【0018】

計算機器108は任意のタイプの計算機器として実現されることができ、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルドまたはラップトップラップトップ機器、携帯機器（例えば携帯電話、パーソナルデジタルアシスタント（PDA）、メディアプレーヤなど）、マルチプロセッサシステム、消費者向け電子製品、小型コンピュータ、大型

50

コンピュータ、上記システムまたは機器のいずれかを含む分散式計算環境などを含むがこれらに限定されない。サーバは、クラウドサーバであってもよく、クラウド計算サーバまたはクラウドホストとも呼ばれ、クラウド計算サービスシステム中のホスト製品として、従来の物理ホストとVPSサービス(「Virtual Private Server」、または「VPS」と略称される)における、管理の難度が高く、業務拡張性が弱いという欠陥を解決する。サーバは、分散式システムのサーバであってもよいし、ブロックチェーンと組み合せられたサーバであってもよい。

【0019】

計算機器108は、端末104から取得された音声信号を処理することで、回答のための出力音声信号および出力ビデオを生成する。

10

この方法により、インタラクションの内容の範囲を著しく増加させ、マンマシンインタラクションの品質とレベルを向上させ、ユーザ体験を向上させることができる。

【0020】

上記の図1は、本開示の複数の実施形態を実現することができる環境100の概略図を示す。以下、図2によってマンマシンインタラクションのための方法200の概略図を説明する。この方法200は、図1における計算機器108または任意の適当な計算機器によって実現することができる。

【0021】

図2に示すように、計算機器108は、受信した音声信号202を取得する。次に、計算機器108は、受信した音声信号を音声識別(ASR)して入力テキスト204を生成する。ここでは、計算機器108は、任意の適当な音声識別アルゴリズムを用いて入力テキスト204を取得することができる。

20

【0022】

計算機器108は、回答用の回答テキスト206を取得するために、取得された入力テキスト204を対話モデルに入力する。この対話モデルはトレーニングされた機械学習モデルであり、そのトレーニングプロセスはオフラインで行うことができる。代替的または付加的には、この対話モデルはニューラルネットワークモデルであり、以下、図4および図5Aと図5Bに関連してこの対話モデルのトレーニングプロセスを紹介する。

【0023】

その後、計算機器108は、音声合成技術(TTS)により回答テキスト206を利用して回答音声信号208を生成するとともに、回答テキスト206に基づいて、現在の回答に使用されている表情および/または動作の標識210をさらに識別することができる。いくつかの実施形態では、この標識は表情および/または動作ラベルであってもよい。いくつかの実施形態では、この標識は表情および/または動作のタイプである。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。

30

【0024】

計算機器108は取得された表情および/または動作の標識に基づいて、出力ビデオ212を生成する。次に、回答音声信号208と出力ビデオ212を、端末上で同期して再生されるように端末に送信する。

【0025】

40

上記の図2は、本開示の複数の実施形態によるマンマシンインタラクションのためのプロセス200の概略図を示す。以下、図3に関連して、本開示のいくつかの実施形態によるマンマシンインタラクションのための方法300のローチャートを説明する。図3の方法300は、図1の計算機器108または任意の適当な計算機器によって実行することができる。

【0026】

ブロック302において、受信した音声信号に基づいて、音声信号に対する回答の回答テキストを生成する。例えば、図2に示すように、計算機器108は、受信した音声信号202に基づいて、受信した音声信号202に対する回答テキスト206を生成する。

【0027】

50

いくつかの実施形態では、計算機器 108 は、受信した音声信号を識別して入力テキスト 204 を生成する。入力テキストを取得するために任意の適当な音声識別技術を採用して音声信号を処理することができる。続いて、計算機器 108 は、入力テキスト 204 に基づいて、回答テキスト 206 を取得する。この方法によって、ユーザから受信された音声の回答テキストを迅速かつ効率的に取得することができる。

【0028】

いくつかの実施形態では、計算機器 108 は、回答テキスト 206 を取得するために、入力テキストと仮想オブジェクトの人格属性を用いて回答テキストを生成する機械学習モデルである対話モデルに入力テキスト 204 と仮想オブジェクトの人格属性を入力する。代替的または付加的には、この対話モデルはニューラルネットワークモデルである。いくつかの実施形態では、この対話モデルは任意の適当な機械学習モデルであってもよい。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。この方法によって、回答テキストを迅速かつ正確に確定することができる。

10

【0029】

いくつかの実施形態では、対話モデルは、仮想オブジェクトの人格属性および入力テキストサンプルと回答テキストサンプルとを含む対話サンプルトを利用してレーニングすることで得られる。この対話モデルは計算機器 108 によってオフラインでトレーニングすることで得られてもよい。計算機器 108 は、まず仮想オブジェクトの人格属性を取得し、人格属性は仮想オブジェクトの、性別、年齢、星座などの、人と関連する特性を説明する。次に、計算機器 108 は、人格属性および入力テキストサンプルと回答テキストサンプルとを含む対話サンプルに基づいて、対話モデルをトレーニングする。トレーニングするとき、人格属性と入力テキストサンプルを入力とし、回答テキストサンプルを出力としてトレーニングする。いくつかの実施形態では、対話モデルは他の計算機器によってオフラインでトレーニングしてもよい。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。この方法によって、対話モデルを迅速的に取得することができる。

20

【0030】

以下、図 4 と図 5 A および図 5 B に関連してこの対話モデルのレーニングを紹介する。図 4 は、本開示のいくつかの実施形態による対話モデルをトレーニングするための方法 400 のフローチャートを示す。図 5 A および図 5 B は本開示のいくつかの実施形態による対話モデルネットワーク構造および用いられるマスクテーブルの一例を示す。

30

【0031】

図 4 に示すように、プレトレーニング段階 404 において、例えば 10 億レベルの人間対話コーパスなどのソーシャルプラットフォーム上で自動的にマイニングされたコーパス 402 を用いて、モデルが基礎的なオープンドメイン対話能力を備えるように、対話モデル 406 をトレーニングする。次に、例えば 5 万レベルの特定の人格属性を有する対話コーパスなどの手動ラベル付け対話コーパス 410 を取得し、人格適合段階 408 において、指定の人格属性を用いて対話する能力を備えるように、対話モデル 406 をさらにトレーニングする。この指定の人格属性は、マンマシンインタラクションで使用しようとする仮想人物の、性別、年齢、趣味、星座などの人格属性である。

40

【0032】

図 5 A は対話モデルのモデル構造を示し、それは入力 504、モデル 502 およびさらなる回答 512 を含む。このモデルはディープラーニングモデルにおける Transformer モデルを用いており、モデルを使用するたびに、回答中の 1 つの単語を生成する。このプロセスは、具体的には、人格情報 506、入力テキスト 508、および回答 510 に既に生成された部分（例えば、単語 1 & 2）をモデルに入力して、さらなる回答 512 の次の単語（3）を生成し、このように再帰して、完全な回答文を生成する。モデルトレーニング時に、効率を向上させるために図 5 B におけるマスクテーブル 514 を用いて、回答の生成にバッチ（Batch）処理の操作を行う。

【0033】

50

ここで、図3に戻り、ブロック304において、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットを含む回答テキストに対応する回答音声信号を生成し、生成された回答音声信号は1セットのテキストユニットに対応する1セットの音声信号ユニットを含む。例えば、計算機器108は、予め記憶された音声信号ユニットとテキストユニットとのマッピング関係を利用して、1セットのテキストユニットを含む回答テキスト206に対応する回答音声信号208を生成し、生成した回答音声信号は該セットのテキストユニットに対応する1セットの音声信号ユニットを含む。

【0034】

いくつかの実施形態では、計算機器108は、回答テキスト206を1セットのテキストユニットに分割する。次に、計算機器108は、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットにおけるテキストユニットに対応する音声信号ユニットを取得する。計算機器108は、音声ユニットに基づいて、回答音声信号を生成する。この方法によって、回答テキストに対応する回答音声信号を迅速かつ効率的に生成することができる。

10

【0035】

いくつかの実施形態では、計算機器108は、1セットのテキストユニットからテキストユニットを選択する。次に、計算機器は、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、音声ライブラリからテキストユニットに対応する音声信号ユニットを検索する。この方式によって、音声信号ユニットを迅速に取得することができ、このプロセスにかかる時間を短縮し、効率を向上させる。

20

【0036】

いくつかの実施形態では、音声ライブラリに音声信号ユニットとテキストユニットとのマッピング関係が記憶され、音声ライブラリにおける音声信号ユニットは、取得された、仮想オブジェクトに関する音声記録データを分割することで取得されるものであり、音声ライブラリにおけるテキストユニットは、分割で得られた音声信号ユニットに基づいて確定されるものである。音声ライブラリは以下の方式によって生成される。まず、仮想オブジェクトに関連する音声記録データを取得する。例えば、仮想オブジェクトに対応する人間の声を録音する。次に、音声記録データを複数の音声信号ユニットに分割する。音声信号ユニットに分割された後、複数の音声信号ユニットに対応する複数のテキストユニットを確定し、ここで、第1音声信号ユニットは1つのテキストユニットに対応する。次に、複数の音声信号ユニットにおける音声信号ユニットと複数のテキストユニットにおける対応するテキストユニットとを関連付けて音声ライブラリに記憶し、それにより音声ライブラリが生成される。この方法により、テキストの音声信号ユニットを取得する効率を高め、取得時間を節約することができる。

30

【0037】

以下、図6に関連して、回答音声信号を生成するプロセスを具体的に説明する。ここで、図6は、本開示のいくつかの実施形態による回答音声信号を生成するための方法600のフローチャートを示す。

【0038】

図6に示すように、機械が人間のチャットをよりリアルにシミュレートするために、仮想キャラクタと一致する人間の声を用いて回答音声信号を生成する。このプロセス600はオフラインとオンラインの2つの部分に分割される。オフライン部分では、ブロック602において、仮想キャラクタと一致する人間の録音録画データを収集する。次に、ブロック604の後に、録音された音声信号を音声ユニットに分割し、対応するテキストユニットとアライメントすることで、単語ごとに対応する音声信号を記憶している音声ライブラリ606を取得する。このオフラインプロセスは、計算機器108または任意の他の適切な装置で行われることができる。

40

【0039】

オンライン部分では、回答テキスト中の単語シーケンスに基づいて音声ライブラリ60

50

6 から対応する音声信号を抽出して出力音声信号を合成する。まず、ブロック 6 0 8 において、計算機器 1 0 8 は回答テキストを取得する。次に、計算機器 1 0 8 は回答テキスト 6 0 8 を 1 セットのテキストユニットに分割する。その後、ブロック 6 1 0 において、音声ライブラリ 6 0 6 からテキストユニットに対応する音声ユニットの抜き取りおよびスライスを行う。次に、ブロック 6 1 2 において、回答音声信号を生成する。したがって、音声ライブラリを利用して回答音声信号をオンラインで取得することができる。

【 0 0 4 0 】

次に、図 3 に戻って引き続き説明し、ブロック 3 0 6 において、回答テキストに基づいて、仮想オブジェクトによって表現される表情および / または動作の標識を確定する。例えば、計算機器 1 0 8 は、回答テキスト 2 0 6 に基づいて、仮想オブジェクト 1 1 0 によって表現される表情および / または動作の標識 2 1 0 を確定する。

10

【 0 0 4 1 】

いくつかの実施形態では、計算機器 1 0 8 は、テキストを用いて表情および / または動作の標識を確定する機械学習モデルである表情および動作識別モデルに回答テキストを入力して、表情および / または動作の標識を取得する。この方法によって、テキストを迅速かつ正確に利用して、使用しようとする表情と動作を確定することができる。

【 0 0 4 2 】

以下、図 7 と図 8 に関連して表情および / または動作の標識および表情および動作の記述を説明する。図 7 は、本開示のいくつかの実施形態による表情および / または動作の例 7 0 0 の概略図を示す。図 8 は、本開示のいくつかの実施形態による表情および動作識別モデルを取得し使用するための方法 8 0 0 のフローチャートを示す。

20

【 0 0 4 3 】

対話において、仮想オブジェクト 1 1 0 の表情と動作は対話内容によって決定され、仮想人物は「私はとても嬉しいです」と答える場合、楽しい表情を用いることができ、「こんにちは」と答える場合、手を振る動作を用いることができ、このため、表情と動作識別は対話モデルにおける回答テキストに基づいて仮想人物の表情と動作ラベルを識別するものである。このプロセスには表情および動作ラベルシステムの設定と識別の 2 つの部分が含まれる。

【 0 0 4 4 】

図 7 において、対話過程に関する高頻度の表情および / または動作に 1 1 個のラベルが設定される。いくつかのシーンでは表情と動作が共同で働くので、システムにおいては、あるラベルが表情であるか動作であるかを厳密に区別していない。いくつかの実施形態では、表情と動作をそれぞれ設定してから、異なるラベルまたは標識を割り当てることができる。回答テキストを利用して表情および / または動作のラベルまたは標識を取得する場合、トレーニングされたモデルによって取得してもよいし、トレーニングされた、表情に対するモデルと動作に対するモデルによって対応する表情ラベルと動作ラベルをそれぞれ取得してもよい。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。

30

【 0 0 4 5 】

表情および動作ラベルの識別プロセスは、図 8 に示すように、オフラインフローとオンラインフローに分けられる。オフラインフローは、ブロック 8 0 2 において、対話テキストの手動ラベル付け表情および動作コーパスを取得する。ブロック 8 0 4 において、BERT 分類モデルをトレーニングし、表情および動作識別モデル 8 0 6 を取得する。オンラインフローでは、ブロック 8 0 8 において回答テキストを取得し、次に回答テキストを表情および動作識別モデル 8 0 6 に入力して、ブロック 8 1 0 において表情および動作識別を行う。次に、ブロック 8 1 2 において、表情および / または動作の標識を出力する。いくつかの実施形態では、この表情および動作識別モデルは、様々な適当なニューラルネットワークモデルなどの任意の適当な機械学習モデルを用いることができる。

40

【 0 0 4 6 】

次に、図 3 に戻って説明を続け、ブロック 3 0 8 において、回答音声信号、表情および

50

／または動作の標識に基づいて、仮想オブジェクトを含む出力ビデオを生成し、出力ビデオは回答音声信号に基づいて確定された、仮想オブジェクトによって表現される唇形シーケンスを含む。例えば、計算機器 108 は、回答音声信号 208、表情および／または動作の標識 210 に基づいて、仮想オブジェクト 110 を含む出力ビデオ 212 を生成する。出力ビデオには、回答音声信号に基づいて確定された、仮想オブジェクトによって表現される唇形シーケンスを含む。このプロセスは、以下、図 9 と図 10 に関連して詳細に説明する。

【0047】

いくつかの実施形態では、計算機器 108 は、回答音声信号 208 と出力ビデオ 212 とを関連付けて出力する。この方法によって、正確なマッチングした音声とビデオの情報を生成することができる。このプロセスでは、回答音声信号 208 と出力ビデオ 212 とを時間的に同期させることによって、ユーザとやり取りをする。

10

【0048】

この方法により、インタラクションの内容の範囲を著しく増加させ、マンマシンインタラクションの品質とレベルを向上させ、ユーザ体験を向上させることができる。

以上、図 3 から図 8 に関連して、本開示のいくつかの実施形態によるマンマシンインタラクションのための方法 300 のローチャートを説明する。以下、図 9 に関連して、回答音声信号、表情および／または動作の標識に基づいて出力ビデオを生成するプロセスについて詳細に説明する。図 9 は、本開示のいくつかの実施形態による出力ビデオを生成するための方法 900 のフローチャートを示す。

20

【0049】

ブロック 902 において、計算機器 108 は回答音声信号を 1 セットの音声信号ユニットに分割する。いくつかの実施形態では、計算機器 108 は、ワード単位で音声信号ユニットを分割する。いくつかの実施形態では、計算機器 108 は、音節単位で音声信号ユニットを分割する。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。当業者は任意の適当な音声サイズで音声ユニットを分割することができる。

【0050】

ブロック 904 において、計算機器 108 は、1 セットの音声信号ユニットに対応する仮想オブジェクトの唇形シーケンスを取得する。計算機器 108 は、対応するデータベースから音声信号ごとに対応する唇形ビデオを検索することができる。音声信号ユニットと唇形の対応関係を生成する場合、まず、仮想オブジェクトに対応する人間の発声ビデオを録画し、次に、ビデオから音声信号ユニットに対応する唇形を抽出する。次に、唇形と音声信号ユニットとを関連付けてデータベースに記憶する。

30

【0051】

ブロック 906 において、計算機器 108 は、表情および／または動作の標識に基づいて、仮想オブジェクトについての対応する表情および／または動作のビデオセグメントを取得する。データベースまたは記憶装置には、表情および／または動作の標識と、対応する表情および／または動作のビデオセグメントとのマッピング関係が事前に記憶される。例えば表情および／または動作のラベルまたはタイプなどの標識を取得した後に、表情および／または動作の標識と、ビデオセグメントとのマッピング関係を利用して、対応するビデオを検索することができる。

40

【0052】

ブロック 908 において、計算機器 108 は、唇形シーケンスをビデオセグメントに結合して出力ビデオを生成する。計算機器は、時系列に、取得された、1 セットの音声信号ユニットに対応する唇形シーケンスをビデオセグメントの各フレームに結合する。

【0053】

いくつかの実施形態では、計算機器 108 は、ビデオセグメントにおける時間軸での所定の時間位置におけるビデオフレームを確定する。次に、計算機器 108 は、唇形シーケンスから所定の時間位置に対応する唇形を取得する。唇形を取得した後、計算機器 108

50

は唇形をビデオフレームに結合して出力ビデオを生成する。この方式により、正確な唇形を含むビデオを迅速に取得することができる。

【0054】

この方法によって、仮想人物の唇形を音声と動作により正確にマッチングすることができ、ユーザの体験を改善する。

以上、図9に関連して、本開示のいくつかの実施形態による出力ビデオを生成するための方法900のフローチャートを示す。以下、図10に関連して、出力ビデオを生成するプロセスについてさらに説明する。図10は、本開示のいくつかの実施形態による出力ビデオを生成するための方法1000のフローチャートを示す。

【0055】

図10においては、生成されたビデオは、回答音声信号と表情動作ラベルに基づいて仮想人物を合成するビデオセグメントを含む。このプロセスは図10に示すように、唇形ビデオの取得、表情動作ビデオの取得およびビデオのレンダリングの三つの部分を含む。

【0056】

唇形ビデオの取得プロセスは、オンラインフローとオフラインフローに分けられる。オフラインフローでは、ブロック1002において、音声および対応する唇形の間人間の撮影を実行する。次に、ブロック1004において、人間の音声と唇形ビデオのアライメントを実行する。このプロセスでは、音声ユニットごとに対応する唇形ビデオを取得する。その後、取得された音声ユニットと唇形ビデオとを関連付けて音声唇形ライブラリ1006に記憶する。オンラインフローでは、ブロック1008において、計算機器108は回答音声信号を取得する。次に、ブロック1010において、計算機器108は回答音声信号を音声信号ユニットに分割し、その後、唇形データベース1006から音声信号ユニットに基づいて対応する唇形を抽出する。

【0057】

表情動作ビデオの取得プロセスもオンラインフローとオフラインフローに分けられる。オフラインフローでは、ブロック1014において、人間の表情動作ビデオを撮影する。次に、ブロック1016において、ビデオを分割して表情および/または動作標識ごとに対応するビデオを取得し、即ち、表情および/または動作をビデオユニットとアライメントする。その後、表情および/または動作ラベルとビデオとを関連付けて表情および/または動作ライブラリ1018に記憶する。いくつかの実施形態では、表情および/または動作ライブラリ1018には、表情および/または動作の標識と、対応するビデオとのマッピング関係を記憶する。いくつかの実施形態では、表情および/または動作ライブラリにおいて、表情および/または動作の標識を用いて、マルチレベルマッピングを利用して対応するビデオを見つける。上記の例は、本開示を説明するためのものに過ぎず、本開示への具体的な限定ではない。

【0058】

オンライン段階のフローでは、ブロック1012において、計算機器108は、入力表情および/動作の標識を取得する。次に、ブロック1020において、表情および/または動作の標識に基づいてビデオセグメントを抽出する。

【0059】

その後、ブロック1022において、唇形シーケンスをビデオセグメントに結合する。このプロセスにおいて、表情と動作ラベルに対応するビデオは時間軸でのビデオフレームによってスッチングされてなり、唇形シーケンスに基づいて、それぞれの唇形を時間軸での同じ位置のビデオフレームにレンダリングし、最終的に組み合わせられたビデオを出力する。次に、ブロック1024において、出力ビデオを生成する。

【0060】

図11は、本開示の実施形態によるマンマシンインタラクションのための装置1100の概略的ブロック図を示す。図11に示すように、装置1100は、受信した音声信号に基づいて、音声信号に対する回答の回答テキストを生成するように構成される回答テキスト生成モジュール1102を含む。装置1100は、音声信号ユニットとテキストユニッ

10

20

30

40

50

トとのマッピング関係に基づいて、1セットのテキストユニットを含む回答テキストに対応する回答音声信号を生成し、生成された回答音声信号は1セットのテキストユニットに対応する1セットの音声ユニットを含むように構成される第1回答音声信号生成モジュール1104をさらに含む。装置1100は、回答テキストに基づいて、仮想オブジェクトによって表現される表情および/または動作の標識を確定するように構成される標識確定モジュール1106をさらに含む。装置1100は、回答音声信号、表情および/または動作の標識に基づいて、仮想オブジェクトを含む出力ビデオを生成し、出力ビデオは回答音声信号に基づいて確定された、仮想オブジェクトによって表現される唇形シーケンスを含むように構成される第1出力ビデオ生成モジュール1108をさらに含む。

【0061】

いくつかの実施形態では、回答テキスト生成モジュール1102は、受信した音声信号を識別して入力テキストを生成するように構成される入力テキスト生成モジュールと、入力テキストに基づいて、回答テキストを取得するように構成される回答テキスト取得モジュールを含む。

【0062】

いくつかの実施形態では、回答テキスト生成モジュールは、回答テキストを取得するために、入力テキストと仮想オブジェクトの人格属性を用いて回答テキストを生成する機械学習モデルである対話モデルに入力テキストと仮想オブジェクトの人格属性を入力するように構成されるモデルに基づく回答テキスト取得モジュールを含む。

【0063】

いくつかの実施形態では、対話モデルは、仮想オブジェクトの人格属性および入力テキストサンプルと回答テキストサンプルとを含む対話サンプルトを利用してレーニングすることで得られるものである。

【0064】

いくつかの実施形態では、第1回答音声信号生成モジュールは、回答テキストを1セットのテキストユニットに分割するように構成されるテキストユニット分割モジュールと、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットにおけるテキストユニットに対応する音声信号ユニットを取得するように構成される音声信号ユニット取得モジュールと、音声ユニットに基づいて回答音声信号を生成するように構成される第2回答音声信号生成モジュールとを含む。

【0065】

いくつかの実施形態では、音声信号ユニット取得モジュールは、音声信号ユニットとテキストユニットとのマッピング関係に基づいて、1セットのテキストユニットからテキストユニットを選択するように構成されるテキストユニット選択モジュールと、音声ライブラリからテキストユニットに対応する音声信号ユニットを検索するように構成される検索モジュールとを含む。

【0066】

いくつかの実施形態では、音声ライブラリには音声信号ユニットとテキストユニットとのマッピング関係が記憶され、音声ライブラリにおける音声信号ユニットは、取得された、前記仮想オブジェクトに関する音声記録データを分割することで取得されるものであり、音声ライブラリにおけるテキストユニットは、分割で得られた音声信号ユニットに基づいて確定されるものである。

【0067】

いくつかの実施形態では、標識判定モジュール1106は、テキストを用いて表情および/または動作の標識を確定する機械学習モデルである表情および動作識別モデルに回答テキストを入力して、表情および/または動作の標識を取得するように構成される表情動作標識取得モジュールを含む。

【0068】

いくつかの実施形態では、第1出力ビデオ生成モジュール1108は回答音声信号を1セットの音声信号ユニットに分割するように構成される音声信号分割モジュールと、1セ

10

20

30

40

50

ットの音声信号ユニットに対応する仮想オブジェクトの唇形シーケンスを取得するように構成される唇形シーケンス取得モジュールと、表情および/または動作の標識に基づいて、仮想オブジェクトについての対応する表情および/または動作のビデオセグメントを取得するように構成されるビデオセグメント取得モジュールと、唇形シーケンスをビデオセグメントに結合して出力ビデオを生成するように構成される第2出力ビデオ生成モジュールとを含む。

【0069】

いくつかの実施形態では、第2出力ビデオ生成モジュールは、ビデオセグメントにおける時間軸での所定の時間位置におけるビデオフレームを確定するように構成されるビデオフレーム確定モジュールと、唇形シーケンスから所定の時間位置に対応する唇形を取得するように構成される唇形取得モジュールと、唇形をビデオフレームに結合して出力ビデオを生成するように構成される結合モジュールとを含む。

10

【0070】

いくつかの実施形態では、装置1100は回答音声信号と出力ビデオとを関連付けて出力するように構成される出力モジュールをさらに含む。

本開示の実施形態によれば、本公開は、電子機器、可読記憶媒体およびコンピュータプログラム製品をさらに提供する。

【0071】

図12は、本開示の実施形態を実施するための例示的な電子機器1200の概略的ブロック図を示す。図1の端末104および計算機器108は、電子機器1200によって実現することができる。電子機器は、ラップトップ型コンピュータ、デスクトップ型コンピュータ、ワークステーション、パーソナルデジタルアシスタント、サーバ、ブレードサーバ、大型コンピュータ、その他の好適なコンピュータなど、様々なデジタルコンピュータを指すことを意図している。電子機器は、例えば、パーソナルデジタル処理、携帯電話、スマートフォン、ウェアラブル機器、その他の類似装置などの様々なモバイル機器を指すこともできる。本明細書に示される部材、それらの接続関係、およびそれらの機能は、ただ一例に過ぎず、本明細書に記載および/または請求の本開示の実現を制限することを意図するものではない。

20

【0072】

図12に示すように、機器1200は、計算ユニット1201を含み、それはリードオンリーメモリ(ROM)1202に記憶されたプログラムまた記憶ユニット1208からランダムアクセスメモリ(RAM)1203にロードされたプログラムによって、種々の適当な操作と処理を実行することができる。RAM1203には、機器1200の動作に必要な種々のプログラムとデータを記憶することもできる。計算ユニット1201、ROM1202およびRAM1203はバス1204によって互いに接続される。入力/出力(I/O)インターフェース1205もバス1204に接続される。

30

【0073】

機器900における複数の部材はI/Oインターフェース1205に接続され、この複数の部材は、例えば、キーボード、マウスなどの入力ユニット1206と、例えば、様々なタイプのディスプレイ、スピーカーなどの出力ユニット1207と、例えば、磁気ディスク、光ディスクなどの記憶ユニット1208と、例えば、ネットワークカード、モデム、無線通信送受信機などの通信ユニット1209と、を含む。通信ユニット1209は、機器1200が例えば、インターネットなどのコンピュータネットワーク及び/又は様々な電気通信ネットワークを介して他の機器と情報/データのやり取りをすることを可能にする。

40

【0074】

計算ユニット1201は処理および計算能力を有する様々な汎用および/または専用の処理コンポーネントであってもよい。計算ユニット1201の例には、中央処理ユニット(CPU)、グラフィックス処理ユニット(GPU)、様々な専用人工知能(AI)計算チップ、様々な機械学習モデルアルゴリズムを実行する計算ユニット、デジタル信号プロ

50

セッサ（DSP）、および任意の適当なプロセッサ、コントローラ、マイクロコントローラなどが含まれるがこれらに限定されない。計算ユニット1201は以上で説明される例えば方法200、300、400、600、800、900および1000のような様々な方法および処理を実行する。例えば、いくつかの実施形態では、方法200、300、400、600、800、900および1000をコンピュータソフトウェアプログラムとして実現することができ、それは記憶ユニット1208などの機械可読媒体に有形的に含まれる。いくつかの実施形態では、コンピュータプログラムの一部または全部は、ROM1202および/または通信ユニット1209を介して機器1200にロードされたりインストールされたりすることができる。コンピュータプログラムがRAM1203にロードされて計算ユニット1201によって実行される場合、以上で説明される方法200、300、400、600、800、900および1000の1つまたは複数のステップを実行することができる。代替的に、他の実施形態において、計算ユニット1201は、他の任意の適当な方法で（例えば、ファームウェアを用いて）、方法200、300、400、600、800、900および1000を実行するように構成される。

【0075】

ここで述べるシステムおよび技術の様々な実施形態は、デジタル電子回路システム、集積回路システム、フィールドプログラマブルゲートアレイ（FPGA）、専用集積回路（ASIC）、専用標準製品（ASSP）、チップ上システムのシステム（SOC）、コンプレックスプログラマブルロジックデバイス（CPLD）、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されてもよい。これら様々な実施形態は、1つまたは複数のコンピュータプログラムに実装され、この1つまたは複数のコンピュータプログラムは、少なくとも1つのプログラマブルプロセッサを含むプログラマブルシステム上で実行することおよび/または解釈することが可能であり、このプログラマブルプロセッサは、専用または汎用のプログラマブルプロセッサであってもよいし、記憶システム、少なくとも1つの入力装置、および少なくとも1つの出力装置からデータおよびコマンドを受信し、この記憶システム、この少なくとも1つの入力装置、およびこの少なくとも1つの出力装置にデータおよびコマンドを送信することが可能である。

【0076】

本開示の方法を実施するためのプログラムコードは、1つまたは複数のプログラミング言語の任意の組み合わせを用いて作成することができる。これらのプログラムコードは、汎用コンピュータ、専用コンピュータ、または他のプログラマブルデータ処理装置のプロセッサまたはコントローラに提供することができ、これによって、プログラムコードがプロセッサまたはコントローラによって実行されると、フローチャートおよび/またはブロック図で規定された機能/操作が実行される。プログラムコードは完全に機械上で実行されても、部分的に機械で実行されても、独立ソフトウェアパッケージとして部分的に機械で実行されかつ部分的に遠隔機械上で実行されても、または、完全に遠隔機械またはサーバー上で実行されてもよい。

【0077】

本開示のコンテキストにおいて、機械可読媒体は、コマンド実行システム、装置、また機器が使用するプログラムまたはコマンド実行システム、装置または機器と組み合わせて使用されるプログラムを含むか記憶することができる有形の媒体であってもよい。機械可読媒体は、機械可読信号媒体または機械可読記憶媒体であってもよい。機械可読媒体は、電子的、磁氣的、光学的、電磁的、赤外線的、または半導体システム、装置や機器、または上記の内容の任意の適当な組み合わせを含むことができるが、これらに限定されない。機械可読記憶媒体のより具体的な例は、1つまたは複数のワイヤに基づく電気接続、ポータブルコンピュータディスク、ハードディスク、ランダムアクセスメモリ（RAM）、リードオンリーメモリ（ROM）、消去可能プログラマブル読み取り専用メモリ（EPROM またはフラッシュメモリ）、光ファイバ、ポータブルコンパクトディスク読み取り専用メモリ（CD-ROM）、光学記憶機器、磁気記憶機器、また上記の内容の任意の適当な組み

10

20

30

40

50

合わせを含むことができる。

【0078】

ユーザとのインタラクションを提供するために、ここで述べたシステムおよび技術をコンピュータ上で実行することができる。このコンピュータは、ユーザに情報を表示するための表示装置（例えば、CRT（Cathode Ray Tube、陰極線管）またはLCD（Liquid Crystal Crystal Display、液晶表示装置）モニター）と、キーボードやポインティング装置を有し、ユーザはこのキーボードやポインティング装置（例えば、マウスやトラックボール）によって入力をコンピュータに提供することができる。他の種類の装置は、さらに、ユーザとのインタラクションを提供するために利用することができる。例えば、ユーザに提供されるフィードバックは、任意の形のセンシングフィードバック（例えば、視覚フィードバック、聴覚フィードバック、または触覚フィードバック）であってもよい。しかも、ユーザからの入力を、任意の形（ボイス入力、音声入力、触覚入力を含む）で受け付けてもよい。

10

【0079】

ここで述べたシステムや技術は、バックステージ部材を含む計算システム（例えば、データサーバとして）や、ミドルウェア部材を含む計算システム（例えば、アプリケーションサーバ）や、フロントエンド部材を含む計算システム（例えば、グラフィカルユーザインタフェースやウェブブラウザを有するユーザコンピュータ、ユーザが、そのグラフィカルユーザインタフェースやウェブブラウザを通じて、それらのシステムや技術の実施形態とのインタラクティブを実現できる）、あるいは、それらのバックステージ部材、ミドルウェア部材、あるいはフロントエンド部材の任意の組み合わせからなる計算システムには実施されてもよい。システムの部材は、任意の形式や媒体のデジタルデータ通信（例えば、通信ネットワーク）により相互に接続されてもよい。通信ネットワークとしては、例えば、LAN（Local Area Network）、WAN（Wide Area Network）、インターネットを含む。

20

【0080】

コンピュータシステムは、クライアントとサーバとを含んでもよい。クライアントとサーバとは、一般に互いに離れ、通常、通信ネットワークを介してやりとりを行う。クライアントとサーバの関係は、対応するコンピュータ上で動作し、かつ、互いにクライアントとサーバの関係を有するコンピュータプログラムにより生成される。

30

【0081】

理解できるように、以上に示した様々な形式のフローを用いて、ステップを再び並び、増加または削除することができる。例えば、本開示に記載された各ステップは、並行して実行されてもよいし、順次実行されてもよいし、異なる順序で実行されてもよいし、本開示に開示された技術的解決手段が所望する結果を実現できれば、本明細書はここでは限定しない。

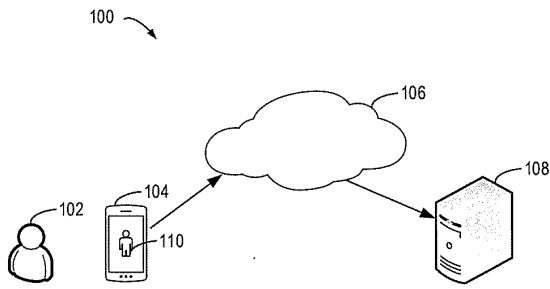
【0082】

上述した具体的な実施形態は、本開示に係る保護範囲に対する制限を構成していない。当業者は、設計要件やその他の要因によって、種々の変更、組み合わせ、サブコンビネーション、代替が可能であることは明らかである。本開示における精神および原則から逸脱することなく行われるいかなる修正、同等物による置換や改良等などは、いずれも本開示の保護範囲に含まれるものである。

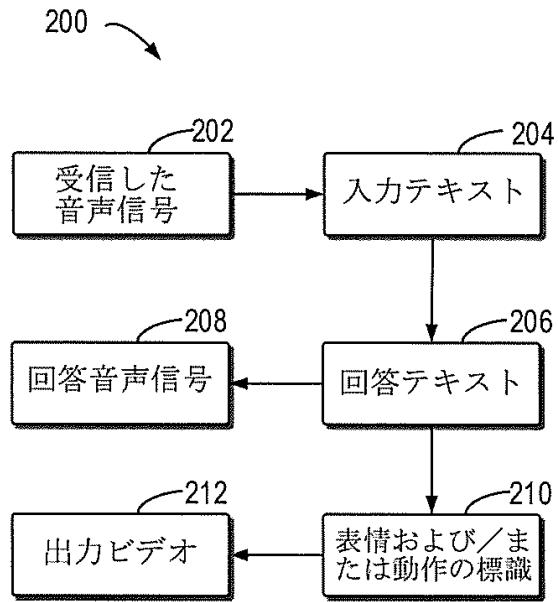
40

【図面】

【図 1】



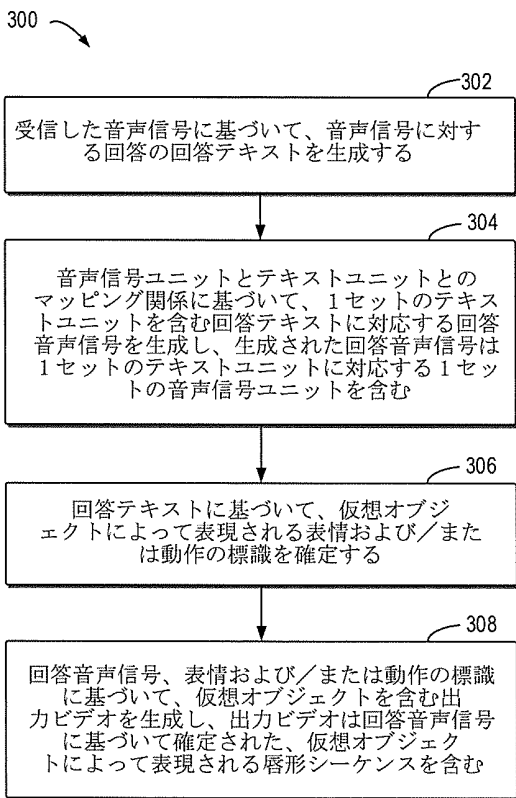
【図 2】



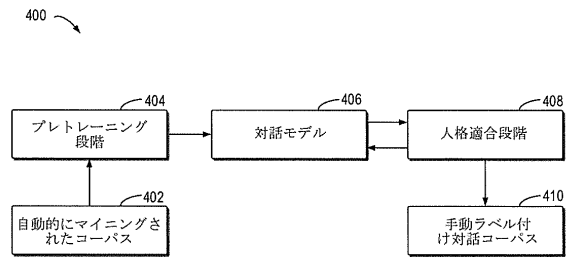
10

20

【図 3】



【図 4】

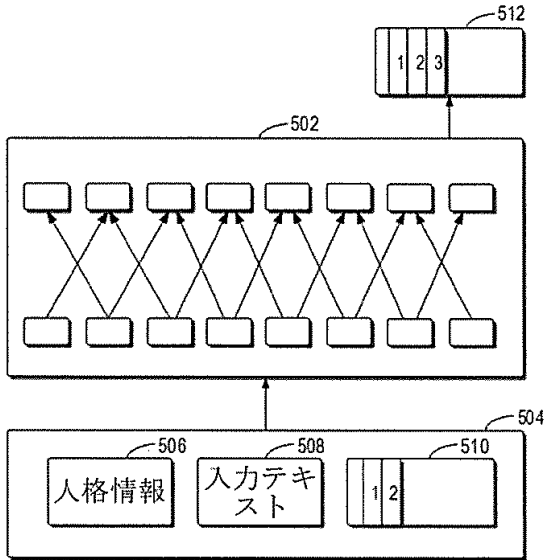


30

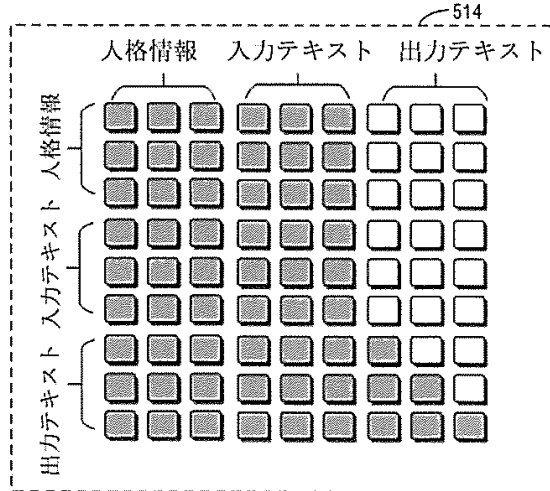
40

50

【図 5 A】

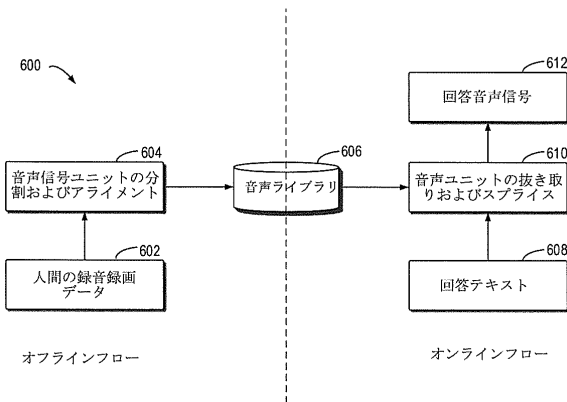


【図 5 B】



10

【図 6】



【図 7】

700

表情およびノミまたは動作レベル	解説
自然な状態—耳を傾けること	両手を体の前に置き、顔の表情が自然で、まばたきやわずかなうなずきの動作
自然な状態—述べる	両手を体の前に置き、片手を外に向け、顔の表情が自然で、まばたきができること
微笑する	主に眉毛や顔全体で、微笑んでいるような表情を表現
口を開けて笑う	口を開けて、はっきりとした口の動きで敬回笑う
しかめっ面	主に眉毛や顔全体で、しかめっ面をする
驚き	主に眉毛、目、顔全体で、驚きの表情を表現
ご挨拶とお別れ	右手を上げて手招きしている
自己紹介	右手を軽く上げ、襟元に向けて軽くタップする
確認—ジェスチャー	右手を上げてOKのジェスチャーをする
否定—頭部	頭を振って否定する (ショートアクション、話す前にアクションを完了させる)
思考と質問	少し首を傾げて、しかめっ面をしているような表情があること

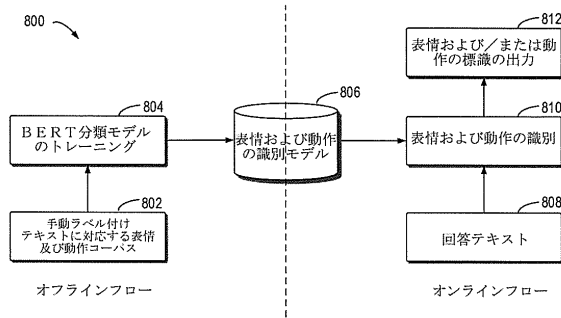
20

30

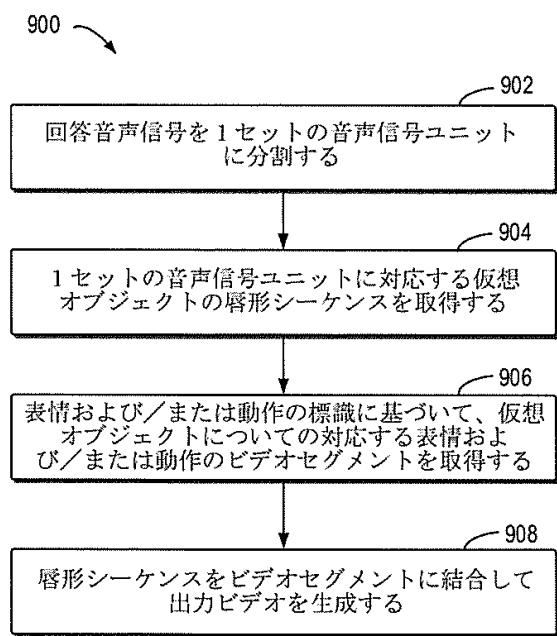
40

50

【 図 8 】



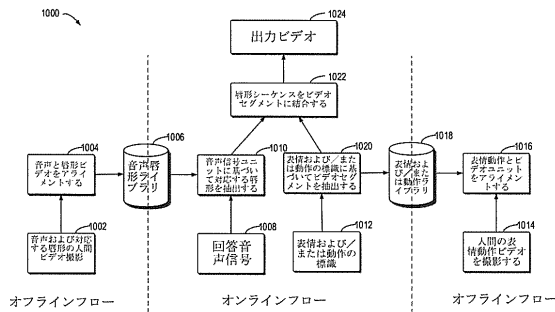
【 図 9 】



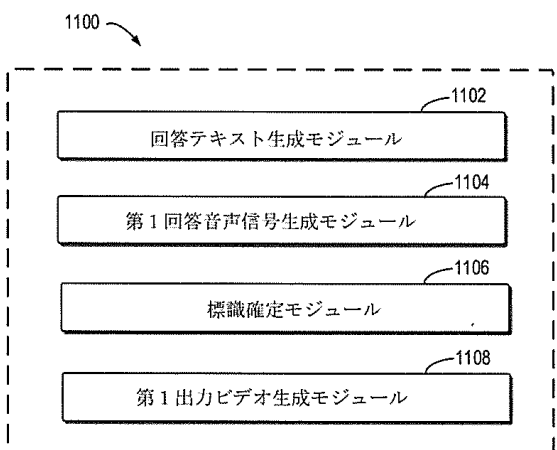
10

20

【 図 10 】



【 図 11 】

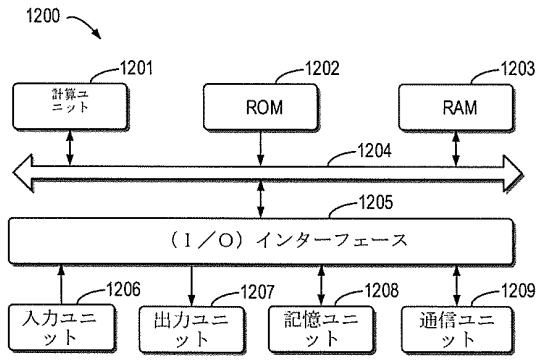


30

40

50

【図 12】



10

20

30

40

50

フロントページの続き

(51)国際特許分類

F I
G 0 6 F 3/01 5 1 0

8 5 , C h i n a

(74)代理人

100118902

弁理士 山本 修

(74)代理人

100106208

弁理士 宮前 徹

(74)代理人

100196508

弁理士 松尾 淳一

(74)代理人

100138759

弁理士 大房 直樹

(72)発明者

ウエンチュエン・ウー

中華人民共和国 ペキン 100085, ハイディアן ディストリクト, シャンディ テンス ストリート, 10番, バイドゥ キャンパス 2階

(72)発明者

フア・ウー

中華人民共和国 ペキン 100085, ハイディアן ディストリクト, シャンディ テンス ストリート, 10番, バイドゥ キャンパス 2階

(72)発明者

ハイフオン・ワーン

中華人民共和国 ペキン 100085, ハイディアן ディストリクト, シャンディ テンス ストリート, 10番, バイドゥ キャンパス 2階

審査官 鈴木 圭一郎

(56)参考文献

特開平09-016800(JP,A)

特開2006-330484(JP,A)

特開2006-099194(JP,A)

特開平11-231899(JP,A)

特開平11-339058(JP,A)

特開2020-160341(JP,A)

特開2004-310034(JP,A)

特開平9-081632(JP,A)

高津 弘明、小林 哲則, 対話エージェントのための性格モデル, 言語処理学会第21回年

次大会 発表論文集 [online], 日本, 言語処理学会, 2015年03月09日, 191~194

(58)調査した分野 (Int.Cl., DB名)

G 0 6 T 1 3 / 4 0

G 1 0 L 1 5 / 0 0

G 1 0 L 1 3 / 0 0

G 0 6 F 3 / 1 6

G 0 6 F 3 / 0 1