

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2020-190930  
(P2020-190930A)

(43) 公開日 令和2年11月26日(2020.11.26)

(51) Int.Cl.		F I	テーマコード (参考)		
<b>G06F</b>	<b>16/58</b>	(2019.01)	G06F 16/58	5E555	
<b>G06T</b>	<b>7/00</b>	(2017.01)	G06T 7/00	350C	5L096
<b>G06F</b>	<b>3/0484</b>	(2013.01)	G06F 3/0484		
<b>G06F</b>	<b>16/56</b>	(2019.01)	G06F 16/56		
<b>G10L</b>	<b>15/22</b>	(2006.01)	G10L 15/22	453	

審査請求 未請求 請求項の数 6 OL (全 24 頁)

(21) 出願番号 特願2019-95922 (P2019-95922)  
(22) 出願日 令和1年5月22日 (2019.5.22)

(71) 出願人 301022471  
国立研究開発法人情報通信研究機構  
東京都小金井市貫井北町4-2-1

(74) 代理人 110001195  
特許業務法人深見特許事務所

(72) 発明者 マガスーバ アリー  
東京都小金井市貫井北町4-2-1 国立  
研究開発法人情報通信研究機構内

(72) 発明者 杉浦 孔明  
東京都小金井市貫井北町4-2-1 国立  
研究開発法人情報通信研究機構内

(72) 発明者 河井 恒  
東京都小金井市貫井北町4-2-1 国立  
研究開発法人情報通信研究機構内

最終頁に続く

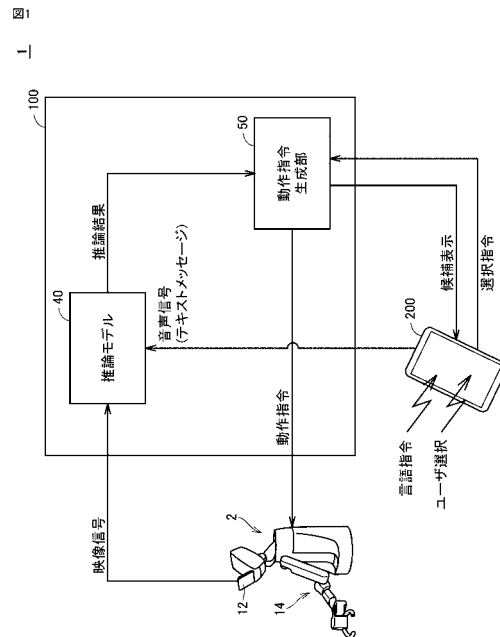
(54) 【発明の名称】 対象物検索システム、対象物検索方法および学習済モデル

(57) 【要約】

【課題】対象となる物体の候補が複数存在するような状況であっても、対象となる物体を特定できる技術を提供する。

【解決手段】対象物検索システムは、特定の対象物に関する命令文を取得する命令文取得部と、命令文に関連付けられた入力画像から、当該入力画像に含まれる個々の物体を示す1または複数の第1の部分画像を抽出する画像抽出部と、命令文と、第1の部分画像の各々と、当該第1の部分画像の画像内環境を示す情報との入力を受けて、第1の部分画像の各々が命令文により特定される対象物である確率を出力する学習済モデルとを含む。学習済モデルは、画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている。

【選択図】 図1



**【特許請求の範囲】****【請求項 1】**

特定の対象物に関する命令文を取得する命令文取得部と、

前記命令文に関連付けられた入力画像から、当該入力画像に含まれる個々の物体を示す 1 または複数の第 1 の部分画像を抽出する画像抽出部と、

前記命令文と、前記第 1 の部分画像の各々と、当該第 1 の部分画像の画像内環境を示す情報との入力を受けて、前記第 1 の部分画像の各々が前記命令文により特定される対象物である確率を出力する学習済モデルとを備え、

前記学習済モデルは、画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている、対象物検索システム。

10

**【請求項 2】**

前記画像抽出部は、前記命令文に関連付けられた前記入力画像から、いずれかの物体が存在する区域を示す 1 または複数の第 2 の部分画像をさらに抽出するように構成されており、

前記学習済モデルは、前記第 1 の部分画像と前記第 2 の部分画像との組み合わせの各々が前記命令文により特定される対象物である確率を出力する、請求項 1 に記載の対象物検索システム。

**【請求項 3】**

前記学習済モデルを規定するパラメータは、前記第 1 の部分画像についてのクロスエントロピー損失関数と、前記第 2 の部分画像についてのクロスエントロピー損失関数とを含むコスト関数に基づいて最適化される、請求項 2 に記載の対象物検索システム。

20

**【請求項 4】**

前記学習済モデルは、

前記命令文から第 1 の特徴量を抽出する第 1 のネットワークと、

前記第 1 の部分画像および当該第 1 の部分画像の画像内環境を示す情報から第 2 の特徴量を抽出する第 2 のネットワークと、

前記第 1 の特徴量および前記第 2 の特徴量に基づいて、前記命令文により特定される対象物である確率を算出する第 3 のネットワークとを含む、請求項 1 ~ 3 のいずれか 1 項に記載の対象物検索システム。

30

**【請求項 5】**

特定の対象物に関する命令文を取得するステップと、

前記命令文に関連付けられた入力画像から、当該入力画像に含まれる個々の物体を示す 1 または複数の第 1 の部分画像を抽出するステップと、

前記命令文と、前記第 1 の部分画像の各々と、当該第 1 の部分画像の画像内環境を示す情報とを学習済モデルに入力して、前記第 1 の部分画像の各々が前記命令文により特定される対象物である確率を出力するステップとを備え、

前記学習済モデルは、画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている、対象物検索方法。

40

**【請求項 6】**

対象物検索システムを構成する学習済モデルであって、

前記学習済モデルは、

特定の対象物に関する命令文と、前記命令文に関連付けられた入力画像に含まれる個々の物体を示す 1 または複数の第 1 の部分画像の各々と、当該第 1 の部分画像の画像内環境を示す情報との入力を受けて、前記第 1 の部分画像の各々が前記命令文により特定される対象物である確率を出力するものであり、

画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている、学習済モデル。

50

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本技術は、命令文の言語理解および言語理解に基づく物体探索に関する。

## 【背景技術】

## 【0002】

お年寄りや体の不自由な人の日常生活を支援するためのニーズに対して、労働力不足などの背景もあり、生活支援ロボットなどが有効な解決手段として提案されている。例えば、家庭向けサービスロボット(DSR: domestic service robot)を標準化するような取り組みも始まっている。

10

## 【0003】

一方で、現時点においては、生活支援ロボットは、言語での対話能力を十分に有しておらず、生活支援ロボットに命令を与えるための手段は極めて限定されている。例えば、対象物検索タスク(object retrieval task)に関して、ユーザがさまざまな言語表現を用いることは難しく、生活支援ロボットは、ある限られた言語表現の範囲内でのみ命令を理解することができるといった程度である。

## 【0004】

画像および言語理解を用いて対象物を推論する技術が提案されている(非特許文献1~3)。これらの技術においては、言語情報と画像情報との間の類似性に基づいて、画像知識および言語知識を関連付けるというアプローチが採用されている。特に、非特許文献1および2は、対象物を把持するタスクに向けられており、非特許文献3は、画像内に含まれる対象物を理解するタスクに向けられている。

20

## 【先行技術文献】

## 【非特許文献】

## 【0005】

【非特許文献1】J. Hatori et al., "Interactively Picking Real-World Objects with Unconstrained Spoken Language Instructions," in IEEE ICRA, 2018, pp. 3774-3781.

【非特許文献2】M. Shridhar and D. Hsu, "Interactive visual grounding of referring expressions for human-robot interaction," in RSS, 2018.

【非特許文献3】Yu L., Tan H., Bansal M. and Berg, T. L., "A joint speaker-listener-reinforcer model for referring expressions," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2017, pp. 7282-7290.

30

【非特許文献4】J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pretraining of deep bidirectional transformers for language understanding," arXiv preprint arXiv:1810.04805, 2018.

【非特許文献5】Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al., "Google's neural machine translation system: Bridging the gap between human and machine translation," arXiv preprint arXiv:1609.08144, 2016.

【非特許文献6】K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," arXiv preprint arXiv:1409.1556, 2014.

40

【非特許文献7】A. Magassouba, K. Sugiura, and H. Kawai, "A Multimodal Classifier Generative Adversarial Network for Carry and Place Tasks From Ambiguous Language Instructions," IEEE RA-L, vol. 3, no. 4, pp. 3113-3120, Oct 2018.

【非特許文献8】K. Sugiura and H. Kawai, "Grounded Language Understanding for Manipulation Instructions Using GAN-Based Classification," IEEE ASRU, 2017.

【非特許文献9】T. Inamura, J. T. C. Tan, K. Sugiura, T. Nagai, and H. Okada, "Development of robocup@home simulation towards long-term large scale hri," in Robot Soccer World Cup. Springer, 2013, pp. 672-680.

## 【発明の概要】

50

## 【発明が解決しようとする課題】

## 【0006】

現実のアプリケーションにおいては、人間が発する言語による命令だけでは、認識対象の物体を一意に特定することはできず、不確実性が残ったものとなり得る。そのため、そのような不確実性にも対応できるシステムが要望されている。

## 【0007】

本技術は、対象となる物体の候補が複数存在するような状況であっても、対象となる物体を特定できる技術を提供することを目的とする。

## 【課題を解決するための手段】

## 【0008】

本発明のある局面に従う対象物検索システムは、特定の対象物に関する命令文を取得する命令文取得部と、命令文に関連付けられた入力画像から、当該入力画像に含まれる個々の物体を示す1または複数の第1の部分画像を抽出する画像抽出部と、命令文と、第1の部分画像の各々と、当該第1の部分画像の画像内環境を示す情報との入力を受けて、第1の部分画像の各々が命令文により特定される対象物である確率を出力する学習済モデルとを含む。学習済モデルは、画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている。

## 【0009】

画像抽出部は、命令文に関連付けられた入力画像から、いずれかの物体が存在する区域を示す1または複数の第2の部分画像をさらに抽出するように構成されていてもよい。学習済モデルは、第1の部分画像と第2の部分画像との組み合わせの各々が命令文により特定される対象物である確率を出力してもよい。

## 【0010】

学習済モデルを規定するパラメータは、第1の部分画像についてのクロスエントロピー損失関数と、第2の部分画像についてのクロスエントロピー損失関数とを含むコスト関数に基づいて最適化されてもよい。

## 【0011】

学習済モデルは、命令文から第1の特徴量を抽出する第1のネットワークと、第1の部分画像および当該第1の部分画像の画像内環境を示す情報から第2の特徴量を抽出する第2のネットワークと、第1の特徴量および第2の特徴量に基づいて、命令文により特定される対象物である確率を算出する第3のネットワークとを含んでいてもよい。

## 【0012】

第3のネットワークは、第1の特徴量および第2の特徴量の入力に対する類似性を評価する識別器と、第1の特徴量と第2の特徴量との連結結果が入力される多層パーセプトロンとを含んでいてもよい。

## 【0013】

第3のネットワークは、第1の特徴量および第2の特徴量が入力される、敵対的生成ネットワークを含んでいてもよい。

## 【0014】

敵対的生成ネットワークは、第2の特徴量についての条件を付して学習されてもよい。第1のネットワークは、命令文に対してサブワード埋め込み処理を行うレイヤと、サブワード埋め込み処理が行われた結果が入力されるリカレントニューラルネットワークとを含んでいてもよい。

## 【0015】

対象物検索システムは、命令文により特定される対象物である確率が相対的に高い複数の第1の部分画像を出力する手段と、出力された複数の第1の部分画像に対するユーザ選択に応答して、選択された第1の部分画像に対応する物体に対して物理的な作用を与えるための動作指令を生成する手段とをさらに含む。

## 【0016】

10

20

30

40

50

本発明の別の局面に従う対象物検索方法は、特定の対象物に関する命令文を取得するステップと、命令文に関連付けられた入力画像から、当該入力画像に含まれる個々の物体を示す1または複数の第1の部分画像を抽出するステップと、命令文と、第1の部分画像の各々と、当該第1の部分画像の画像内環境を示す情報と学習済モデルに入力して、第1の部分画像の各々が命令文により特定される対象物である確率を出力するステップとを含む。学習済モデルは、画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている。

#### 【0017】

本発明のさらに別の局面に従えば、対象物検索システムを構成する学習済モデルが提供される。学習済モデルは、特定の対象物に関する命令文と、命令文に関連付けられた入力画像に含まれる個々の物体を示す1または複数の第1の部分画像の各々と、当該第1の部分画像の画像内環境を示す情報との入力を受けて、第1の部分画像の各々が命令文により特定される対象物である確率を出力するものであり、画像に含まれるいずれかの物体を特定する命令文と、当該命令文により特定される物体を示す部分画像とを含むトレーニングデータセットにより学習されている。

#### 【発明の効果】

#### 【0018】

本技術によれば、対象となる物体の候補が複数存在するような状況であっても、対象となる物体を特定できる。

#### 【図面の簡単な説明】

#### 【0019】

【図1】本実施の形態に従う対象物検索システムのシステム概要を示す模式図である。

【図2】本実施の形態に従う対象物検索システムの処理内容の概略を説明するための図である。

【図3】本実施の形態に従う情報処理装置のハードウェア構成例を示す模式図である。

【図4】本実施の形態に従う対象物検索システムにおいて採用される推論モデルの処理内容を説明するための図である。

【図5】本実施の形態に従う対象物検索システムにおいて採用される推論モデルの改良された処理内容を説明するための図である。

【図6】本実施の形態に従う対象物検索システムが提供する推論結果の一例を示す模式図である。

【図7】本実施の形態に従う対象物検索システムにおいて利用されるトレーニングデータセットの生成手順を示すフローチャートである。

【図8】本実施の形態に従う対象物検索システムにおいて利用される物体認識技術の結果例を示す図である。

【図9】本実施の形態に従う対象物検索システムにおけるトレーニングの処理手順を示すフローチャートである。

【図10】本実施の形態に従う対象物検索システムにおける推論処理の処理手順を示すフローチャートである。

#### 【発明を実施するための形態】

#### 【0020】

本発明の実施の形態について、図面を参照しながら詳細に説明する。なお、図中の同一または相当部分については、同一符号を付してその説明は繰り返さない。

#### 【0021】

##### [A. システム概要]

まず、本実施の形態に従う対象物検索システム1の概要について説明する。図1は、本実施の形態に従う対象物検索システム1のシステム概要を示す模式図である。図1を参照して、対象物検索システム1は、典型的には、ロボット2と、情報処理装置100と、端末装置200とを含む。情報処理装置100は、機能モジュールとして、推論モデル40

10

20

30

40

50

および動作指令生成部 50 を含む。

【0022】

情報処理装置 100 には、ロボット 2 に設けられたカメラ 12 からの映像信号が入力される。ロボット 2 のカメラ 12 の視野は、ユーザの視野と少なくとも一部は重複しているものとする。この状況に応じて、ユーザが端末装置 200 に向けて任意の言語命令（発話による命令文）を与えると、端末装置 200 を介して情報処理装置 100 へ音声信号が入力される。

【0023】

情報処理装置 100 の推論モデル 40 は、学習済モデルであり、端末装置 200 からの音声信号と、ロボット 2 のカメラ 12 からの映像信号との入力を受けて、推論結果を算出する。推論結果は、カメラ 12 により撮像された画像内に含まれる 1 または複数の対象物について、ユーザからの言語命令により指定された対象物である確率を含む。

10

【0024】

情報処理装置 100 の動作指令生成部 50 は、推論モデル 40 からの推論結果を受けて、端末装置 200 に操作対象の対象物の候補を表示するとともに、対象物の候補のうち、ユーザにより選択された対象物を示す選択指令を端末装置 200 から受付ける。動作指令生成部 50 は、選択指令に従って、対象の対象物を決定するとともに、対応する動作指令を生成して、ロボット 2 へ出力する。ロボット 2 は、動作指令に従って対象物に対する物理的な作用を与える作用部 14 を駆動する。

20

【0025】

このように、本実施の形態に従う対象物検索システム 1 においては、ユーザによる言語命令に応じて、画像内に存在する各対象物が言語命令によって指定された対象物である確率を推論する。対象物検索システム 1 は、このような推論結果を用いることで、ロボット 2 に対して、適切な動作指令を与えることができる。

【0026】

次に、本実施の形態に従う対象物検索システム 1 における処理内容の概略を説明する。図 2 は、本実施の形態に従う対象物検索システム 1 の処理内容の概略を説明するための図である。図 2 を参照して、対象物検索システム 1 においては、ユーザが発した音声信号 20 およびロボット 2 に設けられたカメラ 12 などにより撮像された入力画像 30 が取得される。

30

【0027】

推論モデル 40 には、音声信号 20 に対応する命令文 22 が入力される。命令文 22 は、音声信号 20 を公知の音声認識することでテキスト化することで生成できる。なお、音声信号 20 に代えて、ユーザがキーボードといった任意の入力デバイスを用いて、命令文 22 の内容を示すテキストを入力するようにしてもよい。情報処理装置 100 は、命令文 22 を取得する命令文取得機能として、音声認識の機能あるいはテキストベースの命令文 22 を受付ける機能を有している。

【0028】

図 2 には、一例として、「Bring me the toy on the wagon」（そのワゴン上のそのおもちゃを取って）といった命令文 22 を示す。このように、命令文 22 は、特定の対象物に関するものであるとする。

40

【0029】

入力画像 30 は、命令文 22 が発せられたシーンを示すものである。入力画像 30 からは、命令文 22 により操作の対象となり得る物体を示す部分画像（以下、「ターゲット画像」あるいは「ターゲット」とも称す。）と、操作の対象となり得る物体が存在し得る区域を示す部分画像（以下、「ソース画像」あるいは「ソース」とも称す。）とが抽出される。

【0030】

本明細書において、「命令文」は、任意の対象物に対する任意の操作を命令するものを意味する。「操作」の内容としては、例えば、対象物を「把持する」、「取る」、「しま

50

う」、「置く」、「移動する」といった動作が想定される。但し、これら列挙した動作に限らず、任意の操作を対象とし得る。

【0031】

図2には、入力画像30から複数のターゲット32および複数のソース34が抽出されている例を示す。推論モデル40には、入力画像30から抽出された1または複数のターゲット32からなるターゲット候補36と、入力画像30から抽出された1または複数のソース34からなるソース候補38とが入力される。

【0032】

推論モデル40は、命令文22、ターゲット候補36およびソース候補38の入力を受けて、命令文22による操作の対象物である「たしからしさ」(likelihood)を、ターゲット32とソース34との各組み合わせ52について算出する。図2に示す例では、算出される「たしからしさ」は、「そのワゴン上のそのおもちゃ」と指定されている対象物である確率を意味する。

10

【0033】

このように、学習済モデルである推論モデル40は、ターゲット32とソース34との組み合わせの各々が命令文22により特定される対象物である確率を出力する。但し、ソース候補38として1つのソース34のみが含まれる場合には、情報処理装置100は、ターゲット32の各々が命令文22により特定される対象物である確率を出力することになる。

【0034】

さらに、ターゲット32とソース34との組み合わせ52のうち、たしからしさが上位のものだけが推論結果54として出力されてもよい。推論結果54の内容は、端末装置200に表示され、ユーザから最終的な選択を受け付けるようにしてもよい。

20

【0035】

[B. 情報処理装置100のハードウェア構成]

次に、本実施の形態に従う情報処理装置100のハードウェア構成の一例について説明する。図3は、本実施の形態に従う情報処理装置100のハードウェア構成例を示す模式図である。情報処理装置100は、典型的には、汎用コンピュータを用いて実現される。

【0036】

図3を参照して、情報処理装置100は、主要なコンポーネントとして、プロセッサ102と、主メモリ104と、ディスプレイ106と、入力デバイス108と、ネットワークインターフェイス(I/F: interface)110と、光学ドライブ112と、入力インターフェイス(I/F)114と、出力インターフェイス(I/F)116と、二次記憶装置120とを含む。これらのコンポーネントは、内部バス118を介して互いに接続される。

30

【0037】

プロセッサ102は、後述するような各種プログラムを実行することで、後述するような処理および機能を実現する演算主体であり、例えば、1または複数のCPU(Central Processing Unit)やGPU(Graphics Processing Unit)などで構成される。複数のコアを有するようなCPUまたはGPUを用いてもよい。

40

【0038】

主メモリ104は、プロセッサ102がプログラムを実行するにあたって、プログラムコードやワークメモリなどを一時的に格納する記憶領域であり、例えば、DRAM(Dynamic Random Access Memory)やSRAM(Static Random Access Memory)などの揮発性メモリデバイスなどで構成される。

【0039】

ディスプレイ106は、処理に係るユーザインターフェイスや処理結果などを出力する表示部であり、例えば、LCD(Liquid Crystal Display)や有機EL(Electroluminescence)ディスプレイなどで構成される。

【0040】

50

入力デバイス 108 は、ユーザからの命令や操作などを受付けるデバイスであり、例えば、キーボード、マウス、タッチパネル、ペンなどで構成される。また、入力デバイス 108 としては、機械学習に必要な音声を収集するためのマイクロフォンを含んでいてもよいし、機械学習に必要な音声を収集した集音デバイスと接続するためのインターフェイスを含んでいてもよい。

#### 【0041】

ネットワークインターフェイス 110 は、インターネット上またはイントラネット上の任意の情報処理装置などとの間でデータを遣り取りする。ネットワークインターフェイス 110 としては、例えば、イーサネット（登録商標）、無線 LAN（Local Area Network）、Bluetooth（登録商標）などの任意の通信方式を採用できる。

10

#### 【0042】

光学ドライブ 112 は、CD-ROM（Compact Disc Read Only Memory）、DVD（Digital Versatile Disc）などの光学ディスク 112 M に格納されている情報を読み出して、内部バス 118 を介して他のコンポーネントへ出力する。光学ディスク 112 M は、非一過的（non-transitory）な記録媒体の一例であり、任意のプログラムを不揮発的に格納した状態で流通する。光学ドライブ 112 が光学ディスク 112 M からプログラムを読み出して、二次記憶装置 120 などにインストールすることで、コンピュータにより対象物検索システム 1 の機能を提供できるようになる。したがって、本発明の主題は、二次記憶装置 120 などにインストールされたプログラム自体、または、本実施の形態に従う機能や処理を実現するためのプログラムを格納した光学ディスク 112 M などの記録媒体でもあり得る。

20

#### 【0043】

図 3 には、非一過的な記録媒体の一例として、光学ディスク 112 M などの光学記録媒体を示すが、これに限らず、フラッシュメモリなどの半導体記録媒体、ハードディスクまたはストレージテープなどの磁気記録媒体、MO（Magneto-Optical disk）などの光磁気記録媒体を用いてもよい。

#### 【0044】

入力インターフェイス 114 は、カメラなどの外部デバイスと接続され、カメラにより撮像された映像信号を取込む。出力インターフェイス 116 は、ロボット 2 などの外部デバイスと接続され、操作可能性の推論結果およびユーザからの命令などに基づいて、必要な動作指令をロボット 2 へ出力する。入力インターフェイス 114 および出力インターフェイス 116 は、USB（Universal Serial Bus）などの汎用的な通信インターフェイスを用いることができる。

30

#### 【0045】

二次記憶装置 120 は、プロセッサ 102 にて実行されるプログラム、後述するようなモデル（ニューラルネットワーク）をトレーニングするためのトレーニングデータセット、および、モデルを規定するパラメータなどを格納するコンポーネントであり、例えば、ハードディスク、SSD（Solid State Drive）などの不揮発性記憶装置で構成される。

#### 【0046】

より具体的には、二次記憶装置 120 は、図示しない OS（Operating System）の他、音声認識プログラム 121 と、画像抽出プログラム 122 と、トレーニングプログラム 123 と、動作指令生成プログラム 124 と、モデルパラメータ 125 とを格納している。また、二次記憶装置 120 には、トレーニングデータセット 126 が格納されていてもよい。

40

#### 【0047】

音声認識プログラム 121 は、後述するように、音声信号 20 に対応する命令文 22 を生成する。画像抽出プログラム 122 は、入力画像に含まれる部分画像を抽出する（図 4 画像抽出部 403 に対応）。トレーニングプログラム 123 は、126 を用いて、推論モデル 40 を規定するパラメータを最適化する。動作指令生成プログラム 124 は、動作指

50

令生成部 50 (図 1) を実現する。モデルパラメータ 125 は、学習済モデルである推論モデルを規定する 1 または複数のパラメータを含む。トレーニングデータセット 126 は、推論モデルを最適化するための教師データであり、後述するようなデータの組からなる。

#### 【0048】

これらのプログラムをプロセッサ 102 で実行する際に必要となるライブラリや機能モジュールの一部を、OS が標準で提供するライブラリまたは機能モジュールを用いて代替するようにしてもよい。この場合には、各プログラム単体では、対応する機能を実現するために必要なプログラムモジュールのすべてを含むものにはならないが、OS の実行環境下にインストールされることで、必要な機能を実現できる。このような一部のライブラリまたは機能モジュールを含まないプログラムであっても、本発明の技術的範囲に含まれ得る。

10

#### 【0049】

また、これらのプログラムは、上述したようないずれかの記録媒体に格納されて流通するだけでなく、インターネットまたはイントラネットを介してサーバ装置などからダウンロードすることで配布されてもよい。

#### 【0050】

図 3 には、単一のコンピュータが情報処理装置 100 を構成する例を示すが、これに限らず、コンピュータネットワークを介して接続された複数のコンピュータが明示的または黙示的に連携して、情報処理装置 100 を含む対象物検索システム 1 を実現するようにしてもよい。複数のコンピュータが連携する場合、一部のコンピュータがいわゆるクラウドコンピュータと称される、ネットワーク上にある不特定のコンピュータであってもよい。

20

#### 【0051】

コンピュータ (プロセッサ 102) がプログラムを実行することで実現される機能の全部または一部を、集積回路などのハードワイヤード回路 (hard-wired circuit) を用いて実現してもよい。例えば、ASIC (Application Specific Integrated Circuit) や FPGA (Field-Programmable Gate Array) などを用いて実現してもよい。

#### 【0052】

当業者であれば、本発明が実施される時代に応じた技術を適宜用いて、適切なハードウェア構成を採用するであろう。

30

#### 【0053】

[C. 推論モデル 40 (MTCM)]

(c1: 概要)

次に、本実施の形態に従う対象物検索システム 1 において採用される推論モデル 40 の概要について説明する。

#### 【0054】

本実施の形態においては、推論モデル 40 として、MTCM (Multimodal Target-source Classifier Model) と略称するモデルを採用した場合を示す。MTCM は、上述した非特許文献 1 ~ 3 に開示されるような、マルチモーダル類似性ベースの統合アプローチ (multimodal similarity-based integration approach) の改良である。

40

#### 【0055】

図 4 は、本実施の形態に従う対象物検索システム 1 において採用される推論モデル 40 の処理内容を説明するための図である。図 4 を参照して、命令文 22 を処理するネットワークである命令文処理部 401 には、マルチレイヤ双方向 LSTM (Long short-term memory) を採用する。併せて、推論精度を高める目的で、BERT モデル (非特許文献 4 など参照) を用いて、サブワード埋め込み処理を付加している。

#### 【0056】

画像情報を処理するネットワークには、CNN (Convolutional Neural Network) モデル 404 を採用している。

#### 【0057】

50

より具体的には、画像抽出部 403 が入力画像 30 を処理することで、ターゲット候補 36 に含まれる  $i$  番目 ( $i \in \{1, \dots, N\}$ ) のターゲット 32、および、ソース候補 38 に含まれる  $i'$  番目 ( $i' \in \{1, \dots, M\}$ ) のソース 34 の各々について、部分画像および画像内の位置が取得される。入力画像 30 は、命令文 22 に関連付けられたものである。そして、画像抽出部 403 は、入力画像 30 に含まれる個々の物体を示す 1 または複数のターゲット 32 (部分画像) を抽出するとともに、いずれかの物体が存在する区域を示す 1 または複数のソース 34 (部分画像) を抽出する。

【0058】

任意のターゲット 32 について、入力データセット  $x(i)$  を以下の (1) 式のように示すことができる。

【0059】

$$x(i) = \{x_{i n s}(i), x_v(i), x_{r e l}(i)\} \dots (1)$$

但し、 $x_{i n s}$  は言語特徴量を示し、 $x_v$  は画像特徴量を示し、 $x_{r e l}$  は関連性特徴量を示す。以下の説明においては、添え字  $i$  を省略して、「入力データセット  $x$ 」とも記載する。

【0060】

画像特徴量  $x_v$  は、ターゲット 32 として抽出された部分画像に対応する。画像特徴量  $x_v$  は、CNN モデル 404 によって処理される。関連性特徴量  $x_{r e l}$  は、各ターゲット 32 の画像内環境 (例えば、他のターゲットとの相対関係、入力画像内の位置、ソースに対する位置など) を示す情報である。

【0061】

画像特徴量  $x_v$  の処理と並行的に、言語特徴量  $x_{i n s}$  は、埋め込み処理がされた上で、マルチレイヤ双方向 LSTM によりエンコーディングされる。

【0062】

言語特徴量  $x_{i n s}$  および画像特徴量  $x_v$  を処理することで得られるそれぞれの潜在表現 (latent representation) 同士を比較するために、3つの MLP (多層パーセプトロン: multilayer perceptron) が配置されている。

【0063】

最終的に、推論モデル 40 からは、ターゲット 32 毎の「たしからしさ」を示す推論結果が出力される。このように、推論モデル 40 には、命令文 22 と、ターゲット 32 の各々と、ターゲット 32 の画像内環境を示す情報 (関連性特徴量  $x_{r e l}$ ) との入力を受けて、ターゲット 32 の各々が命令文 22 により特定される対象物である確率を出力する。

【0064】

(c2: 命令文 22 を処理するネットワーク)

次に、推論モデル 40 の命令文 22 を処理するネットワークについて説明する。図 4 を参照して、命令文処理部 401 には、音声信号 20 に対応する命令文 22 が入力され、入力された命令文 22 は、リカレントニューラルネットワークの一例としてのマルチレイヤ双方向 LSTM により処理されて、抽出された非音声特徴量  $o_I$  が MLP 402 に入力される。

【0065】

推論モデル 40 においては、マルチレイヤ双方向 LSTM の前段に、サブワード埋め込みモデルである BERT モデルが配置されている。サブワードモデルを用いて、マルチレイヤ双方向 LSTM に入力する埋め込みベクトルを初期化する。すなわち、命令文 22 を処理するネットワークは、命令文 22 に対してサブワード埋め込み処理を行うレイヤ (典型例として、BERT モデル) と、記サブワード埋め込み処理が行われた結果が入力されるリカレントニューラルネットワーク (典型例として、マルチレイヤ双方向 LSTM) とを含む。

【0066】

BERT モデルは、双方向トランスフォーマに基づく言語エンコーディングモデルである。BERT モデルを用いることで、フレキシブル性およびロバスト性を高めることがで

10

20

30

40

50

きる。現在利用できるBERTモデルは、35億個のワードを用いてトレーニングされているため、頻出頻度の少ないワードであってもデータのスパース性は問題にはならない。

【0067】

また、BERTモデルは、ワードベースのトークナイゼーション(ワードトークン)ではなく、サブワードのトークナイゼーション(サブワードトークン)を用いる。サブワードのトークナイゼーションは、ワードの一部を用いたトークナイゼーションを意味する。例えば、以下のTable 1に示すように、頻出頻度の少ないワードやミススペルされたワードに対しても、頻出頻度の高いワードを用いたトークンを生成できる。

【0068】

【表1】

10

Table 1

Expression	(a)ワードトークン	(b)サブワードトークン
topright object	topright, object	top, right, object
sprayer	<UNK>	spray, er
greyis bottle	<UNK>, bottle	grey, is, bottle

【0069】

なお、Table 1において、<UNK>は、トークンを生成できないことを意味する。なお、埋め込みモデル(BERTモデル)は、推論モデル40がトレーニングされるに伴って、微調整されることになる。

20

【0070】

BERTモデルから出力されるトークンがマルチレイヤ双方向LSTMに入力される。マルチレイヤ双方向LSTMは、公知技術であるので、ここでは詳細な説明は行わない。

【0071】

さらに、マルチレイヤ双方向LSTMからの出力は、MLP402に入力される。MLP402からは、入力された命令文22の非音声特徴量oIが出力される。

【0072】

このように、推論モデル40の命令文22を処理するネットワークは、命令文処理部401およびMLP402を含み、命令文22から非音声特徴量oIを抽出する。

【0073】

30

(c3: 画像情報を処理するネットワーク)

次に、推論モデル40の画像情報を処理するネットワークについて説明する。

【0074】

CNNモデル404としては、例えば、非特許文献6に示されるような16層ネットワーク(VGG16)を用いて、画像特徴量をエンコーディングできる。CNNモデル404からの出力は、連結部405において関連性特徴量x<sub>rel</sub>と連結される。

【0075】

ターゲット候補36に含まれるN個のターゲット32の各々と、対応する関連性特徴量x<sub>rel</sub>とについて、連結部405による連結結果が出力される。そして、すべてのターゲット32についての連結結果がMLP406に入力される。MLP406からは、ターゲット候補36に含まれる複数のターゲット32についての画像特徴量oVが出力される。

40

【0076】

このように、推論モデル40の画像情報を処理するネットワークは、ターゲット32およびターゲット32の画像内環境を示す情報である関連性特徴量x<sub>rel</sub>から画像特徴量oVを抽出する。

【0077】

(c4: 推論モデル40の推論結果を生成する出力部410)

推論モデル40の推論結果Yは、以下の(2)式のように示すことができる。

【0078】

50

$$Y = \{ y_{targ}, y_{src} \} \dots (2)$$

但し、 $y_{targ}$  はターゲットについての推論結果を示し、 $y_{src}$  はソースについての推論結果を示す。推論結果  $y_{targ}$  および  $y_{src}$  は、いずれも  $N \times M$  次元のベクトルとして規定される。

【0079】

ターゲットについての推論結果  $y_{targ}$  は、入力された命令文 22 の非音声特徴量  $o_I$  と、ターゲット候補 36 に含まれる複数のターゲット 32 についての画像特徴量  $o_V$  とが類似性識別器 407 に入力されることで算出される。

【0080】

ソースについての推論結果  $y_{src}$  は、入力された命令文 22 の非音声特徴量  $o_I$  と、ターゲット候補 36 に含まれる複数のターゲット 32 についての画像特徴量  $o_V$  とが連結部 408 において連結された結果が M L P 409 に入力されることで算出される。

10

【0081】

このように、出力部 410 は、非音声特徴量  $o_I$  および画像特徴量  $o_V$  に基づいて、各ターゲット 32 が命令文 22 により特定される対象物である確率を算出するネットワークである。より具体的には、このネットワークは、非音声特徴量  $o_I$  および画像特徴量  $o_V$  の入力に対する類似性を評価する類似性識別器 407 と、非音声特徴量  $o_I$  と画像特徴量  $o_V$  との連結結果が入力される M L P 408 とを含む。

【0082】

(c5: 推論モデル 40 のトレーニング)

20

推論モデル 40 のコスト関数 J M T C M は、以下の (3) 式のように定義できる。

【0083】

$$J M T C M = \alpha_1 J_{targ} + \alpha_2 J_{src} \dots (3)$$

但し、 $\alpha_1$  および  $\alpha_2$  は重みパラメータであり、 $J_{targ}$  は、ターゲット 32 についてのクロスエントロピー損失関数であり、 $J_{src}$  は、ソース 34 についてのクロスエントロピー損失関数である。クロスエントロピー損失関数  $J_{targ}$  および  $J_{src}$  は、以下の (4-1) および (4-2) 式のように定義できる。

【0084】

【数 1】

$$J_{targ} = -\sum_n \sum_m y_{targ\_nm}^* \log p(y_{targ\_nm}) \dots (4-1)$$

30

$$J_{src} = -\sum_n \sum_m y_{src\_nm}^* \log p(y_{src\_nm}) \dots (4-2)$$

【0085】

但し、 $y_{targ\_nm}^*$  および  $y_{src\_nm}^*$  は、 $n$  番目のサンプルの  $m$  番目の次元についてのラベル (正解) を示し、 $y_{targ\_nm}$  および  $y_{src\_nm}$  は、 $n$  番目のサンプルの  $m$  番目の次元についての推論結果を示す。

【0086】

図 4 に示される推論モデル 40 は、予め用意されたトレーニングデータセット 126 に対して、上述した (3) 式で定義されるコスト関数 J M T C M が最小になるようにパラメータを最適化することで構成される。このように、推論モデル 40 を規定するパラメータは、ターゲット 32 についてのクロスエントロピー損失関数  $J_{targ}$  と、ソース 34 についてのクロスエントロピー損失関数  $J_{src}$  とを含むコスト関数に基づいて最適化されることになる。

40

【0087】

[D: 推論モデル 40 (M T C M - G A N)]

(d1: 概略)

上述した M T C M からなる推論モデル 40 に対して、敵対的生成ネットワーク (G A N : generative adversarial nets) を付加することで、トレーニングデータを増大させて

50

、識別性能を高めることもできる。以下、MTCMおよびGANからなる推論モデル40（以下、「MTCM-GAN」とも記載する。）について説明する。

【0088】

先に、GANについて概略する。GANフレームワークは、生成器(generator)Gおよび識別器(discriminator)Dの2つの敵対的ネットワークで構成される。生成器Gは、所与の分布データを模倣することで疑似データを生成する。並行して、識別器Dは、入力データが真(real)であるか偽(fake)であるかを識別(推論)する。これらのネットワークの目的として、生成器Gはより真に近いデータを生成するようになり、一方で、識別器Dはその識別能力(推論能力)を向上させる。

【0089】

生成器Gは、正規分布からランダムにサンプルされた複数次元の入力z(ノイズ)を用いて、疑似サンプル $x_{fake}$ を生成する。真正サンプル $x_{real}$ と疑似サンプル $x_{fake}$ とを識別するために、ソースフラグ $S \in \{real, fake\}$ に応じて、入力 $x = x_{real}$ または $x = x_{fake}$ が識別器Dには選択的に入力される。識別器Dの出力は、推論確率 $p_D(S = real | x) = D(x)$ となる。生成器Gの損失関数 $J_G$ および識別器Dの損失関数 $J_D$ は、以下の(5-1)式で定義される損失関数 $J_S$ を用いて、(5-2)および(5-3)式のように定義できる。これらの損失関数を用いて、GANのパラメータが最適化される。

【0090】

【数2】

$$J_S = -\frac{1}{2} E_{x_{real}} \log D(x_{real}) - \frac{1}{2} E_z \log(1 - D(x_{fake})) \quad \dots(5-1)$$

$$J_D = J_S \quad \dots(5-2)$$

$$J_G = -J_S \quad \dots(5-3)$$

【0091】

(d2: MTCM-GAN)

GANのデータ増大特性を利用して、推論モデル40にデータ増大および識別を同時に行う機能を付加した改良例について説明する(非特許文献7および非特許文献8など参照)。

【0092】

図5は、本実施の形態に従う対象物検索システム1において採用される推論モデル40の改良された処理内容を説明するための図である。図5を参照して、MTCM-GANは、図4に示す推論モデル40の推論結果を生成する出力部410に代えて、出力部420を有している。推論モデル40の出力部420以外の部分は、図4と同様であるので、詳細な説明は繰り返さない。

【0093】

出力部420は、連結部421、422と、生成器423と、選択部424と、識別器425とを含む。生成器423および識別器425が敵対的ネットワークを構成する。出力部420においては、生成器423により生成される疑似データは識別器425の識別能力を向上させる。識別器425は、真正サンプル $x_{real}$ と疑似サンプル $x_{fake}$ とを識別するだけでなく、候補となるターゲットの「たしからしさ」を推論することで識別タスクも実行することになる。

【0094】

そのため、出力部420の識別器425は、推論確率 $p_D(S)$ に加えて、第2の出力として、各ターゲットが命令文22による操作の対象物である確率を示す推論確率 $p_D(y_{target})$ を出力する。

【0095】

また、識別器425のコスト関数 $J_D$ は、以下の(6)式のように定義できる。

10

20

30

40

50

$$JD = JS + J \dots (6)$$

但し、 $J$  は重みパラメータであり、 $J$  は上述した (5 - 1) 式において定義したクロスエントロピー損失関数である。

【0096】

MTCMの推論モデル40の初期状態を考慮すると、図5に示される出力部420の識別器425へ入力されるデータセットは、以下の(7)式のように設定できる。

【0097】

$$xGAN = \{ x_{real} = (oV, oI), x_{fake} = (z, oV) \} \dots (7)$$

このように、出力部420は、非音声特徴量 $oI$ および画像特徴量 $oV$ が入力される、敵対的生成ネットワークを含む。敵対的生成ネットワークを用いることで、識別性能を高めることもできる。

【0098】

(d3: 条件付きMTCM-GAN)

図5に示すMTCM-GANにおいて、生成器423に対して画像特徴量 $oV$ についての条件を付してもよい。これは、生成器423および識別器425が全結合ネットワークであるとともに、画像特徴量 $oV$ が、生成器423(真正サンプル $x_{real}$ を通じて)および識別器425の両方に入力されるからである。すなわち、敵対的生成ネットワーク(出力部420)は、画像特徴量 $oV$ についての条件を付して学習されてもよい。

【0099】

より具体的には、識別器425は、入力ソースがいずれであるかを推論するのに加えて、組み合わせ( $oI, oV$ )が正しいものであるか否かを推論する。一方、生成器423は、正解/不正解の組み合わせ( $oI, oV$ )を疑似データとして生成する。そこで、同一のシーンのランダムに選択された不正解のターゲット $j$ を考慮しつつ、各ターゲット $i$ について、正解の特徴量の組み合わせ( $oI(i), oV(i)$ )および不正解の特徴量の組み合わせ( $oI(i), oV(j)$ )を用意し、これを用いてトレーニングを行うようにしてもよい。

【0100】

[E. 推論結果の出力例]

次に、本実施の形態に従う対象物検索システム1による推論結果の出力例について説明する。本実施の形態に従う対象物検索システム1は、ターゲット32とソース34との組み合わせの各々について、命令文22による操作の対象物である「たしからしさ」(すなわち、確率)を算出できる。

【0101】

ターゲット32とソース34との組み合わせについての確率、あるいは、ターゲット32についての確率をユーザに提示するようにしてもよい。

【0102】

図6は、本実施の形態に従う対象物検索システム1が提供する推論結果の一例を示す模式図である。図6を参照して、対象物検索システム1は、端末装置200のスクリーン上などに推論結果を含む結果表示300を提供することができる。結果表示300は、シーン(入力画像30)から抽出されたターゲット32の各々について、命令文22による操作の対象物である確率60が表示されている。

【0103】

図6に示すように、結果表示300にはターゲット32毎の確率60が表示されているので、ユーザは、意図したターゲット32をより容易に選択できる。また、算出される確率を定量的に評価できるので、しきい値などの条件に基づいて、対象となるターゲット32を自動的に選択することができる。

【0104】

[F. 処理手順]

次に、本実施の形態に従う対象物検索システム1における処理手順について説明する。

10

20

30

40

50

## 【0105】

( f 1 : トレーニングデータセット 1 2 6 の生成 )

図 7 は、本実施の形態に従う対象物検索システム 1 において利用されるトレーニングデータセット 1 2 6 の生成手順を示すフローチャートである。図 7 に示す各ステップは、コンピュータにより実行されてもよいし、一部をユーザ自身が実行してもよい。

## 【0106】

図 7 を参照して、ユーザは、シーンを示す 1 または複数の入力画像を取得する ( ステップ S 1 0 0 )。入力画像は、現実の室内を撮像することで取得してもよいし、画像共有サイトなどから任意にダウンロードすることで取得してもよい。

## 【0107】

続いて、取得された 1 または複数の入力画像のうち 1 つを選択し ( ステップ S 1 0 2 )、選択された入力画像から物体を示す 1 または複数の領域 ( ターゲット 3 2 ) を抽出する ( ステップ S 1 0 4 )。抽出された 1 または複数のターゲット 3 2 のうち 1 つを選択し ( ステップ S 1 0 6 )、選択されたターゲット 3 2 が存在する領域を示す部分画像 ( ソース 3 4 ) を抽出する ( ステップ S 1 0 8 )。さらに、選択されたターゲット 3 2 に対して、対応する物体の名称をラベルとして付与する ( ステップ S 1 1 0 )。併せて、選択されたターゲット 3 2 とシーンとの関連性を示す情報 ( 関連性特徴量  $x_{r e l}$  ) を設定する ( ステップ S 1 1 2 )。

## 【0108】

ステップ S 1 0 4 および S 1 1 0 に関して、入力画像に対して公知の物体認識技術を用いて自動的にターゲット 3 2 となり得る領域を抽出するようにしてもよい。

## 【0109】

図 8 は、本実施の形態に従う対象物検索システム 1 において利用される物体認識技術の結果例を示す図である。図 8 に示すように、入力された画像に対して、物体が存在する領域が特定および抽出されるとともに、特定された物体を示すラベル ( 例えば、物品名 ) が自動的に抽出される。このような抽出結果を用いて、ターゲット 3 2 および対応するラベルのデータセットを自動的に生成できる。

## 【0110】

このような物体認識技術としては、SSD ( Single Shot MultiBox Detector ) や YOLO ( You Only Look Once ) などのアルゴリズムを用いることができる。

## 【0111】

あるいは、領域抽出およびラベル付与を手動で行うようにしてもよい。さらには、公知のアルゴリズムを用いて自動的に領域を抽出した上で、手動でラベルを付与するようにしてもよい。

## 【0112】

また、ソース 3 4 については、ターゲット 3 2 を囲むような領域を抽出するようにしてもよい。

## 【0113】

ステップ S 1 1 2 に関して、抽出したターゲット 3 2 とシーンとの関連性を示す情報としては、例えば、「右下」や「左上」といった自然言語表現であってもよいし、位置の情報を示す符号であってもよい。

## 【0114】

さらに、選択されたターゲット 3 2 に関する 1 または複数の命令文 2 2 を取得する ( ステップ S 1 1 4 )。1 または複数の命令文 2 2 は、ユーザが任意に考えて設定してもよい。例えば、いわゆるクラウドワークに対して入力画像を提供するとともに、対応する 1 または複数の命令文 2 2 を応答してもらうような形態が想定できる。

## 【0115】

ステップ S 1 0 6 ~ S 1 1 4 の処理によって、1 つの入力画像に含まれる 1 つのターゲット 3 2 に対応付けられる、ラベル、ソース 3 4、関連性特徴量  $x_{r e l}$ 、命令文 2 2 からなるデータセットを取得できる。

10

20

30

40

50

## 【0116】

選択された入力画像に含まれるターゲット32のすべてについて処理が完了したか否かが判断される(ステップS116)。選択された入力画像に含まれるターゲット32のうち処理が完了していないものがあれば(ステップS116においてNO)、新たなターゲット32が選択され(ステップS118)、ステップS108以下の処理が繰り返される。

## 【0117】

選択された入力画像に含まれるすべてのターゲット32について処理が完了していれば(ステップS116においてYES)、取得された入力画像のすべてについて処理が完了したか否かが判断される(ステップS120)。取得された入力画像のうち処理が完了していないものがあれば(ステップS120においてNO)、新たな入力画像が選択され(ステップS122)、ステップS104以下の処理が繰り返される。

10

## 【0118】

取得されたすべての入力画像について処理が完了していれば(ステップS120においてYES)、ステップS106~S114の処理によって得られるデータセットがトレーニングデータセット126として出力される(ステップS124)。そして、処理は終了する。

## 【0119】

図7に示すトレーニングデータセット126の生成手順によれば、各入力画像に含まれる各ターゲット32に対応付けられる、ラベル、ソース34、関連性特徴量 $x_{rel}$ 、命令文22からなるデータセットを取得できる。各トレーニングデータセット126には、ターゲット32およびソース34の位置および大きさの情報を含めるようにしてもよい。

20

## 【0120】

(f2: 推論モデル40のトレーニング)

図9は、本実施の形態に従う対象物検索システム1におけるトレーニングの処理手順を示すフローチャートである。図9に示す各ステップは、情報処理装置100のプロセッサ102がトレーニングプログラム123を実行することで実現されてもよい。

## 【0121】

図9を参照して、情報処理装置100は、予め用意されたトレーニングデータセット126のうち1つのデータセットを選択し(ステップS200)、選択されたデータセットから入力データセット $x$ および対応する正解ラベル( $y_{target}^*$ および $y_{src}^*$ など)を生成する(ステップS202)。なお、MTCM-GANからなる推論モデル40を採用する場合には、正解ラベルとして、非音声特徴量 $o_I$ および画像特徴量 $x_v$ を用いて、推論確率 $p_D(S)$ および推論確率 $p_D(y_{target})$ を算出してもよい。

30

## 【0122】

情報処理装置100は、予め用意されたトレーニングデータセット126のすべてについての処理が完了したか否かを判断する(ステップS204)。予め用意されたトレーニングデータセット126のうち処理が完了していないものがあれば(ステップS204においてNO)、新たなデータセットが選択され(ステップS206)、ステップS202以下の処理が繰り返される。

40

## 【0123】

予め用意されたトレーニングデータセット126に含まれるすべてのデータセットについて処理が完了していれば(ステップS204においてYES)、情報処理装置100は、生成された入力データセット $x$ を推論モデル40に入力するとともに、算出される推論結果と対応する正解ラベルとの誤差に基づいて、推論モデル40のモデルパラメータを最適化する(ステップS208)。すなわち、学習済モデルである推論モデル40は、入力画像に含まれるいずれかの物体を特定する命令文22と、命令文22により特定されるターゲット32(物体を示す部分画像)とを含むトレーニングデータセット126により学習されることで、生成される。

## 【0124】

50

より具体的には、推論モデル 40 のモデルパラメータの最適化には、上述したようなクロスエントロピー損失関数が用いられる。

【0125】

なお、バッチノーマリゼーションやドロップアウトなどの公知の加速化手法を採用できる。

【0126】

(f3: 推論モデル 40 を用いた推論処理)

図 10 は、本実施の形態に従う対象物検索システム 1 における推論処理の処理手順を示すフローチャートである。図 10 に示す各ステップは、情報処理装置 100 のプロセッサ 102 が各種プログラムを実行することで実現されてもよい。

10

【0127】

図 10 を参照して、情報処理装置 100 は、端末装置 200 から音声信号 20 が入力されると (ステップ S300)、入力された音声信号 20 を音声認識してテキストベースの命令文 22 を取得する (ステップ S302)。このように、情報処理装置 100 は、特定の対象物に関する命令文 22 を取得する。

【0128】

並行して、情報処理装置 100 は、命令文 22 に関連付けられた入力画像 30 から、入力画像 30 に含まれる個々の物体を示す 1 または複数のターゲット 32 (部分画像) を抽出するとともに、いずれかの物体が存在する区域を示す 1 または複数のソース 34 (部分画像) を抽出する。より具体的には、情報処理装置 100 は、入力画像 30 を取得し (ステップ S304)、取得した入力画像 30 から 1 または複数のターゲット 32 および 1 または複数のソース 34 を抽出する (ステップ S306)。

20

【0129】

情報処理装置 100 は、ステップ S306 において抽出したいずれかのターゲット 32 といずれかのソース 34 との組み合わせを選択する (ステップ S308) とともに、選択した組み合わせにおける関連性特徴量  $x_{r_e_1}$  を決定する (ステップ S310)。

【0130】

そして、情報処理装置 100 は、命令文 22 と、ターゲット 32 の各々と、ターゲット 32 の画像内環境を示す情報とを学習済モデルである推論モデル 40 に入力して、ターゲット 32 の各々が命令文 22 により特定される対象物である確率を出力する。すなわち、情報処理装置 100 は、情報処理装置 100 は、命令文 22 (ステップ S302)、選択した組み合わせを構成するターゲット 32 (ステップ S308) および選択した組み合わせにおける関連性特徴量  $x_{r_e_1}$  (ステップ S310) を推論モデル 40 に入力し、選択した組み合わせについての推論結果を算出する (ステップ S312)。

30

【0131】

情報処理装置 100 は、ステップ S306 において抽出したターゲット 32 とソース 34 とのすべての組み合わせについて推論結果の算出が完了したか否かを判断する (ステップ S314)。推論結果の算出が完了していない組み合わせが存在していれば (ステップ S314 において NO)、情報処理装置 100 は、ターゲット 32 とソース 34 との新たな組み合わせを選択し (ステップ S316)、ステップ S310 以下の処理を実行する。

40

【0132】

すべての組み合わせについて推論結果の算出が完了していれば (ステップ S314 において YES)、情報処理装置 100 は、推論結果のスコアが上位の 1 または複数の組み合わせを選択して、選択された組み合わせを含む画面を端末装置 200 に表示する (ステップ S318)。このように、情報処理装置 100 は、命令文 22 により特定される対象物である確率が相対的に高い複数のターゲット 32 (および、対応するソース 34) を出力する。

【0133】

情報処理装置 100 は、端末装置 200 からユーザの選択指令を受けると (ステップ S320)、選択指令により指定されたターゲット 32 を対象として動作指令を生成する (

50

ステップ S 3 2 2 )。このように、情報処理装置 1 0 0 は、出力された複数のターゲット 3 2 ( および、対応するソース 3 4 ) に対するユーザ選択に応答して、選択されたターゲット 3 2 に対応する物体に対して物理的な作用を与えるための動作指令を生成する。なお、生成された動作指令は、ロボット 2 などへ出力されてもよい。

【 0 1 3 4 】

[ G . 評価結果 ]

次に、本実施の形態に従う対象物検索システム 1 の性能を評価した結果例を示す。

【 0 1 3 5 】

( g 1 : P F N - P I C )

まず、P F N - P I C データセット ( 非特許文献 1 ) を用いて、本実施の形態に従う対象物検索システム 1 の性能を評価した。上述したように、推論モデル 4 0 としては、M T C M、M T C M に敵対的生成ネットワーク ( G A N ) を付加した M T C M - G A N、および、M T C M - G A N に対して条件を付したモデルの 3 種類を採用可能である。それぞれのモデルについて、非特許文献 1 に開示される手法との比較を含めて評価を行った。

10

【 0 1 3 6 】

より具体的には、P F N - P I C データセットのうち 8 9 8 6 1 組のデータセットを用いて推論モデル 4 0 をトレーニングするとともに、別の 8 9 8 組のデータセットを用いて評価を行った。

【 0 1 3 7 】

なお、P F N - P I C データセットは、ピックアッププレイスのタスクに向けられたものであり、上から見て 4 つのボックスのいずれかに配置された対象物を見つけることが想定されている。各ボックスをソース 3 4 と見なし、対象物をターゲット 3 2 と見なして評価を行った。

20

【 0 1 3 8 】

命令文 2 2 としては、例えば、"Grab the black mug and put it in the lower right box." ( 黒いマグをつかんで右下のボックスに入れなさい。 ) といった、対象となるターゲット 3 2 およびソース 3 4 を含むものを用いた。

【 0 1 3 9 】

本実施の形態に従う対象物検索システム 1 においては、ターゲット 3 2 とソース 3 4 との組み合わせ毎、あるいは、ターゲット 3 2 毎に、命令文 2 2 による操作の対象物である「たしからしさ」( 確率 ) を算出できる。そのため、複数の物体が命令文 2 2 による操作の対象物になり得る ( すなわち、命令文 2 2 が不確実性を含んでいる ) 場合であっても、算出される確率に基づいて処理が可能である。

30

【 0 1 4 0 】

これに対して、非特許文献 1 に開示される手法は、命令文 2 2 による操作の対象物は 1 つであることが前提となっており、複数の物体が対象物になり得ることは何ら想定されていない。

【 0 1 4 1 】

非特許文献 1 は、マルチモーダル類似性ベースの手法であり、類似性の正確性 ( 類似性についての正答率 ) を評価した。この類似性の正確性は、推論モデル 4 0 として M T C M を採用した場合に、類似性識別器 4 0 7 から出力されるターゲットについての推論結果  $y_{target}$  ( 図 4 参照 ) に相当する。Table 1 においては、類似性の正確性に関しては、非特許文献 1 の手法と M T C M とを比較している。

40

【 0 1 4 2 】

本実施の形態に従う対象物検索システム 1 においては、与えられた命令文 2 2 の対象となるターゲット 3 2 以外に 1 または複数の不正解のターゲット 3 2 を用意した場合の正確性を評価した ( 領域毎の正確性 )。パラメータ は、不正解のターゲット 3 2 の数を示す。

【 0 1 4 3 】

【表 2】

Table 2

手法	類似性の 正確性	領域毎の正確性				
		$\gamma=1$	$\gamma=2$	$\gamma=4$	$\gamma=45$	$\gamma=10$
非特許文献1	88	-	-	-	-	
MTCM	88.8	94.5	95.4	96.1	96.3	97.1
MTCM-GAN (条件無)	-	95.7	96.1	96.2	96.4	96.8
MTCM-GAN (条件付)	-	95.9	96.3	96.5	96.5	97.2

10

## 【0144】

Table 2に示すように、推論モデル40としてMTCMを採用した場合であっても、非特許文献1に開示される手法に比較して改善効果を見ることができる。さらに、条件付きMTCM-GANを推論モデル40として採用することで、最も高い識別性能が発揮されていることが分かる。

## 【0145】

20

(g2:WRS-VS)

次に、World Robot Summit 2018 Virtual Space (以下、「WRS-VS」とも称す。) challengeで利用されたデータセットを用いても、本実施の形態に従う対象物検索システム1の性能を評価した。WRS-VSで用いられたデータセットは、SIGVerseに基づくものである(非特許文献9など参照)。

## 【0146】

より具体的には、WRS-VSデータセットのうち1010組のデータセットを用いて推論モデル40をトレーニングするとともに、別の37組のデータセットを用いて評価を行った結果をTable 3に示す。なお、パラメータは、不正解のターゲット32の数を示す。

30

## 【0147】

## 【表 3】

Table 3

手法	領域毎の正確性	
	$\gamma=1$	$\gamma=2$
MTCM	83.5	85.6
MTCM-GAN (条件無)	89.6	92.6
MTCM-GAN (条件付)	90.7	93.6

40

## 【0148】

Table 3に示すように、推論モデル40として、条件付きMTCM-GANを推論モデル40として採用することで、最も高い識別性能が発揮されていることが分かる。

## 【0149】

[H. 変形例]

本実施の形態に従う対象物検索システム1(推論モデル40)を十分に大きいトレーニングデータセットを用いてトレーニングすることで、汎化性能を高めることができる。こ

50

の場合、トレーニングによって得られたモデルパラメータのみを配布するようにしてもよい。

【 0 1 5 0 】

上述の説明においては、タスクが実行される場所にシステムを配置する、いわゆるオンプレミス環境に適した処理例を示すが、これに限らず、コンピュータネットワーク上に配置された 1 または複数のサーバを用いて、タスクを処理する、いわゆるクラウドサービス環境を採用してもよい。

【 0 1 5 1 】

本実施の形態に従う推論モデル 4 0 は、要求されるタスクの内容や実行環境などに応じて適宜適切な実装が可能である。例えば、推論モデル 4 0 を別のモデルの一部として組み込む、あるいは、推論モデル 4 0 と別のモデルとを組み合わせるといった実装形態が可能である。

10

【 0 1 5 2 】

[ I . まとめ ]

人間が発する言語による命令だけでは、認識対象の物体を一意に特定することはできず、不確実性が残ったものとなり得るが、本実施の形態に従えば、対象となる物体の候補が複数存在するような状況であっても、対象となる物体である確率を評価しつつ、対象となる物体を容易に特定できる。

【 0 1 5 3 】

今回開示された実施の形態は、すべての点で例示であって制限的なものではないと考えられるべきである。本発明の範囲は、上記した実施の形態の説明ではなくて特許請求の範囲によって示され、特許請求の範囲と均等の意味および範囲内でのすべての変更が含まれることが意図される。

20

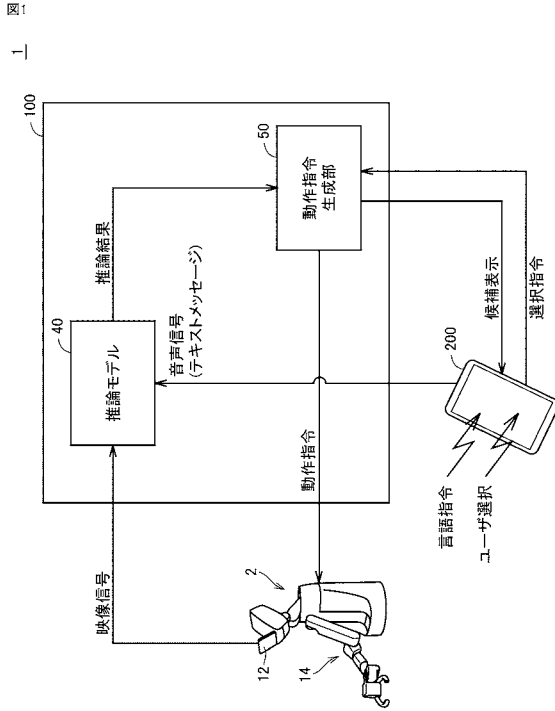
【 符号の説明 】

【 0 1 5 4 】

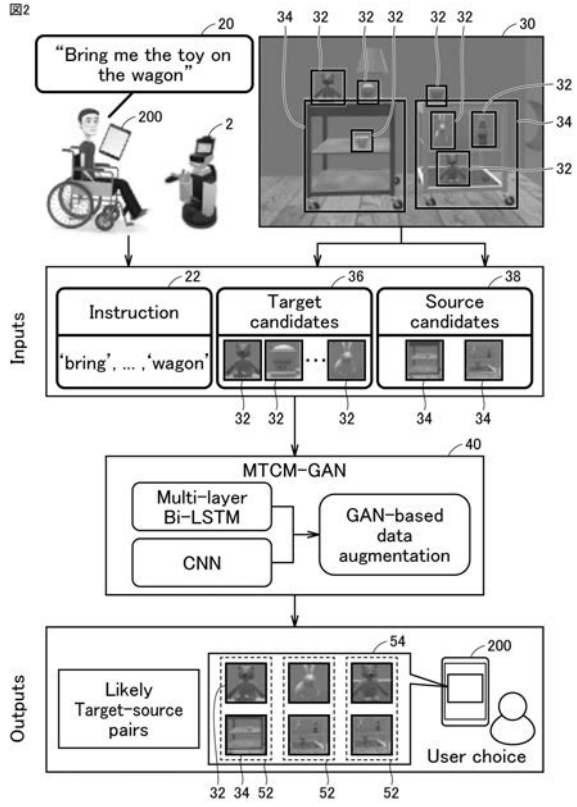
1 対象物検索システム、2 ロボット、12 カメラ、14 作用部、20 音声信号、22 命令文、30 入力画像、32 ターゲット、34 ソース、36 ターゲット候補、38 ソース候補、40 推論モデル、50 動作指令生成部、52 組み合わせ、54 推論結果、60 確率、100 情報処理装置、102 プロセッサ、104 主メモリ、106 ディスプレイ、108 入力デバイス、110 ネットワークインターフェイス、112 光学ドライブ、112 M 光学ディスク、114 入力インターフェイス、116 出力インターフェイス、118 内部バス、120 二次記憶装置、121 音声認識プログラム、122 画像抽出プログラム、123 トレーニングプログラム、124 動作指令生成プログラム、125 モデルパラメータ、126 トレーニングデータセット、200 端末装置、300 結果表示、401 命令文処理部、402, 406, 409 MLP、403 画像抽出部、404 モデル、405, 408, 421, 422 連結部、407 類似性識別器、410, 420 出力部、423 生成器、424 選択部、425 識別器。

30

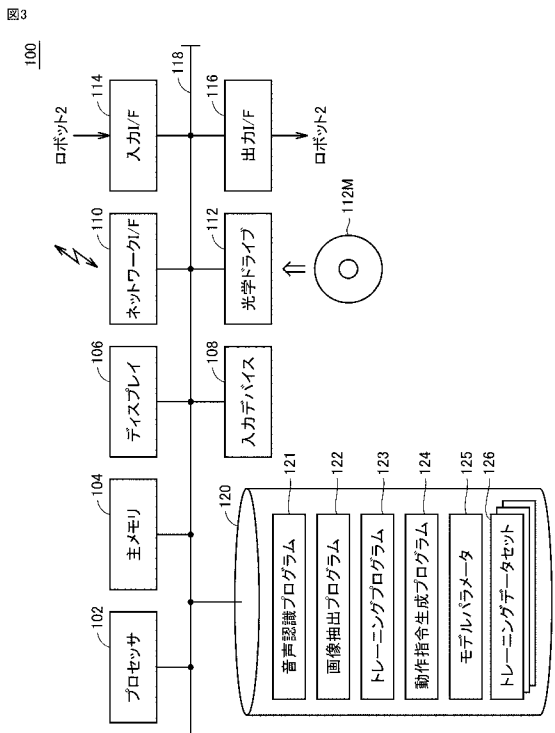
【図1】



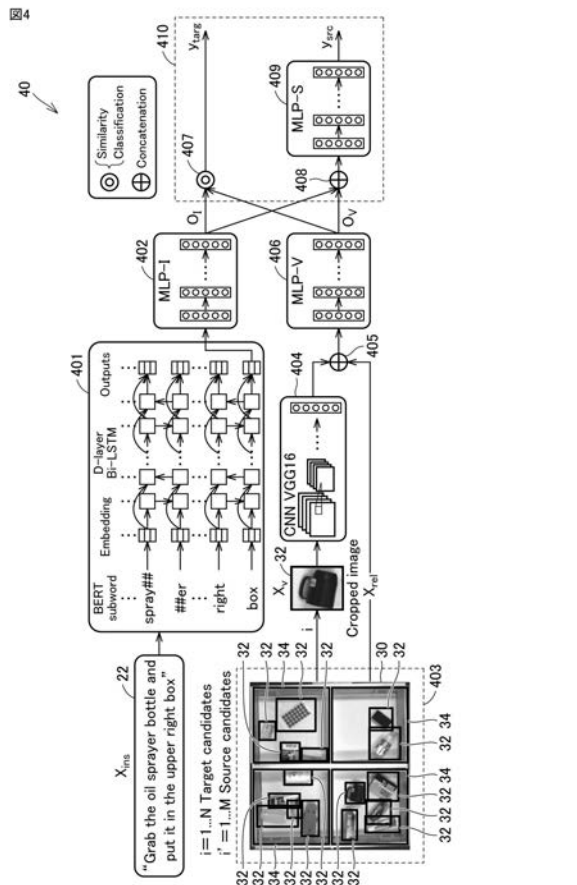
【図2】



【図3】



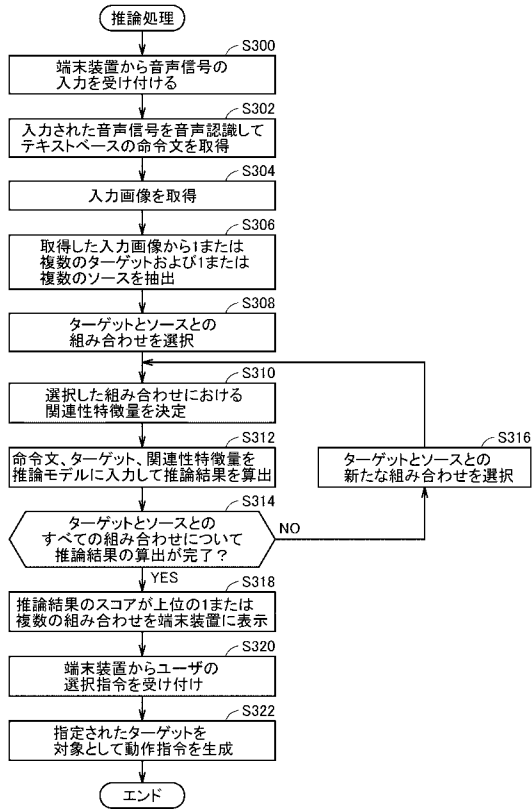
【図4】





【図 10】

図10



---

フロントページの続き

Fターム(参考) 5E555 AA23 AA46 BA01 BA83 BB01 BC04 BC17 BE17 CA42 CA47  
CB45 CB64 CC03 DA31 DB32 DB39 DB53 DC10 DC61 DD11  
EA07 EA22 EA25 EA27 FA00  
5L096 BA05 HA11 JA11 KA04