

TDP-MATに基づく実画像を対象とした物体操作指示理解

○小槻誠太郎, 石川慎太郎, 杉浦孔明 (慶應義塾大学)

1. はじめに

高齢化が進展する現代社会において、介助者の不足が深刻な社会問題となっている。こうした社会問題に対する解決策の一つとして生活支援ロボットが挙げられる。一方、生活支援ロボットが人間と自然な対話をする能力は未だ不十分である。そこで本研究では、生活支援ロボットにおける物体操作指示理解タスクを扱う。

本研究では自然言語による命令文と状況を説明する画像が与えられた際、正しく命令内容を解釈し、対象物体を特定することを目的とする。例えば“Look in the vase that is next to the potted plant”という命令文及び観葉植物と3つの花瓶が写った画像が与えられた場合、観葉植物に一番近い花瓶を対象物体として予測する事が望ましい。一方、人間の発する命令はしばしば曖昧であり、その内容を正確に理解することは容易ではない。例えば上述の例では参照表現を正確に解釈し、3つの花瓶から観葉植物に一番近いものを選ぶ必要がある。こうした参照表現を含む命令文では対象物体の特定が難しく、実際に誤る例が報告されている [1]。

既存手法である Target-Dependent UNITER [2] は物体操作指示理解タスクにおいて UNITER [3] を使用して高い精度を達成したモデルである。一方で Target-Dependent UNITER は精度がばらつきやすく、複雑なシーンを扱う際に十分な性能が安定して得られないという問題がある。

本研究では敵対的学習を導入し、Perceiver [4] によって特徴量間の関係性のモデル化を行う物体操作指示理解モデルである Target-Dependent Perceiver with Moment Adversarial Training (TDP-MAT) を提案する。敵対的学習は頑健なモデルを獲得するための一般的な手法であり、これによって提案手法はモデルの頑健性を向上させると期待される。また、Perceiver によって画像中の領域間の関係性のモデル化、及び画像中の領域群と命令文の特徴量の関係性のモデル化を行うことで、計算量の入力長さへの依存を抑えつつ、モデルの性能が向上することが期待される。提案手法と既存手法との違いは、視覚特徴量に対して、MAT [5] に基づく敵対的摂動を加える点、また Perceiver によって画像中の領域間の関係性のモデル化、及び画像中の領域群と命令文の特徴量の関係性のモデル化を行う点にある。

本研究の貢献を以下に示す。

- Perceiver によって画像中の領域間の関係性のモデル化、及び画像中の領域群と命令文の特徴量の関係性のモデル化を行う TDP-MAT を提案する。
- Perceiver の出力に対して MAT に基づく敵対的摂動を加える構造を導入する。
- REVERIE データセット [6] から命令文とそれに対応する実画像を収集し、REVERIE-fetch データセットを作成する。

2. 関連研究

マルチモーダル言語処理分野のサーベイ論文として、[7] が挙げられる。マルチモーダル言語処理分野のうち言語と画像を扱う分野に分類されるものとして Vision-and-Language Navigation (VLN) が挙げられる。

Multimodal Target-source Classifier Model [8] は自然言語による命令文と視覚的な情報から対象物体を特定するモデルである。HLSM-MAT [5] は VLN において、敵対的な摂動を特徴量空間に加える Moment-based Adversarial Training (MAT) を提案し、サブゴールおよび状態表現の特徴量空間に対して適用した手法である。

VLN 分野における有名なデータセットとして、Room-to-Room (R2R) [9], MatterPort3D [10], REVERIE [6] が挙げられる。R2R は実際の室内環境における視覚情報に基づいた自然言語による移動タスクのためのデータセットである。MatterPort3D は室内環境におけるシーン理解を行うための大規模な深度付き画像データセットである。REVERIE データセットは Remote Embodied Visual Referring Expression in Real Indoor Environments (REVERIE) タスクを行うためのデータセットである。この REVERIE タスクは、環境中に存在する対象物体に関連した自然言語による命令文が与えられ、始めに対象物体が存在する場所をゴールとする複数ステップの行動選択サブタスクを行い、その後対象物体を特定するサブタスクを行う。

3. 問題設定

本論文では以下のように用語を定義する。

- **対象物体:** 命令文が対象としている物体
- **候補物体:** 対象物体であるか否かの候補となる物体
- **コンテキスト物体:** 物体検出器で検出される物体

なお、対象領域、候補領域、コンテキスト領域はそれぞれ対象物体、候補物体、コンテキスト物体を囲むバウンディングボックスを指すものとする。

本論文で扱うタスクは、命令文と、それに対応する画像から物体検出器によって抽出された候補領域、コンテキスト領域が与えられ、候補領域が対象領域であるか否かの2値分類を行うというものである。以降、本タスクを REVERIE-fetch タスクと呼ぶ。REVERIE-fetch タスクにおいて期待される出力は、候補領域が対象領域である確率の予測値 $p(\hat{y})$ である。候補領域と対象領域が一致しているときは $p(\hat{y}) = 1$ 、異なるときは $p(\hat{y}) = 0$ と出力することが望ましい。図1に REVERIE-fetch タスクの代表例を示す。図1においては、画像中央奥に位置する赤いバウンディングボックスで示す枝編み細工の花瓶が対象物体となる。例えば候補領域として緑色のバウンディングボックスが与えられた時は $p(\hat{y}) = 0$ 、候補領域として対象領域に対応する赤いバウンディングボックスが与えられたときは $p(\hat{y}) = 1$ と出力するこ

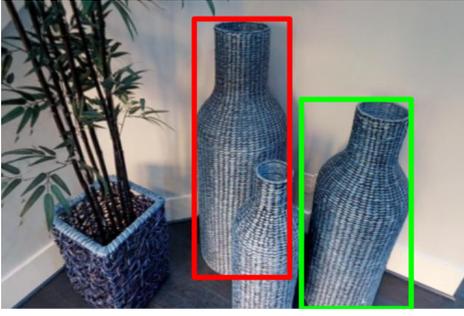


図 1 REVERIE-fetch タスクの例.

命令文: "Look in the left wicker vase that is next to the potted plant on the second floor at the foot of the stairs"

とが期待される. 本タスクの評価尺度には, 候補領域が対象領域であるか否かを分類する精度 Acc を使用する. なお本研究では, 物体検出誤りが十分に少ない物体検出器を利用できることを前提とし, 命令文に対応する画像から候補領域, 及びコンテキスト領域を抽出する過程はこの物体検出器によって行われることを仮定する.

4. 提案手法

提案ネットワークの構造を図 2 に示す. 図 2 において, Instruction は命令文, Context Regions はコンテキスト領域群, Candidate Region は候補領域を表す.

4.1 入力

ネットワークの入力 \mathbf{x} を以下のように定義する.

$$\mathbf{x} = f(\mathbf{X}_{inst}, \mathbf{X}_{cont}, \mathbf{X}_{cand}) \quad (1)$$

$$\mathbf{X}_{inst} = f(\mathbf{x}_{inst}, \mathbf{x}_{pos}) \quad (2)$$

$$\mathbf{X}_{cont} = \left\{ \mathbf{x}_{cont}^{(i)}, \mathbf{x}_{contloc}^{(i)} \mid i = 1, \dots, N_{FRCNN} \right\} \quad (3)$$

$$\mathbf{X}_{cand} = f(\mathbf{x}_{cand}, \mathbf{x}_{candloc}) \quad (4)$$

ここで, \mathbf{X}_{inst} は命令文, \mathbf{X}_{cont} はコンテキスト領域群の特徴量, \mathbf{X}_{cand} は候補領域の特徴量を表す. \mathbf{x}_{inst} は命令文のトークン列, \mathbf{x}_{pos} は命令文の各トークンの先頭からの順序であり, $\mathbf{x}_{cont}^{(i)}$ は i 番目の物体のコンテキスト領域の特徴量, $\mathbf{x}_{contloc}^{(i)}$ はそのコンテキスト領域の位置埋め込み, \mathbf{x}_{cand} は候補領域の特徴量, $\mathbf{x}_{candloc}$ は候補領域の位置埋め込みを表す. また, N_{FRCNN} は Faster R-CNN [11] によって検出した画像中の領域の数である.

\mathbf{X}_{inst} に関しては, SentencePiece [12] によるトークン化で \mathbf{x}_{inst} を獲得する. \mathbf{x}_{cont} に関しては, 全体画像を Faster R-CNN に入力し, Faster R-CNN における ResNet50 の fc6 層の出力によって \mathbf{x}_{cont} を獲得する. 一方で $\mathbf{x}_{contloc}$ は Faster R-CNN によって得たバウンディングボックスの位置情報をエンコードして得る. \mathbf{X}_{cand} は \mathbf{X}_{cont} の中から判定対象として選択する. なお, \mathbf{X}_{cont} 及び \mathbf{X}_{cand} における位置埋め込み $\mathbf{x}_{contloc}$ 及び $\mathbf{x}_{candloc}$ は, バウンディングボックスの左上の頂点の座標を (x_1, y_1) , 右下の座標を (x_2, y_2) とするとき, $[x_1, y_1, x_2, y_2, x_2 - x_1, y_2 - y_1, (x_2 - x_1) / (y_2 - y_1)]^T$ という 7 次元ベクトルを使用する.

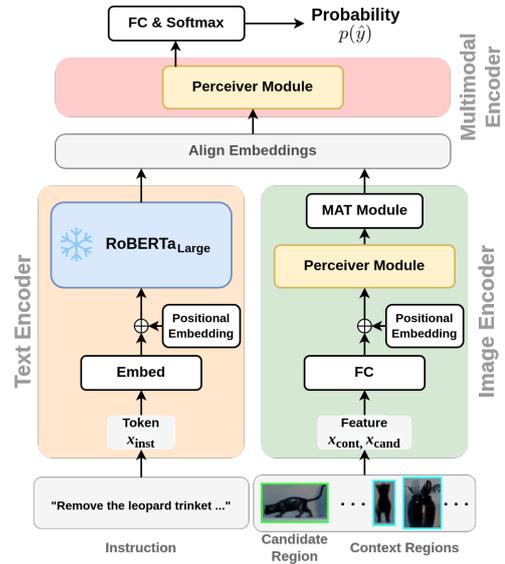


図 2 提案手法のネットワーク構成

4.2 Text Encoder

Text Encoder は \mathbf{x}_{inst} を入力として特徴量 E_T を出力する. 始めに \mathbf{x}_{inst} と \mathbf{x}_{pos} をそれぞれ 1024 次元のベクトル表現に埋め込んでから足し合わせ, 次にパラメータを固定した事前学習済み RoBERTa [13] に入力することで E_T を得る. なお, RoBERTa は事前学習に特化した最終層を取り除いた 23 層の transformer layer で構成される.

4.3 Image Encoder

Image Encoder は 2 つの全結合層, Perceiver モジュール, そして MAT モジュールから構成される. Image Encoder は \mathbf{X}_{cand} 及び \mathbf{X}_{cont} を入力とし, それぞれを全結合層によって 1024 次元のベクトル表現に埋め込んで足し合わせる. 次に足し合わせて得た特徴量を Perceiver モジュールに入力して各ベクトル表現の関係性をモデル化し, さらに MAT モジュールによって敵対的摂動を加え, 特徴量 E_{cand} 及び E_{cont} を出力する.

4.3.1 MAT モジュール

本モジュールは [5] で提案された Moment-based Adversarial Training に基づく敵対的摂動を加えるモジュールである. ここで敵対的摂動の更新方法を示す. 始めに, 交差エントロピー誤差に関する勾配 r_δ を計算する. 続いて, Adam をもとにした 2 種類の移動平均 \mathbf{m}_t , \mathbf{v}_t を導入する.

$$\mathbf{m}_t = \rho_1 \mathbf{m}_{t-1} + (1 - \rho_1) r_\delta E(\mathbf{v}_t), \quad (5)$$

$$\mathbf{v}_t = \rho_2 \mathbf{v}_{t-1} + (1 - \rho_2) (r_\delta E(\mathbf{v}_t))^2 \quad (6)$$

t は現在の摂動更新ステップ, ρ_1 と ρ_2 は各移動平均の平滑化係数である. これらを用いて, 摂動の更新幅 Δ_t を次のように導出する.

$$\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \rho_1^t}, \quad \hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \rho_2^t}, \quad (7)$$

$$\Delta_t = \eta \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}} \quad (8)$$

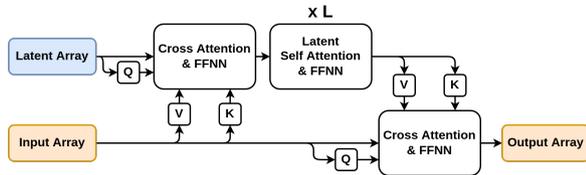


図3 Perceiver モジュールの構造

ここで η は MAT モジュールの学習率, ε はゼロ除算を避けるための小さな値である. 最後に導出した摂動の更新幅 Δ_t を利用し, [5] に従って を更新する.

4.3.2 Perceiver モジュール

本モジュールの構造を図3に示す. 本モジュールは Perceiver [4] を用い, 入力された特徴量間の関係性をモデル化する. Perceiver は, 入力された特徴量を, 潜在変数との cross attention を取ることによって固定サイズの潜在空間にマッピングし, 潜在空間にマッピングした特徴量に対して self attention を繰り返す. これにより, 入力の系列長に対して線形の計算量を達成しつつ入力された特徴量間の関係性をモデル化できる. 出力の際には, 潜在空間にマッピングされた特徴量を, 入力特徴量との cross attention を取って入力と同じ空間にマッピングする. 本モジュールではアテンション機構及び FFNN の初期化に DeepNet [14] で提案された DeepNorm を用い, 同様に DeepNorm に従って残差接続にモジュールの深さに応じた係数を掛ける.

4.4 Multimodal Encoder

本モジュールは Perceiver モジュールで構成される. 本モジュールは Text Encoder 及び Image Encoder の出力を連結した $[E_{\text{cand}}; E_{\text{inst}}; E_{\text{cont}}]$ を入力とし, 候補領域が対象領域である確率の予測値 $p(\hat{y})$ を出力する. モデルの学習に際して, 損失関数には交差エントロピー誤差を使用した.

5. 実験設定

5.1 データセット

我々は REVERIE データセット [6] から命令文及び各命令文に対応する実画像を収集し, 1 枚の画像から対象物体を特定するタスクのための REVERIE-fetch データセットを作成した. なお, 本実験では, REVERIE-fetch データセットの画像を Faster R-CNN [11] に入力し, 各画像について候補領域及びコンテキスト領域群を抽出してサンプルを作成した.

REVERIE-fetch データセットは, 画像と英語の命令文のペアからなるサンプルを 20990 サンプル含む. 語彙サイズは 2576, 全単語数は 399476, 平均文長は 19.0 である. REVERIE-fetch データセットは訓練集合に 20440 サンプル, 検証集合に 246 サンプル, そしてテスト集合に 304 サンプルのデータを含んでいる.

5.2 パラメータ設定

Image Encoder 及び Multimodal Encoder で使用される Perceiver [4] は, Latent Self Attention 層の数を 3, 潜在空間中の特徴量の次元数を 512, Attention の Head 数を 16 とした. 最適化には AdamW を使用し, 学習率は 10^{-3} , ステップ数は 32000, バッチサイズは

表1 REVERIE-fetch データセットにおける定量的結果

Method	Condition	Acc [%]
Baseline [2]		83.2 ± 2.25
Ours	full	85.5 ± 1.21
	(i)-a	50.0 ± 0.00
	(i)-b	83.8 ± 1.53
	(ii)	82.3 ± 1.50
	(iii)	83.2 ± 2.05

32, ドロップアウト率は 0.1 とした. なお, 1 ステップは 1 つのバッチの処理を意味する.

TDP-MAT の総パラメータ数は 510M, 学習可能パラメータ数は 168M である. 学習にはメモリ 24GB 搭載の GeForce RTX 3090 及びメモリ 64GB 搭載の Intel Core i9-10900KF を使用した. 学習には約 3 時間を要し, 推論には約 4.1ms/sample を要した.

6. 実験結果

6.1 定量的結果

表1に REVERIE-fetch データセットにおける各手法の精度を示す. なお, 精度の欄には 5 回の試行における平均値と標準偏差を示す. 本実験ではベースライン手法として, 類似のタスクにおいて良好な結果が報告されている Target-Dependent UNITER [2] を用いた. 表1より, 提案手法およびベースライン手法の精度はそれぞれ 85.5% および 83.2% であり, 提案手法がベースライン手法を 2.3 ポイント上回っている.

6.2 Ablation Study

Ablation Study には, 以下の 4 条件を定めた.

- (i)-a W/o MAT: Image Encoder において MAT モジュールによって敵対的摂動を加える場合と加えない場合で, 性能への影響を調べた.
- (i)-b W/o MAT+Smaller Learning Rate: 学習率を 10^{-4} に設定した上で, Image Encoder において MAT モジュールによって敵対的摂動を加える場合と加えない場合で, 性能への影響を調べた.
- (ii) W/o Perceiver Module in Image Encoder: Image Encoder において, Perceiver モジュールによって領域間の関係性をモデル化する場合としない場合を比較し, 性能への影響を調べた.
- (iii) W/o RoBERTa: Text Encoder においてパラメータを固定した事前学習済み RoBERTa を使用する場合としない場合を比較し, 性能への影響を調べた.

各条件での定量的結果を表1に示す. まず MAT モジュールを取り除いた場合, 同様の学習率では学習が失敗した. 一方, 学習率を 10^{-4} に設定した上で MAT モジュールを取り除いた場合は 1.7 ポイント精度が低下した. また, Text Encoder からパラメータを固定した事前学習済み RoBERTa を取り除いた場合は 2.3 ポイント, Image Encoder 中の Perceiver モジュールを取り除いた場合は 3.2 ポイント精度が低下した. 即ち Image Encoder において Perceiver モジュールによって領域間の関係性をモデル化することが性能に最も寄与していると言える. また, MAT の導入によって高い学習率の条件下で学習の発散を防げることが示された.

