

Prototypical Contrastive Transfer Learning for Multimodal Language Understanding

Seitaro Otsuki¹, Shintaro Ishikawa¹ and Komei Sugiura¹

Abstract—Although domestic service robots are expected to assist individuals who require support, they cannot currently interact smoothly with people through natural language. For example, given the instruction “Bring me a bottle from the kitchen,” it is difficult for such robots to specify the bottle in an indoor environment. Most conventional models have been trained on real-world datasets that are labor-intensive to collect, and they have not fully leveraged simulation data through a transfer learning framework. In this study, we propose a novel transfer learning approach for multimodal language understanding called Prototypical Contrastive Transfer Learning (PCTL), which uses a new contrastive loss called Dual ProtoNCE. We introduce PCTL to the task of identifying target objects in domestic environments according to free-form natural language instructions. To validate PCTL, we built new real-world and simulation datasets. Our experiment demonstrated that PCTL outperformed existing methods. Specifically, PCTL achieved an accuracy of 78.1%, whereas simple fine-tuning achieved an accuracy of 73.4%.

I. INTRODUCTION

In our aging society, the demand for daily care and support is increasing, which is leading to a shortage of home care workers. Domestic service robots (DSRs) are gaining popularity as a solution because of their ability to physically assist individuals. However, DSRs currently lack the capability to smoothly interact with people through natural language. To train their language comprehension models, it is desirable to use data collected in real-world environments. However, collecting and annotating such real-world data can be labor-intensive. By contrast, collecting training data using a simulator is much more cost-effective. Hence, it is advantageous to leverage simulation data through a transfer learning framework.

In this study, we focus on the task of identifying the target object in a given scenario using natural language instructions for object manipulation. For instance, given the instruction “Bring me the book closest to the lamp,” and a scene in which several books are near the lamp, the robot is expected to specify the book closest to the lamp as the target object. It is not easy to understand the meaning of human instructions correctly because such instructions are often ambiguous. In the above example, the robot should identify the book closest to the lamp among all the observed books by correctly comprehending the referring expression in the given instruction. Magassouba et al. [1] report cases in which a robot fails to comprehend instructions containing referring expressions.

¹The authors are with Keio University, 3-14-1 Hiyoshi, Kohoku, Yokohama, Kanagawa 223-8522, Japan. {otsu8sei14, shin.0116, komei.sugiura}@keio.jp

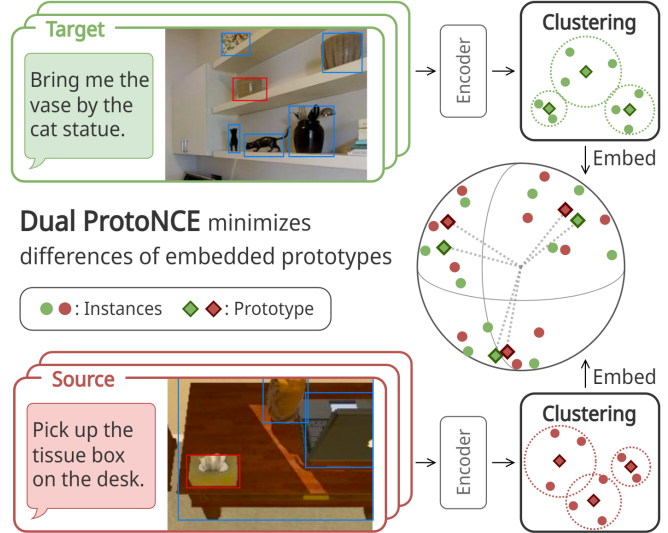


Fig. 1: Summary of our method PCTL. Given data from the source and target domains, PCTL aims to alleviate the influence of the domain gap by minimizing Dual ProtoNCE.

In this study, we aim to transfer experience gained from *simulation* data to *real-world* data to improve performance in our target task. In the task of multimodal language understanding for object manipulation, most conventional models have been trained only on real-world data [2]–[4]. However, such approaches have difficulty in increasing dataset size because building real-world datasets is labor-intensive. By contrast, collecting training data using a simulator is a significantly more cost-effective approach. Consequently, we expect that leveraging simulation data within a transfer learning framework will effectively enhance model performance.

In this paper, we propose Prototypical Contrastive Transfer Learning (PCTL), which is a novel transfer learning approach for the multimodal language understanding task. PCTL performs contrastive learning between the source and target domain data using our new contrastive loss called Dual ProtoNCE. We expect that PCTL will alleviate the influence of the gap between the two domains by minimizing Dual ProtoNCE. A summary of our method is shown in Fig. 1.

To design Dual ProtoNCE, we extended ProtoNCE [5] to transfer learning. ProtoNCE does not simultaneously handle source and target domains because it is not designed for transfer learning. Unlike ProtoNCE, Dual ProtoNCE is designed as a contrastive loss between the source and target domains. By defining a new contrastive loss between the source and target domains, we expect that Dual ProtoNCE

will enable contrastive representation learning to bridge the gap between the two domains.

Our key contributions are as follows:

- We introduce transfer learning to the task of identifying target objects in domestic environments according to free-form natural language instructions.
- We propose PCTL, which is a novel transfer learning approach for the multimodal language understanding task.
- Within PCTL, we develop Dual ProtoNCE, which is a novel contrastive loss generalized for transfer learning.

II. RELATED WORK

A. Multimodal Language Understanding

Many surveys have been conducted on multimodal language understanding and vision-language pretraining (VLP) [6]–[10]. Uppal et al. [7] present an overview of the latest trends in research on multimodal language understanding. They consider task formulations, evaluation metrics, model architecture, and other topics, for example, bias and fairness, and adversarial attacks. Long et al. [9] describe the general task definition and architecture of recent VLP models. They also discuss vision and language data encoding methods, and the mainstream model structure. They further summarize several essential pretraining and fine-tuning strategies.

Several studies have tackled referring expression comprehension (REC), which is one of the multimodal language understanding tasks [11]–[16]. In the REC task, models should ground a target object in an image described by a referring expression. Our task formulation is slightly more flexible than that of REC. In detail, we can address cases in which more than one or no target object exists in a given scene by formulating the task as the binary classification of whether or not the candidate object is the target. Some studies have attempted to build models for specifying the target object using natural language instructions and visual information [3], [4]. Specifically, Target-Dependent UNITER [4] (TDU) uses UNITER-based transformer architecture to model the relationship between text and visual features. Thus, TDU is pretrainable on general-purpose datasets. Additionally, Ishikawa et al. [4] formulates the task of identifying the target object in a given scenario using natural language instructions for object manipulation as the Multimodal Language Understanding for Fetching Instruction (MLU-FI) task. Ishikawa et al. [17] proposes Moment-based Adversarial Training (MAT), which is an adversarial training approach for vision-and-language navigation (VLN) tasks.

B. Datasets

Several datasets exist for the MLU-FI task. PFN-PIC [2] is a dataset that consists of images and instructions about objects in the scene. It contains images of approximately 20 commodities in four boxes taken in the real world, with the limitation of a fixed viewpoint. WRS-unialt [4] is a dataset that consists of images and instructions collected using a simulator. Additionally, these images were observed from various viewpoints.

In this study, we target the MLU-FI task in various real-world indoor environments with scenes observed from various viewpoints. However, to the best of our knowledge, no standard real-world dataset exists for this task. Therefore, we built new datasets by collecting data necessary for the task from datasets used in the VLN task. Several standard datasets exist for VLN [18]–[20]. The Room-to-Room (R2R) dataset [18] is a benchmark dataset for VLN in building-scale 3D environments in the real world. In the R2R navigation task, autonomous agents are required to follow navigation instructions in previously unseen indoor environments. The Remote Embodied Visual Referring Expression in Real Indoor Environments (REVERIE) dataset [19] is a standard dataset for the VLN task in real indoor environments. The REVERIE task consists of the subtask of navigating to a location where the target object exists, followed by another subtask of identifying the target object. These VLN datasets were built on the data provided by MatterPort3D [21]. MatterPort3D is a large-scale RGB-D dataset for scene understanding in various indoor environments in the real world. The dataset contains 10,800 panoramic views from 194,400 RGB-D images of 90 building-scale scenes with varied annotation, such as segmentation information.

C. Contrastive Learning

Contrastive learning is an approach used to learn a good data representation in a self-supervised manner. This category of learning strategies aims to align all instances in the embedding space where they are well-separated and locally smooth by leveraging the contrastive loss. Various contrastive learning frameworks have been proposed in the context of representation learning for vision [22]–[24], language [25], and multimodal models [26]–[28].

ProtoNCE [5] is a contrastive loss designed to implicitly encode the semantic structure of data into the embedding space by leveraging data prototypes obtained by clustering on embeddings as positive and negative features. In this study, we develop a novel contrastive loss called Dual ProtoNCE by extending ProtoNCE to transfer learning. Unlike ProtoNCE, we enable contrastive learning across different domains by leveraging data prototypes obtained from clustering for the embedded features of each domain.

III. PROBLEM STATEMENT

A. Preliminaries

The terminology used in this paper is defined as follows:

- **Target object:** object referred to in the natural language instruction.
- **Candidate object:** object that the model predicts whether it matches the target object or not.
- **Context objects:** objects detected by an object detector.

We refer to the bounding boxes of the target, candidate, and context objects as the target, candidate, and context regions, respectively.

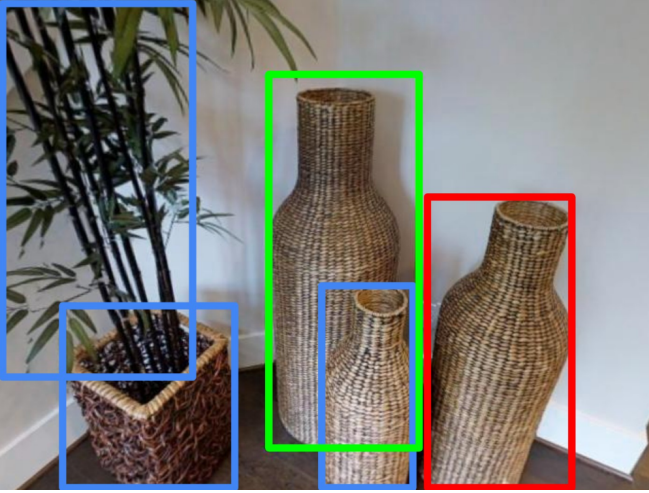


Fig. 2: Typical sample of the target task. The instruction is “Look in the left wicker vase next to the potted plant on the second floor at the foot of the stairs.” The red and green bounding boxes represent the candidate and target regions, respectively. Note that blue bounding boxes represent the context regions.

B. Task Formulation

We focus on the MLU-FI task. In this task, given a natural language instruction, candidate region, and context regions, the model is required to perform the binary classification of whether the candidate object matches the target object or not.

The MLU-FI task is characterized as follows:

- Input: An instruction, candidate region, and context regions.
- Output: Predicted probability $p(\hat{y} = 1)$. y and \hat{y} denote a label and predicted label, respectively. The condition $y = 1$ indicates that the candidate object matches the target object.

Fig. 2 shows a typical sample of the task. In the sample, the target object is the wicker vase enclosed by the green bounding box. In this case, $p(\hat{y} = 1) = 0$ because the given candidate object does not match the target object.

We assume that an object detector is used to extract the candidate region and context regions from the image. It is worth noting that the task is not a multi-class classification task to select a single object from all the objects in the image. The binary classification setting allows us to consider the case in which multiple or no target objects exist in the given image.

For transfer learning, the target samples are collected in the real world, whereas the source samples are collected using a simulator. Every sample consists of a set of an instruction, a candidate region, and context regions. They are collected in indoor environments for the MLU-FI task.

IV. PROPOSED METHOD

In this study, we propose PCTL, which is a novel transfer learning approach for the multimodal language understanding task. Specifically, we introduce Dual ProtoNCE, which is

a contrastive loss generalized for transfer learning. Although we apply our transfer learning approach to the MLU-FI task, our proposed method can be used for transfer learning on other multimodal language understanding tasks.

Fig. 3 shows an overview of our training framework. Our overall framework, referred to as PCTL, has three main modules: Encoder, Momentum Encoder, and Clustering Module.

A. Input

We define the input \mathbf{x} to our model as follows:

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

$$\mathbf{X}_{\text{cont}} = \{\mathbf{x}_{\text{cont}}^{(i)} | i = 1, \dots, N_{\text{det}}\}, \quad (2)$$

where \mathbf{x}_{inst} , \mathbf{x}_{cand} , and $\mathbf{x}_{\text{cont}}^{(i)}$ denote a natural language instruction, candidate region, and i -th context region, respectively. We use Faster R-CNN to detect N_{det} context regions.

It should be noted that we use positional encoding for \mathbf{x}_{cand} , \mathbf{X}_{cont} , and the text features extracted from \mathbf{x}_{inst} . We use a seven-dimensional vector $[x_1, y_1, x_2, y_2, x_2 - x_1, y_2 - y_1, (x_2 - x_1) \cdot (y_2 - y_1)]^T$ as positional encoding for \mathbf{x}_{cand} and \mathbf{X}_{cont} , where (x_1, y_1) and (x_2, y_2) are the coordinates of the top-left and bottom-right corners, respectively. Additionally, we assume that x_1 , x_2 , y_1 , and y_2 are normalized by the width and height of the input image.

B. Encoder

Encoder f_{θ} is parameterized by θ . The structure of f_{θ} follows TDU [4] and has three main parts: Text Embedder, Image Embedder, and Multi-Layer Transformer. The Text Embedder tokenizes \mathbf{x}_{inst} using WordPiece [29] and converts them to text features. The Image Embedder embeds \mathbf{x}_{cand} and $\mathbf{x}_{\text{cont}}^{(i)}$ into visual features. The Multi-Layer Transformer takes text and visual features as input and models the relationship between them. The output is the final hidden vector of the Multi-Layer Transformer corresponding to the input feature, \mathbf{x}_{cand} .

Hereafter, we refer to a sample of the source domain as (\mathbf{x}_s, y_s) . Similarly, (\mathbf{x}_t, y_t) denotes that of the target domain. The input and output of f_{θ} are denoted as follows:

$$\mathbf{u} = f_{\theta}(\mathbf{x}_s) \in \mathbb{R}^{768} \quad (3)$$

$$\mathbf{v} = f_{\theta}(\mathbf{x}_t) \in \mathbb{R}^{768}, \quad (4)$$

where \mathbf{u} and \mathbf{v} denote the feature vector of the source sample \mathbf{x}_s and that of the target sample \mathbf{x}_t . They are used for k -means clustering and the loss function. Classifier g consists of a two-layer MLP and softmax function. g takes the output of f_{θ} and calculates the predicted probability,

$$p(\hat{y} = 1) = g(f_{\theta}(\mathbf{x})). \quad (5)$$

C. Momentum Encoder

The Momentum Encoder $f_{\theta'}$ has the same structure as f_{θ} . $f_{\theta'}$ is parametrized by θ' , which is modeled as a moving average of θ . Specifically, we update θ' as follows:

$$\theta' \leftarrow \gamma \theta' + (1 - \gamma) \theta, \quad (6)$$

samples and target samples independently:

$$\mathcal{L}_{\text{Intra}} = \mathcal{L}_{\text{Target}} + \mathcal{L}_{\text{Source}} \quad (13)$$

$$\mathcal{L}_{\text{Target}} = \sum_{i=1}^n - \left(\log \frac{\exp(\dagger \mathbf{v}_i \cdot \dagger \mathbf{v}'_i / \tau)}{\sum_{j \in J} \exp(\dagger \mathbf{v}_i \cdot \dagger \mathbf{v}'_j / \tau)} \right) + \frac{1}{M} \sum_{m=1}^M \log \frac{\exp(\dagger \mathbf{v}_i \cdot \dagger \mathbf{c}_s^{(m)} / \phi_s^{(m)})}{\sum_{j \in J'} \exp(\dagger \mathbf{v}_i \cdot \dagger \mathbf{c}_j^{(m)} / \phi_j^{(m)})} \quad (14)$$

$$\mathcal{L}_{\text{Source}} = \sum_{i=1}^n - \left(\log \frac{\exp(\dagger \mathbf{u}_i \cdot \dagger \mathbf{u}'_i / \tau)}{\sum_{j \in J} \exp(\dagger \mathbf{u}_i \cdot \dagger \mathbf{u}'_j / \tau)} \right) + \frac{1}{M} \sum_{m=1}^M \log \frac{\exp(\dagger \mathbf{u}_i \cdot \dagger \mathbf{d}_s^{(m)} / \varphi_s^{(m)})}{\sum_{j \in J'} \exp(\dagger \mathbf{u}_i \cdot \dagger \mathbf{d}_j^{(m)} / \varphi_j^{(m)})}, \quad (15)$$

where $\dagger \mathbf{a}$ represents $\mathbf{a} / \|\mathbf{a}\|_2$ and ϕ and φ denote the concentration factors of \mathbf{c} and \mathbf{d} , respectively.

Next, we compute $\mathcal{L}_{\text{Inter}}$ to bridge the gap between the two domains. $\mathcal{L}_{\text{Inter}}$ is defined as follows:

$$\mathcal{L}_{\text{Inter}} = \mathcal{L}_{\text{S2T}} + \mathcal{L}_{\text{T2S}}, \quad (16)$$

where \mathcal{L}_{S2T} is the contrastive loss defined between source domain features \mathbf{u} and the prototypes of target domain \mathbf{c} , and \mathcal{L}_{T2S} is similarly defined between \mathbf{v} and \mathbf{d} . They are expressed as

$$\mathcal{L}_{\text{S2T}} = -\frac{1}{M} \sum_{i=1}^n \sum_{m=1}^M \left(\log \frac{\exp(\dagger \mathbf{u}_i \cdot \dagger \mathbf{c}_s^{(m)} / \phi_s^{(m)})}{\sum_{j \in J'} \exp(\dagger \mathbf{u}_i \cdot \dagger \mathbf{c}_j^{(m)} / \phi_j^{(m)})} \right), \quad (17)$$

$$\mathcal{L}_{\text{T2S}} = -\frac{1}{M} \sum_{i=1}^n \sum_{m=1}^M \left(\log \frac{\exp(\dagger \mathbf{v}_i \cdot \dagger \mathbf{d}_s^{(m)} / \varphi_s^{(m)})}{\sum_{j \in J'} \exp(\dagger \mathbf{v}_i \cdot \dagger \mathbf{d}_j^{(m)} / \varphi_j^{(m)})} \right). \quad (18)$$

Let \mathcal{L}_{CE} and λ be the cross-entropy loss and a hyperparameter, respectively. Our overall loss function \mathcal{L} is defined as,

$$\mathcal{L} = \lambda \mathcal{L}_{\text{DualProtoNCE}} + \mathcal{L}_t + \mathcal{L}_s \quad (19)$$

$$\mathcal{L}_t = \sum_{i=1}^n \left(\mathcal{L}_{\text{CE}}(g(f_\theta(\mathbf{x}_t^{(i)})), y_t^{(i)}) + \frac{1}{M} \sum_{m=1}^M \mathcal{L}_{\text{CE}}(g(\mathbf{c}_s^{(m)}), y_t^{(i)}) \right) \quad (20)$$

$$\mathcal{L}_s = \sum_{i=1}^n \left(\mathcal{L}_{\text{CE}}(g(f_\theta(\mathbf{x}_s^{(i)})), y_s^{(i)}) + \frac{1}{M} \sum_{m=1}^M \mathcal{L}_{\text{CE}}(g(\mathbf{d}_s^{(m)}), y_s^{(i)}) \right), \quad (21)$$

where $\mathbf{c}_s^{(m)}$ and $\mathbf{d}_s^{(m)}$ are the *prototypes* closest to the embedding $f_\theta(\mathbf{x}_t^{(i)})$ and $f_\theta(\mathbf{x}_s^{(i)})$, respectively.

V. EXPERIMENTS

A. Datasets

To validate our model in real-world domestic environments, we built a new dataset called REVERIE-fetch. This is because no standard real-world dataset exists for the MLU-FI task, to the best of our knowledge. To construct such a dataset, we collected images and natural language instructions based on the REVERIE dataset [19], which is a standard dataset for VLN in real-world indoor environments. Note that this dataset is not directly applicable to our task.

We first collected the cubemaps [31] of the goal points provided in the original dataset because the target object is placed at the goal point of the navigation task in the REVERIE task. Then, we extracted images in which target objects existed from the collected cubemaps. Over 1,000 annotators collected the instructions in the REVERIE dataset using Amazon Mechanical Turk. The annotators viewed an animation of the route and a randomly highlighted target object via the interactive 3D WebGL simulator. Then, they were asked to provide instructions to find and manipulate the target object.

Regarding the source-domain datasets, we extracted the source samples from the ALFRED dataset [32] and built the ALFRED-fetch-b dataset. The ALFRED dataset is a standard dataset for VLN with object manipulation. It includes 25,743 English instructions that describe 8,055 expert demonstrations. It contains multiple sequential subgoals that constitute the given goal, an instruction for each subgoal, and images observed from the agent's first-person views at each timestep of ground-truth behavior. The new ALFRED-fetch-b dataset consists of instructions and images from the training set of the original ALFRED dataset. They were collected in a scenario in which the subgoal was to pick up an object. In this study, we gathered the images from scenes just before the picking action.

As a preprocessing step for data extracted from the REVERIE and ALFRED datasets, we extracted the candidate and context regions from the images using Faster R-CNN [33] to create samples. We labeled those with a GIoU [34] greater than 0.80 of their target and candidate regions as

TABLE I: Experimental setup.

PCTL	$\mathcal{L}_{\text{DualProtoNCE}}$	$(r, r') = (32, 32),$ $k^{(1)} = 64, \lambda = 1/32$
	Concentration, ϕ	$\tau' = 0.2, \alpha = 10$
Transformer		#L: 12, #H: 768, #A: 12
Optimizer		SGD w/ momentum 0.9
Learning rate (LR)		8×10^{-4}
LR for Multi-Layer Transformer		8×10^{-5}
Batch size		64
#Epoch		30

TABLE II:
Quantitative results on the REVERIE-fetch dataset.

Method	Accuracy [%]
Target domain only	73.0 \pm 1.87
Fine-tuning	73.4 \pm 11.8
MCDDA+ [35]	74.9 \pm 3.94
Ours	78.1 \pm 2.49

positive samples and those with a GIoU less than 0.45 as negative samples. The target regions were provided in the original datasets. In the test set of REVERIE-fetch, we manually removed inappropriate samples caused by misdetection.

The REVERIE-fetch dataset consists of 10,243 image-instruction pairs with a vocabulary size of 1958 words, a total of 188,965 words, and an average sentence length of 18.4 words. The ALFRED-fetch-b dataset similarly consists of 34,286 samples of image-instruction pairs and a total of 399,964 words. Its vocabulary size and average sentence length are 1,558 and 11.7 words, respectively. The REVERIE-fetch dataset includes 8,302, 994, and 947 samples in the training, validation, and test sets, respectively. We built the training set with data collected from the training set and the seen split of the validation set of the REVERIE dataset. The validation and test sets consist of data collected from the unseen split of the validation set of the REVERIE dataset. The ALFRED-fetch-b dataset includes 27,492; 3,470; and 3,324 samples in the training, validation, and test sets, respectively. We collected these samples from the training set of the ALFRED dataset. It should be mentioned that there were no overlaps among the training, validation, and test sets for either dataset. We used the training set to train our model and the validation set to tune the hyperparameters. We evaluated our model on the test set of the REVERIE-fetch dataset.

B. Experimental Setup

Table I summarizes the experimental setup. Note that #L, #H, and #A denote the number of layers, hidden size, and number of attention heads in the Multi-Layer Transformer, respectively.

Our model had roughly 110 million trainable parameters. We trained our model on a GeForce RTX 3090 with 24GB of memory and an Intel Core i9-10900KF with 64GB of memory. It took 2 hours to train our model. The inference time was approximately 59.3 milliseconds for one sample. We evaluated the value of the loss \mathcal{L}_{CE} of the model for every epoch on the validation set of REVERIE-fetch. We used the test set accuracy of the REVERIE-fetch dataset when the value of \mathcal{L}_{CE} on the validation set of REVERIE-fetch was minimized.

C. Quantitative Results

We conducted experiments to compare the proposed and baseline methods. Table II shows the accuracy on the REVERIE-fetch dataset. The right column shows the means

and standard deviations over five trials. In this experiment, we used accuracy as the evaluation metric because the numbers of positive and negative samples were almost balanced.

We set the following three baseline settings:

- (i) Target domain only: We trained the model only on target samples.
- (ii) Fine-tuning: We performed pretraining on the source samples and fine-tuned the model with the target samples.
- (iii) MCDDA+: We extended the Maximum Classifier Discrepancy for Domain Adaptation [35] (MCDDA) and applied it to a supervised transfer learning setting.

We set up baselines (i) and (ii) to compare our approach with the approach in which no data from the source domain was used and the approach in which pretraining was performed on data from the source domain, respectively. As reported in [35], MCDDA performs well as an unsupervised transfer learning method on image classification. Therefore, we extended MCDDA to a supervised transfer learning setting and used it as baseline (iii). We called this extended method MCDDA+.

As listed in Table II, our method achieved an accuracy of 78.1%, whereas the accuracy of baselines (i), (ii), and (iii) were 73.0%, 73.4%, and 74.9%, respectively. Therefore our method outperformed all the baselines (i), (ii), and (iii) by 5.1, 4.7, and 3.2 points in terms of accuracy, respectively. The performance difference between baseline (i) and our method was statistically significant (p-value was lower than 0.01).

D. Ablation Studies

We conducted ablation studies to investigate the contribution of M and the combination of $k^{(m)}$ to performance. Specifically, we set the following conditions:

- (i)-a $M = 1, k^{(1)} = 33$
- (i)-b $M = 1, k^{(1)} = 64$
- (i)-c $M = 1, k^{(1)} = 128$
- (ii) $M = 3, (k^{(1)}, k^{(2)}, k^{(3)}) = (64, 128, 256)$

It is worth noting that we chose $k^{(1)} = 33$ for condition (i)-a because it is the minimum value of $k^{(1)}$ that allows us to select 32 negative prototypes ($r' = 32$) and one positive prototype from a total of $k^{(1)}$ prototypes.

Table III lists the quantitative results of the ablation study. The accuracy column shows the means and standard deviations over five trials. As described in Table III, PCTL

TABLE III:
Quantitative results of the ablation studies.

Condition	Accuracy [%]
Ours	78.1 \pm 2.49
(i)-a	75.6 \pm 1.96
(i)-b	73.7 \pm 2.92
(i)-c	77.4 \pm 1.96
(ii)	71.7 \pm 10.3



(a) “Go down the stairs to the lower balcony area and turn off the lamp on the dresser.”

(b) “Go to the lounge on the first level where the red carpet is and move the black vase to the right of the mirror.”

(c) “Fluff the light silver pillow on the smaller couch in the living room”

Fig. 4: Qualitative results on the REVERIE-fetch dataset. Panels (a), (b), and (c) show the true positive, true negative, and false positive cases, respectively. The red and green bounding boxes indicate the candidate and target regions, respectively. For panels (a) and (b), our method correctly predicted whether the candidate object matched the target object or not by successfully comprehending the instruction and scene.

achieved accuracy of 78.1%, whereas the accuracy under conditions (i)-a, (i)-b, (i)-c, and (ii) were 75.6%, 73.7%, 77.4%, and 71.7%, respectively. This indicates that our method achieved the highest accuracy under all the ablation conditions for $k^{(m)}$ and M . Moreover, this result indicates that a decrease or increase of $k^{(m)}$ reduced performance.

E. Qualitative Results

Qualitative results are shown in Fig. 4. Fig. 4 (a) shows a successful sample. The instruction was “Go down the stairs to the lower balcony area and turn off the lamp on the dresser.” The target object is the lamp on the dresser. Our method correctly predicted that the candidate object matched the target object, whereas the baseline (i) wrongly predicted that the candidate object did not match the target object. Fig. 4 (b) shows another successful sample. The given instruction was “Go to the lounge on the first level where the red carpet is and move the black vase to the right of the mirror.” The target object was the vase on the right side of the mirror. Our method successfully identified the candidate object as a different object from the target object, whereas baseline (i) failed to do this.

Fig. 4 (c) shows a failed sample. The instruction was “Fluff the light silver pillow on the smaller couch in the living room.” The target object was the pillow on the left side of the couch. Our method incorrectly predicted that the candidate object matched the target object.

F. Error Analysis and Discussion

The results consisted of 371, 65, 126, and 385 samples for true positives, false positives (FP), false negatives (FN) and true negatives, respectively. Thus, there were 191 samples for the failed cases. We randomly selected 50 FP and 50 FN samples to analyze the causes of errors. Table IV categorizes these samples. We classified the causes of errors into the following seven types:

- Comprehension Error (CE): The model failed to process the visual information and instruction correctly.

This class includes cases in which the model failed to comprehend the given referring expression or correctly specify the object to which the textual information in the instruction referred.

- Missing Landmark: The given image did not contain the visual information w.r.t. the referring expression. For example, the model failed to predict the chair nearest the kitchen because the instruction had the referring expression, “nearest the kitchen,” but the given image did not contain a kitchen.
- Small Region: The model failed to specify the target object because the target region was smaller than 1% of the entire image area.
- Ambiguous Instruction: The given instruction was ambiguous; hence, the model failed to specify the target object.
- Annotation Error: Annotation errors occurred in the bounding boxes and/or instructions.
- Severe Occlusion: The target object was severely occluded by other objects.
- Multiple Objects: The candidate region enclosed multiple objects.

As shown in Table IV, the main bottleneck was CE. We could reduce the number of cases by using a huge number of

TABLE IV: Categorization of failed samples.

Error Type	#Error
Comprehension Error	43
Missing Landmark	17
Small Region	14
Ambiguous Instruction	11
Annotation Error	10
Severe Occlusion	3
Multiple Objects	2

source samples or introducing pretrained models [26], [27] that embed language features and visual features into the same embedding space.

VI. CONCLUSIONS

In this study, we proposed PCTL, which is a novel transfer learning approach for the multimodal language understanding task. Specifically, we developed Dual ProtoNCE, which is a new contrastive loss generalized for transfer learning.

Our key contributions are as follows:

- We introduced transfer learning to the MLU-FI task.
- We proposed PCTL, which is a novel transfer learning approach for the multimodal language understanding task.
- Within PCTL, we developed Dual ProtoNCE, which is a novel contrastive loss generalized for transfer learning.
- PCTL outperformed the baselines in terms of the accuracy of the MLU-FI task on the REVERIE-fetch dataset.

In future work, we plan to enrich the source-domain dataset using a simulator and apply the model trained by the PCTL framework to physical robots.

ACKNOWLEDGMENT

This work was partially supported by JSPS KAKENHI Grant Number 20H04269, JST Moonshot, and NEDO.

REFERENCES

- [1] A. Magassouba, K. Sugiura, and H. Kawai, "A Multimodal Target-Source Classifier With Attention Branches to Understand Ambiguous Instructions for Fetching Daily Objects," *RA-L*, vol. 5, no. 2, pp. 532–539, 2020.
- [2] J. Hatori, Y. Kikuchi, S. Kobayashi, *et al.*, "Interactively Picking Real-World Objects with Unconstrained Spoken Language Instructions," in *ICRA*, 2018, pp. 3774–3781.
- [3] A. Magassouba, K. Sugiura, A. T. Quoc, H. Kawai, *et al.*, "Understanding Natural Language Instructions for Fetching Daily Objects Using GAN-Based Multimodal Target–Source Classification," *RA-L*, vol. 4, no. 4, pp. 3884–3891, 2019.
- [4] S. Ishikawa and K. Sugiura, "Target-dependent UNITER: A Transformer-Based Multimodal Language Comprehension Model for Domestic Service Robots," *RA-L*, vol. 6, no. 4, pp. 8401–8408, 2021.
- [5] J. Li, P. Zhou, C. Xiong, and S. Hoi, "Prototypical Contrastive Learning of Unsupervised Representations," in *ICLR*, 2021.
- [6] A. Mogadala, M. Kalimuthu, D. Klakow, *et al.*, "Trends in integration of vision and language research: A survey of tasks, datasets, and methods," *JAIR*, vol. 71, pp. 1183–1317, 2021.
- [7] S. Uppal, S. Bhagat, D. Hazarika, N. Majumder, S. Poria, R. Zimmermann, and A. Zadeh, "Multimodal research in vision and language: A review of current and emerging trends," *Information Fusion*, vol. 77, pp. 149–171, 2022.
- [8] Y. Du, Z. Liu, J. Li, and W. X. Zhao, "A Survey of Vision-Language Pre-Trained Models," in *IJCAI*, 2022.
- [9] S. Long, F. Cao, S. C. Han, and H. Yang, "Vision-and-Language Pretrained Models: A Survey," in *IJCAI*, 2022.
- [10] F.-L. Chen, D.-Z. Zhang, M.-L. Han, X.-Y. Chen, J. Shi, S. Xu, and B. Xu, "VLP: A Survey on Vision-language Pre-training," *Machine Intelligence Research*, vol. 20, no. 1, pp. 38–56, 2023.
- [11] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg, "MAttNet: Modular Attention Network for Referring Expression Comprehension," in *CVPR*, 2018, pp. 1307–1315.
- [12] J. Lu, D. Batra, D. Parikh, and S. Lee, "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks," in *NeurIPS*, vol. 32, 2019.
- [13] Y.-C. Chen, L. Li, L. Yu, A. El Kholly, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, "UNITER: UNiversal Image-TExt Representation Learning," in *ECCV*, 2020, pp. 104–120.
- [14] Z. Gan, Y.-C. Chen, L. Li, C. Zhu, Y. Cheng, and J. Liu, "Large-Scale Adversarial Training for Vision-and-Language Representation Learning," in *NeurIPS*, vol. 33, 2020, pp. 6616–6628.
- [15] A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion, "MDETR - Modulated Detection for End-to-End Multi-Modal Understanding," in *ICCV*, 2021, pp. 1780–1790.
- [16] P. Wang, A. Yang, R. Men, J. Lin, S. Bai, Z. Li, J. Ma, C. Zhou, J. Zhou, and H. Yang, "OFA: Unifying Architectures, Tasks, and Modalities Through a Simple Sequence-to-Sequence Learning Framework," in *ICML*, 2022, pp. 23 318–23 340.
- [17] S. Ishikawa and K. Sugiura, "Moment-based Adversarial Training for Embodied Language Comprehension," in *ICPR*, 2022, pp. 4139–4145.
- [18] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, A. Van Den Hengel, *et al.*, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *CVPR*, 2018, pp. 3674–3683.
- [19] Y. Qi, Q. Wu, P. Anderson, X. Wang, W. Y. Wang, C. Shen, A. v. d. Hengel, *et al.*, "Reverie: Remote embodied visual referring expression in real indoor environments," in *CVPR*, 2020, pp. 9982–9991.
- [20] A. Ku, P. Anderson, R. Patel, E. Ie, and J. Baldridge, "Room-Across-Room: Multilingual Vision-and-Language Navigation with Dense Spatiotemporal Grounding," in *EMNLP*, 2020, pp. 4392–4412.
- [21] A. X. Chang, A. Dai, T. A. Funkhouser, M. Halber, M. Nießner, M. Savva, S. Song, A. Zeng, Y. Zhang, *et al.*, "Matterport3D: Learning from RGB-D Data in Indoor Environments," in *3DV*, 2017, pp. 667–676.
- [22] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A Simple Framework for Contrastive Learning of Visual Representations," in *ICML*, 2020, pp. 1597–1607.
- [23] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum Contrast for Unsupervised Visual Representation Learning," in *CVPR*, 2020, pp. 9729–9738.
- [24] I. Misra and L. v. d. Maaten, "Self-Supervised Learning of Pretext-Invariant Representations," in *CVPR*, 2020, pp. 6707–6717.
- [25] T. Gao, X. Yao, and D. Chen, "SimCSE: Simple Contrastive Learning of Sentence Embeddings," in *EMNLP*, 2021, pp. 6894–6910.
- [26] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning Transferable Visual Models From Natural Language Supervision," in *ICML*, 2021, pp. 8748–8763.
- [27] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision," in *ICML*, 2021, pp. 4904–4916.
- [28] B. Wu, R. Cheng, P. Zhang, T. Gao, J. E. Gonzalez, and P. Vajda, "Data Efficient Language-Supervised Zero-Shot Recognition with Optimal Transport Distillation," in *ICLR*, 2022.
- [29] 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.
- [30] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation Learning with Contrastive Predictive Coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [31] F. Duanmu, Y. He, X. Xiu, P. Hanhart, Y. Ye, and Y. Wang, "Hybrid Cubemap Projection Format for 360-Degree Video Coding," in *Data Compression Conference*, 2018, pp. 404–404.
- [32] M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and D. Fox, "ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks," in *CVPR*, 2020, pp. 10 740–10 749.
- [33] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *IEEE Trans. PAMI*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [34] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, S. Savarese, *et al.*, "Generalized intersection over union: A metric and a loss for bounding box regression," in *CVPR*, 2019, pp. 658–666.
- [35] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *CVPR*, 2018, pp. 3723–3732.