# Prompt-based Learning for Unpaired Image Captioning

Peipei Zhu, Xiao Wang, *Member, IEEE*, Lin Zhu, Zhenglong Sun, *Senior Member, IEEE*
Wei-Shi Zheng, Yaowei Wang, *Member, IEEE*, and Changwen Chen, *Fellow, IEEE*

*Abstract*—**Unpaired Image Captioning (UIC) has been developed to learn image descriptions from unaligned vision-language sample pairs. Existing works usually tackle this task using adversarial learning and visual concept reward based on reinforcement learning. However, these existing works were only able to learn limited cross-domain information in vision and language domains, which restrains the captioning performance of UIC. Inspired by the success of Vision-Language Pre-Trained Models (VL-PTMs) in this research, we attempt to infer the cross-domain cue information about a given image from the large VL-PTMs for the UIC task. This research is also motivated by recent successes of prompt learning in many downstream multi-modal tasks, including image-text retrieval and vision question answering. In this work, a semantic prompt is introduced and aggregated with visual features for more accurate caption prediction under the adversarial learning framework. In addition, a metric prompt is designed to select high-quality pseudo image-caption samples obtained from the basic captioning model and refine the model in an iterative manner. Extensive experiments on the COCO and Flickr30K datasets validate the promising captioning ability of the proposed model. We expect that the proposed prompt-based UIC model will stimulate a new line of research for the VL-PTMs based captioning.**

*Index Terms*—**Prompt-based Learning, Unpaired Image Captioning, Semantic Prompt, Metric Prompt**

## I. INTRODUCTION

The goal of image captioning is to automatically describe visual images with natural languages. This is a cross-modality task that transfers information from the image domain to the language domain [1]–[3]. With the release of large-scale captioning datasets [4], [5] and the advances in deep learning, the performance of image captioning has been continuously improved. It has been widely used in many applications, such as human-robot interaction [6], [7], visual aid for the blind [8]–[10], and automatic driving [11]–[13]. The mainstream image captioning models have followed the encoder-decoder paradigm [14], [15], which encodes the image into feature representation first and then decodes it into a sentence in a word-by-word fashion. Although the performance is good, such supervised learning based captioning models rely on massively labeled vision-language pairs [16]–[18], which is time- and energy-consuming. Also, the models trained on limited samples may have poor generalization ability.

Considering the limitations of the fully-supervised image captioning paradigm, captioning using unpaired vision-language samples draws more and more attention as this approach does not require carefully labeled image-text training pairs. Usually, these models are developed based on *adversarial learning* [19], [20] and the *visual concept reward* based on reinforcement learning [21], [22]. As an early attempt, adversarial learning can only be utilized to guide the optimization of UIC parameters from the perspective of the overall structure of the sentence, while the correlations between the vision domain and the language domain have not been sufficiently explored. The concept reward based UIC models simply restrain their captions to contain the detected visual concepts (such as "dog" and "tree"), therefore, their performance heavily depends on object detectors and very limited cross-domain knowledge is concerned. How to exploit more vision-language knowledge without paired image-text samples for UIC is still a challenging research problem to be resolved.

Recently, the pre-trained giant models [23] have demonstrated their abundant prior knowledge by their superior performance in multiple domains and tasks, including natural language processing, computer vision, and multi-modal. These models carry an extremely large number of parameters and are pre-trained on the super large-scale corpus. For example, the CLIP [23] is pre-trained with 400 million image-text pairs using cosine similarity maximization. Its superior performance on zero/few-shot learning demonstrates that it carries a lot of visual-language prior knowledge. Many other computer vision tasks have proved that the CLIP features further improve their performance significantly [24]–[26]. On the other hand, prompt learning [27] is proposed to better leverage pre-trained models to improve overall performance on downstream tasks, such as PPT [28], CoOp [29], and VPT [30]. These works inspire us to *design new mechanisms for UIC by extracting prior vision-language knowledge from pre-trained big models.*

In this paper, a novel Prompt-based Learning scheme is proposed for UIC, termed PL-UIC, which can extract prior knowledge from the large-scale VL-PTMs. The key insight of

Peipei Zhu is with the School of Science and Engineering, Chinese University of Hong Kong, Shenzhen 518172, China, and also with the Peng Cheng Laboratory, Shenzhen 518066, China (e-mail: peipeizhu@link.cuhk.edu.cn).

Xiao Wang is with the School of Computer Science and Technology, Anhui University, Hefei 230601, China. (email: wangxiaocvpr@foxmail.com)

Lin Zhu is with the School of Computer Science, Beijing Institute of Technology, Beijing 100081, China. (e-mail: linzhu@pku.edu.cn)

Zhenglong Sun is with the School of Science and Engineering, Chinese University of Hong Kong, Shenzhen 518172, China (e-mail: sunzhenglong@cuhk.edu.cn).

Wei-Shi Zheng is the School of Computer Science and Engineering, Sun Yat-Sen University, Guangzhou 510275, China, and also with with the Peng Cheng Laboratory, Shenzhen 518066, China (e-mail: wszheng@ieee.org).

Yaowei Wang is with the Peng Cheng Laboratory, Shenzhen 518066, China (e-mail: wangyw@pcl.ac.cn).

Changwen Chen is with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong 999077, China (e-mail: changwen.chen@polyu.edu.hk).

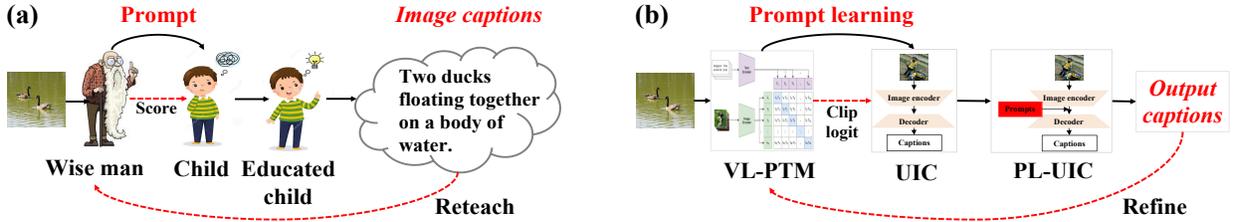Corresponding author: Yaowei Wang, Changwen Chen

Fig. 1: The prompt-based learning for unpaired image captioning. (a) An ignorant child learns knowledge from a wise man to describe an image. (b) The PL-UIC is developed to utilize prompts, learned from a vison-language pre-trained model (VL-PTM), to generate captions for images. These prompts of each image contain abundant contextual information of the matched images and texts, which is the information the previous UIC model does not have but is indispensable. The red dotted lines represent the process of caption scoring and reteaching in (a), which corresponds to the process of caption filtering and UIC model refining in (b).

this idea is similar to coaching a child to describe an image with the help of a wise man, as illustrated in Fig. 1. The child may describe the content of the given image more accurately if the wise man could give some important prompts. Therefore, two kinds of prompts are designed, *i.e.*, the *semantic prompt* and *metric prompt*, to imitate such a learning paradigm. More specifically, the visual images are taken as input to the semantic prompt extraction module, consisting of the pre-trained VL-PTMs (CLIP [23] is used in our experiments) and a feed-forward layer. The predicted prompt vector will be fed into the CLIP model to adjust its context and then align the image and prompt accurately. Then, the semantic prompt is injected into the adversarial learning based UIC framework for more intelligent and robust caption generation. The metric prompt is designed to transform the aforementioned unsupervised captioning optimization into a semi-supervised manner. As the pseudo captions can be obtained using the basic captioning model, and then high-quality samples can be filtered based on the metric prompt to polish the captioning model in an iterative way. As elaborated in Fig. 2, the metric value of an image and a caption obtained from the CLIP model can serve as the metric prompt. This semantic prompt-based learning and metric prompt guided high-quality sample filtering are integrated to form a strong caption generator *without using annotated aligned image-text pairs.*

To sum up, the contributions of this paper can be summarized as the following three aspects:

• We have developed a novel Prompt-based Learning scheme for Unpaired Image Captioning, termed PL-UIC, which can make full use of VL-PTMs for high-performance captioning. To the best of our knowledge, it is the first work to infer the cue information (*i.e.*, the prompt) about a given image that exists in the large VL-PTMs for the UIC task.

• Two types of simple yet effective prompt schemes have been designed for the UIC task, *i.e.*, the semantic prompt and the metric prompt. The semantic prompt has been devised to extract vision-aware prior knowledge via the textual format and taken as input to guide the caption generation. The metric prompt guided pseudo label filter has been designed to help improving the selection of highly-matched image-caption pairs, which enabled us to enhance the proposed UIC model in a semi-supervised way.

• Extensive experiments have been carried out on the widely used COCO and Flickr30k datasets to demonstrate that the proposed prompt-based learning can efficiently boost the performance in caption generation. The design principle proposed in this research can also be applied to other applications that demand prior knowledge.

## II. RELATED WORK

In this section, we review the related works on supervised image captioning, unpaired image captioning models, and prompt learning.

**Image Captioning.** Classical image captioning implements the encoder-decoder architecture, which first encodes images into features and decodes these image features into sentences [4], [31], [32] later. The goal of these models is to maximize the probability of generating the correct captions, relying on tremendous image-caption pairs [33], [34]. To solve the problem of the tight dependence on the costly image-caption pairs, some researchers proposed to use fewer and fewer pairs to complete the task, including novel object captioning [35], [36] and semi-supervised image captioning [37], [38]. Despite the promising captioning reform that has been completed, the costly paired image-caption datasets are indispensable in the training process. Distinct from all these works, we attempt to complete UIC without requiring any image-caption pair.

**Unpaired Image Captioning.** Distinct from the afore-mentioned supervised image captioning, UIC is to generate descriptions for images without requiring any image-caption pairs. Feng *et al.* [21] tackled UIC via adversarial learning and the alignments between images and visual concepts. Although the UIC is achieved, the captioning performance has a big gap between UIC and supervised image captioning due to the weak vision-language correlations. Thus, some researchers put effort into enhancing the weak cross-domain correlations in the task [14], [22], [39]. For example, Laina *et al.* [22] proposed to narrow the domain gap between images and languages by a shared embedding space of images and visual concepts. Also, several works focused on adopting scene graph modeling in UIC to align more textual information with images, including relationships and attributes [40]–[43]. The following methods achieved better captioning performance since much more vision-language alignment is explored in UIC. Despite the

| Images | | | | |
|---|---|---|---|---|
| Captions | a man and a girl playing a video game on a couch . | a group of people standing in front of a building . | a large elephant standing next to a baby elephant . | a man with a backpack standing next to a suitcase . |
| Metrics | 13.4922 | 14.5000 | 28.7031 | 34.5312 |

Fig. 2: The metric prompt of image-caption pairs generated by the prompt-based UIC model. The higher value of the metric value, the higher quality of the image-caption pairs.

enhanced captioning performance, there is still much room for improvement due to the neglect of the majority of vision-language correlations. Different from all these works, we attempt to utilize the prompt-based learning in UIC, which is aided by the pre-trained CLIP model [23] with abundant vision-language prior knowledge.

**Prompt-based Learning.** Prompt-based learning methods are proposed in natural language processing (NLP), which aim to reduce or obviate the requirement for large supervised datasets in the downstream tasks [27]. When learning a language model, task-specific prompt functions are designed to model the probability of the text prompt itself, and this probability is then adopted to tackle the task [44]–[46]. During the training process, these prompt functions are utilized to instruct pre-trained models to perform corresponding tasks conditionally [28], [47]. The developments of prompt-based methods make zero-shot and few-shot learning in NLP tasks come true [48], [49]. Inspired by these developments, researchers tried to extend it into vision tasks, such as image classification [50], visual question answering [51], [52], image captioning [52], etc. In this work, to incorporate the vision-language prior knowledge into the UIC task, we propose a PL-UIC model which is inspired by the prompt-based learning. Distinct from all existing works, we take a two-step procedure: learning the semantic prompt of each image from a large pre-trained VL-PTMs and applying these prompts as additional guidance to generate captions firstly, and then one metric prompt is utilized to design a high-quality pseudo label filter to further enhance the captioning performance.

## III. METHODOLOGY

In this section, we will first introduce an overview of the proposed PL-UIC. Then, the unpaired image captioning is reviewed briefly. After that, the motivation of semantic prompt and metric prompt is discussed carefully. Later, the designed semantic prompt for UIC is introduced in detail. Finally, we discuss the metric prompt guided pseudo label filter for UIC model polish.

### A. Overview

The target of PL-UIC is to train an unpaired image captioning model guided by prompts. As shown in Fig. 3, the overall framework consists of three parts, *i.e.*, the semantic prompt extraction module, the adversarial learning based UIC model, and the metric prompt based generator refining module. Specifically, the semantic prompt extraction module contains

the CLIP model and a feed-forward layer to provide the semantic prompt for each image in the image datasets. Then, the semantic prompt guided UIC model is trained to generate captions by adopting unpaired image-text samples. Later, the caption generator refining module is trained using pseudo labels. These pseudo labels are selected by the designed metric prompt, which can measure the correlations of image-text pairs. We combine the designed two prompt-related modules into the common adversarial learning UIC framework, which achieves promising results on two widely used captioning datasets.

### B. Preliminary: Unpaired Image Captioning

Unpaired image captioning is to describe images without using any aligned vision-text pairs. The image dataset with $N_i$ images is represented by $\mathcal{I} = \{I_1, I_2, ..., I_{N_i}\}$, the unpaired language dataset is represented as $\mathcal{S} = \{S_1, S_2, ..., S_{N_s}\}$ with $N_s$ sentences, and the captions generated by the UIC model are denoted as $\mathcal{C} = \{C_1, C_2, ..., C_{N_c}\}$, where $N_c$ means the number of generated captions. For simplicity, we utilize $I$, $C$, and $S$ to represent an image, a virtual generated caption, and a real sentence, respectively.

The common pipeline of UIC is illustrated in Fig. 4, where a CNN-LSTM network is adopted for the encoder-decoder framework. Firstly, a CNN encoder is utilized to extract the image features $f_I$. Then, an LSTM decoder network $LSTM_g$ is adopted to transform the image features into texts $C = \{c_1, c_2, ..., c_{n_c}\}$, where $n_c$ is the number of words in one caption. And this procedure can be written as:

$$x_0 = FC(f_I), \quad (1)$$

$$x_t = W_e c_t, t \in \{1, ..., n-1\}, \quad (2)$$

$$[o_{t+1}, h_{t+1}] = LSTM_g(x_t, h_t), t \in \{0, ..., n-1\} \quad (3)$$

$$c_t \sim o_t, t \in \{1, ..., n\} \quad (4)$$

where $x_t$, $c_t$, $o_t$, and $h_t$ are the input of the $LSTM_g$ decoder layer, a one-hot vector representation of the outputted word, the probability of every word in the dictionary, and the hidden state of the LSTM at the $t$-th time step, separately. $x_0$ is the initial input of the decoder, and one feed-forward layer $FC$ is utilized to adjust the image features. $h_0$ is a zero vector for the initial hidden states of the decoder. $W_e$ is for word embedding. Besides, an LSTM discriminator $LSTM_d$ is adopted to differentiate a real sentence $S$ and a virtual sentence $C$, which is utilized to make the generated captions as real sentence as possible. Formally,

$$[q_t, h_t] = LSTM_d(x_t, h_{t-1}), t \in \{1, ..., n\}, \quad (5)$$

where $q_t$ represents the probability that the outputted sequential words $\{c_1, ..., c_t\}$ of one caption $C$ or $\{s_1, ..., s_t\}$ from a real sentence $S$ are regarded as a real sentence $S$.

Although the adversarial learning based UIC models work well in some scenarios, however, their overall performance is still limited due to the less cross-domain knowledge. We think this situation can be alleviated by introducing vision-language prior knowledge.
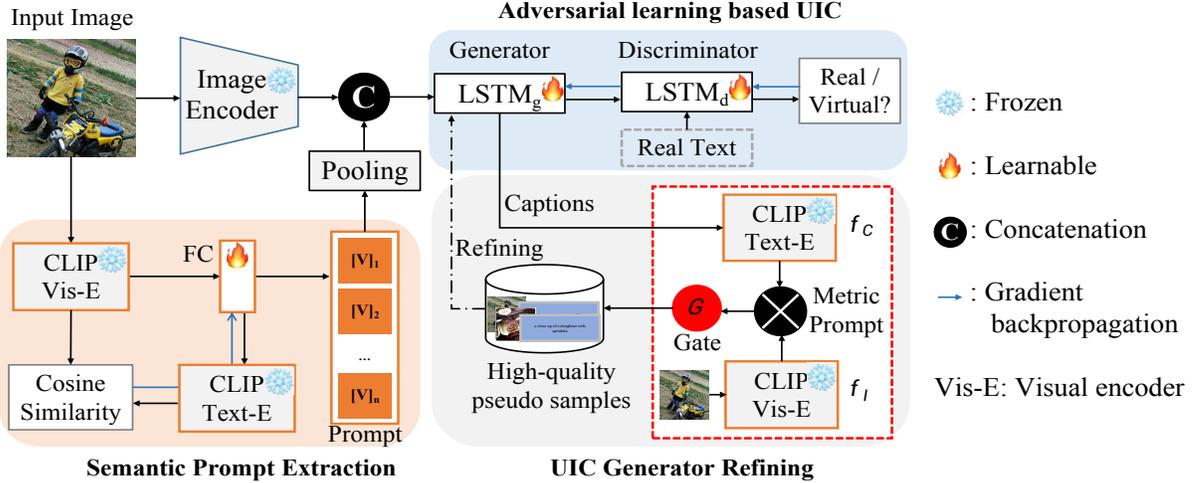
Fig. 3: An overview of the proposed PL-UIC framework, which contains three modules, *i.e.*, the semantic prompt extraction, the adversarial learning based UIC, and the UIC generator refining module. Firstly, the designed semantic prompt extraction module, a feed-forward layer with the frozen CLIP model, is trained to generate the semantic prompt for each image. Then, these prompts are implemented as guidance in the adversarial learning based UIC. Finally, the proposed pseudo label filtering is utilized to re-train the UIC generator iteratively, which selects high-quality image-caption pairs as the pseudo samples through the measurement of the metric prompt.
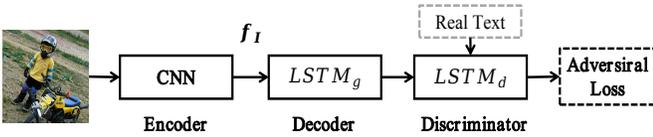


Fig. 4: The adversarial learning for UIC.

### C. Motivation of Semantic Prompt and Metric Prompt

The motivation of the semantic prompt and metric prompt is the vision-language prior knowledge fully reflected in the feature representation of the VL-PTMs. For example, the CLIP model used in this paper performs well in zero/few-shot learning since its abundant language-aware visual features. The CLIP model includes one visual encoder and one sentence encoder. The visual encoder is used to extract the deep image features, which can be implemented with ResNet-101 [53] or ViT [54]. The text encoder is used to learn the sentence features using a Transformer network [55]. Based on a large amount of image-text dataset (400 million) crawled from the Internet, the CLIP is optimized to align features of these two domains using contrastive learning [56]. Specifically, its target is to maximize the cosine similarity between the matched image-caption pairs and minimize the unmatched ones. Therefore, the cross-domain vision-language prior knowledge can be mastered by the CLIP.

To intuitively exhibit the prior knowledge of the CLIP model, two types of examples are visualized in Fig. 5. Fig. 5 (a) is the zero-shot image classification, which reflects the prior knowledge from the vision domain to the language domain. Fig. 5 (b) is about highlighting the most related regions in an image for the given texts without training, which illustrates the prior knowledge from the language domain to the vision domain.

The main target of this paper is to extract the learned knowledge from CLIP for the guidance of unpaired image captioning from three aspects, *i.e.*, the visual encoder, the semantic prompt, and the metric prompt. The visual encoder of the CLIP model is adopted to extract the image features for UIC directly, where abundant language-aware image features are provided to guide the caption generation. The details of the two prompts are introduced in the following sections.

### D. Semantic Prompt for Unpaired Image Captioning

For unpaired image captioning, the performance of existing methods is restrained by the limited cross-domain knowledge. To alleviate the issue, a semantic prompt extraction module is designed to explore more cross-domain prior knowledge from the VL-PTMs. The prompt extracted from the module will be taken as guidance for the caption generation.

**Semantic Prompt Extraction** is used to draw vision-aware textual knowledge from pre-trained vision-language models, *i.e.*, the semantic prompt for UIC task. The semantic prompt is represented by $\hat{P}_i \in \{\hat{P}_1, \hat{P}_2, ..., \hat{P}_{N_p}\}$ for each image $I_i$ in the image datasets. For simplicity, we utilize $\hat{P}$ to represent a semantic prompt. Different from the prior knowledge obtained from the fixed visual encoder, the semantic prompt is obtained by training with the frozen CLIP model via given visual images. In detail, the features from the image encoder of the CLIP model are taken as the input, as shown in Fig. 3. For practical implementations, we adopt one feed-forward layer to achieve this goal. Then, the output of the feed-forward layer $FC$ will be fed into the text encoder for subsequent processing. Formally,

$$p = FC(f_I), \tag{6}$$

$$f_p = TE(p), \tag{7}$$

where $p \in \mathbb{R}^y$ means the output of the $FC$ layer. $y$ denotes the dimension of the semantic prompt, and it will be reshaped into
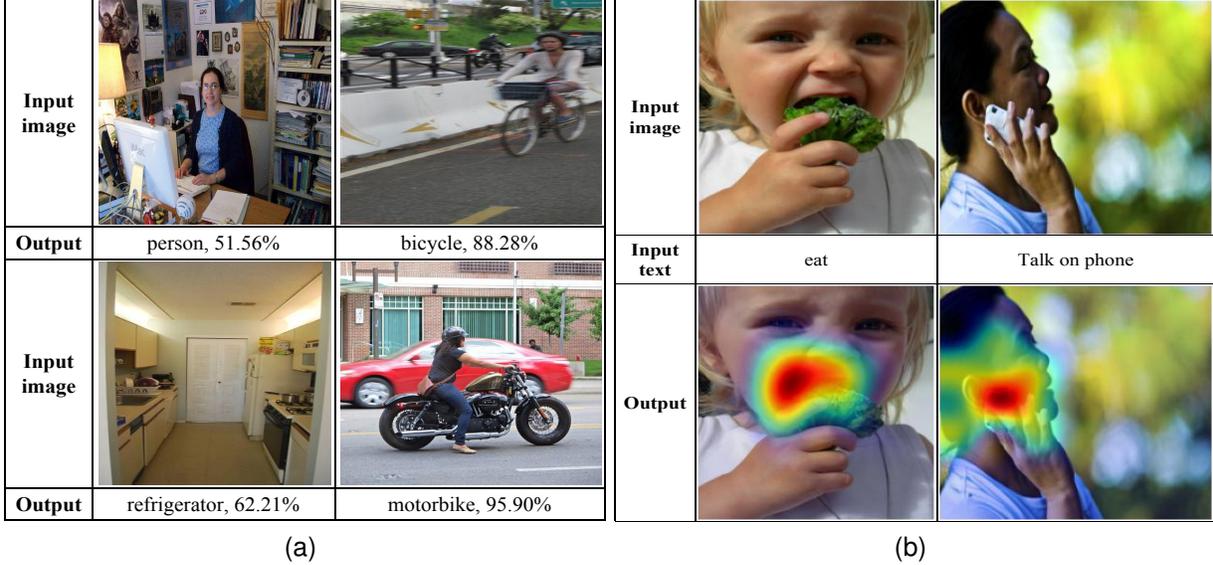
Fig. 5: The illustration of the vision-language prior knowledge in the CLIP model. (a) Examples of the zero-shot image classification, and (b) The highlighted heatmaps for an input caption without training.

the length of the prompt and the dimension of each prompt, such as $8 \times 512$. $TE$ means the CLIP text encoder and $f_p$ means the encoded prompt features. Note that the parameters of the image and text encoder of the CLIP model are frozen. The objective of this module is to maximize the cosine similarity of the image features $f_I$ and the matched prompt features $f_p$, and minimize the cosine similarity of unmatched pairs. The cosine similarity $\hat{S}$ can be written as:

$$\hat{S} = \frac{\sum_i f_I \cdot f_p}{||f_I|| \times ||f_p||}, \tag{8}$$

where $\sum_i$ means the summation of all elements in $f_I \cdot f_p$.

To obtain the semantic prompt $\hat{P}$ for each image, the "SOS", "CLS", and "EOS" embeddings are concatenated at the start and end of the FC layer, which represents the start and end of the prompts, respectively. Moreover, the positional pointers of the CLIP model are added to the former embeddings, and the results $\hat{P}$ are taken as the learned prompt for each image. Formally,

$$\hat{P} = [E_{SOS}, p, E_{CLS}, E_{EOS}] + E_{pos}, \tag{9}$$

where $[\cdot]$ means the concatenation operation, and $E$ represents various types of embeddings in CLIP. The $E_{pos}$ indicates the positional embeddings in the CLIP model. And we utilize $[V]_1, [V]_2, ..., [V]_n$ to represent a semantic prompt $\hat{P}$, where $[V]_1$ denotes the features for a textual word. Note that only the image dataset of UIC is utilized for training this module, and the architecture of the trainable layer in the semantic prompt extraction module can be replaced by other networks for better performance, which we left in the future work.

**Caption Generation** To integrate the semantic prompt $\hat{P}$ into the UIC task, we take the prompts of each image as part of the input of the UIC network. Firstly, a pooling block $Pool$ is adopted to reduce the computation related to the prompts.

Then, the prompts $f_{\hat{P}}$ and image features $f_I$ are concatenated together and are taken as the input features of the caption generator $LSTM_g$ at the first time step $x_0$ to guide the caption generation. Formally,

$$f_{\hat{P}} = Pool(\hat{P}) \tag{10}$$

$$x_0 = [f_I, f_{\hat{P}}]. \tag{11}$$

Then, the sentence decoder $LSTM_g$ is utilized to generate texts based on $x_0$ and the process takes the same principle as in the common UIC model, as shown in Equation 2, 3, and 4. After that, a virtual text will be inputted into the discriminator $LSTM_d$ to differentiate whether it is real or not, shown in Equation 5. Through the training of sentence decoder and discriminator, we will obtain relatively accurate captions for each image $\{I, C\}$ through the trained semantic prompt-based UIC model.

### E. Metric Prompt for UIC Generator Refining

After the training of semantic prompt-based UIC is done, we can get better results than the prompt-free UIC model. Due to the limitation of unsupervised learning of captioning generator, the unpaired image-text samples are still hard to optimize, which bring us sub-optimal performance only. Therefore, we are inspired to transform the unpaired image-text samples into paired ones using generated captions for the training images. Then, these data can be utilized as pseudo labels to train a fully-supervised caption generator.

Considering that the UIC model is still weak in the initial stage, indiscriminate leveraging these data may harm our model. In this section, a pseudo label filtering scheme is proposed to filter out the low-quality image-caption pairs and leave the high-quality ones, which enables us to polish our model in a more stable way. Based on the proposed scheme, the
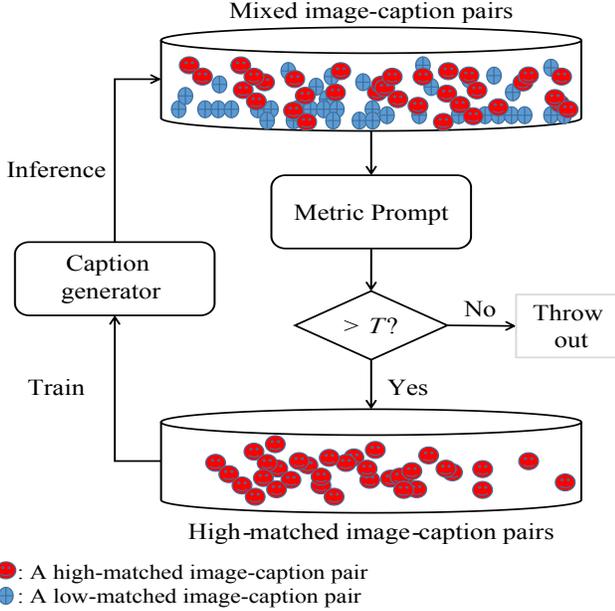
Fig. 6: The process of the caption generator refining. The generator first outputs mixed image-caption pairs, including low-matched ones and high-matched ones. Then, the metric prompt of one image-caption pair is computed and compared with a threshold $T$. Suppose the metric prompt is higher than the $T$, the image-caption pair will be regarded as a high-matched pair and be selected to train the caption generator again.

model refining and label filtering processes can be executed iteratively. The whole process is illustrated in Fig. 6. In the filtering process, the metric prompt is designed to select the high-matched image-caption pairs. Unlike the prior knowledge in the semantic prompt, the metric prompt contains different cross-domain knowledge that indicates whether an image and a caption are matched.

**Metric Prompt.** The metric prompt $L$ is adopted as a criterion to filter out the low-quality image-caption pairs, which represents the cosine similarity between an image and a caption. Specifically, we input the visual image and its caption into the pre-trained CLIP model. Therefore, we can get the prompt via:

$$L = \boldsymbol{w} \times \left( \frac{\sum_i \boldsymbol{f_I} \cdot \boldsymbol{f_C}}{||\boldsymbol{f_I}|| \times ||\boldsymbol{f_C}||} \right), \qquad (12)$$

where $\boldsymbol{w}$ and $\boldsymbol{f_C}$ are the weights and the caption features, respectively.

As illustrated in Fig. 2, the high-quality image-caption pair has a high metric value, and the low-quality image-caption pair has a low metric value. Suppose we have a metric gate $G$ with threshold $T$, if the metric is lower than the threshold $T$, the raw image caption $\{\boldsymbol{I}, \boldsymbol{C}\}$ will be filtered out. Otherwise, the raw image caption pair will be preserved to train the semi-supervised caption generator.

**Iterative Sample Filtering and Caption Generator Refining.** In the experiments, we conduct the captioning generator refining and sample filtering in an iterative way. As the captioning generator will output higher quality captions when a new iterative loop is completed (*i.e.*, the sample selection

and generator re-training), therefore, it can in turn contribute to training a better semi-supervised caption generator. The filtering and re-training procedure can be executed multiple times to further enhance the captioning ability of the generator. We provide the experimental analysis of the influence of the iteration number in Section IV-D4.

### F. Training and Inference

**Training Phase.** In this work, we apply a multi-stage optimization scheme for the network since the proposed framework consists of multiple modules which are required to be optimized one by one.

*Stage-I:* The semantic prompt extraction module, *i.e.*, a single FC layer embedded on the CLIP model, is firstly trained on an image dataset. As a result, the semantic prompt of each image will be obtained, which represents the contextual features of images and texts, and will be utilized as a part of the input of the UIC model.

*Stage-II:* Aided by the aforementioned semantic prompts, the UIC model is trained by the image features and semantic prompts, which is optimized by adversarial learning and visual concept alignments. The trained model can be adopted to generate the captions of images.

*Stage-III:* With the generated image-text pairs, the captioning model is re-trained to enhance the captioning performance. In this stage, metric prompt-based pseudo label filtering is implemented to select the high-quality image-text pairs as the pseudo labels. These labels can be utilized to train a better semi-supervised model. The filtering and re-training strategies loop multiple times to enhance the captioning performance.

**Inference Phase.** Given the testing images, we first feed them into the semantic prompt extraction module to obtain the semantic prompt. Then, these images and corresponding prompts are inputted into the trained semi-supervised model to obtain high-quality captions.

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metrics

**Datasets.** Different from the standard image captioning task, which utilizes paired image-text samples, the UIC task adopts unpaired datasets. In our experiments, five datasets are involved, including two image datasets (*i.e.*, the COCO [57] and Flickr30K [58]) and three language datasets (*i.e.*, the COCO captions, Shutterstock sentences [21], and Conceptual sentences [59]).

A brief introduction to these datasets is given below. 1). COCO [57] is one of the popular benchmarks of the image captioning tasks, which contains 123,287 images. The common Karpathy split [60] is utilized in the experiments, with 113,287 training images, 5,000 validation images, and 5,000 testing images. 2). Flickr30K dataset consists of 31,783 images with 29,783 training images, 1,000 validation images, and 1,000 testing images, which is the common split [58] for the experiments on image captioning. 3). COCO caption dataset is annotated manually. Each image is annotated with 5 captions. For training, the COCO images and captions are utilized in an unpaired way. A word vocabulary is built by words

TABLE I: Comparisons of the proposed PL-UIC with related methods on the testing split of COCO dataset for UIC. SP: semantic prompt. * represents the re-trained version of the related works with the same powerful CLIP visual encoder. "(finetuning SP)" means that all parameters in the semantic extraction module are trainable. "(3)" represents three loops of metric prompt based model refining.

| Method | B4 | M | R | C | S |
|---|---|---|---|---|---|
| Feat2sen [21]* | 15.7 | 18.0 | 41.5 | 48.7 | 10.8 |
| Recons-Align [21]* | 22.3 | 21.5 | 47.1 | 71.9 | 14.4 |
| WS-UIC [43]* | 22.6 | 20.9 | 46.6 | 69.2 | 14.1 |
| PL-UIC (finetuning SP) (3) | 24.3 | 22.3 | 48.7 | 75.6 | 14.5 |
| PL-UIC (3) | 24.9 | 22.5 | 49.3 | 77.3 | 14.9 |
| PL-UIC | **25.0** | **22.6** | **49.4** | **77.9** | **15.1** |

TABLE II: Comparisons with related methods for UIC using independent datasets. * represents the re-trained version of the related works with the same powerful CLIP visual encoder as the image encoder in UIC.

| Method | B4 | M | R | C | S |
|---|---|---|---|---|---|
| **COCO images + Shutterstock sentences** | | | | | |
| Feat2sen [21]* | 5.5 | 12.7 | 28.8 | 24.8 | 8.4 |
| Recons-Align [21]* | 6.6 | 12.9 | 28.0 | 31.3 | 8.9 |
| WS-UIC [43]* | 6.4 | 12.4 | 29.2 | 26.7 | 7.9 |
| PL-UIC | **10.0** | **16.2** | **35.8** | **45.8** | **11.5** |
| **Flickr30k images + Conceptual sentences** | | | | | |
| Feat2sen [21]* | 7.8 | **11.4** | 29.7 | 9.5 | **6.5** |
| Recons-Align [21]* | 5.2 | 10.1 | 26.2 | **11.4** | 5.7 |
| WS-UIC [43]* | 6.2 | 10.3 | 32.4 | 7.6 | 4.7 |
| PL-UIC | **9.7** | 10.9 | **32.7** | 8.8 | **6.5** |

TABLE III: Ablation study to verify the effectiveness of each module in the proposed PL-UIC. W/o P: without using VL-PTMs. VP: vision encoder of VL-PTMs. SP: semantic prompt. MP: metric prompt. MP(1): metric prompt with 1 iteration. MP(3): metric prompt with 3 iterations.

| Method | | | | | Evaluation Metric | | | | |
|---|---|---|---|---|---|---|---|---|---|
| W/o P | VP | SP | MP (1) | MP (3) | B4 | M | R | C | S |
| ✓ | | | | | 18.4 | 18.3 | 43.3 | 56.7 | 11.5 |
| ✓ | | ✓ | | | 19.0 | 18.9 | 43.9 | 58.6 | 11.6 |
| | ✓ | | | | 20.7 | 20.1 | 45.9 | 64.4 | 12.7 |
| | ✓ | ✓ | | | 21.6 | 20.6 | 46.1 | 67.2 | 13.4 |
| | ✓ | | ✓ | | 24.0 | 22.1 | 48.0 | 74.8 | 14.0 |
| | ✓ | ✓ | ✓ | | 24.2 | 22.1 | 48.7 | 75.1 | 14.7 |
| | ✓ | ✓ | | ✓ | **24.9** | **22.5** | **49.3** | **77.3** | **14.9** |

presented no less than 4 times in captions of the training images. 4). Shutterstock sentence dataset is crawled from the Shutterstock website by [21], which contains 2,282,444 different image captions. 5). Conceptual sentences [59] are automatically collected from webs, which include 3.3 million image-caption pairs. The captions are utilized as a sentence dataset.

**Evaluation Metrics.** Five evaluation metrics, including BLEU-4 (B4) [61], METEOR (M) [62], ROUGE (R) [63], CIDEr (C) [64] and SPICE (S) [65], are adopted to evaluate the captioning performance, whose values are computed based on the ground-truth captions of test images. Among these metrics, BLEU-4 [61] was proposed for machine translation by computing the overlapping units of the machine translations and human translations in 2002; language-specific evaluation was brought to evaluate machine translation of any language in 2014 and the metric is called METER [62]; ROUGE [63] was mainly designed for automatic abstracting by counting the overlapping units of computer-outputted summaries and human labeled summaries in 2004; in 2015, CIDEr [64] was proposed for evaluating image captioning by capturing human judgment of consensus; and the semantic propositional content of human caption evaluation was considered in SPICE [65] for evaluating image captioning in 2016. The higher the metric values, the better the experimental performance of UIC.

### B. Implementation Details

In UIC, we utilize the pretrained ViT-b16 CLIP model as the image encoder. The object concepts are extracted by Faster-RCNN [66] related object detection model. The LSTM with 512 hidden states is utilized for the sentence decoder and discriminator. The learning rate of the caption generator is set as 0.001 at the training **Stage II** in Section III-F. At the training **Stage III**, the learning rate is set as 0.00001. The architecture of the semi-supervised image captioning model is the same as the generator of the UIC model. The whole PL-UIC model is trained on 2 V100 GPUs.

For the semantic prompt extraction and pseudo label filtering, the ViT-b16 visual encoder [23] is utilized in the CLIP model. For pseudo label filtering of COCO image datasets, the threshold of metric prompt is set as 30. For the Flickr30k

image dataset, the threshold of the metric prompt is set as 26 due to the smaller scale of the image dataset.

### C. Comparison on Benchmarks

In this section, two kinds of experimental comparisons are carried out to validate the effectiveness of the proposed approach. The **first** one is to compare the PL-UIC with other approaches on the COCO dataset using unpaired image and caption samples. The **second** setting is to compare PL-UIC with algorithms adopting independent datasets, *i.e.*, COCO images with Shuttershock sentences, and Flickr30k images with Conceptual captions. We compare the PL-UIC with the related methods Feat2sen [21], Recons-Align of [21], and WS-UIC [43]. Feat2sen relies on pseudo labels, outputted by an object-to-sentence model [21], to train a fully-supervised caption generator, where only sparse knowledge of the text domain and the image domain is concerned. Recons-Align depends on adversarial learning and the visual concept reward, both of which cannot be implemented to explore the contextual vision-language information thoroughly. WS-UIC relies on one more unrecognized object loss to improve the alignment between the objects and images, which is still weak to explore plentiful of cross-domain knowledge. As a result, the experimental performance of these methods is severely constrained. Different from all these works, we adopt the semantic prompt and metric prompt of each image, containing full of vision-language prior knowledge, as additional guidance for caption generation.
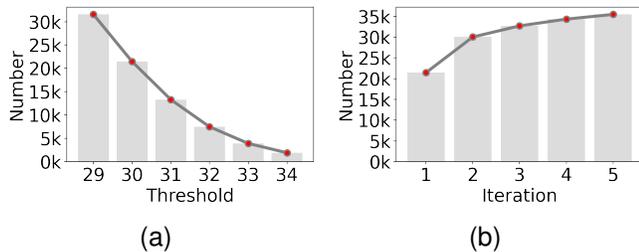
Fig. 7: The number of pseudo labels in (a) different thresholds of the metric gate and (b) different iterations of the metric prompt related model refining. The threshold of the metric gate is set as 30 in (b).

For the first setting, we report the experimental results of ours and the related UIC models in Table I. * means the re-trained version of methods with the same powerful image encoder of PL-UIC. Besides the method Feat2sen, Recons-Align, and WS-UIC, we compare with the PL-UIC (finetuning SP) (3) and PL-UIC (3). The PL-UIC (finetuning SP) (3) means training PL-UIC with finetuned semantic prompts and 3 iterations of the generator refining, where the semantic prompt are extracted by training the whole model of the semantic prompt extraction instead of one feed-forward layer. PL-UIC (3) means that the experiment is carried out by 3 iterations of the caption generator refining. Generally speaking, the experimental results of the proposed PL-UIC achieve the best captioning performance over all five evaluation metrics. For example, the C value of PL-UIC is 6% higher than the C value of Recons-Align [14]. The results basically demonstrate that the extracted prompts, with plentiful contextual vision-language information, can provide efficient guidance for captioning. And the iterative pseudo label filtering scheme can provide high-quality labels for generator refining. Compared the PL-UIC (finetuning SP) (3) with PL-UIC (3), the higher experimental results of PL-UIC (3) verify that the training of only one feed-forward layer of the designed semantic prompt extraction method is more effective than training all parameters of the whole module.

For the second setting, we compare the proposed approach with related methods by utilizing independent datasets. From Table II, it is easy to find that our proposed PL-UIC outperforms all the compared methods by a significant margin in the experiments with COCO images and Shutterstock sentences. More in detail, the C value of PL-UIC is 45.8. Meanwhile, the value of Recons-Align is only 31.3, which is 14.5 lower than PL-UIC. For the experiments with Flickr30k images and Conceptual sentences, the overall captioning performance of PL-UIC exceeds all the compared methods, although it has weak performance on metrics M and C. These experimental results fully demonstrate that our PL-UIC model has the strongest captioning ability in the independent dataset settings since the vision-language prior knowledge learned from VL-PTMs.

To compare the captioning ability of all these methods more intuitively, two images with the corresponding captions generated by these methods are elaborated in Fig. 8. From the figure, we can observe that the captions outputted by the

proposed PL-UIC are better than the other compared methods. Take the first image as an example, the Feat2sen* method outputs "a young boy with a drink in front of a person holding a snack in his hand.", where only the concept "a young boy" is accurate. The method Recons-Align* and WS-UIC* generate incorrect concept "donut" and "holding", respectively. PL-UIC (finetuning SP) (3) generates a totally wrong caption "a person riding a bike down a street next to a building". The captions outputted by the proposed method PL-UIC (3) and PL-UIC are semantically correct with "a little boy", "eating", and "a piece of cake", which contain almost all the important concepts as in the ground-truth captions. Visualizing these captions intuitively demonstrates the superiority of the proposed PL-UIC with plentiful cross-domain prior knowledge.

### D. Ablation Studies

*1) Component Analysis:* To facilitate readers to have a more thorough understanding of the proposed approach, extensive experiments for the component analysis have been conducted on the testing split of the COCO dataset. As shown in Table III, the following modules are implemented or not for the comparisons: "W/o P" means utilizing an image encoder from a relatively small off-the-shell model instead of the VL-PTMs; "VP" is adopting the visual encoder of CLIP as the image encoder of UIC; "SP" means the semantic prompt; "MP" denotes the metric prompt; "(1)" represents 1 iteration of the generator refining; and "(3)" means 3 iterations of the generator refining.

As shown in Table III, we can find that the model "W/o P" only achieves 18.4, 18.3, 43.3, 56.7, and 11.5 on the B4, M, R, C, and S metric, respectively. The experimental results of the method "W/o P" are obviously lower than the method "VP", which demonstrates the effectiveness of the language-aware visual prior knowledge from the CLIP model. When introducing the semantic prompt to the "W/o P", the overall results can be improved to 19.0, 18.9, 43.9, 58.6, and 11.6, respectively. Also, when the other two methods "VP" and "VP" with "MP(1)" are combined with "SP", the experimental performances are also better than the methods without "SP". These results validate that our designed semantic prompt based on VL-PTMs is beneficial for the UIC task due to the abundant vision-aware language prior knowledge. When we integrate the metric prompt-related module, our results can be further improved. Specifically, the model with "MP(1)" is better than the "non-MP(1)" version. For instance, the B4, M, R, C, and S can be improved up by +3.3, +2.0, +2.1, +10.4, and +1.3, respectively, when one iteration of the metric prompt-related module is used with "VP". Obviously, more iterations can bring better results. These results fully validate the effectiveness of our proposed modules for the UIC task due to the prior metric knowledge of the VL-PTMs. We hope our model can bring new insights to exploring contextual vision-language information for the VL-PTMs or prompt learning-based UIC.

*2) Effect of Different Image Encoders:* To investigate the effects of the language-aware prior knowledge from different vision backbones, we conduct the UIC experiments with three image encoders from different vision backbones, *i.e.*, ResNet-101, ViT-b32, and ViT-b16, illustrated in Table IV. Moreover,

| Image | Method | Caption | Image | Method | Caption |
|---|---|---|---|---|---|
| | Feat2sen* | a young boy with a drink in front of a person holding a snack in his hand . | | Feat2sen* | a group of people standing around a boat . |
| | Recons-Align* | a little boy eating a donut with a face | | Recons-Align* | a boat in the water on a sunny day . |
| | WS-UIC* | a little boy holding a piece of cake | | WS-UIC* | a boat in the water near a tree |
| | PL-UIC (finetuning SP) (3) | a person riding a bike down a street next to a building | | PL-UIC (finetuning SP) (3) | a wooden boat sitting next to a body of water . |
| | PL-UIC (3) | a little boy eating a piece of cake . | | PL-UIC (3) | a boat sitting on top of a body of water . |
| | PL-UIC | a little boy eating a piece of cake . | | PL-UIC | a boat sitting on top of a body of water . |
| | Ground-truth | Baby boy at the table eating cake frosting off his hand. | | Ground-truth | A blue boat docked on a green lush shore. |

Fig. 8: Examples of the captions generated by different methods on the COCO dataset with ground-truth captions.

the experiments with Inception-V4, not related to the VL-PTMs, are carried out for comparison. In the table, we can observe that all these three VL-PTMs related experiments achieve more promising captioning performance than the non-VL-PTMs experiments. For example, the VL-PTMs related ResNet-101 experiments perform 20.2, 20.1, 45.0, 63.3, and 12.8 over metric B4, M, R, C, and S, respectively, while Inception-V4 related experiments only achieve 19.0, 18.9, 43.9, 58.6, and 11.6 over the same metrics, respectively. These experiments show that the language-aware vision prior knowledge in the image encoders is effective and efficient for the UIC task, no matter what type of vision backbone it is.

*3) Effect of Semantic Prompt:* **Prompt Length.** To explore the influence of the length of the semantic prompt, three lengths are utilized to carry out the experiments, *i.e.*, 4, 8, and 16. As exhibited in Table VI, these three experiments have comparable experimental performance over evaluation metrics B4, M, R, and S. For example, these three experiments obtain 21.4, 21.6, and 21.6 over metric B4, respectively. As for metric C, the experiment of the semantic prompt with length 8 has a slight advantage over the other two experiments. In detail, experiments with prompt length 8 achieve 67.2 on metric C, while the experiments with prompt length 4 and 16 obtain only 66.8 and 65.5 on the same metric, respectively. We can conclude that the prompt length has a relatively bigger influence on the evaluation metric C than other metrics. And we choose the prompt length 8 in other types of semantic prompt-related experiments.

**Prompt from Different Vision Backbones.** To verify the effectiveness of the semantic prompt learned from different vision backbones, three different backbones, including ResNet-101, ViT-b32, and ViT-b16, are utilized in the experiments to extract the semantic prompt, respectively. For comparison, we carry out the experiments with the same image encoder but without the semantic prompt (W/o Prompt). As elaborated in Table V, all these semantic prompt-related experiments elaborate promising captioning performance over five evaluation metrics than the non-prompt experiments, demonstrating the generality of vision backbones for the semantic prompt extraction. For instance, the experimental results with the semantic prompt learned from the ResNet-101 backbone of the CLIP model achieve 21.1, 20.4, 45.7, 66.8, and 13.5 on metrics B4, M, R, C, and S, respectively, while the experiments without semantic prompt only perform 20.7, 20.1, 45.9, 64.4, and 12.7 over the same evaluation metrics, respectively.

*4) Effect of Metric Prompt:* **Iterative Generator Re-training.** The multiple iterations of the generator re-training are indispensable for training a better generator since higher-

TABLE IV: The effect of the image encoder from different vision backbones.

| Method | B4 | M | R | C | S |
|---|---|---|---|---|---|
| InceptionV4 | 19.0 | 18.9 | 43.9 | 58.6 | 11.6 |
| ResNet-101 | 20.2 | 20.1 | 45.0 | 63.3 | 12.8 |
| ViT-b32 | 19.8 | 19.2 | 44.8 | 58.6 | 12.0 |
| ViT-b16 | 21.3 | 20.4 | 45.8 | 66.5 | 13.5 |

TABLE V: The effect of the semantic prompt from different vision backbones.

| Method | B4 | M | R | C | S |
|---|---|---|---|---|---|
| W/o Prompt | 20.7 | 20.1 | 45.9 | 64.4 | 12.7 |
| ResNet-101 | 21.1 | 20.4 | 45.7 | 66.8 | 13.5 |
| ViT-b32 | 21.1 | 20.7 | 46.0 | 66.4 | 13.6 |
| ViT-b16 | 21.3 | 20.4 | 45.8 | 66.5 | 13.5 |

quality labels will be collected in the next loop. From Fig. 7 (b), we can observe that the number of pseudo labels increases with the number of iterations. As more high-quality labels can be used in the supervised learning phase, the captioning performance is getting better and better, as shown in Fig. 9. The values are higher and higher over all five evaluation metrics, including B4, M, R, C, and S. These two figures demonstrate the usefulness of the iteration scheme.

**Threshold.** The pseudo labels for UIC generator refining are generated by comparing the metric prompt of an image and a text with a predefined threshold. To achieve better performance, we set different thresholds in this experiment, as shown in Table VII. The higher the threshold, the less the number of pseudo labels, as illustrated in Fig. 7 (a). If the threshold equals 0, it means utilizing all generated image-caption pairs as the pseudo labels to refine the caption generator without filtering. Compared to the experiments with non-zero thresholds 29, 30, 31, and 32, the experimental results with 0 threshold are much worse. Since the threshold 33 and 34 are so high, the number of pseudo labels become very small with less dataset variety which limits the refining performance. Obviously, we can find that when the threshold is set as 30, the best captioning results can be obtained as 24.2, 22.1, 48.7, 75.1, and 14.7 for evaluation metrics B4, M, R, C, and S, respectively.

TABLE VI: The effect of prompt length.

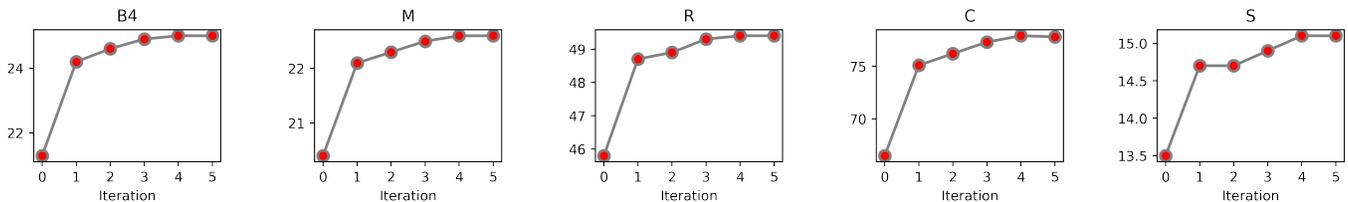| Length | B4 | M | R | C | S |
|---|---|---|---|---|---|
| 4 | 21.4 | **20.6** | 46.0 | 66.8 | **13.6** |
| 8 | **21.6** | **20.6** | **46.1** | **67.2** | 13.4 |
| 16 | **21.6** | 20.3 | 45.9 | 65.5 | 13.0 |

Fig. 9: The iteration influence of the pseudo label filtering on different evaluation metrics, including the B4, M, R, C, and S.

| Images |  |  |  |  |
|---|---|---|---|---|
| Captions | a man is shaving a sheep in his mouth | a group of people riding bike down a street | a black and white photo of a police officer on a motorcycle. | a group of people sitting at a table in a restaurant |
| Metrics | 18.5938 | 19.2500 | 17.3281 | 19.5469 |
| Images |  |  |  |  |
| Captions | a man standing in a kitchen preparing food. | a young boy and a girl sitting on a toilet. | a motorcycle parked in front of a building | a little girl sitting on a bench in a park. |
| Metrics | 34.4375 | 37.2812 | 34.1250 | 32.5000 |

Fig. 10: The metric value of image-caption pairs.

TABLE VII: The effect of different thresholds of the metric prompt.

| Threshold | B4 | M | R | C | S |
|---|---|---|---|---|---|
| 0 | 22.7 | 20.7 | 46.6 | 67.2 | 13.3 |
| 29 | **24.2** | 22.0 | 48.5 | 74.2 | 14.5 |
| 30 | **24.2** | **22.1** | **48.7** | **75.1** | **14.7** |
| 31 | 23.8 | 22.0 | 48.6 | 74.4 | 14.6 |
| 32 | 23.6 | 21.8 | 48.3 | 72.8 | 14.2 |
| 33 | 22.0 | 20.8 | 47.0 | 67.6 | 13.3 |
| 34 | 20.6 | 20.0 | 45.9 | 61.7 | 12.4 |

### E. Qualitative Results

**Visualization of the Metric Prompt of the Image-caption Pairs.** To better demonstrate the usefulness of the metric prompt in pseudo label filtering, we visualize the metric value of several image-caption pairs. As displayed in Fig. 10, we can find that the low-quality image-caption pairs have lower metric prompt values, and the high-quality image-caption pairs have higher metric prompt values. Take the first image-caption pair as an example, the image is about "a man makes a call in a telephone hall", but the generated caption is "a man is shaving a sheep in his mouth", which is mismatched except for the concept "a man". Thus, the metric value has a much lower value "18.5938". All these examples fully demonstrate that utilizing metric prompts to filter pseudo labels is reasonable.

**Visualization of the Distance between Semantic Prompts.** The square distances between similar images and different images are illustrated in Fig. 11, respectively. From the figure, we can observe that the square distance between similar images is smaller than the distance between different images. For example, the square distance between two images with a similar "toilet room" scene is only 6.67, while the distance is up to 112.74 between the image "birds flying in the sky" and the image "a lot of people with many cupcakes". The phenomenon of smaller prompt distance existing between images with similar scenes and larger prompt distance existing between images with different scenes demonstrates that the semantic prompt obtained the reasonable semantics of the images.

**Visualization of the Generated Captions.** Fig. 12 shows several representative captions outputted by multiple methods, *i.e.*, VP, VP + SP, PL-UIC(1), and PL-UIC. We can observe that these various UIC models can output reasonable descriptions by using the semantic prompt and the metric prompt. Let us take

| | | |
|---|---|---|
| **Similar images** |  |  |
| **Semantic prompt distance** | 6.670418 | 5.823758 |
| **Different images** |  |  |
| **Semantic prompt distance** | 111.639206 | 112.73506 |

Fig. 11: The distance of semantic prompts between similar images and different images, respectively.

| | | | | |
|---|---|---|---|---|
| **Images** |  |  |  |  |
| **VP** | a man surfing in the ocean on a sunny day . | a man in a suit and tie sitting on a chair . | a woman sitting at a table in a restaurant . | a man sitting at a table with a laptop . |
| **VP + SP** | a man skiing down a snow covered slope . | a man sitting on a bench holding a cell phone . | a woman sitting at a table eating food | a man sitting at a table with a bottle of wine . |
| **PL-UIC(1)** | a man in a baseball uniform holding a bat . | a man holding a cell phone in his hand . | a woman standing in a kitchen preparing food . | a man sitting at a table in a kitchen . |
| **PL-UIC** | a man in a baseball uniform swinging a bat . | a man holding a cell phone in his hand . | a woman standing in a kitchen preparing food . | a man sitting at a table in a kitchen . |
| **Ground-truth** | A baseball player taking a swing at a ball | A man stands and talks on his cell phone. | A man preparing food on a large old oven. | A man sits in a wooden kitchen at a table. |

Fig. 12: Examples of the generated captions on the COCO dataset with ground-truth captions.

the first image as an example, the concept "man" is described accurately in the VP method, but "surfing in the ocean on a sunny day" is imprecise. As for the VP + SP method, the key concept "slop" has been correctly recognized. Besides the formerly mentioned concept "man", the "baseball uniform" and "bat" are promisingly described in PL-UIC(1) and PL-UIC, respectively. Especially, the action "swinging" is captured accurately by the method PL-UIC. These qualitative samples strongly verify the advantages of the proposed semantic prompt,

the metric prompt, and the extraordinary ability of the iteration strategy in caption generator refining.

## V. CONCLUSION

We have presented a novel Prompt-based Learning scheme for Unpaired Image Captioning (PL-UIC). By introducing the vision-language pre-trained model (*i.e.*, CLIP), the proposed PL-UIC has for the first time leveraged vision-language prior knowledge in the unpaired image captioning task. To take

the advantage of the vision-language alignment in the pre-trained model, two types of prompts are designed, *i.e.*, semantic prompt and metric prompt. The semantic prompt was devised to guide the caption generation for unpaired image captioning in an unpaired supervision fashion. The boosted experimental performance demonstrated that the vision-aware language prior knowledge is effective for generating captions in the UIC task. To further explore the vision-language prior knowledge, the metric prompt was designed to filter pseudo image-caption pairs for the UIC generator refinement in a paired supervision fashion, so the performance of PL-UIC was greatly enhanced. Extensive experiments demonstrated that the guidance of vision-language prior knowledge is extremely helpful to the UIC task. It also indicates that it is worthy of focusing on the research of paired image-text transformation for UIC. Overall, we hope that our work will shed light on the development of more effective UIC models.

## REFERENCES

[1] J. Wu, T. Chen, H. Wu, Z. Yang, G. Luo, L. Lin, Fine-grained image captioning with global-local discriminative objective, IEEE Transactions on Multimedia (2020).

[2] J. Zhang, K. Mei, Y. Zheng, J. Fan, Integrating part of speech guidance for image captioning, IEEE Transactions on Multimedia 23 (2020) 92–104.

[3] L. Yang, H. Wang, P. Tang, Q. Li, Captionnet: A tailor-made recurrent neural network for generating image descriptions, IEEE Transactions on Multimedia 23 (2020) 835–845.

[4] O. Vinyals, A. Toshev, S. Bengio, D. Erhan, Show and tell: A neural image caption generator, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3156–3164.

[5] O. Vinyals, A. Toshev, S. Bengio, D. Erhan, Show and tell: Lessons learned from the 2015 mscoco image captioning challenge, IEEE transactions on pattern analysis and machine intelligence 39 (4) (2016) 652–663.

[6] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, D. Batra, Visual dialog, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 326–335.

[7] Y. Hong, Q. Wu, Y. Qi, C. Rodriguez-Opazo, S. Gould, Vln bert: A recurrent vision-and-language bert for navigation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 1643–1653.

[8] S. Wu, J. Wieland, O. Farivar, J. Schiller, Automatic alt-text: Computer-generated image descriptions for blind users on a social network service, in: Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing, 2017, pp. 1180–1192.

[9] D. Gurari, Q. Li, A. J. Stangl, A. Guo, C. Lin, K. Grauman, J. Luo, J. P. Bigham, Vizwiz grand challenge: Answering visual questions from blind people, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 3608–3617.

[10] D. Gurari, Y. Zhao, M. Zhang, N. Bhattacharya, Captioning images taken by people who are blind, in: European Conference on Computer Vision, Springer, 2020, pp. 417–434.

[11] J. Kim, A. Rohrbach, T. Darrell, J. Canny, Z. Akata, Textual explanations for self-driving vehicles, in: Proceedings of the European conference on computer vision, 2018, pp. 563–578.

[12] X. Wang, W. Xiong, H. Wang, W. Y. Wang, Look before you leap: Bridging model-free and model-based reinforcement learning for planned-ahead vision-and-language navigation, in: Proceedings of the European Conference on Computer Vision, 2018, pp. 37–53.

[13] D. Omeiza, H. Webb, M. Jirotka, L. Kunze, Explanations in autonomous driving: A survey, IEEE Transactions on Intelligent Transportation Systems (2021).

[14] H. Ben, Y. Pan, Y. Li, T. Yao, R. Hong, M. Wang, T. Mei, Unpaired image captioning with semantic-constrained self-learning, IEEE Transactions on Multimedia (2021).

[15] R. Del Chiaro, B. Twardowski, A. Bagdanov, J. Van de Weijer, Ratt: Recurrent attention to transient tasks for continual image captioning, Advances in Neural Information Processing Systems 33 (2020) 16736–16748.

[16] M. Yang, W. Zhao, W. Xu, Y. Feng, Z. Zhao, X. Chen, K. Lei, Multitask learning for cross-domain image captioning, IEEE Transactions on Multimedia 21 (4) (2018) 1047–1061.

[17] Z. Zhang, Q. Wu, Y. Wang, F. Chen, High-quality image captioning with fine-grained and semantic-guided visual attention, IEEE Transactions on Multimedia 21 (7) (2018) 1681–1693.

[18] X. Li, S. Jiang, Know more say less: Image captioning based on scene graphs, IEEE Transactions on Multimedia 21 (8) (2019) 2117–2130.

[19] T.-H. Chen, Y.-H. Liao, C.-Y. Chuang, W.-T. Hsu, J. Fu, M. Sun, Show, adapt and tell: Adversarial training of cross-domain image captioner, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 521–530.

[20] J. Donahue, P. Krähenbühl, T. Darrell, Adversarial feature learning, arXiv preprint arXiv:1605.09782 (2016).

[21] Y. Feng, L. Ma, W. Liu, J. Luo, Unsupervised image captioning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 4125–4134.

[22] I. Laina, C. Rupprecht, N. Navab, Towards unsupervised image captioning with shared multimodal embeddings, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 7414–7424.

[23] 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: International Conference on Machine Learning, PMLR, 2021, pp. 8748–8763.

[24] C. Wang, M. Chai, M. He, D. Chen, J. Liao, Clip-nerf: Text-and-image driven manipulation of neural radiance fields, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 3835–3844.

[25] R. Zhang, Z. Guo, W. Zhang, K. Li, X. Miao, B. Cui, Y. Qiao, P. Gao, H. Li, Pointclip: Point cloud understanding by clip, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 8552–8562.

[26] M. Narasimhan, A. Rohrbach, T. Darrell, Clip-it! language-guided video summarization, Advances in Neural Information Processing Systems 34 (2021) 13988–14000.

[27] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, G. Neubig, Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing, arXiv preprint arXiv:2107.13586 (2021).

[28] Y. Gu, X. Han, Z. Liu, M. Huang, Ppt: Pre-trained prompt tuning for few-shot learning, arXiv preprint arXiv:2109.04332 (2021).

[29] K. Zhou, J. Yang, C. C. Loy, Z. Liu, Learning to prompt for vision-language models, International Journal of Computer Vision 130 (9) (2022) 2337–2348.

[30] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, S.-N. Lim, Visual prompt tuning, arXiv preprint arXiv:2203.12119 (2022).

[31] L. Wu, M. Xu, J. Wang, S. Perry, Recall what you see continually using gridlstm in image captioning, IEEE Transactions on Multimedia 22 (3) (2019) 808–818.

[32] N. Xu, H. Zhang, A.-A. Liu, W. Nie, Y. Su, J. Nie, Y. Zhang, Multi-level policy and reward-based deep reinforcement learning framework for image captioning, IEEE Transactions on Multimedia 22 (5) (2019) 1372–1383.

[33] M. Z. Hossain, F. Sohel, M. F. Shiratuddin, H. Laga, A comprehensive survey of deep learning for image captioning, ACM Computing Surveys (CsUR) 51 (6) (2019) 1–36.

[34] L. Guo, J. Liu, S. Lu, H. Lu, Show, tell, and polish: Ruminant decoding for image captioning, IEEE Transactions on Multimedia 22 (8) (2019) 2149–2162.

[35] T. Yao, Y. Pan, Y. Li, T. Mei, Incorporating copying mechanism in image captioning for learning novel objects, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 6580–6588.

[36] L. A. Hendricks, S. Venugopalan, M. Rohrbach, R. Mooney, K. Saenko, T. Darrell, Deep compositional captioning: Describing novel object categories without paired training data, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 1–10.

[37] W. Chen, A. Lucchi, T. Hofmann, A semi-supervised framework for image captioning, arXiv preprint arXiv:1611.05321 (2016).

[38] D.-J. Kim, J. Choi, T.-H. Oh, I. S. Kweon, Image captioning with very scarce supervised data: Adversarial semi-supervised learning approach, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), 2019, pp. 2012–2023.

[39] D. Guo, Y. Wang, P. Song, M. Wang, Recurrent relational memory network for unsupervised image captioning, in: International Joint Conferences on Artificial Intelligence, 2021, pp. 920–926.

[40] J. Gu, S. Joty, J. Cai, H. Zhao, X. Yang, G. Wang, Unpaired image captioning via scene graph alignments, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 10323–10332.

[41] F. Liu, M. Gao, T. Zhang, Y. Zou, Exploring semantic relationships for image captioning without parallel data, in: IEEE International Conference on Data Mining, IEEE, 2019, pp. 439–448.

[42] S. Cao, G. An, Z. Zheng, Q. Ruan, Interactions guided generative adversarial network for unsupervised image captioning, Neurocomputing 417 (2020) 419–431.

[43] P. Zhu, X. Wang, Y. Luo, Z. Sun, W.-S. Zheng, Y. Wang, C. Chen, Unpaired image captioning by image-level weakly-supervised visual concept recognition, IEEE Transactions on Multimedia (2022).

[44] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. Dy, T. Pfister, Learning to prompt for continual learning, arXiv preprint arXiv:2112.08654 (2021).

[45] S. Hu, N. Ding, H. Wang, Z. Liu, J. Li, M. Sun, Knowledgeable prompt-tuning: Incorporating knowledge into prompt verbalizer for text classification, arXiv preprint arXiv:2108.02035 (2021).

[46] O. Rubin, J. Herzig, J. Berant, Learning to retrieve prompts for in-context learning, arXiv preprint arXiv:2112.08633 (2021).

[47] N. Ding, S. Hu, W. Zhao, Y. Chen, Z. Liu, H.-T. Zheng, M. Sun, Openprompt: An open-source framework for prompt-learning, arXiv preprint arXiv:2111.01998 (2021).

[48] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models are few-shot learners, Advances in neural information processing systems 33 (2020) 1877–1901.

[49] T. Schick, H. Schütze, Exploiting cloze questions for few shot text classification and natural language inference, arXiv preprint arXiv:2001.07676 (2020).

[50] K. Zhou, J. Yang, C. C. Loy, Z. Liu, Learning to prompt for vision-language models, arXiv preprint arXiv:2109.01134 (2021).

[51] Y. Yao, A. Zhang, Z. Zhang, Z. Liu, T.-S. Chua, M. Sun, Cpt: Colorful prompt tuning for pre-trained vision-language models, arXiv preprint arXiv:2109.11797 (2021).

[52] W. Jin, Y. Cheng, Y. Shen, W. Chen, X. Ren, A good prompt is worth millions of parameters: Low-resource prompt-based learning for vision-language models, in: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics, 2022, pp. 2763–2775.

[53] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[54] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, arXiv preprint arXiv:2010.11929 (2020).

[55] N. Parmar, A. Vaswani, J. Uszkoreit, L. Kaiser, N. Shazeer, A. Ku, D. Tran, Image transformer, in: International Conference on Machine Learning, PMLR, 2018, pp. 4055–4064.

[56] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for contrastive learning of visual representations, in: International conference on machine learning, PMLR, 2020, pp. 1597–1607.

[57] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: European conference on computer vision, Springer, 2014, pp. 740–755.

[58] J. Lu, C. Xiong, D. Parikh, R. Socher, Knowing when to look: Adaptive attention via a visual sentinel for image captioning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 375–383.

[59] P. Sharma, N. Ding, S. Goodman, R. Soricut, Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2018, pp. 2556–2565.

[60] K. Andrej, F.-F. Li, Deep visual-semantic alignments for generating image descriptions, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3128–3137.

[61] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, Bleu: A method for automatic evaluation of machine translation, in: Proceedings of the 40th annual meeting of the Association for Computational Linguistics, 2002, pp. 311–318.

[62] M. Denkowski, A. Lavie, Meteor universal: Language specific translation evaluation for any target language, in: Proceedings of the ninth workshop on statistical machine translation, 2014, pp. 376–380.

[63] C.-Y. LIN, Looking for a good metrics: Rouge and its evaluation, in: Proc. of the 4th NTCIR Workshops (open submission), 2004, 2004, pp. 1–8.

[64] R. Vedantam, C. Lawrence Zitnick, D. Parikh, Cider: Consensus-based image description evaluation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 4566–4575.

[65] P. Anderson, B. Fernando, M. Johnson, S. Gould, Spice: Semantic propositional image caption evaluation, in: European conference on computer vision, Springer, 2016, pp. 382–398.

[66] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, IEEE transactions on pattern analysis and machine intelligence 39 (6) (2016) 1137–1149.