

# Text Data-Centric Image Captioning with Interactive Prompts

Yiyu Wang, Hao Luo, Jungang Xu\*, Yingfei Sun, Fan Wang

**Abstract**—Supervised image captioning approaches have made great progress, but it is challenging to collect high-quality human-annotated image-text data. Recently, large-scale vision and language models (e.g., CLIP) and large-scale generative language models (e.g., GPT-2) have shown strong performances in various tasks, which also provide some new solutions for image captioning with web paired data, unpaired data or even text-only data. Among them, the mainstream solution is to project image embeddings into the text embedding space with the assistance of consistent representations between image-text pairs from the CLIP model. However, the current methods still face several challenges in adapting to the diversity of data configurations in a unified solution, accurately estimating image-text embedding bias, and correcting unsatisfactory prediction results in the inference stage. This paper proposes a new Text data-centric approach with Interactive Prompts for image Captioning, named TIPCap. 1) We consider four different settings which gradually reduce the dependence on paired data. 2) We construct a mapping module driven by multivariate Gaussian distribution to mitigate the modality gap, which is applicable to the above four different settings. 3) We propose a prompt interaction module that can incorporate optional prompt information before generating captions. Extensive experiments show that our TIPCap outperforms other weakly or unsupervised image captioning methods and achieves a new state-of-the-art performance on two widely used datasets, i.e., MS-COCO and Flickr30K.

**Index Terms**—image captioning, weakly or unsupervised approaches, modality gap, interactive prompt.

## I. INTRODUCTION

Image captioning is a typical vision-language task, which aims to automatically generate textual descriptions for given images. In the past few years, although image captioning has made great progress, the models [1]–[11] are generally trained on human-annotated image-text data like MS-COCO [12] and Flickr30K [13], which are challenging to collect. Some unsupervised works [14], [15] try to mitigate this issue using unpaired image-text data, but still need complex pseudo-training or adversarial training to ensure the semantic alignment between decoder and image. Recently, the large foundation models such as BERT [16], GPT-2 [17], T5 [18], CLIP [19], ALIGN [20] and BLIP [21], etc, have provide some new solutions for the vision-language tasks, including image captioning. Based on their strong generalization ability and image-text representation alignment, some methods [22]–[27] use low-cost paired image-text data collected from web (abbreviated as *web data* in this paper) or even text-only data to train models which can even exhibit zero-shot caption capability.

Specifically, there are two main research lines showing promising potential in image captioning. First, some foundation models (e.g. BLIP [21], BLIP2 [28], OFA [29], and

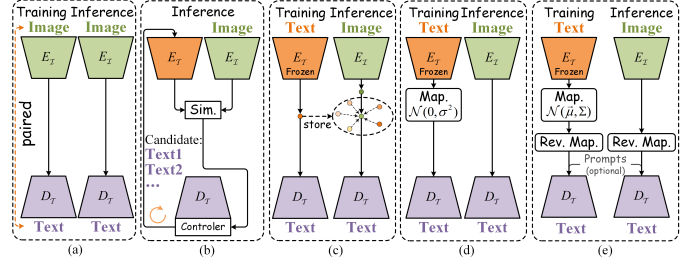


Fig. 1. Comparison of different methods,  $E_{\mathcal{I}}$ ,  $E_{\mathcal{T}}$  and  $D_{\mathcal{T}}$  indicate image encoder, text encoder and text decoder respectively. (a) supervised method. (b) ZeroCap and MAGIC. (c) DeCap. (d) CapDec and CLOSE. (e) our approach.

SEEM [30], etc.) trained on large-scale web data successfully unify multiple vision-language understanding and generation tasks, allowing them to directly predict captions for images. However, due to the low-quality labels and significant noise in web data, these models still need adding high-quality image-text datasets such as MS-COCO in the training/fine-tuning data to achieve comparable performance. Second, another line is to leverage the image-text representation alignment capability of the CLIP model to reduce the cost of acquiring training data. It is generally easier to obtain a textual corpus than to obtain high-quality image-text data. Since CLIP can provide consistent feature representations for image-text pairs, some works train caption models with text-only data or little additional paired data.

This paper focuses on the second research line, because it is a more caption-specific and resource-friendly solution. Additionally, aligning image-text representations rather than directly predict captions can reduce the requirements of the foundation model on the label quality of image-text data. Among the related works shown in the Fig. 1, ZeroCap [23] and MAGIC [24] generate multiple text and then determining the predicted caption according to the feature similarity between image and text calculated by the CLIP model. However, the untrained generation process and the frequent CLIP text encoder forward limit both the model performance and efficiency. DeCap [26] proposes another efficient and text-only required method, which builds a support memory using all text features. In the inference stage, the image embedding can be projected into the text embedding space. However, the support memory restricts its ability to scale up to extensive data. Latest CapDec [25] and CLOSE [27], the closest paradigm to ours, make a assumption that the feature bias between a image-text pair in the CLIP embedding space can be estimated with a Gaussian distribution of  $\mathcal{N}(0, \sigma^2)$ . The value of  $\sigma$  is estimated

from few paired MS-COCO data in CapDec, but it is set as a hyper-parameter in CLOSE. In this way, they need text-only training data because the text embedding can be projected into the image embedding space.

Although the above methods propose some ingenious solutions, we still point out some crucial issues here. 1) Above methods are usually only compatible with one or two specific data configurations. However, in real-world applications, users have very different data configurations. For example, apart from the text corpus, some web data, few high-quality image-text data like MS-COCO or some web image data can also be provided. How to propose a unified solution to deal with different data configurations? 2) The popular assumption that image-text feature bias is an independent Gaussian distribution may be sub-optimal because correlations exist between different feature dimensions. How to propose a better approximation? 3) These methods inevitably output unsatisfactory results. Can we allow the model to handle user-provided prompts (such as the object in the image) to improve predictions?

Based on the above motivations, we propose a new approach TIPCap that is text data-centric with interactive prompts for image captioning, as shown in Fig. 1 (e). Specifically, our TIPCap combines CLIP and GPT-2 to fully leverage the advantages of pre-trained models and contains three extra key modules: a mapping module, a reverse mapping module, and a prompt interaction module. Firstly, we have taken four different data settings into account, which almost cover the vast majority of data configuration scenarios. A unified solution is proposed to estimate text-to-image embedding maps. Secondly, taking into account the correlation between feature dimensions, the mapping module is driven by multivariate Gaussian distribution instead of independent Gaussian distribution, which aims to mitigate the modality gap by performing a simple projection from CLIP text embedding space to CLIP image embedding space. The reverse mapping module performs a weak projection from CLIP image embedding space back to CLIP text embedding space for stronger robustness. During inference, our TIPCap no longer needs mapping module but directly inputs CLIP image embedding into reverse mapping module and follow-up modules to generate captions. Thirdly, the prompt interaction module endows TIPCap with the ability to fuse additional prompt information to generate higher-quality descriptions. With these modules, TIPCap can be trained on text-centric data, and predict captions which can be further improved with manual prompts for a given image.

To evaluate our approach, we conduct extensive experiments on two commonly used datasets: MS-COCO [12] and Flickr30K [13]. The results demonstrate that our approach significantly outperforms existing weakly or unsupervised approaches, and achieves a new state-of-the-art performance.

Our major contributions can be summarized as follows:

- (1) We propose a new approach TIPCap for image captioning, which provides a unified solution for four settings with different data configurations;
- (2) The mapping module utilizes multivariate Gaussian distribution to mitigate the modality gap effectively and outperforms independent Gaussian distribution; our model is

able to handle prompt information, which further enhances the flexibility;

- (3) Extensive experiments demonstrate the effectiveness of TIPCap and achieve a new state-of-the-art performance.

## II. RELATED WORK

### A. Image Captioning

1) *Supervised Approaches*: Inspired by the development of deep learning methods in machine translation [31], [32], most of existing models utilize encoder-decoder framework. Earlier works [1], [2], [33], [34] adopt CNN to extract image features and decode them into sentence by LSTM [35]. Xu *et al.* [2] introduce an attention mechanism which can dynamically focus on salient regions of the given image. After that, Anderson *et al.* [3] propose to use Faster R-CNN [36] as encoder and achieve significant improvement. Some subsequent works [4], [5], [8], [37], [38] follow this paradigm. Recently, transformer-based models have demonstrated excellent performance in image captioning task [5], [7], [39]–[41]. Although supervised methods have achieved impressive results, high-quality human-annotated paired image-text data is essential.

2) *Zero-shot Approaches*: Zero-shot image captioning aims to generate description without human-annotated data. While ZeroCap [23] and MAGIC [24] realize zero-shot image captioning by combining CLIP and GPT-2, and both of them introduce weak visual control cues through the cosine similarity between generated text and given image. Specifically, ZeroCap relies on gradient descent to update the context cache of GPT-2 to make the output match given images; and MAGIC proposes a new decoding strategy to regularize the generated word to be close to given image and previously generated context. However, frequent CLIP text encoder forward slows the inference speed significantly. Wang *et al.* [42] argue that the above methods are prone to fall into harmful contextual language prior and ignore the visual information of given image. DeCap [26] contains a frozen CLIP and a lightweight text decoder. During training, DeCap takes the CLIP text embedding as input to reconstruct its textual sequence, and stores all training CLIP text embeddings as a support memory. During inference, they project the image embedding into CLIP text embedding space by calculating a weighted sum of all embeddings in support memory based on the cosine similarity. CapDec [25] and CLOSE [27] also perform zero-shot image captioning using text-only data and have a similar paradigm, estimating the modality gap between image and text by an independent Gaussian distribution  $\mathcal{N}(0, \sigma^2)$ .

### B. Vision-Language Models

Inspired by BERT [16] and its task-agnostic pre-training paradigm, a line of works [43]–[46] have extended it to Vision-Language (VL) for learning joint representations of image content and natural language, these models directly rely on pre-trained object detector to extract image region features and employ a multimodal encoder to fuse multi-modal features by solving tasks such as masked language modeling (MLM) and image-text matching (ITM). Another line of works (*e.g.* CLIP [19] and ALIGN [20]) construct unimodal encoders

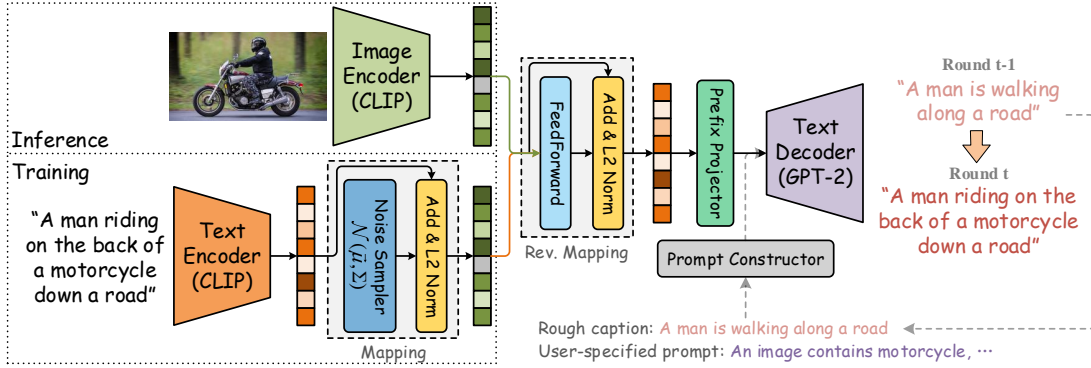


Fig. 2. The overall framework of our approach. Our approach TIPCap is based on a pre-trained CLIP model and a pre-trained GPT-2 model. During training, we first exploit CLIP to extract CLIP text embedding and project it into CLIP image embedding space by a mapping module; then we reconstruct text embedding by a reverse mapping module and inject optional prompt information; finally, GPT-2 generates description. In the inference stage, we no longer need the mapping module but directly feed CLIP image embedding into reverse mapping module and follow-up modules to generate captions.

for image and text separately without densely cross-domain connection, and perform pre-training from scratch on web-scale noisy dataset using only a contrastive loss. ALBEF [47] believes that the learning of image-text interactions and alignment are both important, so introduces a contrastive loss to align image and text before multimodal encoder. The above models are all encoder-based and are not easy to directly transfer to text generation tasks. Considering this issue, BLIP [21] proposes a unified model for VL understanding and generation. Specifically, BLIP constructs unimodal encoders and text decoder for different VL tasks and applies parameter sharing between text encoder and decoder except for the self-attention layers.

### III. METHOD

The overall framework of our approach called TIPCap is shown in Fig. 2. We first introduce model architecture, and four settings with different data configurations, then the details of interactive prompts, and our training objectives.

#### A. Model Architecture

As shown in Fig. 2, we utilize the frozen CLIP [19] and GPT-2 [17] following existing methods. In addition, our TIPCap contains a mapping module, a reverse mapping module, and a prefix projector.

Given a text sequence  $T$ , we first extract the CLIP text embedding  $T_e$  and project it into CLIP image embedding space by a mapping module:

$$T_e = \text{CLIP}_{\text{text}}(T) \quad (1)$$

$$T_e^I = \text{Mapping}(T_e, \mathcal{N}(\bar{\mu}, \Sigma)) \quad (2)$$

where  $\mathcal{N}(\bar{\mu}, \Sigma)$  indicates a multivariate Gaussian distribution with mean  $\bar{\mu}$  and covariance  $\Sigma$ . Briefly, the mapping module performs a simple projection by adding a noise  $\epsilon$  obeying  $\mathcal{N}(\bar{\mu}, \Sigma)$  into  $T_e$ .

Then, a reverse mapping module re-projects  $T_e^I$  back into CLIP text embedding space:

$$T_e^T = \text{Reverse-Mapping}(T_e^I) \quad (3)$$

TABLE I  
TAKING THE EXPERIMENTS ON MS-COCO AS EXAMPLE, A COMPARISON OF DIFFERENT DATA SETTINGS.

	base data	external data
Setting 1	COCO text	1% COCO paired data.
Setting 2	COCO text	100K YFCC15M paired data.
Setting 3	COCO text	5K YFCC15M image data.
Setting 4	COCO text	None

The prefix projector module projects CLIP embedding from CLIP dimension to GPT-2 dimension. Finally,  $P_e$  with optional prompt embedding  $P_{t_e}$  are incorporated as input of GPT-2 to generate captions:

$$P_e = \text{Prefix-Projector}(T_e^T) \quad (4)$$

$$T^I = \text{GPT-2}(P_e, P_{t_e}) \quad (5)$$

The inference forward is slightly different from training forward described above: We no longer need the mapping module but directly extract CLIP image embedding and feed it into reverse mapping module and follow-up modules.

#### B. Estimation of $\bar{\mu}$ and $\Sigma$

The mapping module is driven by a multivariate Gaussian distribution  $\mathcal{N}(\bar{\mu}, \Sigma)$  to mitigate the modality gap, which brings up a crucial issue: how to obtain the mean  $\bar{\mu}$  and covariance  $\Sigma$  that can effectively characterize the true modality bias  $\delta$  between CLIP image embedding and text embedding?

In this paper, we consider four settings with different data configurations, which can cover the majority of real-world scenarios (text corpus  $T_{\text{corpus}}$  is available in each setting):

- (1) Setting 1: Few human-annotated high-quality paired data  $\langle I_{\text{human}}, T_{\text{human}} \rangle$  is available, where  $T_{\text{human}}$  is **homologous** to  $T_{\text{corpus}}$ ;
- (2) Setting 2: Low-quality paired web data  $\langle I_{\text{web}}, T_{\text{web}} \rangle$  is available, where  $T_{\text{web}}$  is **heterologous** to  $T_{\text{corpus}}$ ;
- (3) Setting 3: No paired data, but a few source-agnostic image data  $I_{\text{any}}$  is available;
- (4) Setting 4: No any extra data, only  $T_{\text{corpus}}$  is available.

**Setting 1.** We first calculate the embedding difference of human annotated paired data:

$$\delta \simeq I_{e,human} - T_{e,human} \quad (6)$$

where  $I_{e,human}$  and  $T_{e,human}$  indicate CLIP image embedding and CLIP text embedding respectively.

As mentioned before,  $\mathcal{N}(\vec{\mu}, \Sigma)$  aims to characterize the modality bias. Here  $T_{human}$  is paired with  $I_{human}$ , and also homologous to  $T_{corpus}$ , so we can directly adopt the mean  $\vec{\mu}_\delta$  and covariance  $\Sigma_\delta$  of  $\delta$  as a tight estimation of  $\vec{\mu}$  and  $\Sigma$ .

**Setting 2.** Similar to setting 1, we can calculate the embedding difference of web data, and use the mean and covariance of embedding difference as an estimation. But it performs worse, due to  $T_{web}$  and  $T_{corpus}$  are heterologous. We propose to apply a simple correction to alleviate this issue:

$$\begin{aligned} \delta &\simeq I_{e,web} - \text{Correct}(T_{e,web}) \\ &= I_{e,web} - (T_{e,web} + \delta_{w \rightarrow c}) \\ &= \delta_{web} - \delta_{w \rightarrow c} \end{aligned} \quad (7)$$

where  $\delta_{web} \sim \mathcal{N}(\vec{\mu}_{web}, \Sigma_{web})$  and  $\delta_{w \rightarrow c} \sim \mathcal{N}(\vec{\mu}_{w \rightarrow c}, \Sigma_{w \rightarrow c})$ , subscript  $w \rightarrow c$  indicates ‘‘web data to corpus data’’, which aims to achieve a domain alignment from  $T_{web}$  to  $T_{corpus}$ . Since there is no pairwise relationship between  $T_{corpus}$  and  $T_{web}$ ,  $\text{Correct}(\cdot)$  is just a global and rough estimation of the domain alignment.

**Setting 3.** In this setting, we extend the mapping module to be trainable instead of pre-defined parameters. Specifically, the covariance can be denoted as  $\Sigma = LL^\top$  by cholesky decomposition, where  $L$  is a lower triangular matrix with the same size as  $\Sigma$ . Through reparameterization trick [48], the noise of  $\epsilon \sim \mathcal{N}(\vec{\mu}, \Sigma)$  can be re-formulated as follows:

$$\epsilon = L\epsilon' + \vec{\mu}, \quad \epsilon' \sim \mathcal{N}(\vec{0}, I) \quad (8)$$

Therefore, we can carry out the mean  $\vec{\mu}$  and matrix  $L$  as trainable parameters.

The difficulty lies in how to drive the training of  $\vec{\mu}$  and  $L$  toward the correct direction. For this purpose, we introduce a few source-agnostic image data  $I_{any}$ , and calculate the embedding difference:

$$\delta \simeq \delta_{any} = I_{e,any} - T_{e,corpus} \quad (9)$$

where  $\delta_{any} \sim \mathcal{N}(\vec{\mu}, \Sigma)$  if  $I_{any}$  is paired with  $T_{corpus}$ . Here we relax this requirement and assume that  $\delta_{any} \sim \mathcal{N}(\vec{\mu}, \Sigma)$  is also roughly satisfied even though image and text are unpaired, to derive our training objective.

Again we apply reparameterization trick, and have:

$$\delta_{any} \sim \mathcal{N}(\vec{\mu}, \Sigma) \simeq L\epsilon + \vec{\mu}, \quad \epsilon \sim \mathcal{N}(\vec{0}, I) \quad (10)$$

then we apply a simple transformation to  $\delta_{any}$ , and result in:

$$\epsilon = L^{-1}(\delta_{any} - \vec{\mu}) \sim \mathcal{N}(\vec{0}, I) \quad (11)$$

From the above, our training goal is to make  $\epsilon$  more close to a standard Gaussian distribution of  $\mathcal{N}(\vec{0}, I)$ :

$$\mathcal{L}_{\text{Map}} = KL(L^{-1}(\delta_{any} - \vec{\mu}) \| \mathcal{N}(\vec{0}, I)) \quad (12)$$

**Setting 4.** Setting 4 explores a more extreme data configuration, where only text data is available. We follow the paradigm

similar to setting 3 and apply  $\mathcal{L}_{\text{Map}}$  to optimize trainable parameters  $\vec{\mu}$  and  $L$ . The calculation of  $\mathcal{L}_{\text{Map}}$  relies on the modality bias between CLIP image embedding and CLIP text embedding. However, if we use the bias between  $T_e^I$  and  $T_e$  to optimize the loss  $\mathcal{L}_{\text{Map}}$ , the model is prone to falling into trivial solutions (collapse), because the modality bias between the output and input of mapping module, *i.e.*  $T_e^I$  and  $T_e$ , is directly sampled from  $\mathcal{N}(\vec{\mu}, \Sigma)$ . To address such issue, we propose several specific designs.

Firstly, we follow other works to apply several asymmetric designs to enhance the robustness and avoid mode collapse. Different from the mapping module consisting of an addition operation between the input and a sampled multivariate Gaussian distribution noise, the reverse mapping module is designed as a FeedForward layer, and aimed at re-project the output of mapping module back into CLIP text embedding space. Then, we do not use the original text embedding  $T_e$  but the reconstructed one  $T_e^T$  to optimize  $\mathcal{L}_{\text{Map}}$ . Then,  $\mathcal{L}_{\text{Map}}$  can be calculated as follows:

$$\mathcal{L}_{\text{Map}} = KL(L^{-1}(\delta_{pseudo} - \vec{\mu}) \| \mathcal{N}(\vec{0}, I)) \quad (13)$$

$$\delta_{pseudo} = T_e^I - T_e^T \quad (14)$$

where  $T_e^I$  and  $T_e^T$  indicate the output of mapping module and reverse mapping module respectively. Since the reverse mapping module does not perform strict reconstruction of CLIP text embedding, which also introduces asymmetry to obtain a more robust performance. In order to ensure the unity of our TIPCap, we also apply the reverse mapping module in other settings except Setting 4.

Secondly, since the lack of real image data to introduce corresponding latent prior information, we employ a relational knowledge distillation loss:

$$\mathcal{L}_{\text{Disti}} = KL(\mathcal{S}_{T_e^I} \| \mathcal{S}_{T_e}) \quad (15)$$

where  $\mathcal{S}_{T_e^I}$  and  $\mathcal{S}_{T_e}$  indicate internal cosine similarity matrix of  $T_e^I$  and  $T_e$  respectively. Employing  $\mathcal{L}_{\text{Disti}}$  is aiming to encourage  $T_e^I$  to have a similar internal cosine similarity to  $T_e$ , which provides a constraint to ensure that  $T_e^I$  and  $T_e$  are semantically related and avoids mode collapse.

### C. Interactive Prompts

Inevitably, image captioning models output unsatisfactory sentences, sometimes with factual errors or missing objects. Based on this issue, we hope to endow our model with the ability to deal with additional prompt information to generate information-enhanced captions.

Inspired by the supervised fine-tuning of InstructGPT [49], we construct prompts as textual sentences and serve as a part of input of GPT-2 model, as shown in Fig. 2. One full prompt contains two parts: a rough caption predicted by the model and a user-specified prompt sentence correcting the caption. The main difficulty focuses on collecting user-specified prompts sentences during training, which should contain information ignored in rough captions but existed in ground truth captions, to introduce positive information for training guidance.

We divide the training into two stages: 1) We perform the first stage of training without introducing prompts and get a

base captioning model  $model_{base}$ , which aims to endow our model with the ability to generate rough captions for the second stage training. 2) Then, we initialize  $model_{prompt}$  with parameters of  $model_{base}$  to perform the second stage of training with introducing prompts. To avoid complex and expensive manual annotation, we explore a simple strategy to generate user-specified prompts. At each training step,  $model_{base}$  is frozen to generate rough captions  $\langle c_r \rangle$ . For user-specified prompt, we do extract nouns or noun phrases set  $\{p_i\}_{i=1}^N$  by part-of-speech tagging from ground truth captions  $\langle c_{gt} \rangle$ . Furthermore, aiming to preserve positive prompt information, we remove nouns or noun phrases that appear in  $\langle c_r \rangle$  from  $\{p_i\}_{i=1}^N$  to get filtered set  $\{p'_i\}_{i=1}^{N'}$ .

Taking Fig. 2 as example, for the round  $t$  generation, we have  $\langle c_r \rangle = \text{"A man is walking along a road."}$ ,  $\{p'_i\}_{i=1}^{N'} = \{\text{"motorcycle"}\}$ , and  $\langle c_{gt} \rangle = \text{"A man riding on the back of a motorcycle down a road."}$  The corresponding full prompt sentence  $P_t$  is constructed by prompt constructor as follows:

**Reference:** *A man is walking along a road.*  
**Prompt:** *An image contains motorcycle.*  
**Prediction:** *A man riding on the back of a motorcycle down a road.*

Note that we hope our TIPCap keeps the ability to perform general captioning task (*i.e.* generate captions without “reference” and “prompt”), thus prompt information is **NOT** always essential. Furthermore, when the generated caption is good and descriptive enough, it is also not necessary to introduce prompts and perform the second inference.

Based on the above considerations, we replace the prompt sentences with padding tokens with a probability of  $p = 0.1$  when performing the stage 2 training (refer to “Implementation details” in Section IV-A); in addition, we also replace the prompts with padding when generated caption is the same as the ground truth caption. Examples of constructed full prompt sentences are shown in Fig. 3.

#### D. Objectives

**Language Modeling Loss.** For a mini-batch of  $n$  texts  $\{T^i\}_{i=1}^n$ , where  $T^i = \{T_1^i, \dots, T_L^i\}$ , we optimize our model by applying Maximum Likelihood Estimation (MLE):

$$\mathcal{L}_{MLE} = -\frac{1}{n \times L} \sum_{i=1}^n \sum_{j=1}^L \log p_{\theta}(T_j^i | T_{1:j-1}^i) \quad (16)$$

where  $\theta$  denotes the parameters of our model.

**Reverse Mapping Reconstruction Loss.** For our reverse mapping module, we define  $\mathcal{L}_{Recons}$  to constrain the reconstruction relationship of  $T_e^T$  to  $T_e$ :

$$\mathcal{L}_{Recons} = \mathcal{L}_{Cosine} + \mathcal{L}_{CL} \quad (17)$$

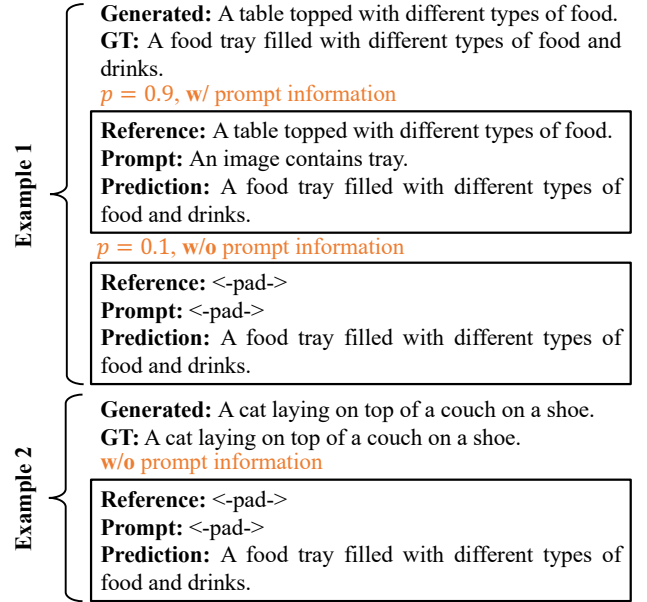


Fig. 3. Examples of constructed full prompt sentences during stage 2 training.

where  $\mathcal{L}_{Cosine}$  and  $\mathcal{L}_{CL}$  denote cosine embedding loss and contrastive loss [19] respectively, and are defined as follows:

$$\mathcal{L}_{Cosine} = \frac{1}{n} \sum_{i=1}^n (1 - \mathcal{S}_{i,i}), \mathcal{L}_{CL} = -\frac{1}{2} (\mathcal{L}_{CL}^{\mathcal{S}} + \mathcal{L}_{CL}^{\mathcal{S}^T}) \quad (18)$$

$$\mathcal{L}_{CL}^{\mathcal{S}} = \frac{1}{n} \sum_{i=1}^n \log \frac{\exp(\mathcal{S}_{i,j}/\tau)}{\sum_{j=1}^n \exp(\mathcal{S}_{i,j}/\tau)} \quad (19)$$

where  $\tau = 0.1$ ,  $\mathcal{S} \in \mathbb{R}^{n \times n}$  indicates the cosine similarity matrix between  $T_e^T$  and  $T_e$ . Intuitively,  $\mathcal{L}_{Cosine}$  is introduced for hard reconstruction, and  $\mathcal{L}_{CL}$  aims to relax the reconstruction constraint. Because we hope that  $T_e^T$  has similar semantics to  $T_e$  without perfect reconstruction.

## IV. EXPERIMENTS

### A. Experimental Setting

1) **Datasets:** For fair performance comparison, we conduct experiments on both MS-COCO [12] and Flickr30K [13] datasets, and follow “Karpathy” split [50]. Furthermore, YFCC15M [51] which has been used to train CLIP is adopt as low-quality paired web data for setting 2 and 3.

**MS-COCO** is a widely used dataset in image captioning task, which consists of 123,287 images, and each is paired with 5 human-annotated captions. **Flickr30K** is another human-annotated dataset similar to MS-COCO, where each image is also paired with 5 reference captions, but only contains 31,000 images. **YFCC15M** is a subset of the large-scale web-collected dataset YFCC100M, where each image is annotated with a weakly paired alt-text.

Taking experiments on MS-COCO as examples, TABLE I shows a comparison of the available data under four different settings, which cover the majority of real-world scenarios. Specifically, we sample 1% paired MS-COCO data for setting 1, about 100K paired YFCC data for setting 2, and 5K YFCC image data for setting 3.

TABLE II

IN-DOMAIN CAPTIONING RESULTS ON MS-COCO AND FLICKR30K. THE SUPERSCRIPT \* INDICATES THAT RESULTS ARE FROM MAGIC [24]. FOR WEAKLY OR UNSUPERVISED APPROACHES, THEY USE CLIP WITH DIFFERENT BACKBONE AS ENCODER. DECAP [26] CONSTRUCTS A LIGHTWEIGHT TRANSFORMER DECODER (T.D.) INSTEAD OF PRE-TRAINED GPT-2. CLOSE [27] USES T5 MODEL [18]. SETTING 1-4 ARE ABBREVIATED AS S1-4.

Method	Encoder	Decoder	MS-COCO						Flickr30k						
			B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S	
<b>Fully Supervised:</b>															
BUTD			77.2	36.2	27.0	56.4	113.5	20.3	-	27.3	21.7	-	56.6	16.0	
UniVLP			-	36.5	28.4	-	116.9	21.2	-	30.1	23.0	-	67.4	17.0	
ClipCap			-	33.5	27.5	-	113.1	21.1	-	-	-	-	-	-	
Oscar			-	36.5	30.3	-	123.7	23.1	-	-	-	-	-	-	
<b>Weakly or Unsupervised:</b>															
ZeroCap*	ViT-B/32	GPT-2	49.8	7.0	15.4	31.8	34.5	9.2	44.7	5.4	11.8	27.3	16.8	6.2	
MAGIC	ViT-B/32	GPT-2	56.8	12.9	17.4	39.9	49.3	11.3	44.5	6.4	13.1	31.6	20.4	7.1	
DeCap	ViT-B/32	T.D.	-	24.7	25.0	-	91.2	18.7	-	21.2	21.8	-	56.7	15.2	
CapDec	RN50x4	GPT-2	69.2	26.4	25.1	51.8	91.8	-	55.5	17.7	20.0	43.9	39.1	-	
CLOSE	ViT-L/14	T5	-	22.1	23.7	-	81.2	17.7	-	-	-	-	-	-	
CLOSE w/Tuned Noise	ViT-L/14	T5	-	28.6	25.2	-	95.4	18.1	-	-	-	-	-	-	
TIPCap (S1)	RN50x4	GPT-2	74.6	30.7	26.7	54.2	106.7	20.3	71.1	25.6	22.5	49.1	63.7	16.2	
TIPCap (S2)	RN50x4	GPT-2	72.7	28.6	25.6	52.5	100.6	19.6	68.0	23.7	21.3	47.3	57.8	15.2	
TIPCap (S3)	RN50x4	GPT-2	73.0	30.4	26.5	53.8	104.5	20.0	68.4	24.2	21.9	48.1	61.2	16.1	
TIPCap (S4)	RN50x4	GPT-2	71.3	29.8	26.2	53.4	102.1	19.4	67.5	24.0	21.7	47.7	59.4	15.9	
TIPCap (S4)	ViT-L/14	GPT-2	73.3	31.4	26.9	54.2	106.6	20.2	69.6	26.1	23.0	49.3	65.7	17.0	

TABLE III

CROSS-DOMAIN CAPTIONING RESULTS.  $X \implies Y$  MEANS THAT THE MODEL IS TRAINED ON DATASET X BUT EVALUATED ON DATASET Y.

Method	Flickr30k $\implies$ MS-COCO						MS-COCO $\implies$ Flickr30k					
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
MAGIC	41.4	5.2	12.5	30.7	18.3	5.7	46.4	6.2	12.2	31.3	17.5	5.9
CapDec	43.3	9.2	16.3	36.7	27.3	-	60.2	17.3	18.6	42.7	35.7	-
DeCap	-	12.1	18.0	-	44.4	10.9	-	16.3	17.9	-	35.7	11.1
TIPCap (S1)	59.8	16.7	19.4	42.3	56.0	12.5	66.9	19.8	19.8	45.3	48.2	13.7
TIPCap (S2)	55.9	14.5	17.8	40.4	47.8	11.4	63.5	17.3	18.4	43.2	41.6	12.3
TIPCap (S3)	55.9	14.2	18.4	40.6	48.7	11.9	63.7	18.6	19.2	44.0	42.8	13.0
TIPCap (S4)	55.8	14.3	18.4	40.5	48.5	11.9	63.8	18.7	19.2	44.1	42.4	12.9

2) *Evaluation metrics*: Following the common paradigm, we adopt five widely used metrics for evaluation, including BLEU- $N$  [31], METEOR [52], ROUGE-L [53], CIDEr [54], and SPICE [55], which are denoted as B- $N$ , M, R, C, and S for simplicity.

3) *Implementation details*: Our TIPCap model is implemented in PyTorch [56] and trained on 4 Nvidia Tesla V100 (32GB) GPUs. We first train our model **without** introducing interactive prompts for 10 epochs, where the learning rate is warmed-up to  $5e^{-4}$  during 1250 steps and decayed linearly; then we froze parameters of the mapping module and reverse mapping module and train our model **with** interactive prompts for another 5 epochs with learning rate of  $1e^{-6}$ , in this training stage. For optimization, we set the batch size on each GPU to 32 and adopt AdamW optimizer [57] with a weight decay of 0.1 in both above stages and beam size is set to 5 during inference. Trained models and source code will be released.

## B. Comparison with Existing Models

1) *In-domain captioning*: TABLE II shows the in-domain captioning results on MS-COCO and Flickr30K. We compare our proposed TIPCap with several approaches with different supervision levels. 1) Fully supervised approaches, which rely on human-annotated paired data for model training, including BUTD [3], UniVLP [58], ClipCap [59], Oscar [46]; and 2) weakly or unsupervised approaches, which adopt pre-trained

foundation models (e.g. CLIP [19], GPT-2 [17] and T5 [18]) to perform image captioning on unpaired or text-only data, including ZeroCap [23], MAGIC [24], DeCap [26], CapDec [25] and CLOSE [27]. Our proposed TIPCap belongs to the second category, weakly or unsupervised approaches.

As shown in TABLE II, we report the performance results of four different data settings. As expected, TIPCap (Setting 1) performs better than the other three settings, as it estimates the distribution parameter of mapping module from high-quality paired data, which introduces strong and credible prior information of modality bias and can be regarded as an upper bound of our proposed approach. TIPCap (Setting 2) uses weakly paired web data but performs worse, because the alt-text is low-quality and has a large margin with training text corpus. TIPCap (Setting 3) and TIPCap (Setting 4) have less data available but also achieve superior performances that significantly outperforms other state-of-the-art models, which shows the effectiveness of our trainable mapping module. Especially, our TIPCap uses CLIP with RN50x4 backbone and GPT-2 model but outperforms the strong competitor CLOSE that adopts CLIP with ViT-L/14 backbone and T5 model.

Overall, our proposed TIPCap achieves state-of-the-art performance and shows substantial performance improvement compared with recent weakly or unsupervised methods, demonstrating the effectiveness and advantages.

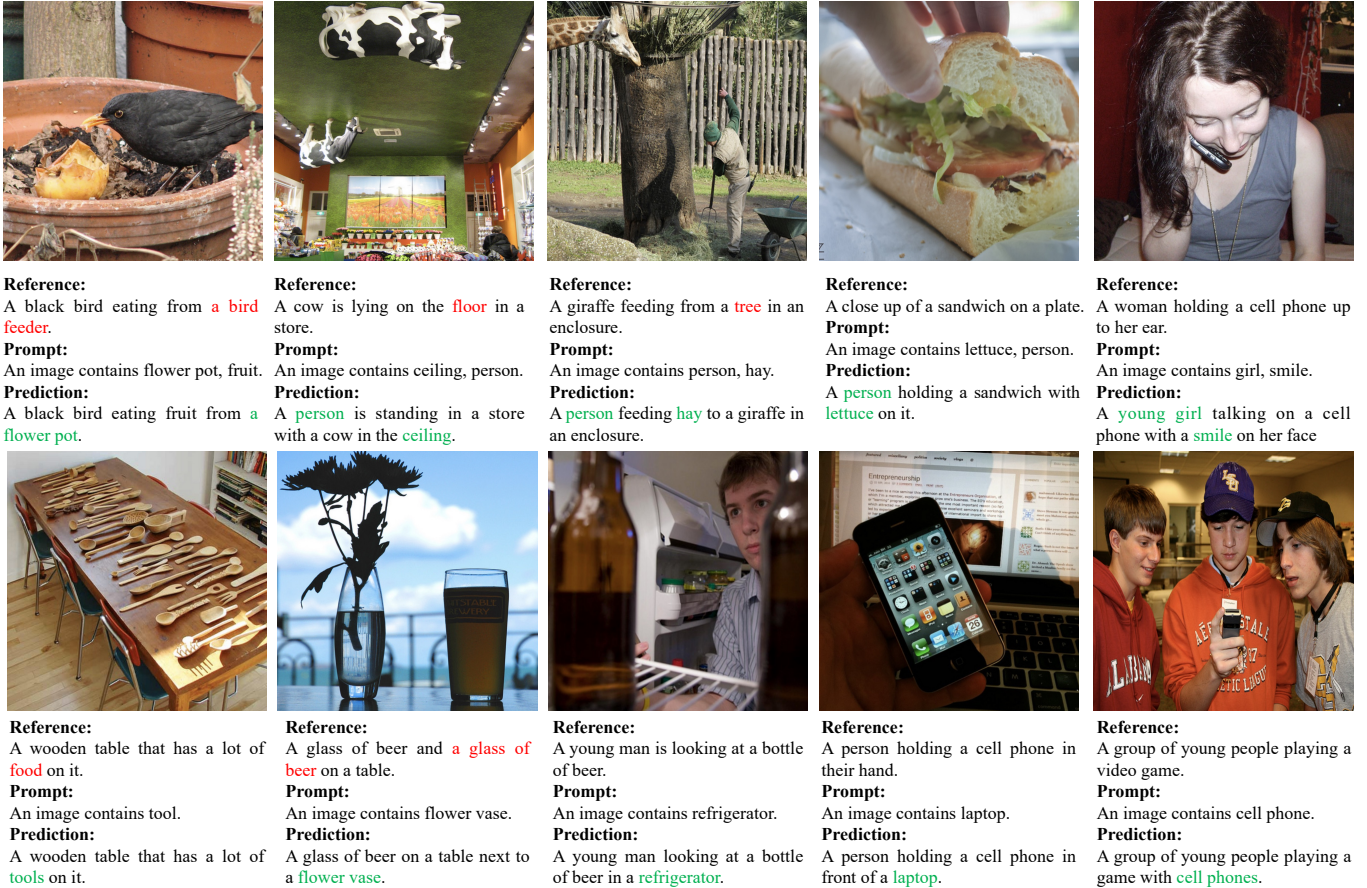


Fig. 4. Examples of captions generated by TIPCap with simulated interactive prompts, images come from MS-COCO karpathy test split. “Reference” indicates the generated caption without prompt information; “Prompt” indicates the simulated user-specified prompt information; “Prediction” shows the new generated caption with prompt information.

TABLE IV  
PERFORMANCE COMPARISON WITH SIMULATED INTERACTIVE PROMPT INFORMATION IN TIPCAP.

		B-1	B-4	M	R	C	S
CLIP RN50x4	S1	77.27 $\pm$ .15	32.73 $\pm$ .06	27.57 $\pm$ .06	55.53 $\pm$ .12	113.17 $\pm$ .35	21.97 $\pm$ .06
	S2	75.93 $\pm$ .06	30.77 $\pm$ .12	26.70 $\pm$ .10	54.03 $\pm$ .06	107.93 $\pm$ .06	21.60 $\pm$ .00
	S3	78.23 $\pm$ .06	33.90 $\pm$ .10	28.50 $\pm$ .00	56.33 $\pm$ .06	116.27 $\pm$ .06	24.27 $\pm$ .06
	S4	77.80 $\pm$ .10	33.60 $\pm$ .17	28.30 $\pm$ .00	56.20 $\pm$ .00	115.33 $\pm$ .06	23.73 $\pm$ .06
CLIP ViT-L/14	S1	77.53 $\pm$ .06	33.03 $\pm$ .12	28.00 $\pm$ .10	55.80 $\pm$ .00	115.87 $\pm$ .23	22.40 $\pm$ .00
	S2	76.30 $\pm$ .10	31.30 $\pm$ .20	27.03 $\pm$ .06	54.23 $\pm$ .06	110.23 $\pm$ .42	22.03 $\pm$ .06
	S3	78.57 $\pm$ .15	34.40 $\pm$ .20	28.40 $\pm$ .00	56.27 $\pm$ .06	116.87 $\pm$ .21	24.10 $\pm$ .00
	S4	78.27 $\pm$ .15	33.90 $\pm$ .17	28.37 $\pm$ .06	56.07 $\pm$ .06	116.57 $\pm$ .42	24.13 $\pm$ .06

2) *Cross-domain captioning*: As shown in TABLE III, we conduct cross-domain experiments to further explore the generalization ability of our approach. Specifically, we train our TIPCap on the source dataset (e.g. MS-COCO), but perform inference on a different target dataset (e.g. Flickr30K). We compare TIPCap with several text-only methods, including MAGIC [24], CapDec [25], and DeCap [26]. Our TIPCap still outperforms all compared methods, demonstrating the superiority of our approach in generalization ability.

### C. Ablation Studies

1) *Impact of Interactive Prompts*: Due to the dynamic nature of the prompt information and the lack of relevant

benchmarks, it is not easy to give accurate quantitative results when introducing interactive prompts.

To verify the effectiveness of interactive prompts, we perform the evaluation on MS-COCO with simulated user-specified prompt information. Specifically, we first apply TIPCap to generate a caption without prompt information as the rough caption, then we do second-time inference to introduce the prompt information, in which the simulated user-specified prompt information is sampled from the nouns or noun key phrases (ignored by rough caption) extracted from corresponding ground truth caption. We perform three times evaluations in each setting and report the average performance with the standard deviation, as shown in TABLE IV. The results demonstrate that the introduction of prompt information

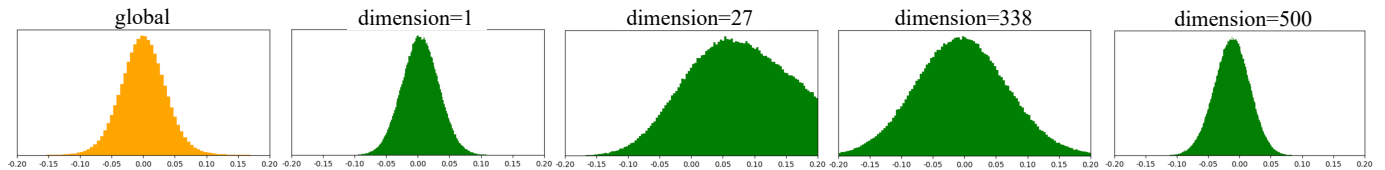


Fig. 5. Histogram visualization of the CLIP image and text embedding difference of MS-COCO training set, where orange and green indicate the histogram statistics on all dimensions (global) and specific dimensions (local) separately.

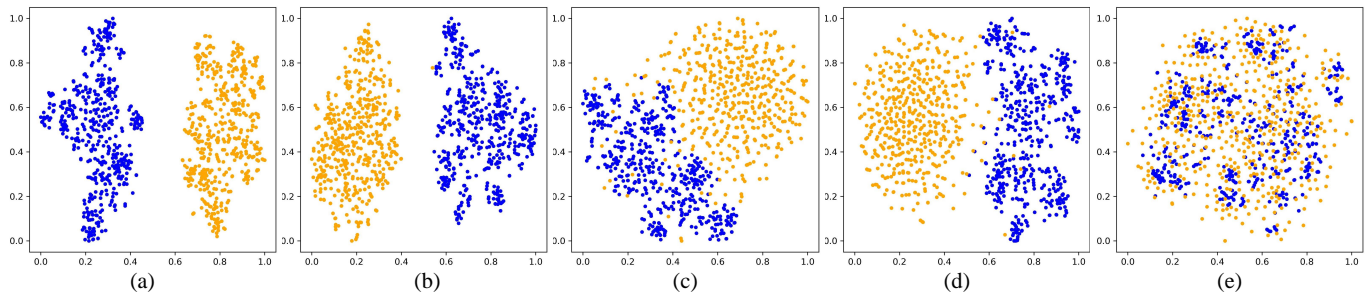


Fig. 6. t-SNE visualization of CLIP image and text embeddings from MS-COCO training set, where blue and yellow points indicate image and text embeddings, respectively. (a) without mapping module; (b) mapping module driven by univariate Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu \simeq 0.0009$  and  $\sigma \simeq 0.0440$  are estimated from MS-COCO paired data; (c) mapping module driven by  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu = 0$  and  $\sigma = \sqrt{0.016}$ , which is adopted in CapDec [25]; (d) mapping module driven by  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu = 0$  and  $\sigma = 0.08$ , which is adopted in CLOSE [27]; (e) mapping module driven by multi-variate Gaussian distribution  $\mathcal{N}(\bar{\mu}, \Sigma)$ , where  $\bar{\mu}$  and  $\Sigma$  are estimated from coco paired data.

TABLE V

ABLATION STUDY OF MAPPING MODULE AND REVERSE MAPPING MODULE. (A) OUR FULL TIPCAP MODEL; (B) TIPCAP WITH THE MAPPING MODULE DRIVEN BY INDEPENDENT GAUSSIAN DISTRIBUTION; (C) TIPCAP MODEL WITH REVERSE MAPPING MODULE REMOVED. “N/A” REFERS TO “NOT APPLICABLE”.

	CLIP RN50x4					CLIP ViT-L/14						
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
(A) TIPCap, $\mathcal{N}(\bar{\mu}, \Sigma)$												
S1	74.6	30.7	26.7	54.2	106.7	20.3	75.4	31.2	27.2	54.5	109.7	20.9
S2	72.7	28.6	25.6	52.5	100.6	19.6	73.4	29.1	26.0	52.7	103.5	20.2
S3	73.0	30.4	26.5	53.8	104.5	20.0	73.7	31.2	26.7	53.9	107.5	20.5
S4	71.3	29.8	26.2	53.4	102.1	19.4	73.3	31.4	26.9	54.2	106.6	20.2
(B) TIPCap, $\mathcal{N}(\mu, \sigma^2)$												
S1	68.5	24.8	24.1	49.7	89.5	18.6	66.7	22.5	23.2	46.8	87.3	18.1
S2	69.4	25.5	24.4	50.2	91.6	18.7	68.0	23.5	23.5	47.7	90.1	18.3
S3	69.9	26.2	24.6	50.7	92.9	18.9	70.3	25.3	24.7	50.2	94.1	19.0
S4	70.0	25.9	24.1	50.5	91.5	18.5	69.9	25.3	24.4	49.9	93.8	18.8
(C) TIPCap, $\mathcal{N}(\bar{\mu}, \Sigma)$ , w/o reverse mapping module												
S1	74.1	30.2	26.2	53.7	105.0	19.8	75.3	31.3	26.9	54.4	109.2	20.8
S2	72.4	28.3	25.1	52.0	99.2	19.1	73.7	29.2	25.9	52.8	103.1	19.8
S3	71.4	30.0	26.3	53.4	102.4	19.4	73.2	30.8	26.6	53.9	105.5	20.1
S4				N/A						N/A		

brings positive influences on performance.

Furthermore, in order to make a more intuitive qualitative comparison, we give some examples of captions generated by TIPCap when introducing simulated interactive prompts, as shown in Fig. 4. In which, “**Reference**” caption is generated by the same TIPCap model without prompt information (*i.e.* rough caption); “**Prompt**” indicates the simulated user-specified prompt information constructed by sampling from the ignored ground truth nouns or phrases; “**Prediction**” gives the new caption generated by our TIPCap.

2) *Impact of Mapping Module*: This paper makes an assumption that the feature bias between image-text paired CLIP embeddings can be estimated as a multivariate Gaussian distribution, which is a natural extension of the assumption adopted in CapDec [25] and CLOSE [27] that use independent Gaussian distribution.

In this part, we conduct experiments to investigate the effect of the above two strategies. The results are reported in TABLE V (A) and (B), we can see that the performances are better when using multivariate Gaussian distribution. The results of setting 1 when using  $\mathcal{N}(\mu, \sigma^2)$  perform worst and intuitively show that the independent Gaussian distribution is a suboptimal strategy, which can not effectively characterize the feature bias. Moreover, the results of setting 3 and setting 4 in TABLE V (B) demonstrate that our TIPCap is also effective when using independent Gaussian distribution and can achieves competitive performances.

Fig. 5 shows simple histogram visualization of the CLIP image and text embedding difference of MS-COCO training set, it can be seen that the distribution of embedding differences in different dimensions are not consistent. And it is why we apply multivariate Gaussian distribution instead of independent Gaussian distribution.

As shown in Fig. 6, we randomly sample 500 paired image-text data from MS-COCO training set and visualize the CLIP image and text embeddings. Fig. 6 (a) shows the clear modality gap between CLIP image embeddings and CLIP text embeddings. Fig. 6 (b), (c) and (d) show the influence of mapping module driven by univariate Gaussian distribution. Fig. 6 (b) indicates that the modality gap still exists after applying the mapping module driven by  $\mathcal{N}(\mu, \sigma^2)$ , even if mean  $\mu$  and standard deviation  $\sigma$  are estimated from MS-COCO paired data. CapDec [25] and CLOSE [27] also use



TABLE VI

PERFORMANCE COMPARISON WITH DIFFERENT PREFIX LENGTH  $L$  AND DIFFERENT LAYER NUMBER OF  $N$  IN THE PREFIX PROJECTOR MODULE. EXPERIMENTS ARE CONDUCTED UNDER SETTING 4.

	CLIP RN50x4						CLIP ViT-L/14						
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S	
L	2	71.1	29.3	26.0	53.2	101.1	19.2	74.0	31.6	26.5	54.2	106.7	20.0
	4	<b>71.3</b>	<b>29.8</b>	<b>26.2</b>	<b>53.4</b>	<b>102.1</b>	<b>19.4</b>	<b>73.3</b>	<b>31.4</b>	<b>26.9</b>	<b>54.2</b>	<b>106.6</b>	<b>20.2</b>
	5	71.1	29.5	26.3	53.3	101.7	19.5	73.3	31.3	26.8	54.2	106.7	20.3
	10	71.0	29.3	26.0	53.1	100.5	19.2	72.0	30.7	27.3	54.4	106.5	20.7
	20	71.1	29.5	26.1	53.1	101.5	19.4	73.4	30.9	26.7	53.9	106.7	20.2
	40	70.1	29.2	26.2	53.0	99.9	19.1	72.3	30.4	27.1	54.3	106.3	20.5
N	1	71.2	29.9	26.1	53.2	101.5	19.1	72.9	31.1	26.5	53.7	105.2	19.9
	2	71.7	30.1	26.3	53.5	101.9	19.4	73.2	31.2	26.7	54.1	106.0	20.1
	3	<b>71.3</b>	<b>29.8</b>	<b>26.2</b>	<b>53.4</b>	<b>102.1</b>	<b>19.4</b>	<b>73.3</b>	<b>31.4</b>	<b>26.9</b>	<b>54.2</b>	<b>106.6</b>	<b>20.2</b>
	4	71.6	29.9	26.2	53.4	102.1	19.3	73.8	31.1	26.7	54.1	107.0	20.5
	5	71.6	30.0	26.2	53.5	102.3	19.5	74.1	31.4	26.8	54.1	108.2	20.6
	6	71.2	29.8	26.3	53.5	102.4	19.6	73.7	31.5	27.0	54.1	108.5	20.8

TABLE VII

PERFORMANCE COMPARISON WITH DIFFERENT RATIO OF PAIRED DATA (*i.e.* MS-COCO) USED IN SETTING 1.

r	CLIP RN50x4						CLIP ViT-L/14					
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
0.2%	73.3	29.4	26.1	53.4	103.0	19.7	74.7	30.3	26.7	53.9	107.5	20.6
0.5%	74.2	30.5	26.6	54.1	105.8	20.1	75.3	31.2	27.1	54.7	109.2	20.9
1%	<b>74.6</b>	<b>30.7</b>	<b>26.7</b>	<b>54.2</b>	<b>106.7</b>	<b>20.3</b>	<b>75.4</b>	<b>31.2</b>	<b>27.2</b>	<b>54.5</b>	<b>109.7</b>	<b>20.9</b>
5%	74.4	30.7	26.7	54.3	106.5	20.3	75.5	31.5	27.3	54.8	109.9	20.9
20%	74.5	31.0	26.7	54.4	107.1	20.3	75.7	31.9	27.3	55.0	110.5	21.0
50%	74.6	30.9	26.8	54.4	107.3	20.3	75.6	31.8	27.5	55.0	111.0	21.1
100%	74.4	30.8	26.7	54.4	106.6	20.3	75.8	31.7	27.3	54.9	110.6	21.0

$\mathcal{N}(\mu, \sigma^2)$  but have a larger standard deviation as shown in Fig. 6 (c) and (d), which mitigate the modality gap but not significantly. Fig. 6 (e) shows that our mapping module driven by multi-variate Gaussian distribution can effectively reduce the modality gap.

3) *Impact of Reverse Mapping Module*: To explore the effectiveness of our proposed reverse mapping module, we report the performances of our TIPCap with the reverse mapping module removed, as shown in TABLE V (C). Comparing TABLE V (A) and (C), it can be seen that the reverse mapping module brings slight performance improvement over all metrics. More importantly, the reverse mapping module is essential for our approaches in setting 4, which makes our proposed TIPCap still applicable when only text data is available.

4) *Impact of prefix projector*: The prefix projector aims to convert CLIP embedding to prefix embeddings as the input of GPT-2 model. Following ClipCap [59] and CapDec [25], we use a transformer-based architecture as prefix projector. To explore the effect of prefix projector, we perform ablation studies on prefix length of  $L$  and layer number of  $N$  respectively, as shown in Table VI. The bold results denote our default implementation.

From the results, TIPCap shows good robustness when increasing  $L$  or  $N$ . It brings an advantage that we do not need a heavy module to perform the indispensable projection from CLIP embedding space to GPT-2 embedding space.

5) *Impact of available external data*: In setting 1, the parameter  $\vec{\mu}$  and  $\Sigma$  are estimated from human-annotated high-quality paired data. We perform estimation using different

TABLE VIII

PERFORMANCE COMPARISON WITH DIFFERENT NUMBER OF IMAGES USED IN SETTING 3. ALL IMAGES ARE RANDOMLY SAMPLED FROM YFCC15M.

#img	CLIP RN50x4						CLIP ViT-L/14					
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
500	73.0	30.3	26.3	53.8	104.0	19.9	73.8	31.1	26.7	54.0	107.3	20.4
1K	72.9	30.3	26.3	53.8	104.2	19.9	73.9	31.2	26.7	54.0	107.3	20.4
5K	<b>73.0</b>	<b>30.4</b>	<b>26.5</b>	<b>53.8</b>	<b>104.5</b>	<b>20.0</b>	<b>73.7</b>	<b>31.2</b>	<b>26.7</b>	<b>53.9</b>	<b>107.5</b>	<b>20.5</b>
10K	72.9	30.4	26.4	53.7	104.0	19.9	73.8	31.2	26.7	53.9	107.3	20.4

TABLE IX

PERFORMANCE COMPARISON OF THREE DIFFERENT SAMPLED DATA UNDER SETTING 1. EXPERIMENTS ARE CONDUCTED ON MS-COCO.

	CLIP RN50x4						CLIP ViT-L/14					
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
sample 1	74.6	30.7	26.7	54.2	106.7	20.3	75.4	31.2	27.2	54.5	109.7	20.9
sample 2	74.8	31.0	26.7	54.3	107.2	20.5	75.3	31.4	27.1	54.5	109.9	20.8
sample 3	74.7	30.9	26.7	54.3	107.2	20.3	75.4	31.6	27.3	54.7	110.7	21.0

TABLE X

PERFORMANCE COMPARISON OF DIFFERENT IMAGE SOURCE (“YFCC15M” v.s. “MS-COCO”) UNDER SETTING 3. EXPERIMENTS ARE CONDUCTED ON MS-COCO.

	CLIP RN50x4						CLIP ViT-L/14					
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
YFCC15M	73.0	30.4	26.5	53.8	104.5	20.0	73.7	31.2	26.7	53.9	107.5	20.5
MS-COCO	72.9	30.6	26.5	53.8	104.9	20.1	74.1	31.5	26.7	54.1	107.9	20.4

ratios of paired data to study its effect as shown in Table VII. Theoretically, the estimated parameters can more accurately characterize the modality bias when more paired data is available. From the results, we can see that the performance tends to stabilize when more than 1% paired data is available, which indicates that we can obtain sufficiently accurate estimated parameters without all data.

In setting 3, we use few source-agnostic image data to introduce the latent prior information of CLIP image embedding space. To explore its effect, we sample different numbers of image data from YFCC15M dataset for training. The results are reported in Table VIII. When using only 500 images, the performance is still advantageous enough, which indicates the data efficiency of our TIPCap.

6) *Impact of different sampled data under setting 1*: Under Setting 1, the parameters of our mapping module (*i.e.* mean and covariance) are estimated from sampled paired data (*e.g.* MS-COCO or Flickr30K), which brings a question: whether different sampled data will influence the performance. To verify the influence of different sampling data, we conduct experiments on MS-COCO dataset. We sample 1% paired MS-COCO data for the estimation of mean and covariance, and perform three times sampling, resulting in three different sets of parameters. Then we train our TIPCap with these three different estimations of  $\mathcal{N}(\vec{\mu}, \Sigma)$ .

The results are reported in Table IX, and show that the performances are robust to different sampling data, which further demonstrates the effectiveness and robustness of applying multivariate Gaussian distribution to mitigate the modality gap.

7) *Impact of image source under setting 3*: As described in Section 4.1, we use 5K YFCC images for our TIPCap training under setting 3, which aims to prove that heterologous images

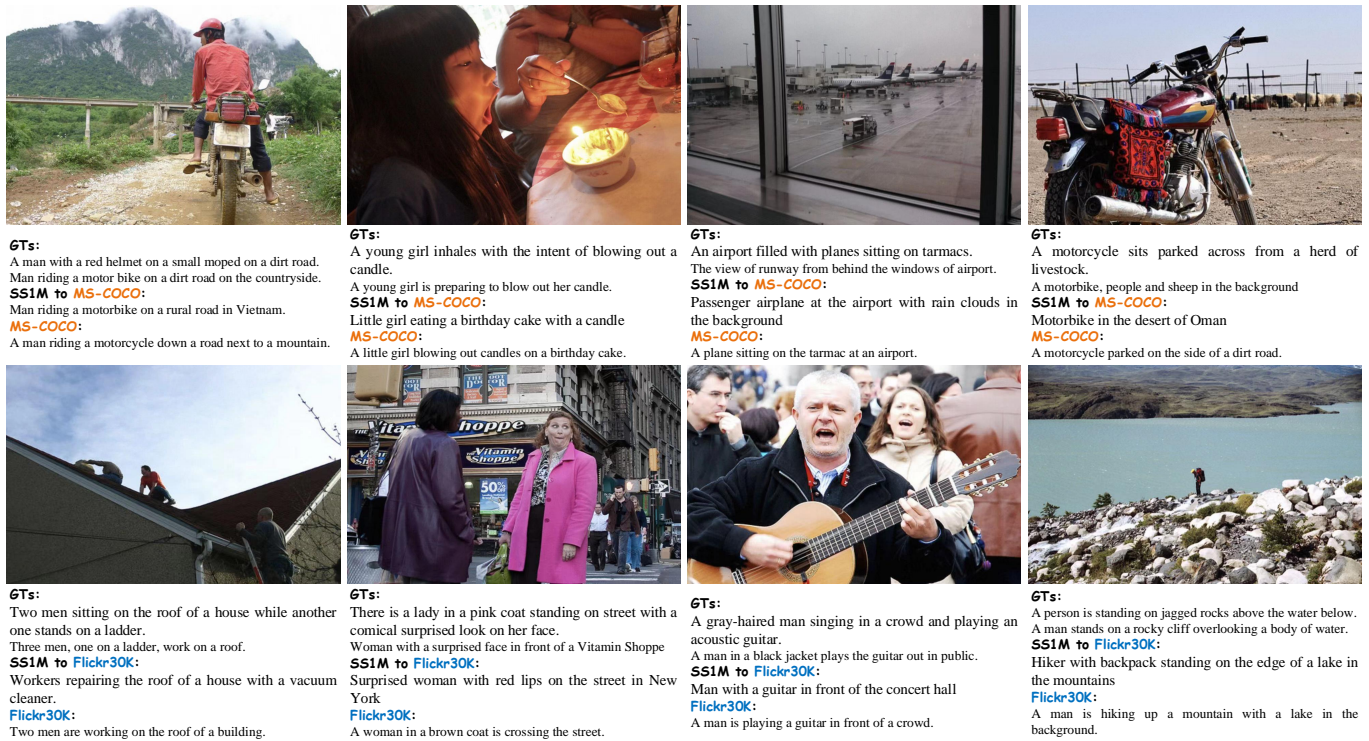


Fig. 7. Examples of captions generated by TIPCap. The first two images comes from MS-COCO dataset and the second two images comes from Flickr30K dataset. “SS1M to MS-COCO / Flickr30K” indicate that captions is generated by TIPCap trained on SS1M. “MS-COCO” and “Flickr30K” denote TIPCap trained on MS-COCO and Flickr30K, respectively.

are also applicable and effective. We also conduct experiments using homologous images on MS-COCO dataset, the results are shown in Table X. Specifically, we adopt the 5K images of MS-COCO karpathy validation split, which are homologous to the training corpus but are not paired and do not affect the performance evaluation on karpathy test split.

From the results, both YFCC images and MS-COCO images achieve comparable performance under setting 3, which indicates that the requirement for image data is source-agnostic. Combining the results shown in Table 8, it indicates that although available images are few and heterologous with training text corpus, TIPCap can still achieve competitive performance under setting 3. Compared under setting 4, TIPCap under setting 3 can achieve better performance because we introduce prior information of real image for training. Although it is difficult to collect high-quality paired image data for text corpus in real-world scenarios, the image data-efficiency under setting 3 makes the collection cost of image data affordable.

#### D. More Generalization Analysis

To further quantitatively evaluate the generalization performance of our proposed TIPCap, we conduct experiments on two new datasets: NoCaps [60] and SS1M [14].

1) *Performance on NoCaps dataset:* MS-COCO is limited to 80 classes, NoCaps [60] provides a benchmark to measure the open-set capability for novel objects (unseen classes in MS-COCO dataset). The NoCaps dataset contains only

validation and test sets, and is divided into three parts: *in-domain* contains images portraying only MS-COCO classes; *near-domain* contains both MS-COCO and novel classes; *out-of-domain* contains only novel classes. Following DeCap [26], we evaluate our TIPCap using the validation set on official evaluation server<sup>1</sup>. Table XII shows the performance comparison on MS-COCO karpathy test split and NoCaps validation split. Compared to DeCap, TIPCap achieves significant improvement on all metrics.

2) *Performance using SS1M dataset:* SS1M [14] is a web-collected text corpus, which uses the name of 80 MS-COCO classes as keywords to crawl image descriptions from Shutterstock<sup>2</sup>, resulting in 2,322,628 distinct image descriptions in total that is larger than MS-COCO and Flickr30K. Fig. XI gives a comparison of some text examples between SS1M and MS-COCO.

We explore using SS1M corpus to train our TIPCap under setting 3, where the image data is sampled from YFCC15M, identical to the experiments conducted on the MS-COCO. After training, we evaluate TIPCap on MS-COCO and Flickr30K test datasets as shown in Table XIII (middle 2 rows), the results indicate that TIPCap can be easily extended to larger datasets and also performs good zero-shot performance. Fig. 7 shows some examples of captions generated by TIPCap trained on SS1M under setting 3 (see “SS1M to MS-COCO / Flickr30K”), which can correctly describe images, even with more descriptive details.

<sup>1</sup><https://eval.ai/web/challenges/challenge-page/355/overview>

<sup>2</sup><https://www.shutterstock.com>

TABLE XI  
COMPARISON OF TEXT BETWEEN SS1M AND MS-COCO.

<p><b>SS1M text examples:</b></p> <ol style="list-style-type: none"> <li>1. Bus Lane Sign in Urban Setting in Black and White Sepia Tone</li> <li>2. France, Paris, 04/04/2015, parc de bercy, a pathway with lush green vegetation and a park bench</li> <li>3. Cat looking at sea in Santorini, Greece</li> <li>4. Yellow bus handles inside the vehicle with blur background</li> <li>5. Many people in the street on the bicycles, holiday</li> <li>6. Concept fighting, friendship, promise, success, Two people put their hands together and raised</li> </ol>
<p><b>MS-COCO text examples:</b></p> <ol style="list-style-type: none"> <li>1. Modern living room interior with many green plants</li> <li>2. A couple of chairs that are at a table</li> <li>3. A flat screen TV mounted to a wall.</li> <li>4. A mid sized bathroom with toilette, shower and vanity mirror above a sink.</li> <li>5. A dock area with various toilets and a television on it.</li> <li>6. There is a room with distinctive things in the picture.</li> </ol>

TABLE XII  
RESULTS ON MS-COCO KARPATHY TEST SPLIT AND NOCAPS VALIDATION SPLIT, WHERE “IN”, “NEAR” AND “OUT” INDICATE “*in-domain*”, “*near-domain*” AND “*out-of-domain*” RESPECTIVELY. ALL MODELS ARE TRAINED ON MS-COCO TEXT CORPUS.

	MS-COCO						NoCaps val (CIDEr)			
	B-1	B-4	M	R	C	S	in	near	out	overall
DeCap	-	24.7	25.0	-	91.2	18.7	65.2	47.8	25.8	45.9
<b>CLIP RN50x4:</b>										
TIPCap (S1)	74.6	30.7	26.7	54.2	106.7	20.3	77.6	63.3	44.8	61.6
TIPCap (S2)	72.7	28.6	25.6	52.5	100.6	19.6	71.1	56.6	39.0	55.1
TIPCap (S3)	73.0	30.4	26.5	53.8	104.5	20.0	74.6	59.2	37.1	56.9
TIPCap (S4)	71.3	29.8	26.2	53.4	102.1	19.4	74.3	59.2	36.4	56.7
<b>CLIP ViT-L/14:</b>										
TIPCap (S1)	75.4	31.2	27.2	54.5	109.7	20.9	81.7	67.5	49.8	65.9
TIPCap (S2)	73.4	29.1	26.0	52.7	103.5	20.2	76.1	61.3	44.5	60.0
TIPCap (S3)	73.7	31.2	26.7	53.9	107.5	20.5	77.3	62.7	40.6	60.9
TIPCap (S4)	73.3	31.4	26.9	54.2	106.6	20.2	80.2	62.3	39.6	60.3

TABLE XIII  
CROSS-DOMAIN IMAGE CAPTIONING PERFORMANCE. WE TRAIN OUR TIPCAP USING SS1M DATASET UNDER SETTING 3 AND SETTING 1, AND PERFORM EVALUATION ON MS-COCO AND FLICKR30K DATASETS.

	SS1M $\Rightarrow$ MS-COCO						SS1M $\Rightarrow$ Flickr30K					
	B-1	B-4	M	R	C	S	B-1	B-4	M	R	C	S
DeCap	-	8.9	17.5	-	50.6	13.1	-	-	-	-	-	-
TIPCap (S3, CLIP RN50x4)	54.7	11.5	18.0	35.7	53.0	13.9	52.2	9.1	15.4	32.7	32.2	11.0
TIPCap (S3, CLIP ViT-L/14)	56.5	12.3	18.6	37.2	57.3	14.5	56.6	10.7	16.2	35.3	36.3	11.3
TIPCap (S1, CLIP RN50x4)	62.9	16.5	20.6	41.3	69.2	16.3	59.2	13.0	17.1	37.2	38.7	11.8
TIPCap (S1, CLIP ViT-L/14)	65.1	18.3	21.5	43.2	73.9	17.0	62.6	14.4	18.1	39.3	44.0	12.5

Another interesting results are also shown in Table XIII (bottom 2 rows), we train our TIPCap on setting 1, where the required parameters  $\vec{\mu}$  and  $\Sigma$  are directly estimated from 1% MS-COCO paired data. It can be seen that this paradigm achieves better performance, even though the distribution parameters are estimated from another dataset. We attribute this phenomenon to two reasons: 1) As shown in Table XI, the text corpus in SS1M have a good fluency and language accuracy; 2) Comparing SS1M and MS-COCO, the styles of their text corpus are not exactly same but similar. Their text both consists of a main object and additional descriptive attribution like state and position.

## V. CONCLUSION

In this paper, we propose a unified solution TIPCap with text-centric training data for image captioning, which can almost cover the vast majority of data configurations in real-world scenarios. In TIPCap, the mapping module, which is driven by a multivariate Gaussian distribution and aims to mitigate the modality gap between image and text embedding. Additionally, TIPCap can incorporate optional prompt information, which is proposed to improve generated captions. Extensive experiments also demonstrate the effectiveness of our TIPCap. We believe this study can give a new paradigm and benefit the community for image captioning.

## REFERENCES

- [1] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 3156–3164.
- [2] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 2048–2057.
- [3] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 6077–6086.
- [4] L. Huang, W. Wang, J. Chen, and X. Wei, "Attention on attention for image captioning," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 4633–4642.
- [5] Y. Pan, T. Yao, Y. Li, and T. Mei, "X-linear attention networks for image captioning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 10968–10977.
- [6] Y. Luo, J. Ji, X. Sun, L. Cao, Y. Wu, F. Huang, C. Lin, and R. Ji, "Dual-level collaborative transformer for image captioning," in *Proc. AAAI Conf. Artif. Intell.*, 2021, pp. 2286–2293.
- [7] Y. Wang, J. Xu, and Y. Sun, "End-to-end transformer based model for image captioning," in *Proc. AAAI Conf. Artif. Intell.*, 2022, pp. 2585–2594.
- [8] L. Wu, M. Xu, L. Sang, T. Yao, and T. Mei, "Noise augmented double-stream graph convolutional networks for image captioning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 31, no. 8, pp. 3118–3127, 2021. [Online]. Available: <https://doi.org/10.1109/TCSVT.2020.3036860>
- [9] W. Jiang, W. Zhou, and H. Hu, "Double-stream position learning transformer network for image captioning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 11, pp. 7706–7718, 2022. [Online]. Available: <https://doi.org/10.1109/TCSVT.2022.3181490>
- [10] S. Cao, G. An, Z. Zheng, and Z. Wang, "Vision-enhanced and consensus-aware transformer for image captioning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 10, pp. 7005–7018, 2022. [Online]. Available: <https://doi.org/10.1109/TCSVT.2022.3178844>
- [11] J. Zhang, Y. Xie, W. Ding, and Z. Wang, "Cross on cross attention: Deep fusion transformer for image captioning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 33, no. 8, pp. 4257–4268, 2023. [Online]. Available: <https://doi.org/10.1109/TCSVT.2023.3243725>
- [12] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 740–755.
- [13] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *Trans. Assoc. Comput. Linguistics*, vol. 2, pp. 67–78, 2014.
- [14] Y. Feng, L. Ma, W. Liu, and J. Luo, "Unsupervised image captioning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 4125–4134.
- [15] I. Laina, C. Rupprecht, and N. Navab, "Towards unsupervised image captioning with shared multimodal embeddings," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 7413–7423.
- [16] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proc. Conf. N. Am. Chapter Assoc. Comput. Linguistics: Hum. Lang. Technol.*, 2019, pp. 4171–4186.
- [17] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," 2019.
- [18] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.
- [19] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proc. Int. Conf. Mach. Learn.*, 2021, pp. 8748–8763.
- [20] C. Jia, Y. Yang, Y. Xia, Y. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *Proc. Int. Conf. Mach. Learn.*, 2021, pp. 4904–4916.
- [21] J. Li, D. Li, C. Xiong, and S. C. H. Hoi, "BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation," in *Proc. Int. Conf. Mach. Learn.*, 2022, pp. 12888–12900.
- [22] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 3558–3568.
- [23] Y. Tewel, Y. Shalev, I. Schwartz, and L. Wolf, "Zerocap: Zero-shot image-to-text generation for visual-semantic arithmetic," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 17897–17907.
- [24] Y. Su, T. Lan, Y. Liu, F. Liu, D. Yogatama, Y. Wang, L. Kong, and N. Collier, "Language models can see: Plugging visual controls in text generation," *ArXiv*, vol. abs/2205.02655, 2022.
- [25] D. Nukrai, R. Mokady, and A. Globerson, "Text-only training for image captioning using noise-injected CLIP," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2022, pp. 4055–4063.
- [26] W. Li, L. Zhu, L. Wen, and Y. Yang, "Decap: Decoding clip latents for zero-shot captioning via text-only training," in *Proc. Int. Conf. Learn. Representations*, 2023.
- [27] S. Gu, C. Clark, and A. Kembhavi, "I can't believe there's no images! learning visual tasks using only language data," *ArXiv*, vol. abs/2211.09778, 2022.
- [28] J. Li, D. Li, S. Savarese, and S. C. H. Hoi, "BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models," *ArXiv*, vol. abs/2301.12597, 2023.
- [29] 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 *Proc. Int. Conf. Mach. Learn.*, 2022, pp. 23318–23340.
- [30] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Gao, and Y. J. Lee, "Segment everything everywhere all at once," *ArXiv*, vol. abs/2304.06718, 2023.
- [31] K. Papineni, S. Roukos, T. Ward, and W. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proc. Annu. Meet. Assoc. Comput. Linguist.*, 2002, pp. 311–318.
- [32] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2014, pp. 1724–1734.
- [33] H. Chen, G. Ding, Z. Lin, S. Zhao, and J. Han, "Show, observe and tell: Attribute-driven attention model for image captioning," in *Proc. Int. Joint Conf. Artif. Intell.*, 2018, pp. 606–612.
- [34] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo, "Image captioning with semantic attention," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 4651–4659.
- [35] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [36] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [37] X. Dong, C. Long, W. Xu, and C. Xiao, "Dual graph convolutional networks with transformer and curriculum learning for image captioning," in *Proc. ACM Int. Conf. Multimed.*, 2021, pp. 2615–2624.
- [38] W. Nie, J. Li, N. Xu, A. Liu, X. Li, and Y. Zhang, "Triangle-reward reinforcement learning: A visual-linguistic semantic alignment for image captioning," in *Proc. ACM Int. Conf. Multimed.*, 2021, pp. 4510–4518.
- [39] S. Herdade, A. Kappeler, K. Boakye, and J. Soares, "Image captioning: Transforming objects into words," in *Proc. Adv. neural inf. proces. syst.*, 2019, pp. 11135–11145.
- [40] J. Yu, J. Li, Z. Yu, and Q. Huang, "Multimodal transformer with multi-view visual representation for image captioning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 12, pp. 4467–4480, 2020. [Online]. Available: <https://doi.org/10.1109/TCSVT.2019.2947482>
- [41] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, "Meshed-memory transformer for image captioning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 10575–10584.
- [42] J. Wang, Y. Zhang, M. Yan, J. Zhang, and J. Sang, "Zero-shot image captioning by anchor-augmented vision-language space alignment," *ArXiv*, vol. abs/2211.07275, 2022.
- [43] H. Tan and M. Bansal, "LXMERT: learning cross-modality encoder representations from transformers," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2019, pp. 5099–5110.
- [44] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Proc. Adv. neural inf. proces. syst.*, 2019, pp. 13–23.
- [45] Y. Chen, L. Li, L. Yu, A. E. Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, "UNITER: universal image-text representation learning," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 104–120.
- [46] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei, Y. Choi, and J. Gao, "Oscar: Object-semantics aligned

- pre-training for vision-language tasks,” in *Proc. Eur. Conf. Comput. Vis.*, ser. Lecture Notes in Computer Science, vol. 12375, 2020, pp. 121–137.
- [47] J. Li, R. R. Selvaraju, A. Gotmare, S. R. Joty, C. Xiong, and S. C. Hoi, “Align before fuse: Vision and language representation learning with momentum distillation,” in *Proc. Adv. neural inf. proces. syst.*, 2021, pp. 9694–9705.
- [48] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” in *Proc. Int. Conf. Learn. Representations*, 2014.
- [49] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. Lowe, “Training language models to follow instructions with human feedback,” in *Proc. Adv. neural inf. proces. syst.*, 2022.
- [50] A. Karpathy and L. Fei-Fei, “Deep visual-semantic alignments for generating image descriptions,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 3128–3137.
- [51] B. Thomee, D. A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, and L. Li, “YFCC100M: the new data in multimedia research,” *Commun. ACM*, vol. 59, no. 2, pp. 64–73, 2016.
- [52] M. J. Denkowski and A. Lavie, “Meteor universal: Language specific translation evaluation for any target language,” in *Proc. ACL Workshop Statistical Machine Translation*, 2014, pp. 376–380.
- [53] C.-Y. Lin, “ROUGE: A package for automatic evaluation of summaries,” in *Proc. ACL Workshop Text Summarization Branches Out*, 2004, pp. 74–81.
- [54] R. Vedantam, C. L. Zitnick, and D. Parikh, “Cider: Consensus-based image description evaluation,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 4566–4575.
- [55] P. Anderson, B. Fernando, M. Johnson, and S. Gould, “Spice: Semantic propositional image caption evaluation,” in *Proc. Eur. Conf. Comput. Vis.*, vol. 9909, 2016, pp. 382–398.
- [56] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Proc. Adv. neural inf. proces. syst.*, 2019, pp. 8024–8035.
- [57] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *Proc. Int. Conf. Learn. Representations*. OpenReview.net, 2019.
- [58] L. Zhou, H. Palangi, L. Zhang, H. Hu, J. J. Corso, and J. Gao, “Unified vision-language pre-training for image captioning and VQA,” in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 13 041–13 049.
- [59] R. Mokady, A. Hertz, and A. H. Bermano, “Clipcap: CLIP prefix for image captioning,” *ArXiv*, vol. abs/2111.09734, 2021.
- [60] H. Agrawal, P. Anderson, K. Desai, Y. Wang, X. Chen, R. Jain, M. Johnson, D. Batra, D. Parikh, and S. Lee, “nocaps: novel object captioning at scale,” in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 8947–8956.