

# Enhancing Visual Grounding in Vision-Language Pre-training with Position-Guided Text Prompts

Alex Jinpeng Wang, Pan Zhou, Mike Zheng Shou, Shuicheng Yan, *Fellow, IEEE*

**Abstract**—Vision-Language Pre-Training (VLP) has demonstrated remarkable potential in aligning image and text pairs, paving the way for a wide range of cross-modal learning tasks. Nevertheless, we have observed that VLP models often fall short in terms of visual grounding and localization capabilities, which are crucial for many downstream tasks, such as visual reasoning. In response, we introduce a novel Position-guided Text Prompt (*PTP*) paradigm to bolster the visual grounding abilities of cross-modal models trained with VLP. In the VLP phase, *PTP* divides an image into  $N \times N$  blocks and employs a widely-used object detector to identify objects within each block. *PTP* then reframes the visual grounding task as a fill-in-the-blank problem, encouraging the model to predict objects in given blocks or regress the blocks of a given object, exemplified by filling “[*P*]” or “[*O*]” in a PTP sentence such as “The block [*P*] has a [*O*]”. This strategy enhances the visual grounding capabilities of VLP models, enabling them to better tackle various downstream tasks. Additionally, we integrate the second-order relationships between objects to further enhance the visual grounding capabilities of our proposed PTP paradigm. Incorporating *PTP* into several state-of-the-art VLP frameworks leads to consistently significant improvements across representative cross-modal learning model architectures and multiple benchmarks, such as zero-shot Flickr30k Retrieval (+5.6 in average recall@1) for ViLT baseline, and COCO Captioning (+5.5 in CIDEr) for the state-of-the-art BLIP baseline. Furthermore, *PTP* attains comparable results with object-detector-based methods and a faster inference speed, as it discards its object detector during inference, unlike other approaches. Our code and pre-trained models are available at <https://github.com/sail-sg/ptp>.

**Index Terms**—Vision-Language Pre-Training, Position-guided Text Prompt, Fill-in-the-blank, Visual Grounding

## 1 INTRODUCTION

Vision-and-language pre-training (VLP) models, such as CLIP [41], ALIGN [20], and CoCa [59], have significantly improved cross-modal tasks like visual question answering [4], natural language visual reasoning [47], and image captioning [1], [9]. The success of these models can be attributed to a two-stage learning process: pre-training on large image-caption data to enhance generalization capabilities, followed by fine-tuning on downstream tasks for seamless adaptation. This efficient pre-training and fine-tuning paradigm has established VLP models as a dominant force in the multi-modal research domain, highlighting their potential for further development and continued performance improvements across various cross-modal applications.

In VLP, visual grounding plays a crucial role in various tasks, as evidenced by previous research [3], [57]. Traditional VLP models [3], [31], [63], depicted at the top of Figure 1 (a), utilize a Faster R-CNN [43] pre-trained on the 1600-class Visual Genome [23] to extract salient region features and bounding boxes. These models then take both the bounding box and object feature as input, allowing them to identify objects within the salient region and determine their locations. However, by using region features as input, the model selectively attends to information within bounding boxes while neglecting contextual data beyond their boundaries [17]. This limitation can result in suboptimal performance on downstream tasks, necessitating the use of additional

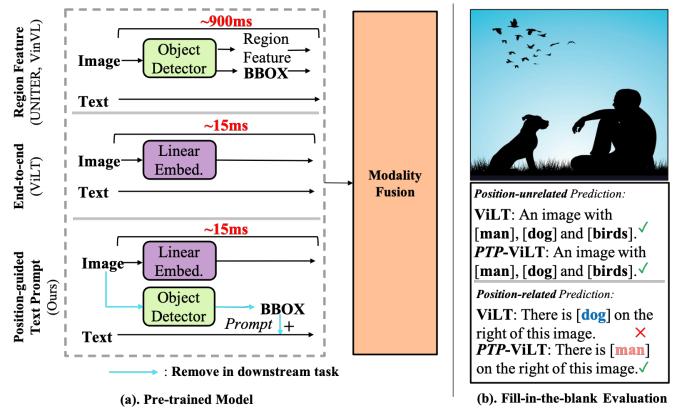


Fig. 1: Comparison of three VLP learning frameworks and their performance. (a) compares region feature-based VLP (RF-VLP), end-to-end VLP (E2E-VLP), and our position-guided text prompt-based VLP (PTP-VLP). Our PTP-VLP requires only about 15ms for inference, which is on par with E2E-VLP and significantly faster than RF-VLP. (b) On position-aware questions that are commonly encountered in many downstream tasks, both RF-VLP and PTP-VLP can accurately predict objects with masked text and image input. In contrast, E2E-VLP struggles to pinpoint the position information of the object in the image.

M. Shou is the corresponding author.

- A. Wang and M. Shou are with the Show Lab, National University of Singapore.
- P. Zhou and S. Yan are with the Sea AI Lab, Singapore.

object detectors to extract objects and consequently causing significantly slower inference speeds.

To get rid of region feature for higher efficiency, recent works [17], [22] (the middle of Figure 1 (a)) adopt raw image as input instead of region features, and train the model with Image Text Matching [10] and Masked Language Modeling

[12] loss end-to-end. Despite their faster speed, these models cannot well learn the object positions and also their relations. As demonstrated in Figure 1 (b), a well-trained ViLT model [22] can successfully identify objects in an image. However, it fails to precisely learn object positions. For instance, it incorrectly predicts “*the dog is on the right of this image.*” During fine-tuning evaluation, downstream tasks necessitate object position information for a comprehensive understanding of the image. This gap significantly hinders performance on downstream tasks, emphasizing the need for improved object position and relation learning.

In this work, our goal is to address the position learning issue in end-to-end (e2e) models while maintaining fast inference times for downstream tasks. Drawing inspiration from recent prompt learning methods [21], [34], [42], [58], we introduce a novel and effective **Position-guided Text Prompt (PTP)** paradigm (depicted at the bottom of Figure 1 (a)) for VLP. The core insight is that by incorporating position-based co-referential markers into both image and text, visual grounding can be transformed into a fill-in-the-blank problem, significantly simplifying the learning of object information. To establish a connection between language expressions and image, PTP comprises two components: (1) block tag generation, which divides the image into  $N \times N$  blocks and identifies objects within each block, and (2) text prompt generation, which embeds the query text into a position-based text query template. This innovative approach facilitates more accurate position learning while retaining the efficiency advantages of e2e models.

Integrating position information into the pre-training phase, our PTP significantly enhances the visual grounding capabilities of VLP models. We also investigate second-order relations between objects to improve reasoning ability. Importantly, our approach maintains fast inference times, as we do not rely on object detectors for downstream tasks. Experimental results reveal that our method substantially outperforms its counterparts, particularly in the zero-shot setting. Our proposed model, PTP-BLIP, demonstrates exceptional performance on the zero-shot image-to-text retrieval Recall@1 task on the COCO dataset, achieving a 6.9% absolute accuracy gain over CoCa [59] while using considerably less training data (4M vs. 3B) and a smaller model size (220M vs. 2.1B). Moreover, PTP’s effectiveness extends beyond retrieval, as evidenced by its success in other visual language tasks such as visual grounding and image captioning. These results underline the potential of our position-guided approach for advancing the state of the art in cross-modal learning.

Our contributions can be summarized as follows: 1). We introduce a novel pre-training paradigm for vision-language models, called cross-modal prompt-based pre-training, which explicitly incorporates position information into the prompt. To the best of our knowledge, this represents the first attempt at employing such an approach for pre-training vision-language models. 2). We design and assess multiple configurations of high-quality cross-modal prompts for our proposed model, PTP, demonstrating the versatility and adaptability of our approach. 3). We carry out comprehensive experiments using four backbone models, showcasing the effectiveness of PTP across a range of vision-language tasks. We expand our approach to encompass

data at the billion-level scale and demonstrate its efficacy with a potent Large Language Model, further highlighting its potential to advance the state of the art in cross-modal learning.

This journal paper extends our previous work [51] in several ways: *First*, we propose a novel second-order prompt to further enhance the understanding of object relationships. This innovation leads to improved state-of-the-art results in various downstream tasks, such as visual-text retrieval, visual question answering, and image captioning. *Second*, we expand our evaluation by including more downstream tasks, specifically visual grounding on the RefCOCO [60] and RefCOCO+ [60] datasets. Additionally, we explore video question answering on MSVD [56], TVQA [25], and TGIF [19], providing a more comprehensive comparison with VLP-related works. *Third*, we provide an in-depth analysis of the visual grounding ability, enriched visualizations of masked blocks, and a thorough ablation study. These elements together showcase the effectiveness and robustness of our second-order PTP approach in advancing the state of the art in cross-modal learning. *Finally*, we have trained the PTP on the extensive DataComp-1B [14] dataset at the billion-level and integrated it with popular Large Language Models, thereby showcasing the versatility and broad applicability of PTP.

## 2 RELATED WORK

### 2.1 Vision-language Pre-training Architectures

Existing VLP models can be roughly grouped into three categories based on their architectures: one-stream models, dual-stream models, and dual-stream with fusion encoder models. We provide an overview of all three architectures:

1) *One-stream*: (e.g., UNITER [10], ViLT [22]) as shown in Figure 2 (a), operates on a concatenation of image and text inputs. 2) *Dual-stream*: (e.g., CLIP [41]) depicted in Figure 2 (b), employs separate and equally expensive transformer encoders for each modality. The two modalities are not concatenated at the input level, with interaction between the pooled image and text vectors occurring at a shallow layer. 3) *Dual-stream with Fusion Encoder*: (e.g., BLIP [27]) illustrated in Figure 2 (c), is a combination of one-stream and dual-stream models that allows for intermediate interaction.

In this work, without loss of generality, we focus on prompting all three types of VLP models due to their prevalence and adaptability to various downstream tasks.

### 2.2 Prompt Learning for Computer Vision

Prompt learning was initially designed for probing knowledge in pre-trained language models and adapting them to specific downstream tasks [34], [42]. In recent years, there has been a growing interest in studying prompt tuning for vision tasks, including multi-modal learning and image understanding. The pioneering work Color Prompt [58] introduces color prompts on images and text color descriptions for visual grounding. Most related to our work is Multi-modality Prompt [21], which presents multi-modal prompt tuning for vision-language pre-training models, achieving promising results on various vision-language tasks.

However, these efforts, akin to early NLP research, focus on prompt engineering during the fine-tuning stage, leaving

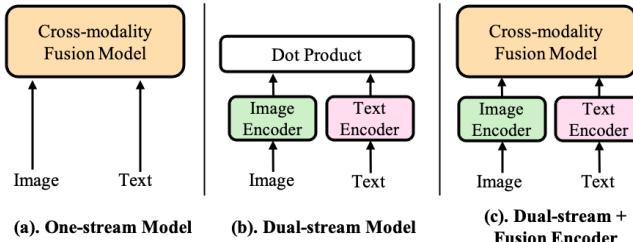


Fig. 2: Three widely-used categories of vision-and-language models. The primary distinction lies in the stage at which cross-modality information fusion occurs. One-stream models perform fusion at an early stage, while dual-stream models fuse information at a late stage. Lastly, the third type of models integrates information at a middle stage, striking a balance between the other two approaches.

the pre-training phase unaffected. In contrast, the goal of using prompt design in our work is to equip the model with the ability to understand semantic concepts at a finer level during the pre-training stage, laying a stronger foundation for downstream tasks.

### 2.3 Learning Position Information in VLP

The grounding ability has proven to be essential for multiple cross-modal tasks [29], [35]. To introduce this ability into VLP models, bottom-up and top-down [3] approaches and their follow-up works [10], [31] concatenate region features and bounding box vectors together as input signals. However, object extraction is time-consuming during inference for downstream tasks. Recently, some works [29], [35], [62] propose training VLP models with additional object localization loss or word patch alignment loss. These methods, however, are difficult to extend as they are specifically designed for particular frameworks. In contrast, we aim to propose a general framework for learning position information. To this end, we introduce a simple text prompt that can be easily integrated into existing frameworks, providing a versatile solution for capturing position information in VLP.

## 3 POSITION-GUIDED TEXT PROMPT

In this section, we first provide a detailed explanation of our proposed Position-guided Text Prompt paradigm (PTP for short). Following this, we demonstrate how to incorporate it into existing vision-language pre-training (VLP) frameworks to enhance their visual grounding capabilities. We use the classical and popular models VILT [22], CLIP [41], and BLIP [27] as examples to showcase the integration.

### 3.1 PTP Paradigm

To enhance the visual grounding ability of cross-modal models trained using VLP, we propose a novel and effective Position-guided Text Prompt (PTP) that helps a cross-modal model perceive objects and align them with the relevant text. PTP differs from conventional vision-language alignment methods, such as [3], [10], [31], [63], which concatenate object features and bounding boxes as input to learn the

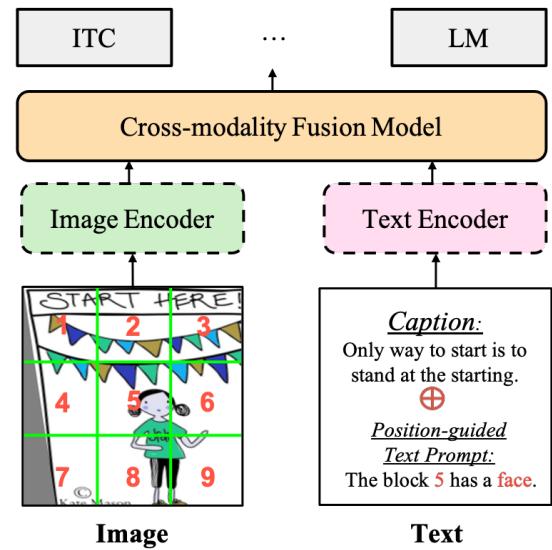


Fig. 3: The overview of our pipeline. Our PTP can be seamlessly integrated with any pre-training framework, including one-stream, dual-stream, and dual-stream with fusion encoder models, as well as most pre-training objectives. The dashed line in the figure indicates that certain models may not be present. For downstream tasks, we remove the text prompt and evaluate the model using standard procedures, ensuring that the integration of PTP does not interfere with task-specific evaluation metrics.

alignment between objects and relevant text. This alternative approach offers several advantages, as discussed in Section 3.2. As illustrated in Figure 3, PTP consists of two steps: 1) block tag generation, which divides an input image into several blocks and identifies objects in each block; and 2) text prompt generation, which reformulates the visual grounding task into a fill-in-the-blank problem based on the object position information from step 1). By solving the fill-in-the-blank problem in PTP, one can easily integrate PTP into a VLP model. We provide details on these steps below.

#### 3.1.1 Block Tag Generation

As shown in Figure 3, for each image-text pair in the training phase, we evenly divide the input image into  $N \times N$  blocks. Then we identify the object in each block using one of the following two methods:

**(1) Object Detector.** First, we adopt a strong Faster-RCNN [43] used in VinVL [63] to extract all objects for each image. This Faster-RCNN version is based on ResNeXt152 and trained on the 1600-class Visual Genome [23]. Then we select the top- $K$  objects, denoted by  $\mathcal{O} = o_{i=1}^K$ , with the highest prediction confidence, where  $o_i = (z_i, q_i)$  denotes an object with a 4-dimensional region position vector  $z$  and object category  $q$ . For each block, we select the objects whose region centers are in that block. Finally, the block tag for this block is  $q$  of the selected objects. In this work, we generate object tags using an object detector by default.

**(2) CLIP Model.** As an alternative to using a heavy object detector, recent works [65], [66] have attempted to generate region supervision based on CLIP [41] due to its efficiency and effectiveness. Inspired by these works,

PTP can also generate block-wise object supervision using the CLIP (ViT-B) model<sup>1</sup>. First, we extract  $M$  (3000 by default) keywords/phrases that are most frequent in the text corpus<sup>2</sup>. These keywords/phrases form our vocabulary  $V$ . Then, we extract the text feature  $e_i, i \in [1, \dots, M]$  of all these  $M$  text embedding using the CLIP text encoder.

Furthermore, we take the image embedding  $h$  from each block and compute the similarity across every text feature. The keyword/phrase with the highest similarity score is selected as the object tag for this particular block. Formally, the index of the object tag per block is computed as:

$$I = \text{argmax}_{y \in [1, \dots, M]} \left( \frac{\exp(h^T e_y)}{\sum w \in V \exp(h^T e_w)} \right), \quad (1)$$

where  $h$  is the visual feature embedding of the selected block. Compared to the object detector, the CLIP model has two advantages. First, instead of pre-defined object categories, more diverse object tags are produced. Second, generating block tags is much faster with the CLIP model than with an object detector, e.g., it is  $40\times$  faster than the Faster-RCNN (ResNeXt152) model.

### 3.1.2 First-order Text Prompt Generation

For the input image of each training pair, Section 3.1.1 already generate the object tags and positions which allows us to design a simple text prompt as follows:

*"The block [P] has a [O]."*

Here,  $P \in 1, \dots, N^2$  denotes the index of the selected block, and  $O$  denotes the object tag generated for that block. If there are multiple objects in a block, we randomly select one object tag for each prompt. This way, each prompt contains object position and text, which can help the model better understand the relationships between objects and text. More prompt design is explored in Section 4.4.

### 3.1.3 Second-order Text Prompt Generation

The first-order relations primarily focus on identifying the position of individual objects within an image. While this approach provides a foundation for understanding object placement, it falls short in capturing the more complex, challenging relationships that exist between objects. For instance, in the example shown in Figure 3, a *flag* is positioned on top of a *girl*, illustrating a higher-order relationship that a first-order prompt would not capture.

To address this limitation and further enhance the learning process, we propose exploring second-order relations by incorporating more sophisticated prompts. As below:

*"The block [P] has a [O] and a [O<sub>2</sub>] in [R] of this block."*

In this format,  $O_2$  is another randomly selected object, and  $R$  represents the relative position, with  $R \in \{\text{top}, \text{bottom}, \text{left}, \text{right}\}$ . By employing this structure, the model is challenged to not only recognize individual objects but also understand their interrelationships and relative positions within the image.

We have named this method *PTP2R*, with  $2R$  representing second-order relations. This approach requires the

1. <https://huggingface.co/openai/clip-vit-base-patch16>  
2. NLTK: (<https://github.com/nltk/nltk>)

TABLE 1: Statistics of the pre-training datasets.

	Dataset	# Images	# Captions	# BBox
Small	COCO	0.11M	0.55M	-
	Visual Genome	0.10M	-	-
	SBU	0.86M	0.86M	-
	CC-3M	2.8M	2.8M	2.69M
Base	4M	4.0M	5.1M	2.69M
	CC-12M	10.2M	10.2M	7M
Large	DataComp-1B	1.17B	1.17B	-
	MMC4	324M	324M	-
	LAION400M	375.3M	375.3M	-

model to conduct reasoning over all objects in the image, allowing it to develop a more comprehensive understanding of the visual scene. By delving into these higher-order relationships, we aim to improve the model's ability to recognize and interpret complex object interactions, ultimately enhancing its performance in VLP tasks.

### 3.2 Pre-training with PTP

In this work, we integrate our *PTP* into mainstream VLP frameworks, leading to *PTP-ViLT* [22], *PTP-CLIP* [41] and *PTP-BLIP* [27]. Following receipt of the *PTP*, we have two options for training these models:

**Integrate into existing tasks.** The simplest method for using text prompt is to change the text input. As shown in Figure 3, the prompted text and original caption were simply padded together. Formally, the input caption  $x$  of our method is represented as:

$$x = [w, q], \quad (2)$$

where  $w$  is text and  $q$  is our generated text prompt. Then we train the VLP models end-to-end with conventional objectives. Following [22], [27], [41], we employ Language Modeling (LM) loss, Image-text Matching (ITM), and Image-text Contrastive (ITC) loss for our *PTP-BLIP*; we use ITM and Masked Language Modeling (MLM) loss to train our *PTP-ViLT*; we only use ITC loss to train our *PTP-CLIP*. We use this method as default for all experiments because of its good performance.

**As a new pretext task.** Alternatively, we explore the position prediction as an additional language modeling task. Formally, if  $D$  is the pretraining data and  $y_1, \dots, y_T$  is a training token sequence of our generated text prompt  $q$ , then at the timestep  $t$ , we devise our model to predict a probability distribution  $p(t) = p(\cdot | y_1, \dots, y_{t-1})$ . Then we regressively try to maximize the probability of being the correct token. The object prediction loss is computed as:

$$\mathcal{L}_{\text{PTP}}(\theta) = -\mathbb{E}_{y \sim D} \left[ \sum_{t=1}^T \log P_\theta(y_t | y_{<t}) \right], \quad (3)$$

where  $\theta$  is the trainable parameters of the model. In this way, the model is asked to predict *which block P has objects and what object O is in this block*.

**Discussion.** Notably, our method does not need to modify the base network and can be applied to any VLP models without bells and whistles. The model is designed to learn position information from raw-pixel image. Note that

**TABLE 2: Comparison with state-of-the-art image captioning methods on NoCaps and COCO Caption.** C: CIDEr, S: SPICE, B@4: BLEU@4. Notice that VinVL $\ddagger$  and LEMON $\ddagger$  require high resolution ( $800 \times 1333$ ) input images.

Method	#Images	Param.	NoCaps validation						COCO Caption		
			in-domain		near-domain		out-domain		Overall	Karpathy test	
			CIDEr	SPICE	CIDEr	SPICE	CIDEr	SPICE	CIDEr	SPICE	B@4
OSCAR [31]	4M	155M	79.6	12.3	66.1	11.5	45.3	9.7	80.9	11.3	37.4
VinVL $\ddagger$ [63]	5.7M	347M	103.1	14.2	96.1	13.8	88.3	12.1	95.5	13.5	38.2
BLIP $\ddagger$ [27]	4M	220M	106.5	14.4	99.3	13.6	95.6	13.0	98.8	14.2	37.0
PTP-BLIP	4M	220M	$108.3 \pm 1.8$	$14.9 \pm 0.5$	$105.0 \pm 5.7$	$14.2 \pm 0.6$	$105.6 \pm 10.0$	$14.2 \pm 1.2$	$106.0 \pm 8.3$	$14.7 \pm 0.5$	$38.6 \pm 1.6$
PTP2R-BLIP	4M	220M	$109.4 \pm 2.9$	$14.9 \pm 0.5$	$105.2 \pm 5.9$	$14.4 \pm 0.8$	$105.6 \pm 10.0$	$14.3 \pm 1.3$	$106.6 \pm 8.9$	$14.9 \pm 0.7$	$38.8 \pm 1.8$
Enc-Dec [8]	15M	—	92.6	12.5	88.3	12.1	94.5	11.9	90.2	12.1	—
BLIP [27]	14M	220M	111.3	15.1	104.5	14.4	102.4	13.7	105.1	14.4	38.6
PTP-BLIP	14M	220M	$112.8 \pm 1.5$	$15.2 \pm 0.1$	$107.3 \pm 2.8$	$14.9 \pm 0.5$	$108.1 \pm 6.7$	$14.3 \pm 0.6$	$106.3 \pm 0.8$	$14.7 \pm 0.3$	$40.1 \pm 1.5$
PTP2R-BLIP	14M	220M	$113.1 \pm 1.8$	$15.2 \pm 0.1$	$107.5 \pm 3.0$	$14.9 \pm 0.5$	$108.6 \pm 7.2$	$14.4 \pm 0.7$	$106.5 \pm 1.0$	$14.8 \pm 0.4$	$40.3 \pm 1.7$
SimVLM <sub>H</sub> [55]	1.8B	1.2B	113.7	—	110.9	—	115.2	—	112.2	—	40.6
LEMON <sub>H</sub> $\ddagger$ [16]	200M	675M	118.0	15.4	116.3	15.1	120.2	14.5	117.3	15.0	42.6

TABLE 3: Hyper-parameters for BLIP baseline.

Task	VQA	Retrieval	NLVR2	Captioning
AdamW with Weight Decay				
Optimizer	1.0			
Gradient clip	Cosine Schedule Decaying to Zero			
LR decay schedule	0.05			
Weight decay rate	2,5	2,5	2,5	2,5
RandAugment	10	6	5	5
Train epochs	64	24	128	16
Train batch size	2e-5	1e-5	3e-5	1e-5

only during the pre-training stage, we would require the object's position information; yet on downstream tasks, we evaluate model in normal end-to-end ways without object information to get rid of the heavy object feature extraction.

## 4 EXPERIMENTS

In this section, we empirically evaluate PTP on multiple downstream tasks and present a comprehensive study.

### 4.1 Experimental Settings

We first describe the pre-training experimental conditions, including the datasets, training configurations, evaluation procedures, and baseline models used in our studies.

**Statistics of the Pre-training Datasets.** In this work, we explore beginning from both 4M (Small) and 14M (Base) setting. As in earlier studies [31], [63], we begin by using a 4M setup made up of four popular pre-training datasets (COCO [32], VG [23], SBU [38] and CC3M [45]). The 14M setting is a combination of 4M setting and CC-12M. Following recent work [27], we also explore 14M setting, which includes additional CC12M [8] (actually only 10M image urls available) dataset besides 4M datasets. We report the data statistics in Table 1. As the URLs for the CC3M and CC12M datasets are derived from the Internet, and given that a portion of them are no longer valid, we have downloaded a total of 2.8 million data instances for CC3M and 10.2 million data instances for CC12M, respectively. It should be noted that the BLIP baseline employs 3 million data instances for CC3M, slightly more than our model. In terms of the number of images that contain bounding boxes, there are 2.69 million such images for CC3M and 7 million for CC12M. These bounding boxes are used in our PTP. For quick evaluation, we pre-train the BLIP model for 50K steps rather than the 200K in [22], [62].

TABLE 4: Hyper-parameters for ViLT baseline.

Task	Dataset	VQA	Retrieval	NLVR2
		VQAV2	COCO	F30K
Optimizer	AdamW with Weight Decay			
Gradient clip	1.0			
LR decay schedule	Cosine Schedule Decaying to Zero			
RandAugment	2,9			
Weight decay rate	0.05			
Train epochs	10	10	5	10
Train batch size	256	256	256	128
LR	1e-4	3e-4	1e-4	1e-4
Warm-up steps	1500	2500	1000	500

Additionally, to assess the effectiveness of our method at a Large scale, we conducted experiments on extensive datasets, encompassing billions of data points, which include DataComp-1B [14], MMC4 [68], and LAION400M [44].

**Training Settings.** Our models are implemented in PyTorch [39] and pre-trained on 8 NVIDIA A100 GPUs. To ensure a fair comparison, we adopt the optimizer and training hyperparameters from the original implementation in the baseline works. Additionally, we investigate the use of RandAugment [11] for data augmentation and utilize all of the original policies, except for color inversion, as color information is crucial for our tasks. Furthermore, we apply affine transformations to augment the bounding boxes in a manner similar to that used for image rotation. We take random image crops of  $224 \times 224$  during pre-training resolution, and increase the image resolution to  $384 \times 384$  for downstream task finetuning.

**Hyper-parameters for Downstream Tasks.** 3. The final decoder outputs of the encoder-decoder model BLIP can be utilized for both multimodal understanding and generation. Therefore, we evaluate its performance on popular vision-language benchmarks, employing the same setup as introduced in the BLIP paper [27]. Specifically, we use the AdamW optimizer for all tasks and train the retrieval task for only 6 epochs to increase efficiency. We believe that increasing the number of epochs would yield better results.

In the case of the ViLT baseline, we primarily focus on three tasks: vision-question answering, image-text retrieval, and natural language visual reasoning. The hyper-parameters for ViLT on these downstream tasks are reported in Table 4. Finally, for the CLIP baseline, we use the same hyper-parameter settings as those employed in BLIP.

**Baselines.** We evaluate three variants of pre-training

**TABLE 5: Results of zero-shot image-text retrieval on Flickr30K and MSCOCO datasets.** The methods that utilize significantly larger models or train on larger corpora have been grayed out. The symbol  $\dagger$  denotes models that were implemented and trained on the same dataset, as the original datasets were either inaccessible or not trained on these splits. The Avg metric represents the mean of all image-to-text and text-to-image recalls.

Method	#Images	Param.	MSCOCO (5K test set)									Flickr30K (1K test set)								
			Image $\rightarrow$ Text			Text $\rightarrow$ Image			Avg	Image $\rightarrow$ Text			Text $\rightarrow$ Image			Avg				
			R@1	R@5	R@10	R@1	R@5	R@10		R@1	R@5	R@10	R@1	R@5	R@10					
Unicoder-VL [26]	4M	170M	—	—	—	—	—	—	64.3	85.8	92.3	48.4	76.0	85.2	75.3					
ImageBERT [40]	4M	170M	44.0	71.2	80.4	32.3	59.0	70.2	59.5	70.7	90.2	94.0	54.3	79.6	87.5	79.4				
ViLT [22]	4M	87M	41.3	79.9	87.9	37.3	67.4	79.0	65.5	69.7	91.0	96.0	53.4	80.7	88.8	79.9				
PTP-ViLT	4M	87M	55.1	82.3	89.1	43.5	70.2	81.2	70.2 <sup>+4.7</sup>	74.5	93.7	96.5	60.3	85.5	90.4	83.5 <sup>+3.6</sup>				
PTP2R-ViLT	4M	87M	55.8	82.7	89.6	44.2	71.3	81.8	70.9 <sup>+5.4</sup>	75.3	94.1	96.8	60.6	85.8	90.5	83.9 <sup>+4.0</sup>				
BLIP + [27]	4M	220M	57.4	81.1	88.7	41.4	66.0	75.3	68.3	76.0	92.8	96.1	58.4	80.0	86.7	81.7				
PTP-BLIP	4M	220M	72.3	91.8	95.7	49.5	75.9	84.2	77.3 <sup>+9.0</sup>	86.4	97.6	98.9	67.0	87.6	92.6	88.4 <sup>+6.7</sup>				
PTP2R-BLIP	4M	220M	72.9	92.3	96.1	49.9	76.2	84.4	77.5 <sup>+9.2</sup>	86.9	97.8	99.0	67.3	87.8	92.7	88.6 <sup>+6.9</sup>				
BLIP + [27]	14M	220M	65.5	86.4	92.3	48.4	73.3	83.5	74.9	83.3	95.8	98.0	70.4	88.3	93.1	88.2				
PTP-BLIP	14M	220M	73.2	92.4	96.1	53.6	79.2	87.1	78.6 <sup>+3.7</sup>	87.1	98.4	99.3	73.1	91.0	94.8	90.3 <sup>+2.1</sup>				
PTP2R-BLIP	14M	220M	73.6	92.6	96.5	53.5	79.0	87.1	78.8 <sup>+3.9</sup>	87.4	98.5	99.3	73.2	91.1	94.8	90.4 <sup>+2.2</sup>				
CLIP [41]	300M	173M	58.4	81.5	88.1	37.8	62.4	72.2	66.7	88.0	98.7	99.4	68.7	90.6	95.2	90.1				
ALIGN [20]	1.8B	820M	58.6	83.0	89.7	45.6	69.8	78.6	70.9	88.6	98.7	99.7	75.7	93.8	96.8	92.2				
FILIP [57]	340M	787M	61.3	84.3	90.4	45.9	70.6	79.3	72.0	89.8	99.2	99.8	75.0	93.4	96.3	92.3				
Flamingo [2]	2.1B	80B	65.9	87.3	92.9	48.0	73.3	82.1	74.9	89.3	98.8	99.7	79.5	95.3	97.9	93.4				
CoCa [32]	3B	2.1B	66.3	86.2	91.8	51.2	74.2	82.0	75.3	92.5	99.5	99.9	80.4	95.7	97.7	94.3				

**TABLE 6: Finetuning results of image-to-text retrieval and text-to-image retrieval on COCO and Flickr30K.** Notice that UNITER [10], OSCAR [31] and VinVL [63] all use bounding box and object feature.

Method	#Images	Param.	MSCOCO (5K test set)									Flickr30K (1K test set)								
			Image $\rightarrow$ Text			Text $\rightarrow$ Image			Avg	Image $\rightarrow$ Text			Text $\rightarrow$ Image			Avg				
			R@1	R@5	R@10	R@1	R@5	R@10		R@1	R@5	R@10	R@1	R@5	R@10					
UNITER [10]	4M	155M	65.7	88.6	93.8	52.9	79.9	88.0	78.2	87.3	98.0	99.2	75.6	94.1	96.8	91.8				
OSCAR [31]	4M	155M	70.0	91.1	95.5	54.0	80.8	88.5	—	—	—	—	—	—	—	—				
VinVL [63]	4M	157M	74.6	92.6	96.3	58.1	83.2	90.1	82.5	—	—	—	—	—	—	—				
ViLT [22]	4M	87M	61.8	86.2	92.6	41.3	72.0	82.5	72.7	81.4	95.6	97.6	61.9	86.8	92.8	86.0				
PTP-ViLT	4M	87M	67.1	90.5	94.3	45.3	79.1	88.4	77.5 <sup>+4.8</sup>	85.2	96.9	98.5	68.8	91.4	95.3	89.4 <sup>+3.4</sup>				
PTP2R-ViLT	4M	87M	68.3	91.2	94.7	45.6	80.6	89.2	78.3 <sup>+5.5</sup>	85.7	97.0	98.2	69.3	91.7	95.6	89.6 <sup>+3.6</sup>				
BLIP $\dagger$ [27]	4M	220M	75.2	93.3	96.3	57.4	82.1	89.5	82.3	94.0	99.1	99.7	82.5	96.4	98.2	95.0				
PTP-BLIP	4M	220M	83.7	97.0	98.7	68.1	89.4	94.2	88.5 <sup>+6.2</sup>	96.1	99.8	100.0	84.2	96.6	98.6	95.9 <sup>+0.9</sup>				
PTP2R-BLIP	4M	220M	84.1	97.2	98.8	69.2	89.9	94.5	89.0 <sup>+6.7</sup>	96.3	99.9	100.0	84.4	96.8	98.9	96.1 <sup>+1.1</sup>				
ALBEF [28]	14M	210M	77.6	94.3	97.2	60.7	84.3	90.5	84.1	95.9	99.8	100.0	85.6	97.5	98.9	96.3				
BLIP [27]	14M	220M	80.6	95.2	97.6	63.1	85.3	91.1	85.5	96.6	99.8	100.0	87.2	97.5	98.8	96.7				
PTP-BLIP	14M	220M	84.2	97.3	98.8	68.8	89.5	94.2	88.8 <sup>+3.3</sup>	97.0	99.9	100.0	87.7	98.2	99.3	97.0 <sup>+0.3</sup>				
PTP2R-BLIP	14M	220M	84.6	97.5	98.8	69.5	90.1	94.8	89.3 <sup>+3.8</sup>	97.1	99.9	100.0	87.9	98.3	99.4	97.1 <sup>+0.4</sup>				
ALIGN [20]	1.8B	820M	77.0	93.5	96.9	59.9	83.3	89.8	83.4	95.3	99.8	100.0	84.9	97.4	98.6	96.0				
FILIP [57]	340M	787M	78.9	94.4	97.4	61.2	84.3	90.6	84.5	96.6	100.0	100.0	87.1	97.7	99.1	96.8				
Florence [61]	900M	893M	81.8	95.2	—	63.2	85.7	—	—	97.2	99.9	—	87.9	98.1	—	—				

frameworks, including one-stream ViLT [22], dual-encoder CLIP [41], and fusion-encoder BLIP [27], for their superior performance. For fair comparisons, we adopt the ViT-B/16 [13] as base vision encoder and use same dataset.

## 4.2 Main Results

In this section, we integrated our *PTP* into existing networks and compare to existing VLP methods on a wide range of vision-language downstream tasks. Include five image-text tasks and two video-text tasks.

**Image Captioning.** This task asks the model to describe the input image. We consider two datasets for image captioning: No-Caps [1] and COCO Captioning [32], [49], both evaluated using the model finetuned on COCO with the LM loss. Similar to BLIP, we start each caption with the phrase “a picture of,” which yields marginally better results. Since image captioning needs a decoding head, dual stream CLIP and one-stream ViLT architectures cannot test this task directly due to missing decoding head. We do not pre-

train with COCO to avoid information leakage. For No-Caps dataset, following BLIP, we adopts a zero-shot setting.

As shown in Table 2, related works utilizing a comparable quantity of pre-training data perform significantly worse than *PTP-BLIP*. The results of our method are closed to the VinVL [63] with fewer training samples and smaller image. Finally, with 14M setting, our method leads to close result with LEMON, which trained on billions data and requires two times higher resolution image.

**Image-Text Retrieval.** We evaluate *PTP* for both image-to-text retrieval (TR) and text-to-image retrieval (IR) on COCO and Flickr30K benchmarks. For *PTP-BLIP*, following original implementation, we adopt an additional re-ranking strategy. Specifically, we first select  $k$  candidates based on the image-text feature similarity, and then rerank the selected candidates based on their pairwise ITM scores. We set  $k = 256$  for COCO and  $k = 128$  for Flickr30K.

We first report zero-shot retrieval result on both image-to-text and text-to-image setting in Table 5. Mainstream methods in VLP, e.g. BUTD [3], OSCAR [31] and

**TABLE 7: Comparison with State-of-the-Art Methods on VQA and NLVR<sup>2</sup>:** It is worth noting that VinVL [63] employs a larger vision backbone and object features extracted using Faster R-CNN. Additionally, ALBEF [28] performs an extra pre-training step for NLVR<sup>2</sup>, while BeIT-3 [54] uses an additional 160GB text corpus.

Method	#Images	Para.	VQA		NLVR <sup>2</sup>	
			test-dev	test-std	dev	test-P
UNITER [10]	4M	155M	72.70	72.91	77.18	77.85
OSCAR [31]	4M	155M	73.16	73.44	78.07	78.36
UNIMO [30]	5.6M	307M	75.06	75.27	-	-
VinVL <sub>L</sub> [63]	5.6M	347M	<b>76.52</b>	<b>76.60</b>	<b>82.67</b>	<b>83.98</b>
ViLT [22]	4M	87M	70.33	-	74.41	74.57
PTP-ViLT	4M	87M	72.13	74.36	76.52	77.83
PTP2R-ViLT	4M	87M	73.44	76.16	77.31	78.50
BLIP † [27]	4M	220M	73.92	74.13	77.52	77.63
PTP-BLIP	4M	220M	75.47	75.88	80.73	81.24
PTP2R-BLIP	4M	220M	75.81	76.03	81.22	81.40
ALBEF [28]	14M	210M	75.84	76.04	82.55	83.14
BLIP [27]	14M	220M	<b>77.54</b>	<b>77.62</b>	<b>82.67</b>	<b>82.30</b>
PTP-BLIP	14M	220M	78.44	78.33	84.55	83.17
PTP2R-BLIP	14M	220M	<b>78.85</b>	<b>78.66</b>	<b>84.92</b>	<b>83.41</b>
SimVLM [55]	1.8B	1.2B	<b>77.87</b>	78.14	81.72	81.77
GIT [52]	0.8B	0.7B	-	78.81	-	-
Beit-3 [54]	35M+	1.9B	84.19	84.03	91.51	92.58

UNITER [10], often use object detector, e.g. Faster-RCNN that is also pretrained on Visual Genome. Moreover, the compared methods in Table 1~4 mostly use object detector. E.g. we use Faster-RCNN adopted by VinVL, but our model (220M) has better performance than VinVL (347M). So the comparison is fair. We find PTP significantly improves baselines on all metrics. For example, for ViLT [22] baseline, PTP leads to 13.8 % absolute improvement (from 41.3 % to 55.1 %) over Recall@1 of image to text retrieval on MSCOCO. In addition, our PTP-BLIP even outperforms CoCa [59] on most recalls of MSCOCO with much less data.

A summary comparison of the fine-tuned settings between different models is presented in Table 6. It is observed that: (1) PTP outperforms the BLIP and ViLT baselines by a large margin in both datasets. For instance, PTP-ViLT achieves an impressive 5.3% improvement on R@1 of TR in MSCOCO. (2) With the strong BLIP baseline, PTP-BLIP achieves state-of-the-art performance at the same scale. Notably, the training cost of PTP remains the same as that of the BLIP baseline, since we train PTP with the same settings as the baseline without increasing the maximum input text token. Moreover, we can even reduce the gap between the 4M setting and ALBEF [28] (14M setting) with a similar dual stream with fusion encoder architecture.

From all these results above, we point out UNITER [10], OSCAR [31], VinVL [63], ImageBERT [40] all use faster-rcnn as we used. However, our PTP leads to much better results than these related works. Besides, we only use object detector in pre-training stage. This indicates *object detector is not the secret for success and how to leverage the position information is essential important for VLP models*.

**Visual Question Answering:** In the context of visual question answering, VQA [4] requires a model to predict an answer based on an image and a corresponding question. For PTP-ViLT, we approach VQA as a multi-answer classification task. On the other hand, for PTP-BLIP, we follow the approach used in [27], [28] and consider VQA as an

**TABLE 8: Comparisons with related works for zero-shot text-to-video retrieval on the 1k test split of the MSRVTT.**

Method	R1↑	R5↑	R10↑	MdR↓
ActBERT [67]	8.6	23.4	33.1	36.0
MIL-NCE [37]	9.9	24.0	32.4	29.5
Frozen-in-time [6]	18.7	39.5	51.6	10.0
OA-Trans [53]	23.4	47.5	55.6	8.0
ViLT† [22]	22.6	46.9	53.2	8.0
PTP-ViLT	27.9	52.5	56.3	7.0
PTP2R-ViLT	<b>28.4</b>	<b>53.1</b>	<b>56.6</b>	<b>7.0</b>

**TABLE 9: Comparison with Related Works on Visual Grounding:** It is worth noting that these related methods were pre-trained on object feature.

Method	Ref-COCO			Ref-COCO+		
	val	testA	testB	val	testA	testB
VL-BERT <sub>L</sub> [46]	-	-	-	72.59	78.57	62.30
ViLBERT [36]	-	-	-	72.34	78.52	62.61
UNITER [10]	81.41	87.04	74.17	75.90	81.45	66.70
VILLA [15]	<b>82.39</b>	<b>87.48</b>	<b>74.84</b>	76.17	81.54	66.84
BLIP† [27]	77.45	83.32	71.33	74.13	79.37	64.41
PTP-BLIP	79.95	85.33	72.46	74.49	80.13	66.34
PTP2R-BLIP	81.83	86.44	74.30	<b>76.65</b>	<b>82.14</b>	<b>67.38</b>

answer generation task to facilitate open-vocabulary VQA for improved performance.

The performance results are reported in Table 7. Our proposed model, PTP, shows a significant improvement over the ViLT baseline, with a gain of 1.8% on both splits. Furthermore, with the 14M setting, PTP-BLIP outperforms SimVLM [55], which utilizes a ViT-Large based vision backbone and 1.8 billion training samples.

**Visual Reasoning.** Natural Language Visual Reasoning (NLVR<sup>2</sup>) [47] task is a binary classification task given triplets of two images and a question in natural language. This task relies on position information heavily. As shown in Table 7, SimVLM [55] is outperformed by PTP-BLIP, which has a reasonable model size and was pretrained on fewer instances. Meanwhile, our method is also closed to VinVL<sub>large</sub> model that adopt larger model and use object feature from strong object detector instead of raw-pixel image as input.

**Visual Grounding.** We follow the approach used in ViLBERT [36] to evaluate our model's visual grounding capabilities in the Referring Expression task, which involves using text to locate image regions. Table 9 presents the results obtained for the 4M setting, which demonstrate the strong grounding ability of our proposed PTP models. Notably, we observe that PTP outperforms several related works that rely on object features extracted using the heavy Faster R-CNN architecture. Moreover, PTP achieves significant improvements over BLIP.

**Video-Text Retrieval.** In this experiment, we test the generalization ability of our method to video-language tasks. Specifically, we perform zero-shot transfer to text-to-video retrieval in Table 8, where we directly evaluate the models trained on COCO-retrieval. To process video input, we simply sample 8 frames uniformly from each video and average the frame features into a single sequence.

Our method outperforms OA-Trans [53], which is a retrieval-focused method, demonstrating the generality ca-

TABLE 10: Comparing with the Open-Flamingo baseline across various language model scales.

Method	Shots	Captioning (CIDEr)		VQA				Classification HatefulMemes	Average
		COCO	FLICKR	ok-vqa	textvqa	vizwiz	vqav2		
Open-flamingo(3B) [5]	0	74.9	52.3	28.2	24.2	23.7	44.6	51.2	42.7
	4	77.3	57.2	30.3	27.0	27.0	45.8	50.6	45.0
	32	93.0	61.1	31.0	28.3	44.1	47.0	50.2	50.7
PTP-Open-flamingo(3B)	0	78.9	53.3	27.4	25.1	29.3	44.2	51.5	44.2
	4	88.4	61.2	28.4	25.7	31.2	46.3	54.0	47.9
	32	94.4	63.5	30.8	25.6	42.3	43.9	52.2	50.3
Open-flamingo(9B) [5]	0	79.5	59.5	28.2	24.2	23.7	44.6	51.6	44.5
	4	89.0	65.8	40.1	28.2	34.1	54.8	54.0	52.3
	32	99.5	61.3	42.4	23.8	44.0	53.3	53.8	54.1
PTP-Open-flamingo(9B)	0	80.2	60.4	27.9	22.5	25.3	46.6	51.5	44.9
	4	91.7	67.2	39.3	28.4	39.0	55.5	53.2	53.4
	32	95.4	68.8	41.8	27.6	42.3	54.9	53.7	55.0

TABLE 11: Comparisons with state-of-the-art methods for video question answering on MSVD, TVQA and TGIF-Frame datasets. Notably, the training samples used in AllInOne [50] are 20 times larger than those used in our method.

Method	MSVD Test	TVQA Val	TGIF-FrameQA
ClipBERT [24]	-	-	59.4
AllInOne [50]	46.5	69.8	62.5
ViLT [22]	45.7	70.4	65.4
PTP-ViLT	48.8	73.4	68.7
PTP2R-ViLT	49.1	73.9	68.7
BLIP† [27]	47.1	71.3	66.4
PTP-BLIP	50.3	72.4	70.2
PTP2R-BLIP	50.8	72.7	70.6

pability of PTP. Also, note that this simple approach not well explored temporal information.

**Video Question Answering.** We report the video question answering results in Table 11. Following All-in-one [50], we explore three widely used benchmarks: MSVD-QA [56], TVQA [25] and TGIF [19]. TGIF FrameQA and MSVD-AQ are open-ended VQA tasks and TVQA is a multiple-choice VQA task. Similar to video-text retrieval, we sample 8 frames for each video. We observe that PTP-BLIP performs well in both multiple-choice and open-ended settings. For example, out PTP2R-BLIP outperform All-in-one [50] by 4.3% on MSVD-QA and 8.3% on TGIF-FrameQA.

### 4.3 Scale-up Experiments

#### 4.3.1 Data Scale Up

In this section, we expand our research methodology to encompass the extensive DataComp-1B dataset [14], which is recognized for its exceptional data quality when compared to the LAION dataset [44]. The DataComp-1B dataset originally comprised 1.4 billion samples; however, it is worth noting that certain URLs have since become unavailable, resulting in the download of 1.17 billion data samples.

Given the substantial computational resources demanded by the BLIP and ViLT architectures, our primary evaluation is centered around the CLIP (ViT B-16) model, chosen for its efficiency in training. Leveraging the inherent simplicity of PTP, which obviates the need for intricate hyperparameter selection and operates seamlessly at the data level, we enable a comprehensive comparison at this

TABLE 12: PTP pre-training with one billion data. We report the image to text Recall1 results on MSCOCO and Flickr30K dataset.

Method	Dataset	MSCOCO R@1	Flickr30K R@1	ImageNet Acc(%)	VTAB Acc(%)
CLIP† [41]	DataComp-1B	68.9	78.2	66.2	53.4
PTP-CLIP	DataComp-1B	71.4	81.2	67.1	54.6

considerable scale. Moreover, in alignment with the Datacompr evaluation framework, we subject the model to a comprehensive array of image classification benchmarks, employing a zero-shot evaluation methodology.

Recognizing the significant time investment associated with deploying an object detector model, we expedite the training process by utilizing data generated by the CLIP model, as described in the methods section. We summarize the comparative performance between the original CLIP model and PTP-CLIP in Table 12. Notably, upon closer scrutiny of the enhanced version, we find that the experimental results remain promising even when dealing with data at the billion-level scale.

#### 4.3.2 Language Model Scale Up

In this experiment, we vary the language model size from BERT-base(200M) [12] to LLAMA(7B) model [48] and analyze the impact of Large Language Models (LLMs) in multi-modality learning. Specifically, we implement our methodology within the open-flamingo architecture [5]. Staying true to the original implementation, we integrate our approach with OPT1.7B [64] and LLAMA-7B [48] as language models, while utilizing the Open-CLIP [18] ViT-L model as the vision encoder. The total parameter is 3B and 9B.

It's noteworthy that due to unforeseen issues with some image URLs from LAION400M [44] and MMC4 [68] datasets, preventing access for reevaluation, the sample size is marginally smaller than open-flamingo and we report results with our implementation. Given the document nature of MMC4, we exclusively introduce text prompts in the image-text pairs.

The model undergoes rigorous training over a span of 6.3 days for 3B model and 17 days for 9B models, leveraging the computational power of 64 Nvidia V100 GPUs. The comprehensive results are presented in Table 10. Our focus primarily centers around few-shot tasks, providing a

TABLE 13: **The ablation on different architectures under 4M setting.** We report the i2t and t2i results on MSCOCO (5K test set). As we do not used object detector in downstream tasks, PTP is 20 times faster than object-feature based model.

Method	Time	MSCOCO (5K test set)										Flickr30K (1K test set)										
		Image → Text			Text → Image			Image → Text			Text → Image											
		R@1	R@5	R@10	R@1	R@5	R@10	Avg	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	Avg	
ViLT [22]	~15	61.8	86.2	92.6	41.3	72.0	82.5	72.7	81.4	95.6	97.6	61.9	86.8	92.8	86.0	<i>One-stream Models</i>						
	~15	<b>67.1</b>	<b>90.5</b>	<b>94.3</b>	<b>45.3</b>	<b>79.1</b>	<b>88.4</b>	<b>77.5</b> <sub>+4.8</sub>	<b>85.2</b>	<b>96.9</b>	<b>98.5</b>	<b>68.8</b>	<b>91.4</b>	<b>95.3</b>	<b>89.4</b> <sub>+3.4</sub>							
CLIP† [41]	~27	64.9	83.2	90.1	50.4	76.3	84.7	74.9	77.5	92.1	94.5	58.6	81.4	88.9	82.2	<i>Dual-stream Models</i>						
	~27	<b>68.3</b>	<b>86.4</b>	<b>92.7</b>	<b>54.1</b>	<b>80.1</b>	<b>86.8</b>	<b>78.1</b> <sub>+3.2</sub>	<b>83.5</b>	<b>94.1</b>	<b>96.3</b>	<b>63.4</b>	<b>88.5</b>	<b>91.7</b>	<b>86.3</b> <sub>+4.1</sub>							
BLIP † [27]	~33	75.2	93.3	96.3	57.4	82.1	89.5	82.3	94.0	99.1	99.7	82.5	96.4	98.2	95.0	<i>Dual-stream + Fusion encoder Models</i>						
	~33	<b>83.7</b>	<b>97.0</b>	<b>98.7</b>	<b>68.1</b>	<b>89.4</b>	<b>94.2</b>	<b>88.5</b> <sub>+6.2</sub>	<b>96.1</b>	<b>99.8</b>	<b>100.0</b>	<b>84.2</b>	<b>96.6</b>	<b>98.6</b>	<b>95.9</b> <sub>+0.9</sub>							
VinVL [63]	~650	74.9	92.6	96.3	58.1	83.2	90.1	82.5	-	-	-	-	-	-	-	<i>Object-feature Based Models</i>						

TABLE 14: **Text prompt vs. additional pretext head.** The last column is COCO captioning task.

Method	COCO TR@1	F30K TR@1	NLVR Acc(%)	Captioning CIDEr
Baseline	70.6	53.4	76.1	121.2
Pretext	72.3 ( <b>1.7↑</b> )	54.7 ( <b>2.3↑</b> )	76.9 ( <b>0.8↑</b> )	123.5 ( <b>2.3↑</b> )
Prompt	<b>73.2</b> ( <b>2.6↑</b> )	<b>55.4</b> ( <b>2.0↑</b> )	<b>77.9</b> ( <b>1.8↑</b> )	<b>127.2</b> ( <b>6.0↑</b> )

nuanced understanding of the model’s representation capabilities, as reflected in the results showcased in the aforementioned table. It’s essential to highlight the importance of data prompt style even for large models, as evidenced by the downstream tasks.

#### 4.4 Ablation & Design Choices

In this section, we first evaluate our method on retrieval task over three well-known baselines under 4M setting for comparison. Then we train a BLIP model on CC3M as baseline and perform various ablations.

**Exploration of Diverse Architectures.** In this research, we conduct experiments utilizing three distinct baseline models, specifically ViLT, CLIP, and BLIP, to examine the influence of PTP on a range of performance indicators. The outcomes of these experiments are detailed in Table 13, demonstrating the performance on both the COCO 5K test set and the Flickr30K 1K test set. The data acquired from the baseline models indicate that our proposed approach, *PTP*, substantially improves image-to-text (i2t) and text-to-image (t2i) performance, underscoring its adaptability and suitability for a variety of visual-language tasks.

Furthermore, we assess the execution time of our model relative to the baseline models. As we do not employ an object detector or prompts for downstream tasks, the computational overhead aligns with that of the baseline models. Impressively, *PTP* is determined to be 20 times more rapid than VinVL [63], which depends on object features. Importantly, in spite of the considerable decrease in execution time, our model produces results of a quality comparable to VinVL, highlighting its efficiency and efficacy.

**Comparing Text Prompts and Additional Pretext Tasks.** This research examines the effects of incorporating *PTP* as a

TABLE 15: **Case study of text prompt on image-text retrieval.** We use the **O** to represent objects and **P** to represent positions. The symbol  $\theta$  represents learnable parameter.

Prompt	TR@1	IR@1
Baseline	70.6	53.4
The [O] is in the block [P].	72.7 ( <b>2.1↑</b> )	54.1 ( <b>0.7↑</b> )
The block [P] looks like [O].	73.3 ( <b>2.7↑</b> )	53.9 ( <b>0.5↑</b> )
The [O] is in which block? In [P].	72.3 ( <b>1.7↑</b> )	54.9 ( <b>1.5↑</b> )
The [O] is located in block [P].	72.3 ( <b>1.7↑</b> )	54.2 ( <b>0.8↑</b> )
(X1, Y1, W, H) has a [O].	72.5 ( <b>1.9↑</b> )	54.3 ( <b>0.9↑</b> )
The block in [NP] has a [O].	73.0 ( <b>2.4↑</b> )	55.1 ( <b>1.7↑</b> )
$\theta_1 \theta_2$ [P] $\theta_3 \theta_4$ [O].	73.1 ( <b>2.5↑</b> )	55.2 ( <b>1.8↑</b> )
The block [P] has a [O].	73.2 ( <b>2.6↑</b> )	55.4 ( <b>2.0↑</b> )
Mixed	72.3 ( <b>1.7↑</b> )	54.7 ( <b>1.2↑</b> )

supplementary pretext task during the pre-training phase of vision-language models. By implementing this strategy, the pretext task does not conflict with other pre-training objectives, such as ITM and ITC, even though it may increase computational expenses. Conversely, the prompt design modifies the textual input, influencing all pre-training goals.

The results are presented in Table 14. We notice that both Pretext and Prompt approaches enhance the baseline performance across all four tasks. Nonetheless, the prompt method is distinctly more advantageous than the pretext approach, particularly for COCO captioning CIDEr scores (127.2 vs 123.5). Consequently, we employ the prompt design by default in this study, owing to its superior efficiency.

**Exploring Various Text Prompts.** In this experiment, we investigate six distinct types of prompts: i. The [O] is in block [P]. ii. The block [P] resembles [O]. iii. In which block is the [O]? In [P]. iv. The [O] is situated in block [P]. v.  $(X_1, Y_1, W, H)$  contains a [O].  $(X_1, Y_1)$  represents the top-left point, while  $W, H$  denote the width and height of the bounding box. vi. The block [P] features a [O]. vii. The block [NP] includes a [O]. NP refers to the usage of nouns to describe block positions, e.g., from upper left to bottom right. The results are reported in Table 15 and we observe:

A precise position does not yield superior results compared to a block, possibly because precise positioning is challenging to learn. Furthermore, we find that utilizing block IDs (e.g., 0) or nouns (e.g., upper left) produces similar

TABLE 16: Other variations of text prompt. [O] is short for object and [P] is short for position.

Prompt	Multipy Position	Multipy Tags	Prompt	COCO Retrieval TR@1	IR@1	NLVR Acc	COCO Captioning CiDER
Baseline	-	-	-	70.6	53.4	76.0	122.6
The object in region [P] looks like [O]. The block [P] has objects [O <sub>1</sub> ], [O <sub>2</sub> ], [O <sub>3</sub> ]. The [O] is located in which region? In [P <sub>1</sub> ], [P <sub>2</sub> ] and [P <sub>3</sub> ].		✓	✓ ✓ ✓	72.5 (1.9↑) 71.9 (0.9↑) 70.7 (0.1↑)	54.3 (0.9↑) 54.7 (0.9↑) 53.6 (0.2↑)	77.8 (1.8↑) 76.8 (0.9↑) 77.1 (1.1↑)	127.4 (4.8↑) 124.5 (1.9↑) 125.2 (2.6↑)
ColorPrompt [58]	✓	✓	✓	70.4 (0.2↓)	53.6 (0.2↑)	75.1 (0.9↓)	120.3 (2.3↓)

TABLE 17: COCO text-to-image retrieval results with the second-order relation over objects.

Prompt	R@1	R@5	R@10
The block [P] has a [O]	73.1	91.9	96.0
+, and the block [P <sub>2</sub> ] has a [O <sub>2</sub> ].	71.5	91.5	95.9
+, and a [O <sub>2</sub> ] on the [P <sub>2</sub> ].	72.2	91.5	95.8
+, and a [O <sub>2</sub> ] on the [R] of this block.	73.6	92.3	96.3
The image has a [O] and a [O <sub>2</sub> ] on the [R].	71.4	91.5	95.6

outcomes. Ultimately, we discover that the hybrid version does not generate the best results. We also note that a single-word change can significantly impact performance, a common issue in prompt learning, as observed in GPT-3 [7]. Our work does not primarily focus on addressing this problem. Additionally, Table 15 demonstrates that using prompts to predict the exact position (four coordinates) of an object’s bounding box results in inferior performance compared to predicting the block.

**Investigating Second-order Prompts.** Table 17 delves into the examination of object relation prompts, where R encompasses terms such as “left,” “right,” “top,” and “bottom.” We explore three distinct variations of PTP2R, including the simple repetition of first-order relations, the relative position of the object, and the relative position of the selected block. O<sub>2</sub>/P<sub>2</sub> means different objects/positions. In addition, we also explore the second-order text prompts solely based on objects and positional relationships.

Incorporating relation prompts leads to a noticeable improvement in COCO text-to-image (t2i) retrieval performance. This enhancement may be ascribed to the heightened complexity linked to learning relation prompts, as indicated by the increase in language mask loss from 1.29 to 1.47. The elevated language mask loss implies that the model encounters greater difficulty in grasping the subtleties of object relations, thus pushing it to develop a more refined comprehension of the relationships between objects within the visual domain. Consequently, the model becomes more adept at handling intricate tasks involving object relations, ultimately resulting in its enhanced performance in the COCO t2i retrieval task. We have noted a challenge encountered by the model when attempting to determine the precise location of the first object in cases where positional information is absent from the initial second-order prompt. This issue primarily arises due to the presence of recurring object classes among the top K largest objects.

**Exploring Additional Prompt Designs.** In our approach, an object may span multiple blocks, and each category may encompass multiple objects. To address this, we also predicts all blocks or objects and investigate several other

TABLE 18: The position information is essential for prompt design. Different variations of object prediction prompt design and evaluate on coco retrieval.

Object Tags	Prompt	Position	TR@1	IR@1
-	-	-	70.6	53.4
✓			70.2 (0.4↓)	52.7 (0.7↓)
✓	✓		70.3 (0.3↓)	52.9 (0.5↓)
✓		✓	70.8 (0.3↓)	52.4 (1.0↓)
✓	✓	✓	73.3 (2.7↑)	55.4 (2.0↑)

prompt. The model is trained on CC3M and evaluated on three downstream tasks. Specifically, we explore the following approaches: i. *Multiple Tags*. We note that a block may contain multiple objects in many cases. We attempt to refine the text prompt as *The block [P] has objects [O<sub>1</sub>], [O<sub>2</sub>], and [O<sub>3</sub>]*. It is important to remember that each block contains a varying number of objects. ii. *Multiple Positions*. We create a multiple position setup, considering that one object could appear in several blocks. We refine the prompt using question-answer pairs. iii. *Synonymous Substitution*. We substitute “block” with “region” and “is” with “looks like.” iv. *CPT* [58] Following this work, we color the detected region for each tag, assigning a unique color to each region.

The results are reported in Table 16. We observe that incorporating multiple objects or positions does not significantly improve the model’s performance on downstream tasks, and the language modeling loss is higher than the baseline. This suggests that the assignment is too difficult for the model to learn. We also find that the outcome of simple synonymous substitution remains consistent with the original text prompt’s outcome. Modeling location information only requires a straightforward prompt. CPT is designed for downstream visual grounding by coloring region proposals for identification, while PTP is for pretraining, which only pretrains a VLP model for numerous downstream tasks. Since most downstream tasks do not have region proposals available (e.g., VQA), CPT cannot generate color prompts to boost grounding, while PTP can, as it improves grounding during the pretraining phase. In fact, we attempted CPT for pretraining and observed worse performance, e.g., 75.1 for CPT vs. 77.8 for PTP on NLVR. For CPT, downstream tasks like NLVR do not have color prompts (CP), while its pretraining uses CP, leading to inconsistent phases.

We also find that selecting the top-1 predicted object and using its corresponding bounding box provides the best performance. It should be noted that the bounding box is rectangular, while the actual object may have various shapes. One possible explanation for this observation is that other regions or blocks may contain excessive background

TABLE 19: The different ways to get grid pseudo label and its corresponding running time. We report the image-to-text retrieval result on the COCO dataset for reference.

Method	Time	R1	R5	R10
baseline	-	70.6	91.3	95.4
Faster-RCNN (ResNet101)	10d	72.7	91.8	95.7
Faster-RCNN (ResNeXt152)	14d	73.3	92.0	96.1
CLIP Similarity	8h	72.9	92.0	96.6

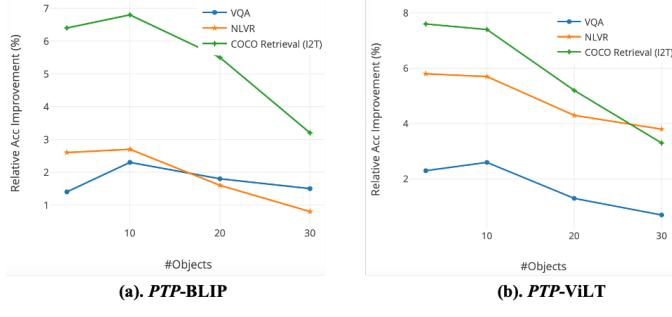


Fig. 4: We varying the number of selecting objects from 3 to 30. We report the result on downstream tasks over BLIP and ViLT baselines.

noise, making them challenging to identify.

**Assessing the Role of Position in Text Prompts.** In this experiment, we investigate the effectiveness of prompting our *PTP* for information at various granularities, such as without Positional. We simply use *[P]* has *[O]* when removing the prompt. The results are listed in Table 18. From these results we observe: i. It is interesting to see that each component is crucial. Without any one component, the downstream performance progressively deteriorates. ii. Although OSCAR [31] found that using object tags as supplementary input improved results when area features were used as input, we demonstrate that object tags are ineffective when raw pixel images are used. This highlights the importance of devising a functional prompt for learning alignment between object tags and image regions.

**Exploring the Number of Blocks.** We investigate whether more fine-grained position information benefits our *PTP*. In Figure 5, we vary the number of blocks from  $1 \times 1$  (removing position information in *PTP*) to  $4 \times 4$  and report the relative performance over BLIP/ViLT models. As can be seen, the results for both backbones improve when the number of blocks is more than 1. However, when there are 16 blocks, all downstream tasks experience a relative decline in performance. The reason may be that the predicted bounding box deviates from the localization of the actual object, resulting in a mesh that is too small and may not contain the selected object. Consequently, we opt for using  $3 \times 3$  blocks, as this configuration offers better accuracy.

**Determining the Optimal Number of Objects.** To generate object tags, one approach is to use Faster-RCNN and detect at least 10 objects from an image. We varied the number of objects from 5 to 30, exploring the effects of different object counts. The results are shown in Figure 4.

We observe that the BLIP baseline exhibits a slight increase in performance at the beginning, emphasizing the

TABLE 20: Part samples with position information. Under 14M setting, we test the result with different amount of pre-training samples with objects.

Method	COCO Retrieval TR@1	NLVR IR@1	COCO Captioning Acc	COCO Captioning CiDER
0%	79.5	62.4	80.5	129.5
19.3%	82.1	66.3	81.4	133.1
68.6%	84.6	69.1	83.1	134.6

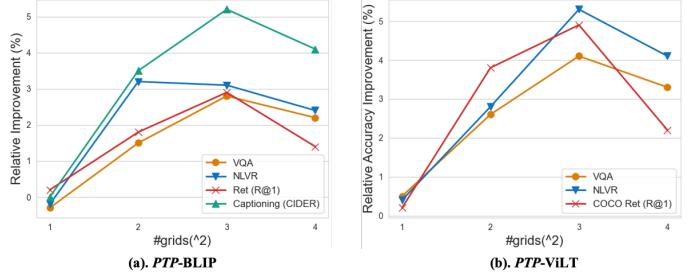


Fig. 5: The relation between the number of blocks and the relative accuracy improvement. We explore two baselines and show the improvements over four different tasks.

significance of data diversity for subsequent tasks. However, the results are less promising with a larger number of objects due to the high likelihood of false predictions resulting from a substantial number of objects with low confidence scores. *PTP* selects objects based on the confidence scores predicted by Faster-RCNN. A higher number of selected objects leads to lower confidence scores and increased noise in the tags, thus requiring a trade-off between the number of selected objects and the associated tag noise. In this study, we set the default number of selected objects to 10.

**Partial Bounding Box Annotation.** Since some URLs for the CC3M dataset are no longer valid and certain images are in an incorrect format, we have opted to extract objects from 2.7 million data points in the CC3M dataset and 7 million data points in the CC12M dataset. Consequently, only 9.7 million of the pre-training samples have available objects. We also report results for the 14 million setting, wherein we utilize the original text without text prompts in cases where objects are not available.

The outcomes for different sampling subsets are presented in Table 20. Our analysis reveals that the 68.6% object availability results in a CiDER value of 134.6 for COCO Captioning and an accuracy of 83.2 for the NLVR test-P. These findings demonstrate that an increased number of annotated samples contributes to better overall performance. Additionally, this observation supports the notion that the *PTP* approach is well-suited for large-scale pre-training, further validating its applicability and potential in the development of more advanced models.

## 5 DISCUSSION

**Evaluating the Necessity of an Object Detector.** In this work, part of the predicted bounding box information comes from Faster-RCNN [43]. To verify the expressive power of objects, we also consider two variations: i. Pure

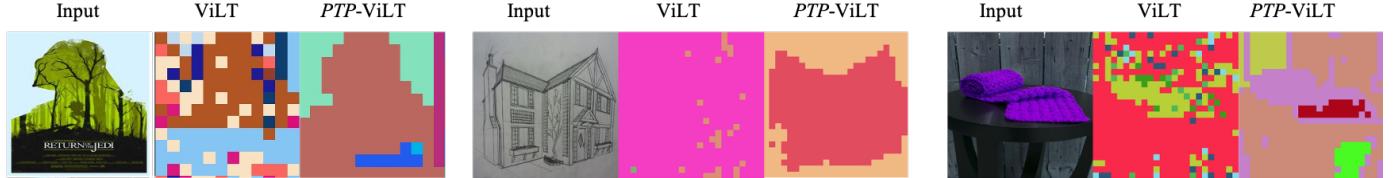


Fig. 6: **Token cluster visualization.** We train ViLT and *PTP-ViLT* with ViT-B/32 model on CC3M train set. We show the token cluster result with KMeans algorithm from CC3M test set [45]. *PTP-ViLT* shows preferable clusters.

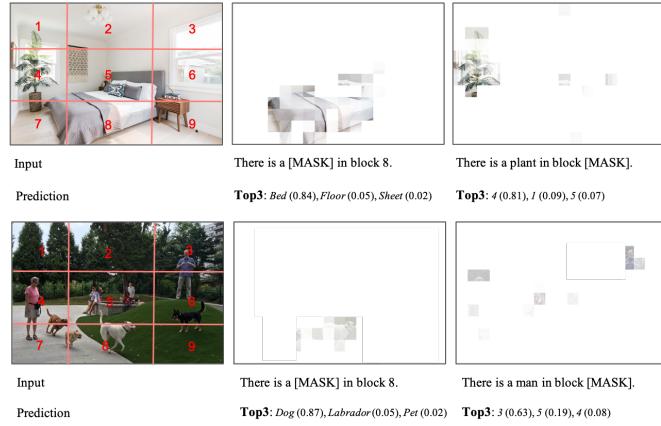


Fig. 7: **Evaluation of the full-in-the-blank task**, where the model is required to predict the objects contained within a given block and identify the blocks having a specific object.

CLIP similarity. This design choice is adapted mainly for efficiency reasons, as utilizing an object detector is time-consuming and not always easy to access. ii. In addition to the powerful ResNext152-based object detector, we also use a smaller ResNet101 based Faster-RCNN.

The results are reported in Table 19. We also report the overall feature extracting time on 8 NVIDIA V100 GPUs. Mainstream works in VLP use object detectors, and slow preprocessing is a common problem. The community tolerates this because objects only need to be extracted once and saved on disk (our features are downloaded from VinVL). To reduce costs further, PTP uses CLIP as a feature extractor, which is 42 times faster than Faster-RCNN. On 8 A100 GPUs, under the 4M setting, BLIP needs 19.4 hours, while PTP with CLIP requires  $8 + 19.4 = 27.4$  hours to train from scratch, which is affordable.

As can be seen from the table, we find that using a stronger detector leads to better results but incurs a substantial computational cost. Moreover, we observe that the result of CLIP embedding is very close to Faster-RCNN (ResNeXt152). Additionally, it takes only around 2.3% of the time of the Faster-RCNN version to extract pseudo labels for each grid. We conclude that a CLIP model is a suitable alternative to an object detector in PTP.

**Position Information Exploration.** To explore whether model training with the *PTP* framework indeed learns position information, we design a fill-in-the-blank evaluation experiment in this section. Following ViLT [22], we mask some key words and ask the model to predict the masked words and show its corresponding heatmap. We design two text prompts: one given the noun to predict the localization

TABLE 21: Comparison with BLIP baseline on two splits of VQA dataset. We report accuracy(%).

Method	Position-related	Position-unrelated	All
BLIP	72.5	74.2	73.8
<i>PTP-BLIP</i>	78.4	74.6	75.6

TABLE 22: COCO text-to-image retrieval. The second-order relation over objects.

Method	ITC	ITM	MLM	TR@1	IR@1	NLVR
Baseline	2.098	0.106	1.731	70.6	53.4	76.1
w GT Label	2.044	0.105	1.722	71.8	54.6	77.4
w <i>PTP</i>	1.883	0.093	1.4290	73.2	55.4	77.9

and the other given the localization to predict the missing noun. We show the top-3 predictions for reference.

The results are shown in Figure 7. Our findings reveal that *PTP-ViLT* is capable of making accurate predictions by utilizing both the block position information and its corresponding visual concepts. Moreover, when we mask only the position noun, we still observe a high probability of correct block prediction. For instance, as depicted in the bottom part of Figure 7, our model accurately identifies all image patches resembling the object of “man”. Based on these experiments and the insights presented in Figure 1, we conclude that *PTP* is an effective tool for facilitating the learning of position information in a vision-language model, using a simple yet powerful text prompt.

Furthermore, we cluster the token-level features with the K-Means algorithm for ViLT and *PTP-ViLT*. Intuitively, tokens with similar semantics should be clustered together. We show the visualization results in Figure 6. Comparing with the ViLT baseline, we observe that our method can cluster similar patches more accurately. This illustrates that our *PTP* has fairly accurately learned semantic information.

**What Kind of Samples Does *PTP* Help?** We undertake a comprehensive analysis of the visual-question answering task. Our investigation reveals that a significant portion of the VQA dataset samples contain position-related words, such as “top” and “sitting in.” To further examine this observation, we construct a vocabulary comprising 30 commonly occurring position words. Subsequently, we categorize the VQA dataset into two subsets: the position-related subset (approximately 27%) and the position-unrelated subset (roughly 73%), based on whether the text contains words from the aforementioned vocabulary.

Table 21 demonstrates the effectiveness of the proposed categorization scheme and its corresponding performance on the test-dev set. Our analysis reveals that *PTP* achieves



Fig. 8: Our text prompt (in red color) and its corresponding bounding box's mask. The block index spans a range from 0 to 8. Furthermore, we employ data augmentation techniques on the bounding box to ensure its alignment with the transformations applied to the input image.

### Captioning



A large bus **sitting next to** a very tall building.

### VQA



Q: What is on the desk **in front of** the boys?

A: Laptop.

### Retrieval



The boy is **sitting in** the fridge.

### (a). Downstream Tasks



Who is wearing glasses?



Is the umbrella **upside down**?



Where is the child **sitting**?

**BLIP** Man (0.25)

Yes (**0.43**)

Arm (0.14)

**PTP-BLIP** Man (0.87)

Yes (0.95)

Arm (0.52)

### (b). Some Examples on VQA

Fig. 9: (a). Position information is essential for mainstream downstream vision-language tasks. (b). In the context of the VQA [4] dataset, our PTP model provides improved predictions for position-related examples.

an accuracy of 78.4% on the position-related subset, which is 5.9% higher than the BLIP baseline. This result highlights the significant contribution of PTP towards enhancing the model's ability to efficiently learn position information and underscores its robust visual grounding capability. Thus, our proposed model can serve as a valuable asset in addressing visual question-answering tasks.

**Comparison with Direct Object Regression.** Learning position knowledge from an object detector is indeed challenging. In this section, we consider the coordinates of predicted object bounding boxes as ground truth labels and regress them during pre-training. Specifically, we add

TABLE 23: Existing object-centric datasets have a limited number of significant objects.

Dataset	i=2	i=3	i=5
COCO [33]	34%	49%	73%
CC12M [8]	41%	58%	84%
DataComp1B(subset) [14]	39%	56%	81%

a lightweight detection head after the BLIP image encoder. Table 22 shows that PTP performs better than the ground truth label approach in terms of pre-training loss (ITC, ITM, MLM) and downstream task performance (TR, IR, NLVR).

We find that the bounding boxes provided by Faster-RCNN are not very precise, and enforcing the model to regress coordinates could bias its grounding ability. Regressing objects directly requires predicting coordinates, which are not very precise. Such an implementation focuses only on the image encoder and does not improve the text decoder model. Moreover, this implementation is complex and involves many tricks to explore, making it difficult to extend to other frameworks easily. In contrast, our block position representation for objects is more accurate, ensuring that the model learns correct position information. With the position-guided text prompt (e.g., giving position/block to predict object), the model learns which blocks contain objects and what objects are in each block. This way, the model implicitly learns visual grounding, as experimentally demonstrated in the previous sections.

**Exploring the Feasibility of Higher-Order PTP Integration.** The consideration of incorporating higher-order relations into image captions presents a significant point of discussion. This deliberation is rooted in the potential extension of caption lengths, which necessitates scrutiny. Notably, conventional Visual-Language Pretraining (VLP) models adhere to specific maximum caption length constraints, typically set at 32 tokens, and this constraint bears relevance in our examination.

It is pivotal to recognize that the majority of images within the corpus exhibit a rather limited diversity of dis-

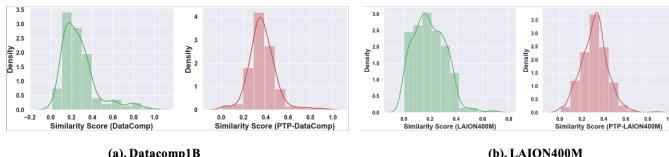


Fig. 10: Image-text similarity score distribution for original and revised captions using *PTP*.

tinct objects. To elucidate this context, we undertook an extensive analysis to quantify the unique object classes within each image, eliminating frequent, non-discriminatory labels such as "sky," "road," and "tree," alongside objects exceeding 1/10th of the image's width. The results of this rigorous analysis are documented in Table 23 for reference.

Our findings unveil a noteworthy statistic; approximately 40% of samples sourced from prominent pre-training datasets, such as DataComp1B [14], encompass two or fewer unique objects. Given this empirical insight, the endeavor to generate third-order relations within these object-centric datasets emerges as a formidable challenge. Therefore, the judicious selection of second-order relations assumes a pragmatic stance within the scope of our inquiry.

#### Erasing Misalignment through PTP Implementation.

A primary challenge observed in existing large-scale vision-language pre-training datasets is the significant misalignment between image and text pairs [14]. In other words, there is often a substantial lack of alignment between the provided images and their corresponding textual descriptions, making it difficult for the model to learn effectively.

In this experiment, we set out to tackle this issue by assessing dot product similarity scores between LAION400M [44] data and Datacomp1B [14]. Specifically, we randomly selected 10% subsets from both datasets for analysis. Utilizing the pre-trained CLIP model, available at the following link<sup>3</sup>, we calculated the dot product similarity scores for image-text pairs. The comparison results are visually presented in Figure 10, demonstrating the efficacy of our method in erasing the impact of noisy image-text pairs.

It is worth noting that in each epoch, our approach generates distinct captions due to the random selection of positions and objects. This not only erases the misalignment issue but also contributes to the enhancement of the overall representation from the data aspect.

## 6 VISUALIZATION

### 6.1 Case Analysis

In this experiment, we showcase several cases involving position information in Figure 9. We observe that position information is crucial for various downstream tasks, including captioning, VQA, and retrieval. To comprehend these tasks, the trained model needs to learn position information.

Since a large number of samples in VQA tasks typically include position information, we evaluate our model on VQA tasks and select some representative samples. Specifically, we display the prediction probability and predicted nouns at the bottom of this figure. We observe that *PTP*

provides accurate predictions in most cases, illustrating that our *PTP* learns position information more effectively.

### 6.2 Bounding Box Visualization

In this section, we present the object detection results obtained using our generated text prompts. Specifically, we randomly select an object from the set  $V$  and then visualize the original image along with the corresponding bounding box mask. It is important to note that we apply the same affine transformation to these bounding boxes as we do to the original image, ensuring consistency between the image and the corresponding bounding box mask.

We randomly select some samples from the overall dataset, and the results are reported in Figure 8. We also observe that the bounding box may be very large and span multiple blocks in some examples (e.g., the first case in the third row). Since we use RandAugment [11] in this work, some objects may be outside the border of the input image. For such situations, we simply replace the specific position with [X], and the final *PTP* is *The block [X] has a [O]*. We also find that some masks may not be square, as seen in the last example in the third row.

## 7 LIMITATIONS AND CONCLUSION

Initially, we attempted to leverage position information from existing object detectors or trained models to enhance Visual-Language Pre-training models using simple prompts. To aid in prompt engineering, we developed a successful practice of cross-modal prompt settings. In addition to the first-order relation between objects and position, we also explored more complicated second-order relations between objects. Through rigorous experiments, we demonstrated that *PTP* serves as a general-purpose pipeline and improves the learning of position information without incurring significant extra computational costs.

Although the current version of *PTP* has demonstrated significant progress in its ability to process and interpret various input data, it is essential to acknowledge the limitations that still persist. One primary concern is that, at this time, *PTP* does not possess the capability to effectively handle instances wherein an incorrect object tag is presented. Furthermore, the present scope of this research has not delved deeply into the intricacies associated with more complex prompts. An in-depth exploration of such prompts would facilitate a better understanding of the model's strengths and weaknesses, paving the way for further refinements and enhancements in future iterations. Looking ahead, it is crucial to broaden the research horizons to evaluate *PTP*'s performance across a diverse range of vision-language tasks.

## ACKNOWLEDGEMENT

This project is supported by the National Research Foundation, Singapore under its NRFF Award NRF-NRFF13-2021-0008, and Mike Zheng Shou's Start-Up Grant from NUS.

3. [https://github.com/mlfoundations/open\\_clip](https://github.com/mlfoundations/open_clip)

## REFERENCES

- [1] Agrawal, H., Desai, K., Wang, Y., Chen, X., Jain, R., Johnson, M., Batra, D., Parikh, D., Lee, S., Anderson, P.: Nocaps: Novel object captioning at scale. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8948–8957 (2019) [1, 6](#)
- [2] Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al.: Flamingo: a visual language model for few-shot learning. arXiv preprint arXiv:2204.14198 (2022) [6](#)
- [3] Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and visual question answering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6077–6086 (2018) [1, 3, 6](#)
- [4] Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: Vqa: Visual question answering. In: Proceedings of the IEEE international conference on computer vision. pp. 2425–2433 (2015) [1, 7, 13](#)
- [5] Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., Jitsev, J., Kornblith, S., Koh, P.W., Ilharco, G., Wortsman, M., Schmidt, L.: Openflamingo: An open-source framework for training large autoregressive vision-language models. arXiv preprint arXiv:2308.01390 (2023) [8](#)
- [6] Bain, M., Nagrani, A., Varol, G., Zisserman, A.: Frozen in time: A joint video and image encoder for end-to-end retrieval. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1728–1738 (2021) [7](#)
- [7] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems **33**, 1877–1901 (2020) [10](#)
- [8] Changpinyo, S., Sharma, P., Ding, N., Soricut, R.: Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3558–3568 (2021) [5, 13](#)
- [9] Chen, X., Fang, H., Lin, T.Y., Vedantam, R., Gupta, S., Dollár, P., Zitnick, C.L.: Microsoft coco captions: Data collection and evaluation server. arXiv preprint arXiv:1504.00325 (2015) [1](#)
- [10] Chen, Y.C., Li, L., Yu, L., El Kholy, A., Ahmed, F., Gan, Z., Cheng, Y., Liu, J.: Uniter: Learning universal image-text representations (2019) [1, 2, 3, 6, 7](#)
- [11] Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.V.: Randaugment: Practical automated data augmentation with a reduced search space. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops. pp. 702–703 (2020) [5, 14](#)
- [12] Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018) [2, 8](#)
- [13] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) [6](#)
- [14] Gadre, S.Y., Ilharco, G., Fang, A., Hayase, J., Smyrnis, G., Nguyen, T., Marten, R., Wortsman, M., Ghosh, D., Zhang, J., et al.: Datacomp: In search of the next generation of multimodal datasets. arXiv preprint arXiv:2304.14108 (2023) [2, 5, 8, 13, 14](#)
- [15] Gan, Z., Chen, Y.C., Li, L., Zhu, C., Cheng, Y., Liu, J.: Large-scale adversarial training for vision-and-language representation learning. Advances in Neural Information Processing Systems **33**, 6616–6628 (2020) [7](#)
- [16] Hu, X., Gan, Z., Wang, J., Yang, Z., Liu, Z., Lu, Y., Wang, L.: Scaling up vision-language pre-training for image captioning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 17980–17989 (2022) [5](#)
- [17] Huang, Z., Zeng, Z., Huang, Y., Liu, B., Fu, D., Fu, J.: Seeing out of the box: End-to-end pre-training for vision-language representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12976–12985 (2021) [1](#)
- [18] Ilharco, G., Wortsman, M., Wightman, R., Gordon, C., Carlini, N., Taori, R., Dave, A., Shankar, V., Namkoong, H., Miller, J., Hajishirzi, H., Farhadi, A., Schmidt, L.: Openclip (Jul 2021). <https://doi.org/10.5281/zenodo.5143773>, <https://doi.org/10.5281/zenodo.5143773>, if you use this software, please cite it as below. [8](#)
- [19] Jang, Y., Song, Y., Yu, Y., Kim, Y., Kim, G.: Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2758–2766 (2017) [2, 8](#)
- [20] Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: International Conference on Machine Learning. pp. 4904–4916. PMLR (2021) [1, 6](#)
- [21] Khattak, M.U., Rasheed, H., Maaz, M., Khan, S., Khan, F.S.: Maple: Multi-modal prompt learning. arXiv preprint arXiv:2210.03117 (2022) [2](#)
- [22] Kim, W., Son, B., Kim, I.: Vilt: Vision-and-language transformer without convolution or region supervision. In: International Conference on Machine Learning. pp. 5583–5594. PMLR (2021) [1, 2, 3, 4, 5, 6, 7, 8, 9, 12](#)
- [23] Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.J., Shamma, D.A., et al.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. International journal of computer vision **123**(1), 32–73 (2017) [1, 3, 5](#)
- [24] Lei, J., Li, L., Zhou, L., Gan, Z., Berg, T.L., Bansal, M., Liu, J.: Less is more: Clipbert for video-and-language learning via sparse sampling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7331–7341 (2021) [8](#)
- [25] Lei, J., Yu, L., Berg, T.L., Bansal, M.: Tvqa+: Spatio-temporal grounding for video question answering. arXiv preprint arXiv:1904.11574 (2019) [2, 8](#)
- [26] Li, G., Duan, N., Fang, Y., Gong, M., Jiang, D.: Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In: Proceedings of the AAAI Conference on Artificial Intelligence (2020) [6](#)
- [27] Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: ICML (2022) [2, 3, 4, 5, 6, 7, 8, 9](#)
- [28] Li, J., Selvaraju, R., Gotmare, A., Joty, S., Xiong, C., Hoi, S.C.H.: Align before fuse: Vision and language representation learning with momentum distillation. Advances in neural information processing systems **34**, 9694–9705 (2021) [6, 7](#)
- [29] Li, L.H., Zhang, P., Zhang, H., Yang, J., Li, C., Zhong, Y., Wang, L., Yuan, L., Zhang, L., Hwang, J.N., et al.: Grounded language-image pre-training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10965–10975 (2022) [3](#)
- [30] Li, W., Gao, C., Niu, G., Xiao, X., Liu, H., Liu, J., Wu, H., Wang, H.: Unimo: Towards unified-modal understanding and generation via cross-modal contrastive learning. arXiv preprint arXiv:2012.15409 (2020) [7](#)
- [31] Li, X., Yin, X., Li, C., Zhang, P., Hu, X., Zhang, L., Wang, L., Hu, H., Dong, L., Wei, F., et al.: Oscar: Object-semantics aligned pre-training for vision-language tasks. In: European Conference on Computer Vision. pp. 121–137. Springer (2020) [1, 3, 5, 6, 7, 11](#)
- [32] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014) [5, 6](#)
- [33] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V **13**. pp. 740–755. Springer (2014) [13](#)
- [34] Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. arXiv preprint arXiv:2107.13586 (2021) [2](#)
- [35] Liu, Z., Stent, S., Li, J., Gideon, J., Han, S.: Loctex: Learning data-efficient visual representations from localized textual supervision. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2167–2176 (2021) [3](#)
- [36] Lu, J., Batra, D., Parikh, D., Lee, S.: Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. Advances in neural information processing systems **32** (2019) [7](#)
- [37] Miech, A., Alayrac, J.B., Smaira, L., Laptev, I., Sivic, J., Zisserman, A.: End-to-end learning of visual representations from uncurated instructional videos. In: Proceedings of the IEEE/CVF Conference

- on Computer Vision and Pattern Recognition. pp. 9879–9889 (2020) 7
- [38] Ordóñez, V., Kulkarni, G., Berg, T.: Im2text: Describing images using 1 million captioned photographs. Advances in neural information processing systems **24** (2011) 5
- [39] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems **32** (2019) 5
- [40] Qi, D., Su, L., Song, J., Cui, E., Bharti, T., Sacheti, A.: Imagebert: Cross-modal pre-training with large-scale weak-supervised image-text data. arXiv preprint arXiv:2001.07966 (2020) 6, 7
- [41] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML. pp. 8748–8763. PMLR (2021) 1, 2, 3, 4, 6, 8, 9
- [42] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., et al.: Exploring the limits of transfer learning with a unified text-to-text transformer. J. Mach. Learn. Res. **21**(140), 1–67 (2020) 2
- [43] Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems **28** (2015) 1, 3, 11
- [44] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems **35**, 25278–25294 (2022) 5, 8, 14
- [45] Sharma, P., Ding, N., Goodman, S., Soricut, R.: 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). pp. 2556–2565 (2018) 5, 12
- [46] Su, W., Zhu, X., Cao, Y., Li, B., Lu, L., Wei, F., Dai, J.: Vi-bert: Pre-training of generic visual-linguistic representations. arXiv preprint arXiv:1908.08530 (2019) 7
- [47] Suhr, A., Zhou, S., Zhang, A., Zhang, I., Bai, H., Artzi, Y.: A corpus for reasoning about natural language grounded in photographs. arXiv preprint arXiv:1811.00491 (2018) 1, 7
- [48] Touvron, H., Larvil, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambrø, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023) 8
- [49] Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: Lessons learned from the 2015 ms coco image captioning challenge. IEEE transactions on pattern analysis and machine intelligence **39**(4), 652–663 (2016) 6
- [50] Wang, A.J., Ge, Y., Yan, R., Ge, Y., Lin, X., Cai, G., Wu, J., Shan, Y., Qie, X., Shou, M.Z.: All in one: Exploring unified video-language pre-training. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022) 8
- [51] Wang, A.J., Zhou, P., Shou, M.Z., Yan, S.C.: Position-guided text prompt for vision language pre-training. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023) 2
- [52] Wang, J., Yang, Z., Hu, X., Li, L., Lin, K., Gan, Z., Liu, Z., Liu, C., Wang, L.: Git: A generative image-to-text transformer for vision and language. arXiv preprint arXiv:2205.14100 (2022) 7
- [53] Wang, J., Ge, Y., Cai, G., Yan, R., Lin, X., Shan, Y., Qie, X., Shou, M.Z.: Object-aware video-language pre-training for retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3313–3322 (2022) 7
- [54] Wang, W., Bao, H., Dong, L., Bjorck, J., Peng, Z., Liu, Q., Aggarwal, K., Mohammed, O.K., Singhal, S., Som, S., et al.: Image as a foreign language: Beit pretraining for all vision and vision-language tasks. arXiv preprint arXiv:2208.10442 (2022) 7
- [55] Wang, Z., Yu, J., Yu, A.W., Dai, Z., Tsvetkov, Y., Cao, Y.: Simvlm: Simple visual language model pretraining with weak supervision. arXiv preprint arXiv:2108.10904 (2021) 5, 7
- [56] Xu, D., Zhao, Z., Xiao, J., Wu, F., Zhang, H., He, X., Zhuang, Y.: Video question answering via gradually refined attention over appearance and motion. In: Proceedings of the 25th ACM international conference on Multimedia. pp. 1645–1653 (2017) 2, 8
- [57] Yao, L., Huang, R., Hou, L., Lu, G., Niu, M., Xu, H., Liang, X., Li, Z., Jiang, X., Xu, C.: Filip: Fine-grained interactive language-image pre-training. arXiv preprint arXiv:2111.07783 (2021) 1, 6
- [58] Yao, Y., Zhang, A., Zhang, Z., Liu, Z., Chua, T.S., Sun, M.: Cpt: Colorful prompt tuning for pre-trained vision-language models. arXiv preprint arXiv:2109.11797 (2021) 2, 10
- [59] Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., Wu, Y.: Coca: Contrastive captioners are image-text foundation models. arXiv preprint arXiv:2205.01917 (2022) 1, 2, 7
- [60] Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14. pp. 69–85. Springer (2016) 2
- [61] Yuan, L., Chen, D., Chen, Y.L., Codella, N., Dai, X., Gao, J., Hu, H., Huang, X., Li, B., Li, C., et al.: Florence: A new foundation model for computer vision. arXiv preprint arXiv:2111.11432 (2021) 6
- [62] Zeng, Y., Zhang, X., Li, H.: Multi-grained vision language pre-training: Aligning texts with visual concepts. arXiv preprint arXiv:2111.08276 (2021) 3, 5
- [63] Zhang, P., Li, X., Hu, X., Yang, J., Zhang, L., Wang, L., Choi, Y., Gao, J.: Vinvl: Making visual representations matter in vision-language models. arXiv preprint arXiv:2101.00529 1(6), 8 (2021) 1, 3, 5, 6, 7, 9
- [64] Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X.V., et al.: Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068 (2022) 8
- [65] Zhong, Y., Yang, J., Zhang, P., Li, C., Codella, N., Li, L.H., Zhou, L., Dai, X., Yuan, L., Li, Y., et al.: Regionclip: Region-based language-image pretraining. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16793–16803 (2022) 3
- [66] Zhou, C., Loy, C.C., Dai, B.: Denseclip: Extract free dense labels from clip. arXiv preprint arXiv:2112.01071 (2021) 3
- [67] Zhu, L., Yang, Y.: Actbert: Learning global-local video-text representations. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8746–8755 (2020) 7
- [68] Zhu, W., Hessel, J., Awadalla, A., Gadre, S.Y., Dodge, J., Fang, A., Yu, Y., Schmidt, L., Wang, W.Y., Choi, Y.: Multimodal c4: An open, billion-scale corpus of images interleaved with text. arXiv preprint arXiv:2304.06939 (2023) 5, 8



**Alex Jinpeng Wang** is currently a Third-year PHD student in National University of Singapore. Before that, he obtained the bachelor and master degree in Sun Yat-Sen University (SYSU). His research focuses on Large-scale Visual Language Pre-training and Data-centric AI. He has published relevant papers in the top-tier conferences: CVPR, ICCV, AAAI, NeurIPS and ECCV. He also serves as reviewer for: CVPR, ICCV, ECCV, ICLR and NeurIPS.



**Pan Zhou** is currently a senior Research Scientist at Sea AI Lab (SAIL) of Sea group. Before that, he worked in Salesforce as a research scientist. He completed his Ph.D. at National University of Singapore (NUS) and Master at Peking University. His research interests include deep learning theory and applications, nonconvex/convex optimization. He has published papers in ICLR, ICML, NeurIPS, CVPR, ICCV, ECCV, AAAI, IJCAI and journals: TPAMI, TIP. He serves as reviewer for top conferences: ICML, NeurIPS, CVPR, ICCV, AAAI and journals: TPAMI, IJCV, TIP, TNNLS and TCSV. He is awarded the Microsoft Research Asia Fellowship.



**Mike Zheng Shou** is a tenure-track Assistant Professor at National University of Singapore. He holds a PhD degree from Columbia University in the City of New York. He received the best paper finalist at CVPR'22 and the best student paper nomination at CVPR'17. His team won 1st place in multiple international challenges including ActivityNet 2017, EPIC-Kitchens 2022, Ego4D 2022 & 2023. He is a Fellow of the National Research Foundation (NRF) Singapore and has been named on the Forbes 30 Under 30

Asia list.



**Shuicheng Yan** (Fellow, IEEE) is currently a visiting professor at BAAI, Beijing, China. Previously, he was the director of Sea AI Lab (SAIL) and Group Chief Scientist of Sea. He is an Fellow of Academy of Engineering, Singapore, IEEE Fellow, ACM Fellow, AAAI Fellow and IAPR Fellow. His research areas include computer vision, machine learning and multimedia analysis. Till now, he has published over 600 papers in top international journals and conferences, with Google Scholar Citation over 90,000 times and

H-index 135. He had been among "Thomson Reuters Highly Cited Researchers" in 2014, 2015, 2016, 2018, 2019. Dr. Yan's team has received winner or honorable-mention prizes for 10 times of two core competitions, Pascal VOC and ImageNet (ILSVRC), which are deemed as "World Cup" in the computer vision community. Also his team won over 10 best paper or best student paper prizes and especially, a grand slam in ACM MM, the top conference in multimedia, including Best Paper Award, Best Student Paper Award and Best Demo Award.