

# ZERO: A LARGE-SCALE CHINESE CROSS-MODAL BENCHMARK WITH A NEW VISION-LANGUAGE FRAMEWORK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Vision-language pre-training (VLP) on large-scale datasets has shown premier performance on various downstream tasks. In contrast to plenty of available benchmarks with English corpus, large-scale pre-training datasets and downstream datasets with Chinese corpus remain largely unexplored. In this paper, we build a large-scale Chinese cross-modal benchmark from ZERO, which is named for our database publicly available for the research community to build VLP models. We release a pre-training dataset and five fine-tuning datasets for downstream tasks, and also develop a pre-training framework of pre-**R**anking + **R**anking with target-guided **D**istillation and feature-guided **D**istillation (R2D2) for cross-modal learning. In specific, a global contrastive pre-ranking is first introduced to learn the individual representations of images and texts. We then fuse the representations in a fine-grained ranking manner via an image-text cross encoder and a text-image cross encoder. To further enhance the capability of our method, a two-way distillation strategy is used with target-guided distillation and feature-guided distillation. We achieve state-of-the-art performance on eleven downstream datasets from four broad categories of tasks including image-text retrieval, image-text matching, image caption, and text-to-image generation.

## 1 INTRODUCTION

Vision-language pre-training (VLP) mainly learns the semantic correspondence between vision and natural language. Previous works (Li et al., 2019b; Singh et al., 2022; Li et al., 2021a) explore the VLP model and achieve significant improvement on various vision-language (V+L) tasks. These methods are supported by massive data (Schuhmann et al., 2021), excellent architectures such as Transformer (Vaswani et al., 2017), and cross-modal models such as CLIP (Radford et al., 2021).

There are plenty of available benchmarks with English corpus, such as Conceptual Captions (Sharma et al., 2018), SBU Captions (Ordonez et al., 2011), and LAION (Schuhmann et al., 2021). Differently, large-scale pre-training datasets and downstream datasets with Chinese corpus are relatively few. M6-Corpus (Lin et al., 2021) is a multi-modal pre-training dataset in Chinese but not publicly available. Wukong (Gu et al., 2022) is a newly published pre-training dataset. Most existing downstream Chinese datasets mainly focus on retrieval tasks, such as Flickr30k-CN (Lan et al., 2017) and AIC-ICC (Wu et al., 2019), which are not sufficient for a complete evaluation of VLP models. Besides, Flickr30k-CN tries to translate English cross-modal downstream datasets into Chinese, however, fails to cover Chinese idioms and often causes translation errors.

In this paper, we introduce a large-scale Chinese cross-modal benchmark (called ZERO), including a pre-training dataset (ZERO-Corpus) and five downstream datasets. Specifically, ZERO-Corpus consists of 23 million image-text pairs, which are collected from the search engine with images and corresponding textual descriptions, by filtering from 5 billion image-text pairs by user click-through rate (CTR). Compared to existing pre-training datasets, ZERO-Corpus is high-quality due to the collection way of user CTR and with diverse textual information for each image. Together with the pre-training dataset, we provide 5 high-quality downstream datasets. Two of them are the largest Chinese V+L downstream datasets and first proposed for Chinese image-text matching task, which is also important for evaluating VLP models. For the image-text retrieval task, we provide 3 datasets,

especially our Flickr30k-CNA, which is a more comprehensive and accurate human-annotated dataset than Flickr30k-CN (Lan et al., 2017). We build a leaderboard on the five downstream test datasets.

From the perspective of cross-modal learning, existing methods can be categorized as single-stream and dual-stream. Most single-stream methods (*e.g.*, (Chen et al., 2020; Li et al., 2021b; Qi et al., 2020)) employ an extra object detector to extract the patch embedding and then align patches and words. As illustrated in Li et al. (2021a), object detectors are annotation-expensive and computing-expensive, because they require bounding box annotations during pre-training and high-resolution (*e.g.*,  $600 \times 1000$ ) images during inference. On the other hand, for dual-stream architectures (*e.g.*, (Radford et al., 2021; Wei et al., 2021)), it is non-trivial to model the fine-grained associations between image and text, since the corresponding representations reside in their own semantic space.

To address these limitations, we introduce two strategies in a new cross-model learning framework. We first omit the object detection module from our network to avoid expensive annotations and computation. We then combine dual-stream architecture with single-stream architecture, where the single-stream architecture consists of two cross-encoders. The cross-encoders are able to learn image-to-text and text-to-image interactions in a fine-grained manner. During pre-training, we design a global contrastive pre-**R**anking loss to obtain image-text representations and fine-grained **R**anking loss to further improve model performance, inspired by industrial technology such as recommend systems (Cheng et al., 2016; Wang et al., 2020) and online advertising (Tan et al., 2021). We also introduce a two-way distillation method, consisting of target-guided **D**istillation and feature-guided **D**istillation. The target-guided distillation increases the robustness when learning from noisy labels, while feature-guided distillation aims to improve the generalization performance. We call this pre-training framework **R2D2**. To summarize, our main contributions are as follows:

- We introduce a Chinese cross-modal benchmark, including a large-scale pre-training dataset, which is high-quality due to the collection way of user CTR and with diverse textual information for each image. We also provide five human-annotated downstream train/val/test sets, two of which are currently the largest Chinese V+L downstream datasets.
- We introduce a VLP framework named **R2D2** for image-text cross-modal learning. Specifically, we propose a pre-**R**anking + **R**anking strategy to learn powerful vision-language representations and a two-way distillation method (*i.e.*, target-guided **D**istillation and feature-guided **D**istillation) to further enhance the learning capability.
- Our proposed method achieves state-of-the-art performance on eleven downstream datasets from four broad categories of V+L tasks, showing the superior ability of our pre-trained model.

## 2 CURATED VLP BENCHMARK: ZERO

### 2.1 PRE-TRAINING DATASETS

Existing public pre-training datasets suffer from two major limitations. First, the image-text pairs are collected usually by their co-occurrence relationship coarsely from third-party search engines or websites. Thus, the collected pairs are inherently noisy. Second, the text corpus lacks diversity as each image usually has one corresponding text description. To overcome these drawbacks, we collect a new dataset for Chinese image-text pre-training, called ZERO-Corpus. Specifically, we extract 23 million samples from 5 billion image-text pairs collected by an image search engine. The key point is filtering the candidates by the highest user CTR, which means users have clicked the most on an image searched by the same query. Moreover, we remove inappropriate images and harmful textual descriptions to keep only the most relevant and high-quality image-text pairs. We also provide diverse textual information for each image, *i.e.*, “Title”, “Content”, and “ImageQuery”. More details about the pre-training datasets can be found in Appendix A.

### 2.2 DOWNSTREAM DATASETS

We use an image search engine to construct 4 Chinese image-text datasets. In these datasets, each image has one corresponding text. We divide the training set, validation set, and test set with a ratio of 8:1:1. 15 human annotators carefully label all the image-text pairs. We also translate all data of Flickr30k (Young et al., 2014) by 6 professional linguists. The details of each dataset are as follows.

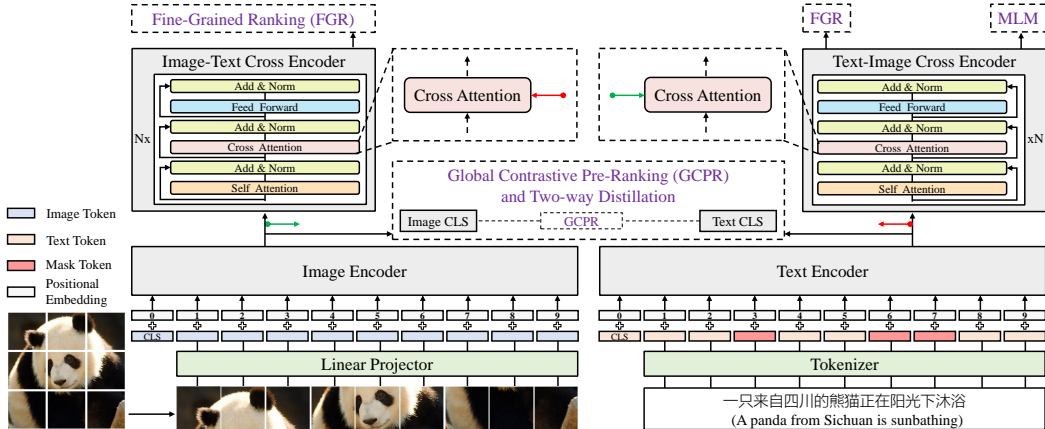


Figure 1: The overall architecture of the proposed framework. The image encoder and the text encoder aim to learn individual features of image and text, respectively. Then, the image features (green circled arrow) are fed into the text-image cross encoder. Similarly, the text features (red circled arrow) are fed into the image-text cross encoder. During pre-training, we apply global contrastive pre-ranking (GCPR) and fine-grained ranking (FGR) as pre-training objectives. Moreover, we introduce mask language modeling (MLM) and two-way distillation to obtain remarkable performance.

**Image-Caption Matching Dataset (ICM).** ICM is collected for the image-text matching task. Each image has a corresponding caption text, which describes the image in detail. We first use CTR to select the most relevant pairs. Then, human annotators manually perform a 2nd round manual correction, obtaining 400,000 image-text pairs, including 200,000 positive cases and 200,000 negative cases. We keep the ratio of positive and negative pairs consistent in each of the train/val/test sets.

**Image-Query Matching Dataset (IQM).** This is a dataset also for the image-text matching task. Different from ICM, we use the search query instead of detailed description text. In this dataset, we randomly select image-query pairs in the candidate set after performing the cleaning process described in Section 2.1. Similarly, IQM contains 200,000 positive cases and 200,000 negative cases. ICM and IQM are currently the largest Chinese vision-language downstream datasets.

**Image-Caption Retrieval Dataset (ICR).** In this dataset, we collect 200,000 image-text pairs under the rules described in ICM. It contains image-to-text retrieval and text-to-image retrieval tasks.

**Image-Query Retrieval Dataset (IQR).** IQR is also proposed for the image-text retrieval task. We randomly select 200,000 queries and the corresponding images as the annotated image-query pairs similar to IQM. We show examples of the above four datasets in Appendix B.

**Flickr30k-CNA Dataset.** Former Flickr30k-CN (Lan et al., 2017) translates the training and validation sets of Flickr30k (Young et al., 2014) using machine translation, and manually translates the test set. We check the machine-translated results and find two kinds of problems. (1) Some sentences have language problems and translation errors. (2) Some sentences have poor semantics. In addition, the different translation ways between the training set and test set prevent the model from achieving accurate performance. We gather 6 professional English and Chinese linguists to meticulously re-translate all data of Flickr30k and double-check each sentence. We name this dataset as Flickr30k-Chinese All (Flickr30k-CNA). We show some cases of the difference between Flickr30k-CN and Flickr30k-CNA in Appendix B.

### 3 APPROACH

#### 3.1 MODEL ARCHITECTURE

From Figure 1, R2D2 contains a text encoder, an image encoder, and two cross encoders. The text encoder and image encoder transform texts and images into sequences of hidden states separately. Then the text and image hidden states interact in the two cross encoders through cross-attention.

**Text Encoder.** We employ a BERT (Devlin et al., 2019) encoder as our text encoder. Given a textual sequence, we first tokenize it using the tokenizer of RoBERTa-wwm-ext (Cui et al., 2020). Here, a special  $[CLS]$  token is appended to the head of the tokenized text. Then, we feed the tokenized text into the text encoder.

**Image Encoder.** For the image encoder, we adopt the Vision Transformer (ViT) (Dosovitskiy et al., 2021). We firstly rescale the input image into a standard size and split the image into patches. Each of the patches is then linearly projected and a position embedding is added. Additionally, a learnable  $[CLS]$  token is concatenated with the patch vectors. The sequential vectors are finally fed into a standard Transformer model to obtain a list of image hidden state vectors.

**Cross Encoder.** The image and text hidden vectors are fused and fed into the cross encoders. Specifically, the linear projection layer is used to change the dimensionality of each text feature and image feature to make them consistent. The multi-layer transformers fuse the feature information of both modalities with the help of cross-attention and produce the final cross-modal outputs.

### 3.2 PRE-TRAINING METHODS

To explore the matching relationship between image and text pairs, we design a mechanism of pre-ranking + ranking, named global contrastive pre-ranking (GCPR) and fine-grained ranking (FGR). We adopt masked language modeling (MLM) to learn the representation of cross-modal models.

**Global Contrastive Pre-Ranking.** Traditional contrastive learning aims to align the representation of multi-modal data (*e.g.*, paired image and text). It maximizes the similarity score of the positive pairs and minimizes the score of the negative pairs. In practice, we use global contrastive learning to accomplish the pre-ranking task. We perform full back-propagation across  $k$  GPUs. For each image  $I_i$  and the corresponding text  $T_i$ , the softmax-normalized similarity score of image-to-text and text-to-image can be defined as:

$$s(I_i, T_i) = \frac{\exp(\text{sim}(I_i, T_i)/\tau)}{\sum_{j=1}^{n \times k} \exp(\text{sim}(I_i, T_j)/\tau)}, \quad s(T_i, I_i) = \frac{\exp(\text{sim}(T_i, I_i)/\tau)}{\sum_{j=1}^{n \times k} \exp(\text{sim}(T_i, I_j)/\tau)}, \quad (1)$$

where  $n$  is the batch size of one GPU,  $k$  is the number of GPUs,  $\tau$  is a learnable temperature parameter, and  $\text{sim}(\cdot, \cdot)$  denotes the cosine similarity between a pair of image-text. Let  $\mathcal{D}$  denote the training data and  $\mathbf{y}(\cdot, \cdot)$  denote the ground-truth one-hot label. The global contrastive pre-ranking loss is calculated by the cross-entropy loss  $\mathcal{L}_c(\cdot)$ , as shown in Equation (2).

$$\mathcal{L}_{\text{GCPR}} = \frac{1}{2} \mathbb{E}_{(I, T) \sim \mathcal{D}} [\mathcal{L}_c(\mathbf{s}(I, T), \mathbf{y}(I, T)) + \mathcal{L}_c(\mathbf{s}(T, I), \mathbf{y}(T, I))]. \quad (2)$$

**Fine-Grained Ranking.** As aforementioned, we apply global contrastive pre-ranking to obtain the individual representations of images and texts, respectively. Relying on these representations, we next perform Fine-Grained Ranking (FGR) loss to conduct a fine-grained ranking task. To be specific, this is a binary classification task, aiming to predict whether an image-text is matched. Formally, we denote  $h_{I_{[CLS]}}$  and  $h_{T_{[CLS]}}$  as the output representations of two cross encoders. Given an image representation  $h_{I_{[CLS]}}$  and a text representation  $h_{T_{[CLS]}}$ , we feed the representations into a fully-connected layer  $g(\cdot)$  to get the predicted probabilities respectively. Let  $\mathbf{y}$  denote the ground-truth label of binary classification, we then compute the FGR loss as:

$$\mathcal{L}_{\text{FGR}} = \frac{1}{2} \mathbb{E}_{(I, T) \sim \mathcal{D}} [\mathcal{L}_c(g(h_{I_{[CLS]}}), \mathbf{y}) + \mathcal{L}_c(g(h_{T_{[CLS]}}), \mathbf{y})]. \quad (3)$$

**Masked Language Modeling.** We apply a masked language modeling loss to the text-image cross encoder to improve the ability to model the relationship between image and text at the token level. 15% of the text tokens are masked in the input. All of these tokens are replaced with the  $[MASK]$  token. For the MLM task (Devlin et al., 2019), the forward operations are executed individually in most VLP models (Chen et al., 2020; Li et al., 2021a), increasing the computational cost of pre-training. In our model, the MLM task utilizes masked text and corresponding images together for denoising, which enhances the interaction between text and images. Since FGR relies heavily on this interaction ability, we propose enhanced training (ET), which integrates the MLM task into the FGR forward operations for positive image-text pairs. Experiments in Section 4.3 show that ET can reduce the computational cost of R2D2 while maintaining the accuracy of the model. For simplicity,  $\mathcal{L}_{\text{MLM}}$  denotes the loss of the MLM task.

### 3.3 TWO-WAY DISTILLATION

Most image-text pre-training data are collected by a semi-automatic program, which may create noisy and inaccurate samples. Imprecise labels are problematic, since they may mislead the model. To address this, we propose target-guided distillation (TgD), a teacher-student paradigm with soft targets. To further improve the generalization performance of the pre-trained model, we introduce feature-guided distillation (FgD), another teacher-student based distillation. For convenience, we call the combination of these two distillations as two-way distillation (TwD).

**Target-guided Distillation.** To decrease the risk of learning from noisy labels, we propose to adopt soft targets generated by momentum-updated encoders. Here, the momentum-updated encoder is the teacher model of distillation, which contains the exponential-moving-average weights. We combine the similarity score  $s(\cdot, \cdot)$  with one-hot labels  $\mathbf{y}(\cdot, \cdot)$  via coefficient  $\alpha$  to generate the final soft targets. Let  $\hat{\mathbf{y}}(I, T)$  and  $\hat{\mathbf{y}}(T, I)$  denote the final soft targets. Taking  $\hat{\mathbf{y}}(I, T)$  as the example, we define it as:

$$\hat{\mathbf{y}}(I, T) = \alpha s(I_m, T) + (1 - \alpha)\mathbf{y}(I, T), \quad (4)$$

where  $I_m$  represents that the images  $I$  are fed into the momentum-updated encoder. During training, we also introduce a queue mechanism and replace  $\hat{\mathbf{y}}(I, T)$  with  $\hat{\mathbf{y}}(I, T_q)$ . In practice, the text queue with a fixed size aims to maintain the recent text representations. We then concatenate the text queue and the text representations of current mini-batch to compute  $s(I_m, T_q)$  and  $\mathbf{y}(I, T_q)$ . Similarly, we perform the same process when constructing  $\hat{\mathbf{y}}(T, I_q)$ .

Considering the effectiveness of features in the queue decreases with increasing time steps, we also maintain a weighted queue  $w$  to mark the reliability of the corresponding position features. Specifically, we decay each element in the queue by a factor of 0.99 per iteration, except for the new incoming item. Further, we replace  $\mathcal{L}_c(\cdot)$  with weighted cross-entropy loss  $\mathcal{L}_w(\cdot)$  in Equation 2. With the target-guided distillation, the  $\mathcal{L}_{\text{GCPR}}^{\text{TgD}}$  is defined as:

$$\mathcal{L}_{\text{GCPR}}^{\text{TgD}} = \frac{1}{2} \mathbb{E}_{(I, T) \sim \mathcal{D}} [\mathcal{L}_w(s(I, T_q), \hat{\mathbf{y}}(I, T_q); w) + \mathcal{L}_w(s(T, I_q), \hat{\mathbf{y}}(T, I_q); w)]. \quad (5)$$

**Feature-guided Distillation.** Similar to TgD, we use a teacher-student paradigm to conduct feature-guided distillation. Taking the text encoder as the example below, the teacher character is the momentum-updated text encoder and the student is the text encoder. Here, the weights of the teacher are updated by all past text encoders via exponential-moving-average. To further improve the capability of the model, we apply a masking strategy to the inputs. In practice, we feed complete inputs into the teacher and masked inputs into the student. Relying on the momentum mechanism, we aim to make the features of the student closer to that of the teacher. Formally, the predicted distributions (*i.e.*,  $\mathcal{P}_t(T)$ ,  $\mathcal{P}_s(T)$ ) of the teacher and the student are defined as follows, respectively.

$$\mathcal{P}_t(T) = \frac{\exp((f_t(T) - \mu)/\tau_t)}{\sum_{i=1}^d \exp((f_t(T)^{(i)} - \mu^{(i)})/\tau_t)}, \quad \mathcal{P}_s(T) = \frac{\exp(f_s(T)/\tau_s)}{\sum_{i=1}^d \exp(f_s(T)^{(i)}/\tau_s)}, \quad (6)$$

where  $f_t(\cdot)$  and  $f_s(\cdot)$  denote the networks of the teacher and the student, respectively. Moreover,  $\mu$  is a momentum-updated mean of  $f_t(\cdot)$ , and  $d$  is the dimension of the features.  $\tau_t$  and  $\tau_s$  are the temperature parameters of the teacher and the student, respectively, which can sharpen the distribution of the features. Note that we do not use  $\mu$  for  $\mathcal{P}_s$  to avoid collapse in feature-guided distillation. We can obtain similar formulations for  $\mathcal{P}_s(I)$  and  $\mathcal{P}_t(I)$ . We perform the feature-guided distillation by the cross-entropy loss, and the loss  $L_{\text{FGD}}$  is defined as:

$$\mathcal{L}_{\text{FGD}} = \frac{1}{2} \mathbb{E}_{(I, T) \sim \mathcal{D}} [\mathcal{L}_c(\mathcal{P}_s(I), \mathcal{P}_t(I)) + \mathcal{L}_c(\mathcal{P}_s(T), \mathcal{P}_t(T))]. \quad (7)$$

Our model is trained with the full objective:

$$\mathcal{L} = \mathcal{L}_{\text{GCPR}}^{\text{TgD}} + \mathcal{L}_{\text{FGR}} + \mathcal{L}_{\text{FGD}} + \mathcal{L}_{\text{MLM}}. \quad (8)$$

## 4 EXPERIMENTS

### 4.1 IMPLEMENTATION DETAILS

The number of transformer layers for the text encoder, and the two cross encoders are 12, 6, and 6, respectively. The text encoder is initialized from RoBERTa-wwm-ext (Cui et al., 2020) while

Table 1: Comparisons with state-of-the-art models on image-text retrieval task.

Dataset	Method	Image-to-Text Retrieval			Text-to-Image Retrieval			R@M
		R@1	R@5	R@10	R@1	R@5	R@10	
Flickr30k-CN	Wukong <sub>ViT-B</sub>	83.9	97.6	99.0	67.6	89.6	94.2	88.7
	Wukong <sub>ViT-L</sub>	92.7	99.1	99.6	77.4	94.5	97.0	93.4
	R2D2 <sub>ViT-B</sub>	92.6	99.1	99.8	78.3	94.6	97.0	93.6
	R2D2 <sub>ViT-L</sub>	<b>95.0</b>	<b>99.7</b>	<b>100.0</b>	<b>83.4</b>	<b>95.9</b>	<b>98.1</b>	<b>95.4</b>
COCO-CN	Wukong <sub>ViT-B</sub>	65.8	90.3	96.6	67.0	91.4	96.7	84.6
	Wukong <sub>ViT-L</sub>	73.3	94.0	98.0	74.0	94.4	98.1	88.6
	R2D2 <sub>ViT-B</sub>	76.1	95.3	98.5	75.1	94.2	98.1	89.6
	R2D2 <sub>ViT-L</sub>	<b>77.4</b>	<b>96.3</b>	<b>98.7</b>	<b>78.1</b>	<b>95.3</b>	<b>98.5</b>	<b>90.7</b>
AIC-ICC	WenLan	45.6	68.0	76.3	34.1	58.9	69.1	58.7
	Wukong <sub>ViT-B</sub>	47.5	70.6	78.6	36.7	36.7	71.7	57.0
	Wukong <sub>ViT-L</sub>	61.6	<b>80.5</b>	<b>86.1</b>	48.6	72.5	80.2	71.6
	R2D2 <sub>ViT-B</sub>	55.9	76.0	82.1	47.1	72.8	80.5	69.1
	R2D2 <sub>ViT-L</sub>	<b>64.4</b>	80.4	85.0	<b>56.8</b>	<b>78.2</b>	<b>83.6</b>	<b>74.7</b>
MUGE	Wukong <sub>ViT-B</sub>	-	-	-	39.2	66.9	77.4	61.2
	Wukong <sub>ViT-L</sub>	-	-	-	52.7	77.9	85.6	72.1
	R2D2 <sub>ViT-B</sub>	-	-	-	47.4	75.1	83.5	68.7
	R2D2 <sub>ViT-L</sub>	-	-	-	<b>53.8</b>	<b>79.6</b>	<b>86.9</b>	<b>73.4</b>
Flickr30k-CNA	R2D2 <sub>ViT-B</sub>	93.2	99.4	99.7	79.6	95.2	97.5	94.1
	R2D2 <sub>ViT-L</sub>	<b>96.5</b>	<b>99.7</b>	<b>100.0</b>	<b>83.6</b>	<b>96.5</b>	<b>98.4</b>	<b>95.8</b>
ICR	R2D2 <sub>ViT-B</sub>	43.4	69.8	78.4	42.2	69.4	77.8	63.5
	R2D2 <sub>ViT-L</sub>	<b>50.6</b>	<b>76.0</b>	<b>82.9</b>	<b>50.1</b>	<b>75.7</b>	<b>82.7</b>	<b>69.6</b>
IQR	R2D2 <sub>ViT-B</sub>	27.9	54.5	64.4	27.4	53.4	63.6	48.5
	R2D2 <sub>ViT-L</sub>	<b>32.9</b>	<b>59.8</b>	<b>69.6</b>	<b>32.6</b>	<b>59.2</b>	<b>68.6</b>	<b>53.8</b>

the two cross encoders are randomly initialized. Following Wukong (Gu et al., 2022), we use the image encoder of 12-layers ViT-Base and 24-layers ViT-Large initialized from CLIP (Radford et al., 2021), and freeze it during pre-training. The resolution of the input image is  $224 \times 224$  in pre-training and fine-tuning. The dimension of the feature vectors of both image and text is 768. We pre-train models with 30 epochs using a batchsize of 32 per GPU. We set  $\tau=0.07$  in Equation 1,  $\tau_s=0.1$ ,  $\tau_t=0.04$  in Equation 6, and  $\alpha=0.4$  in Equation 4. Moreover, the momentum is set as  $m=0.995$ , and the queue size is 36,864. R2D2 is pre-trained with the mixed-precision technique for 4 days using 64 A100 GPUs. The pre-trained model is adapted to four V+L downstream tasks: image-text retrieval, image-text matching, image caption, and text-to-image generation. More details about the downstream datasets and fine-tuning strategy can refer to Appendix C.

## 4.2 COMPARISONS WITH STATE-OF-THE-ART

For both image-to-text retrieval and text-to-image retrieval tasks, we report Recall@1 (R@1), Recall@5 (R@5), Recall@10 (R@10), and Mean Recall (R@M). The results of WenLan (Huo et al., 2021) and Wukong (Gu et al., 2022) are excerpted from their paper. From Table 1, our models outperform state-of-the-art on all datasets. Moreover, R2D2<sub>ViT-L</sub> outperforms R2D2<sub>ViT-B</sub>. These results indicate that our framework is able to learn better fine-grained associations between image and text. We report the results of Flickr30k-CNA on the test set of Flickr30k-CN for a fair comparison. R2D2 fine-tuned on Flickr30k-CNA outperforms that on Flickr30k-CN, since the quality of human-translated Flickr30k-CNA is much higher than that of machine-translated Flickr30k-CN.

Table 2 reports the comparison with existing methods on other V+L understanding tasks. Unlike the image-text retrieval task, there are few datasets for the Chinese image-text matching (ITM) task. Thus, we introduce image-caption matching dataset (ICM) and image-query matching dataset (IQM) for the Chinese ITM task and show the corresponding results. Also, we evaluate Wukong and WenLan on these datasets for the ITM task. We use Area Under Curve (AUC) as the metric. For the image captioning task, fine-tuning is conducted on the training split of AIC-ICC (Wu et al., 2019). We adopt four widely-used evaluation metrics: BLEU, METEOR, ROUGE-L, and CIDEr following WenLan. Table 2 also presents text-to-image generation results on ECommerce-T2I dataset<sup>1</sup>. The metric of

<sup>1</sup><https://tianchi.aliyun.com/muge>

Table 2: Comparison with state-of-the-art models on downstream vision-language tasks.

Method	Image-Text Matching		Image Caption			Text-to-Image Generation	
	AUC (ICM)	AUC (IQM)	BLEU	METEOR	ROUGE-L	CIDEr	FID
UNITER	-	-	62.8	38.7	69.2	199.7	-
WenLan	61.9	57.6	66.1	41.1	71.9	220.7	-
Wukong <sub>VIT-B</sub>	79.2	75.1	66.7	71.2	72.2	224.2	23.7
Wukong <sub>VIT-L</sub>	81.8	78.1	68.9	74.5	72.3	243.1	18.8
R2D2 <sub>VIT-B</sub>	86.5	82.0	67.1	74.7	72.4	227.2	19.2
R2D2 <sub>VIT-L</sub>	<b>88.1</b>	<b>83.6</b>	<b>70.6</b>	<b>76.5</b>	<b>74.3</b>	<b>246.5</b>	<b>15.7</b>

Table 3: Effect of the proposed VLP framework and pre-training dataset. We compare different VLP frameworks on the same pre-training dataset (rows 1-3), and compare same VLP framework on different pre-training datasets (rows 3-4). We report R@M, AUC, CIDEr, and FID for four V+L downstream tasks respectively.

Method	Pre-training Dataset	Image-Text Retrieval		Image-Text Matching		Image Caption	Text-to-Image Generation
		Flick30k-CN	COCO-CN	ICM	IQM	AIC-ICC	ECommerce-T2I
WenLan	Wukong (100M)	79.5	75.1	68.2	65.1	225.6	-
Wukong <sub>VIT-L</sub>	Wukong (100M)	93.4	88.6	81.8	78.1	243.1	18.8
R2D2 <sub>VIT-L</sub>	Wukong (100M)	95.2	90.1	86.5	81.5	245.8	16.4
R2D2 <sub>VIT-L</sub>	ZERO-Corpus (23M)	<b>95.4</b>	<b>90.7</b>	<b>88.1</b>	<b>83.6</b>	<b>246.5</b>	<b>15.7</b>

Fréchet Inception Distance (FID) is reported. Our model achieves state-of-the-art performance on these V+L downstream tasks, showing the superior capabilities of R2D2.

**Effect of the proposed VLP framework.** Different from existing methods that are unfairly evaluated based on different pre-training datasets, we conduct comparative experiments using the same pre-training data in a fairer way. From Table 3, we show the results (rows 1-3) of different VLP frameworks pre-trained on the Wukong dataset and then fine-tuned on six downstream datasets. Our R2D2 outperforms competitors on six datasets when pre-trained with the same 100M dataset, which demonstrates the superiority of our VLP framework.

**Effect of the proposed pre-training dataset.** Similarly, we provide comparison results (rows 3-4 of Table 3) of our R2D2 framework pre-trained on the 100M Wukong dataset and the proposed ZERO-corpus which contains 23M image-text pairs, respectively. R2D2 pre-trained on the 23M ZERO-corpus achieves better results than on the much larger 100M Wukong dataset. This improvement comes from the data quality of our ZERO-corpus dataset, which is filtered by user click-through rate and provides diverse text descriptions along with each image.

### 4.3 ABLATION STUDY

**Effect of Fine-Grained Ranking (FGR).** We conduct ablation studies on the first 10% of ZERO-Corpus. For simplicity, we define R2D2<sub>VIT-L</sub> as R2D2 in the ablation study. We first train a restricted version of R2D2 using only the global contrastive pre-ranking and the two-way distillation strategy. We denote it as PRD2. This restricted setting is conceptually similar to CLIP (Radford et al., 2021) that involves a local contrastive loss. R2D2 outperforms PRD2 on the downstream tasks, indicating the effectiveness of the proposed pre-ranking + ranking framework.

**Effect of Enhanced Training (ET).** In this experiment, we demonstrate the effectiveness of enhanced training. From the third row of Table 4, R2D2 (with ET) performs slightly better than R2D2 w/o ET on all downstream tasks. Another advantage is that R2D2 uses less computational resources than R2D2 w/o ET. R2D2 requires 154.0 GFLOPs and can run at 1.4 iterations per second (Iter/s), while without ET we get 168.8 GFLOPs and 1.1 Iter/s. This indicates that ET is able to both reduce the computational cost and improve the capability of the learning process.

**Effect of Masked Language Modeling (MLM).** Compared to R2D2 w/o MLM, R2D2 obtain better performance on all downstream tasks. MLM allows R2D2 to learn robust representations by masking data during training. These results indicate that MLM is indeed effective for downstream tasks.

**Effect of Two-way Distillation (TwD).** The proposed two-way distillation is composed of target-guided distillation (TgD) and feature-guided distillation (FgD). By analyzing the two components

Table 4: Effect of different components of R2D2. Note that we conduct ablation studies and report the average results on all downstream datasets. R@\* denotes the result for the image-text retrieval task. We report AUC, CIDEr, and FID for image-text matching, image caption and text-to-image generation tasks respectively.

Method	Image-to-Text Retrieval			Text-to-Image Retrieval			R@M	AUC	CIDEr	FID
	R@1	R@5	R@10	R@1	R@5	R@10				
PRD2	53.61	75.13	81.60	43.62	70.79	79.88	66.71	73.49	239.71	19.29
R2D2	<b>64.20</b>	<b>80.63</b>	<b>85.55</b>	<b>56.23</b>	<b>77.81</b>	<b>84.06</b>	<b>74.06</b>	<b>80.51</b>	<b>243.26</b>	<b>17.65</b>
R2D2 w/o ET	63.25	78.49	85.09	55.56	77.42	83.08	73.53	80.27	243.07	17.80
R2D2 w/o MLM	63.92	80.21	85.33	55.91	77.52	83.79	73.68	80.03	242.95	17.89
R2D2 w/o TwD	63.19	79.62	84.91	54.85	76.77	83.58	73.10	80.31	242.83	17.96
R2D2 w/o TgD	63.98	80.59	85.50	56.13	77.29	83.44	73.71	80.42	243.06	17.83
R2D2 w/o FgD	63.47	79.99	85.24	55.05	76.98	83.63	73.29	80.40	242.98	17.90

Table 5: Zero-shot results of different methods on downstream vision-language tasks.

Dataset	Method	Image-to-Text Retrieval			Text-to-Image Retrieval			R@M	AUC
		R@1	R@5	R@10	R@1	R@5	R@10		
Flickr30k-CN	Wukong <sub>vIT-L</sub>	76.1	94.8	97.5	51.7	78.9	86.3	80.9	-
	R2D2 <sub>vIT-L</sub>	70.2	94.1	97.6	55.9	83.5	90.6	<b>82.0</b>	-
COCO-CN	Wukong <sub>vIT-L</sub>	55.2	81.0	90.6	53.4	80.2	90.1	75.1	-
	R2D2 <sub>vIT-L</sub>	58.1	86.8	93.3	55.0	83.1	92.5	<b>78.1</b>	-
MUGE	Wukong <sub>vIT-L</sub>	-	-	-	42.7	69.0	78.0	<b>63.2</b>	-
	R2D2 <sub>vIT-L</sub>	-	-	-	41.0	67.8	76.6	61.8	-
AIC-ICC	Wukong <sub>vIT-L</sub>	18.2	34.5	42.4	8.8	20.3	27.3	20.7	-
	R2D2 <sub>vIT-L</sub>	22.0	36.8	42.1	10.6	21.7	27.0	<b>26.7</b>	-
ICR	Wukong <sub>vIT-L</sub>	35.1	58.2	66.3	33.7	58.0	66.5	53.0	-
	R2D2 <sub>vIT-L</sub>	46.8	72.6	79.3	44.6	70.2	76.0	<b>64.9</b>	-
IQR	Wukong <sub>vIT-L</sub>	26.1	48.9	58.1	24.9	48.1	57.7	44.0	-
	R2D2 <sub>vIT-L</sub>	31.7	58.1	67.1	30.0	54.8	62.1	<b>50.6</b>	-
ICM	R2D2 <sub>vIT-L</sub>	-	-	-	-	-	-	-	<b>87.4</b>
IQM	R2D2 <sub>vIT-L</sub>	-	-	-	-	-	-	-	<b>81.8</b>

of TwD, we see that performing feature alignment is important, since the model w/o FgD shows a more noticeable drop in performance. Although milder, removing TgD also causes a reduction in performance. These results indicate that both components are relevant and TwD is an effective way to improve the generalization performance of the pre-trained model.

#### 4.4 FURTHER EXPERIMENTS

**Zero-shot Tasks.** To demonstrate the generalization performance of our method, we conduct zero-shot transfer experiments. From Table 5, compared with current state-of-the-art Wukong<sub>vIT-L</sub>, our R2D2<sub>vIT-L</sub> achieves comparable or even better performance on Flickr30k-CN, COCO-CN, MUGE, AIC-ICC, ICR, and IQR. We show more results on our proposed downstream tasks *i.e.*, ICM, and IQM. Note that the results of R2D2<sub>vIT-L</sub> on Flickr30k-CNA are the same as that of Flickr30k-CN, since we use the same test set for a fair comparison. In this way, we do not report the results of R2D2<sub>vIT-L</sub> on Flickr30k-CNA.

**Entity-conditioned Image Visualization.** In this experiment, we visualize the attention map of images on COCO-CN. Specifically, we first extract an entity from the Chinese text and calculate the attention score of an image-entity pair. Here, we select the third layer of the text-image cross encoder following (Li et al., 2021a). Figure 2 illustrates the visual explanations of four images over four different entities. It shows that R2D2 learns well to align text with the correct content inside the image. More analysis is shown in Appendix D.





Figure 2: Entity-conditioned image visualization.

## 5 RELATED WORK

### 5.1 VISION-LANGUAGE DATASETS

Chinese vision-language benchmark requires images and high-quality Chinese texts, which are hard to obtain and still rare for the research community’s reach. To this end, existing public datasets (Lan et al., 2017; Li et al., 2019c) use machine translation to adapt their English versions (Chen et al., 2015; Young et al., 2014) to Chinese, but the data quality is sacrificed due to machine translation errors. Newly reported datasets with Chinese texts (Fei et al., 2021; Gu et al., 2022; Lin et al., 2021) are proposed for Chinese VLP. However, they are either not publicly available or lack sufficient downstream tasks. In this paper, we propose a Chinese vision-language benchmark that covers a large-scale pre-training dataset and five high-quality downstream datasets.

### 5.2 VISION-LANGUAGE PRE-TRAINING LEARNING

**Vision-Language Architecture.** The vision-language pre-training architectures can be categorized as: single-stream and dual-stream. Most existing single-stream models (Chen et al., 2020; Li et al., 2020; 2019b; Lu et al., 2020; Qi et al., 2020) concatenate image and text as a single input to model the interactions between image and text within a transformer model (Vaswani et al., 2017). On the other hand, popular dual-stream models (Faghri et al., 2018; Jia et al., 2021; Li et al., 2019a; Lu et al., 2019; Radford et al., 2021; Wei et al., 2021) aim to align image and text into a unified semantic space via contrastive learning. Besides, some works (Li et al., 2021a)(Li et al., 2022) align the individual features of images and texts in a dual-stream architecture, and then fuse the features in a unified semantic space via a single-stream architecture. However, they ignore supervised signals from images. In addition, they use traditional masked language modeling (MLM) and local contrastive learning to conduct pre-training tasks, leading to potential inferior model performance. In this paper, we explore the effective signals via an image-text cross encoder and a text-image cross encoder while also maintaining the bottom dual-stream architecture. Moreover, we improve MLM with enhanced training and apply global contrastive learning to further improve performance.

**Knowledge Distillation.** The general purpose of knowledge distillation is to improve the student model’s performance by simulating the output of the teacher network (Hinton et al., 2014; Xie et al., 2020; Yue et al., 2020). Compared to previous works (Li et al., 2021a; Yue et al., 2020), we propose target-guided distillation with a weighted momentum queue and feature-guided distillation to stabilize the model representations for vision-language pre-training.

## 6 CONCLUSION

In this paper, we introduce a Chinese vision-language benchmark called ZERO and a vision-language method, namely R2D2. ZERO includes a large-scale pre-training dataset and five human-annotated downstream datasets. Our R2D2 adopts a framework of pre-ranking + ranking for cross-modal learning. To alleviate the risk from noise and improve the model capability, we propose a two-way distillation strategy. We achieve state-of-the-art results on eleven downstream datasets of four V+L tasks. A limitation of our method is that it only treats Chinese text, and we will adopt R2D2 to English VLP learning in future work. We will release all datasets and models to promote the development of vision-language learning. We expect that the good cross-modal benchmark and framework will encourage a plethora of engineers to develop more effective methods in specific real-world scenarios.

## REPRODUCIBILITY STATEMENT

The implementation code can be found in supplemental materials. The code platform (PyTorch) we use is public. The pre-training settings are in Section 4.1 of the main paper. In addition, we also provide detailed fine-tuning strategies and experiment parameters in Appendix C.

## REFERENCES

- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European Conference on Computer Vision*, pp. 104–120, 2020.
- Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*, pp. 7–10, 2016.
- Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. Revisiting pre-trained models for chinese natural language processing. In *Conference on Empirical Methods in Natural Language Processing*, pp. 657–668, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Human Language Technology Conference of the NAACL*, pp. 4171–4186, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. Vse++: Improving visual-semantic embeddings with hard negatives. In *Proceedings of the British Machine Vision Conference*, 2018.
- Nanyi Fei, Zhiwu Lu, Yizhao Gao, Guoxing Yang, Yuqi Huo, Jingyuan Wen, Haoyu Lu, Ruihua Song, Xin Gao, Tao Xiang, et al. Wenlan 2.0: Make ai imagine via a multimodal foundation model. *arXiv preprint arXiv:2110.14378*, 2021.
- Jiaxi Gu, Xiaojun Meng, Guansong Lu, Lu Hou, Minzhe Niu, Hang Xu, Xiaodan Liang, Wei Zhang, Xin Jiang, and Chunjing Xu. Wukong: 100 million large-scale chinese cross-modal pre-training dataset and a foundation framework. *arXiv preprint arXiv:2202.06767*, 2022.
- Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *Advances in Neural Information Processing Systems Workshop*, 2014.
- Yuqi Huo, Manli Zhang, Guangzhen Liu, Haoyu Lu, Yizhao Gao, Guoxing Yang, Jingyuan Wen, Heng Zhang, Baogui Xu, Weihao Zheng, et al. Wenlan: Bridging vision and language by large-scale multi-modal pre-training. *arXiv preprint arXiv:2103.06561*, 2021.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pp. 4904–4916, 2021.
- Weiyu Lan, Xirong Li, and Jianfeng Dong. Fluency-guided cross-lingual image captioning. In *Proceedings of the 25th ACM international conference on Multimedia*, pp. 1549–1557, 2017.
- Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 11336–11344, 2020.

- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems*, 34, 2021a.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, 2022.
- Kunpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. Visual semantic reasoning for image-text matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4654–4662, 2019a.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019b.
- Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. Unimo: Towards unified-modal understanding and generation via cross-modal contrastive learning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2592–2607, 2021b.
- Xirong Li, Chaoxi Xu, Xiaoxu Wang, Weiyu Lan, Zhengxiong Jia, Gang Yang, and Jieping Xu. Cocrn for cross-lingual image tagging, captioning, and retrieval. *IEEE Transactions on Multimedia*, 21(9):2347–2360, 2019c.
- Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. M6: A chinese multimodal pretrainer. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3251–3261, 2021.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10437–10446, 2020.
- Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in Neural Information Processing Systems*, 24, 2011.
- Di Qi, Lin Su, Jia Song, Edward Cui, Taroon Bharti, and Arun Sacheti. Imagebert: Cross-modal pre-training with large-scale weak-supervised image-text data. *arXiv preprint arXiv:2001.07966*, 2020.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *Advances in Neural Information Processing Systems Workshop*, 2021.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2556–2565, 2018.

- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15638–15650, 2022.
- Shulong Tan, Meifang Li, Weijie Zhao, Yandan Zheng, Xin Pei, and Ping Li. Multi-task and multi-scene unified ranking model for online advertising. In *2021 IEEE International Conference on Big Data*, pp. 2046–2051, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- Zhe Wang, Liqin Zhao, Biye Jiang, Guorui Zhou, Xiaoqiang Zhu, and Kun Gai. Cold: Towards the next generation of pre-ranking system. In *2nd Workshop on Deep Learning Practice for High-Dimensional Sparse Data with KDD*, 2020.
- Fangyun Wei, Yue Gao, Zhirong Wu, Han Hu, and Stephen Lin. Aligning pretraining for detection via object-level contrastive learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipai Zhou, Guosen Lin, Yanwei Fu, et al. Ai challenger: A large-scale dataset for going deeper in image understanding. In *IEEE International Conference on Multimedia and Expo*, 2019.
- Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10687–10698, 2020.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- Kaiyu Yue, Jiangfan Deng, and Feng Zhou. Matching guided distillation. In *European Conference on Computer Vision*, pp. 312–328, 2020.

## A DETAILS OF ZERO-CORPUS

We illustrate several representative examples of ZERO-Corpus in Figure A. Each sample contains one image and its corresponding attributes. For ease of understanding, we add an English translation version after each Chinese text. There are 3 types of text fields associated with each image: “Title”, “Content” and “ImageQuery”. “Title” and “Content” come from the source webpage containing the

	<p><b>Title:</b> 五大地缝奇观欣赏 (View of the five fissure wonders)</p> <p><b>Content:</b> 奉节地缝亦称天井峡地缝, 全长有37公里, 最大深度有229米, 而最窄处仅2米、而峡谷高度达900米, 形成气势宏伟的“一线天”, 被岩溶专家称作“世界喀斯特峡谷奇中之稀”。峡谷上段较为开阔, 但愈往下愈狭窄, 上部宽10至30米, 谷底宽仅1至30米, 悬崖最深处达300米 (Fengjie fissure, also known as Tianjingxia fissure, has a total length of 37 kilometers and a maximum depth of 229 meters. The narrowest point is only 2 meters and the height of the canyon is 900 meters, forming a magnificent “one-line sky”. The Fengjie fissure is called “the rarest karst canyon in the world” by karst experts. The upper part of the fissure is relatively open, but it becomes narrower as it goes down. The upper part is 10 to 30 meters wide, the bottom of the valley is only 1 to 30 meters wide, and the deepest cliff is 300 meters.)</p> <p><b>ImageQuery:</b> 天井峡地缝 (TianJingXia fissure)</p> <p><b>ImageUrl:</b> <a href="https://n.sinaimg.cn/sinakd20200518ac/219/w490h529/20200518/b5dc-itvqcaz9771277.jpg">https://n.sinaimg.cn/sinakd20200518ac/219/w490h529/20200518/b5dc-itvqcaz9771277.jpg</a></p>
	<p><b>Title:</b> 游览大沼国立公园, 这里山清水秀白云蓝天, 大沼、小沼、莼菜沼三个高山湖皆属于大沼国立公园 (Onuma National Park is with clear waters, white clouds and a blue sky. All three alpine lakes (i.e., Onuma, Konuma, and Uzbekistan) belong to Onuma National Park.)</p> <p><b>Content:</b> 游览大沼国立公园, 这里山清水秀白云蓝天, 大沼、小沼、莼菜沼三个高山湖皆属于大沼国立公园。大沼是由驹岳火山喷发后生成的面积24平方公里的湖泊, 有大小126个岛屿、32湖湾所组成, 这些岛屿由18座桥梁连接的景象十分秀美, 富有欧洲风味的风景。 (Onuma National Park is with clear waters, white clouds and a blue sky. All three alpine lakes (i.e., Onuma, Konuma, and Uzbekistan) belong to Onuma National Park. Onuma is a lake with an area of 24 square kilometers formed after the eruption of the Komagatake volcano. It consists of 126 islands and 32 bays. The view of these islands connected by 18 bridges is very beautiful, full of European-style scenery.)</p> <p><b>ImageQuery:</b> 蓝天山清水秀 (Blue sky, beautiful scenery)</p> <p><b>ImageUrl:</b> <a href="https://img1.qunarzz.com/travel/d5/1510/fb/7336bbdc82ce34f7.jpg_r_720x480x95_8bdcd811.jpg">https://img1.qunarzz.com/travel/d5/1510/fb/7336bbdc82ce34f7.jpg_r_720x480x95_8bdcd811.jpg</a></p>
	<p><b>Title:</b> 英宠物狗戴墨镜穿潮装, 百变时装造型受热捧 (British pet dogs wear sunglasses and trendy clothes. The ever-changing fashion styles are popular.)</p> <p><b>Content:</b> 一只名叫托斯特(Toast)的查尔斯王小猎犬不用拥有专属于自己的漂亮手提包 (A King Charles Spaniel named Toast doesn't have its own fancy handbag.)</p> <p><b>ImageQuery:</b> 戴墨镜的狗, 戴墨镜的人, 狗戴墨镜, 墨镜狗狗, 戴墨镜的狗狗图片, 宠物戴墨镜, 漂亮的宠物狗造型, 宠物戴墨镜和围巾, 橙色的宠物狗, 小猎犬戴墨镜, 舔脚, 时装造型, 狗狗舔脚, 小狗戴墨镜, 狗狗戴墨镜 (Dog with sunglasses)</p> <p><b>ImageUrl:</b> <a href="http://www.people.com.cn/mediafile/pic/20140610/32/6255024078836657304.jpg">http://www.people.com.cn/mediafile/pic/20140610/32/6255024078836657304.jpg</a></p>
	<p><b>Title:</b> 猴子捞月 (Monkey fishing for the moon)</p> <p><b>Content:</b> 猴子捞月 (Monkey fishing for the moon)</p> <p><b>ImageQuery:</b> 猴子捞月 (Monkey fishing for the moon)</p> <p><b>ImageUrl:</b> <a href="https://youer.chazidian.com/uploadfile/image/20121114/3921c63abc1decc364602880cb88c962.jpg">https://youer.chazidian.com/uploadfile/image/20121114/3921c63abc1decc364602880cb88c962.jpg</a></p>
	<p><b>Title:</b> 零基础学绘画-彩铅《紫红色百合花》 (Zero Basic Learning Painting - Color Lead "Fuchsia Lily")</p> <p><b>Content:</b> 最终的效果如图, 能出这样的效果, 真的是一层层涂出来的 (The final view is shown in the figure. To achieve such a view, it is painted layer by layer.)</p> <p><b>ImageQuery:</b> 彩铅百合, 彩铅百合绘画大全 (Color lead lily, color lead lily painting Daquan)</p> <p><b>ImageUrl:</b> <a href="http://5b0988e595225.cdn.sohucs.com/q_70,c_zoom,w_640/images/20180404/aff256df13914b918522acf66094856d.jpeg">http://5b0988e595225.cdn.sohucs.com/q_70,c_zoom,w_640/images/20180404/aff256df13914b918522acf66094856d.jpeg</a></p>
	<p><b>Title:</b> 宾夕法尼亚州立大学 (The Pennsylvania State University)</p> <p><b>Content:</b> 宾夕法尼亚州立大学 (The Pennsylvania State University)</p> <p><b>ImageQuery:</b> 宾夕法尼亚大学校徽, 宾州州立大学 (The school badge of The Pennsylvania State University, The Pennsylvania State University)</p> <p><b>ImageUrl:</b> <a href="http://pic.baik.e.soso.com/p/20140414/bki-20140414115942-921021200.jpg">http://pic.baik.e.soso.com/p/20140414/bki-20140414115942-921021200.jpg</a></p>

Figure A: Examples of ZERO-Corpus.

image, and the latter is also termed as surrounding text in other works. “ImageQuery” is the associated query string for the corresponding image. The average length of “Title”, “Content”, and “ImageQuery” is 5, 18, and 29, respectively. During pre-training, we randomly select one from the 3 text fields to construct an image-text pair, ensuring the data diversity. In this way, the pre-trained model can flexibly fit different text lengths on various downstream tasks. For instance, the text length of AIC-ICC is about 18 words while the text length of MUGE is less than 10 words.

We apply a series of filtering strategies to construct the ZERO-Corpus. For images, we filter out images with both dimensions smaller than 100 pixels or aspect ratio out of the range [1/4, 4]. In addition, we filter images that contain sensitive information, such as sexual, violent scenes, etc. For texts, we remove texts shorter than 2 words or longer than 128 words. Also, we remove texts that contain sensitive as in image filtering. We hope this dataset will bring help to the research community.

## B EXAMPLES OF THE PROPOSED DOWNSTREAM DATASETS

We illustrate examples of ICM, IQM, ICM, and IQR in Figure B. Moreover, Figure C highlights some cases of the difference between Flickr30k-CN and our proposed Flickr30k-CNA.



Figure B: Image-text examples of ICM, IQM, ICR and IQR from left to right.

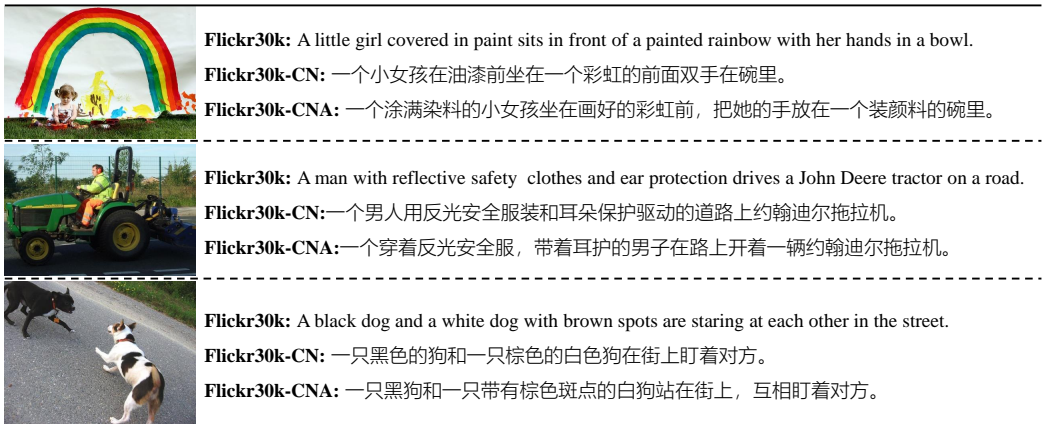


Figure C: Comparisons of Flickr30k, Flickr30k-CN and our proposed Flickr30k-CNA.

## C DETAILS OF FINE-TUNING STRATEGY

**Fine-tuning Strategy of Image-Text retrieval.** We jointly optimize the GCPR loss (Equation 2) and the FGR loss (Equation 3). We extract the individual features of images and texts via our dual-stream encoder and compute the similarity of all image-text pairs. Then we take the top-K candidates and use two cross encoders to further calculate the corresponding similarity scores for ranking during



inference. We use a mean operation for the outputs of the two cross encoders. Here, we adjust the K on different downstream datasets. We fine-tune the pre-trained model with 20 epochs on 7 downstream datasets, including Flickr30k-CN, COCO-CN, AIC-ICC, MUGE, ICR, IQR, and Flickr30k-CNA. K is set as 128, 256, 32, 64, 64, 64, 128, respectively. The batchsize is 32 and the learning rate is  $1e^{-5}$ .

For both image-to-text retrieval (TR) and text-to-image retrieval (IR) tasks, we report Recall@1 (R@1), Recall@5 (R@5), Recall@10 (R@10), and Mean Recall (R@M). For AIC-ICC and MUGE, we report their results on the validation sets, since their test sets are not released. For ICR and IQR, we also report the results on the validation sets in this paper, since we use the corresponding test sets to build a leaderboard. The test set of Flickr30k-CNA is also added to the leaderboard. For Flickr30k-CNA, we show the performance on the test set of Flickr30k-CN for a fair comparison in the main paper. For the remaining downstream datasets, we report the results on the test sets. Following (Gu et al., 2022), we select the first 10,000 images with the corresponding 50,000 texts when testing on AIC-ICC. In particular, we only provide IR scores on MUGE since it only has IR settings.

**Fine-tuning Strategy of Image-Text Matching.** This task predicts whether an image-text pair is matched or not. During fine-tuning, we only apply the FGR loss (Equation 3). We fine-tune the models with 5 epochs using a batchsize of 64. The initial learning rate is  $1e^{-5}$ . Additionally, we report the results on the validation sets of ICM and IQM.

**Fine-tuning Strategy of Image Caption.** Given an image, the goal of the image-caption task is to generate a caption to describe the image. Similar to Transformer(Vaswani et al., 2017), the image-caption model consists of an encoder and a decoder, where the encoder aims to extract the embedding of the given image and the decoder generates tokens of the caption. In specific, we use the image encoder and the text-image cross encoder of R2D2 to initialize the image-caption encoder and decoder, respectively. We fine-tune the image-caption model on the training split of AIC-ICC Wu et al. (2019) with 20 epochs. The batchsize is 128 and the learning rate is  $1e^{-4}$ .

**Fine-tuning Strategy of Text-to-Image Generation.** Text-to-image generation requires the model to generate an image corresponding to the input text. Following DALL-E 2 (Ramesh et al., 2022), we build a generation model, including a CLIP-based module, a prior module and a decoder module. Specifically, the dual-stream weights of R2D2 are used to initialize the CLIP-based module. We fine-tune the CLIP-based module and fix it in the next step. Then, we train the prior module to generate image embeddings for given texts. Finally, we fix two former modules and train a diffusion decoder to invert the image embeddings to generate images. All three components of the generation model are fine-tuned on the ECommerce-T2I dataset with 20 epochs, respectively. The batchsize is 16 and the learning rate is  $1e^{-4}$ . The statistics of all downstream datasets are present in Table A.

Table A: Statistics and comparison of different datasets.

Dataset	#Image/#Text		
	Train	Val	Test
Flickr30k-CN (Lan et al., 2017)	29.7K/148.9K	1K/5K	1K/5K
COCO-CN (Li et al., 2019c)	18.3K/20K	1K/1.1K	1K/1K
AIC-ICC (Wu et al., 2019)	210K/1.05M	30K/150K	30K/150K
MUGE (Lin et al., 2021)	129.4K/248.8K	29.8K/5K	30.4K/5K
ICM	320K/320K	40K/40K	40K/40K
IQM	320K/320K	40K/40K	40K/40K
ICR	160K/160K	20K/20K	20K/20K
IQR	160K/160K	20K/20K	20K/20K
Flickr30k-CNA	29.7K/148.9K	1K/5K	1K/5K
ECommerce-T2I	9K/9K	5K/5K	5K/5K

## D MORE CASES ABOUT ENTITY-CONDITIONED IMAGE VISUALIZATION

In this experiment, we provide more cases of image visualization given an entity. From Figure D, our proposed R2D2 is able to align the entities with patches of images. Especially, R2D2 has the ability to capture the salient areas when given an image with complex backgrounds, such as the images of

“A train” and “A bull”. Our R2D2 also precisely locates different objects within the same image, as shown in the images of “A cup of yogurt” and “Banana”. Moreover, we analyze some bad cases in Figure E. We find that the attention score is disturbed when two adjacent entities are present in an image. This phenomenon is particularly evident for objects with similar colors or categories.

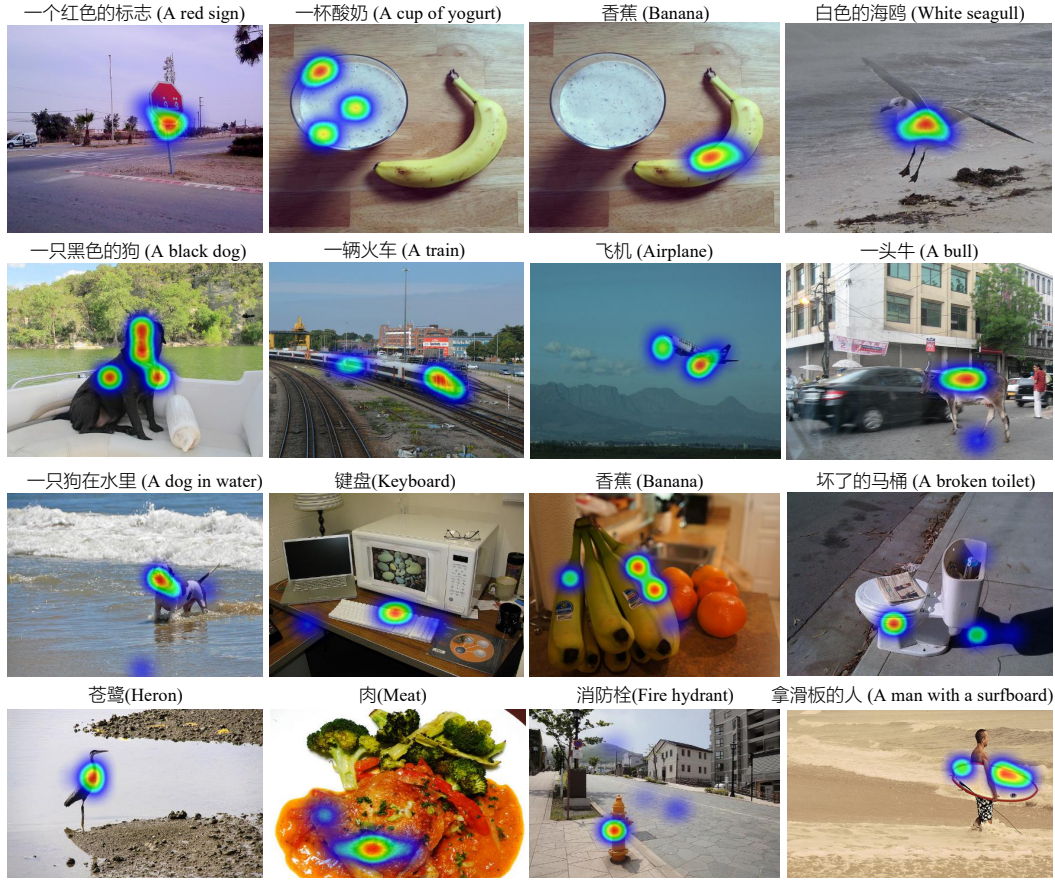


Figure D: More Examples of entity-conditioned image visualization.



Figure E: Bad cases of entity-conditioned image visualization.