

# CCMB: A Large-scale Chinese Cross-modal Benchmark

Chunyu Xie  
360 AI Research  
Beijing, China  
xiechunyu@360.cn

Heng Cai\*  
360 AI Research  
Beijing, China  
caiheng1@360.cn

Jincheng Li\*  
360 AI Research  
Beijing, China  
lijincheng@360.cn

Fanjing Kong  
360 AI Research  
Beijing, China  
kongfanjing@360.cn

Xiaoyu Wu  
360 AI Research  
Beijing, China  
wuxiaoyu1@360.cn

Jianfei Song  
360 AI Research  
Beijing, China  
songjianfei@360.cn

Henrique Morimitsu  
360 AI Research  
Tsinghua University  
Beijing, China  
henrique.morimitsu@mail.tsinghua.edu.cn

Lin Yao  
360 AI Research  
Beijing, China  
yaolin@360.cn

Dexin Wang  
360 Search Department  
Beijing, China  
wangdexin@360.cn

Xiangzheng Zhang  
360 Search Department  
Beijing, China  
zhangxiangzheng@360.cn

Dawei Leng  
360 AI Research  
Beijing, China  
lengdawei@360.cn

Baochang Zhang  
Beihang University  
Beijing, China  
bczhang@buaa.edu.cn

Xiangyang Ji  
Tsinghua University  
Beijing, China  
xyji@tsinghua.edu.cn

Yafeng Deng<sup>†</sup>  
360 AI Research  
Tsinghua University  
Beijing, China  
dengyafeng@gmail.com

## 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 work, we build a large-scale high-quality Chinese Cross-Modal Benchmark named CCMB for the research community, which contains the currently largest public pre-training dataset Zero and five human-annotated fine-tuning datasets for downstream tasks. Zero contains 250 million images paired with 750 million text descriptions, plus two of the five fine-tuning datasets are also currently the largest ones for Chinese cross-modal downstream tasks. Along with the CCMB,

we also develop a VLP framework named R2D2, applying a pre-Ranking + Ranking strategy to learn powerful vision-language representations and a two-way distillation method (i.e., target-guided Distillation and feature-guided Distillation) to further enhance the learning capability. With the Zero and the R2D2 VLP framework, we achieve state-of-the-art performance on twelve downstream datasets from five broad categories of tasks including image-text retrieval, image-text matching, image caption, text-to-image generation, and zero-shot image classification. The datasets, models, and codes are available at <https://github.com/yuxie11/R2D2>

## CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence.**

## KEYWORDS

large-scale datasets, vision-language pre-training

## ACM Reference Format:

Chunyu Xie, Heng Cai, Jincheng Li, Fanjing Kong, Xiaoyu Wu, Jianfei Song, Henrique Morimitsu, Lin Yao, Dexin Wang, Xiangzheng Zhang, Dawei Leng, Baochang Zhang, Xiangyang Ji, and Yafeng Deng. 2023. CCMB: A Large-scale Chinese Cross-modal Benchmark. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3581783.3611877>

\*Both are second authors

<sup>†</sup>Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*MM '23, October 29–November 3, 2023, Ottawa, ON, Canada*

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 979-8-4007-0108-5/23/10...\$15.00  
<https://doi.org/10.1145/3581783.3611877>

## 1 INTRODUCTION

Vision-language pre-training (VLP) mainly learns the semantic correspondence between vision and natural language. Previous works [4, 21, 34, 38] explore the VLP model and achieve significant improvement on various vision-language (V+L) tasks. These methods are supported by massive data [32], excellent architectures such as Transformer [37], and cross-modal models such as CLIP [30].

There are plenty of available benchmarks with English corpus, such as Conceptual Captions [33], SBU Captions [28], and LAION [32]. Differently, large-scale pre-training datasets and downstream datasets with Chinese corpus are relatively few. M6-Corpus [24] is a multi-modal pre-training dataset in Chinese but not publicly available. BriVL (also called WenLan) [11] constructs a vision-language dataset called WSCD, but only releases 5M image-text pairs. Wukong [12] is a newly published pre-training dataset with 100M image-text pairs. Most existing downstream Chinese datasets mainly focus on retrieval tasks, such as Flickr30k-CN [16] and COCO-CN [23], which are not sufficient for a complete evaluation of VLP models. Besides, Flickr30k-CN tries to translate English cross-modal downstream datasets into Chinese, which, however, fails to cover Chinese idioms and often causes translation errors.

In this paper, we introduce a large-scale Chinese cross-modal benchmark called CCMB, including a pre-training dataset (Zero) and five downstream datasets. Specifically, Zero consists of 250 million images and 750 million descriptive texts, which is the largest public Chinese V+L pre-training dataset. Zero is collected from the search engine with images and corresponding textual descriptions, by filtering from 5 billion image-text data by user click-through rate (CTR). Compared to existing pre-training datasets, Zero is high-quality due to the user CTR filtering method and the diverse textual information for each image. Table 1 shows an overview of V+L pre-training datasets. Together with the pre-training dataset, we provide 5 high-quality human-annotated downstream datasets. To the best of our knowledge, two of them are the first proposed datasets for the Chinese image-text matching task, which is also important for evaluating VLP models. They are also the largest Chinese V+L downstream datasets. 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 [16]. The statistics of the public and our proposed downstream datasets are shown in Table 2.

From the perspective of cross-modal learning, existing methods are mainly categorized as single-stream and dual-stream. Most single-stream methods (e.g., [3, 22, 29]) employ an extra object detector to extract the patch embedding and then align patches and words. As illustrated in [19], object detectors are annotation-expensive and computing-expensive, because they require bounding box annotations during pre-training and high-resolution (such as  $600 \times 1000$ ) images during inference. On the other hand, for dual-stream architectures (e.g., [11, 30, 39]), it is non-trivial to model the fine-grained associations between image and text, since the corresponding representations reside in their own semantic space. Some works [18, 19, 34] omit the object detection module and combine dual-stream architecture with single-stream architecture, showing powerful performance on multi-modal downstream tasks.

Inspired by this line of work, we introduce a VLP framework called R2D2, a combination architecture of dual-stream and single-stream. We apply global contrastive pre-ranking to obtain image-text representations and fine-grained ranking to further improve model performance. Besides, we introduce a two-way distillation method into the model, consisting of target-guided distillation and feature-guided distillation. The target-guided distillation increases the robustness when learning from noisy labels, while feature-guided distillation aims to improve the generalization performance. We apply masked language modeling with enhanced training, which improves the capability of the model while reducing the training cost. To summarize, our main contributions are as follows:

- We construct the largest public Chinese vision-language pre-training dataset, containing 250 million images and 750 million corresponding texts. It is high-quality due to the filtering method by user CTR and the diverse textual information for each image. We provide five human-annotated cross-modal downstream datasets, two of which are currently the largest Chinese vision-language downstream datasets.
- We introduce a vision-language pre-training framework named R2D2 for cross-modal learning. Our proposed method achieves state-of-the-art performance on twelve downstream datasets from five broad categories of vision-language tasks, showing the superior ability of our pre-trained model.

## 2 RELATED WORK

### 2.1 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 [16, 23] use machine translation to adapt their English versions [2, 43] to Chinese, but the data quality is sacrificed due to machine translation errors. Newly reported datasets with Chinese texts [11, 12, 24] 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 downstream datasets.

### 2.2 Vision-Language Pre-training Learning

The vision-language pre-training architectures can be categorized as: single-stream and dual-stream. Most existing single-stream models [3, 17, 21, 27, 29] concatenate image and text as a single input to model the interactions between image and text within a transformer model [37]. On the other hand, popular dual-stream models [10, 14, 20, 26, 30, 39, 41] aim to align image and text into a unified semantic space via contrastive learning. Besides, some works [18, 19, 34] 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. These works show that a combined architecture of dual-stream and single-stream achieves better performance than only one. R2D2 explores 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 masked language modeling with enhanced training and propose a two-way distillation to stabilize the model representations for vision-language pre-training.

**Table 1: Statistics of the vision-language pre-training datasets. The details of Zero can refer to Section 3.1.**

Dataset	Language	Availability	#Image	#Text
Visual Genome [15]	English	Yes	108K	5.4M
SBU Captions [28]	English	Yes	875K	875K
CC3M [33]	English	Yes	3.1M	3.1M
CC12M [1]	English	Yes	12M	12M
RedCaps [7]	English	Yes	12M	12M
WIT [35]	Multilingual	Yes	11.5M	37.6M
YFCC100M [36]	English	Yes	100M	200M
LAION-400M [32]	English	Yes	400M	400M
WSCD [11]	Chinese	Yes	5M	5M
M6-Corpus [24]	Chinese	No	60.5M	60.5M
Wukong [12]	Chinese	Yes	100M	100M
Zero	Chinese	Yes	250M	750M

### 3 CCMB

#### 3.1 Pre-training Dataset

Existing public pre-training datasets suffer from two 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.

To this end, we first collect 5 billion image-text data from an image search engine. We try to mitigate the noise of these image-text pairs via user click-through rate (CTR) and obtain about 250 million images and 750 million corresponding texts, namely Zero. In other words, each image in Zero is with about 3 textual descriptions, *i.e.*, “Title”, “Content”, and “ImageQuery”. “Title” and “Content” come from the source webpage containing the image. “Title” is the title of the webpage, and “Content” represents the surrounding text of the image in the webpage. “ImageQuery” is the user search query for the corresponding image. The average length of “Title”, “Content”, and “ImageQuery” is 18, 29, and 5, respectively. We show an example in Figure 1 and more examples in Appendix.

**How to remove the irrelevant content?** We apply a series of filtering strategies to construct the Zero. For images, we filter out images with dimensions smaller than 100 pixels or aspect ratio out of the range [1/4, 4]. We then 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. Moreover, we remove image-text pairs that contain sensitive words and personal names in the text.

**Why use CTR instead of random selection?** After collecting 5 billion image-text data, we can randomly select a part of them to conduct pre-training experiments under the consideration of computational resources and time cost. However, random selection brings noises in pre-training, which may degrade model performance. To address this, we use an inherent metric CTR in search engine data. The CTR indicates the number of times that users click on an image for a given text. We observe that image-text pairs with low CTR are irrelevant in most cases. That is, an image-text pair

**Table 2: Statistics of the vision-language downstream datasets. Our downstream datasets can refer to Section 3.2.**

Dataset	Annotation	Image-Text Pairs		
		Train	Val	Test
Flickr30k-CN [16]	Machine Translation	29K	1K	1K
COCO-CN [23]	Human Annotation	18K	1K	1K
AIC-ICC [40]	Human Annotation	210K	30K	30K
MUGE [24]	-	129K	29K	30K
ECommerce-T2I [24]	-	9K	5K	5K
Flickr30k-CNA	Human Annotation	29K	1K	1K
ICR	Human Annotation	160K	20K	20K
IQR	Human Annotation	160K	20K	20K
ICM	Human Annotation	320K	40K	40K
IQM	Human Annotation	320K	40K	40K

with high CTR is strongly correlated due to user interaction. We then rank all 5 billion image-text data by CTR and filter the top 250 million images and the corresponding texts as Zero.

**Why use three types of text instead of one?** Each image has three kinds of textual descriptions in raw data. During pre-training, we construct an image-text pair per iteration by randomly selecting one of them. Besides, we also conduct a variety of combinations of different types of textual information, such as without “Title”. We find the best mode is to use all three types of text instead of other combinations, which brings more data diversity and potentially improves the model performance.

**Why is Zero large-scale and high-quality?** As shown in Table 1, the proposed Zero is a large-scale Chinese pre-training dataset with 250 million images and 750 million corresponding texts. To the best of our knowledge, Zero is the largest Chinese pre-training image-text dataset. Relying on the CTR filtering method and the diverse textual descriptions, Zero is high-quality because a small amount of pre-training data (about 10% of Zero, 23M) surpasses the previous state-of-the-art [12], which uses 100M image-text pairs. More details can be found in Section 5.3.

#### 3.2 Downstream Dataset

We perform human annotation in downstream datasets, *i.e.*, ICM, IQM, ICR, IQR, and Flickr30k-CNA. For the first four datasets, the selection strategy of image-text pairs is the same as the collection of the pre-training dataset Zero. We divide the training set, validation set, and test set with a ratio of 8:1:1. We check that the downstream data do not appear in the pre-training dataset via the hash value of images. Then, 15 human annotators carefully label the image-text pairs. Specifically, the human annotators verify whether an image-text pair is relevant. That is, the human annotators mark them as positive or negative pairs until a pre-defined data size is reached. Note that we do not rewrite the caption or query for each image in these downstream datasets. For Flickr30k-CNA, we gather 6 professional English and Chinese linguists to meticulously translate all data of Flickr30k [43] and double-check each sentence. The details of each dataset are as follows.

**Image-Caption Matching Dataset (ICM).** ICM is collected for the image-text matching task. The image-text matching task is

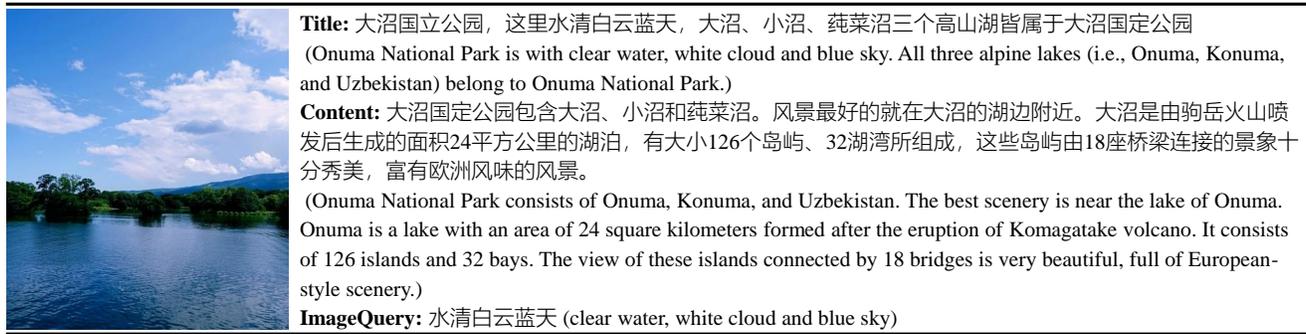


Figure 1: An example of Zero. More samples can be found in Appendix.

a binary classification task, aiming to predict whether an image-text pair is matched. Each image has a corresponding caption text. We first select image-text pairs beyond the 5 billion data. 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. Similarly, IQM contains 200,000 positive cases and 200,000 negative cases. ICM and IQM are the largest Chinese vision-language downstream datasets.

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

**Image-Query Retrieval Dataset (IQR).** IQR is also proposed for the image-text retrieval task. We collect 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 Figure B in Appendix.

**Flickr30k-CNA Dataset.** Former Flickr30k-CN [16] translates the training and validation sets of Flickr30k [43] 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 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.

## 4 METHODOLOGY

### 4.1 Model Architecture

From Figure 2, the architecture contains a text encoder, an image encoder, and two cross encoders. The text encoder is a BERT [8] using the tokenizer of RoBERTa-wwm-ext [5]. For the image encoder, we adopt the Vision Transformer (ViT) [9]. The two cross encoders are multi-layer transformers. 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.

### 4.2 Pre-training Objectives

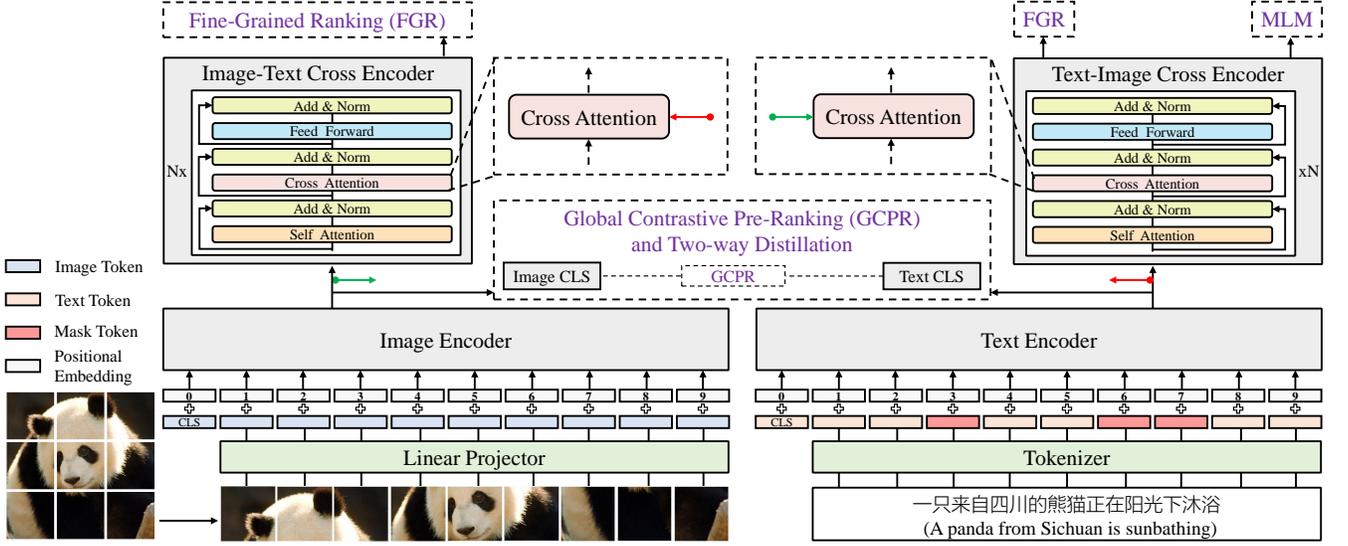
We jointly optimize R2D2 with the following four objectives. To fully explore the matching relationship between image and text pairs, we design a mechanism of pre-ranking + ranking, named global contrastive pre-ranking (GCPR) and fine-grained ranking (FGR). To further enhance the capability of the model, we propose a two-way distillation (TwD) strategy consisting of target-guided distillation and feature-guided distillation. We adopt masked language modeling (MLM) with enhanced training (ET) to efficiently learn the representation of cross-modal models. We conduct the ablation study to verify the effectiveness of each pre-training strategy in Section 5.3.

**Global Contrastive Pre-Ranking.** Our global contrastive pre-ranking method is similar to that of CLIP [30], aiming to align the representation of multi-modal data (e.g., paired image and text). The open-source CLIP implementation [30] only performs back-propagation of the contrastive loss from the local GPU, where negative samples are not fully utilized. Instead, We back-propagate the gradients across all  $k$  GPUs. Inspired by MoCo [13], we also introduce a queue mechanism. In practice, two queues with a fixed size  $M$  aim to maintain the recent image and text representations from the momentum-updated encoders, respectively. 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:

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

$$s(T_i, I_i) = \frac{\exp(\text{sim}(T_i, I_i)/\tau)}{\sum_{j=1}^{n \times k + M} \exp(\text{sim}(T_i, I_j)/\tau)},$$
(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. 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. Let  $\mathcal{D}$  denote the training data and  $y(\cdot, \cdot)$



**Figure 2: 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), fine-grained ranking (FGR), two-way distillation (TwD), and mask language modeling (MLM) with enhanced training (ET) as pre-training objectives.**

denote the ground-truth one-hot label. The global contrastive pre-ranking (GCPR) loss is calculated by the weighted cross-entropy loss  $\mathcal{L}_w(\cdot)$ , as shown in Equation (2).

$$\begin{aligned}\mathcal{L}_{i2t}^w(I, T) &= \mathcal{L}_w(s(I, T), y(I, T); w), \\ \mathcal{L}_{t2i}^w(T, I) &= \mathcal{L}_w(s(T, I), y(T, I); w), \\ \mathcal{L}_{GCPR} &= \frac{1}{2} \mathbf{E}_{(I, T) \sim \mathcal{D}} [\mathcal{L}_{i2t}^w(I, T) + \mathcal{L}_{t2i}^w(T, I)].\end{aligned}\quad (2)$$

**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 be specific, this is a binary classification task, aiming to predict whether an image-text pair 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 by the cross-entropy loss  $\mathcal{L}_c(\cdot)$  as:

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

The selection strategy of negative pairs is in Appendix.

**Two-way Distillation.** Relying on the momentum-updated encoders in contrastive learning, we introduce target-guided distillation (TgD) to decrease the risk of learning from noisy labels, and feature-guided distillation (FgD) to improve the generalization performance of the pre-trained model. We conduct target-guided

distillation to learn from pseudo-targets generated by the momentum model following ALBEF [19]. In practice, we replace the target in Equation (2) with the pseudo-targets. More details about the training process of TgD can be found in Appendix. Besides, target-guided distillation and feature-guided distillation both adopt a teacher-student paradigm. For convenience, we call the combination of TgD and FgD as two-way distillation (TwD).

Below are the details of FgD. 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.

$$\begin{aligned}\mathcal{P}_t(T) &= \frac{\exp((f_t(T) - \mu)/\tau_t)}{\sum_{i=1}^d \exp((f_t(T)^{(i)} - \mu^{(i)})/\tau_t)}, \\ \mathcal{P}_s(T) &= \frac{\exp(f_s(T)/\tau_s)}{\sum_{i=1}^d \exp(f_s(T)^{(i)}/\tau_s)},\end{aligned}\quad (4)$$

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:

$$\mathcal{L}_{\text{FGD}} = \frac{1}{2} \mathbb{E}_{(I,T) \sim \mathcal{D}} [\mathcal{L}_c(\mathcal{P}_s(I), \mathcal{P}_t(I)) + \mathcal{L}_c(\mathcal{P}_s(T), \mathcal{P}_t(T))]. \quad (5)$$

Through experiments in Section 5.3, we observe a noticeable performance gain by performing FGD.

**Masked Language Modeling with Enhanced Training.** 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. The forward operations of MLM [8] and FGR are executed individually in most VLP models [3, 19, 34], 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), applying FGR and MLM loss on text tokens with masking simultaneously. Experiments in Section 5.3 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 with enhanced training. Our model is trained with the full objective:

$$\mathcal{L} = \mathcal{L}_{\text{GCPR}}^{\text{TgD}} + \mathcal{L}_{\text{FGD}} + \mathcal{L}_{\text{FGR}} + \mathcal{L}_{\text{MLM}}. \quad (6)$$

## 5 EXPERIMENTS

### 5.1 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-wvm-ext [5] while the two cross encoders are randomly initialized. Following Wukong [12], we use the image encoder of 12-layers ViT-Base and 24-layers ViT-Large initialized from CLIP [30], and freeze it during pre-training. The resolution of the input image is 224×224 in pre-training and fine-tuning. The dimension of the feature vectors of both image and text is 768. We pre-train models with 15 epochs using a batch size of 4096 on 128 NVIDIA A100 GPUs.  $\tau$  in Equation 1 is initialized to 0.07 and can be learned during the training process. We set  $\tau_s = 0.1$  and  $\tau_t = 0.04$  in Equation 4. Moreover, the momentum is set as  $m = 0.995$ , and the queue size is 36,864. We adopt the Adam optimizer and the cosine learning rate schedule with a linear warmup [25]. The pre-trained model is adapted to five vision-language downstream tasks: image-text retrieval, image-text matching, image caption, text-to-image generation, and zero-shot image classification. More details can refer to Appendix.

### 5.2 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 BriVL [11] and Wukong [12] are excerpted from their paper. Wukong reproduces the CLIP-style [30] and FILIP-style [42] models. Their results are also included. From Table 3, our models outperform state-of-the-art on all datasets.

**Table 3: Comparisons with state-of-the-art models on image-text retrieval task. CNA represents our Flickr30k-CNA.**

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	CLIP <sub>ViT-B</sub>	87.1	97.7	98.8	69.0	90.3	95.0	89.7
	CLIP <sub>ViT-L</sub> [30]	91.6	99.1	99.7	77.3	94.4	97.2	93.2
	FILIP <sub>ViT-B</sub>	72.1	91.3	95.8	57.5	84.3	90.6	81.9
	FILIP <sub>ViT-L</sub> [42]	90.6	98.8	99.6	76.9	94.9	97.4	93.0
	Wukong <sub>ViT-B</sub>	83.9	97.6	99.0	67.6	89.6	94.2	88.7
	Wukong <sub>ViT-L</sub> [12]	92.7	99.1	99.6	77.4	94.5	97.0	93.4
	R2D2 <sub>ViT-B</sub>	93.2	99.2	99.8	79.2	95.2	97.3	94.0
	R2D2 <sub>ViT-L</sub>	<b>95.6</b>	<b>99.8</b>	<b>100.0</b>	<b>84.4</b>	<b>96.7</b>	<b>98.4</b>	<b>95.8</b>
COCO-CN	CLIP <sub>ViT-B</sub>	68.7	93.6	97.5	68.9	93.3	97.3	86.6
	CLIP <sub>ViT-L</sub> [30]	68.3	93.0	97.3	70.1	92.2	96.4	86.2
	FILIP <sub>ViT-B</sub>	52.7	81.3	88.3	56.2	86.8	94.3	76.6
	FILIP <sub>ViT-L</sub> [42]	69.1	91.3	96.9	72.2	92.4	97.2	86.5
	Wukong <sub>ViT-B</sub>	65.8	90.3	96.6	67.0	91.4	96.7	84.6
	Wukong <sub>ViT-L</sub> [12]	73.3	94.0	98.0	74.0	94.4	98.1	88.6
	R2D2 <sub>ViT-B</sub>	78.1	96.2	98.6	76.0	94.9	98.3	90.3
	R2D2 <sub>ViT-L</sub>	<b>79.3</b>	<b>97.1</b>	<b>98.7</b>	<b>79.1</b>	<b>96.5</b>	<b>98.9</b>	<b>91.6</b>
AIC-ICC	BriVL [11]	45.6	68.0	76.3	34.1	58.9	69.1	58.7
	CLIP <sub>ViT-B</sub>	50.5	73.0	80.2	38.1	63.7	73.3	63.1
	CLIP <sub>ViT-L</sub> [30]	59.1	79.5	85.2	46.2	70.7	78.6	69.9
	FILIP <sub>ViT-B</sub>	42.5	67.2	76.0	32.9	58.4	68.8	57.6
	FILIP <sub>ViT-L</sub> [42]	54.1	75.8	82.8	44.9	69.0	77.5	67.4
	Wukong <sub>ViT-B</sub>	47.5	70.6	78.6	36.7	66.7	71.7	57.0
	Wukong <sub>ViT-L</sub> [12]	61.6	<b>80.5</b>	<b>86.1</b>	48.6	72.5	80.2	71.6
	R2D2 <sub>ViT-B</sub>	56.8	76.2	82.1	47.6	72.8	80.2	69.3
R2D2 <sub>ViT-L</sub>	<b>65.4</b>	80.3	84.7	<b>57.3</b>	<b>78.1</b>	<b>83.0</b>	<b>74.8</b>	
MUGE	CLIP <sub>ViT-B</sub>	-	-	-	43.5	71.7	80.6	65.3
	CLIP <sub>ViT-L</sub> [30]	-	-	-	50.1	76.9	84.9	70.6
	FILIP <sub>ViT-B</sub>	-	-	-	30.6	58.2	70.2	53.0
	FILIP <sub>ViT-L</sub> [42]	-	-	-	43.5	71.5	80.9	65.3
	Wukong <sub>ViT-B</sub>	-	-	-	39.2	66.9	77.4	61.2
	Wukong <sub>ViT-L</sub> [12]	-	-	-	52.7	77.9	85.6	72.1
	R2D2 <sub>ViT-B</sub>	-	-	-	53.4	78.1	86.0	72.5
	R2D2 <sub>ViT-L</sub>	-	-	-	<b>60.1</b>	<b>82.9</b>	<b>89.4</b>	<b>77.5</b>
CNA	R2D2 <sub>ViT-B</sub>	93.6	99.5	99.8	80.5	95.6	97.7	94.5
	R2D2 <sub>ViT-L</sub>	<b>96.9</b>	<b>99.8</b>	<b>100.0</b>	<b>84.9</b>	<b>97.0</b>	<b>98.6</b>	<b>96.2</b>
ICR	R2D2 <sub>ViT-B</sub>	53.4	75.4	83.4	52.1	73.3	82.0	69.9
	R2D2 <sub>ViT-L</sub>	<b>61.5</b>	<b>82.9</b>	<b>87.7</b>	<b>60.7</b>	<b>82.0</b>	<b>86.9</b>	<b>77.0</b>
IQM	R2D2 <sub>ViT-B</sub>	37.0	62.1	70.9	35.8	61.2	70.5	56.3
	R2D2 <sub>ViT-L</sub>	<b>41.9</b>	<b>67.8</b>	<b>75.9</b>	<b>41.3</b>	<b>67.6</b>	<b>75.4</b>	<b>61.7</b>

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 4 reports the comparison with existing methods on other V+L downstream 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 BriVL 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 [40]. We adopt four widely-used evaluation metrics: BLEU, METEOR, ROUGE-L, and

**Table 4: Comparison with state-of-the-art models on downstream vision-language tasks.**

Method	Image-Text Matching		Image Caption				Text-to-Image Generation	Zero-shot Image Classification	
	AUC (ICM)	AUC (IQM)	BLEU	METEOR	ROUGE-L	CIDEr	FID	Top-1 Acc.	Top-5 Acc.
BriVL [11]	61.9	57.6	66.1	41.1	71.9	220.7	-	24.3	56.8
Wukong <sub>ViT-B</sub>	79.2	75.1	66.7	71.2	72.2	224.2	23.7	49.1	74.2
Wukong <sub>ViT-L</sub> [12]	81.8	78.1	68.9	74.5	72.3	243.1	18.8	55.0	80.5
R2D2 <sub>ViT-B</sub>	88.6	84.9	68.3	76.3	73.2	230.2	18.9	50.6	78.1
R2D2 <sub>ViT-L</sub>	<b>90.6</b>	<b>86.7</b>	<b>71.8</b>	<b>78.2</b>	<b>75.3</b>	<b>247.9</b>	<b>14.4</b>	<b>56.9</b>	<b>83.3</b>

**Table 5: Comparison with state-of-the-arts which combine dual-stream and single-stream architectures. Classification represents zero-shot image classification. We report R@M, AUC, CIDEr, FID, and Top-1 accuracy for five V+L downstream tasks respectively.**

Method	Image-Text Retrieval		Image-Text Matching		Image Caption	Text-to-Image Generation	Classification
	Flick30k-CN	COCO-CN	ICM	IQM	AIC-ICC	ECommerce-T2I	ImageNet
ALBEF[19]	90.1	84.9	79.5	74.7	226.0	21.4	35.9
FLAVA[34]	91.4	85.1	80.1	75.6	226.2	21.0	37.2
R2D2 <sub>ViT-B</sub>	<b>92.2</b>	<b>86.3</b>	<b>81.1</b>	<b>76.3</b>	<b>226.8</b>	<b>20.9</b>	<b>37.5</b>

**Table 6: Effect of the proposed pre-training dataset. Classification represents zero-shot image classification. We report R@M, AUC, CIDEr, FID, and Top-1 accuracy for five V+L downstream tasks respectively.**

Method	Pre-training Dataset	Image-Text Retrieval		Image-Text Matching		Image Caption	Text-to-Image Generation	Classification
		Flick30k-CN	COCO-CN	ICM	IQM	AIC-ICC	ECommerce-T2I	ImageNet
R2D2	Wukong (100M) [12]	95.2	90.1	86.5	81.5	245.8	16.4	55.6
R2D2	Zero (23M)	95.4	90.7	88.1	83.6	246.5	15.7	55.7
R2D2	Zero (250M)	<b>95.8</b>	<b>91.6</b>	<b>90.6</b>	<b>86.7</b>	<b>247.9</b>	<b>14.4</b>	<b>56.9</b>

**Table 7: Effect of different components of R2D2. Note that we conduct ablation studies and report the average results on all downstream datasets. Generation and classification represent text-to-image generation and zero-shot image classification, respectively. R@\* denotes the result for the image-text retrieval task. We report AUC, CIDEr, FID, and Top-1 accuracy for image-text matching, image caption, text-to-image generation, and zero-shot image classification tasks respectively.**

Method	Image-to-Text Retrieval			Text-to-Image Retrieval			R@M	Image-Text Matching	Image Caption	Generation	Classification
	R@1	R@5	R@10	R@1	R@5	R@10		AUC	CIDEr	FID	Top-1 Acc.
PRD2	53.92	75.67	82.01	43.97	71.19	80.28	67.14	73.89	239.91	19.21	32.27
R2D2	<b>64.51</b>	<b>81.02</b>	<b>85.92</b>	<b>56.63</b>	<b>78.22</b>	<b>84.49</b>	<b>74.45</b>	<b>80.82</b>	<b>243.29</b>	<b>17.58</b>	<b>39.96</b>
R2D2 w/o ET	64.14	78.48	84.96	55.32	77.38	83.01	73.81	80.31	243.01	17.82	39.72
R2D2 w/o MLM	63.72	80.19	85.14	55.73	77.29	83.70	73.57	80.01	242.90	17.91	39.54
R2D2 w/o TwD	63.08	79.51	84.69	54.69	76.74	83.53	73.03	79.98	242.64	18.16	38.76
R2D2 w/o TgD	63.87	80.43	85.39	55.97	77.02	83.23	73.52	80.39	243.01	17.90	39.29
R2D2 w/o FgD	63.39	79.86	85.01	54.92	76.83	83.45	73.11	80.28	242.83	18.01	38.92

CIDEr following BriVL. Table 4 also presents text-to-image generation results on ECommerce-T2I dataset<sup>1</sup> [24]. The metric of Fréchet Inception Distance (FID) is reported. We evaluate our pre-trained models on ImageNet[6] for the zero-shot image classification task. Class labels are translated from English. Top-1 and Top-5 accuracy are reported. Our model achieves state-of-the-art performance on these V+L downstream tasks.

ALBEF [19] and FLAVA [34] also combine dual-stream unimodal encoders and single-stream multimodal encoders. They are pre-trained with English Corpus, lacking the ability to perform Chinese downstream tasks. To make a comparison with these methods, We pre-train ALBEF, FLAVA, and R2D2 on the first 1% of Zero. We replace the text encoder and tokenizer of these baselines with the same as ours. Considering that ALBEF and FLAVA use ViT-Base as the image encoder, we show the comparative performance of

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

R2D2<sub>VIT-B</sub> in Table 5. In summary, the results on various tasks demonstrate the superiority of our framework.

### 5.3 Ablation Study

**Effect of the Proposed Pre-training Dataset.** To demonstrate the effectiveness of our proposed pre-training dataset, we provide comparison results of our R2D2 framework pre-trained on the 100M Wukong dataset [12] and the proposed Zero in Table 6. Wukong is the previous largest publicly available Chinese image-text pre-training dataset. For simplicity, we define R2D2<sub>VIT-L</sub> as R2D2 in the ablation study. R2D2 pre-trained on the 23M pre-training dataset (a subset of Zero) achieves better results than the ones on the much larger 100M Wukong dataset. This improvement verifies the high quality of our Zero dataset, which is filtered by user click-through rate and provides diverse text descriptions along with each image, compared to previous datasets. Moreover, we achieve the best results on the whole pre-training dataset, *i.e.*, Zero with 250M high-quality image-text pairs.

**Effect of Fine-Grained Ranking (FGR).** We conduct subsequent ablation studies on the first 1% of Zero. 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 [30]. R2D2 outperforms PRD2 on the downstream tasks, indicating that the two cross encoders can effectively interact with image and text information through cross-attention.

**Effect of Enhanced Training (ET).** From the third row of Table 7, R2D2 (with ET) performs slightly better than R2D2 w/o ET. Furthermore, 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 obtains better performance on all downstream tasks. MLM allows R2D2 to learn robust representations by masking data. 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 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.

### 5.4 Further Experiments

**Zero-shot Tasks.** In this section, we conduct zero-shot transfer experiments. From Table 8, our R2D2<sub>VIT-L</sub> achieves the best performance on Flickr30k-CN, COCO-CN, MUGE, AIC-ICC, ICR, and IQR. For example, R2D2<sub>VIT-L</sub> achieves 80.5% R@M on COCO-CN, an absolute 5.3% gain over the previous best performance. These results demonstrate sound generalization ability of R2D2. The results of R2D2<sub>VIT-L</sub> on Flickr30k-CNA are the same as that of Flickr30k-CN,

**Table 8: Zero-shot results on image-text retrieval task.**

	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	CLIP <sub>VIT-L</sub> [30]	75.0	94.5	97.7	51.8	78.6	85.9	80.6
	FILIP <sub>VIT-L</sub> [42]	<b>78.9</b>	96.2	98.1	55.7	81.2	87.9	83.0
	Wukong <sub>VIT-L</sub> [12]	76.1	94.8	97.5	51.7	78.9	86.3	80.9
	R2D2 <sub>VIT-L</sub>	77.6	<b>96.7</b>	<b>98.9</b>	<b>60.9</b>	<b>86.8</b>	<b>92.7</b>	<b>85.6</b>
COCO-CN	CLIP <sub>VIT-L</sub> [30]	51.0	80.0	89.7	48.7	76.8	86.4	72.1
	FILIP <sub>VIT-L</sub> [42]	56.9	82.4	90.9	52.7	79.9	88.6	75.2
	Wukong <sub>VIT-L</sub> [12]	55.2	81.0	90.6	53.4	80.2	90.1	75.1
	R2D2 <sub>VIT-L</sub>	<b>63.3</b>	<b>89.3</b>	<b>95.7</b>	<b>56.4</b>	<b>85.0</b>	<b>93.1</b>	<b>80.5</b>
MUGE	CLIP <sub>VIT-L</sub> [30]	-	-	-	43.3	69.2	78.4	63.6
	FILIP <sub>VIT-L</sub> [42]	-	-	-	37.6	63.4	73.6	58.2
	Wukong <sub>VIT-L</sub> [12]	-	-	-	42.7	69.0	78.0	63.2
	R2D2 <sub>VIT-L</sub>	-	-	-	<b>49.5</b>	<b>75.7</b>	<b>83.2</b>	<b>69.5</b>
AIC-ICC	CLIP <sub>VIT-L</sub> [30]	16.8	32.0	39.8	9.7	21.1	27.5	24.5
	FILIP <sub>VIT-L</sub> [42]	20.6	37.0	45.4	11.3	24.3	31.4	28.3
	Wukong <sub>VIT-L</sub> [12]	18.2	34.5	42.4	8.8	20.3	27.3	25.3
	R2D2 <sub>VIT-L</sub>	<b>30.7</b>	<b>47.2</b>	<b>52.9</b>	<b>14.9</b>	<b>28.1</b>	<b>33.4</b>	<b>34.5</b>
ICR	CLIP <sub>VIT-L</sub> [30]	30.3	52.9	61.6	29.0	51.9	60.9	47.8
	FILIP <sub>VIT-L</sub> [42]	27.3	49.6	58.3	25.4	48.5	57.7	44.5
	Wukong <sub>VIT-L</sub> [12]	35.1	58.2	66.3	33.7	58.0	66.5	53.0
	R2D2 <sub>VIT-L</sub>	<b>58.0</b>	<b>80.5</b>	<b>85.2</b>	<b>55.9</b>	<b>78.2</b>	<b>82.4</b>	<b>73.4</b>
IQR	CLIP <sub>VIT-L</sub> [30]	24.3	47.1	56.2	22.2	45.2	54.8	41.6
	FILIP <sub>VIT-L</sub> [42]	21.9	43.2	52.8	19.9	42.0	52.0	38.6
	Wukong <sub>VIT-L</sub> [12]	26.1	48.9	58.1	24.9	48.1	57.7	44.0
	R2D2 <sub>VIT-L</sub>	<b>38.4</b>	<b>64.8</b>	<b>72.8</b>	<b>37.4</b>	<b>62.6</b>	<b>69.0</b>	<b>57.5</b>

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. In addition, the AUC scores of R2D2<sub>VIT-L</sub> on ICM and IQM are 89.8% and 84.5%, respectively.

**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 [19]. Figure D in Appendix shows that R2D2 learns well to align text with the correct content inside the image.

## 6 CONCLUSION

In this paper, we introduce a large-scale Chinese cross-modal benchmark called CCMB and a vision-language framework named R2D2. CCMB includes a high-quality pre-training dataset Zero, which is the largest Chinese cross-modal dataset, and five human-annotated downstream datasets, two of which are the largest Chinese vision-language downstream datasets and the first proposed datasets for the Chinese image-text matching task. R2D2 adopts a framework of pre-ranking + ranking for cross-modal learning, boosted with feature-guided distillation, target-guided distillation, and enhanced training. After pre-training, R2D2 achieves state-of-the-art results on fine-tuning and zero-shot settings on twelve downstream datasets of five vision-language tasks. 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.

**Acknowledgement.** This work was supported by the National Key Research and Development Program of China (No.2018AAA0100400).

## REFERENCES

- [1] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 2021. 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*. 3558–3568.
- [2] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325* (2015).
- [3] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European Conference on Computer Vision*. 104–120.
- [4] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. 2021. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*. 1931–1942.
- [5] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2020. Revisiting Pre-Trained Models for Chinese Natural Language Processing. In *Conference on Empirical Methods in Natural Language Processing*. 657–668.
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 248–255.
- [7] Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. 2021. RedCaps: Web-curated image-text data created by the people, for the people. *Advances in Neural Information Processing Systems Track on Datasets and Benchmarks* (2021).
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Human Language Technology Conference of the NAACL*. 4171–4186.
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*.
- [10] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. 2018. Vse++: Improving visual-semantic embeddings with hard negatives. In *Proceedings of the British Machine Vision Conference*.
- [11] Nanyi Fei, Zhiwu Lu, Yizhao Gao, Guoxing Yang, Yuqi Huo, Jingyuan Wen, Haoyu Lu, Ruihua Song, Xin Gao, Tao Xiang, et al. 2022. Towards artificial general intelligence via a multimodal foundation model. *Nature Communications* 13, 1 (2022), 1–13.
- [12] Jiayi Gu, Xiaojun Meng, Guansong Lu, Lu Hou, Minzhe Niu, Hang Xu, Xiaodan Liang, Wei Zhang, Xin Jiang, and Chunjing Xu. 2022. Wukong: 100 Million Large-scale Chinese Cross-modal Pre-training Dataset and A Foundation Framework. *arXiv preprint arXiv:2202.06767* (2022).
- [13] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9729–9738.
- [14] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*. 4904–4916.
- [15] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision* 123, 1 (2017), 32–73.
- [16] Weiyu Lan, Xirong Li, and Jianfeng Dong. 2017. Fluency-guided cross-lingual image captioning. In *Proceedings of the 25th ACM international conference on Multimedia*. 1549–1557.
- [17] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 11336–11344.
- [18] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*.
- [19] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems* (2021).
- [20] Kumpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. 2019. Visual semantic reasoning for image-text matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4654–4662.
- [21] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557* (2019).
- [22] Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. 2021. 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)*. 2592–2607.
- [23] Xirong Li, Chaoxi Xu, Xiaoxu Wang, Weiyu Lan, Zhengxiong Jia, Gang Yang, and Jieping Xu. 2019. COCO-CN for cross-lingual image tagging, captioning, and retrieval. *IEEE Transactions on Multimedia* 21, 9 (2019), 2347–2360.
- [24] Junyang Lin, Rui Men, An Yang, Chang Zhou, Ming Ding, Yichang Zhang, Peng Wang, Ang Wang, Le Jiang, Xianyan Jia, et al. 2021. M6: A chinese multimodal pretrainer. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3251–3261.
- [25] Ilya Loshchilov and Frank Hutter. 2016. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* (2016).
- [26] Jiaseen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in Neural Information Processing Systems* (2019).
- [27] Jiaseen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 2020. 12-in-1: Multi-task vision and language representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10437–10446.
- [28] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. 2011. Im2text: Describing images using 1 million captioned photographs. *Advances in Neural Information Processing Systems* (2011).
- [29] Di Qi, Lin Su, Jia Song, Edward Cui, Taroon Bharti, and Arun Sacheti. 2020. Imagebert: Cross-modal pre-training with large-scale weak-supervised image-text data. *arXiv preprint arXiv:2001.07966* (2020).
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*. 8748–8763.
- [31] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* (2022).
- [32] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *Advances in Neural Information Processing Systems Workshop* (2021).
- [33] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. 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)*. 2556–2565.
- [34] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. FLAVA: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15638–15650.
- [35] Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. 2021. Wit: Wikipedia-based image text dataset for multimodal multi-lingual machine learning. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2443–2449.
- [36] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. 2016. YFCC100M: The new data in multimedia research. *Commun. ACM* 59, 2 (2016), 64–73.
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems* (2017).
- [38] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*. 23318–23340.
- [39] Fangyu Wei, Yue Gao, Zhirong Wu, Han Hu, and Stephen Lin. 2021. Aligning pretraining for detection via object-level contrastive learning. *Advances in Neural Information Processing Systems* (2021).
- [40] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipeng Zhou, Guosen Lin, Yanwei Fu, et al. 2019. Ai challenger: A large-scale dataset for going deeper in image understanding. In *IEEE International Conference on Multimedia and Expo*.
- [41] An Yang, Junshu Pan, Junyang Lin, Rui Men, Yichang Zhang, Jingren Zhou, and Chang Zhou. 2022. Chinese CLIP: Contrastive Vision-Language Pretraining in Chinese. *arXiv preprint arXiv:2211.01335* (2022).
- [42] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguang Li, Xin Jiang, and Chunjing Xu. 2022. FLIP: Fine-grained Interactive Language-Image Pre-Training. In *International Conference on Learning Representations*.
- [43] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics* 2 (2014), 67–78.

## A DETAILS OF ZERO

We illustrate several representative examples of Zero in Figure A. There are 3 types of text fields associated with each image: “Title”, “Content” and “ImageQuery”.

## B EXAMPLES OF THE PROPOSED DOWNSTREAM DATASETS

Figure B illustrates examples of ICM, IQM, ICR, and IQR. Figure C highlights some cases of the difference between Flickr30k-CN and our proposed Flickr30k-CNA.

## C MORE IMPLEMENTATION DETAILS

**Training process of target-guided distillation.** We use target-guided distillation following the settings of the momentum distillation in ALBEF [19]. Our goal is to generate soft targets to replace the ground-truth labels in Equation (2). During training, we perform the following processes. i. For image and text, we use a momentum-updated encoder as the teacher model respectively, which contains the exponential-moving-average weights. ii. We use the teacher models to obtain the corresponding image-text features and compute their similarity scores. iii. We combine the above similarity scores with ground-truth labels via a coefficient parameter to generate the final soft targets. iv. We replace the ground-truth labels in Equation (2) with the generated soft targets.

**Selection Strategy of Negative Pairs in Image-text Matching.** We obtain hard negative samples by sampling in a mini-batch. Given an image in the mini-batch, we select the corresponding negative text by ranking the contrastive scores of the current batch. We choose the higher score except for the original positive text of the image. In this way, we construct one image-text negative pair for image-text matching loss. The negative images of each text are similar to the description above.

**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. During inference, we use the top-K strategy to rank the scores in the two cross encoders. We extract the individual features of images and texts via dual-stream encoders. For each image feature, we select the top-K candidate text features and construct K image-text pairs. We feed the K image-text pairs into two cross encoders to calculate similarity scores. We obtain two K-dimensional score matrices and average them to obtain a final K-dimensional score matrix for ranking. 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. 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 [12], 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[37], 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 [40] 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 [31], 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  $1e^{-4}$ .

**Fine-tuning Strategy of Zero-shot Image Classification.** Given an image, the zero-shot image classification task aims to predict the corresponding class label. Following [12], we use R2D2 to conduct zero-shot image classification task on ImageNet[6]. All the class labels in ImageNet are translated into Chinese.

## D ENTITY-CONDITIONED IMAGE VISUALIZATION

In this experiment, we visualize the attention map of images on COCO-CN. From Figure D, 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”. 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.

---



**Title:** 五大地缝奇观欣赏 (View of the five fissure wonders)

**Content:** 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.)

**ImageQuery:** 天井峡地缝 (TianJingXia fissure)

---



**Title:** 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.)

**Content:** 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.)

**ImageQuery:** 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)

---



**Title:** 美呆了!25万盆鲜花齐聚小榄菊花展

(Stunningly beautiful! 250,000 pots of flowers gathered at the Xiaolan chrysanthemum exhibition.)

**Content:** 大立菊、盆景菊、悬崖菊

(Dali chrysanthemum, bonsai chrysanthemum, cliff chrysanthemum)

**ImageQuery:** 大立菊 (Dali chrysanthemum)

---



**Title:** 零基础学绘画-彩铅《紫红色百合花》 (Zero Basic Learning Painting - Color Lead "Fuchsia Lily")

**Content:** 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.)

**ImageQuery:** 彩铅百合,彩铅百合绘画大全 (Color lead lily, color lead lily painting Daquan)

---



**Title:** 茶百戏, 一种能使茶汤纹脉形成物象的民间艺术

(Tea Baixi, a folk art that can make the veins of tea soup form objects.)

**Content:** 乌龙茶汤显现的茶百戏图

(Tea Baixi shown in Oolong tea soup)

**ImageQuery:** 茶百戏 (Tea Baixi )

---

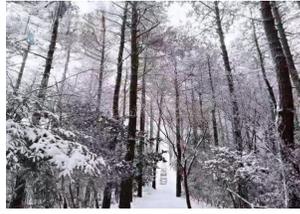
Figure A: Examples of Zero.



这么晴好的天，当然开得快！大家一定要抓住机会，去欣赏洛阳市这一一年一度的杏花满山。  
On such a sunny day, of course it drives fast! Everyone must seize the opportunity to appreciate the annual apricot blossoms in Luoyang City.



恩施民族服饰  
Enshi National Costume



这场雨雪天气将持续到今天早上，预计平原地区的积雪将达到1-4cm。  
The rain and snow will continue until this morning, with 1-4cm of snow expected in the plains.



紫乐用什么花盆  
What flower pot does Zi Le use

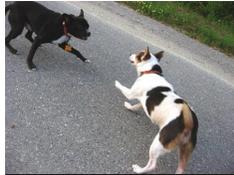
Figure B: Image-text examples of ICM, IQM, ICR and IQR from left to right.



**Flickr30k:** A little girl covered in paint sits in front of a painted rainbow with her hands in a bowl.  
**Flickr30k-CN:** 一个小女孩在油漆前坐在一个彩虹的前面双手在碗里。  
**Flickr30k-CNA:** 一个涂满染料的小女孩坐在画好的彩虹前，把她的手放在一个装颜料的碗里。



**Flickr30k:** A man with reflective safety clothes and ear protection drives a John Deere tractor on a road.  
**Flickr30k-CN:** 一个男人用反光安全服装和耳朵保护驱动的道路上约翰迪尔拖拉机。  
**Flickr30k-CNA:** 一个穿着反光安全服，带着耳护的男子在路上开着一辆约翰迪尔拖拉机。



**Flickr30k:** A black dog and a white dog with brown spots are staring at each other in the street.  
**Flickr30k-CN:** 一只黑色的狗和一只棕色的白色狗在街上盯着对方。  
**Flickr30k-CNA:** 一只黑狗和一只带有棕色斑点的白狗站在街上，互相盯着对方。

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

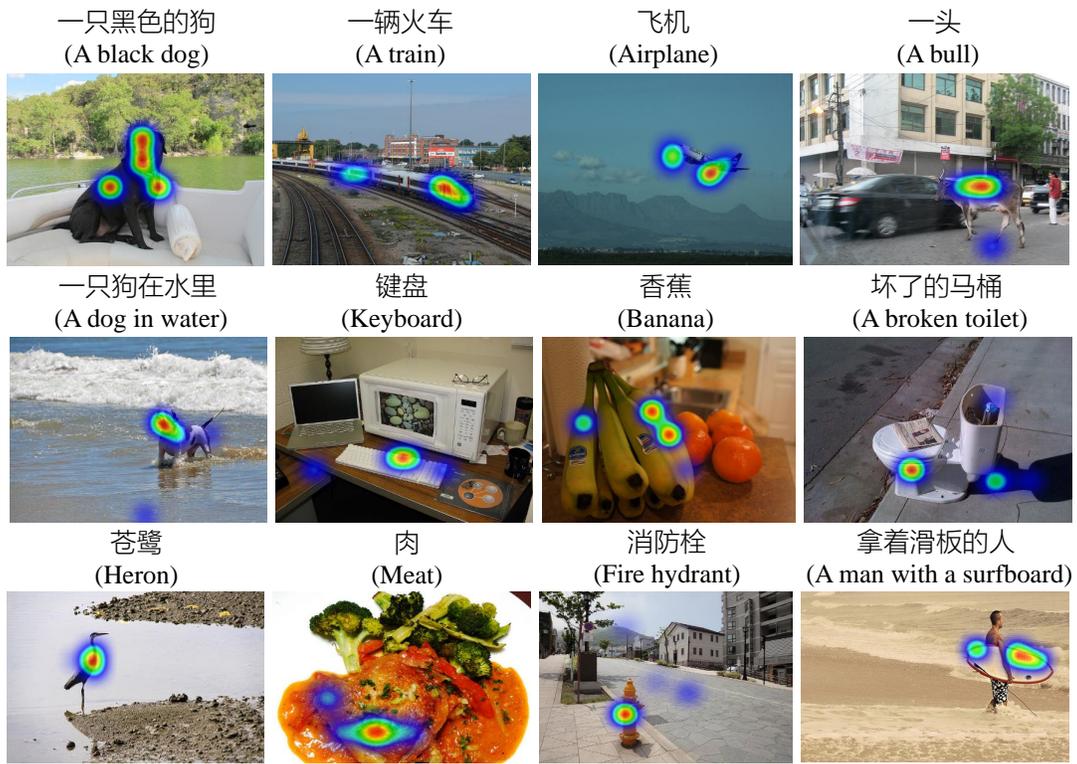


Figure D: More Examples of entity-conditioned image visualization.

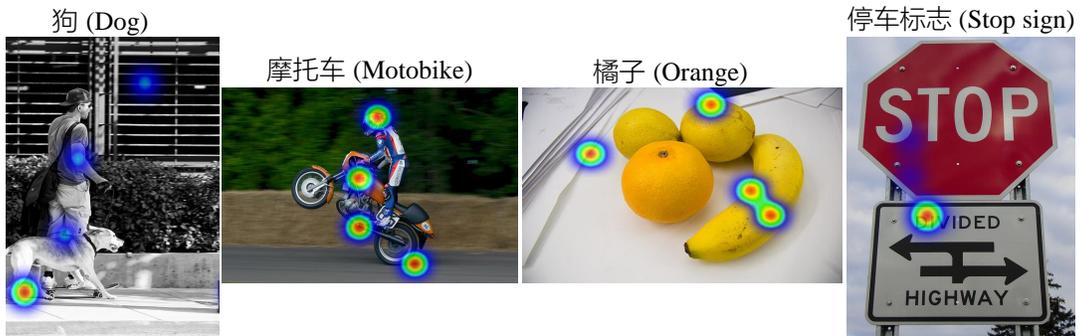


Figure E: Bad cases of entity-conditioned image visualization.