

MaskCLIP, which incorporates masked self-distillation into contrastive language-image pretraining to advance the transferable visual encoder. The core idea of masked self-distillation is to distill representation from a full image to the representation predicted from a masked image. Such incorporation enjoys two vital benefits. First, masked self-distillation targets local patch representation learning, which is complementary to vision-language contrastive focusing on text-related representation. Second, masked self-distillation is also consistent with vision-language contrastive from the perspective of training objective as both utilize the visual encoder for feature aligning, and thus is able to learn local semantics getting indirect supervision from the language. Symmetrically, they also introduce the local semantic supervision into the text branch, which further improves the pretraining performance.

1) Firstly, the learned feature representation shall characterize local patches, serving as a complementary for global representation in VL contrastive. Inspired by the recent success of MIM in learning patch representations, we also randomly mask the input image with a large portion to force the visual encoder to focus on the remaining visible patches. (2) Secondly, the learned representation for local patches shall possess semantic meanings, being consistent with the global representation receiving semantic text supervision. We bring mean teacher self-distillation to supervise the learned patch representations with the visual feature representations, enabling implicit supervision from natural language. The resulting objective is denoted as masked self-distillation where the student model and the teacher model come from the same neural networks and the knowledge is distilled from the full image (fed to the teacher model) to the masked image (fed to student model)

Self knowledge distillation not pretrained knowledge distillation and this is for pretraining approach

Masked image modeling is able to learn representations for local patches. We argue that the image encoder only pays attention to the text-described objects under VL contrastive due to sparse text description and to the centric objects under image contrastive due to central-crop augmentation. In contrast, masked image modeling forces the image encoder to focus on local patches using token-wise objectives by mandatorily masking a large portion of patches

MaskCLIP: Masked Self-Distillation Advances Contrastive Language-Image Pretraining

Xiaoyi Dong^{1,†}, Jianmin Bao^{2,*}, Yinglin Zheng³, Ting Zhang², Dongdong Chen^{4,†}, Hao Yang², Ming Zeng³, Weiming Zhang¹, Lu Yuan⁴, Dong Chen², Fang Wen², Nenghai Yu¹

¹University of Science and Technology of China ²Microsoft Research Asia

³Xiamen University ⁴Microsoft Cloud + AI

{dlight@mail., zhangwm@, ynh@}.ustc.edu.cn cddlyf@gmail.com

{jianbao, ting.zhang, luyuan, doch, fangwen}@microsoft.com

{zhengyinglin@stu., zengming@}xmu.edu.cn yanghao.alexis@foxmail.com

Abstract

This paper presents a simple yet effective framework **MaskCLIP**, which **incorporates a newly proposed masked self-distillation into contrastive language-image pretraining**. The **core idea of masked self-distillation is to distill representation from a full image to the representation predicted from a masked image**. Such incorporation enjoys **two vital benefits**. First, **masked self-distillation targets local patch representation learning, which is complementary to vision-language contrastive focusing on text-related representation**. Second, **masked self-distillation is also consistent with vision-language contrastive from the perspective of training objective as both utilize the visual encoder for feature aligning, and thus is able to learn local semantics getting indirect supervision from the language**. We provide specially designed experiments with a comprehensive analysis to validate the two benefits. **Symmetrically, we also introduce the local semantic supervision into the text branch, which further improves the pretraining performance**. With extensive experiments, we show that MaskCLIP, when applied to various challenging downstream tasks, **achieves superior results in linear probing, finetuning, and zero-shot performance with the guidance of the language encoder**. Code will be release at <https://github.com/LightDXY/MaskCLIP>.

1. Introduction

Vision-language (VL) contrastive learning [34, 56] has shown remarkable success in pretraining for various tasks. With large-scale image-text pairs available on the Internet, the model composed of a simple dual encoder design learns

strong semantic prior by aligning between image and text. The resulting visual encoder not only exhibits excellent linear probing and finetuning performance, but also enables impressive zero-shot performance with the guidance of the language encoder, showing the generality of natural language and its ability to supervise a wide range of visual concepts.

Nonetheless, the associated language description, though providing richer information than mere class labels, still can hardly describe all the information in the corresponding image, as images are continuous signals with fine-grained details and complex semantics. As a result, the VL contrastive by aligning global representations may only focus on the text-described objects and ignore the rest which might be useful for downstream tasks.

In this paper, we are interested in how to fully leverage the image itself to facilitate the VL contrastive to further improve the transfer capability. (1) Firstly, the learned feature representation shall characterize local patches, serving as a complementary for global representation in VL contrastive. Inspired by the recent success of masked image modeling [4, 22, 29, 56, 65, 66] in learning patch representations, we also randomly mask the input image with a large portion to force the visual encoder to focus on the remaining visible patches. (2) Secondly, the learned representation for local patches shall possess semantic meanings, being consistent with the global representation receiving semantic text supervision. We bring mean teacher self-distillation [28, 62] to supervise the learned patch representations with the visual feature representations, enabling implicit supervision from natural language. The resulting objective is denoted as **masked self-distillation** where the student model and the teacher model come from the same neural networks and the knowledge is distilled from the full image (fed to the teacher model) to the masked image (fed to student model). To this end, we introduce **MaskCLIP** by **incorporating masked self-**

*Equal contribution, † Corresponding Author

†Work done during an internship at Microsoft Research Asia

distillation into VL contrastive to advance the transferable visual encoder.

There are several recent attempts [54, 75] also exploring the capability of the visual encoder under natural language supervision. The common approach is to introduce contrastive learning or masked image modeling on the vision side together with contrastive language-image pretraining. However, the performance indeed improves based on CLIP but does not as well as our masked self-distillation. We argue that (1) the contrastive learning objective based on central crop augmentation actually learns global representations for salient objects while lack of attention on the surrounding backgrounds [12]; and (2) masked image modeling usually needs to remap the learned representation to pixels [29] or discrete tokens [4]. Such low-level prediction target is inefficient for semantic feature learning and thus also conflicts with high-level language supervision in VL contrastive. A brief illustration is presented in Figure 1. In the experiments, we conduct comprehensive ablations to analyze the difference and provide numerical and visual evidence for better understanding.

Symmetrically, we argue that local semantic supervision on the text branch is also helpful for the text encoder and eventually beneficial for zero-shot performance. So we introduce the same mask-data-modeling format supervision into the text branch as well. Different from images where the pixel is low-level signal, the words crafted by human beings are already highly semantic, so we use the tokenized word piece as the prediction target directly, following the well-studied mask language modeling method BERT. Meanwhile, to reduce the output conflicts between contrastive learning and mask language modeling, we introduce a small decoder for the mask language modeling branch.

We train our MaskCLIP on a subset of a publicly available image-text pairs dataset, YFCC [63], and thoroughly evaluate the transfer ability of visual representations on several vision benchmarks: ImageNet-1K [20] for classification, ADE20K [76] for semantic segmentation, MS-COCO [44] for detection and segmentation, as well as a batch of other classification benchmarks. When it comes to ImageNet-1K [20] classification, MaskCLIP achieves +6.9%, +7.2%, +1.3% higher than CLIP for zero-shot transfer, linear probing, and finetuning respectively. For vision downstream tasks, we reach +2.7 mIoU on ADE20K [76] and +1.8 AP^b, +1.4 AP^m on MS-COCO [44]. For vision-language tasks, MaskCLIP achieves +6.1% average zero-shot accuracy on 20 datasets, and +17.2%, +12.8% rank@1 improvement on the Flickr30K [74] image-text retrieval. In the recent Image Classification in the Wild challenge academic track, our MaskCLIP gets the 1st result with 48.9% TOP-1 average accuracy, surpassing the second team with 3.4%.

In summary, the major contributions of this work are:

1. We present a novel vision-language pretraining

framework MaskCLIP, by introducing masked self-distillation objective to facilitate VL contrastive for better transferable visual models.

2. We present extensive ablation studies on MaskCLIP variants and provide in-depth analysis numerically and visually to help understand how the proposed masked self-distillation assists VL contrastive.
3. We demonstrate our MaskCLIP on tens of benchmarks, showing the superiority under all three settings: zero-shot, linear probing, and finetuning.

2. Related Work

Vision-language pretraining Recent years have seen rapid progress made in vision-language pretraining [16, 21, 36, 38–42, 49–51, 55, 60, 61, 79]. Several multiple cross-modality loss functions have been proposed for the training objective, such as image-text matching [16, 40, 49, 61, 69], masked language modeling [16, 40, 49, 60, 61], masked image modeling [16, 49, 60, 61], contrastive loss [38, 41, 42]. These objects are often mixed with each other to form a compound objective. While a variety of approaches have been proposed, few works investigate the performance on visual representation learning for image classification. Recently, CLIP [56] and ALIGN [34] show that the image-text contrastive learning objective achieves promising performance for visual representation learning. There are many following works proposed to further improve the pretraining performance, DeCLIP [77], SLIP [54], COTS [48], ViCHA [59], CYCLIP [27] use additional uni/multi-modality supervision to improve the model capability, and PyramidCLIP [26], KLITE [58], IDEA [33] seek to external knowledge from pre-trained models or datasets as the additional guidance. FILIP [72] and LOUPE [37] introduce fine-grained alignment to the model. Focusing on this research direction, we analyze the desired properties of supervision which could be complementary to CLIP, and propose the masked self-distillation objective incorporated with the image-text contrastive loss to further improve pretraining performance for various visual understanding tasks.

Self-supervised learning Self-supervised visual representation learning has attracted increasing attention over the past few years. The objective of the self-supervised learning is mainly divided into two categories: contrastive and generative [45]. The contrastive methods, such as MOCO [13, 30], SimCLR [10, 11], BYOL [28], SimSiam [14], and DINO [6] measure the similar and dissimilar samples by contrastive loss. Their success heavily depends on the strong data augmentation. The generative methods, such as BEiT [4], MAE [29], PeCo [22], BEVT [65], BootMAE [23] and MaskFeat [66] leverage masked image modeling to reconstruct the remaining masked part of its original input from the given visible parts. The generative methods show more

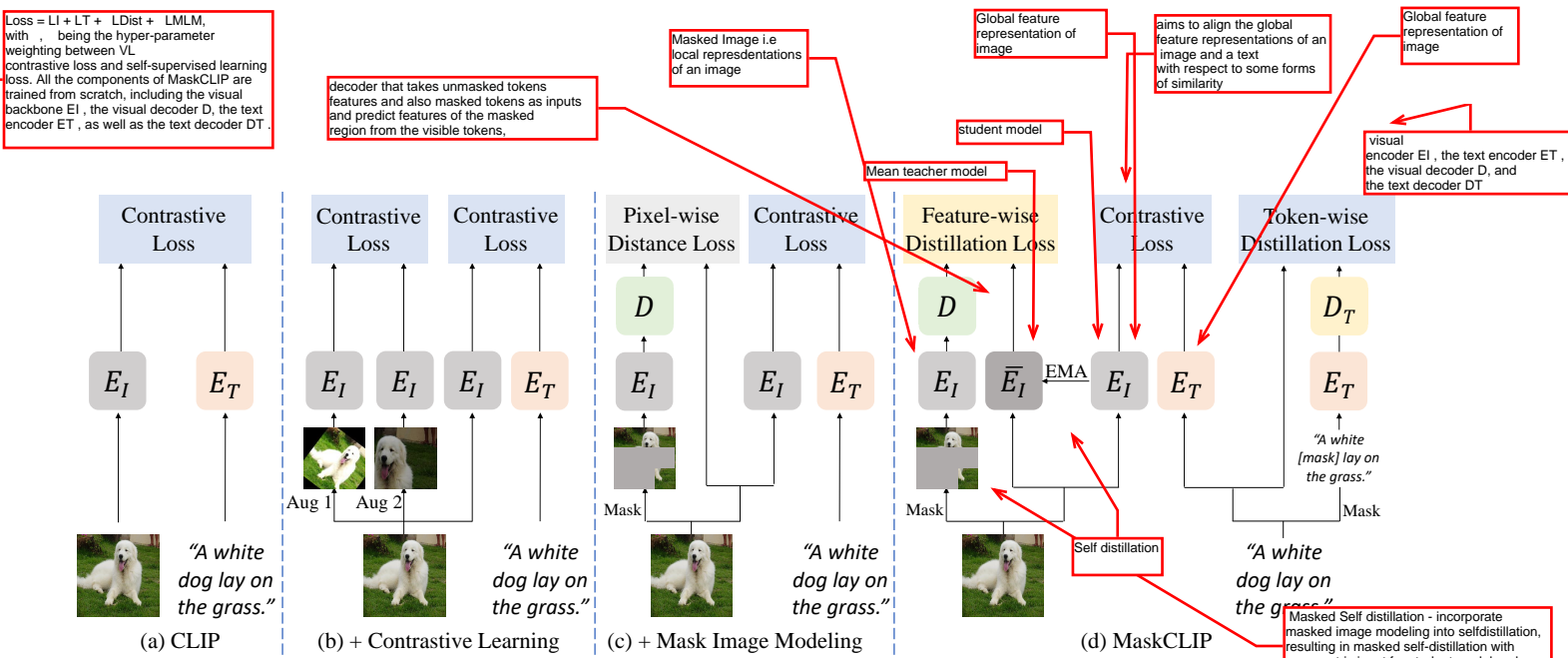


Figure 1. Pipeline comparison between combination CLIP with different vision self-supervised learning methods. (a) Vanilla CLIP. (b) CLIP + contrastive learning. (c) CLIP + pixel prediction mask image modeling. (d) CLIP + mask self-distillation, i.e. MaskCLIP. The E_T , E_I is the text encoder and image encoder respectively, and all the E_I , E_T within each pipeline share the weight. \bar{E}_I is the mean-teacher model, whose weight is updated by the exponential moving average of E_I and does not require gradient.

promising transfer performance than the contrastive methods, as generative objective learns patch representations while contrastive objective focuses on learning centric global representations [12].

Self-knowledge distillation Self-knowledge distillation [35] aims to distill the knowledge in a model itself and uses it for training the model. Instead of distilling knowledge from a pretrained teacher model [32], self-knowledge distillation regards a temporal ensemble of the student model as the teacher. It means that a student model becomes a teacher model itself, which gradually utilizes its own knowledge for softening the hard targets to be more informative during training. Self-knowledge distillation has been explored in semi-supervised learning [62], contrastive learning [17, 38], self-supervised learning [3, 7]. In this paper, we use visual features supervised by natural language for guidance in masked self-distillation which naturally fit VL contrastive to learn more transferable visual representations.

3. MaskCLIP

We introduce MaskCLIP, a novel framework that learns visual representations. The core part of MaskCLIP is its backbone image encoder, denoted by E_I as shown in Figure 1. It obtains the transferable capability during pretraining that could benefit downstream vision tasks. Following recent self-supervised approaches [4, 15, 29, 54], we implement the backbone E_I as a Vision Transformer (ViT) [25]. The prediction results from E_I given an input image I then should be a collection of visual feature tokens, represented as

$$E_I(I) = \{f_{cls}, f_1, f_2, \dots, f_N\}. \quad (1)$$

Here cls is short for class token. $1, \dots, N$ are the indexes of the non-class tokens.

The rest of this section starts with the utilization of language supervision. More shall be emphasized on the masked self-distillation, which we deem crucial for vision-language learning.

3.1. Vision-language Contrastive

Following [34, 56], we introduce a Transformer-based text encoder E_T to leverage language knowledge. It aims to align the global feature representations of an image and a text with respect to some forms of similarity. Precisely, consider a given image-text pair $\{I, T\}$, besides extracting the visual feature representation $E_I(I)$ using the vision backbone as shown by Equation 1, we additionally use the text encoder E_T to extract linguistic features from the text T .

The mean feature of the two branches are regarded as the global representations and are fed into a projection head (implemented as a fully-connected layer) respectively to obtain the metric embeddings e^T and e^I . Image-text contrastive loss is employed to align them during pretraining. The loss can be formulated as $\mathcal{L}_T + \mathcal{L}_I$, with

$$\begin{aligned} \mathcal{L}_I &= -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(e_i^I e_i^T / \sigma)}{\sum_{j=1}^B \exp(e_i^I e_j^T / \sigma)} \\ \mathcal{L}_T &= -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(e_i^T e_i^I / \sigma)}{\sum_{j=1}^B \exp(e_i^T e_j^I / \sigma)}, \end{aligned} \quad (2)$$

where B stands for the number of image-text pairs within a training mini-batch, i, j are indexes within the batch; σ stands for the temperature for the loss functions, which is learned together with all other parameters during training.

To reduce the output conflict between the global image text contrastive learning and the local mask language modeling, we further introduce a small text decoder, which shares the same architecture as the encoder but with only a few layers. So that the global prediction and local prediction are conducted at different layers.

3.2. Masked Self-distillation for Visual Encoder

Knowledge distillation is a learning paradigm where a student model is trained to match the output of a given teacher model, so that the student model can be improved by the teacher. Instead of bringing in an external teacher, self-distillation methods such as [7, 28, 62] proposes using a *mean teacher* model that is derived from the student itself. In specific, the teacher shares the same structure with the student, while the parameters of the teacher are exponential moving averages (EMA) of the parameters from the student. In the following, we would use the term “EMA model” to represent such mean teacher model constructed from the student.

MaskCLIP leverages the mean teacher self-distillation to enhance its vision representations. Let \bar{E}_I be the EMA model of the backbone encoder E_I . θ_t and $\bar{\theta}_t$ are the parameters of E_I and \bar{E}_I at training step t . $\bar{\theta}_t$ is updated with

$$\bar{\theta}_t = \alpha \bar{\theta}_{t-1} + (1 - \alpha) \theta_t, \quad (3)$$

where α is a hyper-parameter for smoothing updates. We propose to incorporate masked image modeling into self-distillation, resulting in *masked self-distillation* with asymmetric input for student model and teacher model.

In specific, considering a given input image I , we first feed it to the EMA model \bar{E}_I (teacher model) to obtain the distillation targets. These target features can be represented as

$$\bar{E}_I(I) = \{\bar{f}_{cls}, \bar{f}_1, \bar{f}_2, \dots, \bar{f}_N\}. \quad (4)$$

In the meantime, we randomly mask a large portion of the input image patches and then feed it into the original backbone E_I (student model). Following [29], we only feed the visible (unmasked) patches, denoted by I' , into the original backbone E_I to speed up computation and save memory. Let \mathcal{M} be the indexes of all the masked tokens. These encoded features corresponding to visible tokens can then be denoted as $E_I(I') = \{f'_{cls}\} \cup \{f'_{k \notin \mathcal{M}}\}$. They are then joined with a shared and learnable feature vector, denoted as m , that represents mask tokens, to form a complete set of features $\{f'_{cls}, f'_1, f'_2, \dots, f'_N\}$, with $f'_{i \in \mathcal{M}} = m$. We attach positional embeddings onto all these tokens, and append a small Transformer D as a decoder to predict features of the masked region from the visible tokens, which could be formulated as

$$\begin{aligned} (D \circ E_I)(I') &= D(\{f'_{cls}, f'_1, f'_2, \dots, f'_N\}) \\ &= \{f''_{cls}, f''_1, f''_2, \dots, f''_N\}. \end{aligned} \quad (5)$$

Inspired by [78], we use an online quantizer $h(\cdot)$ to transform the output features into a soft codewords distribution, and minimize the cross-entropy between the target features and the predicted features. Formally,

$$\mathcal{L}_{Dist} = \frac{1}{|\mathcal{M}|} \sum_{k \in \mathcal{M}} -\bar{h}(\bar{f}_k)^T \log h(f''_k). \quad (6)$$

here the parameter of the teacher quantizer $\bar{h}(\cdot)$ is also EMA updated by the online quantizer, similar to the teacher model.

3.3. Local Semantic Learning for Text Encoder

Besides the local semantic supervision for the visual encoder, we argue it is also helpful for the text encoder. So we introduce the BERT pretraining into the text branch. For the text $T = \{t_{sos}, t_1, t_2, \dots, t_M, t_{eos}\}$, we denote the masked input as $T' = \{t'_{sos}, t'_1, t'_2, \dots, t'_M, t'_{eos}\}$, where $t'_{i \in \mathcal{M}_T} = m_t$ and $t'_{i \notin \mathcal{M}_T} = t_i$, and \mathcal{M}_T be the indexes of all the masked text tokens. The output feature of the encoder is $E_T(T')$.

To reduce the output conflict between the global image-text contrastive learning and the local mask language modeling, we further introduce a small text decoder, which shares the same architecture as the encoder but with only a few layers. So that the global prediction and local prediction are conducted at different layers. We denote the output feature as: $(D_T \circ E_T)(T') = \{t''_{sos}, t''_1, t''_2, \dots, t''_M, t''_{eos}\}$ and the loss could be formulated as:

$$\mathcal{L}_{MLM} = \frac{1}{|\mathcal{M}_T|} \sum_{k \in \mathcal{M}_T} -t_k^T \log t''_k. \quad (7)$$

3.4. Overall Loss Functions

Finally, we pretrain MaskCLIP with all these losses combined:

$$\mathcal{L}_I + \mathcal{L}_T + \lambda \mathcal{L}_{Dist} + \beta \mathcal{L}_{MLM}, \quad (8)$$

with λ, β being the hyper-parameter weighting between VL contrastive loss and self-supervised learning loss. All the components of MaskCLIP are trained from scratch, including the visual backbone E_I , the visual decoder D , the text encoder E_T , as well as the text decoder D_T .

4. Experiments

4.1. Setup

Model architecture. Our framework consists of the visual encoder E_I , the text encoder E_T , the visual decoder D , and the text decoder D_T . We adopt the widely used Transformer ViT-B/16 [25] for a fair comparison. It is composed of 12 layers, 768 width, and 12 head. The input image is 224×224 resolution and is further split into 14×14 patches with size 16×16 . A learnable cls token is prepended to the 196 embeddings. For the text encoder, we adopt a 12-layer, 512-width, and 8-head Transformer following CLIP [56], and the text decoder has 4 layers. The number of text tokens is fixed to 77 with necessary truncations or paddings. For the image decoder, we directly use a one-layer Vision Transformer.

Pretraining details. We train our proposed MaskCLIP from scratch for 25 epochs, the batch size is fixed to 4096 for all the experiments. The masks used in the mask self-distillation

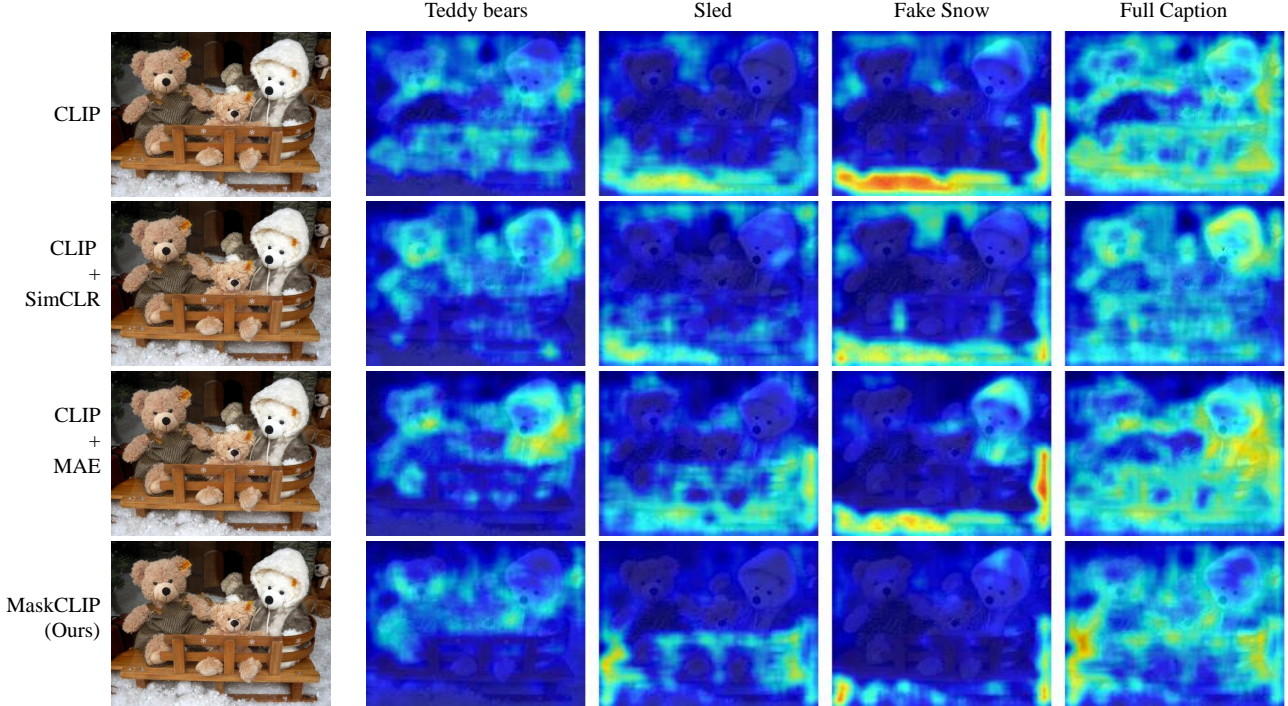


Figure 2. Visualization of the similarity between text and image features. The images and captions are from the MS-COCO val set. Here we show the image feature similarity with both full caption and different objects in it. The caption is “Three teddy bears sit in a sled in snow”. More results could be found in the supplemental materials.

	Training		IN-1K			Flicker30K	
	Memory	Time	0-shot	Linear	Finetune	I2T	T2I
CLIP	14G	1.00×	37.6	66.5	82.3	52.9	32.8
CLIP+SimCLR	30G	2.67×	42.8	72.1	82.6	58.6	41.3
CLIP+MAE	16G	1.30×	42.1	68.5	83.2	57.3	41.1
MaskCLIP	19G	1.75×	44.5	73.7	83.6	70.1	45.6

Table 1. Results of boosting CLIP with different kinds of vision self-supervised learning methods.

branch and mask language modeling branch are random mask with a mask ratio of 75% and 20%. We pretrain all the models with the commonly used YFCC15M dataset, which is flited from the YFCC100M [63] dataset by [56].

Downstream details. We evaluated MaskCLIP on several downstream datasets, including ImageNet-1K [20], ADE20K [76], MS-COCO [44], Flicr30K [74] *et al.* For ImageNet-1K, we report zero-shot, linear probing, and finetuning performance. The zero-shot is conducted following the label prompt setting in SLIP [54]. For linear probing, we fix the backbone and train a new linear classifier for 90 epochs. For finetuning, we follow the setting in BEiT [4] and finetune the model for 100 epochs with a layer-decayed learning rate. See supplemental materials for more details.

4.2. Analysis

We first present our analysis by studying different ways of boosting CLIP. The baseline is CLIP [56] trained on the

YFCC-15M. Besides the introduced masked self-distillation, we consider two other popular methods: (1) SimCLR [10], a representative contrastive method; and (2) MAE [29] the state-of-the-art masked image modeling approaches. All the compared methods are trained on the YFCC-15M for a fair comparison. We have the following observations.

Vision self-supervision helps VL contrastive. We evaluate the models on both vision task ImageNet-1K [20] classification and vision-language task image-text retrieval on Flicker30K [74] and present the comparison in Table 1. All the added vision self-supervision, regardless of contrastive or generative, improves the baseline CLIP. Among them, our proposed MaskCLIP achieves the best results in terms of all the evaluation metrics, outperforming CLIP with +6.9%, +7.2%, +1.3% on ImageNet-1K classification for zero-shot, linear probing, and finetuning respectively, and +17.2%, +12.8% on Flicker30K for image-to-text retrieval and text-to-image retrieval. We also report the training GPU memory usage and time-consuming cost in Table 1. It is worth noting that the contrastive model (CLIP+SimCLR) compares two additional views of the input image, resulting in larger GPU memory usage and longer training time.

Masked image modeling is able to learn representations for local patches. We argue that the image encoder only pays attention to the text-described objects under VL contrastive due to sparse text description and to the centric objects under image contrastive due to central-crop augmenta-

Method	Objective	ADE20K mIoU	Pascal mIoU
CLIP	Global	7.2	13.5
CLIP + SimCLR	Global + Global	6.3	11.9
CLIP + MAE	Global + Pixel-wise Local	8.3	16.4
MaskCLIP	Global + Token-wise Local	10.2	17.2

Table 2. Annotation-free zero-shot segmentation results on ADE20K and Pascal Context.

tion. In contrast, masked image modeling forces the image encoder to focus on local patches using token-wise objectives by mandatorily masking a large portion of patches. Here, we provide numerical comparisons for evidence. We conduct an “Annotation-free zero-shot segmentation” experiment to test the zero-shot segmentation. The results on such a dense prediction task would better reveal the ability of local patch representations than global classification. Following the design in DenseCLIP [77], we use the prompted label feature as the linear classification weight to realize segmentation, without any training procedure. We evaluate the performance on two widely used datasets: ADE20K [76] and Pascal Context [53]. The results are shown in Table 2. We can see that equipped with masked image modeling, our MaskCLIP as well as CLIP+MAE achieves better results than CLIP and CLIP+SimCLR, validating our hypothesis.

Masked self-distillation learns semantic representations for local patches. Our masked self-distillation predicts visual features dynamically outputted by the visual encoder and thus implicitly gets supervision from the text side via VL contrastive. While MAE predicts fixed low-level pixels, making it inefficient to learn semantic representations (as the objective may force the representation to memorize low-level details) and thus causing conflict with VL contrastive. To show this, we select images from MS-COCO [44] and calculate the feature similarity between image features and their corresponding caption features. We also select objects in the caption, prompt it to a new caption, such as “a photo of teddy bears”, and calculate the similarities. An example is shown in Figure 5 (More can be found in the supplementary material). Comparing MaskCLIP with CLIP+MAE in the fourth column, we can see that CLIP+MAE uses color as evidence and fails to distinguish the white teddy bear from the white snow. While our MaskCLIP successfully differentiates the two objects, suggesting ours learn more semantic features. On the other hand, the superior results of MaskCLIP shown in Table 1 and Table 2 also validate this. It is worth mentioning that CLIP and CLIP+SimCLR fail to have a correct response partition for different single objects like MaskCLIP, further justifying our second observation.

4.3. Comparison with Previous Methods

To show the effectiveness of MaskCLIP as a general vision-language pretrain method, we conduct experiments on both vision tasks and vision-language tasks. For vision

Method	Epoch	IN-1K			ADE20K mIoU	MS-COCO	
		0-Shot	Lin	FT		AP ^b	AP ^m
DeiT [64]	300*	–	–	81.8	47.4	44.1	39.8
SimCLR [10]	25	–	64.0	82.5	48.0	44.6	40.2
MAE [29]	25	–	56.2	82.5	46.5	43.2	39.1
CLIP [56]	25	37.6	66.5	82.3	47.8	43.6	39.5
SLIP [54]	25	42.8	72.1	82.6	48.5	44.0	40.3
MaskCLIP	25	44.5	73.7	83.6	50.5	45.4	40.9

Table 3. Comparison with previous methods, including supervised baselines, self-supervised pretraining methods, and vision-language pretraining methods. * is the epoch of the supervised baseline on ImageNet-1K.

tasks, we report results on ImageNet-1K [20] classification, MS-COCO [44] object detection, and ADE20K [76] semantic segmentation. For vision-language tasks, we report zero-shot results on recent challenging ICinW 20 datasets benchmark and image-text retrieval results on Flickr30K [74] and MS-COCO [44]. In the following, we compare with the supervised baseline DeiT [64], self-supervised methods SimCLR [10] and MAE [29], and vision-language methods CLIP [56] and SLIP [54]. For a fair comparison, we train SimCLR and MAE on YFCC-15M [63] with the same epochs.

Classification on ImageNet-1K. As shown in Table 3, MaskCLIP benefits from the advantages of both VL pretraining and image mask self-distillation that shows strong performance on all the metrics. For zero-shot tasks, MaskCLIP outperforms CLIP by +6.9% with 25 epoch training and achieves +1.7% higher than the recent work SLIP. When it comes to finetune, MaskCLIP reaches 83.6% top-1 accuracy, and outperforms CLIP by +1.3%.

Semantic segmentation on ADE20K. Then we apply our MaskCLIP to the semantic segmentation task. Here we use the UperNet [68] framework with 512×512 input and end-to-end training for 160K iterations. The evaluation metric is the mean Intersection of Union (mIoU) and we report single-scale evaluation results here. The results are given in Table 3. Our method achieves 50.5 mIoU, +2.7 mIoU than our baseline method CLIP, and +2.0 mIoU than SLIP. This verifies the effectiveness of our introduced incorporation.

Object detection and instance segmentation on MS-COCO. We further investigate our transfer performance on object detection and instance segmentation in Table 3. Here we use Mask-RCNN [31] framework with single-scale input and $1 \times$ schedule (12 epochs). Our method achieves 45.4 box AP and 40.9 mask AP, +1.8/1.4 better than CLIP, and +1.4/0.6 better than SLIP.

Zero-shot on small datasets. We also report zero-shot performance on 20 small datasets under the ICinW setting (see the introduction below) in Table 4. We find that all the methods perform poorly on some datasets such as Aircraft(1% acc for random guessing, we omit the description in the

	Average	Caltech-101	CIFAR-10	CIFAR-100	Country211	DTD	EuroSAT	FER-2013	Aircraft	Food-101	GTSRB	Memes	KittiDis	MNIST	Flowers	Pets	PatchCam	SST2	RESISC45	Cars	Voc2007
<i>Pretraining on YFCC-15M</i>																					
CLIP	34.0	58.6	68.5	36.9	10.8	21.4	30.5	16.9	5.1	51.6	6.5	51.1	25.9	5.0	52.7	28.6	51.7	52.5	22.4	4.5	79.1
SLIP	37.8	70.9	82.6	48.6	11.8	26.6	19.8	18.1	5.6	59.9	12.6	51.8	29.4	9.8	56.3	31.4	55.3	51.5	28.5	5.4	80.5
MaskCLIP	40.1	72.0	80.2	57.5	12.6	27.9	44.0	20.3	6.1	64.9	8.5	52.0	34.3	4.9	57.0	34.3	50.1	49.9	35.7	6.7	82.1
<i>Pretraining on ICinW Academic Track Stting: YFCC-15M, GCC3M+12M, ImageNet-21K(ImageNet-1K is removed)</i>																					
1st MaskCLIP	48.9	86.4	95.3	78.3	11.6	33.0	57.7	18.8	8.0	78.9	17.3	52.8	16.0	7.3	74.2	74.4	52.1	46.2	54.3	26.5	82.3
2nd KLITE*	45.5	87.4	92.7	68.8	8.2	32.2	27.9	17.4	4.3	72.4	11.4	48.4	31.1	12.8	75.6	65.9	50.6	52.9	44.4	10.2	82.3
3rd YT-CLIP	44.5	77.8	83.5	58.4	11.9	31.9	40.7	27.1	6.9	68.7	18.8	52.3	9.1	18.8	53.1	69.3	51.5	50.3	52.7	19.7	79.3
4th UniCL†	44.0	84.8	90.2	67.8	6.7	25.4	35.3	30.8	3.5	68.3	11.1	51.0	17.9	11.3	71.7	44.9	52.1	49.5	41.4	24.2	81.3
5th Gramer*	43.2	83.9	92.9	69.5	7.3	25.5	24.4	30.4	2.7	71.0	9.0	52.6	12.4	10.1	70.4	52.4	50.6	50.1	44.8	13.8	81.3

Table 4. Zero-shot evaluation on ICinW classification benchmarks. Best results in **bold**. * indicates using Swin-B as the backbone, † indicates using Focal-B as the backbone.

	Training Epoch	Flickr30K						MS-COCO					
		Image-to-text			Text-to-image			Image-to-text			Text-to-image		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
CLIP [56]	25	52.9	79.6	87.2	32.8	60.8	71.2	27.5	53.5	65.0	17.7	38.8	50.5
SLIP [54]	25	58.6	85.1	91.7	41.3	68.7	78.6	33.4	59.8	70.6	21.5	44.4	56.3
MaskCLIP	25	70.1	90.3	95.3	45.6	73.4	82.1	41.4	67.9	77.5	25.5	49.7	61.3

Table 5. Results of zero-shot image-text retrieval on Flickr30K and MS-COCO datasets. Best results in **bold**.

following), Fer(24.7%), Country211(0.5%), GTSRB(5.9%), Cars(0.8%). This might be caused by the data domain gap that the YFCC-15M contains few related images and descriptions. For the rest of the datasets, all the methods get reasonable performance and our MaskCLIP gets the best performance on most datasets.

Image Classification in the Wild (ICinW) Challenge The ICinW challenge [1] is a newly proposed visual pretraining benchmark, which contains 20 diverse downstream classification datasets, measuring the ability of pre-training models on both the prediction accuracy and their transfer efficiency in a new task. The pretraining is limited to three datasets: YFCC-15M [63], GCC3M [57]+12M [8] and ImageNet-21K [20] (ImageNet-1K data is excluded). We pretrain our MaskCLIP on it and get the **1_{st}** result in the zero-shot track [2] (we submit the results anonymously). As shown in Table 4, the **2_{nd}** team KLITE uses a strong Swin-B [46] as the backbone and additional knowledge from GPT-3 [5] and Wiktionary [52], and the **4_{th}** use the strong Focal-B [70] as the backbone, while our MaskCLIP greatly outperforms these methods with a simple ViT-B backbone and no additional knowledge.

Zero-shot on text-image retrieval. We further report the zero-shot text-image retrieval results on two benchmark datasets, Flickr30K [74] and MS-COCO [44]. We find that the text without any prefixes or suffixes works well for all the models. Table 5 shows the results. We can see that MaskCLIP exhibits a strong zero-shot performance. For example, with 25 epochs training, MaskCLIP reaches 41.4% Rank@1 image-to-text accuracy on MS-COCO, outperform-

ing CLIP with 13.9%, and 25.5% Rank@1 text-to-image accuracy, +7.8% higher than CLIP.

4.4. Ablations

We compare our default settings with other alternatives to justify the efficacy of our model designs.

Training objectives ablation. As shown in Table.6a, when we remove the mask language modeling loss \mathcal{L}_{MLM} , the performance of the image-text task drops, including the zero-shot accuracy and retrieval performance. While benefiting from the distillation loss, the finetuning performance on ImageNet-1K is not influenced. When we remove the distillation loss \mathcal{L}_{Dis} , we observe a performance drop on all tasks, especially the finetuning results.

Distillation loss format. Different from previous methods [3, 24, 29] that calculate the per-element distance as the loss function, we use an online tokenizer to map the feature to soft codewords and use the cross-entropy loss as the supervision. Here we study their difference in Table.6b. We find that although they get similar fine-tuning performance, the CE loss gets better zero-shot and linear probing performance. The reason may be that the per-element MSE loss leads the model to fit some unnecessary details of the target feature, while the CE loss with soft tokenizer helps the model to focus more on the important feature.

Distillation & MLM loss weight. Here we set the loss weight of the CLIP branch as 1 and study the loss weight of the two additional branches. As shown in Table.6e and Table.6f, setting $\lambda = 1$ or $\beta = 1$ emphasize too much on new tasks, which mislead the model to a wrong converge

Model	0-Shot	FT	I2T/T2I
MaskCLIP	44.5	83.6	70.1/45.6
w/o \mathcal{L}_{MLM}	42.8	83.6	65.0/41.6
w/o \mathcal{L}_{Dis}	42.0	82.4	65.4/40.5

(a) **Training Objectives ablation.** Both is necessary for MaskCLIP.

Loss	0-Shot	Lin	FT
MSE	43.8	73.2	83.6
CE	44.5	73.7	83.6

(b) **Distillation loss format.** The online tokenizer with cross-entropy loss works slightly better than MSE loss.

Depth	0-Shot	Lin	FT
1	44.5	73.7	83.6
2	43.7	72.9	83.4
4	43.5	72.5	83.3

(c) **Visual decoder Depth.** A shallow decoder gets better performance.

Depth	0-Shot	I2T/T2I
0	43.5	65.2/44.1
1	44.3	70.4/45.3
2	44.3	70.2/45.4
4	44.5	70.1/45.6
8	44.2	67.5/44.7

(d) **Text decoder depth.** The decoder is necessary and a shallow one works better.

Weight	0-Shot	Lin	FT
1	38.5	68.2	82.5
0.1	44.4	73.5	83.5
0.05	44.5	73.7	83.6
0.01	43.6	73.0	83.4

(e) **Distillation loss weight.** A small loss weight works well for MaskCLIP.

Weight	0-Shot	I2T/T2I
1	36.5	51.7/32.1
0.1	44.3	69.2/45.9
0.05	44.5	70.1/45.6
0.01	43.2	70.6/45.6

(f) **MLM loss weight.** A small loss weight works better.

Table 6. **MaskCLIP ablation experiments** with YFCC-15M dataset. We report zero-shot(0-Shot), fine-tuning (FT), and linear probing (Lin) accuracy (%) for image-encoder-related ablation. And zero shot image-to-text, text-to-image retrieval (I2T/T2I) for text encoder-related ablations. Default settings are marked in gray .

direction, resulting in poor performance. When we reduce the loss weight by $10\times$, the two additional tasks are helpful for the model and show a consistent gain on all the metrics. We suspect this is because the CLIP loss requires two different capabilities: understanding the input content and aligning them into a shared feature space. And the goal of the two additional self-supervised learning tasks is to facilitate understanding.

Image & Text decoder depth. Then we study the influence of the decoder depth for both image and text decoders. As shown in Table.6c, we find the image decoder with only one layer works well, increasing the decoder depth leads to worse performance on all metrics. Similarly, Table.6d shows that the text branch benefits from a shallow decoder design. We argue that a too-deep decoder would make the encoder lazy, relying on the strong decoder to resolve the challenging mask feature/language modeling tasks. And the different depth choice between the image and text branches is caused by the framework difference: the image branch sees the mask tokens at the decoder, while the text branch takes the mask tokens as the encoder input. Note that if we remove the text decoder, the performance gets worse. We think this is largely caused by the output conflict that the global recognition feature aggregation and local word prediction are conducted at the same layer.

Single-Stage v.s. two-Stage. Our MaskCLIP learns the VL contrastive and masked self-distillation simultaneously and jointly in a single stage. One possible variant is to first train CLIP and then use CLIP feature from the first stage to train masked image modeling as in [66, 67]. We report results on three datasets in Table 7. We can see that the second stage achieves better finetuning results compared with results from the stage one, showing the effectiveness of masked image modeling. Nonetheless, such two-stage training requires longer training time and loses the transfer capability in a zero-

Method	Epoch	IN-1K		Flicker30K		ADE20K	
		0-shot	FT	I2T	T2I	0-shot	FT
Two-Stage	Stage1	25	37.6 82.3	52.9 32.8	7.2 47.8		
	Stage2	25	— 83.4	— —	— 48.2		
MaskCLIP	25	44.5	83.6	70.1	45.6	10.2	50.5

Table 7. Comparison between two-stage method and our single-stage MaskCLIP.

shot setting. In contrast, our MaskCLIP achieves superior results under all settings with fewer epochs.

5. Conclusion

We present MaskCLIP, a new VL pretraining framework that incorporates masked self-distillation into VL contrastive. We point out that masked self-distillation learns local semantics, fitting nicely to the VL contrastive that aims to learn global semantics, and this is supported with comprehensively designed experiments. We also utilize mask language modeling to enhance the text encoder which is critical for zero-shot performance. The resulting visual encoder shows strong transfer capability across widely adopted benchmarks for linear probing, fine-tuning, and also zero-shot evaluation.

References

- [1] <https://eval.ai/web/challenges/challenge-page/1832/overview>.
- [2] <https://eval.ai/web/challenges/challenge-page/1832/leaderboard/4298>.
- [3] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. Data2vec: A general framework for self-supervised learning in speech, vision and language. *arXiv preprint arXiv:2202.03555*, 2022.

- [4] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.
- [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.
- [8] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 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*, pages 3558–3568, 2021.
- [9] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [11] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020.
- [12] Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised representation learning. *arXiv preprint arXiv:2202.03026*, 2022.
- [13] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [14] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning, 2020.
- [15] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021.
- [16] 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*, pages 104–120. Springer, 2020.
- [17] Ruizhe Cheng, Bichen Wu, Peizhao Zhang, Peter Vajda, and Joseph E Gonzalez. Data-efficient language-supervised zero-shot learning with self-distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3119–3124, 2021.
- [18] MMSegmentation Contributors. Mmsegmentation, an open source semantic segmentation toolbox. <https://github.com/open-mmlab/mmssegmentation>, 2020.
- [19] Yufeng Cui, Lichen Zhao, Feng Liang, Yangguang Li, and Jing Shao. Democratizing contrastive language-image pre-training: A clip benchmark of data, model, and supervision. *arXiv preprint arXiv:2203.05796*, 2022.
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [21] Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11162–11173, 2021.
- [22] Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Peco: Perceptual codebook for bert pre-training of vision transformers. *arXiv preprint arXiv:2111.12710*, 2021.
- [23] Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Bootstrapped masked autoencoders for vision bert pretraining. *arXiv preprint arXiv:2207.07116*, 2022.
- [24] Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Bootstrapped masked autoencoders for vision bert pretraining. *arXiv preprint arXiv:2207.07116*, 2022.
- [25] 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. *arXiv preprint arXiv:2010.11929*, 2020.
- [26] Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke Li, and Chunhua Shen. Pyramidclip: Hierarchical feature alignment for vision-language model pretraining. *arXiv preprint arXiv:2204.14095*, 2022.
- [27] Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan A Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *arXiv preprint arXiv:2205.14459*, 2022.
- [28] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33:21271–21284, 2020.
- [29] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.
- [30] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.

- [31] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [32] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015.
- [33] Xinyu Huang, Youcai Zhang, Ying Cheng, Weiwei Tian, Ruiwei Zhao, Rui Feng, Yuejie Zhang, Yaqian Li, Yandong Guo, and Xiaobo Zhang. Idea: Increasing text diversity via online multi-label recognition for vision-language pre-training. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4573–4583.
- [34] 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*, pages 4904–4916. PMLR, 2021.
- [35] Kyungyul Kim, ByeongMoon Ji, Doyoung Yoon, and Sangheum Hwang. Self-knowledge distillation: A simple way for better generalization. *arXiv preprint arXiv:2006.12000*, 3, 2020.
- [36] 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*, volume 34, pages 11336–11344, 2020.
- [37] Juncheng Li, Xin He, Longhui Wei, Long Qian, Linchao Zhu, Lingxi Xie, Yueting Zhuang, Qi Tian, and Siliang Tang. Fine-grained semantically aligned vision-language pre-training. *arXiv preprint arXiv:2208.02515*, 2022.
- [38] 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, 2021.
- [39] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. Hero: Hierarchical encoder for video+ language omni-representation pre-training. *arXiv preprint arXiv:2005.00200*, 2020.
- [40] Liunan 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*, 2019.
- [41] 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. *arXiv preprint arXiv:2012.15409*, 2020.
- [42] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020.
- [43] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. In *International Conference on Learning Representations*, 2022.
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [45] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [46] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.
- [47] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [48] Haoyu Lu, Nanyi Fei, Yuqi Huo, Yizhao Gao, Zhiwu Lu, and Ji-Rong Wen. Cots: Collaborative two-stream vision-language pre-training model for cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15692–15701, 2022.
- [49] 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.
- [50] 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*, pages 10437–10446, 2020.
- [51] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. Univl: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*, 2020.
- [52] Christian M Meyer and Iryna Gurevych. *Wiktionary: A new rival for expert-built lexicons? Exploring the possibilities of collaborative lexicography*. na, 2012.
- [53] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 891–898, 2014.
- [54] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. *arXiv preprint arXiv:2112.12750*, 2021.
- [55] 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.
- [56] 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*, pages 8748–8763. PMLR, 2021.

- [57] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypemymed, 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)*, pages 2556–2565, 2018.
- [58] Sheng Shen, Chunyuan Li, Xiaowei Hu, Yujia Xie, Jianwei Yang, Pengchuan Zhang, Anna Rohrbach, Zhe Gan, Lijuan Wang, Lu Yuan, et al. K-lite: Learning transferable visual models with external knowledge. *arXiv preprint arXiv:2204.09222*, 2022.
- [59] Mustafa Shukor, Guillaume Couairon, and Matthieu Cord. Efficient vision-language pretraining with visual concepts and hierarchical alignment. *arXiv preprint arXiv:2208.13628*, 2022.
- [60] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. V1-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.
- [61] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- [62] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017.
- [63] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016.
- [64] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [65] Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Yu-Gang Jiang, Luwei Zhou, and Lu Yuan. Bevt: Bert pretraining of video transformers. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2022)*, 2022.
- [66] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. *arXiv preprint arXiv:2112.09133*, 2021.
- [67] Longhui Wei, Lingxi Xie, Wengang Zhou, Houqiang Li, and Qi Tian. Mvp: Multimodality-guided visual pre-training. *arXiv preprint arXiv:2203.05175*, 2022.
- [68] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 418–434, 2018.
- [69] Haiyang Xu, Ming Yan, Chenliang Li, Bin Bi, Songfang Huang, Wenming Xiao, and Fei Huang. E2e-vlp: End-to-end vision-language pre-training enhanced by visual learning. *arXiv preprint arXiv:2106.01804*, 2021.
- [70] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal self-attention for local-global interactions in vision transformers. *arXiv preprint arXiv:2107.00641*, 2021.
- [71] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Bin Xiao, Ce Liu, Lu Yuan, and Jianfeng Gao. Unified contrastive learning in image-text-label space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19163–19173, 2022.
- [72] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *arXiv preprint arXiv:2111.07783*, 2021.
- [73] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.
- [74] 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.
- [75] Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Yangyu Huang, Lu Yuan, Dong Chen, Ming Zeng, and Fang Wen. General facial representation learning in a visual-linguistic manner. *arXiv preprint arXiv:2112.03109*, 2021.
- [76] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017.
- [77] Chong Zhou, Chen Change Loy, and Bo Dai. Denseclip: Extract free dense labels from clip. *arXiv preprint arXiv:2112.01071*, 2021.
- [78] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. Image bert pre-training with online tokenizer. In *International Conference on Learning Representations*, 2021.
- [79] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13041–13049, 2020.

A. More Experiment

Comparison over small model and small dataset. As some baselines report ViT-B/32 instead of ViT-B/16, in order to compare, we further experiment MaskCLIP with a smaller model ViT-B/32 and report the zero-shot performance on ImageNet-1K. As shown in Table 8 left, our MaskCLIP outperforms the combination [19] of two recent strong methods DeCLIP [43] and FILIP [72]. We also investigate the performance on a smaller dataset CC3M [57] (we use ViT-B/16 here in coherency with previous experiments). Table 8(right part) shows that MaskCLIP achieves consistent gain.

ViT-B/32	0-Shot	Lin	FT	CC3M	0-Shot	Lin	FT
CLIP	26.1	60.5	74.3	CLIP	17.1	53.3	78.5
DeFILIP	36.4	—	—	SLIP	23.0	65.4	81.4
MaskCLIP	38.5	69.1	79.2	MaskCLIP	24.4	66.1	82.5

Table 8. Results of zero-shot performance on ImageNet-1K when pretrained with ViT-B/32 model(left) or CC3M dataset(right) .

Ablation on distillation loss. Here we further study the effectiveness of each component in the distillation loss. We start from CLIP+MAE and add three components of the distillation loss one by one. We find that 1) using the feature as the prediction target improves all metrics; 2) using EMA model gets better performance; 3) the MLM loss improves all the vision-language tasks.

	0-Shot	FT	Seg	Det	I2T/T2I
CLIP+MAE (baseline)	42.1	83.2	49.1	44.5/40.4	57.3/41.1
+ Feature prediction	42.6	83.4	49.9	45.1/40.6	62.3/41.4
+ EMA model	42.8	83.6	50.4	45.5/40.9	65.0/41.6
+ MLM loss	44.5	83.6	50.5	45.4/40.9	70.1/45.6

Table 9. Component ablation of the distillation loss.

B. Experiment detail

Pre-training We train our proposed MaskCLIP from scratch and training for 25 epochs, the batch size is fixed to 4096 for all the experiments. We use 32 V100 for training with 128 samples per GPU. We use the AdamW [47] optimizer with weight decay 0.1. The learning rate is set to $1e^{-3}$ with one epoch warm-up and decay to $1e^{-5}$ followed by a cosine schedule. The masks used in the mask self-distillation branch are random mask with a mask ratio of 75%. The EMA weight is set to 0.999 and linearly increases to 0.9999 during the training. We pretrain all the models with the commonly used YFCC15M dataset, which is flited from the YFCC100M [63] dataset by [56].

For the ICinW academic track experiment, we pre-train the model with three datasets: YFCC-15M [63], GCC3M [57]+12M [8] and ImageNet-21K [20] (ImageNet-1K data is excluded). Here we use the UniCL [71] to utilize the ImageNet-22k dataset in the pretraining with a unified

format. We train the model for 32 epochs and 16384 batch size, the rest settings are the same as the YFCC15M setting.

Zero-shot ImageNet-1K classification. For zero-shot on ImageNet-1K, we follow the prompt setting in [54] to convert the labels to text features, which contains 7 prompt templates and we use the average feature as the final label feature. We calculate the similarity between image feature and all the label features to get its zero-shot classification result.

Linear-probing ImageNet-1K classification. For linear probing, we fix the backbone and train a new linear classifier for 90 epochs. Following the setting in MAE [29], we add a batch-norm layer without learnable affine parameters before the classifier to avoid adjusting the learning rate for each model. We set the batch size to 16384 and use the LARS [73] optimizer with weight decay 0 and momentum 0.9. The learning rate is set to 6.4 and decays to 0 following the cosine schedule.

Fine-tuning ImageNet-1K classification. When fine-tuning on the ImageNet-1K dataset, we average pool the output of the last transformer of the encoder and feed it to a softmax-normalized classifier. We fine-tune 100 epochs for all the experiments, the learning rate is warmed up to 0.0006 for 20 epochs and decay to $1e^{-6}$ following the cosine schedule. Similar to recent works, we also apply the layer decayed learning rate used in [4] and we set the decay factor as 0.7. Note that we use the pure ViT architecture, *without* the techniques used in [4], such as layer scale and relative position embedding. The evaluation metric is top-1 validation accuracy of a single 224×224 crop.

Zero-shot Semantic segmentation. Here we follow the setting in DenseCLIP [77] based on the implementation from mmsegmentaion [18]. For ADE20K and MS-COCO, we report the single-scale test result with 512×512 input. For Pascal Context, we use 480×480 input. To avoid the influence of position embedding caused by changing input size, we use sliding inference with 224×224 input and stride 112. To convert the labels to text embedding, we use 85 prompt templates and use the average feature as the final label feature.

ADE20K Semantic segmentation. Here we use: UperNet [68] based on the implementation from mmsegmentaion [18]. For UperNet, we follow the settings in [4] and use AdamW [47] optimizer with initial learning rate $2e^{-4}$, weight decay of 0.05 and batch size of 16 (8 GPUs with 2 images per GPU) for 160K iterations. The learning rate warmups with 1500 iterations at the beginning and decays with a linear decay strategy. We use the layer decay [4] for the backbone and we set it as 0.6. As the ViT architecture outputs features with the same size, here we add four different scale FPNs to scale the feature map into different size. Specifically, we upsample the output feature of the 4th block $4\times$, upsample the output feature of the 6th block

$2\times$, keep the output feature of the $8th$ block unchanged and downsample the output feature of the $12th$ block $2\times$. We use the default augmentation setting in mmsegmentation including random horizontal flipping, random re-scaling (ratio range $[0.5, 2.0]$) and random photo-metric distortion. All the models are trained with input size 512×512 . The stochastic depth is set to 0.1. When it comes to testing, we report single-scale test result.

COCO Object Detection and Instance Segmentation.

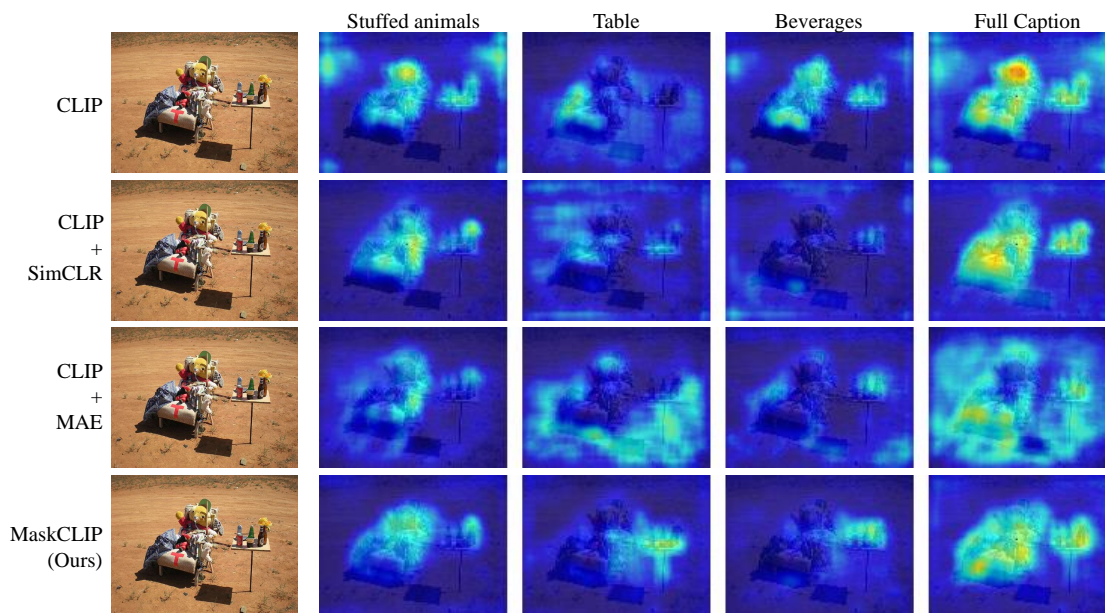
We use the classical object detection framework Mask R-CNN [31] based on the implementation from mmdetection [9]. We train it the $1\times$ schedule with single-scale input (image is resized so that the shorter side is 800 pixels, while the longer side does not exceed 1333 pixels) for 12 epochs. We use AdamW [47] optimizer with a learning rate of $1e^{-4}$, weight decay of 0.05 and batch size of 16. We also use the layer decay [4] for the backbone and we set it as 0.75. The learning rate declines at the $8th$ and $11th$ epoch with decay rate being 0.1. The stochastic depth is set to 0.1. Similar to the implementation of semantic segmentation above, we also use four different scale FPNs to scale the feature map into different size.

C. More visualization results.

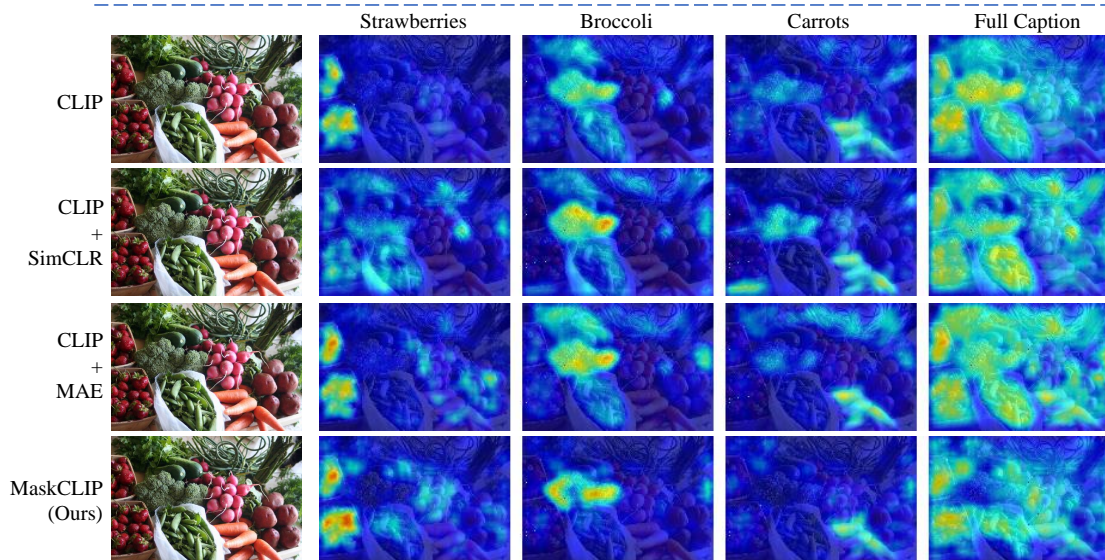
Here we provide more visualization results on the MSCOCO val set. In most cases, our MaskCLIP gets a better feature alignment performance between image and text.

D. Societal impacts

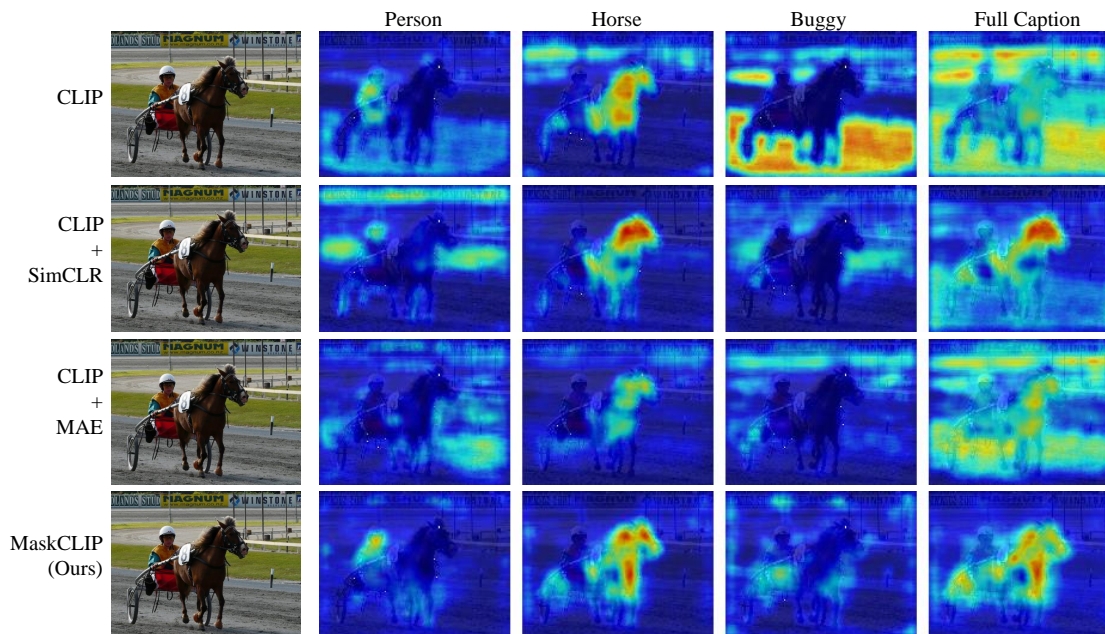
MaskCLIP is an improvement of CLIP, so it has the same societal impacts of CLIP, including some malicious usages and positive applications. Meanwhile, CLIP and MaskCLIP may suffer from some unwanted data bias, as the data used for training are roughly collected from the Internet.



Large stuffed animal posed outdoors as if sitting in a chair with beverages on a table.



various fruits and vegetables are on display close together.



A person in a buggy drawn by a horse.

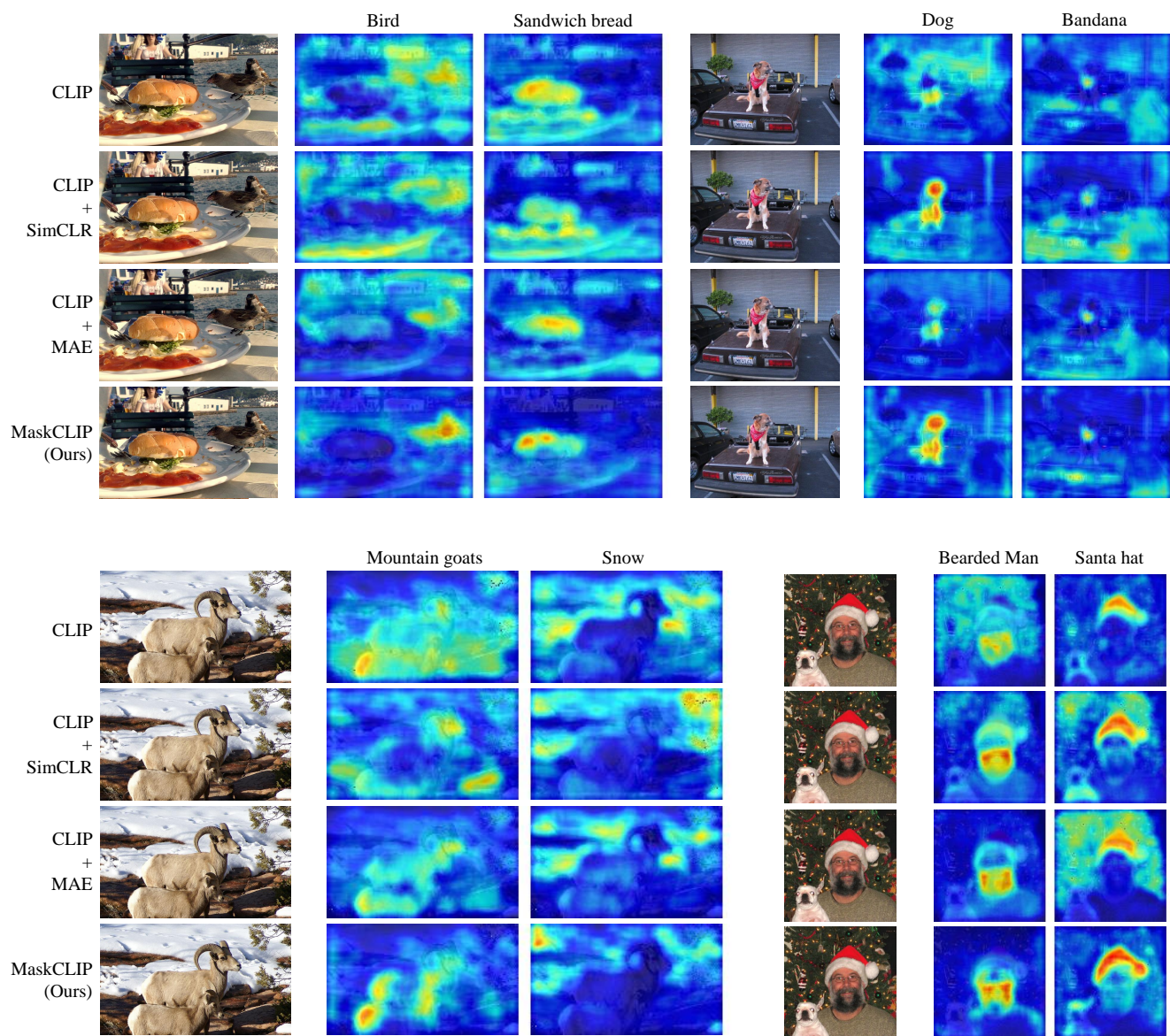


Figure 4. Visualization of the similarity between words and image features. The images and captions are from the MS-COCO val set.

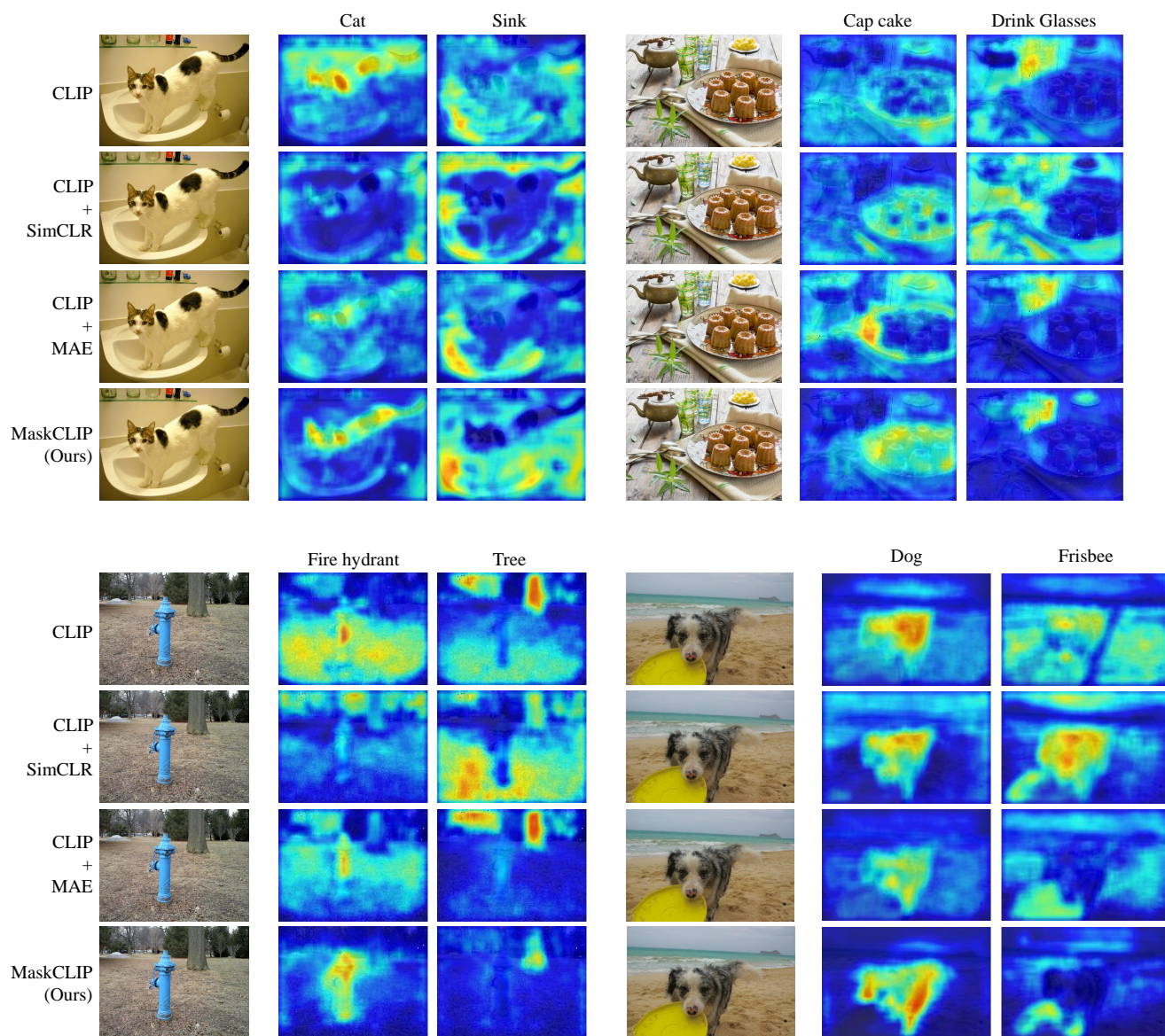


Figure 5. Visualization of the similarity between words and image features. The images and captions are from the MS-COCO val set.