

Zero-Shot Detection via Vision and Language Knowledge Distillation

Xiuye Gu* Tsung-Yi Lin Weicheng Kuo Yin Cui

Google Research

{xiuyegu, tsungyi, weicheng, yincui}@google.com

Abstract

*Zero-shot image classification has made promising progress by training the aligned image and text encoders. The goal of this work is to advance zero-shot object detection, which aims to detect novel objects without bounding box nor mask annotations. We propose **ViLD**, a training method via **Vision and Language knowledge Distillation**. We distill the knowledge from a pre-trained zero-shot image classification model (e.g., CLIP [33]) into a two-stage detector (e.g., Mask R-CNN [17]). Our method aligns the region embeddings in the detector to the text and image embeddings inferred by the pre-trained model. We use the text embeddings as the detection classifier, obtained by feeding category names into the pre-trained text encoder. We then minimize the distance between the region embeddings and image embeddings, obtained by feeding region proposals into the pre-trained image encoder. During inference, we include text embeddings of novel categories into the detection classifier for zero-shot detection. We benchmark the performance on LVIS dataset [15] by holding out all rare categories as novel categories. ViLD obtains 16.1 mask AP_r with a Mask R-CNN (ResNet-50 FPN) for zero-shot detection, outperforming the supervised counterpart by 3.8. The model can directly transfer to other datasets, achieving 72.2 AP₅₀, 36.6 AP and 11.8 AP on PASCAL VOC, COCO and Objects365, respectively.*

1. Introduction

Consider the image in Figure 1, can we design object detectors beyond recognizing only *base categories* (e.g., *toy*) present in training labels and expand the vocabulary to detect *novel categories* (e.g., *toy elephant*)?

Existing object detection algorithms often learn to detect only the categories present in the training data. A common approach to increase the detection vocabulary is by collecting images with more labeled categories. The research community has recently collected new object detection datasets

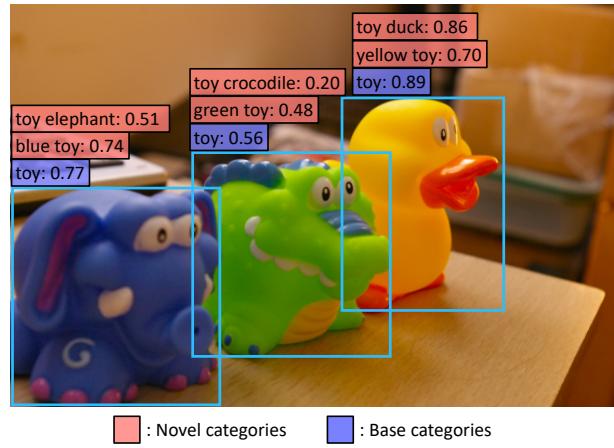


Figure 1: **An example of our zero-shot object detection with free-form text classifiers.** After training our zero-shot detector on base categories (purple), we can use the text embeddings of novel categories (pink) to detect novel object categories that do not exist in the training dataset.

with large vocabularies [15, 25]. LVIS [15] is a milestone of these efforts by building a dataset with 1203 categories. With such a rich vocabulary, it becomes quite challenging to collect enough training examples for all categories. By Zipf’s law, object categories naturally follow a long-tailed distribution. To find sufficient training examples for rare categories, exponentially more data is needed [15]. This implies creating a large vocabulary detection dataset is expensive to scale up.

On the other hand, paired images and texts are abundant on the Internet. Recently, Radford *et al.* [33] train a joint vision and language model using 400 million image and text pairs and demonstrate impressive zero-shot recognition abilities on over 30 computer vision datasets. Despite the great success on learning image-level representations for zero-shot classification, learning object-level representations for zero-shot object detection is still challenging. In this work, we consider the idea of borrowing the knowledge from a pre-trained zero-shot image classification model to enable zero-shot object detection.

*Work done as a member of the Google AI Residency program.

We begin with a straightforward approach to use a zero-shot image classification model. We turn zero-shot detection into two sub-problems: 1) generalized object proposal and 2) zero-shot image classification. This approach can be seen as a variant of R-CNN [13]. We train a region proposal model using examples in the base categories. Then we use the pre-trained zero-shot image classification model to classify cropped object proposals, which contain both base and novel categories. We benchmark the performance on LVIS [15] by holding out all rare categories as novel categories and treat other categories as base categories. To our surprise, the performance on the novel categories is already on par with the supervised counterpart. However, this approach is very slow for inference because it requires to feed all object proposals one-by-one into the classification model. Moreover, the performance on the base categories is much worse.

We propose **ViLD** (**Vision and Language knowledge Distillation**) for training two-stage zero-shot detectors. ViLD consists of two components: learning with the text embeddings (ViLD-text) and the image embeddings (ViLD-image) inferred by a zero-shot image classification model (*e.g.*, CLIP [33]). In **ViLD-text**, we obtain the text embeddings by feeding the category text prompts into the pre-trained text encoder. We then replace the trainable detection classifier with the fixed text embeddings. Similar approaches have been used in prior zero-shot detection work [3, 35]. However, these existing methods only use text embeddings learned from a language corpus, *e.g.*, [32]. In contrast, we find text embeddings learned jointly with visual data, *e.g.*, [33], can better represent the visual similarity between text prompts. For example, “bucket” and “trash can” can be far away in GloVe, but close in CLIP due to the shape similarity. Using CLIP text embeddings attains **10.1** AP on novel categories in LVIS, significantly outperforming the **3.0** AP of using GloVe. In **ViLD-image**, we obtain the image embeddings by feeding the object proposals, which contain objects with or without category labels, into the pre-trained image encoder. We align the region embeddings in a two-stage object detector to these image embeddings. In contrast to ViLD-text, ViLD-image provides supervision to object proposals that do not have category labels.

We show that ViLD, a joint learning from text and image embeddings, achieves **16.1** AP on novel categories in LVIS, surpassing the supervised learning counterpart by **3.8**. We obtain the best performance of **22.6** AP by ensembling predictions of ViLD-text and CLIP on the proposals. In addition, the detector trained with ViLD on LVIS can transfer to other datasets including Objects365, COCO, PASCAL VOC Detection. Finally, we qualitatively demonstrate the model is able to detect objects with novel fine-grained categories and attributes as shown in Figure 1.

We summarize our contributions as follows:

- We propose ViLD, a zero-shot detection method that distills the knowledge in a zero-shot image classification model.
- ViLD is the first zero-shot detection method being evaluated on the challenging LVIS dataset. It surpasses the supervised learning counterpart on novel categories.
- Our findings suggest a new direction for handling large vocabulary detection, apart from scaling up datasets with long-tailed instance annotations.

2. Related Work

Zero-shot visual recognition: The exploration of zero shot learning in the vision community starts from recognition. Zero-shot visual recognition aims to classify images whose categories are never seen during training. Prior to the breakthrough of deep learning, earlier work study the use of attributes to encode categories as vectors of binary attributes [22, 10, 38] and learn label embeddings [2]. Researchers have also explored class hierarchy, class similarity, and object parts as discriminative features to aid the knowledge transfer [38, 1, 8, 23, 4, 44]. Another direction focuses on learning the representation of semantic classes to serve as a visual classifier. DeVise [11] learns a deep visual-semantic embedding space, and performs nearest neighbor search in the word embedding space. ConSE [31] builds upon the idea of DeVise, but instead uses an existing image classifier and a word embedding model to build the embedding system. Wang *et al.* [43] distills information from both word embeddings and knowledge graphs. In a very recent work, CLIP [33] pushes the limit by first collecting 400 million image-text pairs and then training a joint image-text embedding model via contrastive learning. Equipped with vast image-text knowledge from the web, CLIP advances state-of-the-art zero-shot recognition on a suite of visual recognition benchmarks. CLIP is closely related to our work as we distill its knowledge for zero-shot detection. While these work focus on image-level recognition, ours focuses on detecting objects of novel categories in complex scenes.

Large vocabulary object detection: There has been growing interest in large vocabulary detection [25, 15]. Natural object categories follow Zipf’s law, and thus rare categories have much fewer examples than common ones. Existing object detection approaches [37, 29, 17] does not perform well for rare categories. Specialized methods are proposed to address the data imbalance issue, including two-stage training [24, 26], loss re-weighting [19, 6, 5, 40, 26], and data re-sampling [30, 15, 42]. Unlike these methods, our proposed approach performs general zero-shot detection, which doesn’t require any annotation for novel categories. Therefore, our method can be directly applied to long-tailed categories without using any specially designed techniques.

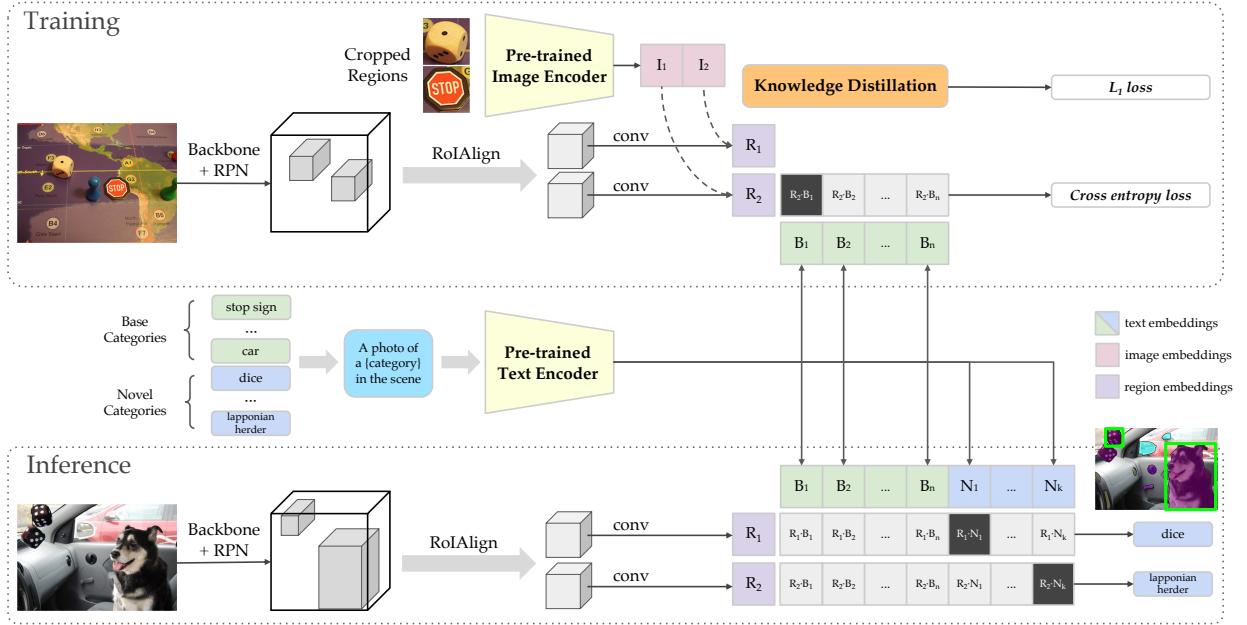


Figure 2: **An overview of using ViLD for zero-shot object detection.** ViLD distills the knowledge from a pre-trained zero-shot image classification model. First, the category text embeddings and the image embeddings of cropped object proposals are computed using the text and image encoders in the classification model. Then, ViLD uses the text embeddings as the region classifier (ViLD-text) and minimizes the distance of the region embedding to the image embedding for each proposal (ViLD-image). During inference, text embeddings of novel categories are used to enable zero-shot detection.

Zero-shot object detection: The field of zero-shot detection (ZSD) is less-explored. Most ZSD approaches leverage visual and text embeddings. Bansal *et al.* [3] adapt visual-semantic embeddings for ZSD and introduce two background-aware approaches. In addition to visual-semantic embeddings, Rahman *et al.* [35] design a loss for both max-margin class separation and semantic clustering. Demirel *et al.* [7] introduce a hybrid model using convex combination of embeddings. Rahman *et al.* [34] adopt a transductive setting to reduce the domain-shift against unseen classes. Hayat *et al.* [16] propose to generate and separate visual features using class semantics. Zheng *et al.* [45] propose a cascade model and learn background classes. There is also a different line of work [46, 21] focusing only on generating region proposals for novel categories. Unlike prior work, to the best of our knowledge, the proposed zero-shot detection method ViLD is the first to distill knowledge from a zero-shot image classification model. We are also the first to benchmark ZSD on the challenging large-scale LVIS dataset [15] and compare with supervised learning.

3. Method

In this section, we introduce three main components in our approach: localization of novel objects, learning from pre-trained text embeddings (ViLD-text), and learning from pre-trained image embeddings (ViLD-image). Figure 2 il-

lustrates an overview of our method. We also study model ensembling approaches for best detection performance.

Settings for zero-shot detection: We have a set of base categories C_B that the model can be trained on, and another set of novel categories C_N that are never seen during training. The main goal of zero-shot detection is to achieve accurate object detection on novel categories.

Zero-shot image classification model: We achieve zero-shot detection by leveraging an off-the-shelf pre-trained zero-shot image classification model (*e.g.* CLIP [33]). The model has a text encoder $\mathcal{T}(\cdot)$ and an image encoder $\mathcal{V}(\cdot)$, which are pre-trained by joint image-text contrastive learning. In ViLD, we do not update these encoders during training, as shown in Figure 2.

3.1. Object proposals for novel categories

The first challenge for zero-shot detection is to localize novel objects. We modify a two-stage object detector (*e.g.*, Mask R-CNN [17]) to detect object proposals with bounding boxes and masks. We replace class-specific localization modules, *i.e.*, the second stage bounding box regression and mask prediction layers, with class-agnostic modules for general object proposals. For each region of interest, these modules only predict a single bounding box and a single mask for all possible categories, instead of one prediction for each category.

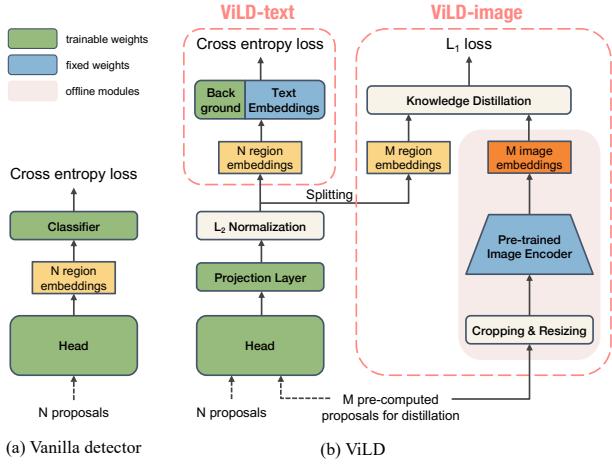


Figure 3: Model architecture and training objectives. The classification head of (a) a vanilla two-stage detector, e.g., Mask R-CNN, and (b) ViLD. ViLD replaces the vanilla classifier with fixed category text embeddings while keeping a learnable background embedding, and use an \mathcal{L}_1 distillation loss to enforce the region embeddings to be similar to the image embeddings. A projection layer and a normalization layer are introduced to adjust the norm and dimension of region embeddings in order to be compatible with the fixed text embeddings.

3.2. Zero-shot detection with cropped regions

Once object candidates are localized, a straightforward approach is to use a pre-trained zero-shot image classifier to classify each region. In particular, we filter out the proposals with high background probabilities and apply non-maximum suppression (NMS) to obtain the top-k regions. We crop and resize the proposals ($\text{crop}(I, r)$) as inputs to the classifier, where I is the image and r is the bounding box of a proposal. The image embeddings inferred by the image encoder is $\mathcal{V}(\text{crop}(I, r))$. We ensemble the image embeddings from $1\times$ and $1.5\times$ crops for each proposal, as the $1.5\times$ crop provides more context cues. The ensembled embedding is then renormalized to unit norm:

$$\begin{aligned} \mathbf{v} &= \mathcal{V}(\text{crop}(I, r_{1\times}) + \mathcal{V}(\text{crop}(I, r_{1.5\times})) \\ \mathcal{V}(\text{crop}(I, r_{\{1\times, 1.5\times\}})) &= \frac{\mathbf{v}}{\|\mathbf{v}\|}. \end{aligned} \quad (1)$$

After that, we apply per class NMS and take regions with top K confidence scores. This method has a slow inference time because every cropped region proposal is fed into the image classification model. In addition, it does not make use of base category labels in the detection dataset.

3.3. Replacing classifiers with text embeddings

In **ViLD**, we learn region embeddings in a two-stage detector to represent each proposal. We define region embeddings as $\mathcal{R}(\phi(I), r)$, where $\phi(\cdot)$ is the backbone model and

\mathcal{R} is a lightweight model to generate region embeddings for each proposal r . Specifically, we take outputs of the layer before the detection classifier as region embeddings. Our goal is to train the region embeddings such that they can be classified with the text embeddings encoded by $\mathcal{T}(\cdot)$.

We first introduce **ViLD-text**. Figure 3(b) shows the model architecture and the training objective. ViLD-text replaces the learnable classifier in Figure 3(a) with the text embeddings. For training, we generate the text embeddings $\mathcal{T}(C_B)$ by feeding the text prompts of base categories (C_B), e.g., “a photo of {category} in the scene”, into the text encoder. For the proposals that do not match to any labeled base categories, they are assigned to the “background” category. Since the text prompt of “background” does not well represent these unmatched proposals, we allow the “background” category to learn its own embedding (\mathbf{e}_{bg}). We compute the cosine similarity between each region embedding $\mathcal{R}(\phi(I), r)$ and all category embeddings, including $\mathcal{T}(C_B)$ and \mathbf{e}_{bg} . Then we apply softmax activation with a temperature τ to these cosine similarities to compute the cross entropy loss. Let $\text{sim}(\mathbf{a}, \mathbf{b}) = \mathbf{a}^\top \mathbf{b} / (\|\mathbf{a}\| \|\mathbf{b}\|)$, \mathbf{t}_i denote elements in $\mathcal{T}(C_B)$, y_r denote the class label of the region r , and \mathcal{L}_{CE} denote the cross entropy loss. The loss function for ViLD-text can be written as:

$$\begin{aligned} \mathbf{e}_r &= \mathcal{R}(\phi(I), r) \\ \mathbf{z}(r) &= [\text{sim}(\mathbf{e}_r, \mathbf{e}_{bg}), \text{sim}(\mathbf{e}_r, \mathbf{t}_1), \dots, \text{sim}(\mathbf{e}_r, \mathbf{t}_{|C_B|})] \\ \mathcal{L}_{\text{ViLD-text}} &= \mathcal{L}_{CE}(\text{softmax}(\mathbf{z}(r)/\tau), y_r). \end{aligned} \quad (2)$$

During inference, we include novel categories (C_N) and generate $\mathcal{T}(C_B \cup C_N)$ (sometimes $\mathcal{T}(C_N)$ only) for zero-shot detection (Figure 2). Our hope is that the model learned from labeled C_B can generalize to novel categories C_N .

3.4. Distilling image embeddings

We then introduce **ViLD-image**, which aims to align region embeddings $\mathcal{R}(\phi(I), r)$ to image embeddings $\mathcal{V}(\text{crop}(I, r))$, introduced in Section 3.2. The goal is to distill the knowledge in the teacher image encoder \mathcal{V} into the student detector. We extract proposals \tilde{r} offline from the training images using a region proposal network pre-trained on base categories. These proposals may contain objects in both C_B and C_N . In contrast, ViLD-text can only learn from C_B . We apply an \mathcal{L}_1 loss between the region and image embeddings to minimize their distance (we use the same ensembled image embeddings as in Section 3.2):

$$\mathcal{L}_{\text{ViLD-image}} = \|\mathcal{V}(\text{crop}(I, \tilde{r}_{\{1\times, 1.5\times\}})) - \mathcal{R}(\phi(I), \tilde{r})\|_1. \quad (3)$$

The training loss of **ViLD** is simply a weighted sum of both objectives:

$$\mathcal{L}_{\text{ViLD}} = \mathcal{L}_{\text{ViLD-text}} + w \cdot \mathcal{L}_{\text{ViLD-image}}, \quad (4)$$

where w is a hyperparameter weight for distilling the image embeddings. Figure 3(b) shows the model architecture and training objectives.

The difference between ViLD-text and ViLD-image is only in the training, where ViLD-text and ViLD-image are trained with $\mathcal{L}_{\text{ViLD-text}}$ and $\mathcal{L}_{\text{ViLD-image}}$ respectively. During inference, ViLD-image, ViLD-text and ViLD share the same model architecture for zero-shot detection.

3.5. Model ensembling

In this section, we explore model ensembling for the best detection performance over base and novel categories.

First, we combine the predictions of a ViLD-text detector with the zero-shot image classification model. The intuition is that ViLD-image learns to approximate the predictions of its zero-shot teacher model, and therefore, we assume using the teacher model directly may improve performance. We use a trained ViLD-text detector to obtain a set of candidate regions and their category confidence scores. We then filter out background regions and apply NMS to obtain the top-k proposals. We use $p_{i,\text{ViLD-text}}$ to denote the confidence scores for each proposal i . We then feed $\text{crop}(I, r)$ to the zero-shot classification model to obtain confidence scores $p_{i,\text{cls}}$. Since we know the two models have different performance on base and novel categories, we introduce a weighted geometric average in the ensemble:

$$p_{i,\text{ensemble}} = \begin{cases} p_{i,\text{ViLD-text}}^\lambda \cdot p_{i,\text{cls}}^{(1-\lambda)}, & \text{if } i \in C_B \\ p_{i,\text{ViLD-text}}^{1-\lambda} \cdot p_{i,\text{cls}}^\lambda, & \text{if } i \in C_N \end{cases} \quad (5)$$

λ is to 2/3, which weighs the prediction of ViLD-text more for base categories and vice versa. Note this approach has a similar slow inference speed as the method in Section 3.2.

Further, we introduce a different ensembling approach to mitigate the inference speed issue of the method above. In ViLD, the cross entropy loss of ViLD-text and the \mathcal{L}_1 distillation loss of ViLD-image is applied to the same region embedding, which may cause contentions. Here, instead, we learn two embeddings for the same region by using two independent heads with the same backbone architecture for ViLD-text and ViLD-image. Text embeddings are applied to these two regions embeddings to obtain confidence scores $p_{i,\text{ViLD-text}}$ and $p_{i,\text{ViLD-image}}$, which are then ensembled using the same way as in Eq. 5 with $p_{i,\text{ViLD-image}}$ replacing $p_{i,\text{cls}}$. We name this approach **ViLD-ensemble**.

4. Experiments

Implementation details: We benchmark on the Mask R-CNN [17] with ResNet [18] FPN [27] backbone and use the same settings for all models unless explicitly specified. All models use 1024×1024 as input image size, large-scale jittering augmentation [12], synchronized batch normalization [20, 14] of batch size 256, weight decay of 4e-5, and an

initial learning rate of 0.32. We train the model from scratch for 180,000 iterations, and divide the learning rate by 10 at $0.9\times$, $0.95\times$, and $0.975\times$ of total iterations. For all our experiments involving a zero-shot classification model, we use the pre-trained CLIP model that is publicly available¹.

4.1. Benchmark settings

The zero-shot detection is a relatively new area. We compare our setting with Bansal *et al.* [3], which use COCO (80 categories) as the major benchmark for evaluation.

Our setting: We train and benchmark zero-shot detection on LVIS [15] (1203 categories). LVIS contains a much larger and diverse set of vocabulary that is more suitable for zero-shot detection. We take its frequent and common categories as the base categories C_B , and hold out rare categories as the novel categories C_N . AP_r , the AP of rare categories, is used to measure the performance.

Comparison with Bansal *et al.* [3]: Bansal *et al.* manually re-sample and divide COCO into new train and validation splits. They use 48 categories for training and 17 categories for validation, removing 15 categories that does not have a synset in the WordNet hierarchy. They then create the training dataset with images only containing objects in the base categories. In contrast, we directly use the training and validation splits in the original LVIS dataset. In our setting, the novel categories C_N are rare categories that are harder to find training examples. On the contrary, in [3], the C_B and C_N are both common object categories from COCO. We remove all the annotations in C_N from the training set. This means training images can contain both objects in C_B and C_N , but only objects in C_B are annotated. We find our setup is more realistic for large datasets and easier to compare the zero-shot performance to the supervised detection baseline.

4.2. Learning object proposals with base categories

We first study whether a detector can predict bounding boxes of novel categories when only trained on base categories. We train the region proposal networks in Mask R-CNN for this experiment. Table 1 shows the average recall (AR) on novel categories. Training with only base categories performs slightly worse by around 2 AR at 100, 300, and 1000 proposals, compared to using both base and novel categories. This experiment demonstrates that, even without seeing novel categories during training, region proposal networks only suffers a small performance drop. We believe better region proposal networks focusing on unseen category generalization should further improve the performance, and leave this for future research.

¹<https://github.com/openai/CLIP>, ViT-B/32.

Supervision	$AR_r@100$	$AR_r@300$	$AR_r@1000$
base	39.3	48.3	55.6
base + novel	41.1	50.9	57.0

Table 1: **Training with only base categories achieves comparable average recall (AR) for novel categories in LVIS.** We compare region proposal networks trained with only base categories vs. base + novel categories and report the bounding box AR.

Method	AP_r	AP_c	AP_f	AP
Supervised (base class only)	0.0	22.6	32.4	22.5
CLIP on cropped regions	13.0	10.6	6.0	9.2
Supervised (base + novel)	4.1	23.5	33.2	23.9
Supervised-RFS (base + novel)	12.3	24.3	32.4	25.4

Table 2: **Using CLIP for zero-shot detection achieves high detection performance on novel categories.** We apply CLIP to classify cropped region proposals and report the mask average precision (AP). The performance on novel categories AP_r is on par with supervised learning approaches. However, the overall performance is still behind supervised methods.

(a) Mask R-CNN R50-FPN				
Method	AP_r	AP_c	AP_f	AP
CLIP on cropped regions	13.0	10.6	6.0	9.2
GloVe baseline	3.0	20.1	30.4	21.2
ViLD-text	10.1	23.9	32.5	24.9
ViLD-image	9.6	8.5	7.8	8.4
ViLD ($w = 0.5$)	16.1	20.0	28.3	22.5
ViLD-ensemble ($w = 0.5$)	16.6	24.6	30.3	25.5
ViLD-text + CLIP \dagger	22.6	24.8	29.2	26.1
Supervised-RFS (base + novel)	12.3	24.3	32.4	25.4

(b) Mask R-CNN R152-FPN				
Method	AP_r	AP_c	AP_f	AP
ViLD-text	11.7	25.8	34.4	26.7
ViLD-image	10.8	10.0	8.7	9.6
ViLD ($w = 1.0$)	18.7	21.1	28.4	23.6
ViLD-ensemble ($w = 2.0$)	18.7	24.9	30.6	26.0
Supervised-RFS (base + novel)	14.4	26.8	34.2	27.6

Table 3: **ViLD outperforms the supervised learning counterpart on novel categories.** The best performance is achieved by ensembling ViLD-text and CLIP but its runtime is **630 \times** of all other methods. Results are mask AP averaged over 3 runs.

4.3. Classifying proposals with CLIP

We explore the potential of a zero-shot classification model on zero-shot detection. Specifically, we use CLIP to classify detected region proposals.

Image embeddings: We crop and resize 1000 object proposals to 224×224 image crops, and then feed them into the CLIP model to obtain region embeddings. In this experiment, we treat the second-stage refined boxes as proposals and use the corresponding masks to report mask AP.

Text classifiers: For each of the 1203 categories in LVIS, we generate 63 text prompts and then feed them into CLIP’s text encoder and ensemble them to obtain the text embedding. We use the synonyms provided in the LVIS dataset. We compare it with the Supervised baseline trained on both base and novel categories, and Supervised-RFS [30, 15] that oversamples images containing infrequent categories. We set $t = 0.001$ in RFS following [15]. As shown in Table 2, its AP_r outperforms the Supervised baseline by a large margin and is on par with Supervised-RFS. However, CLIP performs much worse for AP_c and AP_f . This indicates directly using CLIP for object detection is still challenging.

4.4. Vision and language knowledge distillation

In this section, we evaluate the performance of detectors trained with ViLD and its variants ViLD-text, ViLD-image, and ViLD-ensemble. These models are significantly faster compared to the method in Section 4.3. Finally, we combine ViLD-text and CLIP to demonstrate our highest AP_r on novel categories. Table 3 summarizes the results.

Text embeddings as classifiers (ViLD-text): We evaluate the performance of ViLD-text (Section 3.3), which uses text-embeddings generated by CLIP. We compare its performance with using GloVe [32] text embeddings trained on large-scale text corpuses. Table 3 shows a Mask R-CNN R50-FPN model trained by ViLD-text attains 10.1 AP_r , which is significantly better than 3.0 AP_r using GloVe. This demonstrates the importance of using text embeddings that are jointly trained with visual data. ViLD-text achieves much higher AP_c and AP_f compared to CLIP on cropped regions (Section 4.3). This is because ViLD-text uses instance annotations in the base categories to align region embeddings and CLIP’s text embeddings. The AP_r is worse, showing that only using 866 base categories in LVIS does not generalize as well as CLIP to novel categories.

Distilling image embeddings (ViLD-image): We study the effectiveness of visual distillation in Section 3.4 and compare it with applying CLIP to classify object proposals. We only use the distillation loss to train our detectors without the classification loss. Ideally, the detection performance of ViLD-image and CLIP on cropped regions would be similar because we train region embeddings inferred by our detector to approximate the embeddings inferred by CLIP on the same proposals. As shown in Table 3, the performance gap is 3.4 AP_r . This demonstrates that visual distillation works for zero-shot detection but yields a small decrease in performance.

Text + visual embeddings (ViLD): Next, we investigate learning from the combination of text and image embeddings in CLIP. Table 4 shows the parameter sweep of different distillation weights using \mathcal{L}_1 and \mathcal{L}_2 losses. Overall, using both text and image embeddings yields good performance across AP_r , AP_c , and AP_f . We find \mathcal{L}_1 loss can better improve the AP_r performance with the trade-off to AP_c and AP_f . This suggests there is a competition between ViLD-text and ViLD-image. In Table 3, we compare ViLD with other ViLD variants and baselines. ViLD has already outperformed Supervised-RFS by 3.8, showing our zero-shot detection approach is better than supervised models on rare categories. Its AP_r is also 6.0 higher than ViLD-text, indicating that additionally learning from the image embeddings boosts the zero-shot performance.

Model ensembling: We study methods discussed in Section 3.5 to reconcile the conflict of joint training with ViLD-text and ViLD-image. We use two ensembling approaches: 1) ensembling ViLD-text with CLIP (**ViLD-text + CLIP**); 2) ensembling ViLD-text and ViLD-image using two different heads (**ViLD-ensemble**). As shown in Table 3, ViLD-ensemble improves performance over ViLD, mainly on AP_c and AP_r . This demonstrates using two different heads and ensembling reduce the competition of ViLD-text and ViLD-image. ViLD-text + CLIP obtains significantly higher AP_r , outperforming ViLD by 6.5, and maintains good performance on $AP_{c,f}$. Note the ViLD-text + CLIP model is slow and impractical for real world applications. This experiment is meant for showing the potential of using a zero-shot image classification model for zero-shot detection.

4.5. Transfer to novel datasets

Using CLIP text embeddings as the detection classifier not only achieves impressive zero-shot performance on LVIS, but also makes it finetuning-free to transfer to other datasets. We simply replace the detection classifier with the category text embeddings of a new dataset. For simplicity, we use the same background embedding trained on LVIS. This is all you need to transfer detectors trained with ViLD. We evaluate the transferability of ViLD (ResNet-50 FPN) on PASCAL VOC [9], COCO [28], and Objects365 [39]. Because these three datasets have much smaller vocabularies than LVIS, category overlap is unavoidable and images can be shared across datasets, *e.g.*, COCO and LVIS. Our experiments do not measure the zero-shot detection performance. Instead, we want to evaluate the generalization ability of our model trained on LVIS.

As shown in Table 5, ViLD achieves better generalization ability than ViLD-text. In PASCAL and COCO, the gap is large. This improvement should be credited to visual distillation, which better aligns region embeddings with the free-from text classifier. We also compare with finetuning the linear classifier and training from scratch with super-

	Distill loss	Distill weight w	AP_r	AP_c	AP_f	AP
No distill	0.0	10.4	22.9	31.3	24.0	
	0.5	13.7	21.7	31.2	24.0	
\mathcal{L}_2 loss	1.0	12.4	22.7	31.4	24.3	
	2.0	13.4	22.0	30.9	24.0	
\mathcal{L}_1 loss	0.05	12.9	22.4	31.7	24.4	
	0.1	14.0	20.9	31.2	23.8	
	0.5	16.3	19.2	27.3	21.9	
	1.0	17.3	18.2	25.1	20.7	

Table 4: **Hyperparameter sweep for visual distillation in ViLD.** \mathcal{L}_1 loss is better than \mathcal{L}_2 loss. For \mathcal{L}_1 loss, there is a trend that AP_r increases as the weight increases, while $AP_{f,c}$ decrease. For all parameter combinations, ViLD outperforms ViLD-text on AP_r . We use ResNet-50 FPN backbone and shorter training iterations (84,375), and report mask AP in this table.

vised learning. Although ViLD has 3-6 AP gap compared to the finetuning and larger gaps compared to the supervised learning across datasets, it is the first time that we can directly transfer a trained detector to different datasets.

4.6. Qualitative results

Qualitative examples: In Figure 4, we visualize ViLD’s detection results in various scenes. It illustrates ViLD is able to detect objects of both novel and base categories, with high-quality bounding boxes and masks prediction for novel objects, *e.g.*, it well separates banana slices from the crapes (novel categories). We also show qualitative results on COCO and Objects365 and find ViLD generalizes well. In Figure 6, we show two failure cases of ViLD. The most common failure cases are the missed detection. A less common mistake is misclassifying the object category.

On-the-fly interactive object detection: We tap the potential of ViLD by using free-form text to interactively recognize fine-grained categories and attributes. After obtaining detection results on base categories, we extract the region embedding and compute its cosine similarity with a small set of on-the-fly free-form texts describing attributes and/or fine-grained categories; we apply a softmax with temperature on top of the similarities. To our surprise, though never trained on fine-grained dog breeds (Figure 5), it correctly distinguishes husky from shiba inu. It also works well on identifying object colors (Figure 1). The results demonstrate knowledge distillation from a zero-shot image classification model helps ViLD to gain understanding of concepts not present in the training. Of course, ViLD does not work all the time, *e.g.*, it fails to recognize poses of animals.

Systematic expansion of dataset vocabulary: In addition, we propose to systematically expand the dataset vocabulary ($\mathbf{v} = \{v_1, \dots, v_p\}$) with a set of attributes ($\mathbf{a} = \{a_1, \dots, a_q\}$)

Method	PASCAL VOC [†]		COCO						Objects365					
	AP ₅₀	AP ₇₅	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l	AP	AP ₅₀	AP ₇₅	AP _s	AP _m	AP _l
ViLD-text	40.5	31.6	28.8	43.4	31.4	11.8	35.5	52.0	10.4	15.8	11.1	4.5	11.4	20.4
ViLD	72.2	56.7	36.6	55.6	39.8	20.7	39.2	52.6	11.8	18.2	12.6	5.5	13.5	21.2
Finetuning	78.9	60.3	39.1	59.8	42.4	21.0	41.7	55.0	15.2	23.9	16.2	7.3	17.2	26.1
Supervised	78.5	49.0	46.5	67.6	50.9	27.1	67.6	77.7	25.6	38.6	28.0	16.0	28.1	36.7

Table 5: **Generalization ability of the detector trained with ViLD on LVIS.** We evaluate the trained model on PASCAL VOC 2007 test set, COCO validation set, and Objects365 V1 validation set. Simply replacing the text embeddings, our approach is able to transfer to various detection datasets of different vocabularies. This finetuning-free transfer achieves performance slightly worse than finetuning the classifier. The Supervised baselines of COCO and Objects365 are trained from scratch using the same setting as the LVIS Supervised baseline; [†]: due to the small size of PASCAL VOC, we train its Supervised baseline from an ImageNet pre-trained checkpoint and report the best result of a hyperparameter sweep. All results are box AP for Mask-RCNN ResNet-50 FPN backbone.

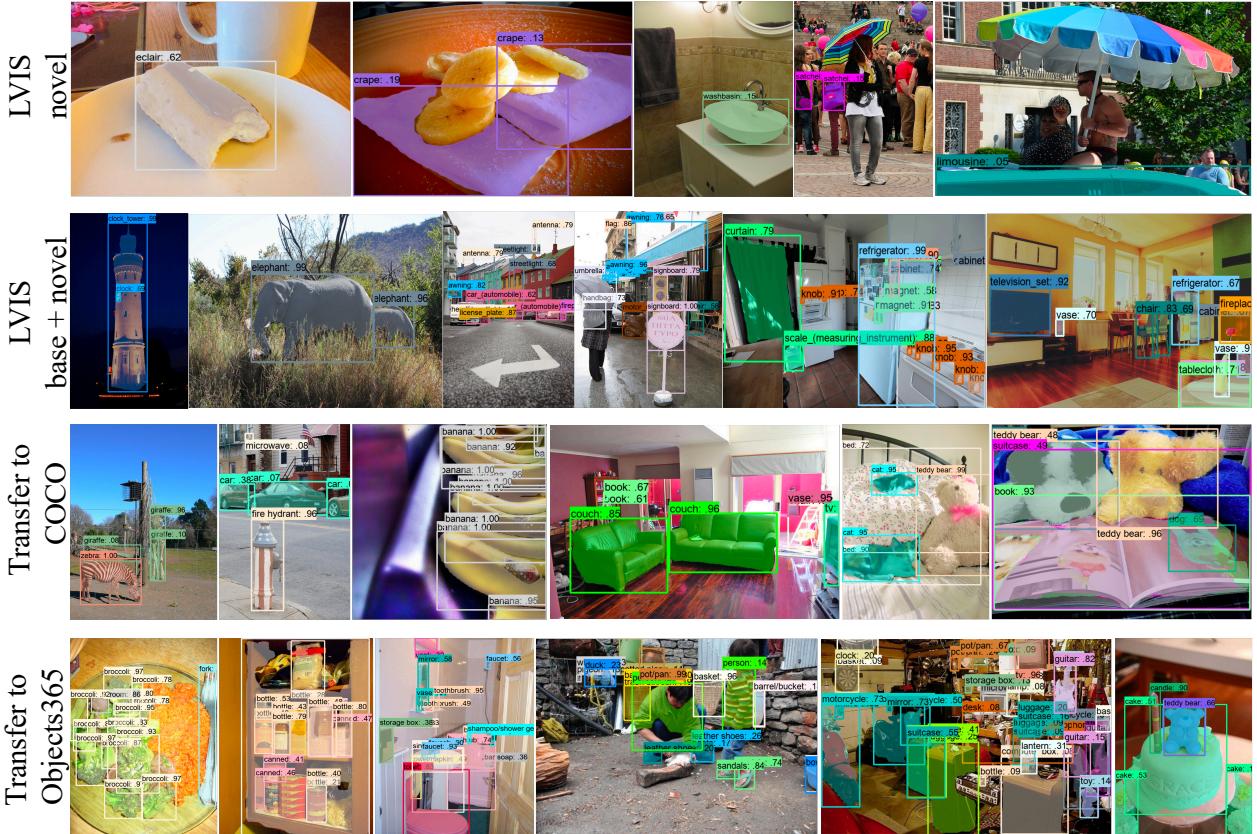


Figure 4: **Qualitative results on LVIS, COCO, and Objects365.** First row: ViLD is able to correctly localize and recognize objects in novel categories. For clarity, we only show the detected novel objects in the scene. Second row: The detected instances on base+novel categories. The performance on base categories are not degraded with ViLD. Last two rows: ViLD can directly transfer to COCO and Objects365 without further fine-tuning.

as follows:

$$\begin{aligned} \Pr(v_i, a_j | \mathbf{e}_r) &= \Pr(v_i | \mathbf{e}_r) \cdot \Pr(a_j | v_i, \mathbf{e}_r) \\ &= \Pr(v_i | \mathbf{e}_r) \cdot \Pr(a_j | \mathbf{e}_r), \end{aligned} \quad (6)$$

where \mathbf{e}_r denotes the region embedding. We assume the event the object belongs to category v_i is conditionally independent to the event it has attribute a_j given the region

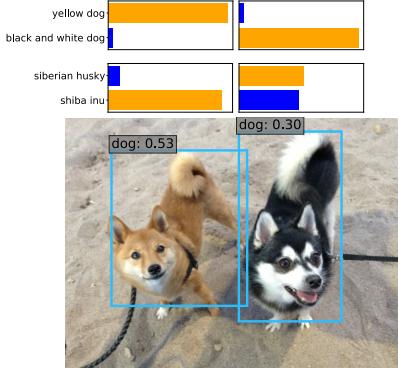
embedding \mathbf{e}_r , i.e., $v_i \perp\!\!\!\perp a_j | \mathbf{e}_r$.

Let τ denote the temperature used for softmax and \mathcal{T} denote the text encoder as in Eq. 2. Then

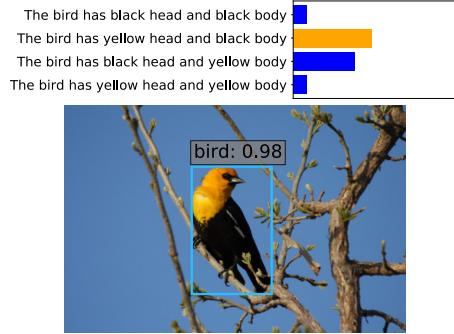
$$\Pr(v_i | \mathbf{e}_r) = \text{softmax}(\text{sim}(\mathbf{e}_r, \mathcal{T}(\mathbf{v}))/\tau), \quad (7)$$

$$\Pr(a_j | \mathbf{e}_r) = \text{softmax}_j(\text{sim}(\mathbf{e}_r, \mathcal{T}(\mathbf{a}))/\tau). \quad (8)$$

In this way, we are able to expand p vocabularies into a



(a) Fine-grained breeds and colors.



(b) Colors of body parts.

Figure 5: **On-the-fly interactive object detection.** One application of ViLD is using on-the-fly free-form text embeddings to further recognize more details of the detected objects, e.g., fine-grained categories and various attributes.

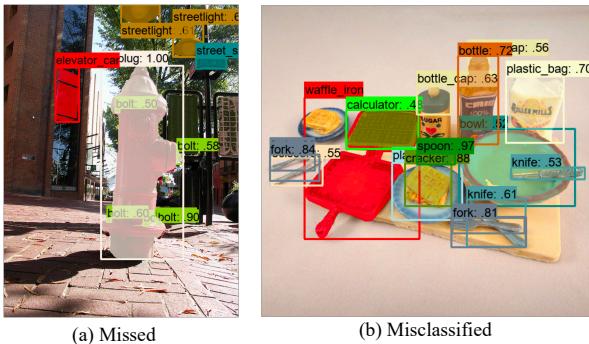


Figure 6: **Failure cases on LVIS novel categories.** The red bounding boxes indicate the groundtruths of failure detections. **(a)** A common failure type where the novel objects are missing, e.g., the elevator car is not detected. **(b)** A less common failure where (part of) the novel objects are misclassified, e.g., half of the waffle iron is detected as a calculator due to visual similarity.

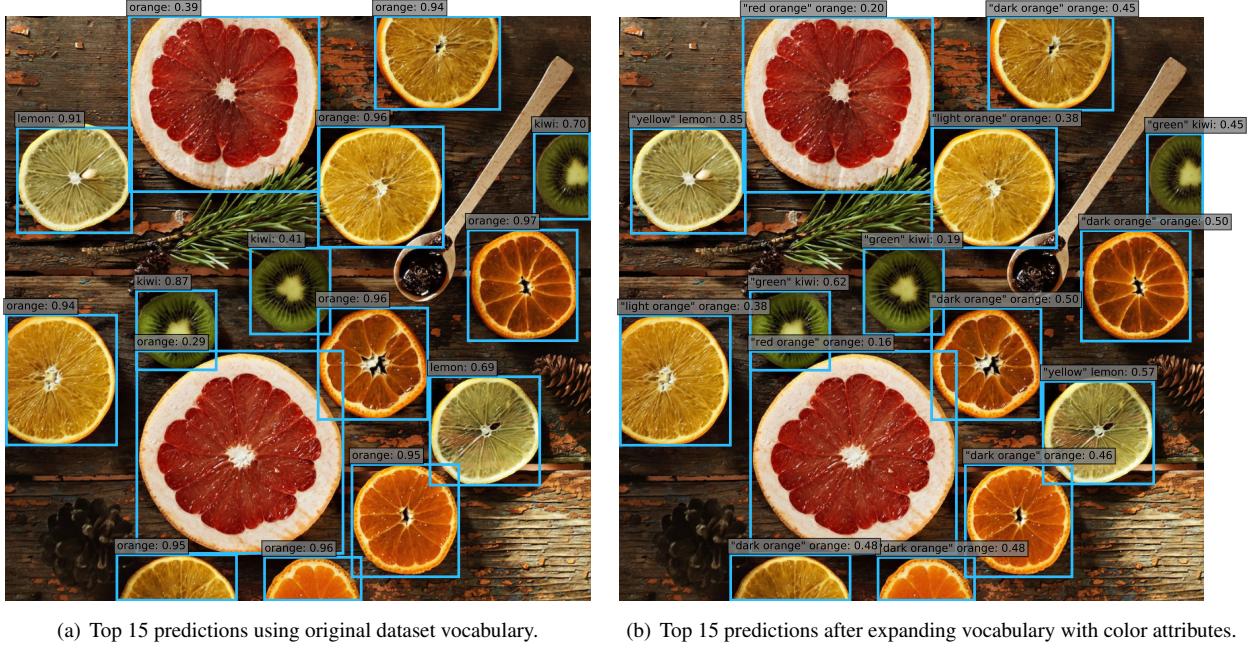
new set of $p \times q$ vocabularies with more attributes. The conditional probability approach is similar to YOLO9000 [36]. We show a qualitative example of this approach in Figure 7, where we use a color attribute set as a . Our zero-shot detector successfully detects fruits with color attributes.

We further expand our vocabularies to fine-grained categories by using all 200 bird species from CUB-200-2011 [41]. Figure 8 shows successful and failure examples of our zero-shot fine-grained detection on CUB-200-2011 images. In general, our model is able to detect visually distinctive species, but fails at other ones.

5. Conclusion

We present ViLD, a zero-shot detection method by distilling knowledge in zero-shot image classification models. ViLD attains 16.1 AP_r on novel categories in LVIS, which surpasses the supervised learning counterpart at the same

inference speed. With model ensembling with CLIP, the performance can be further improved to 22.6, showing the potential of our approach. We demonstrate the detector learned from LVIS can be directly transferred to 3 other detection datasets. Moreover, ViLD can use free-from text to detect fine-grained categories and objects with attributes. Finally, ViLD shows a scalable alternative approach for detecting long-tailed categories instead of collecting expensive annotations for object detection.



(a) Top 15 predictions using original dataset vocabulary.

(b) Top 15 predictions after expanding vocabulary with color attributes.

Figure 7: Systematic expansion of dataset vocabulary with colors. We add 11 color attributes (*red orange*, *dark orange*, *light orange*, *yellow*, *green*, *cyan*, *blue*, *purple*, *black*, *brown*, *white*) to LVIS categories, which expand the vocabulary size by $11\times$. Above we show an example of detection results on an image with fruits. Our zero-shot detector is able to assign the correct color to each object. A class-agnostic NMS with threshold 0.9 is applied.

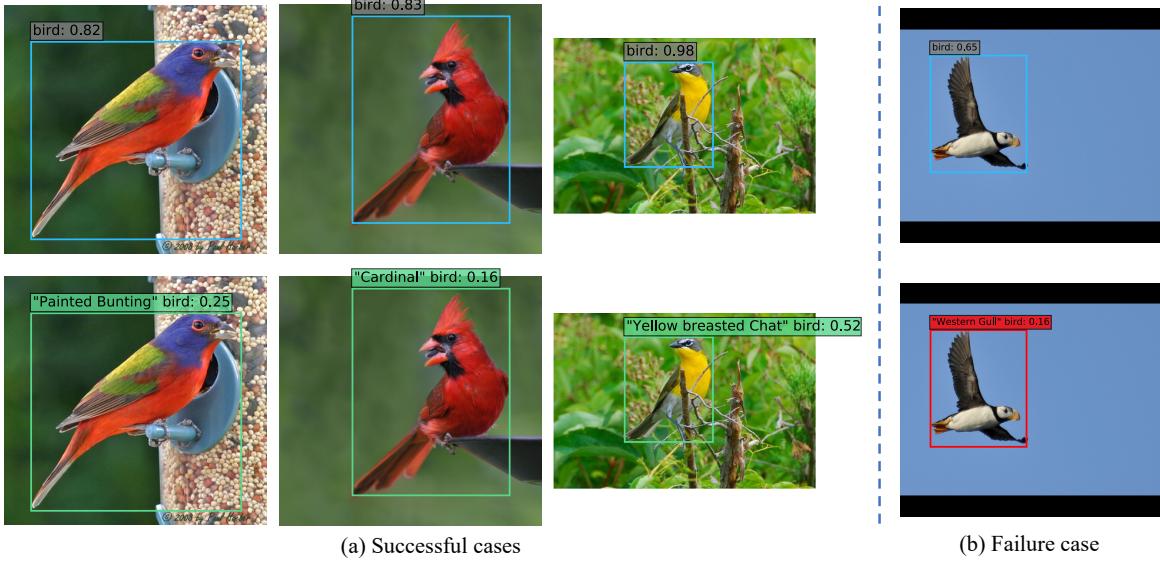


Figure 8: Systematic expansion of dataset vocabulary with fine-grained categories. We use the systematic expansion method to detect 200 fine-grained bird species in CUB-200-2011. **(a)**: Our zero-shot detector is able to perform fine-grained detection **(bottom)** conditioned on the detector trained on LVIS **(top)**. **(b)**: It fails at recognizing visually non-distinctive species. It incorrectly assigns “Western Gull” to “Horned Puffin” due to visual similarity.

References

- [1] Zeynep Akata, Mateusz Malinowski, Mario Fritz, and Bernt Schiele. Multi-cue zero-shot learning with strong supervision. In *CVPR*, 2016. 2
- [2] Zeynep Akata, Florent Perronnin, Zaid Harchaoui, and Cordelia Schmid. Label-embedding for image classification. *TPAMI*, 2015. 2
- [3] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. In *ECCV*, 2018. 2, 3, 5
- [4] Yannick Le Cacheux, Herve Le Borgne, and Michel Crucianu. Modeling inter and intra-class relations in the triplet loss for zero-shot learning. In *ICCV*, 2019. 2
- [5] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *NeurIPS*, 2019. 2
- [6] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019. 2
- [7] Berk Demirel, Ramazan Gokberk Cinbis, and Nazli Ikizler-Cinbis. Zero-shot object detection by hybrid region embedding. In *BMVC*, 2018. 3
- [8] Mohamed Elhoseiny, Yizhe Zhu, Han Zhang, and Ahmed Elgammal. Link the head to the “beak”: Zero shot learning from noisy text description at part precision. In *CVPR*, 2017. 2
- [9] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 7, 13
- [10] Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In *CVPR*, 2009. 2
- [11] Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Marc’Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *NeurIPS*, 2013. 2
- [12] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *CVPR*, 2021. 5
- [13] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 2
- [14] Ross Girshick, Ilya Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018. 5
- [15] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 1, 2, 3, 5, 6, 13
- [16] Nasir Hayat, Munawar Hayat, Shafin Rahman, Salman Khan, Syed Waqas Zamir, and Fahad Shahbaz Khan. Synthesizing the unseen for zero-shot object detection. In *ACCV*, 2020. 3
- [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 1, 2, 3, 5
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [19] Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. Learning deep representation for imbalanced classification. In *CVPR*, 2016. 2
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5
- [21] Ayush Jaiswal, Yue Wu, Pradeep Natarajan, and Premkumar Natarajan. Class-agnostic object detection. In *WACV*, 2021. 3
- [22] Dinesh Jayaraman and Kristen Grauman. Zero shot recognition with unreliable attributes. *arXiv preprint arXiv:1409.4327*, 2014. 2
- [23] Zhong Ji, Yanwei Fu, Jichang Guo, Yanwei Pang, Zhongfei Mark Zhang, et al. Stacked semantics-guided attention model for fine-grained zero-shot learning. In *NeurIPS*, 2018. 2
- [24] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *ICLR*, 2020. 2
- [25] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallochi, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020. 1, 2
- [26] Yu Li, Tao Wang, Bingyi Kang, Sheng Tang, Chunfeng Wang, Jintao Li, and Jiashi Feng. Overcoming classifier imbalance for long-tail object detection with balanced group softmax. In *CVPR*, 2020. 2
- [27] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 5
- [28] 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 *ECCV*, 2014. 7
- [29] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 2
- [30] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018. 2, 6
- [31] Mohammad Norouzi, Tomas Mikolov, Samy Bengio, Yoram Singer, Jonathon Shlens, Andrea Frome, Greg S Corrado, and Jeffrey Dean. Zero-shot learning by convex combination of semantic embeddings. *arXiv preprint arXiv:1312.5650*, 2013. 2
- [32] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *EMNLP*, 2014. 2, 6
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 1, 2, 3, 13

- [34] Shafin Rahman, Salman Khan, and Nick Barnes. Transductive learning for zero-shot object detection. In *ICCV*, 2019. 3
- [35] Shafin Rahman, Salman Khan, and Fatih Porikli. Zero-shot object detection: Learning to simultaneously recognize and localize novel concepts. In *ACCV*, 2018. 2, 3
- [36] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *CVPR*, 2017. 9
- [37] Shaoqing Ren, Kaiming He, Ross B Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 2
- [38] Marcus Rohrbach, Michael Stark, and Bernt Schiele. Evaluating knowledge transfer and zero-shot learning in a large-scale setting. In *CVPR*, 2011. 2
- [39] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *ICCV*, 2019. 7
- [40] Jingru Tan, Changbao Wang, Buyu Li, Quanquan Li, Wanli Ouyang, Changqing Yin, and Junjie Yan. Equalization loss for long-tailed object recognition. In *CVPR*, 2020. 2
- [41] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 9
- [42] Tao Wang, Yu Li, Bingyi Kang, Junnan Li, Junhao Liew, Sheng Tang, Steven Hoi, and Jiashi Feng. The devil is in classification: A simple framework for long-tail object detection and instance segmentation. In *ECCV*, 2020. 2
- [43] Xiaolong Wang, Yufei Ye, and Abhinav Gupta. Zero-shot recognition via semantic embeddings and knowledge graphs. In *CVPR*, 2018. 2
- [44] Guo-Sen Xie, Li Liu, Fan Zhu, Fang Zhao, Zheng Zhang, Yazhou Yao, Jie Qin, and Ling Shao. Region graph embedding network for zero-shot learning. In *ECCV*, 2020. 2
- [45] Ye Zheng, Ruoran Huang, Chuanchi Han, Xi Huang, and Li Cui. Background learnable cascade for zero-shot object detection. In *ACCV*, 2020. 3
- [46] Pengkai Zhu, Hanxiao Wang, and Venkatesh Saligrama. Don't even look once: Synthesizing features for zero-shot detection. In *CVPR*, 2020. 3

A. Additional Qualitative Results

Failure case of mask prediction: See Figure 9. We notice this occasionally happens when the object is large and has clear surface color boundaries. It seems that the class-agnostic mask head sometimes makes predictions based on low-level appearance rather than semantics.



Figure 9: **An example of ViLD on PASCAL VOC showing a mask of poor quality.** The class-agnostic mask prediction head occasionally predicts masks based on low-level appearance rather than semantics, and thus fail to obtain a complete instance mask.

Transfer to PASCAL VOC: In Figure 10, we show qualitative results of transferring a zero-shot detector trained on LVIS [15] to PASCAL VOC Detection (2007 test set) [9], without finetuning (Section 4.5 in the main paper). Results demonstrate that the transferring works well.

B. Additional Experiment Details

Text prompts: Since the zero-shot classification model we use is trained on full sentences, we use the text embeddings of “This is a photo of a {category}” and ensemble various prompts, *e.g.*, “This is a photo of a small / medium / large {category}”, following [33]. In addition, we also include prompts containing the phrase “in the scene”: “There is a / the {category} in the scene” and “This is a / the / one {category} in the scene”.

Model used for qualitative results: For all qualitative results, we use a Mask R-CNN R152-FPN model (ViLD $w = 1.0$ in Table 3(b)).

Class-agnostic supervised baselines: For a fair comparison, we train the second stage box/mask prediction heads of Supervised and Supervised-RFS baselines in a class-agnostic manner (Section 3.1), the same manner as ours.

Details for R-CNN style experiments: For CLIP on cropped regions in Section 4.2 and ViLD-text + CLIP in Section 4.3, when obtaining the region proposals from Mask R-CNN, we apply a class-agnostic NMS with 0.9 threshold, and output a maximum of 1000 proposals. After

getting results from CLIP or ensembling, we apply a class-specific NMS with a threshold of 0.6, and output a maximum of 300 detection results. We use 300 as the maximum number of detections for all experiments.

C. Analysis of CLIP on Cropped Regions

In this section, we analyze some common failure cases of CLIP on cropped regions and discuss possible ways to mitigate these problems.

Visual similarity: This confusion is common for any classifiers and detectors, especially on large vocabularies. In Figure 11(a), we show two failure examples due to visual similarity. Since we only use a relatively small ViT-B/32 CLIP model, potentially we can improve the performance by using a higher-capacity pre-trained model.

Aspect ratio: This issue is introduced by the pre-processing of inputs in CLIP. We use the ViT-B/32 CLIP with a fixed input resolution of 224×224 . It resizes the shorter edge of the image to 224, and then uses a center crop. However, since region proposals can have more extreme aspect ratios than the training images for CLIP, and some proposals are tiny, we directly resize the proposals to that resolution, which might cause some issues. For example, the thin structure in Figure 11(b) right will be highly distorted with the pre-processing. And the oven and fridge can be confusing with the distorted aspect ratio. There might be some simple remedies for this, *e.g.*, pasting the cropped region with original aspect ratio on a black background.

Multiple objects in a bounding box: Multiple objects in a region interfere CLIP’s classification results, see Figure 11(c), where a corner of an aquarium dominates the prediction. This is due to CLIP’s pre-training method, which pairs an entire image with its caption. The caption is usually about salient objects in the image. It’s hard to mitigate this issue at the zero-shot classification model’s end. On the other hand, a supervised detector are trained to recognize the major object in the bounding box. So when distilling knowledge from a zero-shot image classification model, keeping training a supervised detector on base categories could help, as can be seen from the qualitative results of ViLD (Figure 4 in the main paper).

Confidence scores predicted by CLIP do not reflect the localization quality: For example, in Figure 12(a), CLIP correctly classifies the object, but gives highest scores to partial detection boxes. CLIP is not trained to measure the quality of bounding boxes. Nonetheless, in object detection, it is important for the higher-quality boxes of the same object to have higher scores. In Figure 12(c), we simply re-score by taking the geometric mean of CLIP probabilities and the Mask R-CNN score of the bounding box, which



Figure 10: **Transfer to PASCAL VOC.** ViLD correctly localizes and recognizes objects when transferring to PASCAL VOC, where images usually have lower resolution than LVIS (our training set). In the third picture, our detector is able to locate tiny bottles, though it fails to detect the person.

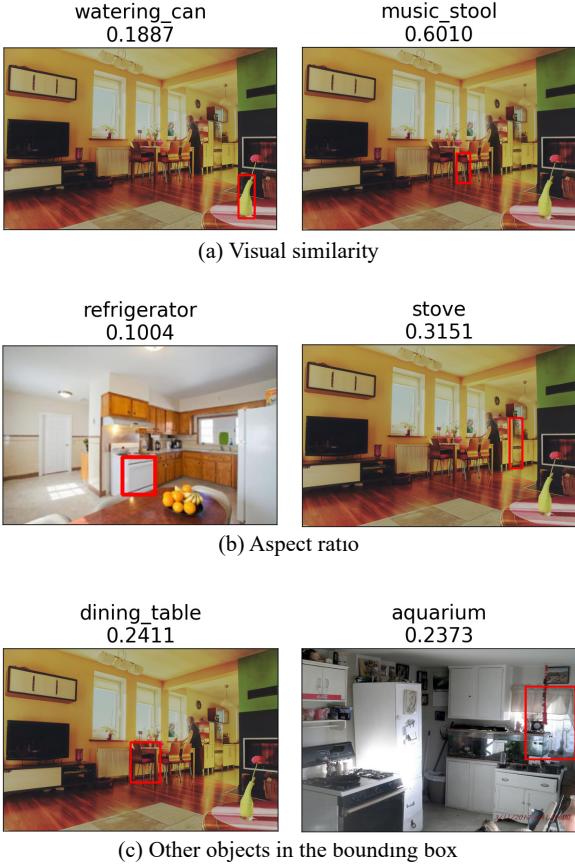


Figure 11: **Typical errors of CLIP on cropped regions.** (a): The prediction and the groundtruth have high visual similarity. (b): Directly resizing the cropped regions changes the aspect ratios, which may cause troubles. (c): CLIP’s predictions are sometimes affected by other objects appearing in the region, rather than predicting what the entire bounding box is.

yields much better top predictions. In Figure 12(b), we show top predictions of the Mask R-CNN model used in (c). Its top predictions have good bounding boxes, while the

predicted categories are wrong. This experiment shows that it’s important to have both a zero-shot classification model for better recognition, as well as a supervised detector participating in scoring for better localization.

For the four issues we discussed in this section, ViLD could also have the first two as we distill visual embeddings into the detector. The last two issues, however, should be mitigated in ViLD, since it keeps the supervised detection training on base categories and uses a detection model after knowledge distillation.

D. Box APs

Table 6 shows the corresponding box AP of Table 3 in the main paper. In general, box AP is slightly higher than mask AP, and the trend/relative performance is consistent between the two. ViLD-ensemble achieves the best box and mask AP_r.

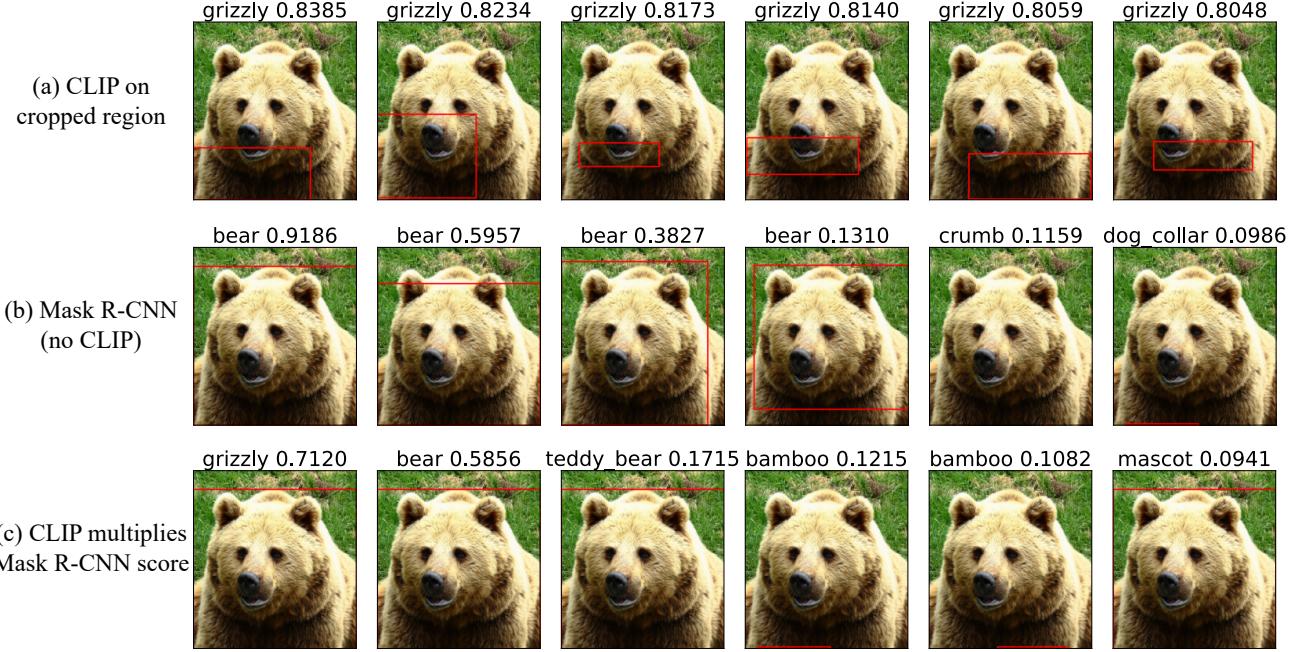


Figure 12: **The prediction scores of CLIP do not reflect the quality of bounding box localization.** (a): Top predictions of CLIP on cropped region. Boxes of poor qualities receive high scores, though the classification is correct. (b): Top predictions of a vanilla Mask R-CNN model. Box qualities are good while the classification is wrong. (c): We take the geometric mean of CLIP classification score and Mask R-CNN score, and use it to rescore (a). In this way, a high-quality box as well as the correct category rank first.

Model	Method	Box				Mask			
		AP _r	AP _c	AP _f	AP	AP _r	AP _c	AP _f	AP
Mask R-CNN (R50-FPN)	CLIP on cropped regions	12.3	9.0	5.0	8.0	13.0	10.6	6.0	9.2
	GloVe baseline	3.2	22.0	34.9	23.8	3.0	20.1	30.4	21.2
	ViLD-text	10.6	26.1	37.4	27.9	10.1	23.9	32.5	24.9
	ViLD-image	8.4	8.2	7.2	7.8	9.6	8.5	7.8	8.4
	ViLD ($w = 0.5$)	16.3	21.2	31.6	24.4	16.1	20.0	28.3	22.5
	ViLD-ensemble ($w = 0.5$)	16.7	26.5	34.2	27.8	16.6	24.6	30.3	25.5
Mask R-CNN (R152-FPN)	ViLD-text + CLIP [†]	23.8	26.7	32.8	28.6	22.6	24.8	29.2	26.1
	Supervised-RFS (base + novel)	13.0	26.7	37.4	28.5	12.3	24.3	32.4	25.4
	ViLD-text	12.3	28.3	39.7	30.0	11.7	25.8	34.4	26.7
	ViLD-image	9.7	9.5	8.2	9.0	10.8	10.0	8.7	9.6
	ViLD ($w = 1.0$)	19.1	22.4	31.5	25.4	18.7	21.1	28.4	23.6
	ViLD-ensemble ($w = 2.0$)	19.8	27.1	34.5	28.7	18.7	24.9	30.6	26.0
	Supervised-RFS (base + novel)	16.2	29.6	39.7	31.2	14.4	26.8	34.2	27.6

Table 6: **Performance of ViLD variants.** This table shows additional box AP for models in Table 3.