# Aligning Bag of Regions for Open-Vocabulary Object Detection

Size Wu[1]    Wenwei Zhang[1]    Sheng Jin[2,3]    Wentao Liu[3,4]    Chen Change Loy[1*]

[1]Nanyang Technological University    [2] The University of Hong Kong
[3] SenseTime Research and Tetras.AI    [4] Shanghai AI Laboratory

{size001, wenwei001, ccloy}@ntu.edu.sg    {jinsheng, liuwentao}@sensetime.com

## Abstract

*Pre-trained vision-language models (VLMs) learn to align vision and language representations on large-scale datasets, where each image-text pair usually contains a bag of semantic concepts. However, existing open-vocabulary object detectors only align region embeddings individually with the corresponding features extracted from the VLMs. Such a design leaves the compositional structure of semantic concepts in a scene under-exploited, although the structure may be implicitly learned by the VLMs. In this work, we propose to align the embedding of* bag of regions *beyond individual regions. The proposed method groups contextually interrelated regions as a bag. The embeddings of regions in a bag are treated as embeddings of words in a sentence, and they are sent to the text encoder of a VLM to obtain the bag-of-regions embedding, which is learned to be aligned to the corresponding features extracted by a frozen VLM. Applied to the commonly used Faster R-CNN, our approach surpasses the previous best results by 4.6 box $AP_{50}$ and 2.8 mask AP on novel categories of open-vocabulary COCO and LVIS benchmarks, respectively. Code and models are available at* https://github.com/wusize/ovdet.

## 1. Introduction

A traditional object detector can only recognize categories learned in the training phase, restricting its application in the real world with a nearly unbounded concept pool. Open-vocabulary object detection (OVD), a task to detect objects whose categories are absent in training, has drawn increasing research attention in recent years.

A typical solution to OVD, known as the distillation-based approach, is to distill the knowledge of rich and unseen categories from pre-trained vision-language models (VLMs) [23, 40]. In particular, VLMs learn aligned image and text representations on large-scale image-text pairs (Fig. 1(a)). Such general knowledge is beneficial for OVD. To extract the knowledge, most distillation-based ap-
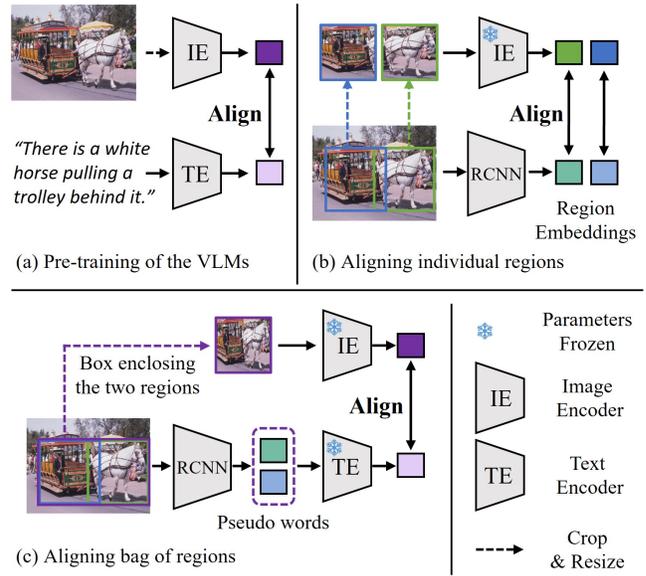


Figure 1. **(a)** Typical vision-language models (VLMs) learn to align representations of images and captions with rich compositional structure. **(b)** Existing distillation-based object detectors align each *individual* region embedding to features extracted by the frozen image encoder of VLMs. **(c)** Instead, the proposed method aligns the embedding of *bag of regions*. The region embeddings in a bag are projected to the word embedding space (dubbed as pseudo words), formed as a sentence, and then sent to the text encoder to obtain the *bag-of-regions* embedding, which is aligned to the corresponding image feature extracted by the frozen VLMs.

proaches [11, 16, 55] align each *individual* region embedding to the corresponding features extracted from the VLM (Fig. 1(b)) with some carefully designed strategies.

We believe VLMs have implicitly learned the inherent compositional structure of multiple semantic concepts (*e.g.*, co-existence of stuff and things [1, 25]) from a colossal amount of image-text pairs. A recent study, MaskCLIP [60], leverages such a notion for zero-shot segmentation. Existing distillation-based OVD approaches, however, have yet to fully exploit the compositional structures encapsulated in VLMs. Beyond the distillation of individual region embed-

1

ding, we propose to align the embedding of *BAg of RegiONs*, dubbed as BARON. Explicitly learning the co-existence of visual concepts encourages the model to understand the scene beyond just recognizing isolated individual objects.

BARON is easy to implement. As shown in Fig. 1(c), BARON first samples contextually interrelated regions to form a 'bag'. Since the region proposal network (RPN) is proven to cover potential novel objects [16, 61], we explore a neighborhood sampling strategy that samples boxes around region proposals to help model the co-occurrence of a bag of visual concepts. Second, BARON obtains the bag-of-regions embeddings by projecting the regional features into the word embedding space and encoding these pseudo words with the text encoder (TE) of a frozen VLM [40]. By projecting region features to pseudo words, BARON naturally allows TE to effectively represent the co-occurring semantic concepts and understand the whole scene. To retain the spatial information of the region boxes, BARON projects the box shape and box center position into embeddings and add to the pseudo words before feeding them to TE.

To train BARON, the bag-of-regions embeddings are learned to be aligned to the embeddings obtained by feeding the image crops that enclose the bag of regions to the teacher, *i.e.*, the image encoder (IE) of the VLM. We adopt a contrastive learning approach [49] to learn the pseudo words and the bag-of-regions embeddings. Consistent with the VLMs' pre-training (*e.g.*, CLIP [40]), the contrastive loss pulls close corresponding student (the detector) and teacher (IE) embedding pairs and pushes away non-corresponding pairs.

We conduct extensive experiments on two challenging benchmarks, OV-COCO and OV-LVIS. The proposed method consistently outperforms existing state-of-the-art methods [11, 55, 61] in different settings. Combined with Faster R-CNN, BARON achieves a 34.0 (4.6 increase) box $AP_{50}$ of novel categories on OV-COCO and 22.6 (2.8 increase) mask mAP of novel categories on OV-LVIS. It is noteworthy that BARON can also distill knowledge from caption supervision – it achieves 32.7 box $AP_{50}$ of novel categories on OV-COCO, outperforming previous approaches that use COCO caption [14, 56, 59, 61].

## 2. Related Work

**Vison-Language Pretraining and Its Applications.** Vision-language pre-training aims to learn aligned image and text representations [13, 22–24, 40] on large-scale image-text pairs. There are many studies [24, 27, 28, 35] that pre-train vision-language models (VLMs) to improve the performance of downstream recognition and generation tasks. There are also studies that learn aligned vision-language representation so that the images can be classified with arbitrary texts [13, 37]. Recent attempts [23, 40, 57] push forward this direction by conducting contrastive

learning in VLMs on billion-scale image-text pairs. These models show impressive zero-shot performance when they are transferred to image classification tasks.

Inspired by the success of VLMs [40], some works try to exploit the alignment of vision-language representations for dense prediction tasks, *e.g.* segmentation [26, 42, 60] and detection [11, 16, 29, 46]. In particular, MaskCLIP [60] shows that the image encoder in VLMs [40] captures the stuff and things in a complex scene, where the pixel embeddings of each concept are naturally aligned with the corresponding text representations, although the original CLIP [40] model does not explicitly learn this target. This implies that VLMs, after trained on a massive amount of image-text pairs, have implicitly learned the compositional structure of multiple semantic concepts, which naturally exist in image-text pairs. This motivates us to explore the representation alignment between bag of regions and bag of words, different from previous works [16, 26] that focus on aligning the representation of individual pixels, regions, or words in VLMs.

**Open-Vocabulary Object Detection.** Traditional object detectors [3, 5, 34, 44, 62] are limited to pre-defined object categories. To detect objects of unseen categories, zero-shot object detection (ZSD) [2, 9, 19, 41, 58] is proposed to align individual region embeddings with the text embeddings of categories through different strategies. Recent attempts further explore open-vocabulary object detection (OVD) [56], a more general form of ZSD that leverages weak supervisions like visual grounding data [29], image captions [14, 56, 59], and image labels [61]. Large-scale pre-trained VLMs [40] are also exploited for their remarkable zero-shot recognition ability. These VLMs can generate conditional queries [55] or serve as a good teacher for knowledge distillation [16]. Specifically, distillation-based approaches [11, 16] extract embeddings on pre-computed region proposals and *individually* align them to the corresponding features obtained from the VLMs. To our best knowledge, BARON is the first attempt to lift the learning from *individual* regions to the *bag of regions* for OVD.

## 3. Method

Our method makes the first attempt to align embedding of *bag of regions* beyond *individual* regions for OVD. We call our method BARON. In this work, we instantiate the idea of BARON based on the commonly used Faster-RCNN [44] and modify it for OVD (Sec. 3.1). We design a simple strategy to form the bag of regions from the region proposals (Sec. 3.2). The bag of regions is treated as a bag of words to obtain the bag-of-regions embedding (Sec. 3.3), which are then aligned with the corresponding features from VLMs (Sec. 3.4). BARON is general that it can align the bag-of-regions embeddings to not only image representations but also text representations (Sec. 3.5).
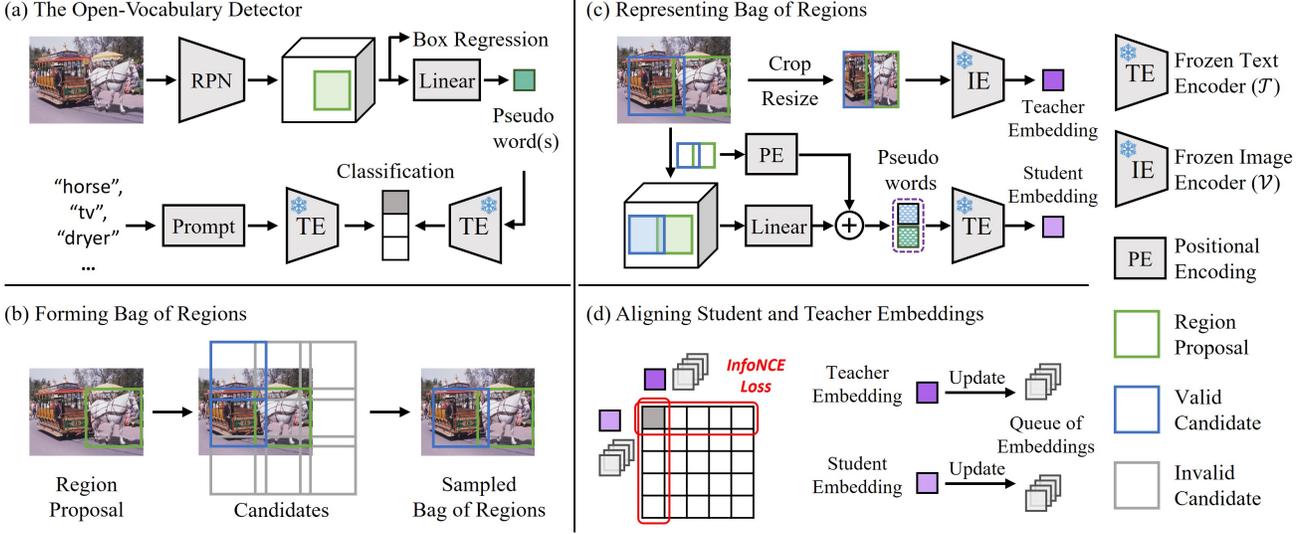
Figure 2. Overview of BARON. **(a)** BARON is based on a Faster R-CNN whose classifier is replaced by a linear layer to map region features into pseudo words. **(b)** BARON takes region proposals and its surrounding boxes to form bags of regions. **(c)** BARON obtains student and teacher embeddings for the bag of regions from the pre-trained VLMs. **(d)** BARON learns the alignment by the InfoNCE loss and maintains queues of embeddings to provide sufficient negative pairs for InfoNCE loss.

## 3.1. Preliminaries

In this paper, we instantiate the idea upon Faster R-CNN [44] for simplicity. The idea can also be used for other architectures [3, 32] applicable for OVD. To enable Faster R-CNN to detect objects from arbitrary vocabularies, we replace the original classifier with a linear layer that projects the region features into the word embedding space (dubbed as pseudo words) (Fig. 2(a)). In practice, the linear layer maps a region feature to multiple pseudo words to represent the rich semantic information of each object, similar to those category names consisting of multiple words (*e.g.*, horse-driven trolley). Finally, we feed these pseudo words into the text encoder and then calculate the similarity with the category embeddings to obtain final classification results.

As shown in Fig. 2(a), given $C$ object categories, we obtain the embedding $f_c$ for the $c$-th category by feeding the category names with a prompt template, *e.g.*, 'a photo of {category} in the scene' to the text encoder $\mathcal{T}$. For a region and its pseudo words $w$, the probability of the region to be classified as the $c$-th category is

$$p_c = \frac{\exp(\tau \cdot \langle \mathcal{T}(w), f_c \rangle)}{\sum_{i=0}^{C-1} \exp(\tau \cdot \langle \mathcal{T}(w), f_i \rangle)}, \quad (1)$$

where $\langle, \rangle$ denotes the cosine similarity and $\tau$ is the temperature to re-scale the value.

During training, only the boxes of base categories are annotated and the learning on base categories follows the convention of Faster R-CNN with regression and classification losses [44]. To learn to detect novel categories that do not

have box annotations in training, previous distillation-based approaches [11, 16] *individually* align region embeddings (*e.g.*, $\mathcal{T}(w)$) to the corresponding features obtained from the VLMs. To further exploit the power of VLMs that capture the compositional structure of multiple concepts, we lift the learning from *individual* regions to the *bag of regions*.

## 3.2. Forming Bag of Regions

Our framework is inspired by existing OVD approaches that distill knowledge from the image encoder of VLMs. Specifically, we choose the image encoder of the VLMs as the teacher and expect it to teach the detector. But different from existing approaches, we wish the detector to learn the co-existence of multiple concepts, especially the potential existence of novel objects. To effectively and efficiently learn such knowledge from the VLMs, we consider the following two properties of regions inside a bag: 1) the regions need to be close to each other because an image crop enclosing distanced regions will include a larger proportion of redundant image contents, distracting the image encoder from representing the bag of regions; 2) the regions should have similar sizes, as an imbalanced size ratio among regions will make the image representations dominated by the largest region.

According to these two requirements, we adopt a simple neighborhood sampling strategy to form the bag of regions based on the region proposals predicted by the region proposal network (RPN) [44]. Concretely, for each region proposal, we take its surrounding eight boxes (neighbors) as the candidates, which are spatially close to each other as shown in Fig. 2(b). We also allow these candidates to
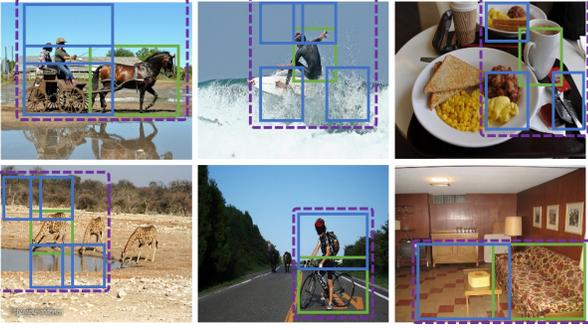
Figure 3. Visualization of sampled bags of regions. Green boxes denote the region proposals and blue boxes are sampled neighbors (candidates). Areas exceeding image boundary are cropped out.

overlap slightly at a specific Intersection over Foreground (IOF) to improve the continuity of regional representations. To balance the sizes of regions inside the bag, we simply let the eight candidate boxes have the same shape as the region proposal. In practice, candidates that exceed the boundary of the image over a certain proportion, *e.g.*, more than $\frac{2}{3}$ of their area out of the image, will be discarded. The remaining candidates are independently sampled, and the sampled boxes, together with the region proposal, form a bag of regions.

We sample $G$ groups for each region proposal to obtain rich bag-of-regions representations. At each time of sampling, the probability of a candidate being sampled is adjusted to prevent the box enclosing the grouped regions from having an extreme aspect ratio. Assuming a base probability $p_b$ and a scaling factor $\alpha$, the probability to sample the left and right candidates is $p = p_b \times min\{(\frac{H}{W})^\alpha, 1\}$ and the probability to sample upper and bottom candidates is $p = p_b \times min\{(\frac{W}{H})^\alpha, 1\}$, where $H$ and $W$ are the height and width of the region proposal box. $\alpha$ is 3.0 by default.

We show some sampling results of bags of regions on COCO dataset [33] in Fig 3. The grouped regions in a bag can cover objects (in blue) co-occurring with the object in the region proposal (in green). The context of the co-occurring objects leads BARON to get a scene-level understanding of the regions, *e.g.* 'a horse pulling the carriage' and 'a man surfing on the wave'. Novel object categories occasionally appear in the bag of regions, *e.g.* the 'carriage' in the top left image of Fig 3 and the 'cup' in the top right. In the following Sec 3.3 and Sec 3.4, these potential novel object categories would be learned in the context of the bag of regions.

### 3.3. Representing Bag of Regions

With the sampled bags of regions, BARON obtains the bag-of-regions embeddings from both the student (*i.e.*, the open-vocabulary object detector) and the teacher (*i.e.*, VLMs). We denote the $j$-th region in the $i$-th group as $b_j^i$ and the pseudo words after the projection layer as $w_j^i$. For the pre-trained VLM, we use $\mathcal{T}$ to denote the text encoder and $\mathcal{V}$ to denote the image encoder.

**Student Bag-of-Regions Embedding.** Because the region features are projected to word embedding space and learned to be aligned with the text embedding of category, a straightforward way to obtain the embedding of a bag of regions is to concatenate the pseudo words and feed them to the text encoder $\mathcal{T}$. However, the spatial information of the regions will be lost in such a process, including relative box center positions and relative box shapes. The center position and shape indicate the spatial relationship among the regions in a bag, which is essential to induce a sentence-like interpretation for the bag of regions. And they are also encoded in the teacher (image encoder of VLMs) through positional embeddings in the input [40]. Therefore, BARON encodes the spatial information into positional embeddings $p_j^i$ that have the same dimension with $w_j^i$, following the practices of Transformers [10, 50]. The positional embeddings are added to the pseudo words before concatenation. Assuming the group size is $N^i$, this representation can be formulated as $f_t^i = \mathcal{T}(w_0^i + p_0^i, w_1^i + p_1^i, \ldots, w_{N^i-1}^i + p_{N^i-1}^i)$.

**Teacher Bag-of-Regions Embedding.** The image embedding of the grouped regions can be obtained by feeding the image crop that encloses the regions to the image encoder $\mathcal{V}$. The image crop may contain redundant contents that are outside the grouped regions; we mask out them in the attention layers of $\mathcal{V}$. The image feature can be formulated as $f_v^i = \mathcal{V}(b_0^i, b_1^i, \ldots, b_{N_i-1}^i)$.

### 3.4. Aligning Bag of Regions

BARON aligns the bag-of-regions embeddings from student and teacher to make the student learn to encode the co-existence of multiple regions, which potentially contain multiple concepts. We adopt the contrastive learning approach used in vision-language pre-training [40]. Specifically, given $G$ bags of regions, the alignment InfoNCE loss [38] between bag-of-regions embeddings is calculated as

$$\mathcal{L}_{\text{bag}} = -\frac{1}{2}\sum_{k=0}^{G-1}(\log(p_{t,v}^k) + \log(p_{v,t}^k)). \quad (2)$$

The $p_{t,v}^k$ and $p_{v,t}^k$ are calculated as

$$p_{t,v}^k = \frac{\exp(\tau' \cdot \langle f_t^k, f_v^k \rangle)}{\sum_{l=0}^{G-1}\exp(\tau' \cdot \langle f_t^k, f_v^l \rangle)} \quad (3)$$

$$p_{v,t}^k = \frac{\exp(\tau' \cdot \langle f_v^k, f_t^k \rangle)}{\sum_{l=0}^{G-1}\exp(\tau' \cdot \langle f_v^k, f_t^l \rangle)}, \quad (4)$$

respectively, where $\tau'$ is the temperature to re-scale the cosine similarity. The loss pulls positive pairs $\{f_t^k, f_v^k\}$ close to each other and pushes away negative pairs $\{f_t^k, f_v^l\}(k \neq l)$.

In practice, the number of groups $G$ is small for a single image. We maintain two queues to save the image and text embeddings in previous iterations during training to provide sufficient negative pairs [20].

Table 1. Comparison with state-of-the-art methods on OV-COCO benchmark. We separately compare our approach with methods distilling knowledge from CLIP and approaches using COCO caption. † means using proposals produced by MAVL [36].

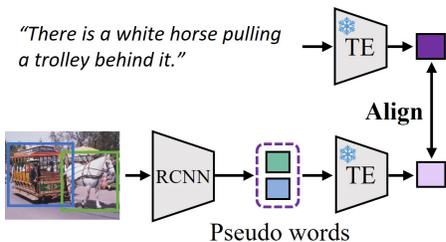| Method | Supervision | Backbone | Detector | $AP_{50}^{novel}$ | $AP_{50}^{base}$ | $AP_{50}$ |
|---|---|---|---|---|---|---|
| ViLD [16] | CLIP | ResNet50-FPN | FasterRCNN | 27.6 | 59.5 | 51.2 |
| OV-DETR [55] | CLIP | ResNet50 | DeformableDETR | 29.4 | 61.0 | 52.7 |
| BARON (Ours) | CLIP | ResNet50-FPN | FasterRCNN | **34.0** | 60.4 | 53.5 |
| OVR-CNN [56] | Caption | ResNet50-C4 | FasterRCNN | 22.8 | 46.0 | 39.9 |
| RegionCLIP [59] | Caption | ResNet50-C4 | FasterRCNN | 26.8 | 54.8 | 47.5 |
| Detic [61] | Caption | ResNet50-C4 | FasterRCNN | 27.8 | 51.1 | 45.0 |
| PB-OVD [14] | Caption | ResNet50-C4 | FasterRCNN | 30.8 | 46.1 | 42.1 |
| VLDet [30] | Caption | ResNet50-C4 | FasterRCNN | 32.0 | 50.6 | 45.8 |
| BARON (Ours) | Caption | ResNet50-C4 | FasterRCNN | **33.1** | 54.8 | 49.1 |
| Rasheed et al. [43]† | CLIP + Caption | ResNet50-C4 | FasterRCNN | 36.6 | 54.0 | 49.4 |
| BARON (Ours)† | CLIP + Caption | ResNet50-C4 | FasterRCNN | **42.7** | 54.9 | 51.7 |



Figure 4. The caption-version BARON. We align the text embeddings of the bag of regions to the caption embeddings.

**Aligning Individual Regions.** The alignment between individual regions' student and teacher embeddings is complementary to that of a bag of regions. Therefore, we also adopt the individual-level distillation in our implementation. For computational efficiency, we obtain the teacher embeddings from the feature map of the image encoder's last attention layer by RoiAlign [21] instead of repeatedly passing image crops to the image encoder. Similarly, the student embeddings are extracted from the text encoder's last attention layer by averaging pseudo-word embeddings of the same region. We apply the InfoNCE loss and keep queues of embeddings to calculate the individual-level loss $\mathcal{L}_{\text{individual}}$.

### 3.5. Caption Supervision

It is noteworthy that BARON can be applied to caption supervision. As shown in Fig. 4, the core idea is to replace the image embedding $f_v$ with embedding obtained by feeding the image captions to the text encoder $\mathcal{T}$. To obtain the region groups, we randomly sample some region proposals generated by the RPN. As an image can have multiple captions, we follow the practice in UniCL [53] to apply a soft cross entropy loss that simultaneously aligns the student bag-of-regions embedding to multiple caption embeddings. The alignment of individual regions is discarded since we cannot get the correspondence between pseudo words of regions and the actual words in a caption without grounding. In this way, BARON now learn to align the bag-of-regions

embedding of the image caption, which also describes the existence of multiple concepts of the image.

## 4. Experiments

**Datasets.** We evaluate our method on the two popular object detection datasets, i.e., COCO [33] and LVIS [17]. For the COCO dataset, we follow OV-RCNN [56] to split the object categories to 48 base categories and 17 novel categories. For the LVIS dataset, we follow ViLD [16] to split the 337 rare categories into novel categories and the rest common and frequent categories into base categories. For brevity, we denote the open-vocabulary benchmarks based on COCO and LVIS as OV-COCO and OV-LVIS.

**Evaluation Metrics.** We evaluate the detection performance on both base and novel categories for completeness. For OV-COCO, we follow OV-RCNN [56] to report the box AP at IoU threshold 0.5, noted as $AP_{50}$. For OV-LVIS, we report both the mask and box AP averaged on IoUs from 0.5 to 0.95, noted as mAP. The $AP_{50}$ of novel categories ($AP_{50}^{novel}$) and mAP of rare categories ($AP_r$) are the main metrics that evaluate the open-vocabulary detection performance on OV-COCO and OV-LVIS, respectively.

**Implementation Details.** We build BARON on Faster R-CNN [44] with ResNet50-FPN [31]. For a fair comparison with existing methods, we initialize the backbone network with weights pre-trained by SOCO [51] and apply synchronized Batch Normalization (SyncBN) [39] following DetPro [11]. We choose the $2\times$ schedule (180, 000 iterations) for the main experiments on COCO and LVIS [6,52]. For the pre-trained VLM, we choose the CLIP [40] model based on ViT-B/32 [10]. For the adaptation to caption supervision, we base our method on Faster R-CNN with ResNet50-C4 backbone [44] and adopt the $1\times$ schedule. For the prompt of category names, we use the hand-crafted prompts in ViLD [16] for all our experiments by default. We use learned prompt only when comparing with DetPro [11].

Table 2. Comparison with state-of-the-art methods on OV-LVIS. * denotes the re-implemented ViLD [16] reported in DetPro [11].

| Method | Ensemble | Learned Prompt | Object Detection | | | | Instance segmentation | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | $AP_r$ | $AP_c$ | $AP_f$ | AP | $AP_r$ | $AP_c$ | $AP_f$ | AP |
| ViLD [16] | - | - | 16.3 | 21.2 | 31.6 | 24.4 | 16.1 | 20.0 | 28.3 | 22.5 |
| OV-DETR [55] | - | - | - | - | - | - | 17.4 | 25.0 | 32.5 | 26.6 |
| BARON (Ours) | - | - | **17.3** | 25.6 | 31.0 | 26.3 | **18.0** | 24.4 | 28.9 | 25.1 |
| ViLD [16] | ✓ | - | 16.7 | 26.5 | 34.2 | 27.8 | 16.6 | 24.6 | 30.3 | 25.5 |
| ViLD* [16] | ✓ | - | 17.4 | 27.5 | 31.9 | 27.5 | 16.8 | 25.6 | 28.5 | 25.2 |
| BARON (Ours) | ✓ | - | **20.1** | 28.4 | 32.2 | 28.4 | **19.2** | 26.8 | 29.4 | 26.5 |
| DetPro [11] | ✓ | ✓ | 20.8 | 27.8 | 32.4 | 28.4 | 19.8 | 25.6 | 28.9 | 25.9 |
| BARON (Ours) | ✓ | ✓ | **23.2** | 29.3 | 32.5 | 29.5 | **22.6** | 27.6 | 29.8 | 27.6 |

Table 3. Comparison of the transfer ability of the model trained on OV-LVIS. * denotes the re-implemented ViLD [16] reported in DetPro [11]. ‡ denotes that we use hand-crafted prompts for a fair comparison with ViLD.

| Method | Pascal VOC | | COCO | | | | | | Objects365 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $AP_{50}$ | $AP_{75}$ | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ |
| Supervised [11] | 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 |
| ViLD* [16] | 73.9 | 57.9 | 34.1 | 52.3 | 36.5 | 21.6 | 38.9 | 46.1 | 11.5 | 17.8 | 12.3 | 4.2 | 11.1 | 17.8 |
| BARON (Ours)‡ | **74.5** | **57.9** | **36.3** | **56.1** | **39.3** | **25.4** | **39.5** | **48.2** | **13.2** | **20.0** | **14.0** | **4.8** | **12.7** | **20.1** |
| DetPro [11] | 74.6 | 57.9 | 34.9 | 53.8 | 37.4 | 22.5 | 39.6 | 46.3 | 12.1 | 18.8 | 12.9 | 4.5 | 11.5 | 18.6 |
| BARON (Ours) | **76.0** | **58.2** | **36.2** | **55.7** | **39.1** | **24.8** | **40.2** | **47.3** | **13.6** | **21.0** | **14.5** | **5.0** | **13.1** | **20.7** |

## 4.1. Benchmark Results

**OV-COCO.** We report the comparison with previous methods in Table 1. BARON surpasses previous state of the arts with either pre-trained VLMs or COCO captions [7], indicating its effectiveness and flexibility. Note that OV-DETR is based on the Deformable DETR [62], which is stronger than the Faster R-CNN [44] with higher performance on base categories. But BARON still outperforms OV-DETR by 4.6 $AP_{50}$ on novel categories. When using caption supervision, BARON even outperforms PB-OVD [14] that uses sophisticated pseudo-labeling. Combining CLIP image features, COCO captions, and MAVL [36] proposals, BARON significantly outperforms Rasheed *et al.* [43] by a large margin.

**OV-LVIS.** We compare BARON with other methods on the OV-LVIS benchmark in Table 2. Because ViLD [16] is trained with large-scale jittering [15] and a prohibitive 32× schedule, DetPro [11] re-implemented it with backbone weights pre-trained by SOCO [51] and a regular 2× schedule. DetPro also proposes learned prompts for the category's names. Besides, an ensembling strategy for classification scores is adopted in ViLD and DetPro. For fair comparison, we respectively implement our method on OV-LVIS with and without the these tricks. BARON achieves the best performance in all the scenarios and can even surpass ViLD that adopts the ensembling strategy without these tricks.

**Transfer to Other Datasets.** We transfer the open-vocabulary detector trained on OV-LVIS to three other datasets, including Pascal VOC 2007 [12] test set, COCO [33] validation set and Objects365 [45] v2 validation set. We compare our results with DetPro [11] and the

ViLD [16] implemented in DetPro. The comparison is fair since all the models are based on the same object detector and training schedule. As shown in Table 3, our approach exhibits better generalization ability on all of the three datasets.

## 4.2. Ablation Study

In this section, we ablate the effectiveness of components in BARONon OV-COCO benchmark.

**Effectiveness of Aligning Bag of Regions.** We start from a baseline that only uses the individual-level loss $\mathcal{L}_{\text{individual}}$. As shown in Table 4(#1), the individual-level baseline achieves 25.7 $mAP_{50}$ on novel categories. We then replace $\mathcal{L}_{\text{individual}}$ with the loss for bag of regions $\mathcal{L}_{\text{bag}}$ (Table 4(#2)). Without considering the spatial information of region boxes, the performance on novel categories is compatible with aligning individual regions. When adding the positional embedding of the region boxes' spatial information, the performance on novel categories (Table 4(#3)) dramatically increases by 7.1 $mAP_{50}$. This means that the spatial information is essential to effectively exploit the compositional structure of co-occurring visual concepts in a bag of regions. Finally, in Table 4(#4), we find that the individual-level loss is complementary to the bag-of-regions alignment, which brings 1.2 $mAP_{50}$ performance gain on novel categories.

**Sampling Strategies.** We explore two baselines to sample bags of regions to support the rationale of our neighborhood sampling strategy. The first is to equally split an image into grids (dubbed as grid sampling) like the pre-training stage in OVR-CNN [56] such that the fixed grids form a bag of regions. And the second is to randomly sample region proposals to form a bag of regions (dubbed as random

Table 4. Effectiveness of main components of BARON

| # | $\mathcal{L}_{\text{individual}}$ | $\mathcal{L}_{\text{bag}}$ | PE | $\text{AP}_{50}^{\text{novel}}$ | $\text{AP}_{50}^{\text{base}}$ | $\text{AP}_{50}$ |
|---|---|---|---|---|---|---|
| 1 | ✓ | - | - | 25.7 | 59.6 | 50.6 |
| 2 | - | ✓ | - | 25.7 | 59.4 | 50.5 |
| 3 | - | ✓ | ✓ | 32.8 | 60.1 | 53.0 |
| 4 | ✓ | ✓ | ✓ | **34.0** | 60.4 | 53.5 |

Table 5. Exploring sampling strategies to obtain bag of regions

| # | Sampling strategy | $\text{AP}_{50}^{\text{novel}}$ | $\text{AP}_{50}^{\text{base}}$ | $\text{AP}_{50}$ |
|---|---|---|---|---|
| 1 | Grid | 25.4 | 58.0 | 49.5 |
| 2 | Random | 27.3 | 53.3 | 46.5 |
| 3 | Ours (reduced) | **32.2** | 58.3 | 51.5 |
| 4 | Ours | **34.0** | 60.4 | 53.5 |

Table 6. Overlap (IOF) between regions

| # | Overlap | $\text{AP}_{50}^{\text{novel}}$ | $\text{AP}_{50}^{\text{base}}$ | $\text{AP}_{50}$ |
|---|---|---|---|---|
| 1 | -0.1 | 32.5 | 59.7 | 52.6 |
| 2 | 0.0 | 33.6 | 60.2 | 53.2 |
| 3 | 0.1 | **34.0** | 60.4 | 53.5 |
| 4 | 0.2 | 33.8 | 59.8 | 53.0 |
| 5 | 0.3 | 33.7 | 60.0 | 53.1 |

Table 7. Ablation study on the sampling probability $p_b$

| $p_b$ | $\text{AP}_{50}^{\text{novel}}$ | $\text{AP}_{50}^{\text{base}}$ | $\text{AP}_{50}$ |
|---|---|---|---|
| 0.1 | 33.7 | 59.8 | 52.9 |
| 0.3 | **34.0** | 60.4 | 53.5 |
| 0.5 | 33.2 | 60.0 | 53.0 |

Table 8. Number of sampled bags per region proposal

| #bags | $\text{AP}_{50}^{\text{novel}}$ | $\text{AP}_{50}^{\text{base}}$ | $\text{AP}_{50}$ |
|---|---|---|---|
| 1 | 32.6 | 58.5 | 51.7 |
| 3 | **34.0** | 60.4 | 53.5 |
| 5 | 33.2 | 60.0 | 53.0 |

Table 9. Number of pseudo words

| # | #words | $\text{AP}_{50}^{\text{novel}}$ | $\text{AP}_{50}^{\text{base}}$ | $\text{AP}_{50}$ |
|---|---|---|---|---|
| 1 | 2 | 31.6 | 59.5 | 52.2 |
| 2 | 4 | 33.1 | 60.1 | 53.0 |
| 3 | 6 | **34.0** | 60.4 | 53.5 |
| 4 | 8 | 33.5 | 59.9 | 53.0 |

sampling). For a fair comparison, we keep the number of sampled regions in each sampling strategy roughly the same. Concretely, we split the images to $3 \times 3 = 9$ grids for the grid sampling strategy and sample 9 region proposals for the random sampling strategy. For these two strategies, we take 4 permutations of the regions to obtain richer bag-of-regions embeddings so that there would be 36 regions for each image. For the neighborhood sampling, we introduce a reduced version of the strategy that restricts the number of region proposals per image to 12 and takes 1 bag per proposal. This is because we record that the average number of regions in a bag is 3 with $p_b = 0.3$ and $\alpha = 3.0$, meaning the average sampled regions for each image in the neighborhood sampling strategy is close to $3 \times 12 = 36$.

The grid sampling strategy, whose fixed regions may either contain too many objects or only small parts of an object, achieves 25.4 mAP$_{50}$ on the novel categories as shown in Table 5(#1). For the random sampling strategy, regions in a bag have different sizes and shapes and the distance between the regions can be large, which hinders the image encoder to exactly represent the bag of regions. It achieves 27.3 mAP$_{50}$ as shown in Table 5(#2). While our neighborhood sampling strategy utilizes the same amount of regions, we record 32.2 mAP on the novel categories in Table 5(#3). Compared to these two baselines, our neighborhood sampling strategy captures potential objects in the vicinity of region proposals and ensures the teacher embedding exactly represents the bag of regions by sampling the neighboring boxes of region proposals. Note that our result in Table 5(#4) is achieved with 3 groups per region proposal and no limitation on the number of region proposals. More details on how we develop the sampling strategy are in the appendix.

**Box Overlap between Regions in a Bag.** We let the sampled regions in a group overlap at a certain IOF. In Table 6, we show how the overlap between boxes affects the performance. The scalar smaller than 0.0 in Table 6(#1) means that we keep an interval between the sampled boxes. The best performance (34.0 mAP$_{50}$) comes with an overlap of

0.1. Table 6(#1−3) indicate that the regions in a group need to have a continuity of semantics. Table 6(#3 − 5) indicate that the regions also need to cover diverse image contents.

**Sampling Probability.** We study the effect of the probability $p_b$ to sample the candidate boxes, which affects the number of regions in a bag. The details are in Sec. 3.2. Note that the $p_b$ will be adjusted by the aspect ratio scaled with the scale factor $\alpha$. We fix $\alpha$ as 3.0 and sample 3 bags for each region proposal. We observe that the best result on novel categories (34.0 mAP$_{50}$) is achieved with $p_b = 0.3$ in Table 7.

**Number of Bags Per Proposal.** We show the effect of the number of sampled bags (#bags) for each region proposal in Table 8. For details, please refer to Sec 3.2. We fix $\alpha$ as 3.0 and the sampling probability $p_b$ as 0.3. We observe that the best result on novel categories (34.0 mAP$_{50}$) is achieved with three bags of regions for each region proposal in Table 8.

**Number of Pseudo Words Per Region.** As an object category often needs many words to reach a precise description, we study the number of pseudo words (#wordss) predicted for each region of interest. Table 9(#1 − 3) show that the performance on novel categories increases with the number of pseudo words, indicating that stacking more pseudo words to a certain extent can strengthen the detector's ability to distinguish object categories. However, the results in Table 9(#1 − 3) show that further increasing the number of pseudo words does not bring performance gain, and redundant words can even do harm to the performance.

### 4.3. Further Analysis

We show qualitative results in this section to further analyze the effectiveness of our method.

**Co-occurrence of Objects.** We first illustrate how a pretrained VLM [40] captures the co-occurrence of objects by comparing the image-text similarity in Fig. 6. We use the CLIP-ViT-B/32 model pre-trained on more than 400 million image-text pairs to obtain the image and text embeddings. For each image, we incrementally add the object categories that appear to the text description. We observe that the simi-
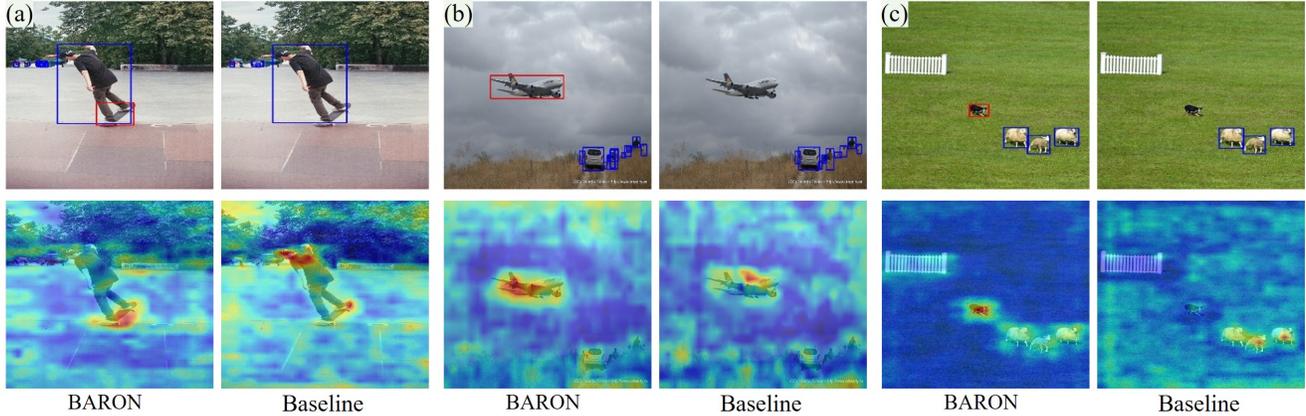
Figure 5. Qualitative comparisons of BARON and the individual-level baseline in Table 4(#1). **Top:** Red boxes are for the novel categories and blue for the base categories. **Bottom:** Feature map's responses to the queried novel object categories. From (a) to (c), the queried novel categories are 'skateboard', 'airplane' and 'dog', respectively. BARON detects objects of novel categories that are missed by the baseline.



*"There is a desk."* (0.265)
*"There is a desk with a monitor."* (0.277)
*"There is a desk with a monitor and keyboard."* (0.283)
*"There is a desk with a monitor, keyboard and mouse."* (0.294)

*"There is a black motorcycle."* (0.272)
*"There is a black motorcycle parked on the road."* (0.279)
*"There are a black motorcycle and a car parked on the road."* (0.295)
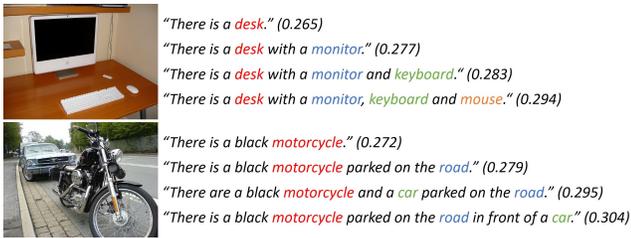*"There is a black motorcycle parked on the road in front of a car."* (0.304)

Figure 6. For each image, we incrementally add the object categories that appear in it to the text description. The similarity score between the image and text embeddings increases as the text description becomes more complete and precise.

larity score between the image and text embeddings increases as the text description includes more concepts as shown in Fig. 6. The examples reveal that the large-scale VLMs could capture the co-occurrence of multiple concepts in an image, although they are not explicitly trained to do so. We believe this is because each image-text pair naturally contains multiple concepts and the VLMs could implicitly learn the underlying connections when training on massive-scale image-text pairs. We also find simple relationship between objects can also be captured by the VLMs, *e.g.*, the similarity increases when we add the relation 'in front of' to the text description in the second example of Fig 6.

**Visualization.** We further visualize the predictions of detectors learned through BARON and the individual-level baseline in Fig. 5. The images are from COCO's validation set. We visualize the feature map's response to the novel categories using the Grad-CAM++ [4]. We find that the model learned through BARON generates responses at locations of novel categories while the individual-level baseline induces weaker, incomplete or diffused responses. We also notice that semantically related objects can respond to the queried object category. For example, in Fig 5(c) the location of a flock of sheep chased by a dog responds to the queried ob-

ject 'dog'. Note the 'dog' and the 'flock of sheep' are not in neighboring regions and the 'dog' is much smaller in size. Even though we form bags of regions by neighboring regions of equal size in training, the model has the ability to capture the relationship between objects with imbalanced sizes or large distance. This phenomenon resembles the generalization ability of human language, where learned concepts can be applied to describe or recognize new things.

## 5. Discussion and Conclusion

This paper goes beyond the learning of individual regions to bag of regions in OVD, exploring the ability of large-scale VLMs to represent the compositional structure of multiple concepts that naturally exists in image-text pairs. We develop a neighborhood sampling strategy to group contextually related regions into a bag and adopt the contrastive learning approach to align the bag-of-regions representations of the detector and pre-trained VLMs, which achieves new state-of-the-art performance on multiple OVD benchmarks.

The compositional structure explored in this paper is mainly about the co-occurrence of objects, and behaves like bag of words [48, 54]. The more complex compositional structure in the language is still under-explored and whether modern pre-trained vision-language models capture such structure still remains an open problem for the community. We look forward to further unveiling the behavior of the VLMs, and more importantly endowing the VLMs with human-like compositional representation to move to more generalized intelligence.

## A1. Implementation Details

We provide more details of the implementation of BARON on OV-COCO [33] and OV-LVIS [18] benchmarks. **Sampling.** For neighborhood sampling strategy, we obtain top $K$ region proposals from the RPN and filter out those with an objectness score lower than 0.85. We also discard regions with an aspect ratio smaller than 0.25 or larger than 4.0. And regions with an area ratio smaller than 0.01 are also discarded. Then we apply NMS on the region proposals with IOU threshold 0.1. The region proposals after NMS are used for neighborhood sampling. We sample $G$ bags of regions for each region proposal with a probability 0.3 to sample each surrounding candidate box. For OV-COCO, we set $K = 300$ and $G = 3$. For OV-LVIS, we set $K = 500$ and $G = 4$ due to the denser spatial distribution of object boxes in the LVIS dataset.

**Classification Loss.** We use CE loss as the classification loss $\mathcal{L}_{\text{cls}}$ on base categories. Given $C$ object categories, we obtain the embedding $f_i$ for the name of the $i$-th category by the text encoder ($\mathcal{T}$) of the VLM. We also learn a background embedding for non-object regions. If a region is labeled as the $c$-th category, the classification loss is

$$\mathcal{L}_{\text{cls}} = -\log \frac{\exp(\tau_{\text{cls}} \cdot \langle \mathcal{T}(w), f_c \rangle)}{\sum_{i=0}^{C} \exp(\tau_{\text{cls}} \cdot \langle \mathcal{T}(w), f_i \rangle)}, \quad (5)$$

where $\tau_{\text{cls}}$ is the temperature to re-scale the cosine similarity, $f_C$ is the background embedding and $w$ is the embedding (pseudo words) of the region. On OV-COCO, we set $\tau_{\text{cls}} = 50.0$. And on OV-LVIS, we set $\tau_{\text{cls}} = 100.0$ since there are orders of magnitude more categories defined in the LVIS dataset.

**Alignment Loss.** Assuming there are $G$ bags of regions and the image (teacher) and text (student) embeddings for the $k$-th bag of regions are $f_v^k$ and $f_t^k$, the alignment loss $\mathcal{L}_{\text{bag}}$ on bag of regions is calculated as

$$\mathcal{L}_{\text{bag}} = -\frac{1}{2} \sum_{k=0}^{G-1} (\log(p_{t,v}^k) + \log(p_{v,t}^k)). \quad (6)$$

The $p_{t,v}^k$ and $p_{v,t}^k$ are calculated as

$$p_{t,v}^k = \frac{\exp(\tau_{\text{bag}} \cdot \langle f_t^k, f_v^k \rangle)}{\sum_{l=0}^{G-1} \exp(\tau_{\text{bag}} \cdot \langle f_t^k, f_v^l \rangle)} \quad (7)$$

$$p_{v,t}^k = \frac{\exp(\tau_{\text{bag}} \cdot \langle f_v^k, f_t^k \rangle)}{\sum_{l=0}^{G-1} \exp(\tau_{\text{bag}} \cdot \langle f_v^k, f_t^l \rangle)}, \quad (8)$$

respectively, where $\tau_{\text{bag}}$ is the temperature to re-scale the cosine similarity. Assuming there are totally $N$ regions and the image (teacher) and text (student) embeddings for the $k$-th region are $g_v^k$ and $g_t^k$, the alignment loss $\mathcal{L}_{\text{individual}}$ on individual regions is calculated as

Table A1. Number of linear layers (#Layers) mapping region features to pseudo-words

| #Layers | $AP_{50}^{novel}$ | $AP_{50}^{base}$ | $AP_{50}$ |
|---|---|---|---|
| 1 | 34.0 | 60.4 | 53.5 |
| 2 | 33.9 | 60.5 | 53.5 |
| 3 | 34.1 | 60.8 | 53.8 |

$$\mathcal{L}_{\text{individual}} = -\frac{1}{2} \sum_{k=0}^{N-1} (\log(q_{t,v}^k) + \log(q_{v,t}^k)). \quad (9)$$

The $q_{t,v}^k$ and $q_{v,t}^k$ are calculated as

$$q_{t,v}^k = \frac{\exp(\tau_{\text{individual}} \cdot \langle g_t^k, g_v^k \rangle)}{\sum_{l=0}^{N-1} \exp(\tau_{\text{individual}} \cdot \langle g_t^k, g_v^l \rangle)} \quad (10)$$

$$q_{v,t}^k = \frac{\exp(\tau_{\text{individual}} \cdot \langle g_v^k, g_t^k \rangle)}{\sum_{l=0}^{N-1} \exp(\tau_{\text{individual}} \cdot \langle g_v^k, g_t^l \rangle)}, \quad (11)$$

respectively, where $\tau_{\text{individual}}$ is the temperature to re-scale the cosine similarity.

On OV-COCO, we set $\tau_{\text{bag}} = 30.0$ and $\tau_{\text{individual}} = 50.0$. Since there are finer-grained definition of categories and denser distribution of object boxes in the LVIS dataset, we set $\tau_{\text{bag}} = 20.0$ and $\tau_{\text{individual}} = 30.0$ on OV-LVIS to make the contrastive learning harder.

**Mapping Region Features to Pseudo-words.** In our implementation, we used a single linear layer to map region features from the detector to pseudo-words. In Table A1, we show that adding more linear layers (#Layers) brings no noticeable improvements. This observation is also in line with Maaz *et al.* [36] that visual properties can be transferred to language models (LMs) by linearly mapping visual features to the input space of LMs.

**Random Word Dropout.** As we apply two different supervision to the pseudo words, the training can lead certain words to overfit to certain losses. To alleviate overfitting, we borrow the idea of Dropout [47] in neural networks where neurons are randomly dropped during training to avoid overfitting to specific neurons. We randomly discard pseudo words for each region with a probability $p_{\text{drop}}$. By default, we set $p_{\text{drop}} = 0.5$ for training on both OV-COCO and OV-LVIS.

**Suppression on Novel Categories.** On OV-COCO, we observe a tendency to overfit on base categories due to the smaller number of categories. And compared with OV-LVIS where the tail categories act as the novel categories, the distribution of novel and base categories on COCO is more balanced. We adopt the following strategies to alleviate suppression on novel categories: (1) detach the objectness prediction branch so that the suppression onto novel categories

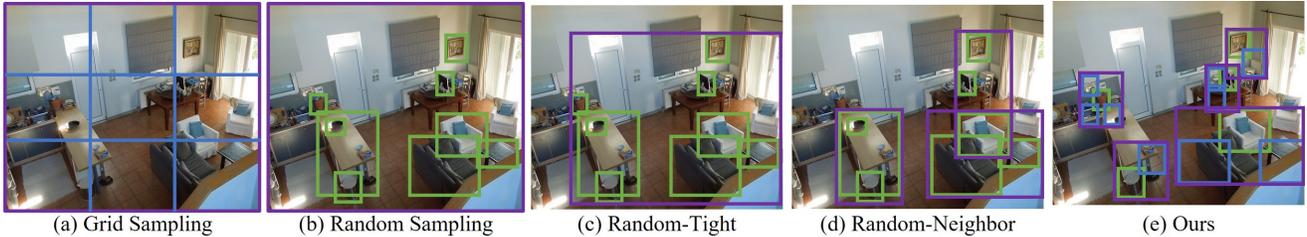(a) Grid Sampling   (b) Random Sampling   (c) Random-Tight   (d) Random-Neighbor   (e) Ours

Figure A1. Comparisons of different sampling strategies. Green boxes denote the region proposals. Blue boxes stand for sampled region boxes. The purple box represents the image crop of a bag of regions (a region group).

Table A2. Different sampling strategies

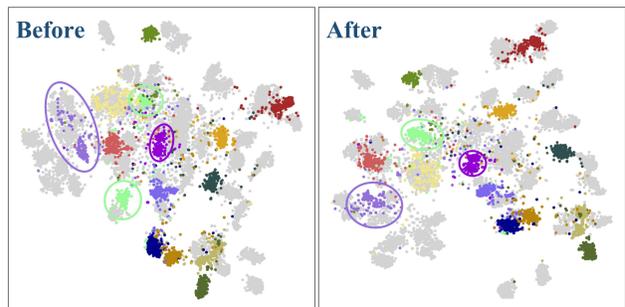| # | Strategy | $AP_{50}^{novel}$ | $AP_{50}^{base}$ | $AP_{50}$ | #regions |
|---|----------|-------|------|------|----------|
| 1 | Grid | 25.4 | 58.0 | 49.5 | 36 |
| 2 | Random | 27.3 | 53.3 | 46.5 | 36 |
| 3 | Random-Tight | 29.5 | 56.9 | 49.7 | 36 |
| 4 | Random-Neighbor | 30.7 | 56.9 | 50.0 | 36 |
| 5 | Ours (reduced) | **32.2** | 58.3 | 51.5 | 36 |
| 6 | Ours | **34.0** | 60.4 | 53.5 | 216 |



Figure A2. tSNE visualization of embeddings on COCO categories. **Left:** the region features *before* being projected to pseudo words. **Right:** embeddings *after* sending pseudo words to the text encoder.

would not be back-propagated to the backbone; (2) save the sampled region proposals into a cache so that regions covering potential novel categories detected in certain iteration can be preserved throughout the training phase; (3) use the output of the second last layer of the VLM (CLIP) for classification and the final output for aligning bag of regions to reduce the competition between the two types of losses.

## A2. Sampling Strategy

We have introduced two baseline sampling strategies, *i.e.* grid sampling and random sampling. The grid sampling strategy is to equally split an image into grids like the pre-training stage in OVR-CNN [56]. And the random sampling strategy is to randomly sample region proposals to form a bag of regions. These two baseline strategies let the bag of regions represent the whole image. We add two other strategies to shift the focus to neighboring (local) regions.

We start from the random sampling strategy and let the bag of regions represent the image crop that tightly encloses them instead of the whole image (dubbed as Random-Tight). Then we move to the neighborhood centered on region proposals (dubbed as Random-Neighbor). For each center region proposal, we randomly sample 2 nearby region proposals with GIOU larger than 0.5 to make a bag of regions. We randomly take 12 region proposals as centers so that the total number of regions is 36, ensuring a fair comparison with other strategies. Table A2 shows the performance of these strategies.

In Fig A1, we show how these sampling strategies differ

and how it gradually develops to our final option. In (a), we find the equally split grids may either contain too many objects or only small parts of an object. From (b) to (c), the bag of regions gradually shift to representing neighboring local regions from representing the whole image. However, we observe that there is always box size imbalance such as the left bottom bag of regions in (d). And there are also large area of redundant image contents between the regions in a bag as shown in (c). The box size imbalance and the redundant image contents hinder the image encoder of a VLM to effectively represent a bag of regions. As shown in (e), our sampling strategy obtains a bag of neighboring regions of equal size while capturing potential objects. Although we still observe image contents between sampled regions that do not belong to a bag of regions, they only account for a small portion of the image crop enclosing the bag of regions.

## A3. Pseudo Word Encoding

Projecting visual features to word embedding space is common in region-based visual-language representation learning methods [8, 35]. In BARON, we project region features into pseudo words to fully exploit the inherent compositional structure of multiple semantic concepts and obtain more distinctive feature embeddings. In Fig A2, we show the tSNE visualization of the region features *before* being pro-
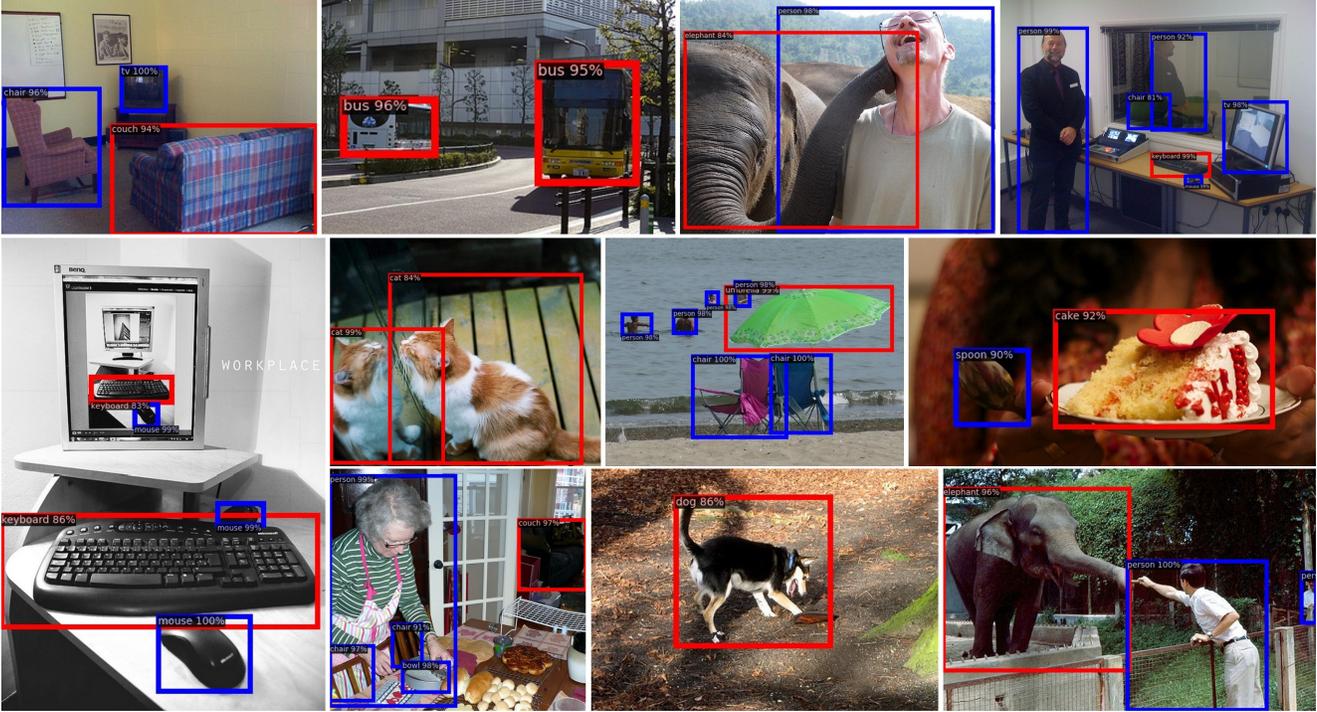
Figure A3. Visualization of detection results on OV-COCO. Red boxes are for novel categories, while blue boxes are for base categories.

jected to pseudo words and embeddings *after* sending pseudo words to the text encoder (TE), *i.e.* $\mathcal{T}(w)$. Gray points represent base categories while chromatic points represent novel categories. With pseudo words encoded by TE, the categories are split into clusters of a more diverse distribution and distinct boundaries.

## A4. Image-Guided Inference

We further examine the generalization ability of our method by using images to guide the inference of the detector. We use the image encoder of CLIP [40] to encode the reference image. And the detector used in this experiment is trained on the LVIS dataset. Given a reference image, our detector is able to detect the object in the reference image as shown in Fig A4. Our detector can even recognize the cartoon characters in the reference images ('pikachu' and 'winnie pooh').

## A5. Detection Results

We show more detection results of our method in Fig A3 and Fig A5. On COCO dataset, BARON correctly detects novel categories including bus, keyboard, couch and so on. On LVIS dataset, BARON detects rare categories like salad plate, fedora hat, gas mask and so on. We also visualize the results when transferring the LVIS-trained detector to Objects365 [45] dataset in Fig A6. We find that the LVIS-



Figure A4. Image-guided inference of the detector trained on LVIS dataset. BARON can even recognize the cartoon characters in the reference images ('pikachu' and 'winnie pooh').

trained detector is able to correctly recognize a wide range of object concepts defined in Objects365 dataset, exhibiting impressive generalization ability.

## A6. Potential Negative Societal Impacts

Our models have learned knowledge from vision-language models (VLMs) that are pre-trained on large-scale web image-text pairs. They potentially inherit and even re-inforce harmful biases and stereotypes in the pre-trained VLMs. We suggest scrupulous probing before applying our models for any purpose.
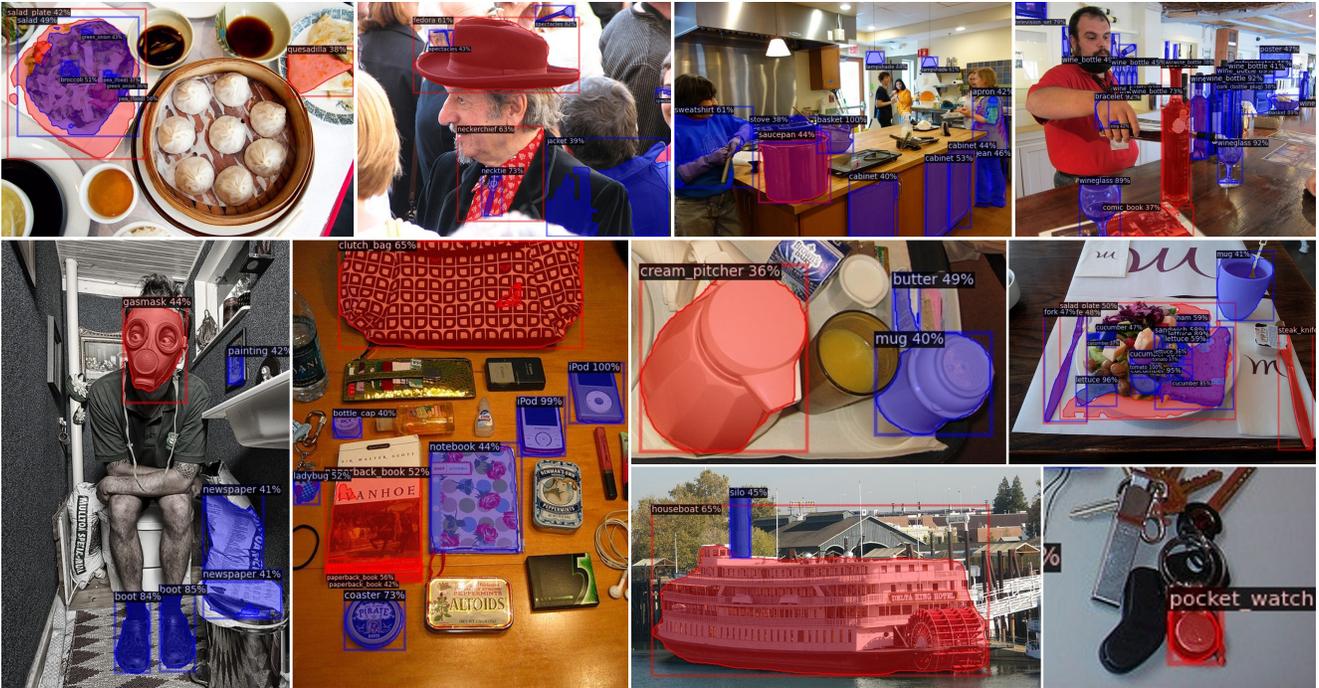
Figure A5. Visualization of detection results on OV-LVIS dataset. Red boxes and masks are for novel (rare) categories, while blue boxes and masks are for base categories.



Figure A6. Visualization of transfer detection results on Objects365 dataset.

# References

[1] Edward H. Adelson. On seeing stuff: the perception of materials by humans and machines. In Bernice E. Rogowitz and Thrasyvoulos N. Pappas, editors, *Human Vision and Electronic Imaging VI*, SPIE Proceedings, pages 1–12, 2001. 1

[2] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. In *Eur. Conf. Comput. Vis.*, 2018. 2

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Eur. Conf. Comput. Vis.*, 2020. 2, 3

[4] Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. In *IEEE Winter App. Comput. Vis.*, 2018. 8

[5] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Hybrid task cascade for instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 2

[6] 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. 5

[7] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 6

[8] 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 *Eur. Conf. Comput. Vis.*, pages 104–120, 2020. 10

[9] Berkan Demirel, Ramazan Gokberk Cinbis, and Nazli Ikizler-Cinbis. Zero-shot object detection by hybrid region embedding. In *Brit. Mach. Vis. Conf.*, 2018. 2

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. Learn. Represent.*, 2021. 4, 5

[11] Yu Du, Fangyun Wei, Zihe Zhang, Miaojing Shi, Yue Gao, and Guoqi Li. Learning to prompt for open-vocabulary object detection with vision-language model. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 1, 2, 3, 5, 6

[12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vis.*, 88(2):303–338, 2010. 6

[13] Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Marc'Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. *Adv. Neural Inform. Process. Syst.*, 2013. 2

[14] Mingfei Gao, Chen Xing, Juan Carlos Niebles, Junnan Li, Ran Xu, Wenhao Liu, and Caiming Xiong. Open vocabulary object detection with pseudo bounding-box labels. *Eur. Conf. Comput. Vis.*, 2022. 2, 5, 6

[15] 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 *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 6

[16] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *Int. Conf. Learn. Represent.*, 2021. 1, 2, 3, 5, 6

[17] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5356–5364, 2019. 5

[18] Agrim Gupta, Piotr Dollár, and Ross B. Girshick. LVIS: A dataset for large vocabulary instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 9

[19] Nasir Hayat, Munawar Hayat, Shafin Rahman, Salman H. Khan, Syed Waqas Zamir, and Fahad Shahbaz Khan. Synthesizing the unseen for zero-shot object detection. In *Asian Conf. Comput. Vis.*, 2020. 2

[20] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 4

[21] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Int. Conf. Comput. Vis.*, 2017. 5

[22] Dinesh Jayaraman and Kristen Grauman. Zero-shot recognition with unreliable attributes. In *Adv. Neural Inform. Process. Syst.*, 2014. 2

[23] 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 *Int. Conf. Mach. Learn.*, 2021. 1, 2

[24] Wonjae Kim, Bokyung Son, and Ildoo Kim. ViLT: Vision-and-language transformer without convolution or region supervision. In *Int. Conf. Mach. Learn.*, 2021. 2

[25] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9404–9413, 2019. 1

[26] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *Int. Conf. Learn. Represent.*, 2022. 2

[27] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *Int. Conf. Mach. Learn.*, 2022. 2

[28] Junnan Li, Ramprasaath R. Selvaraju, Akhilesh Gotmare, Shafiq R. Joty, Caiming Xiong, and Steven Chu-Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. In *Adv. Neural Inform. Process. Syst.*, 2021. 2

13

[29] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[30] Chuang Lin, Peize Sun, Yi Jiang, Ping Luo, Lizhen Qu, Gholamreza Haffari, Zehuan Yuan, and Jianfei Cai. Learning object-language alignments for open-vocabulary object detection. *Int. Conf. Learn. Represent.*, 2023. 5

[31] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. Feature pyramid networks for object detection. In *Int. Conf. Comput. Vis.*, 2017. 5

[32] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Int. Conf. Comput. Vis.*, 2017. 3

[33] 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 *Eur. Conf. Comput. Vis.*, 2014. 4, 5, 6, 9

[34] Li Liu, Wanli Ouyang, Xiaogang Wang, Paul W. Fieguth, Jie Chen, Xinwang Liu, and Matti Pietikäinen. Deep learning for generic object detection: A survey. *Int. J. Comput. Vis.*, 2020. 2

[35] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Adv. Neural Inform. Process. Syst.*, 2019. 2, 10

[36] Muhammad Maaz, Hanoona Rasheed, Salman Khan, Fahad Shahbaz Khan, Rao Muhammad Anwer, and Ming-Hsuan Yang. Class-agnostic object detection with multi-modal transformer. In *Eur. Conf. Comput. Vis.*, pages 512–531, 2022. 5, 6, 9

[37] Mohammad Norouzi, Tomás Mikolov, Samy Bengio, Yoram Singer, Jonathon Shlens, Andrea Frome, Greg Corrado, and Jeffrey Dean. Zero-shot learning by convex combination of semantic embeddings. In *Int. Conf. Learn. Represent.*, 2014. 2

[38] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 4

[39] Chao Peng, Tete Xiao, Zeming Li, Yuning Jiang, Xiangyu Zhang, Kai Jia, Gang Yu, and Jian Sun. MegDet: A large mini-batch object detector. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 5

[40] 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 *Int. Conf. Mach. Learn.*, 2021. 1, 2, 4, 5, 7, 11

[41] Shafin Rahman, Salman H. Khan, and Nick Barnes. Transductive learning for zero-shot object detection. In *Int. Conf. Comput. Vis.*, 2019. 2

[42] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. DenseCLIP: Language-guided dense prediction with context-aware prompting. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[43] Hanoona Abdul Rasheed, Muhammad Maaz, Muhammd Uzair Khattak, Salman Khan, and Fahad Khan. Bridging the gap between object and image-level representations for open-vocabulary detection. In *Adv. Neural Inform. Process. Syst.*, 2022. 5, 6

[44] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Adv. Neural Inform. Process. Syst.*, 2015. 2, 3, 5, 6

[45] 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 *Int. Conf. Comput. Vis.*, 2019. 6, 11

[46] Hengcan Shi, Munawar Hayat, Yicheng Wu, and Jianfei Cai. ProposalCLIP: Unsupervised open-category object proposal generation via exploiting clip cues. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[47] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 2014. 9

[48] Ajinkya Tejankar, Maziar Sanjabi, Bichen Wu, Saining Xie, Madian Khabsa, Hamed Pirsiavash, and Hamed Firooz. A fistful of words: Learning transferable visual models from bag-of-words supervision. *CoRR*, abs/2112.13884, 2021. 8

[49] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Eur. Conf. Comput. Vis.*, 2020. 2

[50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Adv. Neural Inform. Process. Syst.*, 2017. 4

[51] Fangyun Wei, Yue Gao, Zhirong Wu, Han Hu, and Stephen Lin. Aligning pretraining for detection via object-level contrastive learning. *Adv. Neural Inform. Process. Syst.*, 2021. 5, 6

[52] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. https://github.com/facebookresearch/detectron2, 2019. 5

[53] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Bin Xiao, Ce Liu, Lu Yuan, and Jianfeng Gao. Unified contrastive learning in image-text-label space. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 5

[54] Mert Yüksekgönül, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? *CoRR*, abs/2210.01936, 2022. 8

[55] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary DETR with conditional matching. *Eur. Conf. Comput. Vis.*, 2022. 1, 2, 5, 6

[56] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 2, 5, 6, 10

[57] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. LiT: Zero-shot transfer with locked-image text tuning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[58] Ye Zheng, Ruoran Huang, Chuanqi Han, Xi Huang, and Li Cui. Background learnable cascade for zero-shot object detection. In *Asian Conf. Comput. Vis.*, 2020. 2

[59] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. RegionCLIP: Region-based language-image pretraining. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2, 5

[60] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *Eur. Conf. Comput. Vis.*, 2022. 1, 2

[61] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Phillip Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. *Eur. Conf. Comput. Vis.*, 2022. 2, 5

[62] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable transformers for end-to-end object detection. *Int. Conf. Learn. Represent.*, 2021. 2, 6