

Adapting CLIP For Phrase Localization Without Further Training

Jiahao Li Greg Shakhnarovich Raymond A. Yeh
Toyota Technological Institute at Chicago
{jiahao, greg, yehr}@ttic.edu

Abstract

Supervised or weakly supervised methods for phrase localization (textual grounding) either rely on human annotations or some other supervised models, *e.g.*, object detectors. Obtaining these annotations is labor-intensive and may be difficult to scale in practice. We propose to leverage recent advances in contrastive language-vision models, CLIP, pre-trained on image and caption pairs collected from the internet. In its original form, CLIP only outputs an image-level embedding without any spatial resolution. We adapt CLIP to generate high-resolution spatial feature maps. Importantly, we can extract feature maps from both ViT and ResNet CLIP model while maintaining the semantic properties of an image embedding. This provides a natural framework for phrase localization. Our method for phrase localization requires no human annotations or additional training. Extensive experiments show that our method outperforms existing no-training methods in zero-shot phrase localization, and in some cases, it even outperforms supervised methods. Code is available at <https://github.com/pals-ttic/adapting-CLIP>.

1 Introduction

Phrase Localization (a.k.a. textual grounding) is the task of localizing bounding boxes referred by textual phrases in a given image. It has many down-stream applications, *e.g.*, visual question answering, image caption, and human computer interaction.

With the advancement in deep learning models, supervised training of models has emerged as an dominant approach to build effective phrase localization systems [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]. Critically, these methods require human annotations specific to textual grounding, *i.e.*, triplets of (image, phrases, bounding boxes). The annotation process to collect such a training set is expensive, labor intensive, and error-prone, especially, when a large dataset is necessary to effectively train these deep models. To address this, some efforts [11, 12, 13, 14, 15] have considered weakly-supervised/unsupervised methods that rely on pre-trained supervised models, *e.g.*, object detectors, from other tasks [16, 17, 18].

Such pre-trained components are naturally biased towards the few commonly seen object categories, such as person, car, *etc.* These biases from the pre-trained model potentially lead to poor performance for uncommon or unseen objects in the training/pre-training set. To study these biases, ZSGNet [19] propose zero-shot phrase localization, where the train and test sets have “non-overlapping” nouns hence evaluating models’ zero-shot capability. In more details, they proposed different levels of “non-overlapping”, at the word-level or at the category level. For example, “Toyota” and “Honda” would be non-overlapping at the word level, but considered overlapping at the category level. However, ZSGNet remains a supervised approach which requires human annotations on phrase localization.

To address the aforementioned issues, we propose a method that **does not rely on** any textual grounding, image classification or bounding box annotations. Instead, we leverage recent advances in large-scale contrastive language-vision model, CLIP [20], trained on a dataset of 400M image-text pairs collected from the internet. We propose a method to adapt/re-purpose CLIP for phrase localization. Importantly, our method can do so **without any extra supervision or training**, *i.e.*, our method immediately improves with more advanced CLIP models.



Figure 1: Illustration of different zero-shot phrase localization splits. ZSGNet [19] proposes different level of semantic overlaps, specifically, in Flickr-S0 only the query phrase needs to be unseen *e.g.*, “a used car” and “a cab” can both refer to a car. However, this is not allowed in Flickr-S1, where the object category is also required to be unseen during training, *e.g.*, “red bucket” is not seen during training; both the category and the phrase is novel.

At a high-level, our method extracts high-resolution pixel-wise semantic features from images. By construction, these features lie in the same semantic space as the text embedding extracted using CLIP. As a result, we can compute per-pixel similarity scores with a given text query to obtain a heatmap. Localizing the described object in the text query becomes a score maximization problem, *i.e.*, finding a bounding box which achieves the highest score characterized by the heatmap.

We demonstrate the effectiveness of our model on zero-shot phrase localization using Flickr30k Entities [21] and Visual Genome (VG) [22] datasets following the zero-shot setting proposed in ZSGNet [19]. Our approach outperforms ZSGNet by an absolute 5% on three out of four zero-shot splits over Flickr30k and VG, despite having never seen any textual grounding or object detection annotations. We also achieve comparable performance in long-tailed object categories compared to no-training methods that utilize pre-trained object detectors.

Summary of our contributions:

- We propose a method for textual grounding entirely from a pre-trained language-vision model (CLIP) without any bounding box supervision.
- We design methods for extracting high-resolution spatial feature maps from CLIP, for both ViT and ResNet architectures.
- We conduct extensive experiments and ablation studies demonstrating the effectiveness of our approach in zero-shot phrase localization.

2 Related Works

Phrase Localization. To study and evaluate the progress of phrase localization, numerous datasets, such as Flickr30k Entities [21], Visual Genome [23] and ReferItGame [24] have been proposed. These datasets contain a rich set of annotations covering a diverse set of objects and phrases.

Numerous approaches have been proposed to tackle the task of phrase localization [1, 2, 3, 4, 5, 7, 8, 9, 10]. These methods can be roughly divided into two groups: (a) Two-stage methods based on the the classical proposal-classification paradigm of object detection [3, 25, 26]. These take object proposals from the image and associated them with the corresponding query text, *e.g.*, ranking these proposals based on embedding similarity to the text query. (b) In contrast, one-stage methods [27, 28] build an end-to-end pipeline without an intermediate proposal stage and are directly trained via bounding box regression.

Another line of work in phrase localization aims to reduce the requirement for textual grounding annotations [12, 14] in a weakly supervised manner. A few efforts exist [13, 15] to build phrase localization models without any phrase localization data at all. Specifically, they utilized various off-the-shelf components such as detectors, word embedding models to create such systems. However, we note that these off-the-shelf detectors, *etc.*, are themselves trained with human-annotated bounding boxes.

Finally, most relevant to our work is ZSGNet [19] where they propose to evaluate the zero-shot generalization of phrase localization systems. Specifically, they re-split Flickr30K Entities and Visual Genome such that they have different levels of semantic overlap; see an illustrating example in Fig. 1.

This helps quantify how well the model can generalize to completely unseen objects. The proposed ZSGNet remains a supervised model, *i.e.*, they still train on examples of (image, phase, bounding box). Different from ZSGNet, we propose to adapt a pre-trained CLIP for phrase localization. Our method does not require any annotations nor training (besides optional tuning of hyper-parameters on a validation set), while still achieving zero-shot phrase localization capacities.

Vision Language Models. Recently, large-scale contrastive models across language and vision domains have received a lot of attention, *e.g.*, CLIP [20] and ALIGN [29]. Trained on $10^8 - 10^9$ of associated text and image pairs, these models aim to learn a joint embedding space between an image and natural language. These joint embeddings have demonstrated impressive capabilities of transferring to down-stream classification tasks without additional fine-tuning, *i.e.*, zero-shot classification. Hence, much of very recent work aims to use this capability beyond classification such as text conditioned image generation [30], open-vocabulary and zero-shot detection [31, 32, 33] and segmentation [34, 35, 36].

Closely related to our work is DenseCLIP [35], a concurrent work on ArXiv. They propose to extract spatial features from CLIP’s ResNet image encoder. These features are then used to construct pseudo-labels for training another deep-net, to predict high-resolution semantic segmentation. While our method also generates spatial features from CLIP, our approach and goals differ: (a) we devise methods to extract spatial features from both the ResNet and (the better performing) ViT architectures, while DenseCLIP is limited to ResNet; (b) our method achieves high-resolution maps without the need of distillation or training; (c) we develop a framework for using the extracted spatial features for zero-shot phrase localization.

3 Preliminaries: CLIP and Attention

We provide a brief overview of CLIP and review Attention Pooling in detail to establish our notations. Familiarity with these concepts is necessary to understand our approach.

CLIP. Recently, contrastive pre-training has emerged as a promising approach to learning effective image representations. CLIP proposes to train a contrastive model on 400 million image and text pairs [20]. Specifically, given a paired image and text $\{(I_n, T_n)\}$, CLIP’s objective is to find image embedding $e_{\text{img}}(I_n) \in \mathbb{R}^D$ and text embedding $e_{\text{txt}}(T_n) \in \mathbb{R}^D$, such that the cosine similarity between e_{img} and e_{txt} is maximized for a given input. With the embedding trained, CLIP demonstrates zero-shot transfer to downstream image classification tasks, typically by using the text prompt “A photo of a {label}.”, where “{label}” is replaced with the category labels from the task.

Extracting Image Embedding with Attention Pooling. Given an image $I \in \mathbb{R}^{H \times W \times 3}$, CLIP uses convolution layers to extract patch features $F \in \mathbb{R}^{H' \times W' \times C}$, where the height and width is reduced to $H' \times W'$ and flattened into the first dimension. *I.e.*, $f_i \in \mathbb{R}^C$ is a feature vector of patch i , commonly referred as a *patch token*.

To create the image embedding, $e_{\text{img}} \in \mathbb{R}^D$, an additional *class token* c is introduced. This class token is passed into attention layers along with the patch tokens F , and its corresponding output is the image embedding. For readability, we denote the class token as the 0th patch token, *i.e.*, $f_0 \triangleq c$. The computation for one attention layer is as follows:

$$e_{\text{img}} = \sum_{i=0}^{H'W'} \alpha_i(\mathbf{W}^Q c, \mathbf{W}^K F) \cdot \mathbf{W}^V f_i, \quad (1)$$

where $\mathbf{W}^{Q/K/V}$ corresponds to a linear transformation to create the queries, keys and values of an attention layer. Next, α_i denotes the i^{th} attention weight:

$$\alpha_i(\mathbf{q}, \mathbf{K}) = \frac{\exp(\mathbf{q}^\top \mathbf{k}_i / \tau)}{\sum_j \exp(\mathbf{q}^\top \mathbf{k}_j / \tau)}, \quad (2)$$

where τ controls the temperature of the attention. We point out that Eq. (1) can be viewed as *an operation over the set of patch tokens*. Specifically, shuffling the patch tokens’ index results in the same e_{img} ; as along as the summation is over the same set.

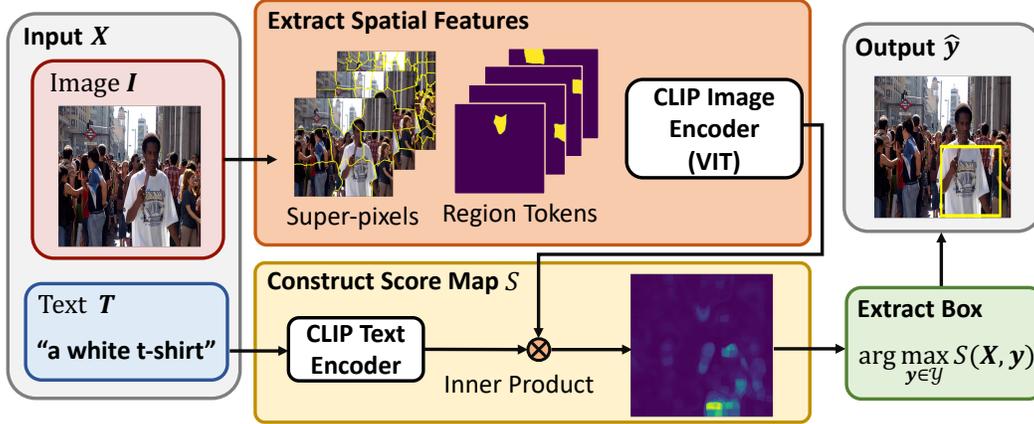


Figure 2: Proposed approach overview. Given a text query and image as input, we extract spatial features (per-pixel) and a text embedding using pre-trained CLIP. For each pixel location, we compute the cosine similarity between the spatial feature and text embedding resulting in a score map which we use to predict a bounding box formulated as a search over all possible bounding boxes.

CLIP’s ResNet architecture uses multiple convolution layers followed by an attention layer. On the other hand, CLIP’s ViT architecture uses a single convolution layer followed by several attention layers.

4 Our Approach: Region-based Attention Pooling

Our goal is to adapt a pre-trained CLIP model for the task of phrase localization. Recall, CLIP is trained to output an embedding vector for a given image or text phrase, hence the image embedding cannot be directly applied to phrase localization which requires spatial reasoning.

To obtain spatial information, we propose to extract spatial features from CLIP for both ViT (Sec. 4.3) and ResNet (Sec. 4.4) architectures. The key is to maintain the semantic meaning (*i.e.*, alignment with language embeddings) of these spatial features to be the same as the original image embedding. After obtaining these spatial features, for each pixel location, we compute the inner product between the spatial feature and the text embedding extracted from CLIP to obtain a score map. Finally, we predict the bounding box that have the largest score according to the extracted map. An overview of our proposed method is depicted in Fig. 2.

4.1 Problem Formulation

Given the input image and text query pair $\mathbf{X} = (\mathbf{I}, \mathbf{T})$, the task of phrase localization is to output the corresponding bounding box $\mathbf{y} = (y_1, y_2, y_3, y_4)$, where a bounding box is defined by its top left (y_1, y_2) and bottom-right (y_3, y_4) corners. We formulate bounding box prediction as a score maximization problem:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}} S(\mathbf{X}, \mathbf{y}), \quad (3)$$

where \mathcal{Y} denotes the set of all possible bounding boxes.

Our score function $S(\mathbf{X}, \mathbf{y})$ is composed from a score map $\phi(\mathbf{X}) \in \mathbb{R}^{H \times W}$ extracted using CLIP. The score of a given bounding box \mathbf{y} is the aggregated values within \mathbf{y} on $\phi(\mathbf{X})$, *i.e.*,

$$S(\mathbf{X}, \mathbf{y}) = \sum_{i=y_1}^{y_3} \sum_{j=y_2}^{y_4} \phi(\mathbf{X})_{ij} - R(\mathbf{y}), \quad (4)$$

where $R(\mathbf{y}) \triangleq \lambda \cdot \text{BoxArea}(\mathbf{y})$ is a penalty term regularizing the box size.

In the remaining of this section, we will discuss how to effectively design this score map $\phi(\mathbf{X})$ from a pre-trained CLIP model without any additional training or data.

4.2 Score Map Design

Our score map $\phi(\mathbf{X})$ is constructed using features extracted from CLIP. Given an input $\mathbf{X} = (\mathbf{I}, \mathbf{T})$, we extract a text embedding from CLIP, $\mathbf{e}_{\text{txt}} \in \mathbb{R}^D$. We also extract a “spatial”, per-pixel, image embedding $\mathbf{E}_{\text{img}} \in \mathbb{R}^{H \times W \times D}$ using CLIP (details deferred to Sec. 4.3 and Sec. 4.4).

As in CLIP, we compute an inner-product between the text embedding and the per-pixel image embedding to relate text to image, *i.e.*,

$$\phi(\mathbf{X})_{ij} = \exp(\mathbf{e}_{\text{txt}} \cdot (\mathbf{E}_{\text{img}})_{ij} / \sigma). \quad (5)$$

where σ denotes the temperature scaling parameter in CLIP.

With the score map defined, we will next show how to extract a spatial per-pixel image embedding \mathbf{E}_{img} from both CLIP’s ViT and ResNet architecture. Recall, CLIP’s image encoder outputs a vector of dimension \mathbb{R}^D without any spatial resolution.

4.3 Spatial Embedding from ViT

Interpreting Attention Layers as Pooling. As reviewed in Sec. 3, CLIP’s ViT encoder consists of L attention layers. At layer l , the class token embedding $\mathbf{c}^{(l+1)}$ is computed as

$$\mathbf{c}^{(l+1)} = \sum_{i=0}^{H'W'} \alpha_i (\mathbf{W}^Q \mathbf{c}^{(l)}, \mathbf{W}^K \mathbf{F}^{(l)}) \cdot \mathbf{W}^V \mathbf{f}_i^{(l)}, \quad (6)$$

recall, for readability, the class token is also denoted as the 0th patch token, *i.e.*, $\mathbf{f}_0^{(l)} \triangleq \mathbf{c}^{(l)}$.

As can be seen, each attention layer outputs a weighted sum over all the patches, *i.e.*, $i \in \{1 \dots, H'W'\}$. In other words, the class token can be viewed as “pooling” information from all the patch tokens and itself. Hence, a forward pass of CLIP’s ViT image encoder can be interpreted as an iterative pooling process as CLIP uses the final output $\mathbf{e}^{(L)}$ as the embedding vector for the entire image.

Modifying Attention for Spatial Features. As our goal is to extract spatial image embedding, we need to perform pooling over *image regions* instead of the entire image. For example, consider an image consists of a tree and a car. The attention layers would aggregate over both the tree and a car. This would lead to an embedding mixing the semantics of both objects, with spatial information lost. On the other hand, if we can just aggregate over the tree region and the car region separately, then the embedding would be spatially dependent.

The key question is how can we extract an embedding for a region, such that the embedding remains aligned with the text embedding \mathbf{e}_{txt} ? We propose to modify Eq. (6) such that the class token only pools information from patches that are inside a region. Just as a class token \mathbf{c} which aggregates over the entire image, we introduce a **region token** $\mathbf{r}^{(l)}$ to aggregate over a region \mathcal{R} . An input region token is initialized from the pre-trained class token, *i.e.*, $\mathbf{r}^{(1)} = \mathbf{c}^{(1)}$. Different from class tokens, we only update $\mathbf{r}^{(l)}$ by aggregating over the patches within \mathcal{R} :

$$\mathbf{r}^{(l+1)} = \left(\sum_{i \in \mathcal{R}} \alpha_i (\mathbf{W}^Q \mathbf{r}^{(l)}, \mathbf{W}^V \mathbf{F}^{(l)}) \cdot \mathbf{W}^V \mathbf{f}_i^{(l)} \right) + \alpha_r \cdot \mathbf{W}^V \mathbf{r}^{(l)}, \quad (7)$$

where \mathcal{R} denotes a set of patch indices covered by the region. See illustration in Fig. 3. Note that the patch tokens and class tokens are updated according to the standard CLIP, *i.e.*, those are unaffected by the introduced region tokens. The operation in Eq. (7) can be implemented as masked self-attention and computed efficiently in parallel.

To extract a $H \times W \times D$ feature map with per-pixel embedding vectors of size D . First, we divide the image into a set of \mathcal{M} non-overlapping spatial regions \mathcal{R}_m using simple linear iterative clustering (SLIC) [37], a classic super-pixel method. Next, we initialize a region token $\mathbf{r}_m^{(1)}$ for each region in \mathcal{M} . Following the aggregation in Eq. (7), we obtain the final embedding vectors $\mathbf{r}_m^{(L)}$, $m = 1, 2, \dots, M$ for each region. Finally, to create a feature map of size $H \times W \times D$, pixel locations with a region is assigned the same final embedding vector, *i.e.*,

$$(\mathbf{E}_{\text{img}})_{ij} = \mathbf{r}_m^{(L)} \quad \forall (i, j) \in \mathcal{R}_m. \quad (8)$$

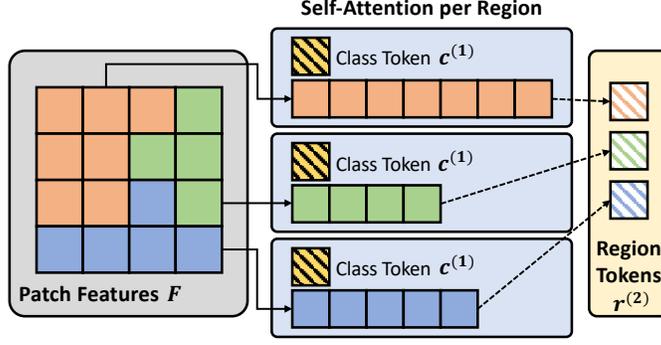


Figure 3: Illustration of the modified attention layer. Here, we illustrate the first attention layer in CLIP’s ViT model. Specifically, we duplicate the class token $c^{(1)}$ as the initialization for region tokens $r^{(1)}$. Attention is performed over the patch features in each of the regions (color-coded) to produce an updated token $r^{(2)}$ for each region.

This approach may be sensitive to number of super-pixels chosen in SLIC. To address this issue, we run SLIC at different “resolution”, *i.e.*, using different numbers of super-pixels to extract multiple feature maps following Eq. (8) and then perform an pixel-wise average.

Increasing Spatial Resolution. In CLIP’s ViT architecture, the patch size is pretty large (32 or 16). As a result, the feature map is extracted at a substantially lower resolution $H' \times W'$. This can be problematic if the super-pixel covers a smaller area than a patch. To address this issue, we propose to modify the first convolution layer of ViT by reducing its stride. This can be viewed as extracting patch features with a sliding window, which leads to an increase in the number of patches, *i.e.*, increase in spatial resolution¹. The remaining attention layers remain unchanged besides processing more patch features. As reviewed in Sec. 3, an attention layer is a “set operation” and the layer can process a larger set without any retraining.

4.4 Spatial Embedding from ResNet

Modifying Attention for Spatial Features. CLIP’s ResNet architecture consists of multiple convolution layers followed by a single attention layer for global pooling. An important difference between the two architectures is how the class token c is trained. In ViT, the class token c is randomly initialized and trained in an end-to-end manner. On the other hand, ResNet’s class token is “defined” to be the average of all patch tokens, *i.e.*,

$$c \triangleq \frac{1}{H'W'} \sum_{i=1}^{H'W'} f_i. \quad (9)$$

Recall, for ResNet architecture, we refer to each vector (per spatial location) from the final convolution layer as a patch token.

Due to this difference, the construction of spatial features (per-patch) is more straight-forward in ResNet. Once again, in a standard attention layer, we can view a class token as aggregating information over all the patches. To obtain spatial features per patch, we need to aggregate over an individual patch. To do so, we introduce a “class token” for each patch $r_m^{(1)} \triangleq f_m$ (analogous to the region token described in Sec. 4.3) and only aggregate over this patch:

$$r_m^{(2)} = \sum_{i \in \{0, m\}} \alpha_i (\mathbf{W}^Q r_m^{(1)}, \mathbf{W}^K [f_0, f_m]) \cdot \mathbf{W}^V f_i. \quad (10)$$

Recall the 0th patch token corresponds to the “class token”, *i.e.*, for a single patch $f_0 \triangleq f_m$. Let’s simplify Eq. (10) by substituting in $r_m^{(1)} = f_m$ and $f_0 = f_m$ to obtain (Note, $\{m, m\}$ denotes a set

¹The positional encoding of these new patches are bilinear interpolation of ones from CLIP.

of duplicate elements):

$$\mathbf{r}_m^{(2)} = \sum_{i \in \{m, m\}} \alpha_i(\mathbf{W}^Q \mathbf{f}_m, \mathbf{W}^K[\mathbf{f}_m, \mathbf{f}_m]) \cdot \mathbf{W}^V \mathbf{f}_i \quad (11)$$

$$= 2 \cdot \alpha_m(\mathbf{W}^Q \mathbf{f}_m, \mathbf{W}^K[\mathbf{f}_m, \mathbf{f}_m]) \cdot \mathbf{W}^V \mathbf{f}_m \quad (12)$$

$$= \mathbf{W}^V \mathbf{f}_m. \quad (13)$$

The last step uses the fact

$$\alpha_m(\mathbf{W}^Q \mathbf{f}_m, \mathbf{W}^K[\mathbf{f}_m, \mathbf{f}_m]) = \frac{\exp((\mathbf{W}^Q \mathbf{f}_m)^\top (\mathbf{W}^K \mathbf{f}_m))}{\sum_{i \in \{m, m\}} \exp((\mathbf{W}^Q \mathbf{f}_i)^\top (\mathbf{W}^K \mathbf{f}_i))} = 0.5. \quad (14)$$

We note that this result in Eq. (13) has also been discovered by a concurrent work, DenseCLIP [35]. Their discovery, however, is based on a hypothesis and empirical validation. Here, we present a mathematical justification of this result. Finally, Eq. (13) enables the extraction of embedding vectors per patch, however, the spatial resolution is low due to large patch sizes.

Increasing Spatial Resolution. To address the low resolution issue when using Eq. (13), we aim to extract a higher resolution feature map without extra training, *e.g.*, distillation in DenseCLIP [35]. The main cause of low-resolution feature map is the down-sampling/pooling layers in the Res-Net architecture which reduce the spatial resolution.

Our method is inspired by Chen *et al.* [38], where they remove the last few max pooling layers from a pre-trained and use atrous (dilated) convolution to more densely sample the output. Following this idea, we replace all convolution layers in CLIP’s image encoder with dilated convolution. Specifically, to get a high-resolution feature map, we increase the dilation factor by a multiple of two for every stride-two down-sampling/pooling in the original model. Pre-trained weights of the model is unmodified.

4.5 Practical Considerations

In Eq. (3), we formulated bounding box prediction as a search over all possible bounding boxes. Naively, using brute force search is infeasible. Existing branch and bound methods, *e.g.*, efficient sub-window search [39] is directly applicable. To further speedup this search, we devise a (greedy) hierarchical search strategy. This involves iteratively searching on a downsampled score map to find a coarse bounding box, then searching again by zooming in to this coarse box. At each iteration, a brute-force search on all, low-resolution, bounding boxes based on integral images is computed using a GPU, which gives significant speedup. While this strategy does not guarantee the optimal solution, empirically we found it to perform well.

5 Experiments

5.1 Experimental Setup and Details

Datasets. We follow ZSGNet’s [19] zero-shot setup on Flickr30k [3] and Visual Genome [23] dataset for evaluation. Note that we only use these datasets for evaluation (and tuning two hyper-parameters) as our models do not require any training on phrase localization annotations. In more details, ZSGNet proposes four **zero-shot** splits based on Flickr30k and Visual Genome:

- **Flickr-S0.** Phrases in the test set are not seen in the training set. But this split does not exclude the possibility that there are phrases in the training set describe objects in the same **object category**. For example, if “man” is in the training set then “woman” is allowed to be in the test set, even though both of these refer to a broader category of “people”. As a result, this split is zero-shot for phrases, but **not for object categories**. This split corresponds to the Case 0 in ZSGNet [19].
- **Flickr-S1.** Phrases in the test set are not seen in the training set **and** there are no phrases in the training set belonging to the same object category of any text phrase. Flickr30k has several common object categories (*e.g.* “people”, “animals”) and one “other” category. ZSGNet uses all the examples in “other” as validation and test sets, and examples of the remaining categories are used in the training set. This split corresponds to the Case 1 in ZSGNet [19].

Table 1: Quantitative Results on various splits of Flickr and Visual Genome datasets. Here, we report the $\text{Acc}@thr$ metric for each of the methods (the higher the better). The CLIP’s architecture is in parentheses; RN50 corresponds to a ResNet architecture and ViT-L/14 corresponds to biggest Vit architecture released by CLIP.

Method	Supervision	Flickr		Zero-Shot Flickr		Zero-Shot VG	
		All	Other	S0	S1	S0	S1
ZSGNet [19]	Full	63.39	45.53	43.02	31.23	19.95	20.77
Wang <i>et al.</i> [13]	BBox	50.49	34.71	-	-	-	-
Parcalabescu <i>et al.</i> [15]	BBox	57.08	38.52	-	-	-	-
Crop & Rank (ViT-L/14)	CLIP	4.21	6.34	7.17	7.51	17.24	16.97
DenseCLIP (RN50) [35]	CLIP	34.26	25.87	31.08	27.78	14.71	15.68
Ours (RN50)	CLIP	34.25	25.93	31.11	27.30	19.49	20.72
Ours (ViT-L/14)	CLIP	43.80	35.42	40.37	36.10	24.47	25.50

- **VG-S0.** In this split of VG, phrases in the training and test sets are from different synsets [40]. Also, no test images contain objects in the training synsets. This split corresponds to the Case 2 in ZSGNet [19].
- **VG-S1.** This is similar to VG-S0, except that each test image contains, in addition to the object to which the phrase refers, an object belonging to the training synsets. This split corresponds to Case 3 in ZSGNet [19].

As not all prior works use these dataset splits, we also report results on standard, *i.e.*, non zero-shot, split of the Flickr30k dataset:

- **Flickr-All.** The original Flickr30k train-val-test split.
- **Flickr-Other.** The original Flickr30k subset split only containing objects belonging to the “other” category. Note that this not the same as Flickr-S1 because ZSGNet reconstruct the splits over the entire train, validation, test splits.

Baselines. First, we compare against a straight-forward baseline based on CLIP, which we named, “Crop & Rank”. For an image-query pair, we first use a traditional bounding box proposal method (selective search [41]) to get the top 200 bounding box proposals. For each proposal, we crop the corresponding image region and resize it to 224×224 . Next, we use CLIP to extract image embedding for each bounding box and the query’s text embedding. The model predicts the bounding box where its cosine similarity with the text embedding is the highest out of all proposals.

Next, we also compared with DenseCLIP [35] by using their method to extract spatial features. Different from our ResNet method, the resolution of the extract feature map is lower. The bounding box extraction procedure is identical to ours. Additionally, we also report performance on the following prior works:

- ZSGNet [19] is a fully supervised method that reports on all of the standard and zero-shot splits.
- Wang *et al.* [13] is a method that assembles off-the-shelf strongly supervised object detectors and word embedding models into a textual grounding model without extra training.
- Parcalabescu *et al.* [15] is an improved model of Wang *et al.* [13] by using additional structure information on image and text.

For completeness, we report performance on fully supervised methods and methods using fully supervised pre-trained object detectors. However, this is not a fair comparison. Specifically, Wang *et al.* [13] and Parcalabescu *et al.* [15] use supervised object detectors which naturally biased them to perform better on object categories that overlaps with the phrase localization dataset.

Evaluation Metric. We adopt the widely used $\text{Acc}@thr$ metric for evaluation, where *thr* refers to the threshold of intersection over union (IoU). A predicted bounding box is considered correct if its IoU with the ground-truth box is greater than the threshold. $\text{Acc}@thr$ reports the accuracy over the dataset, *i.e.*, the fraction of correct predictions over a dataset. We follow ZSGNet [19] to use a threshold of 0.5 for Flickr30k and 0.3 for VG.

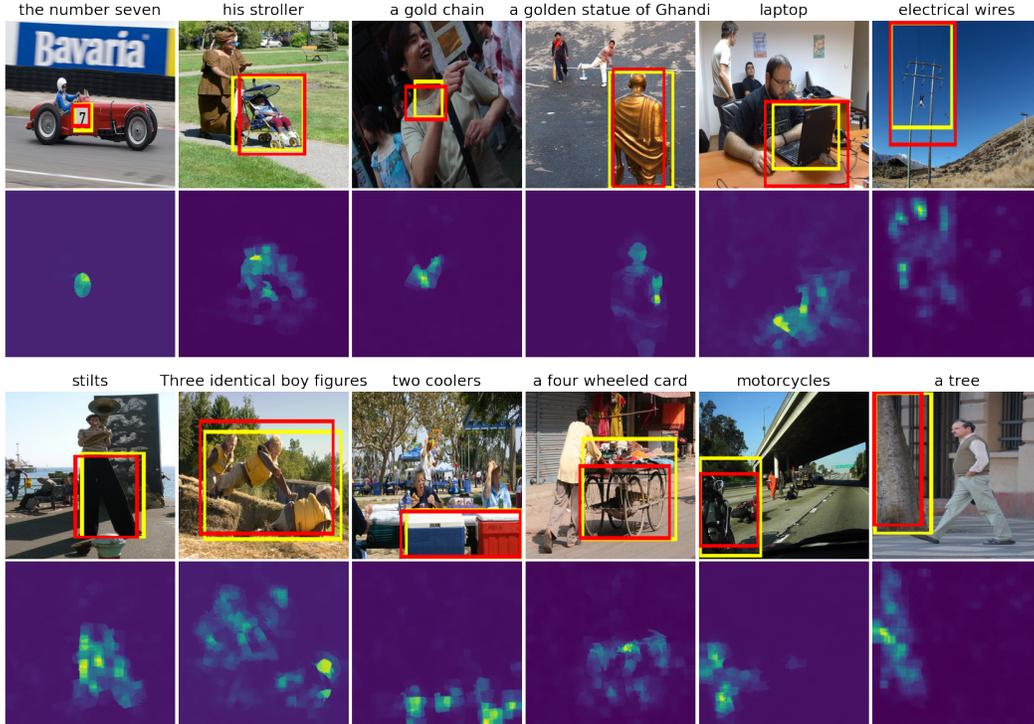


Figure 4: Qualitative results from our method. The query text is shown at the top. We visualize ground-truth bounding boxes in yellow and our predicted boxes in red. For each image, we also visualize the score map extracted from our method using CLIP’s ViT-L14 architecture.

Implementation Details. We use the latest-released CLIP model (ViT with patch size 14 and a $4\times$ scaled ResNet50) with our method². We obtain regions by pooling SLIC maps with 100, 150, ..., 600 superpixels. The two parameters that have the biggest impact on the performance is the weighting factor λ in the bounding box score Eq. (4) and the temperature term σ in Eq. (5). We tune these hyperparameters for different variants of our method on the validation sets of Flickr30k and VG respectively.

5.2 Quantitative Results

We report quantitative comparisons with the baselines in Tab. 1.

Zero-Shot Results. Without any supervision from phrase localization datasets, our method outperforms all compared baselines using CLIP supervision by a large margin. Additionally, our method even outperforms the fully supervised methods (ZSGNet [19]) on Flickr-S1, VG-S1, VG-S2. We hypothesize that ZSGNet is negatively impacted by bias introduced by the training data distribution and have limited generalization capabilities.

For Flickr-S0, our model achieves comparable performance to ZSGNet. We emphasize that Flickr-S0 is **not a truly zero-shot split**, because the same object category can appear in both the train and test set, while only the phrases are “unseen”. This is more favorable to supervised methods because they can learn the visual features from the training data to better localize objects.

Non Zero-Shot Results. For non-zero shot splits (Flickr-All and Flickr-Other), not surprisingly, supervised methods give the best performance. Our method achieves comparable performance to prior work that used supervised object detectors [13, 15]. Specifically, observe that the performance of ZSGNet on Flickr-All is much better than that of Flickr-Other. This again shows the bias introduced by the dominance of common object categories. On the other hand, the gap between ZSGNet and

²Pre-trained models are available at <https://github.com/openai/CLIP>

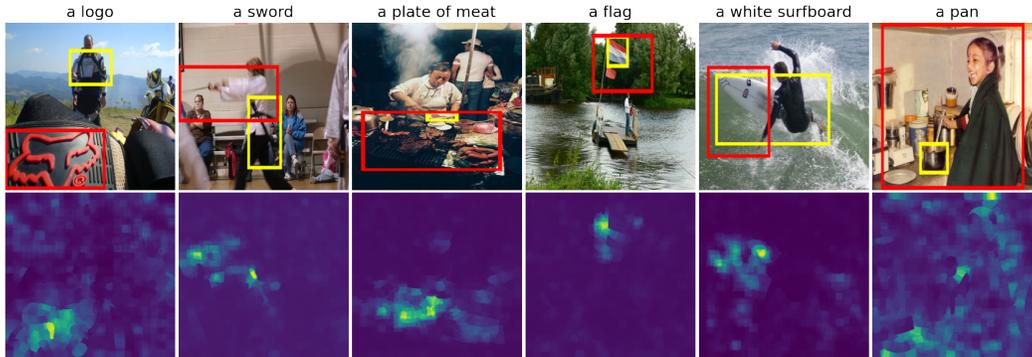


Figure 5: Failure cases of our method. Ground-truth in yellow and predicted box in red.

Table 2: Ablation study comparing different backbone architectures. We observe that the phrase localization performance improves as the CLIP model size increases.

Architecture	Flickr		Zero-Shot Flickr		Zero-Shot VG	
	All	Other	S0	S1	S0	S1
RN50	34.30	25.84	31.07	27.28	17.51	18.28
RN50×4	34.25	25.93	31.11	27.30	19.49	20.72
ViT-B/32	29.49	24.72	27.48	25.00	13.43	14.21
ViT-B/16	41.54	35.29	37.82	35.70	18.23	19.19
ViT-L/14	43.80	35.42	40.37	36.10	24.47	25.50

ours on Flickr-All is larger than on Flickr-Other. This shows that our method, which is not fine-tuned on any textual grounding dataset, is less prone to bias in object categories.

5.3 Qualitative Results

Beyond the quantitative results, we further study the quality of the extracted score map. In Fig. 4, we visualize some examples of our localization method on the Flickr30k dataset, including the score maps and the predicted bounding boxes. Ground-truth bounding boxes are shown in yellow and our predicted bounding boxes are shown in red. We observe that our method can correctly localize uncommon phrases such as “stilts” and “electrical wires”. This demonstrates our method’s ability to handle all kinds of open vocabulary/world objects, including those rarely seen in fully-annotated object detection datasets.

We also observe typical failure cases of our method, visualized in Fig. 5. First, annotations from the dataset (Flickr30k) may be ambiguous (see col. 1 and 2) Second, our method tends to have challenges in localizing small objects (see col. 3 and 4). Specifically, for the image in column 4, the heatmap correctly identified the flag, however, the bounding box extraction procedure did not predict a small bounding box. Third, we observe that our method may have challenges when handling occlusions (see col. 5), the score map only partially identified the surfboard. Finally, there are examples that the score map is not interpretable (see col. 6) and more investigation and understanding of CLIP embedding is necessary.

5.4 Ablation Study

We conduct ablation studies on two aspects of our method to show their effects on the phrase localization performance: (a) how important is the quality of CLIP model; (b) how important is the spatial resolution of feature maps.

CLIP Backbones. In Tab. 2, we show qualitative results for our method with different pre-trained CLIP architectures. It can be seen that larger models consistently outperform smaller models. This is most obvious for the ViT architecture. With smaller patch size, ViT patch tokens can capture

Table 3: Ablation study on spatial resolution of feature maps.

Architecture	Flickr S0		Flickr S1		VG S0		VG S1	
	1×	2×	1×	2×	1×	2×	1×	2×
ViT-B/32	27.05	27.48	23.80	25.00	12.79	13.43	13.46	14.21
ViT-B/16	34.10	37.82	29.26	35.70	13.51	18.23	13.87	19.19
ViT-L/14	34.67	40.37	30.30	36.10	18.92	24.47	19.54	25.50

more delicate spatial details of the scene, and therefore can better compute the features for some given region, even if it is small. Additionally, we observe that the best ViT CLIP model outperforms ResNet, validating our contribution over DenseCLIP [35] which is restricted to ResNet architectures.

Feature Map Resolution. In Tab. 3, we report an ablation study on the spatial resolution of feature maps. We report results with and without the upsampling technique for ViT. We consistently observe that upsampling to higher resolution (2×) feature maps results in better phrase localization performance.

6 Conclusion

We present a method for visual grounding by leveraging pre-trained language-vision model, CLIP, without any extra training. Our method generates high-resolution pixel-wise feature vectors from ViT and ResNet architectures of CLIP and computes per-pixel similarity scores with the embedding of the text query to create a score map. We then search over the score map to find the best box that localizes the object. Experiments on two datasets show that our method can effectively localize a wide variety of objects with complex text queries in the zero-shot setting. This method is a first step towards visual grounding based on large-scale pre-trained language-vision models without extra training, reducing potential biases caused by the limited size of human annotated datasets.

References

- [1] L. Wang, Y. Li, and S. Lazebnik, “Learning deep structure-preserving image-text embeddings,” in *CVPR*, 2016.
- [2] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell, “Natural language object retrieval,” in *Proc. CVPR*, 2016.
- [3] B. A. Plummer, A. Mallya, C. M. Cervantes, J. Hockenmaier, and S. Lazebnik, “Phrase localization and visual relationship detection with comprehensive image-language cues,” in *Proc. ICCV*, 2017.
- [4] L. Wang, Y. Li, J. Huang, and S. Lazebnik, “Learning two-branch neural networks for image-text matching tasks,” *IEEE TPAMI*, 2018.
- [5] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach, “Multimodal compact bilinear pooling for visual question answering and visual grounding,” in *EMNLP*, 2016.
- [6] R. A. Yeh, J. Xiong, W.-M. Hwu, M. Do, and A. G. Schwing, “Interpretable and globally optimal prediction for textual grounding using image concepts,” in *Proc. NeurIPS*, 2017.
- [7] B. A. Plummer, P. Kordas, M. H. Kiapour, S. Zheng, R. Piramuthu, and S. Lazebnik, “Conditional image-text embedding networks,” in *Proc. ECCV*, 2018.
- [8] S. Yang, G. Li, and Y. Yu, “Dynamic graph attention for referring expression comprehension,” in *Proc. ICCV*, 2019.
- [9] Z. Yang, T. Chen, L. Wang, and J. Luo, “Improving one-stage visual grounding by recursive sub-query construction,” in *Proc. ECCV*, 2020.
- [10] Z. Mu, S. Tang, J. Tan, Q. Yu, and Y. Zhuang, “Disentangled motif-aware graph learning for phrase grounding,” in *Proc. AAAI*, 2021.
- [11] A. Rohrbach, M. Rohrbach, R. Hu, T. Darrell, and B. Schiele, “Grounding of textual phrases in images by reconstruction,” in *ECCV*, 2016.

- [12] R. A. Yeh, M. N. Do, and A. G. Schwing, “Unsupervised textual grounding: Linking words to image concepts,” in *Proc. CVPR*, 2018.
- [13] J. Wang and L. Specia, “Phrase localization without paired training examples,” in *Proc. ICCV*, 2019.
- [14] L. Wang, J. Huang, Y. Li, K. Xu, Z. Yang, and D. Yu, “Improving weakly supervised visual grounding by contrastive knowledge distillation,” in *Proc. CVPR*, 2021.
- [15] L. Parcalabescu and A. Frank, “Exploring phrase grounding without training: Contextualisation and extension to text-based image retrieval,” in *Proc. CVPRW*, 2020.
- [16] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Proc. CVPR*, 2009.
- [17] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The PASCAL visual object classes (VOC) challenge,” *IJCV*, 2010.
- [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Proc. ECCV*, 2014.
- [19] A. Sadhu, K. Chen, and R. Nevatia, “Zero-shot grounding of objects from natural language queries,” in *Proc. ICCV*, 2019.
- [20] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, “Learning transferable visual models from natural language supervision,” in *Proc. ICML*.
- [21] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, “Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models,” in *Proc. ICCV*, 2015.
- [22] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, “From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions,” *TACL*, 2014.
- [23] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, *et al.*, “Visual genome: Connecting language and vision using crowdsourced dense image annotations,” *IJCV*, 2017.
- [24] S. Kazemzadeh, V. Ordonez, M. Matten, and T. L. Berg, “ReferItGame: Referring to objects in photographs of natural scenes,” in *Proc. EMNLP*, 2014.
- [25] L. Wang, Y. Li, and S. Lazebnik, “Learning two-branch neural networks for image-text matching tasks,” *arXiv preprint arXiv:1704.03470*, 2017.
- [26] R. Hu, M. Rohrbach, J. Andreas, T. Darrell, and K. Saenko, “Modeling relationships in referential expressions with compositional modular networks,” in *Proc. CVPR*, 2017.
- [27] K. Endo, M. Aono, E. Nichols, and K. Funakoshi, “An attention-based regression model for grounding textual phrases in images,” in *Proc. IJCAI*, 2017.
- [28] K. Chen*, R. Kovvuri*, and R. Nevatia, “Query-guided regression network with context policy for phrase grounding,” in *Proc. ICCV*, 2017. * equal contribution.
- [29] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, “Scaling up visual and vision-language representation learning with noisy text supervision,” in *Proc. ICML*, 2021.
- [30] X. Liu, D. H. Park, S. Azadi, G. Zhang, A. Chopikyan, Y. Hu, H. Shi, A. Rohrbach, and T. Darrell, “More control for free! image synthesis with semantic diffusion guidance,” *arXiv preprint arXiv:2112.05744*, 2021.
- [31] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, “Open-vocabulary object detection via vision and language knowledge distillation,” *arXiv preprint arXiv:2104.13921*, 2021.
- [32] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, “Zero-shot detection via vision and language knowledge distillation,” *arXiv e-prints*, pp. arXiv–2104, 2021.
- [33] A. Zareian, K. D. Rosa, D. H. Hu, and S.-F. Chang, “Open-vocabulary object detection using captions,” in *Proc. CVPR*, 2021.
- [34] G. Ghiasi, X. Gu, Y. Cui, and T.-Y. Lin, “Open-vocabulary image segmentation,” *arXiv preprint arXiv:2112.12143*, 2021.

- [35] C. Zhou, C. C. Loy, and B. Dai, “DenseCLIP: Extract free dense labels from clip,” *arXiv preprint arXiv:2112.01071*, 2021.
- [36] M. Xu, Z. Zhang, F. Wei, Y. Lin, Y. Cao, H. Hu, and X. Bai, “A simple baseline for zero-shot semantic segmentation with pre-trained vision-language model,” *arXiv preprint arXiv:2112.14757*, 2021.
- [37] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, “Slic superpixels compared to state-of-the-art superpixel methods,” *IEEE TPAMI*, 2012.
- [38] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *TPAMI*, 2017.
- [39] C. H. Lampert, M. B. Blaschko, and T. Hofmann, “Efficient Subwindow Search: A Branch and Bound Framework for Object Localization,” *PAMI*, 2009.
- [40] G. A. Miller, “WordNet: a lexical database for English,” *Communications of the ACM*, 1995.
- [41] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders, “Selective search for object recognition,” *IJCV*, 2013.

Appendix

A Additional Results

Architecture comparison. In Fig. A1, we show comparisons of heatmaps extracted from different CLIP architectures. In general, we observe ViT architectures to perform better than the ResNet architecture, with ViT-L/14 performing the best.

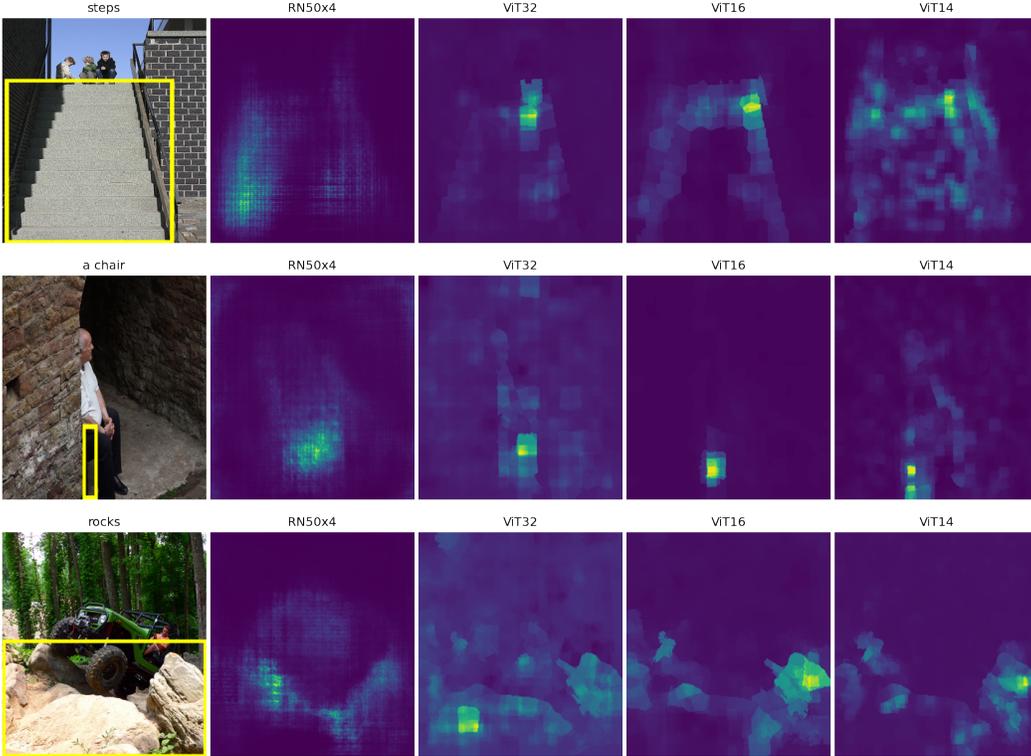


Figure A1: Heatmaps extracted from different architectures

Resolution comparison. In Fig. A2, we visualize our approach with and without using our technique for increasing the spatial resolution. For both ResNet and ViT architectures, we observe better localization of the object when using a higher resolution heatmap.

Comparisons on bounding box prediction. In Sec. 4.5, we describe a (greedy) hierarchical search strategy to find the best bounding box. While this greedy approach does not guarantee finding the best bounding box, empirically, we found it to perform similar to efficient sub-window search (ESS) [39]. In Fig. A3, we observe that ESS and the hierarchical search results in nearly identical box predictions. In Tab. A1, we report the running time (mean and standard deviation over 100 examples) for the hierarchical search and ESS. We observe that the hierarchical search results in faster running time. Note the ESS implementation is in Python and runs only on CPU.

Table A1: Comparison of bounding box search time

Method	Hardware	Mean Time (s)	Standard Deviation (s)
Hierarchical	RTX6000	2.3×10^{-3}	2.6×10^{-5}
ESS (Python)	CPU	3.1×10^{-1}	4.1×10^{-1}

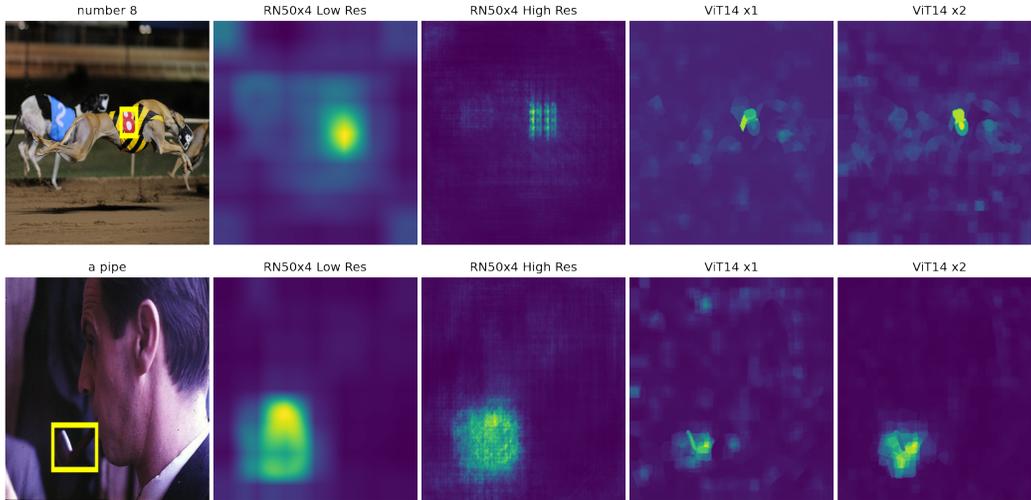


Figure A2: Heatmaps with different resolutions

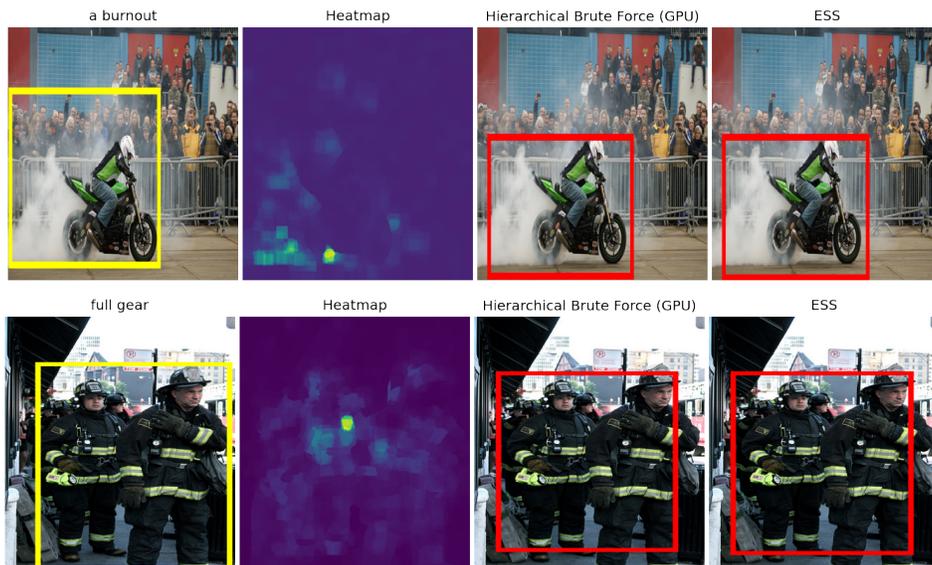


Figure A3: Qualitative results with different box search methods

Additional Qualitative Result. We provide more qualitative results in Fig. A4 and Fig. A5. We observe that our approach successfully localizes a diverse set of phrases, *e.g.*, “vending machine”, “water spray”, “rainbow flags”, *etc.*

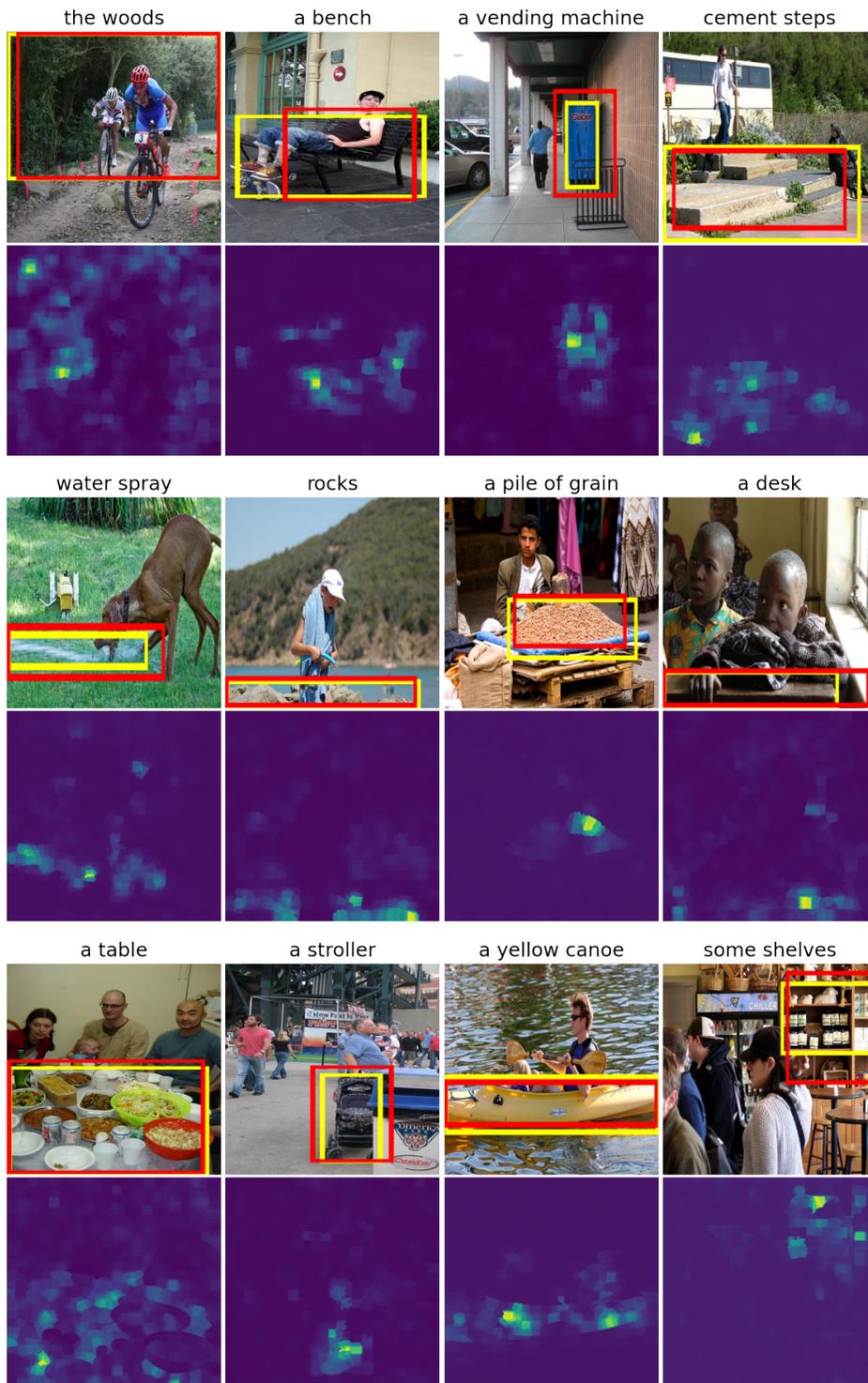


Figure A4: Additional qualitative results.

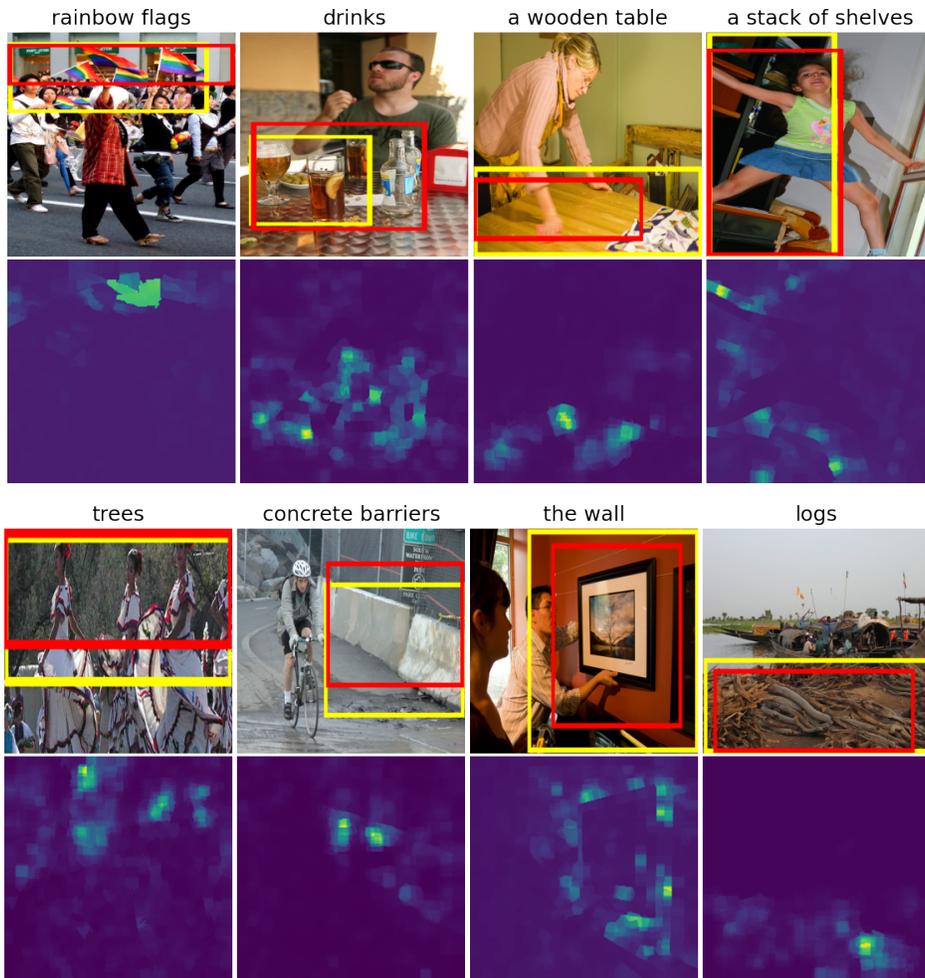


Figure A5: Additional qualitative results.