# Your Large Vision-Language Model Only Needs
# A Few Attention Heads For Visual Grounding

Seil Kang    Jinyeong Kim    Junhyeok Kim    Seong Jae Hwang
Yonsei University
{seil, jinyeong1324, timespt, seongjae}@yonsei.ac.kr

## Abstract

*Visual grounding seeks to localize the image region corresponding to a free-form text description. Recently, the strong multimodal capabilities of Large Vision-Language Models (LVLMs) have driven substantial improvements in visual grounding, though they inevitably require fine-tuning and additional model components to explicitly generate bounding boxes or segmentation masks. However, we discover that a few attention heads in frozen LVLMs demonstrate strong visual grounding capabilities. We refer to these heads, which consistently capture object locations related to text semantics, as localization heads. Using localization heads, we introduce a straightforward and effective training-free visual grounding framework that utilizes text-to-image attention maps from localization heads to identify the target objects. Surprisingly, only three out of thousands of attention heads are sufficient to achieve competitive localization performance compared to existing LVLM-based visual grounding methods that require fine-tuning. Our findings suggest that LVLMs can innately ground objects based on a deep comprehension of the text-image relationship, as they implicitly focus on relevant image regions to generate informative text outputs. All the source codes will be made available to the public.*

## 1. Introduction

Visual grounding is a task that, given textual descriptions, identifies and localizes relevant objects within an image, producing outputs such as bounding boxes [39, 70] or segmentation masks [18]. Recently, this vision-language task, which inherently requires a deep understanding of the relationship between images and text, has seen significant advancements with the emergence of powerful Large Vision-Language Models (LVLMs) [28, 35, 36, 56]. However, since LVLMs are primarily designed to generate text outputs, directly leveraging them as a vision-language tool to identify and localize objects within an image (*i.e.*, visual grounding) presents technical challenges. Inevitably, cur-
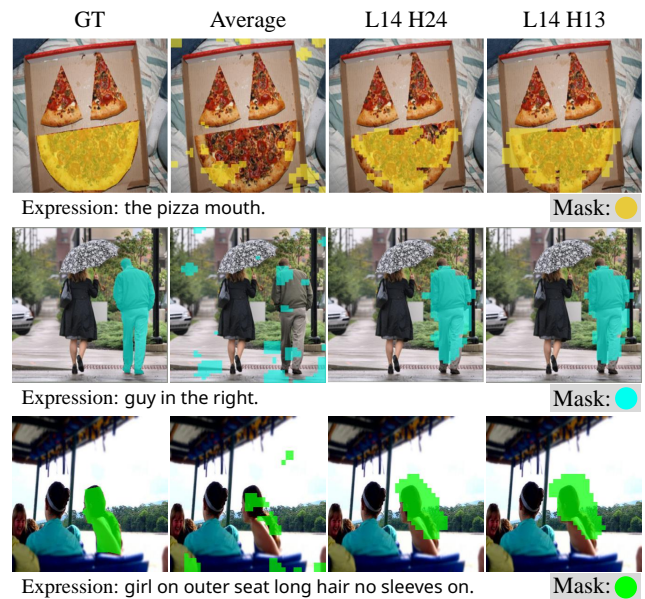


Figure 1. Visualization of the text-to-image attention maps from LLaVA-1.5-7B [35]. While the average attention map initially seems uninformative for localization, a closer examination reveals that LVLM possesses built-in *localization heads* that consistently capture key areas of an image corresponding to the referring text, regardless of sample variations. L14 H24 refers to the 24th attention head in the 14th layer of the LVLM.

rent LVLM-based visual grounding methods require explicit fine-tuning of LVLMs with additional visual grounding datasets and modifications to model components to enable the generation of bounding boxes [6, 61, 68] or segmentation masks [26, 45, 65, 74].

Despite the interesting integration of LVLMs in previous visual grounding works, a fundamental question remains: *since LVLMs generate text outputs that imply an understanding of specific image regions, is it possible to explicitly observe this mechanism in action?* In other words, we ask whether we can extract how the LVLMs "focus" on specific image regions corresponding to given text descriptions for visual grounding. A natural first approach to addressing this question might be to examine the text-to-image atten-

tion maps, which reveal how a text description attends to different image patches. To explore this, we visualize the average attention maps of LVLMs across various layers and heads—a common method in ViTs [11, 73] and diffusion models (DMs) [4, 17, 54]—anticipating that they would capture the regions associated with the referring text. However, unlike the interpretable attention patterns observed in ViTs and DMs, the text-to-image attention maps in LVLMs appear sparse and contain significant noise, as illustrated in the second column of Fig. 1. This suggests that the current use of LVLM attention maps may struggle to accurately pinpoint relevant objects for visual grounding.

However, interestingly, our work reveals that not the average of the attention maps, but *some* small subset of attention heads are capable of providing tangible and precise text-image attention maps. In particular, we find that a few attention heads in LVLMs consistently capture regions in images corresponding to the referred text, regardless of the samples. We refer to these heads as *localization heads*. For example, as presented in the third and fourth columns of Fig. 1, the attention maps of the 24th head of the 14th layer (L14 H24) and the 13th head of the 14th layer (L14 H13) in LLaVA-1.5-7B [35] consistently highlight the regions of interest based on the referred text.

In this work, we introduce how we systematically identify such localization heads based on two explicit criteria. (1) We measure how much each attention head focuses on the image by calculating the *attention sum* and only select the heads that dominantly attend to the image. (2) Among these heads, the ones that specifically pay attention to a certain region of the image, which is measured by *spatial entropy* [2], are considered to effectively localize the referred object. We validate that the selected localization heads consistently capture objects closely associated with the text.

With our localization heads, we introduce a simple yet effective training-free visual grounding framework. The attention maps from the localization heads are assembled to predict the bounding box or mask of the referred object. Notably, only three localization heads are enough to localize the referred object within the image, suggesting that they are highly specialized to attend to relevant image regions. As shown in Fig. 2, in contrast to existing fine-tuning based methods, our framework is training-free, eliminating the need for additional fine-tuning LVLMs for visual grounding tasks.

We validate our approach across ten different LVLMs with varying parameter counts, architectures, and training datasets, demonstrating its broad applicability. Our framework outperforms the existing training-free methods by significant margins. Furthermore, our method performs comparably to specially fine-tuned LVLMs for visual grounding tasks (*e.g.*, LISA [26]). The results indicate that LVLMs can serve as effective text-referring localizers, intrinsically
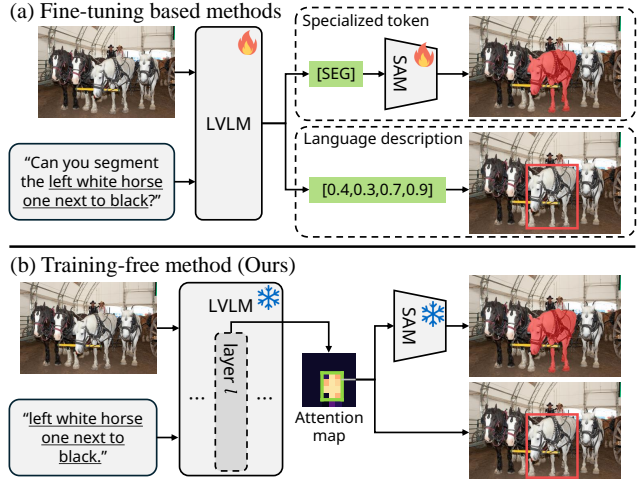


Figure 2. Comparison of LVLM frameworks for visual grounding. (a) Existing methods generally fine-tune a LVLM to leverage specialized tokens (*e.g.*, [SEG]) or language descriptions for visual grounding. (b) Our framework utilizes the attention maps of only a few localization heads from frozen LVLMs.

identifying regions that are relevant and coherent with the text expression. To the best of our knowledge, we are the first to identify the localization properties of specific attention heads in LVLMs.

In summary, our contributions are as follows:
- We discover that the specific attention heads in LVLMs have the capability for visual grounding, which we refer to as *localization heads*.
- We propose a simple yet effective framework for LVLM-based training-free visual grounding with localization heads. The attention maps from a few localization heads are utilized to predict the bounding box or mask of the referred object.
- We evaluate our approach across various LVLMs. Our framework demonstrates superior performance by a large margin compared to other training-free methods and even performs comparably to fine-tuned methods.

## 2. Related Works

**Visual Grounding.** Visual grounding aims to identify the region in the image based on a free-form natural language expression [5], which has expanded the scope of detection and segmentation tasks to a more realistic scenario [50, 66]. Two prominent tasks within visual grounding are Referring Expression Comprehension (REC) [39, 70] and Referring Expression Segmentation (RES) [18, 34]. REC focuses on localizing a referred object in an image and generating a bounding box, while RES further requires a pixel-level segmentation mask. In order to address these tasks, numerous studies have been conducted to explore effective methods that consider both text and visual information simultaneously [23, 29, 33, 42, 49, 62, 67, 75].

**Application of LVLMs in Grounding Tasks.** Recently, visual grounding has been significantly advanced by leveraging the outstanding vision-language processing capabilities of LVLMs. To incorporate LVLMs into visual grounding tasks, existing methods include visual grounding datasets in the training process and implement additional components to extract localization information. For example, LISA [26] introduces [SEG] token as a mask embedding and generates a segmentation mask using additional mask decoder [24]. F-LMM [64] leverages the attention weights of frozen LVLMs, but still requires training its mask refinement modules on visual grounding datasets. In contrast, we propose a training-free visual grounding method that directly utilizes LVLMs.

**Training-Free Visual Grounding.** Given the high performance of multimodal foundation models across diverse vision-language tasks, training-free visual grounding emerges as a new research direction. Existing training-free methods typically apply internal features or attention maps from CLIP [44] or Text-to-Image Diffusion Models (DMs) [47]. CLIP-based methods typically employ off-the-shelf models [24, 46] to generate region proposals and select the most relevant bounding box [52] or mask [53, 71] based on the CLIP similarity score with the text query. On the other hand, DM-based methods utilize the residue of the text-to-image diffusion process (*e.g.*, the attention map) to predict the segmentation mask [3, 40]. Our work advances this line of research by introducing the first LVLM-based training-free visual grounding framework.

## 3. Background

**Notation.** Large Vision-Language Models (LVLMs) typically consist of three main components: a vision encoder, a projector, and a large language model. For an input image $X_v$, the vision encoder and the projector transform the image into a sequence of visual embedding $Z_v \in \mathbb{R}^{P^2 \times d}$, where $P^2$ is the number of flattened image tokens and $d$ is the hidden dimension. Similarly, an input text $X_t$ is converted into a sequence of token embeddings $Z_t \in \mathbb{R}^{L \times d}$, where $L$ is the number of tokens in the text. The visual and textual embeddings are concatenated as $Z^0 = [Z_v; Z_t] \in \mathbb{R}^{(P^2+L) \times d}$ and fed into the large language model (LLM) as the input embeddings.

**Multi-Head Self-Attention.** The input embeddings $Z^0$ pass through a series of decoder blocks, which consists of multi-head self-attention and feed-forward neural network module. Specifically, we focus on the attention heads, as these are the only components where tokens interact. In layer $\ell$ and head $h$, the hidden state from the previous layer $Z^{\ell-1}$ is projected into query $Q$, key $K$, and value $V \in \mathbb{R}^{(P^2+L) \times d_h}$ matrices, where $d_h$ is the hidden dimension of the attention head. Then, the attention head com-

putes the attention weights as:

$$\text{Attn}^{\ell,h}(Z^{\ell-1}) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right). \qquad (1)$$

Note that the attention weights reflect the similarity between the query $Q$ and key $K$ matrices.

**Investigation of Image-Text Interaction.** Considering that LLM decoding operates in an auto-regressive manner, information flows from preceding tokens to subsequent ones, resulting in the final token to encapsulate the context of the entire sentence [20, 59]. Thus, we posit that the query vector of the last input text token $q_\text{txt}$ serves as a representative query for the whole sentence. For example, in the sentence "`the pizza mouth.`" in Fig. 1, the query vector of the last token [.] is utilized in our experiments. To investigate image-text interactions, we examine the attention weights of where the query is $q_\text{txt}$ and keys are image tokens. Specifically, considering a slight modification of Eq. (1), for the attention weights $a^{\ell,h}$ at layer $\ell$ and head $h$ with $q_\text{txt}$ as a query token:

$$a^{\ell,h} = \text{softmax}\left(\frac{q_\text{txt}K^\top}{\sqrt{d_h}}\right) \in \mathbb{R}^{P^2+L}, \qquad (2)$$

we focus on the first $P^2$ components, $a^{\ell,h}[1:P^2]$, for our analysis. In the following sections of this paper, this will also be denoted as L$\ell$ H$h$ for simplicity. For example, L5 H3 refers to the third attention head in the fifth layer of the LVLM.

## 4. Towards Discovering Localization Heads

Recent studies [9, 57, 76] have shown that the attention heads exhibit distinct characteristics, motivating us to find specific heads possessing the potential to serve as effective *text referring localizers*. In this section, we propose attention sum and spatial entropy in Sec. 4.1 as two criteria for selecting such heads. Through experiments in Sec. 4.2, we validate that the heads capturing objects corresponding to the text description can be successfully identified based on the proposed criteria. Note that the first two layers of the LLM are consistently excluded in our analyses, as the early layers are known to operate differently from the other layers [25]. To demonstrate the generalizability of our findings, we conduct experiments across various LVLMs [7, 8, 30, 35, 36, 38, 69] and datasets [18, 22]. Details of the experimental setup and more results are provided in the Appendix Sec. A. and C., respectively.

### 4.1. Criteria to Find Localization Heads

Our final goal is to identify heads that excel in text referring. To achieve this, we propose two criteria in this section.

**Criterion 1: Attention Sum.** To identify heads that predominantly focus on the overall image, we first consider
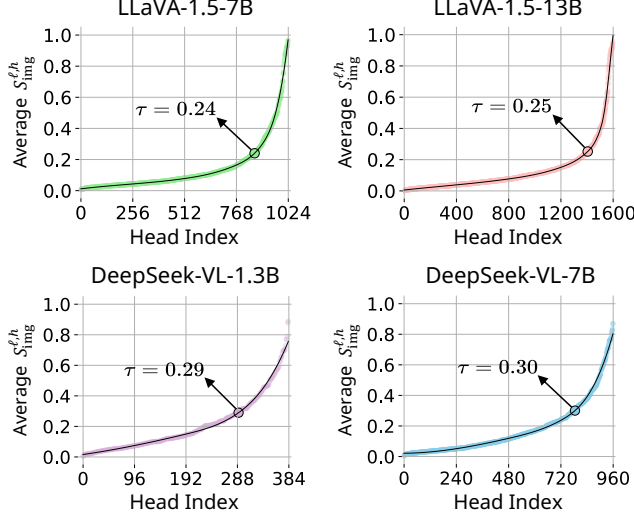
Figure 3. Average $S_{\text{img}}^{\ell,h}$ values for each attention head. We sort the heads in ascending order of $S_{\text{img}}^{\ell,h}$. Attention heads with $S_{\text{img}}^{\ell,h} \geq \tau$ are considered to effectively attend to the image, where $\tau$ is the threshold determined by the maximum curvature in the graph.

attention sum $S_{\text{img}}^{\ell,h} = \sum_{i=1}^{P^2} \boldsymbol{a}^{\ell,h}[i]$, which quantifies the relevance of image information to $\boldsymbol{q}_{\text{txt}}$ within individual attention heads. Then, the average $S_{\text{img}}^{\ell,h}$ for each head is computed across 1,000 random samples from RefCOCO [22] training set.

As shown in Fig. 3, most attention heads exhibit low $S_{\text{img}}^{\ell,h}$ values, indicating that relatively few heads contribute significantly to the model's text-image interaction. To distinguish heads with high $S_{\text{img}}^{\ell,h}$ from those with low values, we set the threshold $\tau$ at the point of the maximum curvature in the graph (e.g., $\tau = 0.24$ in LLaVA-1.5-7B [35]). We deem the heads with $S_{\text{img}}^{\ell,h} \geq \tau$ to effectively attend to image. While we adopt the maximum curvature as a practical choice, we note that our analysis remains robust across a range of reasonable $\tau$ values. For analyses using alternative $\tau$ values, please refer to Appendix Sec. C.

**Criterion 2: Spatial Entropy.** For an attention head to be considered effective at focusing on objects, it must not only have a high attention sum value for the image but also concentrate its attention specifically around the objects. Since it is reasonable to assume that the object patches tend to stay near each other [51, 58, 72], we evaluate how locally a cluster is formed in each attention map through spatial entropy [2, 41] to identify localization heads.

Fig. 4 presents an example of how spatial entropy is calculated. First, we reshape the attention weights $\boldsymbol{a}^{\ell,h}[1:P^2]$ into a $P \times P$ attention map $\boldsymbol{A}^{\ell,h}$. The attention map is binarized by assigning a value of 1 to elements above the mean and 0 to those below it [41]. Next, we identify connected components $C_i$ [14], defined as a set of coordinates connected via 8-neighbors. Then, for the set of $N$ connected components $\{C_i\}_{i=1}^{N}$, the spatial entropy $H$ is calculated
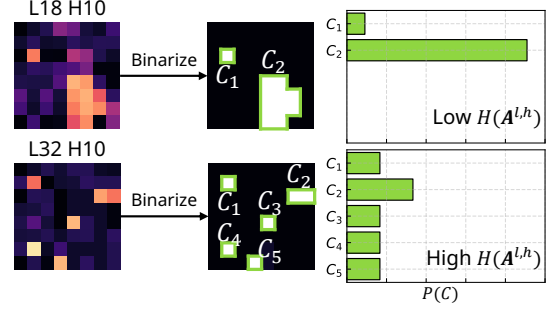


Figure 4. Illustration of the process for calculating spatial entropy. The attention map is binarized, and the spatial entropy is computed based on the sizes of its connected components $\{C_i\}_{i=1}^{N}$.



Figure 5. Overview of finding localization heads. We first identify heads with high attention sum. Then, we evaluate spatial entropy for each head and select 10 heads with the lowest spatial entropy. We repeat this process for 1,000 image-text pairs and calculate the selection frequency of each head.

as:

$$H(\boldsymbol{A}^{\ell,h}) = -\sum_{i=1}^{N} P(C_i) \log P(C_i), \quad (3)$$

where $P(C_i) = |C_i| / \sum_{i=1}^{N} |C_i|$. As a result, an attention map $\boldsymbol{A}^{\ell,h}$ is considered effectively localized if it exhibits low spatial entropy. For more mathematical details on spatial entropy, please refer to the Appendix Sec. B.

## 4.2. Finding Localization Heads via Criteria

In this section, we utilize the two criteria described earlier to select a small subset of attention heads. Then, we demonstrate that the selected heads effectively capture objects relevant to the text.

4

Figure 6. (a) Selection frequency of individual heads. Only a few heads exhibit high selection frequency, suggesting that their attention maps are consistently 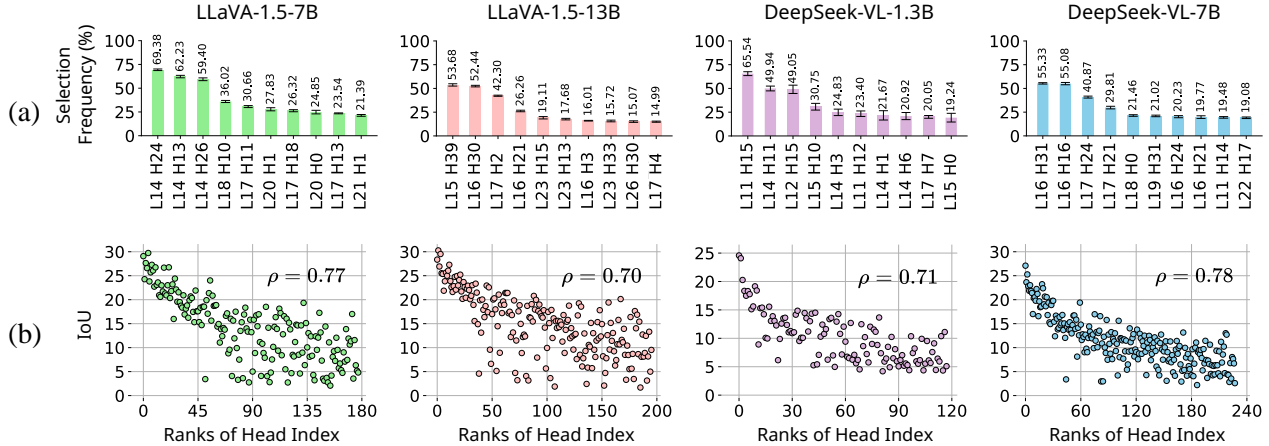well-localized. We calculate the selection frequency five times and report the average and standard deviation. (b) Scatter plot illustrating the relationship between selection frequency rank and each head's average IoU. Heads with higher selection frequency tend to show higher IoU values, indicating that they capture text semantics more effectively. The Spearman correlation coefficient ($\rho$) between rank and IoU is displayed in the top-right corner. The results of the Spearman correlation are statistically significant ($p < 0.001$).

To begin with, we rank all attention heads in order of how well they meet our criteria. Specifically, for 1,000 random image-text samples from the RefCOCO [22] training set, we retain all the heads that satisfy $S_{\text{img}}^{\ell,h} \geq \tau$. Among these heads, we calculate the frequency with which each head exhibits the 10-lowest spatial entropy across the samples to identify heads consistently exhibiting low spatial entropy. We refer to this metric as the selection frequency. The overall process is illustrated in Fig. 5, and the results are reported in Fig. 6(a). Now, we assign ranks to each head based on their selection frequency, with higher-ranked heads being those with high selection frequency. For example, in Fig. 6(a), with LLaVA-1.5-7B [35], head L14 H24 ranks first, followed by head L14 H13 in second place.

Finally, we aim to demonstrate that higher-ranked heads are more effective at capturing objects relevant to the text. To this end, we binarize the attention maps of each head to obtain pseudo-masks and measure the IoU between these pseudo-masks and the ground truth (GT) masks. Then, we visualize the relationship between head ranks, derived from Fig. 6(a), and their IoU values as a scatter plot, shown in Fig. 6(b). Note that only the heads with a selection frequency of at least 1% are considered in this analysis.

As visualized in Fig. 6(b), attention heads with higher selection frequency tend to exhibit higher average IoU. We also calculate the Spearman correlation coefficient to quantitatively evaluate the relationship between the selection frequency and IoU. The correlation coefficients are above 0.7 for all LVLMs, indicating strong positive correlations. This trend becomes increasingly evident for heads with higher ranks, leading us to conclude that a small number of top-ranked heads strongly capture semantic information. We refer to these heads as *localization heads*. Since the



Figure 7. Our training-free visual grounding framework. Attention maps of localization heads are assembled into a combined map, which is then used to define the bounding box or segmentation mask.

trend consistently appears across various LVLMs (see Appendix Sec. C. for trends across more LVLMs), we claim that localization heads are an innate property of LVLMs.

## 5. Visual Grounding with Localization Heads

In the previous section, we demonstrated that our criteria effectively identifies text-referring localization heads. Building on this, we propose a simple yet effective method to solve visual grounding tasks using these localization heads.

Specifically, our objective is to perform visual grounding tasks, given an LVLM. To achieve this, the localization heads of the LVLM must first be identified. Following the process we described in Sec. 4.2 and Fig. 5, we rank the heads based on the selection frequency and select the heads with the $k$-highest rank. Subsequently, an image-text pair for which a mask is to be generated is fed into the LVLM, and attention maps are extracted from the localization heads.

As illustrated in Fig. 7, Gaussian smoothing is applied to each attention map of the localization head to preserve detailed localization information while minimizing

5

Table 1. Comparison of our method with existing fine-tuning based and training-free methods on the REC (Referring Expression Comprehension) task. All fine-tuning based methods are trained on the training set of the corresponding datasets. Best performance is colored in red for fine-tuning and in blue for training-free methods.

| Method | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val | test |
| *Fine-tuning based methods* | | | | | | | | |
| MDETR [21] | 86.8 | 89.6 | 81.4 | 79.5 | 84.1 | 70.6 | 81.6 | 80.9 |
| SeqTR [77] | 87.0 | 90.2 | 83.6 | 78.7 | 84.5 | 71.9 | 82.7 | 83.4 |
| G-DINO [37] | 89.2 | 91.9 | 86.0 | 81.1 | 87.4 | 74.7 | 84.2 | 84.9 |
| ONE-PEACE [60] | 92.6 | 94.2 | 89.3 | 88.8 | 92.2 | 83.2 | 89.2 | 89.3 |
| UNINEXT [31] | 92.6 | 94.3 | 91.5 | 85.2 | 89.6 | 79.8 | 88.7 | 89.4 |
| *Fine-tuning based methods w/ LVLMs* | | | | | | | | |
| Shikra-7B [6] | 87.0 | 90.6 | 80.2 | 81.6 | 87.4 | 72.1 | 82.3 | 82.2 |
| Ferret-7B [68] | 87.5 | 91.4 | 82.5 | 80.8 | 87.4 | 73.1 | 83.9 | 84.8 |
| Shikra-13B [6] | 87.8 | 91.1 | 81.8 | 82.9 | 87.8 | 74.4 | 82.6 | 83.2 |
| Ferret-13B [68] | 89.5 | 92.4 | 84.4 | 82.8 | 88.1 | 75.2 | 85.8 | 86.3 |
| CogVLM-17B [61] | 92.8 | 94.8 | 89.0 | 88.7 | 92.9 | 83.4 | 89.8 | 90.8 |
| *Training-free methods* | | | | | | | | |
| ReCLIP [52] | 45.8 | 46.1 | 47.1 | 47.9 | 50.1 | 45.1 | 59.3 | 59.0 |
| Han et al. [15] | 49.4 | 47.8 | 51.7 | 48.9 | 50.0 | 46.9 | 61.0 | 60.0 |
| GroundVLP [48] | 65.0 | 73.5 | 55.0 | 68.8 | 78.1 | 57.3 | 74.7 | 75.0 |
| *Training-free methods w/ LVLMs (Ours)* | | | | | | | | |
| DeepSeek-VL-1.3B | 73.2 | 77.7 | 70.7 | 62.0 | 66.7 | 57.1 | 65.2 | 69.3 |
| Mini-Gemini-2B | 74.0 | 77.5 | 71.1 | 62.5 | 67.8 | 59.3 | 65.1 | 69.3 |
| InternVL-6B | 85.2 | 86.4 | 78.5 | 78.0 | 83.3 | 71.9 | 81.1 | 80.5 |
| Yi-VL-6B | 85.1 | 86.8 | 78.4 | 78.9 | 84.2 | 72.2 | 80.5 | 80.9 |
| DeepSeek-VL-7B | 85.3 | 87.2 | 81.0 | 77.8 | 83.9 | 73.5 | 81.1 | 82.8 |
| ShareGPT4V-7B | 86.1 | 87.1 | 80.5 | 79.7 | 86.2 | 71.3 | 82.4 | 82.9 |
| LLaVA-7B | 80.3 | 83.5 | 77.4 | 74.5 | 80.2 | 69.3 | 77.5 | 77.1 |
| LLaVA-1.5-7B | 86.5 | 89.8 | 80.2 | 80.1 | 86.3 | 71.9 | 82.3 | 83.0 |
| LLaVA-13B | 82.8 | 85.3 | 79.8 | 79.3 | 82.4 | 73.0 | 79.8 | 79.5 |
| LLaVA-1.5-13B | 87.2 | 90.0 | 83.3 | 82.7 | 88.5 | 74.0 | 84.3 | 85.5 |

Table 2. Comparison of our method with existing fine-tuning based and training-free methods on the RES (Referring Expression Segmentation) task. All fine-tuning based methods, except for LISA [26] and GSVA [65], are trained on the training set of the corresponding datasets. Red and blue colors are used as in Tab. 1.

| Method | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val | test |
| *Fine-tuning based methods* | | | | | | | | |
| LAVT [67] | 72.7 | 75.8 | 68.8 | 62.1 | 68.4 | 55.1 | 61.2 | 62.1 |
| ReLA [34] | 73.8 | 76.5 | 70.2 | 66.0 | 71.0 | 57.7 | 65.0 | 66.0 |
| UniRef++ [63] | 79.1 | 82.1 | 77.5 | 68.4 | 74.0 | 61.5 | 71.4 | 72.8 |
| UNINEXT [31] | 82.2 | 83.4 | 81.3 | 72.5 | 76.4 | 66.2 | 74.4 | 76.4 |
| *Fine-tuning based methods w/ LVLMs* | | | | | | | | |
| LISA-7B [26] | 74.1 | 76.5 | 71.1 | 62.4 | 67.4 | 56.5 | 66.4 | 68.5 |
| GSVA-7B [65] | 76.4 | 77.4 | 72.8 | 64.5 | 67.7 | 58.6 | 71.1 | 72.0 |
| LISA-13B [65] | 73.4 | 76.2 | 69.5 | 62.3 | 66.6 | 56.3 | 68.2 | 68.5 |
| GSVA-13B [65] | 77.7 | 79.9 | 74.2 | 68.0 | 71.5 | 61.5 | 73.2 | 73.9 |
| GLaMM [45] | 79.5 | 83.2 | 76.9 | 75.9 | 78.7 | 68.8 | 76.8 | 78.4 |
| PSALM [74] | 83.6 | 84.7 | 81.6 | 72.9 | 75.5 | 70.1 | 73.8 | 74.4 |
| *Training-free methods* | | | | | | | | |
| Yu et al. [71] | 24.9 | 23.6 | 24.7 | 26.2 | 24.9 | 25.8 | 31.1 | 31.0 |
| TAS [53] | 29.5 | 30.3 | 28.2 | 33.2 | 38.8 | 28.0 | 35.8 | 36.2 |
| Ref-Diff [40] | 35.2 | 37.4 | 34.5 | 35.6 | 38.7 | 31.4 | 38.6 | 37.5 |
| *Training-free methods w/ LVLMs (Ours)* | | | | | | | | |
| DeepSeek-VL-1.3B | 56.3 | 57.0 | 52.7 | 51.2 | 55.5 | 49.2 | 52.3 | 55.8 |
| Mini-Gemini-2B | 59.8 | 60.3 | 55.5 | 56.3 | 59.9 | 51.8 | 55.1 | 60.3 |
| InternVL-6B | 62.1 | 65.8 | 60.9 | 62.2 | 65.5 | 55.5 | 63.5 | 65.4 |
| Yi-VL-6B | 62.5 | 65.8 | 60.7 | 61.0 | 65.3 | 56.0 | 64.0 | 67.0 |
| DeepSeek-VL-7B | 73.9 | 76.6 | 70.7 | 63.1 | 66.1 | 56.5 | 64.0 | 68.9 |
| ShareGPT4V-7B | 73.5 | 76.7 | 70.1 | 59.4 | 63.8 | 55.9 | 60.7 | 65.1 |
| LLaVA-7B | 65.4 | 66.2 | 61.1 | 59.9 | 63.2 | 52.7 | 59.7 | 63.3 |
| LLaVA-1.5-7B | 74.2 | 76.5 | 70.4 | 62.5 | 65.2 | 56.0 | 64.2 | 68.1 |
| LLaVA-13B | 66.8 | 68.0 | 63.7 | 62.3 | 66.9 | 57.3 | 65.0 | 68.2 |
| LLaVA-1.5-13B | 76.1 | 78.9 | 72.8 | 64.1 | 67.1 | 57.3 | 67.7 | 69.0 |

potential random noise. The resulting maps are assembled through element-wise summation to produce the combined map. This combined map is then binarized to produce the pseudo-mask. Finally, the largest rectangle encompassing the pseudo-mask is identified and can be used as a bounding box. Additionally, this bounding box can serve as a prompt for SAM [24] to address the segmentation task. Additional details on the algorithm used to find the bounding box are provided in Appendix. Sec. B, and the ablation study on Gaussian smoothing is presented in Appendix. Sec. D.

# 6. Experiments

In this section, we verify whether the localization head discovered through our selection process ensures robust performance on well-known visual grounding benchmarks. Additionally, we conduct ablation studies to validate the settings of our method.

## 6.1. Experimental Setup

**Models.** We apply our approach across ten LVLMs to validate its broad applicability. The main experiments include DeepSeek-VL [38], Mini-Gemini [30], InternVL [8], Yi-VL [69], ShareGPT4V [7], LLaVA [36], and LLaVA-1.5 [35], with model sizes ranging from 1.3B to 13B. The number of localization heads is fixed to $k = 3$ for all models.

**Benchmarks.** To assess visual grounding capabilities, we conduct experiments on Referring Expression Comprehension (REC) and Referring Expression Segmentation (RES) tasks. REC requires the model to predict the bounding box of the referred object, while RES requires the segmentation mask. We use the RefCOCO, RefCOCO+ [22], and RefCOCOg [18] datasets. We further evaluate the performance of our method on the more challenging scenario, Reasoning Segmentation (ReasonSeg) [26], which requires complex reasoning or world knowledge. For the REC task, we report the performance using Acc@0.5 metric, which is the standard detection metric for REC. For the RES and ReasonSeg task, cIoU is used as the evaluation metric.

**Baselines.** We compare our method with existing fine-tuning based and training-free approaches. The fine-tuning based methods include visual grounding specialist models [21, 31, 34, 37, 63, 67, 77], along with fine-tuned LVLMs for object localization [6, 61, 68] or segmentation tasks [26, 45, 65, 74]. The training-free methods include
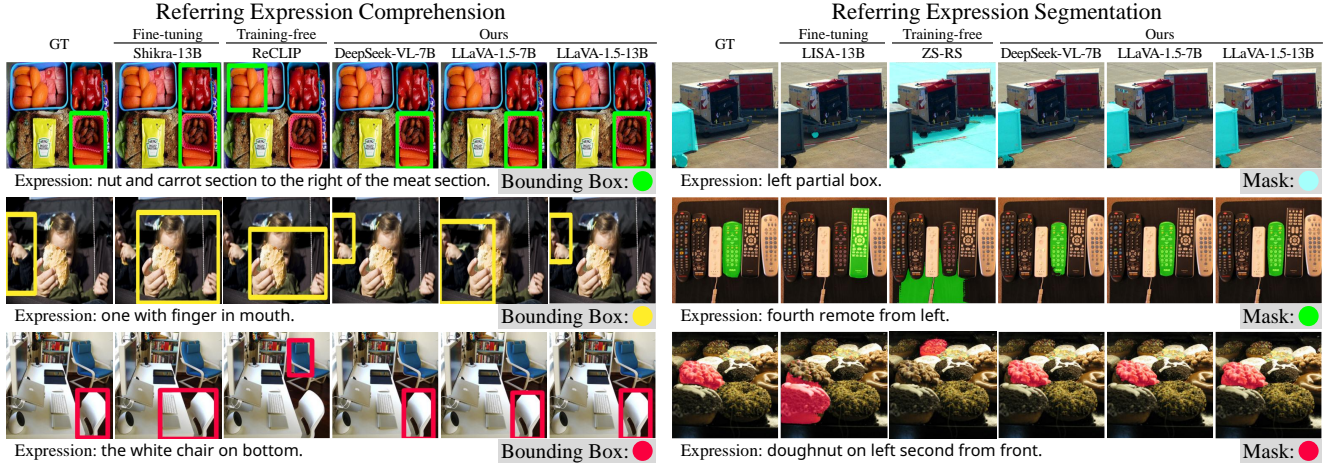
Figure 8. Qualitative results of our framework with the baseline models. LVLMs successfully localize the referred objects in various challenging scenarios including multiple similar objects, non-salient objects, and complex spatial relations.

Table 3. Comparison of our framework with LISA [26] on the ReasonSeg (Reasoning Segmentation) benchmark.

| Method | val | test | | |
|---|---|---|---|---|
| | overall | short | long | overall |
| LISA-7B [26] | 52.3 | 48.5 | 48.9 | 48.8 |
| LISA-13B [26] | 60.3 | 50.0 | 50.9 | 50.8 |
| LLaVA-1.5-7B (Ours) | 52.4 | 48.0 | 49.1 | 48.7 |
| LLaVA-1.5-13B (Ours) | 60.5 | 48.7 | 51.0 | 49.9 |

CLIP-based methods [15, 48, 52, 53, 71] and DM-based method [40]. More details on the experimental setup are provided in the Appendix Sec. A.

## 6.2. Main Results

**REC and RES.** Tab. 1 and Tab. 2 present the results of our method and the baseline models on the REC and RES tasks, respectively. Our framework achieves substantial improvements over the existing training-free methods. Surprisingly, our method performs comparably to the fine-tuned LVLMs, even though our method does not require additional training. For example, in the REC task, the best performance of our approach achieves results on par with Shikra [6] and Ferret [68], which share the same base LLMs as LLaVA-1.5 [35], but are fine-tuned for localization tasks. A similar finding is observed with LISA [26] in the RES task. The results indicate that frozen LVLMs can effectively localize the referred object without any additional training, due to the presence of localization heads.

Notably, the visual grounding capability is enhanced as the model evolves. First, performance consistently improves as model size increases (1.3B to 13B). Second, updates in architecture and training data (*e.g.*, LLaVA to LLaVA-1.5) also boost performance. This observation suggests that the grounding ability of LVLMs could be further enhanced with larger models and more diverse training data.

Fig. 8 compares the qualitative results of our method with those of the baseline models. The results demonstrate that LVLMs can accurately identify the correct object regions, even in challenging scenarios where multiple similar objects are present, or when the referred object is not prominently centered in the image. According to [52], CLIP-based methods struggle to interpret orientation descriptors (*e.g.*, "left"). Therefore, they have to manually decompose the referring expression into multiple components [71] or rely on post-processing steps that use the object's spatial information [53]. In contrast, our framework can directly predict the bounding box or segmentation mask of the referred object without carefully designed post-processing steps, with the help of the strong text comprehension capabilities of LVLMs. More qualitative results are provided in the Appendix Sec. E.

**Reasoning Segmentation.** Tab. 3 shows the results of our method and LISA [26] on the ReasonSeg. For a fair comparison, we compare both methods using the same backbone model, LLaVA-1.5 [35]. Our method performs comparably to LISA and sometimes outperforms it. The results suggest that the localization heads in LVLMs are generalizable to various visual grounding tasks, including those that require complex reasoning or world knowledge.

## 6.3. Ablation Studies

**Number of Localization Heads.** In our main experiments, we set the number of localization heads to $k = 3$. Here, we investigate the effect of varying $k$ on visual grounding performance. Tab. 4 presents the results of our framework with different $k$ values. We observe that the performance generally improves as $k$ increases from 1 to 3, indicating that top-3 heads complement each other to provide more accurate localization. However, increasing $k$ further does not guarantee better performance, implying that additional heads may

Table 4. Ablation study on the number of localization heads ($k$) on the RefCOCO validation set for the RES task.

| Method | $k$ (# of Localization Heads) | | | | |
|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 |
| DeepSeek-VL-1.3B | 55.1 | 56.3 | 56.3 | 55.3 | 51.2 |
| MiniGemini-2B | 58.0 | 58.5 | 59.8 | 59.1 | 54.2 |
| InternVL-6B | 61.3 | 61.8 | 62.1 | 61.0 | 55.7 |
| Yi-VL-6B | 61.8 | 62.1 | 62.5 | 62.6 | 55.4 |
| DeepSeek-VL-7B | 70.1 | 72.2 | 73.9 | 73.0 | 65.3 |
| ShareGPT4V-7B | 70.3 | 72.4 | 73.5 | 73.5 | 60.8 |
| LLaVA-7B | 62.7 | 63.1 | 65.4 | 65.3 | 57.7 |
| LLaVA-1.5-7B | 70.3 | 73.1 | 74.2 | 74.1 | 65.4 |
| LLaVA-13B | 63.5 | 64.7 | 66.8 | 66.4 | 57.8 |
| LLaVA-1.5-13B | 71.7 | 75.7 | 76.1 | 76.0 | 65.7 |
| Average | 64.5 | 66.0 | 67.1 | 65.4 | 58.9 |

Table 5. Ablation study on the validation of criteria and selection methods for localization heads. The results are reported on the RefCOCO validation set and LLaVA-1.5-13B.

| Criteria | | Selection | | REC | RES |
|---|---|---|---|---|---|
| $S_{\text{img}}^{\ell,h}$ | $H(\boldsymbol{A}^{\ell,h})$ | Fixed | Greedy | | |
| ✓ | | | ✓ | 23.7 | 18.8 |
| | ✓ | | ✓ | 29.8 | 21.5 |
| ✓ | ✓ | | ✓ | 67.4 | 63.8 |
| ✓ | | ✓ | | 23.9 | 19.3 |
| | ✓ | ✓ | | 31.3 | 25.7 |
| ✓ | ✓ | ✓ | | 87.2 | 76.1 |

introduce noise or redundancy. It is worth noting that the optimal $k$ trend remains consistent across different LVLMs. The results suggest that similar numbers of attention heads are responsible for localization of referred objects in various LVLMs, even though the total number of heads and model architectures differ.

**Validation of Criteria and Selection Methods for Localization Heads.** In Sec. 4.1, we propose two criteria, attention sum $S_{\text{img}}^{\ell,h}$ and spatial entropy $H(\boldsymbol{A}^{\ell,h})$, to identify localization heads. Then, we select the fixed top-$k$ heads based on the selection frequency, as described in Sec. 4.2. We ablate the effectiveness of each criterion and validate selection methods. For criterion ablation, we evaluate the performance of our method using each criterion individually: (1) selecting heads with either the highest $S_{\text{img}}^{\ell,h}$ or (2) the lowest $H(\boldsymbol{A}^{\ell,h})$ only. For selection validation, we compare the performance of our method (denoted as the 'fixed' method for comparison) with 'greedy' selection, where the top-$k$ heads are selected and aggregated per sample. Further details regarding the settings are provided in Appendix Sec. A.

Tab. 5 shows the results of these ablation studies. The performance drops significantly when only one criterion is used, indicating that both criteria are essential for identifying localization heads. Furthermore, the greedy selection
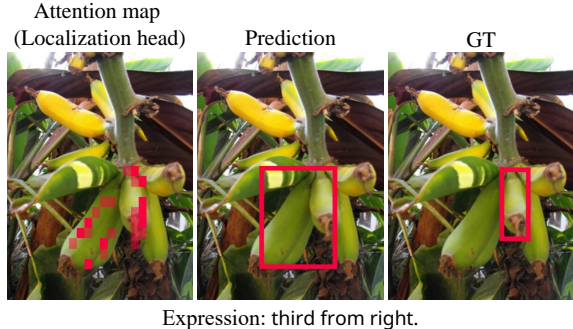


Expression: third from right.

Figure 9. Failure case of the LLaVA-1.5-13B [35] in visual grounding. The text-to-image attention map from a localization head (L15 H39) shows where the model focuses, helping to understand the model's failure.

method shows worse results than the fixed method. While our criteria identify attention maps exhibiting apparent clusters, they do not ensure that these clusters are formed around text semantics. As a result, the greedy method may select heads that are localized but not text-referred. In contrast, our method involves a statistical analysis (*i.e.*, selection frequency). This ensures that the localization heads are genuinely text-referred, consistently focusing on text-related regions rather than arbitrarily clustered areas.

### 6.4. Understanding LVLMs When They Fail

Here, we briefly discuss how localization heads may also help us better understand LVLMs. Specifically, localization heads allow us to identify where LVLMs focus when they fail to ground the correct object. Fig. 9 illustrates an example where the model fails to predict the correct object, the third banana from the right. As shown in the first column of Fig. 9, the text-to-image attention map from a localization head focuses on both the third and fourth bananas from the right. This observation suggests that LVLMs struggle with pinpointing the exact location of objects. These findings show the localization head's potential to provide a transparent understanding of where the LVLMs focus.

### 7. Conclusion

In this work, we identify *localization heads* within various LVLMs via criteria, which exhibit strong visual grounding capabilities in response to textual queries. We then propose a simple yet effective training-free framework that assembles the text-to-image attention maps from a few localization heads to predict bounding boxes and segmentation masks for text-relevant regions in the image. Our approach achieves competitive performance compared to fine-tuning based methods. Therefore, we conclude that LVLMs can act as text-referring localizers for visual grounding tasks with their inherent property under the attention mechanisms. We hope that our work opens up new possibilities for analyzing and utilizing the attention mechanisms of LVLMs.

# References

[1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. 2023. 16

[2] Michael Batty. Spatial entropy. *Geographical Analysis*, 6(1): 1–31, 1974. 2, 4, 13

[3] Ryan Burgert, Kanchana Ranasinghe, Xiang Li, and Michael S Ryoo. Peekaboo: Text to image diffusion models are zero-shot segmentors. *arXiv preprint arXiv:2211.13224*, 2022. 3

[4] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22560–22570, 2023. 2

[5] Khyathi Raghavi Chandu, Yonatan Bisk, and Alan W Black. Grounding 'grounding' in NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4283–4305, Online, 2021. Association for Computational Linguistics. 2

[6] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 1, 6, 7

[7] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023. 3, 6, 14, 16

[8] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198, 2024. 3, 6, 14, 16

[9] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? an analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy, 2019. Association for Computational Linguistics. 3

[10] Explosion AI. spaCy: Industrial-strength Natural Language Processing in Python. 15

[11] Wei Gao, Fang Wan, Xingjia Pan, Zhiliang Peng, Qi Tian, Zhenjun Han, Bolei Zhou, and Qixiang Ye. Ts-cam: Token semantic coupled attention map for weakly supervised object localization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2886–2895, 2021. 2

[12] Cristina González, Nicolás Ayobi, Isabela Hernández, José Hernández, Jordi Pont-Tuset, and Pablo Arbeláez. Panoptic narrative grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1364–1373, 2021. 15

[13] Ronald L. Graham. An efficient algorithm for determining the convex hull of a finite planar set. *Info. Proc. Lett.*, 1: 132–133, 1972. 14

[14] Costantino Grana, Daniele Borghesani, and Rita Cucchiara. Optimized block-based connected components labeling with decision trees. *IEEE Transactions on Image Processing*, 19 (6):1596–1609, 2010. 4

[15] Zeyu Han, Fangrui Zhu, Qianru Lao, and Huaizu Jiang. Zero-shot referring expression comprehension via structural similarity between images and captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14364–14374, 2024. 6, 7

[16] Shuting He, Henghui Ding, Chang Liu, and Xudong Jiang. Grec: Generalized referring expression comprehension. *arXiv preprint arXiv:2308.16182*, 2023. 15

[17] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2

[18] Ronghang Hu, Marcus Rohrbach, and Trevor Darrell. Segmentation from natural language expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 108–124. Springer, 2016. 1, 2, 3, 6, 13

[19] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR, 2021. 16

[20] Nick Jiang, Anish Kachinthaya, Suzie Petryk, and Yossi Gandelsman. Interpreting and editing vision-language representations to mitigate hallucinations. *arXiv preprint arXiv:2410.02762*, 2024. 3

[21] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1780–1790, 2021. 6

[22] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. 3, 4, 5, 6, 13

[23] Seoyeon Kim, Minguk Kang, Dongwon Kim, Jaesik Park, and Suha Kwak. Extending clip's image-text alignment to referring image segmentation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4611–4628, 2024. 2

[24] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 3, 6

[25] Vedang Lad, Wes Gurnee, and Max Tegmark. The remarkable robustness of llms: Stages of inference? *arXiv preprint arXiv:2406.19384*, 2024. 3

[26] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation

via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 1, 2, 3, 6, 7, 13

[27] Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text documents, 2023. 16

[28] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 1, 16

[29] Ruiyu Li, Kaican Li, Yi-Chun Kuo, Michelle Shu, Xiaojuan Qi, Xiaoyong Shen, and Jiaya Jia. Referring image segmentation via recurrent refinement networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2018. 2

[30] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024. 3, 6, 14, 15, 16

[31] Fangjian Lin, Jianlong Yuan, Sitong Wu, Fan Wang, and Zhibin Wang. Uninext: Exploring a unified architecture for vision recognition. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 3200–3208, 2023. 6

[32] 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 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 13

[33] Chenxi Liu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, and Alan Yuille. Recurrent multimodal interaction for referring image segmentation. In *Proceedings of the IEEE international conference on computer vision*, pages 1271–1280, 2017. 2

[34] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 23592–23601, 2023. 2, 6, 15

[35] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 1, 2, 3, 4, 5, 6, 7, 8, 15, 16, 23, 24

[36] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 1, 3, 6, 14, 16

[37] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 6

[38] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 3, 6, 15, 16

[39] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016. 1, 2

[40] Minheng Ni, Yabo Zhang, Kailai Feng, Xiaoming Li, Yiwen Guo, and Wangmeng Zuo. Ref-diff: Zero-shot referring image segmentation with generative models. *arXiv preprint arXiv:2308.16777*, 2023. 3, 6, 7

[41] Elia Peruzzo, Enver Sangineto, Yahui Liu, Marco De Nadai, Wei Bi, Bruno Lepri, and Nicu Sebe. Spatial entropy as an inductive bias for vision transformers. *Machine Learning*, 113(9):6945–6975, 2024. 4, 13

[42] Koutilya Pnvr, Bharat Singh, Pallabi Ghosh, Behjat Siddiquie, and David Jacobs. Ld-znet: A latent diffusion approach for text-based image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4157–4168, 2023. 2

[43] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*. 16, 24

[44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3

[45] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018, 2024. 1, 6, 15

[46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149, 2016. 3

[47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3

[48] Haozhan Shen, Tiancheng Zhao, Mingwei Zhu, and Jianwei Yin. Groundvlp: Harnessing zero-shot visual grounding from vision-language pre-training and open-vocabulary object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4766–4775, 2024. 6, 7

[49] Hengcan Shi, Hongliang Li, Fanman Meng, and Qingbo Wu. Key-word-aware network for referring expression image seg-

mentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 38–54, 2018. 2

[50] Mohit Shridhar, Dixant Mittal, and David Hsu. Ingress: Interactive visual grounding of referring expressions. *The International Journal of Robotics Research*, 39(2-3):217–232, 2020. 2

[51] Oriane Siméoni, Gilles Puy, Huy V Vo, Simon Roburin, Spyros Gidaris, Andrei Bursuc, Patrick Pérez, Renaud Marlet, and Jean Ponce. Localizing objects with self-supervised transformers and no labels. *arXiv preprint arXiv:2109.14279*, 2021. 4

[52] Sanjay Subramanian, William Merrill, Trevor Darrell, Matt Gardner, Sameer Singh, and Anna Rohrbach. ReCLIP: A strong zero-shot baseline for referring expression comprehension. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5198–5215, Dublin, Ireland, 2022. Association for Computational Linguistics. 3, 6, 7

[53] Yucheng Suo, Linchao Zhu, and Yi Yang. Text augmented spatial aware zero-shot referring image segmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1032–1043, Singapore, 2023. Association for Computational Linguistics. 3, 6, 7

[54] Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin, and Ferhan Ture. What the daam: Interpreting stable diffusion using cross attention. *arXiv preprint arXiv:2210.04885*, 2022. 2

[55] OpenGVLab Team. Internvl2, 2024. 16

[56] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024. 1

[57] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808, Florence, Italy, 2019. Association for Computational Linguistics. 3

[58] Feng Wang, Jieru Mei, and Alan Yuille. Sclip: Rethinking self-attention for dense vision-language inference. In *European Conference on Computer Vision*, pages 315–332. Springer, 2025. 4

[59] Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. Label words are anchors: An information flow perspective for understanding in-context learning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9840–9855, Singapore, 2023. Association for Computational Linguistics. 3

[60] Peng Wang, Shijie Wang, Junyang Lin, Shuai Bai, Xiaohuan Zhou, Jingren Zhou, Xinggang Wang, and Chang Zhou. One-peace: Exploring one general representation model toward unlimited modalities. *arXiv preprint arXiv:2305.11172*, 2023. 6

[61] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. 1, 6

[62] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11686–11695, 2022. 2

[63] Jiannan Wu, Yi Jiang, Bin Yan, Huchuan Lu, Zehuan Yuan, and Ping Luo. Uniref++: Segment every reference object in spatial and temporal spaces. *arXiv preprint arXiv:2312.15715*, 2023. 6

[64] Size Wu, Sheng Jin, Wenwei Zhang, Lumin Xu, Wentao Liu, Wei Li, and Chen Change Loy. F-lmm: Grounding frozen large multimodal models. *arXiv preprint arXiv:2406.05821*, 2024. 3, 15

[65] Zhuofan Xia, Dongchen Han, Yizeng Han, Xuran Pan, Shiji Song, and Gao Huang. Gsva: Generalized segmentation via multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3858–3869, 2024. 1, 6

[66] Li Xu, Mark He Huang, Xindi Shang, Zehuan Yuan, Ying Sun, and Jun Liu. Meta compositional referring expression segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19478–19487, 2023. 2

[67] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Lavt: Language-aware vision transformer for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18155–18165, 2022. 2, 6

[68] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023. 1, 6, 7

[69] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024. 3, 6, 14, 16

[70] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016. 1, 2

[71] Seonghoon Yu, Paul Hongsuck Seo, and Jeany Son. Zero-shot referring image segmentation with global-local context features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19456–19465, 2023. 3, 6, 7

[72] Sukmin Yun, Hankook Lee, Jaehyung Kim, and Jinwoo Shin. Patch-level representation learning for self-supervised vision transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8354–8363, 2022. 4

[73] Bowen Zhang, Zhi Tian, Quan Tang, Xiangxiang Chu, Xiaolin Wei, Chunhua Shen, et al. Segvit: Semantic segmentation with plain vision transformers. *Advances in Neural Information Processing Systems*, 35:4971–4982, 2022. 2

11

[74] Zheng Zhang, Yeyao Ma, Enming Zhang, and Xiang Bai. Psalm: Pixelwise segmentation with large multi-modal model. *arXiv preprint arXiv:2403.14598*, 2024. 1, 6

[75] Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-to-image diffusion models for visual perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5729–5739, 2023. 2

[76] Zifan Zheng, Yezhaohui Wang, Yuxin Huang, Shichao Song, Bo Tang, Feiyu Xiong, and Zhiyu Li. Attention heads of large language models: A survey. *arXiv preprint arXiv:2409.03752*, 2024. 3

[77] Chaoyang Zhu, Yiyi Zhou, Yunhang Shen, Gen Luo, Xingjia Pan, Mingbao Lin, Chao Chen, Liujuan Cao, Xiaoshuai Sun, and Rongrong Ji. Seqtr: A simple yet universal network for visual grounding. In *European Conference on Computer Vision*, pages 598–615. Springer, 2022. 6

# Appendix

## A. Experimental Details

**Experimental Setting.** All experiments and evaluations are conducted on a single NVIDIA GeForce RTX A6000 48GB GPU. We only use the inference stage of the models without any fine-tuning or training.

**Analysis Setting.** In Sec. 4, we identify and analyze localization heads in various LVLMs. We use the RefCOCO training set to prevent validation set leakage. To calculate the selection frequency of individual heads, we randomly select 1,000 image-text pair samples from the RefCOCO training set and average the results over five trials to validate consistency. When analyzing the selection frequency and IoU, we binarize the attention weights by assigning 1 above the mean value and 0 below it and calculate the IoU between the binarized attention weights and the ground-truth mask. We repeat this process for 1,000 image-text pairs and average the IoU scores.

**Dataset Details.** We evaluate our method on the following datasets:

- RefCOCO, RefCOCO+ [22], and RefCOCOg [18] datasets, sourced from MS-COCO [32], offer a collection of referring expressions and associated images. RefCOCO consists of 19,994 images paired with 142,210 expressions, while RefCOCO+ includes 19,992 images and 141,564 expressions. RefCOCOg, on the other hand, contains 26,771 images and 104,560 expressions. The expressions in RefCOCO and RefCOCO+ are generally concise, with an average of 1.6 nouns and 3.6 words per expression. In contrast, RefCOCOg features more descriptive expressions, averaging 2.8 nouns and 8.4 words.

- ReasonSeg: The dataset and benchmark for reasoning segmentation were first introduced in LISA [26]. The resulting ReasonSeg benchmark consists of 1,218 image-instruction-mask data samples, which are further divided into three splits: training (239 samples), validation (200 samples), and test (779 samples).

**Main Experiments Setting.** We evaluate our method on the following tasks:

- Referring Expression Comprehension (REC) and Referring Expression Segmentation (RES): The datasets evaluated for the main results in Sec.6.2 include RefCOCO (validation, test-A, test-B), RefCOCO+ (validation, test-A, test-B), and RefCOCOg (validation, test). All evaluations were conducted using the UNC split.

- Reasoning Segmentation (ReasonSeg): Reasoning Segmentation was first introduced in LISA [26]. This task shares a similar formulation with the referring expression segmentation task but is considerably more challenging. The key distinction lies in the complexity of the query text in reasoning segmentation. Rather than simple phrases (e.g., "the blue mug"), the queries involve more nuanced

descriptions (e.g., "the container used for drinking, located next to the plate") or longer sentences (e.g., "Find the item on the table that someone would use to hold liquid, often paired with a saucer"). These queries demand advanced reasoning and a deeper understanding of contextual and world knowledge. All reasoning segmentation results were evaluated using the ReasonSeg benchmark, which includes both the validation set and test set. Performance was measured across short queries, long queries, and overall, following the same experimental setup as LISA [26] to ensure consistency in comparisons.

**Ablation Studies Setting.** In Sec. 6.3, we ablate the effectiveness of each criterion and validate the selection methods. In this section, we provide details of the ablation studies.

For criterion ablation, we consider two approaches: (1) selecting heads based solely on the highest $S_{\text{img}}^{\ell,h}$ values, or (2) selecting heads based solely on the lowest $H(\boldsymbol{A}^{\ell,h})$ values. In approach (1), we select the 10 heads with the highest $S_{\text{img}}^{\ell,h}$ values and calculate their selection frequency. Similarly, in approach (2), we select the 10 heads with the lowest $H(\boldsymbol{A}^{\ell,h})$ values and calculate their selection frequency.

For selection validation, we introduce the 'greedy' selection method, which selects the top-$k$ heads per sample without considering the overall selection frequency. When applying the greedy selection method and criterion (1) simultaneously, we select the top-$k$ heads with the highest $S_{\text{img}}^{\ell,h}$ values for each sample. Criterion (2) is applied in a similar manner, simultaneously selecting the top-$k$ heads with the lowest $H(\boldsymbol{A}^{\ell,h})$ values for each sample.

## B. Detailed Description of Algorithms

### B.1. Spatial Entropy

Spatial entropy [2] adjusts the probability of attention being focused in a region by factoring in the size of that region, ensuring fair comparison across areas of different sizes. Note that, our spatial entropy calculation is based on the previous work [41] which validated the effectiveness of spatial entropy in image attention maps within vision transformer. We begin by computing the image attention map $\boldsymbol{A}^{\ell,h}$ as follows:

$$\boldsymbol{A}^{\ell,h} = \text{ReLU}\left(\text{reshape}(\boldsymbol{a}^{\ell,h}) - m\right), \quad (4)$$

where the ReLU function is applied after reshaping by $P \times P$, and it retains only those values in $\boldsymbol{a}^{\ell,h}$ that are greater than the mean $m$. Next, we identify the connected components $C_{\boldsymbol{A}^{\ell,h}} = \{C_1, C_2, \ldots, C_n\}$ from $\boldsymbol{A}^{\ell,h}$:

$$C_{\boldsymbol{A}^{\ell,h}} = \text{ConnectedComponents}(\boldsymbol{A}^{\ell,h}), \quad (5)$$

where the connected components are determined by applying an 8-connectivity relation among the non-zero elements of $\boldsymbol{A}^{\ell,h}$. Each connected component $C_n$ (with

$1 \leq n \leq N$) in $C_{\boldsymbol{A}^{\ell,h}}$ is defined as the set of coordinates $C_n = \{(x_1, y_1), (x_2, y_2), \ldots, (x_{k_n}, y_{k_n})\}$ for the $n$-th component, where $k_n = |C_n|$ represents the cardinality, or the number of elements, in $C_n$. Finally, we calculate the spatial entropy $H(\boldsymbol{A}^{\ell,h})$ as follows:

$$H(\boldsymbol{A}^{\ell,h}) = -\sum_{n=1}^{N} P(C_n) \log P(C_n), \qquad (6)$$

where this entropy is computed using Shannon's entropy formula. Here, $P(C_n)$ represents the probability of observing each connected component $C_n$ within $\boldsymbol{A}^{\ell,h}$. The probability $P(C_n)$ for each component $C_n$ is defined as:

$$P(C_n) = \frac{|C_n|}{\sum_{n=1}^{N} |C_n|}, \qquad (7)$$

where $P(C_n)$ is calculated by dividing the area of $C_n$ by the total area of all components in $\boldsymbol{A}^{\ell,h}$. This provides a normalized measure of spatial focus. The resulting spatial entropy $H(\boldsymbol{A}^{\ell,h})$ ranges from 0 to 1. A value of 0 indicates that attention is completely focused on a single region, while a value of 1 suggests that attention is evenly distributed across the image. This measure thus enables us to evaluate the dispersion of the model's attention across different regions within the image.

## B.2. Details of Our Framework

In this section, we provide a detailed description of our framework, described in Sec. 5 of the main paper.

**Binarization of the Attention Map.** The attention map is binarized by setting values above the mean to 1. This approach effectively highlights the most significant regions of the attention map.

**Gaussian Smoothing.** Gaussian smoothing is applied using a kernel size of $k = 7$ and a standard deviation of $\sigma = 1.0$. These parameters ensure a balance between smoothing effects and detail preservation.

**Convex Hull Algorithm for Bounding Box.** To determine the bounding box in an assembled attention map from the localization heads, we employ the convex hull algorithm [13]. In cases where multiple convex hulls are present within the same attention map, we retain only the largest convex hull. Subsequently, we calculate the smallest tight bounding box that encloses the retained convex hull and we use it as the final bounding box.

## C. More Analysis on Localization Heads

### C.1. Extended Analysis Across More LVLMs

In this section, we extend the analysis of localization heads in Sec. 4 of the main paper to more LVLMs, including InternVL [8], LLaVA [36], Mini-Gemini [30], ShareGPT4V [7], and Yi-VL [69].

**Average Attention Sum in More LVLMs.** We extend Fig. 3 in the main paper to demonstrate that relatively few attention heads significantly contribute to the model's text-image interaction. As shown in Fig. 11, the trend of the average $S_{\text{img}}^{\ell,h}$ values remains consistent across different LVLMs.

**Selection Frequency and IoU in More LVLMs.** Similar to the above, we extend Fig. 6 in the main paper to cover additional LVLMs. Fig. 12 presents the selection frequency and a scatter plot of selection frequency rank versus IoU for each attention head across various LVLMs. The results confirm that our observations hold consistently across different LVLMs.

### C.2. Robustness of Localization Head Selection

In this section, we validate the robustness of our localization head selection method across different threshold values ($\tau$) and the number of selected heads ($N$). The experiments below indicate that localization head selection is not sensitive to the choice of $\tau$ or $N$.

**Threshold $\tau$.** Fig. 3 in the main paper presents the average $S_{\text{img}}$ values for each attention head, setting the threshold $\tau$ at the point where the maximum curvature is observed. We select maximum curvature as the threshold to reduce the need for manual tuning; however, other $\tau$ values can also be considered. Therefore, we further validate that plausible $\tau$ values can give consistent results with the maximum curvature. To this end, we calculate the selection frequency of the heads based on different $\tau$ values and compare them with the results obtained using the maximum curvature. The results are presented in Fig. 13. We observe that the same localization heads are consistently selected across different $\tau$ values, indicating that our analysis results are robust to the choice of $\tau$.

**Number of Heads $N$.** In Fig. 6(a) of the main paper, we select the 10 heads with the lowest $H(\boldsymbol{A}^{\ell,h})$ values and repeat the process for 1,000 image-text pairs to calculate the selection frequency. We also investigate the effect of selecting different numbers of heads ($N$) on the selection frequency. We conduct experiments from $N = 1$ to $N = 14$ and compare the results with the selection frequency obtained using $N = 10$ (default setting). As shown in Fig. 14, we can obtain the same top-3 localization heads consistently across different $N$ values, suggesting that the selection of localization heads is robust to the choice of $N$.

## D. More Experiments

### D.1. Comparison with Baseline Models

Most LVLMs, including the LLaVA [36] family, likely encode localization knowledge in their pretrained weights, possibly due to pretraining with bounding box coordinates or visual instruction prompts [36]. In Tab. 6, we compare baseline models and our proposed method, revealing the

Table 6. Comparison performance to baseline models

| REC (RefCOCOg) | DeepSeekVL-1.3B | LLaVA-1.5-7B | LLaVA-1.5-13B |
|---|---|---|---|
| Baseline | 1.5 | 2.92 | 5.28 |
| Ours | 65.2 | 82.3 | 84.3 |

Table 7. Performance comparison between F-LMM [64] and our method on the RES task. We note that F-LMM models are trained on the training set of Referring Expression Segmentation datasets.

| Method | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val | test |
| *F-LMM (Fine-tuning on RES)* | | | | | | | | |
| DeepSeek-VL-1.3B | 75.0 | 78.1 | 69.5 | 62.8 | 70.8 | 56.3 | 68.2 | 68.5 |
| Mini-Gemini-2B | 75.0 | 78.6 | 69.3 | 63.7 | 71.4 | 53.3 | 67.3 | 67.4 |
| DeepSeek-VL-7B | 76.1 | 78.8 | 72.0 | 66.4 | 73.2 | 57.6 | 70.1 | 70.4 |
| LLaVA-1.5-7B | 75.2 | 79.1 | 71.9 | 63.7 | 71.8 | 54.7 | 67.1 | 68.1 |
| *Ours (Training-free)* | | | | | | | | |
| DeepSeek-VL-1.3B | 56.3 | 57.0 | 52.7 | 51.2 | 55.5 | 49.2 | 52.3 | 55.8 |
| Mini-Gemini-2B | 59.8 | 60.3 | 55.5 | 56.3 | 59.9 | 51.8 | 55.1 | 60.3 |
| DeepSeek-VL-7B | 73.9 | 76.6 | 70.7 | 63.1 | 67.1 | 56.5 | 64.0 | 68.9 |
| LLaVA-1.5-7B | 74.2 | 76.5 | 70.4 | 62.5 | 65.2 | 56.0 | 64.2 | 68.1 |

baseline models' poor localization accuracy, likely due to their focus on describing objects rather than precise localization. Moreover, the localization head might provide only indirect support when text generation unfolds in its usual course. As a result, it becomes difficult for the model to directly output accurate object or region coordinates required for visual grounding, unless the information from this head is explicitly scrutinized. Thus, the localization head's practical value can be realized as long as it is integrated with our proposed method.

## D.2. Comparision with F-LMM

We compare our method with F-LMM [64], which also leverages the attention weights of frozen LVLMs for visual grounding. The differences between F-LMM and our method are as follows. First, F-LMM still requires fine-tuning its mask decoder modules on visual grounding datasets (i.e., referring expression segmentation datasets). Second, F-LMM uses all attention heads without considering the relative importance of each, leaving the decoder modules to interpret the entire set of attention weights. In contrast, our approach requires no fine-tuning and directly utilizes a few selected attention heads that are particularly useful for localizing objects in the image. Furthermore, our framework provides a transparent understanding of where the model focuses through localization heads, which is not available in F-LMM.

Tab. 7 presents the performance comparison between F-LMM and our method on the RES task. In smaller LVLMs (*e.g.*, DeepSeek-VL-1.3B [38] and Mini-Gemini-2B [30]), F-LMM outperforms our method. However, in relatively larger LVLMs (e.g., DeepSeek-VL-7B [38] and LLaVA-1.5-7B [35]), our method demonstrates performance comparable to F-LMM, with only a slight gap. This result sug-

Table 8. Ablation study on Gaussian smoothing parameters ($\sigma$ and $\kappa$). The performance is evaluated using the RefCOCO validation set (UNC split) with the LLaVA-1.5-13B [35].

| Task | $\sigma$ (standard deviation) | | | | | |
|---|---|---|---|---|---|---|
| | 0.0 | 0.4 | 0.8 | 1.0 | 1.4 | 1.8 |
| REC | 85.5 | 86.8 | 87.2 | 87.2 | 86.8 | 84.3 |
| RES | 74.3 | 75.2 | 76.1 | 76.1 | 75.2 | 72.7 |

| Task | $\kappa$ (kernel size) | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 3 | 5 | 7 | 9 | 11 |
| REC | 85.5 | 86.5 | 86.5 | 87.2 | 87.2 | 87.2 |
| RES | 74.3 | 75.2 | 75.2 | 76.1 | 76.1 | 76.1 |

Table 9. Performance comparison with F-LMM [64] on the PNG [12] benchmark.

| PNG (all) | Ours | F-LMM |
|---|---|---|
| DeepSeekVL-7B | 66.7 | 65.7 |

gests that the localization heads have competitive potential with the specialized mask decoder modules for visual grounding tasks, especially in relatively larger LVLMs.

## D.3. Gaussian Smoothing Ablation

When assembling the attention map in the localization head (see Sec. 5 of the main paper), we apply Gaussian smoothing to the attention map to minimize potential random noise. In this section, we conduct an ablation study on the parameters of Gaussian smoothing to better understand the robustness of our framework across different values of standard deviation $\sigma$ and kernel size $\kappa$. For the experiments, LLaVA-1.5-13B [35] was evaluated using the RefCOCO validation set (UNC split).

The results are presented in Tab. 8. Regardless of the selected $\sigma$ and $\kappa$, Gaussian smoothing consistently enhances performance in almost all cases. The findings highlight that the framework is robust to varying choices of $\sigma$ and $\kappa$. Furthermore, even when using the basic attention map of localization heads without Gaussian smoothing ($\sigma = 0$ or $\kappa = 1$), the performance remains competitive, with only a 1.9% drop compared to the best case. This demonstrates that Gaussian smoothing only serves as an auxiliary postprocessing step for refining the attention map from localization heads.

## D.4. Multi-Object Grounding Tasks

Beyond single-object tasks, our pipeline also suggests promise for multi-object grounding. We utilize spaCy [10] to extract noun tokens for generating attention maps (see Fig. 10), obtaining comparable results on the PNG benchmark [12], with improvements observed relative to F-LMM (see Tab. 9). Similarly, we believe this approach holds promise for extension to other various tasks [16, 34, 45].

## E. More Qualitative Results

We present more qualitative results of our framework, including the performance of 10 LVLMs [7, 8, 30, 35, 36, 38, 69], with parameter numbers ranging from 1.3B to 13B, on visual grounding tasks. Fig. 15, Fig. 16, and Fig. 17 present the qualitative results of our method on the Referring Expression Comprehension (REC), Referring Expression Segmentation (RES), and Reasoning Segmentation tasks, respectively. The results demonstrate that only a few selected localization heads are sufficient to accurately localize objects in the image based on the text query. Our method effectively localizes objects in various scenarios.

## F. Applications

### F.1. Real World Application

Fig. 18 illustrates that the localization heads effectively capture the region or object of interest in images from the real world, based on the provided expressions. This result demonstrates the robustness of the localization heads across various types of data.

### F.2. Image Editing

Fig. 19 presents the results of image inpainting performed by integrating Stable Diffusion XL (SDXL) [43]. The frozen LVLM generates a segmentation mask corresponding to the expression, and this mask, along with an additional text prompt, is used as input to the diffusion model to generate the desired image. These results demonstrate that the segmentation mask corresponding to the referred text, output by a small number of localization heads from the frozen LVLM, can serve as guidance for diffusion models. This compatibility enables its application in image editing tasks.

## G. Limitations

We propose a simple yet effective framework for training-free visual grounding, which leverages the localization heads of LVLMs. Our framework successfully localizes objects in images based on text queries without requiring any fine-tuning and achieves superior performance compared to existing training-free methods. However, our method still has some limitations that could be addressed in future work.

First, our work, as illustrated in Fig. 10, reveals the potential for multi-object grounding; however, the establishment of a formalized pipeline or the development of a more streamlined implementation remains limited. The task of rendering the identified localization head more practical, user-friendly, and adaptable across a diverse range of applications continues to pose a significant challenge. This presents a compelling avenue for future research.



Figure 10. Multi-object segmentation results from the localization heads of DeepSeekVL-7B, along with the corresponding raw attention maps.

Second, our method is less suitable for LVLMs or methods that do not preserve spatial information in images (*e.g.*, pooling) [1, 19, 27, 28, 55]. These methods make it challenging to explicitly obtain image attention maps. To collect the attention map, a reverse computation is required to determine the order in which image tokens were input during processing. We leave the application of our framework to these methods for future exploration.

Figure 11. Average $S_{\text{img}}^{\ell,h}$ values for each attention head in more LVLMs. $\tau$ is set at the point where the maximum curvature is observed.



Figure 12. Selection frequency of individual heads and scatter plot of selection frequency rank versus each head's average IoU in more LVLMs. The Spearman correlation coefficient ($\rho$) between the selection frequency rank and the average IoU is displayed in the top-right corner of each plot. The observed Spearman correlation are statistically significant ($p < 0.001$) for all LVLMs.

17

Figure 13. Selection frequency of individual heads across different $\tau$ values. $\tau$ represents the threshold for the sum of each head's attention map. Our analy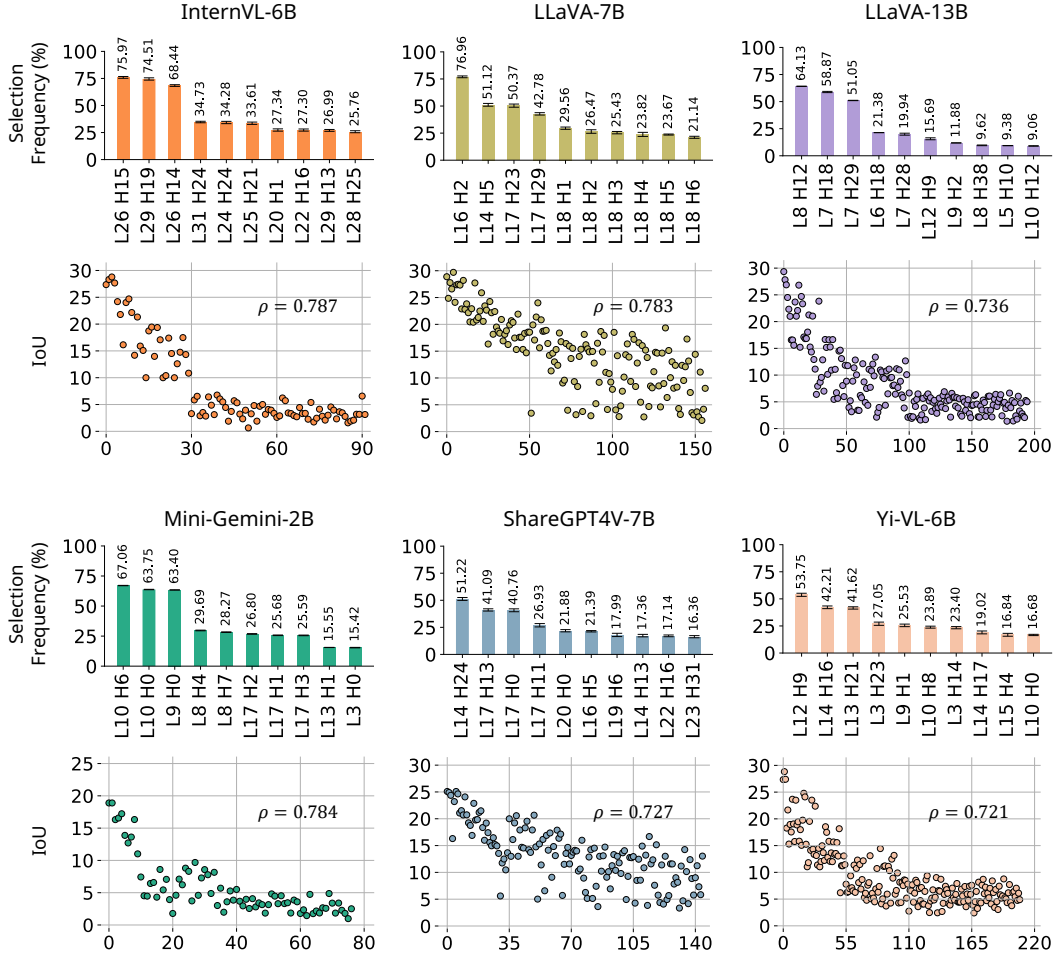sis focuses on heads with attention map sums greater than $\tau$, which are selected as targets for selection frequency evaluation. In the main paper, we select the threshold where the maximum curvature is observed. The top-3 localization heads remain consistent across different $\tau$ values, demonstrating the robustness of our analysis to variations in $\tau$.

18

Figure 14. Selection frequency of individual heads across different $N$ values. $N$ refers to the number of selected heads based on the lowest $H(\boldsymbol{A}^{\ell,h})$ values. Default setting is $N = 10$. The top-3 localization heads are consistent across different $N$ values, indicating the robustness of localization head selection to the choice of $N$.

Figure 15. Qualitative results of Referring Expression Comprehension.

Figure 16. Qualitative results of Referring Expression Segmentation.

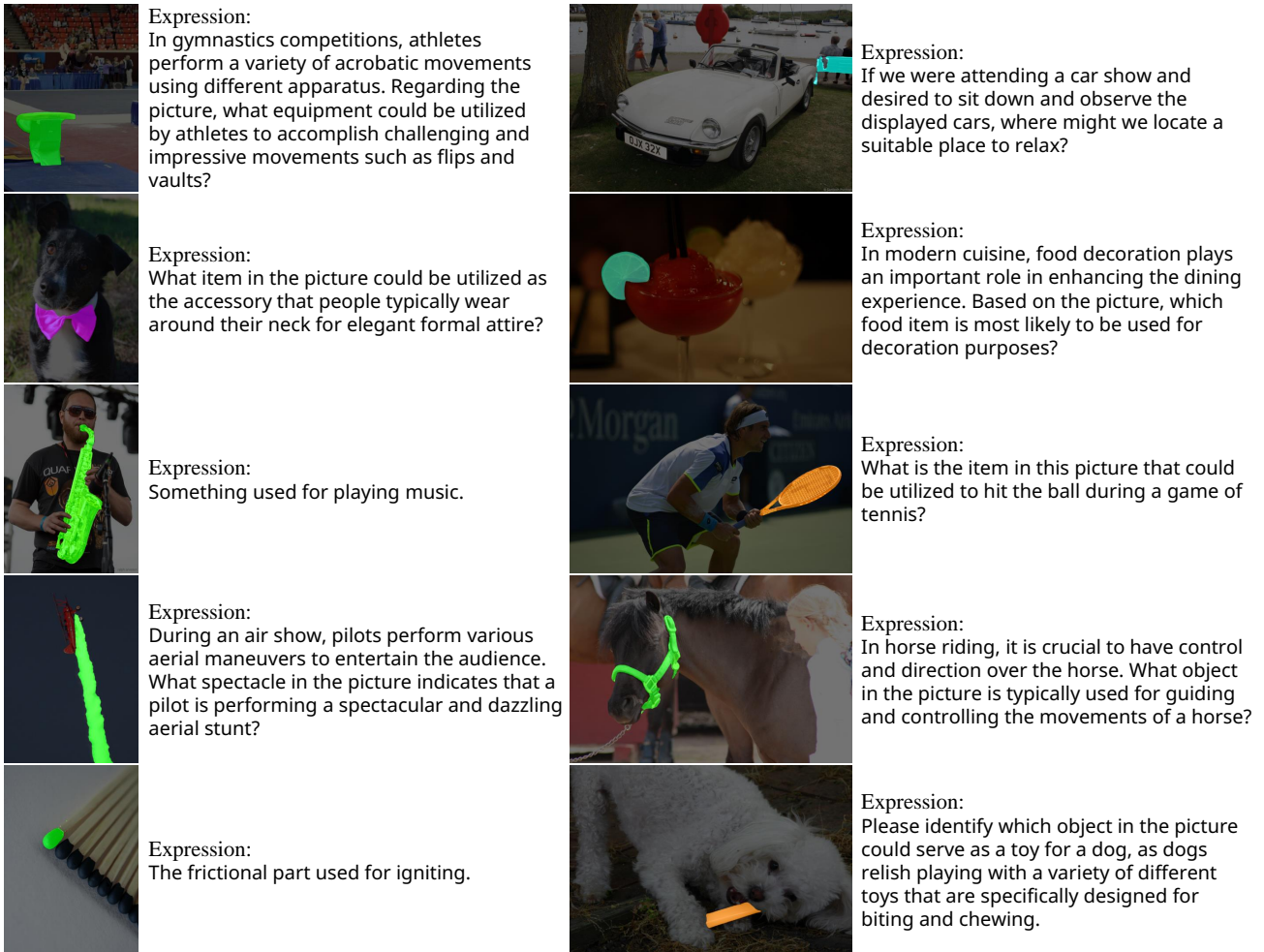LLaVA-1.5-13B with only three attention heads (L15 H39, L16 H30, L7 H2)



Expression:
In gymnastics competitions, athletes perform a variety of acrobatic movements using different apparatus. Regarding the picture, what equipment could be utilized by athletes to accomplish challenging and impressive movements such as flips and vaults?

Expression:
If we were attending a car show and desired to sit down and observe the displayed cars, where might we locate a suitable place to relax?

Expression:
What item in the picture could be utilized as the accessory that people typically wear around their neck for elegant formal attire?

Expression:
In modern cuisine, food decoration plays an important role in enhancing the dining experience. Based on the picture, which food item is most likely to be used for decoration purposes?

Expression:
Something used for playing music.

Expression:
What is the item in this picture that could be utilized to hit the ball during a game of tennis?

Expression:
During an air show, pilots perform various aerial maneuvers to entertain the audience. What spectacle in the picture indicates that a pilot is performing a spectacular and dazzling aerial stunt?

Expression:
In horse riding, it is crucial to have control and direction over the horse. What object in the picture is typically used for guiding and controlling the movements of a horse?

Expression:
The frictional part used for igniting.

Expression:
Please identify which object in the picture could serve as a toy for a dog, as dogs relish playing with a variety of different toys that are specifically designed for biting and chewing.

Figure 17. Qualitative results of Reasoning Segmentation.

LLaVA-1.5-13B with only three attention heads (L15 H39, L16 H30, L7 H2)



| | |
|---|---|
| Original Image | Expression: a father and the youngest. |



| | |
|---|---|
| Original Image | Expression: Ohtani Shohei. |



| | | |
|---|---|---|
| Original Image | Expression: Bruno Mars. | Expression: Blackpink Rose. |



| | | |
|---|---|---|
| Original Image | Expression: an Electric-type Pokémon. | Expression: a Pokémon Trainer. |

Figure 18. Qualitative results of real-world image segmentation. LLaVA-1.5-13B [35] uses only three attention heads (L15 H39, L16 H30, L7 H2) as localization heads to produce a precise segmentation masks related to the text expressions. The whitened regions in the images represent the segmentation mask output by the model.

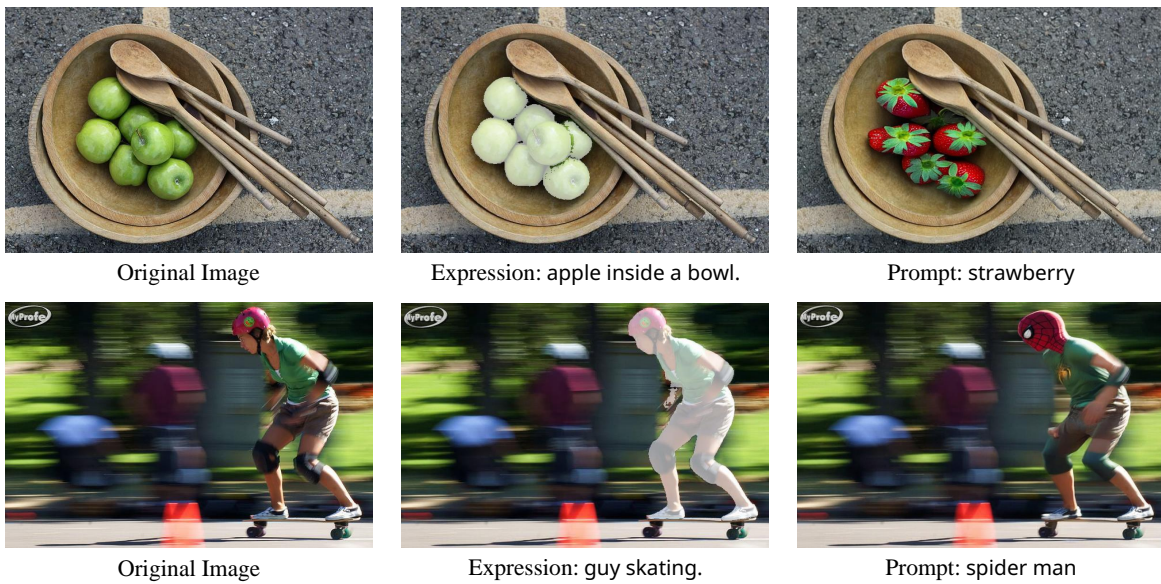| Original Image | Expression: apple inside a bowl. | Prompt: strawberry |
| Original Image | Expression: guy skating. | Prompt: spider man |

Figure 19. Qualitative results of generating the desired image through integration with a diffusion model [43]. Given an original image, our method generates a mask from the LVLM based on the text describing the desired modifications. This mask is then used as guidance for a diffusion model to perform image editing. Using the segmentation mask obtained through the localization head of the frozen LVLM [35], it is possible to generate semantic objects that align with the prompt at the specified mask locations.