

X-SAM: From Segment Anything to Any Segmentation

Hao Wang^{1,2}, Limeng Qiao³, Zequn Jie³, Zhijian Huang¹, Chengjian Feng³,
Qingfang Zheng², Lin Ma³, Xiangyuan Lan^{2*}, Xiaodan Liang^{1*}

¹Sun Yat-sen University, ²Peng Cheng Laboratory, ³Meituan Inc.

{wanghao9610, xdliang328}@gmail.com, lanxy@pcl.ac.cn

Abstract

Large Language Models (LLMs) demonstrate strong capabilities in broad knowledge representation, yet they are inherently deficient in pixel-level perceptual understanding. Although the Segment Anything Model (SAM) represents a significant advancement in visual-prompt-driven image segmentation, it exhibits notable limitations in multi-mask prediction and category-specific segmentation tasks, and it cannot integrate all segmentation tasks within a unified model architecture. To address these limitations, we present X-SAM, a streamlined Multimodal Large Language Model (MLLM) framework that extends the segmentation paradigm from *segment anything* to *any segmentation*. Specifically, we introduce a novel unified framework that enables more advanced pixel-level perceptual comprehension for MLLMs. Furthermore, we propose a new segmentation task, termed Visual GrounDed (VGD) segmentation, which segments all instance objects with interactive visual prompts and empowers MLLMs with visual grounded, pixel-wise interpretative capabilities. To enable effective training on diverse data sources, we present a unified training strategy that supports co-training across multiple datasets. Experimental results demonstrate that X-SAM achieves state-of-the-art performance on a wide range of image segmentation benchmarks, highlighting its efficiency for multimodal, pixel-level visual understanding. Code is available at <https://github.com/wanghao9610/X-SAM>.

1 Introduction

Multi-modal Large Language Models (MLLMs) have exhibited substantial advancements alongside the rapid development of Large Language Models (LLMs) (Bai et al. 2023; Touvron et al. 2023a,b; Abdin et al. 2024) and multi-modal pre-training methods (Radford et al. 2021; Jia et al. 2021; Zhai et al. 2023). These models have shown remarkable effectiveness in a wide range of applications, including image captioning (Xu et al. 2015), VQA (Antol et al. 2015), and visual editing (Chen et al. 2018). Nevertheless, a significant impediment to developing a truly generalized model remains: contemporary MLLMs are restricted to generating solely textual outputs. This limitation poses a considerable challenge in directly addressing tasks that require pixel-level comprehension of visual data, such as image segmentation, which is the most critical task in the field of computer vision.

*Corresponding author.

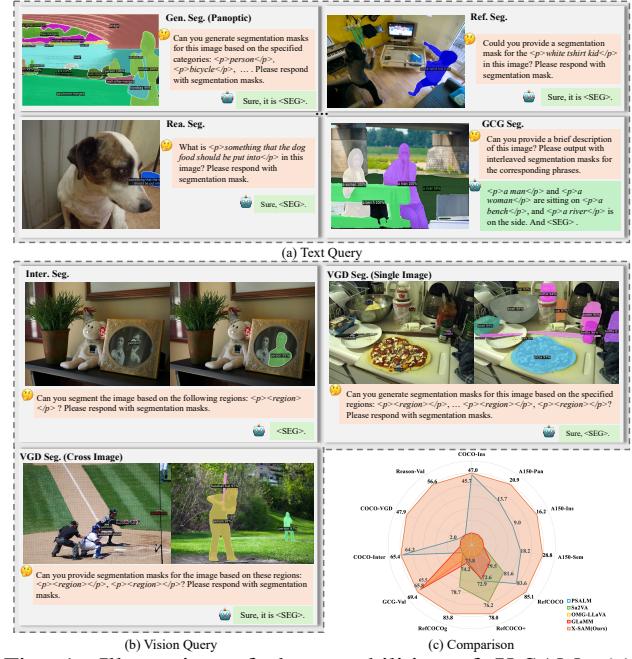


Fig. 1: Illustration of the capabilities of X-SAM. (a). Text query tasks: Generic(Gen.), Referring(Ref.), Reasoning(Rea.), and Grounded Conversation Generation(GCG) segmentation, etc.. (b). Vision query tasks: Interactive(Inter.) and Visual GrounDed(VGD) segmentation for single and cross-image. (c). X-SAM outperforms existing MLLMs on all segmentation benchmarks. Best viewed in zoom.

The Segment Anything Model (SAM) represents a foundational segmentation model that demonstrates exceptional efficacy in generating dense segmentation masks and has inspired the development of various segmentation tasks, such as high-quality segmentation (Ke et al. 2023), matching anything (Li et al. 2024c), and tracking anything (Rajić et al. 2025). Nevertheless, SAM’s architecture is fundamentally constrained by its dependency on visual prompts, which significantly limits its direct applicability to a wide range of image segmentation tasks, including generic (semantic, instance, panoptic) segmentation, referring segmentation, and open-vocabulary (OV) segmentation, among others. Achieving a unified framework capable of addressing various image segmentation tasks remains a challenging problem.

In this work, we introduce X-SAM, an innovative framework that unifies diverse image segmentation tasks, expanding the segmentation paradigm from *segment anything* to *any segmentation*. To accomplish this objective, our approach addresses three critical technical challenges: (1) *Task formulation*: Transforming SAM into a versatile segmentation architecture with cross-task applicability. (2) *Modality enhancement*: Augmenting LLMs with multimodal input processing capabilities. (3) *Unified framework*: Developing a cohesive approach to effectively facilitate comprehensive segmentation applications across diverse domains.

First, we develop a unified segmentation MLLM architecture that incorporates a unified mask decoder capable of generating segmentation masks suitable for generalized image segmentation tasks. Second, we expand the multimodal capabilities of MLLMs to process not only textual queries but also visual queries. Specifically, we introduce a novel task termed Visual Grounded (VGD) segmentation, which segments all instance objects with interactive visual prompts in an image. This task introduces visual guide modalities into large language models (LLMs). Furthermore, we propose a unified input format and training methodology that reformulates segmentation tasks within a unified framework, thus optimizing the adaptation of MLLMs to diverse image segmentation tasks.

As shown in Fig. 1 and Tab. 1, we present the comprehensive capabilities of X-SAM and compare them with those of other methods. Our proposed framework exhibits capabilities in processing text query-based tasks, such as generic segmentation and referring segmentation, while simultaneously accommodating vision query-based tasks such as interactive segmentation (Zhang et al. 2024d) and our novel VGD segmentation, which functions effectively in both single-image and cross-image contexts. Furthermore, X-SAM leverages the reasoning and generative capacities of LLMs, thereby enabling advanced reasoning segmentation and Grounded Conversation Generation (GCG) (Rasheed et al. 2024) segmentation.

X-SAM undergoes co-training with a diverse range of datasets. We perform a comprehensive evaluation on more than twenty segmentation datasets across seven distinct image segmentation tasks, even including the image conversion task. X-SAM achieves the state-of-the-art performance across all image segmentation benchmarks, and establishes a robust new baseline for unified pixel-level image understanding. In summary, our contributions are as follows:

- We introduce X-SAM, a novel unified framework that extends the segmentation paradigm from *segment anything* to *any segmentation*. Our approach formulates diverse image segmentation tasks into a standardized segmentation format.
- We propose a new image segmentation benchmark, Visual GrounDed (VGD) Segmentation, which provides visual grounded prompts for MLLMs to segment instance objects in images. The benchmark introduces user-friendly inputs to ground the segmentation objects and guide the MLLMs to output the segmentation masks.
- We present a unified multi-stage training strategy to co-

train X-SAM with a diverse range of datasets, and conduct extensive evaluations on more than twenty image segmentation benchmarks, achieving state-of-the-art performance on all of them. This establishes a new strong baseline for unified pixel-level perceptual understanding in MLLMs.

2 Related Work

Multi-modal Large Language Model. Multi-modal learning has evolved from early models focused on task-specific fusion and feature extraction (Li et al. 2022b), to leveraging large language models (Brown et al. 2020; Touvron et al. 2023a,b) for generalized, instruction-tuned multi-task benchmarks (Li et al. 2023a; Liu et al. 2024c; Hudson et al. 2019). LLaVA (Liu et al. 2024a,b, 2023a) introduced visual feature tokenization, inspiring advances in visual representation (Yuan et al. 2024b), specialized vision extensions (Lai et al. 2024; Lin et al. 2023; Dong et al. 2024a,b; Zhang et al. 2023; Ren et al. 2024; Zhang et al. 2024a; Zang et al. 2025), and language-guided segmentation (Li et al. 2024e; Zhang et al. 2024b). However, most progress remains task-specific. To our knowledge, we are the first to successfully implement a comprehensive approach, opening new directions for image segmentation.

Multi-modal Grounded Segmentation. Recent works (Pan et al. 2024; Sun et al. 2024; Zhou et al. 2022a; Bar et al. 2022; Wang et al. 2023a,b) explore visual initiation methods in vision, including learnable tokens (Zhou et al. 2022a), mask-visual-modeling (Wang et al. 2024; Fang et al. 2023; Wang et al. 2023b), and visual prompting encoders (Yuan et al. 2024a; Wang et al. 2023c,a). SAM (Kirillov et al. 2023b) and its extensions (Xu et al. 2024; Yuan et al. 2024a) introduce visual grounding signals to segmentation models, greatly improving performance. Interactive segmentation (Li et al. 2024e) further enhances user-guided segmentation for MLLMs. However, existing methods cannot freely treat grounded input as textual input for segmentation. To address this, we propose Visual GrounDed (VGD) segmentation, enabling more diverse multi-modal grounded segmentation.

Unified Segmentation Model. Vision transformers (Mehta et al. 2021; Dosovitskiy et al. 2020; Carion et al. 2020) have advanced universal segmentation, with recent works (Zhou et al. 2022b; Xu et al. 2024; Li et al. 2024d; Zhou et al. 2024; Sun et al. 2023; Yang et al. 2021; Cheng et al. 2022b; Wang et al. 2021a) developing end-to-end mask classification frameworks that outperform earlier models (Zhou et al. 2022c; Li et al. 2021, 2020; Chen et al. 2019) across various applications. Research has expanded to open-world and open-vocabulary segmentation (Wu et al. 2024; Yuan et al. 2024a; Qi et al. 2022a,b), as well as unified architectures for multiple tasks (Athar et al. 2023; Gu et al. 2023; Yan et al. 2023b; Xu et al. 2024; Jain et al. 2023; Li et al. 2024e). However, most methods focus solely on visual segmentation and lack interactive textual and visual prompts found in MLLMs. To address this, we combine SAM with MLLMs, extending SAM from *segment anything* to *any segmentation*, and introduce a unified framework adaptable to all image segmentation tasks, establishing a new strong baseline.

Tab. 1: Comparison of Capability. We compare different methods on both segmentation-specific (Gray) and MLLM-based.

| Method | Text Query | | | | | | | Vision Query | |
|---------------------------------|------------|---------|-----------|-----------|----------|-------------|----------|--------------|---|
| | Gen. Seg. | OV Seg. | Ref. Seg. | Rea. Seg. | GCG Seg. | Inter. Seg. | VGD Seg. | | |
| SAM(Kirillov et al. 2023a) | ✓ | | | | | | ✓ | | |
| Mask2Former(Cheng et al. 2022a) | ✓ | | | | | | | | |
| ODISE(Xu et al. 2023) | ✓ | ✓ | | | | | | | |
| UNINEXT(Yan et al. 2023a) | ✓ | | | ✓ | | | | ✓ | |
| SEEM(Zou et al. 2023b) | ✓ | ✓ | | ✓ | | | | ✓ | |
| OMG-Seg(Li et al. 2024e) | ✓ | ✓ | | | | | ✓ | | |
| LISA(Lai et al. 2024) | | | | ✓ | ✓ | | | | |
| GLaMM(Rasheed et al. 2024) | | | | ✓ | | | ✓ | | |
| PixelLM(Zhang et al. 2024d) | | | | ✓ | | | | | |
| OMG-LLaVA(Zhang et al. 2024c) | ✓ | | | ✓ | | | ✓ | | |
| Sa2VA(Yuan et al. 2025) | | | | ✓ | ✓ | | ✓ | | |
| PSALM(Zhang et al. 2024d) | ✓ | ✓ | | ✓ | | | ✓ | | |
| X-SAM (Ours) | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ |

3 Method

To achieve unified image segmentation, we present X-SAM, a novel multi-modal segmentation MLLM. We design a versatile input format and a unified framework to integrate diverse segmentation tasks into a single model. Additionally, we introduce an innovative training strategy that enables SAM to handle any segmentation task. The following sections detail our methodology.

3.1 Formulation

The development of a unified segmentation model is fraught with challenges stemming from the diverse nature of segmentation tasks and the variability in input format. To address these issues, we introduce a versatile input format tailored to support a wide range of image segmentation tasks, laying the groundwork for the unified framework of X-SAM. We delineate the input format into two primary categories: text query input and vision query input. The text query input consists exclusively of linguistic prompts derived from user requests, the vision query input integrates both linguistic prompts and visual prompts provided by the user.

Text Query Input. The majority of existing image segmentation tasks can be conceptualized as text query inputs, including generic segmentation (Kirillov et al. 2019), referring segmentation, open-vocabulary (OV) segmentation (Li et al. 2022a), GCG segmentation (Rasheed et al. 2024), and reasoning segmentation (Lai et al. 2024). A text query input encapsulates the user’s request along with the specific category or object to be segmented, which may be embedded within the user’s prompt or generated by a large language model (LLM). To facilitate the GCG segmentation task, inspired by GLaMM (Rasheed et al. 2024), we incorporate two special phrase tokens, $\langle p \rangle$ and $\langle /p \rangle$, into the tokenizer to denote the beginning and end of a phrase, respectively. For each category in generic segmentation and GCG segmentation, phrase in referring segmentation, or sentence in reasoning segmentation, the format is standardized as “ $\langle p \rangle$ category/phrase/sentence $\langle /p \rangle$ ”. Specifically, the $\langle p \rangle$

and $\langle /p \rangle$ tokens are not only encoded in the input tokens but also generated in the output tokens, ensuring consistency across different tasks. Additionally, for the output, we introduce a special token $\langle \text{SEG} \rangle$ into the tokenizer to signify the segmentation result, following the approach in (Lai et al. 2024).

Vision Query Input. Beyond text query inputs, some tasks necessitate vision query inputs, such as interactive segmentation (Zhang et al. 2024d) and the Visual Grounded segmentation proposed in this work. In contrast to text query inputs, vision query inputs incorporate a visual prompt from the user, which may take the form of points, scribbles, boxes, or masks. To denote the visual prompt, we employ a dedicated token, $\langle \text{region} \rangle$, within the input format. Analogous to the text query input, the visual prompt is formatted as “ $\langle p \rangle \langle \text{region} \rangle \langle /p \rangle$ ” and the segmentation output is similarly indicated by the $\langle \text{SEG} \rangle$ token. The $\langle \text{region} \rangle$ token serves as a placeholder for the visual prompt and will be replaced by the region feature extracted from the segmentation encoder.

Unified Formulation. The latent language embeddings between the $\langle p \rangle$ and $\langle /p \rangle$ tokens are used as the condition embeddings for the segmentation decoder to compute the classification scores. Based on this formulation, we achieve a unified framework for all image segmentation tasks. Given an input image $\mathbf{X}_v \in \mathbb{R}^{H \times W \times 3}$ and a language instruction $\mathbf{X}_q \in \mathbb{R}^{P \times 1}$, the model takes the image and language instruction as inputs and outputs a language response $\mathbf{Y}_q \in \mathbb{R}^{L \times 1}$ and a segmentation mask $\mathbf{Y}_m \in \mathbb{R}^{H \times W}$. Here, P is the length of the input text tokens, and L is the total length of the input and output text tokens. H and W denote the height and width of the image, respectively. Detailed input format examples can be found in Fig. 1 (a) and (b).

3.2 Architecture

In this section, we propose X-SAM, a unified segmentation MLLM for any segmentation. As shown in Fig. 2, it includes dual encoders, dual projectors, an LLM, a segmentation connector, and a segmentation decoder.

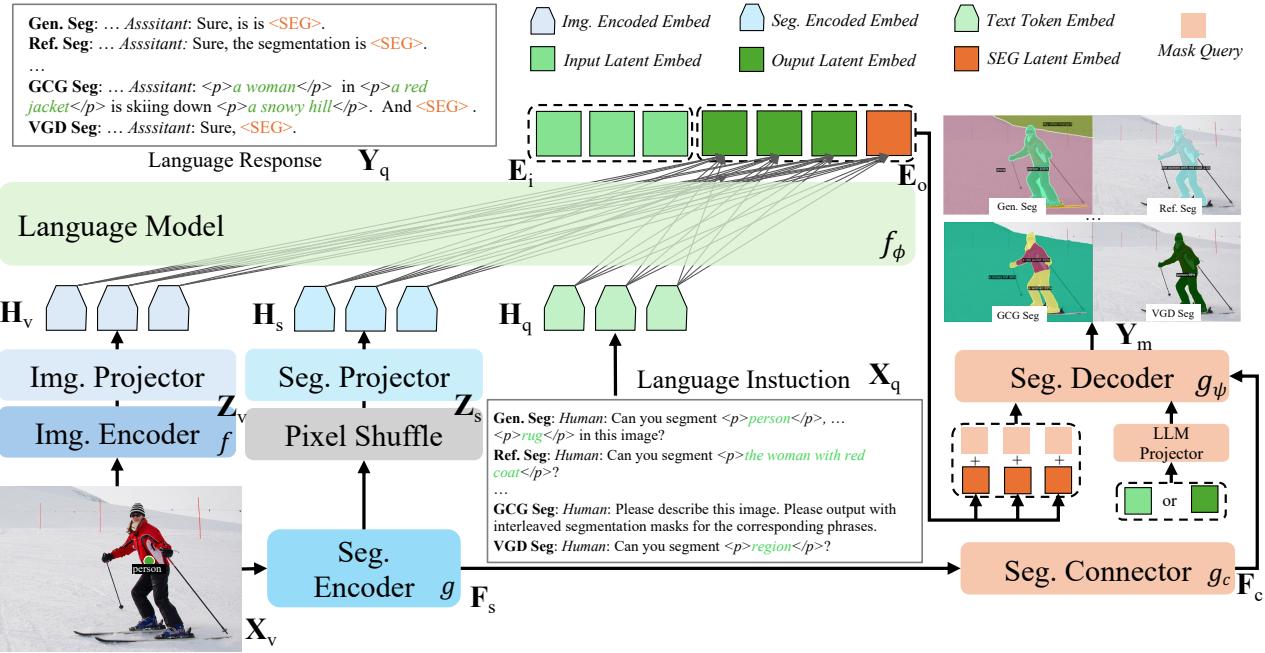


Fig. 2: The Overview of X-SAM. X-SAM comprises dual encoders, dual projectors, a language model, a segmentation connector, and a segmentation decoder. The dual encoders process the image and project features to match text embedding dimensions, which are then input to the language model with tokenized text for instruction-guided understanding. The SAM features are connected to the segmentation decoder, which uses the LLM’s <SEG> token to generate segmentation masks.

Dual Encoders. There are two encoders in X-SAM, an image encoder and a segmentation encoder. The image encoder f is used to extract the global image feature $Z_v = f(X_v)$, while the segmentation encoder g extracts the fine-grained image feature $Z_s = g(X_v)$. The feature from the image encoder is global and benefits image understanding tasks, whereas the feature from the segmentation encoder is fine-grained and benefits image segmentation tasks. We adopt SigLIP2-so400m (Tschannen et al. 2025) as the image encoder and SAM-L (Ke et al. 2023) as the segmentation encoder.

Dual Projectors. To enhance the LLM’s understanding of the image, we concatenate the features from the image encoder and the segmentation encoder before passing them to the LLM. Specifically, the feature from the segmentation encoder is too large to be processed directly by the LLM, so we utilize a pixel-shuffle operation to reduce its spatial size. We then project the reduced feature into the language embedding space H_q via an MLP projector W_s . For the feature from the image encoder, we directly project it into the language embedding space via an MLP projector W_i , such that $H_v = W_i \cdot Z_v$ and $H_q = W_s \cdot Z_s$. We then concatenate the features from dual projectors and the language embeddings, and input them into the LLM f_ϕ .

Segmentation Connector. For image segmentation tasks, fine-grained multi-scale features are crucial for the segmentation decoder to accurately predict segmentation masks. The output of the segmentation encoder in SAM is single-scale (1/16) with reduced spatial resolution. To obtain multi-scale features, we design a segmentation connector g_c , to bridge the segmentation encoder and decoder. As shown in

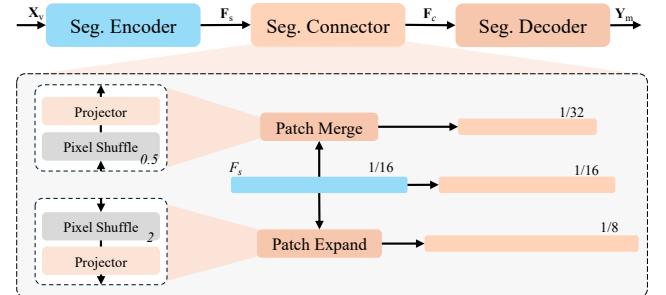


Fig. 3: The Architecture of Segmentation Connector.

Fig. 3, we perform patch-merge using a pixel-shuffle (Chen et al. 2024) with a scale of 0.5 to reduce the spatial size of the last feature in the encoder to a smaller scale (1/32). We also perform patch-expand with a pixel-shuffle of scale 2.0 to increase the spatial size of the last feature to a larger scale (1/8), resulting in multi-scale features for the segmentation decoder.

Segmentation Decoder. The Segment Anything Model (SAM) can segment a single object based on input text or visual prompts, but it fails to segment all objects in a single inference. To segment all objects at once, we replace its original segmentation decoder with a new decoder, following the approach in (Cheng et al. 2022a; VS et al. 2024). The segmentation decoder g_p predicts masks and their category probabilities from either the input latent embedding E_i or the output latent embedding E_o , multi-scale segmentation features \mathbf{F}_c , and a set of mask query tokens plus the <SEG>

Tab. 2: Comprehensive Performance Comparison. We compare X-SAM to segmentation-specific models (Gray) and MLLMs. “ \times ” denotes unsupported tasks. “–” indicates unreported results. X-SAM achieves state-of-the-art performance across all segmentation tasks with a single model. Best results are in **bold**, second-best are underlined.

| Method | Gen. Seg. | OV Seg. | Ref. Seg. | Rea. Seg. | GCG Seg. | Inter. Seg. | VGD Seg. |
|-----------------------------------|---------------------------|---------------------------|---------------------------|--------------------|--------------------|--------------------|--------------------|
| | Pan. / Ins. / Sem. | Pan. / Ins. / Sem. | RefCOCO / + / g | Val / Test | Val / Test | Point / Box | Point / Box |
| SAM-L(Kirillov et al. 2023a) | \times | \times | \times | \times | \times | 51.8 / 76.6 | \times |
| Mask2Former-L(Cheng et al. 2022a) | 57.8 / 48.6 / 67.4 | \times | \times | \times | \times | \times | \times |
| SEEM-B(Zou et al. 2023b) | 56.1 / 46.4 / 66.3 | \times | - / - / 65.6 | \times | \times | 47.8 / 44.9 | \times |
| ODISE(Xu et al. 2023) | 55.4 / 46.0 / 65.2 | 22.6 / 14.4 / 29.9 | \times | \times | \times | \times | \times |
| OMG-Seg(Li et al. 2024e) | 53.8 / - / - | \times | \times | \times | \times | \times | \times |
| LISA-7B(Lai et al. 2024) | \times | \times | 74.9 / 65.1 / 67.9 | <u>52.9 / 47.3</u> | \times | \times | \times |
| GLaMM(Rasheed et al. 2024) | \times | \times | 79.5 / 72.6 / 74.2 | \times | <u>65.8 / 64.6</u> | \times | \times |
| PixelLM-7B(Ren et al. 2024) | \times | \times | 73.0 / 66.3 / 69.3 | \times | \times | \times | \times |
| OMG-LLA-VA-7B(Zhang et al. 2024c) | 53.8 / - / - | \times | 78.0 / 69.1 / 72.9 | \times | <u>65.5 / 64.7</u> | \times | \times |
| Sa2VA-8B(Yuan et al. 2025) | \times | \times | 81.6 / <u>76.2 / 78.7</u> | - / - | - / - | \times | \times |
| PSALM(Zhang et al. 2024d) | 55.9 / 45.7 / 66.6 | <u>13.7 / 9.0 / 18.2</u> | <u>83.6 / 72.9 / 73.8</u> | \times | \times | 64.3 / 67.3 | 2.0 / 3.7 |
| X-SAM (Ours) | <u>54.7 / 47.0 / 66.5</u> | 20.9 / 16.2 / 28.8 | 85.1 / 78.0 / 83.8 | 56.6 / 57.8 | 69.4 / 69.0 | 65.4 / 70.0 | 47.9 / 49.5 |

token embedding, which bridges the LLM output with the segmentation decoder. Notably, we introduce a latent background embedding to represent the “ignore” category for all tasks, thereby unifying all image segmentation tasks with one model.

3.3 Training

To improve the performance on diverse image segmentation tasks, we propose a novel multi-stage training strategy. The training strategy consists of three stages: segmentor fine-tuning, alignment pre-training, and mixed fine-tuning.

Stage 1: Segmentor Fine-tuning. As the segmentation decoder is redesigned, we need to train the segmentor to adapt to segment all objects in a single forward pass. We follow the training pipeline in (Cheng et al. 2022a), which trains the model on the popular COCO-Panoptic (Kirillov et al. 2019) dataset. To enable faster convergence during training, we unfreeze all the parameters in the segmentor while training the segmentation encoder with a lower learning rate. The training objective, \mathcal{L}_{seg} , is the same as in (Cheng et al. 2022a), and is defined as the sum of the classification loss \mathcal{L}_{cls} , the mask loss $\mathcal{L}_{\text{mask}}$, and the dice loss $\mathcal{L}_{\text{dice}}$:

$$\mathcal{L}_{\text{seg}} = \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{mask}} + \mathcal{L}_{\text{dice}} \quad (1)$$

Stage 2: Alignment Pre-training. To align the language embeddings and visual embeddings, we perform alignment pre-training on the LLaVA-558K dataset, following (Liu et al. 2023b). We keep the dual encoders and the LLM parameters frozen, and only train the dual projectors. In this way, the image embeddings and segmentation embeddings can be aligned with the pre-trained LLM word embeddings. The training objective for alignment pre-training is an auto-regressive loss $\mathcal{L}_{\text{regressive}}$:

$$\mathcal{L}_{\text{regressive}} = - \sum_{i=1}^N \log p_\theta \left(\mathcal{Y}_q^{[P+i]} | \mathcal{Y}_q^{[:i-1]}, \mathcal{X}_q^{[:i-1]} \right), \quad (2)$$

where \mathcal{X}_q is the input sequence $\mathcal{X}_q = [x_1, x_2, \dots, x_p] \in \mathbb{R}^{P \times D}$, \mathcal{Y}_q is the output sequence $\mathcal{Y}_q = [y_1, y_2, \dots, y_l] \in$

$\mathbb{R}^{L \times D}$, where $L = P + N$ represents the length of output sequence, D represents the hidden size of LLM. θ is a trainable parameter in LLM, and we only calculate the loss for the generated text.

Stage 3: Mixed Fine-tuning. X-SAM is co-trained on multiple datasets across diverse tasks in an end-to-end manner. For the image conversation task, we adopt the auto-regressive loss $\mathcal{L}_{\text{regressive}}$ as is common in MLLM training. For the segmentation tasks, we not only use the segmentation loss as in segmentor training, but also add the auto-regressive loss to the training objective. Benefiting from the unified formulation and simple training objective, end-to-end mixed fine-tuning across diverse tasks can be performed within a unified framework. The training objective for mixed fine-tuning can be formulated as:

$$\mathcal{L}_{\text{total}} = \begin{cases} \mathcal{L}_{\text{regressive}}, & \text{conversation} \\ \mathcal{L}_{\text{regressive}} + \mathcal{L}_{\text{seg}}, & \text{segmentation} \end{cases} \quad (3)$$

4 Experiments

4.1 Experiment Settings

Datasets and Tasks. For segmentor fine-tuning, we train on the COCO-Panoptic (Kirillov et al. 2019) dataset. For alignment pre-training, we utilize the LLaVA-558K (Liu et al. 2023b) dataset. For end-to-end mixed fine-tuning, we incorporate one image conversation dataset and five types of image segmentation datasets into the training process. To balance the training data across these diverse datasets, we set the training epoch to 1 and adjust the resampling rates of different datasets using dataset balance resampling. After training, X-SAM is capable of performing a variety of tasks, including Image Conversation, Generic, Referring, Reasoning, GCG, Interactive, and VGD Segmentation. Additionally, X-SAM supports Open-Vocabulary (OV) (OV-semantic, OV-instance, OV-panoptic) segmentation, enabling it to segment all objects defined by the input prompt, even those never seen before. Note that COCO-VGD is our proposed VGD segmentation dataset, which is built on the COCO2017 dataset. Details of the datasets are presented in Appendix A.1.

Tab. 3: Comparison of Referring Segmentation. We evaluate methods on referring segmentation benchmarks by (M)LLMs.

| Method | (M)LLM | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|--------------------------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | val | testA | testB | val | testA | testB | val | test |
| SEEM-L (Zou et al. 2023b) | - | - | - | - | - | - | - | 65.6 | - |
| UNINEXT-L (Yan et al. 2023a) | - | 80.3 | 82.6 | 77.8 | 70.0 | 74.9 | 62.6 | 73.4 | 73.7 |
| UNINEXT-H (Yan et al. 2023a) | - | 82.2 | 83.4 | 81.3 | 72.5 | 76.4 | 66.2 | 74.7 | 76.4 |
| GLaMM (Rasheed et al. 2024) | Vicuna-7B | 79.5 | 83.2 | 76.9 | 72.6 | 78.7 | 64.6 | 74.2 | 74.9 |
| OMG-LLaVA (Zhang et al. 2024c) | InternLM-7B | 77.2 | 79.8 | 74.1 | 68.7 | 73.0 | 61.6 | 71.7 | 71.9 |
| Sa2VA(Yuan et al. 2025) | InternVL2-8B | 81.6 | - | - | 76.2 | - | - | 78.7 | - |
| PSALM (Zhang et al. 2024d) | Phi-1.5-1.3B | 83.6 | 84.7 | 81.6 | 72.9 | 75.5 | 70.1 | 73.8 | 74.4 |
| X-SAM (Ours) | Phi-3-3.8B | 85.1 | 87.1 | 83.4 | 78.0 | 81.0 | 74.4 | 83.8 | 83.9 |

Tab. 4: Comparison of GCG Segmentation. † indicates pretraining with the Grand dataset (Rasheed et al. 2024).

| Methods | Val | | | | Test | | | |
|---|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | METEOR | CIDEr | AP50 | mIoU | METEOR | CIDEr | AP50 | mIoU |
| Kosmos-2(Peng et al. 2023) | 16.1 | 27.6 | 17.1 | 55.6 | 15.8 | 27.2 | 17.2 | 56.8 |
| LISA-7B(Lai et al. 2024) | 13.0 | 33.9 | 25.2 | 62.0 | 12.9 | 32.2 | 24.8 | 61.7 |
| GLaMM-7B [†] (Rasheed et al. 2024) | 15.2 | 43.1 | 28.9 | 65.8 | 14.6 | 37.9 | 27.2 | 64.6 |
| OMG-LLaVA-7B(Zhang et al. 2024c) | 14.9 | 41.2 | 29.9 | 65.5 | 14.5 | 38.5 | 28.6 | 64.7 |
| X-SAM (Ours) | 15.4 | 46.3 | 33.2 | 69.4 | 15.1 | 42.7 | 32.9 | 69.0 |

Evaluation Metrics. We conduct extensive experiments to evaluate the performance of X-SAM. For generic segmentation and open-vocabulary segmentation, we use PQ, mIoU, and mAP as the main metrics for panoptic, semantic, and instance segmentation, respectively. For referring segmentation and reasoning segmentation, we adopt cloU and gloU as metrics, following (Zhang et al. 2024d). For GCG segmentation, we use M, C, AP50, and mIoU as metrics, following (Rasheed et al. 2024). For interactive segmentation, we use mIoU and cloU, also following (Zhang et al. 2024d). For VGD segmentation, we use AP and AP50. For image conversation, we adopt scores from common MLLM benchmarks as the main metrics, following (Liu et al. 2023b).

Implementation Details. We adopt the XTuner (Contributors 2023) codebase for training and evaluation. During segmentor fine-tuning, we train all parameters, set the batch size to 64, and use a learning rate of 1e-5 for the SAM encoder and 1e-4 for the other parameters. The number of training epochs is set to 36. For alignment pre-training, we train only the dual projector parameters, with a batch size of 256, a learning rate of 1e-3, and one training epoch. For end-to-end mixed fine-tuning, we train all parameters, set the batch size to 64, and use a learning rate of 4e-6 for the dual encoders and 4e-5 for the other parameters, with one training epoch. All training is conducted on 16 A100 GPUs. For image conversation evaluation, we use the VLMEvalKit (Duan et al. 2024) codebase to evaluate performance on MLLM benchmarks. For segmentation task evaluation, we follow the settings described in the corresponding papers and repositories. More implementation details are provided in Appendix A.3.

4.2 Main Results

We conduct extensive evaluation on seven segmentation tasks, including Generic, Open-Vocabulary, Referring, Rea-

soning, GCG, Interactive, and VGD Segmentation.

Overall. In Tab. 2, we compare X-SAM with current segmentation-specific models and MLLMs. X-SAM demonstrates the most comprehensive capabilities. It achieves performance comparable to state-of-the-art in generic segmentation, and achieves the best performance on other benchmarks, with a single model. X-SAM sets a new state-of-the-art record for image segmentation benchmarks. Detailed results for each task are discussed below.

Referring Segmentation. We evaluate X-SAM on RefCOCO, RefCOCO+, and RefCOCOg, with the results shown in Tab. 3. X-SAM outperforms PSALM (Zhang et al. 2024d) by 1.5% cloU, 5.1% cloU, and 10.0% cloU on the validation sets of RefCOCO, RefCOCO+, and RefCOCOg, respectively. Compared to Sa2VA-8B (Yuan et al. 2025), X-SAM achieves better results with a smaller model size. It shows performance improvements of 3.5% cloU, 1.8% cloU, and 5.1% cloU on RefCOCO, RefCOCO+, and RefCOCOg, respectively.

GCG Segmentation. Grounded conversation generation demands detailed image and pixel-level understanding, requiring MLLMs to link captioned objects to their segmentation masks. As shown in Tab. 4, X-SAM achieves a significant performance improvement compared to previous methods and obtains the best results on both the Val and Test sets. In terms of image-level understanding, X-SAM outperforms GLaMM (Rasheed et al. 2024) by 0.2% METEOR and 3.2% CIDEr on the Val set, and by 0.5% METEOR and 4.8% CIDEr on the Test set. In terms of pixel-level understanding, X-SAM outperforms OMG-LLaVA (Zhang et al. 2024c) by 3.3% AP and 3.9% mIoU on the Val set, and by 4.3% AP and 4.3% mIoU on the Test set.

VGD Segmentation. Visual grounded segmentation demands vision query understanding, requiring MLLMs to

Tab. 5: Comparison of VGD Segmentation. † indicates evaluation results following X-SAM setting.

| Method | Point | | Scribble | | Box | | Mask | |
|-----------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | AP | AP50 | AP | AP50 | AP | AP50 | AP | AP50 |
| PSALM† (Zhang et al. 2024d) | 2.0 | 3.3 | 2.8 | 4.4 | 3.7 | 5.8 | 2.3 | 3.3 |
| X-SAM (Ours) | 47.9 | 72.5 | 48.7 | 73.4 | 49.5 | 74.7 | 49.7 | 74.9 |

Tab. 6: Ablation on Fine-Tuning(FT).

| FT | COCO-Pan | A150-OV | RefCOCO | Reason-Val |
|----------|-------------|-------------|-------------|-------------|
| | PQ | AP | cIoU | gIoU |
| Specific | 55.3 | 16.4 | 81.0 | 48.2 |
| Mixed | 54.5(↓ 0.8) | 22.4(↑ 6.0) | 85.4(↑ 4.4) | 57.1(↑ 8.9) |

Tab. 7: Ablation on Dual Encoders. Swin† is initialized from Mask2Former (M2F) (Cheng et al. 2022a).

| Encoder | COCO-Pan | A150-OV | GCG-Val | COCO-VGD |
|-----------|-------------|-------------|-------------|-------------|
| | PQ | AP | mIoU | AP |
| ViT - | 54.5 | 16.4 | 64.8 | 40.7 |
| ViT Swin† | 56.2(↑ 1.7) | 18.6(↑ 2.2) | 62.5(↓ 2.3) | 48.6(↑ 7.9) |
| ViT SAM | 54.7(↑ 0.2) | 20.9(↑ 4.5) | 69.4(↑ 4.6) | 47.9(↑ 7.2) |

comprehend the visual modality and segment all related instances. Tab. 5 presents the VGD segmentation results. As VGD segmentation is our newly proposed task, we evaluate PSALM (Zhang et al. 2024d) following X-SAM’s settings. X-SAM outperforms PSALM by 45.9% AP, 45.9% AP, 45.8% AP, and 47.4% AP on Point, Scribble, Box, and Mask visual prompts, respectively.

More results and discussions for other segmentation and conversation benchmarks are provided in Appendix A.5.

4.3 Abaltions

We conduct ablation studies on mixed fine-tuning, dual encoders, multi-stage training, and segmentor architecture, presenting selected benchmark results due to space limitations.

Mixed Fine-tuning. We ablate the impact of mixed fine-tuning on X-SAM’s performance. As shown in Tab. 6, mixed fine-tuning improves performance on out-of-domain COCO benchmarks, demonstrating X-SAM’s robust segmentation capabilities—for example, a 6.0% AP increase on A150-OV and an 8.9% gIoU increase on Reason-Val. However, it results in a 0.8% PQ decrease on COCO-Pan due to the challenge of balancing performance in multisource training.

Dual Encoders. We ablate the design of the dual encoders in X-SAM. As shown in Tab. 7, dual encoders with either a SAM or Swin encoder benefit VGD segmentation, achieving 7.2% AP and 7.9% AP on COCO-VGD, respectively. Moreover, dual encoders with a SAM encoder consistently improve performance on GCG-Val and A150-OV, while the Swin encoder, which lacks robust segmentation capabilities, provides only a small improvement on A150-OV and even has a negative impact on GCG-Val.

Multi-stage Training. We ablate the impact of the multi-stage training strategy. As shown in Tab. 8, the S1 segmentor

Tab. 8: Ablation on Multi-Stage(M-Stage) Training. S1: Stage 1, S2: Stage 2, S3: Stage 3, Conv.: convolution.

| M-Stage | COCO-Pan | A150-OV | GCG-Val | Conv.-MMB |
|------------|-------------|-------------|-------------|-------------|
| | PQ | AP | mIoU | Acc |
| S3 | 45.2 | 19.4 | 60.6 | 67.2 |
| S1, S3 | 54.5(↑ 9.3) | 20.9(↑ 1.5) | 65.4(↑ 4.8) | 67.4(↑ 0.2) |
| S1, S2, S3 | 54.7(↑ 9.5) | 20.9(↑ 1.5) | 69.4(↑ 8.8) | 69.3(↑ 2.1) |

Tab. 9: Ablation on Segmentor Architecture. Conn.: connector, M-Scale: multi-scale, Con.: convolution, M2F: Mask2Former.

| Conn. | Decoder | M-Scale | PQ | AP | mIoU |
|-------|---------|---------|--------------|--------------|--------------|
| - | SAM | ✗ | 40.9 | 26.3 | 49.5 |
| MLP | M2F | ✗ | 50.1(↑ 9.2) | 38.9(↑ 12.6) | 60.2(↑ 10.7) |
| Con. | M2F | ✗ | 50.3(↑ 9.4) | 39.1(↑ 12.8) | 60.6(↑ 11.1) |
| Con. | M2F | ✓ | 51.6(↑ 10.7) | 41.5(↑ 15.2) | 61.6(↑ 12.1) |

fine-tuning phase boosts the segmentation capability, producing a notable improvement of 9.3% PQ in COCO-Pan and 1.5% AP in the A150-OV datasets. Meanwhile, the S2 alignment pre-training phase enhances image understanding capabilities, contributing an additional 2.1% Accuracy on Conv.-MMB. By integrating these stages, X-SAM demonstrates robust advances in image segmentation and comprehension, establishing its effectiveness in addressing complex visual tasks.

Segmentor Architecture. We ablate the impact of segmentor architecture by performing segmentor fine-tuning for 12 epochs. As shown in Tab. 9, M2F decoder brings a large improvement with 9.2% PQ as the effective design of M2F. The convolution connector performs better than the MLP connector, as the convolution spatial-awareness benefits segmentation, and multi-scale further improves the performance(10.7% PQ) with more diverse scale features.

More ablation results can be found in the Appendix A.6.

5 Conclusion

In this work, we propose X-SAM, a unified segmentation MLLM that extends the segmentation paradigm from *segment anything* to *any segmentation*, integrating all image segmentation tasks into a single model. Our method can process various multimodal inputs in MLLMs, including both text and visual queries. Moreover, to equip MLLMs with visual grounded perception capabilities, we introduce a new segmentation task, Visual GrounDed (VGD) segmentation, further extending the capabilities of the unified segmentation model. We conduct extensive experiments across all image segmentation tasks, and X-SAM achieves state-of-the-art performance on each task with a single model.

A Technical Appendices and Supplementary Material

In the Appendix, we first provide more details on the dataset, model architecture, and implementation of the proposed method. Then, we present additional experimental results on more benchmarks to demonstrate the effectiveness of our approach. Next, we include ablation studies on dataset balance resampling and the image encoder. Following that, we provide further visualization results for different tasks. Finally, we discuss the limitations and future work.

A.1 More Dataset Details

Training Datasets. In Tab. 10, we show the datasets used in our multi-stage training. For the segmentor fine-tuning stage, we fine-tune the segmentor on generic segmentation datasets. For the alignment pre-training stage, we pre-train the dual projectors on the alignment LLaVA 558K (Liu et al. 2023b) dataset. For the mixed fine-tuning stage, we fine-tune the whole model on mixed datasets, including both segmentation and conversation datasets. There are six types of datasets in total, including one image-level dataset and five segmentation datasets.

Building COCO-VGD Dataset. The COCO-VGD dataset is built on the images and annotations of the COCO2017 instance segmentation dataset, which provides instance-level segmentation masks for each object in an image. We automatically generate four types of visual prompts: point, scribe, box, and mask, for each instance in the image, following (Zhang et al. 2024d). We randomly sample one type of visual prompt for each category as the visually grounded prompt during training and evaluation.

Dataset Balance Resampling. As shown in Tab. 10, the dataset sizes vary, ranging from 0.2K to 665K. The data ratio among the datasets is crucial to the model’s performance during mixed fine-tuning. To balance the different dataset sizes, we propose a dataset balance resampling strategy following (Gupta et al. 2019). For each dataset d , let f_d be the frequency of dataset d in the mixed datasets. We define the dataset-level repeat factor as $r_d = \max(1, \sqrt{t/f_d})$, where t is a hyperparameter controlling the oversampling ratio. We then repeat dataset d with a repeat factor of r_d .

A.2 More Model Details

Model Framework. For the segmentor, we adopt SAM-L (Ke et al. 2023) as the segmentation encoder and the Mask2Former head (Cheng et al. 2022a) as the decoder. To reduce the number of connector parameters, we employ a bottleneck architecture, which first reduces the dimension of the segmentation feature to 512 via a 1×1 convolution, then further refines the feature via a 3×3 convolution, and finally expands the dimension to a value determined by the spatial scale of the pixel shuffle (Chen et al. 2024) operation using another 1×1 convolution. For the MLLM, we use SigLIP2-so400m (Tschannen et al. 2025) as the image encoder, an MLP as the image projector, another MLP with a pixel shuffle operation as the segmentation projector, and Phi-3-mini-4k-instruct (Abdin et al. 2024) as the LLM. The total number of parameters in X-SAM is about 5B.

Region Sampling. To sample region features from the vision query, we adopt the region sampling strategy from (You et al. 2023). Specifically, we first convert the vision query into a binary mask, then perform point sampling on the segmentor-encoded features to obtain the region features, and finally apply mean pooling to produce the final region features. These region features are placed in the corresponding position of `<region>` in the language instruction, serving as vision query categories for the segmentation decoder.

A.3 More Training Details

Stage1: Segmentor Fine-tuning. During segmentor fine-tuning, we unfreeze all parameters of the segmentor, including the SAM encoder, segmentation connector, and segmentation decoder. The learning rate is set to 1e-4 for all components except the segmentation encoder, which uses a learning rate of 1e-5. We set the batch size to 64 and train for 36 epochs. The SAM encoder is initialized with pre-trained weights, while the segmentation connector and decoder are initialized with random weights. Additionally, we apply random scale augmentation to the images during training, with a scale of [0.1, 2.0].

Stage2: Alignment Pre-training. During alignment pre-training, we train only the parameters of the dual projectors and keep all other parameters fixed. The learning rate is set to 1e-3, and the batch size is set to 256. The dual projectors are initialized with random weights, the segmentation encoder is initialized with the pre-trained weights from the segmentor fine-tuning stage, and the image encoder and LLM are initialized with their official pre-trained weights. Training is conducted for 1 epoch.

Stage3: Mixed Fine-tuning. During mixed fine-tuning, we fine-tune all parameters of the model. The learning rate for the dual encoders is set to 4e-6, while the learning rate for the other modules is set to 4e-5. The batch size is set to 64, and training is conducted for 1 epoch. The segmentation encoder, segmentation connector, and segmentation decoder are initialized with pre-trained weights from the segmentor fine-tuning stage, and the image encoder is initialized with official pre-trained weights. The dual projectors are initialized with pre-trained weights from the alignment pre-training stage. Additionally, to make training more stable, we ensure that all data within a global batch come from the same source.

The hyper-parameters in the multi-stage training are shown in Tab. 11. A simplified illustration of the multi-stage training is shown in Fig. 4.

A.4 More Evaluation Details

PSALM COCO-VGD Evaluation. PSALM (Zhang et al. 2024d) is a segmentation MLLM that supports generic segmentation, open-vocabulary segmentation, referring segmentation, interactive segmentation, and more. To evaluate its performance on our proposed COCO-VGD dataset, we follow the same evaluation process as for X-SAM. We randomly sample some instance annotations, as done with X-SAM, and then feed them to PSALM to obtain instance-level predictions. PSALM is trained on the COCO-Interactive dataset, which shares the same source as COCO-VGD but

Tab. 10: The Datasets Specification. The datasets used in multi-stage training. # Sample is the total number of samples in this task.

| Task | Datasets | # Sample |
|--|---|--|
| Generic Segmentation | Stage 1: Segmentor Fine-tuning COCO Panoptic (Kirillov et al. 2019) (118K) | 118K |
| Image Conversation | Stage 2: Alignment Pre-training LLaVA (Liu et al. 2023b) (558k) | 558k |
| Image Conversation Generic Segmentation VGD Segmentation Referring Segmentation GCG Segmentation Reasoning Segmentation | Stage 3: Mixed Fine-tuning LLaVA-1.5 (Liu et al. 2023b) (665K) COCO Panoptic (Kirillov et al. 2019) (118K) COCO-VGD (117K) RefCOCO(17K), RefCOCO+(17K), RefCOCOg(22k) Grand-f (Rasheed et al. 2024) (1K), RefCOCOg (19K), PSG (28K), Flickr (148K) LISA ReasonSeg (Lai et al. 2024) (0.2K) | 665K 118K 117K 301K 195K 0.2K |

Tab. 11: The Hyper-parameters in Multi-stage Training of X-SAM.

| Item | Stage 1 | Stage 2 | Stage 3 |
|---------------------|----------------------------------|------------------------|-------------------|
| | Segmentor Fine-tuning | Alignment Pre-training | Mixed Fine-tuning |
| batch size | 64 | 256 | 64 |
| training epochs | 36 | 1 | 1 |
| lr of dual encoders | 1e-5 | - | 4e-6 |
| lr of other modules | 1e-4 | 1e-3 | 4e-5 |
| optimizer | AdamW (Loshchilov et al. 2017) | | |
| optimizer momentum | $\beta_1 = 0.9, \beta_2 = 0.999$ | | |
| weight decay | 0.05 | 0 | 0.05 |
| warmup ratio | 0.03 | 0.03 | 0.03 |
| clip max norm | 0.01 | 1 | 1 |
| image augmentation | random scale[0.1, 2.0] | - | - |

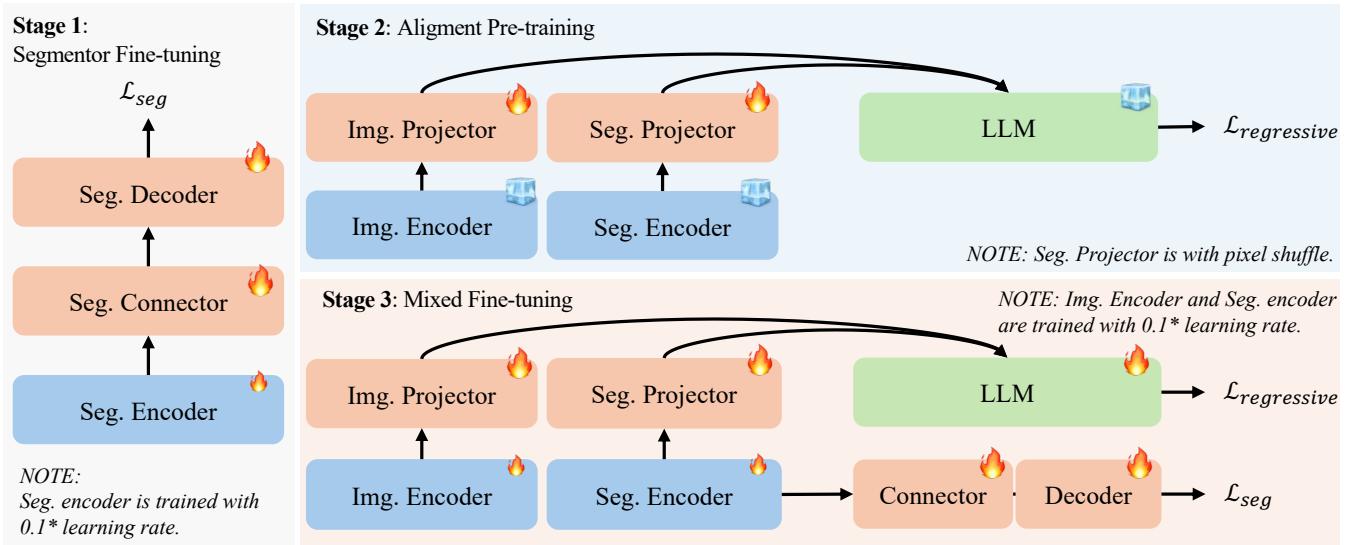


Fig. 4: The Multi-stage Training of X-SAM. X-SAM performs a multi-stage training process, including segmentor fine-tuning, alignment pre-training, and mixed fine-tuning. *Segmentor fine-tuning*: train the segmentor on the segmentation datasets to obtain a generalized segmentor. *Alignment pre-training*: train the dual projectors to align the vision features and the LLM features. *Mixed fine-tuning*: fine-tune the dual projectors, the segmentation decoder, and the LLM on the mixed datasets.

Tab. 12: Comparison of Generic Segmentation. We compare different methods on the generic segmentation benchmarks.

| Method | Encoder | PQ | AP | mIoU |
|---------------------------------|--------------|-------------|-------------|-------------|
| Mask2Former(Cheng et al. 2022a) | Swin-L | 57.8 | 48.6 | 67.4 |
| X-Decoder(Zou et al. 2023a) | DaViT-B | 56.2 | 45.8 | 66.0 |
| SEEM(Zou et al. 2023b) | DaViT-B | 56.1 | 46.4 | 66.3 |
| OMG-LLaVA(Zhang et al. 2024c) | ConvNeXt-XXL | 53.8 | - | - |
| PSALM(Zhang et al. 2024d) | Swin-B | 55.9 | <u>45.7</u> | 66.6 |
| X-SAM (Ours) | SAM-L | <u>54.7</u> | 47.0 | <u>66.5</u> |

Tab. 14: Comparison of Reasoning Segmentation. We compare X-SAM with other methods on the reasoning segmentation benchmark.

| Method | val | | test | | | | | |
|--------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | overall | | short query | | long query | | overall | |
| | gIoU | cloU | gIoU | cloU | gIoU | cloU | gIoU | cloU |
| OVSeg (Liang et al. 2023) | 28.5 | 18.6 | 18.0 | 15.5 | 28.7 | 22.5 | 26.1 | 20.8 |
| SEEM (Zou et al. 2023b) | 25.5 | 21.2 | 20.1 | 11.5 | 25.6 | 20.8 | 24.3 | 18.7 |
| LISA-7B (Lai et al. 2024) | 44.4 | 46.0 | 37.6 | 34.4 | 36.6 | 34.7 | 36.8 | 34.1 |
| LISA-7B (ft) (Lai et al. 2024) | <u>52.9</u> | 54.0 | <u>40.6</u> | <u>40.6</u> | <u>49.4</u> | 51.0 | <u>47.3</u> | 48.4 |
| X-SAM (Ours) | 56.6 | 32.9 | 47.7 | 48.1 | 56.0 | <u>40.8</u> | 57.8 | 41.0 |

Tab. 15: Comparison of Interactive Segmentation. We compare X-SAM with other methods on the interactive segmentation benchmark.

| Method | Point | | Scribble | | Box | | Mask | |
|------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | mIoU | cloU | mIoU | cloU | mIoU | cloU | mIoU | cloU |
| SAM-B(Kirillov et al. 2023a) | 48.7 | 33.6 | - | - | 73.7 | 68.7 | - | - |
| SAM-L(Kirillov et al. 2023a) | 51.8 | 37.7 | - | - | 76.6 | 71.6 | - | - |
| SEEM-B(Zou et al. 2023b) | 47.8 | 57.8 | 43.0 | 44.0 | 44.9 | 42.1 | 48.4 | 65.0 |
| OMG-Seg(Li et al. 2024e) | 59.3 | - | - | - | - | - | - | - |
| PSALM(Zhang et al. 2024d) | <u>64.3</u> | 74.0 | 66.9 | 80.0 | <u>67.3</u> | 80.9 | <u>67.6</u> | 82.4 |
| X-SAM (Ours) | 65.4 | <u>62.9</u> | 66.9 | <u>75.7</u> | 69.6 | <u>75.4</u> | 69.7 | <u>77.0</u> |

performs classification for each interactive visual prompt. As a result, the AP metric for VGD segmentation is poor because the instance-level predictions are of low quality. This may explain why PSALM lacks the capability for instance-level visual grounding segmentation.

X-SAM COCO-Interactive Evaluation. X-SAM is the first unified segmentation MLLM, capable of adapting to all image segmentation tasks, including interactive segmentation (Zhang et al. 2024d). To evaluate its performance on the COCO-Interactive dataset, we first filter the instance-level predictions using a threshold of 0.5. We then calculate the IoU score between each remaining prediction and the visual prompt mask. Finally, we select the instance prediction with the highest IoU score as the final interactive segmentation result.

A.5 More Experimental Results

Generic Segmentation. Tab. 12 presents the results of generic segmentation. Thanks to our segmentor design and

Tab. 13: Comparison of OV Segmentation. We compare different methods on the A150-OV segmentation benchmarks.

| Method | PQ | AP | mIoU |
|-----------------------------|-------------|-------------|-------------|
| MaskCLIP (Ding et al. 2022) | 15.1 | 6.0 | 23.7 |
| ODISE(Xu et al. 2023) | 22.6 | 14.4 | 29.9 |
| PSALM(Zhang et al. 2024d) | 13.7 | <u>9.0</u> | 18.2 |
| HyperSeg(Wei et al. 2024b) | <u>16.1</u> | - | <u>22.3</u> |
| X-SAM (Ours) | 20.9 | 16.2 | 28.8 |

fine-tuning, X-SAM can adapt to generic segmentation and achieves competitive performance on the COCO-Panoptic dataset.

Open-Vocabulary Segmentation. Tab. 13 shows the results of open-vocabulary segmentation. Benefiting from the robust mask generation of SAM and our mixed fine-tuning strategy, X-SAM achieves the best performance on open-vocabulary segmentation tasks.

Reasoning Segmentation. In Tab. 14, we present the results of reasoning segmentation on the reasoning segmentation benchmark. We report the performance of our method on both the validation set and the test set, following (Lai et al. 2024). X-SAM achieves the best gIoU metric on both the validation and test sets, even though it is not specifically designed for reasoning segmentation. While the cloU metric is not the best, it remains comparable to state-of-the-art methods. As the number of samples in the validation and test sets is limited, the results on this benchmark may not be stable.

Interactive Segmentation. Tab. 15 shows the results of

Tab. 16: Comparison of Closed-set Segmentation. We compare X-SAM with other methods on the COCO-Panoptic benchmark.

| Method | Encoder | # Params. | Epochs | PQ | PQ^{Th} | PQ^{St} | AP_{pan}^{Th} | mIoU pan |
|---------------------------------|---------|-----------|--------|------|-----------|-----------|-----------------|---------------|
| Max-DeepLab(Wang et al. 2021b) | Max-L | 451M | 216 | 51.1 | 57.0 | 42.2 | - | - |
| MaskFormer(Cheng et al. 2022c) | Swin-L | 212M | 300 | 52.7 | 58.5 | 44.0 | 40.1 | 64.8 |
| K-Net(Zhang et al. 2021) | Swin-L | - | 36 | 54.6 | 60.2 | 46.0 | - | - |
| Mask2Former(Cheng et al. 2022a) | Swin-L | 216M | 100 | 57.8 | 64.2 | 48.1 | 48.6 | 67.4 |
| X-SAM (Ours) | SAM-L | 364M | 36 | 54.2 | 60.0 | 45.4 | 44.2 | 63.8 |

Tab. 17: Comparison of Image-level Benchmarks. We compare X-SAM with other methods on the image-level benchmarks, including MME (Fu et al. 2024), MMBench (Liu et al. 2024c), SEED-Bench (Li et al. 2024a), POPE (Li et al. 2023b), and AI2D (Kembhavi et al. 2016).

| Method | MME | MMBench | SEED-Bench | POPE | AI2D |
|--------------------------------|-------------------|-------------|-------------|-------------|-------------|
| LLaVA-1.5 (Liu et al. 2024a) | 1510 / - | 64.3 | 58.6 | 87.3 | - |
| LLaVA-OV (Li et al. 2024b) | 1580 / 418 | 80.8 | 75.4 | - | 81.4 |
| LISA (Lai et al. 2024) | 1 / 1 | 0.4 | - | 0.0 | 0.0 |
| PixelLM (Ren et al. 2024) | 309 / 135 | 17.4 | - | 0.0 | 0.0 |
| LaSagnA (Wei et al. 2024a) | 0 / 0 | 0.0 | - | 0.0 | 0.0 |
| GLaMM (Rasheed et al. 2024) | 14 / 9 | 36.8 | - | 0.94 | 28.2 |
| OMG-LLaVA (Zhang et al. 2024c) | 1177 / 235 | 47.9 | 56.5 | 80.0 | 42.9 |
| X-SAM (Ours) | 1374 / 312 | 69.3 | 69.3 | 89.3 | 62.6 |

the interactive segmentation. Since interactive segmentation shares similar data with VGD segmentation, we exclude the training data for interactive segmentation and perform a process similar to that used for VGD segmentation to obtain the interactive segmentation results. X-SAM achieves the best or second-best performance on interactive segmentation, even without being trained on the specific data.

Closed-Set Segmentation. In Tab. 16, we present the results of segmentor fine-tuning on the closed-set COCO-Panoptic benchmark. To preserve the robust generalization ability of SAM in mask prediction, we fine-tune the SAM encoder for only 36 epochs. Compared with other methods, our approach achieves comparable performance to the SAM-L encoder with just 36 epochs of fine-tuning.

Image-level Benchmarks. In Tab. 17, we present the results of image-level benchmarks, including MME (Fu et al. 2024), MMBench (Liu et al. 2024c), SEED-Bench (Li et al. 2024a), POPE (Li et al. 2023b), and AI2D (Kembhavi et al. 2016). When jointly co-training with segmentation and conversation datasets, X-SAM achieves the best performance compared to other segmentation MLLMs on these benchmarks. Compared to LISA (Lai et al. 2024), PixelLM (Ren et al. 2024), and GLaMM (Rasheed et al. 2024), our method achieves significant improvements, demonstrating its effectiveness—even outperforming the previous best, OMG-LLaVA (Zhang et al. 2024c). On the POPE benchmark, X-SAM even surpasses LLaVA-V1.5 (Liu et al. 2023b), which is designed specifically for image-level conversation.

A.6 More Ablation Studies

Dataset Balance Resampling. In Tab. 18, we ablate the oversampling ratio t from 0 to 0.1 in dataset balance re-

sampling. When $t = 0$, no dataset is oversampled. Otherwise, datasets are oversampled with a repeat factor of $r_d = \max(1, \sqrt{t/f_d})$. We find that the performance on some small datasets is sensitive to the oversampling ratio t , especially for the reasoning segmentation dataset, where the gIoU score improves from 44.1% to 56.6% as t increases from 0 to 0.1. Meanwhile, larger datasets are not sensitive to the oversampling ratio t . As a result, the overall performance improvement reaches its highest when $t = 0.1$. Therefore, we set t to 0.1 in the final experiment.

Image Encoder. In Tab. 19, we ablate the image encoder of X-SAM by replacing it with CLIP (Radford et al. 2021), SigLIP-so400m (Zhai et al. 2023), and SigLIP2-so400m (Tschannen et al. 2025). It can be observed that using more powerful image encoders improves the image content understanding ability of X-SAM, especially on the image conversation and GCG segmentation benchmarks, and even enhances performance on the generic segmentation benchmark. Although SigLIP achieves the best performance on the reasoning segmentation benchmark, it does not have a performance advantage on the other benchmarks. Meanwhile, SigLIP2 demonstrates more robust and consistently better performance across all benchmarks. Therefore, we adopt SigLIP2-so400m as the image encoder for the final experiment.

A.7 More Visualization Results

Generic Segmentation. Fig. 5 shows the visualization results of X-SAM in generic segmentation, including semantic, instance and panoptic segmentation, which needs both the semantic and instance level understanding of the image. X-SAM can generate accurate and complete masks for the

Tab. 18: Ablation of Dataset Balance Resampling. $\sum \Delta$ represents the total performance improvement compared to the first row.

| t | COCO-Pan PQ | A150-OV AP | RefCOCO cIoU | Reason-Val gIoU | GCG-Val mIoU | COCO-VGD AP | Conv.-MMB Acc | $\sum \Delta$ |
|-------|----------------|---------------|-----------------|--------------------|-----------------|----------------|------------------|-----------------|
| 0 | 53.5 | 20.6 | 84.2 | 44.1 | 62.8 | 47.4 | 68.6 | 0.0 |
| 0.001 | 53.9 | <u>20.9</u> | 85.5 | 44.1 | 62.7 | 47.7 | 67.3 | $\uparrow 0.9$ |
| 0.01 | 54.3 | 21.0 | 84.6 | <u>46.5</u> | <u>63.1</u> | 48.1 | 67.4 | $\uparrow 3.8$ |
| 0.1 | 54.7 | <u>20.9</u> | <u>85.1</u> | 56.6 | 69.4 | 47.9 | 69.3 | $\uparrow 22.7$ |

Tab. 19: Ablation of Image Encoder. We ablate the impact of different image encoders, including CLIP (Radford et al. 2021), SigLIP (Zhai et al. 2023), and SigLIP2 (Tschanne et al. 2025), on the performance of X-SAM.

| Img. Enc. | COCO-Pan PQ | A150-OV AP | RefCOCO cIoU | Reason-Val gIoU | GCG-Val mIoU | COCO-VGD AP | Conv.-MMB Acc | Conv.-MMB Acc |
|--------------------------------|----------------|---------------|-----------------|--------------------|-----------------|----------------|-------------------|------------------|
| | | | | | | | | |
| CLIP (Radford et al. 2021) | <u>53.7</u> | <u>21.2</u> | <u>84.8</u> | 55.6 | <u>62.6</u> | 48.0 | 1327 / 277 | 68.5 |
| SigLIP (Zhai et al. 2023) | 53.2 | 21.7 | 83.3 | 61.6 | 61.8 | 48.0 | 1331 / 280 | 68.4 |
| SigLIP2 (Tschanne et al. 2025) | 54.7 | 20.9 | 85.1 | <u>56.6</u> | 69.4 | 47.9 | 1374 / 312 | 69.3 |

objects in the image.

Open-Vocabulary Segmentation. Fig. 6 shows the visualization results of X-SAM in open-vocabulary (OV) segmentation, including OV-semantic, OV-instance, and OV-panoptic segmentation, which require segmenting objects that may not exist in the training set. X-SAM can segment objects that are not in the training set, demonstrating the robust generalization ability of the proposed method.

GCG Segmentation. Fig. 7 shows the visualization results of X-SAM in GCG segmentation, which needs to describe the image and output the corresponding mask. X-SAM can not only effectively understand the image and generate the language description, but also generate the segmentation masks for the corresponding phrases.

Referring Segmentation. Fig. 8 shows the visualization results of X-SAM in referring segmentation, which needs to segment objects that are referred to by natural language. X-SAM can effectively understand the referring expression and segment the objects that are referred to by natural language.

Reasoning Segmentation. Fig. 9 shows the visualization results of X-SAM in reasoning segmentation, which needs to segment the object that is related to the question. X-SAM can effectively understand the complex question and then generate the corresponding mask for the question.

Interactive Segmentation. Fig. 10 shows the visualization results of X-SAM in interactive segmentation, which needs to segment the individual objects with which the user interacts. X-SAM can generate the corresponding mask for the user’s interactive visual prompt.

VGD Segmentation. Fig. 11 shows the visualization results of X-SAM in VGD segmentation for a single image, which needs to segment all the objects in the *single image* that is grounded with the user’s visual prompt. X-SAM can effectively segment all the objects in the image given by the user’s visual grounded prompt. In addition, X-SAM can perform VGD segmentation on the *cross image*, which needs to segment the objects grounded in another image. Fig. 12 shows

the visualization results of X-SAM in VGD segmentation for cross image, which demonstrates the effectiveness of X-SAM for VGD segmentation both in single image and cross image.

A.8 Further Discussions

Limitation. Although X-SAM achieves unified segmentation by extending *segment anything* to *any segmentation*, there is still considerable room for improvement. First, joint co-training with segmentation datasets and conversation datasets negatively impacts performance on some segmentation datasets, a phenomenon also observed in (Zhang et al. 2024c) and (Rasheed et al. 2024). This challenge may be addressed by designing a more balanced dataset mixture. Second, the performance of X-SAM is not optimal across all tasks, which has also been observed in other unified segmentation methods (Zhang et al. 2024c; Rasheed et al. 2024; Lai et al. 2024). This challenge remains a major obstacle for unified models and may be addressed by scaling up the model size and the amount of training data.

Future Work. Several avenues for future work can be explored with our novel unified framework. We highlight two potential directions. The first is to integrate X-SAM with SAM2 (Ravi et al. 2024), a unified model for segmentation in images and videos. This integration would further extend the application of X-SAM to video segmentation. The second direction is to extend VGD segmentation to the video domain, which would constitute an interesting video segmentation task and introduce visually grounded temporal information to segmentation. We plan to explore these directions in the future, provided that more computational resources become available.

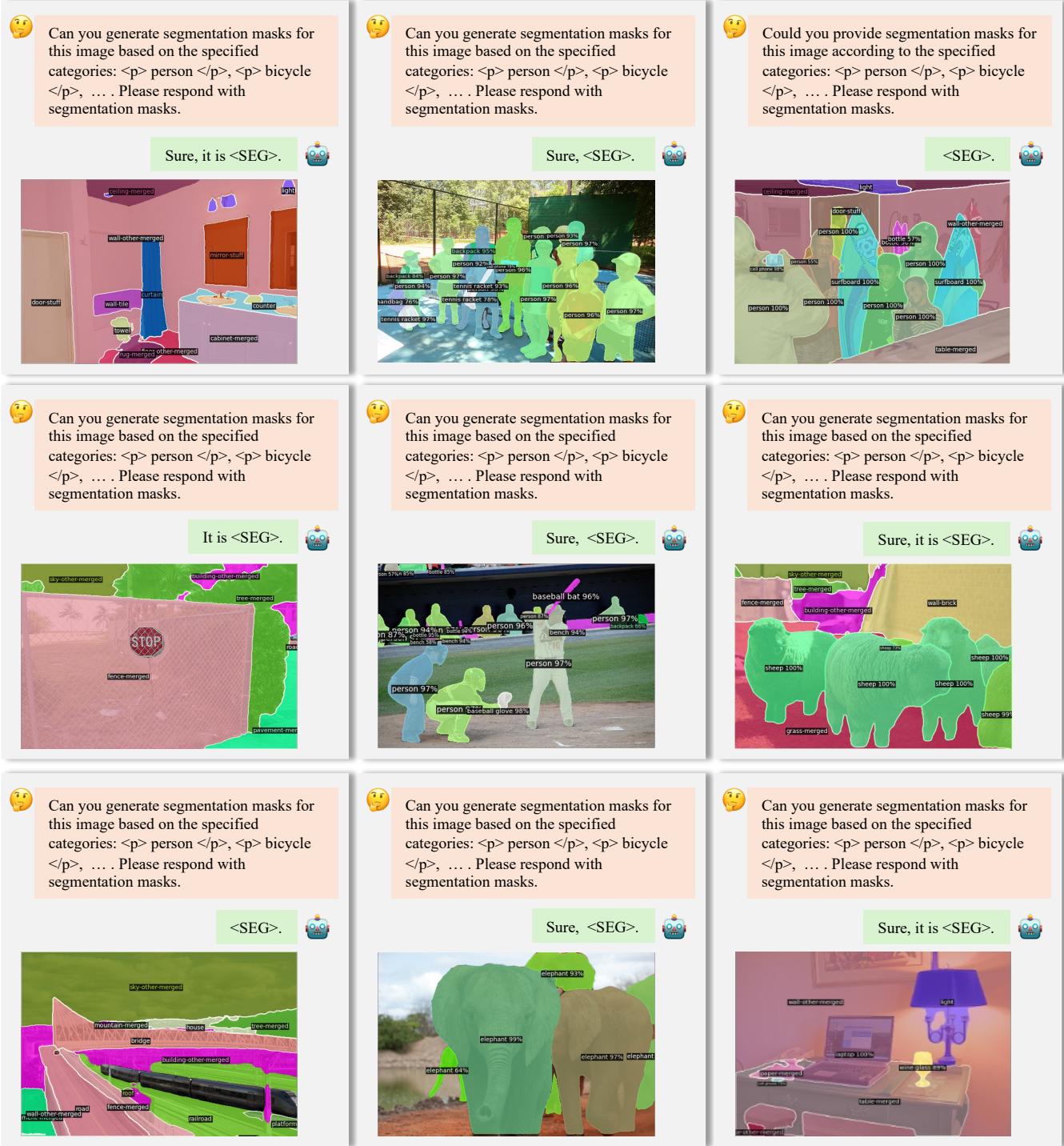


Fig. 5: Visualization Results of Generic Segmentation. Visualized images are sampled from the COCO2017 Val set. More category names are omitted for better visualization.



Fig. 6: Visualization Results of Open-Vocabulary (OV) Segmentation. Visualized images are sampled from the ADE20K Val set. More category names are omitted for better visualization.

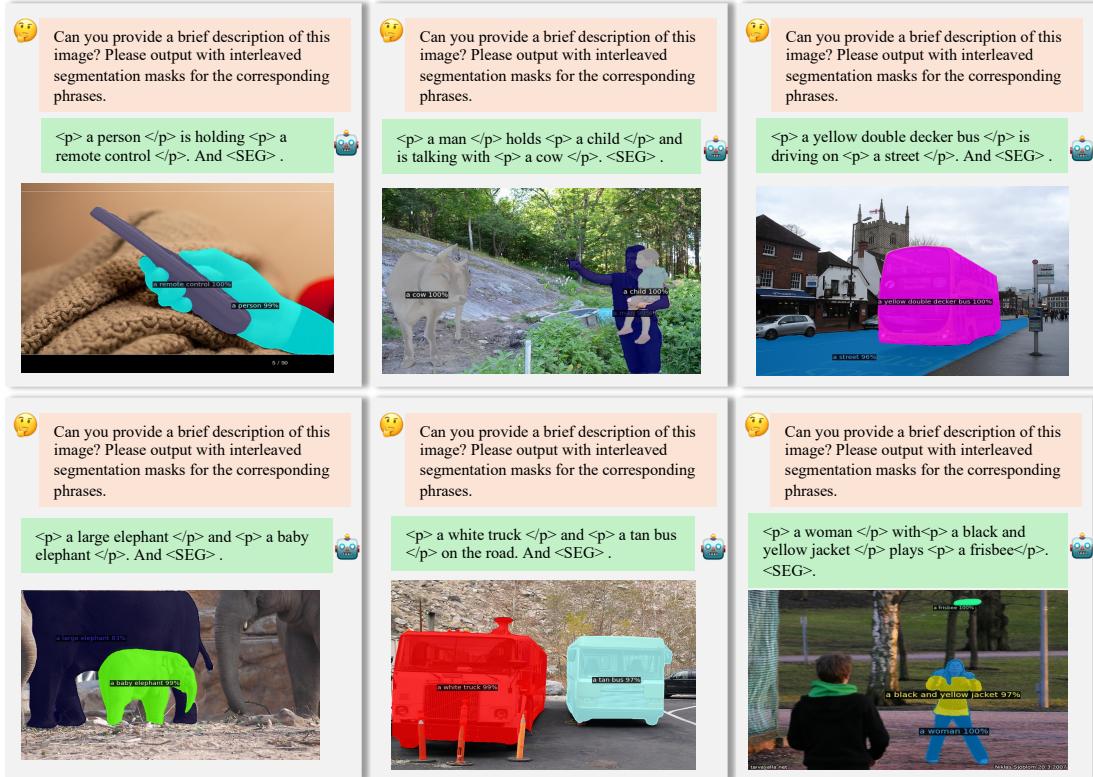


Fig. 7: Visualization Results of GCG Segmentation. Visualized images are sampled from the Open-PSG Val set.

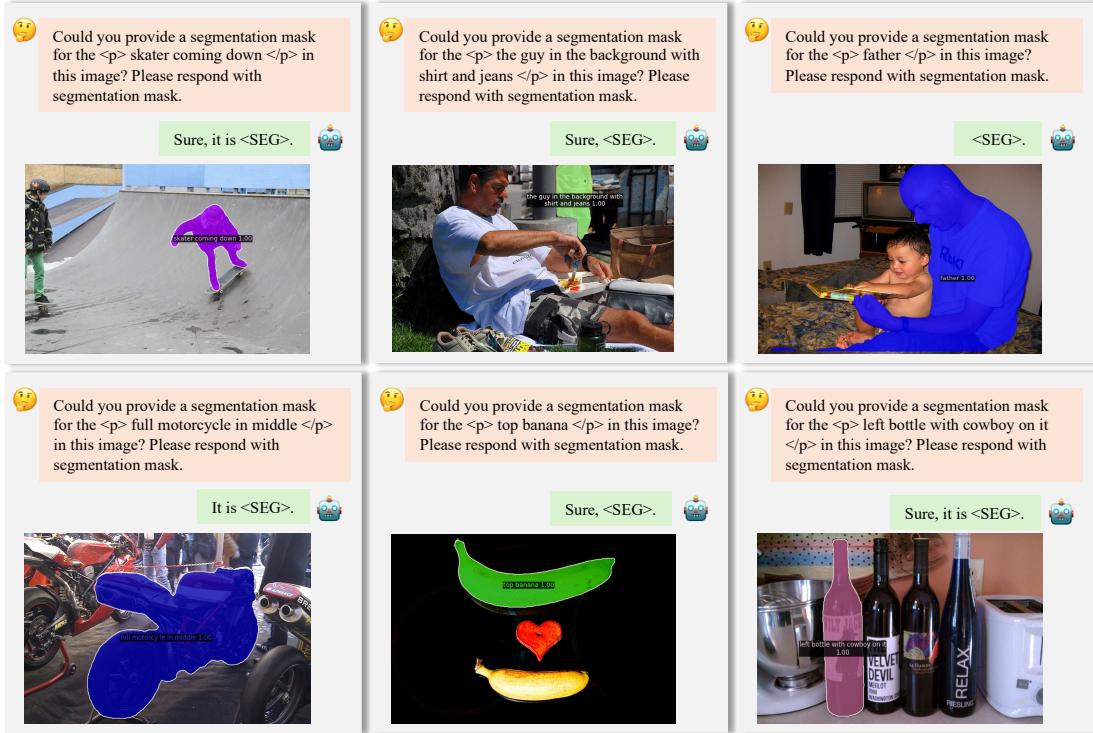


Fig. 8: Visualization Results of Referring Segmentation. Visualized images are sampled from the RefCOCO Val set.

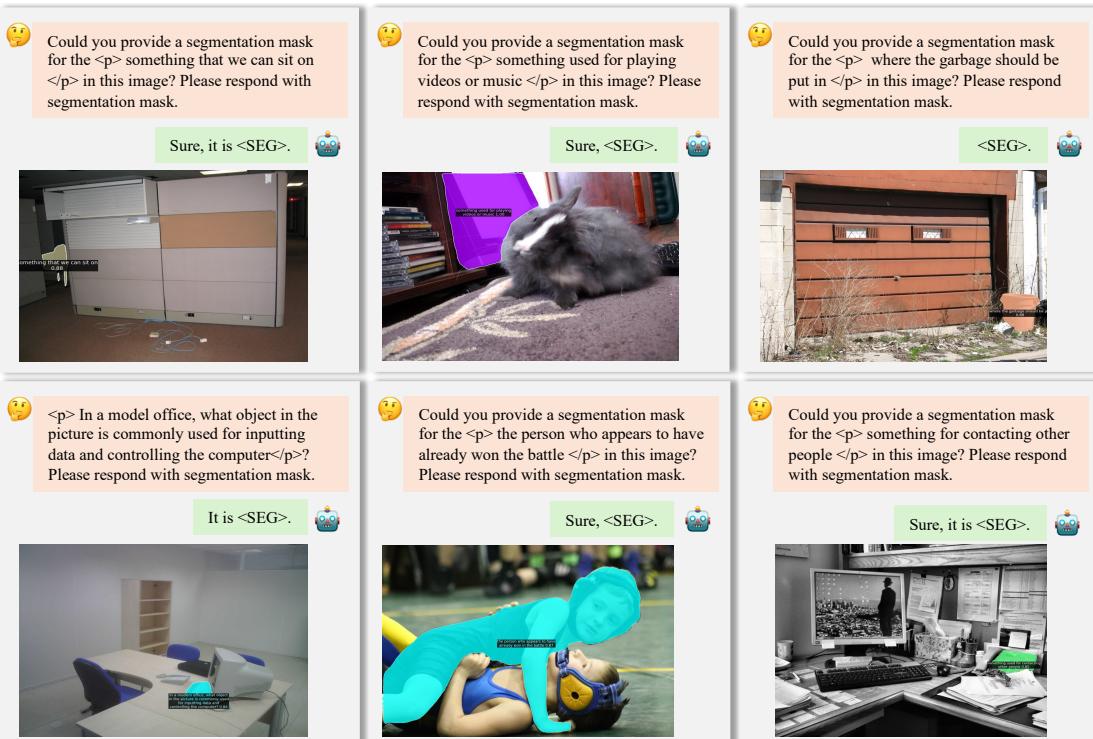


Fig. 9: Visualization Results of Reasoning Segmentation. Visualized images are sampled from the reasoning segmentation Val set.



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

Sure, it is <SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

<SEG>. 🤖

Point Vision Query



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

Sure, it is <SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

Sure, <SEG>. 🤖

Scribble Vision Query



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

Sure, <SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

Sure, <SEG>. 🤖

Box Vision Query



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

<SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>?
Please respond with segmentation masks.

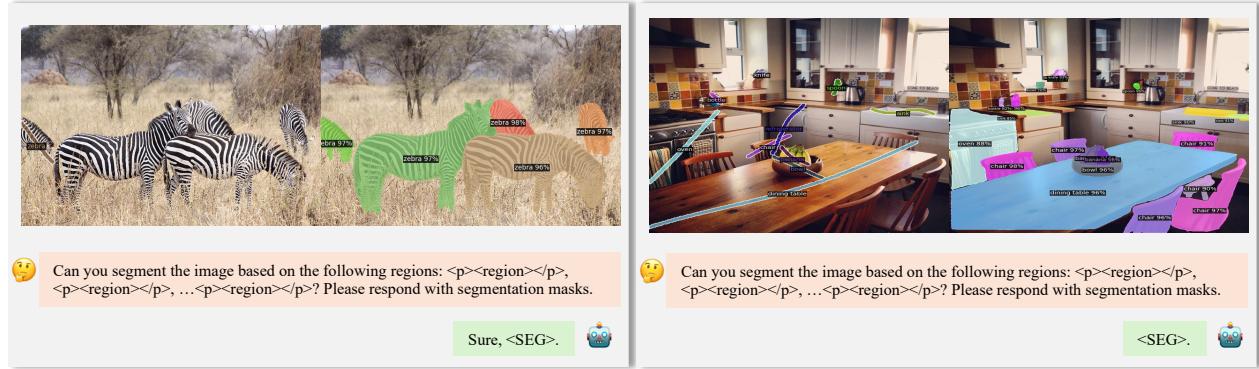
Sure, <SEG>. 🤖

Mask Vision Query

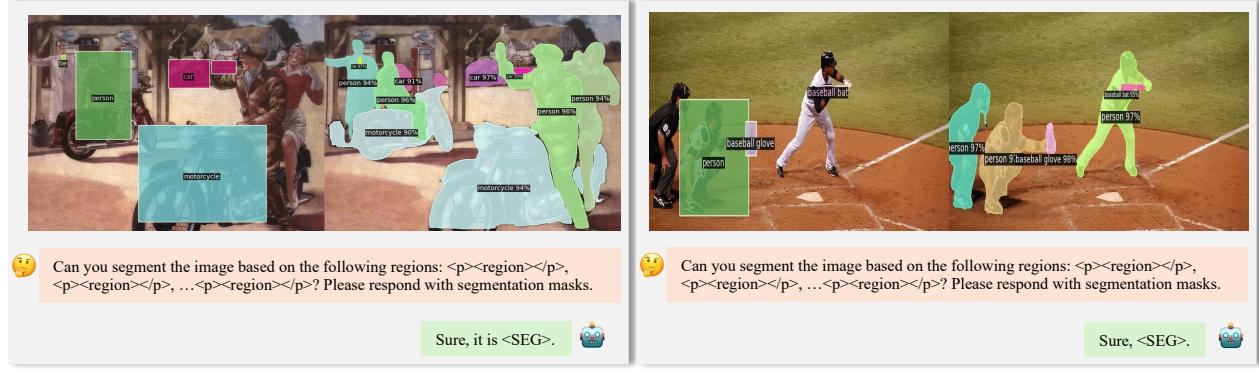
Fig. 10: Visualization Results of Interactive Segmentation. Visualized images are sampled from the COCO-Interactive Val set. In each sample, the left image is the original image with the interactive prompts, and the right image is the visualized result of our method. From top to bottom, there are four types of visual interactive prompts: point, scribe, box, and mask. We visualize the category name of the interactive prompt for better visualization.



Point Vision Query



Scribble Vision Query



Box Vision Query



Mask Vision Query

Fig. 11: Visualization Results of VGD Segmentation (Single Image). Visualized images are sampled from COCO-VGD Val set. In each sample, the left image is the original image with the visual grounded prompts, and the right image is the visualized result of our method. From top to bottom, there are four types of visual grounded prompts: point, scribble, box, and mask. We visualize the category name of the grounded prompt for better visualization.



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

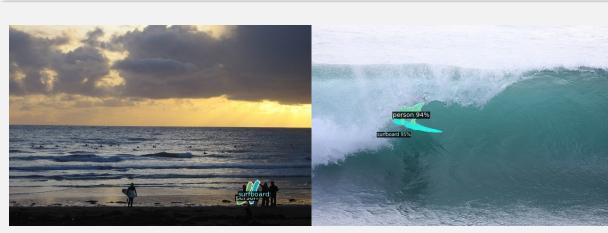
<SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

Sure, it is <SEG>. 🤖

Point Vision Query



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

Sure, <SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

<SEG>. 🤖

Scribble Vision Query



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

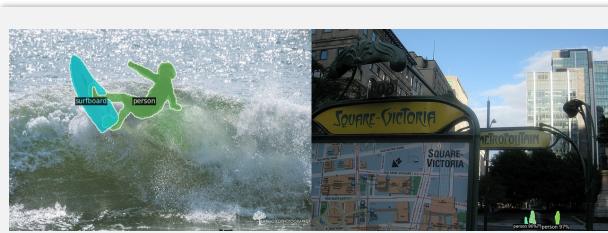
Sure, it is <SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

Sure, <SEG>. 🤖

Box Vision Query



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

Sure, <SEG>. 🤖



🤔 Can you segment the image based on the following regions: <p><region></p>, <p><region></p>, ...<p><region></p>? Please respond with segmentation masks.

<SEG>. 🤖

Mask Vision Query

Fig. 12: Visualization Results of VGD Segmentation (Cross Image). Visualized images are sampled from the COCO-VGD Val set. In each sample, the left image is a different image with the visual grounded prompts, and the right image is the visualized result of our method. From top to bottom, there are four types of visual grounded prompts: point, scribe, box, and mask. We visualize the category name of the grounded prompt for better visualization.

References

- Abdin, M.; Aneja, J.; Awadalla, H.; Awadallah, A.; Awan, A. A.; Bach, N.; Bahree, A.; Bakhtiari, A.; Bao, J.; Behl, H.; et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Antol, S.; Agrawal, A.; Lu, J.; Mitchell, M.; Batra, D.; Zitnick, C. L.; and Parikh, D. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, 2425–2433.
- Athar, A.; Hermans, A.; Luiten, J.; Ramanan, D.; and Leibe, B. 2023. Tarvis: A unified approach for target-based video segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18738–18748.
- Bai, J.; Bai, S.; Chu, Y.; Cui, Z.; Dang, K.; Deng, X.; Fan, Y.; Ge, W.; Han, Y.; Huang, F.; et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Bar, A.; Gandelsman, Y.; Darrell, T.; Globerson, A.; and Efros, A. 2022. Visual prompting via image inpainting. *Advances in Neural Information Processing Systems*, 35: 25005–25017.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Carion, N.; Massa, F.; Synnaeve, G.; Usunier, N.; Kirillov, A.; and Zagoruyko, S. 2020. End-to-end object detection with transformers. In *European conference on computer vision*, 213–229. Springer.
- Chen, J.; Shen, Y.; Gao, J.; Liu, J.; and Liu, X. 2018. Language-based image editing with recurrent attentive models. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 8721–8729.
- Chen, K.; Pang, J.; Wang, J.; Xiong, Y.; Li, X.; Sun, S.; Feng, W.; Liu, Z.; Shi, J.; Ouyang, W.; et al. 2019. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 4974–4983.
- Chen, Z.; Wang, W.; Tian, H.; Ye, S.; Gao, Z.; Cui, E.; Tong, W.; Hu, K.; Luo, J.; Ma, Z.; et al. 2024. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12): 220101.
- Cheng, B.; Misra, I.; Schwing, A. G.; Kirillov, A.; and Girdhar, R. 2022a. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 1290–1299.
- Cheng, B.; Misra, I.; Schwing, A. G.; Kirillov, A.; and Girdhar, R. 2022b. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 1290–1299.
- Cheng, B.; Misra, I.; Schwing, A. G.; Kirillov, A.; and Girdhar, R. 2022c. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 1290–1299.
- conference on computer vision and pattern recognition, 1290–1299.
- Contributors, X. 2023. XTuner: A Toolkit for Efficiently Fine-tuning LLM. <https://github.com/InternLM/xtuner>.
- Ding, Z.; Wang, J.; and Tu, Z. 2022. Open-vocabulary universal image segmentation with maskclip. *arXiv preprint arXiv:2208.08984*.
- Dong, X.; Zhang, P.; Zang, Y.; Cao, Y.; Wang, B.; Ouyang, L.; Wei, X.; Zhang, S.; Duan, H.; Cao, M.; et al. 2024a. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*.
- Dong, X.; Zhang, P.; Zang, Y.; Cao, Y.; Wang, B.; Ouyang, L.; Zhang, S.; Duan, H.; Zhang, W.; Li, Y.; et al. 2024b. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *Advances in Neural Information Processing Systems*, 37: 42566–42592.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Duan, H.; Yang, J.; Qiao, Y.; Fang, X.; Chen, L.; Liu, Y.; Dong, X.; Zang, Y.; Zhang, P.; Wang, J.; et al. 2024. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, 11198–11201.
- Fang, Z.; Li, X.; Li, X.; Buhmann, J. M.; Loy, C. C.; and Liu, M. 2023. Explore in-context learning for 3d point cloud understanding. *Advances in Neural Information Processing Systems*, 36: 42382–42395.
- Fu, C.; Chen, P.; Shen, Y.; Qin, Y.; Zhang, M.; Lin, X.; Yang, J.; Zheng, X.; Li, K.; Sun, X.; Wu, Y.; and Ji, R. 2024. MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models. *arXiv preprint arXiv:2404.08506*.
- Gu, X.; Cui, Y.; Huang, J.; Rashwan, A.; Yang, X.; Zhou, X.; Ghiasi, G.; Kuo, W.; Chen, H.; Chen, L.-C.; et al. 2023. Dataseg: Taming a universal multi-dataset multi-task segmentation model. *Advances in Neural Information Processing Systems*, 36: 67329–67354.
- Gupta, A.; Dollar, P.; and Girshick, R. 2019. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5356–5364.
- Hudson, D. A.; and Manning, C. D. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 6700–6709.
- Jain, J.; Li, J.; Chiu, M. T.; Hassani, A.; Orlov, N.; and Shi, H. 2023. Oneformer: One transformer to rule universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2989–2998.

- Jia, C.; Yang, Y.; Xia, Y.; Chen, Y.-T.; Parekh, Z.; Pham, H.; Le, Q.; Sung, Y.-H.; Li, Z.; and Duerig, T. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, 4904–4916. PMLR.
- Ke, L.; Ye, M.; Danelljan, M.; Tai, Y.-W.; Tang, C.-K.; Yu, F.; et al. 2023. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36: 29914–29934.
- Kembhavi, A.; Salvato, M.; Kolve, E.; Seo, M.; Hajishirzi, H.; and Farhadi, A. 2016. A diagram is worth a dozen images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV* 14, 235–251. Springer.
- Kirillov, A.; He, K.; Girshick, R.; Rother, C.; and Dollár, P. 2019. Panoptic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9404–9413.
- Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; et al. 2023a. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, 4015–4026.
- Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; et al. 2023b. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, 4015–4026.
- Lai, X.; Tian, Z.; Chen, Y.; Li, Y.; Yuan, Y.; Liu, S.; and Jia, J. 2024. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9579–9589.
- Li, B.; Ge, Y.; Ge, Y.; Wang, G.; Wang, R.; Zhang, R.; and Shan, Y. 2024a. Seed-bench: Benchmarking multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 13299–13308.
- Li, B.; Weinberger, K. Q.; Belongie, S.; Koltun, V.; and Ranftl, R. 2022a. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*.
- Li, B.; Zhang, Y.; Guo, D.; Zhang, R.; Li, F.; Zhang, H.; Zhang, K.; Zhang, P.; Li, Y.; Liu, Z.; et al. 2024b. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022b. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, 12888–12900. PMLR.
- Li, S.; Ke, L.; Danelljan, M.; Piccinelli, L.; Segu, M.; Van Gool, L.; and Yu, F. 2024c. Matching anything by segmenting anything. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18963–18973.
- Li, X.; Ding, H.; Yuan, H.; Zhang, W.; Pang, J.; Cheng, G.; Chen, K.; Liu, Z.; and Loy, C. C. 2024d. Transformer-based visual segmentation: A survey. *IEEE transactions on pattern analysis and machine intelligence*.
- Li, X.; You, A.; Zhu, Z.; Zhao, H.; Yang, M.; Yang, K.; Tan, S.; and Tong, Y. 2020. Semantic flow for fast and accurate scene parsing. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I* 16, 775–793. Springer.
- Li, X.; Yuan, H.; Li, W.; Ding, H.; Wu, S.; Zhang, W.; Li, Y.; Chen, K.; and Loy, C. C. 2024e. Omg-seg: Is one model good enough for all segmentation? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 27948–27959.
- Li, X.; Zhang, L.; Cheng, G.; Yang, K.; Tong, Y.; Zhu, X.; and Xiang, T. 2021. Global aggregation then local distribution for scene parsing. *IEEE Transactions on Image Processing*, 30: 6829–6842.
- Li, Y.; Du, Y.; Zhou, K.; Wang, J.; Zhao, W. X.; and Wen, J.-R. 2023a. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Li, Y.; Du, Y.; Zhou, K.; Wang, J.; Zhao, W. X.; and Wen, J.-R. 2023b. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Liang, F.; Wu, B.; Dai, X.; Li, K.; Zhao, Y.; Zhang, H.; Zhang, P.; Vajda, P.; and Marculescu, D. 2023. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 7061–7070.
- Lin, B.; Ye, Y.; Zhu, B.; Cui, J.; Ning, M.; Jin, P.; and Yuan, L. 2023. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*.
- Liu, H.; Li, C.; Li, Y.; and Lee, Y. J. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26296–26306.
- Liu, H.; Li, C.; Li, Y.; Li, B.; Zhang, Y.; Shen, S.; and Lee, Y. J. 2024b. Llavanext: Improved reasoning, ocr, and world knowledge.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023a. Visual instruction tuning. *Advances in neural information processing systems*, 36: 34892–34916.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023b. Visual instruction tuning. *Advances in neural information processing systems*, 36: 34892–34916.
- Liu, Y.; Duan, H.; Zhang, Y.; Li, B.; Zhang, S.; Zhao, W.; Yuan, Y.; Wang, J.; He, C.; Liu, Z.; et al. 2024c. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, 216–233. Springer.
- Loshchilov, I.; and Hutter, F. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Mehta, S.; and Rastegari, M. 2021. Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer. *arXiv preprint arXiv:2110.02178*.
- Pan, T.; Tang, L.; Wang, X.; and Shan, S. 2024. Tokenize anything via prompting. In *European Conference on Computer Vision*, 330–348. Springer.
- Peng, Z.; Wang, W.; Dong, L.; Hao, Y.; Huang, S.; Ma, S.; and Wei, F. 2023. Kosmos-2: Grounding multimodal

- large language models to the world. *arXiv preprint arXiv:2306.14824*.
- Qi, L.; Kuen, J.; Guo, W.; Shen, T.; Gu, J.; Jia, J.; Lin, Z.; and Yang, M.-H. 2022a. High-quality entity segmentation. *arXiv preprint arXiv:2211.05776*.
- Qi, L.; Kuen, J.; Wang, Y.; Gu, J.; Zhao, H.; Torr, P.; Lin, Z.; and Jia, J. 2022b. Open world entity segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(7): 8743–8756.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.
- Rajić, F.; Ke, L.; Tai, Y.-W.; Tang, C.-K.; Danelljan, M.; and Yu, F. 2025. Segment anything meets point tracking. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 9302–9311. IEEE.
- Rasheed, H.; Maaz, M.; Shaji, S.; Shaker, A.; Khan, S.; Cholakkal, H.; Anwer, R. M.; Xing, E.; Yang, M.-H.; and Khan, F. S. 2024. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 13009–13018.
- Ravi, N.; Gabeur, V.; Hu, Y.-T.; Hu, R.; Ryali, C.; Ma, T.; Khedr, H.; Rädle, R.; Rolland, C.; Gustafson, L.; et al. 2024. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*.
- Ren, Z.; Huang, Z.; Wei, Y.; Zhao, Y.; Fu, D.; Feng, J.; and Jin, X. 2024. Pixellm: Pixel reasoning with large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26374–26383.
- Sun, Q.; Cui, Y.; Zhang, X.; Zhang, F.; Yu, Q.; Wang, Y.; Rao, Y.; Liu, J.; Huang, T.; and Wang, X. 2024. Generative multimodal models are in-context learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14398–14409.
- Sun, S.; Wang, W.; Howard, A.; Yu, Q.; Torr, P.; and Chen, L.-C. 2023. Remax: Relaxing for better training on efficient panoptic segmentation. *Advances in Neural Information Processing Systems*, 36: 73480–73496.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Tschannen, M.; Gritsenko, A.; Wang, X.; Naeem, M. F.; Alabdulmohsin, I.; Parthasarathy, N.; Evans, T.; Beyer, L.; Xia, Y.; Mustafa, B.; et al. 2025. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*.
- VS, V.; Borse, S.; Park, H.; Das, D.; Patel, V.; Hayat, M.; and Porikli, F. 2024. PosSAM: Panoptic Open-vocabulary Segment Anything. *arXiv:2403.09620*.
- Wang, H.; Zhu, Y.; Adam, H.; Yuille, A.; and Chen, L.-C. 2021a. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5463–5474.
- Wang, H.; Zhu, Y.; Adam, H.; Yuille, A.; and Chen, L.-C. 2021b. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5463–5474.
- Wang, X.; Fang, Z.; Li, X.; Li, X.; Chen, C.; and Liu, M. 2024. Skeleton-in-context: Unified skeleton sequence modeling with in-context learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2436–2446.
- Wang, X.; Wang, W.; Cao, Y.; Shen, C.; and Huang, T. 2023a. Images speak in images: A generalist painter for in-context visual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6830–6839.
- Wang, X.; Zhang, X.; Cao, Y.; Wang, W.; Shen, C.; and Huang, T. 2023b. Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*.
- Wang, Z.; Jiang, Y.; Lu, Y.; He, P.; Chen, W.; Wang, Z.; Zhou, M.; et al. 2023c. In-context learning unlocked for diffusion models. *Advances in Neural Information Processing Systems*, 36: 8542–8562.
- Wei, C.; Tan, H.; Zhong, Y.; Yang, Y.; and Ma, L. 2024a. Lasagna: Language-based segmentation assistant for complex queries. *arXiv preprint arXiv:2404.08506*.
- Wei, C.; Zhong, Y.; Tan, H.; Liu, Y.; Zhao, Z.; Hu, J.; and Yang, Y. 2024b. HyperSeg: Towards Universal Visual Segmentation with Large Language Model. *arXiv preprint arXiv:2411.17606*.
- Wu, J.; Li, X.; Xu, S.; Yuan, H.; Ding, H.; Yang, Y.; Li, X.; Zhang, J.; Tong, Y.; Jiang, X.; et al. 2024. Towards open vocabulary learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(7): 5092–5113.
- Xu, J.; Liu, S.; Vahdat, A.; Byeon, W.; Wang, X.; and De Mello, S. 2023. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2955–2966.
- Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A.; Salakhudinov, R.; Zemel, R.; and Bengio, Y. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, 2048–2057. PMLR.
- Xu, S.; Yuan, H.; Shi, Q.; Qi, L.; Wang, J.; Yang, Y.; Li, Y.; Chen, K.; Tong, Y.; Ghanem, B.; et al. 2024. Rapsam: Towards real-time all-purpose segment anything. *arXiv preprint arXiv:2401.10228*.
- Yan, B.; Jiang, Y.; Wu, J.; Wang, D.; Luo, P.; Yuan, Z.; and Lu, H. 2023a. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15325–15336.

- Yan, B.; Jiang, Y.; Wu, J.; Wang, D.; Luo, P.; Yuan, Z.; and Lu, H. 2023b. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15325–15336.
- Yang, Z.; Wei, Y.; and Yang, Y. 2021. Collaborative video object segmentation by multi-scale foreground-background integration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9): 4701–4712.
- You, H.; Zhang, H.; Gan, Z.; Du, X.; Zhang, B.; Wang, Z.; Cao, L.; Chang, S.-F.; and Yang, Y. 2023. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*.
- Yuan, H.; Li, X.; Zhang, T.; Huang, Z.; Xu, S.; Ji, S.; Tong, Y.; Qi, L.; Feng, J.; and Yang, M.-H. 2025. Sa2VA: Marrying SAM2 with LLaVA for Dense Grounded Understanding of Images and Videos. *arXiv preprint arXiv:2501.04001*.
- Yuan, H.; Li, X.; Zhou, C.; Li, Y.; Chen, K.; and Loy, C. C. 2024a. Open-vocabulary SAM: Segment and recognize twenty-thousand classes interactively. In *European Conference on Computer Vision*, 419–437. Springer.
- Yuan, Y.; Li, W.; Liu, J.; Tang, D.; Luo, X.; Qin, C.; Zhang, L.; and Zhu, J. 2024b. Osprey: Pixel understanding with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 28202–28211.
- Zang, Y.; Li, W.; Han, J.; Zhou, K.; and Loy, C. C. 2025. Contextual object detection with multimodal large language models. *International Journal of Computer Vision*, 133(2): 825–843.
- Zhai, X.; Mustafa, B.; Kolesnikov, A.; and Beyer, L. 2023. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, 11975–11986.
- Zhang, H.; Li, H.; Li, F.; Ren, T.; Zou, X.; Liu, S.; Huang, S.; Gao, J.; Leizhang; Li, C.; et al. 2024a. Llava-grounding: Grounded visual chat with large multimodal models. In *European Conference on Computer Vision*, 19–35. Springer.
- Zhang, P.; Dong, X.; Wang, B.; Cao, Y.; Xu, C.; Ouyang, L.; Zhao, Z.; Duan, H.; Zhang, S.; Ding, S.; et al. 2023. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*.
- Zhang, T.; Li, X.; Fei, H.; Yuan, H.; Wu, S.; Ji, S.; Loy, C. C.; and Yan, S. 2024b. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. *Advances in Neural Information Processing Systems*, 37: 71737–71767.
- Zhang, T.; Li, X.; Fei, H.; Yuan, H.; Wu, S.; Ji, S.; Loy, C. C.; and Yan, S. 2024c. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. *Advances in Neural Information Processing Systems*, 37: 71737–71767.
- Zhang, W.; Pang, J.; Chen, K.; and Loy, C. C. 2021. K-net: Towards unified image segmentation. *Advances in Neural Information Processing Systems*, 34: 10326–10338.
- Zhang, Z.; Ma, Y.; Zhang, E.; and Bai, X. 2024d. Psalm: Pixelwise segmentation with large multi-modal model. In *European Conference on Computer Vision*, 74–91. Springer.
- Zhou, K.; Yang, J.; Loy, C. C.; and Liu, Z. 2022a. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9): 2337–2348.
- Zhou, Q.; Feng, Z.; Gu, Q.; Pang, J.; Cheng, G.; Lu, X.; Shi, J.; and Ma, L. 2022b. Context-aware mixup for domain adaptive semantic segmentation. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(2): 804–817.
- Zhou, Q.; Li, X.; He, L.; Yang, Y.; Cheng, G.; Tong, Y.; Ma, L.; and Tao, D. 2022c. TransVOD: End-to-end video object detection with spatial-temporal transformers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6): 7853–7869.
- Zhou, Y.; Zhang, T.; Ji, S.; Yan, S.; and Li, X. 2024. Improving video segmentation via dynamic anchor queries. In *European Conference on Computer Vision*, 446–463. Springer.
- Zou, X.; Dou, Z.-Y.; Yang, J.; Gan, Z.; Li, L.; Li, C.; Dai, X.; Behl, H.; Wang, J.; Yuan, L.; et al. 2023a. Generalized decoding for pixel, image, and language. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 15116–15127.
- Zou, X.; Yang, J.; Zhang, H.; Li, F.; Li, L.; Wang, J.; Wang, L.; Gao, J.; and Lee, Y. J. 2023b. Segment everything everywhere all at once. *Advances in neural information processing systems*, 36: 19769–19782.