

Open-Vocabulary Camouflaged Object Segmentation with Cascaded Vision Language Models

Kai Zhao^{1,2}, Wubang Yuan¹, Zheng Wang¹, Guanyi Li¹,
Xiaoqiang Zhu^{1✉}, Deng-Ping Fan³, Dan Zeng¹

¹ Shanghai University ² UCLA ³ Nankai University

kz@kaizhao.net, fdp@nankai.edu.cn
{yuanwubang, zhengwang, kwunyatlee, xqzhu, dzeng}@shu.edu.cn

Abstract. Open-Vocabulary Camouflaged Object Segmentation (OVCOS) seeks to segment and classify camouflaged objects from arbitrary categories, presenting unique challenges due to visual ambiguity and unseen categories. Recent approaches typically adopt a two-stage paradigm: first segmenting objects, then classifying the segmented regions using Vision Language Models (VLMs). However, these methods (1) suffer from a domain gap caused by the mismatch between VLMs’ full-image training and cropped-region inference, and (2) depend on generic segmentation models optimized for well-delineated objects, making them less effective for camouflaged objects. Without explicit guidance, generic segmentation models often overlook subtle boundaries, leading to imprecise segmentation. In this paper, we introduce a novel VLM-guided cascaded framework to address these issues in OVCOS. For segmentation, we leverage the Segment Anything Model (SAM), guided by the VLM. Our framework uses VLM-derived features as explicit prompts to SAM, effectively directing attention to camouflaged regions and significantly improving localization accuracy. For classification, we avoid the domain gap introduced by *hard* cropping. Instead, we treat the segmentation output as a *soft* spatial prior via the alpha channel, which retains the full image context while providing precise spatial guidance, leading to more accurate and context-aware classification of camouflaged objects. The same VLM is shared across both segmentation and classification to ensure efficiency and semantic consistency. Extensive experiments on both OVCOS and conventional camouflaged object segmentation benchmarks demonstrate the clear superiority of our method, highlighting the effectiveness of leveraging rich VLM semantics for both segmentation and classification of camouflaged objects. The code and models are open-sourced at <https://github.com/intcomp/camouflaged-vlm>.

1 Introduction

Open-Vocabulary Camouflaged Object Segmentation (OVCOS) is a challenging task that requires segmenting and classifying camouflaged objects of novel categories not seen during training [36]. Compared to traditional semantic segmentation [7, 40, 5], OVCOS faces additional challenges because it requires recognizing novel categories in visually ambiguous scenes, where camouflage leads to low contrast, indistinct boundaries, and high similarity between objects and their backgrounds. These challenges are particularly relevant in real-world applications such as medical image analysis [13] and agricultural monitoring [29], where annotations are scarce and target categories are often seen.

Several existing open-vocabulary segmentation approaches [10, 24, 3, 39] utilize vision-language models (VLMs), e.g. CLIP [37], to directly classify each pixel across the entire input image, thereby improving semantic generalization. These approaches operate under a one-stage framework. However, VLMs are pre-trained for image-level understanding, creating a granularity mismatch that hinders effective visual-semantic alignment and limits semantic transfer, often leading to suboptimal performance [44].

To mitigate this gap, recent works [36, 44, 11, 12, 27] first perform class-agnostic segmentation and then classify the segmented regions using VLM. This pipeline forms a two-stage framework. This decoupling of segmentation and classification partially alleviates the granularity mismatch [44]. However, in the segmentation stage, as shown in Figure 1(a), many existing approaches typically rely on the generic segmentation architectures [44, 11, 12, 27, 7] to spot the target region. These generic segmentation models are primarily

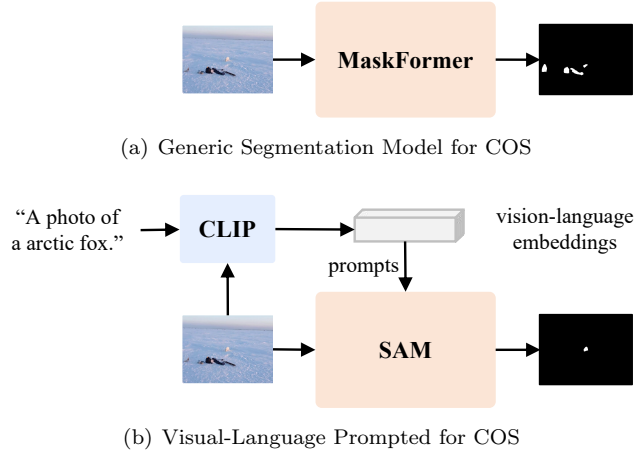


Figure 1: Different Camouflaged Object Segmentation (COS) paradigms in two-stage OVCOS. (a) Generic segmentation models, such as MaskFormer [7], typically operate directly on the input image without target-specific guidance, and are primarily designed to segment salient foreground objects. (b) Our segmentation model leverages vision-language embeddings from CLIP as prompts to guide the SAM model, directing attention to the camouflaged area.

tailored for well-delineated objects and might fail to generalize effectively to camouflaged scenarios, where targets are subtle, indistinct, and visually embedded in complex backgrounds. The lack of alignment between the pretraining objectives and the demands of camouflaged segmentation leads to imprecise localization. In addition, most existing methods do not incorporate explicit edge-aware mechanisms, which are crucial for accurately delineating objects with weak or ambiguous boundaries.

Recent advances in foundation models such as the Segment Anything Model (SAM) [22] have shown remarkable generalization across various segmentation tasks, largely due to their ability to perform prompt-guided segmentation. By using prompts to specify target regions, SAM can adapt its attention to user-defined areas, making it particularly effective for specialized tasks such as camouflaged object segmentation. We propose an adapted SAM architecture tailored for camouflaged object segmentation. As illustrated in Figure 1(b), we integrate CLIP-derived visual and textual embeddings as prompts into the SAM mask decoder, providing task-specific semantic guidance that enhances the ability of the model to focus on the camouflaged targets. Additionally, we enhance the mask decoder with conditional multi-way attention and an edge-aware refinement module to improve boundary precision, effectively handling the indistinct contours characteristic of camouflage.

In the classification stage, most existing methods crop the segmented regions for classification [36, 44, 11] (Figure 2(a)), introducing a domain gap since CLIP is pre-trained on full images. To mitigate the domain gap, we adopt a region-aware classification strategy that replaces hard cropping with a soft spatial prior derived from the segmentation mask, applied via the image’s alpha channel. Our approach preserves the full image context while providing explicit spatial guidance. The predicted segmentation mask serves as a soft spatial prior and is fused with the input image via a lightweight integration module before being processed by the CLIP [37] image encoder. Figure 2 compares the *hard* and *soft* spatial guidance. Additionally, we fine-tune CLIP using a multi-modal prompting strategy similar to [20], jointly optimizing both visual and textual prompts. This enhances semantic alignment and task-specific adaptability, enabling region-aware classification without disrupting global semantics.

Building on these ingredients, we introduce the **C**ascaded **O**pen-vocabulary **C**amouflaged **U**nder**S**tanding network (COCUS), a novel two-stage framework for the OVCOS task that explicitly decouples the process into *segment* and *classify*. In the first stage (segmentation), we use CLIP [37] to extract visual and textual features. These features serve as prompts to the SAM [22] for segmentation. This prompt-based guidance allows SAM to focus more precisely on camouflaged target regions, enhancing localization in visually ambiguous scenes. In the second stage (classification), the segmentation output serves as spatial guidance to refine the integration with the original image, allowing CLIP to perform open-vocabulary classification with improved focus on target regions. By disentangling segmentation and classification, our method enables more accurate

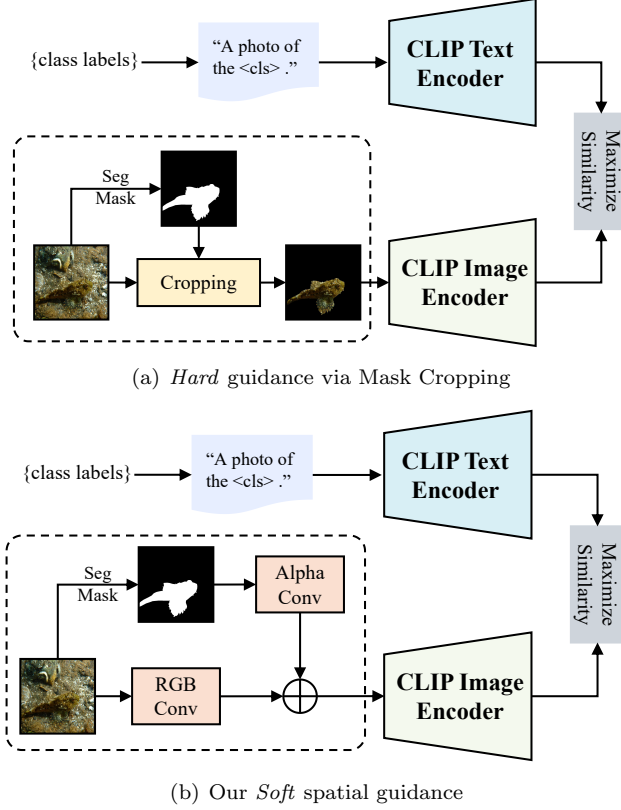


Figure 2: Comparison of mask-guided classification strategies. (a) Mask cropping strategy: applies the segmentation mask to crop the input image before feeding it into the CLIP image encoder. (b) Ours: fuses the segmentation mask with the original image for region-aware classification while retaining full-image context.

semantic interpretation of camouflaged objects through prompt-based guidance segmentation and region-aware classification.

Extensive experiments on the OVCamo [36] benchmark demonstrate the effectiveness of the proposed framework for the OVCOS task. Compared to the strong baseline OVCoser [36], we achieve consistent improvements across all major evaluation metrics, establishing a new state-of-the-art on this challenging benchmark. Moreover, the adapted SAM [22] demonstrates strong performance on the conventional COS task, confirming that CLIP-derived embeddings prompting and edge-aware refinement remains effective in standard closed-set scenarios.

The main contributions of this work are as follows:

- We propose a novel two-stage framework for OVCOS that explicitly decouples segmentation and classification. Our approach employs a prompt-guided segmentation model to generate a mask, which serves as a *soft* spatial guidance for the classification stage while preserving full-image context.
- We propose an adapted SAM as the segmentation model, enhanced for camouflaged object localization by injecting CLIP-derived textual and visual embeddings as prompts. This design provides rich semantic guidance that steers attention toward visually ambiguous regions. Furthermore, we improve SAM’s mask decoder with conditional multi-way attention and edge-aware refinement, improving both spatial accuracy and boundary delineation.
- Extensive experiments on the OVCamo benchmark demonstrate that our method achieves state-of-the-art performance. Moreover, the adapted SAM exhibits strong generalization on the conventional COS task, validating the effectiveness of our framework across both open- and closed-set camouflaged segmentation scenarios.

The remainder of this paper is organized as follows. Section 2 reviews recent advances in open-vocabulary segmentation and camouflaged object understanding. Section 3 presents the proposed framework, detailing its cascaded design, CLIP fine-tuning pipeline and adapted SAM segmentation model. Section 4 presents implementation details, including training settings and architectural configurations, followed by comprehensive experimental results and ablation studies.

2 Related Work

2.1 Vision-Language Models

Vision-language models (VLMs) are neural architectures that learn joint visual-textual representations by embedding both image and text inputs into a shared semantic space. A seminal model in this domain, CLIP [37], jointly learns image and text representations via contrastive learning on large-scale web data, demonstrating strong generalization across open-vocabulary tasks such as object detection [47, 15, 48, 26] and segmentation [10, 44, 11, 12, 27, 43, 46, 30, 42]. However, vanilla CLIP often performs poorly in downstream tasks without task-specific adaptation. To address this limitation, researchers have proposed a variety of fine-tuning approaches. Alpha-CLIP [38] introduces spatially adaptive attention to enhance focus on semantically relevant image regions. CoOp [52] and Co-CoOp [51] optimize textual prompts for better few-shot performance and generalization, respectively. Visual prompt tuning [1] further enhances adaptability by injecting fine-grained prompts into the vision branch. To overcome the limitations of single-modality tuning, recent works [20, 21, 45] adopt multi-modal strategies. FGVP [45] learns patch-level visual prompts to improve alignment across diverse tasks. MaPLe [20] jointly tunes prompts in both visual and textual encoders, preserving CLIP’s generality while enabling task-specific adaptation. In this work, we adopt a multi-modal prompt tuning framework similar to MaPLe to fine-tune CLIP, enhancing semantic alignment for OVCOS.

2.2 Camouflaged Object Segmentation

Camouflaged Object Segmentation (COS) has emerged as a significant research focus in computer vision, with the aim of segmenting objects that visually blend into their surroundings. Unlike traditional tasks such as salient object detection [2, 35, 18, 28, 25] and semantic segmentation [17, 50], COS is inherently more challenging than traditional segmentation tasks due to low object-background contrast, ambiguous boundaries, and high background similarity. It holds practical value in domains such as medical image analysis [13] and agricultural monitoring [29]. COS is typically formulated as a class-agnostic task, focusing on segmenting camouflaged regions within complex visual scenes. The available work [13, 14, 31, 33, 34, 19, 16] to date has demonstrated strong performance on established benchmark datasets [13, 31, 23]. Recent advances have introduced several SAM-based methods [22, 6, 32] adapted for COS, which use prompt tuning and architectural modifications to improve segmentation performance in complex scenes.

2.3 Open-Vocabulary Camouflaged Object Segmentation

Open-Vocabulary Camouflaged Object Segmentation (OVCOS) is a specialized subtask of open-vocabulary segmentation, in which the goal is to segment and recognize camouflaged objects belonging to arbitrary textual categories. Open-vocabulary segmentation aims to align visual and textual representations in a shared embedding space, enabling pixel-level segmentation for unseen or novel categories. Early methods [49] used semantic hierarchies and concept graphs to bridge word concepts and semantic relations. With the rise of VLMs like CLIP [37], recent open-vocabulary segmentation methods have shifted toward leveraging pretrained VLMs to directly connect visual regions with text queries. These approaches follow one-stage and two-stage paradigms. One-stage methods such as MaskCLIP [12] adapt CLIP for segmentation without additional training. SAN [43] enhances feature representations via adapters. CAT-Seg [10] introduces cost aggregation between image and text embeddings. And FC-CLIP [46] employs hierarchical feature fusion. However, these methods often suffer from suboptimal alignment due to CLIP’s image-level representations. The two-stage methods address this by decoupling segmentation and classification. For example, SimSeg [44] uses a cascaded design with MaskFormer [7] for class-agnostic mask generation and CLIP for classification. OVSeg [27] fine-tunes CLIP on diverse and noisy data to improve generalization. In [41], text-to-image diffusion model is

employed for mask generation. While these two-stage framework methods work well on generic objects, they fall short in camouflaged scenarios. OVCOS is especially difficult because low contrast visuals, ambiguous edges, and visually similar backgrounds all contribute to degraded segmentation and classification results. OVCoser [36] is the first to address this task by combining a dedicated camouflaged segmentation model with a CLIP-based classifier in a two-stage pipeline. However, it relies on cropped inputs for classification and does not fully exploit VLM semantics in segmentation.

3 Methodology

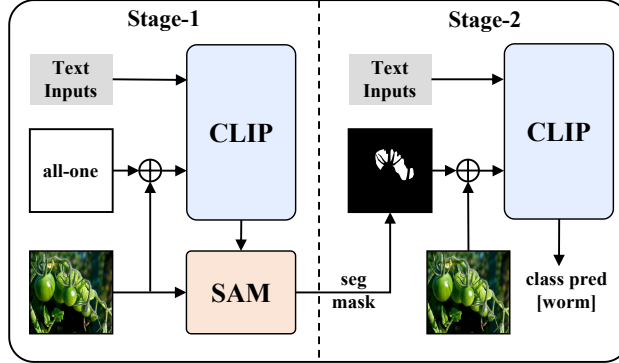


Figure 3: Overview of the cascaded *segment and classify* framework. In Stage-1, the adapted SAM model generates a class-agnostic camouflaged segmentation mask using textual and visual embeddings as prompts. In Stage-2, we use the generated segmentation mask to enable region-aware open-vocabulary classification.

3.1 Problem Definition

Open-Vocabulary Camouflaged Object Segmentation (OVCOS) aims to *segment and classify* camouflaged objects belonging to novel categories unseen during training. Formally, let $\mathcal{C}_{\text{seen}}$ denote the set of categories available during training, and $\mathcal{C}_{\text{unseen}}$ represent the disjoint set of target categories at inference time, such that $\mathcal{C}_{\text{seen}} \cap \mathcal{C}_{\text{unseen}} = \emptyset$. Given an input image I and novel class labels $\mathcal{C}_{\text{unseen}}$, the model is required to produce a segmentation mask M highlighting the camouflaged object and predicts its corresponding class label $\hat{y} \in \mathcal{C}_{\text{unseen}}$.

To address this task, we adopt a *segment-and-classify* strategy. In the first stage, a class-agnostic segmentation model localizes camouflaged regions guided by visual and textual semantics. In the second stage, a vision-language model performs open-vocabulary classification by comparing the visual representation of the segmented regions with textual embeddings of the novel class labels, supporting recognition in an open-set setting.

3.2 Overview

Figure 3 demonstrates the proposed two-stage framework for the OVCOS. During inference, the first stage generates a class-agnostic camouflaged segmentation mask, while the second stage performs open-vocabulary classification based on the segmented regions. We use the same CLIP model for both stages. Our CLIP model accepts a triplet $\{I \in \mathbb{R}^{H \times W \times 3}, M, \text{text}\}$ as input, where I and M are image and mask, and text is a description of the input, with the format of ‘a photo of *< something >*’. The CLIP model outputs visual and textual embeddings, E_v and E_t , which serve as prompts to guide segmentation in the first stage and are used for similarity-based open-vocabulary classification in the second stage. Notably, to ensure a consistent input format across stages, we use an all-one mask in the first stage, while in the second stage, the predicted segmentation mask is used as input.

In the first stage, as shown in Figure 3 (left), we perform segmentation guided by textual and visual embeddings. The inputs consist of an RGB image $I \in \mathbb{R}^{H \times W \times 3}$ and a set of class labels $\mathcal{C} = \{c_1, c_2, \dots, c_N\}$,

where N denotes the number of candidate classes. They are processed by the CLIP [37] model to produce textual embedding E_t and visual embedding E_v optimized for camouflaged object understanding. These embeddings serve as prompts and, together with the image I , are input into the adapted SAM model to guide the prediction of a class-agnostic camouflaged segmentation mask $M \in [0, 1]^{H \times W \times 1}$, effectively localizing the camouflaged object.

As shown in Figure 3 (right), in the second stage, we perform open-vocabulary classification guided by the segmented result. The inputs include the same RGB image I and class labels $\mathcal{C}_{\text{unseen}}$ from the first stage, and the predicted segmentation mask M which is used as an additional input to the CLIP model as spatial guidance. These inputs are processed by the CLIP model as stage one, which now focuses more precisely on the localized object area. The model then outputs a predicted class label $\hat{y} \in \mathcal{C}_{\text{unseen}}$, identifying the category of the camouflaged object. Let $E_t^N \in \mathbb{R}^{N \times d}$, and $E_v \in \mathbb{R}^{1 \times d}$ be the textual and visual embeddings, where $d = 768$ is the feature dimension, we first calculate the similarity scores $S \in \mathbb{R}^N$:

$$S = E_t^N \cdot (E_v)^T. \quad (1)$$

During training, we first fine-tune our CLIP [37] model by optimizing learnable prompts in both the language and vision branches to enhance its sensitivity to camouflaged objects, while keeping all encoder parameters fixed. Figure 4 illustrates the fine-tuning pipeline of our CLIP. After fine-tuning, we freeze the CLIP model as a feature extractor and train the SAM [22] model using visual-textual features from CLIP as prompts. The details of the CLIP fine-tuning process are provided in Section 3.3, and the architecture of the adapted SAM is described in Section 3.4.

3.3 CLIP Fine-Tuning Pipeline

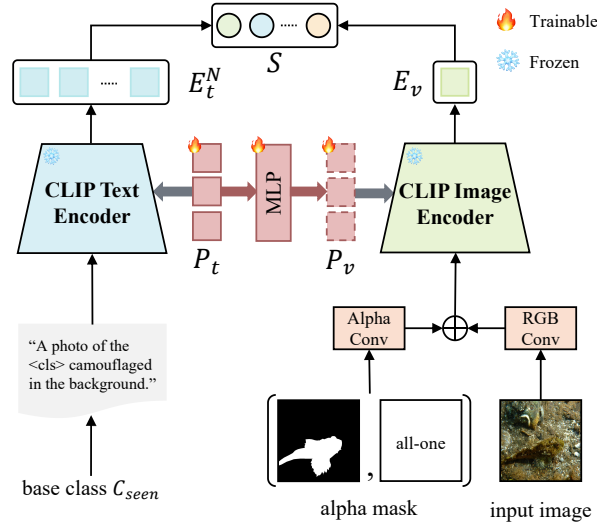


Figure 4: The CLIP fine-tuning pipeline. The language branch encodes base class labels $\mathcal{C}_{\text{seen}}$ with a camouflage-specific prompt template and learnable textual prompts P_t to obtain textual embeddings E_t^N . The vision branch fuses features from the input image and alpha mask, combined with visual prompts P_v injected via an MLP, and passes them to the frozen CLIP image encoder to obtain visual embedding E_v . Similarity scores S are computed by aligning E_t^N and E_v in a shared space.

We fine-tune the CLIP model using a multi-modal prompting strategy to enhance its ability to capture subtle semantic cues for camouflaged object segmentation, as shown in Figure 4. Our CLIP variant is a modified version of Alpha-CLIP [38]. Previous prompting strategies in CLIP [37] typically operate on visual or textual modality. Language-only prompt tuning methods [44, 52, 51] optimize learnable prompts solely in the language branch, while visual-only approaches [12, 27, 45] inject prompts exclusively into the vision branch. In this work, we adopt a multi-modal prompting strategy, following [20], which jointly optimizes both textual and visual prompts to enhance multi-modal alignment and better adapt to task-specific objectives.

Specifically, we append learnable textual prompts P_t to the language branch and generate the corresponding visual prompts P_v , which are produced by conditioning on the textual prompts through an MLP injector. The language and textual prompts P_t and P_v are shown in middle of Figure 4. During fine-tuning, only the textual prompts and injector parameters are updated, while the rest of CLIP model remains frozen. This lightweight strategy promotes efficient adaptation and enables improved semantic alignment across modalities. Below, we outline the fine-tuning pipeline of CLIP model.

The fine-tuning pipeline begins with the language branch, where the base class labels $\mathcal{C}_{\text{seen}}$ are formatted using the prompt template “A photo of the $\langle class \rangle$ camouflaged in the background.” and enriched with learnable textual prompts P_t . These are processed by the frozen CLIP text encoder to produce textual embeddings $E_t^N \in \mathbb{R}^{N \times 768}$, where N is the number of the base classes.

Currently, in the vision branch, the input RGB image $I \in \mathbb{R}^{H \times W \times 3}$ is combined with an auxiliary alpha mask $A \in \mathbb{R}^{H \times W \times 1}$. The alpha mask A is randomly selected as either the all-one mask A_J or the ground-truth segmentation mask A_{gt} , each with equal probability. This enables the CLIP model to optionally accept a mask as input, defaulting to an all-one matrix when no mask is provided—for example, during the first stage of segmentation.

The image I and the alpha mask A are separately processed through dedicated convolutional layers, e.g. AlphaConv and RGBConv in Figure 4, to extract modality-specific features, which are then fused to form the visual representation. This fused representation, along with the injected visual prompts P_v generated by a lightweight MLP-based injector, is fed into the frozen CLIP image encoder to obtain the visual embedding $E_v \in \mathbb{R}^{1 \times 768}$.

Finally, the textual and visual embeddings are used to compute the similarity score as defined in Equation (1), which is used to calculate a cross-entropy loss against the ground-truth class labels.

3.4 Adapted SAM

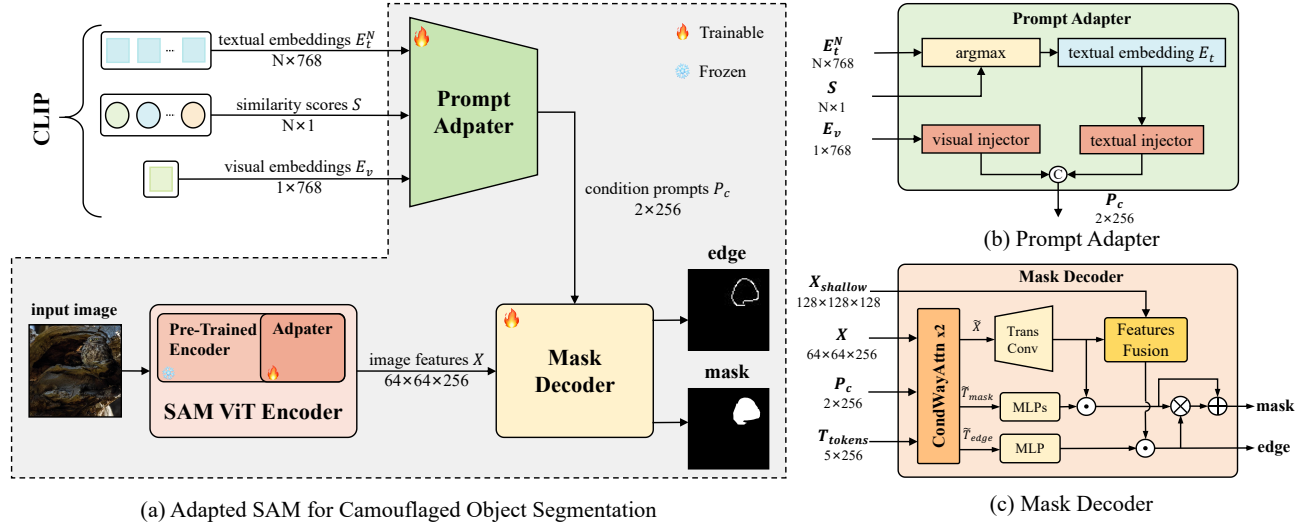


Figure 5: **Overview of the adapted SAM framework.** (a) *Adapted SAM for COS*: Our fine-tuned CLIP provides textual embeddings E_t^N , visual embedding E_v , and similarity scores S , which are projected into condition prompts P_c via a Prompt Adapter. Image features X extracted by SAM ViT encoder are refined by adapters. The Mask Decoder integrates X and P_c to predict the segmentation mask M and edge map E , enabling precise localization. (b) *Prompt Adapter*: Selects the most relevant textual embedding based on S , and projects both E_t and E_v into a unified condition space via lightweight MLPs to guide the decoder. (c) *Adapted Mask Decoder*: Combines image features X , condition prompts P_c , and output tokens T_{tokens} to produce accurate masks and edge maps, improving segmentation in camouflaged scenes.

We build upon the SAM [22] to address the unique challenges of COS. While SAM excels at general-purpose segmentation, it struggles with the subtle visual cues and semantic ambiguities inherent to camouflaged objects. To overcome these limitations, as shown in Figure 5(a), we adapt SAM by incorporating

textual and visual embeddings guidance and edge-aware enhancements for improved segmentation.

Specifically, we integrate our fine-tuned CLIP model with SAM to provide semantic context. The CLIP model produces textual embeddings $E_t^N \in \mathbb{R}^{N \times 768}$, a visual embedding $E_v \in \mathbb{R}^{1 \times 768}$, and similarity scores $S \in \mathbb{R}^{N \times 1}$. These embeddings are further processed by a prompt adapter, which projects them into condition prompts $P_c \in \mathbb{R}^{2 \times 256}$, providing high-level semantic guidance into the segmentation pipeline.

In parallel, the SAM ViT encoder extracts image features $X \in \mathbb{R}^{64 \times 64 \times 256}$ from the input image. To adapt SAM to camouflage-specific cues, we introduce lightweight adapter modules that refine the image features X while keeping the backbone frozen.

Finally, the refined image features X and condition prompts P_c are fused within a mask decoder, which outputs a segmentation mask $M \in \mathbb{R}^{H \times W \times 1}$ and an edge map $E \in \mathbb{R}^{H \times W \times 1}$. The integration of refined image features and condition prompts within the decoder ensures accurate object localization and precise boundary delineation.

3.4.1 Prompt Adapter

The Prompt Adapter refines textual and visual embeddings from our fine-tuned CLIP to generate condition prompts for segmentation guidance, as shown in Figure 5(b). Given textual embeddings $E_t^N = \{e_t^1, e_t^2, \dots, e_t^N\}$, visual embedding E_v , and similarity scores $S = \{s_1, s_2, \dots, s_N\}$, the adapter selects the textual embedding corresponding to the highest similarity score:

$$i^* = \arg \max_i s_i, \quad E_t = e_t^{i^*}. \quad (2)$$

The selected textual embedding E_t and visual embedding E_v are projected into a shared 256-dimensional condition space using lightweight MLP-based injectors. The resulting condition prompts $P_c \in \mathbb{R}^{2 \times 256}$ provide high-level semantic and visual guidance to the segmentation mask decoder, enhancing object localization and boundary accuracy. Formally, this is defined as:

$$P_t = \text{MLP}_{\text{text}}(E_t), \quad P_v = \text{MLP}_{\text{vis}}(E_v), \quad (3)$$

$$P_c = [P_t, P_v] \in \mathbb{R}^{2 \times 256}, \quad (4)$$

where $\text{MLP}_{\text{text}}(\cdot)$ and $\text{MLP}_{\text{vis}}(\cdot)$ denote the projection functions for textual and visual features, respectively.

3.4.2 Mask Decoder

We adapt the original SAM [22] mask decoder to address the specific challenges of camouflaged object segmentation by introducing semantic conditioning and edge-aware enhancements. The modified decoder integrates multi-level image features X , condition prompts P_c , and output tokens T_{tokens} , including mask tokens T_{mask} and an edge token T_{edge} , to localize objects precisely and refine boundaries accurately, as shown in Figure 5(c).

We first apply two Conditional Multi-Way Attention (CondWayAttn $\times 2$) modules to model the interactions among image features, condition prompts, and tokens. Each block enables dense bidirectional information flow between the image features, condition prompts, and output tokens. Specifically, it includes image-to-token and image-to-condition attention to incorporate visual context, token-to-condition and token-to-image attention to align output tokens with semantic and spatial cues, as well as token self-attention and an MLP layer to capture intra-token dependencies and perform feature transformation. The enhanced outputs are computed as:

$$\tilde{X}, \tilde{T}_{\text{mask}}, \tilde{T}_{\text{edge}} = \text{CondWayAttn}(X, P_c, T_{\text{token}}). \quad (5)$$

The attention-enhanced features \tilde{X} are first upsampled using a transposed convolution to restore spatial resolution. To recover fine-grained details, these features are fused with shallow image features X_{shallow} through a fusion block defined as:

$$X_{\text{fusion}} = \text{TConv}(\tilde{X}) + \text{Conv}(\text{ReLU}(\text{Norm}(\text{Conv}(X_{\text{shallow}})))). \quad (6)$$

The attention-enhanced mask and edge tokens are then projected via task-specific MLPs. The coarse segmentation mask is computed by element-wise multiplication of the mask token with the upsampled features:

$$M_{\text{coarse}} = \text{MLPs}(\tilde{T}_{\text{mask}}) \odot \text{TConv}(\tilde{X}). \quad (7)$$

Similarly, the edge map is predicted by interacting the edge token with the fused feature map:

$$E = \text{MLP}(\tilde{T}_{\text{edge}}) \odot X_{\text{fusion}}. \quad (8)$$

The final refined mask incorporates edge guidance by multiplying the coarse mask with the edge map, followed by a residual addition:

$$M_{\text{fine}} = M_{\text{coarse}} + (M_{\text{coarse}} \otimes E). \quad (9)$$

This edge-guided refinement enhances boundary accuracy while preserving regional consistency, effectively handling low-contrast and subtle camouflaged structures. The effectiveness of this module is evidenced by the ablation study in Section 4.3.2.

4 Experiments

4.1 Implement Details

4.1.1 Datasets

We evaluate our method on two tasks: Open-Vocabulary Camouflaged Object Segmentation (OVCOS) and Camouflaged Object Segmentation (COS).

For the OVCOS task, all experiments are conducted on the OVCamo [36] dataset, a benchmark specifically curated for this setting. It comprises 11,483 images sourced from various publicly available datasets, covering 75 camouflaged object categories embedded in complex natural scenes. To enable open-vocabulary evaluation, the dataset is divided into two disjoint subsets by category: the training set $\mathcal{D}_{\text{train}}$ includes 7,713 images from 14 seen categories, while the test set $\mathcal{D}_{\text{test}}$ contains 3,770 images from 61 unseen categories, following an approximate 7:3 split.

For the COS task, we evaluate on three widely used benchmarks: CAMO [23], COD10K [13], and NC4K [31]. A total of 4,040 images from CAMO and COD10K are used for training. We conduct evaluation on the remaining images from these datasets, as well as the entire NC4K set. The detailed statistics of all datasets, covering training/testing splits, are presented in Table 1.

Table 1: Summary of datasets used for OVCOS and COS tasks.

Dataset	Task	Total	Train	Test	Categories
OVCamo [36]	OVCOS	11,483	7,713	3,770	75 (14/61)
CAMO [23]	COS	1,250	1,000	250	–
COD10K [13]	COS	5,066	3,040	2,026	–
NC4K [31]	COS	4,121	–	4,121	–

4.1.2 Evaluation Metrics

To ensure fair and comprehensive evaluation of OVCOS performance, we adopt a set of evaluation metrics tailored for OVCOS, which are adapted from those originally proposed in the camouflaged scene understanding task [13, 8]. Specifically, we use six metrics: class structure measure cS_m , class weighted F-measure cF_β^w , class mean absolute error $cMAE$, class standard F-measure cF_β , class enhanced alignment measure cE_m , and class intersection over union $cIoU$. These metrics are standard in the open-vocabulary segmentation literature [10, 43, 46, 44, 27, 53], jointly assessing classification accuracy and segmentation quality for a balanced evaluation of model performance.

For the COS task, we follow established protocols [13] and adopt four commonly used metrics: structure measure S_α , enhanced alignment measure E_ϕ , weighted F-measure F_β^ω , and mean absolute error MAE .

Among the four standard COS metrics, S_α , E_ϕ , and F_β^ω evaluate structural and region-aware similarity between predictions and ground truth, where higher values indicate better performance. Conversely, MAE measures pixel-wise error, with lower values indicating better accuracy.

Table 2: Comparison of our method with state-of-the-art CLIP-based open-vocabulary segmentation approaches and the baseline model OVCoser on the OVCamo [36] dataset. The bolded values indicate the results of our method, which achieves the best overall performance. The second best is underlined.

Model	VLM	Train Set	Finetune	$cS_m \uparrow$	$cF_\beta^\omega \uparrow$	$cMAE \downarrow$	$cF_\beta \uparrow$	$cE_m \uparrow$	$cIoU \uparrow$
SimSeg [44]	CLIP-ViT-B/16 [37]	COCO-Stuff [4]	OVCamo [36]	0.098	0.071	0.852	0.081	0.128	0.0
OVSeg [27]	CLIP-ViT-L/14 [37]	COCO-Stuff [4]	OVCamo [36]	0.164	0.131	0.763	0.147	0.208	0.123
ODISE [41]	CLIP-ViT-L/14 [37]	COCO-Stuff [4]	OVCamo [36]	0.182	0.125	0.691	0.219	0.309	0.189
SAN [43]	CLIP-ViT-L/14 [37]	COCO-Stuff [4]	OVCamo [36]	0.321	0.216	0.550	0.236	0.331	0.204
FC-CLIP [46]	CLIP-ConvNeXt-L [9]	COCO-Stuff [4]	OVCamo [36]	0.124	0.074	0.798	0.088	0.162	0.072
CAT-Seg [10]	CLIP-ViT-L/14 [37]	COCO-Stuff [4]	OVCamo [36]	0.185	0.094	0.702	0.110	0.185	0.088
OVCoser [36]	CLIP-ConvNeXt-L [9]	OVCamo [36]	–	<u>0.579</u>	<u>0.490</u>	<u>0.336</u>	<u>0.520</u>	<u>0.616</u>	<u>0.443</u>
Ours	Our Fine-Tuned CLIP	OVCamo [36]	–	0.668	0.615	0.265	0.631	0.697	0.568

Table 3: Quantitative COS task comparison results on three benchmark datasets. The best performance per metric is highlighted in bold, and the second best is underlined.

Method	CAMO [23]				COD10K [13]				NC4K [31]			
	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	$MAE \downarrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	$MAE \downarrow$	$S_\alpha \uparrow$	$E_\phi \uparrow$	$F_\beta^\omega \uparrow$	$MAE \downarrow$
SINet [14]	0.751	0.771	0.606	0.100	0.771	0.806	0.551	0.051	0.808	0.883	0.768	0.058
RankNet [31]	0.712	0.791	0.583	0.104	0.767	0.861	0.611	0.045	0.840	0.904	0.802	0.048
PFNet [33]	0.782	0.852	0.695	0.085	0.800	0.868	0.660	0.040	0.829	0.887	0.784	0.053
SINetV2 [13]	0.820	0.882	0.743	0.070	0.815	0.887	0.680	0.037	0.847	0.903	0.770	0.048
ZoomNet [34]	0.820	0.892	0.752	0.066	0.838	0.911	0.729	0.029	0.853	0.912	0.784	0.043
SegMaR [19]	0.815	0.872	0.742	0.071	0.833	0.895	0.724	0.033	0.841	0.905	0.781	0.046
DGNet [16]	0.839	0.901	0.769	0.057	0.822	0.896	0.693	0.033	0.857	0.911	0.784	0.042
SAM [22]	0.684	0.687	0.606	0.132	0.783	0.798	0.701	0.050	0.767	0.776	0.696	0.078
SAM-Adapter [6]	<u>0.847</u>	0.873	0.765	0.070	<u>0.883</u>	<u>0.918</u>	<u>0.801</u>	<u>0.025</u>	–	–	–	–
MedSAM [32]	0.820	0.904	<u>0.779</u>	<u>0.065</u>	0.841	0.917	0.751	0.033	<u>0.866</u>	<u>0.929</u>	<u>0.821</u>	<u>0.041</u>
Ours	0.865	<u>0.902</u>	0.789	0.057	0.905	0.947	0.845	0.019	0.904	0.933	0.852	0.031

4.1.3 Training Details

All experiments are conducted on a system equipped with two NVIDIA RTX 3090Ti GPUs running Ubuntu 20.04. Our framework is implemented on PyTorch with CUDA 11.8 acceleration.

During the CLIP model fine-tuning, we adopt a multi-modal prompting strategy on the pre-trained ViT-L/14 Alpha-CLIP model [38]. The model is trained on the OVCamo [36] dataset for 10 epochs using stochastic gradient descent (SGD) with a learning rate of 0.0035 and a batch size of 8 on a single GPU, following the setup in [20]. Additionally, the input alpha mask is randomly selected as either an all-one mask or the ground truth segmentation mask with equal probability, balancing global context encoding and localized focus.

During adapted SAM training, the fine-tuned CLIP is integrated into our adapted SAM architecture, based on the ViT-H variant of SAM [22]. The network is trained for 20 epochs using the Adam optimizer with an initial learning rate of 2×10^{-4} , decayed via cosine annealing. Training is conducted on two GPUs with a batch size of 2 and completes in approximately 24 hours.

4.2 Experimental Results

4.2.1 Quantitative Results on OVCOS

To thoroughly evaluate the effectiveness of our proposed framework, we compare it with recent state-of-the-art open-vocabulary segmentation methods, including CAT-Seg [10], SAN [43], SimSeg [44], OVSeg [27], FC-CLIP [46], ODISE [41] and the baseline OVCoser [36]. For a fair comparison, all models are trained or

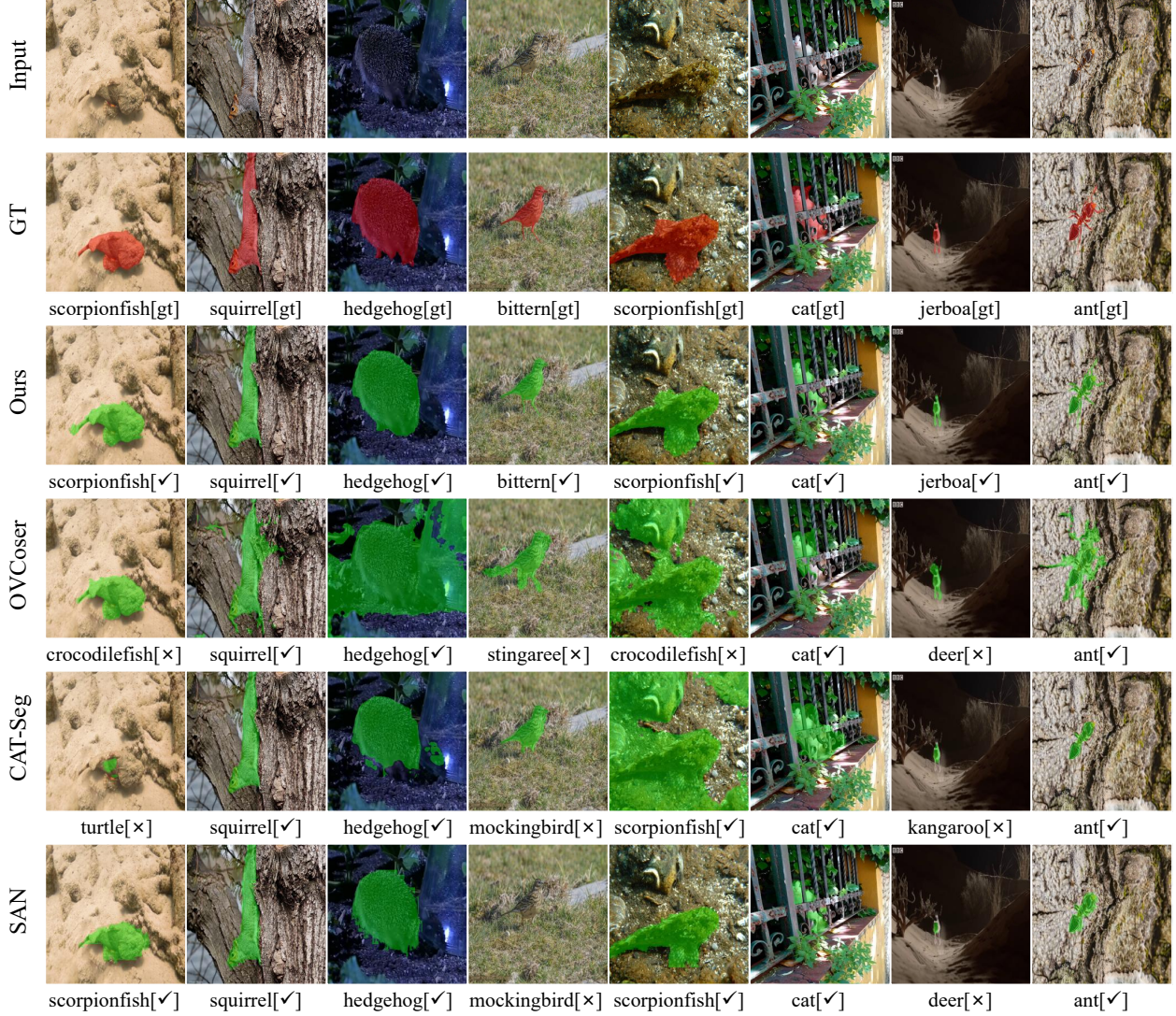


Figure 6: Qualitative comparison between our method and CLIP-based baselines on OVCamo [36]. Each column depicts the input image, segmentation result, and predicted label. Predicted label is shown below each segmentation result, where [✓] indicates correct prediction and [✗] denotes an incorrect one.

fine-tuned on the OVCamo [36] dataset. We adopt the large variants of such approaches wherever available, except for SimSeg, which is only released in its base form. As shown in Table 2, our method consistently outperforms all competitors across multiple evaluation metrics.

While open-vocabulary segmentation methods such as SAN [43], OVSeg [27], and CAT-Seg [10] benefit from large-scale pretraining, they lack task-specific adaptation, resulting in limited performance on OVCOS (e.g., OVSeg: 0.164 cS_m , 0.123 $cIoU$). The baseline OVCoser [36] improves results by integrating camouflage segmentation with CLIP-based classification, achieving 0.579 cS_m and 0.443 $cIoU$, but it does not fine-tune vision-language embeddings or incorporate semantic guidance into segmentation.

Compared to the results of existing methods shown in Table 2, our framework leverages fine-tuned CLIP and a task-adapted SAM to enhance both segmentation and classification. Ours achieves state-of-the-art results, surpassing the baseline OVCoser [36] by notable margins across all metrics: +8.9% in cS_m , +12.5% in $cIoU$, +12.5% in cF_β^w , +11.1% in cF_β , +8.1% in cE_m , and a reduction of 7.1% in $cMAE$. These results highlight the effectiveness of our cascaded design and cross-modal semantic conditioning in tackling the OVCOS challenge.

Table 4: Classification performance comparison of different CLIP models on the OVCamo [36] test set.

Model	Alpha	Top-1↑	Top-5↑
CLIP-ConvNeXt-L [9]	–	0.6944	0.8918
CLIP-ViT-L/14 [37]	–	0.7040	0.8915
Alpha-CLIP [38]	all one	0.6934	0.8849
Alpha-CLIP [38]	gt	0.7467	0.9456
Ours	all one	0.7462	0.9003
Ours	gt	0.7859	0.9497

4.2.2 Quantitative Results on COS

As shown in Table 3, our adapted SAM model achieves competitive performance across three widely used COS benchmarks: CAMO [23], COD10K [13], and NC4K [31]. Compared to both traditional non-SAM-based methods [13, 14, 31, 33, 34, 16, 19] and recent SAM-based approaches [22, 6, 32], our model consistently outperforms across all datasets.

Specifically, the adapted SAM ranks first on 11 out of 12 evaluation metrics and second on the remaining one, demonstrating strong generalization and robustness in diverse camouflage scenarios. Ours achieves notable improvements in structure-aware metrics (S_α , E_ϕ), region-aware precision (F_β^w), and pixel-level accuracy (MAE), particularly on the COD10K and NC4K datasets. These results highlight the effectiveness of our edge-enhanced architecture and prompt-guided segmentation in capturing fine-grained boundaries and ensuring semantic consistency.

4.2.3 Qualitative Results of OVCOS

To further validate our quantitative findings, we present qualitative comparisons in Figure 6. Ours consistently delivers superior segmentation quality, accurately delineating camouflaged objects with well-preserved shapes and precise boundaries—even in low-contrast and cluttered backgrounds. Compared to other methods, our approach better maintains object integrity and minimizes background leakage, demonstrating enhanced robustness in challenging camouflage scenarios.

In terms of classification, ours reliably predicts correct categories across diverse samples, outperforming prior methods that often misclassify visually ambiguous targets. This classification accuracy improvement stems from our region-aware classification strategy, which integrates segmentation masks as alpha masks into the fine-tuned CLIP model. Combined with multi-modal prompting and edge-aware decoding, ours achieves high fidelity in both localization and recognition under open-vocabulary conditions.

Table 5: Ablation results showing the effectiveness of our fine-tuned CLIP on OVCOS performance.

Model	VLM	$cS_m \uparrow$	$cF_\beta^w \uparrow$	$cMAE \downarrow$	$cF_\beta \uparrow$	$cE_m \uparrow$	$cIoU \uparrow$
COCOSIP-ConvNeXt-L [9]		0.567	0.518	0.375	0.534	0.591	0.481
COCOSLIP-ViT-L/14 [37]		0.580	0.536	0.353	0.551	0.617	0.503
COCUSAlpha-CLIP [38]		0.639	0.589	0.299	0.603	0.668	0.545
COCUSr Fine-Tuned CLIP		0.668	0.615	0.265	0.631	0.697	0.568

Table 6: Ablation study of Conditional Multi-Way Attention (CMA) and Edge Enhancement (EDE) in the adapted mask decoder.

Model	$cS_m \uparrow$	$cF_\beta^w \uparrow$	$cMAE \downarrow$	$cF_\beta \uparrow$	$cE_m \uparrow$	$cIoU \uparrow$
Baseline	0.644	0.599	0.281	0.610	0.651	0.549
+ EDE	0.650	0.605	0.278	0.615	0.666	0.554
+ CMA	0.652	0.607	0.273	0.621	0.683	0.551
+ CMA, EDE	0.668	0.615	0.265	0.631	0.697	0.568

4.3 Ablation Study

4.3.1 Effectiveness of the Fine-Tuned CLIP

To assess the effectiveness of our fine-tuned CLIP model, we conduct comparative experiments on classification accuracy using several CLIP [37] variants, including CLIP-ConvNeXt-L [9], CLIP-ViT-L/14 [37], the original Alpha-CLIP [38], and ours. Evaluation results on the OVCamo [36] test set are presented in Table 4, where *all one* refers to using an all-one alpha mask, and *gt* denotes the ground-truth segmentation mask.

The results clearly show that Alpha-CLIP [38] with *gt* alpha mask significantly outperforms both CLIP-ConvNeXt-L [9] and CLIP-ViT-L/14 [37], the backbone models used in OVCoser [36], demonstrating the advantage of incorporating an auxiliary alpha mask for open-vocabulary classification. Moreover, ours further improves classification accuracy, reducing the performance gap between the *all one* and *gt* settings. Ours achieves the highest accuracy when using the *gt* mask, emphasizing the benefit of integrating task-specific semantic information during inference.

Additionally, we evaluate the impact of integrating the fine-tuned CLIP into ours framework. As shown in Table 5, our fine-tuned CLIP model consistently outperforms the original Alpha-CLIP across all key metrics, confirming its effectiveness in enhancing semantic representation and improving overall model performance on the OVCOS task.

4.3.2 Impact of Adapted Mask Decoder

As shown in Table 6, we perform ablation studies on the OVCamo [36] dataset to assess the effectiveness of the proposed Conditional Multi-Way Attention (CMA) and Edge Enhancement (EDE) modules in our adapted mask decoder. Beginning with a baseline SAM [22] model equipped with lightweight adapters—corresponding to the original SAM mask decoder without either enhancement—we observe consistent performance improvements when incorporating CMA or EDE individually.

Specifically, adding the EDE module (+ *EDE*) leads to notable gains in contour-sensitive metrics such as cF_{β}^w and $cIoU$, indicating that explicit edge modeling enhances boundary precision. In contrast, introducing the CMA module (+ *CMA*) results in stronger improvements in semantic-aware metrics like cE_m and cF_{β} , demonstrating that conditional attention effectively enriches textual-visual feature fusion.

When both modules are combined, our method achieves the best performance across all metrics, underscoring the complementary strengths of semantic conditioning and edge-aware refinement. These findings confirm the importance of both enhancements in improving segmentation quality for challenging camouflaged object scenarios.

5 Conclusion

In this paper, we present COCUS, a two-stage framework for OVCOS that explicitly decouples segmentation and classification. In the first stage, visual and textual embeddings are extracted using our fine-tuned CLIP model. These embeddings guide an adapted SAM with a redesigned mask decoder to enhance object localization and boundary precision. In the second stage, the predicted segmentation mask is fused with the input image to guide the attention of the model toward the target regions, enabling region-aware classification without relying on cropped inputs. Extensive experiments on both OVCOS and COS benchmarks show that ours outperforms existing open-vocabulary segmentation methods. The adapted SAM also achieves superior results on the COS benchmarks. These experiments confirm the benefits of our two-stage framework and edge-aware enhancements in complex camouflage scenarios.

References

- [1] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modifying pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 3(11-12):3, 2022.
- [2] Ali Borji, Ming-Ming Cheng, Qibin Hou, Huaizu Jiang, and Jia Li. Salient object detection: A survey. *Computational Visual Media*, 5:117–150, 2019.

- [3] Maxime Bucher, Tuan-Hung VU, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [4] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2018.
- [5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision*, September 2018.
- [6] Tianrun Chen, Lanyun Zhu, Chaotao Deng, Runlong Cao, Yan Wang, Shangzhan Zhang, Zejian Li, Lingyun Sun, Ying Zang, and Papa Mao. Sam-adapter: Adapting segment anything in underperformed scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 3367–3375, October 2023.
- [7] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *Advances in Neural Information Processing Systems*, volume 34, pages 17864–17875, 2021.
- [8] Xuelian Cheng, Huan Xiong, Deng-Ping Fan, Yiran Zhong, Mehrtash Harandi, Tom Drummond, and Zongyuan Ge. Implicit motion handling for video camouflaged object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13864–13873, June 2022.
- [9] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829, June 2023.
- [10] Seokju Cho, Heeseong Shin, Sunghwan Hong, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4113–4123, June 2024.
- [11] Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin Dai. Decoupling zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11583–11592, June 2022.
- [12] Zheng Ding, Jieke Wang, and Zhuowen Tu. Open-vocabulary universal image segmentation with maskclip. In *International Conference on Machine Learning*, pages 8090–8102, 2023.
- [13] Deng-Ping Fan, Ge-Peng Ji, Ming-Ming Cheng, and Ling Shao. Concealed object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6024–6042, 2022.
- [14] Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Camouflaged object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2020.
- [15] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *International Conference on Learning Representations*, 2022.
- [16] Ge-Peng Ji, Deng-Ping Fan, Yu-Cheng Chou, Dengxin Dai, Alexander Liniger, and Luc Van Gool. Deep gradient learning for efficient camouflaged object detection. *Machine Intelligence Research*, 20(1):92–108, 2023.
- [17] Wei Ji, Jingjing Li, Cheng Bian, Zongwei Zhou, Jiaying Zhao, Alan L Yuille, and Li Cheng. Multispectral video semantic segmentation: A benchmark dataset and baseline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1094–1104, 2023.
- [18] Wei Ji, Ge Yan, Jingjing Li, Yongri Piao, Shunyu Yao, Miao Zhang, Li Cheng, and Huchuan Lu. Dmra: Depth-induced multi-scale recurrent attention network for rgb-d saliency detection. *IEEE Transactions on Image Processing*, 31:2321–2336, 2022.

- [19] Qi Jia, Shuilian Yao, Yu Liu, Xin Fan, Risheng Liu, and Zhongxuan Luo. Segment, magnify and reiterate: Detecting camouflaged objects the hard way. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4713–4722, June 2022.
- [20] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19113–19122, June 2023.
- [21] Muhammad Uzair Khattak, Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Self-regulating prompts: Foundational model adaptation without forgetting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15190–15200, October 2023.
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, October 2023.
- [23] Trung-Nghia Le, Tam V. Nguyen, Zhongliang Nie, Minh-Triet Tran, and Akihiro Sugimoto. Anabran network for camouflaged object segmentation. *Computer Vision and Image Understanding*, 184:45–56, 2019.
- [24] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2022.
- [25] Jingjing Li, Wei Ji, Miao Zhang, Yongri Piao, Huchuan Lu, and Li Cheng. Delving into calibrated depth for accurate rgb-d salient object detection. *International Journal of Computer Vision*, 131(4):855–876, 2023.
- [26] Shuai Li, Minghan Li, Pengfei Wang, and Lei Zhang. Opensd: Unified open-vocabulary segmentation and detection. *arXiv preprint arXiv:2312.06703*, 2023.
- [27] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7061–7070, June 2023.
- [28] Jiang-Jiang Liu, Qibin Hou, Zhi-Ang Liu, and Ming-Ming Cheng. Poolnet+: Exploring the potential of pooling for salient object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):887–904, 2022.
- [29] Liu Liu, Rujing Wang, Chengjun Xie, Po Yang, Fangyuan Wang, Sud Sudirman, and Wancai Liu. Pest-net: An end-to-end deep learning approach for large-scale multi-class pest detection and classification. *IEEE Access*, 7:45301–45312, 2019.
- [30] Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. SegCLIP: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 23033–23044, 23–29 Jul 2023.
- [31] Yunqiu Lv, Jing Zhang, Yuchao Dai, Aixuan Li, Bowen Liu, Nick Barnes, and Deng-Ping Fan. Simultaneously localize, segment and rank the camouflaged objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11591–11601, June 2021.
- [32] Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024.
- [33] Haiyang Mei, Ge-Peng Ji, Ziqi Wei, Xin Yang, Xiaopeng Wei, and Deng-Ping Fan. Camouflaged object segmentation with distraction mining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8772–8781, June 2021.

- [34] Youwei Pang, Xiaoqi Zhao, Tian-Zhu Xiang, Lihe Zhang, and Huchuan Lu. Zoom in and out: A mixed-scale triplet network for camouflaged object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2160–2170, June 2022.
- [35] Youwei Pang, Xiaoqi Zhao, Lihe Zhang, and Huchuan Lu. Multi-scale interactive network for salient object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2020.
- [36] Youwei Pang, Xiaoqi Zhao, Jiaming Zuo, Lihe Zhang, and Huchuan Lu. Open-vocabulary camouflaged object segmentation. In *European Conference on Computer Vision*, pages 476–495, 2024.
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763, 18–24 Jul 2021.
- [38] Zeyi Sun, Ye Fang, Tong Wu, Pan Zhang, Yuhang Zang, Shu Kong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. Alpha-clip: A clip model focusing on wherever you want. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13019–13029, June 2024.
- [39] Yongqin Xian, Subhabrata Choudhury, Yang He, Bernt Schiele, and Zeynep Akata. Semantic projection network for zero- and few-label semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2019.
- [40] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *Advances in Neural Information Processing Systems*, volume 34, pages 12077–12090, 2021.
- [41] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2955–2966, June 2023.
- [42] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu Qiao, and Weidi Xie. Learning open-vocabulary semantic segmentation models from natural language supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2935–2944, June 2023.
- [43] Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xiang Bai. Side adapter network for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2945–2954, June 2023.
- [44] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language model. In *European Conference on Computer Vision*, pages 736–753, 2022.
- [45] Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang, and Jian Yang. Fine-grained visual prompting. In *Advances in Neural Information Processing Systems*, volume 36, pages 24993–25006, 2023.
- [46] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- [47] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary detr with conditional matching. In *European conference on computer vision*, pages 106–122, 2022.
- [48] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14393–14402, June 2021.
- [49] Hang Zhao, Xavier Puig, Bolei Zhou, Sanja Fidler, and Antonio Torralba. Open vocabulary scene parsing. In *Proceedings of the IEEE International Conference on Computer Vision*, Oct 2017.

- [50] Xiaoqi Zhao, Youwei Pang, Jiaxing Yang, Lihe Zhang, and Huchuan Lu. Multi-source fusion and automatic predictor selection for zero-shot video object segmentation. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 2645–2653, 2021.
- [51] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16816–16825, June 2022.
- [52] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.
- [53] Chaoyang Zhu and Long Chen. A survey on open-vocabulary detection and segmentation: Past, present, and future. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):8954–8975, 2024.