

# RSPrompter: Learning to Prompt for Remote Sensing Instance Segmentation based on Visual Foundation Model

Keyan Chen<sup>1</sup>, Chenyang Liu<sup>1</sup>, Hao Chen<sup>2</sup>, Haotian Zhang<sup>1</sup>, Wenyuan Li<sup>3</sup>, Zhengxia Zou<sup>1</sup>, and Zhenwei Shi<sup>1,\*</sup>

<sup>1</sup>Beihang University, <sup>2</sup>Shanghai AI Laboratory, <sup>3</sup>The University of Hong Kong

**Abstract**—Leveraging the extensive training data from SA-1B, the Segment Anything Model (SAM) demonstrates remarkable generalization and zero-shot capabilities. However, as a category-agnostic instance segmentation method, SAM heavily relies on prior manual guidance, including points, boxes, and coarse-grained masks. Furthermore, its performance in remote sensing image segmentation tasks remains largely unexplored and unproven. In this paper, we aim to develop an automated instance segmentation approach for remote sensing images, based on the foundational SAM model and incorporating semantic category information. Drawing inspiration from prompt learning, we propose a method to learn the generation of appropriate prompts for SAM. This enables SAM to produce semantically discernible segmentation results for remote sensing images, a concept we have termed RSPrompter. We also propose several ongoing derivatives for instance segmentation tasks, drawing on recent advancements within the SAM community, and compare their performance with RSPrompter. Extensive experimental results, derived from the WHU building, NWPU VHR-10, and SSDD datasets, validate the effectiveness of our proposed method. The code for our method is publicly available at <https://kyanchen.github.io/RSPrompter>.

**Index Terms**—Remote sensing images, foundation model, SAM, instance segmentation, prompt learning.

## I. INTRODUCTION

Instance segmentation is a pivotal task in the analysis of remote sensing images, facilitating a semantic-level understanding of each instance within the images. This process yields crucial information regarding the location (where), category (what), and shape (how) of each object [1–9]. The accurate perception and comprehension of surfaces in remote sensing images significantly propel the advancement of remote sensing for earth observation. The applications of this technology span a multitude of domains, including but not limited to national defense, land surveying, disaster monitoring, and traffic planning [7, 10–15].

Deep learning algorithms have shown significant promise in instance segmentation for remote sensing images, demonstrating their capacity to extract deep, discernible features from raw data [16–19]. Currently, the predominant instance segmentation algorithms include two-stage R-CNN series algorithms (e.g., Mask R-CNN [20], Cascade Mask R-CNN [21], Mask Scoring R-CNN [22], HTC [23], and HQ-ISNet [1]), as well as one-stage algorithms (e.g., YOLACT [24], BlendMask [25], EmbedMask [26], Condinst [27], SOLO [28], and Mask2Former [29]). Nevertheless, the intricacy of remote sensing image backgrounds and the diversity of scenes

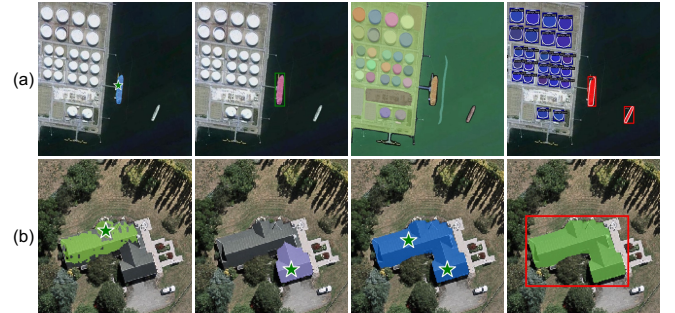


Fig. 1. Comparative instance segmentation results. (a) demonstrates the instance segmentation results derived from point-based prompts, box-based prompts, SAM’s “everything” mode (which segments all objects within the image), and RSPrompter. Note that SAM executes category-agnostic instance segmentation, dependent on manually provided prior prompts. (b) showcases the segmentation results of point-based prompts originating from various locations, a two-point-based prompt, and a box-based prompt. It is evident that the type, location, and number of prompts significantly impact the results produced by SAM.

pose challenges to the generalization and adaptability of these algorithms. Consequently, the development of instance segmentation models capable of accommodating broad remote sensing scenarios is of crucial importance for the interpretation of remote sensing imagery.

In recent years, notable advancements have been made in foundational models, such as GPT-4 [30], Flamingo [31], and SAM [32], which have significantly contributed to societal progress. Despite the inherent big data attributes of remote sensing since its inception [11, 33], there has yet to be the development of foundational models specifically tailored to this field. The primary aim of this paper is not to create a universal foundational model for remote sensing but rather to investigate the applicability of the SAM foundational model, originating from the computer vision domain, to instance segmentation in remote sensing imagery. We foresee that such foundational models will catalyze continued advancement and growth within the remote sensing field.

The SAM model, trained on over a billion masks, demonstrates exceptional generalization capabilities, enabling it to segment any object in any image without additional training [32]. This advancement opens new horizons for intelligent image analysis and understanding [34–39]. Nevertheless, SAM’s interactive framework necessitates a prior prompt, such as a point, box, or mask, to be provided in conjunction with the

input image. It operates as a category-agnostic segmentation method, as illustrated in Fig. 1 (a). These constraints render SAM unsuitable for the fully automatic understanding of remote-sensing images. Furthermore, we note that the complex background interference and the absence of well-defined object edges in remote sensing image scenarios pose significant challenges to SAM’s segmentation capabilities. SAM’s performance in segmenting remote-sensing image objects is heavily dependent on the type, location, and quantity of prompts. In most cases, refined manual prompts are indispensable to achieve the desired outcomes, as shown in Fig. 1 (b). This observation indicates that SAM suggests considerable limitations when directly applied to instance segmentation in remote-sensing images.

To enhance the remote sensing image instance segmentation of the foundation model, we propose a novel approach called RSPrompter, for learning how to generate prompts that can enhance the capabilities of the SAM framework. Our research is primarily motivated by the SAM, wherein each group of prompts can yield an instance-specific mask via the mask decoder. We postulate that if we could autonomously generate groups of category-related prompts, the SAM decoder would be capable of generating multiple instance-level masks, each with their respective category labels. However, this process presents two significant challenges: (i) How can we source category-related prompts? (ii) What type of prompts should be selected for the mask decoder?

Our research is primarily centered on the SAM framework, a category-agnostic interactive segmentation model. We propose a lightweight feature enhancer to collect features from the SAM encoder’s intermediate layers for the subsequent prompter. The prompter can generate prompts with semantic categories. Furthermore, the prompts in the SAM model can take the form of points (foreground/background points), boxes, or masks. The generation of point coordinates necessitates a search within the original SAM prompt’s manifold, which significantly constrains the prompter’s optimization space. To circumvent this limitation, we propose a more flexible representation of prompts, *i.e.*, prompt embeddings, rather than the original coordinates. Those can be interpreted as the embeddings of points or boxes. This strategy also mitigates the challenge of gradient flow from high-dimensional to low-dimensional features and vice versa, *i.e.*, from high-dimension image features to point coordinates and subsequently to positional encodings.

In addition to our primary research (RSPrompter in Fig. 2 (d)), we have conducted an extensive review of the current advancements and derivatives in the SAM community [34, 36, 40]. These primarily include methods based on the SAM backbone, methods that integrate SAM with classifiers, and methods that combine SAM with detectors. The primary contributions of this paper can be summarized as follows:

(i) We propose a novel prompt learning method that augments the SAM model’s capabilities, thereby facilitating instance segmentation in remote sensing imagery.

(ii) We undertake a comprehensive evaluation of the SAM model’s performance when integrated with other discrimina-

tive models for instance segmentation pertaining to remote sensing imagery.

(iii) We present extensive results on three diverse remote sensing instance segmentation datasets, varying in size, categories, and modalities, to demonstrate the efficacy of the proposed RSPrompter.

The structure of this paper is organized as follows: Sec. II offers a thorough review of the relevant literature. Sec. III provides an in-depth exploration of the extension methods based on the SAM model for instance segmentation, as well as a detailed discussion of the proposed RSPrompter. Sec. IV presents both quantitative and qualitative results, supplemented by ablation studies to further substantiate our findings. Finally, Sec. V concludes the paper and encapsulates the key insights.

## II. RELATED WORKS

### A. Deep Learning based Instance Segmentation

The goal of instance segmentation is to pinpoint the location of each target within an image and provide a corresponding semantic mask. This task is inherently more complex than object detection and semantic segmentation [18, 41]. Current deep learning-based instance segmentation approaches can be broadly divided into two categories: two-stage and single-stage methods. The former primarily builds upon the Mask R-CNN [20] series, which has evolved from the two-stage Faster R-CNN [42] object detector by incorporating a parallel mask prediction branch. As research progresses, a growing number of researchers are refining this framework to achieve enhanced performance. PANet [16] optimizes the information path between features by introducing a bottom-up path based on FPN [43]. In HTC [23], a multi-task, multi-stage hybrid cascade structure is proposed, and the spatial context is augmented by integrating the segmentation branch, resulting in significant performance improvements over Mask R-CNN and Cascade Mask R-CNN [21]. The Mask Scoring R-CNN [22] incorporates a mask IoU branch within the Mask R-CNN framework to evaluate segmentation quality. The HQ-ISNet [1] introduces an instance segmentation method for remote sensing imagery based on Cascade Mask R-CNN, which fully exploits multi-level feature maps and preserves the detailed information contained within high-resolution images.

While two-stage methods can produce refined segmentation results, achieving the desired speed of segmentation remains a challenge. With the increasing popularity of single-stage object detectors, numerous researchers have sought to adapt these detectors for instance segmentation tasks. For instance, YOLACT [24] addresses the instance segmentation task by generating a set of prototype masks and predicting mask coefficients for each instance. CondInst [27] provides a novel perspective on the instance segmentation problem by utilizing a dynamic masking head, outperforming existing methods such as Mask R-CNN in terms of instance segmentation performance. SOLO [28] reframes the instance segmentation problem as predicting semantic categories and generating instance masks for each pixel in the feature map. With the widespread use of Transformers [44], DETR [45] has emerged as a fully end-to-end object detector. Inspired by the task modeling and

training procedures employed in DETR, MaskFormer [46] treats segmentation tasks as mask classification problems, but with a slow convergence speed. Mask2Former [29] introduces masked attention to limit cross-attention to the foreground region, significantly enhancing network training speed.

Instance segmentation and object detection tasks are mutually reinforcing, and their development has reached a plateau at the technical level. Currently, research on foundational segmentation and detection models has become a popular area of focus [32, 47–51]. In this paper, we validate the performance of the SAM foundation model when applied to instance segmentation tasks in remote sensing imagery.

### B. Foundation Model

In recent years, foundation models have sparked a transformative shift in the field of artificial intelligence. Training on vast quantities of data has endowed these models with impressive zero-shot generalization capabilities across a multitude of scenarios [32, 52–54]. Prominent models such as Chat-GPT [55], GPT-4 [30], and Stable Diffusion [56] have further propelled the evolution of artificial intelligence, contributing significantly to the advancement of human civilization and exerting considerable influence across various industries. Motivated by the success of foundational models in Natural Language Processing (NLP), researchers have begun to investigate their potential applicability within the domain of computer vision. While the majority of these models aim to extract accessible knowledge from freely available data [31, 52, 57], the recent SAM model [32] takes an innovative approach by constructing a data engine in which the model is co-developed with model-in-the-loop dataset annotation. SAM uniquely utilizes an extensive collection of masks (11 million images, comprising over 1 billion masks), demonstrating robust generalization capabilities. However, it was originally designed as a task-agnostic segmentation model requiring prompts (*i.e.*, input of prior points, bounding boxes, or masks), and thus, it does not facilitate end-to-end automated segmentation perception. In this paper, we do not focus on the design and training of a foundational remote sensing instance segmentation model. Instead, we investigate the potential of leveraging SAM’s powerful general segmentation capabilities for remote-sensing images, with the aim of inspiring readers and fellow researchers. Moreover, the proposed learn-to-prompt method can be extended to other foundation visual models beyond SAM for task-specific or domain-specific downstream tasks.

### C. Prompt Learning

Historically, machine learning tasks have predominantly centered on fully supervised learning, where task-specific models were trained exclusively on labeled instances of the target task [58, 59]. Over time, the learning paradigm has significantly evolved, transitioning from fully supervised learning to a “pre-training and fine-tuning” approach for downstream tasks. This shift allows models to leverage the general features acquired during pre-training [60–63]. In the recent era of foundation models, a new paradigm, “pre-training and prompting” [57, 64–68], has emerged. In this paradigm, researchers

no longer exclusively adapt the model to downstream tasks. Instead, they redesign the input using prompts to reconstruct downstream tasks to align with the original pre-training task [52, 69, 70]. Prompt learning can aid in reducing semantic discrepancies, bridging the gap between pre-training and fine-tuning, and preventing overfitting of the head. Since the introduction of GPT-3 [71], prompt learning has evolved from traditional discrete [67] and continuous prompt construction [57, 68] to large-scale model-oriented in-context learning [31], instruction-tuning [72–74], and chain-of-thought approaches [75–77]. Current methods for constructing prompts primarily involve manual templates, heuristic-based templates, generation, word embedding fine-tuning, and pseudo tokens [67, 78]. In this paper, we propose a prompt generator that generates SAM-compatible prompt inputs. This prompt generator is category-related and yields semantic instance segmentation results.

## III. METHODOLOGY

In this section, we will introduce our proposed RSPrompter, a learning-to-prompt approach grounded in the SAM framework, specifically designed for instance segmentation in remote sensing imagery. The section will encompass the following aspects: a revisit of the SAM framework, the introduction of both anchor-based and query-based RSPrompter, including their loss functions, and some straightforward extensions adapting SAM to instance segmentation.

### A. A Revisit of SAM

The SAM framework, an interactive segmentation approach predicated on provided prompts such as foreground/background points, bounding boxes, and masks, comprises three primary components: an image encoder ( $\Phi_{i-enc}$ ), a prompt encoder ( $\Phi_{p-enc}$ ), and a mask decoder ( $\Phi_{m-dec}$ ). SAM employs a pre-trained Masked AutoEncoder (MAE) [79] based on the Vision Transformer (ViT) [80] to process images into intermediate features, and encodes the prior prompts as embedding tokens. The mask decoder’s cross-attention mechanism then enables interaction between image features and prompt embeddings, culminating in a mask output. This process is depicted in Fig. 3 and can be expressed as follows:

$$\begin{aligned}
 F_{img} &= \Phi_{i-enc}(\mathcal{I}) \\
 F_{sparse} &= \Phi_{p-enc}(p_{sparse}) \\
 F_{dense} &= \Phi_{p-enc}(p_{dense}) \\
 F_{out} &= \text{Cat}(T_{mc-filter}, T_{IoU}, F_{sparse}) \\
 \mathcal{M} &= \Phi_{m-dec}(F_{img} + F_{dense}, F_{out})
 \end{aligned} \tag{1}$$

where  $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$  denotes the original image,  $F_{img} \in \mathbb{R}^{h \times w \times c}$  represents the intermediate image features,  $p_{sparse}$  encompasses the sparse prompts including foreground/background points and bounding boxes, and  $F_{sparse} \in \mathbb{R}^{k \times c}$  signifies the sparse prompt tokens encoded by  $\Phi_{p-enc}$ . Furthermore,  $p_{dense} \in \mathbb{R}^{H \times W}$  refers to the coarse segmentation mask, and  $F_{dense} \in \mathbb{R}^{h \times w \times c}$  is the dense representation extracted by the prompt encoder  $\Phi_{p-enc}$ , which is an optional

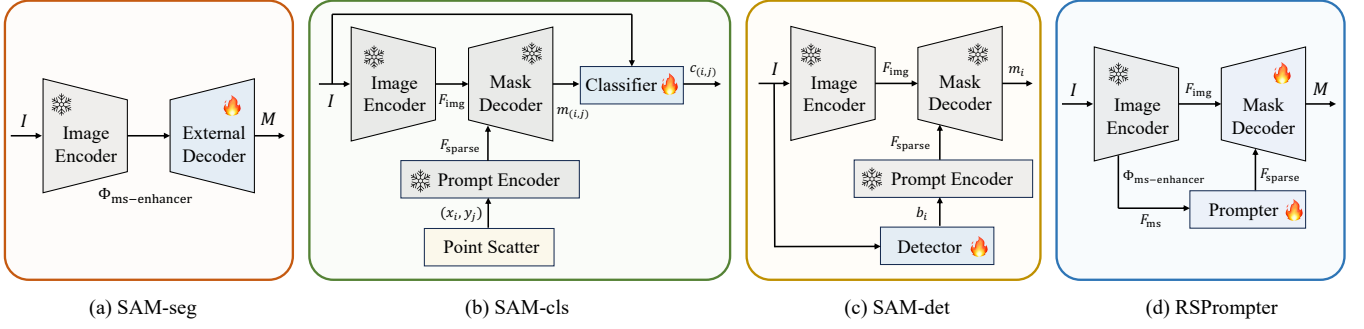


Fig. 2. From left to right, the figure illustrates SAM-seg, SAM-cls, SAM-det, and RSPrompter as alternative solutions for applying SAM to remote sensing image instance segmentation tasks.

(a) An instance segmentation head is added after SAM’s image encoder. (b) SAM’s “everything” mode generates masks for all objects in an image, which are subsequently classified into specific categories by a classifier. (c) Object bounding boxes are first produced by an object detector and then used as prior prompts for SAM to obtain the corresponding masks. (d) The proposed RSPrompter creates category-relevant prompt embeddings for instant segmentation. The snowflake icon in the diagram indicates that the model parameters are maintained in a frozen state, while the fire symbol signifies active training.

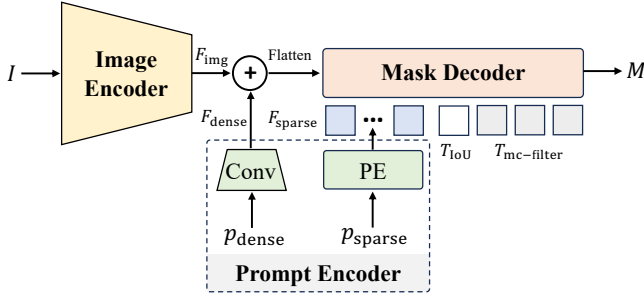


Fig. 3. The figure presents a schematic illustration of the SAM framework, encompassing an image encoder, a prompt encoder, and a mask decoder. The framework generates corresponding object masks with the given input prompts.

input for SAM.  $T_{mc-filter} \in \mathbb{R}^{4 \times c}$  and  $T_{IoU} \in \mathbb{R}^{1 \times c}$  are the pre-inserted learnable tokens representing four different mask filters and their corresponding IoU predictions.  $\mathcal{M}$  denotes the predicted multi-choice masks. In our study, diverse outputs are not required, so we directly select the first mask as the final prediction.

## B. RSPrompter

1) *Overview*: The structure of the proposed RSPrompter is depicted in Fig. 2 (d). Let us consider a training dataset, denoted as  $\mathcal{D}_{train} = \{(\mathcal{I}_1, y_1), \dots, (\mathcal{I}_N, y_N)\}$ , wherein  $\mathcal{I}_i \in \mathbb{R}^{H \times W \times 3}$  signifies an image, and  $y_i = \{b_i, c_i, m_i\}$  corresponds to its respective ground-truth annotations, encompassing the coordinates of  $n$  object bounding boxes ( $b_i \in \mathbb{R}^{n_i \times 4}$ ), their affiliated semantic categories ( $c_i \in \mathbb{R}^{n_i \times C}$ ), and binary masks ( $m_i \in \mathbb{R}^{n_i \times H \times W}$ ). The primary objective is to train a prompter for SAM that is capable of processing any image from a test set ( $\mathcal{I}_k \sim \mathcal{D}_{test}$ ), concurrently localizing the objects and inferring their semantic categories and instance masks,

which can be articulated as follows:

$$\begin{aligned}
 F_{img}, \{F_i\} &= \Phi_{i-enc}(\mathcal{I}_k) \\
 \{F_{ms}^j\} &= \Phi_{ms-enhancer}(\{F_i\}) \\
 \{F_{sparse}^m, c_m\} &= \Phi_{prompter}(\{F_{ms}^j\}) \\
 F_{out} &= \text{Cat}(T_{mc-filter}, T_{IoU}, F_{sparse}^m) \\
 \mathcal{M}_m &= \Phi_{m-dec}(F_{img}, F_{out})
 \end{aligned} \tag{2}$$

where the image ( $\mathcal{I}_k \in \mathbb{R}^{H \times W \times 3}$ ) is processed by the frozen SAM image encoder to generate  $F_{img} \in \mathbb{R}^{h \times w \times c}$  and multiple intermediate feature maps  $\{F_i\} \in \mathbb{R}^{K \times h \times w \times c}$ .  $F_{img}$  is utilized by the SAM decoder to obtain prompt-guided masks, while  $\{F_i\}$  is progressively processed by an efficient multi-scale feature enhancer ( $\Phi_{ms-enhancer}$ ) to obtain the multi-scale feature maps ( $\{F_{ms}^j \in \mathbb{R}^{\frac{H}{2^{j+1}} \times \frac{W}{2^{j+1}} \times c}\}$ ,  $j \in \{1, 2, 3, 4, 5\}$ ) and a prompter ( $\Phi_{prompter}$ ) to acquire multiple groups of prompts ( $F_{sparse}^m \in \mathbb{R}^{K_p \times c}$ ,  $m \in \{1, \dots, N_p\}$ ) and corresponding semantic categories ( $c_m \in \mathbb{R}^C$ ,  $m \in \{1, \dots, N_p\}$ ).  $K_p$  defines the number of prompt embeddings for each mask generation.  $N_p$  is the number of prompt groups to define the number of output instance masks. Two distinct structures, namely anchor-based and query-based, have been employed for the prompt generator ( $\Phi_{prompter}$ ).

It is crucial to note that  $F_{sparse}^m$  only contains foreground target instance prompts, with the semantic category given by  $c_m$ . A single  $F_{sparse}^m$  is a combination of multiple prompts, *i.e.*, representing an instance mask with multiple point embeddings or a bounding box. For the sake of simplicity, the superscript  $k$  in  $\mathcal{I}_k$  will be omitted when describing the following proposed model.

2) *Multi-scale Feature Enhancer*: To extract semantically pertinent and discriminative features without augmenting the computational complexity of the prompter, we propose a lightweight multi-scale feature enhancer. This enhancer comprises a feature aggregator and a feature splitter. The feature aggregator is designed to learn representative semantic features from a multitude of intermediate feature layers within

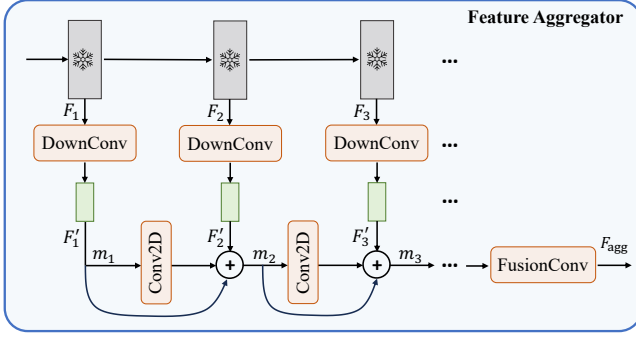


Fig. 4. Illustration of the proposed lightweight feature aggregator, which extracts semantic information from the ViT backbone and performs a lightweight fusion process.

the SAM ViT backbone. Concurrently, the feature splitter is employed to generate multi-scale pyramid feature maps.

**Feature Aggregator:** The feature aggregator, as illustrated in Fig. 4, can be defined recursively as follows:

$$\begin{aligned}
 F'_i &= \Phi_{\text{DownConv}}(F_i) \\
 m_1 &= F'_1 \\
 m_i &= m_{i-1} + \Phi_{\text{Conv2D}}(m_{i-1}) + F'_i \\
 F_{\text{agg}} &= \Phi_{\text{FusionConv}}(m_k)
 \end{aligned} \quad (3)$$

where  $F_i \in \mathbb{R}^{h \times w \times c}$  and  $F'_i \in \mathbb{R}^{h \times w \times 32}$  represent the SAM backbone’s features and down-sampled features produced by  $\Phi_{\text{DownConv}}$ , respectively. The process involves the utilization of a  $1 \times 1$  Convolution-ReLU block to decrease the channels from  $c$  to 32, followed by a  $3 \times 3$  Convolution-ReLU block to augment the spatial information. Given our hypothesis that only rudimentary information about the target location is required, we audaciously further diminish the channel dimension to curtail computational overhead.  $\Phi_{\text{Conv2D}}$  signifies a  $3 \times 3$  Convolution-ReLU block, whereas  $\Phi_{\text{FusionConv}}$  denotes the final fusion convolution layers, consisting of two  $3 \times 3$  convolution layers and a  $1 \times 1$  convolution layer to restore the channel dimension.

**Feature Splitter:** To procure multi-scale features, a straightforward feature splitter ( $\Phi_{\text{f-split}}$ ) is applied to  $F_{\text{agg}}$ . The  $\Phi_{\text{f-split}}$  employs transposed convolutional layers to generate up-sampled features, and max pooling to yield down-sampled features. By leveraging the up/down-sampling layers, we ultimately derive five feature maps of varying scales, *i.e.*,  $\{F_{\text{ms}}^j \in \mathbb{R}^{\frac{H}{2^{j+1}} \times \frac{W}{2^{j+1}} \times c}\}$ , where  $j \in \{1, 2, 3, 4, 5\}$ .

3) *Anchor-based Prompter:* Upon the acquisition of the semantic features, it becomes feasible to utilize the prompter to generate prompt embeddings for the SAM mask decoder. This section will primarily concentrate on the exploration of the anchor-based prompter.

**Architecture:** We begin with generating candidate object boxes utilizing the anchor-based Region Proposal Network (RPN). Subsequently, we extract the unique visual feature representation of each object from the positionally encoded feature map via RoI Pooling [20] with the proposal. This

visual feature is then used to derive three perception heads: the semantic head, the localization head, and the prompt head. The role of the semantic head is to identify a specific object category, whereas the localization head is responsible for establishing the matching criterion between the generated prompt representation and the target instance mask, *i.e.*, greedy matching based on localization (Intersection over Union, IoU). The prompt head, on the other hand, generates the necessary prompt embedding for the SAM mask decoder. A comprehensive illustration of the entire process is provided in Fig. 5 and can be mathematically represented by the following equation:

$$\begin{aligned}
 \{o_i\} &= \Phi_{\text{rpn}}(F_{\text{ms}}) \\
 v_i &= \Phi_{\text{roi-p}}(F_{\text{ms}} + \text{PE}, o_i) \\
 c_i &= \Phi_{\text{cls}}(v_i) \\
 b_i &= \Phi_{\text{reg}}(v_i) \\
 e_i &= \Phi_{\text{prompt}}(v_i) \\
 F_{\text{sparse}}^i &= e_i + \sin e_i
 \end{aligned} \quad (4)$$

where  $F_{\text{ms}}$  denotes the multi-scale features derived from the enhancer, while  $\Phi_{\text{rpn}}$  signifies a lightweight RPN with object proposals ( $\{o_i\}$ ). Given that the  $\Phi_{\text{roi-p}}$  operations have the potential to lose positional information relative to the entire image in prompt generation, we integrate positional encoding (PE) into  $F_{\text{ms}}$ .  $\Phi_{\text{cls}}$ ,  $\Phi_{\text{reg}}$ , and  $\Phi_{\text{prompt}}$  correspond to the semantic head, the localization head, and the prompt head, respectively. To ensure alignment between the generated prompt embeddings and the embeddings from SAM’s prompt encoder, we utilize sine function to directly generate high-frequency information, as opposed to predicting it through the network. This approach is adopted due to the inherent difficulty neural networks face in predicting high-frequency information. The efficacy of this design has been substantiated through experiments.

**Loss Function:** The primary framework of the anchor-based prompter is fundamentally aligned with the structure of Faster R-CNN [42]. The comprehensive loss encompasses several components, namely, objectness loss and localization loss within the RPN, classification loss attributed to the semantic head, regression loss associated with the localization head, and segmentation loss for the SAM decoder. Therefore, the cumulative loss can be mathematically represented as follows:

$$\mathcal{L}_{\text{anchor}} = \frac{1}{M} \sum_i \mathcal{L}_{\text{rpn}}^i + \frac{1}{N} \sum_j (\mathcal{L}_{\text{cls}}^j + \mathbb{1}^j (\mathcal{L}_{\text{reg}}^j + \mathcal{L}_{\text{seg}}^j)) \quad (5)$$

where  $\mathcal{L}_{\text{rpn}}$  represents the region proposal loss.  $\mathcal{L}_{\text{cls}}$  denotes the Cross-Entropy (CE) loss, calculated between the predicted category and the target.  $\mathcal{L}_{\text{reg}}$  is the SmoothL1 loss, computed based on the predicted coordinate offsets and the target offsets between the ground truth and the prior box. Additionally,  $\mathcal{L}_{\text{seg}}$  signifies the binary CE loss between the SAM-decoded mask and the ground-truth instance mask, wherein the IoU of the boxes determines the supervised matching criteria. The indicator function  $\mathbb{1}$  is utilized to validate a positive match.

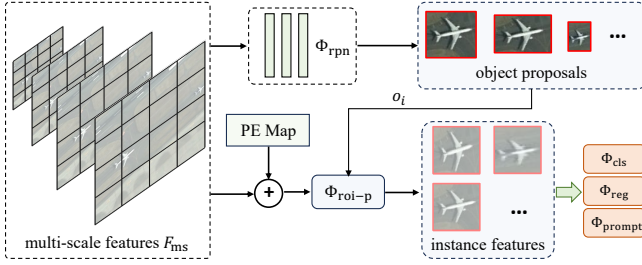


Fig. 5. A diagram of the anchor-based prompter.

4) *Query-based Prompter*: The procedure for the anchor-based prompter is relatively complex, necessitating the use of box information for mask matching and supervised training. To simplify this process, we propose a novel query-based prompter that leverages optimal transport as its foundation.

**Architecture**: The query-based prompter is primarily composed of a Transformer encoder and decoder as shown in Fig. 6. The encoder is utilized to extract high-level gathered semantic features as follows:

$$\begin{aligned} F'_{ms} &= \text{Cat}(\{F'_{ms} + \text{PE}_i + \text{LE}_i\}) \\ \hat{F}_m, \{\hat{F}_i\} &= \Phi_{\text{split}} \circ \Phi_{\text{T-enc}}(F'_{ms}) \end{aligned} \quad (6)$$

where  $\text{PE}_i$  and  $\text{LE}_i$  denote the positional encoding and level encoding of the  $i$ -th level, respectively.  $\text{Cat}(\cdot)$  signifies tensor concatenation along the channel dimension.  $\Phi_{\text{T-enc}}$  symbolizes the Transformer encoder layers, while  $\Phi_{\text{split}}$  denotes the process of partitioning the self-attended aggregated features into  $\{\hat{F}_i\}$ , maintaining their original multi-scale spatial dimensions.  $\hat{F}_m \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times c}$  is the largest-size feature map in  $\{\hat{F}_i\}$ , for coarse mask generation.

The decoder is utilized to transform the preset learnable query into the prompt embedding for SAM and corresponding semantic categories via a cross-attention interaction, as follows:

$$\begin{aligned} F_{\text{query}}^i &= \Phi_{\text{T-dec}}(\hat{F}_{i-1}, F_{\text{query}}^{i-1}) \\ \hat{c}_i &= \Phi_{\text{mlp-cls}}(F_{\text{query}}^i) \\ \hat{f}_i &= \Phi_{\text{mlp-mask}}(F_{\text{query}}^i) \\ \hat{e}_i &= \Phi_{\text{mlp-prompt}}(F_{\text{query}}^i) \\ \hat{m}_{\text{coarse}}^i &= \Phi_{\text{e-sum}}(\hat{F}_m, \hat{f}_i) \\ F_{\text{dense}}^i &= \Phi_{\text{p-enc}}(\hat{m}_{\text{coarse}}^i) \\ F_{\text{sparse}}^i &= e_i + \sin e_i \end{aligned} \quad (7)$$

where  $F_{\text{query}}^0 \in \mathbb{R}^{N_p \times c}$  denotes zero-initialized learnable tokens.  $\Phi_{\text{mlp-cls}}$  and  $\Phi_{\text{mlp-mask}}$  are linear projection layers employed to derive the object class ( $\hat{c}_i \in \mathbb{R}^{N_p \times C}$ ) and mask filter ( $\hat{f}_i \in \mathbb{R}^{N_p \times c}$ ) for the  $i$ -th level, respectively.  $\Phi_{\text{mlp-prompt}}$  is a two-layer MLP designed to acquire the projected prompt embeddings ( $\hat{e}_i \in \mathbb{R}^{N_p \times K_p \times c}$ ). Here,  $N_p$  signifies the number of prompt groups, *i.e.*, the number of instances.  $K_p$  is used to define the number of embeddings per prompt, *i.e.*, the

number of prompts required to represent an instance target. The  $i$ -th level's coarse segmentation mask is achieved by  $\Phi_{\text{e-sum}}$  through linearly weighting  $\hat{F}_m$  using  $\hat{f}_i$ .  $\Phi_{\text{p-enc}}$  is the SAM mask prompt encoder.  $F_{\text{dense}}^i$  and  $F_{\text{sparse}}^i$  can be decoded by the SAM mask decoder as Eq. 2 to yield the fine-grained mask prediction. The overall equation can be computed recurrently along the level  $i$  to obtain multi-semantic results. For inference, only the last layer is considered, and the tight bounding boxes of binary masks are procured through mathematical operations.

**Loss Function**: In the training process for the query-based prompter, two primary steps are undertaken: (i) the matching of  $N_p$  predicted masks to  $K$  ground-truth instance masks (generally,  $N_p > K$ ); (ii) the subsequent implementation of supervised training utilizing the matched labels. During the execution of optimal transport matching, we establish the following matching cost, which incorporates both the predicted category and mask:

$$\begin{aligned} \Omega &= \arg \min_{\omega} \sum_i^N \mathcal{L}_{\text{match}}(\hat{y}_i, y_{\omega(i)}) \\ \mathcal{L}_{\text{match}} &= \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{seg-ce}} + \mathcal{L}_{\text{seg-dice}} \end{aligned} \quad (8)$$

where  $\omega$  denotes the assignment relationship, while  $\hat{y}_i$  and  $y_i$  represent the prediction and the label, respectively. The Hungarian algorithm [81] is utilized to identify the optimal assignment between the  $N_p$  predictions and  $K$  targets. The matching cost takes into account the similarity between predictions and ground-truth annotations. More specifically, it includes the class classification matching cost ( $\mathcal{L}_{\text{cls}}$ ), mask cross-entropy cost ( $\mathcal{L}_{\text{seg-ce}}$ ), and mask dice cost ( $\mathcal{L}_{\text{seg-dice}}$ ).

Upon pairing each predicted instance with its corresponding ground truth, we are then able to apply the supervision terms. These primarily consist of multi-class classification and binary mask classification, as detailed below:

$$\mathcal{L}_{\text{query}} = \frac{1}{N_p} \sum_i^{N_p} (\mathcal{L}_{\text{cls}}^i + \mathbb{1}^i \mathcal{L}_{\text{seg}}^i) \quad (9)$$

where  $\mathcal{L}_{\text{cls}}$  represents the cross-entropy loss calculated between the predicted category and the target.  $\mathcal{L}_{\text{seg}}$  denotes the binary cross-entropy loss between the predicted mask and the matched ground-truth instance mask, encompassing both the predicted coarse and fine-grained masks.  $\mathbb{1}$  is an indicator to confirm a positive match.

### C. Extensions on SAM for Instance Segmentation

We conducted a comprehensive survey of the SAM community. In addition to the proposed RSPrompter, we also introduced three other SAM-based instance segmentation methods for comparative analysis. These methods are illustrated in Fig. 2 (a), (b), and (c), and we hope that they will catalyze future research. The methods include the application of an external instance segmentation head based on SAM's image encoder, the use of an object classifier predicated on SAM's "everything" mode, and the incorporation of an additional

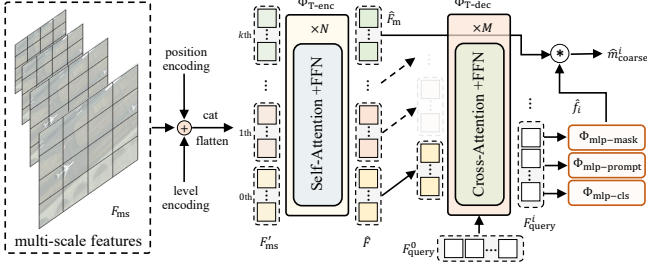


Fig. 6. An illustration of the proposed query-based prompter.

object detector. In the ensuing sections, these methods will be referred to as SAM-seg, SAM-cls, and SAM-det, respectively. It is worth noting that SAM-det enjoys widespread adoption within the SAM community, while the other two methods were conceived and developed by us.

1) *SAM-seg*: In SAM-seg, we leverage the knowledge derived from SAM’s image encoder, while keeping the cumbersome encoder frozen. We utilize the multi-scale feature enhancer, as detailed in Sec. III-B2, to obtain multi-scale features that are subsequently used for the instance segmentation heads, *i.e.*, the mask decoder in Mask R-CNN [20] and Mask2Former [29]. The procedure can be described as follows:

$$\begin{aligned} \{F_i\} &= \Phi_{i\text{-enc}}(\mathcal{I}) \\ \{F_{ms}^j\} &= \Phi_{ms\text{-enhancer}}(\{F_i\}) \\ \mathcal{M} &= \Phi_{ext\text{-dec}}(\{F_{ms}^j\}) \end{aligned} \quad (10)$$

where  $\{F_i\} \in \mathbb{R}^{k \times h \times w \times c}$ ,  $i \in \{1, \dots, k\}$  is the multi-layer semantic feature maps derived from the ViT backbone.  $\Phi_{ms\text{-enhancer}}$  refers to the multi-scale feature enhancer, as detailed in Sec. III-B2.  $\Phi_{ext\text{-dec}}$  denotes the externally inherited instance segmentation head, which could be, Mask R-CNN [20] or Mask2Former [29].

2) *SAM-cls*: In SAM-cls, we initially employ the “everything” mode of SAM to segment all potential instance targets within the image. This is internally accomplished by uniformly distributing points across the image and considering each point as a prompt input for an instance. Upon acquiring the masks for all instances, we can assign label to each mask utilizing a classifier. The procedure can be illustrated as follows:

$$\begin{aligned} F_{img} &= \Phi_{i\text{-enc}}(\mathcal{I}) \\ F_{sparse}^{(i,j)} &= \Phi_{p\text{-enc}}((x_i, y_j)) \\ F_{out}^{(i,j)} &= \text{Cat}(T_{mc\text{-filter}}, T_{IoU}, F_{sparse}^{(i,j)}) \\ m_{(i,j)} &= \Phi_{m\text{-dec}}(F_{img}, F_{out}^{(i,j)}) \\ c_{(i,j)} &= \Phi_{ext\text{-cls}}(\mathcal{I}, m_{(i,j)}) \end{aligned} \quad (11)$$

where  $(x_i, y_j)$  denotes the point prompt. For every image, we employ  $32 \times 32$  points to generate category-agnostic instance masks.  $\Phi_{ext\text{-cls}}$  denotes the external mask classifier, and  $c_{(i,j)}$  refers to the assigned category. For simplicity, we directly utilize the lightweight ResNet18 [63] to label the masks. It

performs classification by processing the original image patch cropped by the mask. During the cropping, we first enlarge the crop box two times and then blur non-mask areas to enhance the discriminative capability. Although the mask’s classification representation could potentially be extracted from the intermediate features of the image encoder for performance improvement, we have opted not to pursue this approach to maintain simplicity in our paper. Alternatively, a pre-trained CLIP model can be leveraged, enabling SAM-cls to function in a zero-shot regime without additional training.

3) *SAM-det*: The SAM-det with straightforward implementation has garnered considerable attention and application within the community. Initially, an object detector is trained to pinpoint the desired targets within the image. Subsequently, the detected bounding boxes are inputted as prompts into the SAM. The comprehensive process can be drawn as follows:

$$\begin{aligned} \{b_i, c_i\} &= \Phi_{ext\text{-det}}(\mathcal{I}) \\ F_{img} &= \Phi_{i\text{-enc}}(\mathcal{I}) \\ F_{sparse}^i &= \Phi_{p\text{-enc}}(b_i) \\ F_{out}^i &= \text{Cat}(T_{mc\text{-filter}}, T_{IoU}, F_{sparse}^i) \\ m_i &= \Phi_{m\text{-dec}}(F_{img}, F_{out}^i) \end{aligned} \quad (12)$$

where  $\{b_i, c_i\}$  denotes the bounding boxes and their corresponding semantic labels, as identified by the externally pre-trained object detector,  $\Phi_{ext\text{-det}}$ . In our study, we take the Faster R-CNN [42] as the detector.

## IV. EXPERIMENTAL RESULTS AND ANALYSES

### A. Experimental Dataset and Settings

We utilize three publicly available remote sensing instance segmentation datasets to validate the effectiveness of the proposed method, including the WHU building extraction dataset [82], the NWPU VHR-10 dataset [2, 83], and the SSDD dataset [1, 84]. These datasets, which vary in terms of size, categories, and modalities, have been extensively employed in the field of remote sensing instance segmentation [1, 2, 4, 85].

**WHU** [82]: We employ the aerial imagery subset from the WHU building extraction dataset. This subset encompasses 8188 non-overlapping RGB images, each of  $512 \times 512$  pixels. These images, captured over Christchurch, New Zealand, exhibit a spatial resolution ranging from 0.0075m to 0.3m. Following the official guidelines, we allocated 4736 images to the training set, 1036 to the validation set, and 2416 to the test set. To procure instance annotations, we utilize the connected component analysis method from the OpenCV library, enabling the transformation of semantic segmentation into an instance segmentation format.

**NWPU** [83]: The NWPU VHR-10 dataset is a remote sensing image object detection dataset comprising ten classes: airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle. It includes 715 optical remote sensing images from Google Earth, with a spatial resolution of 0.5-2m, and 85 pan-sharpened color

infrared images from the Vaihingen dataset, with a spatial resolution of 0.08m. We take 80% of the data for training and the remaining 20% for testing. The instance annotations provided by [2] are utilized for both training and evaluation.

**SSDD** [84]: The SAR Ship Detection Dataset (SSDD) comprises 1160 SAR images with a resolution range spanning from 1 to 15 meters and includes 2540 ship instances. We randomly assign 80% of the images for training and the remaining 20% for testing. The instance masks are annotated as per the guidelines provided by [1].

### B. Evaluation Protocol and Metrics

To evaluate the performance of the proposed method, we employ the widely recognized COCO [86] mean average precision (mAP) metric. This metric is frequently utilized to objectively assess the effectiveness of object detection and instance segmentation methods [9, 42, 57]. A prediction is deemed a true positive when the predicted box or mask of an instance exhibits an intersection over union (IoU) with its corresponding ground truth exceeding a threshold  $T$ , and when its predicted category aligns. In this study, we employ  $AP_{\text{box}}$ ,  $AP_{\text{box}}^{50}$ ,  $AP_{\text{box}}^{75}$ ,  $AP_{\text{mask}}$ ,  $AP_{\text{mask}}^{50}$ , and  $AP_{\text{mask}}^{75}$  for evaluation. AP refers to metrics averaged across all 10 IoU thresholds (0.50 : 0.05 : 0.95) and all categories. A larger AP value denotes more accurate predicted instance masks and, consequently, superior instance segmentation performance.  $AP^{50}$  represents the calculation under the IoU threshold of 0.50, while  $AP^{75}$  embodies a stricter metric corresponding to the calculation under the IoU threshold of 0.75. Consequently,  $AP^{75}$  outperforms  $AP^{50}$  in the evaluation of mask accuracy, with a higher  $AP^{75}$  value indicating more accurate instance masks.

### C. Implementation Details

The proposed method focuses on learning to prompt remote sensing image instance segmentation utilizing the SAM foundation model. In our experimental procedures, we employ the ViT-Huge backbone of SAM, unless specified otherwise.

1) *Architecture Details*: The proposed method focuses on the SAM framework, which generates multiple-choice segmentation masks for a single prompt. However, our method anticipates only a single instance mask for each learned prompt. As a result, we select the first mask as the final output. For each group of prompts, we set the number of prompt embeddings to 5, *i.e.*,  $K_p = 5$ . In the feature aggregator, we opt to take features from every 2 layers, as opposed to every layer, to enhance efficiency. For the anchor-based prompt, the RPN network is derived from Faster R-CNN [42], and other hyper-parameters in the training remain consistent. For the query-based prompt, we only take the last 3 small-size feature maps from the feature splitter as the inputs to improve efficiency. We employ a 3-layer transformer encoder and a 6-layer transformer decoder, implementing multi-scale training from the outputs of the decoder at 3 different levels. We only take the last-level SAM decoded mask as the final prediction. The number of learnable tokens could be easily determined

based on the distribution of object instances in each image for efficiency, *i.e.*,  $N_p = 100, 70, 30$  for the WHU, NWPU, and SSDD datasets, respectively. It is worth noting that a larger number also yields satisfactory results.

2) *Training Details*: In the training phase, we adhere to an image size of  $1024 \times 1024$ , consistent with the original input of the SAM model. To augment the training samples, we employ techniques such as horizontal flipping and large-scale jittering. The image encoder remains frozen throughout the training phase. During the testing, we predicted up to 100 instance masks per image for evaluation. All experiments are conducted using NVIDIA A100 Tensor Core GPUs. For the optimization process, we utilize the AdamW optimizer with an initial learning rate of  $1e-4$  to train our model. We establish a mini-batch size of 16. The total number of training epochs is set at 200/300 for the WHU dataset and 300/500 for both the NWPU and SSDD datasets (RSPrompter-anchor/RSPrompter-query). We implement a Cosine Annealing scheduler [90] with a linear warm-up strategy to decay the learning rate. Our proposed method is developed using the PyTorch framework, and all the additional modules are trained from scratch. It is important to note that, to enhance training efficiency, we have developed and incorporated the Automatic Mixed Precision (AMP) [91] in RSPrompter-anchor training and DeepSpeed ZeRO [92] stage 2 with FP16 in RSPrompter-query.

### D. Comparison with the State-of-the-Art

We compare our proposed method with a range of state-of-the-art instance segmentation methods, encompassing multi-stage approaches such as Mask R-CNN [20], Mask Scoring (MS) R-CNN [22], HTC [23], SCNet [88], CATNet [5], and HQ-ISNet [1]. We also take into account single-stage methods like SOLOv2 [87], CondInst [27], BoxInst [89], and Mask2Former [29]. Within this spectrum, SOLOv2, CondInst, BoxInst, and Mask2Former are categorized as filter-based methodologies, while CATNet and HQ-ISNet are classified as Mask R-CNN-based remote sensing instance segmentation techniques. To augment instance segmentation methods on SAM, we have also proposed SAM-seg (Mask R-CNN) and SAM-seg (Mask2Former), which incorporate Mask R-CNN and Mask2Former heads and training regimes, respectively. SAM-cls is envisaged as a minimalist instance segmentation method that leverages the “everything” mode of SAM to capture all instances within the image and employs an ImageNet-initialized ResNet18 [63] to label all instance masks. SAM-det signifies the initial training of a Faster R-CNN [42] detector to procure boxes, succeeded by the generation of corresponding instance masks by SAM using the box prompts, which have been widely embraced by the community. RSPrompter-anchor and RSPrompter-query denote the anchor-based and query-based promoters, respectively. All the aforementioned methods are implemented following their official publications, utilizing PyTorch as the platform.

1) *Quantitative Results on the WHU Dataset*: The comparative results of RSPrompter with other methods on the WHU dataset are delineated in Table I, with the superior



TABLE I  
COMPARATIVE ANALYSIS OF THE PROPOSED METHODS AND STATE-OF-THE-ART METHODS ON THE WHU DATASET, DEMONSTRATING AP(%) FOR BOXES AND MASKS AT VARIOUS IOU THRESHOLDS.

Method	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
Mask R-CNN [20]	66.4	86.3	76.3	65.6	87.1	76.7
MS R-CNN [22]	67.7	87.2	77.1	66.9	87.5	77.5
HTC [23]	68.4	87.8	77.8	67.7	88.1	78.3
SOLOv2 [87]	-	-	-	65.2	86.7	75.2
SCNet [88]	68.1	87.6	77.7	66.5	87.9	81.2
CondInst [27]	66.7	87.5	76.7	66.6	87.8	78.7
BoxInst [89]	66.7	86.4	75.7	55.0	86.5	63.2
Mask2Former [29]	69.3	87.2	78.0	69.2	88.5	79.3
CATNet [5]	66.7	86.3	76.4	66.1	86.6	76.8
HQ-ISNet [1]	66.1	86.0	75.7	66.5	86.4	78.9
SAM-seg (Mask R-CNN)	70.3	89.8	81.9	70.1	89.9	81.0
SAM-seg (Mask2Former)	70.7	88.4	79.1	71.1	89.5	81.1
SAM-cls	46.8	65.7	53.5	49.3	71.2	57.6
SAM-det	69.1	90.1	79.2	61.8	89.1	71.0
RSPrompter-anchor	71.9	90.9	<b>82.8</b>	70.4	90.0	80.5
RSPrompter-query	<b>72.5</b>	<b>91.0</b>	81.7	<b>72.5</b>	<b>92.0</b>	<b>82.9</b>

TABLE II  
COMPARATIVE ANALYSIS OF THE PROPOSED METHODS AND STATE-OF-THE-ART METHODS ON THE NWPU DATASET, DEMONSTRATING AP(%) FOR BOXES AND MASKS AT VARIOUS IOU THRESHOLDS.

Method	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
Mask R-CNN [20]	62.3	88.3	75.2	59.7	89.2	65.6
MS R-CNN [22]	62.3	88.6	73.1	60.7	88.7	67.7
HTC [23]	63.9	88.9	75.4	60.9	88.6	64.4
SOLO v2 [87]	-	-	-	50.9	77.5	54.1
SCNet [88]	60.0	87.5	69.1	58.1	87.4	62.0
CondInst [27]	62.3	87.8	73.3	59.0	88.5	62.8
BoxInst [89]	64.8	89.3	73.0	47.6	77.2	51.3
Mask2Former [29]	57.4	75.5	63.7	58.8	83.1	63.5
CATNet [5]	63.2	89.0	73.8	60.4	89.6	65.5
HQ-ISNet [1]	63.5	89.9	75.0	60.4	89.6	64.1
SAM-seg (Mask R-CNN)	68.8	92.2	80.1	65.2	92.0	71.6
SAM-seg (Mask2Former)	63.1	86.3	70.6	64.3	89.6	70.1
SAM-cls	40.2	57.1	44.5	44.0	66.0	47.6
SAM-det	64.2	89.6	74.6	51.5	74.8	54.0
RSPrompter-anchor	<b>70.3</b>	<b>93.6</b>	<b>81.0</b>	66.1	<b>92.7</b>	70.6
RSPrompter-query	68.4	90.3	74.0	<b>67.5</b>	91.7	<b>74.8</b>

performance underscored in bold. The task involves executing single-class instance segmentation of buildings in optical RGB band remote sensing imagery. RSPrompter-query exhibits the most commendable performance for both box and mask predictions, achieving AP<sub>box</sub> and AP<sub>mask</sub> values of 72.5/72.5. Notably, SAM-seg (Mask2Former) outperforms the original Mask2Former (69.3/69.2) with 70.7/71.1 on AP<sub>box</sub> and AP<sub>mask</sub>, while SAM-seg (Mask R-CNN) surpasses the original Mask R-CNN (66.4/65.6) with 70.3/70.1. Moreover, both RSPrompter-anchor and RSPrompter-query further enhance the performance to 71.9/70.4 and 72.5/72.5, respec-

tively, outshining SAM-det, which conducts detection before segmentation.

These observations suggest that the learning-to-prompt approach efficaciously adapts SAM for instance segmentation tasks in optical remote sensing imagery. Furthermore, they demonstrate that the SAM backbone, trained on an extensive dataset, can offer invaluable instance segmentation guidance even when it is completely frozen (as observed in SAM-seg).

2) *Quantitative Results on the NWPU Dataset:* To further validate the effectiveness of RSPrompter, comparative experiments were conducted utilizing the NWPU dataset. Despite its

TABLE III  
COMPARATIVE ANALYSIS OF THE PROPOSED METHODS AND STATE-OF-THE-ART METHODS ON THE SSDD DATASET, ILLUSTRATING AP(%) FOR BOXES AND MASKS AT VARIOUS IOU THRESHOLDS.

Method	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
Mask R-CNN [20]	67.7	95.6	84.9	64.3	92.6	80.9
MS R-CNN [22]	67.8	95.3	85.9	64.9	93.3	80.4
HTC [23]	69.3	97.1	85.7	64.1	94.4	80.6
SOLO v2 [87]	-	-	-	58.5	86.2	74.0
SCNet [88]	66.9	92.5	82.5	64.9	92.6	80.1
CondInst [27]	68.1	92.4	85.5	62.5	93.4	81.2
BoxInst [89]	62.8	96.2	74.7	45.2	92.3	35.3
Mask2Former [29]	62.7	90.7	75.6	64.4	93.0	82.4
CATNet [5]	67.5	96.8	80.4	63.9	93.7	80.1
HQ-ISNet [1]	66.6	95.9	80.2	63.4	95.1	78.1
SAM-seg (Mask R-CNN)	68.7	97.2	84.3	66.1	94.5	83.7
SAM-seg (Mask2Former)	63.0	94.9	75.6	66.5	95.0	83.6
SAM-cls	43.2	70.8	48.8	47.5	78.1	57.7
SAM-det	70.0	95.8	85.3	46.0	93.8	37.0
RSPrompter-anchor	<b>70.4</b>	<b>97.7</b>	<b>86.2</b>	66.8	94.7	84.0
RSPrompter-query	66.0	95.6	78.7	<b>67.3</b>	<b>95.6</b>	<b>84.3</b>

smaller size relative to the WHU dataset, the NWPU dataset encompasses a more diverse range of instance categories, comprising 10 classes of remote sensing objects. The comprehensive results of various methodologies applied to this dataset are delineated in Tab. II. The results indicate that RSPrompter-anchor surpasses other methods in terms of box prediction, achieving a score of 70.3 AP<sub>box</sub>. Meanwhile, RSPrompter-query excels in mask prediction, with a score of 67.5 AP<sub>mask</sub>. When compared with Mask R-CNN-based methods, single-stage methods exhibit a marginal decline in performance on this dataset, particularly the Transformer-based Mask2Former. This performance decrement could potentially be ascribed to the relatively diminutive size of the dataset, which may present difficulties for single-stage methods in attaining adequate generalization across the entire data domain. This is especially true for Transformer-based methods, which require a substantial volume of training data. Nonetheless, it is worth noting that the performance of SAM-based SAM-seg (Mask2Former) and RSPrompter-query remains commendable. The performance metrics improved from 57.4/58.8 for Mask2Former to 63.1/64.3 for SAM-seg (Mask2Former), and further to 68.4/67.5 for RSPrompter-query.

These observations suggest that SAM, when trained on a substantial volume of data, can demonstrate significant generalization capabilities even on a smaller dataset. Despite variations in the image domain, the performance of SAM can be enhanced through the learning-to-prompt approach.

3) *Quantitative Results on the SSDD Dataset:* The SSDD dataset was utilized to conduct a comprehensive evaluation of RSPrompter’s proficiency in executing remote sensing image instance segmentation. The SSDD dataset, a single-category SAR ship instance segmentation dataset, presents a distinctly different modality in comparison to the previously discussed datasets. It also exhibits substantial variations in training data

from SAM. The AP values obtained for various methods on this dataset are presented in Tab. III. Upon analysis of the results, it is evident that the SAM-seg (Mask R-CNN) (68.7/66.1) and SAM-seg (Mask2Former) (63.0/66.5), which are predicated on the SAM backbone, yield marginal improvements over the original Mask R-CNN (67.7/64.3) and Mask2Former (62.7/64.4). This indicates a discrepancy between the image domain used for SAM training and the SAR data domain. By liberating the constrained space, both the RSPrompter-anchor and RSPrompter-query further improved the performance, thereby corroborating the effectiveness of the RSPrompter as well.

4) *Qualitative Visual Comparisons:* To facilitate a more comprehensive visual comparison with other methodologies, we conducted a qualitative analysis of the segmentation results derived from SAM-based methods and other state-of-the-art instance segmentation approaches. Fig. 7, 8, and 9 illustrate sample segmentation instances from the WHU dataset, NWPU dataset, and SSDD dataset, respectively. The proposed RSPrompter yields significant visual enhancements in instance segmentation. When compared with alternative methods, RSPrompter produces superior results, characterized by sharper edges, more defined contours, improved completeness, and a closer approximation to the ground-truth references.

#### E. Ablation Study

In this section, we undertake a series of experiments on the NWPU dataset to investigate the significance of each component and parameter setting within our proposed methodology. Unless stated otherwise, all models are trained to utilize the same configuration with the ViT-Huge image encoder. We confine our ablation experiments to RSPrompter-query, given its straightforward design, which led us to consider it as the principal method for this study.

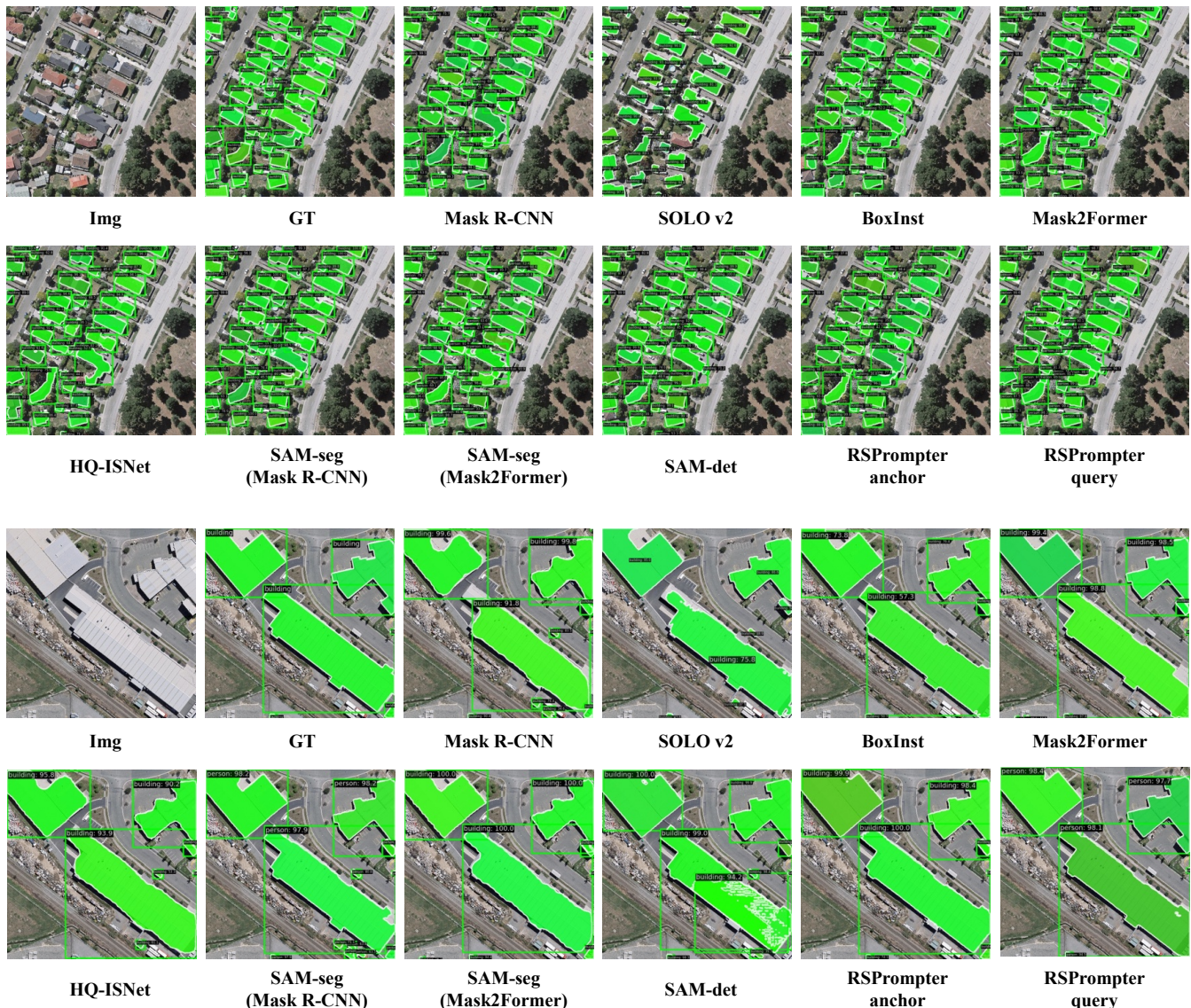


Fig. 7. Comparative visualization of segmentation results for image samples from the WHU dataset.

### 1) Impacts of Various Backbones in the Image Encoder:

Different image encoders not only affect the model’s inference speed but also significantly impact its performance. SAM employs the MAE pre-trained ViT as its image encoder, which is available in three versions: base, large, and huge. We conducted experiments based on RSPrompter-query using backbones of varying parameter sizes, and Tab. IV presents the performance of different backbone versions on the NWPU dataset. As illustrated in the table, the instance segmentation performance exhibits an increasing trend as the size of the backbone network escalates, ranging from 64.4/65.7 to 68.4/67.5 in  $AP_{\text{box}}$  and  $AP_{\text{mask}}$  metrics. Depending on the specific application context, different model sizes can be selected to achieve the optimal balance between segmentation efficiency and effectiveness.

2) Impacts of Varied Multi-scale Semantic Features in the Aggregator: The input for the feature aggregator is sourced from features at disparate levels within the ViT backbone. To

TABLE IV  
SEGMENTATION PERFORMANCE ON THE NWPU DATASET WITH VARIOUS ENCODERS AND THEIR PARAMETER COUNTS.

Encoder	Params.	$AP_{\text{box}}$	$AP_{\text{box}}^{50}$	$AP_{\text{box}}^{75}$	$AP_{\text{mask}}$	$AP_{\text{mask}}^{50}$	$AP_{\text{mask}}^{75}$
ViT-B	86M	64.4	85.3	71.2	65.7	90.1	71.1
ViT-L	307M	66.8	88.8	72.9	66.0	91.1	70.5
ViT-H	632M	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>

maintain optimal efficiency, we refrain from utilizing every layer of the ViT-H as the input for the feature aggregator. To demonstrate the impact of diverse feature layer selections on the ultimate segmentation results, we executed experiments delineated in Tab. V. The notation [start:end:step] signifies the index of the returned feature map, extending from the start to the end with the specified step size. The table demonstrates that the optimal performance and efficiency in prompt generation are achieved by features from every alternate layer.

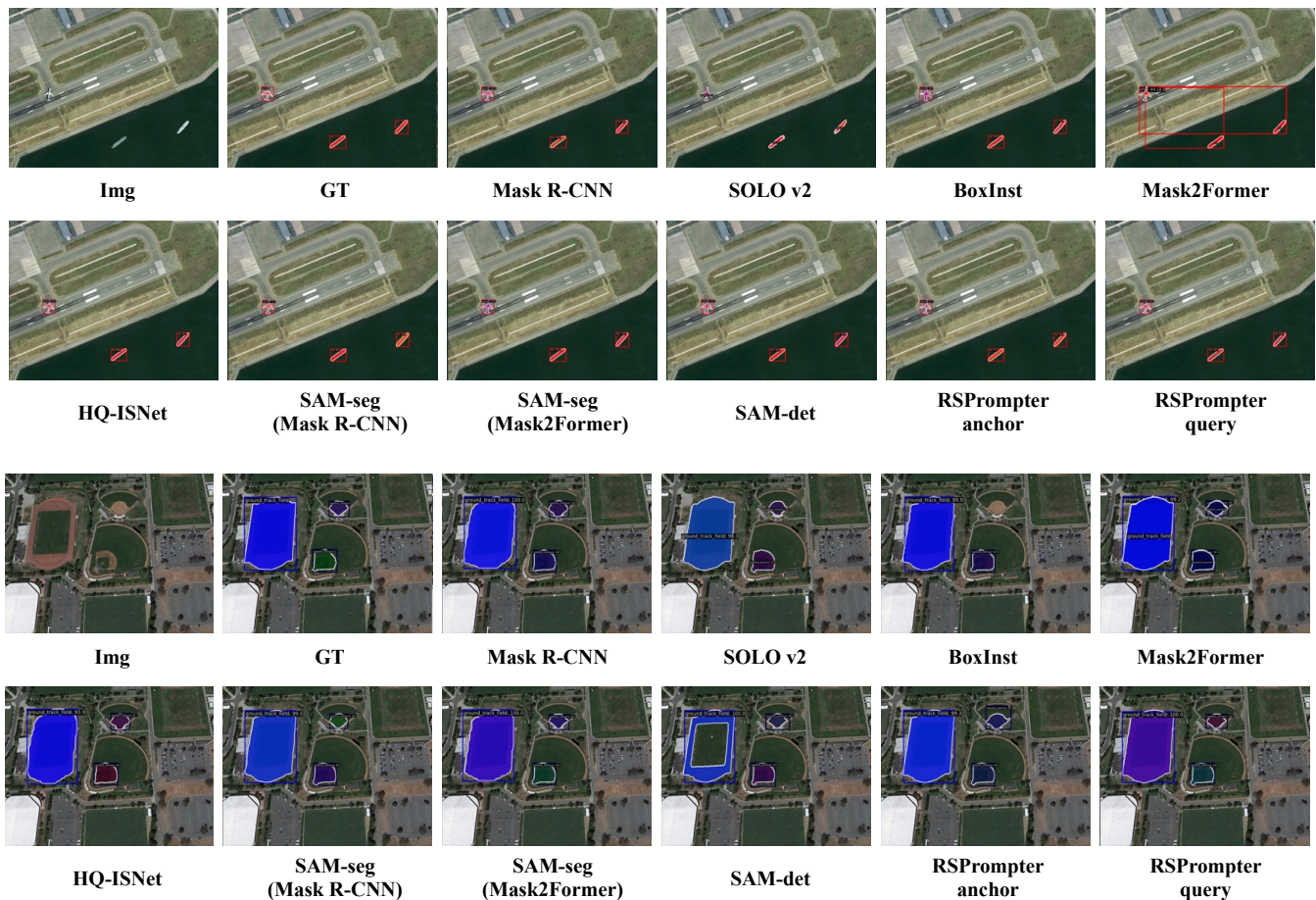


Fig. 8. Comparative visualization of segmentation results for image samples from the NWPU dataset.

TABLE V

SEGMENTATION PERFORMANCE WITH VARIOUS HIERARCHICAL FEATURES TO THE AGGREGATOR. [START:END:STEP] DENOTES THE INDEX OF THE KEPT FEATURE MAPS, RANGING FROM THE START TO THE END WITH A SPECIFIED STEP INTERVAL.

Feat.	N Layer	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
[0 : 32 : 1]	32	67.9	89.7	73.7	67.4	91.5	73.9
[0 : 32 : 2]	16	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>
[0 : 32 : 4]	8	66.3	88.5	72.4	66.3	90.6	72.0
[0 : 16 : 1]	16	65.7	88.0	72.1	65.1	90.3	71.1
[16 : 32 : 1]	16	63.8	84.2	69.8	64.7	88.8	68.7
[31 : 32 : 1]	1	62.2	83.9	68.9	64.1	88.9	67.2

3) *Impacts of DownConv Dimension Reduction in the Aggregator*: In the pursuit of developing a lightweight feature aggregation network, we have strategically minimized the channel dimensions of the features extracted from the image encoder. Specifically, the original dimensions of 768, 1024, and 1280, as associated with ViT-B, ViT-L, and ViT-H respectively, have been reduced to a uniform dimension of 32. As evidenced in Tab. VI, this substantial reduction in channel dimensions does not precipitate a significant decrease in performance. This suggests that it is sufficient to provide a rudimentary prompt to yield relatively precise mask results within SAM.

TABLE VI

IMPACT OF FEATURE DIMENSION REDUCTION IN THE AGGREGATORS ON THE BOX AND MASK PREDICTION.

Dimension	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
16	66.0	88.3	72.2	66.6	90.9	72.0
32	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>
64	66.4	89.2	72.6	67.1	91.1	73.0
128	64.4	85.2	70.1	65.9	89.9	71.5

4) *Impacts of Varied Architectural Designs in the Aggregator*: We have incorporated residual connections between layers to augment the propagation of semantic information within the feature aggregator, as depicted in Fig. 4. As corroborated by the data presented in the Tab. VII, the lack of residual connections considerably affects the final segmentation efficacy, thereby implying that a wholly serial structure may not be suitable for the aggregation of features from disparate layers of ViT. Additionally, we have put forth a parallel feature aggregation structure that concatenates all the extracted feature layers for prompt generation. Nevertheless, this approach did not yield the desired performance.

5) *Impacts of Query Numbers and Prompt Embedding Numbers in the Prompter*: The prompter is designed to generate  $N_p$  sets of prompts for each image, with each

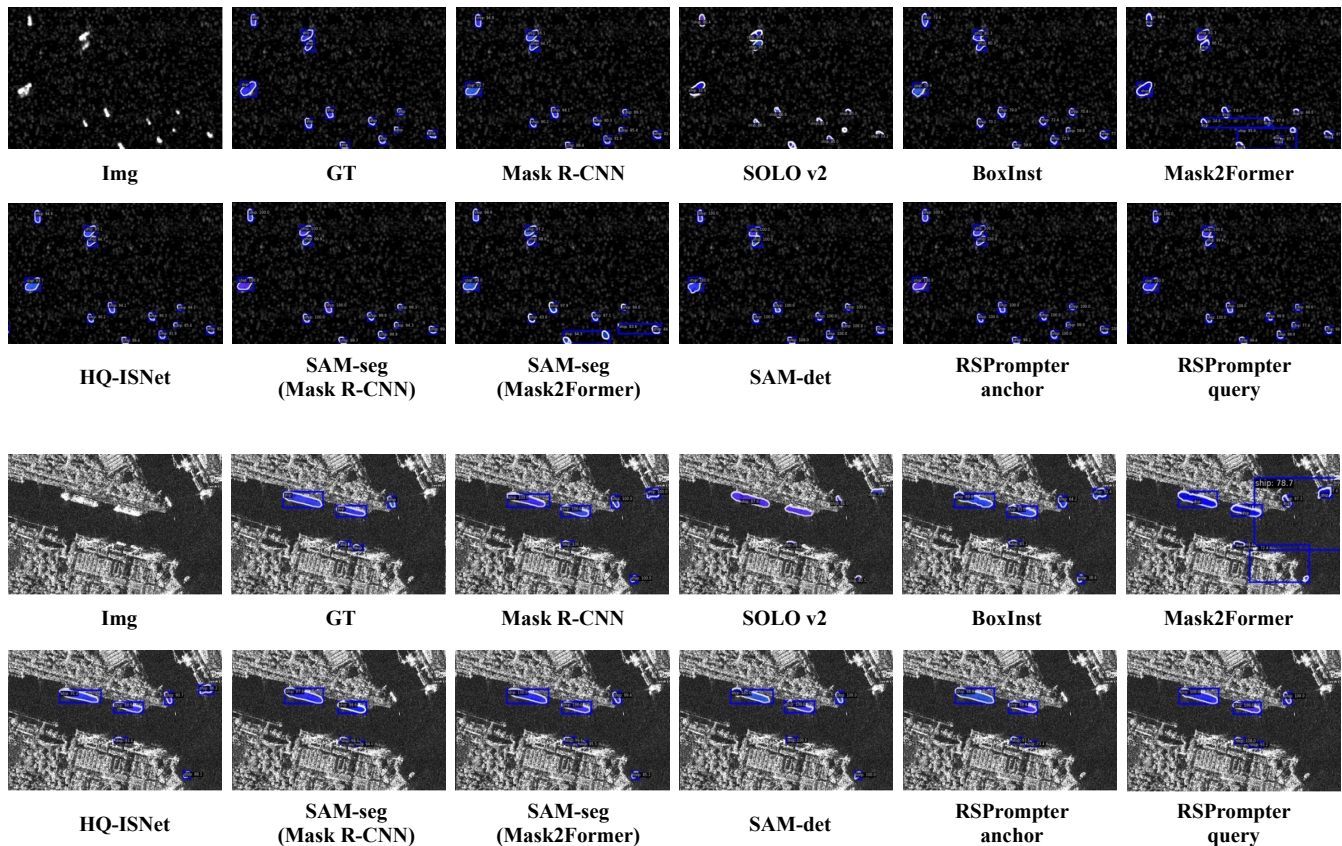


Fig. 9. Comparative visualization of segmentation results for image samples from the SSDD dataset.

TABLE VII

SEGMENTATION PERFORMANCE WITH DIFFERENT ARCHITECTURES IN THE AGGREGATOR.

Architectures design	$AP_{\text{box}}$	$AP_{\text{box}}^{50}$	$AP_{\text{box}}^{75}$	$AP_{\text{mask}}$	$AP_{\text{mask}}^{50}$	$AP_{\text{mask}}^{75}$
Ours	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>
w/o res. connection	62.5	83.5	69.3	63.8	87.3	67.9
w/ para. concatenation	63.7	85.0	70.1	64.9	89.0	69.7

set symbolizing an instance mask. Each set comprises  $K_p$  prompts, where  $K_p$  can be interpreted as the number of points representing a target instance. These two parameters have the potential to influence the ultimate segmentation outcomes, thereby necessitating controlled experiments, as demonstrated in Tab. VIII. We evaluated  $N_p$  values of 50, 70, and 100 while maintaining  $K_p$  at 1, 3, 5, 7, and 9. The findings suggest that optimal performance is achieved when  $N_p = 70$  and  $K_p = 5$ . To further comprehend the design of these parameters, we examined the distribution of instance numbers within the dataset and discovered that it mirrors the value of  $N_p$ . Consequently, we propose that the selection of  $N_p$  should consider the number of targets per image present in the dataset. The choice of  $K_p$  should neither be excessively small nor large: a small value may not adequately represent complex instances, while a large value may deviate from the distribution of the original number of SAM prompts. Similar trends were observed in the other two datasets as well.

TABLE VIII

SEGMENTATION PERFORMANCE WITH VARYING QUERY NUMBERS AND PROMPT EMBEDDING NUMBERS.

$N_p$	$K_p$	$AP_{\text{box}}$	$AP_{\text{box}}^{50}$	$AP_{\text{box}}^{75}$	$AP_{\text{mask}}$	$AP_{\text{mask}}^{50}$	$AP_{\text{mask}}^{75}$
50	5	67.7	88.9	73.6	66.8	91.4	71.7
70	5	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>
100	5	66.2	87.7	72.4	66.6	90.1	73.8
70	1	64.8	86.3	70.3	65.1	89.4	69.8
70	3	66.7	88.6	71.3	66.5	91.5	71.4
70	5	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>
70	7	66.6	88.8	72.2	66.9	89.9	73.3
70	9	65.2	87.7	72.5	65.8	88.6	72.5

6) *Impacts of Applying Sine Regularization for Prompt Embedding Generation in the Prompter*: The original SAM prompt encoder translates a coordinate-based prompt into high-frequency embeddings, which subsequently govern the decoding of masks through Fourier encoding. However, the feature generated by the prompter is inherently smooth due to the intrinsic characteristics of neural networks. To reconcile the prompt embeddings from both sources, we utilize a sine function to directly map the output of the prompter into the high-frequency domain. The efficacy of this design is substantiated in the second row of Tab. IX. The empirical results reveal that in the absence of the sine transformation, the performance metrics decline from 68.4/67.5 to 63.7/64.6 in  $AP_{\text{box}}$  and  $AP_{\text{mask}}$  respectively.

TABLE IX

IMPACTS OF THE SINE REGULARIZATION IN THE PROMPTER, TRAINING WITH SAM’S MASK DECODER FROZEN, AND EMPLOYING A MULTI-SCALE TRAINING REGIME.

Ablation	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
Ours (none)	<b>68.4</b>	<b>90.3</b>	<b>74.0</b>	<b>67.5</b>	<b>91.7</b>	<b>74.8</b>
w/o sine	63.7	85.0	69.8	64.6	88.1	70.2
w/o train dec.	56.8	77.5	62.0	59.5	83.0	63.8
w/o ms loss	59.9	80.2	66.6	62.3	85.8	67.7

TABLE X

IMPACTS OF VARYING KNOWLEDGE BASES ON INSTANCE SEGMENTATION PERFORMANCE. ALL THE BACKBONES ARE FROZEN EXCEPT IMAGENET PRE-TRAINED R50.

Base	Backbone	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
ImageNet [63]	R50*	57.4	75.5	63.7	58.8	83.1	63.5
MAE [79]	ViT-B	49.3	71.3	53.3	51.6	76.7	54.1
SeCo [95]	R50	55.5	74.4	61.0	56.7	79.8	60.1
SatMAE [93]	ViT-L	57.2	78.3	63.2	59.6	84.0	63.6
Scale-MAE [94]	ViT-L	59.6	79.3	66.4	60.1	83.7	63.7
CACo [96]	R50	57.8	79.1	64.1	59.0	84.2	62.7
SAM (Ours)	ViT-H	<b>63.1</b>	<b>86.3</b>	<b>70.6</b>	<b>64.3</b>	<b>89.6</b>	<b>70.1</b>

7) *Impacts of Freezing the Mask Decoder*: The mask decoder in SAM is characterized by a lightweight design, and we have incorporated it into the training process. To investigate the potential enhancements, we conducted an ablation study to exclude it from training, with the results presented in the third row of Tab. IX. The result suggests a significant decline in segmentation performance, thereby indicating that fine-tuning the SAM decoder for a downstream task is recommended.

8) *Impacts of Multi-scale Supervision Provided by the Splitter*: Considering that ViT generates features of uniform scale, we have designed a simplified feature splitter to provide multi-scale feature maps for ensuing multi-scale supervision during the decoding process, as depicted in Eq. 7). However, during the forward pass, we exclusively employ the output from the terminal layer of the decoder. The effectiveness of this design is underscored in the final row of Tab. IX, where the performance metrics exhibited an improvement from 59.9/62.3 to 68.4/67.5.

9) *Impacts of Varying Knowledge Bases on SAM-seg*: We regard models pre-trained on extensive datasets as knowledge bases. In this context, we have executed experiments on SAM-seg (Mask2Former) to evaluate the influence of various pre-trained knowledge bases on remote sensing segmentation. Unless explicitly mentioned, the model architecture and training regimes align with those of SAM-seg (Mask2Former), implying that the image encoder is frozen. We have conducted experiments employing the following methods: ResNet50 pre-trained on ImageNet [63], ViT-B pre-trained with MAE [79], SatMAE [93], Scale-MAE [94], SeCo [95], and CACo [96]. As presented in Tab. X, the experimental results suggest that the latent knowledge offered by the SAM backbone is more apt for segmentation tasks.

10) *Analysis of the Low Performance on SAM-cl*: SAM-cl is entirely dependent on SAM for the instance segmentation of remote sensing images, which is subsequently followed by the training and annotation of a classifier on the predicted segmentation results. As evinced in the preceding experiments,

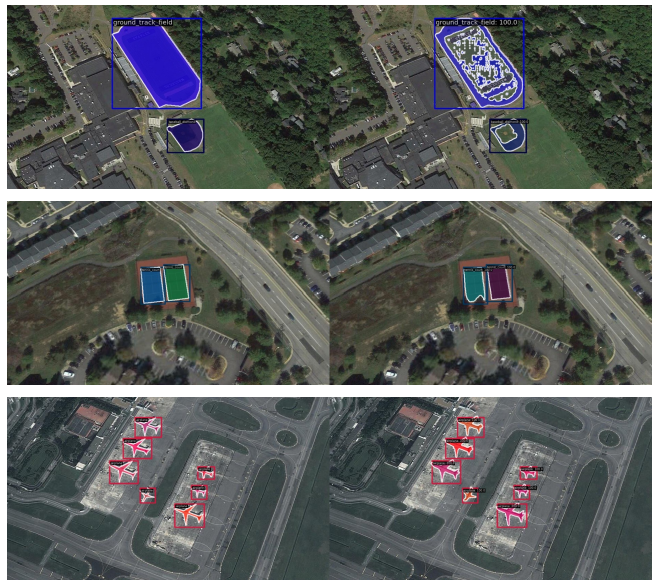


Fig. 10. Some results from SAM-det (oracle). Left: ground-truths; Right: predictions.

it manifests sub-optimal performance. To negate the influence engendered by the performance of the classifier, we exhibit the Top-1 accuracy during classification. The result is 98.5%, which underscores that the classifier is capable of accurately categorizing various classes, which further suggests that the low performance of SAM-cl can be ascribed to the pre-segmentation conducted by SAM due to the domain shift.

11) *Oracle Experimentation on SAM-det*: SAM-det utilizes a detector to procure box-based prompts, which are subsequently incorporated into SAM for mask extraction. Within this framework, we directly input annotated boxes to ascertain the maximum potential of SAM-det. It becomes apparent that despite the accuracy of the categories and the precision of the box-level annotations, SAM-det does not consistently deliver robust segmentation results, as shown in Tab. XI. This observation implies a substantial discrepancy when SAM is directly employed in the perception of remote sensing scenarios. Furthermore, the provision of box-level prompts does not invariably ensure the completeness of the mask. Fig. 10 illustrates examples of various segmentation scenarios.

TABLE XI  
A CONSIDERABLE PERFORMANCE MARGIN, EVEN LEVERAGING THE ANNOTATED BOXES FOR SAM.

Method	AP <sub>box</sub>	AP <sub>box</sub> <sup>50</sup>	AP <sub>box</sub> <sup>75</sup>	AP <sub>mask</sub>	AP <sub>mask</sub> <sup>50</sup>	AP <sub>mask</sub> <sup>75</sup>
box-predicted	64.2	89.6	74.6	51.5	74.8	54.0
box-annotated	-	-	-	52.6	84.6	55.8

## F. Discussions

In this paper, we introduce a prompt-learning methodology based on SAM, designed to enhance the processing of remote sensing images with foundation models. This method is not confined to the SAM model and can be extrapolated to a range

of foundation models. Throughout the experimental phase, we have pinpointed several potential areas for enhancement.

As highlighted by SAM, while the mask decoder of SAM is notably lightweight (4.1M), it does not necessarily denote that its computational demands are similarly minimal. The input token number of the mask decoder transformer is considerable ( $> 64 \times 64$ ), and as a prompt-based interactive segmentation head, a forward calculation is necessitated for each prompt group. Therefore, when addressing 100 instance targets within a single image, the forward calculation must be executed 100 times, which is computationally demanding. Researchers may contemplate reconfiguring this segmentation head for downstream tasks. Additionally, RSPrompter-query, predicated on optimal transport matching, converges at a slower pace than RSPrompter-anchor due to the absence of relatively explicit supervisory information. However, its network structure is more simplistic, lightweight, and exhibits superior performance compared to RSPrompter-anchor. Researchers could investigate strategies to optimize its convergence speed. Finally, the prompt learning methodology proposed herein demonstrates exceptional generalization performance on smaller datasets, significantly outperforming alternative methods. This indicates that in situations where data is insufficient to train or fine-tune an appropriate network, prompt engineering design could be contemplated for foundation models.

## V. CONCLUSION

In this paper, we present RSPrompter, a novel prompt learning methodology for instance segmentation of remote sensing images, which capitalizes on the SAM foundation model. The primary objective of RSPrompter is to learn the generation of prompt inputs for SAM, thereby enabling it to autonomously procure semantic instance-level masks. This stands in contrast to the original SAM, which necessitates additional manually-crafted prompts to attain category-agnostic masks. The design philosophy underpinning RSPrompter is not confined to the SAM model and can be extrapolated to other foundational models. Based on this philosophy, we have developed two specific implementation schemes: RSPrompter-anchor, predicated on pre-set anchors, and RSPrompter-query, which is reliant on queries and optimal transport matching. Each structure harbors its unique merits. Additionally, we have explored and proposed a variety of methods and variants within the SAM community for this task and compared them with our prompt learning approach. The efficacy of each component within the RSPrompter has been validated through comprehensive ablation studies. Simultaneously, experimental results on three public remote sensing datasets substantiate that our method surpasses other state-of-the-art instance segmentation methods, as well as several additional SAM-based methods.

## REFERENCES

- [1] H. Su, S. Wei, S. Liu, J. Liang, C. Wang, J. Shi, and X. Zhang, "Hq-isnet: High-quality instance segmentation for remote sensing imagery," *Remote Sensing*, vol. 12, no. 6, p. 989, 2020.
- [2] H. Su, S. Wei, M. Yan, C. Wang, J. Shi, and X. Zhang, "Object detection and instance segmentation in remote sensing imagery based on precise mask r-cnn," in *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 2019, pp. 1454–1457.
- [3] T. Zhang, X. Zhang, P. Zhu, X. Tang, C. Li, L. Jiao, and H. Zhou, "Semantic attention and scale complementary network for instance segmentation in remote sensing images," *IEEE Transactions on Cybernetics*, vol. 52, no. 10, pp. 10 999–11 013, 2021.
- [4] X. Xu, Z. Feng, C. Cao, M. Li, J. Wu, Z. Wu, Y. Shang, and S. Ye, "An improved swin transformer-based model for remote sensing object detection and instance segmentation," *Remote Sensing*, vol. 13, no. 23, p. 4779, 2021.
- [5] Y. Liu, H. Li, C. Hu, S. Luo, H. Shen, and C. W. Chen, "Catnet: context aggregation network for instance segmentation in remote sensing images," *arXiv preprint arXiv:2111.11057*, 2021.
- [6] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," *Proceedings of the IEEE*, 2023.
- [7] K. Chen, Z. Zou, and Z. Shi, "Building extraction from remote sensing images with sparse token transformers," *Remote Sensing*, vol. 13, no. 21, p. 4441, 2021.
- [8] K. Chen, W. Li, J. Chen, Z. Zou, and Z. Shi, "Resolution-agnostic remote sensing scene classification with implicit neural representations," *IEEE Geoscience and Remote Sensing Letters*, 2022.
- [9] J. Chen, K. Chen, H. Chen, Z. Zou, and Z. Shi, "A degraded reconstruction enhancement-based method for tiny ship detection in remote sensing images with a new large-scale dataset," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, 2022.
- [10] G. Cheng and J. Han, "A survey on object detection in optical remote sensing images," *ISPRS journal of photogrammetry and remote sensing*, vol. 117, pp. 11–28, 2016.
- [11] W. Li, K. Chen, and Z. Shi, "Geographical supervision correction for remote sensing representation learning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–20, 2022.
- [12] K. Li, G. Wan, G. Cheng, L. Meng, and J. Han, "Object detection in optical remote sensing images: A survey and a new benchmark," *ISPRS journal of photogrammetry and remote sensing*, vol. 159, pp. 296–307, 2020.
- [13] H. Chen, H. Zhang, K. Chen, C. Zhou, S. Chen, Z. Zhou, and Z. Shi, "Remote sensing image change detection towards continuous bitemporal resolution differences," *arXiv preprint arXiv:2305.14722*, 2023.
- [14] K. Chen, W. Li, S. Lei, J. Chen, X. Jiang, Z. Zou, and Z. Shi, "Continuous remote sensing image super-resolution based on context interaction in implicit function space," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [15] J. Chen, K. Chen, H. Chen, W. Li, Z. Zou, and Z. Shi, "Contrastive learning for fine-grained ship classification in remote sensing images," *IEEE Transactions on Geo-*

- science and Remote Sensing*, vol. 60, pp. 1–16, 2022.
- [16] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, “Path aggregation network for instance segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8759–8768.
- [17] R. Roscher, M. Volpi, C. Mallet, L. Drees, and J. D. Wegner, “Semcity toulouse: A benchmark for building instance segmentation in satellite images,” *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 5, pp. 109–116, 2020.
- [18] S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos, “Image segmentation using deep learning: A survey,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 7, pp. 3523–3542, 2021.
- [19] F. Fan, X. Zeng, S. Wei, H. Zhang, D. Tang, J. Shi, and X. Zhang, “Efficient instance segmentation paradigm for interpreting sar and optical images,” *Remote Sensing*, vol. 14, no. 3, p. 531, 2022.
- [20] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [21] Z. Cai and N. Vasconcelos, “Cascade r-cnn: high quality object detection and instance segmentation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 5, pp. 1483–1498, 2019.
- [22] Z. Huang, L. Huang, Y. Gong, C. Huang, and X. Wang, “Mask scoring r-cnn,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 6409–6418.
- [23] K. Chen, J. Pang, J. Wang, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Shi, W. Ouyang *et al.*, “Hybrid task cascade for instance segmentation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4974–4983.
- [24] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, “Yolact: Real-time instance segmentation,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9157–9166.
- [25] H. Chen, K. Sun, Z. Tian, C. Shen, Y. Huang, and Y. Yan, “Blendmask: Top-down meets bottom-up for instance segmentation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8573–8581.
- [26] H. Ying, Z. Huang, S. Liu, T. Shao, and K. Zhou, “Embedmask: Embedding coupling for one-stage instance segmentation,” *arXiv preprint arXiv:1912.01954*, 2019.
- [27] Z. Tian, C. Shen, and H. Chen, “Conditional convolutions for instance segmentation,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 2020, pp. 282–298.
- [28] X. Wang, T. Kong, C. Shen, Y. Jiang, and L. Li, “Solo: Segmenting objects by locations,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*. Springer, 2020, pp. 649–665.
- [29] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, “Masked-attention mask transformer for universal image segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1290–1299.
- [30] OpenAI, “Gpt-4 technical report,” 2023.
- [31] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, “Flamingo: a visual language model for few-shot learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 23 716–23 736, 2022.
- [32] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, “Segment anything,” *arXiv preprint arXiv:2304.02643*, 2023.
- [33] W. Li, K. Chen, H. Chen, and Z. Shi, “Geographical knowledge-driven representation learning for remote sensing images,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–16, 2021.
- [34] J. Ma and B. Wang, “Segment anything in medical images,” *arXiv preprint arXiv:2304.12306*, 2023.
- [35] J. Cen, Z. Zhou, J. Fang, W. Shen, L. Xie, X. Zhang, and Q. Tian, “Segment anything in 3d with nerfs,” *arXiv preprint arXiv:2304.12308*, 2023.
- [36] R. Zhang, Z. Jiang, Z. Guo, S. Yan, J. Pan, H. Dong, P. Gao, and H. Li, “Personalize segment anything model with one shot,” *arXiv preprint arXiv:2305.03048*, 2023.
- [37] T. Shaharabany, A. Dahan, R. Giryes, and L. Wolf, “Autosam: Adapting sam to medical images by overloading the prompt encoder,” *arXiv preprint arXiv:2306.06370*, 2023.
- [38] L. P. Osco, Q. Wu, E. L. de Lemos, W. N. Gonçalves, A. P. M. Ramos, J. Li, and J. M. Junior, “The segment anything model (sam) for remote sensing applications: From zero to one shot,” *arXiv preprint arXiv:2306.16623*, 2023.
- [39] X. Hu, X. Xu, and Y. Shi, “How to efficiently adapt large segmentation model (sam) to medical images,” *arXiv preprint arXiv:2306.13731*, 2023.
- [40] L. Ke, M. Ye, M. Danelljan, Y. Liu, Y.-W. Tai, C.-K. Tang, and F. Yu, “Segment anything in high quality,” 2023.
- [41] A. M. Hafiz and G. M. Bhat, “A survey on instance segmentation: state of the art,” *International journal of multimedia information retrieval*, vol. 9, no. 3, pp. 171–189, 2020.
- [42] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *Advances in neural information processing systems*, vol. 28, 2015.
- [43] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, “Feature pyramid networks for object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.



- [45] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 2020, pp. 213–229.
- [46] B. Cheng, A. Schwing, and A. Kirillov, “Per-pixel classification is not all you need for semantic segmentation,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 17 864–17 875, 2021.
- [47] M. Minderer, A. Gritsenko, A. Stone, M. Neumann, D. Weissenborn, A. Dosovitskiy, A. Mahendran, A. Arnab, M. Dehghani, Z. Shen *et al.*, “Simple open-vocabulary object detection with vision transformers,” *arXiv preprint arXiv:2205.06230*, 2022.
- [48] F. Liang, B. Wu, X. Dai, K. Li, Y. Zhao, H. Zhang, P. Zhang, P. Vajda, and D. Marculescu, “Open-vocabulary semantic segmentation with mask-adapted clip,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7061–7070.
- [49] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu *et al.*, “Grounding dino: Marrying dino with grounded pre-training for open-set object detection,” *arXiv preprint arXiv:2303.05499*, 2023.
- [50] X. Wang, X. Zhang, Y. Cao, W. Wang, C. Shen, and T. Huang, “Seggpt: Segmenting everything in context,” *arXiv preprint arXiv:2304.03284*, 2023.
- [51] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Gao, and Y. J. Lee, “Segment everything everywhere all at once,” *arXiv preprint arXiv:2304.06718*, 2023.
- [52] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [53] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, “Scaling up visual and vision-language representation learning with noisy text supervision,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 4904–4916.
- [54] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, “Cnn features off-the-shelf: an astounding baseline for recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2014, pp. 806–813.
- [55] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 27 730–27 744, 2022.
- [56] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [57] K. Chen, X. Jiang, Y. Hu, X. Tang, Y. Gao, J. Chen, and W. Xie, “Ovarnet: Towards open-vocabulary object attribute recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 23 518–23 527.
- [58] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [59] M. Z. Alom, T. M. Taha, C. Yakopcic, S. Westberg, P. Sidike, M. S. Nasrin, B. C. Van Esesn, A. A. S. Awwal, and V. K. Asari, “The history began from alexnet: A comprehensive survey on deep learning approaches,” *arXiv preprint arXiv:1803.01164*, 2018.
- [60] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International journal of computer vision*, vol. 115, pp. 211–252, 2015.
- [61] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [62] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [63] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [64] B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” *arXiv preprint arXiv:2104.08691*, 2021.
- [65] C. Feng, Y. Zhong, Z. Jie, X. Chu, H. Ren, X. Wei, W. Xie, and L. Ma, “Promptdet: Towards open-vocabulary detection using uncurated images,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*. Springer, 2022, pp. 701–717.
- [66] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, “Visual prompt tuning,” in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*. Springer, 2022, pp. 709–727.
- [67] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, “Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing,” *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–35, 2023.
- [68] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, “Learning to prompt for vision-language models,” *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.
- [69] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [70] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [71] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Ka-

- plan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [72] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *arXiv preprint arXiv:2304.08485*, 2023.
- [73] P. Gupta, C. Jiao, Y.-T. Yeh, S. Mehri, M. Eskenazi, and J. P. Bigham, “Improving zero and few-shot generalization in dialogue through instruction tuning,” *arXiv preprint arXiv:2205.12673*, 2022.
- [74] B. Peng, C. Li, P. He, M. Galley, and J. Gao, “Instruction tuning with gpt-4,” *arXiv preprint arXiv:2304.03277*, 2023.
- [75] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. Chi, Q. Le, and D. Zhou, “Chain of thought prompting elicits reasoning in large language models,” *arXiv preprint arXiv:2201.11903*, 2022.
- [76] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, and D. Zhou, “Self-consistency improves chain of thought reasoning in language models,” *arXiv preprint arXiv:2203.11171*, 2022.
- [77] Z. Zhang, A. Zhang, M. Li, and A. Smola, “Automatic chain of thought prompting in large language models,” *arXiv preprint arXiv:2210.03493*, 2022.
- [78] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. Dy, and T. Pfister, “Learning to prompt for continual learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 139–149.
- [79] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 000–16 009.
- [80] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [81] H. W. Kuhn, “The hungarian method for the assignment problem,” *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [82] S. Ji, S. Wei, and M. Lu, “Fully convolutional networks for multisource building extraction from an open aerial and satellite imagery data set,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 1, pp. 574–586, 2018.
- [83] G. Cheng, J. Han, P. Zhou, and L. Guo, “Multi-class geospatial object detection and geographic image classification based on collection of part detectors,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 98, pp. 119–132, 2014.
- [84] T. Zhang, X. Zhang, J. Li, X. Xu, B. Wang, X. Zhan, Y. Xu, X. Ke, T. Zeng, H. Su *et al.*, “Sar ship detection dataset (ssdd): Official release and comprehensive data analysis,” *Remote Sensing*, vol. 13, no. 18, p. 3690, 2021.
- [85] T. Wu, Y. Hu, L. Peng, and R. Chen, “Improved anchor-free instance segmentation for building extraction from high-resolution remote sensing images,” *Remote Sensing*, vol. 12, no. 18, p. 2910, 2020.
- [86] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [87] X. Wang, R. Zhang, T. Kong, L. Li, and C. Shen, “Solov2: Dynamic and fast instance segmentation,” *Advances in Neural information processing systems*, vol. 33, pp. 17 721–17 732, 2020.
- [88] T. Vu, H. Kang, and C. D. Yoo, “Scnet: Training inference sample consistency for instance segmentation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 3, 2021, pp. 2701–2709.
- [89] Z. Tian, C. Shen, X. Wang, and H. Chen, “Boxinst: High-performance instance segmentation with box annotations,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5443–5452.
- [90] I. Loshchilov and F. Hutter, “Sgdr: Stochastic gradient descent with warm restarts,” *arXiv preprint arXiv:1608.03983*, 2016.
- [91] P. Micikevicius, S. Narang, J. Alben, G. Diamos, E. Elsen, D. Garcia, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh *et al.*, “Mixed precision training,” *arXiv preprint arXiv:1710.03740*, 2017.
- [92] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, “Zero: Memory optimizations toward training trillion parameter models,” in *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 2020, pp. 1–16.
- [93] Y. Cong, S. Khanna, C. Meng, P. Liu, E. Rozi, Y. He, M. Burke, D. Lobell, and S. Ermon, “Satmae: Pre-training transformers for temporal and multi-spectral satellite imagery,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 197–211, 2022.
- [94] C. J. Reed, R. Gupta, S. Li, S. Brockman, C. Funk, B. Clipp, K. Keutzer, S. Candido, M. Uyttendaele, and T. Darrell, “Scale-mae: A scale-aware masked autoencoder for multiscale geospatial representation learning,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4088–4099.
- [95] O. Manas, A. Lacoste, X. Giró-i Nieto, D. Vazquez, and P. Rodriguez, “Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9414–9423.
- [96] U. Mall, B. Hariharan, and K. Bala, “Change-aware sampling and contrastive learning for satellite images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5261–5270.