

GRASP: Geospatial pixel Reasoning via Structured Policy learning

Chengjie Jiang¹, Yunqi Zhou², Jiafeng Yan², Jing Li^{3,*}

¹ Tsinghua University, Beijing, China – jiangcj25@tsinghua.org.cn

² Central University of Finance and Economics, Beijing, China – {zhouyunqi, yanjiafeng}@email.cufe.edu.cn

³ East China Normal University, Shanghai, China – lijing2017@cufe.edu.cn

Keywords: Vision-language model, Segmentation, Reinforcement learning.

Abstract

Geospatial pixel reasoning is a nascent remote-sensing task that aims to generate segmentation masks directly from natural-language instructions. Prevailing MLLM-based systems co-train a language model and a mask decoder with dense pixel supervision, which is expensive and often weak on out-of-domain (OOD) data. We introduce GRASP, a structured policy-learning framework. In our design, a multimodal large language model (MLLM) first emits task-relevant bounding boxes and positive points from a vision-language instruction. These outputs are then passed to a pretrained segmentation model, which consumes them as prompts to generate the final mask. Instead of supervised fine-tuning, we optimize the system purely with reinforcement learning: the model is trained solely with GRPO, guided by format rewards and accuracy rewards computed on boxes and points (no mask supervision). This leverages strong priors in foundation models, minimizes trainable parameters, and enables learning from inexpensive annotations. We additionally curate GRASP-1k, which contains reasoning-intensive queries, detailed reasoning traces, and fine-grained segmentation annotations. Evaluations on both in-domain and out-of-domain (OOD) test sets show state-of-the-art results: about 4% improvement in-domain and up to 54% on OOD benchmarks. The experiment results evidence our model’s robust generalization and demonstrate that complex geospatial segmentation behaviors can be learned via RL from weak spatial cues. Code and the dataset will be released open-source.

1. Introduction

Remote sensing technology plays a crucial role in geospatial information analysis, significantly benefiting various applications such as urban planning, environmental monitoring, and disaster relief (Ma et al., 2019, Liao et al., 2024, Li et al., 2025b). Traditional computer vision methods—such as object detection and semantic segmentation—have long supported RS practitioners in rapidly localizing and identifying objects in overhead imagery (Camps-Valls et al., 2006, Kotaridis and Lazaridou, 2021).

Compared with natural images, RS imagery is captured from aerial viewpoints and typically contains substantial background clutter, atmospheric distortions, extreme scale variation, and a low proportion of salient foreground objects. These characteristics make visual grounding more difficult and often lead to suboptimal transfer of methods developed for natural scenes, which tend to rely on prominent objects and limited scale changes. These RS-specific properties motivate RS-specific modeling rather than direct transplantation from the natural-image setting. Consequently, a number of RS-specific grounding models have been proposed to better address the unique challenges of remote sensing imagery (Zhan et al., 2023, Zhou et al., 2024, Choudhury et al., 2025).

At the same time, the earliest grounding formulations in both natural and RS domains were designed to localize all objects of a given category label within a scene. However, realistic geospatial scenarios often involve diverse

and complex objects of interest that defy such simple categorization.

Recent approaches (Liu et al., 2024b, Yuan et al., 2024, Pan et al., 2024, Chen et al., 2025b, Chen et al., 2025a) have introduced referring segmentation into remote sensing tasks, enabling models to handle descriptive labels with positional and attribute information. Consequently, models can precisely localize specific targets instead of indiscriminately segmenting all similar objects within a scene. Nevertheless, such methods still require experts to manually extract and format descriptions according to predefined schemas.

To further enhance model intelligence, it is natural to seek autonomous inference capabilities, enabling models to accurately locate objects based directly on natural language instructions. As a result, SegEarth-R1 (Li et al., 2025a) proposed the task of geospatial pixel reasoning to facilitate implicit querying and reasoning, generating segmentation masks based on natural language prompts. Inspired by models such as LISA (Lai et al., 2024, Yang et al., 2023) and PixelLM (Ren et al., 2024), SegEarth-R1 aligns visual and textual embeddings, which are jointly fed into a Large Language Model to autoregressively produce a <SEG> token. This token is then decoded by a mask decoder to obtain the final segmentation output. GeoPix (Ou et al., 2025) and GeoPixel (Shabbir et al., 2025) concurrently adopted similar frameworks, as illustrated in Fig. 1. In practice, these methods co-train the multimodal LLM (MLLM) and the segmentation decoder and require large amounts of pixel-level mask supervision.

However, this paradigm faces two notable challenges.

* Corresponding author

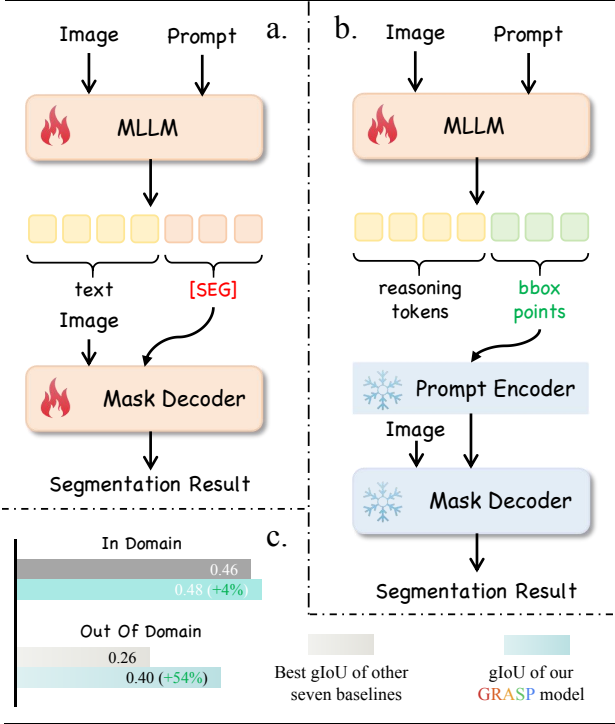


Figure 1. Comparison between previous approaches and our proposed method. (a) The architecture of prior geospatial pixel reasoning methods. (b) The architecture of our method. (c) Quantitative results.

(i) Costly supervision. High-quality remote sensing segmentation annotation demands domain experts to delineate fine-grained polygons at sub-meter resolution, carefully handling occlusions, fuzzy boundaries, and class ambiguity. This process is labor-intensive and time-consuming, making large-scale mask curation prohibitively expensive. (ii) Limited generalization under SFT. Supervised fine-tuning (SFT) inherently fits the training distribution, which can cause training instability or overfitting and degrade performance on out-of-domain (OOD) data. These limitations raise two guiding questions: (1) *Can complex segmentation behavior be learned from cheaper supervision signals—such as bounding boxes and points—to reduce annotation cost?* (2) *Can we employ a more robust training regime that improves generalization, especially on OOD scenarios?*

Motivated by evidence that reinforcement learning (RL) can markedly strengthen reasoning and generalization (Guo et al., 2025), we propose **GRASP**, a structured policy-learning framework for geospatial pixel reasoning. As shown in Fig. 1, we design a cascaded architecture. Given a vision-language instruction, the MLLM directly outputs task-relevant bounding boxes and positive points; these outputs then serve as prompts to a pre-trained segmentation model to produce the final mask. Since both pretrained MLLMs and segmentation models already provide strong reasoning and fine-grained segmentation priors, it is more effective to directly refine their decision-making process with RL rather than re-training large modules from scratch. RL excels at aligning model behavior with task-specific objectives through reward-driven optimization, allowing us to keep trainable parameters minimal while still substantially improv-

ing performance. To this end, we adopt a pure RL paradigm—specifically GRPO (Shao et al., 2024) as the sole optimization algorithm. We further introduce novel reward functions, comprising format rewards and accuracy rewards. Crucially, the accuracy rewards do not rely on pixel masks; instead, they use box and point accuracy, enabling training from inexpensive annotations while avoiding large-scale mask supervision.

In addition, we curate **GRASP-1k**, a finely annotated benchmark to rigorously evaluate generalization. Extensive experiments demonstrate that **GRASP** surpasses state-of-the-art methods on both in-domain and OOD test sets, as illustrated in Fig. 1.

- We show that complex geospatial segmentation can be effectively learned from bounding boxes and points, dramatically reducing annotation cost compared with pixel-level masks via carefully designed RL rewards.
- We introduce **GRASP**, an RL-based cascaded architecture that couples an MLLM with a pretrained segmentation model, yielding interpretable reasoning chains and precise segmentations.
- We release **GRASP-1k** for rigorous evaluation and demonstrate state-of-the-art performance. Our model achieves around 4% improvement on the in-domain test set and up to 54% improvement on out-of-domain benchmarks, highlighting the strong robustness and generalization ability of our framework.

2. Related Works

2.1 Referring Image Segmentation in Remote Sensing

In remote sensing, LGCE (Yuan et al., 2024) first adapted Transformer-based architectures to Referring Image Segmentation (RIS), specifically targeting challenges unique to remote sensing imagery, such as small and sparse objects and scale variation. Subsequent advancements like RMSIN (Liu et al., 2024b) further enhanced segmentation accuracy using rotation-aware convolutions and improved cross-scale feature interactions. DANet (Pan et al., 2024) refined multimodal interactions to tackle inter-domain discrepancies, significantly boosting model generalization capabilities. Despite these improvements, current remote sensing RIS methods still exhibit notable limitations. Primarily, these models possess limited textual comprehension capabilities due to reliance on relatively simple text encoders like BERT (Devlin et al., 2019). Furthermore, these approaches inherently lack advanced reasoning abilities, hindering their capability to directly process and understand complex natural language instructions.

2.2 Geospatial Pixel Reasoning

Motivated by RIS limitations in handling complex implicit queries involving spatial patterns, object relationships, and environmental contexts, SegEarth-R1 (Li et al., 2025a) introduced the Geospatial Pixel Reasoning

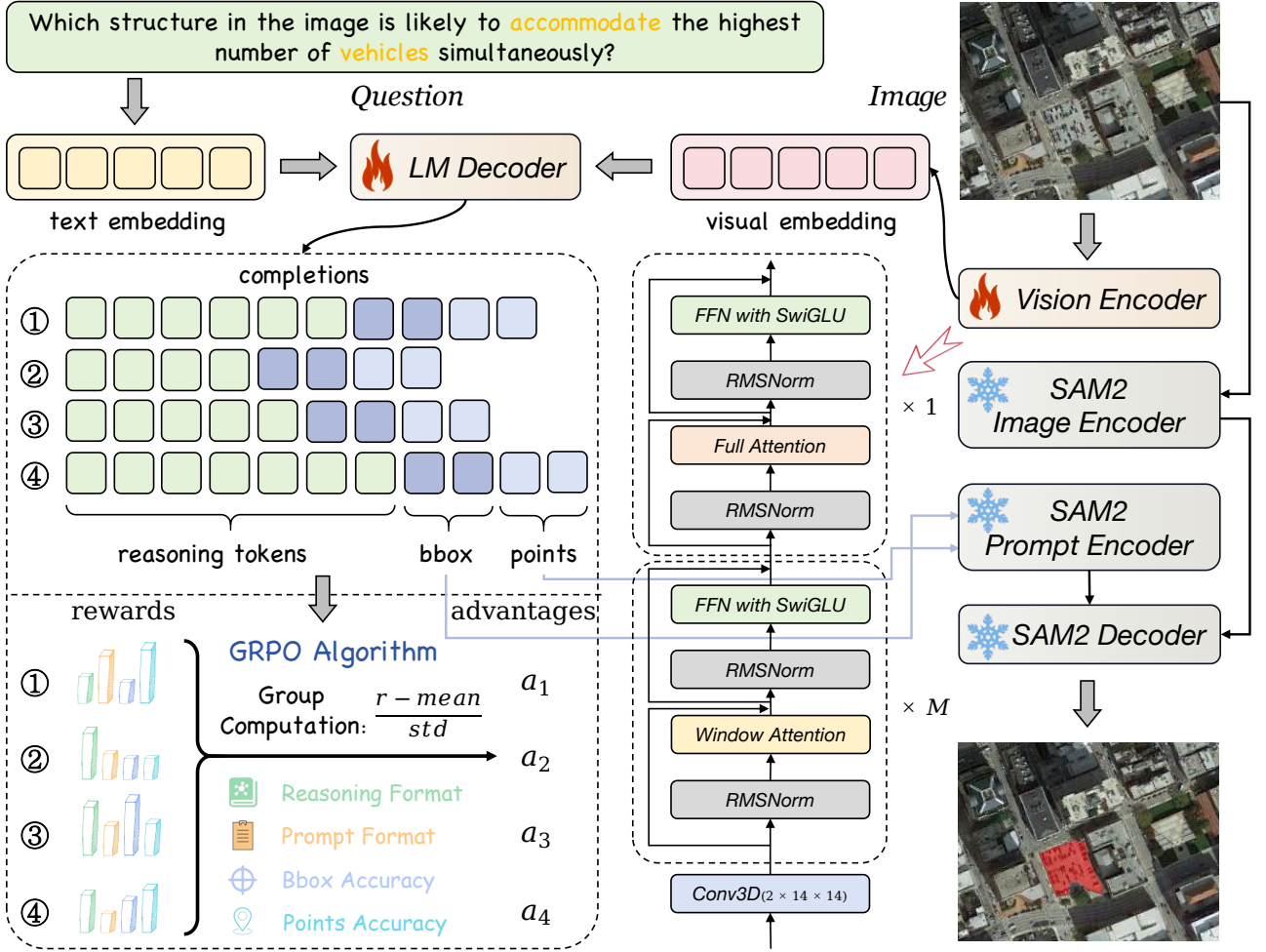


Figure 2. **Overview of the model architecture and the training workflow.** The framework consists of two main components: a multimodal large language model (MLLM) and a segmentation model. The MLLM comprises a vision encoder and an LM decoder. The LM decoder output contains both reasoning tokens and spatial grounding predictions in the form of bounding boxes and positive points. The spatial grounding predictions are fed into the SAM2 prompt encoder as prompts, while the original image is simultaneously input to the SAM2 image encoder. Finally, the SAM2 decoder combines both sources of information to produce the segmentation mask.

task. This task enables implicit querying and reasoning, generating masks directly from natural language descriptions.

Most Geospatial Pixel Reasoning methods in remote sensing adopt the framework introduced by LISA (Lai et al., 2024, Yang et al., 2023). This framework integrates specialized segmentation tokens within multimodal large language models (MLLMs). Upon generating these tokens, the framework employs additional mask decoding mechanisms, typically supported by pretrained visual backbones such as the Segment Anything Model (SAM) (Kirillov et al., 2023, Ravi et al., 2024), to produce precise spatial masks from textual inputs.

Specifically, SegEarth-R1 directly maps descriptive embeddings to segmentation masks through query-mask interactions. GeoPix (Ou et al., 2025) enhances image understanding capabilities at the pixel level, whereas GeoPixel (Shabbir et al., 2025) emphasizes hierarchical visual encoding optimized for handling high-resolution and densely populated target scenarios. Nonetheless, these Geospatial Pixel Reasoning methods currently face

significant limitations. They require finely annotated datasets for effective reasoning and segmentation training, and frequently exhibit overfitting within specific training domains, thereby limiting their generalization to diverse or out-of-domain scenarios.

Recent research in natural image reasoning segmentation, exemplified by SegZero (Liu et al., 2025), has demonstrated that reinforcement learning (RL) alone can significantly enhance the reasoning capabilities of MLLMs, even when trained solely on referring segmentation tasks. This strategy has shown superior performance on out-of-domain datasets, motivating further exploration of RL-based approaches to improve Geospatial Pixel Reasoning in remote sensing contexts.

3. Method

In this section, we present the proposed GRASP model and the corresponding reinforcement learning framework in detail. We first outline our overall approach to geospatial pixel reasoning and the associated processing pipeline in Sec. 3.1. Sec. 3.2 then provides a detailed description

of the **GRASP** model architecture together with its reinforcement learning framework. Finally, Sec. 3.3 introduces the reward functions specifically designed to guide the model training process.

3.1 Pipeline Formulation

Given a remote sensing image I and a reasoning-intensive question Q , the goal of the geospatial pixel reasoning task is to produce a segmentation mask M whose covered region precisely corresponds to the answer to Q .

Solving this task requires both strong visual reasoning capabilities and fine-grained segmentation accuracy. Inspired by (Liu et al., 2025), we decouple the reasoning and segmentation processes. Specifically, we first employ a reinforcement learning approach to train a multimodal large language model (MLLM) that, given Q , generates a reasoning process along with the bounding box B and two positive points P_1 and P_2 of the target object. These spatial grounding cues (B, P_1, P_2) are then provided as prompts to a pretrained segmentation model, which produces the final mask M .

3.2 GRASP Model

The geospatial pixel reasoning task requires a model to possess both strong visual reasoning ability and fine-grained segmentation capability. Recent advances in multimodal large language models (MLLMs) (Liu et al., 2023, Liu et al., 2024a, Chen et al., 2024, Bai et al., 2023, Bai et al., 2025, Lu et al., 2024) have demonstrated powerful vision-language reasoning skills, while pretrained segmentation models (Kirillov et al., 2023, Ravi et al., 2024) can produce high-precision segmentation masks. This motivates us to design a cascaded framework that combines the strengths of both, optimized through an appropriate learning strategy. Based on this idea, we propose **GRASP**, a model architecture that integrates an MLLM with a segmentation model, and we employ pure reinforcement learning to fully exploit the reasoning ability of the MLLM. An overview of the architecture is shown in Fig. 2.

Reasoning via Qwen. For the reasoning stage, we adopt Qwen2.5-VL (Bai et al., 2025), a powerful MLLM capable of generating bounding boxes for target objects given text instructions, but lacking native fine-grained

segmentation capability. To adapt it to the geospatial pixel reasoning domain, we set the parameters of its vision encoder to be trainable, enabling domain transfer to remote sensing imagery. We also make the LM decoder trainable, training it to output not only a bounding box B but also two positive points P_1 and P_2 to enhance localization accuracy. Formally, this process can be expressed as:

$$(B, P_1, P_2) = \text{MLLM}(I, Q; \theta_{\text{vis}}, \theta_{\text{dec}}), \quad (1)$$

where I and Q denote the input image and question, and $\theta_{\text{vis}}, \theta_{\text{dec}}$ are the trainable parameters of the vision encoder and LM decoder, respectively.

Segmentation via SAM. For the segmentation stage, we employ SAM2 (Ravi et al., 2024), a state-of-the-art open-world segmentation model. SAM2 accepts diverse prompts, including dense prompts such as masks and sparse prompts such as bounding boxes and points. Since the MLLM does not natively output dense masks, we use only the predicted bounding box and points as sparse prompts, which are fed into SAM2’s prompt encoder. The SAM2 parameters are frozen during training. We feed the image embedding from the SAM2 image encoder and the prompt embedding from the SAM2 prompt encoder jointly into the SAM2 decoder to produce the final fine-grained segmentation mask M :

$$M = \text{SAM2}(\text{Enc}_{\text{img}}(I), \text{Enc}_{\text{prompt}}(B, P_1, P_2)). \quad (2)$$

Optimization via GRPO. We train the MLLM with Grouped Relative Policy Optimization (GRPO) (Shao et al., 2024). For each question q sampled from $P(Q)$, the policy model $\pi_{\theta_{\text{old}}}$ generates a group of G candidate responses $\{o_1, o_2, \dots, o_G\}$, each represented as a token sequence. Each candidate is assigned a scalar reward $\{r_1, r_2, \dots, r_G\}$ based on the rule-based reward functions described in Sec. 3.3, covering both output format compliance and grounding accuracy. The rewards within each group are standardized to obtain relative advantages:

$$A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\})}. \quad (3)$$

The updated policy π_{θ} is then optimized by maximizing:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{q, \{o_i\}} \left[\frac{1}{G} \sum_{i=1}^G \min \left(\frac{\pi_{\theta}(o_i | q)}{\pi_{\theta_{\text{old}}}(o_i | q)} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i | q)}{\pi_{\theta_{\text{old}}}(o_i | q)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) - \beta \text{D}_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) \right] \quad (4)$$

where ε and β are hyperparameters, and π_{ref} is the frozen base model Qwen2.5-VL used for KL regularization.

3.3 Reward Design for RLVR

Reinforcement Learning with Verifiable Rewards (RLVR) requires clear and well-defined rule-based rewards. We design **five** rewards along two dimensions: (i) output format compliance, which ensures that the model output can be correctly parsed into the reasoning

chain, bounding box, and points; and (ii) localization accuracy, which evaluates the precision of the predicted bounding box and points. The total reward is the sum of the five individual rewards.

Notation. Bounding boxes are axis-aligned. Let the ground-truth (GT) box be $B_g = (x_1^g, y_1^g, x_2^g, y_2^g)$ with width $w_g = x_2^g - x_1^g$ and height $h_g = y_2^g - y_1^g$, and define $s_{\min} = \min(w_g, h_g)$, $s_{\max} = \max(w_g, h_g)$. The predicted bounding box is denoted $B_p = (x_1^p, y_1^p, x_2^p, y_2^p)$, and the corresponding center coordinates and dimensions

User Prompt for GRASP :

" Please find the target object for Question: '{Question}' and mark it with a bounding box and two positive points. "
 " Output the thinking process in <think> </think> and final answer in <answer> </answer> tags. "
 " Output a bounding box and two positive points: one at the center of the object and one outlier point on the object's edge in JSON format. "
 " i.e., <think> thinking process here </think> \n <answer> { 'bbox' : [10,100,200,210], 'points 1' : [30,110], 'points 2' : [35,180] } </answer> "

Figure 3. User prompt for GRASP. {Question} is replaced with geospatial pixel reasoning question Q in both training and inference stage.

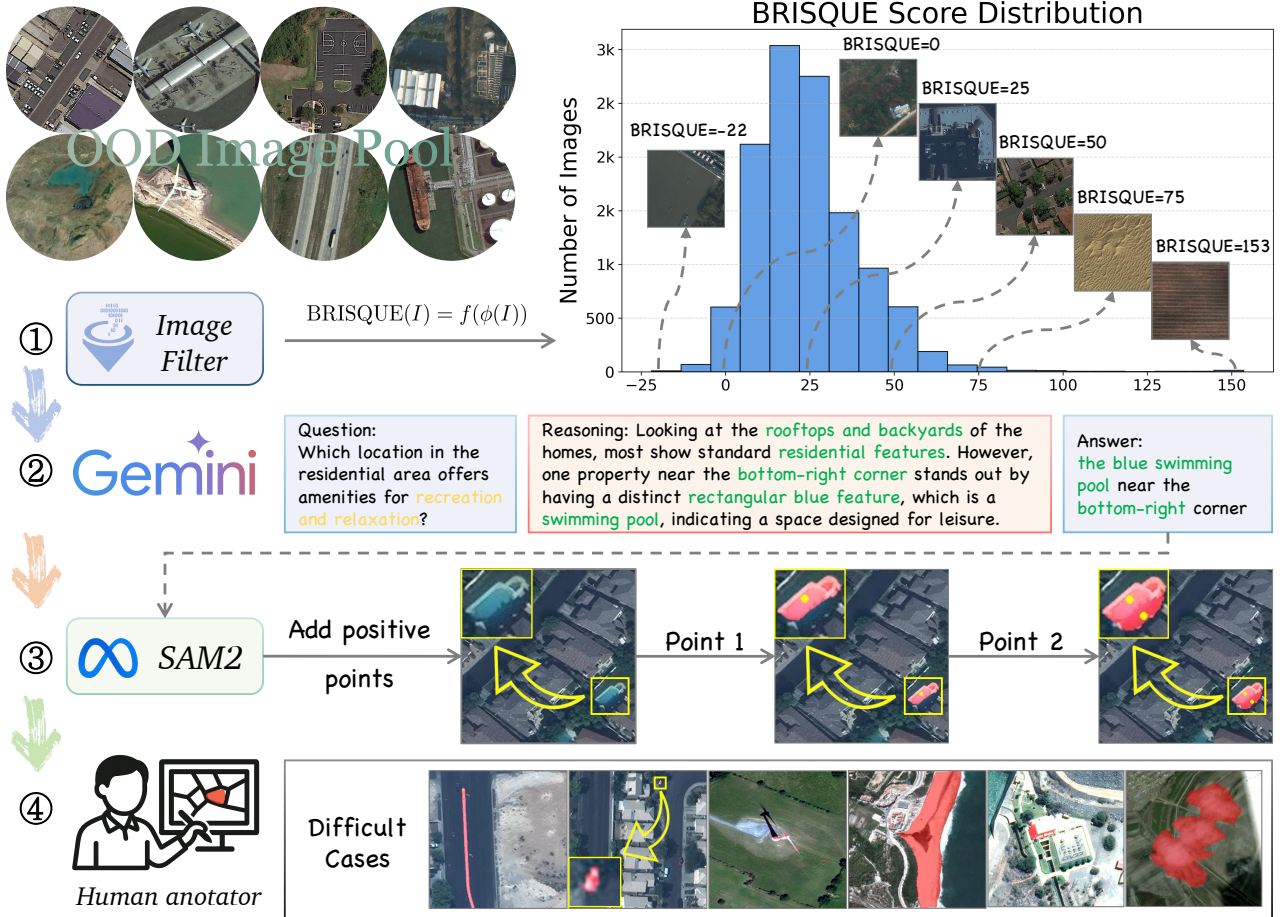


Figure 4. Overview of the GRASP-1k construction pipeline. We first curate seven out-of-domain (OOD) image pools and filter out low-quality images with BRISQUE scores above 50. For the remaining images, Gemini-2.5-Pro is used to generate reasoning-intensive questions, complete with detailed explanations and spatially grounded answers. Human annotators then click on positive points indicated by the answer and leverage SAM2 for rapid segmentation. In challenging cases where SAM2 fails, manual annotations are performed using LabelMe.

are (c_x^p, c_y^p, w^p, h^p) . Predicted points are denoted p_1^p, p_2^p and GT points are p_1^g, p_2^g .

Format Rewards

(1) Reasoning Format Reward. This discrete reward enforces the model to output a structured reasoning process. The output must include a reasoning chain enclosed within <think> and </think> tags, and a grounding answer enclosed within <answer> and </answer> tags. If the output fully meets these format requirements, the reward is 1; otherwise, it is 0.

(2) Prompt Format Reward. This reward regulates the output format of the grounding answer. As illustrated in Fig. 3, the content inside

<answer>...</answer> must strictly contain a bounding box, point 1, and point 2 following the specified schema. Only when the output exactly follows this schema will the reward be 1; otherwise, it is 0.

Localization Rewards

(3) Bbox IoU Reward. We calculate the intersection-over-union between the predicted bounding box B_p and the ground-truth bounding box B_g as

$$\text{IoU}(B_p, B_g) = \frac{\text{area}(B_p \cap B_g)}{\text{area}(B_p \cup B_g)}. \quad (5)$$

If the IoU is greater than 0.5, the reward is 1; otherwise, it is 0.

(4) **Bbox Distance Reward.** While IoU is useful for evaluating overlap, it can be insensitive to small positional errors and produces a binary signal when used with a fixed threshold. To provide a smoother optimization signal and ensure fairness across different object scales, we introduce a scale-normalized distance reward. Let $B_p = (c_x^p, c_y^p, w^p, h^p)$ and $B_g = (c_x^g, c_y^g, w^g, h^g)$ be the center coordinates and dimensions of the predicted and GT boxes. The normalized L1 distance is computed as

$$d_{\text{bbox}} = \frac{1}{2} \left(\frac{|c_x^p - c_x^g|}{w_g} + \frac{|c_y^p - c_y^g|}{h_g} \right), \quad (6)$$

where the width and height terms are omitted because their normalization is symmetric to the center offset terms. The soft reward is defined as

$$R_{\text{bbox-dist}} = \max \left(0, 1 - \frac{d_{\text{bbox}}}{0.5} \right), \quad (7)$$

which decays linearly with the normalized distance.

(5) **Points Accuracy Reward.** We first require both predicted points to lie inside the GT bounding box. The scale-normalized L1 distances for point 1 and point 2 are computed as

$$d_1 = \frac{\|p_1^p - p_1^g\|_1}{s_{\min}}, \quad d_2 = \frac{\|p_2^p - p_2^g\|_1}{s_{\max}}. \quad (8)$$

We take the mean of the two:

$$S = \frac{d_1 + d_2}{2}. \quad (9)$$

If both points are inside the GT bounding box and $S < 0.5$, the reward is 1; otherwise, it is 0.

4. Data Construction

This section introduces the construction of two types of data: (1) the transformation of existing dense segmentation labels into sparse supervision consisting of bounding box and positive points, and (2) the construction pipeline of the **GRASP-1k** benchmark, a fine-grained geospatial pixel reasoning dataset designed for rigorous evaluation under out-of-domain (OOD) conditions.

4.1 Training Data Construction

To ensure fair comparison with existing geospatial pixel reasoning methods, we do not introduce any external grounding datasets. Instead, we reconstruct the supervision format based on existing datasets with dense segmentation annotations.

As shown in Fig. 5, given a binary mask \mathcal{M} representing the dense object segmentation, we first extract the bounding box $\mathbf{b} = [x_{\min}, y_{\min}, x_{\max}, y_{\max}]$ by locating the extreme foreground pixels of \mathcal{M} . Then, to generate two positive points \mathbf{p}_1 and \mathbf{p}_2 as sparse supervision, we adopt the following strategy: - \mathbf{p}_1 is the center of the maximal inscribed circle within the mask, approximated by the point with the largest Euclidean distance to the mask boundary. - \mathbf{p}_2 is a supplementary point sampled

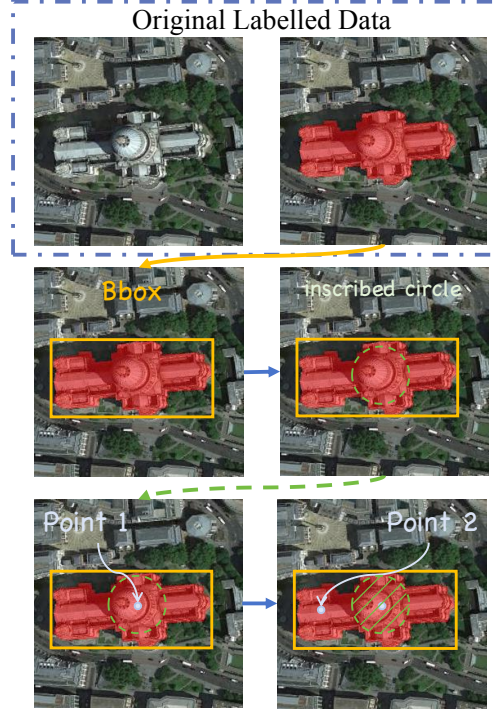


Figure 5. Reconstruction pipeline for training data. We transform dense segmentation masks into sparse supervision comprising a bounding box and two positive points.

from the outer ring of the mask beyond the inscribed circle. If no such candidate exists, we select the farthest point on the mask boundary from \mathbf{p}_1 .

4.2 GRASP-1k Construction Pipeline

To construct the **GRASP-1k** benchmark, we curate an additional set of image pools entirely from out-of-domain (OOD) sources. These images are distinct from those used during model training, which are drawn from Earth-Reason (Million-AID (Long et al., 2021) and DIOR (Li et al., 2020)) and GeoPixInstruct (HRSC2016 (Liu et al., 2017), DOTA-V2.0 (Xia et al., 2018), FAIR1M-2.0 (Sun et al., 2022)). Specifically, we select six OOD datasets and randomly sample 2,000 instances from each, as summarized in Table 1.

Data Source	Image Size	Resolution
CVUSA (Workman et al., 2015)	800	0.08m
NWPU-RESISC45 (Cheng et al., 2017)	256	0.2~30m
CrowdAI (Mohanty, 2018)	—	<0.5m
fMoW (Christie et al., 2018)	74×58~16184×16288	0.5m
CVACT (Liu and Li, 2019)	1200	0.12m
LoveDA (Wang et al., 2021)	1024	0.3m

Table 1. The six OOD image pools used in **GRASP-1k**.

To ensure data quality, we apply the BRISQUE (Blind/Referenceless Image Spatial Quality Evaluator) (Mittal et al., 2012) metric to assess image clarity, contrast, and distortion. BRISQUE computes quality scores based on deviations from natural scene statistics. Formally:

$$\text{BRISQUE}(I) = f(\phi(I)), \quad (10)$$

where $\phi(I)$ denotes a set of natural scene statistics features extracted from locally normalized luminance patches of image I , and $f(\cdot)$ is a learned regression model

that maps these features to perceptual quality scores. Lower scores indicate better perceived quality.

Following this, we discard all images with BRISQUE scores greater than or equal to 50:

$$\mathcal{I}_{\text{clean}} = \{I_i \mid \text{BRISQUE}(I_i) < 50\}. \quad (11)$$

As shown in Fig. 4, images with high scores typically lack visual clarity or contain limited semantic information, making them unsuitable for constructing reasoning-intensive prompts.

For the remaining high-quality images, we employ Gemini-2.5-Pro (Team et al., 2023) to generate reasoning-oriented prompts. Each data point consists of (1) a visually grounded question with logical reasoning requirements, (2) a directional answer referencing specific spatial locations, and (3) a detailed explanation chain. Human annotators then identify positive points based on the answer and apply SAM2 (Ravi et al., 2024) to generate corresponding segmentation masks. For ambiguous or invalid cases, samples are discarded. In rare cases where SAM2 fails to produce satisfactory results, manual segmentation is performed using LabelMe (Wada, n.d.).

As a result, we construct GRASP-1k, a high-quality benchmark of 1071 reasoning samples with fine-grained segmentation annotations based entirely on OOD imagery, offering a rigorous testbed for evaluating geospatial pixel reasoning.

5. Experiments

5.1 Experimental Settings

Implementation Details. We adopt Qwen2.5-VL-7B-Instruct (Bai et al., 2025) as the base reasoning model and SAM2-Large (Ravi et al., 2024) as the base segmentation model. During both the training and inference stages, we provided the Qwen-VL with the user prompt illustrated in Fig. 3 to guide its output. GRASP is trained using the verl (Sheng et al., 2024) reinforcement learning framework, which provides a scalable implementation of GRPO. The training is performed on 8 NVIDIA A100-40G GPUs, with hyperparameters configured as follows: the KL loss coefficient is set to 5.0×10^{-3} , the learning rate is 1.0×10^{-6} , and the micro-batch sizes per device are 8 for policy updates and 4 for experience collection.

Evaluation Datasets. Our method is trained on the EarthReason (Li et al., 2025a) and GeoPixInstruct (Ou et al., 2025) training sets, and we merge their corresponding test sets to form a single in-domain evaluation dataset. In addition, we construct a completely out-of-domain benchmark, GRASP-1k, from a newly curated image pool, which serves as our out-of-domain (OOD) test set.

Evaluation Methods. We compared our model against both natural image reasoning segmentation models (LISA-7b (Lai et al., 2024) and PixelLM-7b (Ren et al., 2024)) and remote sensing referring image segmentation models (LGCE (Yuan et al., 2024) and RMSIN (Liu

et al., 2024b)), as well as three recent models with geospatial pixel reasoning capabilities (GeoPixel (Shabbir et al., 2025), GeoPix (Ou et al., 2025), and SegEarth-R1 (Li et al., 2025a)). To ensure a fair comparison, all baseline models were fine-tuned or trained using supervised fine-tuning (SFT) on the same training dataset as ours.

Evaluation Metrics. According to previous works (Li et al., 2025a, Ou et al., 2025, Shabbir et al., 2025), we evaluate our model using four widely adopted metrics: mIoU, gIoU and cIoU. mIoU measures the mean Intersection-over-Union across all samples, reflecting overall segmentation accuracy. gIoU extends IoU by penalizing non-overlapping predictions through the inclusion of the smallest enclosing box area. cIoU further improves gIoU by incorporating the distance between the centers of predicted and ground-truth boxes as well as aspect ratio consistency.

5.2 Qualitative Results

As shown in Fig. 6, our model is capable of generating both detailed reasoning chains and fine-grained segmentation masks. Benefiting from the strong reasoning capacity of Qwen-VL and the pre-trained fine-grained segmentation capability of SAM2, we did not employ any supervised fine-tuning with segmentation or reasoning data. Instead, we solely used reinforcement learning (RL) to activate the inherent potential of this cascaded architecture and achieve strong performance in both in-domain and out-of-domain scenarios. Moreover, GRASP consistently handles diverse challenges: it remains robust against distractors involving visually similar objects (first column), effectively localizes small objects (fourth column), and successfully performs reasoning and localization tasks that are strongly tied to remote sensing knowledge (sixth column).

5.3 Comparisons on In-domain Datasets

Model	Pub.	mIoU	gIoU	cIoU
LISA-7b (Lai et al., 2024)	CVPR'24	0.28	0.30	0.30
PixelLM-7b (Ren et al., 2024)	CVPR'24	0.17	0.17	0.17
LGCE (Yuan et al., 2024)	TGRS'24	0.29	0.29	0.29
RMSIN (Liu et al., 2024b)	CVPR'24	0.50	0.45	0.42
GeoPixel (Shabbir et al., 2025)	ICML'25	0.20	0.17	0.18
GeoPix (Ou et al., 2025)	GRSM'25	0.31	0.32	0.29
SegEarth-R1 (Li et al., 2025a)	arxiv'25	<u>0.46</u>	<u>0.46</u>	<u>0.47</u>
GRASP(Ours)	—	<u>0.46</u>	0.48	0.48

Table 2. In-domain performance comparison.

As illustrated in Fig. 7, the top three rows provide a qualitative comparison between our model and seven state-of-the-art (SOTA) baselines on in-domain datasets. Most models are able to localize the relevant regions; however, our model produces the most fine-grained segmentation boundaries, which can be attributed to the strong segmentation capability of SAM2. Tab. 2 presents the corresponding quantitative comparison, where our model achieves overall SOTA performance on the in-domain test set. In particular, benefiting from the use of Bbox Distance Reward and Points Accuracy Reward during training, our model ranks first on both gIoU and cIoU, which explicitly account for the degree of prediction deviation.



Figure 6. **Qualitative results of the GRASP model.** The top three rows present cases from the in-domain test set, with the first row derived from GeoPixInstruct and the second and third rows from EarthReason. The bottom three rows show cases from GRASP-1k, an out-of-domain (OOD) test set. Our model demonstrates the ability to generate both detailed reasoning chains and fine-grained segmentation masks.

Model	Pub.	mIoU	gIoU	cIoU
LISA-7b (Lai et al., 2024)	CVPR'24	0.25	0.16	0.14
PixelLM-7b (Ren et al., 2024)	CVPR'24	0.29	0.18	0.20
LGCE (Yuan et al., 2024)	TGRS'24	0.12	0.06	0.08
RMSIN (Liu et al., 2024b)	CVPR'24	0.29	0.14	0.16
GeoPixel (Shabbir et al., 2025)	ICML'25	0.29	0.10	0.16
GeoPix (Ou et al., 2025)	GRSM'25	0.33	0.26	0.28
SegEarth-R1 (Li et al., 2025a)	arxiv'25	0.28	0.12	0.17
GRASP(Ours)	—	0.46	0.40	0.39

Table 3. Out-of-domain performance comparison.

5.4 Comparisons on GRASP-1k

As shown in the bottom three rows of Fig. 7, our model demonstrates substantially superior performance on the out-of-domain test set compared with all other baselines. In complex geospatial pixel reasoning scenarios, our model is able to accurately infer the target objects and produce fine-grained segmentation masks, while most competing models yield low-confidence predictions manifested as either large coarse regions or scattered patches. Tab. 3 further provides the quantitative com-

parison: our model significantly outperforms the others on the OOD test set, achieving improvements of 39% in mIoU, 54% in gIoU, and 39% in cIoU. These results highlight the strong robustness brought by reinforcement learning.

5.5 Ablation Study

Model	mIoU _{id}	gIoU _{id}	cIoU _{id}	mIoU _{ood}	gIoU _{ood}	cIoU _{ood}
GRASP-Zero	0.15	0.14	0.14	0.11	0.09	0.08
GRASP-SFT	0.46	0.48	0.48	0.37	0.32	0.35
GRASP-RL	0.45	0.48	0.48	0.46	0.40	0.39

Table 4. Ablation on training strategy.

SFT vs. RL We compare supervised fine-tuning (SFT) and reinforcement learning (RL) within our cascaded framework (Tab. 4). Subscripts _{id} and _{ood} denote in-domain and out-of-domain test sets, respectively. GRASP-Zero is an untrained cascade that uses Qwen2.5-VL-7B-Instruct to predict bounding boxes from the image and instruction, which are then fed to SAM2 for

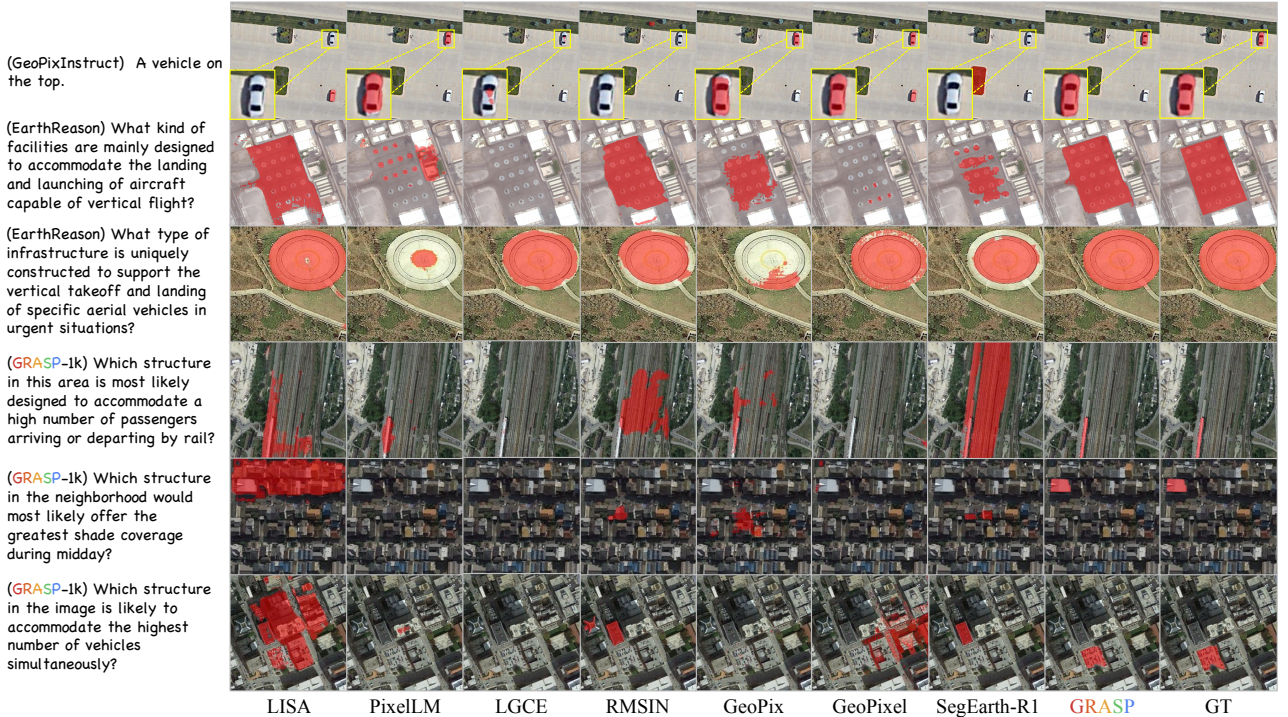


Figure 7. Qualitative comparison between our model and seven other state-of-the-art methods on both in-domain and out-of-domain test sets.

masking. It performs poorly, reflecting the substantial domain gap and the need for task-specific learning. GRASP-SFT follows prior work: Qwen-VL outputs a [SEG] token that is passed through a 1-layer MLP into SAM2’s prompt encoder, and both Qwen and SAM2 are jointly fine-tuned; this yields the best in-domain scores, consistent with SFT’s strong fit to the training distribution. GRASP-RL is our GRPO-trained model. It matches SFT in-domain performance while markedly improving OOD performance. These results indicate that RL enhances generalization without sacrificing in-domain accuracy.

Ablation of rewards We conduct a detailed ablation study on the reward design, as reported in Tab. 5. The format reward enforces both reasoning format and prompt format, the Bbox IoU reward provides strict overlap supervision, the Bbox distance reward penalizes deviations in box localization, and the point accuracy reward evaluates the correctness of predicted points. Without the format reward, the model fails to learn effectively, since properly structured outputs are a prerequisite for computing the other rewards. Comparing the second and third rows shows that box-based accuracy rewards are more effective than point-based rewards, which may be attributed to the pretrained model’s stronger prior for generating bounding boxes and to the fact that boxes provide a more informative prompt for SAM2. The last four rows further demonstrate that each reward contributes positively: all designs improve both in-domain and out-of-domain performance, and the combination of all four achieves the best results with 0.46 mIoU in both settings.

Format	Bbox IoU	Bbox Dis	Points	mIoU _{id}	mIoU _{ood}
×	✓	✓	✓	0.20	0.11
✓	×	×	✓	0.36	0.29
✓	✓	×	×	0.42	0.40
✓	✓	✓	×	0.44	0.43
✓	✓	✓	✓	0.46	0.46

Table 5. Ablation of rewards type.

6. Conclusion

In this work, we addressed the emerging task of geospatial pixel reasoning, which aims to generate fine-grained segmentation masks in remote sensing imagery directly from natural language instructions. While prior approaches rely on pixel-level mask supervision and supervised fine-tuning, they suffer from costly annotation requirements and limited generalization. To overcome these challenges, we introduced GRASP, a structured policy-learning framework that cascades an MLLM with a pretrained segmentation model and leverages reinforcement learning for optimization. By designing novel rewards based on format and grounding accuracy, our method effectively learns complex segmentation behaviors from inexpensive annotations while preserving the strong priors of foundation models.

We further released GRASP-1k, a benchmark containing reasoning-intensive instructions, detailed reasoning traces, and fine-grained segmentation labels. Extensive experiments demonstrate that GRASP achieves state-of-the-art results, improving in-domain accuracy by around 4% and OOD performance by up to 54%, thereby highlighting the robustness and generalization capability of our framework.

Looking forward, we believe that reinforcement learning with structured rewards opens a promising direction for bridging vision–language reasoning and dense prediction tasks in remote sensing, paving the way toward more autonomous and cost-efficient geospatial analysis.

References

- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F. et al., 2023. Qwen technical report. [arXiv preprint arXiv:2309.16609](#).
- Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J. et al., 2025. Qwen2. 5-vl technical report. [arXiv preprint arXiv:2502.13923](#).
- Camps-Valls, G., Gomez-Chova, L., Muñoz-Marí, J., Vila-Francés, J., Calpe-Maravilla, J., 2006. Composite kernels for hyperspectral image classification. *IEEE geoscience and remote sensing letters*, 3(1), 93–97.
- Chen, K., Liu, C., Chen, B., Zhang, J., Zou, Z., Shi, Z., 2025a. RSRefSeg 2: Decoupling Referring Remote Sensing Image Segmentation with Foundation Models. [arXiv preprint arXiv:2507.06231](#).
- Chen, K., Zhang, J., Liu, C., Zou, Z., Shi, Z., 2025b. Rsrefseg: Referring remote sensing image segmentation with foundation models. [arXiv preprint arXiv:2501.06809](#).
- Chen, Z., Wu, J., Wang, W., Su, W., Chen, G., Xing, S., Zhong, M., Zhang, Q., Zhu, X., Lu, L. et al., 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 24185–24198.
- Cheng, G., Han, J., Lu, X., 2017. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10), 1865–1883.
- Choudhury, S., Kurkure, P., Banerjee, B., 2025. Improving visual grounding in remote sensing images with adaptive modality guidance. *ISPRS Journal of Photogrammetry and Remote Sensing*, 224, 42–58.
- Christie, G., Fendley, N., Wilson, J., Mukherjee, R., 2018. Functional map of the world. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6172–6180.
- Devlin, J., Chang, M.-W., Lee, K., Toutanova, K., 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 4171–4186.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X. et al., 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. [arXiv preprint arXiv:2501.12948](#).
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A. C., Lo, W.-Y. et al., 2023. Segment anything. *Proceedings of the IEEE/CVF international conference on computer vision*, 4015–4026.
- Kotaridis, I., Lazaridou, M., 2021. Remote sensing image segmentation advances: A meta-analysis. *ISPRS Journal of Photogrammetry and Remote Sensing*, 173, 309–322.
- Lai, X., Tian, Z., Chen, Y., Li, Y., Yuan, Y., Liu, S., Jia, J., 2024. Lisa: Reasoning segmentation via large language model. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9579–9589.
- Li, K., Wan, G., Cheng, G., Meng, L., Han, J., 2020. Object detection in optical remote sensing images: A survey and a new benchmark. *ISPRS journal of photogrammetry and remote sensing*, 159, 296–307.
- Li, K., Xin, Z., Pang, L., Pang, C., Deng, Y., Yao, J., Xia, G., Meng, D., Wang, Z., Cao, X., 2025a. SegEarth-R1: Geospatial Pixel Reasoning via Large Language Model. [arXiv preprint arXiv:2504.09644](#).
- Li, L., Long, D., Wang, Y., Woolway, R. I., 2025b. Global dominance of seasonality in shaping lake-surface-extent dynamics. *Nature*, 1–8.
- Liao, Z., Yue, C., He, B., Zhao, K., Ciais, P., Alkama, R., Grassi, G., Sitch, S., Chen, R., Quan, X. et al., 2024. Growing biomass carbon stock in China driven by expansion and conservation of woody areas. *Nature Geoscience*, 17(11), 1127–1134.
- Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., Lee, Y. J., 2024a. Lllavanext: Improved reasoning, ocr, and world knowledge.
- Liu, H., Li, C., Wu, Q., Lee, Y. J., 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36, 34892–34916.
- Liu, L., Li, H., 2019. Lending orientation to neural networks for cross-view geo-localization. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5624–5633.
- Liu, S., Ma, Y., Zhang, X., Wang, H., Ji, J., Sun, X., Ji, R., 2024b. Rotated multi-scale interaction network for referring remote sensing image segmentation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26658–26668.
- Liu, Y., Peng, B., Zhong, Z., Yue, Z., Lu, F., Yu, B., Jia, J., 2025. Seg-zero: Reasoning-chain guided segmentation via cognitive reinforcement. [arXiv preprint arXiv:2503.06520](#).
- Liu, Z., Yuan, L., Weng, L., Yang, Y., 2017. A high resolution optical satellite image dataset for ship recognition and some new baselines. *International conference on pattern recognition applications and methods*, 2, SciTePress, 324–331.

- Long, Y., Xia, G.-S., Li, S., Yang, W., Yang, M. Y., Zhu, X. X., Zhang, L., Li, D., 2021. On creating benchmark dataset for aerial image interpretation: Reviews, guidances, and million-aid. IEEE Journal of selected topics in applied earth observations and remote sensing, 14, 4205–4230.
- Lu, H., Liu, W., Zhang, B., Wang, B., Dong, K., Liu, B., Sun, J., Ren, T., Li, Z., Yang, H. et al., 2024. Deepseek-vl: towards real-world vision-language understanding. arXiv preprint arXiv:2403.05525.
- Ma, L., Liu, Y., Zhang, X., Ye, Y., Yin, G., Johnson, B. A., 2019. Deep learning in remote sensing applications: A meta-analysis and review. ISPRS journal of photogrammetry and remote sensing, 152, 166–177.
- Mittal, A., Moorthy, A. K., Bovik, A. C., 2012. No-reference image quality assessment in the spatial domain. IEEE Transactions on image processing, 21(12), 4695–4708.
- Mohanty, S. P., 2018. Crowdai mapping challenge 2018 : Baseline with mask rcnn. <https://github.com/crowdai/crowdai-mapping-challenge-mask-rcnn>.
- Ou, R., Hu, Y., Zhang, F., Chen, J., Liu, Y., 2025. GeoPix: Multi-Modal Large Language Model for Pixel-level Image Understanding in Remote Sensing. arXiv preprint arXiv:2501.06828.
- Pan, Y., Sun, R., Wang, Y., Zhang, T., Zhang, Y., 2024. Rethinking the implicit optimization paradigm with dual alignments for referring remote sensing image segmentation. Proceedings of the 32nd ACM International Conference on Multimedia, 2031–2040.
- Ravi, N., Gabeur, V., Hu, Y.-T., Hu, R., Ryali, C., Ma, T., Khedr, H., Rädle, R., Rolland, C., Gustafson, L. et al., 2024. Sam 2: Segment anything in images and videos. arXiv preprint arXiv:2408.00714.
- Ren, Z., Huang, Z., Wei, Y., Zhao, Y., Fu, D., Feng, J., Jin, X., 2024. Pixellm: Pixel reasoning with large multimodal model. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 26374–26383.
- Shabbir, A., Zumri, M., Bennamoun, M., Khan, F. S., Khan, S., 2025. GeoPixel: Pixel Grounding Large Multimodal Model in Remote Sensing. arXiv preprint arXiv:2501.13925.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y. et al., 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. arXiv preprint arXiv:2402.03300.
- Sheng, G., Zhang, C., Ye, Z., Wu, X., Zhang, W., Zhang, R., Peng, Y., Lin, H., Wu, C., 2024. HybridFlow: A Flexible and Efficient RLHF Framework. arXiv preprint arXiv: 2409.19256.
- Sun, X., Wang, P., Yan, Z., Xu, F., Wang, R., Diao, W., Chen, J., Li, J., Feng, Y., Xu, T. et al., 2022. FAIR1M: A benchmark dataset for fine-grained object recognition in high-resolution remote sensing imagery. ISPRS Journal of Photogrammetry and Remote Sensing, 184, 116–130.
- Team, G., Anil, R., Borgeaud, S., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., Millican, K. et al., 2023. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805.
- Wada, K., n.d. Labelme: Image Polygonal Annotation with Python.
- Wang, J., Zheng, Z., Ma, A., Lu, X., Zhong, Y., 2021. LoveDA: A remote sensing land-cover dataset for domain adaptive semantic segmentation. arXiv preprint arXiv:2110.08733.
- Workman, S., Souvenir, R., Jacobs, N., 2015. Wide-area image geolocalization with aerial reference imagery. Proceedings of the IEEE International Conference on Computer Vision, 3961–3969.
- Xia, G.-S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo, J., Datcu, M., Pelillo, M., Zhang, L., 2018. Dota: A large-scale dataset for object detection in aerial images. The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Yang, S., Qu, T., Lai, X., Tian, Z., Peng, B., Liu, S., Jia, J., 2023. LISA++: An Improved Baseline for Reasoning Segmentation with Large Language Model. arXiv preprint arXiv:2312.17240.
- Yuan, Z., Mou, L., Hua, Y., Zhu, X. X., 2024. Rrsis: Referring remote sensing image segmentation. IEEE Transactions on Geoscience and Remote Sensing.
- Zhan, Y., Xiong, Z., Yuan, Y., 2023. Rsvg: Exploring data and models for visual grounding on remote sensing data. IEEE Transactions on Geoscience and Remote Sensing, 61, 1–13.
- Zhou, Y., Lan, M., Li, X., Ke, Y., Jiang, X., Feng, L., Zhang, W., 2024. Geoground: A unified large vision-language model. for remote sensing visual grounding. arXiv preprint arXiv:2411.11904.