# EarthMind: Towards Multi-Granular and Multi-Sensor Earth Observation with Large Multimodal Models

**Yan Shu**[1]    **Bin Ren**[1,4,5]    **Zhitong Xiong**[3]    **Danda Pani Paudel**[5]
**Luc Van Gool**[5]    **Begüm Demir**[2]    **Nicu Sebe**[1]    **Paolo Rota**[1]
[1]University of Trento    [2]Technische Universität Berlin    [3]Technical University of Munich
[4]University of Pisa    [5]INSAIT, Sofia University "St. Kliment Ohridski"
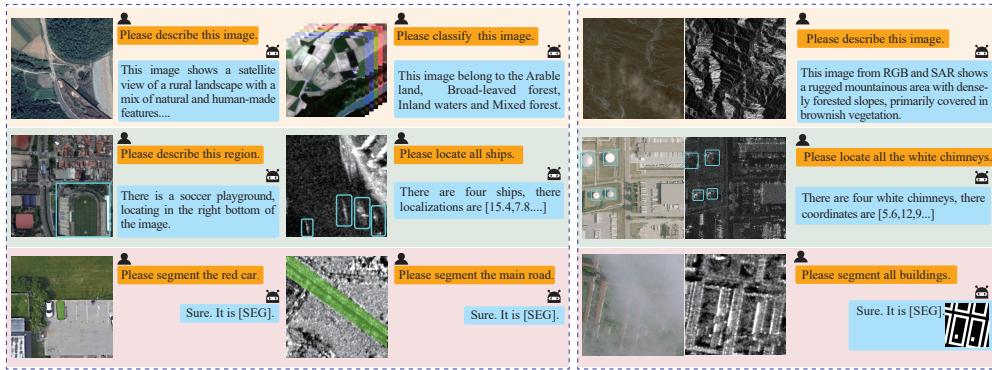https://github.com/shuyansy/EarthMind

Figure 1: The proposed EarthMind supports unified multi-granular understanding for Earth Observation (EO) imagery, including image-level, region-level, and pixel-level tasks. In addition, it enables complementary multi-sensor fusion across Optical and SAR modalities.

## Abstract

Large Multimodal Models (LMMs) have demonstrated strong performance in various vision-language tasks. However, they often struggle to comprehensively understand Earth Observation (EO) data, which is critical for monitoring the environment and the effects of human activity on it. In this work, we present **EarthMind**, a novel vision-language framework for multi-granular and multi-sensor EO data understanding. EarthMind features two core components: (1) *Spatial Attention Prompting (SAP)*, which reallocates attention within the LLM to enhance pixel-level understanding; and (2) *Cross-modal Fusion*, which aligns heterogeneous modalities into a shared space and adaptively reweighs tokens based on their information density for effective fusion. To facilitate multi-sensor fusion evaluation, we propose **EarthMind-Bench**, a comprehensive benchmark with over 2,000 human-annotated multi-sensor image-question pairs, covering a wide range of perception and reasoning tasks. Extensive experiments demonstrate the effectiveness of EarthMind. It achieves state-of-the-art performance on EarthMind-Bench, surpassing GPT-4o despite being only 4B in scale. Moreover, EarthMind outperforms existing methods on multiple public EO benchmarks, showcasing its potential to handle both multi-granular and multi-sensor challenges in a unified framework.

# 1 Introduction

Large multimodal models (LMMs), which integrate large language models (LLMs) [1, 2] with visual encoders, have shown remarkable success across a variety of vision-language tasks, including image captioning [3, 4], visual question answering [5–7], and grounding [8–10]. Among various application domains, Earth Observation (EO) [11–13] is of particular importance, as it allows monitoring the Earth and the effect of human activities on it. However, LMMs trained on general-purpose images often struggle to generalize to EO data, due to a significant domain gap. Recent work has addressed this challenge by constructing large-scale instruction tuning datasets [14–18] specifically tailored for EO, enabling better adaptation of LMMs to this domain.

Despite recent advances, existing LMMs are limited in their understanding of EO data. Firstly, EO tasks span **multiple levels of granularity**, from pixel-level segmentation [19, 20], over region-level semantic understanding [21, 22], up to image-level scene classification [23, 24]. Secondly, EO data comprise **multiple sensing modalities**, including optical imagery (e.g., RGB and Multispectral) and Synthetic Aperture Radar (SAR) [25, 26]. These modalities are inherently complementary: optical images provide rich texture and spectral information under favorable conditions, while SAR captures structural details regardless of weather or illumination. Although various sensor types exist, effective fusion, particularly between SAR and optical data modalities, remains a key challenge for EO understanding. As summarized in Tab. 1, achieving fine-grained, multi-sensor comprehension in EO remains largely unresolved.

To tackle these challenges, we introduce **EarthMind**, the first LMM capable of fusing multi-sensor EO inputs and performing reasoning across multiple semantic levels, as shown in Fig. 1. It achieves this by projecting heterogeneous features from different sensors and scales into a unified semantic space, thus enabling effective interpretation by LLMs. The novelty of EarthMind lies in two key design concepts that enable spatial grounding and cross-modal understanding in EO data. First, *Spatial Attention Prompting (SAP)* enhances pixel-level grounding by explicitly extracting and reallocating attention to regions aligned with queried objects. This overcomes limitations of prior approaches [9, 10] that combine segmentation foundation models [27, 28] with LLMs but degrade in EO settings due to vague boundaries and scale imbalances. Second, a *Cross-modal Fusion* mechanism, built upon token-level contrastive learning, guides the integration of complementary modalities (e.g., RGB and SAR) into a unified semantic space. Equipped with Modality Mutual Attention, EarthMind adaptively selects the most informative features from each modality, thereby facilitating robust autoregressive learning within the LLM.

Additionally, we propose **EarthMind-Bench**, a new benchmark designed to evaluate LMMs in challenging EO scenarios. As shown in Tab. 1, EarthMind-Bench offers several unique features. Firstly, it encompasses *multi-granular tasks*, ranging from coarse-grained image understanding to fine-grained segmentation. Secondly, it introduces *multi-sensor data*, in particular paired RGB-SAR imagery, enabling evaluation of cross-modal fusion capabilities. Thirdly, it covers *multi-level questions*, spanning low-level perception as well as high-level reasoning. EarthMind-Bench consists of more than 2,000 multiple-choice and open-ended questions, providing a comprehensive benchmark to assess the ability of LMMs to interpret and reason over EO data.

We implemented EarthMind based on Qwen-2.5-3B [29], and its effectiveness is demonstrated from three key perspectives. Firstly, it achieves state-of-the-art performance on several downstream tasks, outperforming existing EO-focused LMMs of larger scale, and surpassing specialized segmentation models for the first time on pixel-level tasks. Secondly, it significantly outperforms the baselines on EarthMind-Bench across both RGB and SAR, where it improves performance through the effective fusion mechanism.

# 2 Related Work

**Earth Observation LMMs.** Starting from the excellent foundation of image LMMs, many works [18, 14, 16, 17, 30, 15, 37, 32, 31] try to transfer the success of general image understanding to the EO field. A key challenge is the lack of instruction-tuned EO datasets. To address this, RSGPT [18] proposed the first large-scale EO image-text paired dataset, enabling conversation tasks such as image captioning and VQA. GeoChat [14] and SkyEyeGPT [16] extend LMM capabilities to region-level visual grounding by introducing region-centric instruction data. LHRS-Bot [15] leverages large-

Table 1: **(Left)** Comparison of EarthMind with existing EO LMMs. EarthMind supports both multi-granular and multi-sensor understanding. **(Right)** Comparison of EarthMind-Bench with existing EO benchmarks in terms of multi-sensor support, granularity, and task level. "S" denotes single-modality input; "M" indicates multiple modalities (used independently); "F" represents paired sensor fusion. "MC" and "OE" refer to multiple-choice and open-ended formats, respectively.

| Method | Multi-Granular | | | Multi-Sensor | | | Benchmark | Multi Sensor | Multi Gran. | Multi Level | Task Type |
|--------|:---:|:---:|:---:|:---:|:---:|---|--------|:---:|:---:|:---:|:---:|
| | Image | Region | Pixel | Handling | Fusion | | | | | | |
| RSGPT [18] | ✓ | ✗ | ✗ | ✗ | ✗ | | RSIEval [18] | S | ✗ | ✓ | MC + OE |
| GeoChat [14] | ✓ | ✓ | ✗ | ✗ | ✗ | | HnstD [34] | S | ✗ | ✗ | MC + OE |
| EarthGPT [17] | ✓ | ✓ | ✗ | ✓ | ✗ | | GeoChat-Bench [14] | S | ✗ | ✗ | OE |
| Earthmarker [30] | ✓ | ✓ | ✗ | ✗ | ✗ | | VRSBench [35] | S | ✗ | ✗ | MC + OE |
| LHRS-bot [15] | ✓ | ✗ | ✗ | ✗ | ✗ | | LHRS-Bench [15] | S | ✗ | ✓ | MC |
| SkyEyeGPT [16] | ✓ | ✗ | ✗ | ✗ | ✗ | | VLEO-Bench [36] | S | ✗ | ✓ | MC + OE |
| Skysensegpt [16] | ✓ | ✗ | ✗ | ✗ | ✗ | | FIT-RSRC [37] | S | ✗ | ✓ | MC |
| EarthDial [31] | ✓ | ✗ | ✗ | ✓ | ✗ | | UrBench [38] | S | ✓ | ✓ | MC |
| GeoPixel [32] | ✓ | ✗ | ✓ | ✗ | ✗ | | XLRS-Bench [39] | S | ✓ | ✓ | MC + OE |
| RSUniVLM [33] | ✓ | ✓ | ✓ | ✗ | ✗ | | GEOBench-VLM [40] | M | ✓ | ✓ | MC + OE |
| **EarthMind** | ✓ | ✓ | ✓ | ✓ | ✓ | | **EarthMind-Bench** | **M + F** | ✓ | ✓ | MC + OE |

scale available EO imagery aligned with OpenStreetMap annotations to improve VLM pretraining. To enhance complex reasoning, SkysenseGPT [37] introduces the FIT-RS dataset that focuses on understanding the relationships between spatial entities. GeoPixel [32] further pushes the boundary to pixel-level grounding by constructing an automatic pipeline to generate grounded conversations based on EO imagery. Beyond optical data, EarthDial [31] incorporates a diverse set of modalities, including Multispectral and Hyperspectral imagery and SAR, with the goal of improving generalization across heterogeneous EO data sources. Despite these advances, current EO LMMs remain limited in jointly supporting multi-granular tasks and multi-sensor fusion, which are essential for real applications.

**Earth Observation Multimodal Benchmarks.** The rapid progress of LMMs in the EO domain has also stimulated the development of dedicated evaluation benchmarks. RSIEval [18] provides human-annotated captions and VQA pairs to evaluate VLMs in remote sensing. LHRS-Bench [15] introduces hierarchical taxonomies to assess LMMs across multiple dimensions. VLEO-Bench [36] focuses on real-world applications, including urban monitoring, disaster relief, and land use. Beyond conversational tasks, VRSBench [35] and GeoChat-Bench [14] incorporate grounding tasks to evaluate the localization capabilities of LMMs. HnstD [34] aims to detect model hallucinations, while FIT-RSRC [37] targets reasoning about object relationships. XLRS-Bench [39] leverages ultra-high-resolution imagery for comprehensive evaluation. Concurrently, GEOBench-VLM [40] proposes a multi-task geospatial benchmark covering diverse EO tasks. Despite these advances, a notable gap remains: None of the existing benchmarks explicitly evaluates the ability of LMMs to perform multi-sensor fusion, which is a key capability in real EO scenarios.

## 3  EarthMind

### 3.1  Overview

From dense urban mapping to natural terrain analysis, EO applications require models that can unify multi-modal data to enable rich scene characterization and improved decision-making. At its core, EarthMind integrates a set of vision encoders with a language model to enable flexible and generalizable cross-modal reasoning.

As shown in Fig. 2, for spatial and semantic understanding at different levels of granularity, our model employs three specialized encoders: a visual encoder $\mathbf{E}_v$ for global semantic perception, a region encoder $\mathbf{E}_r$ for object-level understanding, and a grounding encoder $\mathbf{E}_g$ for fine-grained spatial segmentation. These components produce hierarchical representations of the input, which are projected into a shared language space using a vision-language projector (VLP). This yields a token sequence $X^V = \{x_1^v, x_2^v, ..., x_P^v\}$ that captures both global and local visual information.

To further enable dense prediction tasks, we introduce a set of learnable segmentation tokens $X^S = \{x_1^s, x_2^s, ..., x_Q^s\}$, where each token can be interpreted as a query embedding tailored to capture a specific aspect of the spatial layout. The combined sequence of visual tokens $X^V$, segmentation tokens $X^S$, and tokenized language queries is processed by an LLM, which performs joint cross-modal reasoning. The hidden states of "<SEG>" tokens and the output of $\mathbf{E}_g$ are passed to a
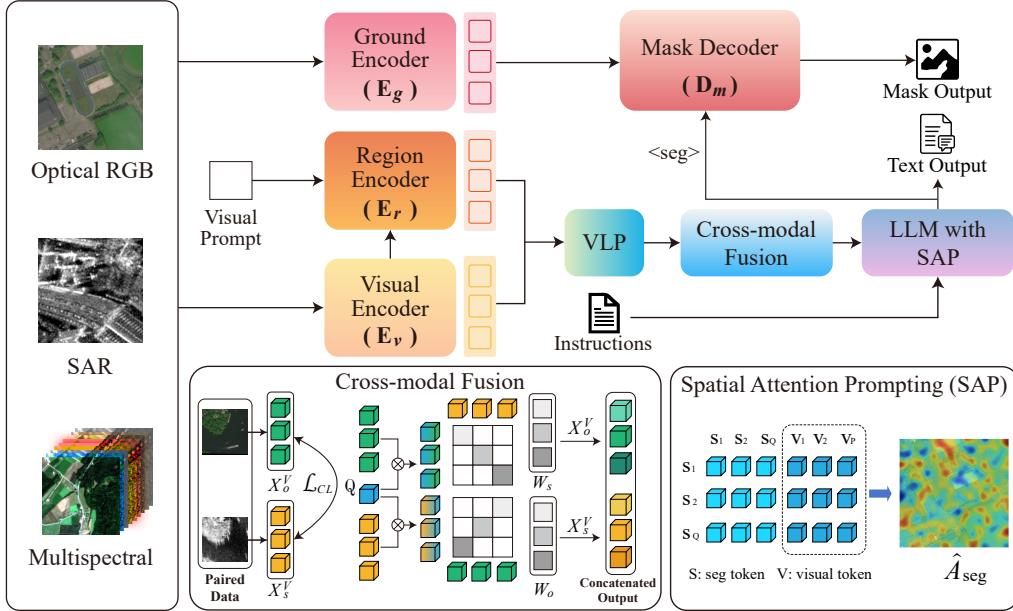
Figure 2: An overview of the EarthMind, which can handle multi-granular and multi-sensor EO data understanding. It facilitates effective interaction between modalities with Cross-modal Fusion and performs accurate pixel-level understanding with Spatial Attention Prompting.

lightweight mask decoder $\mathbf{D}_m$ to produce fine-grained segmentation maps. To further improve spatial grounding in complex EO scenes, we propose *Spatial Attention Prompting* (SAP), detailed in Sec. 3.2.

In parallel, our framework supports multi-sensor fusion by treating heterogeneous EO data as a form of video-like input. Inspired by recent advances in video-language models [41, 42], we adopt a unified data formatting strategy: single- and dual-channel SAR images are padded to form pseudo-RGB frames, while multispectral bands are grouped in triplets to construct temporally-ordered multi-frame sequences. These sequences mimic the temporal structure of video inputs and are processed by shared encoders, enabling our model to exploit cross-frame dependencies and spectral complementarity.

To enhance inter-modal alignment, we propose a *Cross-modal Fusion* module (details in Sec. 3.3) that leverages paired sensor inputs under a unified language conditioning. This module facilitates effective interaction between modalities by attending to complementary cues across spectral and spatial dimensions.

## 3.2 Spatial Attention Prompting (SAP)

In the above architecture, the "<SEG>" token acts as an implicit language query that searches for potential target objects within the image. However, empirical results show a significant performance drop when applying this framework to EO scenarios. To investigate the cause of this degradation, we examine whether the "<SEG>" token effectively captures the semantics of the queried objects aligned with the query in EO imagery.

During LLM inference, tokens interact with each other through multi-head attention layers. A token that receives higher attention weights from others is considered more influential in the model's reasoning process. Motivated by this, we measure the importance of each image token with respect to the "<SEG>" token by computing the cross-attention maps between segmentation tokens and all image tokens across all transformer layers, resulting in $A_{\text{seg}} \in \mathbb{R}^{L \times H \times Q \times P}$, where $L$ is the number of transformer layers and $H$ is the number of attention heads. We compute the attention map by applying softmax over the image token dimension, and average over the heads to obtain a Seg-to-Image attention map:

$$\hat{A}_{\text{seg}} = \frac{1}{H} \sum_{h=1}^{H} \text{Softmax}(A_{\text{seg}}^{(l,h)}), \quad \hat{A}_{\text{seg}} \in \mathbb{R}^{L \times Q \times P} \tag{1}$$

4

As shown in Fig. 2, the attention map reflects the spatial focus of the LLM's segmentation prompts. However, EO imagery often exhibits extreme object scale variations, ambiguous textures, and weak semantic boundaries, which lead to attention drifting away from the target regions. Our attention visualizations further illustrate this phenomenon in Sec. 5.4.

To guide the model toward more semantically meaningful spatial priors, we introduce a supervision signal based on ground-truth masks. Specifically, we downsample the binary mask to match the image token resolution and treat it as a probability distribution over tokens. Then, we minimize the Kullback–Leibler (KL) divergence [43] between the normalized attention map and the target distribution:

$$\mathcal{L}_{\text{KL}} = \sum_{l=1}^{L} \sum_{q=1}^{Q} \text{KL}\left(\hat{A}_{\text{seg}}^{(l,q)} \,\|\, G_q\right) \tag{2}$$

where $G_q \in \mathbb{R}^P$ is the ground-truth token-level distribution corresponding to the $q$-th seg token. This training objective encourages the segmentation tokens to attend more accurately to the queried regions, thereby enhancing pixel-level mask prediction.

## 3.3 Cross-modal Fusion

Given multi-sensor input, we extract cross-modal image tokens before feeding into the LLM: $X_o^V = \{x_1^o, x_2^o, ..., x_P^o\}$ and $X_s^V = \{x_1^s, x_2^s, ..., x_P^s\}$, which represent features from optical (RGB) and non-optical (e.g., SAR) modalities, respectively. To fuse modalities, we propose the concepts of Modality Alignment and Modality Mutual Attention, as detailed below.

**Modality Alignment.** To facilitate effective multi-sensor fusion, we adopt an online contrastive learning strategy to align non-optical features (e.g., SAR) with the optical (RGB) feature space, leveraging the pretrained vision-language projector for optical inputs. Specifically, given a training batch of size $B$, we define the cross-modal contrastive loss as:

$$\mathcal{L}_{CL} = -\frac{1}{B} \sum_{i=1}^{B} \left( \log \frac{f(x_i^o, x_i^s)}{f(x_i^o, x_i^s) + \sum_{k \neq i} f(x_i^o, x_k^s)} + \log \frac{f(x_i^s, x_i^o)}{f(x_i^s, x_i^o) + \sum_{k \neq i} f(x_i^s, x_k^o)} \right) \tag{3}$$

where $f(\cdot)$ represents the cross-modal cosine similarity operation between paired tokens.

**Modality Mutual Attention.** After aligning different modalities into a shared semantic space, we further enable the model to dynamically assess the information content of each modality, thereby emphasizing the most informative representations for downstream LLM reasoning.

Specifically, we first introduce a set of learnable queries $Q \in \mathbb{R}^{k \times D}$ to extract neighborhood-aware features from each modality. For example, the representation of SAR is computed via element-wise multiplication: $\hat{X}_s^V = Q \odot X_s^V$, where $\hat{X}_s^V \in \mathbb{R}^{k \times P \times D}$. RGB features are processed similarly to obtain $\hat{X}_o^V$. We then compute cross-modal importance weights using a dot product:

$$W_s = \text{Softmax}(X_o^V (\hat{X}_s^V)^\top), \quad W_o = \text{Softmax}(X_s^V (\hat{X}_o^V)^\top) \tag{4}$$

where the weights $W_s, W_o \in \mathbb{R}^P$ reflect the relevance of the SAR and RGB tokens in cross-modal contexts, by which we reweight $X_s^V$ and $X_o^V$ accordingly and concatenate them into the LLM for joint understanding.

## 3.4 Training

EarthMind is trained via instruction tuning across diverse supervision formats, including VQA and segmentation tasks. The model learns to generate the target response conditioned on the image tokens, segmentation tokens, and the task instruction. Formally, the generation probability of the next token is defined as:

$$\Pr\left( t_{i+1} \mid \underbrace{\texttt{<img>}_1, \ldots, \texttt{<img>}_P}_{\text{image tokens}}, \underbrace{\texttt{<seg>}_1, \ldots, \texttt{<seg>}_Q}_{\text{segmentation tokens}}, \underbrace{s_1, \ldots, s_M}_{\text{instruction}}, \underbrace{t_1, \ldots, t_i}_{\text{generated response}} ; \boldsymbol{\Theta} \right) \tag{5}$$

where $\boldsymbol{\Theta}$ denotes the trainable parameters of the LMM. For VQA tasks, we apply standard autoregressive training, minimizing the token-level cross-entropy loss over the ground-truth responses.
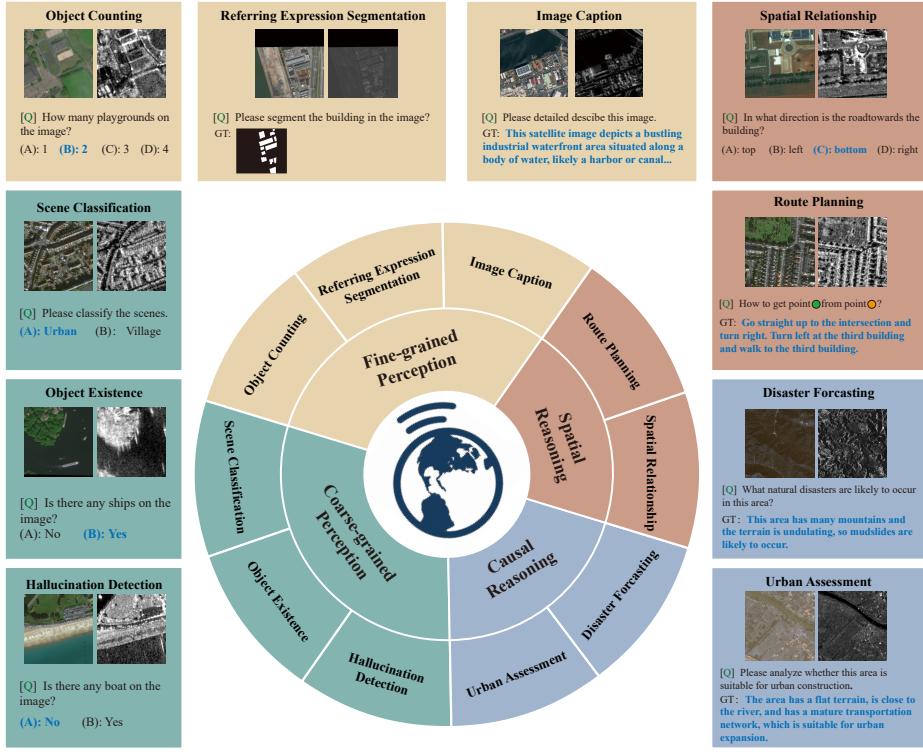
Figure 3: Examples of EarthMind-Bench. There are ten tasks to evaluate the multi-sensor fusion ability of LMMs, covering perception and reasoning levels. The LLMs are asked to solve the problem (with the ground-truth answers marked in blue) based on RGB, SAR or RGB-SAR fusion settings.

For segmentation tasks, we use a combination of pixel-wise cross-entropy Loss and Dice Loss. In addition, the KL loss described in Sec. 3.2 is applied to supervise attention-based spatial priors.

# 4 EarthMind-Bench

We observe that existing EO benchmarks lack support for multi-granular understanding and multi-sensor fusion. To address this limitation, we propose EarthMind-Bench, a new benchmark constructed from high-quality RGB–SAR paired data across various public datasets, including OpenEarthMap-SAR [44], DFC2023 Track2 [45], WHU-OPT-SAR [46], MSAW [47], and MultiResSAR [48]. We curate 2,000 samples from their test sets, and design a suite of 10 tasks spanning perception and reasoning, enabling a comprehensive evaluation of LMMs in EO scenarios, as shown in Fig. 3.

**Perception.** We include six tasks to evaluate both coarse- and fine-grained visual understanding, such as *scene classification*, *object existence*, *object counting*, and *image captioning*. In particular, we highlight two challenging tasks: *1) Hallucination detection*, which requires models to answer some leading questions, and *2) Referring expression segmentation*, where models must produce binary masks conditioned on natural language queries.

**Reasoning.** We organize reasoning tasks into two categories. First, spatial reasoning tasks include *spatial relationship*, where models infer the relative positions of given objects, and *route planning*, which requires generating a feasible path from a specified start point to an end point. Second, causal reasoning tasks consist of *disaster forecasting* and *urban development assessment*, where models are expected not only to make predictions or judgments, but also to provide reasoning evidence.

All tasks are formatted as either multiple-choice (MC) or open-ended (OE) questions. For MC tasks, we report average accuracy across tasks as the evaluation metric. For OE tasks, we adopt GPT-based scoring to assess the quality of generated responses. All annotations are produced by human experts and have undergone a rigorous verification process, which is described in detail in Appendix D.

6

Table 2: Experimental results on EarthMind-Bench, in which different evaluation settings are employed to compare multi-sensor understanding ability. "M-Avg": the average performance of multiple-choice tasks; "O-Avg": the average performance of open-ended tasks. "Referring Segmentation" refers to referring expression segmentation, which is evaluated by mIoU. † denotes proprietary models. Bold means the best performance.

| Model | Size | M-Avg | O-Avg | Scene Class. | Object Exist. | Halluci. Detect. | Object Count. | Spatial Relation. | Referring Segmen. | Image Caption | Disaster Forecast. | Route Plann. | Urban Assess. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Full mark* | – | 100 | 5 | 100 | 100 | 100 | 100 | 100 | 100 | 5 | 5 | 5 | 5 |
| **Evaluation on RGB modality** | | | | | | | | | | | | | |
| GPT-4o† [60] | - | **70.0** | 2.63 | **97.3** | 79.9 | **86.4** | 34.0 | **52.3** | - | **4.58** | 1.75 | **2.01** | 2.18 |
| GPT-4V† [1] | - | 61.5 | 2.12 | 90.7 | 72.9 | 75.9 | 39.0 | 29.2 | - | 3.28 | 1.54 | 1.82 | 1.86 |
| GeoChat [14] | 7B | 41.6 | 1.78 | 71.3 | 51.8 | 46.8 | 18.9 | 19.0 | - | 1.92 | 1.73 | 1.33 | 2.14 |
| LHRS-bot [15] | 7B | 47.4 | 1.84 | 76.0 | 58.3 | 58.3 | 25.2 | 19.4 | - | 2.56 | 1.75 | 1.55 | 1.50 |
| Skysensegpt [37] | 7B | 45.7 | 1.45 | 77.1 | 56.8 | 56.7 | 26.8 | 11.1 | - | 1.68 | 1.40 | 1.18 | 1.55 |
| GeoPixel [32] | 7B | 59.4 | 2.06 | 87.2 | 67.8 | 73.6 | 33.5 | 34.7 | 46.8 | 2.80 | 1.80 | 1.68 | 1.95 |
| **EarthMind** | 4B | 69.0 | **2.82** | 96.5 | **81.2** | 81.6 | **47.3** | 38.4 | **54.0** | 3.35 | **3.37** | **2.01** | **2.55** |
| **Evaluation on SAR modality** | | | | | | | | | | | | | |
| GPT-4o† [60] | - | 55.9 | 2.40 | 75.6 | 71.4 | 72.9 | 22.8 | 36.6 | - | 2.89 | 3.05 | 1.65 | 2.04 |
| GPT-4V† [1] | - | 50.5 | 2.22 | 61.2 | 73.3 | 67.1 | 31.0 | 19.9 | - | 2.63 | 2.98 | 1.40 | 1.85 |
| GeoChat [14] | 7B | 40.2 | 1.45 | 58.1 | 49.8 | 46.8 | 27.9 | 18.5 | - | 1.78 | 1.65 | 1.25 | 1.56 |
| LHRS-bot [15] | 7B | 41.0 | 1.71 | 58.5 | 56.3 | 48.5 | 23.5 | 18.1 | - | 1.86 | 1.70 | 1.45 | 1.83 |
| Skysensegpt [37] | 7B | 39.9 | 1.55 | 51.9 | 53.8 | 49.2 | 33.8 | 10.7 | - | 1.76 | 1.70 | 1.35 | 1.38 |
| GeoPixel [32] | 7B | 53.0 | 1.80 | 76.8 | 59.0 | 65.7 | 30.5 | 32.8 | 35.9 | 2.08 | 1.97 | 1.45 | 1.68 |
| **EarthMind** | 4B | 67.5 | 2.64 | 95.4 | 77.4 | 74.6 | 46.8 | 43.1 | 53.0 | 3.10 | 3.25 | 1.89 | 2.30 |
| **Evaluation on RGB-SAR Fused modality** | | | | | | | | | | | | | |
| GPT-4o† [60] | - | 67.7 | 2.28 | **97.7** | 79.6 | **86.2** | 31.6 | 43.5 | - | 3.68 | 1.59 | 1.82 | 2.03 |
| GPT-4V† [1] | - | 58.2 | 1.93 | 91.1 | 64.8 | 62.4 | 32.8 | 39.8 | - | 2.89 | 1.48 | 1.57 | 1.79 |
| **EarthMind** | 4B | 70.6 | 3.02 | 97.7 | 82.4 | 85.4 | 47.3 | 40.3 | 54.5 | 3.80 | 3.37 | 2.21 | 2.70 |

# 5 Experiments

## 5.1 Implementation Details

EarthMind builds upon InternVL2 [49] and SAM2 [28], adopting a curriculum learning strategy to fine-tune the model across three progressive stages. First, we enhance the instruction-following capability using 1.7M general image-text data, which covers image-level captioning, VQA, region-level object understanding and text-driven segmentation. Second, we introduce 1M EO-specific multimodal data to adapt EarthMind to the remote sensing domain. Third, we utilize our synthesized multi-sensor conversation corpus, and selectively retain examples from earlier stages to mitigate catastrophic forgetting. We fine-tune EarthMind with a learning rate of 4e-5 and a batch size of 2, training only the vision-language projector, the LLM via LoRA [50], and the mask decoder. All experiments are conducted on 8 NVIDIA A100-80G GPUs. More details about training datasets and details can be seen in Appendix F.

## 5.2 Benchmarks

In addition to our proposed EarthMind-Bench, we further evaluate the effectiveness of EarthMind on several widely-used EO benchmarks, assessing its capability in both multi-granularity and multi-sensor understanding. First, we conduct classification and vision-language evaluations on AID [51], UCMerced [52], RSVQA-HRBEN [53] and the VQA task of VRSBench [35] to perform image-level evaluation. AID is a large-scale aerial image classification dataset consisting of 10,000 images from 30 scene categories, collected via Google Earth. We follow the evaluation protocol of [14], using 20% of the dataset for testing. UCMerced contains 2,100 aerial images from 21 land use classes, commonly used for remote sensing scene classification. RSVQA-HRBEN consists of 47k question-answer pairs, in which we follow [14] to test on the "presence" and "Comparison" categories. The VQA part of VRSBench contains 123k samples covering diverse EO scenarios. Second, for region-level evaluation, we evaluate region-level understanding across region-level captioning and visual grounding task, where DIOR-RSVG [54] and visual grounding task of VRSBench are employed. Third, we conduct referring expression segmentation on the test sets of RefSegRS [55] and RRSIS-D [56] for pixel-level grounding, which include 1,817 and 3,481 test instances, respectively. Last, we involve the test set of BigEarthNet [57], SoSAT-LCZ42 [58], which consists of multispectral data and SAR ship detection dataset [59] to prove our multi-sensor understanding ability.

Table 3: Quantitative performance of EarthMind on public benchmarks. For image-level evaluation, we report accuracy on AID, UC-Merced, RSVQA-HRBEN, and the VQA task of VRSBench. For region-level evaluation, we report CIDEr scores on DIOR-RSVG and Acc@0.5 on the Visual Grounding task of VRSBench. For pixel-level tasks, mean Intersection over Union (mIoU) is reported.

| Method | Image-level | | | | Region-level | | Method | Pixel-level | |
|---|---|---|---|---|---|---|---|---|---|
| | AID | UC | RSVQA | VRS-VQA | RSVG | VRS-VG | | RRSIS-D | RefSegRS |
| GPT-4o [60] | 74.7 | 88.8 | - | - | - | - | LAVT [61] | 56.8 | 47.4 |
| LLaVA-1.5 [6] | 72.0 | 84.4 | 63.1 | 76.4 | - | 5.1 | RIS-DMMI [62] | 60.3 | 52.2 |
| GeoChat [14] | 72.0 | 84.4 | 72.3 | 76.0 | 30.9 | 49.8 | Caris [63] | 62.1 | 42.7 |
| EarthGPT [17] | - | - | 72.0 | - | 232.8 | - | RM-SIN [56] | 64.2 | 42.6 |
| EarthMarker [30] | 78.0 | 86.5 | - | - | 379.3 | - | CroBIM [64] | 64.5 | 59.8 |
| EarthDial [31] | 88.8 | 92.4 | 72.5 | - | - | - | GeoPixel [32] | 67.3 | - |
| **EarthMind** | **97.2** | **95.0** | **74.0** | **78.9** | **428.2** | **55.6** | **EarthMind** | **82.2** | **62.6** |

Table 4: Quantitative performance on public multi-sensor benchmarks including Multispectral and SAR domain. Evaluation follows the protocol of [31]. † indicates that the results on benchmarks were reproduced using their official weights.

| Method | Multispectral | | Method | SAR | | | | |
|---|---|---|---|---|---|---|---|---|
| | BigEarthNet | SoSAT-LCZ42 | | Small | Medium | Large | Single | Multiple |
| GPT-4o [60] | 49.0 | 15.5 | GPT-4o [60] | 0.70 | 0.90 | 3.20 | 1.20 | 0 |
| Qwen2.5-VL†[65] | 36.2 | 18.9 | GeoChat†[14] | 2.61 | 4.92 | 6.95 | 2.58 | 1.40 |
| EarthDial [31] | 69.9 | **60.7** | EarthDial [31] | 12.14 | 26.02 | 35.56 | 26.03 | 6.06 |
| **EarthMind** | **70.4** | 58.3 | **EarthMind** | **13.58** | **28.55** | **36.78** | **27.45** | **6.99** |

## 5.3 Main Results

**Results on EarthMind-Bench.** EarthMind-Bench supports evaluation under three settings: RGB only, SAR only, and RGB–SAR fusion. We compare EarthMind with state-of-the-art EO-specific LMMs and proprietary models such as GPT-4V [1] and GPT-4o [60]. Tab. 2 summarizes the results, from which we highlight three key findings: *1) Fine-grained and open-ended tasks remain challenging for existing LMMs.* While coarse-grained tasks like scene classification are largely solved, tasks such as object counting and spatial relationship understanding still exhibit significant performance gaps. In particular, referring segmentation proves difficult for most models due to the lack of pixel-level reasoning ability. *2) Most models generalize poorly to SAR inputs.* Performance under SAR-only settings lags behind RGB settings for nearly all models, likely due to training data limitations. In contrast, EarthMind demonstrates strong generalization to SAR inputs, benefiting from multi-sensor training. *3) Effective fusion requires learning cross-modal complementarity, not just stacking modalities.* As the only open-source model with explicit fusion capability, EarthMind is compared against GPT-4 variants by feeding both RGB and SAR as pseudo-RGB inputs using the official multi-image interface. Results show that GPT-4 models experience performance degradation compared to RGB-only inputs, particularly on fine-grained tasks like Route Planning, Object Counting, and Spatial Relationship Understanding, where SAR data provides robust structural cues under adverse conditions. In contrast, EarthMind effectively captures cross-modal complementarity, yielding consistent improvements over single-modality inputs.

**Results on Public Benchmarks.** We evaluate EarthMind on a series of mainstream EO benchmarks (see Tab. 3). EarthMind consistently delivers strong performance across multiple levels of understanding. On image-level tasks such as scene classification (AID, UC-Merced) and EO visual question answering (RSVQA-HRBEN and VRSBench-VQA), it significantly outperforms previous models, including GPT-4o, despite using only 4B parameters. Moving to region-level tasks, EarthMind achieves a CIDEr score of 428.2 on DIOR-RSVG and 55.6% accuracy on VRSBench visual grounding, outperforming visual-prompt-based methods and highlighting the strength of our region encoder. When it comes to pixel-level benchmarks, EarthMind achieves top results on both RRSIS-D and RefSegRS, surpassing even specialized segmentation models and EO-focused LMMs. Beyond the RGB domain, EarthMind also demonstrates competitive results on multi-sensor understanding, as shown in Tab. 4, demonstrating the strong generalization towards complex EO scenarios.
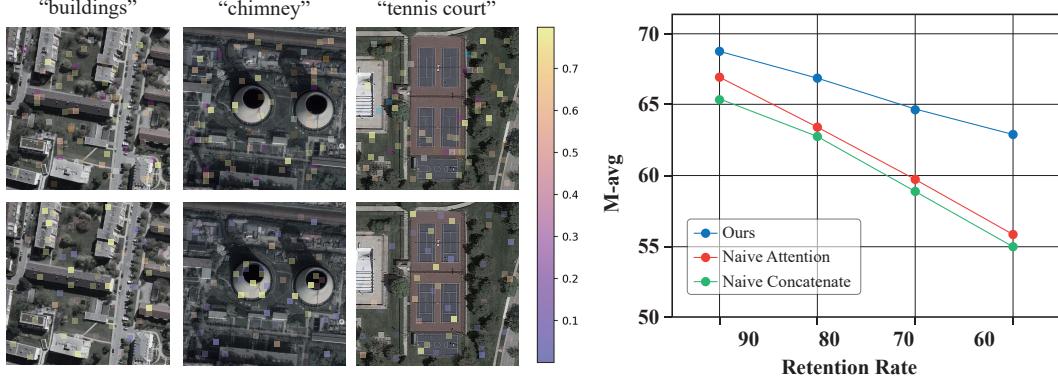
Figure 4: **(Left)** Visualization of attention maps between segmentation tokens and image tokens within LLM. The top row is the results of "w/o SAP" and the bottom row is the results of "w/ SAP", where the attention is reallocated to the region specified by the query. **(Right)** Comparison of our modality-mutual attention with other baselines under the "Token drop" study.

Table 5: Ablation on Spatial Attention Promoting **(Left)** and Cross-Modal Fusion **(Right)**

| Settings | RefSegRS | | RRSIS-D | |
|---|---|---|---|---|
| | mIoU ↑ | oIoU ↑ | mIoU ↑ | oIoU ↑ |
| w/o SAP | 46.7 | 86.0 | 67.5 | 94.8 |
| w/ SAP (first layer) | 46.4 | 86.0 | 67.7 | 95.0 |
| w/ SAP (middle layer) | 48.0 | 86.9 | 71.8 | 95.5 |
| w/ SAP (last layer) | 47.2 | 86.5 | 69.2 | 95.0 |
| w/ SAP (all layers) | **48.4** | **87.4** | **72.0** | **96.0** |

| Settings | M-Avg ↑ | mIoU ↑ |
|---|---|---|
| w/o Modality Alignment | 69.0 | 53.6 |
| Naive Concatenate | 68.4 | 52.7 |
| Naive Attention | 68.6 | 52.5 |
| Ours | **70.6** | **54.5** |

## 5.4 Ablations

EarthMind significantly enhances multi-granularity and multi-sensor EO understanding through the proposed (i) Spatial Attention Prompting (SAP) and (ii) Cross-modal Fusion. We delve into each of these components in the following sections. More discussions can be seen in Appendix G.

**Effectiveness of SAP.** First, we evaluate the effectiveness of our proposed SAP method. To ensure a fair comparison, we train models from scratch on the training splits of RefSegRS and RRSIS-D, and report results on their corresponding test sets. As shown in Tab. 5, incorporating SAP significantly improves segmentation performance. Second, we visualize the attention maps between the <SEG> token and image tokens across transformer layers within the LLM. As illustrated in Fig. 4, SAP effectively suppresses attention drift and helps reallocate attention toward the anchor regions corresponding to the query. Lastly, we investigate how the performance of SAP varies with the supervision layer. Results show that applying supervision at middle layers yields the strongest performance boost, as these layers balance low-level spatial detail and high-level semantics. Importantly, SAP introduces no additional parameters, making it a lightweight yet effective enhancement for pixel-level grounding within LLMs.

**Effectiveness of Cross-modal Fusion.** To evaluate the effectiveness of our cross-modal fusion design, we conduct ablations on the multi-sensor setting with three configurations: (1) w/o Modality Alignment; (2) Naive Concatenate: visual tokens from different modalities are directly concatenated and passed to the LLM; (3) Naive Attention: token importance is computed using cosine similarity between paired SAR and optical features. We report both multiple-choice accuracy and referring expression segmentation results on EarthMind-Bench. As shown in Tab. 5, our full model significantly outperforms all baselines, demonstrating the importance of both alignment and adaptive weighting for robust fusion. Furthermore, to assess the informativeness of modality-selected tokens, we compare different token retention strategies. For baselines, we randomly discard a portion of the input tokens before feeding them into the LLM. In contrast, our method selectively drops tokens with lower learned importance weights. As shown in Fig. 4, even under the same retention ratio, EarthMind retains more informative content, leading to better performance. This demonstrates that our method reserves the most complementary information across them.

9

# 6 Limitations and Conclusion

**Limitations of EarthMind.** Training EarthMind demands considerable computational resources due to its use of multiple visual encoders for multi-level understanding. A promising direction is to optimize the architecture via Mixture-of-Experts or knowledge distillation to reduce redundancy. Additionally, a modality-aligned encoder that jointly embeds heterogeneous sensor inputs into a shared semantic space could further improve efficiency.

In this work, we propose EarthMind, a unified framework for multi-granular and multi-sensor understanding in Earth Observation (EO). To enhance pixel-level performance, we introduce Spatial Attention Prompting (SAP), which guides attention reallocation within the LLM. We also design a Cross-modal Fusion mechanism that aligns heterogeneous modalities into a shared space and adaptively selects informative features. To support evaluation, we curate EarthMind-Bench, the first benchmark tailored for multi-sensor fusion in LMMs. Extensive experiments validate the effectiveness and strong generalization of EarthMind across diverse EO tasks.

# References

[1] OpenAI. Gpt-4 technical report, 2023.

[2] OpenAI. Introducing chatgpt. https://openai.com/blog/chatgpt, 2022.

[3] Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. Pllava: Parameter-free llava extension from images to videos for video dense captioning. *arXiv preprint arXiv:2404.16994*, 2024.

[4] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. In *ECCV*, pages 370–387. Springer, 2024.

[5] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742. PMLR, 2023.

[6] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36:34892–34916, 2023.

[7] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

[8] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *CVPR*, pages 9579–9589, 2024.

[9] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *CVPR*, pages 13009–13018, 2024.

[10] Haobo Yuan, Xiangtai Li, Tao Zhang, Zilong Huang, Shilin Xu, Shunping Ji, Yunhai Tong, Lu Qi, Jiashi Feng, and Ming-Hsuan Yang. Sa2va: Marrying sam2 with llava for dense grounded understanding of images and videos. *arXiv preprint arXiv:2501.04001*, 2025.

[11] Laila Bashmal, Yakoub Bazi, Farid Melgani, Mohamad M Al Rahhal, and Mansour Abdulaziz Al Zuair. Language integration in remote sensing: Tasks, datasets, and future directions. *IEEE Geoscience and Remote Sensing Magazine*, 11(4):63–93, 2023.

[12] Nathan Ratledge, Gabe Cadamuro, Brandon De La Cuesta, Matthieu Stigler, and Marshall Burke. Using machine learning to assess the livelihood impact of electricity access. *Nature*, 611(7936):491–495, 2022.

[13] Jun Yang, Peng Gong, Rong Fu, Minghua Zhang, Jingming Chen, Shunlin Liang, Bing Xu, Jiancheng Shi, and Robert Dickinson. The role of satellite remote sensing in climate change studies. *Nature climate change*, 3(10):875–883, 2013.

[14] Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, and Fahad Shahbaz Khan. Geochat: Grounded large vision-language model for remote sensing. In *CVPR*, pages 27831–27840, 2024.

[15] Dilxat Muhtar, Zhenshi Li, Feng Gu, Xueliang Zhang, and Pengfeng Xiao. Lhrs-bot: Empowering remote sensing with vgi-enhanced large multimodal language model. In *ECCV*, pages 440–457. Springer, 2024.

[16] Yang Zhan, Zhitong Xiong, and Yuan Yuan. Skyeyegpt: Unifying remote sensing vision-language tasks via instruction tuning with large language model. *ISPRS Journal of Photogrammetry and Remote Sensing,*, 221:64–77, 2025.

[17] Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, and Xuerui Mao. Earthgpt: A universal multi-modal large language model for multi-sensor image comprehension in remote sensing domain. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.

[18] Yuan Hu, Jianlong Yuan, Congcong Wen, Xiaonan Lu, Yu Liu, and Xiang Li. Rsgpt: A remote sensing vision language model and benchmark. *ISPRS Journal of Photogrammetry and Remote Sensing,*, 224:272–286, 2025.

[19] Ioannis Kotaridis and Maria Lazaridou. Remote sensing image segmentation advances: A meta-analysis. *ISPRS Journal of Photogrammetry and Remote Sensing,*, 173:309–322, 2021.

[20] Mohammad D Hossain and Dongmei Chen. Segmentation for object-based image analysis (obia): A review of algorithms and challenges from remote sensing perspective. *ISPRS Journal of Photogrammetry and Remote Sensing,*, 150:115–134, 2019.

[21] Libao Zhang and Kaina Yang. Region-of-interest extraction based on frequency domain analysis and salient region detection for remote sensing image. *IEEE Geoscience and Remote Sensing Letters*, 11(5):916–920, 2013.

[22] Germain Forestier, Anne Puissant, Cédric Wemmert, and Pierre Gançarski. Knowledge-based region labeling for remote sensing image interpretation. *Computers, Environment and Urban Systems*, 36(5):470–480, 2012.

[23] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.

[24] Gong Cheng, Xingxing Xie, Junwei Han, Lei Guo, and Gui-Song Xia. Remote sensing image scene classification meets deep learning: Challenges, methods, benchmarks, and opportunities. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 13:3735–3756, 2020.

[25] Yuanzhi Zhang, Hongsheng Zhang, and Hui Lin. Improving the impervious surface estimation with combined use of optical and sar remote sensing images. *Remote Sensing of Environment*, 141:155–167, 2014.

[26] Michael Schmitt, Florence Tupin, and Xiao Xiang Zhu. Fusion of sar and optical remote sensing data—challenges and recent trends. In *IGARSS*, pages 5458–5461. IEEE, 2017.

[27] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*, pages 4015–4026, 2023.

[28] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.

[29] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

[30] Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, Jun Li, and Xuerui Mao. Earthmarker: A visual prompting multi-modal large language model for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.

[31] Sagar Soni, Akshay Dudhane, Hiyam Debary, Mustansar Fiaz, Muhammad Akhtar Munir, Muhammad Sohail Danish, Paolo Fraccaro, Campbell D Watson, Levente J Klein, Fahad Shahbaz Khan, et al. Earthdial: Turning multi-sensory earth observations to interactive dialogues. *arXiv preprint arXiv:2412.15190*, 2024.

[32] Akashah Shabbir, Mohammed Zumri, Mohammed Bennamoun, Fahad S Khan, and Salman Khan. Geopixel: Pixel grounding large multimodal model in remote sensing. *ICML*, 2025.

[33] Xu Liu and Zhouhui Lian. Rsunivlm: A unified vision language model for remote sensing via granularity-oriented mixture of experts. *arXiv preprint arXiv:2412.05679*, 2024.

[34] Chao Pang, Xingxing Weng, Jiang Wu, Jiayu Li, Yi Liu, Jiaxing Sun, Weijia Li, Shuai Wang, Litong Feng, Gui-Song Xia, et al. Vhm: Versatile and honest vision language model for remote sensing image analysis. In *AAAI*, volume 39, pages 6381–6388, 2025.

[35] Xiang Li, Jian Ding, and Mohamed Elhoseiny. Vrsbench: A versatile vision-language benchmark dataset for remote sensing image understanding. *arXiv preprint arXiv:2406.12384*, 2024.

[36] Chenhui Zhang and Sherrie Wang. Good at captioning bad at counting: Benchmarking gpt-4v on earth observation data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7839–7849, 2024.

[37] Junwei Luo, Zhen Pang, Yongjun Zhang, Tingzhu Wang, Linlin Wang, Bo Dang, Jiangwei Lao, Jian Wang, Jingdong Chen, Yihua Tan, et al. Skysensegpt: A fine-grained instruction tuning dataset and model for remote sensing vision-language understanding. *arXiv preprint arXiv:2406.10100*, 2024.

[38] Baichuan Zhou, Haote Yang, Dairong Chen, Junyan Ye, Tianyi Bai, Jinhua Yu, Songyang Zhang, Dahua Lin, Conghui He, and Weijia Li. Urbench: A comprehensive benchmark for evaluating large multimodal models in multi-view urban scenarios. In *AAAI*, volume 39, pages 10707–10715, 2025.

[39] Fengxiang Wang, Hongzhen Wang, Mingshuo Chen, Di Wang, Yulin Wang, Zonghao Guo, Qiang Ma, Long Lan, Wenjing Yang, Jing Zhang, et al. Xlrs-bench: Could your multimodal llms understand extremely large ultra-high-resolution remote sensing imagery? *arXiv preprint arXiv:2503.23771*, 2025.

[40] Muhammad Sohail Danish, Muhammad Akhtar Munir, Syed Roshaan Ali Shah, Kartik Kuckreja, Fahad Shahbaz Khan, Paolo Fraccaro, Alexandre Lacoste, and Salman Khan. Geobench-vlm: Benchmarking vision-language models for geospatial tasks. *arXiv preprint arXiv:2411.19325*, 2024.

[41] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.

[42] Yan Shu, Zheng Liu, Peitian Zhang, Minghao Qin, Junjie Zhou, Zhengyang Liang, Tiejun Huang, and Bo Zhao. Video-xl: Extra-long vision language model for hour-scale video understanding. *arXiv preprint arXiv:2409.14485*, 2024.

[43] Tim Van Erven and Peter Harremos. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820, 2014.

[44] Junshi Xia, Hongruixuan Chen, Clifford Broni-Bediako, Yimin Wei, Jian Song, and Naoto Yokoya. Openearthmap-sar: A benchmark synthetic aperture radar dataset for global high-resolution land cover mapping. *arXiv preprint arXiv:2501.10891*, 2025.

[45] Claudio Persello, Ronny Hänsch, Gemine Vivone, Kaiqiang Chen, Zhiyuan Yan, Deke Tang, Hai Huang, Michael Schmitt, and Xian Sun. 2023 ieee grss data fusion contest: Large-scale fine-grained building classification for semantic urban reconstruction [technical committees]. *IEEE Geoscience and Remote Sensing Magazine*, 11(1):94–97, 2023.

[46] Xue Li, Guo Zhang, Hao Cui, Shasha Hou, Shunyao Wang, Xin Li, Yujia Chen, Zhijiang Li, and Li Zhang. Mcanet: A joint semantic segmentation framework of optical and sar images for land use classification. *International Journal of Applied Earth Observation and Geoinformation*, 106:102638, 2022.

[47] Jacob Shermeyer, Daniel Hogan, Jason Brown, Adam Van Etten, Nicholas Weir, Fabio Pacifici, Ronny Hansch, Alexei Bastidas, Scott Soenen, Todd Bacastow, et al. Spacenet 6: Multi-sensor all weather mapping dataset. In *CVPR Workshops*, pages 196–197, 2020.

[48] Wenfei Zhang, Ruipeng Zhao, Yongxiang Yao, Yi Wan, Peihao Wu, Jiayuan Li, Yansheng Li, and Yongjun Zhang. Multi-resolution sar and optical remote sensing image registration methods: A review, datasets, and future perspectives. *arXiv preprint arXiv:2502.01002*, 2025.

[49] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.

[50] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.

[51] Gui-Song Xia, Jingwen Hu, Fan Hu, Baoguang Shi, Xiang Bai, Yanfei Zhong, Liangpei Zhang, and Xiaoqiang Lu. Aid: A benchmark data set for performance evaluation of aerial scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(7):3965–3981, 2017.

[52] Yi Yang and Shawn Newsam. Bag-of-visual-words and spatial extensions for land-use classification. In *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*, pages 270–279, 2010.

[53] Sylvain Lobry, Diego Marcos, Jesse Murray, and Devis Tuia. Rsvqa: Visual question answering for remote sensing data. *IEEE Transactions on Geoscience and Remote Sensing*, 58(12):8555–8566, 2020.

[54] Yang Zhan, Zhitong Xiong, and Yuan Yuan. Rsvg: Exploring data and models for visual grounding on remote sensing data. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–13, 2023.

[55] Zhenghang Yuan, Lichao Mou, Yuansheng Hua, and Xiao Xiang Zhu. Rrsis: Referring remote sensing image segmentation. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.

[56] Sihan Liu, Yiwei Ma, Xiaoqing Zhang, Haowei Wang, Jiayi Ji, Xiaoshuai Sun, and Rongrong Ji. Rotated multi-scale interaction network for referring remote sensing image segmentation. In *CVPR*, pages 26658–26668, 2024.

[57] Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In *IGARSS*, pages 5901–5904. IEEE, 2019.

[58] Xiao Xiang Zhu, Jingliang Hu, Chunping Qiu, Yilei Shi, Jian Kang, Lichao Mou, Hossein Bagheri, Matthias Häberle, Yuansheng Hua, Rong Huang, et al. So2sat lcz42: A benchmark dataset for global local climate zones classification. *arXiv preprint arXiv:1912.12171*, 2019.

[59] Yuanyuan Wang, Chao Wang, Hong Zhang, Yingbo Dong, and Sisi Wei. A sar dataset of ship detection for deep learning under complex backgrounds. *remote sensing*, 11(7):765, 2019.

[60] OpenAI. Gpt-4o. https://openai.com/index/hello-gpt-4o/, May 2024.

[61] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Lavt: Language-aware vision transformer for referring image segmentation. In *CVPR*, pages 18155–18165, 2022.

[62] Yutao Hu, Qixiong Wang, Wenqi Shao, Enze Xie, Zhenguo Li, Jungong Han, and Ping Luo. Beyond one-to-one: Rethinking the referring image segmentation. In *ICCV*, pages 4067–4077, 2023.

[63] Sun-Ao Liu, Yiheng Zhang, Zhaofan Qiu, Hongtao Xie, Yongdong Zhang, and Ting Yao. Caris: Context-aware referring image segmentation. In *ACM MM*, pages 779–788, 2023.

[64] Zhe Dong, Yuzhe Sun, Yanfeng Gu, and Tianzhu Liu. Cross-modal bidirectional interaction model for referring remote sensing image segmentation. *arXiv preprint arXiv:2410.08613*, 2024.

[65] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[66] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023.

[67] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video understanding. *arXiv preprint arXiv:2406.04264*, 2024.

[68] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, pages 787–798, 2014.

[69] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, pages 69–85. Springer, 2016.

[70] Yuqian Yuan, Wentong Li, Jian Liu, Dongqi Tang, Xinjie Luo, Chi Qin, Lei Zhang, and Jianke Zhu. Osprey: Pixel understanding with visual instruction tuning. In *CVPR*, pages 28202–28211, 2024.

[71] Xiangrui Liu, Yan Shu, Zheng Liu, Ao Li, Yang Tian, and Bo Zhao. Video-xl-pro: Reconstructive token compression for extremely long video understanding. *arXiv preprint arXiv:2503.18478*, 2025.

## A  Overview of Appendix

## B  Discussions and Broader Impact

We discuss the limitation of EarthMind in Sec. 6, here we further discuss the limitation of EarthMind-Bench. Currently, EarthMind-Bench includes only paired optical (RGB) and SAR modalities. To enable a more comprehensive evaluation of LMMs' multi-sensor fusion capabilities, future extensions should incorporate additional sensing modalities such as multispectral, hyperspectral, and infrared imagery. Moreover, benchmarking diverse fusion scenarios involving more than two modalities would further reflect EO challenges.

**Broader Impact.** EarthMind offers a broader impact to the large multimodal model (LMM) community in both methodology and application. On the one hand, it demonstrates how vision-language models can be extended to handle multi-granular tasks via the proposed *Spatial Attention Prompting* mechanism, which reallocates LLM attention to task-relevant regions. This architectural design is not limited to Earth Observation (EO), but can be generalized to other domains requiring fine-grained spatial reasoning, such as medical imaging and autonomous driving.

On the other hand, for the EO community, EarthMind contributes both algorithmically and empirically: it introduces an effective multi-sensor fusion framework and proposes a scalable training pipeline, accompanied by a new benchmark and instruction-tuning dataset. These efforts can benefit various downstream applications in remote sensing, including disaster forecasting, infrastructure monitoring, and route navigation.

## C  Details of EarthMind

EarthMind is built upon the InternVL-2 framework [49]. In InternVL-2, each image is divided into multiple patches at pre-defined scales. Each patch is processed by the visual encoder and encoded into 256 tokens. For instance, an image with 4 patches (plus a global image token) yields $(4 + 1) \times 256$ visual tokens in total. For region-level perception, we adopt the region encoder from GPT4RoI [66], where each detected region is encoded into 256 tokens as well. To support multi-sensor input, non-optical imagery (e.g., SAR or multispectral data) is transformed into a synthetic video-like sequence by stacking pre-processed frames. This sequence is then concatenated before the input query, following the protocol of multi-frame vision-language models.

In addition, we extend the tokenizer by introducing a special "[SEG]" token for pixel-level grounding. The hidden state of the final LLM layer corresponding to the "[SEG]" token serves as the semantic prompt for mask generation. During inference, if the "[SEG]" token is not generated, we interpret it as an indication that the queried object is not present in the image.

## D  Details of EarthMind-Bench

**Overview.** We present a more detailed analysis of EarthMind-Bench in Fig. 5. Subfigures (a–c) illustrate the distribution of task categories, while (d) and (e) report the number of samples per task. Notably, our benchmark includes 438 referring expression segmentation samples paired with corresponding binary masks, highlighting its support for fine-grained pixel-level grounding. We further visualize the word clouds of questions (f) and answers (g), showing that EarthMind-Bench covers a wide variety of object types and semantics, enabling comprehensive evaluation across perception and reasoning tasks.
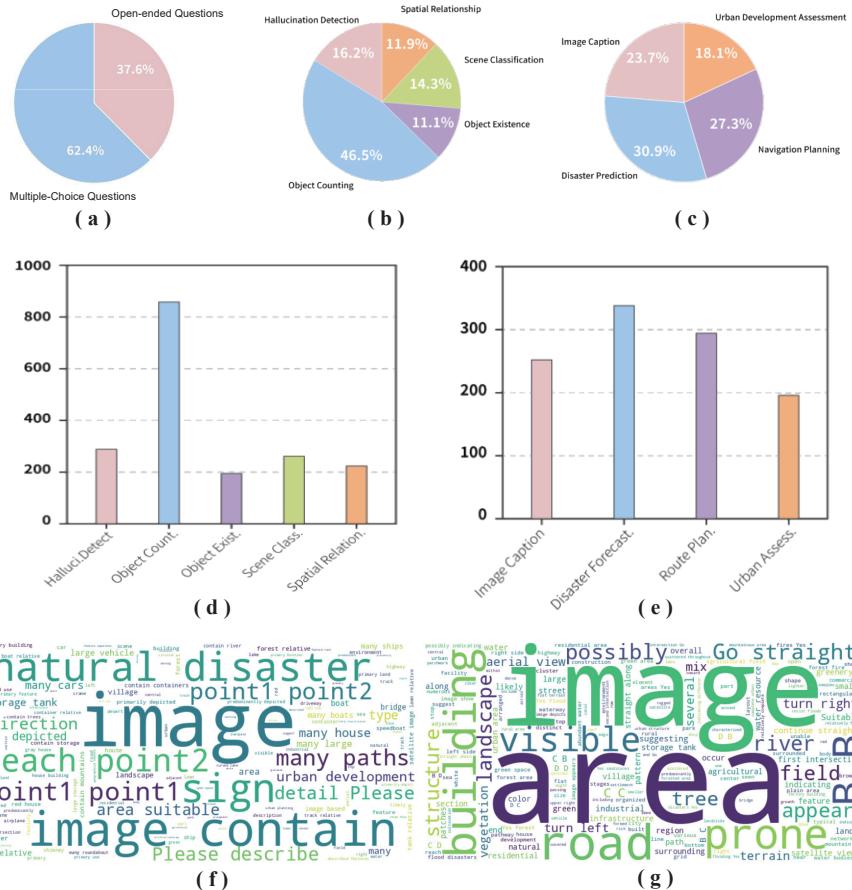
Figure 5: Detailed data statistics of the EarthMind-Bench.

**Data Annotations.** Among the 10 subtasks in EarthMind-Bench, annotations for Scene Classification and Referring Expression Segmentation are directly inherited from their original datasets. For the remaining 8 tasks, we rely on human annotations. We recruit 8 domain experts in geoscience, each assigned to annotate a specific task. After the initial annotation phase, all experts cross-validate each other's work and score the quality of samples. Only high-quality samples with consensus are retained in the final benchmark. To reduce the annotation burden, we first leverage GPT-4o to generate detailed image descriptions, which annotators can reference when constructing question-answer pairs. To ensure consistency and annotation fidelity, we also provide strict task-specific annotation guidelines for all subtasks, which are shown as follows:

- **1. Object Existence** Determine whether a specific object or class exists in the image (e.g., Is there a bridge in the image?"). *Guideline:* The object should be clearly visible and occupy a non-trivial area. If heavily occluded or ambiguous, mark as Not present."

- **2. Hallucination Detection** Detect whether the model falsely recognizes objects that do not exist. *Guideline:* Ask for fine-grained distinctions between semantically similar objects (e.g., "Is there a train or a bus in the image?"). Confirm that the mistaken object does not exist even under partial occlusion.

- **3. Object Counting** Count the number of objects from a given category (e.g., "How many buildings are visible in the image?"). *Guideline:* Only count objects that are visually distinguishable and not clustered into ambiguous shapes. Accept a small margin of error (±1) in complex scenes.

- **4. Spatial Relationship** Identify relative positions between objects (e.g., What is next to the waterbody?"). *Guideline:* Use cardinal or contextual spatial relations (e.g., to the left of," adjacent to," surrounded by"). Annotators should verify object co-existence and relative proximity.

- **5. Route Planning** Generate or select feasible navigation routes from a start to a target location. *Guideline:* Ensure the proposed path avoids obstacles (e.g., rivers, buildings), follows terrain constraints (e.g., avoid steep hills), and adheres to logical movement (e.g., roads preferred over fields).

- **6. Image Captioning** Generate descriptive sentences summarizing the content and layout of the image. *Guideline:* Captions should mention key objects, land types, spatial layout, and human-made structures if visible. Avoid hallucination and be concise yet informative.

- **7. Disaster Forecasting** Assess the likelihood of disaster based on visual evidence (e.g., "Is this area prone to flooding?"). *Guideline:* Use cues such as terrain (lowlands), proximity to water, lack of infrastructure, or signs of previous disaster impact. Avoid speculative answers.

- **8. Urban Development Assessment** Evaluate whether an area is suitable for urban development. *Guideline:* Consider factors such as flat terrain, absence of natural barriers, existing road access, and land cover type. Annotators should justify suitability with visual cues.

**Evaluation Metrics.** For multiple-choice tasks, we evaluate model performance using standard accuracy, measuring the percentage of predictions that exactly match the ground-truth option. For open-ended tasks, inspired by [67], we adopt a GPT-4-based evaluation protocol that assesses the alignment between the model-generated answers and human annotations. As illustrated in Fig. 6, we prompt GPT-4 to score the responses based on semantic similarity and correctness, following a structured rubric to ensure consistent evaluation across tasks.

---

**Evaluation Prompt For EarthMind–Bench open–ended Task**

**###Task Description:** You are required to evaluate a respondent's answer based on a provided question, some scoring points,
and the respondent's answer. You should provide two scores. The first is the accuracy score, which should range from 1 to 5. The
second is the relevance score, which should also range from 1 to 5. Below are the criteria for each scoring category.

**###Scoring Criteria:** Please rate the similarity between the **predicted caption** and the
**ground truth** based on the following criteria:1 - Completely unrelated (content is very different)
2 - Slightly related, but most descriptions do not match
3 - Somewhat similar, with a few common details but also clear differences
4 - Mostly matching, only a few minor differences
5 - Highly consistent, both descriptions describe the same content in detail
**##INSTRUCTION:**
Output Scores in JSON Format: Present the scores in JSON format as follows...

Figure 6: The illustration of open-ended task evaluation prompt.

---

# E   Details of Experimental Settings

We elaborate on the training and inference details of EarthMind. Specifically, we report the hyperparameters in the fine-tuning stage, as shown in Tab. 6.

# F   Training Dataset

EarthMind is trained on large-scale natural image datasets, including LLaVA-665K [6], 56K referring expression data [68, 69], and 214K grounding conversation generation samples [9]. These datasets cover image-level captioning, VQA, and text-driven segmentation. We also incorporate 724K region-level descriptions from the Osprey dataset [70] to improve region-level understanding capacity. Second, we introduce EO-specific multimodal data to adapt EarthMind to the remote sensing domain. This includes 1M VQA data from EarthGPT [17], 142K EO conversations from VRSBench [35], 31K

| Hyperparameter | Value |
|---|---|
| Overall batch size | 64 |
| Learning rate | 4e-5 |
| LR Scheduler | Cosine decay |
| DeepSpeed ZeRO Stage | ZeRO-2 |
| Optimizer | Adam |
| Warmup ratio | 0.3 |
| Epoch | 1 |
| Weight decay | 0 |
| Precision | bf16 |

Table 6: Hyperparameters of EarthMind.

region-level captions from DIOR-RSVG [54], and 21K referring segmentation samples from RRSIS-D [56] and RefSegRS [55]. Moreover, we involve 500k multi-spectral data, including BigEarthNet [57] and SoSAT-LCZ42 [57]. Third, we utilize our synthesized multi-sensor conversation corpus (20K RGB-SAR paired dialogues) and selectively retain examples from earlier stages to mitigate catastrophic forgetting.

**Paired Multi-Sensor Training Data Curation.** We construct a high-quality paired RGB–SAR training corpus from five publicly available datasets: OpenEarthMap-SAR [44], DFC2023 Track2 [45], WHU-OPT-SAR [46], MSAW [47], and MultiResSAR [48]. All RGB–SAR pairs are sampled from their training splits to avoid overlap with EarthMind-Bench. To scale data generation, we design an automatic synthetic pipeline based on GPT-4o. Since GPT-4o performs better on RGB imagery than SAR, we use RGB images and their classification labels as input prompts. The prompt is tailored for large vision-language models (VLMs), comprising the image and a list of detected object categories (e.g., buildings, roads, water bodies). The model is first asked to generate a comprehensive caption summarizing the scene and spatial relationships among objects. Then, it is instructed to synthesize 3–10 diverse question–answer pairs. These cover various reasoning types, such as: (1) Object Existence (e.g., "Is there a river in the image?"), (2) Counting (e.g., "How many buildings are visible?"), (3) Spatial Relation (e.g., "What is next to the forest?"), (4) Object Localization (e.g., "Where is the road located?"), and (5) Scene-Level Understanding (e.g., "Is this area suitable for agriculture?"). All outputs are returned in a structured JSON format for seamless integration into our training pipeline. In total, we curate 20k synthetic multi-sensor samples, which serve as paired RGB–SAR instruction data for adaptation of multi-sensor LMMs.

Table 7: Ablations on the non-rgb data processing.

| Method | Multispectral | | Method | SAR | | | | |
|---|---|---|---|---|---|---|---|---|
| | BigEarthNet | SoSAT-LCZ42 | | Small | Medium | Large | Single | Multiple |
| Three-band grouping | 70.4 | 58.3 | Zero-padding | 12.14 | 26.02 | 35.56 | 26.03 | 6.06 |
| Single-band grouping | 71.2 | 59.2 | Channel replication | 11.08 | 25.99 | 34.38 | 25.99 | 4.32 |

Table 8: Ablation study on Multi-granular joint training.

| Method | Image-level (Accuracy) | | | | Region-level | | Pixel-level (mIoU) | |
|---|---|---|---|---|---|---|---|---|
| | AID | UC | RSVQA | VRS-VQA | RSVG | VRS-VG | RRSIS-D | RefSegRS |
| Only trained on Image data | 97.0 | 95.1 | 73.8 | 77.8 | - | - | - | - |
| Only trained on Region data | - | - | - | - | 379.6 | 49.8 | - | - |
| Only trained on segmentation data | - | - | - | - | - | - | 77.6 | 59.3 |
| **EarthMind** | **97.2** | **95.0** | **74.0** | **78.9** | **428.2** | **55.6** | **82.2** | **62.6** |

# G  More Ablation Studies

**Ablation on Multi-Sensor Data Processing.** One of the key strengths of EarthMind lies in its ability to handle multi-sensor data beyond standard RGB imagery. To better understand the impact of different preprocessing strategies, we conduct a series of experiments focused on the handling

Table 9: Ablation study on joint multi-sensor training.

| Method | EarthMind-Bench | | |
|--------|:---:|:---:|:---:|
| | RGB | SAR | RGB+SAR |
| Only trained on RGB | 68.4 | 30.1 | 28.4 |
| Only trained on SAR | 45.6 | 59.8 | 22.3 |
| **trained on paired RGB-SAR** | **69.0** | **67.5** | **70.6** |

of SAR and multispectral (MS) data. For SAR inputs with fewer than three channels, we compare two strategies: (1) *Zero padding*, where the missing channels are filled with zeros; (2) *Channel replication*, where existing channels are duplicated to reach three channels. For multispectral data with more than three bands, we evaluate: (1) *Three-band grouping*, where every three consecutive bands are grouped to form one RGB-like frame; (2) *Single-band grouping*, where each band is treated as an individual frame, forming a multi-frame sequence. We evaluate each method under the same training settings and report classification accuracy on the test sets of BigEarthNet (for MS) and SAR Ship Detection (for SAR), as summarized in Tab. 7. Our results show that zero padding outperforms channel replication for SAR data, likely because copying channels introduces redundancy and potential noise. For MS data, single-band grouping slightly improves performance due to better spectral resolution, but incurs substantial computational overhead due to the increased number of tokens. Considering both effectiveness and efficiency, we adopt zero padding for SAR and three-band grouping for MS as our default configuration. In the future, we will consider involve token reduction techniques [71] to reduce overhead cost.

**Ablation on Joint Multi-Granular and Multi-Sensor Training.** EarthMind is designed to be jointly trained on both multi-granular and multi-sensor data. To validate the effectiveness of this unified training paradigm, we conduct ablation studies along two axes: (1) Multi-Granular Training. We compare two settings: (i) Joint training with image-level, region-level, and pixel-level data. (ii) Independent training for each granularity using the same total amount of data. (2) Multi-Sensor Training. We also compare: (i) Joint training with all available sensor modalities (e.g., RGB, SAR, MS). (ii) Independent training with each modality separately. As shown in Tab. 8 and Tab. 9, the multi-granular co-training strategy consistently outperforms independently trained counterparts, especially on pixel-level tasks. This suggests that high-level semantic supervision (e.g., image-level QA) can improve fine-grained understanding through shared representation learning. Similarly, multi-sensor co-training improves generalization across modalities. Notably, the performance under SAR-only evaluation is significantly enhanced by leveraging complementary information learned from RGB and MS data during joint training. This highlights EarthMind's ability to exploit cross-modal synergy in both perception and reasoning tasks.

**Ablation on Curriculum Training Strategy.** EarthMind adopts a curriculum learning strategy that gradually adapts the model from general vision-language data to remote sensing tasks. Specifically, the training proceeds in three stages: (1) pretraining on large-scale natural image VQA and captioning datasets, (2) domain adaptation using EO-specific RGB data, and (3) fine-tuning on multi-sensor data (e.g., SAR). To evaluate the effectiveness of this curriculum, we conduct ablation studies by removing or reordering training stages. As shown in Tab. 10 and Tab. 11, pretraining on general image-language data provides strong foundational capabilities, leading to significant performance improvements on EO tasks, particularly in pixel-level segmentation. Furthermore, RGB-domain training not only enhances performance on RGB inputs but also boosts multi-sensor fusion results, demonstrating its role as an effective bridge between general vision and SAR-specific domains.

Table 10: Ablation study on the first stage of curriculum training.

| Method | Image-level (Accuracy) | | | | Region-level | | Pixel-level (mIoU) | |
|--------|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | AID | UC | RSVQA | VRS-VQA | RSVG | VRS-VG | RRSIS-D | RefSegRS |
| w/o pretraining | 96.5 | 94.8 | 73.2 | 77.6 | 406.7 | 49.8 | 75.4 | 55.3 |
| **with pretraining** | **97.2** | **95.0** | **74.0** | **78.9** | **428.2** | **55.6** | **82.2** | **62.6** |

Table 11: Ablation study on the second stage of curriculum learning.

| Method | EarthMind-Bench | | |
| --- | --- | --- | --- |
| | **RGB** | **SAR** | **RGB+SAR** |
| w/o pretraining on RGB EO data | 67.5 | 64.3 | 68.9 |
| **with pretraining on RGB EO data** | **69.0** | **67.5** | **70.6** |

## H More Visualization Results
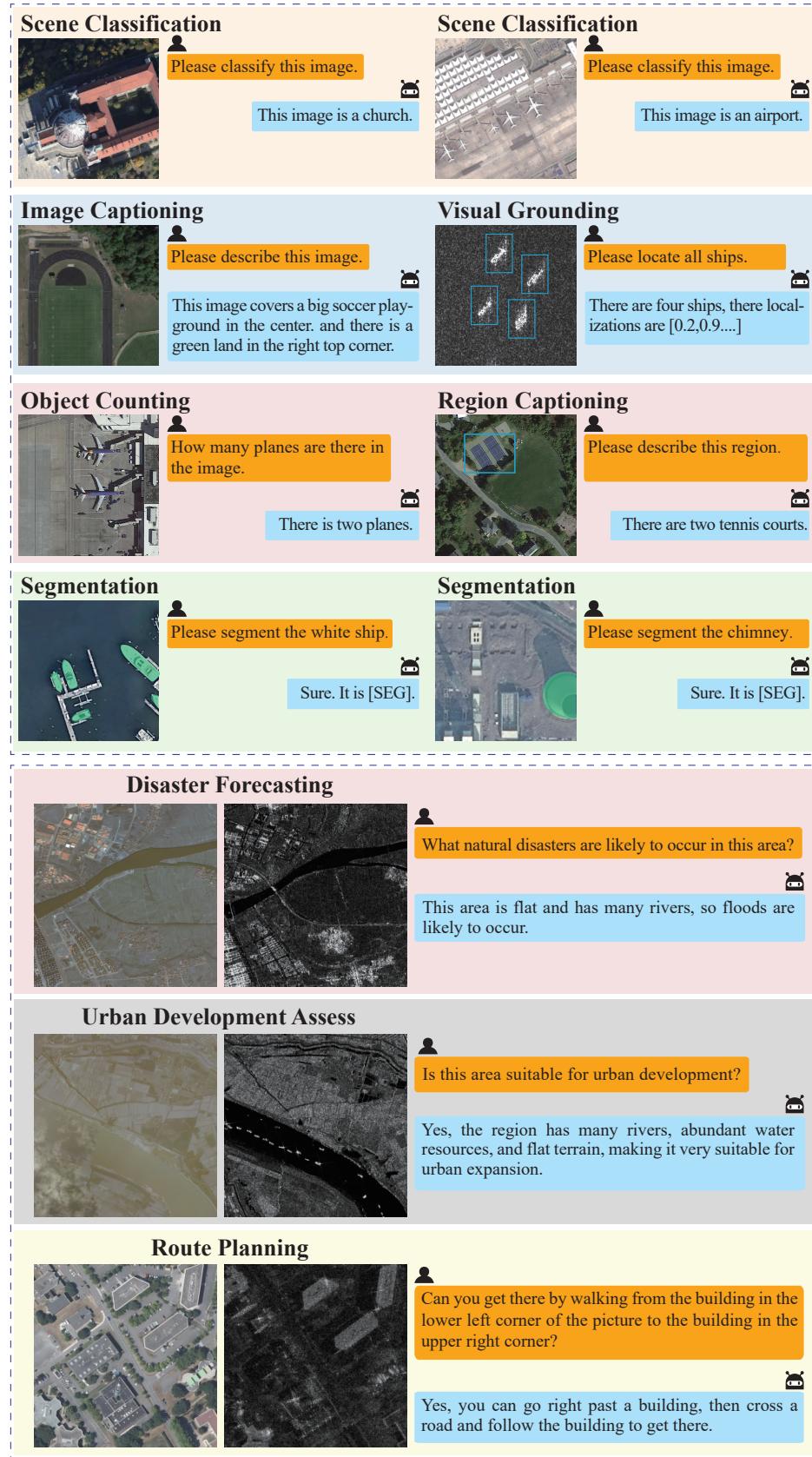
More visualization results of EarthMind can be seen in Fig. 7.

Figure 7: More visualization of EarthMind.