# EarthGPT: A Universal Multi-modal Large Language Model for Multi-sensor Image Comprehension in Remote Sensing Domain

Wei Zhang*, Miaoxin Cai*, *Student Member, IEEE,* Tong Zhang, *Student Member, IEEE,*
Yin Zhuang, *Member, IEEE,* and Xuerui Mao†

*Abstract*—**Multi-modal large language models (MLLMs) have demonstrated remarkable success in vision and visual-language tasks within the natural image domain. Owing to the significant diversities between the natural and remote sensing (RS) images, the development of MLLMs in the RS domain is still in the infant stage. To fill the gap, a pioneer MLLM named EarthGPT integrating various multi-sensor RS interpretation tasks uniformly is proposed in this paper for universal RS image comprehension. Firstly, a visual-enhanced perception mechanism is constructed to refine and incorporate coarse-scale semantic perception information and fine-scale detailed perception information. Secondly, a cross-modal mutual comprehension approach is proposed, aiming at enhancing the interplay between visual perception and language comprehension and deepening the comprehension of both visual and language content. Finally, a unified instruction tuning method for multi-sensor multi-task in the RS domain is proposed to unify a wide range of tasks including scene classification, image captioning, region-level captioning, visual question answering (VQA), visual grounding, object detection, etc. More importantly, a dataset named MMRS-1M featuring large-scale multi-sensor multi-modal RS instruction-following is constructed, comprising over 1M image-text pairs based on 34 existing diverse RS datasets and including multi-sensor images such as optical, synthetic aperture radar (SAR), and infrared. The MMRS-1M dataset addresses the drawback of MLLMs on RS expert knowledge and stimulates the development of MLLMs in the RS domain. Extensive experiments are conducted, demonstrating the EarthGPT's superior performance in various RS visual interpretation tasks compared with the other specialist models and MLLMs, proving the effectiveness of the proposed EarthGPT and offering a versatile paradigm for open-set reasoning tasks. Our code and dataset are available at *https://github.com/wivizhang/EarthGPT*.**

*Index Terms*—**Multi-modal large language model (MLLM), remote sensing (RS), instruction-following, multi-sensor.**

\* The authors contributed equally to this work.

† Corresponding author: Xuerui Mao.

Wei Zhang is with the Advanced Research Institute of Multidisciplinary Sciences, Beijing Institute of Technology, Beijing 100081, China, and also with the School of Mechatronical Engineering, Beijing Institute of Technology, Beijing 100081, China. (e-mail: w.w.zhanger@gmail.com).

Xuerui Mao is with the Advanced Research Institute of Multidisciplinary Sciences, Beijing Institute of Technology, Beijing 100081, China, and with the School of Mechatronical Engineering, Beijing Institute of Technology, Beijing 100081, China, and also with Yangtze Delta Region Academy of Beijing Institute of Technology, Jiaxing 314003, China. (e-mail: xmao@bit.edu.cn).

Miaoxin Cai, Tong Zhang, and Yin Zhuang are with the National Key Laboratory of Science and Technology on Space-Born Intelligent Information Processing, Beijing Institute of Technology, Beijing 100081, China. (e-mail: 3120220667@bit.edu.cn, bit_zhangtong@163.com, yzhuang@bit.edu.cn).

## I. INTRODUCTION

UNIFYING multi-task interpretation in the remote sensing (RS) domain is crucial in practical application as real-world scenarios often demand comprehensive analyses to make informed decisions. Although deep learning methods in RS have been successful in RS visual analysis tasks [1]–[3], current methods mainly follow a one-task-one-architecture paradigm which limits their ability to handle multi-sensor RS images, multiple tasks, and generalize to open-set reasoning.

Most recently, the development of multi-modal large language models (MLLMs) [4]–[7], has shown great success in the natural image domain, as they not only possess robust language interaction capabilities but also exhibit impressive multi-task reasoning skills in real-world scenarios [8], compared to smaller and domain-specific models tailored for particular tasks. These MLLMs showed effectiveness in tasks like detailed image description, question answering, spatial localization guided by language instructions, etc. The remarkable multi-modal reasoning ability allows MLLMs to generalize well in new situations and demonstrate powerful zero-shot capabilities across open-set tasks [9].

However, owing to the significant diversities between the natural and multi-sensor RS images such as the imaging conditions, environments, scales, and object viewpoints, there are significant challenges for the application of MLLMs in the RS domain. Although rare MLLM works such as RSGPT [10] and GeoChat [11] have been proposed to explore the MLLMs for RS and to solve multiple tasks, they still have limitations. Notably, RSGPT has to tune the model for each task individually, leading to the lack of generalizable ability in the open-set domain. RSGPT also struggles with tasks like classification, detection, and visual grounding. In addition, GeoChat is trained on optical RS images and lacks the multi-sensor datasets to realize synthetic aperture radar (SAR), and infrared modalities comprehension. It is clear that the study of MLLM in the RS domain is still in its infancy. To address these problems thoroughly, this paper aims to unify a wide range of RS tasks and concentrates on constructing a large-scale multi-modal dataset containing multi-sensor RS images based on diverse existing RS datasets.

To leverage the robust generalization and emergent ability of large language models (LLMs), a versatile MLLM called EarthGPT is proposed for multi-sensor RS image comprehension, unifying various RS interpretation tasks effectively.

Fig. 1. EarthGPT is a pioneering model designed to seamlessly unify multi-sensor and diverse RS intelligent interpretation tasks in a unified framework, guided by user language instructions, and is versatile at performing visual-language dialogues across optical, SAR, and infrared images. EarthGPT's capabilities extend to a wide range of tasks including scene classification, image description, visual question answering, target description, visual localization, and object detection.

As illustrated in Fig. 1, EarthGPT can serve as an intelligent assistant capable of efficiently handling a wide range of RS tasks through language interaction. In EarthGPT, three key techniques are proposed. Firstly, a visual-enhanced perception mechanism is constructed to mix various visual backbones. Specifically, the mixed backbones refine coarse-scale semantic perception information and fine-scale detailed perception information, thereby enhancing visual comprehension. Secondly, a cross-modal mutual comprehension approach is proposed. By directly concatenating the visual features with language features to generate multi-modal input for LLM, then unfreezing the self-attention and RMSNorm layer of transformer blocks to train on natural common datasets, visual-language alignment is implemented, and the interplay comprehension between visual and language content is deepened. Finally, a unified instruction tuning method for multi-sensor and multi-

task in the RS domain is proposed, by continuing fine-tuning LLM using the bias tuning strategy, MLLM is endowed with the capability of unifying a wide range of tasks including scene classification, image captioning, region-level captioning, visual question answering (VQA), visual grounding, horizontal bounding box (HBB) object detection, and oriented bounding box (OBB) object detection. More importantly, a large-scale dataset called MMRS-1M which is a multi-sensor multi-modal RS instruction-following dataset, is created. MMRS-1M comprises over 1M image-text pairs transformed from 34 existing diverse RS datasets and includes multi-sensor images such as optical, SAR, and infrared. This dataset is tailored to the unique visual modalities and geographical characteristics of the RS domain, effectively mitigating the lack of RS domain expertise of MLLMs. Furthermore, MMRS-1M serves as a catalyst for the advancement of MLLMs in the RS domain.

Extensive experiments are conducted in multiple RS datasets. It shows that the proposed EarthGPT surpasses the state-of-the-art (SOTA) specialist models and MLLMs in most RS tasks. In particular, for the image captioning, VQA, and visual grounding task in the supervised setting, EarthGPT shows a notable improvement in the NWPU-Captions [12], CRSVQA [13], and DIOR-RSVG [14] datasets compared with other specialist models. We also assess generalization ability in the open-set domain for the proposed EarthGPT. For the zero-shot scene classification task, EarthGPT achieves 77.37% and 74.72% accuracy on the CLRS [15] and NaSC-TG2 [16] datasets, respectively, far exceeding other MLLMs. On the MAR20 [17] detection dataset, EarthGPT achieves AP@40% metrics of 90.47% accuracy, outperforming other MLLMs and open-set detection models. In conclusion, experimental results demonstrate EarthGPT's superior performance in a wide range of RS multi-sensor image comprehension tasks, and robust generalization capability in open-set reasoning tasks.

In summary, the contributions of this paper are as follows.

- To our best knowledge, the largest multi-modal multi-sensor RS instruction-following dataset named MMRS-1M is constructed, consisting of over 1M image-text pairs that include optical, SAR, and infrared RS images. This dataset effectively mitigates the lack of RS domain expertise in MLLMs.
- A pioneer MLLM called EarthGPT is proposed for multi-sensor RS image comprehension, and is capable in a wide range of visual-language RS tasks uniformly. It consists of three techniques. The first is a visual-enhanced perception mechanism to refine coarse-scale and fine-scale visual perception information. The second is a cross-modal mutual comprehension approach that realizes the mutual comprehension of visual and language content. The last is a unified instruction tuning method designed for multi-sensor RS image comprehension and a wide range of RS visual tasks.
- Extensive experiments demonstrate the EarthGPT's superior performance in multi-sensor RS visual interpretation tasks compared with the other specialist models and MLLMs in scene classification, image captioning, region-level captioning, VQA, visual grounding, and object detection. EarthGPT therefore represents a significant advance of MLLMs in the field of RS and contributes a versatile paradigm for RS visual-language mutual comprehension as well as the open-set reasoning ability.

## II. RELATED WORK

### A. MLLMs

The continuous emergence of LLMs has fueled rapid advancements in natural language processing, showcasing extraordinary capabilities in language modeling across diverse contexts. Pre-training on massive text corpora, fine-tuning on specific domains, and the continual expansion of model parameters have enabled LLMs to achieve new milestones in various benchmark tests. Notably, OpenAI's ChatGPT [18] addresses users' needs across different scenarios through human-machine dialogues and text generations. The LLaMA-guided fine-tuning models [19]–[21] have gained favor among LLM researchers. In addition, GPT-4-LLM [22] demonstrates significant improvements when provided with high-quality instruction-following datasets. As a fundamental component of general intelligence, multi-modal perception is a crucial step toward achieving universal models. Researchers [4]–[7] are currently devoted to incorporating multimodal data beyond text into LLM to support specific tasks across various modalities. Models like VisualGPT [4], BLIP [5], Flamingo [7], and Kosmos-1 [23] showcase strong multimodal reasoning potential after aligning LLMs with image modalities. To achieve modality alignment, many works, such as MiniGPT-4 [9], LLAMA-Adapter V2 [24], mPLUG-Owl [25], employed projection layers, zero-shot attention mechanism, and intermediate networks to fuse LLaMA and visual modality. With the continuous progress of MLLMs, additional modalities such as audio, video [26], [27], 3D point clouds [28]–[30] are being integrated and aligned. The extraordinary potential of MLLMs is also evident in fields such as LIDAR [31] and robotics [32].

Nevertheless, existing MLLMs are conventionally tailored for the integration of visual and textual elements in natural scenes, missing the capacity to capture the distinctive contextual complexities of the RS domain. Those MLLMs's adaptability to the complex characteristics of RS data is limited, hindering their capability to effectively realize downstream RS tasks reasoning. To address this gap and create an open-set assistant for RS, EarthGPT is proposed to seamlessly integrate multiple RS tasks and visual modalities from multi-sensors.

### B. MLLMs for Remote Sensing

Previous RS large models primarily employ self-supervised methods [33], [34] and have progressed in RS visual tasks. However, these models are pre-trained based on visual modalities, lacking alignment with language and struggling with narrow application scenarios, hindering multi-modal understanding and reasoning. Recently, there has been an emergence of MLLMs for RS. For instance, Remoteclip [35] uses contrastive learning pre-training with image-text pairs from existing datasets, demonstrating strong zero-shot classification and image-text retrieval abilities. However, Remoteclip faces challenges in tasks like image captioning, region-level vision grounding, and visual language response due to constraints in the training pattern. RSGPT [10] fine-tunes Instruct-BLIP [36] on a high-quality image-text pair dataset, showing good image-text caption and VQA capabilities but struggling with tasks like classification, detection, and visual grounding. GeoChat [11] has been proposed to explore RS MLLMs and to solve multiple tasks. Nevertheless, GeoChat and the above models are trained on optical RS images and lack the generality for multi-sensor visual modalities like SAR, infrared, etc. To achieve more universal multi-modal reasoning in the RS domain and address the limitations of existing multi-modal RS models [12], [13], [37]–[40] in open-set dialogue, unifying multiple RS tasks, and multi-sensor image comprehension, this paper focuses on a unified MLLM. The proposed EarthGPT develops MLLM from the natural domain to the RS domain

through instruction tuning. Furthermore, to significantly enhance visual-language alignment and fully adapt to the unique characteristics of the RS domain, a visual-enhanced perception mechanism, a cross-modal mutual comprehension mechanism and a unified instruction tuning method for multi-sensor RS image comprehension have been incorporated into EarthGPT. Simultaneously, the constructed MMRS-1M dataset contains various visual modalities.

### C. Remote Sensing Datasets

RS datasets are the core and foundation of RS intelligent interpretation models. Currently, RS datasets mainly focus on classification, detection, segmentation, image captioning, and VQA tasks. For classification, commonly used datasets include AID [41], EuroSAT [42], NWPU-RESISC45 [43], UCMerced LandUse [44], and WHU-RS19 [45]. For detection, optical datasets like DIOR [46], DOTA [47], FAIR1M [48], HRRSD [49], NWPUVHR10 [50], and SAR datasets like AIR-SARShip-2.0 [51], HRSID [52], and SSDD [53] are used to train expert models. Segmentation datasets typically include Vaihingen, Potsdam, iSAID [54], LoveDA [55], etc. Image captioning datasets include Sydney-Captions [56], RSICD [37], NWPU-Captions [12], RSITMD [57], and UCM-Captions [56] are used for generating image descriptions or image-text retrieval. VQA datasets mainly consist of Floodnet [40], RSVQA-LR [39], RSVQA-HR [39], RSIVQA [38], and CRSVQA [13]. The datasets mentioned above primarily emphasize a single visual modality within an individual task, leading to models with limited generalization capabilities and restricted to specific tasks. MLLMs represent a potential solution to tackle the challenge. However, in the current landscape of the RS domain, there is a scarcity of high-quality image-text instruction-following data. This scarcity poses a significant obstacle when it comes to training multi-modal intelligent dialogue assistants that can achieve a high level of alignment between visual and language, unify various tasks into a framework, and integrate different multi-sensor visual modalities effectively.

To address this challenge, we have constructed a dataset called MMRS-1M containing more than 1M image-text pairs, covering tasks such as classification, detection, image captioning, VQA, visual grounding, etc. MMRS-1M encompasses three visual modalities such as optical, infrared, and SAR. MMRS-1M is designed to promote the continuous development of MLLMs in the RS field.

### III. METHODOLOGY

This section introduces the proposed EarthGPT framework. We first overview the overall model architecture (Section III-A). Then, three fundamental techniques in EarthGPT are elaborated involving the visual-enhanced perception mechanism (Section III-B), the cross-modal mutual comprehension approach (Section III-C), and the unified instruction tuning method for RS (Section III-D).

### A. Overview

Our goal is to break down the differences between natural images and RS images that lead to the obstacles for the application of MLLMs in the RS domain and propose a universal MLLM specifically customized for RS. To this end, a unified MLLM called EarthGPT is proposed for RS imagery comprehension. The overall framework of EarthGPT is shown in Fig. 2 (a). Note that three pivotal techniques are developed. Firstly, to tackle the challenge that RS images commonly contain various disturbances that compromise clarity, the visual-enhanced perception mechanism is proposed, as shown in step 1. Specifically, ViT and CNN backbones are designed to refine global and detailed visual features through multi-layer and multi-scale visual perception. Then, twofold visual tokens are concatenated and projected to ensure dimension alignment with language instruction tokens. The hybrid visual tokens contain both neighboring dependencies and long-range visual interactions, offering subtle local insights and extensive regional correlations. Subsequently, the cross-modal mutual comprehension approach is given in step 2. This part is to equip the LLM with fundamental visual comprehension capabilities and achieve visual-language alignment. Particularly, the visual and language tokens are merged as the multi-modal input for the unfrozen LLM tuning to enhance the interaction between the diverse visual and linguistic contexts. Finally, to address the limitations of the existing RS models are specialist ones, unable to handle multi-source and multi-task in one structure. A unified instruction tuning for RS method is developed in step 3. Concretely, the bias tuning of linear layers in LLM is performed on the newly constructed MMRS dataset, achieving the universal visual interpretation in the RS domain. We will introduce each technique as follows in detail.

### B. Visual-enhanced Perception

The core of the visual-enhanced perception mechanism is to utilize the strengths of diverse image encoding models to enhance and refine multi-granularity critical visual information while simultaneously mitigating various disturbances in RS images. The visual-enhanced perception mechanism consists of two modules as shown below.

**Multi-Layer Visual Perception.** We adopt the vision transformer (ViT) [58] as the image encoder to extract multi-layer features. In ViT, images are represented as sequences, allowing models to learn image structures independently. First, multi-layer visual intermediate feature maps $\{V_a^i\}_{i=1}^n$ are extracted from different encoder layers to capture long-range contextual information, where $n$ represents the number of layers. The original images $I$ are transformed into visual representations $V_a^1$. Then, all the extracted visual features are concatenated along the channel dimension. Such visual feature fusion from different layers is beneficial to capture subtle differences in RS images and provides a more comprehensive understanding of images. The entire process can be formulated as

$$V_a^1 = \text{ViTBlock}_1(I), \tag{1}$$

$$V_a^i = \text{ViTBlock}_i(V_a^{i-1}), \; i = 2, ..., n, \tag{2}$$

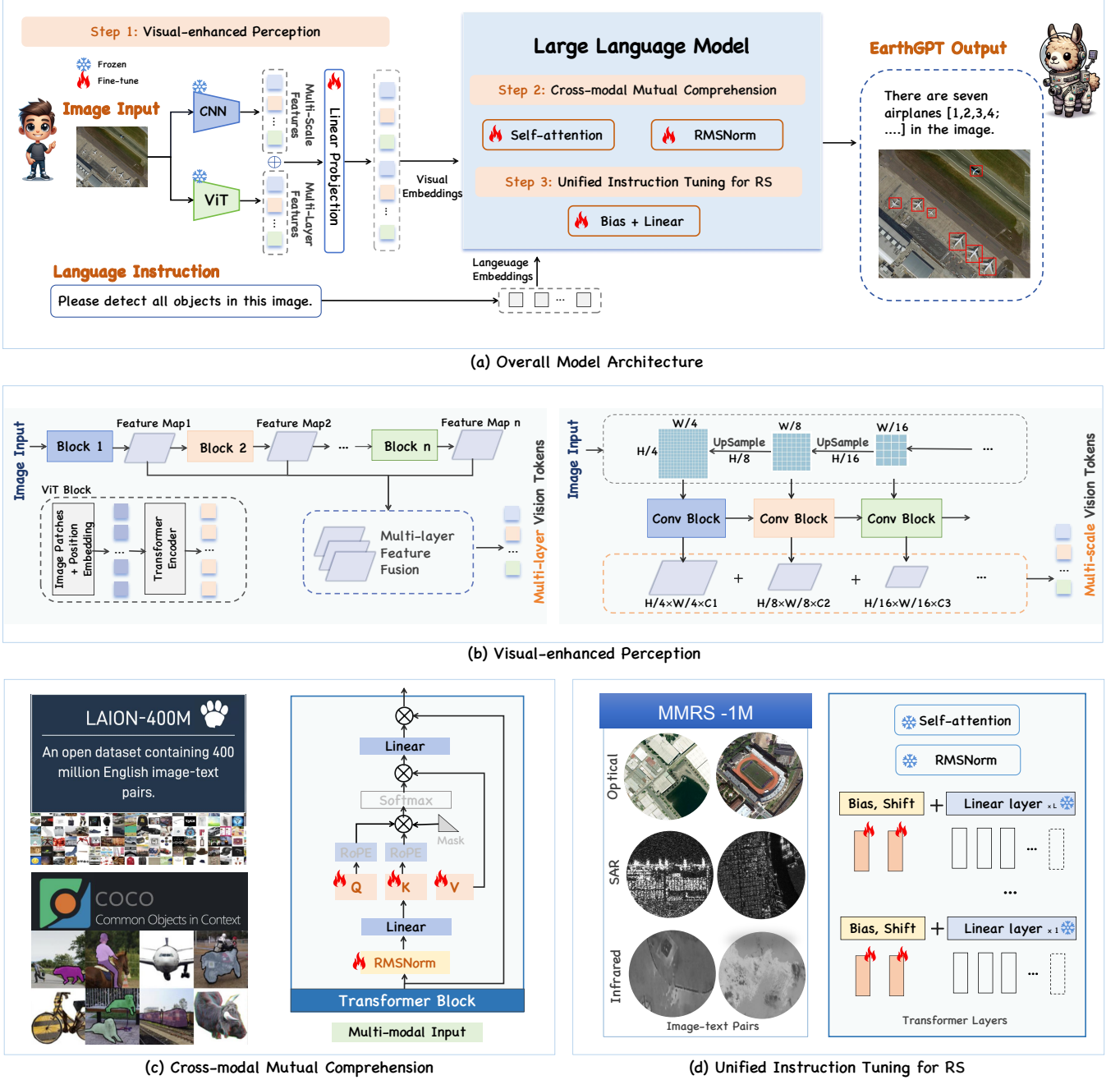$$V_a = \text{Concat}\left[V_a^1, V_a^2, ... V_a^n\right]. \tag{3}$$

Fig. 2. (a) Overall model architecture of EarthGPT. (b) Illustration of the visual-enhanced perception mechanism. (c) Illustration of the cross-modal mutual comprehension approach. (d) Illustration of the unified instruction tuning method for RS.

In this way, the extracted features both involve spatial-aware information from early layers and semantic indicative from later layers. The multi-layer visual perception is shown at the left of Fig. 2 (b).

**Multi-scale Visual Perception.** To integrate multi-scale local details into visual representations, the CNN [59] backbone is designed as the image encoder. Compared to ViT, CNN excels in extracting localized features, such as edges and textures, through its inherent spatial hierarchies and local receptive fields. The extracted multi-scale visual features are denoted as $\{V_b^i\}_{i=1}^m$, where $m$ represents the image scale numbers. As illustrated at the right of Fig. 2 (b), the input of CNN includes

multiple spatial resolutions of $\frac{H}{4} \times \frac{W}{4}$, $\frac{H}{8} \times \frac{W}{8}$, $\frac{H}{16} \times \frac{W}{16}$, etc, where $H \times W$ is the original input image resolution. The convolution blocks convert the inputs into token embeddings $V_b^1 \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C_1}$, $V_b^2 \in \mathbb{R}^{\frac{H}{8} \times \frac{W}{8} \times C_2}$, $V_b^3 \in \mathbb{R}^{\frac{H}{16} \times \frac{W}{16} \times C_3}$, etc. Meanwhile, all scale features are transformed into the same channel dimension and concatenated together channel-wisely. The feature extraction process can be written as

$$V_b^1 = \text{ConvBlock}_1(I_1), \tag{4}$$

$$V_b^i = \text{ConvBlock}_i(V_b^{i-1}), \ i = 2, ..., m, \tag{5}$$

$$V_b = \text{Concat}\left[V_b^1, V_b^2, ... V_b^m\right], \tag{6}$$

where the $I_1$ represents the first scale input. The multi-scale fused visual features encompass information of multi-granularity, enabling the capture of both broad semantics and intricate information.

After extracting visual features using ViT and CNN, multi-layer and multi-scale features are concatenated channel-wisely. Then, a learnable projection layer is used for dimension alignment with language tokens. The new visual tokens after dimension alignment are denoted as $V_p$. This process is formulated as

$$V_p = \text{Projection}(\text{Concat}(V_a, V_b)). \qquad (7)$$

The integration of enhanced visual perception information improves the accuracy and efficiency of the image interpretation.

### C. Cross-modal Mutual Comprehension

The key point of the cross-modal mutual comprehension approach lies in visual and language perception information fusion and delivering the multi-modal input into the LLM for alignment and interaction training.

Visual tokens are obtained based on the aforementioned visual-enhanced perception mechanism. Meanwhile, following most natural language processing (NLP) models, the language tokenizer is employed to segment language instructions into discrete token embeddings. Subsequently, overall visual embeddings $V_p$ are concatenated with the language instruction embeddings $L_p$ to create LLM input tokens $\mathcal{X}$ as

$$\mathcal{X} = \text{Concat}[\overbrace{V_p^1, V_p^2, \ ... \ , V_p^{N_v}}^{\text{visual tokens } V_p}, \underbrace{L_p^1, L_p^2, \ ... \ , L_p^{N_l}}_{\text{language tokens } L_p}], \qquad (8)$$

where $N_v$ represents the token length of visual features, $N_l$ represents the token length of language features, $(V_p^1, V_p^2, \ ... \ , V_p^{N_v})$ are the mixed visual tokens from $V_p$, $(L_p^1, L_p^2, \ ... \ , L_p^{N_l})$ are the language instruction tokens from $L_p$. Then, multi-modal input $\mathcal{X}$ is fed into LLM for fusion and integration.

In current MLLMs [4]–[7], [9], vision-language understanding is achieved through training on a frozen LLM, designed to avoid expensive full-parameter fine-tuning. However, freezing all the LLM's weights significantly limits its potential for comprehensive cross-modal learning. To cope with the challenge of the knowledge gap from different modalities and realize multi-modal mutual comprehension, an unfrozen vision-language alignment strategy is designed. Specifically, the self-attention and normalization layers are unfrozen for training on the most common domain data (LAION-400M [60], COCO Caption [61]). This strategy aids the LLM in thoroughly understanding multi-modal representations and simultaneously preserves the inherent language response capabilities of LLM through the partly frozen modules. In our method, LLaMA-2 [62] with powerful language understanding capability is adopted as the initial LLM. LLaMA-2 is composed of $L$ Transformer blocks. The Transformer block consists of self-attention, RMSNorm, and FFN layers, whose key module self-attention head is composed of key $\mathbf{K}$, query $\mathbf{Q}$, and the value $\mathbf{V}$, as shown in Fig. 2 (c). The $\mathbf{K}$, $\mathbf{Q}$, and $\mathbf{V}$ are implemented via linear layers. The detailed implementation is as follows

$$\mathbf{Q}(x) = \mathbf{W}_q \ x + \mathbf{b}_q, \qquad (9)$$

$$\mathbf{K}(x) = \mathbf{W}_k \ x + \mathbf{b}_k, \qquad (10)$$

$$\mathbf{V}(x) = \mathbf{W}_v \ x + \mathbf{b}_v. \qquad (11)$$

The parameters $\mathbf{W}_q$, $\mathbf{W}_k$, $\mathbf{W}_v$, $\mathbf{b}_q$, $\mathbf{b}_k$, and $\mathbf{b}_v$ are updated during the training. The RMSNorm component of the transformer is also unfrozen, and the scale operation $\gamma$ is set as a trainable parameter. The RMSNorm can be expressed as

$$y = \frac{x}{\sqrt{\text{Mean}(x^2) + \varepsilon}} * \gamma, \qquad (12)$$

where

$$\text{Mean}(x^2) = \frac{1}{N} \sum_{i=1}^{N} x_i^2. \qquad (13)$$

During the training process, the adopted cross-entropy loss function is described as $\mathcal{L}$, the multi-modal input sequence length is denoted as $N_{vl}$ and the parameters of EarthGPT are represented as $\theta$, where $w_i$ denotes the $i$-th word. The cross-entropy loss function $\mathcal{L}$ can be formulated as

$$\mathcal{L} = - \sum_{i=1}^{N_{vl}} \log P(w_i|(w_1, w_2...w_{i-1}; \theta). \qquad (14)$$

After visual-language alignment, the language-only LLM is converted into an MLLM, which can generate the visual interpretation response based on the integrated multi-modal information.

### D. Unified Instruction Tuning for RS

After the cross-modal mutual comprehension training, the proposed EarthGPT has basic multi-modal reasoning and dialogue capabilities in the natural domain but struggles to follow instructions and perform inference for RS downstream tasks in the aforementioned stage. To enhance EarthGPT's ability to follow instructions for various downstream tasks and expand its applicability from the natural domain to the RS one, we develop an extensive and unified instruction-following dataset called MMRS-1M. In the MMRS-1M dataset, we have standardized all downstream tasks into the format of VQA instructions. By utilizing MMRS-1M for tuning, the model excels in accurately understanding and executing various instructions in the RS domain.

Technically, to preserve the visual captioning capability obtained in the previous stage and simultaneously enhance compliance with task instructions, we freeze all the weights from the cross-modal mutual comprehension phase and introduce new learnable parameters into the LLaMA-2 model. Specifically, inspired by LLaMA-Adapter V2 [24], the linear layer has been modified by introducing two learnable parameters including a bias $\beta$ and a shift $\alpha$ into linear layers, as shown in Fig. 2 (d). Given a linear layer $y(x) = \mathbf{W}x$, it can be transformed by incorporating the bias factor and shift factor. The process can be expressed as

$$y(x) = \alpha \cdot (\mathbf{W}x + \beta), \qquad (15)$$

where

$$\alpha \sim \mathcal{N}(\mu, \sigma^2), \ \beta = \text{Init}(0). \tag{16}$$

The bias $\beta$ is initialized with zeros, and the shift $\alpha$ is initialized with a random Gaussian, ensuring training stability and effectiveness. Analogously, the loss function in this stage is given as Eq. (14).

Leveraging the advanced reasoning abilities of LLM and our abundant MMRS-1M dataset, EarthGPT is equipped with versatile skills for multi-sensor visual comprehension in the RS domain guided by the language instructions and shows the potential for real-world practical applications.

## IV. DATASET CONSTRUCTION

Currently, lacking domain-specific datasets for RS hampers the effective application of MLLMs in the intelligent interpretation of geographic information and open-set dialogues.

TABLE I
DETAILS ON THE TRAINING SAMPLES USED FOR THE MMRS-1M.

| Task | Data | Size | Type |
|---|---|---|---|
| Image Captioning | RSICD | 24,333 | optical |
| | UCM-Captions | 9,986 | optical |
| | RSITMD | 20,096 | optical |
| | Sydney-Captions | 2,837 | optical |
| | NWPU-Captions | 141,631 | optical |
| VQA | FloodNet | 4,511 | optical |
| | RSVQA_LR | 67,228 | optical |
| | RSIVQA | 68,625 | optical |
| | CRSVQA | 900 | optical |
| Classification | NWPU-RESISC45 | 6,300 | optical |
| | EuroSAT | 5,400 | optical |
| | UCMerced-Landuse | 420 | optical |
| | WHU-RS19 | 196 | optical |
| | RSSCN7 | 560 | optical |
| | DSCR | 11,951 | optical |
| | FGSCR-42 | 3,878 | optical |
| Detection | DOTA | 163,486 | optical |
| | DIOR | 58,200 | optical |
| | FAR1M | 40,466 | optical |
| | NWPUVHR10 | 3,190 | optical |
| | HRRSD | 23,044 | optical |
| | RSOD | 747 | optical |
| | UCAS-AOD | 1,510 | optical |
| | VisDrone | 186,810 | optical |
| | AIR-SARShip-2.0 | 1,433 | SAR |
| | SSDD | 1,856 | SAR |
| | HRISD | 7,265 | SAR |
| | HIT-UAV | 10,608 | infrared |
| | Sea-shipping | 16,023 | infrared |
| | Infrared-security | 17,234 | infrared |
| | Aerial-mancar | 33,051 | infrared |
| | Double-light-vehicle | 7,922 | infrared |
| | Oceanic ship | 2,505 | infrared |
| Visual Grounding | DIOR-RSVG | 30,820 | optical |
| Region-level Captioning | DIOR-RSVG | 30,820 | optical |

Therefore, a diverse and comprehensive instruction-following dataset for image-text conversation in RS imagery is indispensable. To address this limitation, a large quantity of existing RS datasets are carefully cleaned and transformed, and a new instruction-following dataset called MMRS-1M is created covering five tasks (e.g., classification, detection, image caption, VQA, and visual grounding) and three visual modalities (e.g., optical, SAR, and infrared) is created. EarthGPT is fine-tuned on MMRS-1M to align visual and language modalities and achieve excellent coarse-grained conversation and fine-grained localization capabilities. The dataset construction process is detailed as follows.

### A. Coarse-grained Conversation Scenarios

To endow EarthGPT with image-level coarse-grained question-answering capabilities, ten classification datasets, five image captioning datasets, and four VQA datasets are collected for instruction fine-tuning.

**Category-to-instruction.** The classification datasets include the scene classification datasets AID [41], EuroSAT [42], NWPU-RESISC45 [43], UCMerced-LandUse [44], WHU-RS19 [45], RSSCN7 [63], as well as the ship classification datasets FGSCR-42 [64] and DSCR [65]. For the format conversion of the classification datasets, category-to-instruction is adopted. Specifically, the category of each image is first extracted, and the template "What is the category of this RS image? Answering the question using a single word or phrase. Reference categories include category 1,..., category n" is used to inquire the image. The GPT model provides a response that includes the name of the category to which the image belongs.

**Caption-to-instruction.** Firstly, for image caption datasets, the existing datasets including Syndney-Captions [56], RSICD [37], NWPU-Captions [12], RSITMD [57], and UCM-Captions [56] are cleaned, removing duplicate captions for the same image. Then the instruction "Please provide a one-sentence caption for the provided RS image in detail." is provided and the EarthGPT model is tasked with describing the image. The number of captions per image determines the number of conversation rounds. Putting multiple captions for each image into a multi-turn dialogue can significantly reduce computation costs while ensuring diverse descriptions of images without compromising information leakage. This conversion process is denoted as caption-to-instruction.

**VQA-to-instruction.** For VQA datasets including Floodnet [40], RSVQA-LR [39], RSIVQA [38], and CRSVQA [13], the instruction "Answering the question using a single word or phrase" is added after the original question to control the output format of the answer. The number of question-answer pairs per image determines the number of dialogue rounds in the conversation. This multi-round dialogue allows EarthGPT to extract and interpret sufficient information from the images.

### B. Fine-grained Conversation Scenarios

In order to enable the RS intelligent assistant to perform target region localization and object detection under various imaging conditions, a collection of optical, SAR, and infrared object detection datasets has been assembled. The

## Classification Datasets ★ Category-to-Instruction

DSCR · EuroSAT · FGSCR-42 · NWPU-RESISC45 · RSSCN7 · UCMerced-LandUse · WHU-RS19

Aircraft_carrier · PermanentCrop · Container ship · Church · Parking · Intersection · Bridge

Q: What is the category of this remote sensing image? Answer the question using a single word or phrase.
A: Aircraft_carrier

## VQA Datasets ★ VQA-to-Instruction

FloodNet · CRSVQA · RSIVQA · RSVQA-LR

Q: How many non flooded buildings can be seen in this image? A: 6

Q: What scene of this image shown? A: storage tanks

Q: How many basketball courts are there in this picture? A: 9

Q: Does this picture contain ship? A: no

Q: How many non flooded buildings can be seen in this image? Answer the question using a single word or phrase.
A: 6

## Image Caption Datasets ★ Caption-to-Instruction

NWPU-Captions · RSITMD · RSICD · Syndney-Captions · UCM-Captions

The church with four blue circular roof is surrounded by some g-ray buildings and trees.

Some planes are in front of the terminal building while some others are embraced by greenbelt.

Many buildings and some green trees are around a playground.

Lots of houses with red and orange roofs arranged in lines.

An intersection with many cars parked at the corners.

Q: Please describe this remote sensing image.
A: The church with four blue circular roof is surrounded by some gray buildings and trees.

## Visual Grounding Datasets ★ VG-to-Instruction

DIOR-RSVG

The ground track field is on the right of the small orange baseball field.

Many airplanes are on the right of a small vehicle.

An airport is on the left of the green tiny golf field.

A windmill on the left.

Q: Please provide a short caption for this region in this remote sensing image: [0.58,0.09,0.97,0.96].
A: The ground track field is on the right of the small orange baseball field.

Q: Please provide the horizontal bounding box coordinate of the region which is described as: The ground track field is on the right of the small orange baseball field
A: [0.58,0.09,0.97,0.96]

## Detection Datasets ★ Detection-to-Instruction

### Optical

DOTA HBB OBB · DIOR HBB OBB · VisDrone HBB · NWPUVHR10 HBB · HRRSD HBB · RSOD HBB · UCAS OBB · FAR1M OBB

### SAR

SSDD HBB OBB · HRISD HBB OBB · AIR-SARShip-2.0 HBB

### Infrared

Aerial-mancar HBB · Security HBB · Double-light-vehicle HBB · Sea-shipping HBB · HIT-UAV HBB OBB

#### HBB

Q: Detect all objects shown in the remote sensing image and describe using horizontal bounding boxes.

A: There are four objects shown in the image, including two large-vehicles [0.14,0.55,0.47,0.69;0.52,0.36,0.66,0.73], and two small-vehicles [0.82,0.82,0.89,0.91;0.81,0.12,0.89,0.17].

Q: Detect all large vehicles shown in the image and describe using horizontal bounding boxes.
A: There are two large vehicles [0.14,0.55,0.47,0.69; 0.52,0.36,0.66,0.73] in the image.
Q: Detect all small vehicles shown in the image and describe using horizontal bounding boxes.
A: There are two small vehicles [0.82,0.82,0.89,0.91; 0.81,0.12,0.89,0.17] in the image.

#### OBB

Q: Detect all objects shown in the remote sensing image and describe using oriented bounding boxes.

A: There are four objects shown in the image, including two large-vehicles [0.14,0.65,0.46,0.55,0.47,0.60, 0.15,0.69;0.52,0.38,0.57,0.36,0.66,0.71,0.62,0.73], and two small-vehicles [0.86,0.82,0.89,0.84,0.86,0.91,0.82,0.89; 0.81,0.14,0.88,0.12,0.89,0.15,0.82,0.17].

Q: Detect all large-vehicles shown in the image and describe using oriented bounding boxes.
A: There are two large-vehicles [0.14,0.65,0.46,0.55,0.47,0.60, 0.15,0.69;0.52,0.38,0.57,0.36,0.66,0.71,0.62,0.73] in the image.
Q: Detect all small-vehicles shown in the image and describe using oriented bounding boxes.
A: There are two small-vehicles [0.86,0.82,0.89,0.84,0.86,0.91,0.82, 0.89;0.81,0.14,0.88,0.12,0.89,0.15,0.82,0.17] in the image.

Fig. 3. The construction process of MMRS-1M dataset. MMRS-1M contains three visual modalities from multi-sensor (e.g., optical, SAR, and infrared) and five RS vision tasks data(e.g., classification, detection, image caption, VQA, and visual grounding).

optical datasets include DIOR [46], DOTA [47], FAIR1M [48], HRRSD [49], NWPUVHR10 [50], RSOD [66], UCAS-AOD [67], and VisDrone [68]. Furthermore, SAR object detection datasets like AIR-SARShip-2.0 [51], HRISD [52], SSDD [53] are included, along with infrared object detection datasets for HIT-UAV [69], Sea-shipping [70], Infrared-security [71], Aerial-mancar [72], Double-light-vehicle [73] and oceanic ship [74]. Additionally, the DIOR-RSVG [14] is

collected for visual grounding and region-level caption.

**Detection-to-instruction.** The detection is divided into HBB and OBB formats. The HBB format defines a bounding box as $[x_{min}, y_{min}, x_{max}, y_{max}]$. Here, $(x_{min}, y_{min})$ and $(x_{max}, y_{max})$ represent the corner points of the bounding box closest and farthest from the coordinate origin, respectively. On the other hand, the OBB format is defined as $[x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4]$. In this format, $(x_1, y_1)$ represents the corner point of the bounding box closest to the coordinate origin, and the remaining three points are sorted in ascending order based on their angles with respect to $(x_1, y_1)$. Additionally, for both HBB and OBB detection, the coordinates of the bounding box are normalized. During the dataset conversion process, specific instructions are used to guide the model in predicting either HBB or OBB. This process is referred to as "detection-to-instruction". The instruction "Detect all objects shown in the RS image and describe using horizontal bounding boxes" is used for HBB detection, while "Detect all objects shown in the RS image and describe using oriented bounding boxes" is used for OBB detection. The model then outputs the coordinates and categories of all objects in the image. Furthermore, to achieve referring detection, an additional multi-round conversation session is incorporated to detect objects of each category.

**VG-to-instruction.** For the visual grounding dataset including DIOR-RSVG, two types of conversation formats are constructed via VG-to-instruction. Specifically, in the first format, users describe the region of interest and provide instructions for MLLM to locate the corresponding target. The model then outputs the coordinates of the identified target. The second format involves the user inputting instructions to describe the target region based on its spatial location coordinates, and MLLM generates a description of the corresponding target region. The former format bestows EarthGPT with the ability to perform visual localization, while the latter is employed to accomplish the region caption task.

The aforementioned is a detailed description of the MMRS-1M dataset construction process. Additionally, for the datasets mentioned above, the training set and validation are processed through data transformation and the construction of MMRS-1M after image cutting and dataset splitting. The details on the training samples can be found in Tab. I. The construction flow, as depicted in Fig. 3, provides a more concrete and intuitive understanding of the dataset construction process.

In summary, we clean and convert 34 various existing RS datasets into a uniform visual instruction-following format, covering optical, SAR, and infrared visual modalities, facilitating the application of MLLMs in the RS domain.

## V. EXPERIMENTS AND ANALYSIS

In this section, we conduct extensive experiments to validate the performance of our proposed EarthGPT. In subsection A, we describe the implementation details. Subsequently, we demonstrate EarthGPT's powerful ability in various scenarios such as classification, image captioning, VQA, visual grounding, and object detection, compared to other MLMs, specialist models, and open-set models.

### A. Implementation Details

In the visual-enhanced perception stage, the DINOv2 ViT-L/14 [75] is adopted as the ViT encoder, and the frozen CLIP ConvNeXt-L [76] is utilized as the CNN encoder. The two visual encoders are kept frozen throughout the training. Subsequently, in the cross-modal mutual comprehension stage, the objective is to convert language-only LLM into an MLLM. We trained EarthGPT on LAION-400M [60], COCO Caption [61] datasets, which primarily focus on natural scene image captioning, to develop basic multi-modal alignment from scratch. Furthermore, in the unified multi-task tuning phase, the goal is to equip MLLM with the versatility needed for diverse RS downstream tasks. As mentioned in Section IV, the MMRS-1M dataset is constructed based on diverse and comprehensive RS task-specific datasets. By unifying all the datasets into a multi-modal conversational format, the training process is optimized, reducing costs and enhancing efficiency. During the training, we only train an off-the-shelf language model LLaMA-2 [21] and randomly initialized visual projections. We use the AdamW optimizer $(\beta1, \beta2) = (0.9, 0.95)$, a maximum learning rate of $2 \times 10^{-5}$, a minimum learning rate of 0.

### B. Scene Classification

To evaluate the classification performance, supervised and zero-shot classification assessments are conducted. In the first evaluation, the test sets of NWPU-RESISC45 are utilized. The NWPU-RESISC45 dataset is a large-scale publicly available dataset for RS image scene classification, published by Northwestern Polytechnical University. It is comprised of 45 scene categories, with each category containing 700 RS images. Each image has a size of $256 \times 256$ pixels. 80% of the dataset is used as the test set to evaluate and compare EarthGPT with other specialist models.

TABLE II
SUPERVISED COMPARISON RESULTS ON NWPU-RESISC45 FOR SPECIALIST MODELS AND OUR EARTHGPT.

| Method | Publication Year | Top-1 Acc |
|---|---|---|
| ***Specialist Models*** | | |
| SeCo [77] | ICCV 2021 | 92.71 |
| MGSNet [78] | TGRS 2023 | **94.57** |
| CSDS [79] | JSTARS 2021 | 93.59 |
| T-CNN [80] | TGRS 2022 | 93.05 |
| PSGAN [81] | TGRS 2022 | 88.47 |
| ***MLLM*** | | |
| **EarthGPT(Ours)** | | 93.84 |

TABLE III
ZERO-SHOT COMPARISON RESULTS ON CLRS AND NASC-TG2 FOR OTHER MLLMS AND OUR EARTHGPT.

| Method | Publication Year | CLRS | NaSC-TG2 |
|---|---|---|---|
| Qwen-VL-Chat [82] | Arxiv 2023 | 51.76 | 36.09 |
| LLaVa-1.5 [6] | NeurIPS 2023 | 55.86 | 40.67 |
| Sphinx [83] | Arxiv 2023 | 60.72 | 49.76 |
| **EarthGPT(Ours)** | | **77.37** | **74.72** |

For zero-shot classification, all the images from the CLRS [15] and NaSC-TG2 [16] datasets are employed as the zero-shot testing set. The CLRS dataset includes 15,000 RS images with 25 land-use types such as airports, beaches, etc. These images were extracted from Google Earth, Bing Maps, Google Maps, and Tianditu, and were released by Central South University in 2020. Each image has a size of $256 \times 256 \times 3$, and the resolutions range from 0.26 m to 8.85 m. The NaSC-TG2 dataset consists of 20,000 RS images of 10 land cover types. These images were extracted from the Tiangong-2 satellite and were released by the Center for Space Applications and Engineering of the Chinese Academy of Sciences in 2021. Each image in this dataset has a size of $128 \times 128 \times 3$, and the pixel resolution is 100 m. For the evaluation of supervised classification and zero-shot classification, Top-1 accuracy is reported.

In the classification evaluation, similar to the training phase, we present all categories from the datasets as reference categories. The model is instructed to classify images using only one word or phrase. First, we report the accuracy for EarthGPT on the test sets of NWPU-RESISC45 in Tab. II, comparing its performance with some SOTA specialist models. It is observed that EarthGPT outperforms SOTA specialist models such as CSDS, and T-CNN with improvements of 0.25% and 0.79%, achieving comparable performance with MGSNet. In the zero-shot classification evaluation, we compare EarthGPT with other MLLMs such as QwenVL-Chat, LLaVa-1.5, and Sphinx in Tab. III. It is apparent that EarthGPT shows a notable improvement in the zero-shot evaluation compared to other MLLMs. Specifically, EarthGPT outperforms Qwen-VL-Chat, LLaVaV1.5, and Sphinx in top-1 accuracy on the CLRS dataset, achieving 25.61%, 21.51%, 16.65% respectively. Additionally, EarthGPT brings improvements of 38.63%, 34.05%, and 24.96% on NaSC-TG2 dataset compared with Qwen-VL-Chat, LLaVaV1.5, and Sphinx, respectively. This suggests that the integration of domain knowledge from RS is highly important for generalizing to unknown classification scenarios.

### C. Image Captioning

For the evaluation of image captioning capability, the test set of the NWPU-Caption dataset is utilized to assess and compare EarthGPT with specialist models in supervised setting. The NWPU-Caption dataset is created by Northwestern Polytechnical University, incorporating 31,500 aerial RS images and 157,500 sentences for RS image caption. Following the setting of MLCA-Net, we employ BLEU1, BLEU2, BLEU3, BLEU4, METEOR, ROUGEL, and CIDErD as evaluation metrics.

From Tab. IV, it can be clearly seen that on the NWPU-Captions dataset, compared to other SOTA methods, EarthGPT improves 10.1%, 13.8%, 17.5%, 17.7%, 10.5%, 18.2%, 66.2%, and 3.7% in terms of BLEU-1, BLEU-2, BLEU-3, BLEU-4, METEOR, ROUGE-L, and CIDEr-D. It is evident that EarthGPT provides accurate, detailed, and diverse descriptions of RS images.

### D. Visual Question Answering

The evaluation of VQA is also divided into supervised and zero-shot assessments. For the supervised VQA setting, we adopt the test set of CRSVQA including 1000 image-question-answer pairs. Following the setting of MQVQA, 10% of image-question-answer pairs are adopted as the test set, and the overall accuracy (OA) is reported. For zero-shot VQA setting, RSVQA-HR is utilized to evaluate the performance of EarthGPT and other MLLMs. The RSVQA-HR dataset was created from high-resolution orthophoto (HRO) datasets obtained from USGS. It consists of 10,659 images, with 1,066,316 question-answer pairs. RSVQA-HR dataset is divided into training, validation, test set 1, and test set 2. In the evaluation of zero-shot VQA, test set 2 is adopted for evaluating the model's robustness to diverse locations. Additionally, following the setting of RSGPT and GeoChat, object counting problems are not included in the computation of metrics for RSVQA-HR. For the RSVQA-HR dataset, the accuracy of different question categories and overall accuracy have been reported.

From Tab. V, it is observed that EarthGPT significantly outperforms other specialist models in terms of OA on CRSVQA datasets. Specifically, EarthGPT surpasses MQVQA and SAN with improvements of 11.09% and 20.83%, respectively. In zero-shot evaluation, EarthGPT achieves an average accuracy of 72.05% on RSVQA-HR above or commensurate with other MLLMs, as shown in Tab. VI. Compared with Qwen-VL-Chat, LLava v1.5, minigptV2, and Sphinx, EarthGPT brings the improvement of 3.66%, 9.00%, 25.60%, and 2.27%, respectively. In conclusion, EarthGPT demonstrates exceptional capabilities in answering questions within entirely new and unfamiliar RS contexts, and a profound ability to comprehend the true significance of queries presented in visual data and text instructions.

### E. Visual Grounding

To evaluate the visual grounding performance, we use the test set of DIOR-RSVG containing 758 grounding questions. DIOR-RSVG dataset comprises 38,320 language expressions across 17,402 RS images, uniquely characterizing individual objects within 20 diverse categories. Evaluation metrics include Pr@0.5, Pr@0.6, Pr@0.7, Pr@0.8, Pr@0.9 mean IoU(mIoU), and cum IoU(cIoU). Specialist models ZSGNet, FAOA, ReSc, LBYL-Net, TranVG, VLTVG and MGVLF are adopted to compare with EarthGPT.

Tab. VII shows the performance of EarthGPT and other specialist models on the DIOR-RSVG test set. It can be seen that EarthGPT brings significant performance improvements of 0.11%, 1.19%, 1.30%, and 3.13% on Pr@0.8, Pr@0.9, mIoU and cIoU, respectively. It is evident that EarthGPT have exceptional spatial location and region-level image comprehension ability.

### F. Object Detection

To thoroughly assess the potential and generalization capabilities of the innovative language-guided paradigm of EarthGPT in object detection, we adopt zero-shot setting for comparison with other MLLMs and specialist models. Given that other MLLMs are primarily designed for HBB detection and struggle with OBB detection, we utilize the HBB format of

TABLE IV
SUPERVISED COMPARISON RESULTS ON NWPU CAPTION FOR SPECIALIST MODELS AND OUR MLLM EARTHGPT.

| Method | Publication Year | BLEU1 | BLEU2 | BLEU3 | BLEU4 | METEOR | ROUGE_L | CIDEr(0-5) | SPICE |
|---|---|---|---|---|---|---|---|---|---|
| *Specialist Model* | | | | | | | | | |
| CSMLF [84] | GRSL 2019 | 77.0 | 64.9 | 53.2 | 47.1 | 32.0 | 57.8 | 106.5 | 26.5 |
| Multimodal [56] | CITS 2016 | 72.5 | 60.3 | 51.8 | 45.5 | 33.6 | 59.1 | 117.9 | 27.6 |
| Attention(soft) [37] | TGRS 2017 | 73.1 | 60.9 | 52.5 | 46.2 | 33.9 | 59.9 | 113.6 | 28.5 |
| Attention(hard) [37] | TGRS 2017 | 73.3 | 61.0 | 52.7 | 46.4 | 34.0 | 60.0 | 110.3 | 28.4 |
| FC-Att [85] | TGRS 2019 | 73.6 | 61.5 | 53.2 | 46.9 | 33.8 | 60.0 | 123.1 | 28.3 |
| SM-Att [85] | TGRS 2019 | 73.9 | 61.7 | 53.2 | 46.8 | 33.0 | 59.3 | 123.6 | 27.6 |
| MLCA-Net [86] | TGRS 2022 | 74.5 | 62.4 | 54.1 | 47.8 | 33.7 | 60.1 | 126.4 | 28.5 |
| *MLLM* | | | | | | | | | |
| **EarthGPT(Ours)** | | **87.1** | **78.7** | **71.6** | **65.5** | **44.5** | **78.2** | **192.6** | **32.2** |

TABLE V
SUPERVISED COMPARISON RESULTS ON CRSVQA FOR SPECIALIST MODELS AND OUR EARTHGPT.

| Method | Publication Year | OA |
|---|---|---|
| *Specialist models* | | |
| Qonly [87] | CVPR 2019 | 23.49 |
| RSVQA [39] | TGRS 2020 | 58.96 |
| RSVQA(GRU) [39] | TGRS 2020 | 59.41 |
| SAN [88] | CVPR 2016 | 61.17 |
| MQVQA [13] | TGRS 2023 | 70.91 |
| *MLLM* | | |
| **EarthGPT(Ours)** | | **82.00** |

TABLE VI
ZERO-SHOT COMPARISON RESULTS ON RSVQA-HR FOR OTHER MLLMS AND OUR EARTHGPT.

| Method | Publication Year | Presence | Comparison | OA |
|---|---|---|---|---|
| Qwen-VL-Chat [82] | Arxiv 2023 | **69.83** | 67.29 | 68.40 |
| LLaVa-1.5 [6] | NeurIPS 2023 | 66.44 | 60.41 | 63.06 |
| MiniGPTV2 [89] | Arxiv 2023 | 40.79 | 50.91 | 46.46 |
| Sphinx [83] | Arxiv 2023 | 64.28 | 74.75 | 69.79 |
| GeoChat [11] | Arxiv 2023 | 58.45 | **83.19** | **72.30** |
| **EarthGPT(Ours)** | | 62.77 | 79.53 | 72.06 |

the MAR20 dataset [17] to evaluate the detection performance, employing AP@40 and AP@50 as metrics. The MAR20 dataset is a large-scale RS aircraft target recognition dataset. It comprises 3,842 images with 22,341 instances, and each target instance is annotated with both horizontal bounding boxes and oriented bounding boxes. Comparing methods include MLLMs Qwen-VL-Chat, Sphinx, Lenna, and open-set object detection models GroundingDINO and mm-GroundingDINO. To address the challenge of MLLMs not predicting confidence

scores, we employ clip-score as a confidence logit for noise filtering. Remoteclip [35] weights are adopted to compute the clipscore. For the evaluation of OBB detection, we adopt the OBB format of MAR20 to compare with specialist models trained on DOTA, such as S2A_Net, CFA, Oriented RepPoints, Oriented R-CNN, Sasm, and R3Det.

From Tab. VIII, on the test set of MAR20, EarthGPT exhibits notable enhancements of AP@40, AP@50 compared to other MLLMs and open-set detection models, underscoring its robust generalization capability under unseen scenarios. Specifically, EarthGPT brings improvements of 1.92% and 1.90% in terms of AP@40 and AP@50 compared with the open-set model GroundingDINO. For OBB detection, Tab. IX demonstrates that EarthGPT achieves competitive performance with specialist models trained on DOTA, showcasing the powerful potential of language-based generative detection paradigm.

## VI. VISUALIZATION

In this section, we present the qualitative experimental result of EarthGPT to demonstrate proficiency in multi-turn dialogue and diverse RS tasks under multi-sensor visual modalities. EarthGPT showcases a remarkable capacity for image-level and region-level visual perception, and adeptly interpreting RS complex visual data. Notably, EarthGPT also shows advanced performance in chain-of-thought reasoning, fostering the emergence of cross-task visual comprehension.

### A. Multi-turn Multi-task Dialogue

Fig. 6 provides a detailed depiction of EarthGPT's capabilities in conducting multi-turn dialogues. Diverging from traditional specialist models limited to specific tasks, Earth-GPT leverages language as a medium to achieve diverse abilities in multi-turn dialogue, including classification, image captioning, VQA, object detection, region-level captioning, and visual grounding. In Fig. 6, we demonstrate EarthGPT's proficiency in accurately identifying the scene as a ground track field, along with a detailed spatial layout description.

TABLE VII
SUPERVISED COMPARISON RESULTS ON DIOR-RSVG FOR SPECIALIST MODELS, OTHER MLLMs, AND OUR EARTHGPT.

| Method | Publication Year | Pr@0.5 | Pr@0.6 | Pr@0.7 | Pr@0.8 | Pr@0.9 | mIoU | cIoU |
|---|---|---|---|---|---|---|---|---|
| *Specialist Model* | | | | | | | | |
| ZSGNet [90] | ICCV 2019 | 51.67 | 48.13 | 42.30 | 32.41 | 10.15 | 44.12 | 51.65 |
| FAOA [91] | ICCV 2019 | 70.86 | 67.37 | 62.04 | 53.19 | 36.44 | 62.86 | 67.28 |
| ReSC [92] | ECCV 2020 | 72.71 | 68.92 | 63.01 | 53.70 | 33.37 | 64.24 | 68.10 |
| LBYL-Net [93] | CVPR 2021 | 73.78 | 69.22 | 65.56 | 47.89 | 15.69 | 65.92 | 76.37 |
| TransVG [94] | ICCV 2021 | 72.41 | 67.38 | 60.05 | 49.10 | 27.84 | 63.56 | 76.27 |
| VLTVG [95] | CVPR 2021 | 75.79 | 72.22 | 66.33 | 55.17 | 33.11 | 66.32 | 77.85 |
| MGVLF [14] | TGRS 2023 | **76.78** | **72.68** | **66.74** | 56.42 | 35.07 | 68.04 | 78.41 |
| *MLLM* | | | | | | | | |
| **EarthGPT(Ours)** | | 76.65 | 71.93 | 66.52 | **56.53** | **37.63** | **69.34** | **81.54** |

TABLE VIII
ZERO-SHOT COMPARISON RESULTS ON MAR20 FOR OTHER METHODS AND OUR EARTHGPT.

| Method | Publication Year | AP@40 | AP@50 |
|---|---|---|---|
| *Open-set Model* | | | |
| GroundingDINO [96] | Arxiv 2023 | 88.55 | 88.21 |
| mm-GroundingDINO [97] | Arxiv 2024 | 88.37 | 88.14 |
| *MLLM* | | | |
| Lenna [98] | Arxiv 2023 | 72.69 | 72.11 |
| Qwen-VL-Chat [82] | Arxiv 2023 | 72.56 | 40.04 |
| Sphinx [83] | Arxiv 2023 | 81.03 | 80.48 |
| **EarthGPT(Ours)** | | **90.47** | **90.11** |

TABLE IX
ZERO-SHOT COMPARISON RESULTS FOR OBB DETECTION ON MAR20 FOR OTHER METHODS AND OUR EARTHGPT.

| Method | Publication Year | AP@40 | AP@50 |
|---|---|---|---|
| *Specialist Model* | | | |
| S2A-Net [99] | TGRS 2021 | 90.71 | 90.23 |
| CFA [100] | CVPR 2021 | 90.66 | 90.28 |
| Oriented RepPoints [101] | CVPR 2022 | 90.68 | **90.56** |
| Oriented R-CNN [102] | ICCV 2021 | **90.70** | 90.54 |
| Sasm [103] | AAAI 2022 | 90.19 | 89.83 |
| R3Det [104] | AAAI 2021 | 90.61 | 90.11 |
| *MLLM* | | | |
| **EarthGPT(Ours)** | | 90.53 | 87.86 |

EarthGPT localizes and describes specific areas precisely, such as the ground track field and the nearby small baseball field. Compared to the previous specialist models, EarthGPT enhances readability and comprehension, making it more user-friendly and accessible.

### B. Multi Visual Modality Inference

Previous specialist models in the RS domain mainly concentrate on single visual modality and single task. Thanks to our MMRS-1M's rich multi-sensor information, laying the foundation for training the versatile EarthGPT. EarthGPT has robust generalization capabilities across optical, SAR, and infrared visual modalities. Fig. 4 showcases EarthGPT's remarkable performance across a variety of tasks in different visual modalities.

Although the MMRS-1M dataset involves limited SAR and infrared data compared to the optical data. Note that EarthGPT overcomes the constraints by using language to integrate different visual modalities. As shown in Fig. 4, in optical imagery, EarthGPT can accurately enumerate the number of red cars in the scene, identify all ships on water surfaces, and provide detailed descriptions of dense residential. In SAR imagery, EarthGPT accurately counts storage tanks, evaluates ship presence, and concisely describes the airport. In infrared imagery, EarthGPT precisely locates cars in certain areas, identifies persons on a basketball court, and accurately interprets their activities. Those instances illustrate EarthGPT's exceptional knowledge transfer and generalization abilities.

### C. Chain-of-thought Prompting for Visual Reasoning

We discover that EarthGPT can employ chain-of-thought prompting to enhance visual understanding and reasoning accuracy. For instance, EarthGPT can enhance the OBB detection completeness through the result of object counting as the prompt. As shown in Fig. 5, in the airport detection scene, EarthGPT initially misses two planes, when promptly by the question "How many planes are there in the airport?", after discovering there are 15 planes in total, EarthGPT intelligently uses this newfound knowledge to re-evaluate the aircraft detection task, ultimately successfully identifying all the planes. In addition, EarthGPT is adept at enhancing the accuracy of OBB detection by harnessing HBB detection results. In the scenario involving SAR ships, EarthGPT effectively detects all ships using HBB but misses ships using OBB. To improve the
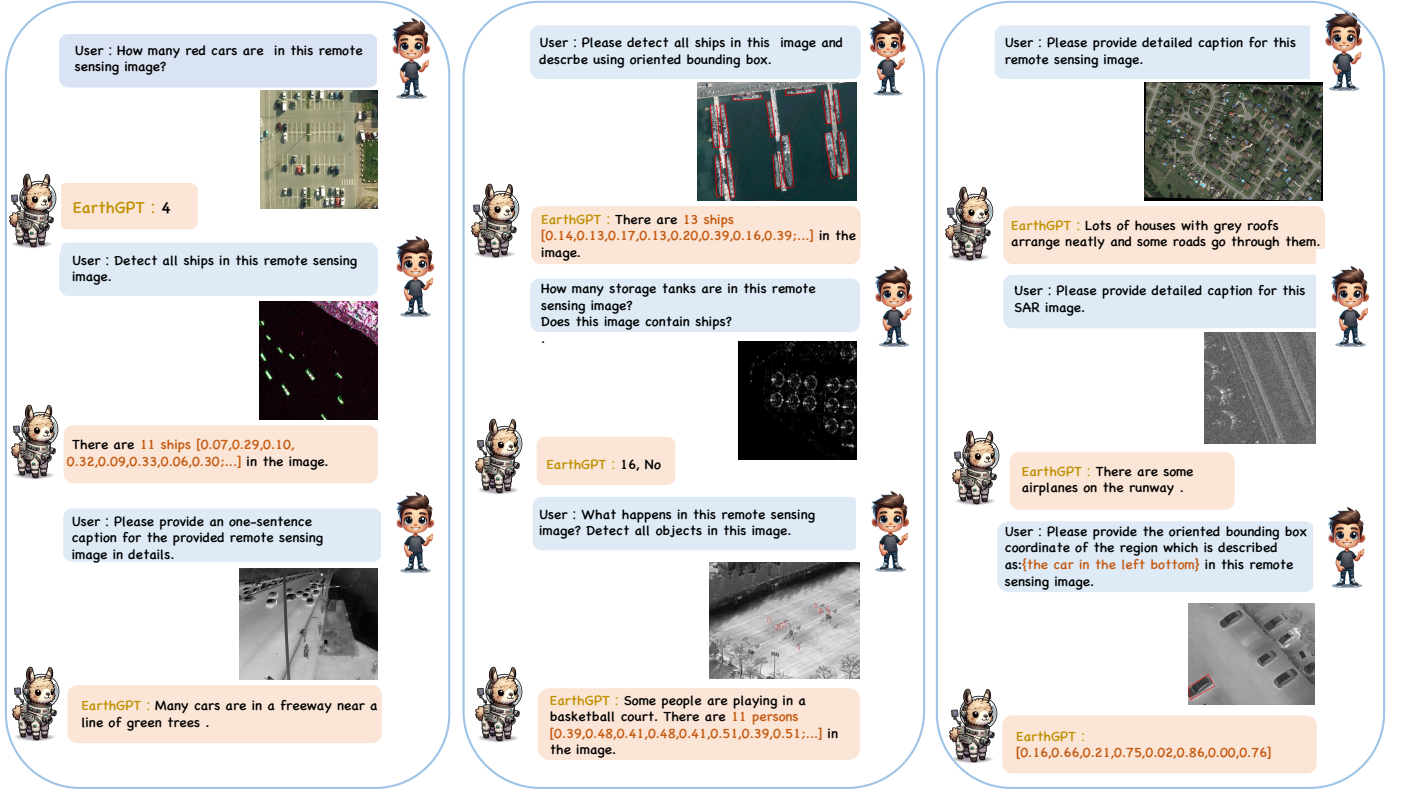
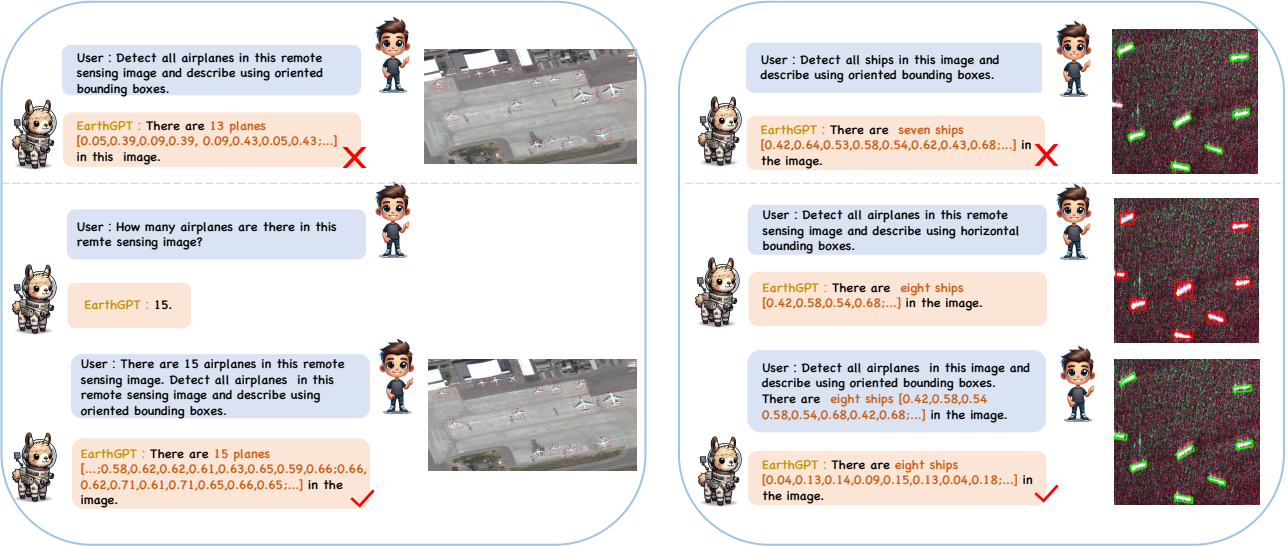Fig. 4. Examples of EarthGPT for different visual modalities inference ability.



Fig. 5. Examples of the chain-of-thought prompting for EarthGPT to perform visual reasoning.

OBB results, EarthGPT utilizes the HBB detection results as a valuable hint to rectify the omissions and deliver more precise OBB results. In conclusion, the chain-of-thought prompting technique facilitates intelligently incorporating insights from related tasks to elevate the performance in various visual interpretation and reasoning tasks.

## VII. CONCLUSION

In this paper, a versatile MLLM EarthGPT, which unifies a wide range of RS tasks and various multi-sensor images,

has been proposed for universal RS image comprehension. EarthGPT integrates coarse-scale and fine-scale visual perception information, bridging the gaps in cross-modal mutual comprehension and vision reasoning. More importantly, EarthGPT is capable of multi-sensor image interpretation and RS downstream tasks including scene classification, image captioning, region-level captioning, VQA, visual grounding, object detection. Furthermore, the MMRS-1M dataset, a multi-sensor multi-modal RS instruction-following dataset compris-

Fig. 6. Examples of EarthGPT for multi-turn various RS tasks dialogue.

ing more than 1M image-text pairs, has been constructed. The MMRS-1M dataset tackles the limitation of MLLMs on RS expert knowledge and encourages the growth of MLLMs tailored specifically for applications in the RS domain. Extensive experiments have demonstrated that EarthGPT surpasses the numerous existing specialist models and MLLMs in various RS visual interpretation tasks, and provides an open-set reasoning capability suitable for multiple RS downstream tasks, both in supervised and zero-shot settings. In the future, we will focus on improving the OBB detection performance of MLLMs and incorporating more modalities into EarthGPT for all-purpose capabilities.

## REFERENCES

[1] Lei Ma, Yu Liu, Xueliang Zhang, Yuanxin Ye, Gaofei Yin, and Brian Alan Johnson. Deep learning in remote sensing applications: A meta-analysis and review. *ISPRS journal of photogrammetry and remote sensing*, 152:166–177, 2019.

[2] Xiao Xiang Zhu, Devis Tuia, Lichao Mou, Gui-Song Xia, Liangpei Zhang, Feng Xu, and Friedrich Fraundorfer. Deep learning in remote sensing: A comprehensive review and list of resources. *IEEE geoscience and remote sensing magazine*, 5(4):8–36, 2017.

[3] Tong Zhang, Peng Gao, Hao Dong, Yin Zhuang, Guanqun Wang, Wei Zhang, and He Chen. Consecutive pre-training: A knowledge transfer learning strategy with relevant unlabeled data for remote sensing domain. *Remote Sensing*, 14(22):5675, 2022.

[4] Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. Visualgpt: Data-efficient adaptation of pretrained language models for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18030–18040, 2022.

[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. *arXiv preprint arXiv:2301.12597*, 2023.

[6] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.

[7] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.

[8] Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. Mmicl: Empowering vision-language model with multi-modal in-context learning. *arXiv preprint arXiv:2309.07915*, 2023.

[9] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

[10] Yuan Hu, Jianlong Yuan, Congcong Wen, Xiaonan Lu, and Xiang Li. Rsgpt: A remote sensing vision language model and benchmark. *arXiv preprint arXiv:2307.15266*, 2023.

[11] Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, and Fahad Shahbaz Khan. Geochat: Grounded large vision-language model for remote sensing, 2023.

[12] Qimin Cheng, Haiyan Huang, Yuan Xu, Yuzhuo Zhou, Huanying Li, and Zhongyuan Wang. Nwpu-captions dataset and mlca-net for remote sensing image captioning. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–19, 2022.

[13] Meimei Zhang, Fang Chen, and Bin Li. Multi-step question-driven visual question answering for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 2023.

[14] 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.

[15] Haifeng Li, Hao Jiang, Xin Gu, Jian Peng, Wenbo Li, Liang Hong, and Chao Tao. Clrs: Continual learning benchmark for remote sensing image scene classification. *Sensors*, 20(4):1226, 2020.

[16] Zhuang Zhou, Shengyang Li, Wei Wu, Weilong Guo, Xuan Li, Guisong Xia, and Zifei Zhao. Nasc-tg2: Natural scene classification with tiangong-2 remotely sensed imagery. *IEEE Journal of Selected Topics*

*in Applied Earth Observations and Remote Sensing*, 14:3228–3242, 2021.

[17] Wenqi Yu, Gong Cheng, Meijun Wang, Yanqing Yao, Xingxing Xie, XW Yao, and JW Han. Mar20: A benchmark for military aircraft recognition in remote sensing images. *National Remote Sensing Bulletin*, 2022.

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

[19] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

[20] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023.

[21] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Alma-hairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

[22] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.

[23] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*, 2023.

[24] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023.

[25] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.

[26] Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bing-shuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*, 2023.

[27] Zineng Tang, Ziyi Yang, Mahmoud Khademi, Yang Liu, Chenguang Zhu, and Mohit Bansal. Codi-2: In-context, interleaved, and interactive any-to-any generation. *arXiv preprint arXiv:2311.18775*, 2023.

[28] Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointllm: Empowering large language models to understand point clouds. *arXiv preprint arXiv:2308.16911*, 2023.

[29] Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023.

[30] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *arXiv preprint arXiv:2307.12981*, 2023.

[31] Senqiao Yang, Jiaming Liu, Ray Zhang, Mingjie Pan, Zoey Guo, Xiaoqi Li, Zehui Chen, Peng Gao, Yandong Guo, and Shanghang Zhang. Lidar-llm: Exploring the potential of large language models for 3d lidar understanding. *arXiv preprint arXiv:2312.14074*, 2023.

[32] Michael Ahn, Debidatta Dwibedi, Chelsea Finn, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Karol Hausman, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, et al. Autort: Embodied foundation models for large scale orchestration of robotic agents. *arXiv preprint arXiv:2401.12963*, 2024.

[33] Xian Sun, Peijin Wang, Wanxuan Lu, Zicong Zhu, Xiaonan Lu, Qibin He, Junxi Li, Xuee Rong, Zhujun Yang, Hao Chang, et al. Ringmo: A remote sensing foundation model with masked image modeling. *IEEE Transactions on Geoscience and Remote Sensing*, 2022.

[34] Yezhen Cong, Samar Khanna, Chenlin Meng, Patrick Liu, Erik Rozi, Yutong He, Marshall Burke, David Lobell, and Stefano Ermon. Sat-mae: Pre-training transformers for temporal and multi-spectral satellite imagery. *Advances in Neural Information Processing Systems*, 35:197–211, 2022.

[35] Fan Liu, Delong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, and Jun Zhou. Remoteclip: A vision language foundation model for remote sensing. *arXiv preprint arXiv:2306.11029*, 2023.

[36] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.

[37] Xiaoqiang Lu, Binqiang Wang, Xiangtao Zheng, and Xuelong Li. Ex-ploring models and data for remote sensing image caption generation. *IEEE Transactions on Geoscience and Remote Sensing*, 56(4):2183–2195, 2017.

[38] Xiangtao Zheng, Binqiang Wang, Xingqian Du, and Xiaoqiang Lu. Mutual attention inception network for remote sensing visual question answering. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–14, 2021.

[39] 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.

[40] Maryam Rahnemoonfar, Tashnim Chowdhury, Argho Sarkar, Debvrat Varshney, Masoud Yari, and Robin Roberson Murphy. Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access*, 9:89644–89654, 2021.

[41] 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.

[42] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.

[43] 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.

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

[45] Dengxin Dai and Wen Yang. Satellite image classification via two-layer sparse coding with biased image representation. *IEEE Geoscience and remote sensing letters*, 8(1):173–176, 2010.

[46] Junwei Han, Dingwen Zhang, Gong Cheng, Lei Guo, and Jinchang Ren. Object detection in optical remote sensing images based on weakly supervised learning and high-level feature learning. *IEEE Transactions on Geoscience and Remote Sensing*, 53(6):3325–3337, 2014.

[47] Jian Ding, Nan Xue, Gui-Song Xia, Xiang Bai, Wen Yang, Michael Yang, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. Object detection in aerial images: A large-scale benchmark and challenges. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.

[48] Xian Sun, Peijin Wang, Zhiyuan Yan, Feng Xu, Ruiping Wang, Wenhui Diao, Jin Chen, Jihao Li, Yingchao Feng, Tao Xu, et al. 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, 2022.

[49] Yuanlin Zhang, Yuan Yuan, Yachuang Feng, and Xiaoqiang Lu. Hierar-chical and robust convolutional neural network for very high-resolution remote sensing object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 57(8):5535–5548, 2019.

[50] Ke Li, Gong Cheng, Shuhui Bu, and Xiong You. Rotation-insensitive and context-augmented object detection in remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 56(4):2337–2348, 2017.

[51] SUN Xian, WANG Zhirui, SUN Yuanrui, DIAO Wenhui, ZHANG Yue, and FU Kun. Air-sarship-1.0: High-resolution sar ship detection dataset. *Journal of Radars*, 8(6):852–863, 2019.

[52] Shunjun Wei, Xiangfeng Zeng, Qizhe Qu, Mou Wang, Hao Su, and Jun Shi. Hrsid: A high-resolution sar images dataset for ship detection and instance segmentation. *Ieee Access*, 8:120234–120254, 2020.

[53] Tianwen Zhang, Xiaoling Zhang, Jianwei Li, Xiaowo Xu, Baoyou Wang, Xu Zhan, Yanqin Xu, Xiao Ke, Tianjiao Zeng, Hao Su, et al. Sar ship detection dataset (ssdd): Official release and comprehensive data analysis. *Remote Sensing*, 13(18):3690, 2021.

[54] Syed Waqas Zamir, Aditya Arora, Akshita Gupta, Salman Khan, Guolei Sun, Fahad Shahbaz Khan, Fan Zhu, Ling Shao, Gui-Song Xia, and Xiang Bai. isaid: A large-scale dataset for instance segmentation in aerial images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 28–37, 2019.

[55] Junjue Wang, Zhuo Zheng, Ailong Ma, Xiaoyan Lu, and Yanfei Zhong. Loveda: A remote sensing land-cover dataset for domain adaptive semantic segmentation. In J. Vanschoren and S. Yeung, editors, *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran Associates, Inc., 2021.

[56] Bo Qu, Xuelong Li, Dacheng Tao, and Xiaoqiang Lu. Deep semantic understanding of high resolution remote sensing image. In *2016 International conference on computer, information and telecommunication systems (Cits)*, pages 1–5. IEEE, 2016.

[57] Zhiqiang Yuan, Wenkai Zhang, Kun Fu, Xuan Li, Chubo Deng, Hongqi Wang, and Xian Sun. Exploring a fine-grained multiscale method for cross-modal remote sensing image retrieval. *arXiv preprint arXiv:2204.09868*, 2022.

[58] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[59] Keiron O'Shea and Ryan Nash. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*, 2015.

[60] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.

[61] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.

[62] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[63] Qin Zou, Lihao Ni, Tong Zhang, and Qian Wang. Deep learning based feature selection for remote sensing scene classification. *IEEE Geoscience and remote sensing letters*, 12(11):2321–2325, 2015.

[64] Yanghua Di, Zhiguo Jiang, and Haopeng Zhang. A public dataset for fine-grained ship classification in optical remote sensing images. *Remote Sensing*, 13(4):747, 2021.

[65] Yanghua Di, Zhiguo Jiang, Haopeng Zhang, and Gang Meng. A public dataset for ship classification in remote sensing images. In *Image and Signal Processing for Remote Sensing XXV*, volume 11155, pages 515–521. SPIE, 2019.

[66] Yang Long, Yiping Gong, Zhifeng Xiao, and Qing Liu. Accurate object localization in remote sensing images based on convolutional neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, 55(5):2486–2498, 2017.

[67] Haigang Zhu, Xiaogang Chen, Weiqun Dai, Kun Fu, Qixiang Ye, and Jianbin Jiao. Orientation robust object detection in aerial images using deep convolutional neural network. In *2015 IEEE International Conference on Image Processing (ICIP)*, pages 3735–3739. IEEE, 2015.

[68] Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Heng Fan, Qinghua Hu, and Haibin Ling. Detection and tracking meet drones challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7380–7399, 2021.

[69] Jiashun Suo, Tianyi Wang, Xingzhou Zhang, Haiyang Chen, Wei Zhou, and Weisong Shi. Hit-uav: A high-altitude infrared thermal dataset for unmanned aerial vehicle-based object detection. *Scientific Data*, 10(1):227, 2023.

[70] InfiRay. Sea-shipping. http://openai.iraytek.com/apply/Sea_shipping.html/, 2021.

[71] InfiRay. Infrared-security. http://openai.iraytek.com/apply/Infrared_security.html/, 2021.

[72] InfiRay. Aerial-mancar. http://openai.raytrontek.com/apply/Aerial_mancar.html/, 2021.

[73] InfiRay. Double-light-vehicle. http://openai.raytrontek.com/apply/Double_light_vehicle.html/, 2021.

[74] Center for Optics Research and Engineering of Shandong University. oceanic-ship. http://www.gxzx.sdu.edu.cn/info/1133/2174.htm/, 2020.

[75] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.

[76] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. If you use this software, please cite it as below.

[77] Oscar Manas, Alexandre Lacoste, Xavier Giró-i Nieto, David Vazquez, and Pau Rodriguez. Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9414–9423, 2021.

[78] Junjie Wang, Wei Li, Mengmeng Zhang, Ran Tao, and Jocelyn Chanussot. Remote-sensing scene classification via multistage self-guided separation network. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–12, 2023.

[79] Xinyu Wang, Liming Yuan, Haixia Xu, and Xianbin Wen. Csds: End-to-end aerial scenes classification with depthwise separable convolution and an attention mechanism. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:10484–10499, 2021.

[80] Weiquan Wang, Yushi Chen, and Pedram Ghamisi. Transferring cnn with adaptive learning for remote sensing scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–18, 2022.

[81] Gong Cheng, Xuxiang Sun, Ke Li, Lei Guo, and Junwei Han. Perturbation-seeking generative adversarial networks: A defense framework for remote sensing image scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–11, 2021.

[82] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.

[83] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv preprint arXiv:2311.07575*, 2023.

[84] Binqiang Wang, Xiaoqiang Lu, Xiangtao Zheng, and Xuelong Li. Semantic descriptions of high-resolution remote sensing images. *IEEE Geoscience and Remote Sensing Letters*, 16(8):1274–1278, 2019.

[85] Xiangrong Zhang, Xin Wang, Xu Tang, Huiyu Zhou, and Chen Li. Description generation for remote sensing images using attribute attention mechanism. *Remote Sensing*, 11(6):612, 2019.

[86] Qimin Cheng, Haiyan Huang, Yuan Xu, Yuzhuo Zhou, Huanying Li, and Zhongyuan Wang. Nwpu-captions dataset and mlca-net for remote sensing image captioning. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–19, 2022.

[87] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019.

[88] Kushal Kafle and Christopher Kanan. Answer-type prediction for visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4976–4984, 2016.

[89] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023.

[90] Arka Sadhu, Kan Chen, and Ram Nevatia. Zero-shot grounding of objects from natural language queries. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4694–4703, 2019.

[91] Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. A fast and accurate one-stage approach to visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4683–4693, 2019.

[92] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and Jiebo Luo. Improving one-stage visual grounding by recursive sub-query construction. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 387–404. Springer, 2020.

[93] Binbin Huang, Dongze Lian, Weixin Luo, and Shenghua Gao. Look before you leap: Learning landmark features for one-stage visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16888–16897, 2021.

[94] Jiajun Deng, Zhengyuan Yang, Tianlang Chen, Wengang Zhou, and Houqiang Li. Transvg: End-to-end visual grounding with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1769–1779, 2021.

[95] Li Yang, Yan Xu, Chunfeng Yuan, Wei Liu, Bing Li, and Weiming Hu. Improving visual grounding with visual-linguistic verification and iterative reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9499–9508, 2022.

[96] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding

dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.

[97] Xiangyu Zhao, Yicheng Chen, Shilin Xu, Xiangtai Li, Xinjiang Wang, Yining Li, and Haian Huang. An open and comprehensive pipeline for unified object grounding and detection. *arXiv preprint arXiv:2401.02361*, 2024.

[98] Fei Wei, Xinyu Zhang, Ailing Zhang, Bo Zhang, and Xiangxiang Chu. Lenna: Language enhanced reasoning detection assistant, 2023.

[99] Jiaming Han, Jian Ding, Jie Li, and Gui-Song Xia. Align deep features for oriented object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2021.

[100] Zonghao Guo, Chang Liu, Xiaosong Zhang, Jianbin Jiao, Xiangyang Ji, and Qixiang Ye. Beyond bounding-box: Convex-hull feature adaptation for oriented and densely packed object detection. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 8792–8801, 2021.

[101] Wentong Li, Yijie Chen, Kaixuan Hu, and Jianke Zhu. Oriented reppoints for aerial object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1829–1838, 2022.

[102] Xingxing Xie, Gong Cheng, Jiabao Wang, Xiwen Yao, and Junwei Han. Oriented r-cnn for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3520–3529, 2021.

[103] Liping Hou, Ke Lu, Jian Xue, and Yuqiu Li. Shape-adaptive selection and measurement for oriented object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.

[104] Xue Yang, Junchi Yan, Ziming Feng, and Tao He. R3det: Refined single-stage detector with feature refinement for rotating object. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3163–3171, 2021.