# VLM-FO1: Bridging the Gap Between High-Level Reasoning and Fine-Grained Perception in VLMs

**Peng Liu**[1], **Haozhan Shen**[3], **Chunxin Fang**[2], **Zhicheng Sun**[1], **Jiajia Liao**[2], **Tiancheng Zhao**[1,2,*]

[1]Om AI Research, [2]Binjiang Institute of Zhejiang University,
[3]College of Computer Science and Technology, Zhejiang University

`tianchez@zju-bj.com`

## Abstract

*Vision-Language Models (VLMs) excel at high-level scene understanding but falter on fine-grained perception tasks requiring precise localization. This failure stems from a fundamental mismatch, as generating exact numerical coordinates is an challenge task for language-centric architectures. In this paper, we introduce VLM-FO1, a novel framework that overcomes this limitation by reframing object-centric perception from a brittle coordinate generation problem into a robust feature retrieval task. Our method operates as a plug-and-play module that integrates with any pre-trained VLM. It leverages a Hybrid Fine-grained Region Encoder (HFRE), featuring a Dual-Vision Encoder, to generate powerful region tokens rich in both semantic and spatial detail. A token-based referencing system then enables the LLM to seamlessly reason about and ground language in these specific visual regions. Experiments show that VLM-FO1 achieves state-of-the-art performance across a diverse suite of benchmarks, demonstrating exceptional capabilities in Object Grounding, Region Generative Understanding, and Visual Region Reasoning. Crucially, our two-stage training strategy ensures these perception gains are achieved without compromising the base model's general visual understanding capabilities. VLM-FO1 establishes an effective and flexible paradigm for building perception-aware VLMs, bridging the gap between high-level reasoning and fine-grained visual grounding.*

## 1. Introduction

The advent of Large Language Models (LLMs) [3, 7, 28, 29, 57, 74, 92, 93, 107] has marked a paradigm shift in artificial intelligence, demonstrating profound capabilities in generation, reasoning, and instruction following. This success has been extended to the visual domain through Vision-Language Models (VLMs) [19, 25, 33, 46, 59, 60, 98, 129], which integrate powerful vision backbones with LLMs to interpret and reason about visual content. By mapping visual features into the language model's embedding space, VLMs have achieved remarkable performance on high-level visual understanding tasks such as visual question answering (VQA) and image captioning. Further advancements, such as the application of reinforcement learning techniques like Group Relative Policy Optimization (GRPO) [86], have continued to enhance their complex reasoning abilities.

Despite these advances, a critical weakness persists: state-of-the-art VLMs struggle with fine-grained visual perception tasks that demand precise spatial localization, such as object detection and grounding. This deficiency severely limits their applicability in real-world scenarios like autonomous robotics, detailed image analysis, and human-computer interaction, where understanding what is in an image is inseparable from knowing where it is. Our evaluations reveal a stark performance gap. On standard benchmarks like COCO [54], specialized detection models routinely achieve a mean Average Precision (mAP) of 50-60. In contrast, even a leading open-source model like Qwen2.5-VL-72B [9] achieves a recall of less than 40%, indicating a fundamental inability to reliably locate all relevant object instances.

This limitation stems from a core architectural mismatch: generating precise numerical coordinates is an "unnatural" task for models fundamentally designed for sequential language generation [82]. The requirement to produce a string of exact floating-point numbers in a specific format is brittle; a single incorrect token can render an entire bounding box prediction invalid. This problem is exacerbated in scenes with multiple instances, where the generation of a long, structured sequence of coordinates challenges the model's attention mechanism, leading to low recall and compounding errors.

To address this perception weakness, several approaches have been explored. Some methods [17, 97, 105] quantize object coordinates into a discrete vocabulary, simplifying the generation task. However, this approach still struggles with multiple instances and suffers from quantization er-
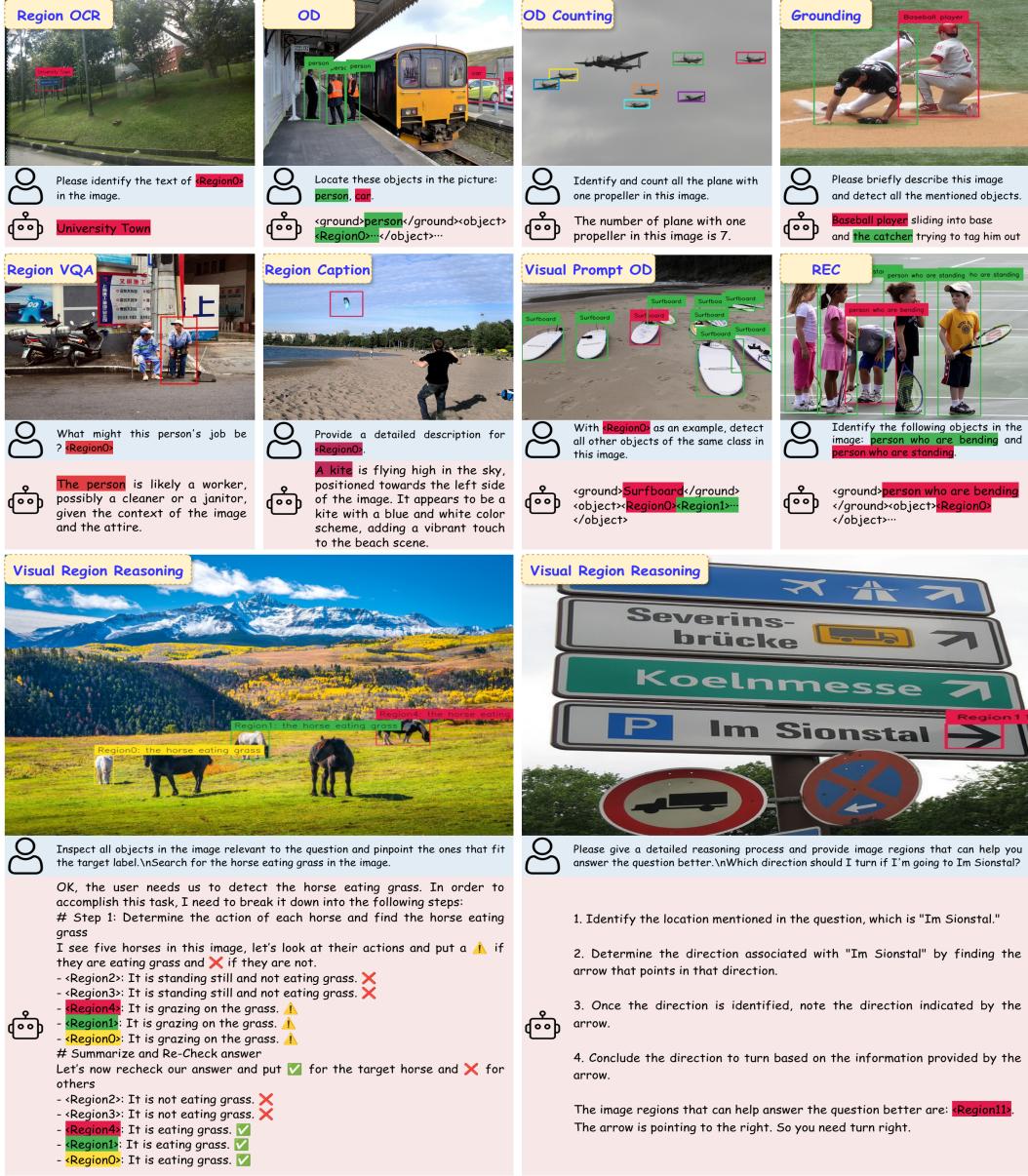
Figure 1. Visualization of VLM-FO1's perception abilities on diverse visual tasks.

rors, particularly on high-resolution images. Another strategy [45, 78, 121] involves appending external prediction heads to the VLM to handle localization. While potentially effective, these modules introduce additional latency and often require complex, task-specific loss functions for optimization. Recently, another direction has emerged that employs agentic, vision-level reasoning. Methods like Zoom-Eye [87] treat the image as a navigable structure, allowing the VLM to dynamically "zoom in" on specific regions to gather fine-grained visual evidence.

A more promising direction [20, 35, 69] reframes the problem by using an external detection model to generate region proposals, effectively converting the difficult generation task into a simpler retrieval task. While innovative, existing methods in this vein have significant drawbacks. They either require joint end-to-end training with the detection model, creating a monolithic and cumbersome system, or they necessitate training a new architecture from scratch on massive datasets. Crucially, both approaches fail to leverage the rich visual understanding and world knowledge already embedded within large-scale, pre-trained VLMs, effectively discarding a powerful and readily available resource.

In response to these challenges, we introduce VLM-FO1,

a novel framework that endows pre-trained VLMs with superior fine-grained perception without compromising their inherent strengths. The core idea is simple: we shift the paradigm from generating box coordinates to directly perceiving the content within them. VLM-FO1 treats any bounding box as a "visual prompt," extracts its features, and converts them into distinct "region tokens" that are fed directly into the LLM. This elegantly transforms object detection into a simple and accurate retrieval task.

Our primary innovations are threefold. First, VLM-FO1 is designed as a plug-and-play enhancement module that can be integrated with any pre-trained VLM, preserving its original capabilities while dramatically improving perception. Second, we introduce a novel Hybrid Fine-grained Region Encoder (HFRE), which features a Dual-Vision Encoder structure. This combines the VLM's original semantic-rich vision encoder with a new perception-enhanced vision encoder, yielding powerful object features that capture both high-level meaning and fine-grained detail. Third, our two-stage decoupled framework separates the training of the proposal model and the VLM. This modularity grants users the flexibility to pair VLM-FO1 with any proposal detector best suited for their specific application.

When combined with our custom-trained Omni Proposal Network (OPN), our lightweight VLM-FO1-3B model achieves 44.4 mAP on COCO, an improvement of over 20 points that places it on par with specialized detectors and far ahead of other VLMs. This strong performance extends to a wide variety of other region-related perception tasks, such as Referring Expression Comprehension (REC), object counting, and OCR, demonstrating the versatility and effectiveness of our approach. In summary, our main contributions are:

- A Flexible and Modular Plug-and-Play Framework: We propose VLM-FO1, a perception enhancement framework whose two-stage, decoupled design allows it to be seamlessly integrated with any pre-trained VLM. This modularity enables practitioners to use off-the-shelf object detectors for proposal generation, enhancing fine-grained perception without requiring full retraining or compromising the VLM's original capabilities.
- A Novel Hybrid Fine-grained Region Encoder (HFRE): We introduce a Dual-Vision Encoder architecture that combines a semantic-rich vision encoder with a perception-enhanced tower to produce region tokens that are rich in both high-level meaning and fine-grained spatial detail.
- State-of-the-Art Performance: We demonstrate the effectiveness of our method by achieving state-of-the-art results across a diverse suite of benchmarks spanning three key perspectives: Object Grounding, Region Generative Understanding, and Visual Region Reasoning, setting a new standard for perception-enhanced VLMs.

## 2. Related Work

### 2.1. Vision-Language Models (VLMs)

Since the emergence of large language models (LLMs)[7, 21, 28, 94], they have achieved remarkable success across a wide range of linguistic applications, which has in turn fostered the development of Vision-Language Models (VLMs). Early pioneering works include [5, 43, 50]. Building on these foundations, LLaVA[59] leveraged GPT-4 [3] to construct training data, achieving strong performance in visual dialogue and reasoning, and inspiring a line of research on visual instruction data [14, 23, 60]. A typical architecture of VLMs encodes visual information through a vision encoder [77, 90, 116] and integrates the resulting visual tokens with textual tokens within the LLM backbone. Today, some of the most widely adopted open-source VLM families include LLaVA[48, 59, 61], QwenVL[8, 9, 98], and InternVL [18, 19, 100, 129].

### 2.2. VLMs with Detection Enhancement

To equip VLMs with detection capability, prior works have explored integrating detection heads [45, 101, 103] or incorporating visual expert models [38, 69, 89, 128]. More generally, most approaches adopt an auto-regressive strategy to sequentially generate the four coordinates of bounding boxes, well aligned with LLM backbones [13, 75, 110]. Building on this paradigm, the Griffon series [117–119] progressively unified localization tasks, introduced high-resolution perception structures, and curated multi-dimensional datasets, extending VLMs to both vision-language and vision-centric settings. More recently, general-purpose VLMs [9, 19, 98, 100, 129] trained with detection data have demonstrated strong performance, and reinforcement learning methods such as GRPO [86] have been employed to further enhance visual reasoning for precise detection [67, 88].

Several prior works, notably Groma [69] and ChatRex [35], share a similar high-level approach by treating detection as a region-token retrieval task with extra detector. However, these methods typically necessitate significant architectural modifications to the base VLM or require joint training with the object detector, mandating a costly retraining process from scratch. In contrast, our VLM-FO1 employs a plug-and-play Hybrid Fine-Grained Region Encoder that seamlessly integrates with existing pre-trained VLMs, preserving their powerful, pre-learned representations while producing more representative and better-aligned region features. Furthermore, a critical distinction lies in the handling of negative categories. While prior models are often limited to detecting only the positive instances present in an image, our training strategy enables VLM-
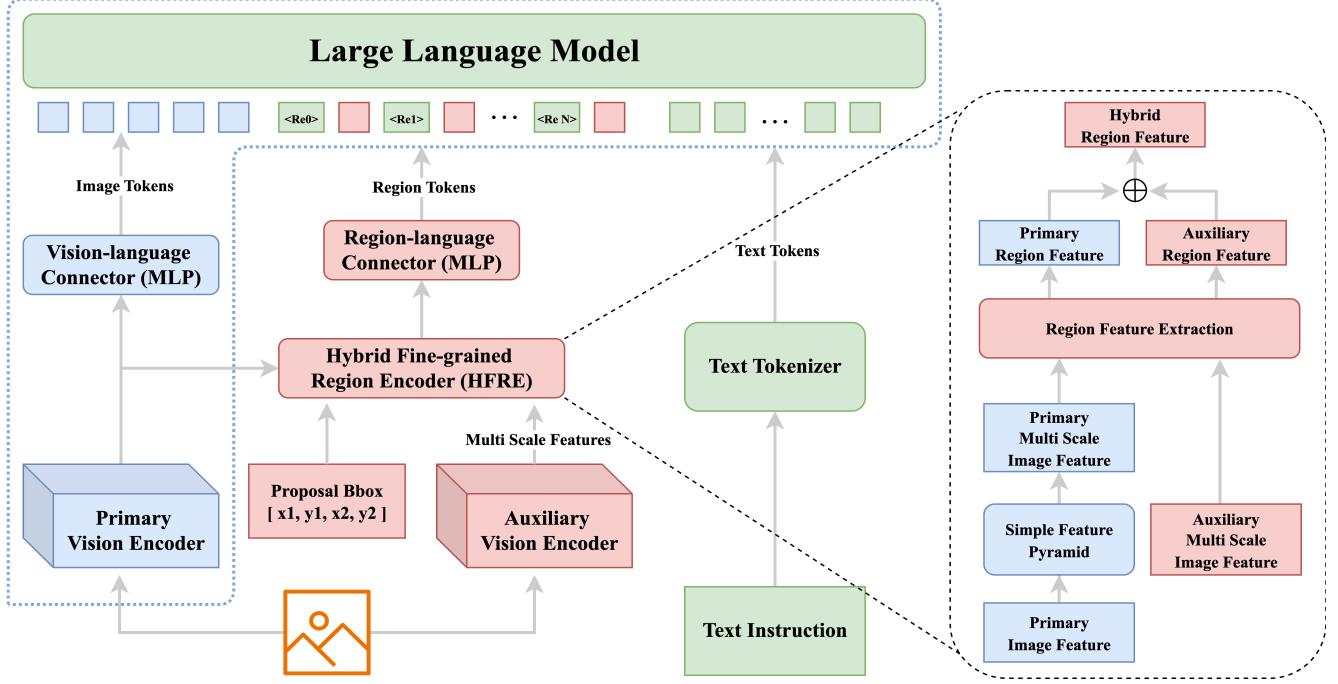
Figure 2. **An overview of our proposed model architecture.** The components enclosed within the blue dotted line represent a standard pretrained VLM, which can be initialized with existing weights to preserve its original performance. Our method introduces external modules, including a dual-vision encoder system and a Hybrid Fine-grained Region Encoder (HFRE), to enrich the VLM's fine-grained perception. These modules process an image to extract and fuse multi-scale region features, which are then fed as new region tokens to the VLM. The inset on the right details the generation of the Hybrid Region Feature.

FO1 to robustly distinguish and reject categories specified in the prompt that are not actually present, a crucial capability for real-world applications.

## 2.3. VLMs with Visual Prompt

Visual prompts take multiple forms. **Bounding-box prompts** provide coarse localization, as in Shikra [13] and InstructDET [24], which directly encode user-specified regions for object grounding. **Marker-based prompts**, such as Set-of-Mark (SoM) [108] and its extension SoM-LLaVA [106], overlay symbolic cues (e.g., circles, arrows, IDs) to highlight target regions. ViP-LLaVA [10] extends this idea by supporting arbitrary free-form markers like scribbles and arrows.

**Pixel-level prompts** enable more fine-grained perception. OMG-LLaVA [123] and VisionLLM [99] tokenize pixel-centric features, while DOrA [104] and CoLLaVO [47] integrate pixel-level annotations to enhance semantic grounding. ControlMLLM [102] further leverages training-free visual prompt learning to model semantic pixel-text relations.

In addition, **soft visual prompts** provide learnable perturbations in the pixel space. Transferable Visual Prompting (TVP) [124] explores prompt generalization across MLLMs, while BlackVIP [72] and ILM-VP [12] design task-specific perturbations to adapt pre-trained models efficiently.

## 3. Methodology

### 3.1. Model Architecture

The overall architecture of VLM-FO1, depicted in Figure 2, is designed as a set of plug-and-play modules that augment a standard pre-trained VLM. Instead of altering the core VLM or training from scratch, our framework introduces specialized components to process region-level visual information. Given an image and a set of bounding box proposals, our model extracts and fuses multi-scale region features, projects them into "region tokens," and feeds them alongside global image and text tokens into the Large Language Model. This allows the LLM to perform reasoning grounded in specific, fine-grained visual evidence. The key components of our architecture are the Dual-Vision Encoder and the Hybrid Fine-grained Region Encoder (HFRE).

**Proposal Regions.** A core tenet of our VLM-FO1 framework is its two-stage, decoupled design, which makes the generation of proposal regions entirely independent from the VLM's perception module. This separation facilitates independent training and optimization of each component. More importantly, it provides exceptional flexibility,

allowing users to switch between different object detectors to generate proposal regions based on specific scenarios, or even manually input regions of interest, all without any additional training. For our experiments and training, we developed an Omni Proposal Network (OPN), a variant based on OmDet-Turbo [126], to serve as a general-purpose detector. The OPN is trained to identify all potential foreground objects in an image, and we use it to generate the proposal regions for all our training and evaluation data.

**Dual-Vision Encoder.** The foundation of our fine-grained understanding capability is the Dual-Vision Encoder, a system engineered to produce region features that are simultaneously rich in semantic meaning and perceptual detail. It synergistically combines two components: the VLM's original Primary Vision Encoder and a new Auxiliary Vision Encoder. The primary encoder, having been co-trained with the LLM, excels at generating semantically-aligned features but lacks spatial precision due to its training on lower-resolution, global scenes. To compensate, the auxiliary encoder acts as a high-resolution detail specialist, processing the image at a higher fidelity to extract multi-scale feature maps rich in the perceptual cues (e.g., edges, textures) necessary for precise localization. While the primary encoder continues to provide global image tokens and semantic context for regions, the auxiliary encoder supplies the critical fine-grained information. The outputs from both are then intelligently fused to create a superior hybrid region feature representation.

**Hybrid Fine-grained Region Encoder (HFRE).** The HFRE is responsible for processing the multi-scale features from the Dual-Vision Encoder and generating the final region tokens for the LLM. This process involves three main stages: multi-scale feature extraction, hybrid feature fusion, and tokenization.

For the auxiliary vision encoder, we select DaViT [26], a vision transformer that combines a CNN-like multi-scale architecture with an efficient dual attention mechanism to capture both fine-grained local details and long-range global context. We extract a set of feature maps, $\{A_1, A_2, A_3, A_4\}$, from its backbone. These maps are upsampled via interpolation to match the spatial dimensions of the largest feature map and then concatenated along the channel dimension to form a dense feature map $A_{concat}$. Given N proposal bounding boxes $B = b_1, ..., b_N$, we use RoIAlign [32] followed by mean pooling to extract the corresponding region features, yielding the auxiliary region features $F_{aux} \in \mathbb{R}^{N \times D_a}$.

For the ViT-based primary vision encoder, which lacks a native feature pyramid, we introduce a Simple Feature Pyramid (SimpleFP) module, inspired by ViTDet [52]. This module takes the last feature map from the encoder and applies a series of convolutions and deconvolutions with strides $\{2, 1, 1/2, 1/4\}$ to construct a feature pyramid. Sim-

ilar to the auxiliary process, we then use RoIAlign to extract region features, resulting in the primary region features $F_{pri} \in \mathbb{R}^{N \times D_p}$.

Finally, the two sets of features are fused by concatenation to form a combined feature representation, $F_{comb} = \text{Concat}(F_{pri}, F_{aux}) \in \mathbb{R}^{N \times (D_p + D_a)}$. To explicitly provide the model with spatial information, we compute sine-cosine positional embeddings $E_{pos}$ from the coordinates of the proposal boxes and add them to the combined features: $F_{hybrid} = F_{comb} + E_{pos}$. This hybrid feature is then passed through a Region-Language Connector, an MLP layer, which projects it into the LLM's embedding space to produce the final region tokens.

## 3.2. Grounding Language to Regions via Token-based Referencing

To ground language in specific visual regions, our framework, inspired by previous works [35, 69], establishes a token-based referencing system that enriches the LLM's input with explicit, addressable region-level information. We augment the standard VLM input of image and text tokens by introducing our new region tokens. To enable the LLM to distinguish between and reference specific regions, we introduce a set of N special tokens, `<region0>`, `<region1>`, `...`, `<regionN-1>`, which serve as unique region index tokens.

The input to the LLM is structured as an interleaved sequence where each region token is preceded by its corresponding index token. This results in a final format of: `<image_tokens>\n<region0><region_token>...<regionN-1><region_token>\n<text_tokens>`. This structure allows the LLM to directly associate the region features with a unique identifier. Consequently, a user can refer to specific regions within a text prompt by simply using the corresponding index token.

For the model's output, we introduce special tokens to handle grounding tasks. The `<ground></ground>` tokens demarcate a noun phrase in the response that requires grounding, and the `<object></object>` tokens enclose the region index tokens that correspond to that phrase. For instance, a valid grounded response would be: "The `<ground>`people`</ground><object><region2><region10></object>` are dancing." This format provides an unambiguous link between textual concepts and their visual referents. For tasks that require simpler referencing without explicit grounding, the model can directly refer the region index token within its natural language response. This structured output format transforms complex localization into a native referencing task for the LLM.

## 3.3. Training Strategy

The training of VLM-FO1 is conducted in two distinct stages designed to efficiently integrate fine-grained per-

| Type | Sub-type | Model | MSCOCO val2017 | ODinW13 | OVDEval |
|------|----------|-------|----------------|---------|---------|
| **Detection Model** | **OD** | **Faster RCNN** [80] | 42.0 | - | - |
| | | **DETR** [11] | 43.3 | - | - |
| | | **DINO** [120] | 49.4 | - | - |
| | **OVD** | **GLIP** [51] | 49.8 | 52.1 | 18.4 |
| | | **Grounding DINO** [63] | 52.5 | 55.7 | 25.3 |
| | | **OmDet-Turbo**[126] | 53.4 | 54.1 | 25.9 |
| **VLM** | **Close-source** | **Gemini 1.5 Pro** [91] | - | 36.7[*] | - |
| | | **GPT-4o** [34] | 3.1 | - | - |
| | **Open-source** | **InternVL2.5-8B** [18] | 12.1 | 20.2[*] | - |
| | | **InternVL2.5-72B** | - | 31.7[*] | - |
| | | **Qwen2.5-VL-7B** [9] | 17.7 | 37.3[*] | - |
| | | **Qwen2.5-VL-72B** | - | 43.1[*] | - |
| | **OBJ-enhanced** | **VLM-R1-7B** [88] | - | - | 31.0 |
| | | **Lumen** [38] | 35.3 | - | - |
| | | **Griffon v2** [118] | 38.5 | - | - |
| | | **Griffon-G-7B** [117] | 40.2 | 43.8[*] | - |
| | | **ChatRex-7B** [35] | 4.3(48.2 reported) | - | - |
| | | **VLM-FO1-3B(Ours)** | **44.4** | **44.0** | **43.7** |

Table 1. **Object Grounding performance on OD benchmarks.** * indicates evaluation under a simplified setting where only ground-truth categories are queried.

ception while preserving the model's extensive pre-trained knowledge.

**Stage 1: Region-Language Alignment Training.** The primary objective of this initial stage is to align the newly introduced region tokens with the LLM's feature space with minimal disruption to the existing VLM weights. To achieve this, we first extend the LLM's vocabulary with our special tokens (e.g., `<RegionN>`, `<ground>`) and freeze the embeddings of the original vocabulary, ensuring that only the new token embeddings are updated. Concurrently, we freeze the parameters of the entire pre-trained VLM, including the primary vision encoder and the LLM itself. Training is focused exclusively on the newly added modules: the HFRE and the Region-Language Connector. This isolated training strategy allows the model to learn a robust mapping from visual regions to token space.

**Stage 2: Perception Instruction Finetuning.** The second stage aims to holistically enhance the model's perception capabilities by fine-tuning the integrated system on a broader set of instruction-based tasks. In this phase, we unfreeze the parameters of the Auxiliary Vision Encoder, the HFRE, the Region-Language Connector, and the LLM. The Primary Vision Encoder remains frozen throughout the entire training process, acting as a stable anchor for the original VLM's semantic understanding. The training dataset is expanded to include a wider variety of perception-focused instruction data.

## 4. Experiments

### 4.1. Experimental Setup

**Model Setup.** Our experiments are built upon the Qwen2.5-VL model, chosen for its excellent baseline performance in visual understanding. For the auxiliary vision encoder, we integrate a pre-trained DaViT-Large model. Within the HFRE module, the auxiliary encoder extracts features from 4 multi-scale layers, resulting in a region feature of dimension 3840. The primary vision encoder utilizes the SimpleFP module to generate 4 multi-scale features, each with a dimension of 512, which are then combined into a 2048-dimension feature. The final hybrid region feature thus has a dimension of 5888. For our two-stage training, we set the learning rate to 1e-3 for Stage 1 and 1e-5 for Stage 2. For each image, we process a maximum of 100 input proposals, selecting the top 100 predictions from our OPN based on their confidence scores.

**Training Data.** Our training data is structured to support our two-stage training strategy, as summarized in Table 10.

- Stage 1 (Region-Language Alignment): In this stage, training is focused on aligning the visual features from the HFRE with the LLM's embedding space. To achieve this, we use a curated collection of datasets centered on region-language tasks. This includes large-scale object detection datasets (COCO [54], O365 [85], V3Det [96]), grounding data (GOLDG [39]), and region caption data (Rexverse-2M [35]).

| Model | LVIS | | PACO | |
|---|---|---|---|---|
| | SS | S-IoU | SS | S-IoU |
| LLaVA-1.5 [59] | 49.0 | 19.8 | 42.2 | 14.6 |
| Kosmos-2 [75] | 39.0 | 8.7 | 32.1 | 4.8 |
| Shikra-7B [13] | 49.7 | 19.8 | 43.6 | 11.4 |
| GPT4RoI-7B | 51.3 | 12.0 | 48.0 | 12.1 |
| Ferret-7B [110] | 63.8 | 36.6 | 58.7 | 26.0 |
| Osprey-7B [113] | 65.2 | 38.2 | 73.1 | 52.7 |
| VisionLLM v2-7B [101] | 68.9 | 46.3 | 67.7 | 44.0 |
| VP-SPHINX-13B [55] | 87.1 | 62.9 | 76.8 | 51.3 |
| DAM-8B [53] | 89.0 | 77.7 | 84.2 | 73.2 |
| PAM-3B [56] | 88.6 | 78.3 | 87.4 | 74.9 |
| ChatRex-7B [35] | 89.8 | 82.6 | **91.4** | **85.1** |
| **VLM-FO1-3B (Ours)** | **92.4** | **86.4** | 88.1 | 77.6 |

Table 2. **Region-level classification performance of VLMs on LVIS and PACO datasets.**

- Stage 2 (Perception SFT): The second stage broadens the model's capabilities by training on a diverse mix of perception-focused instruction datasets. In addition to the data from Stage 1, we incorporate datasets for REC, grounding, region captioning, region reasoning, counting, region QA, and OCR. Furthermore, for detection-related tasks, we introduce rejection samples for 20% of the data, where the model is prompted to find objects that are not present in the image. This strategy encourages the model to be more discerning and avoid hallucinating objects based solely on the text prompt. Crucially, to mitigate catastrophic forgetting, we also include a subset of data from the OmChat-SFT collection (which contains data from LLaVA-1.5 [60], The Cauldron[46], CogVLM[33], etc.). This mix of conventional VLM task data ensures that the model retains its high-level scene interpretation abilities while mastering fine-grained perception.

## 4.2. Main Results

To comprehensively assess the effectiveness of VLM-FO1, we assess its capabilities across three key dimensions: Object Grounding, Region Generative Understanding, and Visual Region Reasoning. We benchmark our model on a diverse suite of tasks within each of these areas.

**Object Grounding.** We first evaluate the model's core ability to ground language in objects through detection tasks. As shown in Table 1, we benchmark VLM-FO1 on standard object detection with COCO [54], open-vocabulary detection in real-world settings with ODinW13 [49], and challenging language-based detection with hard negatives on OVDEval [109].

The results clearly demonstrate VLM-FO1's superiority. On COCO and ODinW13, VLM-FO1 significantly out-

| Model | COCO Text Accuracy(%) |
|---|---|
| ChatSpot-7B [125] | 31.8 |
| PAM-3B [56] | 42.2 |
| VP-LLAVA-8B [55] | 44.8 |
| VP-SPHINX-13B [55] | 45.4 |
| **VLM-FO1-3B (Ours)** | **59.0** |

Table 3. **Regional OCR performance on COCOText benchmark.**

| Model | Ferret Bench Refer. Reasoning |
|---|---|
| LLaVA-7B [59] | 31.7 |
| Kosmos-2 [75] | 33.7 |
| Osprey-7B [113] | 67.8 |
| Ferret-13B [110] | 68.7 |
| Ferret-v2-13B [122] | 79.4 |
| VP-LLAVA-8B [55] | 68.9 |
| VP-SPHINX-13B [55] | 71.4 |
| **VLM-FO1-3B (Ours)** | **80.1** |

Table 4. **Referring Reasoning performance of Ferret Bench.**

performs other VLM-based models, showcasing its powerful perception and high recall. For instance, GPT-4o [34] achieves a mere 3.1 mAP on COCO, confirming that even the most advanced models fail at direct coordinate regression. While ChatRex-7B [35] reports a high mAP of 48.2, this is achieved under a non-standard evaluation protocol where only the ground-truth categories for each image are provided as queries; under standard COCO evaluation, its performance drops to 4.3 mAP, likely due to an inability to handle negative categories. VLM-FO1 successfully overcomes both of these fundamental challenges.

More impressively, on ODinW13, our model achieves the highest score despite being tested under the rigorous, standard mAP protocol. It is important to note that many other VLMs (marked with *) are evaluated on ODinW13 using a simplified setting where only ground-truth categories are fed to the model individually. This easier setting avoids the challenge of distinguishing hard negatives. Even against models tested in this simplified manner, VLM-FO1, under standard evaluation, still comes out on top.

The most compelling results are on OVDEval, which evaluates performance on linguistic labels with hard negatives. Here, VLM-FO1 surpasses not only other VLMs but also specialized detection models like Grounding DINO. This highlights a key advantage of our approach: VLM-FO1 effectively leverages the world knowledge, entity

| Model | Refcoco | | | Refcoco+ | | | Refcocog | | HumanRef | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val | test | P | R | DF1 |
| Gemini 1.5 Pro [91] | 73.2 | 72.9 | 74.6 | 62.5 | 63.9 | 65.0 | 75.2 | 76.2 | - | - | - |
| DINOX [81] | - | - | - | - | - | - | - | - | 33.1 | 75.2 | 23.3 |
| Grounding DINO [63] | 90.6 | 93.2 | 88.2 | 88.2 | 89.0 | 75.9 | 86.1 | 87.0 | 33.1 | 75.2 | 23.3 |
| InternVL-2.5-8B [18] | 90.3 | 94.5 | 85.9 | 85.2 | 91.5 | 78.8 | 86.7 | 87.6 | 37.8 | 29.8 | 31.9 |
| Ferret-7B [110] | 87.5 | 91.4 | 82.5 | 80.8 | 87.4 | 73.1 | 83.9 | 84.8 | 43.2 | 34.4 | 34.3 |
| Groma-7B [69] | 89.5 | 92.1 | 86.3 | 83.9 | 88.9 | 78.1 | 86.37 | 87.01 | 48.7 | 65.9 | 42.1 |
| ChatRex-7B [35] | 91.0 | **94.1** | 87.0 | **89.8** | 91.9 | 79.3 | **89.8** | **90.0** | 72.2 | 50.4 | 55.6 |
| Qwen2.5-VL-7B [9] | 90.0 | 92.5 | 85.4 | 84.2 | 89.1 | 76.9 | 87.2 | 87.2 | 68.5 | 52.5 | 56.2 |
| Molmo-7B-D [25] | - | - | - | - | - | - | - | - | 82.5 | 77.7 | 72.6 |
| RexSeek-7B [37] | - | - | - | - | - | - | 84.0 | 84.4 | 85.8 | 85.9 | 82.3 |
| **VLM-FO1-3B(Ours)** | **91.12** | 93.7 | **87.6** | 86.4 | **91.9** | **80.6** | 88.9 | 88.3 | 87.1 | 83.3 | **82.6** |

Table 5. **Model Performance on Referring Benchmarks**

| Type | Model | CountBench Acc(%) | PixMo Count |
|---|---|---|---|
| Close Source | GPT-4V [1] | 69.9 | 45.0 |
| | GPT-4o-0513 [34] | 87.9 | 59.6 |
| | Gemini 1.5 Pro [91] | 85.8 | 64.3 |
| | Claude-3 Opus [6] | 83.6 | 43.3 |
| | Claude-3.5 Sonnet [6] | 89.7 | 58.3 |
| Open Source | LLaVA-1.5-13B [60] | 47.1 | 35.2 |
| | LLaVA OneVision-72B [48] | 84.3 | 60.7 |
| | InternVL2-8B [18] | 57.8 | 43.9 |
| | InternVL2-Llama-3-76B [18] | 74.7 | 54.6 |
| | InternVL2.5-78B [18] | 72.1 | - |
| | Qwen2-VL-72B [98] | 80.4 | 55.7 |
| | Pixtral-12B [4] | 78.8 | 51.7 |
| | Llama-3.2V-90B-Instruct [28] | 78.5 | 58.5 |
| | Molmo-7B-D [25] | 88.5 | 84.8 |
| | Molmo-72B [25] | **91.2** | 85.2 |
| | **VLM-FO1-3B (Ours)** | 87.8 | **86.0** |

Table 6. **Model Performance on Object Counting Benchmarks**

recognition, and reasoning abilities inherited from its VLM foundation to disambiguate complex and challenging text prompts.

**Region Generative Understanding.** We further evaluate our model's ability to understand and generate accurate textual descriptions based on specific visual regions. This is tested across three diverse tasks: region-level classification (Table 2), region-based OCR (Table 3), and referring reasoning (Table 4). The results unequivocally demonstrate VLM-FO1's SOTA capabilities. On the object-level LVIS and part-level PACO [79] datasets, our model sets a new state-of-the-art for region classification, with our efficient 3B model outperforming significantly larger 8B and 13B models. Our architecture demonstrates a strong capability for generating precise text targeting fine-grained regions. On the COCOText [95] benchmark for regional OCR, VLM-FO1 achieves a staggering 59.0% accuracy, surpassing the next best model by over 13 points. Finally, on the challenging referring reasoning subset of Ferret Bench [110], our model achieves a new SOTA score of 80.1, demonstrating that its strong fine-grained perception directly translates to a more accurate understanding of specific visual regions and their relationships.

**Visual Region Reasoning.** In this section, we evaluate the model's ability to leverage its fine-grained region features to perform complex reasoning. We benchmark this on two challenging task families: Referring Expression Comprehension (Table 5) and Object Counting (Table 6). In REC tasks such as Refcoco/+/g, the model must reason about a natural language description to identify the correct object. Our model achieves consistently top-tier results across these benchmarks. More significantly, on HumanRef [37], a difficult benchmark focusing on people with complex descriptions (attributes, positions, interactions) and hard negatives, VLM-FO1 demonstrates remarkable performance, achieving a new state-of-the-art. This success underscores its robust reasoning capability and its skill in disambiguating between visually similar instances based on nuanced language. Furthermore, in Object Counting, a task notorious for causing failures in large VLMs, our model excels by adopting a "Detect-then-Count" reasoning process. It first localizes all target instances and then aggregates the count, leading to superior accuracy on CountBench [73] and the challenging PixMo-Count [25] benchmark. This methodical approach allows our compact V-FO1-3B to outperform even much larger closed-source

| Model | AVG | MMBench v1.1 [65] | AI2D [42] | MMStar [16] | Hallusion Bench [30] | OCR Bench [66] | MathVista [68] | MMVet [68] | MMMU Val [114] |
|---|---|---|---|---|---|---|---|---|---|
| **Qwen2.5-VL-3B** | 64.5 | 76.8 | 81.4 | 56.3 | 46.6 | 82.8 | 61.2 | 60 | 51.2 |
| **VLM-FO1-3B** | **64.6** | 78.2 | 81.2 | 56.9 | 47.9 | 82.3 | 65.6 | 54.9 | 49.9 |

Table 7. **Comparison of general VLM capabilities on OpenCompass benchmarks.**

| Model | Average Score |
|---|---|
| VLM-FO1-3B | **67.65** |
| VLM-FO1-3B (QwenViT unfrozen) | 66.35 |
| Only Aux. Region Feat. (DaViT unfrozen) | 65.89 |
| Only Prim. Region Feat. (QwenViT frozen) | 65.76 |
| Only Prim. Region Feat. (QwenViT unfrozen) | 66.15 |

Table 8. **Ablation study for HFRE module.** (Aux. Region Feat. denotes Auxiliary Region Feature, and Prim. Region Feat. denotes Primary Region Feature.)

| Model | Average Score |
|---|---|
| Only Prim. Region Feat. with SimpleFP | **66.15** |
| Only Prim. Region Feat. w/o SimpleFP | 64.94 |

Table 9. **Impact of SimpleFP on Primary Region Feature Performance.**

models like GPT-4V and open-source models up to 72B parameters, highlighting the power of grounding complex reasoning in accurate, fine-grained perception.

## 4.3. Preservation of General VLM Capabilities

A critical aspect of our framework is its ability to enhance fine-grained perception without degrading the base model's general visual understanding capabilities. To validate this, we evaluated our VLM-FO1-3B model against the original Qwen2.5-VL-3B on the comprehensive OpenCompass [22] benchmark. The results, shown in Table 7, confirm that our method successfully avoids catastrophic forgetting.

Across a wide range of general VLM benchmarks—including MMBench [65], AI2D [42], and MMStar [16]—our VLM-FO1 model's performance remains virtually identical to the base Qwen2.5-VL model, with a negligible difference in the average score. This demonstrates the effectiveness of our two-stage training strategy; by initially freezing the core VLM during the Region-Language Alignment phase and subsequently mixing in general VLM instruction data during the second stage, our method successfully prevents the degradation of pre-existing knowledge. The results confirm that VLM-FO1 acts as a true enhancement module, allowing pre-trained VLMs to gain superior fine-grained perception while retaining their powerful and general visual understanding abilities.

## 5. Ablation Studies

To validate the effectiveness of the individual components in our model design, we conduct a series of ablation studies. For efficiency, all ablation experiments are conducted by training on a representative subset of our full perception SFT dataset. The "Average Score" reported in this section is an average calculated from the benchmark scores evaluated in the previous sections.

### 5.1. Effectiveness of Components in HFRE

We first analyze the contribution of the different components within the HFRE module. As shown in Table 8, we evaluate the effect of different feature combinations and training strategies. Our full model, VLM-FO1-3B, which combines primary and auxiliary region features and keeps the primary vision encoder (QwenViT) frozen, achieves the highest average score of 67.65. In contrast, using only the auxiliary region feature (65.89) or only the primary region feature (66.15) leads to a drop in performance, demonstrating the importance of combining both semantic and perceptual information from the two vision encoders. Furthermore, we observe that fine-tuning the primary vision encoder (66.35) slightly degrades performance, suggesting that freezing the original, well-aligned vision encoder helps to preserve its valuable semantic priors. These results strongly validate the effectiveness of our HFRE design.

### 5.2. Effectiveness of the Simple Feature Pyramid

Next, we evaluate the effectiveness of introducing the SimpleFP module for the ViT-based primary vision encoder. As shown in Table 9, in a controlled setting using only the primary region feature, we compare the performance with and

without the SimpleFP module. The results show that removing the SimpleFP module causes a significant drop in the average score from 66.15 to 64.94. This performance gap clearly indicates that constructing a multi-scale feature pyramid from the ViT's single-scale feature map is critical for extracting high-quality, information-rich region features.

## 6. Conclusion

In this work, we introduced VLM-FO1, a novel framework that successfully bridges the gap between the high-level reasoning of Vision-Language Models and the demands of fine-grained visual perception. By reframing object-centric tasks from a coordinate generation problem to a feature retrieval task, we circumvent a fundamental limitation of language-centric architectures. Our proposed method, featuring a plug-and-play modular design and an innovative Hybrid Fine-grained Region Encoder, effectively enhances pre-trained VLMs with state-of-the-art perception capabilities. Our extensive experiments demonstrate that VLM-FO1 achieves exceptional performance across a wide range of benchmarks spanning object grounding, region-level understanding, and complex visual reasoning, often outperforming much larger models. Crucially, these gains are achieved without compromising the base model's original general-purpose abilities, validating our training strategy and architectural design. VLM-FO1 establishes a powerful and flexible paradigm for developing the next generation of VLMs, paving the way for models with a deeper, more actionable understanding of the visual world.

## References

[1] Gpt-4v(ision) system card. 2023. 8

[2] Manoj Acharya, Kushal Kafle, and Christopher Kanan. Tallyqa: Answering complex counting questions. In *Proceedings of the AAAI conference on artificial intelligence*, pages 8076–8084, 2019. 2

[3] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1, 3

[4] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. Pixtral 12b. *arXiv preprint arXiv:2410.07073*, 2024. 8

[5] 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. 3

[6] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 1(1):4, 2024. 8

[7] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 1, 3

[8] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 3

[9] 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. 1, 3, 6, 8

[10] Mu Cai, Haotian Liu, Siva Karthik Mustikovela, Gregory P. Meyer, Yuning Chai, Dennis Park, and Yong Jae Lee. Making large multimodal models understand arbitrary visual prompts. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2024. 4

[11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 6

[12] Aochuan Chen, Yuguang Yao, Pin-Yu Chen, Yihua Zhang, and Sijia Liu. Understanding and improving visual prompting: A label-mapping perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19133–19143, 2023. 4

[13] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 3, 4, 7

[14] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023. 3

[15] 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 *European Conference on Computer Vision*, pages 370–387. Springer, 2024. 2

[16] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *Advances in Neural Information Processing Systems*, 37:27056–27087, 2024. 9

[17] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*, 2021. 1

[18] 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. 3, 6, 8

[19] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 1, 3

[20] An-Chieh Cheng, Yang Fu, Yukang Chen, Zhijian Liu, Xiaolong Li, Subhashree Radhakrishnan, Song Han, Yao Lu, Jan Kautz, Pavlo Molchanov, et al. 3d aware region prompted vision language model. *arXiv preprint arXiv:2509.13317*, 2025. 2

[21] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024. 3

[22] OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models. https://github.com/open-compass/opencompass, 2023. 9

[23] 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. 3

[24] Ronghao Dang, Jiangyan Feng, Haodong Zhang, Chongjian Ge, Lin Song, Lijun Gong, Chengju Liu, Qijun Chen, Feng Zhu, Rui Zhao, et al. Instructdet: Diversifying referring object detection with generalized instructions. *arXiv preprint arXiv:2310.05136*, 2023. 4

[25] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv e-prints*, pages arXiv–2409, 2024. 1, 8

[26] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. In *European conference on computer vision*, pages 74–92. Springer, 2022. 5

[27] Dawei Du, Pengfei Zhu, Longyin Wen, Xiao Bian, Haibin Lin, Qinghua Hu, Tao Peng, Jiayu Zheng, Xinyao Wang, Yue Zhang, et al. Visdrone-det2019: The vision meets drone object detection in image challenge results. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pages 0–0, 2019. 2

[28] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407, 2024. 1, 3, 8

[29] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024. 1

[30] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024. 9

[31] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364, 2019. 2

[32] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 5

[33] Wenyi Hong, Weihan Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*, 2024. 1, 7

[34] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 6, 7, 8

[35] Qing Jiang, Gen Luo, Yuqin Yang, Yuda Xiong, Yihao Chen, Zhaoyang Zeng, Tianhe Ren, and Lei Zhang. Chatrex: Taming multimodal llm for joint perception and understanding. *arXiv preprint arXiv:2411.18363*, 2024. 2, 3, 5, 6, 7, 8

[36] Qing Jiang, Xingyu Chen, Zhaoyang Zeng, Junzhi Yu, and Lei Zhang. Rex-thinker: Grounded object referring via chain-of-thought reasoning. *arXiv preprint arXiv:2506.04034*, 2025. 2

[37] Qing Jiang, Lin Wu, Zhaoyang Zeng, Tianhe Ren, Yuda Xiong, Yihao Chen, Qin Liu, and Lei Zhang. Referring to any person. *arXiv preprint arXiv:2503.08507*, 2025. 8

[38] Yang Jiao, Shaoxiang Chen, Zequn Jie, Jingjing Chen, Lin Ma, and Yu-Gang Jiang. Lumen: Unleashing versatile vision-centric capabilities of large multimodal models. *Advances in Neural Information Processing Systems*, 37: 81461–81488, 2024. 3, 6

[39] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1780–1790, 2021. 6, 2

[40] Dimosthenis Karatzas, Lluis Gomez-Bigorda, Anguelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. Icdar 2015 competition on robust reading. In *2015 13th international conference on document analysis and recognition (ICDAR)*, pages 1156–1160. IEEE, 2015. 2

[41] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. 2

[42] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *European conference on computer vision*, pages 235–251. Springer, 2016. 9

[43] Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multimodal inputs and outputs. In *International Conference on Machine Learning*, pages 17283–17300. PMLR, 2023. 3

[44] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017. 2

[45] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 2, 3

[46] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *Advances in Neural Information Processing Systems*, 37:87874–87907, 2024. 1, 7

[47] Byung-Kwan Lee, Beomchan Park, Chae Won Kim, and Yong Man Ro. Collavo: Crayon large language and vision model. *arXiv preprint arXiv:2402.11248*, 2024. 4

[48] 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. 3, 8

[49] Chunyuan Li, Haotian Liu, Liunian Li, Pengchuan Zhang, Jyoti Aneja, Jianwei Yang, Ping Jin, Houdong Hu, Zicheng Liu, Yong Jae Lee, et al. Elevater: A benchmark and toolkit for evaluating language-augmented visual models. *Advances in Neural Information Processing Systems*, 35: 9287–9301, 2022. 7

[50] 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 *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 3

[51] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10975, 2022. 6

[52] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *European conference on computer vision*, pages 280–296. Springer, 2022. 5

[53] Long Lian, Yifan Ding, Yunhao Ge, Sifei Liu, Hanzi Mao, Boyi Li, Marco Pavone, Ming-Yu Liu, Trevor Darrell, Adam Yala, and Yin Cui. Describe anything: Detailed localized image and video captioning. *arXiv preprint arXiv:2504.16072*, 2025. 7

[54] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 6, 7

[55] Weifeng Lin, Xinyu Wei, Ruichuan An, Peng Gao, Bocheng Zou, Yulin Luo, Siyuan Huang, Shanghang Zhang, and Hongsheng Li. Draw-and-understand: Leveraging visual prompts to enable mllms to comprehend what you want. *arXiv preprint arXiv:2403.20271*, 2024. 7, 2

[56] Weifeng Lin, Xinyu Wei, Ruichuan An, Tianhe Ren, Tingwei Chen, Renrui Zhang, Ziyu Guo, Wentao Zhang, Lei Zhang, and Hongsheng Li. Perceive anything: Recognize, explain, caption, and segment anything in images and videos. *arXiv preprint arXiv:2506.05302*, 2025. 7

[57] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 1

[58] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 23592–23601, 2023. 2

[59] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1, 3, 7

[60] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26296–26306, 2024. 1, 3, 7, 8

[61] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 3

[62] Junzhuo Liu, Xuzheng Yang, Weiwei Li, and Peng Wang. Finecops-ref: A new dataset and task for fine-grained compositional referring expression comprehension. *arXiv preprint arXiv:2409.14750*, 2024. 2

[63] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pages 38–55. Springer, 2024. 6, 8

[64] Yuliang Liu, Hao Chen, Chunhua Shen, Tong He, Lianwen Jin, and Liangwei Wang. Abcnet: Real-time scene text spotting with adaptive bezier-curve network. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9809–9818, 2020. 2

[65] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multimodal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024. 9

[66] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024. 9

[67] Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. Visualrft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025. 3

[68] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023. 9

[69] Chuofan Ma, Yi Jiang, Jiannan Wu, Zehuan Yuan, and Xiaojuan Qi. Groma: Localized visual tokenization for grounding multimodal large language models. In *European Conference on Computer Vision*, pages 417–435. Springer, 2024. 2, 3, 5, 8

[70] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016. 2

[71] Nibal Nayef, Yash Patel, Michal Busta, Pinaki Nath Chowdhury, Dimosthenis Karatzas, Wafa Khlif, Jiri Matas, Umapada Pal, Jean-Christophe Burie, Cheng-lin Liu, et al. Icdar2019 robust reading challenge on multi-lingual scene text detection and recognition—rrc-mlt-2019. In *2019 International conference on document analysis and recognition (ICDAR)*, pages 1582–1587. IEEE, 2019. 2

[72] Changdae Oh, Hyeji Hwang, Hee-young Lee, YongTaek Lim, Geunyoung Jung, Jiyoung Jung, Hosik Choi, and Kyungwoo Song. Blackvip: Black-box visual prompting for robust transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24224–24235, 2023. 4

[73] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3170–3180, 2023. 8

[74] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023. 1

[75] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. 3, 7, 2

[76] Birgit Pfitzmann, Christoph Auer, Michele Dolfi, Ahmed S Nassar, and Peter Staar. Doclaynet: A large humanannotated dataset for document-layout segmentation. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 3743–3751, 2022. 2

[77] Khoi Pham, Kushal Kafle, Zhe Lin, Zhihong Ding, Scott Cohen, Quan Tran, and Abhinav Shrivastava. Learning to predict visual attributes in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13018–13028, 2021. 3, 2

[78] Renjie Pi, Jiahui Gao, Shizhe Diao, Rui Pan, Hanze Dong, Jipeng Zhang, Lewei Yao, Jianhua Han, Hang Xu, Lingpeng Kong, et al. Detgpt: Detect what you need via reasoning. *arXiv preprint arXiv:2305.14167*, 2023. 2

[79] Vignesh Ramanathan, Anmol Kalia, Vladan Petrovic, Yi Wen, Baixue Zheng, Baishan Guo, Rui Wang, Aaron Marquez, Rama Kovvuri, Abhishek Kadian, et al. Paco: Parts and attributes of common objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7141–7151, 2023. 8, 2

[80] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 6

[81] Tianhe Ren, Yihao Chen, Qing Jiang, Zhaoyang Zeng, Yuda Xiong, Wenlong Liu, Zhengyu Ma, Junyi Shen, Yuan Gao, Xiaoke Jiang, et al. Dino-x: A unified vision model for open-world object detection and understanding. *arXiv preprint arXiv:2411.14347*, 2024. 8

[82] Karthick Panner Selvam. Why large language models fail at precision regression, 2025. 1

[83] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Unleashing chain-of-thought reasoning in multi-modal language models. *CoRR*, 2024. 2

[84] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123*, 2018. 2

[85] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8430–8439, 2019. 6, 2

[86] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 1, 3

[87] Haozhan Shen, Kangjia Zhao, Tiancheng Zhao, Ruochen Xu, Zilun Zhang, Mingwei Zhu, and Jianwei Yin. Zoomeye: Enhancing multimodal llms with human-like zooming capabilities through tree-based image exploration. *arXiv preprint arXiv:2411.16044*, 2024. 2

[88] Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*, 2025. 3, 6

[89] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36: 38154–38180, 2023. 3

[90] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023. 3

[91] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 6, 8

[92] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities, 2023. 1

[93] 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. 1

[94] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, 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*, 2023. 3

[95] Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. Coco-text: Dataset and benchmark for text detection and recognition in natural images. *arXiv preprint arXiv:1601.07140*, 2016. 8

[96] Jiaqi Wang, Pan Zhang, Tao Chu, Yuhang Cao, Yujie Zhou, Tong Wu, Bin Wang, Conghui He, and Dahua Lin. V3det: Vast vocabulary visual detection dataset. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19844–19854, 2023. 6, 2

[97] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International conference on machine learning*, pages 23318–23340. PMLR, 2022. 1

[98] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1, 3, 8

[99] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, and Jifeng Dai. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks, 2023. 4

[100] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 3

[101] Jiannan Wu, Muyan Zhong, Sen Xing, Zeqiang Lai, Zhaoyang Liu, Zhe Chen, Wenhai Wang, Xizhou Zhu, Lewei Lu, Tong Lu, et al. Visionllm v2: An end-to-end generalist multimodal large language model for hundreds of vision-language tasks. *Advances in Neural Information Processing Systems*, 37:69925–69975, 2024. 3, 7

[102] Mingrui Wu, Xinyue Cai, Jiayi Ji, Jiale Li, Oucheng Huang, Gen Luo, Hao Fei, Guannan Jiang, Xiaoshuai Sun, and Rongrong Ji. Controlmllm: Training-free visual prompt learning for multimodal large language models. *Advances in Neural Information Processing Systems*, 37:45206–45234, 2024. 4

[103] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. *arXiv preprint arXiv:2312.14135*, 2023. 3

[104] Tung-Yu Wu, Sheng-Yu Huang, and Yu-Chiang Frank Wang. Dora: 3d visual grounding with order-aware referring. *CoRR*, 2024. 4

[105] Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4818–4829, 2024. 1

[106] An Yan, Zhengyuan Yang, Junda Wu, Wanrong Zhu, Jianwei Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Julian McAuley, Jianfeng Gao, et al. List items one by one: A new data source and learning paradigm for multimodal llms. *arXiv preprint arXiv:2404.16375*, 2024. 4

[107] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. 1

[108] Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*, 2023. 4

[109] Yiyang Yao, Peng Liu, Tiancheng Zhao, Qianqian Zhang, Jiajia Liao, Chunxin Fang, Kyusong Lee, and Qing Wang. How to evaluate the generalization of detection? a benchmark for comprehensive open-vocabulary detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6630–6638, 2024. 7

[110] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023. 3, 7, 8

[111] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *European conference on computer vision*, pages 69–85. Springer, 2016. 2

[112] Zhihan Yu and Ruifan Li. Revisiting counterfactual problems in referring expression comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13438–13448, 2024. 2

[113] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28202–28211, 2024. 7, 2

[114] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 9

[115] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6720–6731, 2019. 2

[116] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. 3

[117] Yufei Zhan, Hongyin Zhao, Yousong Zhu, Fan Yang, Ming Tang, and Jinqiao Wang. Griffon-g: Bridging vision-language and vision-centric tasks via large multimodal models. *arXiv preprint arXiv:2410.16163*, 2024. 3, 6

[118] Yufei Zhan, Shurong Zheng, Yousong Zhu, Hongyin Zhao, Fan Yang, Ming Tang, and Jinqiao Wang. Griffon v2: Advancing multimodal perception with high-resolution scaling and visual-language co-referring. *arXiv preprint arXiv:2403.09333*, 2024. 6

[119] Yufei Zhan, Yousong Zhu, Zhiyang Chen, Fan Yang, Ming Tang, and Jinqiao Wang. Griffon: Spelling out all object locations at any granularity with large language models. In *European Conference on Computer Vision*, pages 405–422. Springer, 2025. 3

[120] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 6

[121] Hao Zhang, Hongyang Li, Feng Li, Tianhe Ren, Xueyan Zou, Shilong Liu, Shijia Huang, Jianfeng Gao, Leizhang, Chunyuan Li, et al. Llava-grounding: Grounded visual chat with large multimodal models. In *European Conference on Computer Vision*, pages 19–35. Springer, 2024. 2

[122] Haotian Zhang, Haoxuan You, Philipp Dufter, Bowen Zhang, Chen Chen, Hong-You Chen, Tsu-Jui Fu, William Yang Wang, Shih-Fu Chang, Zhe Gan, et al. Ferret-v2: An improved baseline for referring and grounding with large language models. *arXiv preprint arXiv:2404.07973*, 2024. 7

[123] Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Change Loy Chen, and Shuicheng Yan. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. In *NeurIPS*, 2024. 4

[124] Yichi Zhang, Yinpeng Dong, Siyuan Zhang, Tianzan Min, Hang Su, and Jun Zhu. Exploring the transferability of visual prompting for multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26562–26572, 2024. 4

[125] Liang Zhao, En Yu, Zheng Ge, Jinrong Yang, Haoran Wei, Hongyu Zhou, Jianjian Sun, Yuang Peng, Runpei Dong, Chunrui Han, et al. Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning. *arXiv preprint arXiv:2307.09474*, 2023. 7

[126] Tiancheng Zhao, Peng Liu, Xuan He, Lu Zhang, and Kyusong Lee. Real-time transformer-based open-vocabulary detection with efficient fusion head. *arXiv preprint arXiv:2403.06892*, 2024. 5, 6

[127] Tiancheng Zhao, Qianqian Zhang, Kyusong Lee, Peng Liu, Lu Zhang, Chunxin Fang, Jiajia Liao, Kelei Jiang, Yibo Ma, and Ruochen Xu. Omchat: A recipe to train multimodal language models with strong long context and video understanding. *arXiv preprint arXiv:2407.04923*, 2024. 2

[128] Yang Zhao, Zhijie Lin, Daquan Zhou, Zilong Huang, Jiashi Feng, and Bingyi Kang. Bubogpt: Enabling visual grounding in multi-modal llms. *arXiv preprint arXiv:2307.08581*, 2023. 3

[129] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1, 3

# VLM-FO1: Bridging the Gap Between High-Level Reasoning and Fine-Grained Perception in VLMs

## Supplementary Material

## 7. Comprehensive Dataset Overview for VLM-FO1

This section details the comprehensive datasets utilized in our two-stage training methodology in Table 10.

## 8. Additional Visualizations

This section presents a more comprehensive set of visual results to further illustrate the performance and capabilities of our VLM-FO1 model. We provide diverse examples in Figure [3,4,5,6,7,8,9,10,11,12], showcasing key aspects of our method's behavior across various scenarios, offering deeper insights beyond the examples included in the main paper.

| Stage | Task Type | Datasets |
|---|---|---|
| Region-Language Alignment | OD | COCO,O365[85],V3Det[96] |
| | Grounding | GOLDG[39] |
| | Region Caption | Rexverse-2M[35] |
| Perception SFT | OD | COCO,V3Det,O365,VAW[77] VisDrone2019[27],LVIS[31] |
| | REC | Refcoco/+/g[41, 70, 111] finecopsref[62],grefcoco[58] CREC[112] |
| | Grounding | GRIT[75] |
| | Region Caption | PACO[79],VG[44],Osprey[113] shareGPT4v[15],Rexverse-2M |
| | Region Reasoning | VisualCoT[83],HumanRef-CoT[36] |
| | Counting | COCO,LVIS,HumanRef-CoT CrowdHuman[84],TallyQA[2] |
| | Region QA | Osprey,VCR[115] MDVP[55],DoclayNet[76] |
| | OCR | MLT2019[71],ICDAR15[40] CurvedSynText150k[64] |
| | VLM-Instruction | OmChat-SFT[127] |

Table 10. **A summary of the datasets used in our two-stage training process.**



Figure 3. Visualization of VLM-FO1's perception abilities on OD task.

Figure 4. Visualization of VLM-FO1's perception abilities on REC task.



Figure 5. Visualization of VLM-FO1's perception abilities on object counting task.
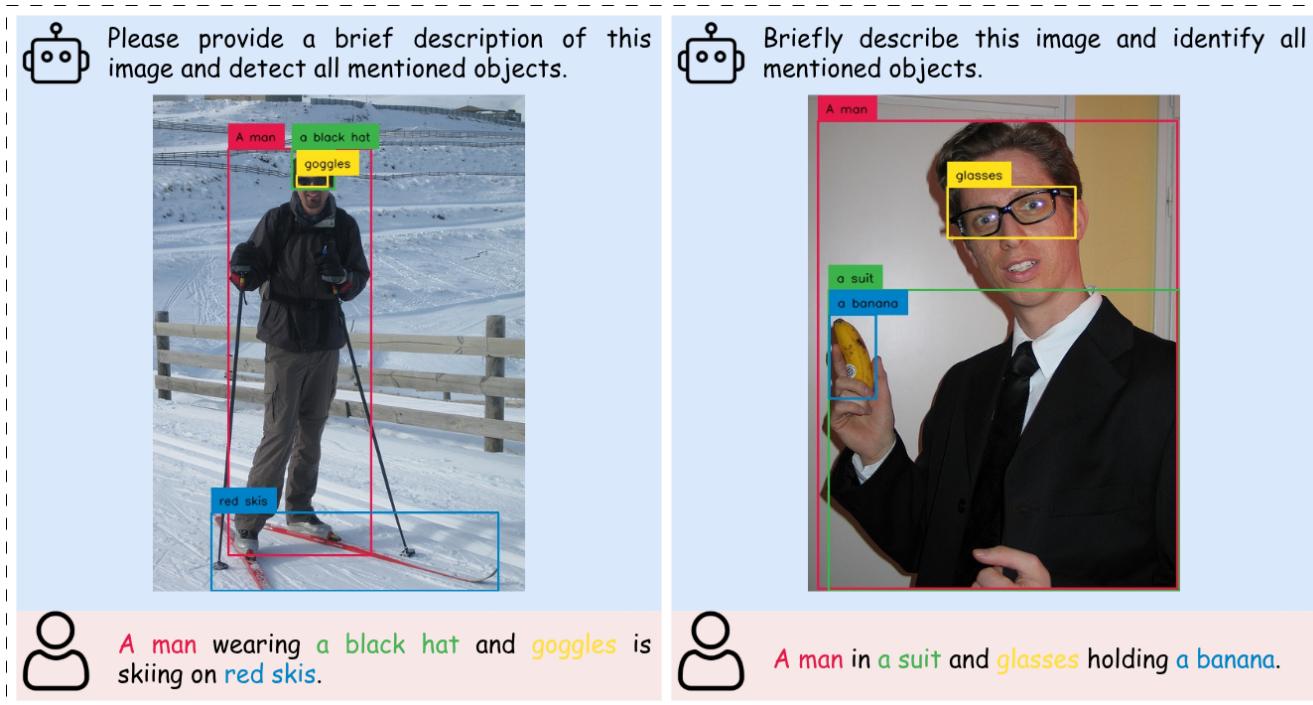
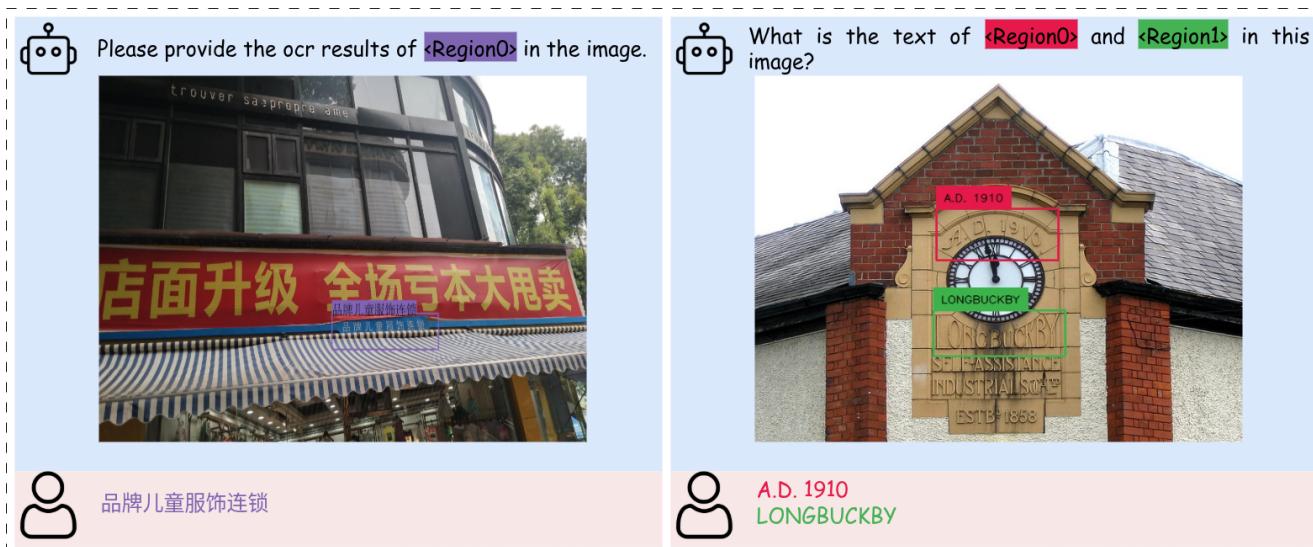Figure 6. Visualization of VLM-FO1's perception abilities on grounding task.



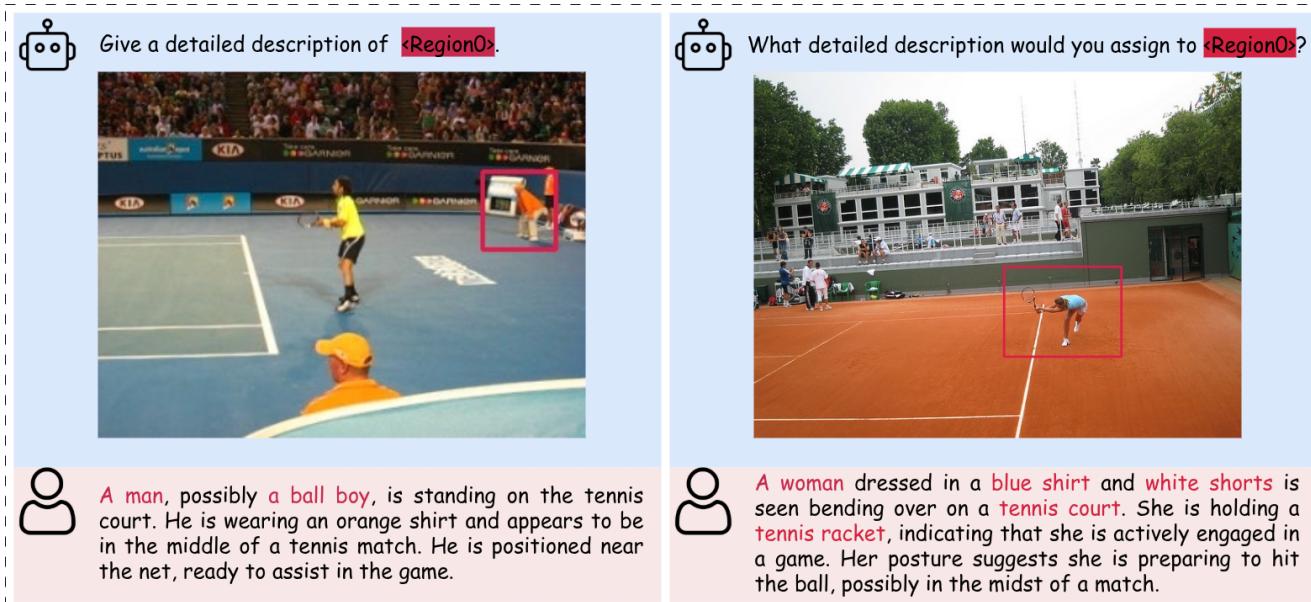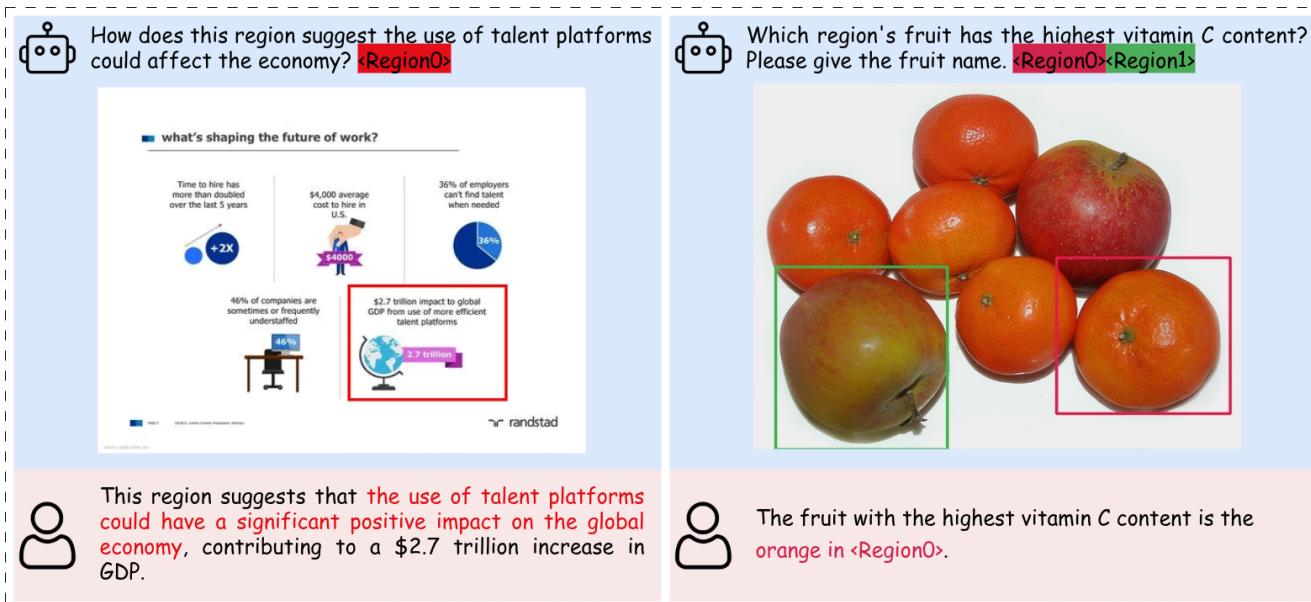Figure 7. Visualization of VLM-FO1's perception abilities on OCR task.

Give a detailed description of <Region0>.

A man, possibly a ball boy, is standing on the tennis court. He is wearing an orange shirt and appears to be in the middle of a tennis match. He is positioned near the net, ready to assist in the game.

What detailed description would you assign to <Region0>?

A woman dressed in a blue shirt and white shorts is seen bending over on a tennis court. She is holding a tennis racket, indicating that she is actively engaged in a game. Her posture suggests she is preparing to hit the ball, possibly in the midst of a match.

Figure 8. Visualization of VLM-FO1's perception abilities on region caption task.



How does this region suggest the use of talent platforms could affect the economy? <Region0>

This region suggests that the use of talent platforms could have a significant positive impact on the global economy, contributing to a $2.7 trillion increase in GDP.

Which region's fruit has the highest vitamin C content? Please give the fruit name. <Region0><Region1>

The fruit with the highest vitamin C content is the orange in <Region0>.

Figure 9. Visualization of VLM-FO1's perception abilities on region VQA task.

Figure 10. Visualization of VLM-FO1's perception abilities on visual prompt OD task.



Figure 11. Visualization of VLM-FO1's perception abilities on visual region reasoning task.

For the reasoning process, review each object in the image associated with the question and identify those that align with the target label.\nPlease detect the person who is not wearing a tie.

OK, the user needs us to detect the person who is not wearing a tie. To accomplish this task, I need to break it down into the following steps:
Step 1: Find all the person wearing a tie.
Step 2: Remove the person from step 1 to get the final answer
# Step 1: Find all the person wearing a tie.
I see four people in this image, let's look at their ties and put a ⚠️ if they are wearing a tie and ❌ if they are not.
- <Region4>: He is not wearing a tie → ❌
- <Region5>: He is wearing a red tie → ⚠️
- <Region12>: He is not wearing a tie → ❌
- <Region8>: He is not wearing a tie → ❌
# Step 2: Remove the person from step 1 to get the final answer
Remove <Region5> and we get <Region4>, <Region12>, <Region8> that meets the description
# Summarize and Re-Check answer
Let's now recheck our answer and put ✅ for the target person and ❌ for others
- <Region4>: He is not wearing a tie → ✅
- <Region5>: He is wearing a tie → ❌
- <Region12>: He is not wearing a tie → ✅
- <Region8>: He is not wearing a tie → ✅

For the reasoning process, carefully examine all objects in the image connected to the question and determine which ones match the target label.\nIdentify the following objects in the image: the plate holding black doughnuts.

OK, the user needs us to detect the plate holding black doughnuts. To accomplish this task, I need to break it down into the following steps:
- Step 1: Find all the plates in the image.
- Step 2: From the plates in Step 1, identify the one holding black doughnuts.
# Step 1: Find all the plates in the image.
There are three plates in the image:
- <Region3>: A white plate with a glazed donut and a brown-colored item on it.
- <Region5>: A white plate with two glazed donuts on it.
- <Region2>: A plate with a black and white checkered pattern, holding a chocolate-dipped donut.
# Step 2: From the plates in Step 1, identify the one holding black doughnuts.
- <Region3>: Holds a glazed donut and a brown-colored item, not black doughnuts → ❌
- <Region5>: Holds two glazed donuts, not black doughnuts → ❌
- <Region2>: Holds a chocolate-dipped donut, which is black in color, so it is the plate holding black doughnuts → ✅
# Summarize and Re-Check answer
Let's now recheck our answer and put ✅ for the target plate and ❌ for others:
- <Region3>: Holds a glazed donut and a brown-colored item, not black doughnuts → ❌
- <Region5>: Holds two glazed donuts, not black doughnuts → ❌
- <Region2>: Holds a chocolate-dipped donut, which is black in color → ✅

Figure 12. Visualization of VLM-FO1's perception abilities on visual region reasoning task.