

Visual-RFT: Visual Reinforcement Fine-Tuning



Ziyu Liu^{1,2*} Zeyi Sun^{1,2*} Yuhang Zang^{2✉} Xiaoyi Dong^{2,3} Yuhang Cao²

Haodong Duan² Dahua Lin^{2,3} Jiaqi Wang^{2✉}

¹Shanghai Jiaotong University ²Shanghai Artificial Intelligence Laboratory

³The Chinese University of Hong Kong

{liuziyu77, szy2023}@sjtu.edu.cn, {zangyuhang, wangjiaqi}@pjlab.org.cn

<https://github.com/Liuziyu77/Visual-RFT>

Abstract

Reinforcement Fine-Tuning (RFT) in Large Reasoning Models like OpenAI o1 learns from feedback on its answers, which is especially useful in applications when fine-tuning data is scarce. Recent open-source work like DeepSeek-R1 demonstrates that reinforcement learning with verifiable reward is one key direction in reproducing o1. While the R1-style model has demonstrated success in language models, its application in multi-modal domains remains under-explored. This work introduces Visual Reinforcement Fine-Tuning (Visual-RFT), which further extends the application areas of RFT on visual tasks. Specifically, Visual-RFT first uses Large Vision-Language Models (LVLMs) to generate multiple responses containing reasoning tokens and final answers for each input, and then uses our proposed visual perception verifiable reward functions to update the model via the policy optimization algorithm such as Group Relative Policy Optimization (GRPO). We design different verifiable reward functions for different perception tasks, such as the Intersection over Union (IoU) reward for object detection. Experimental results on fine-grained image classification, few-shot object detection, reasoning grounding, as well as open-vocabulary object detection benchmarks show the competitive performance and advanced generalization ability of Visual-RFT compared with Supervised Fine-tuning (SFT). For example, Visual-RFT improves accuracy by 24.3% over the baseline in one-shot fine-grained image classification with around 100 samples. In few-shot object detection, Visual-RFT also exceeds the baseline by 21.9 on COCO's two-shot setting and 15.4 on LVIS. Our Visual-RFT represents a paradigm shift in fine-tuning LVLMs, offering a data-efficient, reward-driven approach that enhances reasoning and adaptability for domain-specific tasks.

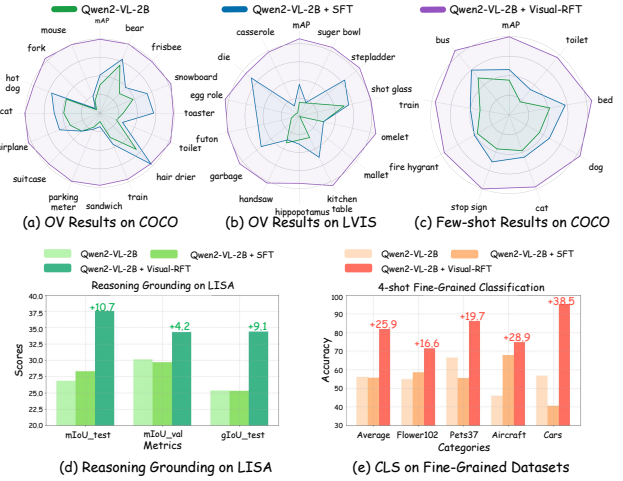


Figure 1. Our **Visual Reinforcement Fine-Tuning (Visual-RFT)** performs better than previous Supervised Fine-Tuning (SFT) on a variety of tasks, such as Open Vocabulary(OV)/Few-shot Detection, Reasoning Grounding, and Fine-grained Classification.

1. Introduction

Large Reasoning Models (LRMs) such as OpenAI o1 [7] represent frontier AI models designed to spend more time “thinking” before answering, and achieving excellent reasoning abilities. One impressive capability of OpenAI o1 is Reinforcement Fine-Tuning (RFT)¹, which efficiently fine-tune the model with merely dozens to thousands of samples to excel at domain-specific tasks. Although the implementation details of o1 are not publicly available, recent open-source studies like DeepSeek R1 [4] reveal one promising direction in reproducing o1 is Verifiable Rewards [4, 12, 37]: the reward score in reinforcement learning is directly determined by pre-defined rules, rather than predicted by the separate reward model [17, 26, 45] trained on preference data.

A primary distinction between the RFT and Previous Supervised Fine-Tuning (SFT) lies in data efficiency. Previous

* Equal contribution. ✉ Corresponding author.

¹<https://openai.com/form/rft-research-program>

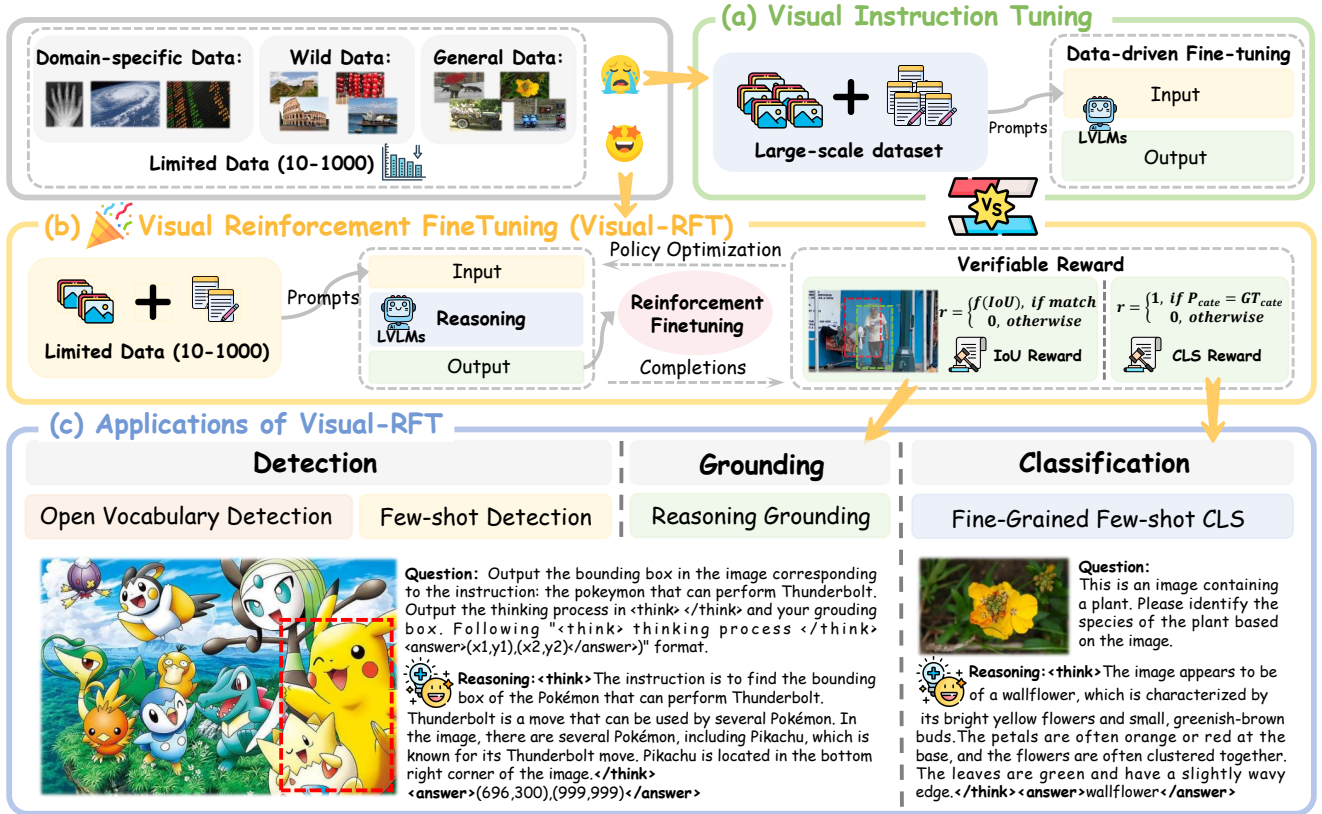


Figure 2. **Overview of Visual-RFT.** Compared to the (a) Visual Instruction Tuning that is data-hungry, (b) our Visual Reinforcement Fine-Tuning (Visual-RFT) is more data efficient with limited data. (c) We successfully empower Large Vision-Language Models (LVLMs) with RFT on a series of multi-modal tasks, and present examples of the model’s reasoning process at the bottom.

SFT paradigm (see Fig. 2 (a)) directly imitates the “ground-truth” answers provided in the high-quality, curated data, thus relying on a large amount of training data. By contrast, RFT evaluates the model’s responses and adjusts based on whether they’re correct, helping it learn through trial and error. Thus, RFT is particularly useful in domains where data is scarce [7, 24]. However, a previous common sense is that RFT is applied merely in tasks like scientific (e.g., mathematics) and code generation. That’s because math and coding exhibit clear and objective final answers or test cases, making their rewards relatively straightforward to verify. In this paper, we demonstrate that RFT can be applied beyond math and code domains to visual perception tasks. Specifically, we introduce Visual Reinforcement Fine-Tuning (Visual-RFT), which successfully extends RFT to empower Large Vision-Language Models (LVLMs) in various multi-modal tasks (see Fig. 1), such as few-shot classification and open-vocabulary object detection.

To extend RFT on visual tasks, we present the implementation details of Visual-RFT in Fig. 2 (b). For each input, Visual-RFT uses Large Vision-Language Models (LVLMs) to generate multiple responses (trajectories) that contain the reasoning tokens and final answers. Crucially, we define task-specific, rule-based verifiable reward functions to

guide policy optimization, such as GRPO [31], in updating the model. For instance, we propose the Intersection over Union (IoU) reward for the object detection task. Our Visual-RFT contrasts with SFT, which relies on memorizing correct answers. Instead, our approach explores different possible solutions and learns to optimize for a desired outcome defined by our verified reward function. It’s about discovering what works best, not just mimicking predefined answers. Our approach shifts the training paradigm from data scaling in SFT to the strategic design of variable reward functions tailored to specific multi-modal tasks. As shown in Fig. 2 (c), the synergistic combination of verifiable rewards and visual perception abilities (e.g., detection, grounding, classification) allows our model to achieve rapid and data-efficient mastery of new concepts, facilitated by a detailed reasoning process.

We validate the effectiveness of Visual-RFT on the following tasks. In fine-grained image classification, the model utilizes its advanced reasoning capabilities to analyze fine-grained categories with high precision. In the one-shot setting with extremely limited data (e.g., around 100 samples), Visual-RFT boosts the accuracy with 24.3% over the baseline, while SFT dropped by 4.3%. In few-shot experiments, Visual-RFT also demonstrates exceptional per-

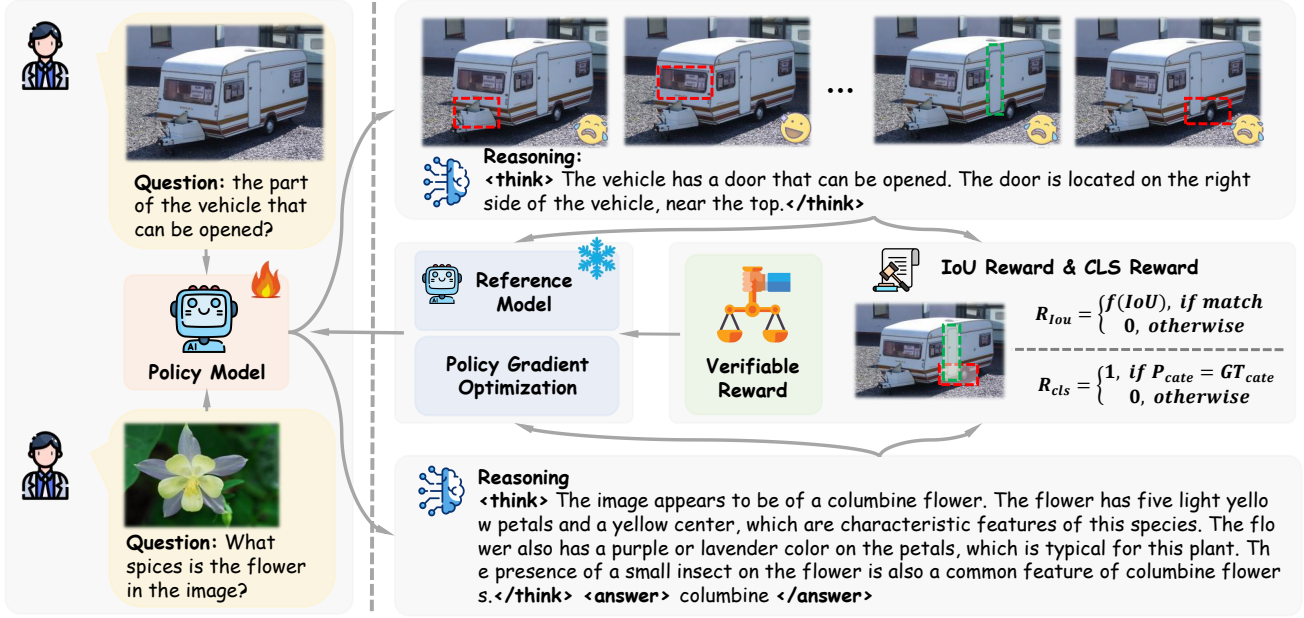


Figure 3. **Framework of Visual-RFT.** Given the question and visual image inputs, the policy model generates multiple responses containing reasoning steps. Then the verifiable reward such as IoU reward and CLS reward is used with the policy gradient optimization algorithm to update the policy model.

formance with minimal training data, showcasing superior few-shot learning capabilities compared to SFT. In reasoning grounding, Visual-RFT excels in the LISA [11] dataset, which heavily relies on reasoning, outperforming specialized models like GroundedSAM [18]. Furthermore, in open vocabulary object detection, Visual-RFT quickly transfers recognition capabilities to new categories, including rare categories in LVIS [5], showing strong generalization. Specifically, the 2B model achieves mAP improvements from 9.8 to 31.3 on new classes of COCO [15] and from 2.7 to 20.7 on selected rare classes of LVIS [5]. These diverse visual perception tasks not only highlight Visual-RFT’s robust generalization capabilities in visual recognition but also underscore the crucial role of reinforcement learning in enhancing visual perception and reasoning.

In summary, our key contributions are as follows:

- (1) We introduce Visual Reinforcement Fine-tuning (Visual-RFT), which extends reinforcement learning with verifiable rewards on visual perception tasks that are effective with limited data for fine-tuning.
- (2) We design different verifiable rewards for different visual tasks that enable efficient, high-quality reward computation at a negligible cost. This allows the seamless transfer of DeepSeek RL’s style reinforcement learning to LVLMs.
- (3) We conduct extensive experiments on various visual perception tasks, including fine-grained image classification, few-shot object detection, reasoning grounding, and open vocabulary object detection. On all the settings, Visual-RFT achieves remarkable performance improvements, significantly surpassing the supervised fine-tuning baselines.

- (4) We fully *open-source* the training code, training data, and evaluation scripts on [Github](#) to facilitate further research.

2. Related Work

Large Vision Language Models (LVLMs) like GPT-4o [23] achieves excellent visual understanding ability by integrating both visual and textual data. This integration enhances the models’ ability to understand complex multi-modal inputs and enables more advanced AI systems [13, 16, 38, 47] capable of processing and responding to both images and text. Generally, the training of LVLMs involves two steps: (a) pre-training and (b) post-training which contains supervised fine-tuning and reinforcement learning. Post-training is crucial in improving the model’s response quality, instruction following, and reasoning abilities. While there has been significant research on using reinforcement learning to enhance LLMs during post-training [1, 3, 25, 28, 32, 33, 36, 40, 44, 52, 53], the progress for LVLMs has been slower. In this paper, we propose Visual-RFT, which used GRPO-based reinforcement algorithms and verifiable reward during the post-training phase to enhance the model’s visual perception and reasoning capabilities.

Reinforcement Learning Recently, with the emergence of reasoning models like OpenAI’s o1 [7], the research focus in Large Language Models (LLMs) has increasingly shifted towards enhancing the models’ reasoning capabilities through reinforcement learning (RL) techniques. Stud-

Table 1. **Prompts used to construct the dataset.** We have listed the detection prompt and classification prompt separately.

Detection Prompt: Detect all objects belonging to the category '{category}' in the image, and provide the bounding boxes (between 0 and 1000, integer) and confidence (between 0 and 1, with two decimal places). If no object belonging to the category '{category}' in the image, return 'No Objects'. Output the thinking process in <think> </think> and final answer in <answer> </answer> tags. The output answer format should be as follows: <think> ... </think><answer>['Position': [x1, y1, x2, y2], 'Confidence': number, ...]</answer> Please strictly follow the format.

Classification Prompt: This is an image containing a plant. Please identify the species of the plant based on the image. Output the thinking process in <think> </think> and final answer in <answer> </answer> tags. The output answer format should be as follows: <think> ... </think> <answer>species name</answer> Please strictly follow the format.

ies have explored improving LLMs' performance in reasoning tasks such as solving mathematical problems [2, 20, 31, 39, 41] and coding [6, 8, 46, 48]. A notable breakthrough in this area is Deepseek-R1-Zero [4], which introduced a new approach to achieving robust reasoning capabilities using RL merely, eliminating the supervised fine-tuning (SFT) stage. However, current research on RL-based reasoning has largely been confined to the language domain, with limited exploration of its application in multi-modal settings. For LVLMs, RL has primarily been used for tasks like mitigating hallucinations and aligning models with human preference [19, 34, 35, 42, 43, 49–51], but there remains a significant gap in research focusing on enhancing reasoning and visual perception of Large Vision Language Models. To address this gap, our work introduces a novel reinforcement fine-tuning strategy Visual-RFT, applying verifiable rewards with GRPO-based [31] RL to a broad range of visual perception tasks. Our approach aims to improve the performance of LVLMs in processing various visual tasks, especially when the fine-tuning data is limited.

3. Methodology

3.1. Preliminary

Reinforcement Learning with Verifiable Rewards. Reinforcement Learning with Verifiable Rewards (RLVR) [4, 12, 37] is a novel training approach designed to enhance language models in tasks with objectively verifiable outcomes, such as math and coding. Unlike previous Reinforcement Learning from Human Feedback (RLHF) [17, 26, 45], which relies on a trained reward model, RLVR instead uses a direct verification function to assess correctness. This approach simplifies the reward mechanism while maintaining strong alignment with the task's inherent cor-

rectness criteria. Given the input question q , the policy model π_θ generates responses o and receives the verifiable reward. More specifically, RLVR optimizes the following objective:

$$\max_{\pi_\theta} \mathbb{E}_{o \sim \pi_\theta(q)} [R_{\text{RLVR}}(q, o)] \quad (1)$$

$$= [R(q, o) - \beta \text{KL}[\pi_\theta(o|q) \parallel \pi_{\text{ref}}(o|q)]], \quad (2)$$

where π_{ref} is the reference model before optimization, R is the verifiable reward function, and β is the hyperparameters to control the KL-divergence. The verifiable reward function R takes the question and output pair (q, o) as inputs, and checks if the ground-truth answer remains the same as the prediction o :

$$R(q, o) = \begin{cases} 1, & \text{if } o = \text{ground truth,} \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

DeepSeek R1-Zero and GRPO. The DeepSeek R1-Zero algorithm eliminates dependence on supervised fine-tuning (SFT) by employing reinforcement learning for training, specifically through its Group Relative Policy Optimization (GRPO) framework. Different from reinforcement learning algorithms such as PPO [30] that require a critic model to evaluate policy performance, GRPO compares groups of candidate responses directly, eliminating the need for an additional critic model. For a given question q , GRPO first generates G distinct responses $\{o_1, o_2, \dots, o_G\}$ from the current policy $\pi_{\theta_{\text{old}}}$. Then GRPO takes actions based on these responses and denotes the obtained rewards as $\{r_1, r_2, \dots, r_G\}$. By computing their mean and standard deviation for normalization, GRPO determines the relative quality of these responses:

$$A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\})}, \quad (4)$$

where A_i represents the relative quality of the i -th answer. GRPO encourages the model to favor better answers with a high reward value within the group.

3.2. Visual-RFT

The framework of Visual-RFT is shown in Fig. 3. The multi-modal input data from the user consists of images and questions. The policy model π_θ outputs a reasoning process and generates a group of responses based on the input. Each response is passed through a verifiable reward function to compute the reward. After group computation of the rewards for each output, the quality of each response is evaluated and used to update the policy model. To ensure the stability of the policy model training, Visual-RFT uses KL divergence to limit the difference between the policy model and the reference model. We will further discuss how to design the verifiable reward for visual tasks in Sec. 3.2.1, and the data preparation steps in Sec. 3.2.2

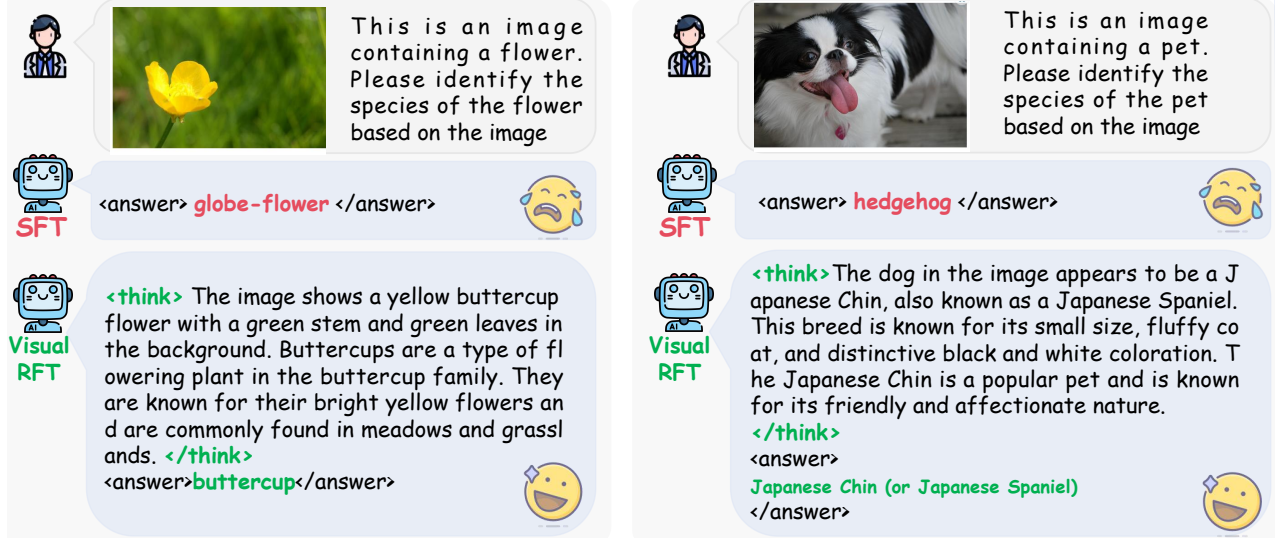


Figure 4. **Qualitative results of Fine-Grained Image Classification.** The thinking process significantly improves the reasoning ability of LVLMs, leading to higher image classification performance.

3.2.1. Verifiable Reward in Visual Perception

The reward model is a key step in reinforcement learning (RL) that aligns models with preference alignment algorithms, which can be as straightforward as a verification function that checks for exact matches between predictions and ground-truth answers. The RL training process in the recent DeepSeek-R1 [4] model achieves a significant improvement in the model’s reasoning ability through the verifiable reward design. To transfer this strategy to the visual domain, we design different rule-based verifiable reward functions for various visual perception tasks.

IoU Reward in Detection Tasks. For the detection task, the model’s output consists of bounding boxes (bbox) and corresponding confidences. The reward function for the detection task should adequately consider the Intersection-over-Union (IoU) metric, which is used to compute the mean Average Precision (mAP) in evaluation. Therefore, we design an IoU and confidence-based reward function R_d . First, for the model’s output box and confidence, we sort these boxes based on their confidence, denoted as $\{b_1, b_2, \dots, b_n\}$. We then match each b_i with the ground truth bbox, $\{b_1^g, b_2^g, \dots, b_m^g\}$, and calculate the IoU, while setting an IoU threshold τ . Bounding boxes with an IoU below this threshold τ are considered invalid, and unmatched bboxes have an IoU of 0. After matching, we obtain the IoU and confidence for each box from the initial set, denoted as $\{iou_1 : c_1, iou_2 : c_2, \dots, iou_n : c_n\}$.

We then use these IoU results and confidence to construct our reward R_d . Our reward R_d consists of three parts, including the IoU reward, Confidence reward, and Format

reward:

$$R_d = R_{IoU} + R_{conf} + R_{format}. \quad (5)$$

The IoU reward R_{IoU} is the average IoU of all the bounding boxes in the model’s output,

$$R_{IoU} = \frac{iou_1 + iou_2 + \dots + iou_n}{n}. \quad (6)$$

The confidence reward R_{conf} is related to IoU. For each bounding box, if the iou_i is non-zero, indicating a successful match, the confidence reward for this single box r_c as the predicted confidence is computed as:

$$r_{ci} = \begin{cases} c_i & , \text{ if } iou_i \neq 0, \\ 1 - c_i & , \text{ if } iou_i = 0. \end{cases} \quad (7)$$

This means that for successfully matched boxes, the higher the confidence, the better. If the iou_i is zero, indicating a failed match, the lower the confidence reward r_c for this box, the better. The overall confidence reward R_{conf} for the model’s output is the average of the confidence rewards of all the bounding boxes in that output,

$$R_{conf} = \frac{\sum_{i=1}^n r_{ci}}{n}. \quad (8)$$

The format reward R_{format} is used to force the model prediction to meet the required HTML tag format of `<think>` and `<answer>` (will detailed in Sec. 3.2.2).

CLS Reward in Classification Tasks. In classification tasks, the reward function we use consists of two parts: accuracy reward R_{acc} and format reward R_{format} . The accuracy

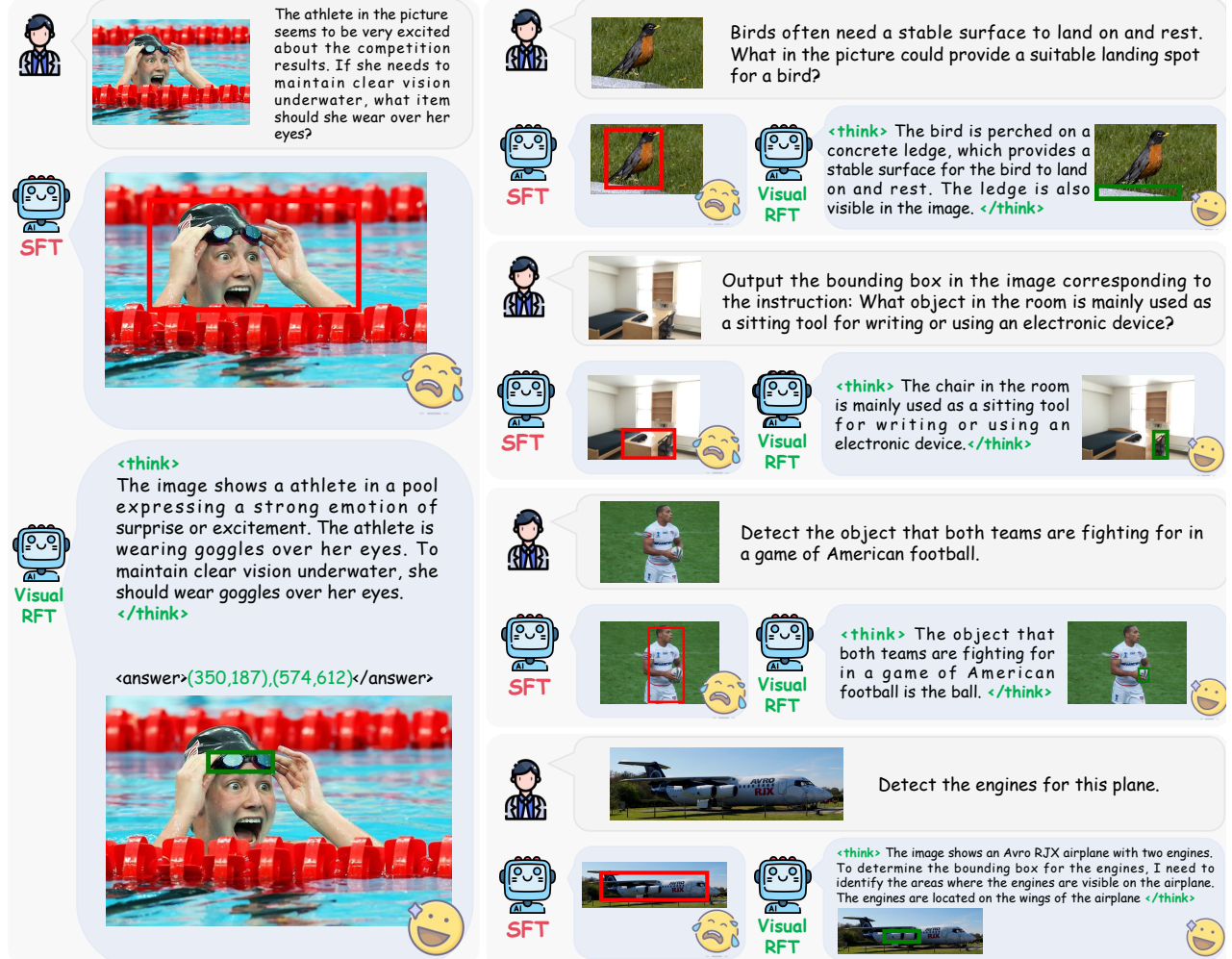


Figure 5. **Qualitative results of reasoning grounding on LISA [11]**. Thinking process significantly improves reasoning grounding ability with Visual-RFT.

reward is determined by comparing the model’s output class with the ground truth class, yielding a value of 1 for correct classification and 0 for incorrect classification:

$$R_{cls} = R_{acc} + R_{format}. \quad (9)$$

3.2.2. Data Preparation

To train the Visual-RFT on various visual perception tasks, we need to construct the multi-modal training dataset. Similar to DeepSeek-R1, to enhance the model’s reasoning ability and apply this ability to improve visual perception, Visual-RFT designed a prompt format to guide the model to output its reasoning process before providing the final answer. The prompts used for detection and classification tasks are shown in Tab 1.

During the training process, we use the format reward to guide the model to output its reasoning process and the

final answer in a structured format. The reasoning process is key to the model’s self-learning and improvement during reinforcement fine-tuning, while the reward determined by the answer directs the model’s optimization.

4. Experiments

4.1. Experimental Setup

Implementation Details Our method is adaptable to various visual perception tasks. We employ a few-shot learning approach, providing the model with a minimal number of samples for training. For the image classification and object detection task, we adopt a few-shot setting to evaluate the model’s fine-grained discriminative and recognition capability, applying reinforcement learning on limited data. Then, for the LISA [11] dataset focusing on reasoning grounding, which demands strong reasoning ability,

Table 2. **Few-shot results on Fine-grained Classification dataset.** We evaluated four fine-grained image classification datasets.

Models	Average	Flower102	Pets37	FGVC	Cars196
Qwen2-VL-2B	56.0	54.8	66.4	45.9	56.8
<i>one-shot</i>					
+ SFT	51.7	56.6	54.7	65.3	30.0
+ Visual-RFT	80.3	70.8	84.1	72.5	93.8
Δ	+24.3	+16.0	+17.7	+26.6	+37.0
<i>2-shot</i>					
+ SFT	58.8	60.3	65.6	68.9	40.2
+ Visual-RFT	83.5	75.8	87.5	75.3	95.4
Δ	+27.5	+21.0	+21.1	+29.4	+38.6
<i>4-shot</i>					
+ SFT	55.6	58.5	55.5	67.9	40.5
+ Visual-RFT	81.9	71.4	86.1	74.8	95.3
Δ	+25.9	+16.6	+19.7	+28.9	+38.5
<i>8-shot</i>					
+ SFT	60.3	59.6	71.4	69.2	40.9
+ Visual-RFT	85.1	77.7	90.2	75.9	96.5
Δ	+29.1	+22.9	+23.8	+30.0	+39.7
<i>16-shot</i>					
+ SFT	64.0	66.8	71.6	76.1	41.5
+ Visual-RFT	85.3	79.2	87.1	79.4	95.3
Δ	+29.3	+24.4	+20.7	+33.5	+38.5

ties, we train the model using Visual-RFT and assess its reasoning and perception performance. Lastly, for open-vocabulary object detection, we evaluate the model’s generalization capability by training the Qwen2-VL-2/7B [38] using Visual-RFT on a subdivided COCO dataset containing 65 base classes. We then test it on 15 novel classes from COCO and 13 rare classes from LVIS [5]. The model’s visual perception and reasoning abilities are assessed in an open-vocabulary detection setting. In our detection experiments, we first prompt the model to check whether the category is present in the image, then predict bound boxes for categories that exist in the images.

4.2. Few-Shot Classification

To demonstrate the extensive generalization ability of Visual-RFT in the visual domain, we conduct few-shot experiments on fine-grained image classification. We selected four datasets: Flower102 [22], Pets37 [27], FGVC-Aircraft [21], and Car196 [10], which contain dozens to hundreds of similar categories, adding significant difficulty

Table 3. **Few-Shot results on COCO dataset of 8 categories.** We conducted one-shot, 2-shot, 4-shot, 8-shot, and 16-shot experiments on 8 categories from the COCO dataset.

Models	mAP	bus	train	fire hydrant	stop sign	cat	dog	bed	toilet
<i>Qwen2-VL-2B</i>									
Baseline	19.6	19.0	15.8	25.8	18.4	29.9	23.2	14.6	9.8
<i>1-shot</i>									
+ SFT	19.5	18.3	17.4	23.1	18.2	28.0	23.4	17.3	10.4
+ Visual-RFT	33.6	23.4	35.7	39.1	23.8	54.3	42.5	19.5	30.8
Δ	+14.0	+4.4	+19.9	+13.3	+5.4	+24.4	+19.3	+4.9	+21.0
<i>2-shot</i>									
+ SFT	21.0	22.1	15.8	29.8	19.0	28.9	26.5	15.5	10.2
+ Visual-RFT	41.5	28.8	39.6	38.2	48.0	63.8	52.7	25.9	34.9
Δ	+21.9	+9.8	+23.8	+12.4	+29.6	+33.9	+29.5	+11.3	+25.1
<i>4-shot</i>									
+ SFT	25.2	25.4	23.6	26.6	21.5	35.6	30.6	18.4	19.9
+ Visual-RFT	40.6	30.0	40.6	45.7	35.0	60.9	44.9	24.6	43.1
Δ	+21.0	+11.0	+24.8	+19.9	+16.6	+31.0	+21.7	+10.0	+33.3
<i>8-shot</i>									
+ SFT	30.2	25.8	35.1	29.4	21.9	44.5	39.0	22.6	23.5
+ Visual-RFT	47.4	36.2	47.9	50.4	47.7	65.2	57.0	30.4	44.0
Δ	+27.8	+17.2	+32.1	+24.6	+29.3	+35.3	+33.8	+15.8	+34.2
<i>16-shot</i>									
+ SFT	31.3	24.0	35.9	32.0	23.6	39.8	40.6	27.5	26.8
+ Visual-RFT	46.8	36.2	44.4	51.4	48.5	66.6	56.2	27.6	43.4
Δ	+27.2	+17.2	+28.6	+25.6	+30.1	+36.7	+33.0	+13.0	+33.6
<i>Qwen2-VL-7B</i>									
Baseline	43.0	35.0	43.3	37.1	36.7	57.3	50.3	37.4	47.1
<i>4-shot</i>									
+ SFT	44.1	41.4	51.7	35.6	30.8	60.5	52.7	38.5	41.5
+ Visual-RFT	54.3	44.3	59.8	52.0	46.0	72.7	62.8	41.9	55.0
Δ	+11.3	+9.3	+16.5	+14.9	+9.3	+15.4	+12.5	+4.5	+7.9

to the classification task.

As shown in Tab. 2, with just one-shot of data, Visual-RFT already delivers a significant performance boost (+24.3%). In contrast, SFT shows a noticeable decline (-4.3%) with the same minimal amount of data. Under the 4-shot setting, the performance of SFT is still slightly lower than the baseline, while the reinforcement fine-tuned model with Visual-RFT achieves an average performance improvement of 25.9. Under the 8-shot and 16-shot settings, as the amount of data increases, SFT’s performance slightly exceeds the baseline. However, SFT still lags significantly behind the performance of the Visual-RFT. In Fig.4, we present some inference cases of the model after reinforcement fine-tuning when handling fine-grained classification tasks. These results not only demonstrate the strong generalization ability of Visual-RFT and its capacity to learn from limited data but also confirm that reinforcement fine-tuning, compared to SFT, leads to a genuine un-

Table 4. **Few-shot results on LVIS dataset of 6 rare categories.** We conducted 10-shot experiments on 6 rare categories from the LVIS dataset.

Models	mAP	horse buggy	die	kitchen table	omelet	papaya	stepladder
Qwen2-VL-2B	4.0	2.9	1.2	13.4	4.7	1.5	0.0
+ SFT	10.0	7.0	7.6	34.1	4.7	6.3	0.0
+ Visual-RFT	19.4	9.1	19.6	42.2	20.4	14.5	10.9
Δ	+15.4	+6.2	+18.4	+29.2	+15.7	+13.0	+10.9
Qwen2-VL-7B	15.4	19.7	21.9	14.5	18.1	18.5	0.0
+ SFT	27.6	26.9	21.9	49.7	29.2	25.2	12.7
+ Visual-RFT	33.8	26.2	27.8	70.6	23.5	21.2	29.3
Δ	+18.4	+6.5	+5.9	+56.1	+5.4	+2.7	+29.3

Table 5. **Few-shot results on MG dataset of 5 categories.** By introducing out-of-domain data, we increased the difficulty of model recognition and reasoning, further demonstrating the strong generalization ability of reinforcement fine-tuning in visual perception tasks.

Models	mAP	bird	feline-or-canid	serpent	slime	wyvern
Qwen2-VL-2B	20.6	12.9	19.8	25.5	18.4	26.4
4-shot						
+ SFT	26.8	19.5	22.4	26.8	33.5	31.8
+ Visual-RFT	61.8	63.9	53.2	70.2	64.5	57.5
Δ	+41.2	+51.0	+33.4	+44.7	+46.1	+31.1
16-shot						
+ SFT	51.3	42.7	44.4	56.4	61.1	52.0
+ Visual-RFT	63.4	59.9	50.8	76.3	71.7	58.1
Δ	+42.8	+47.0	+56.4	+50.8	+53.3	+31.7

derstanding of the task and deeper learning from reasoning.

4.3. Few-Shot Object Detection

Few-shot learning has always been one of the core challenges faced by traditional visual models and large-scale vision-language models (LVLMs). Reinforcement fine-tuning provides a new solution to this problem by enabling the model to quickly learn and understand with a small amount of data. We selected eight categories from the COCO dataset, with 1, 2, 4, 8, and 16 images per category, to construct training sets with limited data. For the LVIS dataset, we select 6 rare categories. Since the training images for these rare categories are very sparse, with each category having between 1 and 10 images, we approximated this as a 10-shot setting. We then train the Qwen2-VL-2/7B model for 200 steps using both reinforcement fine-tuning and SFT, to evaluate the model’s learning ability with lim-

Table 6. **Reasoning Grounding Results on LISA [11].** Visual-RFT surpasses SFT in reasoning grounding with 239 training images.

Model	mIoU _{test}	mIoU _{val}	gIoU _{test}
OV-Seg [14]	28.4	30.5	26.1
X-Decoder [54]	28.5	29.1	24.3
GroundedSAM [18]	26.2	28.6	21.3
Qwen2-VL-2B	26.9	30.1	25.3
+ SFT	28.3	29.7	25.3
+ Visual-RFT	37.6	34.4	34.4
Δ	+10.7	+4.3	+9.1
Qwen2-VL-7B	40.4	45.2	38.0
+ SFT	39.1	43.9	37.2
+ Visual-RFT	43.9	47.1	43.7
Δ	+3.5	+1.9	+5.6

Table 7. **Open Vocabulary Object Detection Results on COCO dataset.** We trained on 65 base categories and tested on 15 novel categories.

Models	mAP_n	mAP_b	mAP_{all}
Qwen2-VL-2B	9.8	6.0	6.7
+ SFT	13.6	7.8	8.9
+ Visual-RFT	31.3	20.6	22.6
Δ	+21.5	+14.6	+15.9
Qwen2-VL-7B	26.3	17.5	19.2
+ SFT	25.7	17.5	19.0
+ Visual-RFT	35.8	25.4	27.4
Δ	+9.5	+7.9	+8.2

ited data.

As shown in Tab. 3 and Tab. 4, although both SFT and reinforcement fine-tuning can improve the model’s recognition accuracy under the few-shot setting, the model after reinforcement fine-tuning consistently outperforms the SFT model by a large margin, maintaining a significant lead. On the COCO [15] categories, as the training data increases, the SFT model reaches an average mAP of approximately 31, while the reinforcement fine-tuned model approaches 47. In the LVIS [5] few-shot experimental results shown in Tab. 4, for the six more challenging rare categories in LVIS, reinforcement fine-tuning still outperforms SFT. The results in Tab. 3 and Tab. 4 clearly demonstrate the exceptional performance of reinforcement fine-tuning in the few-shot setting, where the model achieves a significant improvement in visual perception capabilities through reinforcement learning with only a small amount of data.

Table 8. **Open Vocabulary Object Detection Results on LVIS dataset.** We trained on the 65 base categories of the COCO dataset and tested on the 13 rare categories of the LVIS dataset.

Models	mAP	casserole	die	egg roll	futon	garbage	handsaw	hippopotamus	kitchen table	mallet	omelet	shot glass	stepladder	sugar bowl
GroudingDINO-B [18]	23.9	17.1	0.0	2.4	47.5	27.7	13.4	15.2	92.5	0.0	26.6	16.0	41.0	10.7
Qwen2-VL-2B	2.7	1.6	1.2	0.0	2.4	0.0	10.0	0.0	13.4	0.2	4.7	2.1	0.0	0.0
+ SFT	7.6	3.9	21.2	0.0	0.0	10.7	9.0	11.6	19.4	0.0	11.7	6.3	0.0	5.2
+ Visual-RFT	20.7	24.5	23.4	2.0	16.0	27.7	20.2	14.4	45.8	11.1	22.7	6.0	6.0	40.2
Δ	+18.0	+22.9	+22.2	+2.0	+13.6	+27.7	+10.2	+14.4	+32.4	+10.9	+18.0	+3.9	+6.0	+40.2
Qwen2-VL-7B	15.7	3.7	21.9	0.7	24.5	15.3	19.2	13.1	14.5	11.9	18.1	27.9	0.0	33.8
+ SFT	24.0	20.8	25.4	0.6	41.8	12.2	19.2	18.8	42.5	11.9	15.3	27.9	28.1	47.8
+ Visual-RFT	30.4	19.7	27.8	4.3	41.8	17.4	35.1	20.0	70.6	16.7	23.5	29.8	29.3	59.8
Δ	+14.7	+16.0	+5.9	+3.6	+17.3	+2.1	+15.9	+6.9	+56.1	+4.8	+5.4	+1.9	+29.3	+26.0

We further test on some abstract out-of-domain datasets. We selected the MG (Monster Girls) dataset, which contains different types of anime-style monster girls. By using out-of-domain data, we increased the difficulty of both model recognition and reasoning, and conducted experiments under 4-shot and 16-shot settings. The results, shown in Tab. 5, indicate that reinforcement fine-tuning achieved a significant performance improvement, surpassing supervised fine-tuning (SFT).

4.4. Reasoning Grounding

Another crucial aspect of vision-language intelligence is grounding the exact object according to user needs. Previous specialized detection systems lack reasoning abilities and fail to fully understand the user’s intentions. Pioneered by LISA [11], there have been works done to enable large language models (LLMs) to output control tokens for other models (such as SAM [9]) or directly predict bounding box coordinates [29, 38] through supervised fine-tuning. In our work, we explore the use of Visual-RFT in this task and find that reinforcement learning (RL) leads to significant improvements over supervised fine-tuning.

We finetune Qwen2-VL 2B/7B model [38] using Visual-RFT and supervised fine-tuning (SFT) on the LISA training set, which consists of 239 images with reasoning grounding objects. We follow the same test setting with LISA and compare the results of SFT and our method, both with 500 fine-tuning steps. As shown in Tab. 6, Visual-RFT significantly improves the final results in terms of bounding box IoU compared to SFT. Additionally, we prompt SAM [9] with the Qwen2-VL predicted bounding box to generate the segmentation mask (evaluated using gIoU). Visual-RFT significantly enhances grounding ability and outperforms previous specialized detection systems. Qualitative results

are visualized in Fig. 5, where the thinking process significantly improves the ability to reason and grounding accuracy. Through Visual-RFT, Qwen2-VL learns to think critically and carefully examine the image to produce accurate grounding results.

4.5. Open Vocabulary Object Detection

The advantage of Visual-RFT over SFT arises from the former’s true deep understanding of the task, rather than merely memorizing the data. To further demonstrate the powerful generalization ability of reinforcement fine-tuning, we conduct open vocabulary object detection experiments. We first randomly sampled 6K annotations from the COCO dataset, which included 65 base categories. We perform Visual-RFT and SFT on the Qwen2-VL-2/7B model [38] using this data, and test the model on 15 new categories it has never seen before. To increase the difficulty, we further test 13 rare categories from the LVIS [5] dataset.

As shown in Tab. 7 and Tab. 8, after reinforcement fine-tuning, the Qwen2-VL-2/7B model achieves an average mAP increase of 21.5 and 9.5 on 15 new categories from the COCO dataset. On the more challenging rare categories of the LVIS [5] dataset, mAP increased by 18.0 and 14.7. The Visual-RFT not only transfers its detection capabilities from the COCO base categories to the new COCO categories but also achieves significant improvements on the more challenging rare categories of LVIS. Notably, for some rare LVIS categories in Tab. 8, the original or SFT-trained models cannot recognize these categories, resulting in 0 AP. However, after reinforcement fine-tuning, the model shows a qualitative leap from 0 to 1 in recognizing these previously unidentifiable categories (such as egg roll and futon). This demonstrates that Visual-RFT has a significant impact

on improving the performance and generalization ability in visual recognition for LVLMS.

5. Conclusion

In this paper, we introduce Visual Reinforcement Fine-tuning (Visual-RFT), the first approach to adapt the GRPO-based reinforcement learning strategy for enhancing the visual perception and grounding ability of LVLMS. By using a rule-based verifiable reward system, Visual-RFT reduces the need for manual labeling and simplifies reward computation, achieving significant improvements across various visual perception tasks. Extensive experiments show that Visual-RFT excels in fine-grained classification, open vocabulary detection, reasoning grounding and few-shot learning tasks. It outperforms supervised fine-tuning (SFT) with minimal data and shows strong generalization. This work demonstrates the potential of reinforcement learning to enhance the capabilities of LVLMS, making them more efficient and effective in visual perception tasks.

References

- [1] Marwa Abdulhai, Isadora White, Charlie Snell, Charles Sun, Joey Hong, Yuexiang Zhai, Kelvin Xu, and Sergey Levine. Lmrl gym: Benchmarks for multi-turn reinforcement learning with language models. *arXiv preprint arXiv:2311.18232*, 2023. 3
- [2] Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024. 4
- [3] Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer. Grounding large language models in interactive environments with on-line reinforcement learning. In *ICLR*, 2023. 3
- [4] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 1, 4, 5
- [5] 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. 3, 7, 8, 9
- [6] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024. 4
- [7] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv:2412.16720*, 2024. 1, 2, 3
- [8] Fangkai Jiao, Geyang Guo, Xingxing Zhang, Nancy F Chen, Shafiq Joty, and Furu Wei. Preference optimization for reasoning with pseudo feedback. *arXiv preprint arXiv:2411.16345*, 2024. 4
- [9] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 9
- [10] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV workshops*, 2013. 7
- [11] 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. 3, 6, 8, 9
- [12] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024. 1, 4
- [13] 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

- [14] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7061–7070, 2023. 8
- [15] 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 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 3, 8
- [16] 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. 3
- [17] Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. Skywork-Reward: Bag of tricks for reward modeling in llms. *arXiv preprint arXiv:2410.18451*, 2024. 1, 4
- [18] 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. 3, 8, 9
- [19] Ziyu Liu, Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Haodong Duan, Conghui He, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. Mia-dpo: Multi-image augmented direct preference optimization for large vision-language models. *arXiv preprint arXiv:2410.17637*, 2024. 4
- [20] Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft: Reasoning with reinforced fine-tuning, 2024. 4
- [21] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 7
- [22] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *ICVGIP*, 2008. 7
- [23] OpenAI. Hello gpt-4o, 2024. 3
- [24] OpenAI. Openai o3-mini system card, 2025. 2
- [25] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 3
- [26] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 1, 4
- [27] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*, 2012. 7
- [28] Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hananeh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization. In *ICLR*, 2023. 3
- [29] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018, 2024. 9
- [30] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017. 4
- [31] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 2, 4
- [32] Charlie Victor Snell, Ilya Kostrikov, Yi Su, Sherry Yang, and Sergey Levine. Offline RL for natural language generation with implicit language q learning. In *ICLR*, 2023. 3
- [33] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In *NeurIPS*, 2022. 3
- [34] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023. 4
- [35] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023. 4
- [36] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. In *ACL*, 2024. 3
- [37] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1.5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025. 1, 4
- [38] 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. 3, 7, 9
- [39] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024. 4

- [40] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *ICLR*, 2023. 3
- [41] Huaiyuan Ying, Shuo Zhang, Linyang Li, Zhejian Zhou, Yunfan Shao, Zhaoye Fei, Yichuan Ma, Jiawei Hong, Kuikun Liu, Ziyi Wang, et al. Internlm-math: Open math large language models toward verifiable reasoning. *arXiv preprint arXiv:2402.06332*, 2024. 4
- [42] Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, et al. RIHF-V: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In *CVPR*, 2024. 4
- [43] Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He, Zhiyuan Liu, Tat-Seng Chua, et al. RLAIIF-V: Aligning mllms through open-source ai feedback for super gpt-4v trustworthiness. *arXiv preprint arXiv:2405.17220*, 2024. 4
- [44] Yuhang Zang, Wei Li, Jun Han, Kaiyang Zhou, and Chen Change Loy. Contextual object detection with multimodal large language models. *IJCV*, 2024. 3
- [45] Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Ziyu Liu, Shengyuan Ding, Shenxi Wu, Yubo Ma, Haodong Duan, Wenwei Zhang, et al. A simple yet effective multi-modal reward model. *arXiv preprint arXiv:2501.12368*, 2025. 1, 4
- [46] Kechi Zhang, Ge Li, Yihong Dong, Jingjing Xu, Jun Zhang, Jing Su, Yongfei Liu, and Zhi Jin. Codedpo: Aligning code models with self generated and verified source code. *arXiv preprint arXiv:2410.05605*, 2024. 4
- [47] Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, et al. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024. 3
- [48] Yuxiang Zhang, Shangxi Wu, Yuqi Yang, Jiangming Shu, Jinlin Xiao, Chao Kong, and Jitao Sang. o1-coder: an o1 replication for coding. *arXiv preprint arXiv:2412.00154*, 2024. 4
- [49] Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. Beyond hallucinations: Enhancing lvlms through hallucination-aware direct preference optimization. *arXiv preprint arXiv:2311.16839*, 2023. 4
- [50] Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. Aligning modalities in vision large language models via preference fine-tuning. *arXiv preprint arXiv:2402.11411*, 2024.
- [51] Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. Aligning modalities in vision large language models via preference fine-tuning. *arXiv preprint arXiv:2402.11411*, 2024. 4
- [52] Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. Archer: Training language model agents via hierarchical multi-turn rl. In *ICML*, 2024. 3
- [53] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv:1909.08593*, 2019. 3
- [54] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu Yuan, et al. Generalized decoding for pixel, image, and language. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15116–15127, 2023. 8