

Reinforced Visual Perception with Tools

Anonymous EMNLP submission

Abstract

Visual reasoning, a cornerstone of human intelligence, encompasses complex perceptual and logical processes essential for solving diverse visual problems. While advances in computer vision have produced powerful models for various perceptual tasks, leveraging these for general visual reasoning remains challenging. Prior work demonstrates that augmenting LLMs with vision models via supervised finetuning improves performance, but faces key limitations such as expensive data generation, reliance on careful data filtering, and poor generalization. To address these issues, we propose REVPT to enhance multi-modal LLMs’ abilities to reason about and use visual tools through reinforcement learning. We introduce a novel RL algorithm based on GRPO, designed to train models to reason with a suite of seven visual tools. Our explorative results across models ranging from 3B to 7B parameters show that our method achieves state-of-the-art performance on several perception-heavy benchmarks, including SAT, CV-Bench, BLINK, and BLINK-Hard, significantly outperforming supervised and text-based RL finetuning baselines. We hope our explorative on RL-based visual tool-usage can bring insights to the community.

1 Introduction

Visual reasoning is a core component of human intelligence. It enables us to solve a wide range of problems, from daily activities such as navigation to challenging geometry problems. Unlike verbal reasoning, visual reasoning is a more complex multifaceted process that requires not just straightforward logical reasoning but also sound visual perception, which further relies on atomic abilities such as object recognition, edge detection, depth estimation, *etc.* Due to the complexity of visual perception, the computer vision community has developed specialized models for different perceptual

tasks, such as RecognizeAnything for object recognition, DepthAnything for depth estimation, and SegmentAnything for segmentation (Yang et al., 2024; Liu et al., 2023d; Zhang et al., 2023).

To leverage these advances in perception for enhanced visual reasoning and problem-solving, prior work attempts to augment (multimodal) language models with vision specialists. For example, VisProg first showcases that augmenting large language models (LLMs) with various vision models boosts models’ performance on diverse visual tasks (Gupta and Kembhavi, 2022). Similarly, VisualSketchPad finds that prompting GPT-4o to use sketching tools and depth models improves its performance on challenging perception and visual math benchmarks by large margins (Hu et al., 2024a). More recently, a few works demonstrate that open-source multimodal language models also benefit from using visual tools and reasoning about their outputs after supervised finetuning on tool-use data (Qi et al., 2024; Liu et al., 2023c).

Nonetheless, there are a few limitations to learning to reason with visual tools with supervised finetuning (SFT): first, it relies on expensive data curation. Prior work heavily relies on expensive commercial models like GPT-4 to generate high-quality tool-use and reasoning traces (Qi et al., 2024; Ma et al., 2024; Shao et al., 2024a). Second, it often requires careful data filtering. Previous efforts reveal that aggressive filtering based on data format, answer correctness and other heuristics is crucial to performance gains (Hu et al., 2024b; Ma et al., 2024). Most importantly, supervised finetuning results in limited generalization: it relies on offline trajectories that pre-define one particular set of tools to use, lacks incentives for the model to explore alternative tools or adapt to new environment outputs, and thus limits its generalization to unseen tools or problems (Guo et al., 2025; Chen et al., 2025; Chu et al., 2025).

To mitigate these limitations, we borrow inspi-

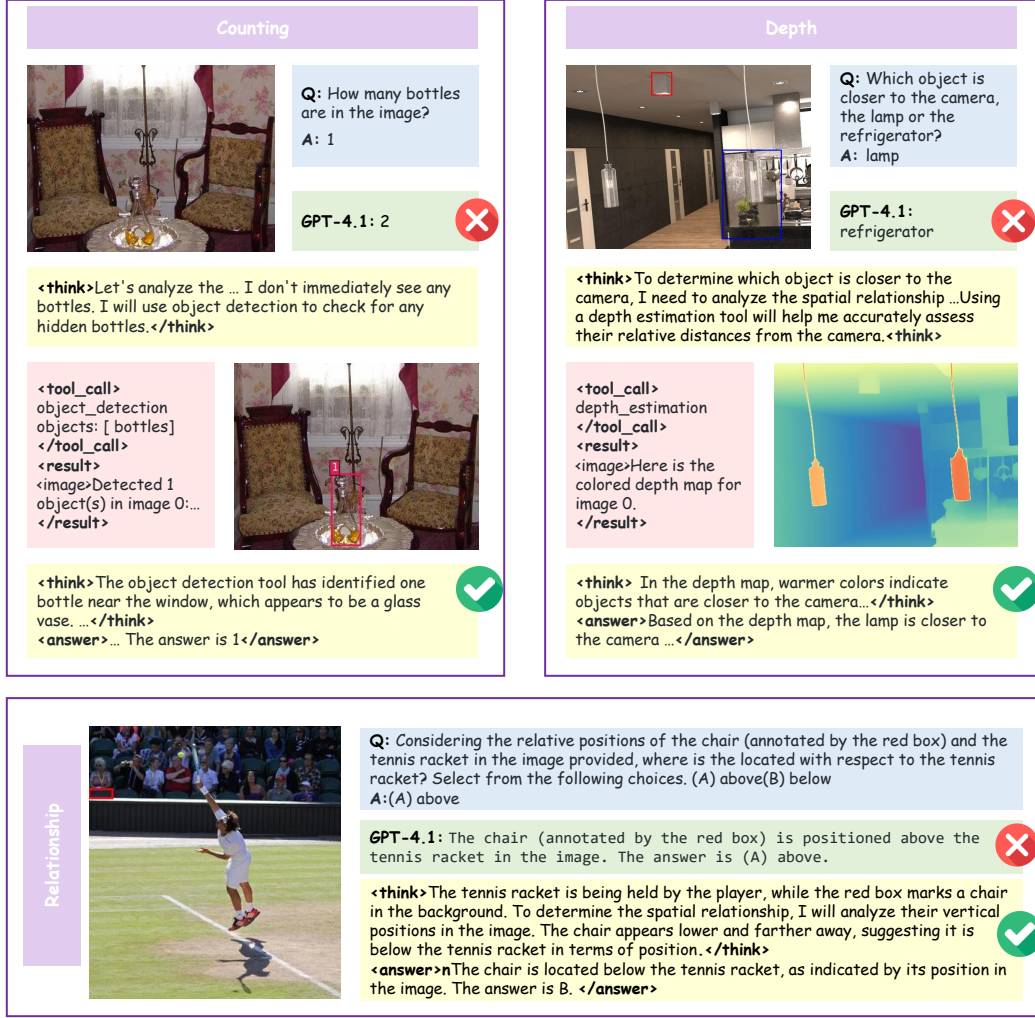


Figure 1: Example outputs of REVPT-3B. REVPT-3B is able to solve problems through thinking and selectively employing visual tools via GRPO training.

ration from recent work on enhancing LLM reasoning with reinforcement learning (RL) (Zhang et al., 2025; Guo et al., 2025; Huang et al., 2025; Zhan et al., 2025; Shen et al., 2025) and propose **Reinforced Visual Perception with Tools** (REVPT), to enhance multimodal language model’s visual problem-solving skills by training them to reason with visual tools via RL. REVPT consists of cold-start and GRPO process to enable efficient reinforcement learning on multimodal language models with 6 visual tools. Unlike SFT, where all the tool-use reasoning traces need to be generated and filtered in advance, we only need to select the appropriate visual questions to perform RL on. Moreover, instead of supervising the model with one correct tool-use trajectory for each question, RL allows the model to explore and learn from different possible solutions and incentivizes

the model to output the best one.

Our explorative experiments with models of different parameter scales—3B and 7B—demonstrate that REVPT-3B enables state-of-the-art performance, consistently outperforming SFT baselines across five perception-intensive benchmarks, including MMVP, CV-Bench, BLINK, and BLINK-Hard. Notably, our approach surpasses the original Instruct version models by significant margins, achieving 6.88% and 4.33% improvements on CV-Bench and MMVP respectively. Furthermore, our models outperform commercial models GPT-4.1 and Gemini-2.0-Flash on the challenging BLINK-Hard benchmark by 3-4 percentage points. We believe that REVPT, alongside our fully open-sourced code, datasets, and evaluation platform, will serve as a valuable resource for the broader research community.

2 Thinking with Images and Tools

In this section, we introduce REVPT, a reinforcement learning framework designed to train multi-modal language models for visual perception tasks. It is built upon the GRPO algorithm (Shao et al., 2024b), aiming to enhance the model’s ability to tackle complex visual problems by integrating visual processing tools as reasoning steps.

Given a multi-modal query, REVPT agent first generates a reasoning process about analyzing the query and deciding whether to call tools or answer directly. Then, by iteratively analyzing the results from the tools, the model generates a final response to the query. The overall architecture and process flow are illustrated in Figure 2.

Through GRPO training, the model has demonstrated the capacity to leverage appropriate tools and accurately interpret processed images, thereby overcoming limitations previously encountered. Our implementation is based on the veRL (Sheng et al., 2024) framework.

2.1 Preliminary: GRPO

Group Relative Policy Optimization (GRPO) presents an alternative approach in the landscape of reinforcement learning algorithms. A key distinction of GRPO is its departure from the actor-critic paradigm, exemplified by algorithms like Proximal Policy Optimization (PPO). Instead, GRPO evaluates policy performance by directly comparing a collection of candidate responses generated by the current policy. The core mechanism of GRPO begins with a given prompt or query, q . For this query, the policy, denoted as π_θ , is used to sample a set of N distinct candidate responses, represented as $\{o_1, o_2, \dots, o_N\}$. Each of these sampled responses o_i is then assessed using a reward function, $R(q, o_i)$, which quantifies the quality or appropriateness of the response o_i in relation to the initial query q . To assess relative quality within the sampled group, GRPO calculates an advantage A_i for each response by normalizing its reward:

$$A_i = \frac{r_i - \text{mean}\{r_1, r_2, \dots, r_N\}}{\text{std}\{r_1, r_2, \dots, r_N\}} \quad (1)$$

where A_i represents the advantage of the candidate response o_i relative to other sampled responses. GRPO encourages the model to generate responses with higher advantages within the group by updating

the policy π_θ using the following objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E} [\{o_i\}_{i=1}^N \sim \pi_{\theta_{old}}(q)] \quad (2)$$

$$= \frac{1}{N} \sum_{i=1}^N \{\min[s_1 \cdot A_i, s_2 \cdot A_i] - \beta \mathbb{D}_{KL}[\pi_\theta || \pi_{ref}]\} \quad (3)$$

$$s_1 = \frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)} \quad (4)$$

$$s_2 = \text{clip}\left(\frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1 + \epsilon, 1 - \epsilon\right) \quad (5)$$

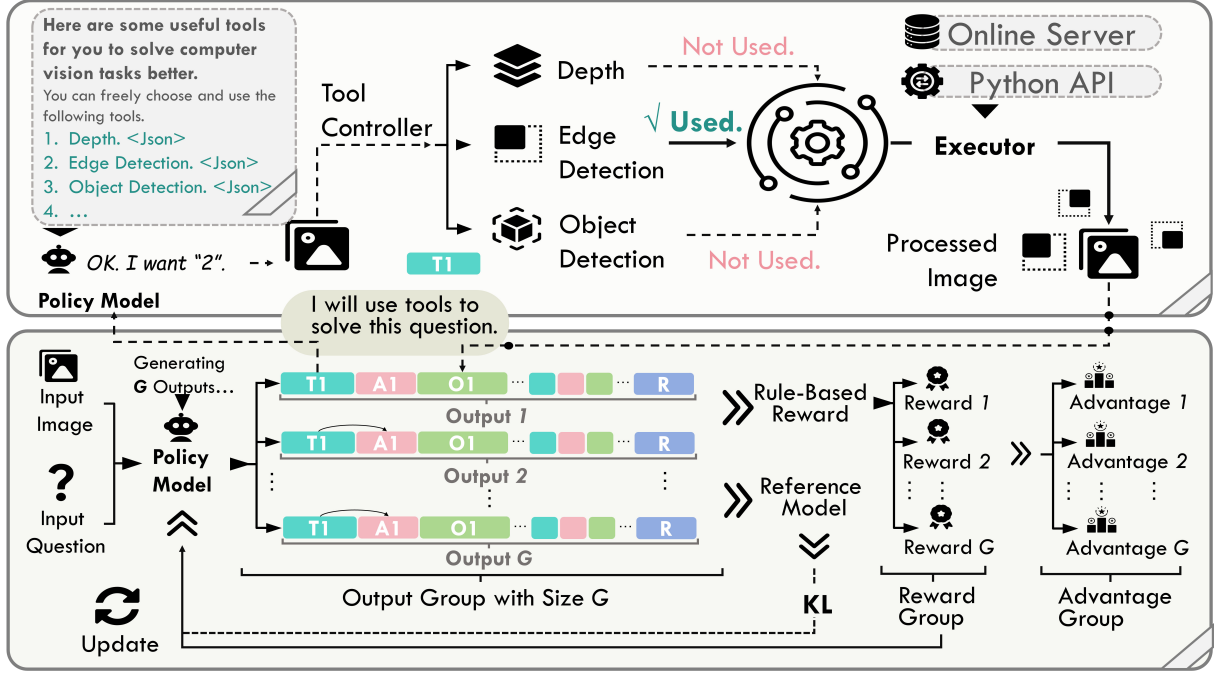
Here, the clip function limits the probability ratio of response o_i under the current and old policies, and \mathbb{D}_{KL} is a Kullback-Leibler divergence term penalizing large deviations from a reference policy π_{ref} , both contributing to stable training.

2.2 Vision Tools

The efficacy of leveraging visual tools for visual reasoning has been well-established. Our framework integrates several such high-performing visual tools, enabling their execution and subsequent result incorporation during the MLLM inference process to construct a comprehensive rollout. While Table 1 offers a detailed summary of each tool’s parameters and specifications, their core functionalities and representative use cases are elaborated upon below.

- **LLMDET** (Fu et al., 2025): This tool is capable of open-vocabulary detection. It takes an input image I_{in} and a textual query q_{text} to locate instances of described objects, submitting their boundaries \mathcal{B}_{out} . It helps model localize objects in the image.
- **SAM** (Kirillov et al., 2023): This tool can generate fine-grained segmentation masks. SAM can be invoked for automatic segmentation, where it takes an input image I_{in} and outputs a set of segmentation masks \mathcal{M}_{out} for all detected objects within the image. Alternatively, it can perform point-prompted segmentation, taking an input image I_{in} and one or more user-defined points \mathcal{P} on the image to output a segmentation mask \mathcal{M}_{out} for the object indicated by the provided point. This tool assists in precisely delineating object boundaries.
- **ZOOMIN**: This tool facilitates focused analysis by taking an input image I_{in} and a specified region of interest \mathcal{R} to output a magnified view

➤ Visual Tools Pool



➤ Visual Tools Training Process Based on GRPO Reinforce Learning.

Figure 2: An overall pipeline of our REVPT. **(Top)**: Model-generated tool requests are managed by a local environment-based Tool Controller, which independently deploys vision tool services (e.g., Depth, Object Detection). These tools’ outputs are then fed back to the LVLm for iterative reasoning. **(Bottom)**: When processing a visual reasoning problem, REVPT employs K-turn rollouts where the model interacts with the tool environment to learn an adaptive policy, culminating in the final model.

\mathcal{I}_{focus} of that region. It helps the model to concentrate on pertinent areas, thereby mitigating the influence of irrelevant information and amplifying salient features.

- **EDGEDETECTION**: This tool is designed to identify and delineate the perimeters of objects and significant textural variations, this tool transforms an input image \mathcal{I}_{in} into a feature map \mathcal{E}_{out} that emphasizes these structural edges. This grants models a sharper perception of object silhouettes, surface patterns, and other intrinsic structural data.
- **DEPTHANYTHING**(Yang et al., 2024): This tool computes spatial depth from a monocular visual input. Given an image \mathcal{I}_{in} , it generates a corresponding depth map \mathcal{D}_{out} which encodes the estimated distance of scene elements from the observer. We use DEPTH_ANYTHING_V2 to serve as the tool. This capability equips models with a more nuanced understanding of three-dimensional arrangements and the relative positioning of objects.
- **DRAWLINE**: This function serves to superim-

pose linear annotations onto an image. Utilizing an input image \mathcal{I}_{in} and precise line parameters \mathcal{L}_{spec} (such as origin and termination points), it produces an augmented image \mathcal{I}_{drawn} featuring these graphical overlays. The lies in visually accentuating particular linear connections, trajectories, or critical measurements for the model’s analytical consideration.

2.3 Cold Start

Initially, our objective is to train a multimodal agent using the R1-Zero method. However, during the training process, we observe a progressive decline in the agent’s propensity to utilize tools. This phenomenon likely stemmed from the fact that solving visual tasks did not inherently require tool usage, and reasoning based on processed images represented a distribution shift from the model’s initial training data. Consequently, we opt to incorporate cold-start data during the initial training phases. This strategy aim to facilitate the model’s initial acquisition of tool utilization skills for solving visual problems.

Table 1: Vision Tools Overview. This table lists vision tools integrated within REVPT, detailing their arguments, result formats, and description for their core function.

Tool	Arguments	Result	Description
LLMDet	Image + Text Prompt	Annotated Image + Boxes	Zero-shot Object Detection
SAM	Image + Point(Optional)	Annotated Image	Auto Segmentation
ZOOMIN	Image + Coordinates + Factor	Cropped Image	Region Cropping
EDGEDETECTION	Image	Edge Map	Edge Detection Using Scharr Algorithm
DEPTHANYTHING	Image	Depth Map	Depth Estimation
DRAWLINES	Image + Point Coordinate + Mode	Annotated Image	Draw Lines

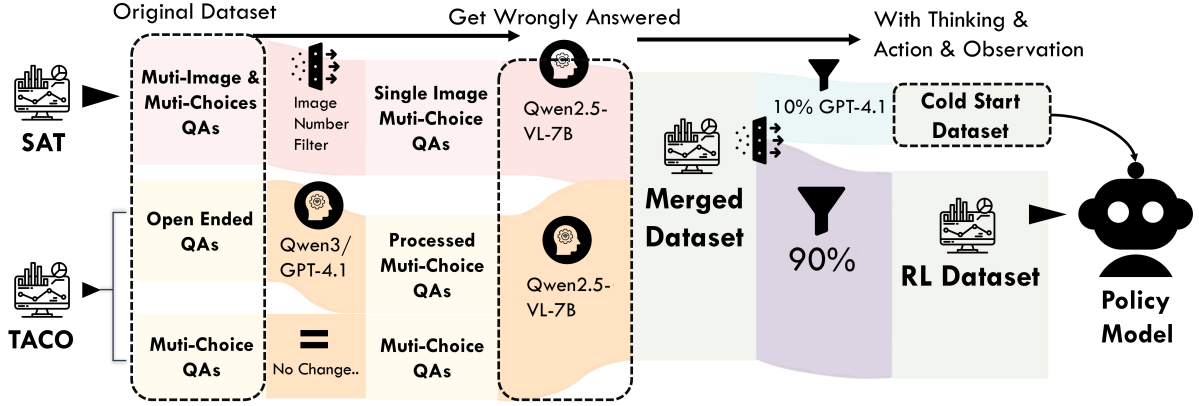


Figure 3: Reinforced visual tool-usage training requires high-quality and verified data. We transform TACO and SAT dataset into multiple-choice quesiton, then filter out easy questions with Qwen2.5-VL-7B.

Currently, we employ GPT-4.1 as the agent to synthesize high-quality data for tool-augmented reasoning in visual problem solving. We require the GPT model to follow a predefined format when answering questions. Specifically, it should first engage in a reasoning process, considering the potential utility of external tools, and then employ a multi-turn generation strategy to address the problem. We explicitly encourage the model in our prompt to generate more tool-assisted rollouts. Subsequently, we filter out rollout trajectories generated by the GPT model that contain errors.

After synthesizing tool-integrated reasoning data, we perform supervised fine-tuning on it. Denote the query as Q^i , the rollout trajectory τ^i as a sequence of actions a_t^i and observations o_t^i . We optimize the cross-entropy loss:

$$\begin{aligned} \mathcal{L}_{\text{SFT}}(\theta) &= -\frac{1}{N} \sum_{i=1}^N \log P_{\theta}(\tau^i | Q^i) \\ &= -\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{n^i} \log P_{\theta}(a_t^i | Q^i, a_{<t}^i, o_{<t}^i) \end{aligned} \quad (6)$$

$$(7)$$

By minimizing \mathcal{L}_{SFT} , the model acquires a robust

Cold-Start policy for sequential vision-tool invocation, providing a solid foundation for the subsequent reinforcement learning phase.

2.4 Reward Design

Vision tasks often possess ground truth data, facilitating the use of rule-based rewards to evaluate a model’s responses. This approach circumvents the introduction of neural network-based reward models, thereby preventing reward hacking.

- **Correctness Checking.** We restructure the dataset into a multiple-choice format. The model is required to put the final answer’s option in the box, enabling reliable rule-based verification of correctness. If the answer match the ground truth, it gets through checking.
- **Format Checking.** In each turn, the response should enclose its thinking process in `<think></think>` and enclose its tool call in `<tool_call></tool_call>` or answer in `<answer></answer>`. If the response matches the pattern, it gets through format checking. Given the query q and the rollout o , the reward is defined as:

$$Reward(q, o) = \begin{cases} 1 & , If FormatCorrect(o) \\ & \wedge AnswerCorrect(o) \\ -1 & , otherwise \end{cases}$$

3 Experiment

3.1 Experiment Setups

Models. We conduct most our experiments on Qwen2.5-VL-3B-Instruct (Wang et al., 2024a) for their strong visual perception and tool-calling capability. We also select an early checkpoint based on 7B model and include some results in Appendix.

Dataset Construction. A high quality cold-start dataset and difficult dataset for RL training is very important for model to learn test-time scaling itself (Yu et al., 2025). Therefore, we filter the SAT dataset (Ray et al., 2024) and Taco (Ma et al., 2024) training set with Qwen2.5-7B-Instruct and retain the samples it answers incorrectly. Then, we split it into 1:9 for code-start data synthetic and training data. To construct high-quality code-start dataset, we leverage GPT-4.1 to synthetic tool trajectory and retain the subset that it answer correctly. Finally, we get 7k cold-start dataset with well-curated reasoning chain and tool trajectory.

Baseline. We compare REVPT against the following models and approaches: **(1) Commercial Models:** We select GPT-4.1 (OpenAI, 2025) and Gemini-2.0-Flash (Google, 2024) as strong multi-modal baselines. Both are evaluated in a zero-shot setting without external tools as strong generalist benchmarks. **(2) TACO:** Learns to invoke 15 external tools (e.g., OCR, depth, etc.) by generating Chain-of-Thought-and-Action (CoTA) sequences via supervised learning on synthetic data. TACO typically executes tools within a single process, contrasting with our RL-based approach and distributed architecture (Ma et al., 2024). **(3) Qwen-Base:** We implement two base models without any tool usage, employing different prompt formats for a robust evaluation. **(4) Qwen-SAT-SFT:** Two models after supervised finetuning with the SAT training set (13k samples) as a strong baseline. We also include an enhanced SFT baseline with rephrased answers generated by the model itself as reported by previous research (Wang et al., 2024b). **(5) Qwen-SFT (cold start):** The model state after cold-start training. **(6) Text-based RL:** An RL agent trained similarly to REVPT but without any visual feedback from tools. This follows the native

GRPO training for MLLM in VisualThinker-R1-Zero (Zhou et al., 2025).

Evaluation. We select 4 multi-modal benchmarks covering diverse capabilities with a focus on visual perception and reasoning tasks. This includes CV-Bench (Tong et al., 2024a), Blink (Fu et al., 2024), MMVP (Tong et al., 2024b), Blink-Hard (Bigverdi et al., 2024). We choose subset with only one image as input.

Experiment Details. We conduct model training using configurations with $8 \times$ NVIDIA A800 GPUs. We leverage LLaMA-Factory (Zheng et al., 2024) as Supervised Fine-tuning platform and Verl as visual tool-base RL platform (Sheng et al., 2024). The training process involved two phases: **(1) Cold-start Period:** Models are trained for 2 epochs with a learning rate of $1e-5$ and a global batch size of 64. **(2) RL Period:** For this phase, models are trained for 100 steps (Loss Curve in Figure 4). We clarify that reinforced learning with visual tool do not reach converge and reward is very unstable.

3.2 Experiment Results

REVPT outperform supervised finetuning and text-based RL. As shown in Table 2, our method significantly outperforms the original model on both CV-Bench and BlinkHard benchmarks, while also demonstrating improvements over SFT cold start models. The performance gains are particularly pronounced in depth and distance tasks, with improvements of 9.5% and 13.23% respectively compared to the base model, indicating that reinforcement learning effectively teaches the model to utilize tools and interpret tool-generated information. Notably, our vision tool-based RL approach substantially outperforms the converged text-based RL by 12.7% on CV-Bench, demonstrating how the model leverages smaller model tools as guidance to acquire more fundamental perception information, thereby enhancing its perceptual capabilities.

Performance Tradeoffs in Perception-specialized Training. Our training on the curated perception subset from SAT and TACO significantly reduced the model’s general performance capabilities. Direct SFT results reveal substantial degradation across several Blink capabilities. More concerning, in our 7B model experiments (detailed in the Table ??), we detect even larger performance declines on the more general MMstar benchmark, with scores approximately 10% below the baseline instruct model. These findings highlight the

Table 2: Comparison performance between different models on vision-centric benchmarks

Model	CV-Bench					BLINK								MMVP	BLINK-HARD		
	Count	Relation	Depth	Distance	Avg.	Counting	IQ Test	Localization	Depth	Reflectance	Relation	Avg.	3		4	5	
Gemini-2.0-Flash	71.95	86.92	87.50	82.17	81.50	75.83	32.67	62.30	79.84	38.81	73.43	59.52	79.34	72.58	66.13	66.13	
GPT-4.1	67.77	92.00	94.50	89.50	84.76	75.20	54.00	73.20	76.80	73.20	54.40	67.80	88.00	71.77	62.90	63.71	
Qwen2.5-VL-3B-Instruct	68.65	74.92	76.00	71.67	72.55	68.33	6.67	49.18	64.52	38.81	83.92	50.95	64.67	61.29	54.03	46.77	
Qwen2.5-VL-3B-SFT	60.53	61.23	79.50	78.33	69.07	40.00	24.67	56.56	50.00	31.34	68.53	44.89	65.42	52.42	46.77	42.74	
Qwen2.5-VL-3B-GRPO	64.85	69.85	72.33	59.33	66.53	60.83	24.00	43.44	49.19	29.85	79.72	47.54	60.33	50.00	48.39	45.97	
R1-3B-cold-start	69.29	86.31	84.17	72.50	77.60	57.50	27.52	51.64	66.94	26.87	73.43	50.06	62.33	63.71	50.81	46.77	
R1-3B	69.92	88.15	85.50	75.50	79.23	54.17	25.33	50.82	66.13	32.09	75.52	50.19	69.00	73.27	66.67	61.22	

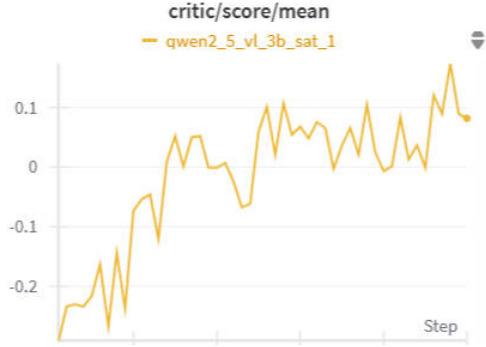


Figure 4: Our reward rapidly upgrade and reach converge.

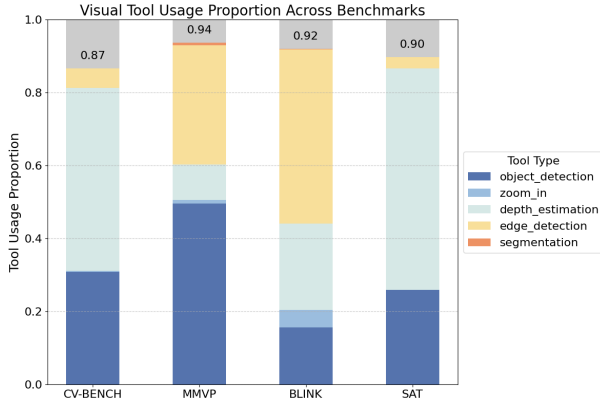


Figure 5: Our model effectively learns tool utilization post-cold-start, enhancing visual perception capabilities. Current RL data shows bias toward object detection and depth estimation over zoom and segmentation tools. Future work will address data balance and generalized perception objectives.

critical relationship between reinforcement learning sample distribution and resulting model capabilities. Our key insight is that developing effective visual tool-using agents requires not only sophisticated interaction environments but also carefully balanced, comprehensive, and high-quality datasets that preserve general capabilities while enhancing targeted skills.

Limitations of Visual Tool Integration. During our experimental process, we observe that these

visual tools often hindered rather than enhanced model performance, sometimes degrading or interfering with the model’s reasoning capabilities. In some cases, the model needed to correct erroneous tool outputs to arrive at accurate answers. For instance, in Figure 6, the model misclassify a mattress as a pillow, resulting in an incorrect response despite the model’s inherent capability to solve the problem correctly. Furthermore, we discover that enhancement of perception capabilities through tool use is significantly influenced by the model’s inherent abilities. In Figure 8, though the depth estimation tool returned a color-coded depth map (where warm colors indicate proximity and cool colors indicate distance), the model—despite acknowledging this information and analyzing it during reasoning—ultimately produced an incorrect answer. We attribute this phenomenon to the model’s original perception capabilities constraining its ability to effectively and accurately utilize external tools, limiting tool-use efficacy and correctness.

4 Related Works

Multimodal language models. Recently, there have been many advances on open-source multimodal models (Awadalla et al., 2023; Chen et al., 2023; Liu et al., 2023b,a, 2024; Dai et al., 2024; Li et al., 2022, 2023; Deitke et al., 2024). These efforts include training multimodal models to take in multiple images, engage in multi-turn conversations, and even understand videos (Liu et al., 2024; Jiang et al., 2024; Li et al., 2024). For example, LLaVA-Next achieves strong multi-image understanding through large-scale interleaved visual instruction tuning with M4-Instruct (Liu et al., 2024). Similarly, Mantis introduces a new large-scale multi-image instruction tuning dataset Mantis-Instruct for multi-image training (Jiang et al., 2024). These efforts pave the foundation for our work on learning vision-language models with image-text interleaved reasoning traces. over perceptual out-

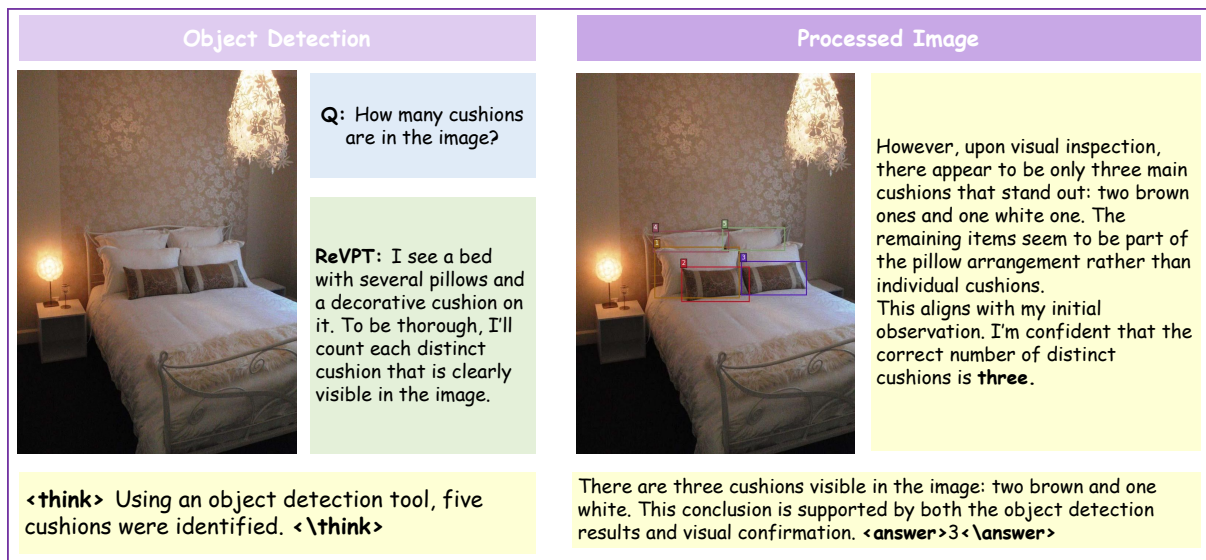


Figure 6: Erroneous outputs from object detection tools disrupt the model’s reasoning chain, ultimately lead to incorrect final predictions.

puts from vision specialists.

Multimodal tool-use. Recently, there has been increasing interest in enhancing multimodal language models with improved tool-use capabilities (Liu et al., 2023c; Qi et al., 2024; Shao et al., 2024a; Ma et al., 2024). LLaVa-Plus demonstrated the feasibility of training such models to utilize vision tools (Liu et al., 2023c). Visual Program Distillation transfers tool-use and reasoning skills into a multimodal model using chain-of-thought (CoT) data derived from programs (Hu et al., 2024b). Similarly, Visual CoT introduces a synthetic CoT dataset aimed at boosting the reasoning abilities of multimodal models (Shao et al., 2024a). More recently, LLaVa-CoT incorporates both perception and reasoning capabilities inspired by GPT-4o (Xu et al., 2025). CogCoM identifies six key manipulation strategies and trains multimodal models with synthetic chain-of-manipulation (CoM) data (Qi et al., 2024), while TACO contributes 273K multimodal reasoning traces derived from the outputs of 15 visual tools (Ma et al., 2024).

Multimodal Language Model Reasoning with RL. DeepSeek-R1 (Guo et al., 2025) has demonstrated that simple rule-based reinforcement learning can effectively induce strong reasoning behaviors. This R1-style reasoning paradigm has shown success in multimodal language models (Shen et al., 2025; Lu et al., 2025; Zhan et al., 2025; Huang et al., 2025; Feng et al., 2025; Li et al., 2025). VLM-R1 (Shen et al., 2025) applies reward-driven

fine-tuning to improve visual reasoning and generalization beyond supervised baselines. Vision-R1 (Zhan et al., 2025) introduces human-free alignment through vision-guided rewards, while another Vision-R1 (Huang et al., 2025) leverages CoT-style tasks and reward shaping to boost multi-step reasoning. UI-R1 (Lu et al., 2025) enhances action prediction in GUI agents via GRPO-based RL, achieving strong performance with compact models. In the video domain, Video-R1 (Feng et al., 2025) improves temporal reasoning through T-GRPO and mixed-modality rewards, and VideoChat-R1 (Li et al., 2025) reinforces spatio-temporal understanding across diverse video question types.

5 Conclusion

In this work, we explore tool-usage reinforcement learning to enable models to utilize external vision tools for test-time scaling. We propose REVPT, a training suite designed to integrate tool-usage with reinforcement learning, enabling models to optimize tool selection and interleaved text-image-tool reasoning through direct interaction and reward feedback. Our experiments across multiple multimodal benchmarks demonstrate that REVPT successfully raises performance beyond the base model’s capabilities, achieving results significantly higher than those obtained through supervised fine-tuning and text-only RL. We hope that REVPT, along with our fully open-sourced code, dataset, and platform, will serve as a valuable resource for the research community.

Limitations and Future Work

Although our method demonstrates exceptional post-training performance on Qwen2.5-VL-3B, numerous opportunities remain for exploring tool design optimization, data distribution ratios, and reward configuration refinements. Furthermore, our training methodology presents additional avenues for investigation—we currently train on the curated dataset for only a single epoch and plan to extend the training duration in future iterations. In our exploratory experiments with the 7B model variant, we observed performance improvements on CV-bench; however, we detected significant performance degradation on the more general MM-star benchmark. We attribute this decline to our dataset composition, which primarily emphasizes perception-oriented samples from SAT and TACO datasets rather than more diverse general-purpose examples, consequently diminishing performance across broader capabilities. Future work will address these data diversity considerations to maintain comprehensive multimodal reasoning abilities.

References

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. *Openflamingo: An open-source framework for training large autoregressive vision-language models*. Preprint, arXiv:2308.01390.

Mahtab Bigverdi, Zelun Luo, Cheng-Yu Hsieh, Ethan Shen, Dongping Chen, Linda G Shapiro, and Ranjay Krishna. 2024. Perception tokens enhance visual reasoning in multimodal language models. *arXiv preprint arXiv:2412.03548*.

Garrick Brazil, Abhinav Kumar, Julian Straub, Nikhila Ravi, Justin Johnson, and Georgia Gkioxari. 2023. Omni3D: A large benchmark and model for 3D object detection in the wild. In *CVPR*, Vancouver, Canada. IEEE.

Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. 2025. Sft or rl? an early investigation into training rl-like reasoning large vision-language models. *arXiv preprint arXiv:2504.11468*.

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning

for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*.

Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. 2025. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2024. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36.

Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, and 1 others. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*.

Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Benyou Wang, and Xiangyu Yue. 2025. Video-rl: Reinforcing video reasoning in mllms. *arXiv preprint arXiv:2503.21776*.

Shenghao Fu, Qize Yang, Qijie Mo, Junkai Yan, Xihan Wei, Jingke Meng, Xiaohua Xie, and Wei-Shi Zheng. 2025. Llm-det: Learning strong open-vocabulary object detectors under the supervision of large language models. *arXiv preprint arXiv:2501.18954*.

Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. 2024. Blink: Multimodal large language models can see but not perceive. *arXiv preprint arXiv:2404.12390*.

Google. 2024. Gemini 2.5 flash. <https://deepmind.google/technologies/gemini/flash/>.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Tanmay Gupta and Aniruddha Kembhavi. 2022. Visual programming: Compositional visual reasoning without training. *ArXiv*, abs/2211.11559.

Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Ranjay Krishna. 2024a. *Visual sketchpad: Sketching as a visual chain of thought for multimodal language models*. Preprint, arXiv:2406.09403.

Yushi Hu, Otilia Stretcu, Chun-Ta Lu, Krishnamurthy Viswanathan, Kenji Hata, Enming Luo, Ranjay Krishna, and Ariel Fuxman. 2024b. Visual program

distillation: Distilling tools and programmatic reasoning into vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9590–9601.

Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. 2025. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint arXiv:2503.06749*.

Dongfu Jiang, Xuan He, Huaye Zeng, Con Wei, Max Ku, Qian Liu, and Wenhui Chen. 2024. Mantis: Interleaved multi-image instruction tuning. *arXiv preprint arXiv:2405.01483*.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment anything. *arXiv:2304.02643*.

Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*.

Xinhao Li, Ziang Yan, Desen Meng, Lu Dong, Xiangyu Zeng, Yinan He, Yali Wang, Yu Qiao, Yi Wang, and Limin Wang. 2025. Videochat-r1: Enhancing spatio-temporal perception via reinforcement fine-tuning. *arXiv preprint arXiv:2504.06958*.

Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. *Microsoft COCO: common objects in context*. *CoRR*, abs/1405.0312.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.

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

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning.

Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. 2023c. *Llava-plus: Learning to use tools for creating multimodal agents*. *Preprint*, arXiv:2311.05437.

Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and 1 others. 2023d. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.

Zhengxi Lu, Yuxiang Chai, Yaxuan Guo, Xi Yin, Liang Liu, Hao Wang, Guanqing Xiong, and Hongsheng Li. 2025. Ui-r1: Enhancing action prediction of gui agents by reinforcement learning. *arXiv preprint arXiv:2503.21620*.

Zixian Ma, Jianguo Zhang, Zhiwei Liu, Jieyu Zhang, Juntao Tan, Manli Shu, Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Caiming Xiong, and 1 others. 2024. Taco: Learning multi-modal action models with synthetic chains-of-thought-and-action. *arXiv preprint arXiv:2412.05479*.

OpenAI. 2025. Gpt-4.1. <https://openai.com/index/gpt-4-1/>. Large language model. Accessed: 2025-05-11.

Ji Qi, Ming Ding, Wei Han Wang, Yushi Bai, Qingsong Lv, Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao Dong, and Jie Tang. 2024. *Cogcom: Train large vision-language models diving into details through chain of manipulations*. *Preprint*, arXiv:2402.04236.

Arijit Ray, Jiafei Duan, Reuben Tan, Dina Bashkirova, Rose Hendrix, Kiana Ehsani, Aniruddha Kembhavi, Bryan A Plummer, Ranjay Krishna, Kuo-Hao Zeng, and 1 others. 2024. Sat: Spatial aptitude training for multimodal language models. *arXiv preprint arXiv:2412.07755*, 3.

Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. 2024a. *Visual cot: Unleashing chain-of-thought reasoning in multi-modal language models*. *Preprint*, arXiv:2403.16999.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and 1 others. 2024b. Deepseek-math: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.

Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, and 1 others. 2025. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*.

Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:2409.19256*.

Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan,

722	Austin Wang, Rob Fergus, Yann LeCun, and Saining Xie. 2024a. Cambrian-1: A fully open, vision-centric exploration of multimodal llms .	776
723		777
724		778
725	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. 2024b. Eyes wide shut? exploring the visual shortcomings of multi-modal llms . <i>Preprint</i> , arXiv:2401.06209.	779
726		780
727		781
728		782
729	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	783
730		784
731		
732		
733		
734		
735		
736		
737	Yifei Wang, Yuyang Wu, Zeming Wei, Stefanie Jegelka, and Yisen Wang. 2024b. A theoretical understanding of self-correction through in-context alignment. <i>arXiv preprint arXiv:2405.18634</i> .	
738		
739		
740		
741	Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. 2025. Llava-cot: Let vision language models reason step-by-step . <i>Preprint</i> , arXiv:2411.10440.	
742		
743		
744		
745	Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. 2024. Depth anything v2. <i>arXiv:2406.09414</i> .	
746		
747		
748	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. <i>arXiv preprint arXiv:2503.14476</i> .	
749		
750		
751		
752		
753	Yufei Zhan, Yousong Zhu, Shurong Zheng, Hongyin Zhao, Fan Yang, Ming Tang, and Jinqiao Wang. 2025. Vision-r1: Evolving human-free alignment in large vision-language models via vision-guided reinforcement learning. <i>arXiv preprint arXiv:2503.18013</i> .	
754		
755		
756		
757		
758	Shaokun Zhang, Yi Dong, Jieyu Zhang, Jan Kautz, Bryan Catanzaro, Andrew Tao, Qingyun Wu, Zhiding Yu, and Guilin Liu. 2025. Nemotron-research-tool-n1: Tool-using language models with reinforced reasoning. <i>arXiv preprint arXiv:2505.00024</i> .	
759		
760		
761		
762		
763	Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, and 1 others. 2023. Recognize anything: A strong image tagging model. <i>arXiv preprint arXiv:2306.03514</i> .	
764		
765		
766		
767		
768	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. LlamaFactory: Unified efficient fine-tuning of 100+ language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , Bangkok, Thailand. Association for Computational Linguistics.	
769		
770		
771		
772		
773		
774		
775		
776	Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. 2017. Scene parsing through ade20k dataset. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> .	
777		
778		
779		
780		
781	Hengguang Zhou, Xirui Li, Ruochen Wang, Minhao Cheng, Tianyi Zhou, and Cho-Jui Hsieh. 2025. R1-zero’s "aha moment" in visual reasoning on a 2b non-sft model. <i>Preprint</i> , arXiv:2503.05132.	
782		
783		
784		

A Dataset Construction Details

We leverage the SAT dataset and the CoTA dataset from TACO for our training.

The SAT (Ray et al., 2024) dataset is a synthetic VQA dataset designed to enhance the spatial reasoning capabilities of Multimodal Large Language Models (MLLMs).

The CoTA dataset from TACO (Ma et al., 2024) is a synthetic dataset comprising Chain-of-Thought-Action data generated by GPT-4o. However, we leverage it solely as a VQA filter for potential tool utilization.

The CoTA data comprises trajectories of thought, action, and observation. We derive QA pairs from this data using Qwen3-30B-A3B. For open-ended questions within this set, we employ Qwen2.5-VL-32B to synthesize multiple-choice options, thereby transforming the entirety of the CoTA data into a MCQA format. Subsequently, we task Qwen2.5-VL-7B with answering all the questions and retain only those that are answered incorrectly. These erroneous examples are more likely to exhibit both correct and incorrect responses during the sampling of rollouts in GRPO training. The resulting dataset is then randomly split into a cold-start dataset and a RL dataset in a 1:9 ratio. For the cold-start dataset, we utilize GPT-4.1 as the agent to synthesize tool-integrated reasoning rollouts, filtering out any rollouts that lead to incorrect answers.

The SAT data is inherently in a MCQA format. We randomly shuffled the answer options and subsequently filtered out data instances for which Qwen2.5-VL-7B provided an incorrect response.

B Experiment Setup Details

B.1 Benchmark and Dataset Details

In this paper, we evaluate five multi-modal benchmarks covering diverse visual reasoning capabilities: CV-Bench (Tong et al., 2024a), Blink (Fu et al., 2024), Blink-Hard (Bigverdi et al., 2024), MMVP (Tong et al., 2024b), and a 1,000 randomly selected subset from SAT (Ray et al., 2024). Our dataset construction incorporates single-image samples from both SAT (Ray et al., 2024) and Taco (Ma et al., 2024). The following sections provide detailed descriptions of these benchmarks and datasets:

- **CV-Bench** (Tong et al., 2024a): A vision-centric benchmark containing 2,638 manually-inspected examples for evaluating MLLMs. It repurposes

standard vision benchmarks (ADE20K (Zhou et al., 2017), COCO (Lin et al., 2014), Omni3D (Brazil et al., 2023)) to assess fundamental 2D and 3D understanding capabilities including spatial relationships, object counting, depth ordering, and relative distance estimation. Questions are programmatically constructed and manually verified for clarity and accuracy.

- **Blink** (Fu et al., 2024): A benchmark with 3,807 multiple-choice questions reformatting 14 classic computer vision tasks to test fundamental visual perception abilities. Despite humans achieving 95.70% accuracy, even top models like GPT-4V and Gemini achieve only 51.26% and 45.72% accuracy respectively. Blink highlights significant gaps between human-level visual perception and current MLLM capabilities.
- **Blink-Hard** (Bigverdi et al., 2024): A challenging benchmark focused on visual perception tasks requiring deeper 3D understanding and spatial reasoning. It evaluates whether MLMs can produce and reason with intermediate representations like depth maps and bounding boxes. The benchmark serves as a testbed for approaches like Perception Tokens that assist models in tackling complex visual reasoning problems.
- **MMVP** (Tong et al., 2024b): Contains 300 questions based on 150 pairs of “CLIP-blind” images that appear similar to CLIP models despite clear visual differences. The benchmark categorizes questions across nine visual patterns and reveals systematic shortcomings in MLLMs’ visual perception abilities, with even top models showing performance gaps of over 50% compared to humans.
- **SAT** (Ray et al., 2024): Contains 218K question-answer pairs covering 22K synthetic scenes testing both static and dynamic spatial reasoning. Unlike benchmarks focused on static relationships, SAT (Ray et al., 2024) incorporates perspective-taking and egocentric action recognition. Generated using a photo-realistic physics engine, it can be scaled and extended to include new scenarios.
- **Taco** (Ma et al., 2024): A framework and dataset with over 1 million synthetic chains-of-thought-and-action traces (filtered to 293K high-quality examples). TACO (Ma et al., 2024) enables models to perform step-by-step reasoning while invoking external tools (OCR, depth estimation,

calculators) to solve complex visual problems, showing performance improvements of up to 15% on challenging multimodal tasks.

B.2 Training Details

We fine-tune Qwen2.5-VL-3B on eight A800 GPUs. The detailed training parameters for cold-start and reinforcement learning are listed in Table 3 and Table 4

Table 3: Training hyperparameters

Name	Value
bf16	True
tf32	True
per_device_train_batch_size	4
gradient_accumulation_steps	2
lr	1e-5
weight_decay	0
warmup_ratio	0.1
lr_scheduler_type	cosine
max_seq_len	16384

Table 4: Training hyperparameters

Name	Value
bf16	True
tf32	True
per_device_train_batch_size	8
ppo_mini_batch_size	128
num_generation	8
kl_loss_coef	1e-3
lr	1e-6
weight_decay	0
warmup_ratio	0.03
lr_scheduler_type	cosine
max_seq_len	16384

The hyperparameters used for the agent are detailed in Table 5. During training with the veRL framework, we observed a degradation in the response quality when batch-inferencing with Qwen2.5-VL by vllm as the batch size increased. Consequently, we imposed a limit on the batch size used for each inference step.

Table 5: Agent config

Name	Value
max_turns	4
max_tokens_per_turn	700
inference_batch_size	8

C Case Study

We show some case studies that our model successfully using tools in Figure 7 and 9. We also include some failure cases that visual tools return wrong output and disturb model’s reasoning chain in 8.

Table 6: Comparison of Performance

Model	CV-Bench					MMStar						
	Count	Relation	Depth	Distance	Avg.	Corse	Fine	Instance	Logical	Math	Science	Avg.
Qwen2.5-VL-7B-Instruct	68.78	84.62	85.83	76.50	78.31	66.80	53.20	65.20	61.60	56.40	40.00	57.20
Qwen2.5-SFT	40.36	65.69	76.17	82.17	64.25	61.20	42.00	58.40	53.20	48.00	33.20	49.33
Qwen2.5-SFT-Cold-Start	70.18	87.08	82.50	73.83	77.98	60.00	53.20	59.20	55.60	55.20	29.60	52.13
TACO-qwen2-7B	63.32	81.08	59.83	57.26	65.52	68.00	50.40	60.40	42.40	45.60	30.40	49.53
ReVPT-7B	71.32	89.23	82.67	74.33	79.00	64.40	46.40	56.00	58.80	66.00	34.40	54.33

Question: Considering the relative positions of the bottle (annotated by the red box) and the fork in the image provided, where is the bottle (annotated by the red box) located with respect to the fork?

(A) left (B) right

Answer: (B)



Assistant:

<think> bottle (red box) is located on the table near the center of the image. The fork is also visible, likely placed on the table as well. To determine their relative positions, I should analyze the locations of both objects. The bottle appears to be more towards the middle-right side of the table, while the fork is closer to the center or slightly left. To confirm, I will use the edge detection tool to better distinguish the boundaries and positions of these objects.</think>

<tool_call>{"name": "edge_detection", "arguments": {"image_id": 0}}</tool_call>

User: <image>

<result>The edge map for image 0.</result>

Assistant:

<think>From the edge detection result, it is clear that the bottle (in the red box) is positioned more towards the center-right area of the table. The fork is visible beside the bottle, further to the right. This confirms that the bottle is to the right of the fork.</think>

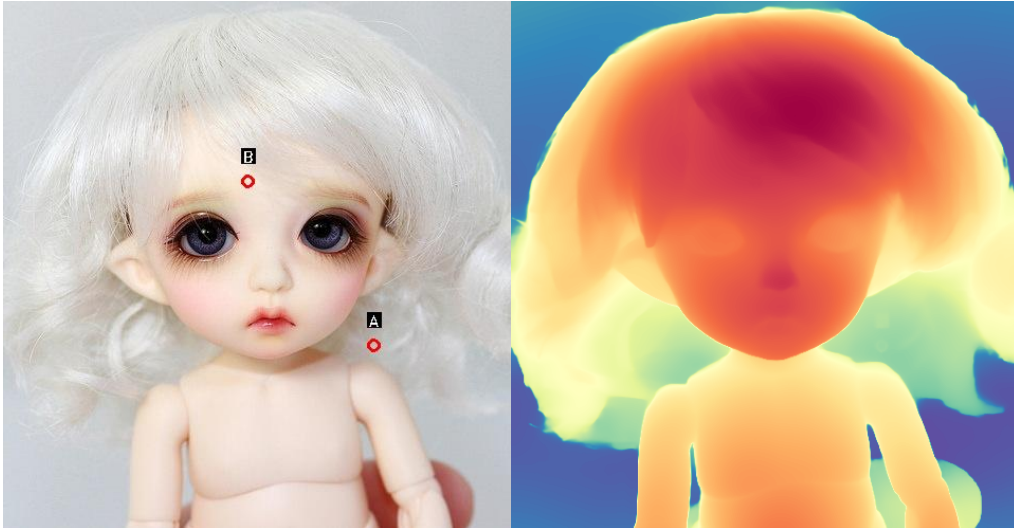
<answer>Based on the image and the edge detection result, the bottle (red box) is located to the right of the fork. Therefore, the answer is \boxed{B}. </answer>

Figure 7: Case study for the tool type: Edge Detection.

Question: Two points are circled on the image, labeled by A and B beside each circle. Which point is closer to the camera? Select from the following choices.

(A) A is closer (B) B is closer

Answer: (B)



Assistant:

<think>Looking at the image, both points A and B are marked near the face of the doll. To determine which one is closer to the camera, I need to analyze their positions relative to the viewer's perspective. Point A appears lower on the doll's face compared to point B, which is higher up on the forehead. However, this could be due to the camera angle or perspective. To be more certain, I should use depth estimation to compare the actual spatial distances of points A and B from the camera.</think>

<tool_call>{"name": "depth_estimation", "arguments": {"image_id": 0}}</tool_call>

User: <image>

<result>The edge map for image 0.</result>

Assistant:

<think>In the depth map, warmer colors indicate areas that are closer to the camera, while cooler colors represent regions farther away. Observing the map, the region around point A (near the lower part of the face) appears much warmer compared to the area around point B (higher on the forehead). This confirms that point A is indeed closer to the camera than point B.</think>

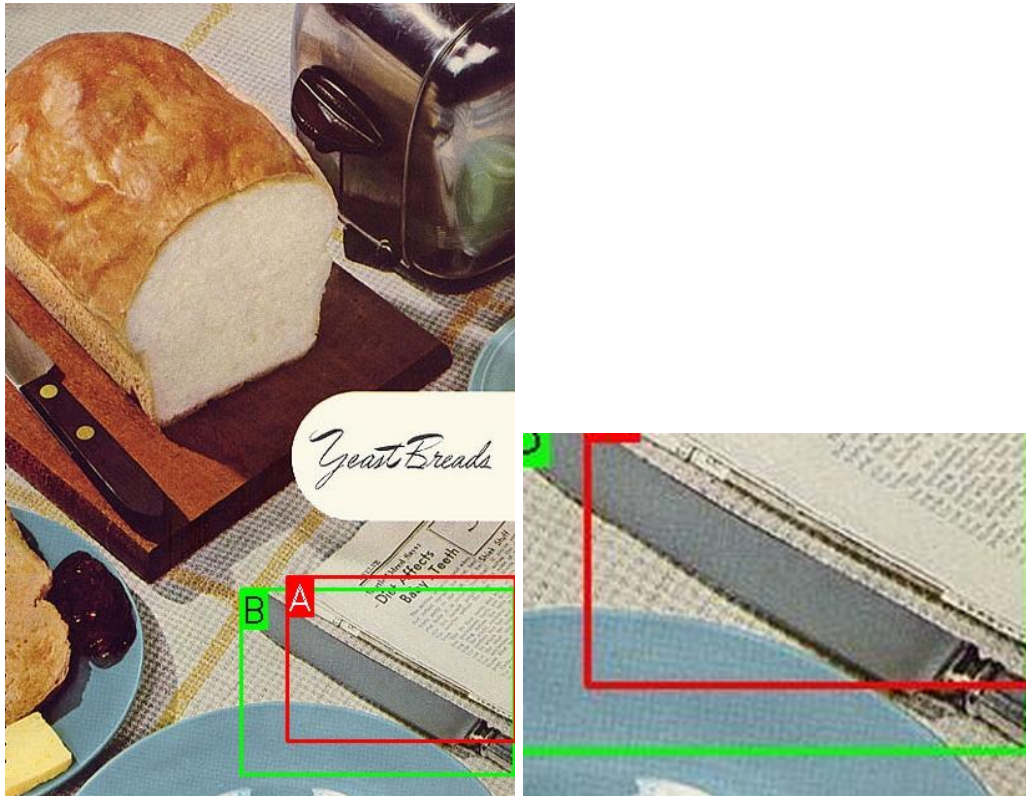
<answer>Based on the depth map, point A is closer to the camera than point B. The warm color at point A supports this conclusion. \boxed{A} </answer>

Figure 8: Case study for the tool error: Depth Estimation.

Question: A bounding box is an annotated rectangle surrounding an object. The edges of bounding boxes should touch the outermost pixels of the object that is being labeled. Given the two bounding boxes on the image, labeled by A and B, which bounding box more accurately localizes and encloses the knife? Select from the following options.

(A) Box A (B) Box B

Answer: (B)



Assistant:

*<think>*Let's carefully examine the positions of the two bounding boxes, A and B. Box A is around what appears to be a book or a stack of papers, while Box B is around a silver object that resembles a knife. The knife is clearly visible and positioned in Box B, whereas Box A does not contain any knife-like object. To confirm, I'll double-check the boundaries of both boxes and their contents.*</think>*

<tool_call>{ "name": "zoom_in", "arguments": { "image_id": 0, "bbox": [200, 490, 480, 720], "factor": 1.5 } }*</tool_call>*

User: *<image>*

*<result>*Zoomed image 0 on [200, 490, 480, 720] with 1.5x magnification.*</result>*

Assistant:

*<think>*The zoomed-in view shows the region where both bounding boxes A and B overlap. By closely inspecting this area, I can see that Box B encompasses the handle and blade of a knife, which is not present in Box A. This confirms my initial assessment. Therefore, Box B is the one that correctly encloses the knife.*</think>*

*<answer>*The bounding box labeled B more accurately localizes and encloses the knife, as it includes the knife's handle and blade. Therefore, the correct answer is *\boxed{B}*. *</answer>*

Figure 9: Case study for the tool type: Zoom In.