# Open Vision Reasoner: Transferring Linguistic Cognitive Behavior for Visual Reasoning

Yana Wei<sup>1,\*</sup>, Liang Zhao<sup>2,\*†</sup>, Jianjian Sun<sup>2,\*</sup>, Kangheng Lin<sup>3</sup>, Jisheng Yin<sup>4</sup>, Jingcheng Hu<sup>5</sup>, Yinmin Zhang<sup>2</sup>, En Yu<sup>6</sup>, Haoran Lv<sup>2</sup>, Zejia Weng<sup>2</sup>, Jia Wang<sup>2</sup>, Qi Han<sup>2</sup>, Zheng Ge<sup>2</sup>, Xiangyu Zhang<sup>2</sup>, Daxin Jiang<sup>2</sup>, Vishal M. Patel<sup>1†</sup>

<sup>1</sup>Johns Hopkins University <sup>2</sup>StepFun <sup>3</sup>BUPT <sup>4</sup>UCAS <sup>5</sup>THU <sup>6</sup>HUST Open-Vision-Reasoner.github.io

#### **Abstract**

The remarkable reasoning capability of large language models (LLMs) stems from cognitive behaviors that emerge through reinforcement with verifiable rewards. This work investigates how to transfer this principle to Multimodal LLMs (MLLMs) to unlock advanced visual reasoning. We introduce a two-stage paradigm built on Qwen2.5-VL-7B: a massive linguistic cold-start fine-tuning, followed by multimodal reinforcement learning (RL) spanning nearly 1,000 steps—surpassing all previous open-source efforts in scale. This pioneering work reveals three fundamental insights: 1) Behavior transfer emerges surprisingly early in cold start due to linguistic mental imagery. 2) Cold start broadly memorizes visual behaviors, while RL critically discerns and scales up effective patterns. 3) Transfer strategically favors high-utility behaviors such as visual reflection. Our resulting model, *Open-Vision-Reasoner* (OVR), achieves state-of-the-art performance on a suite of reasoning benchmarks, including 95.3% on MATH500, 51.8% on MathVision and 54.6% on MathVerse. We release our model, data, and training dynamics to catalyze the development of more capable, behavior-aligned multimodal reasoners.

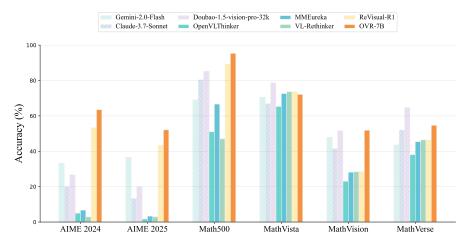


Figure 1: **Performance comparison** with state-of-the-art models on both textual (AIME 2024, AIME 2025 [7], MATH500 [38]) and multimodal (MathVista [64], MathVision [100], MathVerse [133]) math reasoning benchmarks. **Open Vision Reasoner (OVR)** demonstrates superior results among open-source models and performs competitively with commercial counterparts.

<sup>\*</sup>Core contribution, †Corresponding authors: zhaoliang02@stepfun.com, vpatel36@jhu.edu

# 1 Introduction

"The eye sees only what the mind is prepared to comprehend." — Robertson Davies

Shifting Reinforcement Learning from Human Feedback (RLHF) [74] to Reinforcement Learning from Verifiable Reward (RLVR) [34, 91] has endowed LLMs [34, 72] with unexpectedly powerful reasoning across mathematics, code, and general problem-solving. At its core, verifiable reward—where correctness is determined by objective, often rule-based criteria—is inherently less susceptible to "reward hacking" [87, 124] than a learned reward model. This robustness proves instrumental in large-scale RL, enabling the internalization and activation of what recent studies [29, 126, 135] term **cognitive behaviors**—patterns like backtracking and subgoal decomposition that are empirically crucial for advanced reasoning.

The multimodal domain, inherently grounded in verifiable visual facts [10, 125], is uniquely suited for this paradigm. Yet, early multimodal RL efforts paradoxically adopted RLHF, relying on learned reward models to approximate objective correctness [98, 138, 139]. Inspired by the success of RLVR in language models, recent efforts have started exploring rule-based rewards in the multimodal setting. Perception-R1 [125] incorporates supervisions such as IoU and Euclidean distance to enhance the perceptual alignment of MLLMs, while works such as R1-OneVision [116] and VLAA-Thinking [8] construct behavior-rich visual reasoning trajectories through complex pipelines including iterative distillation and synthesizing. Recently, ReVisual-R1 [14] adopts a effective language-only cold start as a foundation for visual reasoning.

Despite this encouraging progress, these approaches still leave a foundational question unanswered: **How can linguistic cognitive behaviors transfer to MLLMs for advanced visual reasoning?** To address this, we build upon the "RL with a cold start" paradigm [34] by conducting large-scale training on Qwen2.5-VL-7B [5], establishing it as a powerful testbed to systematically analyze how such behaviors emerge and scale in the multimodal domain.

To this end, we introduce a robust two-stage methodology designed to first instill linguistic cognitive patterns and then activate them for visual reasoning. Our process begins with a large-scale cold start, fine-tuning Qwen2.5-VL-7B on over 2 million examples to build a strong foundation. This is followed by a prolonged reinforcement learning phase under the Open-Reasoner-Zero [40] framework, leveraging over 0.3 million mixed-modality examples. To the best of our knowledge, this represents the largest open-source RL practice on this model. The resulting model, Open-Vision-Reasoner (OVR), validates our approach by achieving strong performance across both language and multimodal benchmarks. As shown in Fig. 1, it achieves 63.5% on AIME2024 and 95.3% on MATH500 for math reasoning, as well as 51.8% on MathVision and 54.6% on MathVerse for visual reasoning.

To further trace the transfer and evolution of cognitive patterns throughout training, we develop a in-depth **visual cognitive behavior analysis**. Three central insights are worth highlighting: (1) Behavior transfer emerges remarkably early in cold start, driven by linguistic patterns encoding *mental imagery* [26, 45] as illustrated in Fig. 11. (2) Cold start broadly *memorizes* diverse visual cognitive behaviors, while RL critically *discerns* and scales up effective patterns. (3) Transfer follows a *strategic* path, favoring behaviors with high utility such as visual reflection. These findings deepen the understanding on visual intelligence scaffolded by linguistic reasoning [24].

We further examine how this paradigm impacts a foundational capability of MLLMs—visual perception. While linguistic cold start introduces perceptual degradation, our study shows that multimodal RL can effectively *recover* this loss. However, we also observe the *limited scalability* of RL when focused solely on perceptual tasks, as reward signals increase without corresponding growth in reasoning complexity (e.g., token length). This limitation motivates a more deliberate integration of diverse, primitive visual cognitive behaviors. Such efforts represent a promising direction toward unlocking the potential of more advanced RL frameworks—multi-turn or even agentic RL built upon visual manipulation and imagination.

In summary, this paper advances the field through the following three key contributions:

- We construct a two-stage training pipeline consisting of a linguistic cold start followed by large-scale multimodal RL, enabling effective transfer of cognitive behaviors in MLLMs.
- Our *Open Vision Reasoner*, the largest open-source RL practice on Qwen2.5-VL-7B, achieves superior performance on both linguistic and multimodal reasoning benchmarks.

Table 1: **Visual Cognitive Behaviors and Linguistic Counterparts.** We define four key visual cognitive behaviors, providing formal definitions, illustrative examples, and their corresponding linguistic counterparts.

Visual Behavior	Example	Definition	Linguistic Counterpart		
Visual Reflection	"Let me see the image again."	The model explicitly revisits the image after identifying a potential mistake or inconsistency in its reasoning, indicating an effort to correct course.	Backtracking		
Divide-and-Conquer	"Let's first look at the numbers on the left."	r r			
Visual Verification	"I will now verify this against the image."  The model confirms that its intermediate conclusions are visually grounded by cross-referencing with the image before proceeding.		Verification		
Goal-driven Visual Tracing	"To get this answer, I need to find an object that"	The model starts from a desired visual conclusion and reasons backwards to identify relevant image evidence that supports it.	Backward Chaining		

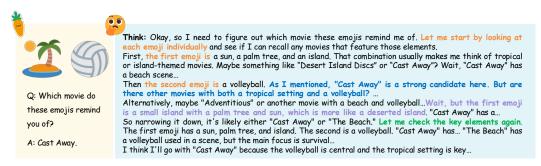


Figure 2: Multiple Cognitive Behaviors in a Single Response. This case shows triggered visual-specific cognitive behaviors like visual divide-and-conquer, reflection, goal-driven visual tracing, along with the linguistic behavior backtracking.

 We conduct an in-depth analysis of visual cognitive behaviors in OVR and provide valuable insights into their transfer and evolution across training stages.

# 2 Cognitive Behavior Preliminaries

Recent studies have highlighted that the emergence of robust reasoning in LLMs is closely tied to the acquisition of certain *cognitive behaviors* [29]. These behaviors reflect structured internal reasoning patterns akin to human problem-solving. Key examples include: (1) Backtracking — revising a previously chosen strategy upon identifying inconsistencies (e.g., "This approach won't work because..."), (2) Verification — checking intermediate steps or partial results, (3) Subgoal Setting — breaking down complex problems into manageable components (e.g., "First, we need to..."), and (4) Backward Chaining — reasoning from the desired outcome back to required inputs (e.g., "To get 75, we need a number divisible by..."). The four patterns form a kind of textual inner monologue that emerges naturally in language-based reasoning.

Based on this, we investigate the *transfer* of cognitive behaviors from language to vision. We define the visual extensions of the aforementioned behaviors—visual reflection, divide-and-conquer, visual verification, and goal-driven visual tracing. Their formal definitions, examples, and corresponding linguistic counterparts are provided in Table 1, while Fig. 2 presents a multimodal example encompassing both linguistic and visual cognitive behaviors. In the following sections, we present a simple yet effective MLLM training pipeline comprising a linguistic cold start followed by multimodal RL (Section 3), and systematically analyze the transfer and scaling of these visual cognitive behaviors (Section 5.2).

# 3 Open Vision Reasoner

In this section, we introduce **Open-Vision-Reasoner** (OVR), a strong multimodal reasoning model build from Qwen2.5-VL-7B [42], from perspectives of training pipeline (Section 3.1), RL algorithm (Section 3.2) and data construction (Section 3.3).

## 3.1 Training Pipeline

To facilitate efficient cognitive development and cross-modal generalization, we employ the popular "RL with a cold start" paradigm [34] with two sequential training stages:

- Stage 1: Linguistic Cold Start. The LLM module is supervised fine-tuned on language-only reasoning datasets distilled from DeepSeek-R1 [34], establishing core cognitive behaviors such as backtracking and subgoal decomposition within a purely linguistic setting.
- Stage 2: Multimodal RL. We apply reinforcement learning with Open-Reasoner-Zero [40] setting on both text and multimodal tasks using verifiable match rewards. This promotes reasoning generalization and aligns previously learned cognitive patterns with visual contexts, enabling effective cross-modal transfer.

# 3.2 RL Algorithm

For the RL stage of our training pipeline, we adopt a lightweight Proximal Policy Optimization (PPO) [82] with Generalized Advantage Estimation (GAE) [81], following the policy and reward design used in *Open-Reasoner-Zero* [40]. We detail the RL for multimodal tasks below:

**Proximal Policy Optimization** For each input, consisting of an image I and a textual prompt q, the policy network  $\pi_{\theta}$  generates n responses  $\{o_1,\ldots,o_n\}$ . Each response  $o_i$  is a trajectory  $\tau_i=(s_0^{(i)},a_0^{(i)},\ldots,s_{T_{i-1}}^{(i)},a_{T_{i-1}}^{(i)})$  of length  $T_i$ . The state  $s_t^{(i)}$  includes q (and potentially encoded I features) and previously generated tokens;  $a_t^{(i)}$  is the token generated at step t. A reward  $r_t^{(i)}$  is computed at each timestep t of trajectory  $\tau_i$ .

We use GAE to balance bias and variance in advantage estimation. The advantage  $\hat{A}_t^{(i)}$  for state-action pair  $(s_t^{(i)}, a_t^{(i)})$  in trajectory  $\tau_i$  is:

$$\hat{A}_{t} = \sum_{l=0}^{T-t-1} (\gamma \lambda)^{l} \delta_{t+l}, \quad \text{where } \delta_{t'} = r_{t'} + \gamma V_{\phi}(s_{t'+1}) - V_{\phi}(s_{t'}). \tag{1}$$

 $V_{\phi}$  is the value function,  $\gamma$ ,  $\lambda$  are discount and GAE factors, and  $V_{\phi}(s_T) = 0$  for terminal states.  $\pi_{\theta}$  is updated by maximizing  $\mathcal{J}_{PPO}(\theta)$  using experiences  $(s_t, a_t, \hat{A}_t)$  sampled under an older policy  $\pi_{old}$ :

$$\mathcal{J}_{\text{PPO}}(\theta) = \hat{\mathbb{E}}_{\pi_{\text{old}}} \left[ \min \left( \rho_t(\theta) \hat{A}_t, \text{clip} \left( \rho_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right]. \tag{2}$$

Here,  $\rho_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\text{old}}(a_t|s_t)}$  and  $\epsilon$  is a clipping parameter (e.g., 0.2).  $\hat{\mathbb{E}}_{\pi_{\text{old}}}$  denotes the empirical average over samples from  $\pi_{\text{old}}$ .  $V_{\phi}$  is trained by minimizing  $\mathcal{J}_{\text{value}}(\phi)$  on samples from  $\pi_{\text{old}}$ , using the empirical discounted returns  $R_t = \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}$ :

$$\mathcal{J}_{\text{value}}(\phi) = \hat{\mathbb{E}}_{\pi_{\text{old}}} \left[ \left( V_{\phi}(s_t) - R_t \right)^2 \right]. \tag{3}$$

**Reward Function.** We adopt the minimalist rule-based reward design, which evaluates only the correctness of model outputs while ignoring formatting or stylistic preferences. Specifically, we extract the predicted answer encapsulated within \boxed{} in the model's output and compare it against the reference answer. A binary reward is assigned—1 for exact matches, and 0 otherwise—enabling a clear, scalable and unhackable reward signal for reinforcement learning.

Table 2: **Comparison on Language Reasoning and General Benchmarks.** Best results are **bold** and the second-best are underlined for *open-source models*. † indicates metrics reproduced by ourselves.

Model	<b>AIME 2024</b>	<b>AIME 2025</b>	MATH500	<b>GPQA Diamond</b>	MMLU	MMLU-Pro
Open-source Models						
Qwen2.5-7B [42]	$6.7^{\dagger}$	$6.7^{\dagger}$	$77.6^{\dagger}$	$32.8^{\dagger}$	$72.6^{\dagger}$	57.5 <sup>†</sup>
Qwen2.5-VL-7B [42]	$6.7^{\dagger}$	$6.7^{\dagger}$	$67.4^{\dagger}$	$31.8^{\dagger}$	$\overline{69.6}^{\dagger}$	51.7 <sup>†</sup>
Open-Reasoner-Zero-7B [40]	17.9	15.6	81.4	36.6	-	-
DeepSeek-R1-Distill-Qwen-7B [34]	55.5	$39.2^{\dagger}$	92.8	49.1	-	-
QwQ-32B-Preview [96]	50.0	33.5	90.6	<u>54.5</u>	-	-
Skywork-R1V-38B [75]	72.0	-	94.0	61.6	-	-
ReVisual-R1 [14]	53.3	43.3	89.2	47.5	-	-
Close-source Models						
Gemini-2.0-Flash [90]	33.4	36.7	69.0	35.4	-	-
OpenAI-o1-mini [72]	63.6	-	90.0	60.0	85.2	80.3
Claude 3.7 Sonnet [3]	20.0	13.3	80.4	61.1	-	80.0
Doubao-1.5-vision-pro-32k [1]	26.7	20.0	85.2	56.1	-	-
OVR-7B	<u>63.5</u>	52.1	95.3	49.8	77.2	67.9

#### 3.3 Dataset Construction

To support cognitive transfer, we carefully curate datasets specifically tailored to each training stage, encompassing both language-only and multimodal domains.

**Data Collection.** We firstly broadly collect prompt-answer pairs to develop both language and multimodal reasoning skills across mathematical, scientific, and logical domains. For language-only scenarios, we utilize public benchmarks including AIME (up to 2023), MATH [38], Numina-Math [52], Tulu3 MATH [48], and OpenR1-Math-220k [2], and other open-source datasets. We also synthesize general logical problems via programmatic generation to further enrich reasoning diversity. Multimodal scenarios incorporate datasets covering geometry problem solving (Geometry3k [61], GeoQA [9], Geos [112]), visual discrimination (IconQA [62], Pixmo [21], ChartQA [67]), visual puzzles (PuzzleVQA [19], AlgoPuzzleVQA [30]), STEM (TQA [46], ScienceQA [63], K12 from [68]) and multimodal math (AtomThink [109], in-house curated math).

**Data Curation.** To refine data quality, we employ a multi-step curation process. *Firstly*, we employ a pre-trained model to automatically filter out samples with high training loss, which typically indicate noise or excessive complexity. *Secondly*, rule-based and model-assisted methods then identify and remove undesirable patterns [52]. *Thirdly*, we apply reweighting to balance coverage, down-weighting overrepresented categories while emphasizing rare but valuable instances. To the end, we distill responses from DeepSeek-R1 [34] to construct approximately 2 *million cold-start data*. To ensure the unhackability and stability during RL, we further exclude problems incompatible with our reward functions (e.g., proof-style questions) and apply difficulty-based heuristic filtering, removing both overly trivial and infeasible samples to ensure well-calibrated learning. This leaves around *300k multimodal RL data*. Further details refer to the appendix.

# 4 Experiments

In this section, we first elaborate our implementation of *Open-Vision-Reasoner* (OVR). Then, we present superior performance across textual benchmarks (Section 4.2) and multimodal scenerios (Section 4.3).

#### 4.1 Implementation Details

Our model is based on Qwen2.5-VL-7B [5] and employs a two-stage training strategy. In the first stage of cold start, we independently fine-tune the LLM module for 5 epochs with a batch size of 640, a sequence length of 64k, and a learning rate of  $2 \times 10^{-4}$  leveraging the default Qwen2.5 configuration [42]. During the subsequent stage of reinforcement learning, following Open-Reasoner-Zero [40], we utilize PPO and configure GAE with  $\gamma=1$  and  $\lambda=1$  to fully capture long-term dependencies crucial for reasoning tasks, enabling stable training. This RL phase proceeds for 900 iterations, during which we adopt a curriculum for the sequence length: it begins at 24k for the first 300 iterations, increases to 32k through iteration 700, and expands to 48k thereafter, with our latest

Table 3: **Evaluation Results on Visual Reasoning Benchmarks.** Best results are **bold** and the second-best are underlined for *open-source models*. † Indicates results reproduced by ourselves.

Model	MathVista	MathVision	hVision MathVerse		WeMath		LogicVista	MMMU-Pro	Cha	rXiv
Wiodei			vision-only	worst	strict	loose			reas.	desc.
SFT Methods										
LLaVA-OneVision-7B [51]	62.6	17.6	17.6	9.0	17.7	-	32.0	24.1	23.6	48.7
InternLM-XComposer2.5 [132]	64.0	17.8	16.2	8.2	14.1	-	34.7	-	-	-
InternVL3-8B [137]	70.5	28.6	33.9	23.0	37.5	-	43.6	-	37.6	73.6
InternVL2.5-8B [15]	64.5	17.0	22.8	9.4	23.5	-	36.0	34.3	32.9	68.6
InternVL2-8B [16]	58.3	20.0	20.4	9.2	20.2	-	33.6	29.0	-	-
Qwen2-VL-7B [101]	61.6	19.2	25.4	11.0	22.3	-	33.3	30.5	34.6	58.0
Qwen2.5-VL-7B [42]	$69.2^{\dagger}$	$25.5^{\dagger}$	41.1	21.8	$31.2^{\dagger}$	$53.1^{\dagger}$	47.9	-	36.4 <sup>†</sup>	$67.3^{\dagger}$
QvQ-72B-Preview [95]	70.3	34.9	48.2	30.7	39.0	-	58.2	-	-	-
Kimi-VL-16B [93]	66.0	21.8	34.1	18.0	32.3	-	42.7	-	-	-
Close-source Models										
Gemini-2.0-Flash [90]	70.4	47.8	43.6	42.1	47.4	-	52.3	-	-	-
OpenAI-GPT-4o [43]	59.9	31.1	40.6	34.5	42.9	-	64.4	-	-	-
Claude 3.7 Sonnet [3]	66.8	41.3	52.0	39.7	58.2	-	49.3	-	-	-
GPT-4o mini [31]	55.1 <sup>†</sup>	27.3 <sup>†</sup>	$30.0^{\dagger}$	$31.6^{\dagger}$	$31.4^{\dagger}$	$48.8^{\dagger}$	$41.4^{\dagger}$	37.6 <sup>†</sup>	$34.10^{\dagger}$	74.92 <sup>†</sup>
doubao-1.5-vision-pro-32k [1]	78.6	51.5	64.7	44.9	64.2	-	65.7			
RL-based Methods										
VLAA-Thinker-Qwen2-7B [8]	59.6	19.8	33.9	15.2	30.5	-	36.0	-	-	-
VLAA-Thinker-Qwen2.5-7B [8]	68.0	26.4	48.2	22.4	41.5	-	48.5	-	-	-
R1-Onevision-7B [116]	64.1	29.9	40.0	-	-	61.8	-	-	-	-
OpenVLThinker-7B [22]	65.3	23.0	38.1	16.8	35.2	-	44.5	-	-	-
MM-Eureka-Qwen-7B [68]	72.6	28.1	45.4	23.0	21.8	-	46.3	-	-	-
MMR1-Math-v0 [76]	69.8	30.7	42.8	17.4	31.9	-	46.8	-	-	-
ThinkLite-7B-VL [102]	71.6	24.6	42.9	16.5	41.8	-	42.7	-	-	-
R1-VL-7B [131]	63.5	24.7	40.0	-	-	-	-	-	-	-
X-REASONER [58]	69.0	29.6	-	-	-	-	-	43.0	-	-
VL-Rethinker-7B [99]	73.7	28.4	46.4	17.8	36.3	-	42.7	41.7	-	-
ReVisual-R1 [14]	73.1	48.8	53.6	27.5	42.0	-	52.3			
WeThink [115]	70.9	27.2	44.7	24.4	48.0	-	53.0			
Skywork-R1V-38B [75]	60.6	42.1	40.4	-	34.1	-	50.6	-	-	-
OVR-7B	72.1	51.8	54.6	33.5	44.6	64.8	54.8	50.2	44.5	73.6

models continuously undergoing this refinement process. We adhere to strict on-policy updates for the policy model and undertake multiple optimization steps for the critic model. Please note that our final model is an **uniform average of several representative intermediate checkpoints**, ensuring balanced and robust performance across various benchmarks. Additional details can be found in the appendix.

# 4.2 Enhanced Language Reasoning and General Capabilities

Our model is initially evaluated on a variety of language benchmarks, which cover mathematical reasoning and general problem-solving skills. Specifically, we include *AIME 2024*, *AIME 2025* [7], *MATH500* [38], *GPQA Diamond* [79], *MMLU* [36], and *MMLU-Pro* [103]. We compare *Open-Vision-Reasoner (OVR)* with strong LLM baselines, including *Qwen2.5-7B* [113], *DeepSeek-R1-Distill-Qwen-7B* [34] and *Open-Reasoner-Zero-7B* [40].

OVR demonstrates exceptional language reasoning capabilities. On the challenging AIME 2024 and 2025 benchmarks, it dramatically **surpasses other 7B open-source models by an average of over 10%**, achieving performance comparable to leading 32B models. This superiority extends to general reasoning tasks, with significant gains of **+4.6%** on MMLU and **+10.4%** on MMLU-Pro over parameter-matched competitors. These results highlight the effectiveness of our curated, high-quality cold-start training data.

#### 4.3 Superior Visual Reasoning Abilities

To evaluate whether the introduced cognitive behavior transfer leads to cross-modal benefits, we further assess the model on a suite of multimodal reasoning benchmarks. These tasks involve image-grounded mathematical reasoning, general multimodal reasoning, and chart understanding. Specifically, we include *MathVista* [64], *MathVision* [100], *MathVerse* [133], *DynaMath* [140], *WeMath* [77], *LogicVista* [110], *MMMU-Pro* [129], and *CharXiv* [105] for evaluation. We compare our model against strong MLLM baselines, including SFT-based methods such as *LLaVA-OneVision* [50] and

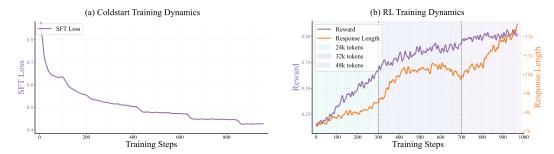


Figure 3: **Training Dynamics.** (a) The cold-start stage shows a step-wise loss decrease. (b) In the RL stage, reward (purple, left axis) and average response length (orange, right axis) grow steadily, with sharp surges after each sequence length expansion.

Qwen2.5-VL [5], as well as recent rule-based RL methods like OpenVLThinker [22], MM-Eureka [68] and ReVisual-R1 [14].

As shown in Table 3, our model sets a new breakthrough for 7B models in visual reasoning. It is **the first post-trained Qwen2.5-VL-7B-based model to surpass the 50% performance on MathVision**, while also achieving state-of-the-art results among 7B models on DynaMath and MathVerse. This strong overall performance is further underscored by a substantial gain on MMMU-Pro (+7.2% over prior SOTA methods). These results indicate that reasoning capabilities acquired through language training can effectively transfer to multimodal tasks, leading to notable improvements in visual reasoning.

### 5 Discussion

# 5.1 Analysis of Training Dynamics

In this section, we present a comprehensive overview of the training dynamics as illustrated in Fig. 3, and provide a detailed analysis of how text and multi-modal reasoning metrics evolve throughout the process as shown in Fig. 4.

During the initial cold-start phase (Fig. 3 (a)), the model's loss rapidly descends to below 0.5. Subsequently, across multiple training epochs, the loss exhibits a step-wise, gradual decrease. In parallel, we observe a corresponding surge in performance across all benchmarks (Fig. 4), which first ascend sharply before transitioning to a phase of slower, more incremental improvement toward their peak. A noteworthy observation is that the *aggressive* training strategy detailed in Section 4.1—employing a large batch size in concert with a high learning rate—proves to be essential. This approach is critical for breaking the model's inherent constraints, thereby successfully imbuing it with new cognitive paradigms and **sculpting a more favorable landscape for reinforcement learning**. It is a prerequisite that enables our model, which originates from an instruction-tuned base, to ultimately achieve text performance that is comparable to, or even surpasses models initialized from base [42] or math-specific checkpoints [114].

Furthermore, Fig. 3 (b) reveals how the model's reward and average token length in the RL phase steadily advance from an initial 7k to exceed 12k. Owing to the stability of the training configuration inherited from *Open-Reasoner-Zero* [40], OVR is successfully trained on a diverse corpus of over 20 multi-modal and language-only datasets without encountering any training collapse or performance degradation. Critically, whenever the token length begins to plateau or even decline, we strategically switch to a longer context length, which invariably catalyzes the next wave of rapid reward growth. Fig. 4 captures the coincident yet unsurprising convergent growth trajectory shared by all eight reasoning benchmarks, spanning both text and multi-modal domains, as they progressively ascend towards their zenith amidst fluctuations.

# 5.2 Multimodal Cognitive Behavior Analysis

Recent studies have highlighted the emergence of cognitive behaviors in MLLMs during visual reasoning tasks—phenomena often dubbed "visual aha moments" [8, 13]. In this work, we move

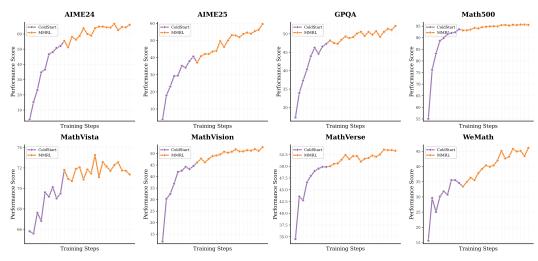


Figure 4: **Performance Evolution on Reasoning Benchmarks.** OVR demonstrates sustained and convergent growth across both linguistic and multi-modal benchmarks throughout the cold start (left) and RL (right) training phases.

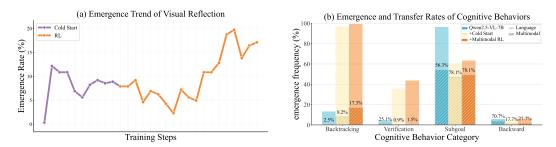


Figure 5: **Multimodal Cognitive Behavior Analysis.** (a) Emergence of visual reflection across the cold start and RL training steps. (b) Emergence and transfer rates of four visual cognitive behaviors across base models and training stages. Numerical values denote the language-to-vision transfer rates for each behavior.

beyond plain observations and systematically investigate how these behaviors are transferred from their linguistic counterparts. Our analysis centers on the four pivotal visual cognitive behaviors introduced in Section 2 which are drawn from foundational research on cognitive patterns [29]. To quantify this process, we employ GPT-40 [31] to analyze the emergence of each behavior within the inference traces of our OVR model.

Visual behaviors emerge remarkably early from cold start. Following Section 5.1, we tracked the dynamics of visual reflection, a significant behavior mentioned in previous studys [47, 107], throughout OVR's training. As depicted in Fig. 5, this vision-specific behavior emerges in significant quantities from the *very beginning* of the cold-start phase and fluctuates throughout subsequent training steps. Strikingly, we observed that even in linguistic problems, DeepSeek-R1's responses [34] frequently exhibited signs of *mental imagery* [26, 45] as shown in Fig. 11(a). The model appeared to construct internal visualizations to support mathematical reasoning, often articulated through phrases such as "*let me visualize*…" or "*let me see the image*." Once this linguistic scaffolding was introduced into our MLLM, these *mental images* were rapidly grounded in actual visual input, enabling their rapid and effective generalization within OVR.

**Cold-start learns broadly, large-scale RL discerns critically.** We further investigate how cognitive behaviors scale during large-scale RL. As shown in Fig. 5(a), after an initial, rapid instillation of patterns during the aggressive cold-start phase, their prevalence is *first suppressed then amplified to unprecedented levels* during multimodal RL. This counter-intuitive dynamic suggests a clear division of labor: the cold-start phase learns broadly, indiscriminately memorizing all available patterns. In

contrast, RL discerns critically, acting as a strategic filter for the crucial tokens [17] and scaling up pivotal behaviors. This process of RL—discarding the dross to select the essence—is significant for achieving superior generalization.

Visual transfer of cognitive behaviors is strategic. To analyze the transition from linguistic to visual cognition, we track the emergence and transfer rates (detailed in Appendix C.1) of four core cognitive behaviors across both language and vision modalities. As shown in Fig. 5(b), the emergence of backtracking and verification steadily increases across training stages, underscoring their growing importance. Among these, the transfer rate of backtracking shows consistent growth—from 2.5% to 17.3%—while verification exhibits near-zero transfer throughout both the cold-start and RL phases. This indicates that transfer is a *strategic* process, for which we posit two potential explanations: (1) Backtracking transfers more readily due to DeepSeek-R1's [34] inherent "mental imagination" capabilities, while verification, lacking a direct linguistic precursor, is more difficult for the MLLM to internalize. (2) Mirroring how humans naturally and instinctively process visual information [107], backtracking is a more *fundamental* component of complex visual reasoning, making its amplification a higher priority during the strategic RL phase. We will investigate these hypotheses in greater depth in our future work.

# 5.3 Beyond Behavior: Visual Perception Analysis and Future Work

Beyond behavioral dynamics, we extend our discussion to a essential capability of MLLMs: **visual perception** under the cold start plus large-scale RL paradigm. In particular, we investigate two key areas of interest—*perceptual hallucination* and *scaling properties*—through a dedicated study on our OVR model.

Cold start impairs perception, while RL enhances. We evaluated both stages of OVR, along with the base model Qwen2.5-VL-7B, on a comprehensive set of multimodal benchmarks targeting visual perception and recognition (MMBench [59], BLINK [28], MMStar [12], HallusionBench [32], POPE [53], RealWorldQA [108], MME [27], MMVet [127]). As shown in Table 4, performance steadily improves across tasks such as MMBench, underscoring the effectiveness of our training paradigm. The cold-start model shows declines on several tasks, notably increased *hallucinations* [56, 118], likely due to token distribution shifts from large-scale linguistic data [53]. However, the regained performance on benchmarks such as MMBench and BLINK demonstrate that long-term multimodal RL can effectively mitigate these issues by **discerning perceptual capabilities** that are critically for multimodal tasks. Looking ahead, degradation from cold start can be mitigated either by incorporating the linguistic data into model pretraining [33, 89], or by introducing more multimodal supervision during the cold start to establish a stronger visual foundation.

The current unscalability of RL for perception policy. Throughout the multimodal RL, we observed a strong correlation between the reward and the average response length in Fig. 3, which is a finding consistent with prior practices [34, 40]. This reinforces response length as an effective reward proxy, indicative of a *scaling property* tied to reasoning depth and computational resources. However, when focusing on specific discriminative perceptual tasks like OCR and counting, we observe a clear divergence. As shown in Fig. 6, while the reward can be effectively increased, the average response length remains largely stagnant.

This unscalable training dynamic on such challenging tasks hints at a more fundamental issue: *the absence of certain core visual cognitive behaviors*. Addressing this *fundamental capability gap* is paramount for achieving robust multimodal scaling. Emerging research offers promising avenues, such as multi-turn RL with agentic *tool-use* (*e.g.*, OpenAI-o3 [73]) and the integration of intrinsic imagining through *mental images* [117, 122]. These approaches hold the potential to bridge current limitations and unlock more scalable multimodal reasoning.

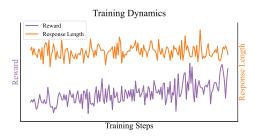


Figure 6: Training Dynamics on perception tasks including OCR and counting.

Table 4: Model Performance on Perception-centric Benchmarks.

Model	MMBench		BLINK	MMStar	HallusionBench	POPE	RealWorldOA	MME	MMVet
	en	cn	DELIVIK	Ministai	HanusionDenen	IOLE	RealWorldQA	MINIE	IVIIVI VCC
Qwen2.5-VL	$85.3^{\dagger}$	$84.6^{\dagger}$	$53.7^{\dagger}$	62.1 <sup>†</sup>	49.0 <sup>†</sup>	$86.1^{\dagger}$	$69.3^{\dagger}$	$1659.7^\dagger$	$63.9^{\dagger}$
+Cold Start	86.1	82.1	51.5	62.4	55.0	82.5	63.1	1549.8	61.8
+Multimodal RL	86.6	84.2	54.1	62.7	53.6	83.2	65.5	1559.1	63.6

# 6 Related Work

Recent breakthroughs like OpenAI's o1 [72] have highlighted the power of RL in unlocking and scaling reasoning capabilities [6, 37, 80] within LLMs. DeepSeek-R1-Zero [34] showed that reasoning capabilities can emerge purely through large-scale RL, leading to complex behaviors like self-verification and reflection. Open-source efforts like Open-Reasoner-Zero [40] further demonstrates that even minimalist RL approaches, such as vanilla PPO [82] with GAE [81] and simple rule-based rewards, can drive scaling in response length and benchmark performance on open-source models [42, 114].

MLLMs [4, 5, 23, 101, 130] have rapidly progressed from basic image captioning [57, 106] to more challenging reasoning tasks [123, 128, 134]. Early efforts primarily relied on supervised fine-tuning with Chain-of-Thought (CoT) datasets [111], while some explored explicit reflection [107] and self-correction [35] mechanisms to emulate human-like reasoning patterns. More recently, methods such as PerPO [139] and MDPO [98] adopt RL-based post-training approaches like DPO [78], where alignment is learned from paired positive/negative responses. These approaches generally follow the RL from Human Feedback (RLHF)[74, 126] or RL from AI Feedback (RLAIF) [49] paradigms, where signals from learned reward models or preference labels are utilized for optimization.

Inspired by the success of RLVR [34] in language models, MLLM research has shifted toward rule-based RL like GRPO [84] into the multimodal domain. This has led to two major lines of efforts: (1) designing task-specific reward objectives [2, 8, 125], and (2) constructing multimodal "thinking" datasets that embed cognitive behaviors within CoT sequences [22, 41, 68, 85]. Additionally, recent powerful MLLMs adopt a language-only cold start [14, 89], using verbal reasoning sequences as a foundation for subsequent multimodal learning. These approaches encourages human-like behaviors [29] or so-called "visual aha moments" in the model responses.

Despite these advances in MLLM, a fine-grained understanding of the underlying reasoning mechanisms remains less explored. In contrast, recent study [29] centered on LLMs posit that effective reasoning is causally linked to the model's acquisition of certain *cognitive behaviors*, such as verification, backtracking, subgoal setting, and backward chaining. The test-time studies have observed that invoking these patterns improve performance [70]. Entropy-based analysis further reveals that regions associated with cognitive tokens are critical for diverse and high-quality reasoning [17]. The multimodal work like Long-Perceptual-Thoughts [54] attempts to explicitly instill such patterns by synthesizing long-form multimodal CoT data.

# 7 Conclusion and Limitations

We propose a two-stage framework to investigate cognitive behavior in the multimodal domain. Combining a linguistic cold start with large-scale multimodal RL enables effective cross-modal transfer and scaling of cognitive patterns. The resulting model, *Open Vision Reasoner*, the largest open-source RL practice built on Qwen2.5-VL-7B, achieves strong performance across linguistic and visual benchmarks. Beyond performance, we systematically analyze visual cognitive behaviors and reveal how they emerge and evolve through training.

While the linguistic cold start provides a strong foundation for reasoning, a trade-off with perception remains. Although multimodal RL partly restores perceptual capability, further work is needed to uncover visual reasoning patterns that better balance or even jointly enhance both aspects. Beyond the observed correlations between behaviors and reasoning, establishing causal insights will further clarify the underlying mechanisms of multimodal reasoning. We hope our findings inspire future research on cognitively aligned multimodal agents and open up new possibilities for scaling vision-language reasoning through behavior-centered learning.

## References

- [1] D. AI. Doubao 1.5 vision pro 32k, 2025. URL https://www.volcengine.com/docs/82379/1553586.
- [2] L. B. Allal, L. Tunstall, A. Lozhkov, E. Bakouch, G. Penedo, and G. M. B. Hynek Kydlicek. Open r1: Evaluating Ilms on uncontaminated math competitions, Feb. 2025. URL https://huggingface.co/blog/open-r1/update-2.
- [3] Anthropic. Claude. https://www.anthropic.com/index/introducing-claude, 2023.
- [4] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [5] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [6] M. Balunovic, J. Dekoninck, I. Petrov, N. Jovanovic, and M. Vechev. Matharena: Evaluating llms on uncontaminated math competitions, february 2025. *URL https://matharena.ai*.
- [7] M. Balunović, J. Dekoninck, and M. V. Ivo Petrov, Nikola Jovanović. Matharena: Evaluating llms on uncontaminated math competitions, Feb. 2025. URL https://matharena.ai/.
- [8] H. Chen, H. Tu, F. Wang, H. Liu, X. Tang, X. Du, Y. Zhou, and C. Xie. Sft or rl? an early investigation into training r1-like reasoning large vision-language models. *arXiv preprint arXiv:2504.11468*, 2025.
- [9] J. Chen, J. Tang, J. Qin, X. Liang, L. Liu, E. P. Xing, and L. Lin. Geoqa: A geometric question answering benchmark towards multimodal numerical reasoning. arXiv preprint arXiv:2105.14517, 2021.
- [10] K. Chen, Z. Zhang, W. Zeng, R. Zhang, F. Zhu, and R. Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv* preprint arXiv:2306.15195, 2023.
- [11] L. Chen, J. Li, X. Dong, P. Zhang, Y. Zang, Z. Chen, H. Duan, J. Wang, Y. Qiao, D. Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024.
- [12] L. Chen, J. Li, X. Dong, P. Zhang, Y. Zang, Z. Chen, H. Duan, J. Wang, Y. Qiao, D. Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*, 2024.
- [13] L. Chen, L. Li, H. Zhao, Y. Song, and Vinci. R1-v: Reinforcing super generalization ability in vision-language models with less than \$3. https://github.com/Deep-Agent/R1-V, 2025. Accessed: 2025-02-02.
- [14] S. Chen, Y. Guo, Z. Su, Y. Li, Y. Wu, J. Chen, J. Chen, W. Wang, X. Qu, and Y. Cheng. Advancing multimodal reasoning: From optimized cold start to staged reinforcement learning, 2025. URL https://arxiv.org/abs/2506.04207.
- [15] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv* preprint arXiv:2412.05271, 2024.
- [16] Z. Chen, W. Wang, H. Tian, S. Ye, Z. Gao, E. Cui, W. Tong, K. Hu, J. Luo, Z. Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
- [17] D. Cheng, S. Huang, X. Zhu, B. Dai, W. X. Zhao, Z. Zhang, and F. Wei. Reasoning with exploration: An entropy perspective, 2025. URL https://arxiv.org/abs/2506.14758.
- [18] D. Cheng, S. Huang, X. Zhu, B. Dai, W. X. Zhao, Z. Zhang, and F. Wei. Reasoning with exploration: An entropy perspective. *arXiv preprint arXiv:2506.14758*, 2025.

- [19] Y. K. Chia, V. T. Y. Han, D. Ghosal, L. Bing, and S. Poria. Puzzlevqa: Diagnosing multimodal reasoning challenges of language models with abstract visual patterns. *arXiv* preprint *arXiv*:2403.13315, 2024.
- [20] T. Chu, Y. Zhai, J. Yang, S. Tong, S. Xie, D. Schuurmans, Q. V. Le, S. Levine, and Y. Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. arXiv preprint arXiv:2501.17161, 2025.
- [21] M. Deitke, C. Clark, S. Lee, R. Tripathi, Y. Yang, J. S. Park, M. Salehi, N. Muennighoff, K. Lo, L. Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. arXiv preprint arXiv:2409.17146, 2024.
- [22] Y. Deng, H. Bansal, F. Yin, N. Peng, W. Wang, and K.-W. Chang. Openvlthinker: An early exploration to complex vision-language reasoning via iterative self-improvement. *arXiv* preprint arXiv:2503.17352, 2025.
- [23] R. Dong, C. Han, Y. Peng, Z. Qi, Z. Ge, J. Yang, L. Zhao, J. Sun, H. Zhou, H. Wei, et al. Dreamllm: Synergistic multimodal comprehension and creation. arXiv preprint arXiv:2309.11499, 2023.
- [24] Y. Du, Z. Liu, Y. Li, W. X. Zhao, Y. Huo, B. Wang, W. Chen, Z. Liu, Z. Wang, and J.-R. Wen. Virgo: A preliminary exploration on reproducing o1-like mllm, 2025. URL https://arxiv.org/abs/2501.01904.
- [25] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- [26] M. Ford. Mental models: towards a cognitive science of language, inference, and consciousness, 1985.
- [27] C. Fu, P. Chen, Y. Shen, Y. Qin, M. Zhang, X. Lin, J. Yang, X. Zheng, K. Li, X. Sun, Y. Wu, and R. Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models, 2024. URL https://arxiv.org/abs/2306.13394.
- [28] X. Fu, Y. Hu, B. Li, Y. Feng, H. Wang, X. Lin, D. Roth, N. A. Smith, W.-C. Ma, and R. Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024.
- [29] K. Gandhi, A. Chakravarthy, A. Singh, N. Lile, and N. D. Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv* preprint *arXiv*:2503.01307, 2025.
- [30] D. Ghosal, V. T. Y. Han, C. Y. Ken, and S. Poria. Are language models puzzle prodigies? algorithmic puzzles unveil serious challenges in multimodal reasoning. *arXiv* preprint *arXiv*:2403.03864, 2024.
- [31] GPT-4o. Hello gpt-4o, 2024. URL https://openai.com/index/hello-gpt-4o/.
- [32] T. Guan, F. Liu, X. Wu, R. Xian, Z. Li, X. Liu, X. Wang, L. Chen, F. Huang, Y. Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024.
- [33] D. Guo, F. Wu, F. Zhu, F. Leng, G. Shi, H. Chen, H. Fan, J. Wang, J. Jiang, J. Wang, et al. Seed1. 5-vl technical report. *arXiv preprint arXiv:2505.07062*, 2025.
- [34] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv* preprint arXiv:2501.12948, 2025.
- [35] J. He, H. Lin, Q. Wang, Y. Fung, and H. Ji. Self-correction is more than refinement: A learning framework for visual and language reasoning tasks. *arXiv preprint arXiv:2410.04055*, 2024.

- [36] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- [37] D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- [38] D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring mathematical problem solving with the math dataset. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021*, 2021.
- [39] J. Hu, X. Wu, Z. Zhu, Xianyu, W. Wang, D. Zhang, and Y. Cao. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework. *arXiv preprint arXiv:2405.11143*, 2024.
- [40] J. Hu, Y. Zhang, Q. Han, D. Jiang, X. Zhang, and H.-Y. Shum. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. arXiv preprint arXiv:2503.24290, 2025.
- [41] W. Huang, B. Jia, Z. Zhai, S. Cao, Z. Ye, F. Zhao, Z. Xu, Y. Hu, and S. Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint arXiv:2503.06749*, 2025.
- [42] B. Hui, J. Yang, Z. Cui, J. Yang, D. Liu, L. Zhang, T. Liu, J. Zhang, B. Yu, K. Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- [43] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [44] N. Jain, K. Han, A. Gu, W.-D. Li, F. Yan, T. Zhang, S. Wang, A. Solar-Lezama, K. Sen, and I. Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- [45] P. N. Johnson-Laird. Mental models in cognitive science. *Cognitive science*, 4(1):71–115, 1980.
- [46] A. Kembhavi, M. Seo, D. Schwenk, J. Choi, A. Farhadi, and H. Hajishirzi. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern recognition*, pages 4999–5007, 2017.
- [47] A. Kumar, V. Zhuang, R. Agarwal, Y. Su, J. D. Co-Reyes, A. Singh, K. Baumli, S. Iqbal, C. Bishop, R. Roelofs, L. M. Zhang, K. McKinney, D. Shrivastava, C. Paduraru, G. Tucker, D. Precup, F. Behbahani, and A. Faust. Training language models to self-correct via reinforcement learning, 2024. URL https://arxiv.org/abs/2409.12917.
- [48] N. Lambert, J. Morrison, V. Pyatkin, S. Huang, H. Ivison, F. Brahman, L. J. V. Miranda, A. Liu, N. Dziri, S. Lyu, Y. Gu, S. Malik, V. Graf, J. D. Hwang, J. Yang, R. L. Bras, O. Tafjord, C. Wilhelm, L. Soldaini, N. A. Smith, Y. Wang, P. Dasigi, and H. Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL https://arxiv.org/abs/2411.15124.
- [49] H. Lee, S. Phatale, H. Mansoor, T. Mesnard, J. Ferret, K. Lu, C. Bishop, E. Hall, V. Carbune, A. Rastogi, and S. Prakash. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback, 2024. URL https://arxiv.org/abs/2309.00267.
- [50] B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, Y. Li, Z. Liu, and C. Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.
- [51] B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, Y. Li, Z. Liu, and C. Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

- [52] J. LI, E. Beeching, L. Tunstall, B. Lipkin, R. Soletskyi, S. C. Huang, K. Rasul, L. Yu, A. Jiang, Z. Shen, Z. Qin, B. Dong, L. Zhou, Y. Fleureau, G. Lample, and S. Polu. Numinamath. [https://huggingface.co/AI-MO/NuminaMath-CoT] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina\_dataset.pdf), 2024.
- [53] Y. Li, Y. Du, K. Zhou, J. Wang, W. X. Zhao, and J.-R. Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- [54] Y.-H. Liao, S. Elflein, L. He, L. Leal-Taixé, Y. Choi, S. Fidler, and D. Acuna. Longper-ceptualthoughts: Distilling system-2 reasoning for system-1 perception. *arXiv preprint arXiv:2504.15362*, 2025.
- [55] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [56] C. Liu, Z. Xu, Q. Wei, J. Wu, J. Zou, X. E. Wang, Y. Zhou, and S. Liu. More thinking, less seeing? assessing amplified hallucination in multimodal reasoning models. *arXiv preprint arXiv:2505.21523*, 2025.
- [57] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
- [58] Q. Liu, S. Zhang, G. Qin, T. Ossowski, Y. Gu, Y. Jin, S. Kiblawi, S. Preston, M. Wei, P. Vozila, et al. X-reasoner: Towards generalizable reasoning across modalities and domains. *arXiv* preprint arXiv:2505.03981, 2025.
- [59] Y. Liu, H. Duan, Y. Zhang, B. Li, S. Zhang, W. Zhao, Y. Yuan, J. Wang, C. He, Z. Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024.
- [60] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [61] P. Lu, R. Gong, S. Jiang, L. Qiu, S. Huang, X. Liang, and S.-C. Zhu. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. arXiv preprint arXiv:2105.04165, 2021.
- [62] P. Lu, L. Qiu, J. Chen, T. Xia, Y. Zhao, W. Zhang, Z. Yu, X. Liang, and S.-C. Zhu. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. *arXiv* preprint arXiv:2110.13214, 2021.
- [63] P. Lu, S. Mishra, T. Xia, L. Qiu, K.-W. Chang, S.-C. Zhu, O. Tafjord, P. Clark, and A. Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022.
- [64] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- [65] M. Luo, S. Tan, J. Wong, X. Shi, W. Tang, M. Roongta, C. Cai, J. Luo, T. Zhang, E. Li, R. A. Popa, and I. Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl, 2025. Notion Blog.
- [66] C. Lyu, S. Gao, Y. Gu, W. Zhang, J. Gao, K. Liu, Z. Wang, S. Li, Q. Zhao, H. Huang, W. Cao, J. Liu, H. Liu, J. Liu, S. Zhang, D. Lin, and K. Chen. Exploring the limit of outcome reward for learning mathematical reasoning, 2025. URL https://arxiv.org/abs/2502.06781.
- [67] A. Masry, D. X. Long, J. Q. Tan, S. Joty, and E. Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. arXiv preprint arXiv:2203.10244, 2022.
- [68] F. Meng, L. Du, Z. Liu, Z. Zhou, Q. Lu, D. Fu, B. Shi, W. Wang, J. He, K. Zhang, et al. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning. *arXiv* preprint arXiv:2503.07365, 2025.

- [69] Y. Min, Z. Chen, J. Jiang, J. Chen, J. Deng, Y. Hu, Y. Tang, J. Wang, X. Cheng, H. Song, W. X. Zhao, Z. Liu, Z. Wang, and J.-R. Wen. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems. *arXiv* preprint arXiv:2412.09413, 2024.
- [70] N. Muennighoff, Z. Yang, W. Shi, X. L. Li, L. Fei-Fei, H. Hajishirzi, L. Zettlemoyer, P. Liang, E. Candès, and T. Hashimoto. s1: Simple test-time scaling, 2025. URL https://arxiv.org/abs/2501.19393.
- [71] OpenAI. GPT-4o Mini: Advancing Cost-Efficient Intelligence. http://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/, 2025. Accessed: 2025-05-16.
- [72] OpenAI. Learning to reason with llms. https://openai.com/index/learning-to-reason-with-llms/, 2025.
- [73] OpenAI. Introducing openai o3 and o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/, 2025.
- [74] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744, 2022.
- [75] Y. Peng, X. Wang, Y. Wei, J. Pei, W. Qiu, A. Jian, Y. Hao, J. Pan, T. Xie, L. Ge, et al. Skywork r1v: Pioneering multimodal reasoning with chain-of-thought. *arXiv preprint arXiv:2504.05599*, 2025.
- [76] Y. Peng, G. Zhang, M. Zhang, Z. You, J. Liu, Q. Zhu, K. Yang, X. Xu, X. Geng, and X. Yang. Lmm-r1: Empowering 3b lmms with strong reasoning abilities through two-stage rule-based rl. *arXiv preprint arXiv:2503.07536*, 2025.
- [77] R. Qiao, Q. Tan, G. Dong, M. Wu, C. Sun, X. Song, Z. GongQue, S. Lei, Z. Wei, M. Zhang, et al. We-math: Does your large multimodal model achieve human-like mathematical reasoning? arXiv preprint arXiv:2407.01284, 2024.
- [78] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [79] D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark. arXiv preprint arXiv:2311.12022, 2023.
- [80] D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- [81] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [82] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [83] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv* preprint arXiv:2402.03300, 2024.
- [84] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv* preprint arXiv:2402.03300, 2024.
- [85] H. Shen, P. Liu, J. Li, C. Fang, Y. Ma, J. Liao, Q. Shen, Z. Zhang, K. Zhao, Q. Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*, 2025.

- [86] H. Shen, T. Wu, Q. Han, Y. Hsieh, J. Wang, Y. Zhang, Y. Cheng, Z. Hao, Y. Ni, X. Wang, et al. Phyx: Does your model have the" wits" for physical reasoning? *arXiv preprint arXiv:2505.15929*, 2025.
- [87] J. Skalse, N. H. R. Howe, D. Krasheninnikov, and D. Krueger. Defining and characterizing reward hacking, 2025. URL https://arxiv.org/abs/2209.13085.
- [88] R. Sutton. The bitter lesson. *Incomplete Ideas (blog)*, 13(1):38, 2019.
- [89] C. Team, Z. Yue, Z. Lin, Y. Song, W. Wang, S. Ren, S. Gu, S. Li, P. Li, L. Zhao, L. Li, K. Bao, H. Tian, H. Zhang, G. Wang, D. Zhu, Cici, C. He, B. Ye, B. Shen, Z. Zhang, Z. Jiang, Z. Zheng, Z. Song, Z. Luo, Y. Yu, Y. Wang, Y. Tian, Y. Tu, Y. Yan, Y. Huang, X. Wang, X. Xu, X. Song, X. Zhang, X. Yong, X. Zhang, X. Deng, W. Yang, W. Ma, W. Lv, W. Zhuang, W. Liu, S. Deng, S. Liu, S. Chen, S. Yu, S. Liu, S. Wang, R. Ma, Q. Wang, P. Wang, N. Chen, M. Zhu, K. Zhou, K. Zhou, K. Fang, J. Shi, J. Dong, J. Xiao, J. Xu, H. Liu, H. Xu, H. Qu, H. Zhao, H. Lv, G. Wang, D. Zhang, D. Zhang, C. Ma, C. Liu, C. Cai, and B. Xia. Mimo-vl technical report, 2025. URL https://arxiv.org/abs/2506.03569.
- [90] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023.
- [91] K. Team. Kimi k1.5: Scaling reinforcement learning with llms. 2025.
- [92] K. Team, A. Du, B. Gao, B. Xing, C. Jiang, C. Chen, C. Li, C. Xiao, C. Du, C. Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [93] K. Team, A. Du, B. Yin, B. Xing, B. Qu, B. Wang, C. Chen, C. Zhang, C. Du, C. Wei, et al. Kimi-vl technical report. *arXiv preprint arXiv:2504.07491*, 2025.
- [94] O. Team. Open Thoughts. https://open-thoughts.ai, Jan. 2025.
- [95] Q. Team. Qvq: To see the world with wisdom, 2024.
- [96] Q. Team. Qwq-32b: Embracing the power of reinforcement learning. *URL: https://qwenlm.github. io/blog/qwq-32b*, 2025.
- [97] L. S. Vygotsky. *Thought and language*, volume 29. MIT press, 2012.
- [98] F. Wang, W. Zhou, J. Y. Huang, N. Xu, S. Zhang, H. Poon, and M. Chen. mdpo: Conditional preference optimization for multimodal large language models. arXiv preprint arXiv:2406.11839, 2024.
- [99] H. Wang, C. Qu, Z. Huang, W. Chu, F. Lin, and W. Chen. VI-rethinker: Incentivizing self-reflection of vision-language models with reinforcement learning. *arXiv* preprint *arXiv*:2504.08837, 2025.
- [100] K. Wang, J. Pan, W. Shi, Z. Lu, H. Ren, A. Zhou, M. Zhan, and H. Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024.
- [101] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024.
- [102] X. Wang, Z. Yang, C. Feng, H. Lu, L. Li, C.-C. Lin, K. Lin, F. Huang, and L. Wang. Sota with less: Mcts-guided sample selection for data-efficient visual reasoning self-improvement. *arXiv* preprint arXiv:2504.07934, 2025.
- [103] Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In *Advances in Neural Information Processing Systems, NeurIPS* 2024, 2024.

- [104] Z. Wang, M. Xia, L. He, H. Chen, Y. Liu, R. Zhu, K. Liang, X. Wu, H. Liu, S. Malladi, A. Chevalier, S. Arora, and D. Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms, 2024. URL https://arxiv.org/abs/2406.18521.
- [105] Z. Wang, M. Xia, L. He, H. Chen, Y. Liu, R. Zhu, K. Liang, X. Wu, H. Liu, S. Malladi, et al. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *Advances in Neural Information Processing Systems*, 37:113569–113697, 2024.
- [106] H. Wei, L. Kong, J. Chen, L. Zhao, Z. Ge, J. Yang, J. Sun, C. Han, and X. Zhang. Vary: Scaling up the vision vocabulary for large vision-language model. In *European Conference on Computer Vision*, pages 408–424. Springer, 2024.
- [107] Y. Wei, L. Zhao, K. Lin, E. Yu, Y. Peng, R. Dong, J. Sun, H. Wei, Z. Ge, X. Zhang, et al. Perception in reflection. *arXiv preprint arXiv:2504.07165*, 2025.
- [108] X.AI. Grok-2 beta release. https://x.ai/blog/grok-2, 2024. Accessed on: 2024-07-02.
- [109] K. Xiang, Z. Liu, Z. Jiang, Y. Nie, R. Huang, H. Fan, H. Li, W. Huang, Y. Zeng, J. Han, et al. Atomthink: A slow thinking framework for multimodal mathematical reasoning. *arXiv* preprint arXiv:2411.11930, 2024.
- [110] Y. Xiao, E. Sun, T. Liu, and W. Wang. Logicvista: Multimodal llm logical reasoning benchmark in visual contexts, 2024. URL https://arxiv.org/abs/2407.04973.
- [111] G. Xu, P. Jin, L. Hao, Y. Song, L. Sun, and L. Yuan. Llava-o1: Let vision language models reason step-by-step. *arXiv preprint arXiv:2411.10440*, 2024.
- [112] L. Xu, Y. Zhao, J. Wang, Y. Wang, B. Pi, C. Wang, M. Zhang, J. Gu, X. Li, X. Zhu, et al. Geosense: Evaluating identification and application of geometric principles in multimodal reasoning. *arXiv* preprint arXiv:2504.12597, 2025.
- [113] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [114] A. Yang, B. Zhang, B. Hui, B. Gao, B. Yu, C. Li, D. Liu, J. Tu, J. Zhou, J. Lin, K. Lu, M. Xue, R. Lin, T. Liu, X. Ren, and Z. Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- [115] J. Yang, F. Ma, Z. Wang, D. Yin, K. Rong, F. Rao, and R. Zhang. Wethink: Toward general-purpose vision-language reasoning via reinforcement learning, 2025. URL https://arxiv.org/abs/2506.07905.
- [116] Y. Yang, X. He, H. Pan, X. Jiang, Y. Deng, X. Yang, H. Lu, D. Yin, F. Rao, M. Zhu, et al. R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization. *arXiv preprint arXiv:2503.10615*, 2025.
- [117] Z. Yang, X. Yu, D. Chen, M. Shen, and C. Gan. Machine mental imagery: Empower multimodal reasoning with latent visual tokens. *arXiv preprint arXiv:2506.17218*, 2025.
- [118] Z. Yao, Y. Liu, Y. Chen, J. Chen, J. Fang, L. Hou, J. Li, and T.-S. Chua. Are reasoning models more prone to hallucination? *arXiv preprint arXiv:2505.23646*, 2025.
- [119] Q. Ye, H. Xu, G. Xu, J. Ye, M. Yan, Y. Zhou, J. Wang, A. Hu, P. Shi, Y. Shi, et al. mplugowl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- [120] E. Yeo, Y. Tong, M. Niu, G. Neubig, and X. Yue. Demystifying long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025.
- [121] B. Yin, Q. Wang, P. Zhang, J. Zhang, K. Wang, Z. Wang, J. Zhang, K. Chandrasegaran, H. Liu, R. Krishna, S. Xie, M. Li, J. Wu, and L. Fei-Fei. Spatial mental modeling from limited views, 2025. URL https://arxiv.org/abs/2506.21458.

- [122] B. Yin, Q. Wang, P. Zhang, J. Zhang, K. Wang, Z. Wang, J. Zhang, K. Chandrasegaran, H. Liu, R. Krishna, et al. Spatial mental modeling from limited views. arXiv preprint arXiv:2506.21458, 2025.
- [123] E. Yu, L. Zhao, Y. Wei, J. Yang, D. Wu, L. Kong, H. Wei, T. Wang, Z. Ge, X. Zhang, et al. Merlin: Empowering multimodal llms with foresight minds. In *European Conference on Computer Vision*, pages 425–443. Springer, 2024.
- [124] E. Yu, K. Lin, L. Zhao, Y. Wei, Z. Zhu, H. Wei, J. Sun, Z. Ge, X. Zhang, J. Wang, et al. Unhackable temporal rewarding for scalable video mllms. arXiv preprint arXiv:2502.12081, 2025.
- [125] E. Yu, K. Lin, L. Zhao, J. Yin, Y. Wei, Y. Peng, H. Wei, J. Sun, C. Han, Z. Ge, et al. Perception-r1: Pioneering perception policy with reinforcement learning. *arXiv* preprint *arXiv*:2504.07954, 2025.
- [126] T. Yu, Y. Yao, H. Zhang, T. He, Y. Han, G. Cui, J. Hu, Z. Liu, H.-T. Zheng, M. Sun, et al. Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13807–13816, 2024.
- [127] W. Yu, Z. Yang, L. Li, J. Wang, K. Lin, Z. Liu, X. Wang, and L. Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.
- [128] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.
- [129] X. Yue, T. Zheng, Y. Ni, Y. Wang, K. Zhang, S. Tong, Y. Sun, B. Yu, G. Zhang, H. Sun, Y. Su, W. Chen, and G. Neubig. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark, 2024. URL https://arxiv.org/abs/2409.02813.
- [130] J. Zhang, R. Dong, H. Wang, X. Ning, H. Geng, P. Li, X. He, Y. Bai, J. Malik, S. Gupta, and H. Zhang. Alphaone: Reasoning models thinking slow and fast at test time. *arXiv* preprint *arXiv*:2505.24863, 2025.
- [131] J. Zhang, J. Huang, H. Yao, S. Liu, X. Zhang, S. Lu, and D. Tao. R1-vl: Learning to reason with multimodal large language models via step-wise group relative policy optimization. *arXiv* preprint arXiv:2503.12937, 2025.
- [132] P. Zhang, X. Dong, Y. Zang, Y. Cao, R. Qian, L. Chen, Q. Guo, H. Duan, B. Wang, L. 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.
- [133] R. Zhang, D. Jiang, Y. Zhang, H. Lin, Z. Guo, P. Qiu, A. Zhou, P. Lu, K.-W. Chang, P. Gao, and H. Li. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems?, 2024. URL https://arxiv.org/abs/2403.14624.
- [134] L. Zhao, E. Yu, Z. Ge, J. Yang, H. Wei, H. Zhou, J. Sun, Y. Peng, R. Dong, C. Han, et al. Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning. *arXiv* preprint arXiv:2307.09474, 2023.
- [135] R. Zhao, A. Meterez, S. Kakade, C. Pehlevan, S. Jelassi, and E. Malach. Echo chamber: Rl post-training amplifies behaviors learned in pretraining. arXiv preprint arXiv:2504.07912, 2025.
- [136] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [137] J. Zhu, W. Wang, Z. Chen, Z. Liu, S. Ye, L. Gu, H. Tian, Y. Duan, W. Su, J. Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.

- [138] K. Zhu, L. Zhao, Z. Ge, and X. Zhang. Self-supervised visual preference alignment. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 291–300, 2024.
- [139] Z. Zhu, L. Zhao, K. Lin, J. Yang, E. Yu, C. Liu, H. Wei, J. Sun, Z. Ge, and X. Zhang. Perpo: Perceptual preference optimization via discriminative rewarding. *arXiv preprint arXiv:2502.04371*, 2025.
- [140] C. Zou, X. Guo, R. Yang, J. Zhang, B. Hu, and H. Zhang. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models, 2024.

# **NeurIPS Paper Checklist**

## 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction accurately reflect the paper's contributions and scope. We assert that this paper pioneers the exploration of cognitive behavior transfer from language to vision. Furthermore, we establish several key insights about the transfer. Driven by these findings, we developed our Open-Vision-Reasoner, which achieves new state-of-the-art (SoTA) results across multiple visual and language tasks.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The discussion of limitation is presented in Section 7.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: the paper does not include theoretical.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

## 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: All reproduce information including model architecture, algorithms and implementation details are provided in this paper.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

We plan to publicly release all the source code and dataset upon paper acceptance.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
  proposed method and baselines. If only a subset of experiments are reproducible, they
  should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We present the experimental setting in section 4.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We do not report error bars in this work.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We report that all experiments are conducted on NVIDIA A100 Tensor Core GPU.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: This work aligns with NeurIPS Code of Ethics.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This work does not have potential negative societal impacts.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: There is no high risk for misuse of our models and datasets.

# Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All code and data are following licenses properly. And all related works are cited.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We introduce new model and data in this paper. While we do not publicly release them at the submission stage to preserve anonymity, we will further release the full code, models, and accompanying documentation publicly.

## Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing experiments and research with human subjects.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing experiments and research with human subjects.

#### Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: We use gpt-4o-mini in our evaluation.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

The appendix includes extended details on data curation (Appendix A), implementation (Appendix B), cognitive behavior evaluation (Appendix C), and additional case studies (Appendix D).

## A Cold-Start Data Curation Details

As mentioned in Section 3.3, a critical component of our initial policy development is the curation of the cold-start SFT dataset. This stage serves as the foundation for subsequent learning, particularly in shaping the model's ability to exhibit structured reasoning and cognitive behavior. To this end, we adopt a multi-stage curation pipeline consisting of data collection, filtering, cleaning, and strategic reweighting.

**Data Acquisition.** We begin by assembling a broad set of prompt-response pairs that span diverse reasoning domains. These include math, science, and general logical reasoning tasks such as puzzles, deductive tasks, and constraint satisfaction problems. Our sources include a mix of public datasets as illustrated in Section 3.3.

**Automated Filtering.** To improve the signal-to-noise ratio, we filter the collected data using a lightweight pretrained LLM as a proxy for quality estimation. Each instance is passed through this model, and those with abnormally high training loss are flagged as noisy, ambiguous, or misaligned. We further apply rule-based and model-assisted pattern detectors to identify and eliminate undesirable data characteristics.

**Difficulty Stratification.** We explicitly incorporate samples from AMC, AIME, Olympiads, and AoPS forums to ensure difficulty levels. We then stratify the collected samples based on their source and inherent problem complexity to balance coverage across easy, intermediate, and challenging reasoning scenarios.

**Reweighting and Balance.** To address imbalances across domains and formats, we apply a reweighting scheme based on coverage and reasoning relevance. Over-represented formats are down-weighted, while rare but cognitively rich categories are given higher sampling probabilities. This ensures a more uniform distribution of reasoning challenges and minimizes overfitting to dominant patterns.

## **B** More Implementation Details

**Model and Optimization Setup** Our model is based on the Qwen2.5-VL [5]. During RL, both the policy and critic networks are initialized from the cold-start model. The value head is initialized from a uniform distribution  $\mathcal{U}(-\sqrt{5},\sqrt{5})$  without bias. The policy and critic networks do not share weights during training. We use the AdamW optimizer with  $\beta = [0.9, 0.95]$  and no weight decay. Learning rates are set to  $1\times 10^{-6}$  for the policy and  $5\times 10^{-6}$  for the critic. We use constant learning rates with a linear warm-up of 50 steps, and employ sample packing for improved throughput. No KL regularization or entropy bonus is used, demonstrating that vanilla PPO remains stable under our setup.

**PPO Training Dynamics** Each PPO update is based on 512 unique prompts, each generating 16 sampled responses using temperature and top-p sampling (both set to 1.0). To ensure training stability, we enforce strict on-policy updates for the policy: each prompt generation corresponds to a single optimization step. In contrast, the critic performs 4 optimization steps per PPO update. We apply batch-level advantage normalization to stabilize training further.

# C Details for Coginitive Behavior Evaluation

In this section, we detail definitions of metrics in cognitive behavior analysis (Section 5.2) and show the prompt for evaluation.

#### **C.1** Behavior Transfer Rate

To quantify how well language-acquired cognitive behaviors generalize to the visual modality, we define the Behavior Transfer Rate (BTR) for each behavior type introduced in Fig. 5. BTR is calculated as the ratio between the emergence rate of visual behaviors and that of their linguistic counterparts. Formally, we compute the Cognitive Behavior Emergence Rate in the visual modality (CBR $_{\nu}$ ) and in the language modality (CBR $_{l}$ ), both evaluated on the multimodal benchmark *MathVision (mini)*. The **BTR** is then defined as:

$$BTR = \frac{CBR_v}{CBR_l}$$

This metric reflects the cross-modal transfer efficiency of cognitive behaviors, with higher values indicating stronger behavioral generalization from language to vision.

#### **C.2** Evaluation Prompt

We design prompts for the LLM-based evaluation. Fig. 7 showcases the prompt template for the cognitive behavior *Backtracking* as an example.

#### Prompt example:

Problem Statement: "<QUESTION>"
Model's Reasoning: "<REASONING CONTENT>"

Evaluate the textual reasoning for Backtracking/Revision. This is when the model, in its text, realizes a chosen path, calculation, or assumption is incorrect or unproductive, and explicitly states it's trying a different approach, redoing calculations with different numbers/operations, or significantly altering its line of textual reasoning.

Example: "That calculation led to a dead end. Let's try subtracting Y from X instead." or "My previous assumption about the sequence was wrong. Starting over with..."

Count distinct instances of such explicit textual backtracking or major revision. Provide the count in <count>X</count> tags. If none, use <count>o</count>.

Figure 7: Prompt design for evaluating *Backtracking*.

# **D** More Case Studies

We present examples of OVR's reasoning process during complex multimodal tasks. In Fig. 10, the model begins by perceiving the image and applying *visual divide-and-conquer* to interpret each emoji individually. It then engages in *subgoal setting*, breaking down the task to explore multiple candidate answers based on each emoji. Throughout the process, the model exhibits *backtracking*, revising earlier hypotheses by *cross-verifying* them with elements in the image. The reasoning process concludes with a final, well-justified prediction.



Ouestion: Based on the paper, which one takes longer to train, RL phase or reward model? Ground Truth: RL phase takes longer.

OVR: </think> Okay, so I need to figure out whether the RL phase or the reward model takes longer to train based on the

information provided in the paper. Let me go through the relevant parts step by step.
First, the paper mentions the Seed1.5-VL model. It talks about a post-training framework using hybrid reinforcement learning with RLHF (Reinforcement Learning with Human Feedback) and RLVF (Reward Learning with Verifier Feedback). The RL phase is mentioned as requiring 24k GPU hours, while the reward model inference uses 60k GPU hours. Wait, but that's for the reward model during inference, right? Because it says "reward model inference use replicas each configured with tensor parallelism." So during training, maybe the reward model is trained for longer?

Wait, the sentence says: "The RL phase of Seed1.5-VL costs 60k GPU hours. The reward model is trained using the same framework as Seed1.5-VL pretraining phase, requiring 24k GPU hours." Hmm, that seems conflicting. Wait, no: "pretraining phase" is separate. Let me parse this sentence again.

Wait, the sentence says: "The RL phase of Seed1.5-VL costs 60k GPU hours. The reward model is trained using the same

war, the sentence says. The RL phase of Secti.5-YL costs own 670 hours. The reward models trained using the same framework as Secti.5-YL pretraining phase, requiring 24k GPU hours." So the RL phase takes 60k hours, and the reward model training (which is part of the post-training framework) uses 24k hours? That would mean the RL phase is longer than the reward model training. But maybe I'm misinterpreting.

Alternatively, perhaps there's a typo or confusion in the description. Let me look again at the exact wording:

"We use the same training framework as Seed1.5-VL pretraining phase for optimization techniques as in the pretraining phase for efficient actor and critic updates, and vLLM [81] for efficient progressive generation of rollouts. Specifically, the RL and RLVF training employs 3-D parallelism [61, 108, 126, 141]; rollout generation and reward/reference model inference use replicas each configured with tensor parallelism [7]. The RL phase of Seed1.5-VL costs 6ok GPU hours. The reward model is trained using the same framework as Seed1.5-VL pretraining phase, requiring 24k GPU hours."

Wait, so here it says that during RL training (which includes both RL and RLVF?), it takes 60k hours. Then separately, the reward model is trained using the same framework as pretraining (which might be different) but takes 24k hours. But if pretraining took longer than post-training phases? Or maybe there's a mix-up here...

But wait, let me check again: Original text from user input:

"We use ... optimization techniques as in the pretraining phase for efficient actor and critic updates ... The RL phase of Seed1.5-VL costs 60k GPU hours ... The reward model is trained using ... same framework as Seed1.5-VL pretraining phase ... requiring 24k GPU hours.'

So during post-training phases

- · RL Phase: 60k
- Reward Model Training: 24k

Thus, RL Phase is longer than Reward Model Training.
Therefore, answer would be that the RL Phase takes longer to train than the Reward Model.

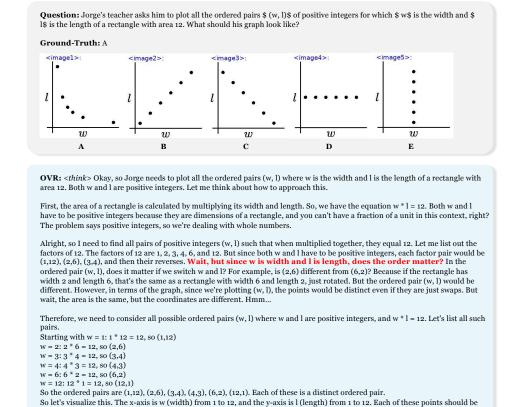
Based on the provided information from the paper:

- RL Phase: Requires \*\*60k GPU hours\*
- Reward Model Training: Requires 24k GPU hours during post-training phases.

Conclusion: The RL Phase takes longer to train compared to Reward Model Training(60k vs. 24k GPU hours).

Answer: RL Phase

Figure 8: Case Study on OCR-based Document Understanding. Red highlights some of the cognitive behaviors.



Looking at the options given:

(A) image1: shows points that seem to go from (1,12) down to (12,1), but maybe in a different arrangement.

Wait, but let me check the images. (image1): a scatter plot with points arranged from left to right along the w-axis increasing, and I decreasing. So like (1,12), (2,6), (3,4), (4,3), (6,2), (12,1). But maybe not all points are there.

(image2): points increasing in both directions? Not sure. (image3): Points decreasing in l as w increases, but maybe with gaps?

(image4): Points all along the w-axis, so l is constant? That can't be right because l would have to be 12 for w=1, 6 for w=2, etc.

(image5): Points all along the l-axis, so w is constant? Similarly, that's not possible.

Wait, but the correct plot should have points at (1,12), (2,6), (3,4), (4,3), (6,2), (12,1). So when plotted on a graph with w on the x-axis and l on the y-axis, these points would form a diagonal line from (1,12) to (12,1), but only at those specific integer points... Let's check the options again.

Assuming image1 has all six points in decreasing order, image2 might have some other arrangement, image3 maybe have some other arrangement, image4 all on w-axis, image5 all on l-axis.

But since the problem says \"positive integers\", and there are exactly six pairs, the correct graph should have six points. So image4 and image5 can be eliminated because they have multiple points on a single axis, which is not possible since for each w, l is determined (except when w=12, l=1 and vice versa). So image4 is points along w-axis (l=1?), but that's not correct. Similarly, image5 would have points along l-axis (w=1?), which is also not correct.

So between image1, image2, image3.

If image1 has points (1,12), (2,6), (3,4), (4,3), (6,2), (12,1), then image1 is correct. If image2 has points in a different order, maybe increasing, but since w is on the x-axis and l on y-axis, increasing w would correspond to decreasing l, so the points should form a decreasing trend...

Therefore, the correct answer is image1, which is option A. </think>

The problem requires plotting all ordered pairs (w, l) where w and l are positive integers such that the area (w \* l = 12). The factor pairs of 12 are (1, 12), (2, 6), (3, 4), (4, 3), (6, 2), and ((12, 1). These pairs form six distinct points on a graph with w on the x-axis and l on the y-axis. The correct graph should show all six points arranged in a decreasing trend from (1, 12) to (12, 1). Among the options, \*\*image1\*\* (Option A) correctly displays all six points in this pattern.

Answer: A

Figure 9: Case Study on Math Reasoning. Red highlights some of the cognitive behaviors.

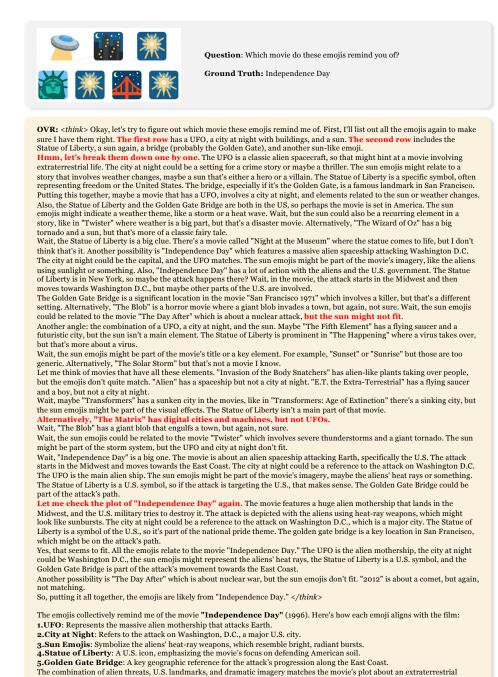


Figure 10: **Case Study on Emoji Quiz.** Red highlights some of the cognitive behaviors. The model first performs *visual decomposition* to interpret each emoji individually. It then engages in *subgoal setting* to explore multiple candidate answers, accompanied by *backtracking* to revise earlier hypotheses. The reasoning concludes with a final, justified prediction.

Answer: Independence Day

**Question:** The number of points equidistant from a circle and two parallel tangents to the circle is: (A) o. (B) 2. (C) 3. (D) 4. (E) infinitely many. (20th American High School Mathematics Examination, 1969)

DeepSeek-R1: <think> \nThe problem is from the 20th American High School Mathematics Examination in 1969. It states: "The number of points equidistant from a circle and two parallel tangents to the circle is\" with options ranging from o to infinitely many.\n\nFirst, I need to understand what is being asked. We have a circle and two parallel tangents to this circle. I need to find points that are equidistant from the circle and from these two tangents. The phrasing is a bit ambiguous. Let me read it again.\n\n\"It says: equidistant from a circle and two parallel tangents.\" I think it means equidistant from the circle itself and from the pair of tangents. But distance to a circle usually means the shortest distance to the circumference, and distance to a line is the perpendicular distance.\n\nBut it says \"from a circle and two parallel tangents.\" It might mean equidistant from the circle and from each of the tangents, but that doesn't make sense because the tangents are lines, and the circle is a curve.  $\n\$ means that the point is equidistant from the circle and from the two tangents considered together, but that seems messy,\n\nAnother interpretation: maybe it means that the distance to the circle is equal to the distance to each tangent, but that would probably be overconstrained.\n\nI recall that in some contexts, \"equidistant from a curve and a line\" means that the shortest distance to the curve equals the shortest distance to the line.\n\nBut here it's a circle and two lines.\n\nPerhaps the problem intends for the distance to the circle to be equal to the distance to the pair of tangents, but since the tangents are parallel, the distance to the pair might be the distance to the closer one or something.\n\nThat doesn't sound right.\n\nLet me visualize this. Let me sketch a circle. Assume it's a circle with center O. Draw two parallel tangents to this circle. Since they are parallel, they must be on opposite sides of the circle. For example, if the circle is centered at the origin, and the tangents are horizontal, say y = r and y = -r, but if the circle has radius r, and center at (0,0), then the tangents would be y = r and y = -r? No.\n\nIf center is at (0,0), radius r, then the tangent lines are  $x = \pm r$  for vertical tangents, but for horizontal, if I want parallel tangents, say both horizontal. So, if I have a circle  $x^2 + y^2 = r^2$ , the horizontal tangents are y = r and y = -r. Is that correct? When y = r,  $x^2 + r^2 = r^2$  implies x = 0, so yes, the point (0,r) is on the circle, and the tangent is horizontal. Similarly for y = -r. So yes, y = r and y = -r are the two parallel tangents.  $\n\$  But in this case, they are symmetric.  $\n\$  Now, I need points P such that the distance from P to the circle is equal to the distance from P to each of the tangents? But that might not be what it means.  $\n$  The problem says \"equidistant from a circle and two parallel tangents

Figure 11: Case Study on DeepSeek-R1. Red highlights the *mental imagery* mentioned in Section 5.2.