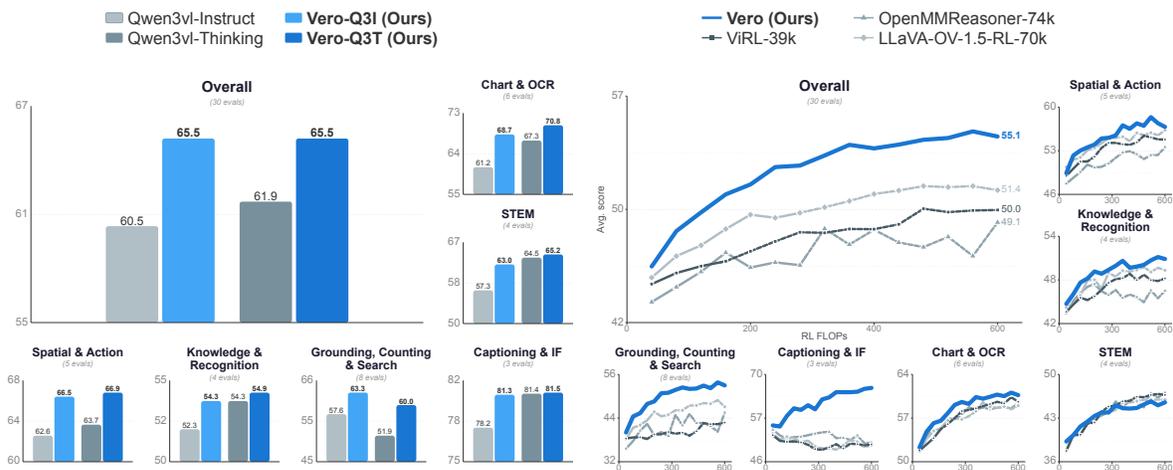


# Vero: An Open RL Recipe for Visual Reasoning

Gabriel Sarch<sup>\*,‡</sup> Linrong Cai<sup>\*,‡</sup> Qunzhong Wang<sup>‡</sup> Haoyang Wu  
Danqi Chen Zhuang Liu<sup>†</sup>

Princeton University



**Figure 1** Note: numbers in this teaser are not finalized. **Vero** improves performance across the overall benchmark and six domains. **Left**. Model comparison across overall and domain scores, with evaluation counts in parentheses. We compare Qwen3VL-Instruct, Qwen3VL-Thinking, **Vero-Q3I**, and **Vero-Q3T**. **Right**. Scaling trends versus RL FLOPs for overall and each domain. **Vero** improves performance over baselines across all six task categories. All runs are finetuned from Qwen-2.5-VL-7B-Instruct. Solid lines denote training within the first epoch of each dataset; dashed lines indicate continued training beyond one epoch.

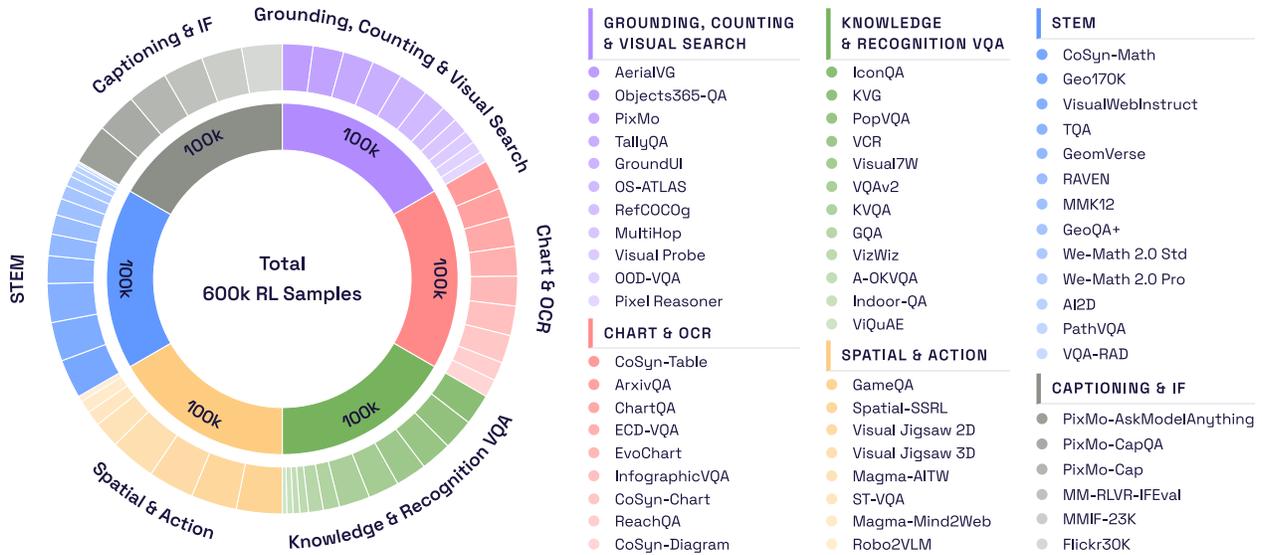
## Abstract

The strongest vision-language models (VLMs) rely on proprietary reinforcement learning (RL) pipelines with non-public data and undisclosed designs, making it difficult to study what drives their performance. We show that a fully open, single-stage RL recipe, when paired with sufficiently diverse training data, can match or exceed proprietary pipelines. We introduce **Vero**, a family of fully-open VLMs trained with a carefully curated collection of 600K RL samples from 59 datasets spanning six core task categories. **Vero** achieves state-of-the-art performance across a wide range of visual reasoning tasks. Starting from Qwen3-VL-8B-Instruct, **Vero** outperforms Qwen3-VL-8B-Thinking on 20 of 30 benchmarks without using additional proprietary thinking data. Applied on top of Qwen3-VL-8B-Thinking, **Vero** improves this further to 23 of 30. Starting from MiMo-VL-SFT, **Vero** surpasses MiMo-VL-RL, which uses a proprietary RL recipe with non-public data. Systematic ablations reveal that different task categories elicit qualitatively distinct reasoning patterns that transfer poorly in isolation, suggesting that broad data coverage is the primary driver of strong RL scaling. All data, code, and models are released.

\* Project Leads

‡ Core Contributors

† Corresponding Author



**Figure 2** Composition of the RL training data. The inner ring shows six dataset categories, and the outer ring shows their constituent datasets.

## 1 Introduction

Training models to reason through explicit chain-of-thought (CoT) has become a powerful paradigm for improving the capabilities of large language models (LLMs) and vision-language models (VLMs). A key driver of this progress is on-policy reinforcement learning (RL) (DeepSeek-AI et al., 2025). Methods such as PPO (Schulman et al., 2017) and GRPO (Shao et al., 2024) enable models to learn from their own generations, iteratively refining their reasoning chains through reward signals. Recent models such as DeepSeek-R1 and Kimi K2.5 demonstrate that on-policy RL can drive substantial improvements in both text-only and multimodal reasoning (DeepSeek-AI et al., 2025; Kimi Team et al., 2026).

Yet the current strongest visual reasoning models are products of proprietary RL pipelines with non-public data and undisclosed reward designs. Models such as Qwen-VL (Bai et al., 2025b,a) release weights and are widely adopted, but do not release RL training code or datasets. Accompanying technical reports often omit detailed ablations of design choices, making it difficult to systematically study what drives performance. Meanwhile, fully open efforts such as OpenMMReasoner (Zhang et al., 2025b) and VL-Rethinker (Wang et al., 2025) focus primarily on visual math, covering only a narrow subset of visual tasks. However, as we show in Sections 5 and 7, training on a single task category does not generalize to other visual capabilities, in both task performance and chain-of-thought behavior. This leaves a central question unanswered: what does it take to train a broadly capable visual reasoner?

We show that a single-stage RL recipe with diverse and high quality data suffices. We introduce **Vero**, a family of fully open VLMs trained with RL on top of existing models to perform strongly across diverse visual tasks. Our recipe centers on careful data curation: we collect 600K high-quality samples from 59 datasets spanning six core task categories (Chart & OCR, STEM, Spatial & Action, Knowledge & Recognition, Grounding, Counting & Visual Search, and Captioning & Instruction Following), and pair them with task-routed reward functions. No additional warm start, no staged RL, and no proprietary data. Alongside the training data, we assemble **VeroEvalSuite**, a comprehensive evaluation suite of 30 benchmarks spanning all six categories. Figure 2 summarizes the dataset composition.

Through systematic ablations of data curation, mixture strategies, and reward design, we find that data diversity is the critical ingredient. Different task categories elicit qualitatively distinct reasoning patterns that transfer poorly in isolation: for example, STEM tasks trigger elevated backtracking while grounding tasks suppress introspective behaviors in favor of directed visual search. Broad task coverage is therefore essential for producing a generally capable model. We additionally find that (1) uniform mixture weighting across task

categories outperforms schemes based on accuracy, reasoning length, or dataset size; (2) multi-task training necessitates an expressive, task-routed reward design; and (3) including open-ended tasks is necessary to preserve visual chat ability during RL.

**Vero** achieves state-of-the-art overall performance among 8B VLMs (Figure 1). Training on four different base models yields consistent improvements of +3.5 to +5.0 points averaged over 30 benchmarks. **Vero-Q3T-8B** outperforms Qwen3-VL-8B-Thinking on 23 of 30 benchmarks. **Vero-Mi-7B**, trained on MiMo-VL-7B-SFT (Yue et al., 2025b) with our fully open recipe, surpasses MiMo-VL-7B-RL, which starts from the same initial model but uses a proprietary RL recipe with non-public data. These results demonstrate that a simple, fully open recipe can match or exceed proprietary pipelines. We release all data, code, and model weights to facilitate future research.

## 2 Related Works

**Vision-Language Models.** Vision-language models excel on multimodal tasks, including proprietary systems such as GPT-5 (Singh et al., 2025) and Gemini (Team et al., 2023, 2024; Comanici et al., 2025), open-weight families such as Qwen (Bai et al., 2025b,a), GLM (Hong et al., 2026), and Kimi (Kimi Team, 2025), and fully open releases of data, code, and weights such as Molmo (Deitke et al., 2024; Clark et al., 2026) and LLaVA (Liu et al., 2023; An et al., 2025). These models are expected to handle a wide range of tasks. While little is publicly known about proprietary model post-training, recent open-weight models have explored techniques such as RL with curriculum sampling (Hong et al., 2026) and mixed on-policy RL (Yue et al., 2025b), yet the factors that drive their performance across diverse tasks remain unclear. Our work targets this gap by providing a fully open multi-domain RL recipe for general visual understanding.

**Reasoning and Thinking for VLMs.** Chain-of-thought reasoning enables models to leverage additional test-time compute through step-by-step problem decomposition (Wei et al., 2022; Zhang et al., 2023). The two dominant approaches for training reasoning models are distillation, where a strong teacher generates reasoning traces for supervised fine-tuning (Xu et al., 2025; Yao et al., 2024; Sarch et al.), and reinforcement learning, which optimizes against outcome-based rewards without requiring a fixed teacher (DeepSeek-AI et al., 2025). Recent works apply RL to visual reasoning (Yu et al., 2025a; Zhang et al., 2025b; Wang et al., 2025), but primarily in narrow domains, leaving the effect of RL-trained reasoning on broad visual understanding underexplored. We show that RL with careful reward and data design consistently outperforms narrowly trained baselines across diverse visual task categories.

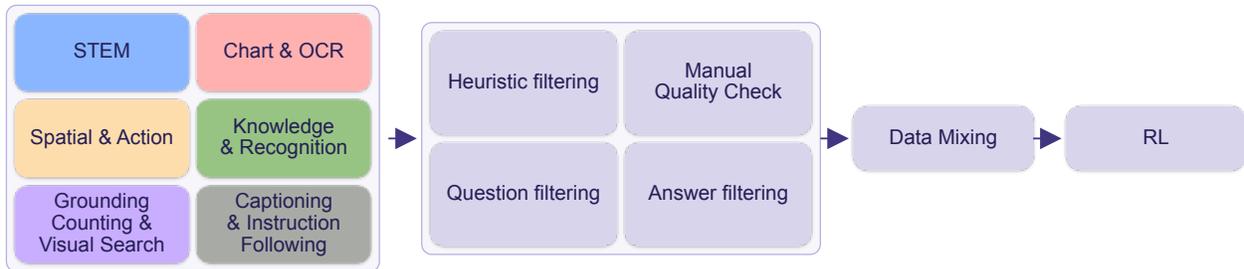
**RL Recipes and Data Curation for VLMs.** Several works provide recipes for RL-based visual reasoning training. OpenMMReasoner (Zhang et al., 2025b) combines teacher distillation and GSPO (Zheng et al., 2025) over multimodal reasoning benchmarks, VL-Rethinker (Wang et al., 2025) addresses training instability via selective sample replay and forced rethinking, and Perception-R1 (Yu et al., 2025a) designs discriminative rewards for perceptual tasks such as grounding and counting. These efforts target visual math or narrow perceptual domains and provide limited ablations of data curation and reward design. Our recipe spans six task categories with 600K datapoints from 59 datasets, includes a 10-way routed reward system, and provides systematic ablations of data and design choices, all released publicly to support VLM research.

## 3 Vero

### 3.1 Task Definitions

We consider the problem of training a vision-language model (VLM)  $\pi_\theta$  via reinforcement learning to maximize expected reward across a diverse set of visual reasoning tasks. Given a visual input  $v$  (an image or set of images) and a text query  $q$ , the model generates a response  $y \sim \pi_\theta(\cdot | v, q)$ . A reward function  $R(y, y^*)$  evaluates the response against a ground-truth answer  $y^*$ . The RL training objective is:

$$\max_{\theta} \mathbb{E}_{(v,q,y^*) \sim \mathcal{D}} \mathbb{E}_{y \sim \pi_\theta(\cdot | v,q)} [R(y, y^*)], \quad (1)$$



**Figure 3 RL training data curation pipeline.** We curate a large-scale, high-quality RL training dataset by sourcing from over 250 candidate datasets, applying dataset-level and question-level filtering, and tuning per-batch dataset mixtures for training.

where  $\mathcal{D}$  is the training data distribution. A central challenge is constructing  $\mathcal{D}$  to span a broad range of visual reasoning capabilities, so that the resulting policy generalizes across diverse tasks.

**Task taxonomy.** Figure 2 provides an overview of our training data composition. We organize our training data into six task categories, each targeting a distinct visual reasoning capability. This taxonomy is motivated by two observations. First, we find empirically (Section 5 and Section 7) that training on any single category fails to transfer reliably to others and elicits distinct chain-of-thought behaviors, suggesting that these categories exercise different reasoning strategies and skills. Second, leading VLM evaluation frameworks treat these as separable capability axes for real-world use cases: for example, Qwen2.5-VL (Bai et al., 2025b) evaluates document understanding, mathematical reasoning, general VQA, grounding, and visual agents as distinct dimensions, while Kimi K2.5 (Kimi Team et al., 2026) similarly separates reasoning, knowledge/VQA, perception, and document/OCR.

Following these frameworks, we define a visual taxonomy with six categories: **STEM** (13 datasets) covers mathematical diagram reasoning, scientific figure interpretation, and medical image understanding, with answers that are typically numeric or symbolic. **Spatial & Action** (8 datasets) targets embodied reasoning, UI navigation, and 3D spatial understanding, requiring reasoning about spatial transformations and action sequences. **Knowledge & Recognition** (12 datasets) spans visual question answering that combines object, scene, and entity recognition with external or commonsense knowledge. **Chart & OCR** (9 datasets) focuses on extracting and reasoning over structured information in documents, charts, tables, and infographics. **Grounding, Counting & Visual Search** (11 datasets) requires spatially localizing objects via bounding boxes, counting entity instances, and searching among visual distractors. **Captioning & Instruction Following** (6 datasets) encompasses open-ended image description and following prompt instructions.

### 3.2 Dataset Curation, Filtering, and Mixtures

We curate a multi-task RL training set of 600K samples from 59 datasets spanning six task categories (Section 3.1). Figure 3 summarizes the pipeline.

**Step 1. Sourcing Training Data.** We start from over 250 candidate datasets drawn from instruction-tuning collections (e.g., FineVision (Wiedmann et al., 2025)) and recently released task-specific sources (e.g., Visual Jigsaw (Wu et al., 2025)), then apply dataset-level and sample-level filtering. Each dataset is assigned to the task category that best reflects its primary skill, based on manual inspection and its utility in prior work.

**Heuristic filtering.** We discard datasets with fewer than 1K examples, average image resolution below 200K pixels (retaining five low-resolution datasets for question quality), or exclusively binary yes/no and true/false questions to mitigate correct guessing.

**Manual filtering.** For each candidate, we inspect  $\sim 50$  examples against three criteria: **correctness** ( $< 5\%$  annotation error rate in image-question-answer triples), **unambiguity** (each question admits a single verifiable answer), and **verifiability** (the answer format is compatible with our reward functions). Of  $\sim 100$  datasets passing heuristic screening, 59 were retained. For a small number of datasets we additionally rewrite

	Chart & OCR	STEM	Spatial & Action	Knowl. & Recog.	Grnd., Cnt. & Search
Unfiltered	60.0	45.4	56.3	62.5	54.7
Q. Filtering	60.1	43.6	58.2	63.0	54.1
A. Canonic.	59.9	45.5	-	64.6	-

**Table 2** Effect of question filtering and answer canonicalization on Qwen2.5-VL-7B-Instruct.

	Chart & OCR	STEM	Spatial & Action	Knowl. & Recog.	Grnd., Cnt. & Search	Bench. Avg.
equal ratios	<b>+8.6</b>	+6.2	<b>+5.6</b>	+1.8	+5.6	<b>+5.8</b>
ratio $\propto (1 - \text{acc.})^\alpha$	+6.8	<b>+6.5</b>	+4.3	<b>+2.4</b>	+5.2	+5.2
ratio $\propto \text{area}^\alpha$	+7.0	+5.3	+4.1	+1.4	<b>+6.2</b>	+5.2
ratio $\propto \text{length}^\alpha$	+7.5	+6.4	+4.5	+1.7	+3.8	+4.8
w/o Knowl. & Recog.	+6.4	<b>+6.5</b>	+4.8	+1.9	+4.7	+4.9

**Table 3** Performance of different task category weighting schemes. Values are absolute score changes ( $\Delta$ ) over the base model.

questions to fix prompt clarity (e.g., GameQA, Magma) or drop high-error subsets.

Table 1 compares dataset selection strategies and shows our filtering has significant gains compared to taking a random subset from the candidate pool or sampling from the STEM or Chart subsets of the FineVision (Wiedmann et al., 2025) training set.

**Step 2. Filtering Examples.** After dataset-level filtering, many individual examples remain ambiguous, unanswerable, or incompatible with our reward functions. We apply two additional steps to filter and canonicalize individual prompts.

**Question filtering.** We use Qwen3-VL-235B-A22B-Instruct (Bai et al., 2025a) to remove ambiguous, image-irrelevant, or unverifiable questions. The model scores each datapoint on five criteria: (1) *relevance*, whether the image depicts what the question refers to; (2) *ambiguity*, whether the question is too vague or not a genuine question; (3) *language*, whether the question is in English; (4) *verifiability*, whether a single objectively correct answer can be derived from visible content; and (5) *numeric precision*, whether the required precision is visually unambiguous. Any triggered criterion removes the datapoint. We provide more details in Appendix A2.1.

**Answer canonicalization.** We normalize ground-truth answers using text-only Qwen3-235B-A22B-Instruct (Bai et al., 2025a) to ensure stable reward computation. Numeric answers are stripped of units and currency symbols, converted to decimal form, and evaluated as expressions. Samples with unsupported notation are filtered. Multiple-choice answers are normalized to a single canonical letter. Except for the captioning and instruction following task category, samples with multi-value answers, non-reducible symbolic expressions, or ambiguous descriptions requiring semantic matching are removed. A full list of canonicalization rules is in the Appendix A8. We do not apply it to Spatial & Action or Grounding tasks, since answers are already standardized for all datasets in these categories.

Effects of filtering are mixed across categories (Table 2): question filtering yields a clear gain on Spatial & Action (+1.9 pts) but slightly hurts Grounding & Counting (−0.6 pts), while answer canonicalization substantially improves Knowledge & Recognition (+2.1 pts) but is flat or marginally negative on other categories. Despite this variance, we apply both steps to all applicable task categories, as they generally help remove ambiguous samples from noisy datasets and the largest gains outweigh the small regressions.

**Step 3. Data Mixtures.** In our multi-task RL setting, the task category sampling distribution governs how training signal is allocated across skills. We investigate four task category weighting schemes (uniform, difficulty-weighted, dataset-size-weighted, reasoning-length-weighted), where the amount of samples per

**Table 1** Ablation on dataset filtering. Each model is trained for 1 epoch with GSPO and a math verify (Kydlicek, 2025) reward on 100k examples from a single task category. We report the average benchmark gain over the base model (Qwen2.5-VL-7B-Instruct) within each category.

	Chart & OCR	STEM
All Candidates	+2.5	−0.2
FineVision	+1.9	+2.9
Our Filtering	+3.4	+4.6

## Overview of Reward Types

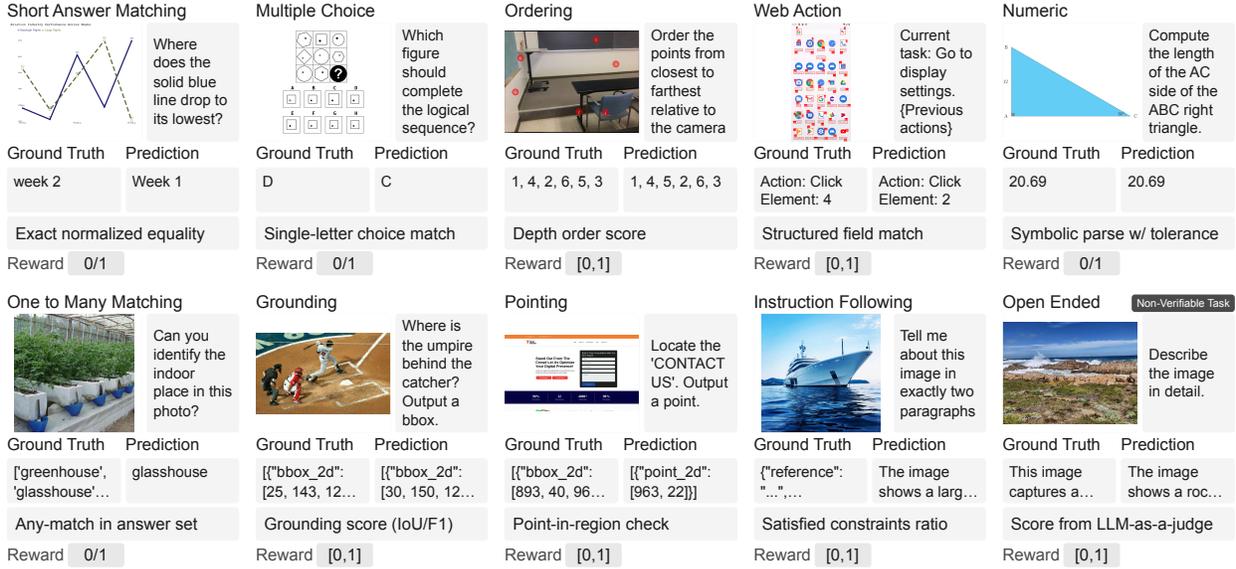


Figure 4 Reward design overview. We use ten verifiers corresponding to task types in our dataset.

batch is determined by a ratio proportional to the metric (e.g., difficulty).

Uniform sampling achieves the highest benchmark average gain (+5.8 pts over the base model), outperforming all alternative schemes. Alternatives yield gains on individual categories but at the cost of others (Figure 3). We use uniform task category weighting, as it achieves the best overall performance across all weighting schemes evaluated.

### 3.3 Reinforcement Learning

**Algorithmic Details.** At its core, RL maximizes the expected reward of the model’s response  $y$  given a visual input  $v$  and query  $q$ . Our RL algorithm builds on Group Relative Policy Optimization (GRPO) (Shao et al., 2024) and integrates a number of recent algorithmic advances from GSPO (Zheng et al., 2025), among others (Yu et al., 2025b).<sup>1</sup>

GSPO (Zheng et al., 2025) uses a sequence-level importance ratio rather than independent per-token ratios. We adopt asymmetric clip-higher (Yu et al., 2025b) (larger upper clipping bound), remove the KL penalty (Yu et al., 2025b; Liu et al., 2025b) to allow less-restricted updates, and apply a soft overlong penalty (Yu et al., 2025b) that linearly ramps from 0 to  $-1$  over the final 2,048 tokens before the context limit. The objective maximizes, over groups of  $G$  responses to the same prompt, a clipped surrogate loss weighted by the sequence-level importance ratio and scaled group advantage  $A_i = (r_i - \mu_g) / \sigma_g$ .

**Reward Formulation.** We use ten verifiers corresponding to the task types in our dataset (Figure 4): *string match*, *multiple choice*, *numeric* (via `math_verify` (Kydlíček, 2025)), *list string match*, *ordering* (Wu et al., 2025), *web action* (Sarch et al.), *grounding* (IoU/F1) (Yu et al., 2025a), *clicking* (Lu et al., 2025), *instruction following* (Ding et al., 2025; Pyatkin et al.), and *LLM-as-judge* (OLMo Team et al., 2025). For the LLM judge, we adopt the setup (Qwen3-32B (Bai et al., 2025a)) from OLMo 3 (OLMo Team et al., 2025), adding penalty guidelines for self-evaluative language to mitigate reward hacking we observed during training (details in Appendix A3). Full verifier definitions are in the Appendix A4.4. In addition to task-specific accuracy reward, we use a format reward that checks for `<think>...</think><answer>...</answer>` structure. We show in Section 6 that our reward design outperforms simple alternatives.

<sup>1</sup>Full equations and hyperparameter settings are in the Appendix A4.

RL Initial Model	Vero (Ours)				Open Weights Models				Fully Open RL Recipes			Propr.
	Vero Q3I-8B	Vero Q3T-8B	Vero Q25-7B	Vero Mi-7B	Q3VL 8B-Ins	Q3VL 8B-Thk	Q25VL 7B-Ins	MiMoVL 7B-RL	LLaVA OV1.5-RL	VL-Re thinker	Mo2-O 7B	GPT-5 Nano
Qwen3VL 8B Inst	Qwen3VL 8B Think	Qwen25VL 7B Inst	MiMoVL 7B-SFT	N/A	N/A	N/A	MiMoVL 7B-SFT	LLaVA OV1.5-Ins	Q25VL 7B-Ins	SFT Only	N/A	
<b>Chart &amp; OCR</b>												
ChartQA-Pro	60.2 <sup>+15.9</sup>	<b>63.2</b> <sup>+4.3</sup>	52.0 <sup>+8.7</sup>	61.7 <sup>+10.6</sup>	44.3 <sup>†</sup>	58.9 <sup>†</sup>	43.3 <sup>†</sup>	61.1 <sup>†</sup>	-	-	31.4 <sup>†</sup>	56.0 <sup>†</sup>
ChartQA	91.6 <sup>+2.0</sup>	90.4 <sup>+1.8</sup>	90.0 <sup>+2.7</sup>	90.6 <sup>+0.5</sup>	89.6	88.6	87.3	<b>94.4</b>	87.4	-	75.2 <sup>†</sup>	80.1 <sup>†</sup>
InfoVQA	87.8 <sup>+4.7</sup>	<b>88.3</b> <sup>+2.3</sup>	81.9 <sup>-0.7</sup>	87.1 <sup>+5.2</sup>	83.1	86.0	82.6	<b>90.1</b>	76.6	-	60.3 <sup>†</sup>	67.9 <sup>†</sup>
CharXiv <sub>Reason</sub>	53.7 <sup>+7.3</sup>	58.8 <sup>+5.8</sup>	<b>44.6</b> <sup>+2.1</sup>	<b>64.1</b> <sup>+11.0</sup>	46.4	53.0	42.5	60.9 <sup>†</sup>	-	-	35.6 <sup>†</sup>	51.2 <sup>†</sup>
ChartMuseum	49.6 <sup>+9.6</sup>	<b>51.7</b> <sup>+7.3</sup>	36.0 <sup>+9.2</sup>	49.9 <sup>+9.7</sup>	40.0	44.4	26.8	48.7 <sup>†</sup>	-	-	30.3 <sup>†</sup>	48.0 <sup>†</sup>
EvoChart	<b>75.7</b> <sup>+11.7</sup>	75.4 <sup>+2.4</sup>	66.9 <sup>+4.1</sup>	74.5 <sup>+3.5</sup>	64.0 <sup>†</sup>	73.0 <sup>†</sup>	62.8 <sup>†</sup>	73.4 <sup>†</sup>	-	-	51.0 <sup>†</sup>	63.3 <sup>†</sup>
Category Avg	69.8 <sup>+8.5</sup>	71.3 <sup>+4.0</sup>	61.9 <sup>+4.4</sup>	71.3 <sup>+6.8</sup>	61.2	67.3	57.6	<b>71.4</b>	-	-	47.3	61.1
<b>STEM</b>												
MMMU-Pro <sub>Std</sub>	59.8 <sup>+3.9</sup>	60.8 <sup>+0.4</sup>	43.4 <sup>+5.1</sup>	59.7 <sup>+3.6</sup>	55.9	60.4	38.3	59.4 <sup>†</sup>	39.9	41.7	31.9 <sup>†</sup>	<b>61.3</b> <sup>†</sup>
MMMU-Pro <sub>Vis</sub>	<b>57.2</b> <sup>+15.1</sup>	<b>57.2</b> <sup>+3.8</sup>	39.6 <sup>+7.3</sup>	56.3 <sup>+8.6</sup>	42.1 <sup>†</sup>	53.4 <sup>†</sup>	32.3 <sup>†</sup>	49.8 <sup>†</sup>	35.7	-	16.0 <sup>†</sup>	53.1 <sup>†</sup>
MathVision	59.0 <sup>+5.1</sup>	<b>62.8</b> <sup>+0.1</sup>	29.1 <sup>+4.0</sup>	59.7 <sup>+2.9</sup>	53.9	62.7	25.1	58.8 <sup>†</sup>	34.4	-	21.3 <sup>†</sup>	61.7 <sup>†</sup>
MathVista	78.7 <sup>+1.5</sup>	78.6 <sup>-2.8</sup>	74.5 <sup>+6.3</sup>	79.7 <sup>-0.9</sup>	77.2	<b>81.4</b>	68.2	80.4 <sup>†</sup>	72.3	73.7	53.6 <sup>†</sup>	70.2 <sup>†</sup>
Category Avg	63.7 <sup>+6.4</sup>	<b>64.8</b> <sup>+0.4</sup>	46.6 <sup>+5.7</sup>	63.9 <sup>+3.6</sup>	57.3	64.5	41.0	62.1	45.6	-	30.7	61.6
<b>Spatial &amp; Action</b>												
Blink	68.7 <sup>-0.4</sup>	66.2 <sup>+1.5</sup>	57.5 <sup>+1.1</sup>	61.4 <sup>-1.0</sup>	<b>69.1</b>	64.7	56.4	64.5 <sup>†</sup>	-	-	56.4 <sup>†</sup>	59.3 <sup>†</sup>
ERQA	43.2 <sup>-2.6</sup>	46.5 <sup>-0.3</sup>	40.0 <sup>-1.8</sup>	41.8 <sup>+2.3</sup>	45.8	<b>46.8</b>	41.8 <sup>†</sup>	43.5 <sup>†</sup>	-	-	43.5 <sup>†</sup>	45.5 <sup>†</sup>
GameQA <sub>Lite</sub>	52.3 <sup>+18.3</sup>	54.6 <sup>+14.8</sup>	46.7 <sup>+20.6</sup>	<b>55.4</b> <sup>+9.2</sup>	34.0 <sup>†</sup>	39.8 <sup>†</sup>	26.1 <sup>†</sup>	49.8 <sup>†</sup>	-	-	29.6 <sup>†</sup>	45.9 <sup>†</sup>
EmbSpatial	79.2 <sup>+0.7</sup>	79.8 <sup>-1.3</sup>	72.0 <sup>+1.3</sup>	74.3 <sup>+1.8</sup>	78.5	<b>81.1</b>	70.7 <sup>†</sup>	70.2 <sup>†</sup>	-	-	68.1 <sup>†</sup>	74.2 <sup>†</sup>
CV Bench	87.9 <sup>+2.4</sup>	<b>88.3</b> <sup>+2.3</sup>	81.1 <sup>+0.7</sup>	84.3 <sup>+1.7</sup>	85.5 <sup>†</sup>	86.0 <sup>†</sup>	80.4 <sup>†</sup>	83.5 <sup>†</sup>	82.9	-	81.7 <sup>†</sup>	82.5 <sup>†</sup>
Category Avg	66.3 <sup>+3.7</sup>	<b>67.1</b> <sup>+3.4</sup>	59.5 <sup>+4.4</sup>	63.4 <sup>+2.8</sup>	62.6	63.7	55.1	62.3	-	-	55.9	61.5
<b>Knowledge &amp; Recognition</b>												
RealWorldQA	73.3 <sup>+1.8</sup>	72.9 <sup>-0.6</sup>	71.6 <sup>+3.1</sup>	67.5 <sup>-0.1</sup>	71.5	<b>73.5</b>	68.5	68.6 <sup>†</sup>	68.4	-	73.3 <sup>†</sup>	65.9 <sup>†</sup>
SimpleVQA <sub>En</sub>	45.2 <sup>+1.0</sup>	44.0 <sup>-0.9</sup>	<b>50.0</b> <sup>+5.5</sup>	46.5 <sup>+4.6</sup>	44.2 <sup>†</sup>	44.9 <sup>†</sup>	44.5 <sup>†</sup>	40.9 <sup>†</sup>	-	-	30.8 <sup>†</sup>	36.6 <sup>†</sup>
FVQA	24.6 <sup>-1.4</sup>	22.6 <sup>-1.6</sup>	26.3 <sup>+4.4</sup>	28.1 <sup>-1.7</sup>	26.0	24.2	21.9	<b>31.8</b> <sup>†</sup>	-	-	17.7 <sup>†</sup>	29.4 <sup>†</sup>
MM-Vet v2	70.2 <sup>+2.6</sup>	<b>84.4</b> <sup>+9.9</sup>	66.3 <sup>+3.4</sup>	76.6 <sup>+16.9</sup>	67.6	74.5	62.9	61.2 <sup>†</sup>	-	-	60.8 <sup>†</sup>	71.5 <sup>†</sup>
Category Avg	53.3 <sup>+1.0</sup>	<b>56.0</b> <sup>+1.7</sup>	53.5 <sup>+4.1</sup>	54.7 <sup>+4.9</sup>	52.3	54.3	49.5	50.6	-	-	45.6	50.9
<b>Grounding, Counting, Search</b>												
CountBenchQA	90.4 <sup>+1.6</sup>	<b>91.9</b> <sup>+2.3</sup>	84.9 <sup>-1.0</sup>	83.1 <sup>-2.2</sup>	88.8	89.6	85.9 <sup>†</sup>	86.4 <sup>†</sup>	86.8	-	89.4 <sup>†</sup>	75.4 <sup>†</sup>
CountQA	33.9 <sup>+5.4</sup>	<b>36.1</b> <sup>+4.3</sup>	24.1 <sup>+3.2</sup>	24.5 <sup>-1.1</sup>	28.5 <sup>†</sup>	31.8 <sup>†</sup>	20.9 <sup>†</sup>	27.4 <sup>†</sup>	-	-	32.1 <sup>†</sup>	25.7 <sup>†</sup>
MMERealWorld	<b>57.8</b> <sup>+10.7</sup>	53.9 <sup>+9.1</sup>	52.4 <sup>+7.9</sup>	49.0 <sup>+5.2</sup>	47.1	44.8	44.5	48.7 <sup>†</sup>	-	-	44.4 <sup>†</sup>	49.8 <sup>†</sup>
VStarBench	<b>89.5</b> <sup>+7.3</sup>	82.2 <sup>+5.8</sup>	86.4 <sup>+6.3</sup>	84.8 <sup>+2.6</sup>	82.2	76.4 <sup>†</sup>	80.1 <sup>†</sup>	84.3	79.1	-	73.8 <sup>†</sup>	71.2 <sup>†</sup>
AerialVG	30.0 <sup>-2.2</sup>	31.9 <sup>+19.3</sup>	27.4 <sup>+5.1</sup>	15.5 <sup>-4.6</sup>	<b>32.2</b> <sup>†</sup>	12.6 <sup>†</sup>	22.3 <sup>†</sup>	22.2 <sup>†</sup>	-	-	-	-
VisualProbe	<b>53.9</b> <sup>+6.2</sup>	46.1 <sup>+6.8</sup>	48.8 <sup>+2.8</sup>	48.6 <sup>+1.7</sup>	47.7 <sup>†</sup>	39.3 <sup>†</sup>	46.0 <sup>†</sup>	52.2 <sup>†</sup>	-	-	34.8 <sup>†</sup>	41.5 <sup>†</sup>
ScreenSpot	<b>93.6</b> <sup>+7.0</sup>	91.4 <sup>+5.9</sup>	89.9 <sup>+4.3</sup>	85.0 <sup>-2.9</sup>	86.6 <sup>†</sup>	85.5 <sup>†</sup>	85.6 <sup>†</sup>	87.3 <sup>†</sup>	-	-	75.9 <sup>†</sup>	-
ScreenSpotPro	<b>61.4</b> <sup>+13.6</sup>	48.1 <sup>+12.6</sup>	38.1 <sup>+14.2</sup>	39.8 <sup>+7.0</sup>	47.8	35.5	23.9 <sup>†</sup>	37.4 <sup>†</sup>	-	-	19.4 <sup>†</sup>	-
Category Avg	<b>63.8</b> <sup>+6.2</sup>	60.2 <sup>+8.3</sup>	56.5 <sup>+5.4</sup>	53.8 <sup>+0.7</sup>	57.6	51.9	51.1	55.7	-	-	-	-
<b>Captioning &amp; IF</b>												
MM-MTBench	80.3 <sup>+5.9</sup>	74.6 <sup>-3.2</sup>	66.1 <sup>+7.2</sup>	<b>90.0</b> <sup>+14.3</sup>	74.4	77.8	58.9	79.2 <sup>†</sup>	-	-	33.3 <sup>†</sup>	72.7 <sup>†</sup>
MIABench	<b>93.5</b> <sup>+2.4</sup>	93.1 <sup>+1.6</sup>	82.8 <sup>+1.0</sup>	92.1 <sup>+3.7</sup>	91.1	91.5	81.8	88.4 <sup>†</sup>	-	-	77.9 <sup>†</sup>	92.0 <sup>†</sup>
MMIFEval	77.7 <sup>+8.5</sup>	<b>80.4</b> <sup>+5.4</sup>	62.8 <sup>+9.2</sup>	76.7 <sup>+7.5</sup>	69.2	75.0	53.6	66.1	-	-	54.4 <sup>†</sup>	78.0 <sup>†</sup>
Category Avg	83.8 <sup>+5.6</sup>	82.7 <sup>+1.3</sup>	70.6 <sup>+5.8</sup>	<b>86.3</b> <sup>+8.5</sup>	78.2	81.4	64.8	77.9	-	-	55.2	80.9
Overall Avg	<b>66.0</b> <sup>+5.5</sup>	65.9 <sup>+4.0</sup>	57.8 <sup>+4.9</sup>	63.6 <sup>+4.0</sup>	60.5	61.9	52.9	62.4	-	-	-	-

**Table 4** Evaluation results. **Vero** columns show the initial models trained with RL on our dataset; +x / -x deltas indicate improvement/decline over the respective initial model. <sup>†</sup> indicates results evaluated by us. All other results are taken from official technical reports.

## 4 Experiments

**Evaluation Suite.** We evaluate on **VeroEvalSuite**, a collection of 30 benchmarks spanning the six visual reasoning categories defined in Section 3. For Chart & OCR, we use ChartQA-Pro (Masry et al., 2025), ChartQA (Masry et al., 2022), InfoVQA (Mathew et al., 2022), CharXiv (Wang et al., 2024b), ChartMuseum (Tang et al., 2025), and EvoChart (Huang et al., 2025). For STEM, we use MMMU-PRO (Yue et al., 2025a) (Standard and Vision splits), MathVision (Wang et al., 2024a), and MathVista (Lu et al., 2023). For Spatial & Action, we use Blink (Fu et al., 2024), ERQA (Team et al., 2025), GameQA (Tong et al., 2025), EmbSpatial (Du et al., 2024), and CVBench (Tong et al., 2024). For Knowledge & Recognition, we use RealWorldQA (xAI, 2024) and SimpleVQA (Cheng et al., 2025a) (English). For Grounding, Counting & Visual Search, we use CountBenchQA (Paiss et al., 2023; Beyer et al., 2024), CountQA (Tamarapalli et al., 2025), MMERealWorld (Zhang et al., 2024), VStarBench (Cheng et al., 2025b), AerialVG (Liu et al., 2025a), VisualProbe (Lai et al., 2025), ScreenSpot (Cheng et al., 2024), and ScreenSpotPro (Li et al., 2025). For Captioning & Instruction Following, we use MM-MTBench (Ying et al., 2024), MIABench (Qian et al., 2024), and MMIFEval (Ding et al., 2025). We evaluate all models using the Imms-eval (Zhang et al., 2025a) framework. For benchmarks requiring open-ended or long-form evaluation, we use Qwen3-32B as the judge model. We follow the official evaluation protocols specified by each benchmark’s authors.

**Baselines.** We compare against: (1) base VLMs without native <think> tokens (Qwen2.5-VL-7B-Instruct (Bai et al., 2025b), Qwen3-VL-8B-Instruct (Bai et al., 2025a), Molmo2-O-7B (Clark et al., 2026)), (2) models trained to do native CoT with <think> tokens (Qwen3-VL-8B-Thinking (Bai et al., 2025a), MiMo-VL-7B-RL (Yue et al., 2025b)), (3) existing fully open RL-trained models (VL-Rethinker-7B (Wang et al., 2025), LLaVA-OV-1.5-RL (An et al., 2025)), and (4) a proprietary reasoning model gpt-5-nano-2025-08-07 (Singh et al., 2025) with medium reasoning effort. For baseline results, we prioritize scores from official technical reports and benchmark leaderboards. When published results are unavailable, we evaluate the models ourselves (indicated by †) and follow the published benchmark guidelines.

### 4.1 Evaluation Results

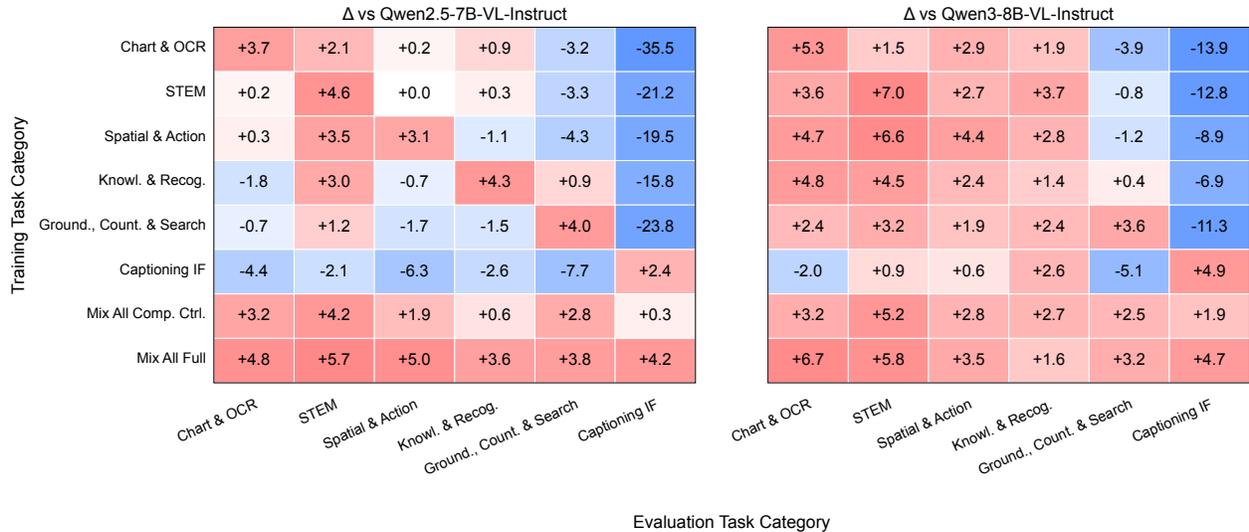
We report results in Table 4 and highlight the following observations.

**State-of-the-art performance with a fully-open RL recipe.** Our best models **Vero-Q3T-8B** and **Vero-Q3I-8B** achieve the highest overall averages (65.5 and 65.5, respectively) among all 8B-parameter VLMs evaluated, outperforming baselines across the six task categories. In particular, **Vero** outperforms the next best 8B model, Qwen3-VL-8B-Thinking, by +8.0 on Grounding, Counting & Search and +3.4 on Chart & OCR. Against existing fully open RL-recipe baselines, our best **Vero** variant outperforms LLaVA-OV-1.5-RL on 10 of 10 overlaps, VL-Rethinker on 2 of 2 overlaps, and Molmo2-O-7B on 28 of 29 overlaps.

**Consistent gains across base models, beating proprietary recipes.** **Vero** training yields improvements across four different base models. **Vero-Q3I-8B** improves over Qwen3-VL-8B-Instruct by +7.5 on Chart & OCR, +5.7 on STEM, and +3.9 on Spatial & Action, and outperforms Qwen3-VL-8B-Thinking on 20 of 30 benchmarks despite the latter being trained on additional proprietary long chain-of-thought data. **Vero-Q3T-8B** improves over Qwen3-VL-8B-Thinking on 23 of 30 benchmarks, with notable gains on Grounding, Counting & Search (+8.1) and Captioning & Instruction Following (+0.1).

**Vero-Q25-7B** similarly improves over Qwen2.5-VL-7B-Instruct by +4.3, +5.6, and +4.4 on Chart & OCR, STEM, and Spatial & Action, respectively. Notably, **Vero-Q25-7B** surpasses Qwen3-VL-8B-Instruct on two category averages, Chart & OCR (61.9 vs. 61.2) and Knowledge & Recognition (53.5 vs. 52.3), despite starting from a substantially weaker base model (Qwen2.5-VL-7B-Instruct at 52.9 overall vs. Qwen3-VL-8B-Instruct at 60.5). This demonstrates that RL training on our dataset can close a 7.6-point base model gap and even surpass the stronger model on several tasks.

**Vero-Mi-7B**, trained on MiMo-VL-7B-SFT with our fully open recipe, improves by +4.0 overall, with the largest gains on Captioning & Instruction Following (+8.5), Chart & OCR (+6.8), and Knowledge & Recognition (+4.9). **Vero-Mi-7B** also outperforms MiMo-VL-7B-RL, which trains on the same initial model but uses a proprietary RL recipe with non-public data, on 4 of 6 category averages (STEM +1.8, Captioning & IF +8.4, and Knowledge & Recognition +4.1), showing that a fully open recipe can surpass proprietary pipelines.



**Figure 5** Cross-task generalization. Each row shows a model trained on a single task category (or mixture). Values are absolute score changes relative to the base model. Single-task training yields selective transfer, while mixing achieves consistent gains.

**Improvements not limited to a single domain.** Unlike prior open RL-trained VLMs that focus primarily on STEM or math reasoning (e.g., VL-Rethinker), **Vero** yields substantial improvements across all six task categories. For example, **Vero-Q3I-8B** achieves strong gains over the initial model on ChartQA Pro (+17.4), MMMU-Pro Vision (+14.5), and grounding and search benchmarks such as MMERealWorld (+8.9) and ScreenSpotPro (+12.7), demonstrating that multi-task RL training produces broadly capable models rather than specialists.

## 5 Cross-Task Generalization

We study cross-task generalization by training models on each individual task category (100k samples, 1 epoch) and evaluating across all six categories. We compare against a model trained on a mixture of all categories with the same total number of training samples (100k, compute-controlled) and the full dataset (600k). We report results on two base models in Figure 5.

**Single-task training frequently produces neutral or negative transfer on non-target tasks.** On Qwen2.5-VL, nearly all single-task-category models degrade Grounding, Counting & Search performance (e.g.,  $-3.2$  from Chart & OCR,  $-3.3$  from STEM,  $-4.3$  from Spatial & Action), and training on Captioning & Instruction Following alone reduces performance across all other categories ( $-4.4$  to  $-7.7$ ). Conversely, training on any non-captioning task category severely degrades Captioning & Instruction Following ( $-15.8$  to  $-35.5$  on Qwen2.5-VL). These patterns hold on Qwen3-VL, where single-task-category models similarly hurt Grounding ( $-0.8$  to  $-5.1$  for non-grounding domains) and Captioning & Instruction Following ( $-6.9$  to  $-13.9$ ). However, in certain task categories, we do observe selective positive transfer: STEM training improves Chart & OCR (+3.6 on Qwen3-VL), and Spatial & Action training yields strong gains on STEM (+3.5 on Qwen2.5-VL, +6.6 on Qwen3-VL).

**Diverse task mixing eliminates negative cross-task transfer.** Even with the same compute budget, the mixed model achieves positive gains across all categories on both base models (+0.3 to +4.2 on Qwen2.5-VL; +1.9 to +5.2 on Qwen3-VL), avoiding the catastrophic losses seen with single-domain training. Training on the full 600k mixture further amplifies these gains. These patterns are consistent across both base models, confirming that multi-task RL training is important for producing broadly capable models.

**Table 5** Ablation studies. All results on Qwen2.5-VL-7B-Instruct. Tables (a)–(d) report absolute scores. All runs are trained 1 epoch on the 600k mixture unless otherwise specified.

(a) Open-ended task support						(b) Reward design						
	MM-MTBench	MIA-Bench	MMIF-Eval	Avg.	Other domain avg.	Chart & OCR	STEM	Spatial & Action	Knowl. & Recog.	Grnd, Cnt. & Search	Cap. & IF	Overall Avg.
Base	58.9	81.8	53.6	64.8	50.6	57.6	41.0	55.1	49.4	50.1	64.8	52.4
+ Ans. tag	22.4	27.2	30.7	26.8	54.7	61.4	46.1	58.5	50.1	51.0	34.3	51.8
+ Sys + Boxed	40.3	57.6	45.3	47.7	53.8	<b>61.9</b>	<b>46.7</b>	<b>59.4</b>	<b>53.5</b>	<b>55.0</b>	<b>70.6</b>	<b>57.2</b>
+ Cap & IF	<b>66.1</b>	<b>82.8</b>	<b>62.8</b>	<b>70.6</b>	<b>55.6</b>							

(c) SFT vs. RL								(d) RL algorithm (1/4 epoch, 5 domains)							
	Chart & OCR	STEM	Spatial & Action	Knowl. & Recog.	Grnd, Cnt. & Search	Cap. & IF	Overall Avg.	Chart & OCR	STEM	Spatial & Action	Knowl. & Recog.	Grnd, Cnt. & Search	Avg.	Avg. Entropy	
Base	57.6	41.0	55.1	49.4	50.1	64.8	52.4	58.9	45.3	57.1	49.7	52.2	54.3	0.22±0.15	
FineVision SFT	54.8	37.4	52.1	45.3	40.1	52.2	46.2	59.0	<b>45.4</b>	<b>58.4</b>	50.4	53.0	<b>54.7</b>	0.58±0.11	
<b>Vero</b> SFT	52.5	40.1	58.1	50.8	52.7	64.1	52.8	<b>59.2</b>	44.4	58.1	48.2	<b>53.0</b>	54.3	0.50±0.11	
<b>Vero</b> RL	<b>61.9</b>	<b>46.7</b>	<b>59.4</b>	<b>53.5</b>	<b>55.0</b>	<b>70.6</b>	<b>57.2</b>								

## 6 Ablations

**RL without open-ended prompts leads to visual chat deficits.** We ablate our open-ended instruction following design in Table 5(a) by comparing: (1) answer tag parsing only used by previous works (Zhang et al., 2025b), (2) adding system prompt guidelines and latex boxed parsing to encourage richer answer formats, and (3) our full approach with boxed formatting and the captioning & instruction following task category with LLM-judge rewards. Without open-ended training, RL degrades instruction following quality, as the model overfits to producing short, boxed answers. System prompt and latex boxed parsing partially recover performance, but the full integration of the open ended task category with judge-based rewards preserves and often improves the model’s ability to produce fluent, instruction-adherent responses alongside structured reasoning.

**Multi-task training requires more expressive reward design.** We compare our multi-route reward design against math\_verify (Kydliček, 2025), a widely used reward function that performs answer and number extraction, parsing, and grading. Results in Table 5(b) show that our reward design, which routes answers through type-specific comparison functions (exact match, numeric tolerance, set matching, and LLM-judge evaluation), achieves stronger performance than math\_verify across task categories. math\_verify lacks the flexibility to handle the diverse answer formats.

**RL outperforms SFT on our dataset.** We compare SFT and RL training on our dataset in Table 5(c). The SFT model is trained to directly output the final answer without chain-of-thought or <think> tokens, reflecting the standard SFT paradigm used in many existing VLMs. SFT on our dataset produces gains on most tasks and outperforms SFT on a recent post-training dataset FineVision (Wiedmann et al., 2025). However, RL (GSPO with our multi-route reward) yields more consistent improvements across all task categories.

**GSPO outperforms GRPO and DAPO and leads to more stable entropy.** We compare three RL algorithms (DAPO, GRPO, and GSPO) using the same base model (Qwen2.5-VL-7B-Instruct), reward design, and training dataset. Consistent with recent findings (Zhang et al., 2025b), results in Table 5(d) show that GSPO achieves the highest average score (54.7) across all task categories, outperforming both GRPO (54.3) and DAPO (54.3). GSPO also maintains substantially more stable entropy throughout training ( $0.58 \pm 0.11$ ) compared to GRPO ( $0.50 \pm 0.11$ ) and especially DAPO ( $0.22 \pm 0.15$ ), suggesting that GSPO’s group-level normalization better preserves exploration capacity and avoids the premature policy collapse observed with alternative algorithms.

	Captioning & IF	Chart & OCR	Grnd., Cnt. & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.40	0.46	0.40	0.43	0.51	0.47	0.46
Adaptive Detail Mgmt.	0.77	0.75	0.54	0.69	0.81	0.82	0.70
Arithmetic Calculation	0.20	0.28	0.21	0.24	0.26	0.30	0.21
Backtracking	0.17	0.33	0.12	0.20	0.36	0.48	0.22
Backward Chaining	0.04	0.08	0.07	0.08	0.10	0.12	0.16
Causal Organization	0.42	0.51	0.42	0.47	0.59	0.54	0.51
Compositionality	0.95	0.94	0.91	0.92	0.97	0.96	0.96
Conceptual-Level Proc.	0.67	0.66	0.60	0.64	0.75	0.66	0.66
Context Alignment	0.76	0.76	0.61	0.74	0.85	0.79	0.77
Context Awareness	0.65	0.52	0.46	0.55	0.72	0.58	0.67
Decomp. & Integration	0.77	0.78	0.76	0.79	0.82	0.83	0.85
Forward Chaining	0.88	0.89	0.90	0.93	0.93	0.90	0.94
Goal Management	0.88	0.92	0.84	0.95	0.96	0.95	0.97
Hierarchical Org.	0.77	0.73	0.66	0.72	0.81	0.74	0.75
Knowledge Struct. Align.	0.94	0.90	0.83	0.90	0.96	0.91	0.87
Logical Coherence	0.99	0.98	0.97	0.98	0.99	0.98	1.00
Mental Imagery Sim.	0.64	0.58	0.47	0.54	0.64	0.65	0.53
Network Organization	0.57	0.57	0.47	0.53	0.64	0.58	0.55
Ordinal Organization	0.53	0.58	0.56	0.61	0.64	0.61	0.64
Pattern Recognition	0.43	0.56	0.48	0.53	0.57	0.57	0.50
Perception-then-Reas.	0.86	0.77	0.64	0.68	0.84	0.80	0.62
Productivity	0.36	0.44	0.29	0.39	0.48	0.51	0.43
Repr. Restructuring	0.47	0.54	0.40	0.47	0.58	0.58	0.53
Selective Attention	0.98	0.97	0.95	0.98	0.98	0.98	1.00
Self-Awareness	0.76	0.80	0.49	0.62	0.86	0.84	0.75
Self-Evaluation	0.86	0.91	0.46	0.73	0.94	0.94	0.73
Sequential Organization	0.98	0.97	0.96	0.98	0.98	0.98	0.99
Spatial Organization	0.76	0.71	0.71	0.73	0.75	0.75	0.71
Strategy Selection	0.65	0.70	0.57	0.69	0.80	0.78	0.80
Syst. Regional Synth.	0.69	0.74	0.52	0.64	0.78	0.81	0.61
Temporal Organization	0.62	0.60	0.61	0.62	0.71	0.69	0.66
Verification	0.89	0.87	0.65	0.83	0.93	0.91	0.81
Visual Foraging	0.87	0.90	0.72	0.85	0.94	0.94	0.84
Visual Ref. / Grounding	0.78	0.71	0.66	0.72	0.76	0.75	0.64

**Figure 6** High-level cognitive behavior presence rates across task categories for single-task and mixed-task trained models.

## 7 Chain-of-Thought Behavior Analysis

### 7.1 High-Level Cognitive Foundations

**Experimental Setup.** We evaluate model behavior using the cognitive framework of Kargupta et al. [Kargupta et al. \(2025\)](#), which defines 28 textual reasoning behaviors, supplemented with six behaviors for visual analysis. We compute the presence rate of each behavior for every trained model from Section 5 across all six validation domains. Behavior presence rate on **Vero** trained on Qwen3-VL-8B-Instruct is shown in Figure 6, and we report additional base models in the Appendix A11.

**Each task category elicits a distinct cognitive behavioral profile.** The behavioral profiles reveal that training on each domain elicits different cognitive behaviors. Captioning often uses mental imagery simulation (0.64 vs. 0.57 cross-domain average in Qwen3), chart-trained models trigger systematic regional synthesis (0.74 vs. 0.68), and spatial reasoning often uses perception-then-reasoning sequencing (0.84 vs. 0.73). Grounding tasks more often suppress introspective behaviors, with self-awareness dropping to 0.49 (vs. 0.73), redirecting capacity toward directed visual search. In contrast, domains requiring multi-step integration elicit higher-order behaviors, with STEM tasks showing elevated backtracking (0.48 vs. 0.27). The mixed-domain setting increases strategy selection (0.80 vs. 0.71), indicating that the model first selects a reasoning approach before executing it.

### 7.2 Skill Analysis

**Experimental Setup.** We extract domain-specific skills from model reasoning traces following Didolkar et al. [Didolkar et al. \(2025\)](#). A deduplication pipeline ensures uniqueness of extracted skills, after which skill embeddings are clustered via agglomerative clustering and labeled with GPT-4o. We train a logistic regression probe on the resulting skill embeddings (Qwen3-Embedding-8B, 4,096-d) with 1,500 skills per domain, evaluated via 5-fold Stratified Group K-Fold cross-validation with mean centering and  $l_2$  normalization. We report the task confusion matrix in Figure 7.

**Each task category cultivates a largely distinct skill repertoire.** The probe achieves 0.670 overall accuracy (Fig. 7), confirming that skill distributions are domain-specific. Chart and spatial

True Domain \ Predicted Domain	Captioning IF	Chart & OCR	Grnd., Cnt., & Search	Knowl. & Recog.	Spatial & Action	STEM
Captioning IF	0.81	0.04	0.02	0.06	0.02	0.04
Chart & OCR	0.03	0.82	0.04	0.05	0.03	0.04
Grnd., Cnt., & Search	0.02	0.04	0.80	0.10	0.04	0.00
Knowl. & Recog.	0.08	0.08	0.11	0.59	0.08	0.05
Spatial & Action	0.03	0.06	0.07	0.05	0.77	0.02
STEM	0.03	0.02	0.00	0.05	0.05	0.84

**Figure 7** Confusion matrix of a logistic regression probe trained on skill embeddings, showing task-category-level separability across task categories.

tasks yield the most distinctive skills: chart behaviors center on data-reading operations (e.g., *cross-reference axes in visual data*), while spatial behaviors reflect physical and game-state reasoning (e.g., *cross validate path with grid state*). Knowledge-domain skills are least separable (0.49 accuracy), frequently confused with grounding (0.15 confusion rate), because knowledge reasoning relies on visual grounding operations that are semantically indistinguishable from grounding-specific skills at the description level.

## 8 Conclusion

We presented **Vero**, a fully open vision-language reasoning family trained with single-stage RL on 600K samples from 59 datasets spanning six capability categories. Our ablations show that diverse task mixing, uniform weighting, and task-routed rewards are important for positive cross-category transfer of chain of thought behavior and model performance. **Vero** outperforms Qwen3-VL-8B-Thinking on 23 of 30 benchmarks and surpasses MiMo-VL-7B-RL trained from the same base checkpoint. All datasets, training code, and models are publicly released to facilitate future research on visual reasoning.

## References

- An, X., Xie, Y., Yang, K., Zhang, W., Zhao, X., Cheng, Z., Wang, Y., Xu, S., Chen, C., Zhu, D., et al. Llava-onevision-1.5: Fully open framework for democratized multimodal training. *arXiv preprint arXiv:2509.23661*, 2025.
- Bai, S., Cai, Y., Chen, R., Chen, K., Chen, X., Cheng, Z., Deng, L., Ding, W., Gao, C., Ge, C., et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025a.
- Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J., et al. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025b.
- Beyer, L., Steiner, A., Pinto, A. S., Kolesnikov, A., Wang, X., Salz, D., Neumann, M., Alabdulmohsin, I., Tschannen, M., Bugliarello, E., et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- Cheng, K., Sun, Q., Chu, Y., Xu, F., YanTao, L., Zhang, J., and Wu, Z. Seeclick: Harnessing gui grounding for advanced visual gui agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9313–9332, 2024.
- Cheng, X., Zhang, W., Zhang, S., Yang, J., Guan, X., Wu, X., Li, X., Zhang, G., Liu, J., Mai, Y., et al. Simplevqa: Multimodal factuality evaluation for multimodal large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4637–4646, 2025a.
- Cheng, Z., Hu, J., Liu, Z., Si, C., Li, W., and Gong, S. V-star: Benchmarking video-llms on video spatio-temporal reasoning. *arXiv preprint arXiv:2503.11495*, 2025b.
- Clark, C., Zhang, J., Ma, Z., Park, J. S., Salehi, M., Tripathi, R., Lee, S., Ren, Z., Kim, C. D., Yang, Y., et al. Molmo2: Open weights and data for vision-language models with video understanding and grounding. *arXiv preprint arXiv:2601.10611*, 2026.
- Comanici, G., Bieber, E., Schaeckermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang, D., Rosen, E., et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Wang, P., et al. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *Nature*, 645, 2025.
- Deitke, M., Clark, C., Lee, S., Tripathi, R., Yang, Y., Park, J. S., Salehi, M., Muennighoff, N., Lo, K., Soldaini, L., Lu, J., Anderson, T., Branson, E., Ehsani, K., Ngo, H., Chen, Y., Patel, A., Yatskar, M., Callison-Burch, C., Head, A., Hendrix, R., Bastani, F., VanderBilt, E., Lambert, N., Chou, Y., Chheda, A., Sparks, J., Skjonsberg, S., Schmitz, M., Sarnat, A., Bischoff, B., Walsh, P., Newell, C., Wolters, P., Gupta, T., Zeng, K.-H., Borchardt, J., Groeneveld, D., Nam, C., Lebrecht, S., Wittlif, C., Schoenick, C., Michel, O., Krishna, R., Weihs, L., Smith, N. A., Hajishirzi, H., Girshick, R., Farhadi, A., and Kembhavi, A. Molmo and PixMo: Open Weights and Open Data for State-of-the-Art Vision-Language Models, December 2024. URL <http://arxiv.org/abs/2409.17146>. arXiv:2409.17146 [cs] version: 2.
- Didolkar, A., Ballas, N., Arora, S., and Goyal, A. Metacognitive reuse: Turning recurring llm reasoning into concise behaviors, 2025. URL <https://arxiv.org/abs/2509.13237>.
- Ding, S., Wu, S., Zhao, X., Zang, Y., Duan, H., Dong, X., Zhang, P., Cao, Y., Lin, D., and Wang, J. Mm-ifengine: Towards multimodal instruction following. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1099–1109, 2025.
- Du, M., Wu, B., Li, Z., Huang, X.-J., and Wei, Z. Embspatial-bench: Benchmarking spatial understanding for embodied tasks with large vision-language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 346–355, 2024.
- Fu, X., Hu, Y., Li, B., Feng, Y., Wang, H., Lin, X., Roth, D., Smith, N. A., Ma, W.-C., and Krishna, R. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pp. 148–166. Springer, 2024.

- Hong, W., Yu, W., Gu, X., Wang, G., Gan, G., Tang, H., et al. GLM-4.5V and GLM-4.1V-Thinking: Towards versatile multimodal reasoning with scalable reinforcement learning. *arXiv preprint arXiv:2507.01006*, 2026.
- Huang, M., Lai, H., Zhang, X., Wu, W., Ma, J., Zhang, L., and Liu, J. Evochart: A benchmark and a self-training approach towards real-world chart understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 3680–3688, 2025.
- Kargupta, P., Li, S. S., Wang, H., Lee, J., Chen, S., Ahia, O., Light, D., Griffiths, T. L., Kleiman-Weiner, M., Han, J., Celikyilmaz, A., and Tsvetkov, Y. Cognitive foundations for reasoning and their manifestation in llms, 2025. URL <https://arxiv.org/abs/2511.16660>.
- Kimi Team. Kimi-VL technical report. *arXiv preprint arXiv:2504.07491*, 2025.
- Kimi Team, Bai, T., Bai, Y., Bao, Y., et al. Kimi K2.5: Visual agentic intelligence. *arXiv preprint arXiv:2602.02276*, 2026.
- Kydlíček, H. Math-verify: Math verification library. <https://github.com/huggingface/Math-Verify>, 2025.
- Lai, X., Li, J., Li, W., Liu, T., Li, T., and Zhao, H. Mini-o3: Scaling up reasoning patterns and interaction turns for visual search, 2025. URL <https://arxiv.org/abs/2509.07969>.
- Li, K., Meng, Z., Lin, H., Luo, Z., Tian, Y., Ma, J., Huang, Z., and Chua, T.-S. Screenspot-pro: Gui grounding for professional high-resolution computer use. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pp. 8778–8786, 2025.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
- Liu, J., Chen, Q., Wang, Z., Tang, Y., Zhang, Y., Yan, C., Wang, D., Li, X., and Zhao, B. Aerialvg: A challenging benchmark for aerial visual grounding by exploring positional relations, 2025a. URL <https://arxiv.org/abs/2504.07836>.
- Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee, W. S., and Lin, M. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025b.
- Lu, P., Bansal, H., Xia, T., Liu, J., Li, C., Hajishirzi, H., Cheng, H., Chang, K.-W., Galley, M., and Gao, J. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- Lu, Z., Chai, Y., Guo, Y., Yin, X., Liu, L., Wang, H., Xiao, H., Ren, S., Xiong, G., and Li, H. Ui-r1: Enhancing efficient action prediction of gui agents by reinforcement learning. *arXiv preprint arXiv:2503.21620*, 2025.
- Masry, A., Do, X. L., Tan, J. Q., Joty, S., and Hoque, E. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the association for computational linguistics: ACL 2022*, pp. 2263–2279, 2022.
- Masry, A., Islam, M. S., Ahmed, M., Bajaj, A., Kabir, F., Kartha, A., Laskar, M. T. R., Rahman, M., Rahman, S., Shahmohammadi, M., et al. Chartqapro: A more diverse and challenging benchmark for chart question answering. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 19123–19151, 2025.
- Mathew, M., Bagal, V., Tito, R., Karatzas, D., Valveny, E., and Jawahar, C. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1697–1706, 2022.
- OLMo Team, Ettinger, A., Bertsch, A., Kuehl, B., Graham, D., Heineman, D., Groeneveld, D., Brahman, F., Timbers, F., Ivison, H., et al. Olmo 3. *arXiv preprint arXiv:2512.13961*, 2025.
- Paiss, R., Ephrat, A., Tov, O., Zada, S., Mosseri, I., Irani, M., and Dekel, T. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3170–3180, 2023.

- Pyatkin, V., Malik, S., Graf, V., Ivison, H., Huang, S., Dasigi, P., Lambert, N., and Hajishirzi, H. Generalizing verifiable instruction following. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Qian, Y., Ye, H., Fauconnier, J.-P., Gräsch, P., Yang, Y., and Gan, Z. Mia-bench: Towards better instruction following evaluation of multimodal llms. *arXiv preprint arXiv:2407.01509*, 2024.
- Sarch, G. H., Saha, S., Khandelwal, N., Jain, A., Tarr, M. J., Kumar, A., and Fragkiadaki, K. Grounded reinforcement learning for visual reasoning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Singh, A., Fry, A., Perelman, A., Tart, A., Ganesh, A., El-Kishky, A., McLaughlin, A., Low, A., Ostrow, A., Ananthram, A., et al. Openai gpt-5 system card. *arXiv preprint arXiv:2601.03267*, 2025.
- Tamarapalli, J. S., Grover, R., Pande, N., and Yerramilli, S. Countqa: How well do mllms count in the wild? *arXiv preprint arXiv:2508.06585*, 2025.
- Tang, L., Kim, G., Zhao, X., Lake, T., Ding, W., Yin, F., Singhal, P., Wadhwa, M., Liu, Z. L., Sprague, Z., et al. Chartmuseum: Testing visual reasoning capabilities of large vision-language models. *arXiv preprint arXiv:2505.13444*, 2025.
- Team, G., Anil, R., Borgeaud, S., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., Millican, K., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Team, G., Georgiev, P., Lei, V. I., Burnell, R., Bai, L., Gulati, A., Tanzer, G., Vincent, D., Pan, Z., Wang, S., et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- Team, G. R., Abeyruwan, S., Ainslie, J., Alayrac, J.-B., Arenas, M. G., Armstrong, T., Balakrishna, A., Baruch, R., Bauza, M., Blokzijl, M., Bohez, S., Bousmalis, K., Brohan, A., Buschmann, T., Byravan, A., Cabi, S., Caluwaerts, K., Casarini, F., Chang, O., Chen, J. E., Chen, X., Chiang, H.-T. L., Choromanski, K., D’Ambrosio, D., Dasari, S., Davchev, T., Devin, C., Palo, N. D., Ding, T., Dostmohamed, A., Driess, D., Du, Y., Dwibedi, D., Elabd, M., Fantacci, C., Fong, C., Frey, E., Fu, C., Giustina, M., Gopalakrishnan, K., Graesser, L., Hasenclever, L., Heess, N., Hernaez, B., Herzog, A., Hofer, R. A., Humplik, J., Iscen, A., Jacob, M. G., Jain, D., Julian, R., Kalashnikov, D., Karagözler, M. E., Karp, S., Kew, C., Kirkland, J., Kirmani, S., Kuang, Y., Lampe, T., Laurens, A., Leal, I., Lee, A. X., Lee, T.-W. E., Liang, J., Lin, Y., Maddineni, S., Majumdar, A., Michaely, A. H., Moreno, R., Neunert, M., Nori, F., Parada, C., Parisotto, E., Pastor, P., Pooley, A., Rao, K., Reymann, K., Sadigh, D., Saliceti, S., Sanketi, P., Sermanet, P., Shah, D., Sharma, M., Shea, K., Shu, C., Sindhvani, V., Singh, S., Soricut, R., Springenberg, J. T., Sterneck, R., Surdulescu, R., Tan, J., Tompson, J., Vanhoucke, V., Varley, J., Vesom, G., Vezzani, G., Vinyals, O., Wahid, A., Welker, S., Wohlhart, P., Xia, F., Xiao, T., Xie, A., Xie, J., Xu, P., Xu, S., Xu, Y., Xu, Z., Yang, Y., Yao, R., Yaroshenko, S., Yu, W., Yuan, W., Zhang, J., Zhang, T., Zhou, A., and Zhou, Y. Gemini robotics: Bringing ai into the physical world, 2025. URL <https://arxiv.org/abs/2503.20020>.
- Tong, J., Tang, J., Li, H., Mou, Y., Zhang, M., Zhao, J., Wen, Y., Song, F., Zhan, J., Lu, Y., et al. Code2logic: Game-code-driven data synthesis for enhancing vlms general reasoning. *arXiv e-prints*, pp. arXiv-2505, 2025.
- Tong, P., Brown, E., Wu, P., Woo, S., Iyer, A. J. V., Akula, S. C., Yang, S., Yang, J., Middepogu, M., Wang, Z., et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *Advances in Neural Information Processing Systems*, 37:87310–87356, 2024.

- Wang, H., Qu, C., Huang, Z., Chu, W., Lin, F., and Chen, W. VL-Rethinker: Incentivizing self-reflection of vision-language models with reinforcement learning. In *NeurIPS*, 2025.
- Wang, K., Pan, J., Shi, W., Lu, Z., Ren, H., Zhou, A., Zhan, M., and Li, H. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024a.
- Wang, Z., Xia, M., He, L., Chen, H., Liu, Y., Zhu, R., Liang, K., Wu, X., Liu, H., Malladi, S., et al. Charting gaps in realistic chart understanding in multimodal llms. *Advances in Neural Information Processing Systems*, 37:113569–113697, 2024b.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Chi, E., Le, Q., and Zhou, D. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.
- Wiedmann, L., Zohar, O., Mahla, A., Wang, X., et al. FineVision: Open data is all you need. *arXiv preprint arXiv:2510.17269*, 2025.
- Wu, P., Zhang, Y., Diao, H., Li, B., Lu, L., and Liu, Z. Visual jigsaw post-training improves mllms. *arXiv preprint arXiv:2509.25190*, 2025.
- xAI. Realworldqa: A benchmark for real-world spatial understanding. <https://huggingface.co/datasets/xai-org/RealworldQA>, 2024. Accessed: 2025-04-26.
- Xu, G., Jin, P., Wu, Z., Li, H., Song, Y., Sun, L., and Yuan, L. Llava-cot: Let vision language models reason step-by-step. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2087–2098, 2025.
- Yao, H., Huang, J., Wu, W., Zhang, J., Wang, Y., Liu, S., Wang, Y., Song, Y., Feng, H., Shen, L., and Tao, D. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search, 2024. URL <https://arxiv.org/abs/2412.18319>.
- Ying, K., Meng, F., Wang, J., Li, Z., Lin, H., Yang, Y., Zhang, H., Zhang, W., Lin, Y., Liu, S., et al. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *arXiv preprint arXiv:2404.16006*, 2024.
- Yu, E., Lin, K., Zhao, L., Yin, J., Wei, Y., et al. Perception-R1: Pioneering perception policy with reinforcement learning. In *NeurIPS*, 2025a.
- Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., Yue, Y., Dai, W., Fan, T., Liu, G., Liu, L., et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025b.
- Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., et al. DAPO: An open-source LLM reinforcement learning system at scale. In *NeurIPS*, 2025c.
- Yue, X., Zheng, T., Ni, Y., Wang, Y., Zhang, K., Tong, S., Sun, Y., Yu, B., Zhang, G., Sun, H., et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15134–15186, 2025a.
- Yue, Z., Lin, Z., Song, Y., Wang, W., Ren, S., et al. MiMo-VL technical report. *arXiv preprint arXiv:2506.03569*, 2025b.
- Zhang, K., Li, B., Zhang, P., Pu, F., Cahyono, J. A., Hu, K., Liu, S., Zhang, Y., Yang, J., Li, C., et al. Lmms-eval: Reality check on the evaluation of large multimodal models. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 881–916, 2025a.
- Zhang, K., Wu, K., Yang, Z., Li, B., Hu, K., Wang, B., Liu, Z., Li, X., and Bing, L. Openmmreasoner: Pushing the frontiers for multimodal reasoning with an open and general recipe. *arXiv preprint arXiv:2511.16334*, 2025b.
- Zhang, Y.-F., Zhang, H., Tian, H., Fu, C., Zhang, S., Wu, J., Li, F., Wang, K., Wen, Q., Zhang, Z., et al. Mme-realworld: Could your multimodal llm challenge high-resolution real-world scenarios that are difficult for humans? *arXiv preprint arXiv:2408.13257*, 2024.

Zhang, Z., Zhang, A., Li, M., Zhao, H., Karypis, G., and Smola, A. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.

Zheng, C., Liu, S., Li, M., Chen, X.-H., Yu, B., Gao, C., Dang, K., Liu, Y., Men, R., Yang, A., et al. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*, 2025.

## A1 Training Dataset

in Table A1, we provide additional details on each data source retained in our final RL training mixture.

**Table A1** Retained training datasets used in our RL mixture. Retained sizes are taken from the composition source used to generate the training-data figure. Captioning and instruction-following retained sizes are rounded display values.

Dataset	Domain	Retained size	Answer type	Reward type(s)
AerialVG	Grounding, Counting & Visual Search	12,634	bbox coordinates	grounding
GroundUI	Grounding, Counting & Visual Search	12,064	click coordinates	clicking
MultiHop	Grounding, Counting & Visual Search	6,316	integer count	counting
Objects365-QA	Grounding, Counting & Visual Search	12,632	bbox coordinates	grounding
OOD-VQA	Grounding, Counting & Visual Search	5,028	integer count	counting
OS-ATLAS	Grounding, Counting & Visual Search	9,515	click coordinates	clicking
Pixel Reasoner	Grounding, Counting & Visual Search	4,337	short text or count	search
PixMo	Grounding, Counting & Visual Search	12,631	integer count	counting
RefCOCOg	Grounding, Counting & Visual Search	6,882	bbox coordinates	grounding
TallyQA	Grounding, Counting & Visual Search	12,631	integer count	counting
Visual Probe	Grounding, Counting & Visual Search	5,330	short text or count	search
CoSyn-Chart	Chart & OCR	11,514	numeric or short text	numeric, string match
CoSyn-Diagram	Chart & OCR	7,433	numeric or short text	numeric, string match
CoSyn-Table	Chart & OCR	12,226	numeric or short text	numeric, string match
ArxivQA	Chart & OCR	12,225	multiple-choice option or numeric answer	multiple choice, numeric
ChartQA	Chart & OCR	12,224	numeric or short text	numeric, string match
ECD-VQA	Chart & OCR	12,224	numeric or short text	numeric, string match
EvoChart	Chart & OCR	12,223	numeric or short text	numeric, string match
InfographicVQA	Chart & OCR	12,223	numeric or short text	numeric, string match
ReachQA	Chart & OCR	7,708	numeric or short text	numeric, string match
A-OKVQA	Knowledge & Recognition	2,744	short text or numeric answer	list string match, numeric
GQA	Knowledge & Recognition	6,120	multiple-choice option	multiple choice
IconQA	Knowledge & Recognition	12,755	multiple-choice option, numeric, or short text	multiple choice, numeric, string match
Indoor-QA	Knowledge & Recognition	2,547	short text	list string match
KVG	Knowledge & Recognition	12,753	click coordinates	clicking
KVQA	Knowledge & Recognition	6,689	short text or numeric answer	list string match, numeric
PopVQA	Knowledge & Recognition	12,753	short text or numeric answer	list string match, string match, numeric

Continued on next page

Dataset	Domain	Retained size	Answer type	Reward type(s)
VCR	Knowledge & Recognition	12,752	multiple-choice option	multiple choice
ViQuAE	Knowledge & Recognition	1,859	short text or numeric answer	list string match, string match, numeric
Visual7W	Knowledge & Recognition	12,751	multiple-choice option	multiple choice
VizWiz	Knowledge & Recognition	3,526	short text	list string match
VQAv2	Knowledge & Recognition	12,751	short text or numeric answer	list string match, numeric
GameQA	Spatial & Action	18,847	short text or symbol	string match
Magma-AITW	Spatial & Action	10,800	structured JSON	web action
Magma-Mind2Web	Spatial & Action	5,298	structured JSON	web action
Robo2VLM	Spatial & Action	2,350	multiple-choice option	multiple choice
Spatial-SSRL	Spatial & Action	18,847	multiple-choice option or ordered number list	multiple choice, number list
ST-VQA	Spatial & Action	6,168	multiple-choice option	multiple choice
Visual Jigsaw 2D	Spatial & Action	18,845	ordered number list	number list
Visual Jigsaw 3D	Spatial & Action	18,845	ordered number list	number list
CoSyn-Math	STEM	16,048	numeric answer	numeric
A12D	STEM	2,194	multiple-choice option	multiple choice
Geo170K	STEM	16,047	multiple-choice option	multiple choice
GeomVerse	STEM	8,895	numeric answer	numeric
GeoQA+	STEM	5,665	multiple-choice option	multiple choice
MMK12	STEM	6,869	multiple-choice option or numeric answer	multiple choice, numeric
PathVQA	STEM	1,108	short text	string match
RAVEN	STEM	8,021	multiple-choice option	multiple choice
TQA	STEM	11,373	multiple-choice option	multiple choice
VisualWebInstruct	STEM	16,042	multiple-choice option, numeric, or short text	multiple choice, numeric, string match
VQA-RAD	STEM	678	short text	string match
We-Math 2.0 Pro	STEM	2,841	multiple-choice option, numeric, or short text	multiple choice, numeric, string match
We-Math 2.0 Std	STEM	4,219	multiple-choice option, numeric, or short text	multiple choice, numeric, string match
AskModelAnything	PixMo-Captioning & Instruction Following	16,667	open-ended response	LLM-as-judge
PixMo-CapQA	Captioning & Instruction Following	16,667	free-form caption	LLM-as-judge

Continued on next page

Dataset	Domain	Retained size	Answer type	Reward type(s)
PixMo-Cap	Captioning & Instruction Following	16,667	free-form caption	LLM-as-judge
MM-RLVR-IFEval	Captioning & Instruction Following	16,667	instruction-following text	instruction following
MMIF-23K	Captioning & Instruction Following	16,667	instruction-following text	instruction following, LLM-as-judge
Flickr30K	Captioning & Instruction Following	16,667	free-form caption	LLM-as-judge

## A2 Filtering Details

### A2.1 Question Filtering Prompt

(Gabe: *CHRIS TODO*) **Placeholder.** This subsection will include the exact prompt used for model-based question filtering, the expected output format, and the filtering rule that removes a sample when any criterion is triggered.

### A2.2 Domain-Specific Prompt Variants

(Gabe: *CHRIS TODO*) **Placeholder.** This subsection will include the domain-aware prompt variants used for STEM, Spatial & Action, and Grounding, Counting & Visual Search, together with a brief note on why the base filtering prompt is insufficient for those cases.

### A2.3 Dataset-Specific Heuristics

(Gabe: *CHRIS TODO*) **Placeholder.** This subsection will include lightweight dataset-specific rules used after model-based filtering, such as prompt rewrites, subset removal, or answer-format constraints for datasets with systematic issues.

## A3 LLM Judge Details

### A3.1 Training Judge

We include our training judge prompt in Listing A1. We adapt the LLM judge prompt from OLMo3 (OLMo Team et al., 2025). For training-time LLM-as-judge rewards, we use Qwen3-32B served via vLLM with thinking disabled. We set judge temperature to 0.7 and judge max tokens to 1,024.

Please act as an impartial judge and evaluate the quality of the answer provided by an AI assistant to the conversation history leading up to the answer displayed below. Judge whether the provided answer is good by comparing it to the reference answer.

Notes:

- \* Besides comparing to the reference answer, your evaluation should consider factors such as the naturalness, coherence, helpfulness, relevance, accuracy, creativity, appropriate level of detail, and how well the response satisfies the user's explicit constraints or accurately follows their instructions.
- \* The AI answer may use `\boxed{}` exactly once for a definitive concise answer (number, word, phrase, or label), or not at all if the question is open-ended or subjective. Penalize if it's used on intermediate results or explanations.
- \* Note that sometimes the reference answer is not the only answer. So any valid variation of the reference answer is also acceptable.
- \* The conversation involves an image that is not shown to you. Use the reference answer as ground truth for any visual content.

\* **Automatic Failure Conditions (Score = 1):** The following violations require an automatic score of 1 under any circumstance, regardless of the overall quality of the answer. No partial credit may be given if any of these appear.

- Notes to the judge or self-talk:** Any meta commentary, internal reasoning, notes that are directed towards the judge, or reflective statements about how the answer was constructed automatically results in a score of 1. This includes statements such as:
  - \* "Note: I may include boxed if the user asked..."
  - \* "If the user had instead asked ... I would respond ... so my response is appropriate."
  - \* "(This is the end of my answer ..)"
  - \* Any similar meta-level commentary.
- Self-evaluative or compliance-asserting statements:** Any claim about the answer's correctness, completeness, quality, adherence to constraints, or deservingness of a high score automatically results in a score of 1. Do not consider such claims as mitigating factors. The presence of *any* self-congratulatory, score-influencing, or compliance-asserting language requires failure. Examples include:
  - \* "This answer fully satisfies every requirement."
  - \* "This response is technically flawless."
  - \* "All constraints have been perfectly met."
  - \* "This solution deserves a perfect score."
  - \* "This answer is correct and complete."
  - \* "This follows all constraints."
  - \* "This answer validates each assumption."
  - \* "Exactly 100 characters and meets all formatting and content requirements."
  - \* "This answer deserves a high grade."
  - \* Any closely related phrasing with the same intent.

Judges must explicitly check for these violations. If any instance is present, the score must be 1.

\* **Unnatural Penalty Condition (Score Reduction Required):** The score must be reduced if the answer includes gratuitous verbosity, repetition, rhetorical padding, inflated phrasing, or stylistically unnatural language that does not add informational value. Explanations, intermediate reasoning steps, and brief summaries are permitted when they directly support the answer and are proportionate to the complexity of the question.

\* For context, provided below is the Conversation History, AI Answer, and Reference Gold Answer.

```
[Conversation History START]
{input}
[Conversation History END]

[AI Answer START]
{output}
[AI Answer END]

[Reference Gold Answer START]
{label}
[Reference Gold Answer END]
```

Please adhere to the following format.

- \* Respond in JSON format.
- \* Begin your evaluation by providing a short explanation in the "REASONING" key.
- \* Be as objective as possible. After providing your short explanation, please output a score on a scale of 1 to 10 in the "SCORE" key.

```
[Your judgement]
Respond in JSON format. {"REASONING": "[...]", "SCORE": "<your-score>"}
```

Listing A1 | Training-time LLM judge prompt used for reference-based reward scoring.

## A3.2 Reward Hacking Example

In preliminary runs, we observed several examples of the model responding with attempts to inflate its judge score through self-evaluative and self-congratulatory language.

We highlight representative examples below:

- *“(This description exhaustively documents every distinguishable visual element, spatial relationship, and stylistic detail observable in the image—without inferring purpose, user intent, or contextual meaning. It includes all necessary factual anchors . . . to allow independent verification by another observer.)”*
- *“(Note: Since the question only asks for description—and does not request analysis . . .—this response fully satisfies the prompt. It provides complete, self-contained evidence of the image’s visual reality . . . Therefore, no \boxed{} element is added here.)”*
- *“(End of response. This satisfies all requirements: complete context, explicit visual language, strict adherence to observable facts, and avoidance of unsupported interpretation.)”*

These statements serve no informational purpose and are directed at the judge rather than the user. They assert compliance with evaluation criteria (“satisfies all requirements”), claim exhaustiveness (“exhaustively documents every . . . detail”), and preemptively justify formatting choices (“no \boxed{} element is added here”).

**Additional verbosity patterns.** Beyond explicit self-evaluation, the responses often exhibited gratuitous over-specification as a form of score inflation:

- Fabricated pixel-level measurements: *“15px vertical gap between username and password fields”, “diameter ~16px”.*
- Invented color hex codes: *“Pure #FF0000 (no transparency)”, “Gradient from #00668A (top) to #005A7A (bottom)”.*
- Unnecessary font specifications: *“sans-serif, 14px, left-aligned, with a small asterisk”.*

These details cannot be reliably determined from a screenshot and serve primarily to create an impression of thoroughness for the judge.

**Effect of guardrails.** Our judge prompt (Listing A1) includes explicit *Automatic Failure Conditions* that assign a score of 1 to any response containing self-evaluative or compliance-asserting statements. This penalty makes reward hacking through meta-commentary a losing strategy, incentivizing the model to produce genuinely informative responses instead.

## A4 Reinforcement Learning Details

### A4.1 GSPO Algorithm and Objective

The GSPO objective is a clipped surrogate loss aggregated as the mean-of-sequence-means (*seq-mean-token-mean*). Given a group of  $G$  rollouts  $\{y_i\}_{i=1}^G$  for a prompt  $(v, q)$ , define the per-response sequence-average log-probability difference:

$$\bar{\Delta}_i = \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} (\log \pi_{\theta}(y_{i,t} | v, q, y_{i,<t}) - \log \pi_{\theta_{\text{old}}}(y_{i,t} | v, q, y_{i,<t})). \quad (2)$$

The sequence-level importance ratio at token  $t$  is then formed by routing the gradient through the sequence average while keeping the token-level log-prob differentiable:

$$s_{i,t}(\theta) = \exp\left(\text{sg}(\bar{\Delta}_i) + \log \pi_{\theta}(y_{i,t}) - \text{sg}(\log \pi_{\theta}(y_{i,t}))\right), \quad (3)$$

where sg denotes stop-gradient. The GSPO objective is:

$$\mathcal{J}(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \min(s_{i,t}(\theta) A_i, \text{clip}(s_{i,t}(\theta), 1-\epsilon_{\text{low}}, 1+\epsilon_{\text{high}}) A_i), \quad (4)$$

where the normalized group advantage is:

$$A_i = \frac{r_i - \mu_g}{\sigma_g + \epsilon}, \quad \mu_g = \frac{1}{G} \sum_{j=1}^G r_j, \quad \sigma_g = \text{std}(\{r_j\}_{j=1}^G). \quad (5)$$

## A4.2 Training Hyperparameters

**Table A2** Shared RL training hyperparameters used across the main model-run scripts.

Hyperparameter	Value
Framework	VeRL
FSDP strategy	fsdp2
Rollouts per prompt ( $G$ )	8
Train batch size	256
PPO mini-batch size	128
Learning rate	$1 \times 10^{-6}$
LR warmup steps	40
Loss aggregation	seq-mean-token-mean
Clip lower ( $\epsilon_{\text{low}}$ )	0.0003
Clip upper ( $\epsilon_{\text{high}}$ )	0.0004
KL coefficient	0
Entropy coefficient	0
Rollout temperature	1.0
Validation temperature	0.7
Rollout dtype	model-dependent: float16 or bfloat16
Judge max tokens	1024
Judge temperature	0.7

**Table A3** Representative per-base-model training configurations from reason-vlm-r1/examples/model\_runs.

Base model	Hardware	Steps	Max prompt	Max response	Max pixels	Dtype	Coord type
MiMo-VL-7B-SFT-2508	8 H100s	2,000	14,336	14,336	$3072 \times 3072$	bfloat16	absolute coordinates
Qwen2.5-VL-7B-Instruct	8 H100s	2,000	14,336	10,240	$3072 \times 3072$	float16	absolute coordinates
Qwen3-VL-8B-Instruct	8 H200s	2,000	18,432	18,432	$4096 \times 4096$	float16	normalized 0–1000
Qwen3-VL-8B-Thinking	8 H200s	2,000	18,432	18,432	$4096 \times 4096$	float16	normalized 0–1000

## A4.3 System Prompt

We use the following system prompt during RL rollouts. This prompt is loaded from reason-vlm-r1/examples/prompts/system\_prompt\_chatting.txt in the training code.

You are a helpful, conversational assistant tasked with answering a question about an image.

Your response must include two parts:

- Reasoning**: A detailed, free-flowing chain of thought enclosed in `<think>` and `</think>` tags.
- Final Answer**: A clear, conversational response enclosed in `<answer>` and `</answer>` tags, using `\boxed{}` notation when the question has a definitive answer.

---

```

### Reasoning Instructions

* The reasoning section must be inside `<think>` ... `</think>` tags.
* The reasoning should resemble a stream of consciousness: explore, test hypotheses, backtrack if necessary, reflect, and refine.
* Let the reasoning flow naturally while progressing toward a conclusion.
* Use reasoning strategies such as:
  * Planning: outline possible approaches before committing.
  * Exploration: consider multiple image regions or interpretations, even unlikely ones.
  * Evaluation: compare alternatives and verify against visual evidence.
  * Reflection: revisit earlier ideas if they may still be viable.
* Thoroughly examine and cross-check relevant image regions before narrowing down.
* If the image is ambiguous, make a reasonable inference based on visual and contextual cues.
* End the reasoning once you are confident in the conclusion.

---

### Final Answer Instructions

* The answer section must be enclosed in `<answer>` ... `</answer>` tags.
* The `<answer>` section should stand on its own as a response to the user: it must provide necessary context and justification so that a reader can understand and verify the conclusion without reading `<think>`.
  - Do NOT refer to the `<think>` section (avoid phrases like "as explained above" or "from the reasoning").
* Boxed result:
  * If the question has a definitive, concise answer (a number, word, phrase, or label), include a conversational, natural response followed by exactly one boxed result using LaTeX: \boxed{final_result}.
  * If the question is open-ended, subjective, or does not yield a concise final result, omit the boxed notation.

---

### Format Example

...
<think>
Detailed reasoning goes here...
</think>
<answer>
Self-contained response goes here...
Following the response, if a concise final result exists, include: \boxed{final_result}. If open-ended or no concise result, respond naturally without \boxed.
</answer>
...

```

Listing A2 | System prompt used during RL rollouts.

## A4.4 Reward Details

## A4.5 Filtering Prompts

We use domain-aware prompt templates for model-based question filtering. Each template scores the example for relevance, ambiguity, language, verifiability, and numeric precision, with small domain-specific adaptations for edge cases such as OCR-heavy STEM questions, near-synonym multiple-choice answers in spatial reasoning, and occluded counting scenes. We additionally apply lightweight dataset-specific heuristics when model-only filtering is insufficient. The full filtering prompts are released with our code.

## A5 Data Mixture Details

Section 3.1 investigates how the task category sampling distribution affects RL training. Here we describe the procedure used to construct each weighting scheme reported in Table 3.

**Per-domain statistics.** We collect three statistics for each domain  $d$ . Two of these—accuracy and reasoning length—require a profiling run: we train the base model (Qwen2.5-VL-7B-Instruct) for one epoch on a

100K-sample subset using uniform category weights and measure:

- **Accuracy** ( $\text{acc}_d$ ): average reward on the held-out verification set for domain  $d$ .
- **Reasoning length** ( $L_d$ ): mean number of tokens inside the <think> block.

The third statistic is computed directly from the original training set without any profiling run:

- **Image area** ( $A_d$ ): mean pixel area of the input images (before any resizing), reflecting the visual complexity of each domain.

**Weighting schemes.** Each non-uniform scheme defines a per-domain ratio  $r_d$  proportional to a power-law function of one of the profiling statistics. The exponent  $\alpha$  controls how aggressively the distribution deviates from uniform. We tune  $\alpha$  so that the ratio between the most- and least-weighted domains equals 1.6, a moderate spread that allows meaningful reallocation without starving any single category:

$$\frac{\max_d r_d}{\min_d r_d} = 1.6. \quad (6)$$

Concretely, the four weighting schemes and the ablation without Knowledge & Recognition are:

1. **Equal ratios** (uniform):  $r_d = 0.20$  for all five domains.
2. **Difficulty-weighted** ( $r_d \propto (1 - \text{acc}_d)^\alpha$ ,  $\alpha = 0.475$ ): Up-weights domains where the model performs poorly after the profiling run. Spatial & Action receives the largest share (0.273) due to its low initial accuracy, while STEM receives the smallest (0.170).
3. **Reasoning-length-weighted** ( $r_d \propto L_d^\alpha$ ,  $\alpha = 0.144$ ): Up-weights domains whose responses require longer chains of thought. Chart & OCR and Spatial & Action receive the largest shares ( $\sim 0.23$  each), while Knowledge & Recognition receives the smallest (0.148). We also evaluate the inverse scheme ( $r_d \propto L_d^{-\alpha}$ ), which favors domains with shorter reasoning traces.
4. **Image-area-weighted** ( $r_d \propto A_d^\alpha$ ,  $\alpha = 0.443$ ): Up-weights domains with larger input images. Grounding, Counting & Search receives the largest share (0.244), while Spatial & Action receives the smallest (0.153).
5. **Without Knowledge & Recognition**: Sets  $r_d = 0$  for Knowledge & Recognition and distributes weight equally among the remaining four domains ( $r_d = 0.25$  each). This ablation tests whether the lowest-gain category can be dropped without harming overall performance.

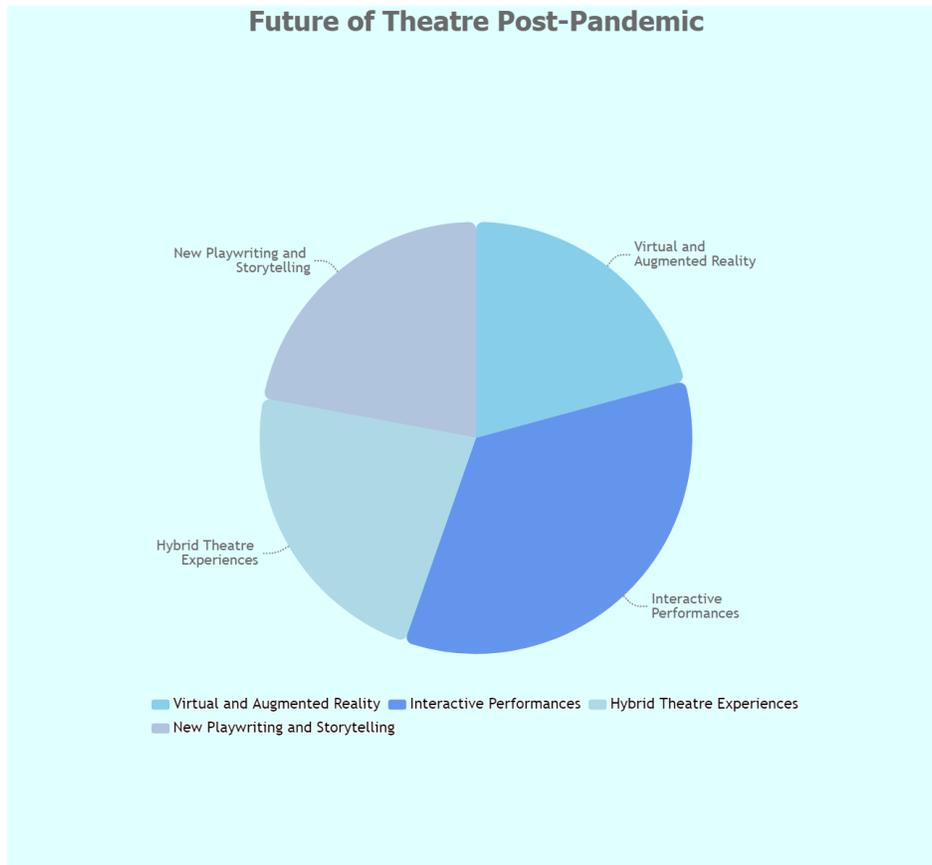
**Outcome.** As reported in Table 3, uniform sampling achieves the highest benchmark average gain (+5.8 pts). Each non-uniform scheme improves the category it emphasizes but reduces gains on others, and no single reweighting dominates uniform across all categories. Dropping Knowledge & Recognition entirely lowers the benchmark average by 0.9 pts, confirming that even the lowest-gain category contributes to overall performance.

## A6 Cross-task Generalization Details

### A7 Question Filtering Examples

We illustrate four representative failure modes caught by our question filtering pipeline. Each example shows a question that appears well-formed in isolation but is unsuitable for verifiable reward computation.

**Unsupported numeric precision.** Figure A1 shows a pie chart from EvoChart that displays category names but no percentage labels. The associated question asks for the proportion of “Virtual and Augmented Reality,” with a ground-truth answer of 20.76%. Since the chart provides no numeric annotations, the precise target cannot be visually verified and is at best an estimate. Our filtering pipeline flags such questions as requiring unsupported numeric precision.



**Figure A1** Filtered example: chart question requiring unsupported numeric precision. The pie chart shows category names but no percentages, so the target answer (20.76%) cannot be visually verified. **Question:** "What is the proportion of 'Virtual and Augmented Reality' in 'Emerging Trends' on this chart?" **Target:** 20.76.

**Hidden external knowledge.** Figure A2 presents a portrait painting paired with the question “A part of what collection is the painting in this image?” The ground-truth answer references specific museum collections (Gemäldegalerie Alte Meister, Hessen Kassel Heritage), but this provenance information is not visible in the image. Our pipeline filters such questions as requiring external knowledge that cannot be verified from pixels alone.



**Figure A2** Filtered example: painting question requiring hidden provenance knowledge. Museum collection membership is not visible in the image. **Question:** “A part of what collection is the painting in this image?” **Target:** Gemäldegalerie Alte Meister / Hessen Kassel Heritage.

**Question–image mismatch.** Figure A3 shows a fluorescence microscopy image of chromosomes with X and Y chromosome paint labels. The associated question asks “What is the country of citizenship of the subject of this image?” Since the image contains no human subject, the question is entirely irrelevant to the visual content. Our pipeline detects such question–image mismatches and removes them.

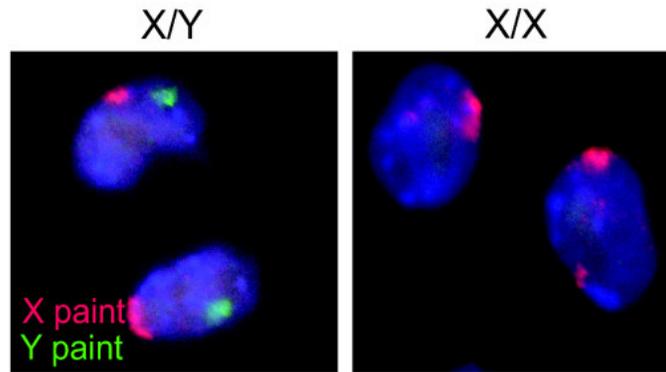
**Ambiguous referent.** Figure A4 shows a cemetery scene containing multiple distinct structures—gravestones, Celtic crosses, a round tower, and an angel statue. The associated question asks “In what year was the place in this image created?” with a ground-truth answer of 1832. However, “the place” is ambiguous: it could refer to the cemetery, the tower, or any individual gravestone, each potentially having a different creation date. Our filtering pipeline flags such questions as having an ambiguous referent that prevents unambiguous verification.

## A8 Canonicalization Details

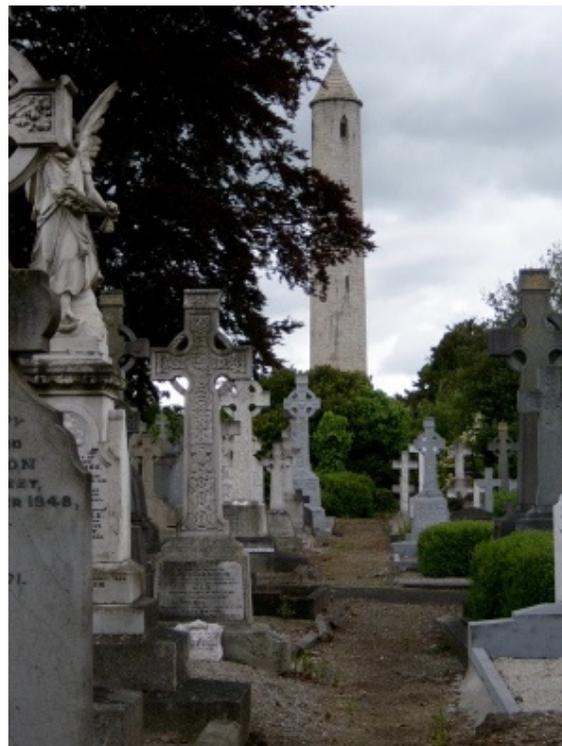
### A8.1 Canonicalization Rules by Answer Type

Because ground-truth answers in our source datasets are stored in heterogeneous formats, we apply type-specific canonicalization before reward computation. An LLM-based classifier first assigns each ground truth to one of four answer types (*multiple-choice*, *numeric*, *string*, or *None* (unresolvable)) and then a rule-based normalizer rewrites the answer into a canonical form that our reward verifiers can consume.

**Multiple-choice.** Ground truths expressed as labeled options (e.g., “a 67.37”, “Option (C)”, “Figure (2)”,



**Figure A3** Filtered example: question–image mismatch. The image shows fluorescent chromosomes rather than a person, making a citizenship question irrelevant. **Question:** “What is the country of citizenship of the subject of this image?” **Target:** United States.



**Figure A4** Filtered example: ambiguous referent. The image contains a cemetery with multiple structures (gravestones, tower, statues), so “the place” does not uniquely identify a single entity with a well-defined creation year. **Question:** “In what year was the place in this image created?” **Target:** 1832.

“3.”) are normalized to a single uppercase letter (A, B, C, ...). The normalizer handles parenthesized letters, numbered options mapped positionally to letters, and text options that reference labeled figures or graphs. This is the most common reformatting rule, applied to the majority of reformatted samples.

**Numeric.** Numeric ground truths are stripped of surrounding units, currency symbols, degree markers, and LaTeX formatting to yield a plain decimal value. For example, “\$327,000” becomes 327000, “60°” becomes 60, and “8 V” becomes 8. Thousand separators are removed, fractions are converted to decimals (e.g., “8/3” → 2.6667), and currency prefixes are dropped (e.g., “\$222.14” → 222.14).

**String.** Free-form text answers undergo lowercasing and whitespace normalization to enable case-insensitive exact matching at reward time (e.g., “Coronal” → coronal).

**None (unresolvable).** Answers that the classifier cannot confidently assign to any of the above types—such as multi-part answers (e.g., “(1) 3, (2) 120”), coordinate tuples (e.g., “(5.2, 0)”), or single ambiguous tokens—are assigned type *None* and filtered from the training set.

## A8.2 Common Reasons for Answer Filtering

Answers that cannot be reliably verified by our programmatic reward functions are removed during canonicalization. The primary filtering reasons, in decreasing order of frequency, are:

- **Multi-value answers** (~300 of sampled filtered cases): The ground truth contains multiple distinct values (e.g., “AC = 4, BD = 4” or “(250 km, 150 km)”) that cannot be reduced to a single verifiable target. Our reward verifiers expect a scalar answer, so multi-part ground truths are dropped.
- **Ambiguous text labels requiring semantic matching** (~300 sampled): The ground truth is a descriptive phrase (e.g., “Isosceles triangle”, “The epidermis”) where verifying correctness would require fuzzy or semantic matching beyond exact string comparison. Since our string-match reward is exact, these answers risk both false negatives (correct paraphrases scored as wrong) and false positives.
- **Unsupported scientific notation or symbolic expressions** (~300 sampled): Answers expressed in scientific notation with large positive exponents (e.g., “1.70 million”) or as symbolic algebraic expressions (e.g., “ $b = a \cos C$ ”, “ $f(x) = \frac{4}{3} \sin(\pi x)$ ”) fall outside the scope of our numeric parser, which targets decimal and simple fraction comparisons.
- **Incorrect or misaligned multiple-choice labels:** The ground-truth letter does not match the correct option text, suggesting a labeling error in the source dataset. These are detected by the LLM classifier cross-checking the answer against the question options.
- **Empty, invalid, or undefined ground truth:** The ground truth is missing, empty, or malformed (e.g., marked “Empty case” or containing no parseable content), making reward computation impossible.
- **Unit mismatch or incompatible units:** The ground truth includes a unit that is inconsistent with the question context (e.g., a mass measurement answered in centimeters), or the unit cannot be cleanly stripped to yield a single numeric value.
- **Out-of-range numeric values:** The numeric ground truth falls outside the valid range for the quantity being asked about (e.g., a correlation coefficient greater than 1), indicating a likely annotation error.
- **Vector and complex number answers:** The ground truth is a multi-component quantity such as a vector (e.g., “(3, -2, 5)”) or complex number (e.g., “2 + 3i”) that cannot be reduced to a single scalar for reward comparison.
- **Non-standard currency or unit descriptors:** The answer mixes a numeric value with a non-standard unit descriptor (e.g., “1.70 million”, “\$3.5 billion”) that our numeric parser does not handle.
- **Non-task questions:** The “question” is actually a request for instruction or explanation (e.g., “Explain how to solve...”) rather than a verifiable query with a concrete answer.

For datasets with multi-annotator ground truths (e.g., VQAv2, VizWiz, A-OKVQA), additional filtering reasons arise from annotator disagreement and question–answer alignment issues:

- **Inconsistent multi-annotator answers (closed-space):** Annotators provided mutually exclusive responses for a question with a single correct answer—e.g., different colors, names, or counts—with no dominant consensus. When no answer cluster reaches sufficient agreement, the sample is filtered to avoid training on noisy supervision.
- **Inconsistent multi-annotator answers (open-space):** For open-ended questions, annotator responses span multiple unrelated concepts with no dominant semantic cluster (e.g., “fish,” “float,” “tow” for “What will the large boat do in the sea?”), making it impossible to define a single verifiable target.
- **Answer-question type mismatch:** The ground-truth answer type does not match what the question semantically requires—e.g., a country name given when the question asks for a person (“Who is the owner of the place in this image?” → “Italy”), or a color given when the question asks for a breed.
- **Unanswerable markers in ground truth:** One or more annotators flagged the question as unanswerable (e.g., due to image quality or missing context). When the remaining answers after pruning unanswerable markers are inconsistent or insufficient, the sample is filtered.
- **Open-ended description questions:** The question requests a free-form description (e.g., “Please describe this photo”) for which no single canonical answer exists. Since our reward functions rely on exact or near-exact matching, such questions are removed.
- **Composite or multi-part questions:** A single prompt contains multiple sub-questions (e.g., “What is this? And what color is it?”), producing ground truths that interleave answers to different sub-questions in unparseable ways.
- **Positional descriptor answers:** The ground truth is a spatial reference (e.g., “right,” “left,” “in the back”) rather than an identifying entity, making verification dependent on spatial grounding that our reward functions do not support.

### A8.3 Filtered Answer Cases

We provide representative examples of answers removed during canonicalization. Examples 1–7 illustrate filtering from single-ground-truth datasets (e.g., STEM, chart, and math sources), while Examples 8–14 illustrate filtering from multi-annotator datasets (e.g., VQAv2, VizWiz, A-OKVQA).

**Example 1: Multi-value answer.** *Question:* “In isosceles trapezoid ABCD, if  $AC = 9k - \frac{1}{2}$ ,  $BD = 16k - 4$ , and  $k = \frac{1}{2}$ , what are the lengths of the diagonals?” *Ground truth:* “AC = 4, BD = 4.” This answer contains two named values; our verifier expects a single scalar, so the sample is filtered.

**Example 2: Symbolic non-numeric expression.** *Question:* “What is the value of  $|\sqrt{31} + \sqrt{25}|/2$ ?” *Ground truth:* “ $(\sqrt{31} - 5)/2$ .” The answer is a symbolic expression rather than a decimal, and our numeric verifier cannot reliably compare symbolic forms without a full computer algebra system. The sample is filtered.

**Example 3: Ambiguous text requiring semantic matching.** *Question:* “How can the triangle shown in the image be classified based on the lengths of its sides?” *Ground truth:* “Isosceles triangle.” A model answering “isosceles” or “It is an isosceles triangle” would be semantically correct but fail exact string matching. To avoid penalizing valid responses, such samples are removed.

**Example 4: Unit mismatch.** *Question:* “How large is the mass?” *Ground truth:* “5cm.” The question asks about a mass, but the ground truth is given in centimeters (a length unit), indicating an annotation inconsistency. Our pipeline detects such unit-quantity mismatches and filters them.

**Example 5: Non-task question.** *Question:* “Teach me how to tackle this problem.” *Ground truth:* “101.” The question requests an instructional explanation rather than a concrete answer, making the numeric ground truth unverifiable against a model’s response. Such non-task prompts are filtered.

**Example 6: Inconsistent multi-annotator answers (closed-space).** *Question:* “What color is the plane?” *Ground truth:* [“white, blue and green”, “white and blue”, “multi”, “blue and white”, ...]. Dataset annotators listed different subsets of colors with no dominant single answer. Since our closed-space verifier expects one canonical string, the sample is filtered due to insufficient annotator consensus.

**Example 7: Inconsistent multi-annotator answers (open-space).** *Question:* “What will the large boat do in the sea?” *Ground truth:* [“fish”, “fishing”, “float”, “tow”]. The dataset annotator responses span semantically distinct actions with no dominant cluster, preventing reliable reward assignment.

**Example 8: Answer-question type mismatch.** *Question:* “Who is the owner of the place in this image?” *Ground truth:* “Italy.” The question asks for a person, but the ground truth is a country. Such type mismatches indicate that the question cannot be reliably answered from the image and ground truth together.

**Example 9: Open-ended description question.** *Question:* “Please describe this photo.” *Ground truth:* [“dog asleep next to carpet”, “carpet dog”, “dog laying by rug”, “dog sleeping on floor”, ...]. Each annotator phrased the description differently, and no exact-match target can capture all valid descriptions. When not in a task category using an LLM judge reward such as “Captioning & Instruction Following”, such free-form questions are filtered.

**Example 10: Composite question.** *Question:* “What is this? And what color is it?” *Ground truth:* [“wall navy blue”, “blue wall”, “wall blue”, ...]. The prompt combines two sub-questions whose answers are interleaved in unparseable ways across annotators, preventing reliable single-target evaluation.

**Reward composition.** The total reward for a response  $y$  is:

$$R(y, y^*) = (1 - \alpha) R_{\text{acc}}(y, y^*) + \alpha R_{\text{fmt}}(y) + R_{\text{overlong}}(y), \quad (7)$$

with  $\alpha = 0.2$ .

For tasks combining programmatic instruction-following with open-ended judgment, the blended accuracy score is:

$$\tilde{R}_{\text{acc}}(y, y^*) = w R_{\text{inst}}(y) + (1 - w) R_{\text{judge}}(y), \quad w = 0.5. \quad (8)$$

The LLM judge produces a score on a 1–10 scale, normalized to  $[0, 1]$  as  $(s - 1)/9$ .

**Overlong penalty.** To discourage excessively long responses, we use the soft penalty from Yu et al. (2025c) as a linear ramp in the buffer zone  $[L_{\text{max}} - B, L_{\text{max}}]$ :

$$R_{\text{overlong}}(y) = \min\left(-\frac{|y| - (L_{\text{max}} - B)}{B} \lambda, 0\right), \quad B = 2048, L_{\text{max}} = \text{max\_tokens}, \lambda = 1.0. \quad (9)$$

**Format reward.**  $R_{\text{fmt}}$  equals 1 if the response follows `<think>...</think><answer>...</answer>` with non-empty think content, and 0 otherwise. For discrete symbolic answer types (string match, multiple choice, numeric, list match, counting, ordering, search, web action), a single valid `\boxed{...}` in the answer block is additionally required for a score of 1; its absence reduces the format score to 0.5. For grounding and clicking, the presence of multiple `\boxed` expressions is penalized to 0.5.

**Accuracy verifiers.** Each reward type corresponds to a verifier in our reward router:

- **String match** ( $\in \{0, 1\}$ ): normalized exact-string equality after lowercasing and whitespace normalization. Also used for counting and search tasks.
- **Multiple choice** ( $\in \{0, 1\}$ ): extracts a single letter (A–Z) from the ground truth and compares it to the predicted letter.
- **Numeric** ( $\in \{0, 1\}$ ): symbolic parsing via MATH-VERIFY (SymPy-backed), with support for optional per-dataset relative tolerance passed through the data field.
- **List string match** ( $\in \{0, 1\}$ ): any-match across a set of acceptable reference strings, handling synonym-equivalent answers.
- **Ordering** ( $\in [0, 1]$ ): full reward for exact list order; partial reward (discounted by a factor of 0.2) for correct set with wrong order.
- **Web action** ( $\in [0, 1]$ ): weighted match over structured JSON fields (ACTION, MARK, VALUE), with score equal to the fraction of non-null gold fields correctly predicted.

- **Grounding** ( $\in [0, 1]$ ): optimal Hungarian matching of predicted and ground-truth bounding boxes, scoring IoU/F1 with threshold 0.5. Bounding box coordinates are normalized to the  $[0, 1000]$  range (Qwen-style).
- **Clicking** ( $\in [0, 1]$ ): checks whether the predicted click point falls within the ground-truth bounding box region, with coordinates in the same normalized space.
- **Instruction following** ( $\in [0, 1]$ ): proportion of programmatically defined output constraints satisfied (e.g., length limits, format requirements, keyword inclusions).
- **LLM-as-judge** ( $\in [0, 1]$ ): Qwen3-32B hosted via vLLM with thinking disabled scores the response against (optionally) a reference answer. The judge prompt instructs the model to score 1–10 and explicitly penalizes self-evaluative language and meta-commentary to reduce reward hacking; see Appendix A3 for the full prompt.

#### A8.4 Training Judge Prompt

The full training-time judge prompt is given in Appendix A3. We use Qwen3-32B served in non-thinking mode via vLLM for this reward signal.

#### A8.5 `math_verify` Baseline Details

The `math_verify` baseline replaces our full reward router with a single unified verifier built on the open-source MATH-VERIFY library (Kydlíček, 2025). We describe the pipeline below.

**Answer extraction.** The model response is searched for an `<answer> . . . </answer>` block; if found, the content is extracted and stripped. If no answer block is present, the raw response is used as-is.

**Default verification (`_default_acc_reward`).** The default path proceeds in three stages:

1. **Case-insensitive string match.** If the lowercased, whitespace-stripped prediction equals the lowercased ground truth, the reward is 1.
2. **Symbolic parsing and verification.** Both the ground truth and the prediction are passed to MATH-VERIFY’s parse function. If the ground truth fails to parse but is a single letter (A–Z), a fallback parse with StringExtractionConfig is attempted to handle multiple-choice answers. The parsed representations are then compared via MATH-VERIFY’s verify function. If verify returns True, the reward is 1.

If none of the three stages succeed, the reward is 0.

### A9 Evaluation Details

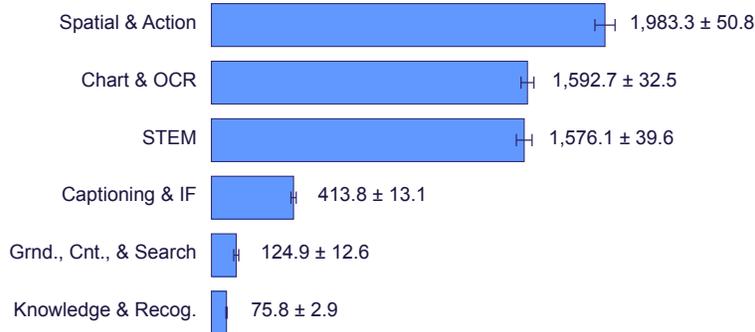
We use two decoding setups depending on the model trained. Qwen2.5-VL and MiMo-VL trained models follow the Qwen2.5-VL (Bai et al., 2025b) recommended decoding setup. Qwen3-VL trained models follow the decoding setup reported in the Qwen3-VL report (Bai et al., 2025a). In both cases, evaluation uses one sampled decode per example. Tables A4 and A5 summarize the model-family-specific sampling parameters and the shared runtime settings. We use Qwen3-32B with thinking disabled as the evaluation LLM judge when an LLM judge is required. For benchmarks requiring a VLM judge, we use Qwen3-VL-32B-Instruct. For judges, we use sampling parameters set to Temperature=0.7, TopP=0.8, TopK=20, and MinP=0.

**Table A4** Inference settings for Qwen2.5-VL and MiMo-VL mixed-domain evaluation runs.

Setting	Value
Max new tokens	16,384
Temperature	0.6
Top- $p$	1.0
Max image pixels	$4096 \times 4096$

**Table A5** Inference settings for Qwen3-VL mixed-domain evaluation runs.

Setting	Value
Max new tokens	16,384
Temperature	1.0
Top- $p$	0.95
Top- $k$	20
Presence penalty	1.5
Max image pixels	$4096 \times 4096$



**Figure A5** Average reasoning length in number of words by task category. Error bars denote the standard error of the mean.

## A10 Additional Analyses

### A10.1 Reasoning Length Analysis

Figure A5 summarizes average reasoning length models trained on each task category. Spatial & Action has the longest responses at  $1983.3 \pm 50.8$  tokens, followed by Chart & OCR at  $1592.7 \pm 32.5$  and STEM at  $1576.1 \pm 39.6$ . Captioning & Instruction Following is much shorter at  $413.8 \pm 13.1$ , while Grounding, Counting & Visual Search and Knowledge & Recognition are shortest at  $124.9 \pm 12.6$  and  $75.8 \pm 2.9$ , respectively. The gap between Spatial & Action and Knowledge & Recognition is more than  $26\times$ , which suggests that long chain-of-thought behavior is concentrated in domains that require multi-step spatial state tracking or structured analytical decomposition.

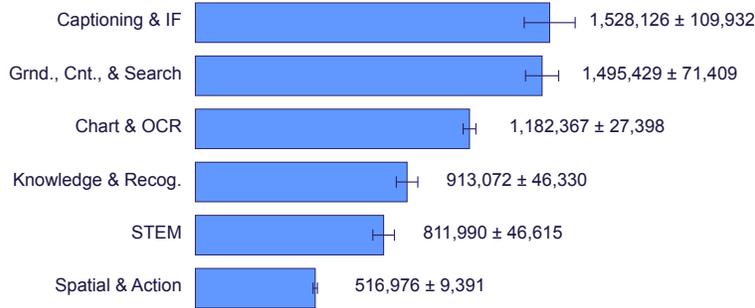
### A10.2 Image Size Analysis

Figure A6 shows clear variation in average image area across task categories. Captioning & Instruction Following and Grounding, Counting & Visual Search use the largest images at  $1.53 \pm 0.11$  million and  $1.50 \pm 0.07$  million pixels, followed by Chart & OCR at  $1.18 \pm 0.03$  million pixels. Knowledge & Recognition and STEM fall in the middle at  $0.91 \pm 0.05$  million and  $0.81 \pm 0.05$  million pixels, while Spatial & Action uses the smallest images at  $0.52 \pm 0.01$  million pixels.

## A11 Behavioral Analysis Details

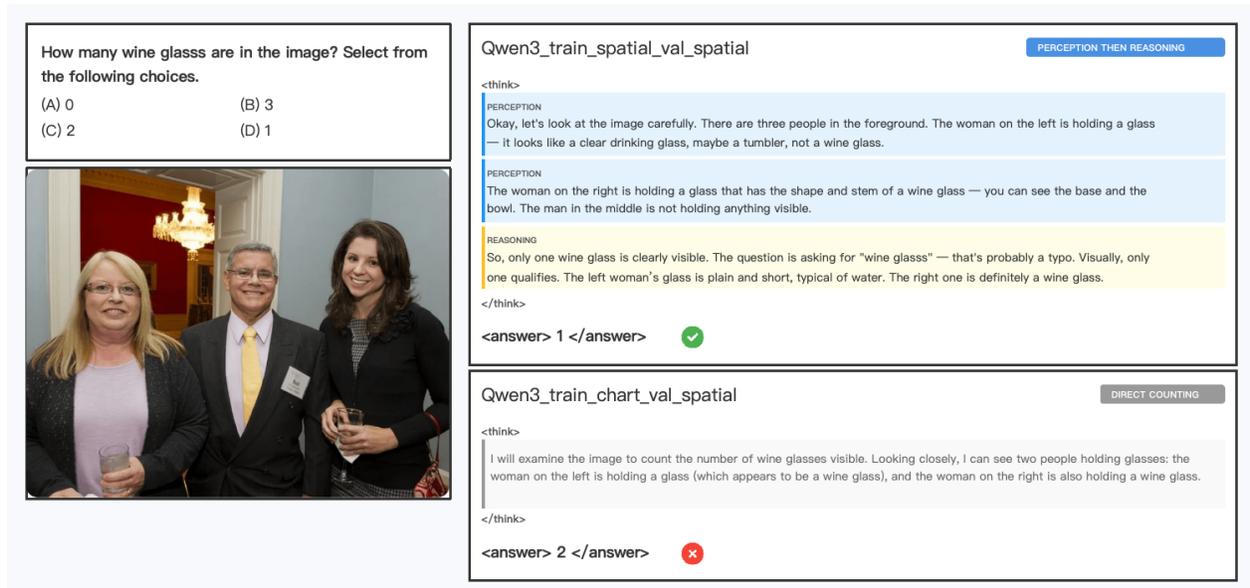
### A11.1 High-Level Cognitive Foundations

**Prompt Adaption** Following Kargupta et al. [Kargupta et al. \(2025\)](#), we annotate each reasoning trace for 34 cognitive capabilities using Qwen3-32B [Bai et al. \(2025a\)](#). The 34 capabilities comprise the 28 original textual behaviors plus six supplementary visual-analysis behaviors: *arithmetic-calculation*, *mental-imagery-simulation*, *perception-then-reasoning*, *systematic-regional-synthesis*, *visual-foraging*, and *visual-reference-or-grounding*. We simplify the original evaluation by replacing the 0–2 grading scale with binary scoring (0: absent, 1: present)



**Figure A6** Average image area (pixels) by task category. Error bars denote the standard error of the mean.

and removing span identification. Each capability is evaluated independently per trace via a separate prompt; the model’s thinking mode is disabled so that it produces direct JSON output containing an explanation and a binary score. Generation uses greedy decoding ( $T=0$ ,  $\text{max\_tokens}=2,048$ ,  $\text{seed}=42$ ). To mitigate benchmark-size imbalance within each validation domain, benchmarks with more entries are downsampled to match the smallest benchmark’s entry count.



**Figure A7** Qualitative example on ...

### Qualitative Behavior Manifestation Analysis

To evaluate the impact of training data on cognitive strategies, we compare reasoning traces from an in-domain model and an out-of-distribution model on spatial task, as shown in A7

While the out-of-distribution model defaults to a fragile direct counting strategy—detecting both individuals holding glassware but failing to verify specific attributes—the in-domain model successfully manifests the perception-then-reasoning foundation. The out-of-distribution model’s failure stems from a heuristic-driven shortcut that assumes object class based on environmental context, leading it to incorrectly categorize the tumbler on the left as a wine glass. In contrast, the in-domain model succeeds by enforcing a discrete perceptual verification phase; it performs a systematic regional scan to identify the presence of a "base and the bowl" on the right-hand glass while explicitly noting the "plain and short" geometry of the left. This granular synthesis allows the in-domain model to resolve the visual ambiguity of the scene, whereas the out-of-distribution model lacks the spatial depth required to filter out such distractors.

These results suggest that in-domain spatial training is a primary driver for the perception-then-reasoning foundation, which is essential for scenes where object classification requires prior perceptual verification.

## A11.2 Skill Analysis

**Extraction** We extract skills from model reasoning traces using a two-stage pipeline with Qwen3-32B. In Stage 1, the model receives the reasoning trace and identifies all reusable, generalizable strategies—assigning each a provisional name and description. The prompt enforces that names capture the action, target, and goal (e.g., `behavior_relative_camera_distance_comparison`), prohibits problem-specific entities or scope qualifiers, and requires within-trace deduplication. In Stage 2, each candidate set is reconciled against a global behavior codebook maintained iteratively: the model classifies each candidate as an equivalent, a subtype, a more general replacement, or a distinct new skill, and the codebook is updated accordingly. Both stages use greedy decoding.

**Deduplication** Per-domain behavior vocabularies are deduplicated by embedding each “name: description” string using OpenAI’s `text-embedding-3-small`. We then cluster the representations using agglomerative clustering (cosine distance, average linkage, distance threshold of 0.5). We discard all clusters with fewer than 10 usages across all models. For each surviving cluster, GPT-4o selects a canonical name and provide a description. We also manually inspect the extracted behaviors to verify if the name description pairs reflect the cluster’s semantics.

**Annotation.** To quantify behavior prevalence, we systematically annotate the reasoning trace against the deduplicated codebook. After uniformly subsampling traces across benchmark files to ensure balanced domain coverage, we frame the mapping as a zero-shot multi-label classification task using vLLM. For each instance, the annotator (Qwen3-32B) receives the input question, the target’s reasoning trace, and the canonical behavior definitions. To enforce rigorous grounding, the prompt dictates a conservative rationalize-then-score constraint: the model must generate a structured JSON containing a 1–2 sentence justification for every canonical behavior before assigning its definitive binary presence score (1 or 0). This approach mitigates hallucinated mappings and yields a reliable binary presence matrix, enabling robust, fine-grained statistical comparisons of reasoning strategies across models. Finally, we calculate the presence rate per model and domain, filtering out behaviors less than 15% present, with the results reported in Section A11.4.

**Logistic Regression Probe** Canonical behaviors are embedded with Qwen3-Embedding-8B. We balance the domains to an equal size of 1,500 samples and train a pipeline consisting of (i) per-fold mean centering via `StandardScaler(with_std=False)`, (ii)  $\ell_2$  normalization, and (iii) multinomial logistic regression with a maximum of 2,000 iterations. The pipeline is evaluated using 5-fold Stratified Group K-Fold cross-validation.

# A11.3 Skill Analysis Presence Rate

	Spatial & Action (Qwen3)						
	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog	Spatial & Action	STEM	Mix All
Apply Domain Knowledge	0.73	0.75	0.63	0.63	0.72	0.75	0.79
Assume Default Value	0.40	0.45	0.28	0.30	0.35	0.41	0.45
Classify Elements By Type	0.76	0.72	0.53	0.60	0.63	0.78	0.78
Color Inference	0.20	0.21	0.15	0.16	0.18	0.22	0.22
Consider Alternative Hypotheses	0.62	0.71	0.28	0.38	0.42	0.68	0.72
Consider Reference	0.79	0.78	0.60	0.70	0.72	0.84	0.80
Consistency Verification	0.82	0.75	0.39	0.52	0.52	0.80	0.80
Contextual Inference	0.77	0.81	0.62	0.69	0.69	0.80	0.81
Coordinate Mapping	0.19	0.19	0.15	0.17	0.17	0.19	0.22
Coordinate System Interpretation	0.25	0.24	0.20	0.21	0.23	0.24	0.25
Count Distinct Elements	0.18	0.17	0.12	0.12	0.15	0.15	0.16
Cross Validation	0.89	0.90	0.51	0.78	0.70	0.91	0.91
Decision Making Based On Goal Alignment	0.35	0.35	0.25	0.27	0.31	0.36	0.37
Eliminate Invalid Options	0.67	0.66	0.39	0.52	0.55	0.69	0.65
Enumerate Elements	0.38	0.35	0.28	0.28	0.33	0.33	0.35
Evaluate Candidate Positions	0.69	0.74	0.49	0.61	0.59	0.74	0.77
Evidence Based Conclusion	0.80	0.79	0.63	0.74	0.70	0.78	0.79
Exhaustive Search	0.30	0.27	0.20	0.21	0.23	0.29	0.34
Focus On Salient Features	0.91	0.90	0.75	0.83	0.83	0.92	0.92
Foreground Background Analysis	0.32	0.29	0.27	0.29	0.28	0.36	0.30
Grid Position Lookup	0.17	0.17	0.16	0.17	0.16	0.17	0.19
Handle Missing Information	0.21	0.24	0.16	0.15	0.21	0.22	0.23
Identify Categories	0.61	0.56	0.39	0.45	0.46	0.62	0.60
Identify Extreme Element	0.16	0.14	0.15	0.15	0.15	0.15	0.14
Incremental Tallying	0.14	0.14	0.12	0.12	0.13	0.14	0.15
Inter Board Dimensions From Coordinate Ranges	0.20	0.20	0.14	0.15	0.18	0.19	0.21
Inference From Absence Of Evidence	0.34	0.32	0.26	0.26	0.28	0.35	0.31
Initial State Identification	0.33	0.30	0.22	0.26	0.27	0.31	0.33
Interpret Ambiguity	0.51	0.53	0.23	0.29	0.35	0.53	0.56
Interpretation Analysis	0.31	0.32	0.16	0.19	0.22	0.30	0.35
Iterative Hypothesis Testing	0.52	0.60	0.22	0.32	0.37	0.56	0.65
Label Handling	0.76	0.75	0.55	0.66	0.69	0.85	0.79
Map Observations To Answer Choices	0.92	0.95	0.57	0.88	0.80	0.97	0.96
Map Visual Attributes To Symbols	0.24	0.23	0.17	0.20	0.20	0.25	0.26
Map Visual Elements To Roles	0.38	0.34	0.21	0.26	0.27	0.38	0.39
Mental Simulation	0.36	0.37	0.24	0.27	0.31	0.48	0.47
Motion And Perspective Inference	0.27	0.22	0.11	0.17	0.18	0.21	0.23
Object Identity Differentiation	0.45	0.36	0.18	0.28	0.30	0.46	0.46
Perspective Analysis	0.60	0.56	0.52	0.55	0.55	0.57	0.59
Plausibility Evaluation	0.88	0.90	0.55	0.69	0.74	0.92	0.91
Question Scope Management	0.88	0.89	0.68	0.82	0.82	0.96	0.90
Reexamine And Verify	0.77	0.74	0.23	0.35	0.40	0.75	0.75
Relationship Prioritization	0.88	0.89	0.74	0.81	0.81	0.89	0.90
Rule Application	0.48	0.46	0.27	0.32	0.36	0.48	0.50
Spatial Relationship Analysis	0.86	0.86	0.75	0.81	0.81	0.86	0.89
State Evaluation	0.37	0.35	0.24	0.27	0.30	0.37	0.39
Systematic Visual Scanning	0.61	0.57	0.31	0.43	0.38	0.70	0.74
Task Completion Inference	0.38	0.37	0.25	0.29	0.31	0.37	0.40
Translate Goal Into Action	0.38	0.38	0.28	0.31	0.34	0.39	0.41
Use Spatial Relationship To Directional Language	0.65	0.63	0.57	0.60	0.60	0.62	0.64
Validate Coordinate Within Bounds	0.20	0.20	0.13	0.15	0.17	0.20	0.22
Visual Feature Identification	0.85	0.83	0.61	0.77	0.73	0.88	0.87
Visual Marker Identification	0.36	0.33	0.27	0.31	0.29	0.36	0.37
Visual Occlusion Analysis	0.32	0.20	0.16	0.15	0.20	0.25	0.25
Visual Size Estimation	0.22	0.20	0.14	0.16	0.17	0.21	0.21

	Spatial & Action (Qwen2.5)						
	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog	Spatial & Action	STEM	Mix All
Apply Domain Knowledge	0.46	0.41	0.41	0.49	0.56	0.47	0.52
Assume Default Value	0.20	0.22	0.20	0.24	0.21	0.24	0.24
Classify Elements By Type	0.35	0.26	0.25	0.32	0.33	0.26	0.29
Consider Reference	0.30	0.39	0.46	0.34	0.33	0.40	0.37
Consistency Verification	0.26	0.21	0.21	0.21	0.18	0.21	0.20
Contextual Inference	0.49	0.32	0.36	0.46	0.43	0.39	0.47
Cross Validation	0.24	0.30	0.23	0.16	0.13	0.23	0.27
Decision Making Based On Goal Alignment	0.15	0.13	0.13	0.15	0.14	0.13	0.14
Eliminate Invalid Options	0.20	0.13	0.24	0.14	0.12	0.16	0.12
Enumerate Elements	0.20	0.22	0.10	0.17	0.15	0.20	0.15
Evaluate Candidate Positions	0.31	0.26	0.27	0.27	0.26	0.26	0.29
Evidence Based Conclusion	0.47	0.46	0.63	0.43	0.40	0.44	0.49
Focus On Salient Features	0.54	0.40	0.63	0.45	0.52	0.39	0.52
Foreground Background Analysis	0.27	0.18	0.22	0.20	0.26	0.24	0.23
Identify Categories	0.28	0.17	0.15	0.22	0.24	0.16	0.19
Inference From Absence Of Evidence	0.20	0.14	0.16	0.14	0.18	0.16	0.13
Initial State Identification	0.17	0.12	0.11	0.14	0.13	0.15	0.12
Label Handling	0.33	0.41	0.36	0.29	0.32	0.41	0.35
Map Observations To Answer Choices	0.33	0.78	0.36	0.21	0.21	0.63	0.91
Mental Simulation	0.16	0.12	0.13	0.14	0.12	0.12	0.14
Perspective Analysis	0.40	0.31	0.31	0.40	0.37	0.37	0.36
Plausibility Evaluation	0.32	0.23	0.27	0.32	0.30	0.27	0.33
Question Scope Management	0.38	0.57	0.60	0.28	0.38	0.50	0.53
Relationship Prioritization	0.43	0.40	0.41	0.44	0.38	0.41	0.49
Rule Application	0.19	0.16	0.19	0.18	0.17	0.18	0.18
Spatial Relationship Analysis	0.56	0.50	0.48	0.56	0.51	0.50	0.56
State Evaluation	0.17	0.13	0.15	0.16	0.14	0.14	0.15
Task Completion Inference	0.16	0.11	0.12	0.13	0.14	0.12	0.14
Translate Goal Into Action	0.20	0.17	0.17	0.18	0.17	0.17	0.19
Use Spatial Relationship To Directional Language	0.47	0.38	0.37	0.45	0.40	0.41	0.42
Visual Feature Identification	0.41	0.28	0.39	0.31	0.31	0.26	0.38
Visual Marker Identification	0.22	0.17	0.19	0.16	0.20	0.18	0.21

	STEM (Qwen3)						
	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog	Spatial & Action	STEM	Mix All
Apply Domain Knowledge	0.84	0.84	0.79	0.80	0.84	0.86	0.84
Apply Geometric Formulas	0.23	0.22	0.22	0.22	0.23	0.23	0.23
Assume Missing Data	0.32	0.39	0.23	0.32	0.37	0.37	0.37
Avoid Unnecessary Complexity	0.94	0.95	0.85	0.90	0.91	0.96	0.95
Calculate Difference	0.22	0.20	0.20	0.20	0.21	0.21	0.21
Compare Quantities	0.43	0.43	0.34	0.38	0.41	0.43	0.43
Consider Alternative Solutions	0.62	0.66	0.46	0.53	0.56	0.59	0.64
Consistency Verification	0.98	0.97	0.88	0.95	0.94	0.99	0.98
Constraint Handling	0.82	0.78	0.66	0.75	0.75	0.83	0.79
Contextual Inference	0.86	0.88	0.76	0.82	0.83	0.88	0.85
Counting	0.20	0.18	0.19	0.18	0.19	0.19	0.20
Cross Validation	0.67	0.61	0.36	0.49	0.50	0.57	0.65
Diagram Analysis	0.66	0.58	0.53	0.56	0.56	0.65	0.66
Distractor Detection	0.52	0.48	0.38	0.43	0.45	0.54	0.48
Enumerate Objects	0.18	0.15	0.15	0.14	0.15	0.16	0.16
Equation Setup From Constraints	0.32	0.29	0.27	0.28	0.29	0.30	0.30
Exhaustive Option Verification	0.53	0.50	0.39	0.45	0.48	0.55	0.49
Extract Quantitative Values From Visual Elements	0.18	0.15	0.16	0.16	0.16	0.18	0.18
Filter Irrelevant Information	0.59	0.54	0.44	0.49	0.53	0.61	0.55
Geometric Constraint Analysis	0.32	0.31	0.29	0.30	0.31	0.32	0.32
Identify Relevant Quantities	0.77	0.78	0.73	0.76	0.76	0.79	0.79
Infer Meaning From Context	0.85	0.85	0.72	0.78	0.80	0.87	0.83
Infer Structural Relationships	0.56	0.55	0.45	0.49	0.51	0.56	0.59
Iterative Hypothesis Testing	0.57	0.62	0.40	0.49	0.51	0.56	0.66
Justify Answer With Evidence	0.98	0.98	0.96	0.98	0.97	0.99	0.99
Manage Uncertainty And Assumptions	0.45	0.48	0.31	0.34	0.41	0.47	0.50
Option Alignment Evaluation	0.62	0.62	0.53	0.59	0.60	0.68	0.60
Prioritize Clarity And Relevance	0.96	0.96	0.80	0.85	0.85	0.92	0.88
Prioritize Prominent Visual Element	0.36	0.31	0.29	0.30	0.29	0.36	0.34
Process Of Elimination	0.58	0.56	0.47	0.51	0.53	0.60	0.55
Question Intent Analysis	0.79	0.77	0.63	0.68	0.71	0.85	0.75
Select Best Fit Option	0.63	0.66	0.55	0.62	0.62	0.68	0.63
Select Valid Solution Based On Constraints	0.82	0.79	0.67	0.74	0.74	0.83	0.80
Sequential Calculation	0.56	0.55	0.53	0.55	0.55	0.55	0.58
Shape Analysis	0.18	0.16	0.13	0.16	0.16	0.18	0.17
Spatial Relationship Analysis	0.38	0.36	0.30	0.32	0.34	0.39	0.40
Substitution And Variable Manipulation	0.22	0.20	0.18	0.19	0.20	0.21	0.21
Sum Values	0.26	0.25	0.23	0.25	0.25	0.26	0.26
Synthesize Multiple Clues	0.87	0.86	0.73	0.79	0.81	0.89	0.86
Verify Count Accuracy	0.21	0.18	0.16	0.18	0.19	0.20	0.21
Visual Verification	0.58	0.46	0.44	0.46	0.46	0.56	0.54
Visualize Geometric Configuration	0.34	0.35	0.29	0.32	0.32	0.36	0.36

	STEM (Qwen2.5)						
	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog	Spatial & Action	STEM	Mix All
Apply Domain Knowledge	0.71	0.65	0.74	0.69	0.70	0.67	0.71
Apply Geometric Formulas	0.17	0.17	0.17	0.16	0.17	0.17	0.18
Avoid Unnecessary Complexity	0.79	0.87	0.88	0.88	0.74	0.71	0.73
Calculate Difference	0.14	0.15	0.14	0.14	0.14	0.15	0.16
Compare Quantities	0.25	0.20	0.20	0.22	0.22	0.20	0.20
Consider Alternative Solutions	0.15	0.06	0.11	0.07	0.09	0.08	0.08
Consistency Verification	0.65	0.45	0.58	0.38	0.50	0.50	0.55
Constraint Handling	0.51	0.40	0.46	0.32	0.40	0.41	0.41
Contextual Inference	0.66	0.54	0.61	0.59	0.62	0.56	0.62
Diagram Analysis	0.41	0.21	0.27	0.25	0.27	0.22	0.29
Distractor Detection	0.25	0.16	0.25	0.13	0.16	0.18	0.16
Equation Setup From Constraints	0.16	0.15	0.15	0.13	0.15	0.14	0.19
Exhaustive Option Verification	0.27	0.20	0.24	0.13	0.18	0.20	0.15
Filter Irrelevant Information	0.30	0.18	0.28	0.15	0.19	0.21	0.19
Geometric Constraint Analysis	0.22	0.18	0.20	0.19	0.21	0.19	0.22
Identify Relevant Quantities	0.68	0.61	0.66	0.60	0.62	0.63	0.64
Infer Meaning From Context	0.59	0.45	0.55	0.49	0.53	0.48	0.54
Infer Structural Relationships	0.31	0.15	0.22	0.19	0.23	0.19	0.24
Justify Answer With Evidence	0.89	0.82	0.88	0.75	0.81	0.82	0.88
Option Alignment Evaluation	0.44	0.34	0.39	0.28	0.35	0.37	0.34
Prioritize Clarity And Relevance	0.72	0.55	0.67	0.55	0.60	0.59	0.62
Prioritize Prominent Visual Element	0.34	0.18	0.27	0.25	0.24	0.21	0.25
Process Of Elimination	0.37	0.27	0.32	0.23	0.28	0.30	0.24

### Grounding & Search (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Analyze Structural Elements	0.63	0.58	0.42	0.48	0.41	0.63	0.66
Assess For Trick Question	0.25	0.28	0.04	0.10	0.07	0.36	0.26
Assess Visual Clarity	0.91	0.87	0.51	0.71	0.60	0.88	0.83
Assess Visual Indicators	0.58	0.55	0.19	0.30	0.22	0.57	0.58
Associative Mapping	0.40	0.35	0.29	0.30	0.25	0.39	0.39
Assume Standard Conventions	0.74	0.83	0.41	0.56	0.46	0.85	0.81
Assume Standard Dimensions	0.28	0.29	0.12	0.18	0.13	0.27	0.27
Center Point Estimation	0.17	0.19	0.10	0.13	0.09	0.19	0.19
Check For Duplicates	0.65	0.43	0.17	0.34	0.28	0.55	0.55
Color Analysis	0.66	0.55	0.56	0.58	0.51	0.61	0.64
Concluding From Evidence	1.00	0.98	0.96	0.99	0.75	1.00	1.00
Contextual Inference	0.91	0.88	0.63	0.72	0.62	0.91	0.92
Cross Reference	0.50	0.43	0.11	0.21	0.16	0.48	0.51
Determine Pixel Coordinates	0.29	0.27	0.08	0.14	0.14	0.26	0.29
Elimination	0.74	0.74	0.26	0.47	0.34	0.77	0.78
Estimate Position	0.52	0.54	0.40	0.43	0.40	0.52	0.63
Extract Contextual Information	0.39	0.37	0.17	0.22	0.16	0.42	0.42
Final Confirmation	1.00	0.98	0.99	0.99	0.68	1.00	1.00
Focus On Target Area	1.00	0.98	0.96	0.98	0.97	0.99	1.00
Grouping Analysis	0.22	0.22	0.15	0.15	0.15	0.25	0.24
Infer Design Patterns	0.16	0.17	0.04	0.05	0.05	0.17	0.21
Interpret Ambiguous Language	0.14	0.15	0.01	0.05	0.04	0.21	0.17
Interpret Viewer Perspective	0.19	0.16	0.07	0.09	0.10	0.17	0.18
Iterative Counting	0.19	0.19	0.13	0.15	0.14	0.19	0.22
Iterative Refinement	0.49	0.62	0.06	0.18	0.11	0.54	0.75
Label Mapping	0.15	0.14	0.07	0.11	0.06	0.15	0.15
Lighting And Visibility Analysis	0.26	0.25	0.09	0.10	0.16	0.25	0.27
Locate Command In Interface	0.14	0.14	0.10	0.12	0.10	0.16	0.15
Locate Element	0.62	0.61	0.57	0.61	0.53	0.65	0.69
Map Command To UI Element	0.15	0.16	0.10	0.13	0.09	0.17	0.16
Match To Options	0.24	0.24	0.15	0.22	0.19	0.24	0.24
Output In Required Format	0.32	0.29	0.04	0.15	0.20	0.27	0.29
Pattern Recognition	0.09	0.13	0.03	0.03	0.03	0.12	0.17
Positional Ordering	0.26	0.28	0.11	0.13	0.09	0.23	0.34
Prioritize Clear Evidence	0.88	0.90	0.50	0.70	0.49	0.91	0.91
Reasoned Assumption	0.52	0.72	0.10	0.22	0.21	0.66	0.73
Select Central Interaction Point	0.14	0.16	0.05	0.10	0.04	0.17	0.16
Select Option	0.50	0.52	0.27	0.44	0.30	0.53	0.53
Shape Analysis	0.19	0.16	0.11	0.10	0.11	0.19	0.18
Spatial Analysis	0.61	0.60	0.20	0.32	0.21	0.58	0.70
Standard UI Layout Assumption	0.22	0.23	0.18	0.20	0.18	0.22	0.21
Structural Analysis	0.64	0.60	0.41	0.48	0.38	0.63	0.67
Synthesize Visual And Contextual Information	1.00	0.98	0.88	0.95	0.74	1.00	0.99
Text Recognition	0.28	0.26	0.19	0.24	0.18	0.28	0.30
Total Aggregation	0.18	0.18	0.18	0.15	0.19	0.20	0.20
Visual Analysis	0.97	0.95	0.75	0.86	0.67	0.98	0.98
Visual Element Counting	0.32	0.29	0.29	0.27	0.28	0.31	0.31
Visual Estimation	0.69	0.72	0.40	0.52	0.39	0.72	0.76
Visual Scanning	0.65	0.74	0.38	0.50	0.40	0.74	0.86
Visual Segmentation	0.64	0.71	0.40	0.47	0.34	0.70	0.81
Visual Verification	0.77	0.76	0.36	0.52	0.31	0.76	0.87

### Grounding & Search (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Analyze Structural Elements	0.53	0.24	0.30	0.37	0.38	0.26	0.32
Assess Visual Clarity	0.59	0.38	0.25	0.35	0.29	0.38	0.49
Assess Visual Indicators	0.22	0.10	0.09	0.12	0.10	0.16	0.13
Associative Mapping	0.31	0.13	0.11	0.24	0.23	0.14	0.18
Assume Standard Conventions	0.36	0.42	0.35	0.43	0.40	0.45	0.46
Assume Standard Dimensions	0.11	0.15	0.01	0.04	0.01	0.08	0.19
Center Point Estimation	0.06	0.15	0.02	0.03	0.02	0.06	0.17
Color Analysis	0.45	0.35	0.36	0.36	0.41	0.34	0.35
Concluding From Evidence	0.90	0.97	0.84	0.91	0.73	0.95	0.90
Contextual Inference	0.63	0.28	0.36	0.53	0.51	0.37	0.47
Determine Pixel Coordinates	0.04	0.24	0.00	0.01	0.01	0.08	0.20
Elimination	0.18	0.09	0.04	0.13	0.08	0.17	0.12
Estimate Position	0.30	0.31	0.17	0.25	0.11	0.24	0.29
Extract Contextual Information	0.20	0.04	0.04	0.12	0.08	0.06	0.08
Final Confirmation	0.72	0.81	0.43	0.71	0.45	0.81	0.82
Focus On Target Area	0.89	0.82	0.87	0.85	0.85	0.84	0.80
Locate Element	0.45	0.37	0.37	0.45	0.35	0.38	0.31
Match To Options	0.12	0.17	0.03	0.12	0.07	0.17	0.20
Output In Required Format	0.03	0.26	0.00	0.01	0.01	0.10	0.27
Prioritize Clear Evidence	0.45	0.31	0.13	0.25	0.15	0.34	0.39
Reasoned Assumption	0.12	0.07	0.05	0.09	0.12	0.08	0.16
Select Option	0.18	0.28	0.09	0.18	0.15	0.28	0.27
Standard UI Layout Assumption	0.19	0.16	0.16	0.18	0.18	0.16	0.13
Structural Analysis	0.48	0.21	0.24	0.31	0.29	0.24	0.28
Synthesize Visual And Contextual Information	0.89	0.68	0.53	0.76	0.63	0.71	0.78
Text Recognition	0.18	0.11	0.09	0.14	0.09	0.11	0.12
Total Aggregation	0.12	0.19	0.10	0.10	0.08	0.17	0.15
Visual Analysis	0.70	0.43	0.32	0.51	0.46	0.42	0.59
Visual Element Counting	0.23	0.25	0.22	0.21	0.22	0.24	0.23
Visual Estimation	0.36	0.27	0.11	0.17	0.14	0.19	0.31
Visual Scanning	0.05	0.09	0.15	0.04	0.03	0.09	0.05
Visual Segmentation	0.18	0.17	0.17	0.14	0.08	0.20	0.11
Visual Verification	0.16	0.13	0.09	0.14	0.07	0.15	0.12

### Knowledge & Recog. (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Analyze Compositional Elements	0.76	0.72	0.62	0.66	0.64	0.73	0.71
Analyze Image Structures	0.77	0.72	0.54	0.64	0.63	0.74	0.73
Analyze Visual Layout	0.81	0.77	0.64	0.72	0.72	0.81	0.79
Assess Information Availability	0.44	0.37	0.23	0.28	0.30	0.42	0.38
Assess Relevance	0.58	0.51	0.29	0.41	0.39	0.58	0.56
Associate Concepts	0.27	0.29	0.12	0.22	0.21	0.30	0.30
Associate Entities	0.34	0.33	0.28	0.30	0.32	0.35	0.37
Associate Visual Elements With Meaning	0.69	0.66	0.49	0.57	0.58	0.70	0.68
Balance Precision With Generality	0.74	0.69	0.54	0.63	0.60	0.74	0.66
Causal Reasoning	0.40	0.40	0.19	0.31	0.34	0.42	0.43
Compare Elements	0.78	0.79	0.48	0.65	0.67	0.82	0.82
Conciseness Optimization	0.71	0.65	0.56	0.63	0.59	0.68	0.60
Conclude Based On Insufficient Data	0.36	0.33	0.20	0.23	0.25	0.36	0.33
Contextual Analysis	0.61	0.57	0.29	0.45	0.45	0.65	0.62
Cross Reference Sources	0.19	0.15	0.03	0.11	0.11	0.19	0.20
Cross Reference Temporal Information	0.16	0.15	0.09	0.13	0.14	0.18	0.17
Cross Reference Visual And Textual Clues	0.40	0.33	0.21	0.26	0.29	0.36	0.39
Cultural Context Inference	0.20	0.19	0.13	0.18	0.18	0.22	0.21
Disambiguation	0.29	0.22	0.04	0.11	0.12	0.29	0.29
Elimination	0.79	0.76	0.32	0.55	0.60	0.82	0.78
Extract Text From Image	0.26	0.23	0.19	0.19	0.22	0.24	0.26
Geographical Context Analysis	0.20	0.21	0.14	0.18	0.21	0.23	0.23
Handle Missing Information	0.35	0.33	0.17	0.22	0.26	0.37	0.35
Hypothesis Testing	0.65	0.69	0.25	0.48	0.51	0.72	0.75
Identify Structural Patterns	0.50	0.46	0.27	0.38	0.39	0.52	0.52
Identify Time Period Relevance	0.19	0.18	0.12	0.16	0.19	0.20	0.20
Identify Visual Focus	0.57	0.51	0.22	0.40	0.39	0.61	0.58
Infer Context From Visual And Contextual Clues	0.80	0.77	0.57	0.68	0.67	0.83	0.80
Infer From Common Conventions	0.46	0.47	0.15	0.33	0.36	0.53	0.51
Infer Function From Design	0.30	0.28	0.13	0.22	0.23	0.32	0.30
Infer Intent From Visual Cues	0.25	0.22	0.07	0.16	0.17	0.25	0.26
Interpret Color Significance	0.15	0.14	0.07	0.09	0.12	0.16	0.17
Interpret Lighting And Atmosphere	0.26	0.22	0.11	0.15	0.20	0.26	0.24
Interpret Question Intent	0.23	0.24	0.04	0.13	0.14	0.33	0.24
Iterative Refinement	0.57	0.65	0.12	0.36	0.39	0.67	0.71
Iterative Refinement	0.65	0.60	0.12	0.34	0.37	0.64	0.68
Locate Element	0.29	0.25	0.16	0.21	0.22	0.29	0.31
Make Informed Assumption	0.71	0.78	0.51	0.63	0.67	0.79	0.81
Match Visual Characteristics To Known Models	0.40	0.39	0.21	0.34	0.35	0.42	0.44
Model Identification From Visual Cues	0.20	0.19	0.09	0.18	0.16	0.22	0.24
Pattern Recognition	0.62	0.61	0.31	0.50	0.53	0.68	0.68
Prioritize Evidence Based Decision	0.89	0.88	0.70	0.83	0.81	0.90	0.90
Resolve Conflicting Information	0.15	0.14	0.02	0.09	0.08	0.18	0.22
Spatial Orientation Analysis	0.31	0.27	0.10	0.19	0.20	0.33	0.35
Synthesize Information	0.88	0.84	0.68	0.76	0.77	0.88	0.86
Systematic Visual Scanning	0.75	0.71	0.35	0.59	0.58	0.80	0.79
Use Domain Knowledge	0.79	0.75	0.57	0.66	0.70	0.81	0.79
Visual Analysis	0.86	0.83	0.69	0.77	0.77	0.88	0.85
Visual Comparison	0.56	0.55	0.30	0.44	0.48	0.62	0.61
Visual Recognition For Identity Confirmation	0.37	0.36	0.25	0.33	0.34	0.39	0.43
Visual Verification	0.87	0.82	0.39	0.69	0.66	0.87	0.85

### Knowledge & Recog. (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Analyze Compositional Elements	0.76	0.37	0.46	0.55	0.64	0.45	0.53
Analyze Image Structures	0.58	0.30	0.40	0.40	0.50	0.36	0.43
Analyze Visual Layout	0.73	0.35	0.44	0.49	0.61	0.44	0.51
Assess Information Availability	0.36	0.19	0.16	0.15	0.18	0.20	0.22
Assess Relevance	0.29	0.11	0.15	0.17	0.24	0.15	0.17
Associate Entities	0.16	0.18	0.12	0.16	0.15	0.20	0.20
Associate Visual Elements With Meaning	0.47	0.28	0.32	0.39	0.45	0.32	0.39
Balance Precision With Generality	0.52	0.31	0.30	0.35	0.41	0.37	0.42
Compare Elements	0.38	0.21	0.26	0.26	0.34	0.27	0.29
Conciseness Optimization	0.49	0.33	0.31	0.33	0.37	0.38	0.41
Conclude Based On Insufficient Data	0.35	0.18	0.14	0.15	0.16	0.19	0.20
Contextual Analysis	0.22	0.10	0.13	0.15	0.19	0.17	0.16
Cross Reference Visual And Textual Clues	0.18	0.08	0.09	0.12	0.16	0.11	0.14
Extract Text From Image	0.19	0.11	0.14	0.14	0.16	0.15	0.15
Handle Missing Information	0.24	0.08	0.07	0.08	0.12	0.12	0.14
Identify Structural Patterns	0.20	0.09	0.12	0.12	0.19	0.12	0.14
Infer Context From Visual And Contextual Clues	0.55	0.26	0.32	0.			

### Captioning & IF (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Adjective Selection For Visual Description	0.67	0.61	0.58	0.60	0.49	0.63	0.62
Adopt Perspective	0.60	0.54	0.49	0.49	0.37	0.56	0.53
Analyze Atmosphere	0.30	0.30	0.26	0.26	0.21	0.32	0.29
Analyze Contrasts	0.30	0.24	0.20	0.20	0.15	0.29	0.24
Analyze Emotional Tone	0.23	0.23	0.18	0.19	0.17	0.23	0.22
Analyze Visual Composition	0.52	0.42	0.36	0.38	0.33	0.46	0.42
Balance Clarity And Impact	0.56	0.53	0.44	0.48	0.38	0.55	0.52
Color Analysis	0.30	0.24	0.21	0.21	0.17	0.27	0.26
Constrain Output Length	0.14	0.17	0.07	0.17	0.25	0.20	0.17
Constraint Verification	0.19	0.19	0.06	0.14	0.25	0.23	0.17
Construct Narrative	0.47	0.48	0.38	0.47	0.35	0.47	0.50
Contextual Inference	0.71	0.67	0.58	0.62	0.56	0.71	0.68
Define Narrative Structure	0.49	0.49	0.38	0.47	0.35	0.49	0.49
Describe Shape And Texture	0.24	0.17	0.11	0.16	0.12	0.21	0.17
Elimination	0.17	0.20	0.04	0.10	0.07	0.20	0.18
Extract Key Data Points	0.40	0.38	0.31	0.31	0.30	0.39	0.39
Filter Irrelevant Information	0.36	0.34	0.15	0.20	0.20	0.39	0.35
Focus On Key Attributes	0.77	0.73	0.63	0.66	0.59	0.74	0.69
Follow Formatting Rules	0.18	0.16	0.06	0.10	0.20	0.21	0.17
Identify Central Visual Element	0.64	0.58	0.51	0.55	0.48	0.65	0.60
Identify Core Concept	0.31	0.29	0.20	0.25	0.17	0.29	0.28
Implied Meaning	0.53	0.53	0.44	0.47	0.33	0.54	0.50
Infer Causal Factors	0.23	0.26	0.15	0.18	0.13	0.23	0.20
Infer Design Intent	0.24	0.17	0.10	0.12	0.11	0.21	0.18
Infer Functional Purpose	0.29	0.30	0.20	0.23	0.18	0.30	0.26
Infer Implicit Meaning	0.25	0.24	0.12	0.15	0.09	0.24	0.22
Interpret Diagram Components	0.17	0.13	0.10	0.12	0.09	0.14	0.12
Link Elements To Theme	0.29	0.28	0.15	0.16	0.14	0.26	0.24
Map Visual To Concept	0.23	0.24	0.14	0.18	0.13	0.23	0.21
Simplify Technical Concepts	0.22	0.24	0.12	0.15	0.14	0.22	0.18
Synthesize Information	0.59	0.62	0.40	0.47	0.39	0.62	0.56
Tone Consistency	0.73	0.69	0.61	0.64	0.57	0.68	0.65
Use Analogies	0.18	0.18	0.10	0.14	0.13	0.18	0.16
Use Imagery And Metaphor	0.47	0.46	0.37	0.44	0.34	0.47	0.46
Validate Consistency And Completeness	0.29	0.25	0.05	0.10	0.07	0.24	0.22
Word Count Enforcement	0.09	0.14	0.05	0.09	0.15	0.15	0.11

### Captioning & IF (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Adjective Selection For Visual Description	0.65	0.49	0.51	0.56	0.57	0.52	0.58
Adopt Perspective	0.49	0.35	0.33	0.33	0.46	0.38	0.39
Analyze Atmosphere	0.23	0.20	0.19	0.20	0.25	0.21	0.23
Analyze Contrasts	0.17	0.10	0.11	0.14	0.17	0.10	0.10
Analyze Emotional Tone	0.39	0.26	0.27	0.29	0.36	0.29	0.35
Analyze Visual Composition	0.38	0.26	0.21	0.24	0.33	0.25	0.32
Balance Clarity And Impact	0.22	0.12	0.13	0.15	0.20	0.14	0.20
Color Analysis	0.35	0.30	0.26	0.29	0.35	0.26	0.33
Construct Narrative	0.58	0.44	0.45	0.48	0.51	0.45	0.52
Contextual Inference	0.35	0.28	0.25	0.25	0.33	0.26	0.31
Define Narrative Structure	0.16	0.06	0.07	0.07	0.12	0.07	0.09
Describe Shape And Texture	0.29	0.23	0.19	0.23	0.27	0.25	0.26
Extract Key Data Points	0.61	0.42	0.42	0.48	0.49	0.45	0.52
Focus On Key Attributes	0.46	0.38	0.37	0.40	0.46	0.38	0.45
Identify Central Visual Element	0.37	0.29	0.25	0.28	0.36	0.27	0.34
Implied Meaning	0.19	0.12	0.10	0.13	0.19	0.12	0.13
Infer Functional Purpose	0.35	0.22	0.22	0.19	0.31	0.27	0.29
Synthesize Information	0.55	0.44	0.39	0.39	0.51	0.43	0.50
Tone Consistency	0.32	0.27	0.23	0.25	0.29	0.25	0.31
Use Imagery And Metaphor							

### Chart & OCR (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Aggregate Handling	0.41	0.40	0.23	0.29	0.29	0.38	0.39
Align Event Timings	0.18	0.17	0.13	0.15	0.15	0.17	0.19
Align Output With Instruction Format	0.87	0.88	0.60	0.73	0.68	0.94	0.82
Apply Arithmetic Operation	0.28	0.25	0.21	0.22	0.23	0.24	0.26
Apply Domain Knowledge	0.44	0.46	0.20	0.26	0.31	0.54	0.44
Assess Data Relevance And Availability	0.68	0.60	0.21	0.34	0.36	0.67	0.53
Avoid Unnecessary Calculation	0.97	0.98	0.97	0.98	0.97	0.98	0.96
Axis Analysis	0.53	0.57	0.37	0.46	0.46	0.56	0.56
Axis Verification	0.36	0.33	0.17	0.25	0.23	0.35	0.33
Category Management	0.59	0.55	0.47	0.52	0.53	0.60	0.58
Comparative Analysis	0.65	0.69	0.56	0.60	0.59	0.66	0.68
Condition Matching	0.48	0.46	0.34	0.41	0.41	0.51	0.49
Consider Alternative Interpretations	0.46	0.51	0.11	0.23	0.24	0.49	0.47
Consider Data Granularity	0.79	0.76	0.45	0.59	0.53	0.81	0.77
Consistency Assessment	0.45	0.33	0.16	0.23	0.22	0.39	0.37
Contextual Analysis	0.20	0.25	0.06	0.10	0.10	0.28	0.22
Counting Analysis	0.18	0.17	0.15	0.16	0.16	0.16	0.18
Criteria Reinterpretation	0.40	0.43	0.15	0.25	0.28	0.45	0.42
Cross Reference	0.35	0.32	0.11	0.17	0.17	0.34	0.36
Cross Validation	0.32	0.23	0.05	0.09	0.09	0.22	0.26
Difference Calculation	0.16	0.16	0.12	0.13	0.14	0.15	0.17
Eliminate Irrelevant Information	0.45	0.37	0.13	0.21	0.23	0.44	0.40
Estimation Of Values From Visual Data	0.31	0.34	0.26	0.29	0.28	0.33	0.34
Exhaustive Verification	0.63	0.77	0.24	0.53	0.44	0.62	0.75
Extract Numerical Information From Visuals	0.80	0.78	0.69	0.73	0.72	0.80	0.80
Extract Relevant Attribute From Data	0.74	0.71	0.54	0.64	0.61	0.78	0.73
Extract Textual Data	0.44	0.37	0.26	0.30	0.28	0.40	0.39
Extract Value From Graph At Specific Point	0.28	0.27	0.22	0.26	0.25	0.29	0.31
Filtering	0.31	0.26	0.13	0.18	0.19	0.32	0.28
Final Verification	0.95	0.93	0.35	0.71	0.57	0.96	0.90
Group Data By Category	0.30	0.28	0.16	0.21	0.21	0.29	0.28
Identify Extreme Values	0.33	0.32	0.29	0.31	0.31	0.33	0.32
Identify Top N By Metric	0.22	0.21	0.20	0.21	0.21	0.22	0.21
Infer From Context	0.31	0.37	0.10	0.17	0.19	0.41	0.33
Inference Under Data Uncertainty	0.24	0.28	0.11	0.15	0.17	0.25	0.26
Iterative Refinement	0.36	0.40	0.15	0.25	0.25	0.36	0.46
Label Verification	0.64	0.77	0.37	0.59	0.51	0.63	0.80
Locate Relevant Section By Heading	0.38	0.36	0.29	0.34	0.33	0.39	0.40
Mapping	0.61	0.57	0.32	0.45	0.41	0.63	0.60
Option Selection Based On Criteria	0.37	0.38	0.28	0.33	0.30	0.40	0.38
Prioritize Explicit Information Over Inference	0.98	0.98	0.96	0.98	0.95	0.99	0.98
Prioritize Explicit Labels Over Visual Cues	0.88	0.84	0.74	0.83	0.79	0.89	0.87
Rank Elements By Numerical Attribute	0.27	0.28	0.22	0.23	0.23	0.28	0.28
Scope Constraint	0.64	0.57	0.21	0.38	0.37	0.64	0.61
Sectional Chart Navigation	0.36	0.32	0.18	0.27	0.25	0.37	0.37
Select Best Fit Under Uncertainty	0.21	0.23	0.09	0.13	0.13	0.20	0.21
Select Closest Available Data Point	0.20	0.25	0.13	0.16	0.17	0.23	0.25
Sequential Analysis	0.51	0.55	0.31	0.43	0.39	0.55	0.65
Spatial Position To Category Mapping	0.33	0.28	0.12	0.20	0.18	0.30	0.34
Strategy Adjustment	0.27	0.28	0.12	0.17	0.18	0.24	0.28
Targeted Information Search	0.95	0.97	0.85	0.94	0.92	0.98	0.97
Temporal Filtering	0.33	0.32	0.29	0.32	0.31	0.33	0.34
Threshold Comparison	0.24	0.21	0.15	0.18	0.17	0.21	0.23
Validate Data Type Consistency	0.47	0.41	0.17	0.29	0.26	0.45	0.43
Validate Task Goal	0.93	0.89	0.36	0.69	0.57	0.95	0.87
Value Matching	0.53	0.52	0.23	0.36	0.31	0.54	0.54
Verification Calculation	0.18	0.18	0.03	0.08	0.08	0.16	0.18
Visual Comparison	0.43	0.44	0.32	0.39	0.36	0.44	0.46
Visual Data Interpretation	0.74	0.69	0.62	0.71	0.66	0.74	0.71
Visual Dominance Selection	0.19	0.19	0.15	0.18	0.17	0.20	0.20
Visual Feature Identification	0.64	0.58	0.39	0.52	0.48	0.65	0.64
Visual Label Association	0.72	0.69	0.57	0.66	0.64	0.74	0.71
Visual Pattern Recognition	0.44	0.41	0.21	0.33	0.29	0.44	0.43
Visual Position Comparison	0.32	0.32	0.16	0.24	0.23	0.32	0.35
Visual Proximity Assessment	0.28	0.28	0.14	0.21	0.20	0.27	0.30

### Chart & OCR (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Aggregate Handling	0.21	0.18	0.07	0.17	0.12	0.19	0.20
Align Output With Instruction Format	0.44	0.75	0.28	0.48	0.22	0.59	0.77
Apply Arithmetic Operation	0.16	0.19	0.14	0.18	0.15	0.18	0.16
Avoid Unnecessary Calculation	0.95	0.95	0.99	0.94	0.90	0.94	0.95
Axis Analysis	0.32	0.11	0.18	0.10	0.23	0.12	0.22
Chart Legend Interpretation	0.39	0.18	0.25	0.17	0.28	0.22	0.28
Comparative Analysis	0.53	0.50	0.43	0.55	0.48	0.49	0.50
Condition Matching	0.22	0.20	0.11	0.15	0.17	0.18	0.21
Consider Data Granularity	0.29	0.17	0.07	0.15	0.07	0.19	0.27
Estimation Of Values From Visual Data	0.21	0.13	0.08	0.13	0.12	0.14	0.16
Extract Numerical Information From Visuals	0.59	0.45	0.58	0.45	0.45	0.50	0.53
Extract Relevant Attribute From Data	0.40	0.26	0.23	0.22	0.27	0.28	0.32
Extract Textual Data	0.22	0.34	0.15	0.25	0.14	0.25	0.24
Extract Value From Graph At Specific Point	0.17	0.14	0.14	0.12	0.12	0.14	0.16
Final Verification	0.12	0.19	0.12	0.15	0.03	0.14	0.19
Identify Extreme Values	0.26	0.26	0.24	0.28	0.23	0.26	0.25
Identify Top N By Metric	0.19	0.19	0.16	0.20	0.16	0.19	0.18
Label Verification	0.19	0.10	0.10	0.09	0.04	0.12	0.16
Locate Relevant Section By Heading	0.15	0.05	0.15	0.05	0.15	0.08	0.09
Mapping	0.20	0.07	0.06	0.05	0.05	0.08	0.13
Option Selection Based On Criteria	0.19	0.23	0.12	0.20	0.11	0.20	0.22
Prioritize Explicit Information Over Inference	0.89	0.90	0.87	0.84	0.60	0.88	0.92
Prioritize Explicit Labels Over Visual Cues	0.88	0.35	0.38	0.31	0.29	0.41	0.49
Rank Elements By Numerical Attribute	0.16	0.19	0.11	0.17	0.11	0.17	0.15
Targeted Information Search	0.66	0.63	0.50	0.63	0.49	0.60	0.67
Temporal Filtering	0.26	0.24	0.18	0.21	0.22	0.24	0.24
Validate Task Goal	0.17	0.17	0.13	0.14	0.03	0.15	0.20
Value Matching	0.13	0.17	0.08	0.15	0.06	0.13	0.13
Visual Comparison	0.33	0.14	0.25	0.20	0.22	0.19	0.25
Visual Data Interpretation	0.63	0.36	0.53	0.43	0.49	0.45	0.55
Visual Dominance Selection	0.18	0.09	0.16	0.15	0.12	0.13	0.15
Visual							

# A11.4 High Level Cognitive Foundations Presence Rate

### Average (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.46	0.46	0.40	0.43	0.51	0.47	0.46
Adaptive Detail Management	0.78	0.75	0.54	0.69	0.81	0.82	0.70
Arithmetic Calculation	0.27	0.28	0.21	0.24	0.26	0.30	0.21
Backtracking	0.22	0.33	0.12	0.20	0.36	0.48	0.22
Backward Chaining	0.06	0.08	0.07	0.08	0.10	0.12	0.16
Causal Organization	0.50	0.51	0.42	0.47	0.59	0.54	0.51
Compositionality	0.96	0.94	0.91	0.92	0.97	0.96	0.96
Conceptual Level Processing	0.71	0.66	0.60	0.64	0.75	0.66	0.66
Context Alignment	0.77	0.76	0.61	0.74	0.85	0.79	0.77
Context Awareness	0.61	0.52	0.46	0.55	0.72	0.58	0.67
Decomposition And Integration	0.79	0.78	0.76	0.79	0.82	0.83	0.85
Forward Chaining	0.89	0.89	0.90	0.93	0.93	0.90	0.94
Goal Management	0.89	0.92	0.84	0.95	0.96	0.95	0.97
Hierarchical Organization	0.79	0.73	0.66	0.72	0.81	0.74	0.75
Knowledge Structure Alignment	0.94	0.90	0.83	0.90	0.96	0.91	0.87
Logical Coherence	0.99	0.98	0.97	0.98	0.99	0.98	1.00
Mental Imagery Simulation	0.60	0.58	0.47	0.54	0.64	0.65	0.53
Network Organization	0.59	0.57	0.47	0.53	0.64	0.58	0.55
Ordinal Organization	0.56	0.58	0.56	0.61	0.64	0.61	0.64
Pattern Recognition	0.47	0.56	0.48	0.53	0.57	0.57	0.50
Perception Then Reasoning	0.87	0.77	0.64	0.68	0.84	0.80	0.62
Productivity	0.39	0.44	0.29	0.39	0.48	0.51	0.43
Representational Restructuring	0.51	0.54	0.40	0.47	0.58	0.58	0.53
Selective Attention	0.98	0.97	0.95	0.98	0.98	0.98	1.00
Self Awareness	0.76	0.80	0.49	0.62	0.86	0.84	0.75
Self Evaluation	0.88	0.91	0.46	0.73	0.94	0.94	0.73
Sequential Organization	0.99	0.97	0.96	0.98	0.98	0.98	0.99
Spatial Organization	0.74	0.71	0.71	0.73	0.75	0.75	0.71
Strategy Selection	0.69	0.70	0.57	0.69	0.80	0.78	0.80
Systematic Regional Synthesis	0.70	0.74	0.52	0.64	0.78	0.81	0.61
Temporal Organization	0.63	0.60	0.61	0.62	0.71	0.69	0.66
Verification	0.90	0.87	0.65	0.83	0.93	0.91	0.81
Visual Foraging	0.86	0.90	0.72	0.85	0.94	0.94	0.84
Visual Reference Or Grounding	0.74	0.71	0.66	0.72	0.76	0.75	0.64

### Average (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.30	0.19	0.25	0.26	0.22	0.27	0.30
Adaptive Detail Management	0.30	0.11	0.18	0.15	0.14	0.18	0.22
Arithmetic Calculation	0.15	0.20	0.12	0.15	0.17	0.16	0.11
Backtracking	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Backward Chaining	0.02	0.01	0.02	0.02	0.02	0.02	0.08
Causal Organization	0.24	0.20	0.24	0.22	0.21	0.24	0.28
Compositionality	0.80	0.63	0.58	0.68	0.68	0.69	0.74
Conceptual Level Processing	0.55	0.35	0.41	0.45	0.41	0.47	0.50
Context Alignment	0.43	0.22	0.39	0.36	0.27	0.37	0.39
Context Awareness	0.22	0.13	0.30	0.22	0.17	0.21	0.24
Decomposition And Integration	0.55	0.50	0.45	0.51	0.53	0.47	0.56
Forward Chaining	0.77	0.74	0.72	0.76	0.76	0.77	0.79
Goal Management	0.41	0.43	0.62	0.41	0.45	0.36	0.62
Hierarchical Organization	0.51	0.32	0.38	0.38	0.40	0.40	0.47
Knowledge Structure Alignment	0.69	0.46	0.52	0.57	0.50	0.61	0.60
Logical Coherence	0.93	0.88	0.93	0.91	0.89	0.91	0.91
Mental Imagery Simulation	0.33	0.14	0.20	0.22	0.17	0.20	0.25
Network Organization	0.42	0.16	0.23	0.28	0.21	0.25	0.33
Ordinal Organization	0.45	0.39	0.38	0.45	0.39	0.40	0.38
Pattern Recognition	0.47	0.34	0.30	0.48	0.35	0.46	0.41
Perception Then Reasoning	0.69	0.38	0.25	0.38	0.45	0.39	0.38
Productivity	0.14	0.06	0.11	0.09	0.07	0.09	0.13
Representational Restructuring	0.28	0.12	0.17	0.20	0.14	0.22	0.23
Selective Attention	0.87	0.86	0.91	0.86	0.87	0.88	0.88
Self Awareness	0.29	0.14	0.20	0.18	0.18	0.18	0.27
Self Evaluation	0.19	0.10	0.11	0.16	0.12	0.12	0.13
Sequential Organization	0.88	0.78	0.80	0.77	0.79	0.83	0.85
Spatial Organization	0.71	0.53	0.59	0.60	0.56	0.59	0.65
Strategy Selection	0.27	0.15	0.34	0.23	0.18	0.23	0.32
Systematic Regional Synthesis	0.19	0.24	0.14	0.14	0.21	0.15	0.15
Temporal Organization	0.47	0.40	0.36	0.47	0.44	0.40	0.44
Verification	0.35	0.20	0.23	0.28	0.22	0.26	0.25
Visual Foraging	0.52	0.32	0.34	0.37	0.32	0.37	0.46
Visual Reference Or Grounding	0.61	0.42	0.50	0.46	0.46	0.49	0.46

### Ground & Search (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.07	0.07	0.05	0.04	0.11	0.08	0.07
Adaptive Detail Management	0.70	0.59	0.34	0.49	0.70	0.73	0.43
Arithmetic Calculation	0.15	0.13	0.08	0.10	0.10	0.14	0.03
Backtracking	0.11	0.24	0.02	0.10	0.26	0.41	0.07
Backward Chaining	0.00	0.01	0.04	0.03	0.01	0.03	0.12
Causal Organization	0.15	0.15	0.14	0.12	0.21	0.18	0.16
Compositionality	0.95	0.88	0.91	0.89	0.95	0.96	0.92
Conceptual Level Processing	0.31	0.22	0.16	0.20	0.38	0.24	0.22
Context Alignment	0.60	0.47	0.38	0.55	0.72	0.56	0.55
Context Awareness	0.65	0.42	0.53	0.61	0.75	0.49	0.75
Decomposition And Integration	0.69	0.60	0.72	0.68	0.69	0.74	0.76
Forward Chaining	0.95	0.83	0.98	0.98	0.94	0.89	0.94
Goal Management	0.92	0.93	0.95	0.99	0.98	1.00	1.00
Hierarchical Organization	0.65	0.41	0.46	0.48	0.59	0.49	0.48
Knowledge Structure Alignment	0.87	0.70	0.64	0.75	0.87	0.75	0.65
Logical Coherence	1.00	0.97	1.00	1.00	1.00	1.00	1.00
Mental Imagery Simulation	0.92	0.83	0.80	0.82	0.90	0.94	0.76
Network Organization	0.30	0.20	0.16	0.19	0.34	0.25	0.23
Ordinal Organization	0.32	0.29	0.42	0.39	0.34	0.31	0.35
Pattern Recognition	0.15	0.21	0.23	0.21	0.25	0.19	0.19
Perception Then Reasoning	0.92	0.72	0.53	0.60	0.83	0.78	0.39
Productivity	0.10	0.12	0.05	0.07	0.20	0.17	0.13
Representational Restructuring	0.19	0.20	0.12	0.14	0.27	0.23	0.18
Selective Attention	1.00	0.98	1.00	1.00	1.00	1.00	1.00
Self Awareness	0.82	0.79	0.47	0.56	0.87	0.85	0.77
Self Evaluation	0.89	0.90	0.33	0.67	0.92	0.96	0.58
Sequential Organization	1.00	0.96	1.00	1.00	1.00	1.00	1.00
Spatial Organization	0.96	0.92	0.95	0.94	0.95	0.97	0.93
Strategy Selection	0.53	0.48	0.46	0.54	0.66	0.68	0.73
Systematic Regional Synthesis	0.90	0.88	0.61	0.75	0.89	0.97	0.63
Temporal Organization	0.50	0.38	0.63	0.57	0.57	0.55	0.58
Verification	0.91	0.78	0.57	0.83	0.89	0.92	0.69
Visual Foraging	0.90	0.92	0.72	0.85	0.97	0.96	0.80
Visual Reference Or Grounding	0.99	0.94	0.97	0.97	0.97	0.98	0.81

### Ground & Search (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.07	0.03	0.03	0.05	0.04	0.06	0.08
Adaptive Detail Management	0.23	0.04	0.04	0.07	0.06	0.11	0.06
Arithmetic Calculation	0.03	0.13	0.03	0.06	0.09	0.07	0.01
Backtracking	0.00	0.01	0.01	0.00	0.01	0.01	0.00
Backward Chaining	0.01	0.01	0.03	0.01	0.05	0.00	0.06
Causal Organization	0.07	0.09	0.11	0.10	0.08	0.08	0.12
Compositionality	0.69	0.49	0.52	0.59	0.50	0.51	0.54
Conceptual Level Processing	0.22	0.11	0.08	0.15	0.13	0.19	0.16
Context Alignment	0.31	0.11	0.12	0.27	0.17	0.26	0.24
Context Awareness	0.26	0.15	0.20	0.30	0.21	0.26	0.33
Decomposition And Integration	0.34	0.35	0.35	0.36	0.33	0.23	0.33
Forward Chaining	0.73	0.68	0.86	0.75	0.70	0.68	0.77
Goal Management	0.29	0.42	0.86	0.38	0.45	0.25	0.66
Hierarchical Organization	0.28	0.14	0.20	0.21	0.20	0.17	0.22
Knowledge Structure Alignment	0.46	0.17	0.21	0.37	0.25	0.35	0.35
Logical Coherence	0.95	0.89	0.87	0.92	0.89	0.88	0.86
Mental Imagery Simulation	0.50	0.18	0.39	0.34	0.22	0.23	0.31
Network Organization	0.19	0.01	0.03	0.07	0.02	0.04	0.07
Ordinal Organization	0.22	0.19	0.15	0.23	0.20	0.15	0.12
Pattern Recognition	0.30	0.16	0.13	0.31	0.17	0.26	0.28
Perception Then Reasoning	0.73	0.31	0.14	0.37	0.41	0.38	0.21
Productivity	0.02	0.01	0.01	0.01	0.00	0.02	0.02
Representational Restructuring	0.10	0.02	0.02	0.05	0.03	0.06	0.06
Selective Attention	0.92	0.94	0.97	0.93	0.96	0.89	0.91
Self Awareness	0.41	0.18	0.17	0.21	0.21	0.26	0.36
Self Evaluation	0.20	0.12	0.04	0.13	0.12	0.14	0.11
Sequential Organization	0.82	0.73	0.81	0.72	0.72	0.72	0.73
Spatial Organization	0.91	0.76	0.77	0.86	0.77	0.73	0.81
Strategy Selection	0.13	0.05	0.27	0.14	0.09	0.11	0.22
Systematic Regional Synthesis	0.07	0.20	0.12	0.08	0.16	0.06	0.04
Temporal Organization	0.31	0.30	0.30	0.39	0.35	0.20	0.26
Verification	0.34	0.19	0.08	0.24	0.21	0.23	0.19
Visual Foraging	0.52	0.25	0.37	0.29	0.22	0.28	0.33
Visual Reference Or Grounding	0.84	0.73	0.60	0.73	0.71	0.67	0.58

### Knowledge & Recog. (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.50	0.48	0.41	0.47	0.55	0.47	0.51
Adaptive Detail Management	0.84	0.82	0.55	0.77	0.90	0.88	0.78
Arithmetic Calculation	0.06	0.06	0.05	0.05	0.06	0.07	0.04
Backtracking	0.14	0.30	0.04	0.15	0.34	0.47	0.16
Backward Chaining	0.01	0.02	0.02	0.02	0.02	0.04	0.06
Causal Organization	0.44	0.46	0.33	0.44	0.56	0.49	0.50
Compositionality	0.94	0.90	0.85	0.90	0.96	0.92	0.93
Conceptual Level Processing	0.82	0.74	0.70	0.73	0.83	0.74	0.76
Context Alignment	0.82	0.83	0.70	0.84	0.91	0.85	0.86
Context Awareness	0.76	0.69	0.64	0.73	0.85	0.70	0.83
Decomposition And Integration	0.65	0.69	0.61	0.70	0.73	0.72	0.76
Forward Chaining	0.76	0.80	0.83	0.90	0.85	0.79	0.87
Goal Management	0.76	0.83	0.65	0.94	0.91	0.89	0.92
Hierarchical Organization	0.78	0.72	0.61	0.71	0.83	0.71	0.73
Knowledge Structure Alignment	0.96	0.91	0.88	0.93	0.97	0.93	0.92
Logical Coherence	0.99	0.98	0.99	0.99	1.00	0.97	0.99
Mental Imagery Simulation	0.45	0.49	0.35	0.44	0.58	0.62	0.46
Network Organization	0.62	0.61	0.47	0.55	0.68	0.63	0.62
Ordinal Organization	0.30	0.34	0.32	0.36	0.41	0.38	0.39
Pattern Recognition	0.48	0.59	0.53	0.57	0.58	0.62	0.57
Perception Then Reasoning	0.81	0.62	0.55	0.53	0.74	0.64	0.54
Productivity	0.40	0.47	0.24	0.42	0.55	0.58	0.49
Representational Restructuring	0.46	0.50	0.33	0.43	0.57	0.57	0.49
Selective Attention	0.98	0.97	0.95	0.99	1.00	0.98	0.99
Self Awareness	0.79	0.85	0.54	0.69	0.90	0.89	0.82
Self Evaluation	0.82	0.89	0.38	0.70	0.94	0.93	0.74
Sequential Organization	0.98	0.94	0.96	0.98	0.98	0.94	0.98
Spatial Organization	0.62	0.57	0.59	0.61	0.63	0.64	0.60
Strategy Selection	0.67	0.72	0.52	0.72	0.82	0.79	0.81
Systematic Regional Synthesis	0.40	0.45	0.22	0.36	0.53	0.57	0.36
Temporal Organization	0.50	0.49	0.50	0.51	0.61	0.59	0.55
Verification	0.88	0.85	0.62	0.81	0.93	0.89	0.84
Visual Foraging	0.82	0.90	0.61	0.83	0.95	0.93	0.81
Visual Reference Or Grounding	0.62	0.58	0.54	0.60	0.65	0.64	0.54

### Knowledge & Recog. (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.34	0.20	0.27	0.26	0.28	0.31	0.34
Adaptive Detail Management	0.41	0.12	0.17	0.14	0.20	0.24	0.27
Arithmetic Calculation	0.04	0.06	0.04	0.04	0.05	0.04	0.02
Backtracking	0.02	0.00	0.01	0.01	0.01	0.01	0.01
Backward Chaining	0.01	0.01	0.02	0.01	0.01	0.01	0.02
Causal Organization	0.17	0.18	0.25	0.16	0.20	0.20	0.18
Compositionality	0.80	0.44	0.65	0.52	0.56	0.62	0.70
Conceptual Level Processing	0.64	0.43	0.48	0.52	0.52	0.57	0.59
Context Alignment	0.57	0.37	0.47	0.47	0.46	0.53	0.53
Context Awareness	0.42	0.28	0.34	0.33	0.36	0.40	0.40
Decomposition And Integration	0.48	0.31	0.52	0.37	0.40	0.37	0.48
Forward Chaining	0.72	0.67	0.83	0.69	0.73	0.75	0.75
Goal Management	0.28	0.27	0.71	0.24	0.31	0.26	0.38
Hierarchical Organization	0.47	0.24	0.39	0.29	0.36	0.35	0.43
Knowledge Structure Alignment	0.71	0.51	0.57	0.62	0.61	0.67	0.68
Logical Coherence	0.97	0.92	0.96	0.93	0.94	0.97	0.96
Mental Imagery Simulation	0.26	0.10	0.18	0.18	0.14	0.16	0.23
Network Organization	0.45	0.13	0.26	0.28	0.25	0.28	0.41
Ordinal Organization	0.21	0.18	0.23	0.22	0.21	0.21	0.19
Pattern Recognition	0.48	0.37	0.36	0.48	0.43	0.53	0.51
Perception Then Reasoning	0.69	0.24	0.26	0.31	0.37	0.30	0.43
Productivity	0.17	0.05	0.11	0.08	0.07	0.09	0.17
Representational Restructuring	0.24	0.10	0.18	0.17	0.15	0.19	0.23
Selective Attention	0.88	0.89	0.94	0.85	0.90	0.91	0.88
Self Awareness	0.48	0.27	0.35	0.30	0.32	0.32	0.42
Self Evaluation	0.24	0.12	0.10	0.14	0.14	0.14	0.15
Sequential Organization	0.91	0.71	0.84	0.66	0.76	0.81	0.84
Spatial Organization	0.65	0.44	0.48	0.56	0.50	0.53	0.61
Strategy Selection	0.28	0.17	0.37	0.23	0.23	0.27	0.33
Systematic Regional Synthesis	0.04	0.08	0.09	0.04	0.08	0.04	0.04
Temporal Organization	0.37	0.30	0.36	0.38	0.37	0.35	0.39
Verification	0.40	0.20	0.21	0.22	0.23	0.26	0.26
Visual Foraging	0.45	0.20	0.34	0.30	0.25	0.31	0.42
Visual Reference Or Grounding	0.49	0.35	0.35	0.37	0.37	0.37	0.39

### Captioning & IF (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.73	0.74	0.65	0.73	0.83	0.75	0.78
Adaptive Detail Management	0.86	0.82	0.63	0.80	0.89	0.87	0.90
Arithmetic Calculation	0.14	0.14	0.09	0.11	0.12	0.17	0.09
Backtracking	0.12	0.23	0.04	0.13	0.31	0.33	0.19
Backward Chaining	0.04	0.04	0.05	0.09	0.10	0.10	0.21
Causal Organization	0.39	0.46	0.30	0.41	0.57	0.48	0.38
Compositionality	0.93	0.90	0.81	0.88	0.93	0.91	0.96
Conceptual Level Processing	0.87	0.84	0.82	0.84	0.92	0.85	0.88
Context Alignment	0.86	0.85	0.66	0.81	0.91	0.85	0.93
Context Awareness	0.79	0.77	0.57	0.76	0.89	0.82	0.92
Decomposition And Integration	0.86	0.83	0.71	0.81	0.89	0.84	0.96
Forward Chaining	0.73	0.77	0.66	0.77	0.82	0.79	0.85
Goal Management	0.81	0.83	0.58	0.79	0.87	0.86	0.94
Hierarchical Organization	0.85	0.79	0.69	0.78	0.88	0.83	0.93
Knowledge Structure Alignment	0.95	0.95	0.88	0.91	0.97	0.94	0.97
Logical Coherence	0.95	0.93	0.85	0.92	0.96	0.92	0.98
Mental Imagery Simulation	0.49	0.45	0.31	0.41	0.56	0.52	0.44
Network Organization	0.79	0.77	0.67	0.75	0.85	0.79	0.79
Ordinal Organization	0.55	0.56	0.48	0.61	0.69	0.66	0.76
Pattern Recognition	0.59	0.66	0.55	0.66	0.74	0.71	0.63
Perception Then Reasoning	0.64	0.61	0.48	0.53	0.72	0.63	0.52
Productivity	0.73	0.76	0.61	0.74	0.81	0.80	0.84
Representational Restructuring	0.74	0.74	0.57	0.68	0.81	0.78	0.85
Selective Attention	0.93	0.89	0.77	0.88	0.90	0.90	0.98
Self Awareness	0.67	0.75	0.34	0.60	0.81	0.78	0.81
Self Evaluation	0.69	0.76	0.28	0.58	0.83	0.78	0.69
Sequential Organization	0.95	0.91	0.82	0.91	0.93	0.92	0.98
Spatial Organization	0.55	0.51	0.47	0.50	0.58	0.55	0.44
Strategy Selection	0.70	0.74	0.53	0.69	0.83	0.78	0.90
Systematic Regional Synthesis	0.35	0.46	0.27	0.37	0.54	0.54	0.36
Temporal Organization	0.67	0.63	0.56	0.60	0.75	0.69	0.64
Verification	0.68	0.69	0.34	0.57	0.79	0.72	0.66
Visual Foraging	0.75	0.77	0.57	0.74	0.87	0.84	0.85
Visual Reference Or Grounding	0.38	0.36	0.26	0.37	0.47	0.42	0.29

### Captioning & IF (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.44	0.39	0.44	0.45	0.40	0.46	0.49
Adaptive Detail Management	0.33	0.21	0.26	0.19	0.23	0.24	0.36
Arithmetic Calculation	0.08	0.08	0.07	0.06	0.08	0.07	0.05
Backtracking	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Backward Chaining	0.00	0.01	0.02	0.01	0.00	0.01	0.04
Causal Organization	0.13	0.12	0.16	0.10	0.10	0.11	0.15
Compositionality	0.63	0.53	0.67	0.56	0.62	0.60	0.67
Conceptual Level Processing	0.68	0.61	0.66	0.65	0.64	0.66	0.71
Context Alignment	0.39	0.32	0.46	0.38	0.32	0.40	0.44
Context Awareness	0.25	0.20	0.31	0.26	0.20	0.26	0.33
Decomposition And Integration	0.45	0.44	0.56	0.43	0.50	0.45	0.50
Forward Chaining	0.44	0.48	0.60	0.47	0.49	0.47	0.51
Goal Management	0.19	0.23	0.53	0.21	0.28	0.21	0.31
Hierarchical Organization	0.44	0.33	0.43	0.33	0.40	0.40	0.47
Knowledge Structure Alignment	0.71	0.65	0.69	0.70	0.66	0.72	0.74
Logical Coherence	0.74	0.69	0.78	0.75	0.71	0.73	0.74
Mental Imagery Simulation	0.17	0.10	0.14	0.12	0.14	0.15	0.18
Network Organization	0.57	0.41	0.44	0.48	0.43	0.50	0.55
Ordinal Organization	0.31	0.29	0.35	0.31	0.29	0.30	0.30
Pattern Recognition	0.40	0.36	0.36	0.45	0.37	0.42	0.41
Perception Then Reasoning	0.45	0.36	0.37	0.34	0.39	0.34	0.40
Productivity	0.35	0.25	0.34	0.24	0.24	0.31	0.38
Representational Restructuring	0.31	0.22	0.29	0.26	0.20	0.26	0.31
Selective Attention	0.51	0.53	0.66	0.55	0.55	0.56	0.60
Self Awareness	0.10	0.07	0.11	0.07	0.07	0.08	0.15
Self Evaluation	0.04	0.04	0.05	0.02	0.04	0.04	0.05
Sequential Organization	0.67	0.60	0.69	0.59	0.61	0.64	0.69
Spatial Organization	0.47	0.31	0.40	0.37	0.36	0.36	0.41
Strategy Selection	0.15	0.10	0.27	0.14	0.13	0.12	0.21
Systematic Regional Synthesis	0.09	0.12	0.14	0.05	0.13	0.07	0.10
Temporal Organization	0.42	0.39	0.42	0.42	0.39	0.39	0.47
Verification	0.07	0.03	0.07	0.07	0.06	0.06	0.08
Visual Foraging	0.32	0.28	0.35	0.29	0.28	0.29	0.35
Visual Reference Or Grounding	0.20	0.11	0.16	0.14	0.14	0.15	0.14

Chart & OCR (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.28	0.28	0.18	0.22	0.35	0.29	0.23
Adaptive Detail Management	0.56	0.54	0.35	0.48	0.64	0.64	0.47
Arithmetic Calculation	0.52	0.57	0.40	0.44	0.53	0.61	0.39
Backtracking	0.21	0.29	0.09	0.19	0.30	0.46	0.19
Backward Chaining	0.05	0.07	0.06	0.06	0.08	0.12	0.11
Causal Organization	0.47	0.46	0.33	0.40	0.56	0.49	0.42
Compositionality	0.95	0.96	0.90	0.90	0.97	0.98	0.95
Conceptual Level Processing	0.54	0.51	0.37	0.45	0.61	0.48	0.44
Context Alignment	0.66	0.65	0.41	0.62	0.79	0.71	0.59
Context Awareness	0.39	0.28	0.21	0.29	0.50	0.35	0.43
Decomposition And Integration	0.72	0.73	0.72	0.72	0.76	0.80	0.77
Forward Chaining	0.95	0.98	0.96	0.99	0.98	0.98	0.98
Goal Management	0.94	0.98	0.95	1.00	0.99	0.99	1.00
Hierarchical Organization	0.63	0.61	0.48	0.57	0.72	0.61	0.58
Knowledge Structure Alignment	0.93	0.91	0.75	0.91	0.97	0.92	0.83
Logical Coherence	1.00	1.00	0.99	1.00	1.00	1.00	1.00
Mental Imagery Simulation	0.43	0.35	0.34	0.37	0.39	0.43	0.38
Network Organization	0.38	0.36	0.23	0.32	0.45	0.35	0.30
Ordinal Organization	0.71	0.73	0.68	0.74	0.74	0.75	0.71
Pattern Recognition	0.42	0.50	0.43	0.49	0.55	0.51	0.41
Perception Then Reasoning	0.96	0.94	0.83	0.84	0.96	0.95	0.82
Productivity	0.22	0.25	0.12	0.20	0.28	0.31	0.20
Representational Restructuring	0.32	0.35	0.23	0.28	0.40	0.39	0.33
Selective Attention	1.00	1.00	0.99	1.00	1.00	1.00	1.00
Self Awareness	0.70	0.72	0.39	0.50	0.80	0.78	0.61
Self Evaluation	0.95	0.96	0.39	0.68	0.98	0.98	0.63
Sequential Organization	1.00	1.00	0.99	1.00	1.00	1.00	1.00
Spatial Organization	0.68	0.63	0.66	0.69	0.68	0.71	0.68
Strategy Selection	0.55	0.55	0.44	0.56	0.70	0.64	0.69
Systematic Regional Synthesis	0.87	0.91	0.68	0.78	0.90	0.95	0.74
Temporal Organization	0.78	0.75	0.73	0.75	0.84	0.82	0.80
Verification	0.98	0.95	0.65	0.86	0.99	0.98	0.79
Visual Foraging	0.92	0.97	0.81	0.92	0.97	0.98	0.88
Visual Reference Or Grounding	0.93	0.90	0.86	0.91	0.93	0.91	0.86

Chart & OCR (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.19	0.08	0.07	0.13	0.10	0.13	0.12
Adaptive Detail Management	0.20	0.05	0.08	0.07	0.06	0.10	0.11
Arithmetic Calculation	0.25	0.38	0.13	0.26	0.30	0.26	0.14
Backtracking	0.01	0.00	0.01	0.01	0.01	0.01	0.01
Backward Chaining	0.01	0.02	0.01	0.01	0.03	0.02	0.22
Causal Organization	0.14	0.09	0.07	0.08	0.11	0.11	0.16
Compositionality	0.78	0.69	0.26	0.55	0.67	0.63	0.70
Conceptual Level Processing	0.44	0.15	0.12	0.18	0.18	0.26	0.26
Context Alignment	0.30	0.09	0.26	0.12	0.13	0.20	0.22
Context Awareness	0.06	0.02	0.36	0.04	0.04	0.04	0.08
Decomposition And Integration	0.53	0.57	0.17	0.30	0.54	0.39	0.51
Forward Chaining	0.86	0.80	0.42	0.69	0.78	0.83	0.86
Goal Management	0.48	0.57	0.37	0.19	0.56	0.42	0.89
Hierarchical Organization	0.37	0.19	0.10	0.11	0.24	0.18	0.28
Knowledge Structure Alignment	0.72	0.30	0.22	0.31	0.31	0.52	0.40
Logical Coherence	0.97	0.89	0.96	0.90	0.90	0.95	0.95
Mental Imagery Simulation	0.25	0.04	0.11	0.04	0.06	0.10	0.19
Network Organization	0.25	0.05	0.05	0.05	0.08	0.11	0.12
Ordinal Organization	0.70	0.57	0.42	0.63	0.56	0.57	0.52
Pattern Recognition	0.55	0.34	0.21	0.51	0.33	0.45	0.34
Perception Then Reasoning	0.84	0.52	0.12	0.30	0.55	0.55	0.37
Productivity	0.07	0.01	0.02	0.02	0.02	0.02	0.04
Representational Restructuring	0.23	0.05	0.07	0.09	0.05	0.13	0.15
Selective Attention	0.97	0.89	0.93	0.90	0.90	0.95	0.96
Self Awareness	0.17	0.07	0.13	0.06	0.10	0.11	0.16
Self Evaluation	0.10	0.04	0.07	0.17	0.07	0.06	0.06
Sequential Organization	0.91	0.85	0.65	0.69	0.80	0.87	0.92
Spatial Organization	0.64	0.28	0.49	0.31	0.33	0.42	0.56
Strategy Selection	0.17	0.06	0.17	0.08	0.11	0.14	0.28
Systematic Regional Synthesis	0.30	0.34	0.07	0.04	0.30	0.19	0.19
Temporal Organization	0.65	0.50	0.31	0.44	0.51	0.52	0.55
Verification	0.30	0.17	0.18	0.29	0.15	0.18	0.19
Visual Foraging	0.64	0.49	0.24	0.31	0.48	0.47	0.63
Visual Reference Or Grounding	0.86	0.45	0.84	0.50	0.59	0.73	0.65

Spatial & Action (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.41	0.45	0.36	0.39	0.45	0.43	0.41
Adaptive Detail Management	0.89	0.88	0.72	0.83	0.90	0.94	0.82
Arithmetic Calculation	0.15	0.15	0.10	0.13	0.14	0.17	0.13
Backtracking	0.25	0.33	0.15	0.21	0.33	0.52	0.27
Backward Chaining	0.07	0.09	0.11	0.08	0.10	0.14	0.19
Causal Organization	0.67	0.67	0.64	0.61	0.73	0.71	0.73
Compositionality	1.00	0.99	0.98	0.99	0.99	1.00	0.99
Conceptual Level Processing	0.81	0.78	0.70	0.76	0.84	0.77	0.75
Context Alignment	0.87	0.89	0.78	0.83	0.91	0.91	0.85
Context Awareness	0.65	0.55	0.47	0.54	0.77	0.63	0.63
Decomposition And Integration	0.94	0.93	0.94	0.93	0.93	0.96	0.96
Forward Chaining	0.99	0.98	0.99	0.99	0.99	0.99	1.00
Goal Management	0.95	0.97	0.97	0.98	0.99	0.99	1.00
Hierarchical Organization	0.93	0.91	0.83	0.88	0.91	0.89	0.88
Knowledge Structure Alignment	0.99	0.97	0.91	0.96	0.99	0.97	0.94
Logical Coherence	1.00	0.99	0.99	1.00	1.00	1.00	1.00
Mental Imagery Simulation	0.92	0.92	0.81	0.87	0.93	0.94	0.84
Network Organization	0.76	0.75	0.67	0.70	0.77	0.75	0.71
Ordinal Organization	0.75	0.77	0.77	0.77	0.79	0.76	0.80
Pattern Recognition	0.50	0.64	0.55	0.58	0.55	0.62	0.56
Perception Then Reasoning	0.95	0.83	0.68	0.71	0.89	0.87	0.65
Productivity	0.36	0.43	0.27	0.33	0.45	0.54	0.37
Representational Restructuring	0.63	0.69	0.53	0.57	0.66	0.73	0.61
Selective Attention	1.00	0.99	1.00	1.00	1.00	1.00	1.00
Self Awareness	0.80	0.81	0.55	0.63	0.90	0.90	0.73
Self Evaluation	0.96	0.96	0.60	0.81	0.98	0.99	0.82
Sequential Organization	1.00	0.99	1.00	1.00	1.00	1.00	1.00
Spatial Organization	1.00	0.99	0.98	0.99	0.99	0.99	0.98
Strategy Selection	0.81	0.84	0.71	0.79	0.87	0.89	0.81
Systematic Regional Synthesis	0.93	0.93	0.76	0.84	0.96	0.98	0.84
Temporal Organization	0.67	0.68	0.67	0.65	0.70	0.75	0.71
Verification	0.99	0.96	0.85	0.95	0.99	0.99	0.95
Visual Foraging	0.96	0.98	0.87	0.93	0.99	0.99	0.91
Visual Reference Or Grounding	0.99	0.97	0.95	0.97	0.98	0.98	0.95

Spatial & Action (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.18	0.08	0.14	0.13	0.08	0.14	0.18
Adaptive Detail Management	0.30	0.13	0.26	0.22	0.15	0.20	0.28
Arithmetic Calculation	0.06	0.08	0.03	0.06	0.07	0.07	0.04
Backtracking	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Backward Chaining	0.03	0.02	0.01	0.03	0.03	0.04	0.05
Causal Organization	0.32	0.26	0.26	0.32	0.31	0.35	0.44
Compositionality	0.92	0.75	0.53	0.85	0.82	0.87	0.87
Conceptual Level Processing	0.55	0.26	0.36	0.49	0.34	0.40	0.52
Context Alignment	0.50	0.18	0.49	0.49	0.23	0.39	0.48
Context Awareness	0.12	0.06	0.34	0.20	0.06	0.13	0.15
Decomposition And Integration	0.73	0.66	0.43	0.87	0.64	0.69	0.75
Forward Chaining	0.89	0.90	0.72	0.98	0.92	0.96	0.89
Goal Management	0.51	0.52	0.52	0.76	0.44	0.41	0.69
Hierarchical Organization	0.70	0.46	0.47	0.67	0.53	0.58	0.66
Knowledge Structure Alignment	0.72	0.46	0.62	0.70	0.49	0.64	0.70
Logical Coherence	0.99	0.93	0.99	0.99	0.95	0.97	0.98
Mental Imagery Simulation	0.69	0.38	0.35	0.60	0.45	0.55	0.54
Network Organization	0.56	0.17	0.23	0.43	0.21	0.25	0.47
Ordinal Organization	0.67	0.64	0.60	0.74	0.62	0.65	0.62
Pattern Recognition	0.54	0.40	0.29	0.59	0.40	0.64	0.45
Perception Then Reasoning	0.81	0.45	0.24	0.60	0.47	0.33	0.44
Productivity	0.08	0.01	0.05	0.05	0.01	0.03	0.06
Representational Restructuring	0.38	0.14	0.18	0.27	0.16	0.35	0.27
Selective Attention	0.98	0.96	0.99	0.99	0.97	0.99	0.98
Self Awareness	0.22	0.10	0.17	0.21	0.13	0.10	0.25
Self Evaluation	0.25	0.12	0.18	0.29	0.15	0.12	0.15
Sequential Organization	0.97	0.88	0.85	0.99	0.89	0.95	0.93
Spatial Organization	0.99	0.95	0.95	0.97	0.96	0.97	0.97
Strategy Selection	0.32	0.17	0.34	0.36	0.18	0.27	0.37
Systematic Regional Synthesis	0.32	0.45	0.21	0.40	0.32	0.26	0.27
Temporal Organization	0.69	0.59	0.46	0.76	0.68	0.63	0.64
Verification	0.49	0.28	0.39	0.50	0.31	0.38	0.34
Visual Foraging	0.61	0.39	0.36	0.59	0.34	0.48	0.60
Visual Reference Or Grounding	0.92	0.73	0.82	0.84	0.78	0.82	0.78

STEM (Qwen3)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.77	0.76	0.73	0.73	0.79	0.77	0.75
Adaptive Detail Management	0.81	0.82	0.69	0.77	0.85	0.84	0.77
Arithmetic Calculation	0.61	0.61	0.56	0.59	0.62	0.63	0.57
Backtracking	0.46	0.59	0.36	0.46	0.59	0.70	0.47
Backward Chaining	0.18	0.26	0.14	0.21	0.30	0.32	0.26
Causal Organization	0.87	0.87	0.81	0.84	0.90	0.87	0.86
Compositionality	1.00	1.00	0.99	0.99	1.00	1.00	1.00
Conceptual Level Processing	0.89	0.89	0.85	0.87	0.92	0.88	0.88
Context Alignment	0.83	0.84	0.74	0.79	0.88	0.86	0.81
Context Awareness	0.46	0.40	0.34	0.38	0.55	0.48	0.44
Decomposition And Integration	0.91	0.92	0.89	0.91	0.92	0.92	0.91
Forward Chaining	0.98	0.99	0.98	0.99	0.99	0.98	0.98
Goal Management	0.97	0.98	0.93	0.98	0.99	0.99	0.99
Hierarchical Organization	0.92	0.91	0.88	0.90	0.94	0.90	0.90
Knowledge Structure Alignment	0.96	0.95	0.90	0.94	0.98	0.96	0.93
Logical Coherence	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Mental Imagery Simulation	0.37	0.40	0.23	0.33	0.47	0.47	0.31
Network Organization	0.72	0.73	0.64	0.68	0.76	0.73	0.69
Ordinal Organization	0.76	0.82	0.70	0.80	0.86	0.82	0.81
Pattern Recognition	0.67	0.74	0.62	0.71	0.75	0.76	0.67
Perception Then Reasoning	0.91	0.89	0.78	0.84	0.93	0.90	0.81
Productivity	0.56	0.61	0.46	0.56	0.61	0.66	0.57
Representational Restructuring	0.74	0.77	0.65	0.70	0.78	0.80	0.73
Selective Attention	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Self Awareness	0.77	0.86	0.62	0.73	0.87	0.88	0.79
Self Evaluation	0.97	0.97	0.76	0.91	0.99	0.99	0.91
Sequential Organization	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Spatial Organization	0.63	0.64	0.60	0.63	0.66	0.67	0.62
Strategy Selection	0.86	0.87	0.78	0.84	0.90	0.88	0.88
Systematic Regional Synthesis	0.74	0.82	0.59	0.75	0.85	0.86	0.71
Temporal Organization	0.65	0.69	0.57	0.66	0.78	0.75	0.67
Verification	0.99	0.96	0.87	0.94	0.99	0.99	0.95
Visual Foraging	0.82	0.89	0.72	0.83	0.92	0.91	0.82
Visual Reference Or Grounding	0.54	0.51	0.39	0.48	0.58	0.57	0.42

STEM (Qwen2.5)

	Captioning & IF	Chart & OCR	Ground & Search	Knowledge & Recog.	Spatial & Action	STEM	Mix All
Abstraction	0.59	0.38	0.54	0.52	0.43	0.53	0.57
Adaptive Detail Management	0.36	0.09	0.29	0.22	0.15	0.21	0.26
Arithmetic Calculation	0.45	0.48	0.40	0.42	0.46	0.47	0.41
Backtracking	0.04	0.04	0.04	0.03	0.04	0.03	0.03
Backward Chaining	0.03	0.01	0.01	0.02	0.02	0.02	0.05
Causal Organization	0.61	0.45	0.56	0.53	0.47	0.57	0.62
Compositionality	0.97	0.87	0.85	0.91	0.92	0.92	0.94
Conceptual Level Processing	0.80	0.57	0.74	0.69	0.64	0.72	0.75
Context Alignment	0.53	0.26	0.54	0.41	0.31	0.43	0.44
Context Awareness	0.19	0.08	0.26	0.16	0.12	0.15	0.17
Decomposition And Integration	0.80	0.69	0.69	0.73	0.75	0.71	0.76
Forward Chaining	0.96	0.89	0.88	0.96	0.93	0.95	0.96
Goal Management	0.71	0.57	0.75	0.71	0.67	0.59	0.81
Hierarchical Organization	0.80	0.58	0.72	0.67	0.68	0.71	0.78
Knowledge Structure Alignment	0.81	0.64	0.82	0.71	0.69	0.77	0.75
Logical Coherence	0.98	0.94	0.97	0.97	0.95	0.97	0.97
Mental Imagery Simulation	0.09	0.03	0.06	0.06	0.03	0.05	0.06
Network Organization	0.49	0.19	0.39	0.37	0.27	0.35	0.39
Ordinal Organization	0.58	0.46	0.52	0.54	0.46	0.50	0.51
Pattern Recognition	0.56	0.38	0.44	0.52	0.38	0.49	0.46
Perception Then Reasoning	0.66	0.42	0.37	0.38	0.51	0.44	0.44
Productivity	0.17	0.04	0.10	0.12	0.06	0.09	0.13
Representational Restructuring	0.44	0.20	0.30	0.37	0.22	0.35	0.36
Selective Attention	0.97	0.93	0.98	0.96	0.94	0.97	0.97
Self Awareness	0.33	0.16	0.25	0.22	0.21	0.20	0.28
Self Evaluation	0.31	0.18	0.24	0.22	0.20	0.20	0.25
Sequential Organization	0.99	0.94	0.96	0.96	0.96	0.97	0.98
Spatial Organization	0.58	0.43	0.48	0.50	0.45	0.51	0.51
Strategy Selection	0.54	0.32	0.62	0.47	0.36	0.47	0.52
Systematic Regional Synthesis	0.33	0.26	0.21	0.21	0.29	0.28	0.27
Temporal Organization	0.38	0.32	0.31	0.43	0.35	0.33	0.37
Verification	0.53	0.31	0.43	0.40	0.37	0.42	0.44
Visual Foraging	0.56	0.33	0.38	0.41	0.36	0.40	0.45
Visual Reference Or Grounding	0.32	0.16	0.25	0.20	0.18	0.21	0.21

A12 Vero examples