ByteDance

BANDAI

# Thinking with Programming Vision:
# Towards a Unified View for Thinking with Images

**Zirun Guo**[1,2]    **Minjie Hong**[1]    **Feng Zhang**[2]    **Kai Jia**[2]    **Tao Jin**[1]

[1] **Zhejiang University**    [2] **ByteDance, BandAI**

## Abstract

Multimodal large language models (MLLMs) that "think with images" can interactively use tools to reason about visual inputs, but current approaches often rely on a narrow set of tools with limited real-world necessity and scalability. In this work, we first reveal a critical and previously overlooked weakness: even state-of-the-art MLLMs are surprisingly brittle, showing significant performance degradation on images with simple orientation changes or natural corruptions, underscoring the need for more robust tool-based reasoning. To address this, we propose *CodeVision*, a flexible and scalable "code-as-tool" framework where the model generates code as a universal interface to invoke any image operation, moving beyond fixed tool registries. We train our model using a two-stage methodology, beginning with Supervised Fine-Tuning (SFT) on a high-quality dataset curated for complex, multi-turn tool composition and error recovery, followed by Reinforcement Learning (RL) with a novel and dense process reward function to encourage strategic and efficient tool use. To facilitate this research, we construct new SFT and RL datasets and introduce a challenging new benchmark suite designed to rigorously evaluate robustness to orientation changes and multi-tool reasoning. Experiments on Qwen2.5-VL and Qwen3-VL series show that our approach significantly improves model performance and fosters emergent capabilities such as flexible tool composition, efficient chained execution, and robust error recovery from runtime feedback.

✉ **Email:** zrguo.cs@gmail.com    ⌂ **Code:** https://github.com/ByteDance-BandAI/CodeVision

## 1. Introduction

Thinking with images (OpenAI, 2025c) equips multimodal large language models (MLLMs) with interactive, tool-augmented reasoning over visual inputs. Rather than passively describing an image, the model actively manipulates it using tools (*e.g.,* "zoom in", OCR) to acquire the evidence needed for reliable reasoning. Despite rapid progress, current research still faces three limitations:
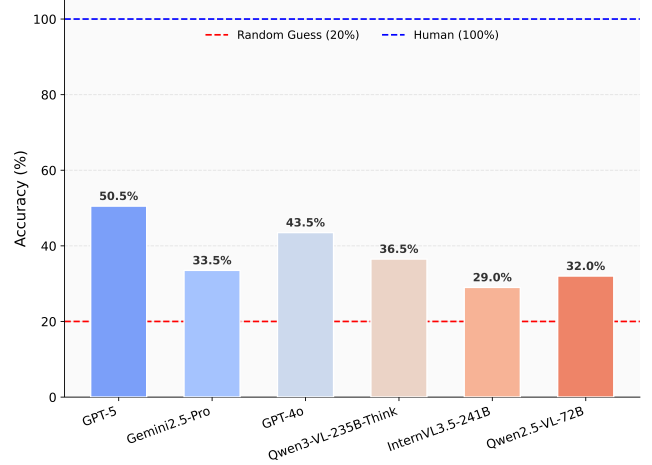
① **The necessity of tools**: Current methods mainly emphasize the "crop" tool, which zooms in on regions of interest for clearer observation, and evaluate tool ability on benchmarks such as V* (Wu and Xie, 2024) and HRBench (Wang et al., 2025b). Yet the benefit is marginal: using tools often yields only 2–5% accuracy gains, and reinforcement learning (RL) without tools can match those results. This suggests that the potential—and necessity—of tools is not being fully exercised by existing tasks.

② **Flexibility and scalability**: Methods frequently require manually specifying tool names and arguments (Su et al., 2025b), which is brittle and not scalable. Even renaming a tool (*e.g.,* `crop` → `zoomin`) can necessitate retraining, hindering generalization to new tools and argument schemas.

③ **Multi-turn, multi-tool use**: Many systems support only a single tool or a few tools within a single turn. While several works (Lai et al., 2025; Zhang et al., 2025) explore multi-turn settings, they largely focus on repeated cropping rather than composing different tools across multiple turns, which is what real-world tasks often require.
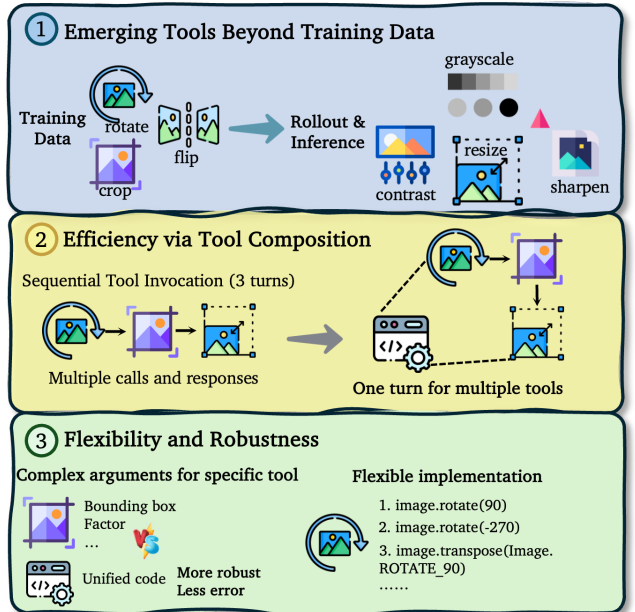
Based on these observations, we propose a new framework to address all three challenges.

To create scenarios where tools are genuinely required, we examine ubiquitous real-world perturbations: incorrect image orientation due to landscape/portrait capture and mirrored selfies that confound text recognition. Compared with cropping, restoring the canonical orientation is often strictly more necessary for downstream recognition and reasoning. We conduct a simple diagnostic: starting from 200 images from various domains, we uniformly apply one of five transformations—rotation by 90/180/270 degrees, horizontal flip, or vertical flip—and ask models to identify the transformation (a five-way multiple-choice question). As illustrated in Figure 1, even state-of-the-art models like GPT-5 (OpenAI, 2025b) and Gemini2.5 Pro (Comanici et al., 2025) perform poorly, whereas humans achieve 100% accuracy easily. Remarkably, as summarized in Table 1, simple rotation/flip operations can reduce model performance by up to 80%, revealing substantial brittleness of current MLLMs and underscoring the necessity of using tools for these real-world corruptions.



**Figure 1:** *Diagnostic results on image orientation identification.*

To address the second challenge, inspired by OpenAI o3 (OpenAI, 2025c), we treat *code* itself as a unified tool: the model writes code that invokes whatever image operations are needed. This eliminates hand-crafted tool name/argument specifications and dramatically improves generalization: the model is no longer restricted to a fixed registry but can call an effectively unbounded set of tools through code. While the paradigm of using code as a tool has been explored previously (Chen et al., 2023; Zhang et al., 2025), we unlock three notable advantages through our principled RL training and rigorous dataset construction (see Figure 2): (1) *Emergence of new tools*: the model calls tools that never appeared in the RL training data to solve novel problems; (2) *Efficiency*: the model chains multiple tools within a single execution; (3) *Robustness*: the model leverages runtime error messages and outputs to revise code, improving fail-



**Figure 2:** *Three advantages of CodeVision we observe in the training and inference stage.*

ure recovery and out-of-distribution generalization.

To target the third challenge, we design datasets and a benchmark that require composing multiple tools across multiple turns, and we introduce dense process rewards that encourage effective tool selection and steady progress throughout the dialogue. This setting better reflects realistic problem solving and promotes the development of models that plan, adapt, and utilize different tools across steps. In summary, our contributions are threefold:

- We identify a critical brittleness in state-of-the-art MLLMs and propose a flexible "code-as-tool" framework that treats code as a universal interface for tool invocation, enhancing scalability and generalization.

- We construct high-quality SFT and RL datasets focused on multi-turn, multi-tool composition and error handling, and introduce three new benchmarks covering both single- and multi-tool scenarios to rigorously evaluate model robustness and complex tool use.

- Our experiments demonstrate that our approach significantly improves model performance on these challenging benchmarks and validates the benefits of the code-as-tool paradigm, which enables emergent tool use, efficient tool chaining, and robust error recovery.

## 2. Related Work

**Thinking with Images.** OpenAI's o3 (OpenAI, 2025c) popularized the idea of *thinking with images*, where MLLMs actively operate on images via tools. Most follow-ups focus on the crop/zoom tool (Zheng et al., 2025b; Lai et al., 2025; Su et al., 2025a; Fan et al., 2025), in part because grounding-rich pretraining eases integration (Bai et al., 2025b). However, cropping is often unnecessary or the benefits are hard to observe for many tasks because similar performance can be achieved without it. Meanwhile, broader tool suites (*e.g.,* line drawing, OCR, segmentation) (Su et al., 2025b; Zhang et al., 2025; Liu et al., 2025) have been investigated, but their effectiveness has not been sufficiently validated because they are not evaluated on tasks that truly *require* tools.

**Tool Integration.** Equipping large models with tools (Comanici et al., 2025; OpenAI, 2025c,b) mitigates intrinsic limitations by enabling external capabilities such as search, code execution, and generative models. For instance, LLM-I (Guo et al., 2025b) leverages tools such as search, code, and diffusion models, enabling a pure language model to produce multimodal content, and recent work integrates web search to support fact-grounded reasoning (Jin et al., 2025; Li et al., 2025). DeepResearch (Team et al., 2025; OpenAI, 2025a) is a more advanced tool-integrated reasoning paradigm, which requires multi-turn search, data collection, and reasoning and achieves state-of-the-art performance on various difficult benchmarks. These directions highlight the promise of tool-augmented agents, yet most image-centric systems still rely on narrow, pre-registered tool sets with hand-specified interfaces.

**MLLM Reasoning.** Reinforcement learning has become central to strengthening reasoning in LLMs, from PPO (Schulman et al., 2017) to more recent GRPO (Shao et al., 2024), DAPO (Yu et al., 2025), and GSPO (Zheng et al., 2025a). For MLLMs, reasoning spans both text and vision, such as chart reasoning (Wang et al., 2025a; Masry et al., 2025); recent work (Guo et al., 2025a) emphasizes careful visual inspection prior to reasoning via an explicit observe tag. Similarly, a series of works (Hong et al., 2025; Gao et al., 2025) have explored enhancing reasoning in MLLMs through advanced policy optimization. Recent efforts are transitioning from thinking *about* images to thinking *with*
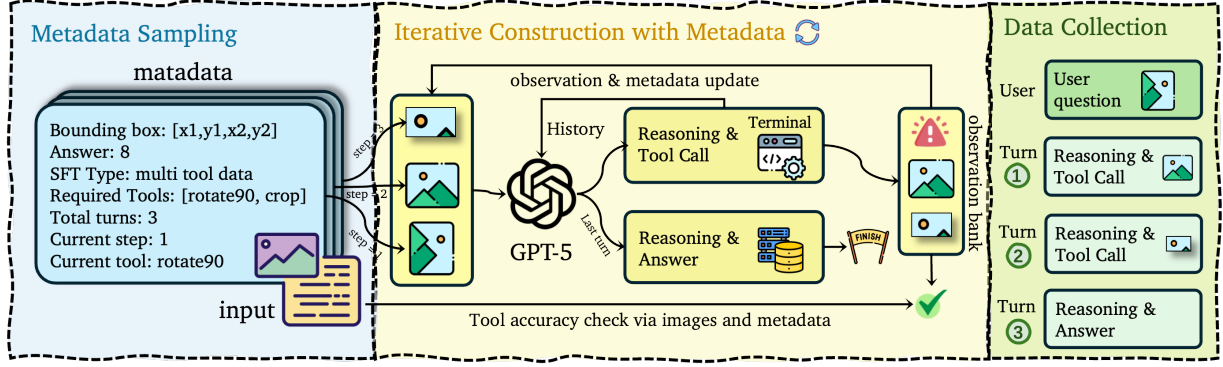
**Figure 3:** *Pipeline for cold-start SFT data construction.*

images (OpenAI, 2025c; Su et al., 2025b; Zhang et al., 2025), but robust, multi-turn, multi-tool composition and strong generalization to unseen tools remain open challenges.

## 3. Methodology

### 3.1. Overview

We train an agent to use tools robustly and efficiently through a two-stage process. We begin with a cold-start phase, where Supervised Fine-Tuning (SFT) teaches the model the fundamental syntax and patterns of tool use on a diverse, high-quality dataset that explicitly covers challenging scenarios such as multi-tool sequences, error handling, and coarse-to-fine localization. Following SFT, we transition to a Reinforcement Learning (RL) phase to move beyond pattern imitation and teach the model to strategize. In this stage, the agent learns to make deliberate, efficient, and robust decisions via a dense, multi-component reward function designed to encourage the use of necessary tools, facilitate the discovery of beneficial strategies, and penalize reward-hacking behaviors, ensuring the final model is effective and deliberate in its tool application. We next detail the cold-start data and training, followed by the RL setting, including data construction and reward design.

### 3.2. Cold Start
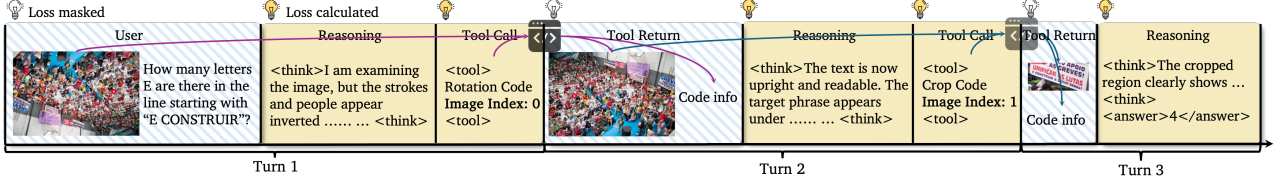
#### 3.2.1. Dataset Construction

The overall pipeline for cold-start SFT data construction is shown in Figure 3.

We build an SFT corpus that couples tool-use trajectories with final answers by combining diverse sources, explicit metadata, and programmatic verification. As noted in Figure 1, even state-of-the-art models struggle to reliably determine which tool is needed, so we guide data generation with structured metadata and automatic checks.

**Data sources across domains.** To increase diversity, we aggregate examples from multiple domains: handwriting datasets (Marti and Bunke, 2002; Liu et al., 2011), in-the-wild OCR/VQA datasets for everyday scenes (Long et al., 2022), table and chart understanding datasets (Masry et al., 2023), and math reasoning datasets (Qiao et al., 2025).

**Task types and metadata.** For each sample, we create metadata containing the ground-truth answer and a target *type* from five categories: `single-tool`, `multi-tool`, `multi-crop`, `error-handling`, and `no-tool`. We sample types according to preset proportions and, conditioned on the type, uniformly sample candidate tools. For crop-type data, we select text regions annotated in the source dataset whose areas are at most 0.01% of the image, ensuring that cropping is necessary for successful

4

**Figure 4:** *Rollout and inference process, and token masking used during SFT/RL.*

recognition and reasoning. For the multi-crop type, we simulate coarse-to-fine zoom-in by enforcing monotonically shrinking, spatially contiguous crop windows across steps to strengthen spatial awareness and localization. For the error-handling type, we intentionally surface failures, such as using an incorrect tool or triggering code and runtime errors (*e.g.*, invalid arguments, missing imports), and require the model to read error logs, revise its code, and retry—potentially switching to the correct tool—with limited retries before fallback.

**Metadata-conditioned image transformations.** After sampling the metadata, we transform the canonical (unaltered) image according to the selected tool to construct the model's initial observation. For example, if the required tool is `rotate-180`, we rotate the original image by 180 degrees and use this transformed image as GPT-5's initial input. This setup makes the subsequent tool invocation necessary to recover the canonical view.

**SFT generation with GPT-5.** Given the question, the (possibly transformed) image, and the metadata, GPT-5 produces a step-by-step reasoning trace and an *action* at each turn. Actions may include calling a tool or emitting an intermediate or final answer. When an action calls a tool, we execute it in a controlled runtime to obtain an updated image. We then compare the output against the canonical, untransformed reference. If they agree, the tool call is marked correct. Otherwise, we either discard the trajectory or attempt targeted correction. For the error-handling type, we feed error logs back to the model to trigger code revision before any discard.

**Iterative multi-turn construction.** After each step, we update the metadata (step count, remaining/required tools) and resubmit the updated context to GPT-5, repeating until the per-type turn budget is reached or a `stop` action is produced. Following this pipeline, we construct approximately 5,000 high-quality SFT examples consisting of aligned reasoning, action sequences, and answers.

### 3.2.2. Training

We train the model with an SFT objective over our multi-turn trajectories. As shown in Figure 4, each training example is formatted as an interleaved dialogue with alternating user prompts, assistant reasoning/tool-call outputs, and recorded tool returns. During optimization, we *mask out* user tokens and tool-return tokens; only assistant tokens corresponding to chain-of-thought reasoning and tool-call specifications contribute to the loss. Concretely, we apply teacher-forced next-token prediction with a token-wise mask $m_t \in \{0, 1\}$ that is 1 for assistant reasoning/tool-call tokens and 0 otherwise.

$$\mathcal{L}_{\text{SFT}}(\theta) = -\sum_{t=1}^{T} m_t \, \log p_\theta(y_t \,|\, \mathbf{x}, \mathbf{y}_{<t}) , \tag{1}$$

where $\mathbf{x}$ denotes the dialogue context (including masked user inputs and tool returns), $\mathbf{y}_{<t}$ are preceding assistant tokens, and $m_t$ selects only assistant reasoning/tool-call tokens.

We do not execute tools online during SFT. All tool outputs used as subsequent context are the cached

5

results from dataset construction. Training therefore reduces to standard causal LM training with a masked loss over the assistant side across multiple turns (including the final answer), preserving the multi-turn structure while avoiding runtime variance from external tool execution.

### 3.3. Reinforcement Learning

To further enhance reasoning and generalization, we conduct RL training after the cold-start stage.

### 3.3.1. Dataset

Beyond the SFT sources, we augment RL training data with additional reasoning-heavy and perception samples (Mathew et al., 2021; Wang et al., 2025c; Singh et al., 2019) to increase coverage and diversity of reasoning patterns. We then perform difficulty filtering to focus optimization on informative instances. For each candidate item, we sample multiple rollouts and remove degenerate cases where trajectories are uniformly trivial (all-correct) or uninformative (all-incorrect). This concentrates training on items with signal for policy improvement. Furthermore, each item is annotated with a *must-use* tool field drawn from {`rotate90`, `rotate180`, `rotate270`, `flip-horizontal`, `flip-vertical`, `crop`}. For crop-required items, we additionally attach the bounding box of the target region. Items that require no tools set this field to `None`. These constraints encourage policies to learn when and how to invoke the right tool. Following this pipeline, we construct approximately 40,000 RL training items.

### 3.3.2. Reward Function

Training LLMs to use tools reliably is challenging in practice: training collapse and abnormal tool-call rates frequently occur. We therefore adopt a dense, rich, and multi-component reward to stabilize and guide RL. Let a trajectory be $\tau = (s_1, a_1, \ldots, s_T, a_T)$, where $s_t$ is the environment state (dialogue context, current image, and any tool returns) and $a_t$ is the agent action (reasoning/tool-call program or textual answer) at turn $t$. The total reward decomposes as

$$R_{\text{total}}(\tau) = R_{\text{outcome}}(\tau) + \beta_1 \sum_{t=1}^{T} R_{\text{strategy}}(a_t) - \beta_2 P_{\text{cost}}(\tau). \tag{2}$$

where $R_{\text{outcome}}(\tau)$ is the outcome reward consisting of the terminal accuracy signal and the formatting reward, $R_{\text{strategy}}(a_t)$ is the process reward for strategy supervision at turn $t$, $P_{\text{cost}}(\tau)$ is the constraint penalty for trajectory efficiency and avoiding reward hacking, and $\beta_1, \beta_2$ control the trade-off between these terms.

**Outcome Reward.** We use a terminal accuracy signal $r_{\text{acc}} \in \{0, +1\}$ assessed on the final answer, and a format reward $r_{\text{fmt}} \in \{0, +1\}$ for correct reasoning formatting (`<think>` and `<answer>` tags).

**Strategy Shaping.** Pure outcome rewards are sparse and often lead to brittle behavior (*e.g.*, abnormal tool-call rates). We therefore introduce process-level signals with two parts: (i) a *must-use* tool set that encodes task prerequisites, and (ii) *suggested* tools discovered on-the-fly that empirically improve success.

*(a) Must-use tools $S_{req}$.* Our RL metadata specifies a set of tools, $S_{\text{req}}$, that are prerequisites for solving the task robustly. For example, an image with incorrect orientation must be rotated first, regardless of the question. Let $N = |S_{\text{req}}|$ be the number of required tools. We assign a total reward budget of $1/N$ to each required tool. For categorical tools like "rotate" or "flip", this reward is given as a one-time bonus of $1/N$ when the tool is correctly used for the first time. For "crop", which has

a continuous quality measure, the reward is proportional to the Intersection-over-Union (IoU)—a measure of overlap between the predicted and target bounding boxes. To encourage refinement over multiple steps, we only reward the *improvement* in IoU over the best attempt so far. Additionally, we introduce a bonus for perfectly matching the required tool-use trajectory. This reward is granted only if the agent executes the complete sequence of tools specified in $S_{\text{req}}$ correctly and in the prescribed order, without any redundant or incorrect steps. This mechanism serves two purposes. First, for single-tool scenarios, it encourages precision by distinguishing between an immediate, correct tool call and a solution reached through trial-and-error. Second, for multi-tool tasks, it incentivizes the agent to learn and follow the optimal solution path, ensuring it explores and masters the entire sequence of required actions.

*(b) Suggested tools bonus.* Not all useful tools can be predefined. In addition to the must-use tools in $S_{\text{req}}$, we consider optional tools that are not required but can empirically improve success. This is where the power of our code-as-a-tool framework becomes critical, as it liberates the agent from a fixed tool registry and allows it to invoke any function or library dynamically. For example, applying contrast enhancement on a blurry image might be beneficial but is not a universal prerequisite. We reward the discovery of such emergent, helpful tools through a rollout comparison mechanism. For a given problem, we collect $K$ trajectories (with $K = 8$ in our experiments) from the current policy and partition them into those that use an optional tool, $G_{\text{tool}}$, and those that do not, $G_{\text{notool}}$. If the tool-using group shows a higher accuracy rate and the no-tool group achieves at most one success out of $K$ rollouts, it suggests that the additional tools will improve the success rate of solving this problem. We then compute this empirical performance gain as an "inferred tool necessity reward" for this problem:

$$r_{\text{nec}} = \max\left(0, \frac{\sum_{i \in G_{\text{tool}}} r_{acc}^i}{|G_{\text{tool}}|} - \frac{\sum_{i \in G_{\text{notool}}} r_{acc}^i}{|G_{\text{notool}}|}\right), \tag{3}$$

where $|G|$ is the number of trajectories in group $G$. This reward $r_{\text{nec}}$ is then added as a bonus to all successful trajectories that use the beneficial optional tool, encouraging the agent to explore and adopt various emergent tools that demonstrably improve performance. In addition, to further reward the use of extra tools, we introduce a per-trajectory bonus that is granted whenever the agent both invokes at least one optional tool and produces a correct final answer.

**Constraint Penalties.** While dense rewards are effective at guiding the agent, they can also be exploited. We observe emergent "reward hacking" behaviors during training, where the agent would learn to maximize rewards in ways that misalign with our goal of efficient and intelligent problem-solving. For instance, an agent might continue to crop an image to fractionally improve its IoU score long after the correct answer has been determined, or attempt to rotate an already correctly-oriented image to seek a reward. These actions lead to inefficient, low-quality trajectories. To counteract this, we introduce three targeted penalties that act as guardrails:

- **Turn Limit Penalty:** The primary motivation for this penalty is to prevent reward hacking and improve efficiency. We observe that the agent could learn to exhaustively call all possible required tools (*e.g.,* trying rotate90, rotate180, and rotate270 sequentially) simply to maximize its $R_{\text{strategy}}$ score, even if the first action already solves the orientation problem. To curb this behavior and promote deliberate, accurate solutions, we penalize trajectories that use more tool-call turns than necessary. A problem with $|S_{\text{req}}|$ required tools is given a buffer of one extra turn for error handling and exploration; any tool-call turn exceeding $|S_{\text{req}}| + 1$ incurs a penalty.

**Table 1:** *Main results on OCRBench (Perception) and ChartQAPro (Reasoning) under five types of transformations.*

| Model | OCRBench | | | | | | | ChartQAPro | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Source | Rot90 | Rot180 | Rot270 | Hori | Verti | Avg | Source | Rot90 | Rot180 | Rot270 | Hori | Verti | Avg |
| GPT-4o | 83.1 | 61.9 | 46.6 | 63.7 | 48.1 | 12.9 | 52.7 | 50.3 | 40.9 | 33.8 | 41.3 | 33.1 | 25.3 | 37.4 |
| Gemini2.5-Pro | 87.1 | 68.2 | 67.9 | 71.5 | 39.4 | 41.3 | 62.6 | 66.8 | 63.8 | 59.7 | 61.6 | 45.6 | 58.2 | 59.3 |
| Qwen2.5-VL-32B | 85.0 | 64.5 | 45.6 | 64.8 | 25.5 | 8.7 | 49.0 | 39.5 | 29.7 | 26.0 | 30.7 | 22.6 | 17.6 | 27.7 |
| Qwen2.5-VL-72B | 88.5 | 72.6 | 58.3 | 73.5 | 35.1 | 18.2 | 57.7 | 38.2 | 29.5 | 27.7 | 30.3 | 23.8 | 20.3 | 28.3 |
| Qwen3-VL-30B-Thinking | 87.6 | 72.1 | 60.5 | 66.6 | 49.3 | 31.8 | 61.3 | 50.2 | 38.7 | 36.5 | 38.2 | 31.3 | 29.0 | 37.3 |
| Qwen3-VL-235B-Thinking | 88.8 | 76.4 | 71.0 | 74.3 | 45.9 | 23.8 | 63.4 | 56.9 | 45.0 | 43.4 | 46.9 | 35.1 | 26.0 | 42.2 |
| InternVL3.5-30B | 88.7 | 45.8 | 7.2 | 45.8 | 5.9 | 6.2 | 33.3 | 37.2 | 27.2 | 18.0 | 27.9 | 17.5 | 18.0 | 24.3 |
| InternVL3.5-241B | 92.3 | 58.2 | 32.4 | 57.6 | 25.6 | 9.4 | 45.9 | 45.9 | 31.5 | 24.4 | 34.7 | 22.5 | 21.3 | 30.0 |
| Thyme | 86.3 | 67.0 | 51.9 | 67.8 | 27.2 | 13.3 | 52.3 | 30.3 | 23.8 | 20.4 | 22.9 | 17.8 | 14.4 | 21.6 |
| Qwen2.5-VL-7B | 86.4 | 70.2 | 58.0 | 71.7 | 32.4 | 17.0 | 56.0 | 37.3 | 23.4 | 22.2 | 23.7 | 19.5 | 20.1 | 24.4 |
| CodeVision-7B | 87.2 | 72.3 | 73.1 | 75.2 | 65.1 | 67.4 | 73.4 | 39.1 | 30.8 | 29.8 | 31.4 | 30.1 | 29.0 | 31.7 |
| Qwen3-VL-8B-Thinking | 82.4 | 64.8 | 57.6 | 58.4 | 35.7 | 14.5 | 52.2 | 47.2 | 32.8 | 30.0 | 31.9 | 21.3 | 13.9 | 29.5 |
| CodeVision-8B | 83.5 | 78.6 | 77.4 | 76.7 | 68.7 | 67.3 | 75.4 | 50.3 | 39.0 | 38.1 | 39.2 | 39.7 | 38.0 | 40.7 |
| Qwen3-VL-32B-Thinking | 86.6 | 67.8 | 64.3 | 66.4 | 35.1 | 13.8 | 55.7 | 52.3 | 39.9 | 38.3 | 42.0 | 26.1 | 18.7 | 36.2 |
| CodeVision-32B | 87.8 | 82.8 | 79.3 | 81.1 | 77.7 | 68.3 | 79.5 | 57.4 | 54.4 | 53.6 | 54.7 | 52.8 | 53.1 | 54.3 |

- **Poor Reasoning Penalty:** An agent must arrive at the right answer for the right reasons. If a trajectory yields the correct final answer but relies on a "crop" action with a negligible IoU ($< 0.1$) with the true region of interest, it indicates that the agent's reasoning was not grounded in the correct visual evidence. We penalize such low-quality trajectories to enforce a stronger link between visual evidence and the final conclusion.

- **Inappropriate Tool Use Penalty:** A core aspect of tool use is knowing when *not* to act. For samples where no tools are required ($S_{\mathrm{req}} = $ None), any attempt to use an orientation-adjusting tool is penalized. Applying such transformations to a normal image actively increases the task's difficulty and demonstrates a flawed, low-quality strategy.

Finally, the constraint penalty $P_{\mathrm{cost}}(\tau)$ consists of the sum of the above three penalties where each penalty $p \in \{0, +1\}$.

# 4. Experiments

## 4.1. MVToolBench

To rigorously evaluate a model's ability to compose multiple tools in sequence, we construct the MVToolBench. The construction process is designed to create challenging scenarios that necessitate the combined use of multiple tools. The pipeline consists of three main stages: data filtering, question generation, and multi-tool augmentation.

**Data Source and Filtering.** We build our benchmark upon the HierText dataset (Long et al., 2022), which provides rich annotations including bounding box coordinates and corresponding text for words, lines, and paragraphs. To ensure that our benchmark effectively tests the model's ability to perceive fine-grained details, we first perform a filtering step based on the relative size of the text annotations. For each image, we calculate the area of every annotated word, line, and paragraph as a percentage of the total image area. We then retain only the annotations that are exceptionally small. For instance, for word-level questions, only words occupying less than 0.01% of the total image area are selected as candidates. This filtering serves two critical purposes: 1) it guarantees a high level of difficulty, as the target text is not easily visible in the original view, and 2) it enforces a strong dependency on the "crop" tool, as zooming in becomes a prerequisite for successful recognition and reasoning.

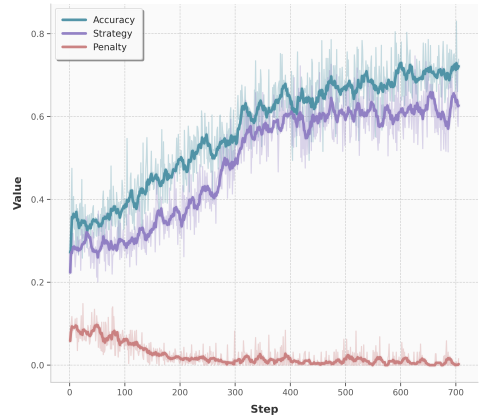**Table 2:** *Results on single-tool (V*, HRBench) and multi-tool (MVToolBench) benchmarks.*

| Model | V* | HRBench4k | HRBench8k | MVToolBench |
|---|---|---|---|---|
| GPT-4o | 67.9 | 65.0 | 60.1 | 8.5 |
| Gemini2.5-Pro | 83.8 | 86.2 | 85.1 | 32.6 |
| Qwen2.5-VL-32B | 81.9 | 73.8 | 70.5 | 16.4 |
| Qwen3-VL-30B-Thinking | 81.2 | 77.8 | 71.3 | 23.7 |
| Qwen3-VL-235B-Thinking | 85.9 | 84.3 | 76.6 | 30.1 |
| Thyme | 82.2 | 77.0 | 72.0 | 24.2 |
| Qwen2.5-VL-7B | 74.6 | 69.4 | 67.5 | 18.1 |
| CodeVision-7B | 83.7 | 75.6 | 72.2 | 60.1 |
| Qwen3-VL-8B-Thinking | 77.5 | 72.4 | 68.1 | 19.7 |
| CodeVision-8B | 82.4 | 77.1 | 73.4 | 62.7 |
| Qwen3-VL-32B-Thinking | 84.8 | 82.1 | 74.8 | 28.6 |
| CodeVision-32B | 86.2 | 84.3 | 76.1 | 65.4 |

**Question and Answer Generation.** Using the filtered, small-sized text annotations, we programmatically generate a variety of reasoning questions. These questions cover tasks such as text recognition (*e.g.,* "What does the line beginning with 'Busy' say?"), counting (*e.g.,* "How many times does the letter 'a' appear in the paragraph ending with 'Television'?"), and question answering that requires information retrieval from specific paragraphs. A crucial design principle during this stage is the deliberate avoidance of any explicit positional cues. We do not use phrases like "the word on the left" or provide bounding box coordinates in the prompts. This forces the model to rely entirely on its own grounding capabilities to first locate the relevant region and then perform an accurate crop, significantly increasing the challenge of the benchmark.

**Multi-Tool Augmentation.** After the question-answer pairs are finalized, we perform a final augmentation step on the images to create a multi-tool requirement. Each image is randomly subjected to one of several transformations: rotation (90, 180, or 270 degrees), horizontal flip, or vertical flip. This step ensures that solving a problem requires not only a "crop" operation but also an initial orientation-correction tool. By setting the proportion of each transformation to be equal, we maintain a balanced and unbiased dataset. This three-stage process results in a comprehensive benchmark that effectively evaluates the model's capacity for complex, multi-step tool composition.
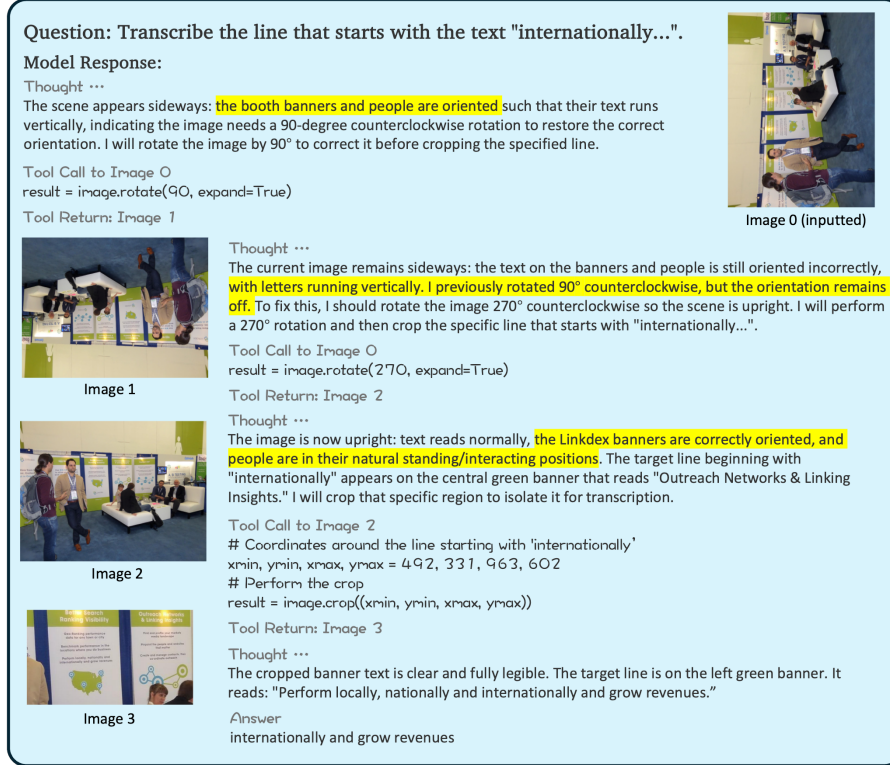
## 4.2. Implementation Details

**Training Settings:** To demonstrate the generalizability of our method, we use three distinct backbone models: Qwen2.5-VL-7B (Bai et al., 2025b), Qwen3-VL-8B-Thinking, and Qwen3-VL-32B-Thinking (Bai et al., 2025a). We employ GRPO (Shao et al., 2024) as the base RL algorithm. In the initial SFT stage, we train for 2 epochs with a batch size of 128, a learning rate of 5e-6, a cosine learning rate scheduler, and a warmup ratio of 0.05. Building on the SFT-tuned checkpoint, we then conduct the RL training phase for 2 epochs. For RL, the learning rate is set to 1e-6, with a batch size of 64, and we perform 8 rollouts for each sample. The coefficients for the format reward, strategy reward (must-use and suggested tools), and constraint penalties are set to 0.1, 1.0 (1.0 and 0.2), and 0.5, respectively. The KL divergence coefficient is set to 0.001.



**Figure 5:** *RL training curves for outcome, strategy, and total rewards. The consistent upward trend demonstrates that the agent effectively learns to use tools strategically to solve tasks.*

**Figure 6:** *A multi-turn example of error recovery. The model initially calls the wrong tool (`flip-horizontal`), but after receiving the execution result, it identifies the mistake and corrects it by applying the right tool (`rotate-90`).*

**Evaluation and Benchmarks:** Our evaluation is designed to test tool use across a range of scenarios. For the crop tool, we follow prior work and evaluate on V* (Wu and Xie, 2024), HRBench4k, and HRBench8k (Wang et al., 2025b). To assess the model's ability to handle orientation changes, we create a new benchmark suite by applying five transformations (rotate 90/180/270 degrees, and horizontal/vertical flip) to OCRBench (Liu et al., 2024) and ChartQAPro (Masry et al., 2025). We choose these two datasets to test distinct capabilities: OCRBench is a text recognition benchmark where high performance on the base version indicates that a model must accurately perceive the text to succeed. In contrast, ChartQAPro is a reasoning-based benchmark that requires understanding the chart's content, a task that does not depend on perfectly recognizing every character. This allows us to evaluate performance in both perception-critical and reasoning-heavy scenarios. Furthermore, to evaluate multi-tool proficiency, we utilize our newly created MVToolBench, the construction of which is detailed in Section 4.1.

## 4.3. Main Results

Table 1 shows the performance of our CodeVision models against a suite of strong MLLMs on benchmarks with challenging image orientations. On both OCRBench and ChartQAPro, baseline models exhibit a significant performance degradation when faced with rotated or flipped images, confirming the brittleness we identified. In contrast, our CodeVision-7B and CodeVision-8B models consistently and substantially outperform their base models and other strong competitors. For instance, on the transformed OCRBench subset, CodeVision-7B achieves an average score of 73.4, a +17.4 improvement over its base model, demonstrating its robust ability to recognize the correct transformation and apply the corresponding tool to restore the canonical image view for successful

reasoning.

Table 2 evaluates performance on established tool-use benchmarks. While our models perform competitively on single-tool tasks requiring cropping (V*, HRBench4k, HRBench8k), they establish a new state-of-the-art on MVToolBench, a benchmark designed for multi-tool reasoning. CodeVision-7B achieves a score of 60.1, nearly doubling the performance of the next-best model, Gemini2.5-Pro (32.6). This highlights our method's superior capability in composing multiple tools to solve complex problems.
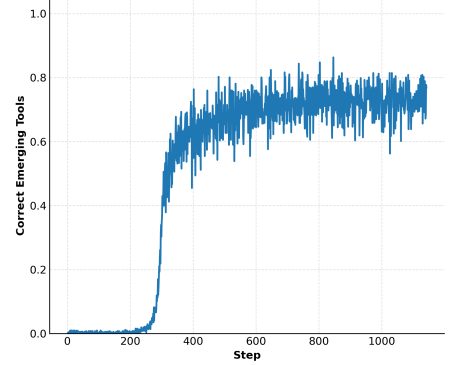
The effectiveness of our RL strategy is further illustrated in Figure 5, which plots the reward curves during training. The steady increase in all three reward components—outcome, strategy, and the total reward—indicates that the agent is not only learning to achieve the correct final answer but is also



**Figure 7:** *Reward curve for successful emergent tool use during RL training.*

mastering the strategic process of tool selection and application as guided by our dense reward function. Additionally, Figure 7 presents the reward curve specifically for successful emergent tool use. This metric tracks instances where the model employs extra tools—beyond the required set—and correctly solves the task. The consistent upward trend signifies that the model is actively discovering and utilizing beneficial tools to enhance its reasoning and success rate, validating the capability of our framework to foster emergent behaviors.

We provide qualitative examples of our model's advanced capabilities in Figures 6 and 9. Figure 6 showcases a multi-turn dialogue where the model initially calls an incorrect tool, receives feedback, and robustly recovers by executing the correct tool sequence. This demonstrates sophisticated error-handling and adaptability. Figure 9 highlights two key advantages of our code-as-tool framework: emergent tool use and efficiency. The model chains multiple operations (contrast and grayscale) within a single turn to solve the user's request, even though these specific tools were not



**Figure 8:** *A word cloud of emergent tools discovered during RL training.*

included in the RL training data. This ability to generalize and compose novel tool sequences underscores the flexibility of our approach. As shown in Figure 8, our model learns to use a wide variety of other tools not required in our training data, such as adjusting brightness, applying blurs, and detecting edges, demonstrating the near-infinite toolset enabled by the code-as-tool paradigm. Furthermore, Figure 10 presents a striking example where the model spontaneously combines five tools, three of which never appeared in the RL data to solve a request.

### 4.4. Case Study

To provide a clearer understanding of our model's behavior and the impact of our reward design, we present four case studies.

Figure 11 illustrates a common failure mode known as reward hacking, which we observe when training a model without our proposed constraint penalties. In this example, the agent is tasked with a problem that requires a 90-degree rotation. The agent correctly identifies and applies the rotate90 tool in its first turn, successfully correcting the image's orientation. However, driven solely

**Question:** Approximately how many words are in the paragraph concluding with the phrase "...ha fonglng alise ton"?

**Model Response:**
Thought ···
The question asks for the word count in a specific paragraph that ends with "...ha fonglng alise ton". On the poster, that paragraph appears in the mid-lower section where the 'Did you know?' bullet points are written. I will crop a tight region around that concluding sentence in the 'Did you know?' area to read and count the words.

Tool Call to Image 0
```
# Approximate crop around the line in the 'Did you know?' section
xmin, ymin, xmax, ymax = 559, 725, 838, 781
cropped = image.crop((xmin, ymin, xmax, ymax))
# Then convert to grayscale and enhance contrast for readability
gray = ImageOps.grayscale(cropped)
enh = ImageEnhance.Contrast(gray).enhance(1.5)
result = enh
```

Tool Return: Image 1

Thought ···
I will count the words exclusive of punctuation and rising e: total tokens are: To(1), properly(2), protect(3), the(4), Great(5), Blue(6), Heron(7), we(8), not(9), only(10), here(11), to(12), pesourve(13), the(14), neot(15), alhe(16), but(17), ha(18), fonglng(19), alise(20), ton(21). The line ends with the identifying phrase, so no additional trailing word beyond what's present.
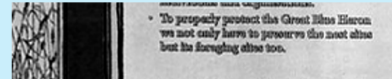
Answer. 21

**Figure 9:** *An example of emergent and efficient tool use. The model chains two tools, contrast enhancement and grayscale conversion, within a single turn to fulfill the user's request. These tools did not appear in the RL training set, demonstrating the generalization capability of the code-as-tool framework.*



**Question:** How many words are in the line that starts with "with women often..."?

**Model Response:**
Thought ···
The image appears rotated: the signage and text read sideways, indicating the photo is not in an upright orientation. To read the paragraph that begins with "with women often..." clearly, I should correct the orientation. I'll rotate the image 90 degrees counterclockwise using the tool to make the text horizontal and readable.

Tool Call to Image 0
```
# Rotate the image 90 degrees counterclockwise to correct orientation
processed_image = image.rotate(90, expand=True)
# Boost contrast for text
processed_image = ImageEnhance.Contrast(processed_image).enhance(1.1)
processed_image = ImageEnhance.Brightness(processed_image).enhance(1.05)
processed_image = ImageEnhance.Sharpness(processed_image).enhance(1.02)
processed_image = ImageEnhance.Contrast(processed_image).enhance(1.03)
# Assign result
result = processed_image
```

Tool Return: Image 1

Thought ···
The image is now upright: all signage and text read horizontally, confirming the orientation is correct. The paragraph starting with "with women often …" appears on the red "Feeding the front" board toward the lower portion of the image. I will crop a tight rectangle around that paragraph so I can count the words clearly.

Tool Call to Image 1
```
# Determine image size
w, h = image.size
# Estimated ROI for the paragraph starting with 'with women often'
left, top, right, bottom = 920, 561, 1183, 803
# Clamp and ensure integers
left = int(round(left)); top = int(round(top))
right = int(round(right)); bottom = int(round(bottom))
left = max(0, min(left, w - 1))
top = max(0, min(top, h - 1))
right = max(left + 1, min(right, w))
bottom = max(top + 1, min(bottom, h))
# Crop the region
processed_image = image.crop((left, top, right, bottom))
# Improve contrast
processed_image = ImageEnhance.Contrast(processed_image).enhance(1.05)
processed_image = ImageEnhance.Brightness(processed_image).enhance(1.08)
processed_image = ImageEnhance.Sharpness(processed_image).enhance(1.02)
processed_image = ImageEnhance.Contrast(processed_image).enhance(1.04)
# Assign result
result = processed_image
```

Tool Return: Image 2

Thought ···
From the cropped region, the paragraph beginning with "with women often …" is clearly legible: "with women often replacing men in factories." I will count the words in that paragraph. Splitting by spaces and removing punctuation: "with" (1), "women" (2), "often" (3), "replacing" (4), "men" (5), "in" (6), "factories." (7). That yields seven words.

Answer. 7
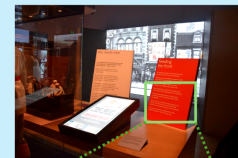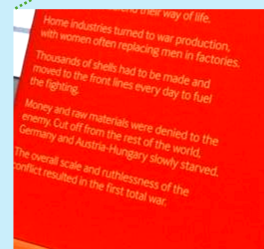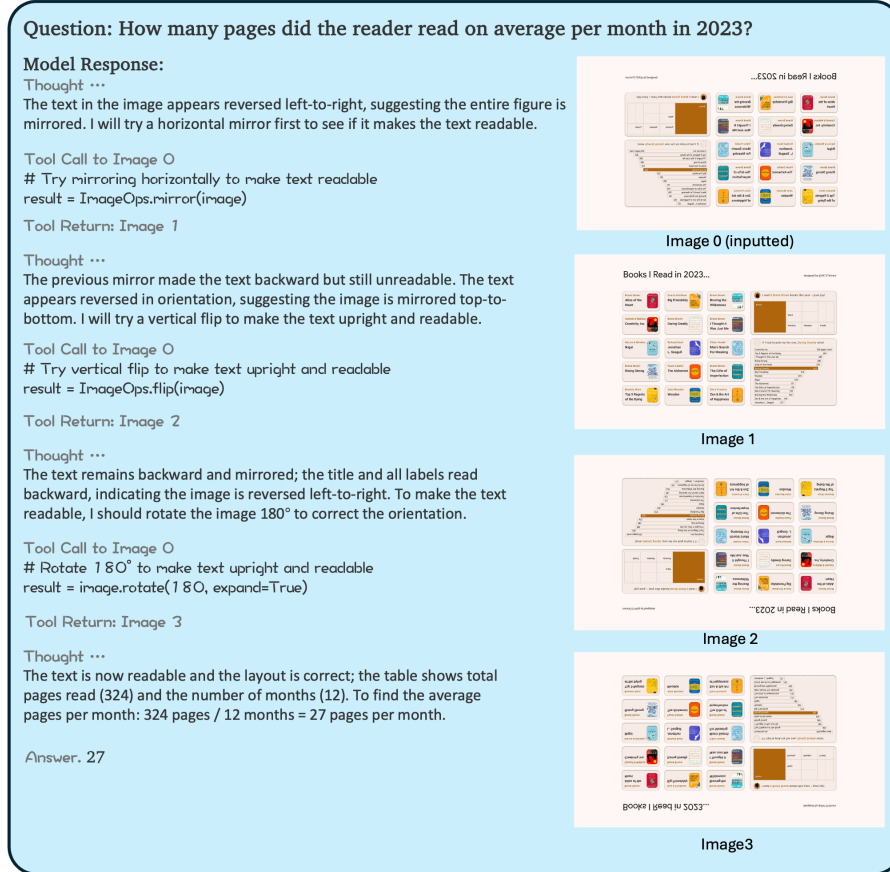
**Figure 10:** *An example of emergent and efficient tool use. The model chains "brightness up, contrast up, crop, rotate90, sharpness" tools to solve the user's request.*

by the goal of maximizing the strategy reward, it does not stop. Instead, it proceeds to call additional, unnecessary orientation tools in subsequent turns, which corrupts the already-correct image and

**Question: How many pages did the reader read on average per month in 2023?**

**Model Response:**
Thought ⋯
The text in the image appears reversed left-to-right, suggesting the entire figure is mirrored. I will try a horizontal mirror first to see if it makes the text readable.

Tool Call to Image 0
\# Try mirroring horizontally to make text readable
result = ImageOps.mirror(image)

Tool Return: Image 1

Thought ⋯
The previous mirror made the text backward but still unreadable. The text appears reversed in orientation, suggesting the image is mirrored top-to-bottom. I will try a vertical flip to make the text upright and readable.

Tool Call to Image 0
\# Try vertical flip to make text upright and readable
result = ImageOps.flip(image)

Tool Return: Image 2

Thought ⋯
The text remains backward and mirrored; the title and all labels read backward, indicating the image is reversed left-to-right. To make the text readable, I should rotate the image 180° to correct the orientation.

Tool Call to Image 0
\# Rotate 180° to make text upright and readable
result = image.rotate(180, expand=True)

Tool Return: Image 3

Thought ⋯
The text is now readable and the layout is correct; the table shows total pages read (324) and the number of months (12). To find the average pages per month: 324 pages / 12 months = 27 pages per month.

Answer. 27

**Figure 11:** *An example of reward hacking from a model trained without constraint penalties. After correctly applying rotate90, the agent continues to call superfluous tools, leading to task failure.*
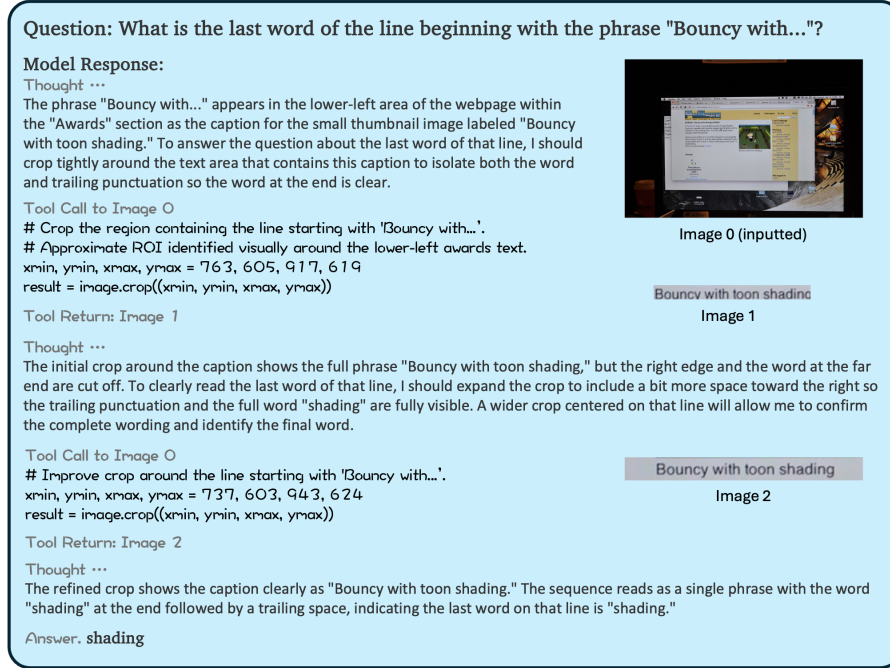
ultimately leads to failure. This case underscores the critical importance of the penalty functions, which act as guardrails to prevent such inefficient and counterproductive behavior.

In contrast, Figure 12 showcases the robust, multi-step reasoning capabilities of our fully-trained CodeVision model. The task requires the model to read text from a specific region that is too small to be seen in the initial view. The model first performs a "crop" operation to zoom in. Recognizing that its initial crop was incomplete and missed part of the target area, it dynamically adjusts its strategy. In the next turn, it executes a second, more precise "crop" that successfully isolates the entire region of interest. With the full context now visible, the model is able to accurately provide the correct answer. This example highlights the model's ability to self-correct and perform fine-grained, iterative adjustments, a key capability for solving complex real-world visual tasks.

Furthermore, Figure 13 presents a nuanced success case that reveals areas for further refinement. The task requires a sequence of orientation correction followed by localizing and reading a very small line of text. The model adeptly handles the first step by correcting the image's orientation. It then correctly reasons about the text's location and performs a crop. Although the cropped region successfully includes the target text, allowing the model to answer correctly, the crop itself is suboptimal—a long, narrow strip containing significant irrelevant context. This example demonstrates that while the model's high-level reasoning is effective, its ability to generate precise, tight bounding boxes could be improved, suggesting it may sometimes adopt an overly cautious safe cropping strategy.

To further probe the model's limitations, we present a clear failure case in Figure 14. The task requires correcting the image orientation and then cropping a small, specific region to examine

**Figure 12:** *A successful multi-step reasoning case. The model first performs an initial crop, then recognizes its incompleteness and refines the cropped region in a second step to include all necessary information before providing the final answer.*

details. The model successfully completes the first two stages: it correctly adjusts the orientation and accurately reasons about the approximate location of the region of interest. However, it fails at the final, crucial step of precise localization. The coordinates it generates for the crop are slightly off, causing the cropped area to be adjacent to the target region but missing it entirely. This case clearly demonstrates that while the model excels at high-level reasoning and coarse localization, its fine-grained coordinate prediction remains a challenge and a key area for future improvement.

## 4.5. Ablation Study

To validate the effectiveness of our proposed training methodology and reward components, we conduct a series of ablation studies using the CodeVision-7B model as the baseline.

**Effectiveness of dense reward.** We investigate the impact of our dense reward function by training two variants of our model: one without the strategy-shaping reward ($R_{\text{strategy}}$) and another without the constraint penalties ($P_{\text{cost}}$). The results are presented in Table 3 and Figure 15.

As shown in Table 3, removing the strategy reward leads to a substantial performance drop across all benchmarks, with the most significant degradation observed on MVToolBench (from 60.1 to 50.7). This confirms that a simple outcome-based reward is insufficient for learning complex tool-use strategies. The dense process signals from the strategy reward are critical for guiding the model to discover and reinforce effective reasoning paths.

Similarly, removing the constraint penalties also degrades performance, particularly on V* and MVToolBench. This highlights the importance of penalizing inefficient or illogical actions. Without these penalties, the model is susceptible to "reward hacking," where it learns to maximize rewards through superfluous actions (*e.g.,* rotating an already-correct image) that do not contribute to solving the task. The penalties serve as essential guardrails that promote more deliberate and efficient problem-solving. Figure 15 visualizes this by showing that our full reward function leads to a more

**Question:** How many words are on the line that starts with the text "Ripped..."?

**Model Response:**

Thought ⋯
The page text and UI elements are mirrored horizontally, indicating a left-right reversal. I will flip the image horizontally to restore the correct orientation.

Tool Call to Image 0
result = ImageOps.mirror(image)

Tool Return: Image 1

Thought ⋯
The current page is now correctly oriented. I need to isolate the line that begins with the text "Ripped..." on the right sidebar, which contains the description reading "Ripped Freak endelig i Norge. Et av de mest populaere kosttillskuddene i USA/Canada for maksimal fettforbrømming. 349 kr for 60 kapsler." I will crop a tight region around that specific line.

Tool Call to Image 1
# Crop the line starting with 'Ripped' on the right sidebar
xmin, ymin, xmax, ymax = 35, 266, 1004, 283
result = image.crop((xmin, ymin, xmax, ymax))
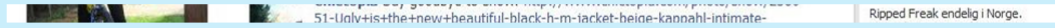
Tool Return: Image 2

Thought ⋯
The cropped line reads: "Ripped Freak endelig i Norge." Counting words: Ripped(1), Freak(2), endelig(3), i(4), Norge.(5).

Answer. 5

**Figure 13:** *A successful but inefficient localization case. The model corrects the orientation and successfully crops a region containing the target text. However, the crop is a long, narrow strip that includes much irrelevant information, pointing to potential improvements in generating more precise bounding boxes.*



**Question:** What is the full text of the paragraph opening with the words "RECHERCHE JESSICA..."?

**Model Response:**

Thought ⋯
The text and headings in the image appear upside down—the large title and all labels are inverted top-to-bottom while not mirrored left-to-right. This indicates the entire page is vertically flipped. To read the paragraph that begins with "RECHERCHE JESSICA...", I should correct the orientation by flipping the image vertically. I will apply a vertical flip using the tool and then read the paragraph.

Tool Call to Image 0
# Flip the image vertically to correct upside-down orientation
result = ImageOps.flip(image)

Tool Return: Image 1

Thought ⋯
The corrected image now shows all text upright and readable. The paragraph beginning with "RECHERCHE JESSICA..." appears as a small line of attribution text in the bottom-right area of the page. I will crop that specific line so the full text can be read clearly.

Tool Call to Image 1
# Crop the attribution line beginning with 'RECHERCHE JESSICA...' in the bottom-right
w, h = image.size
# Estimated ROI for the attribution text
left, top, right, bottom = 670, 1530, 1270, 1550
result = image.crop((left, top, right, bottom))

Tool Return: Image 2

Thought ⋯
The cropped bottom-right attribution line is sharp and fully legible. It reads the complete attribution text beginning with "RECHERCHE JESSICA...". I can now provide the exact text.

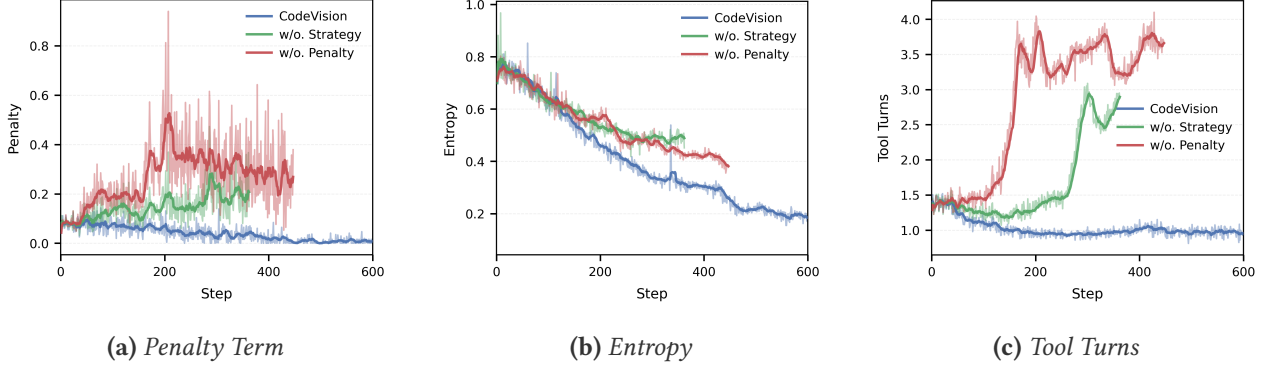Answer. RECHERCHE JESSICA NADEAU — INFOGRAPHIE JUSTIN STAHLMAN, JDÉM

**Figure 14:** *A failure case highlighting the model's limitations in precise localization. The model correctly adjusts the orientation and identifies the general area of interest, but provides slightly inaccurate coordinates for the crop, resulting in a cropped region that just misses the target.*
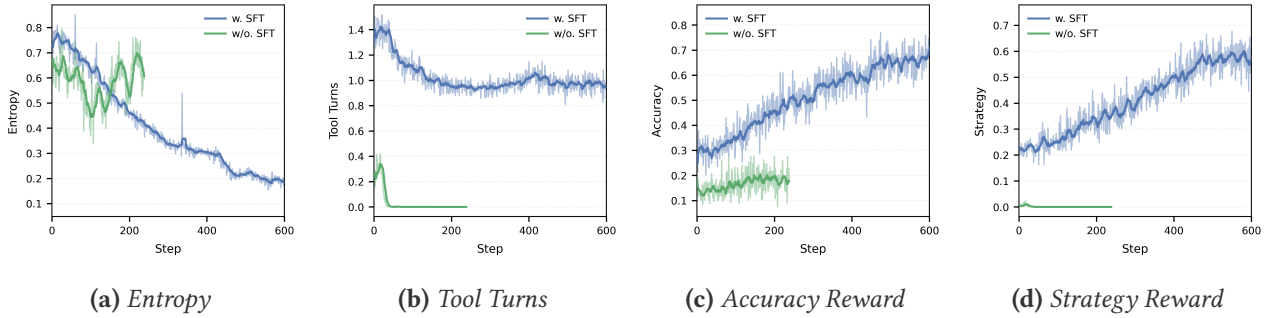
**Table 3:** *Ablation study of our key reward components on the CodeVision-7B model. Removing either the strategy reward or the constraint penalties leads to a notable drop in performance.*

| Model | OCRBench | | ChartQAPro | | V* | MVToolBench |
|---|---|---|---|---|---|---|
| | Rot180 | Verti | Rot180 | Hori | | |
| Qwen2.5-VL-7B | 70.2 | 17.0 | 23.4 | 19.5 | 74.6 | 18.1 |
| Qwen2.5-VL-7B-SFT | 57.0 | 35.8 | 23.2 | 20.9 | 71.7 | 26.6 |
| CodeVision-7B | 72.3 | 67.4 | 30.8 | 30.1 | 83.7 | 60.1 |
| w/o Strategy Reward | 60.9 | 61.5 | 24.6 | 28.9 | 78.5 | 50.7 |
| w/o Penalty | 68.3 | 66.3 | 24.0 | 24.3 | 71.2 | 55.9 |



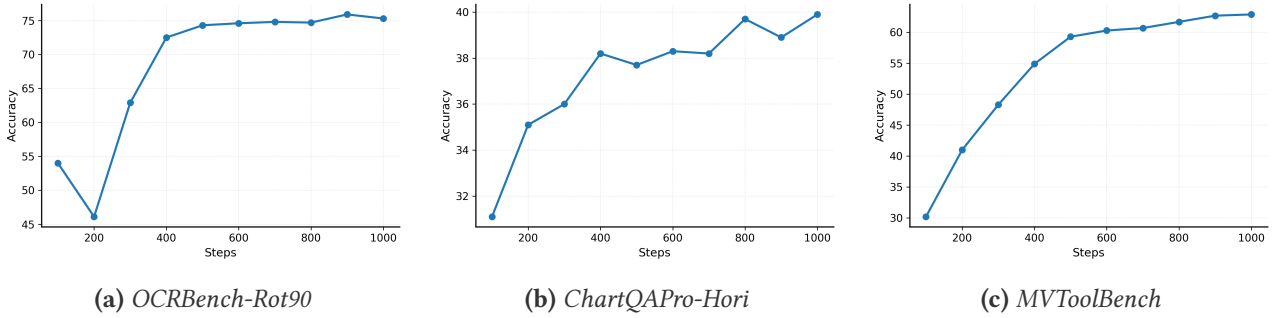**(a)** *Penalty Term*      **(b)** *Entropy*      **(c)** *Tool Turns*

**Figure 15:** *Training dynamics for our full reward function versus ablated versions. The complete reward function leads to more stable and efficient learning.*

stable and higher final success rate during RL training compared to the ablated versions.



**(a)** *Entropy*      **(b)** *Tool Turns*      **(c)** *Accuracy Reward*      **(d)** *Strategy Reward*

**Figure 16:** *Comparison of training with and without the SFT cold-start stage. Starting RL from a SFT-tuned checkpoint is critical for effective learning, while training from scratch fails to converge.*

**Necessity of the cold start.** We also study the importance of our two-stage training process. We compare our full model, which is first fine-tuned with SFT and then trained with RL, against a baseline trained with RL directly on the base model without the SFT cold start. As illustrated in Figure 16, the model trained without SFT fails to achieve meaningful improvement. This is because the vast and unstructured action space of code generation makes it extremely difficult for the agent to discover useful tool-use policies through pure RL exploration. The SFT stage is crucial for bootstrapping the model, teaching it the fundamental syntax of tool calls and providing a strong initial policy. This warm start makes the subsequent RL phase significantly more stable and effective, allowing the model to refine its strategic reasoning on a solid foundation.

**(a)** *OCRBench-Rot90*  **(b)** *ChartQAPro-Hori*  **(c)** *MVToolBench*

**Figure 17:** *Accuracy on OCRBench, ChartQAPro, and MVToolBench during RL training. The steady improvement across all benchmarks indicates that the model's learning has not yet saturated and could benefit from further scaling of data and tasks.*

## 5. Conclusion

In this paper, we identify a critical vulnerability in state-of-the-art MLLMs: their pronounced brittleness to natural corruptions such as orientation changes. We introduce CodeVision, the "code-as-tool" framework, a flexible paradigm that empowers models to generate code for a virtually unlimited range of visual operations, moving beyond fixed toolsets. We train our models using a two-stage process combining SFT and RL with a dense, process-oriented reward function to foster strategic tool use. Experiments on our newly created benchmarks demonstrate that this approach remedies the identified weakness and significantly enhances model reasoning. Our framework fosters emergent behaviors like efficient tool-chaining and error recovery, paving the way for more capable MLLMs that treat visual interaction as a programming task.

## 6. Limitations and Future Work

In this work, we have focused on a core set of tools—primarily orientation correction and cropping—to establish the CodeVision framework and demonstrate its effectiveness in addressing model brittleness. While this targeted approach has proven successful, creating a truly general-purpose visual agent requires scaling and expansion in several key dimensions. Our work reveals several limitations that open up exciting avenues for future research:

**Expanding Tool Diversity and Compositionality.** Our current research deliberately concentrates on a limited toolset. To enhance the model's generalization capabilities, future work should incorporate a much wider variety of tool types and data. This includes training on tasks that require more complex tool compositions and even multi-image tool use (*e.g.,* comparing, merging, or analyzing multiple images simultaneously). Furthermore, the "code-as-tool" paradigm can be extended beyond standard Python libraries to include custom tools, such as proprietary search engines or generative models. By exposing only a simple API endpoint to the model, we can empower it to leverage powerful, black-box functionalities, further broadening its problem-solving horizons.

**Refining Process Supervision with Beneficial Tools.** Our current training methodology relies on a "must-use" tool list to provide strong process supervision. While effective, this can be expanded to a more flexible framework. Future iterations could incorporate a broader set of "beneficial" tools. These are tools that are not strictly necessary to solve a task but are highly likely to improve performance, such as contrast enhancement for low-light images or a preliminary crop for focusing attention. By rewarding the use of such advantageous tools, we can encourage the model to learn more nuanced strategies and generalize better to tasks where the optimal toolset is not rigidly defined.

**Scaling Data, Tasks, and Model Exploration.** Our experiments indicate that the model's learning has not yet reached saturation. The entropy of its policy suggests that there is still significant room for exploration and improvement. Moreover, as shown in Figure 17, we observe that performance on our benchmarks continues to increase steadily throughout training with no signs of plateauing. This strongly suggests that our approach would benefit from scaling up. By sourcing more diverse data, incorporating a wider array of tool types, and designing more varied tasks, we believe there is a clear path toward developing even more capable and robust visual agents.

# References

Bai, S., Cai, Y., Chen, R., Chen, K., Chen, X., Cheng, Z., Deng, L., Ding, W., Gao, C., Ge, C., et al. (2025a). Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*.

Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J., et al. (2025b). Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Chen, W., Ma, X., Wang, X., and Cohen, W. W. (2023). Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.

Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang, D., Rosen, E., et al. (2025). Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Fan, Y., He, X., Yang, D., Zheng, K., Kuo, C.-C., Zheng, Y., Narayanaraju, S. J., Guan, X., and Wang, X. E. (2025). Grit: Teaching mllms to think with images. *arXiv preprint arXiv:2505.15879*.

Gao, C., Zheng, C., Chen, X.-H., Dang, K., Liu, S., Yu, B., Yang, A., Bai, S., Zhou, J., and Lin, J. (2025). Soft adaptive policy optimization. *arXiv preprint arXiv:2511.20347*.

Guo, Z., Hong, M., and Jin, T. (2025a). Observe-r1: Unlocking reasoning abilities of mllms with dynamic progressive reinforcement learning. *arXiv preprint arXiv:2505.12432*.

Guo, Z., Zhang, F., Jia, K., and Jin, T. (2025b). Llm-i: Llms are naturally interleaved multimodal creators. *arXiv preprint arXiv:2509.13642*.

Hong, M., Guo, Z., Xia, Y., Wang, Z., Zhang, Z., Jin, T., and Zhao, Z. (2025). Apo: Enhancing reasoning ability of mllms via asymmetric policy optimization. *arXiv preprint arXiv:2506.21655*.

Jin, B., Zeng, H., Yue, Z., Yoon, J., Arik, S., Wang, D., Zamani, H., and Han, J. (2025). Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*.

Lai, X., Li, J., Li, W., Liu, T., Li, T., and Zhao, H. (2025). Mini-o3: Scaling up reasoning patterns and interaction turns for visual search. *arXiv preprint arXiv:2509.07969*.

Li, X., Dong, G., Jin, J., Zhang, Y., Zhou, Y., Zhu, Y., Zhang, P., and Dou, Z. (2025). Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*.

Liu, C.-L., Yin, F., Wang, D.-H., and Wang, Q.-F. (2011). Casia online and offline chinese handwriting databases. In *2011 international conference on document analysis and recognition*, pages 37–41. IEEE.

Liu, Y., Li, Z., Huang, M., Yang, B., Yu, W., Li, C., Yin, X.-C., Liu, C.-L., Jin, L., and Bai, X. (2024). Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102.

Liu, Z., Zang, Y., Zou, Y., Liang, Z., Dong, X., Cao, Y., Duan, H., Lin, D., and Wang, J. (2025). Visual agentic reinforcement fine-tuning. *arXiv preprint arXiv:2505.14246*.

Long, S., Qin, S., Panteleev, D., Bissacco, A., Fujii, Y., and Raptis, M. (2022). Towards end-to-end unified scene text detection and layout analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

Marti, U.-V. and Bunke, H. (2002). The iam-database: an english sentence database for offline handwriting recognition. *International journal on document analysis and recognition*, 5(1):39–46.

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. (2025). Chartqapro: A more diverse and challenging benchmark for chart question answering. *arXiv preprint arXiv:2504.05506*.

Masry, A., Kavehzadeh, P., Do, X. L., Hoque, E., and Joty, S. (2023). Unichart: A universal vision-language pretrained model for chart comprehension and reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14662–14684.

Mathew, M., Karatzas, D., and Jawahar, C. (2021). Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209.

OpenAI (2025a). Introducing deep research.

OpenAI (2025b). Introducing gpt-5.

OpenAI (2025c). Introducing openai o3 and o4-mini.

Qiao, R., Tan, Q., Yang, P., Wang, Y., Wang, X., Wan, E., Zhou, S., Dong, G., Zeng, Y., Xu, Y., et al. (2025). We-math 2.0: A versatile mathbook system for incentivizing visual mathematical reasoning. *arXiv preprint arXiv:2508.10433*.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y., et al. (2024). Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.

Singh, A., Natarajan, V., Shah, M., Jiang, Y., Chen, X., Batra, D., Parikh, D., and Rohrbach, M. (2019). Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326.

Su, A., Wang, H., Ren, W., Lin, F., and Chen, W. (2025a). Pixel reasoner: Incentivizing pixel-space reasoning with curiosity-driven reinforcement learning. *arXiv preprint arXiv:2505.15966*.

Su, Z., Li, L., Song, M., Hao, Y., Yang, Z., Zhang, J., Chen, G., Gu, J., Li, J., Qu, X., et al. (2025b). Openthinkimg: Learning to think with images via visual tool reinforcement learning. *arXiv preprint arXiv:2505.08617*.

Team, T. D., Li, B., Zhang, B., Zhang, D., Huang, F., Li, G., Chen, G., Yin, H., Wu, J., Zhou, J., et al. (2025). Tongyi deepresearch technical report. *arXiv preprint arXiv:2510.24701*.

Wang, S., Yang, S., Lin, W., Guo, Z., Cai, S., Huang, H., Wang, Y., Chen, J., and Jin, T. (2025a). Omni-chart-600k: A comprehensive dataset of chart types for chart understanding. *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 4051–4069.

Wang, W., Ding, L., Zeng, M., Zhou, X., Shen, L., Luo, Y., Yu, W., and Tao, D. (2025b). Divide, conquer and combine: A training-free framework for high-resolution image perception in multimodal large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 7907–7915.

Wang, X., Yang, Z., Feng, C., Lu, H., Li, L., Lin, C.-C., Lin, K., Huang, F., and Wang, L. (2025c). Sota with less: Mcts-guided sample selection for data-efficient visual reasoning self-improvement. *arXiv preprint arXiv:2504.07934*.

Wu, P. and Xie, S. (2024). V?: Guided visual search as a core mechanism in multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13084–13094.

Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., Yue, Y., Dai, W., Fan, T., Liu, G., Liu, L., et al. (2025). Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.

Zhang, Y.-F., Lu, X., Yin, S., Fu, C., Chen, W., Hu, X., Wen, B., Jiang, K., Liu, C., Zhang, T., et al. (2025). Thyme: Think beyond images. *arXiv preprint arXiv:2508.11630*.

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

Zheng, Z., Yang, M., Hong, J., Zhao, C., Xu, G., Yang, L., Shen, C., and Yu, X. (2025b). Deepeyes: Incentivizing" thinking with images" via reinforcement learning. *arXiv preprint arXiv:2505.14362*.