

MindWatcher: Toward Smarter Multimodal Tool-Integrated Reasoning

MindGPT-ov Team
Li Auto Inc

Traditional workflow-based agents exhibit limited intelligence when addressing real-world problems requiring tool invocation. Tool-integrated reasoning (TIR) agents capable of autonomous reasoning and tool invocation are rapidly emerging as a powerful approach for complex decision-making tasks involving multi-step interactions with external environments. In this work, we introduce MindWatcher, a TIR agent integrating interleaved thinking and multimodal chain-of-thought (CoT) reasoning. MindWatcher can autonomously decide whether and how to invoke diverse tools and coordinate their use, without relying on human prompts or workflows. The interleaved thinking paradigm enables the model to switch between thinking and tool calling at any intermediate stage, while its multimodal CoT capability allows manipulation of images during reasoning to yield more precise search results. We implement automated data auditing and evaluation pipelines, complemented by manually curated high-quality datasets for training, and we construct a benchmark, called MindWatcher-Evaluate Bench (MWE-Bench), to evaluate its performance. MindWatcher is equipped with a comprehensive suite of auxiliary reasoning tools, enabling it to address broad-domain multimodal problems. A large-scale, high-quality local image retrieval database, covering eight categories including cars, animals, and plants, endows model with robust object recognition despite its small size. Finally, we design a more efficient training infrastructure for MindWatcher, enhancing training speed and hardware utilization. Experiments not only demonstrate that MindWatcher matches or exceeds the performance of larger or more recent models through superior tool invocation, but also uncover critical insights for agent training, such as the genetic inheritance phenomenon in agentic RL. The agent reasoning framework, MWE-Bench, three smaller-scale agent models (2B, 3B, and 4B) distilled from MindWatcher 32B, and related resources will be open-sourced.

- 📅 **Last Update Date:** December 29, 2025
- ✉️ **Correspondence:** chenwei10@lixiang.com
- 🔗 **Code:** <https://github.com/TIMMY-CHAN/MindWatcher>
- 🤗 **Hugging Face:** <https://huggingface.co/datasets/Lost-Cloud/MWE-Bench>

1 Introduction

Large language models (LLMs) [27, 1, 14, 42, 9, 37, 25] have achieved remarkable progress in recent years, demonstrating strong capabilities in language understanding, knowledge acquisition, and complex reasoning tasks. However, despite the powerful world knowledge and

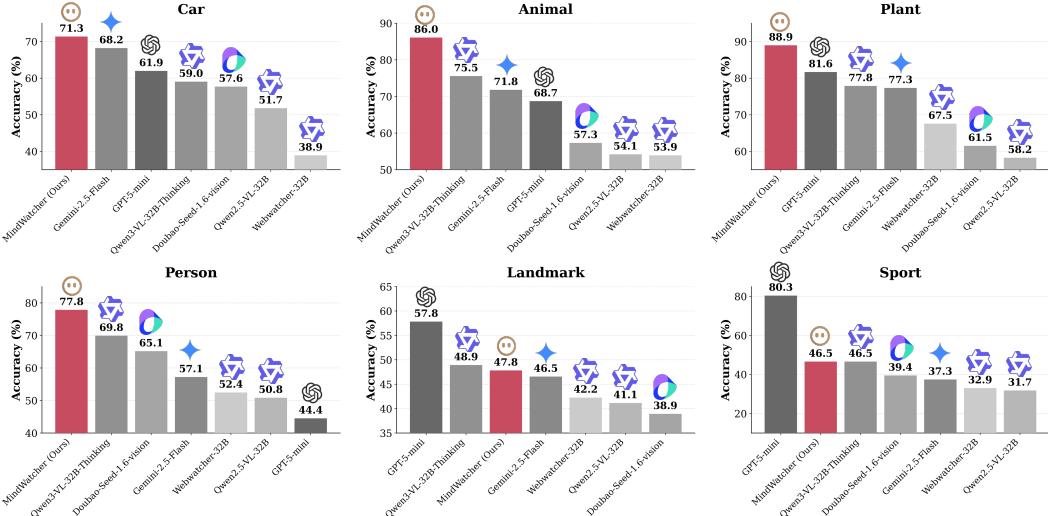


Figure 1: MWE-Bench Performance of MindWatcher.

multimodal capabilities of the latest models such as Gemini 2.5 Pro [9], most LLMs remain fundamentally constrained by the limits of their parametric knowledge: they struggle to cover long-tail information and fine-grained domain-specific knowledge [7], and they cannot directly access real-time information that emerges after training. These structural bottlenecks hinder their reliability in many real-world applications, especially those requiring external knowledge, multi-step information integration, or cross-modal reasoning. Equipping LLMs with external tools has therefore become a key strategy to overcome these limitations. By connecting models with retrieval engines [19, 29], computation tools, or code interpreters, the boundary of problem-solving capabilities can be substantially extended.

Traditional tool-augmented approaches typically rely on manually designed workflows [17, 33] to orchestrate tool invocation, yet such methods exhibit limited adaptability when confronted with the diversity and uncertainty inherent in open-domain environments, which become even more fragile when handling cross-modal demands. Multi-agent systems [41, 31, 22, 21] partially alleviate these issues: a powerful planner agent is responsible for decision-making, while tool-specialized agents execute designated subtasks. This architecture has become highly popular in the industry and significantly improves system flexibility and scalability. However, it also introduces new complexity and overhead, including redundant model deployment and latency caused by chained interactions, which limits its further expansion. With the emergence of thought-augmented models [38, 14], the research community increasingly recognizes that intelligent systems need not rely on multi-component designs: a single unified language model can assume both planning and acting roles. This has led to the rise of Tool-Integrated Reasoning (TIR) methods [26], exemplified by the ReAct [43] paradigm. The core idea is to let the model explicitly generate intermediate thoughts, autonomously invoke tools, and iteratively make decisions based on environmental feedback. TIR agents can dynamically plan multi-step operations in open-world tasks and achieve end-to-end problem solving, making them a promising path toward more general-purpose agents.

However, current TIR systems still fall short of being truly practical and general intelligent agents, with significant limitations across several key dimensions. From an application perspective, existing TIR agents [26, 22, 11] are predominantly focused on text-based tasks, particularly DeepSearch-style reasoning centered on retrieval. Only a small number of works [13, 44] attempt to introduce visual capabilities, and most rely solely on image search tools without enabling the agent to directly manipulate images or perform fine-grained cross-modal reasoning to support problem solving. This severely limits their performance on multimodal tasks and prevents them from tackling the many vision-driven decision-making scenarios found in real-world environments.

From a training methodology perspective, TIR agents face a triple challenge across data, algorithms, and training frameworks. High-quality reasoning trajectories involving multiple tools and multi-step interactions are extremely difficult to construct manually. SFT-based training [24, 6] often causes models to “imitate” the thought-action format rather than truly “learn” the underlying strategy—manifested in excessive, redundant tool calls on simple problems and substantial performance degradation on general benchmarks. Moreover, existing training frameworks lack fine-grained supervision over the interleaved process of thinking, tool invocation, and subsequent reasoning, preventing models from forming stable and reliable tool-use behaviors and exacerbating issues such as tool misuse and unnecessary calls. From the perspective of tool ecosystems, many core retrieval capabilities, especially visual retrieval, rely on expensive external APIs. Their high cost under frequent invocation further constrains the practical deployment of TIR agents in local or enterprise settings.

To address the challenges outlined above, we introduce MindWatcher, a TIR agent capable of autonomous planning and execution, multimodal perception, and robust tool coordination. Leveraging an interleaved thinking paradigm and a multimodal Chain-of-Thought (CoT) mechanism, MindWatcher can flexibly alternate between internal thinking and external tool invocation at any stage of the reasoning process. By integrating fine-grained visual operations into the reasoning chain, the agent achieves precise region-level visual perception and more accurate cross-modal information retrieval.

To avoid the drawbacks of conventional SFT, such as rigid imitation of reasoning formats and redundant tool calls on simple tasks, MindWatcher abandons standard SFT and instead adopts a continuous reinforcement learning (RL) strategy conducted in both real and offline environments. We develop two automated image–text pair construction pipelines to reduce data generation costs. In parallel, we equip MindWatcher with a comprehensive set of tools that cover core multimodal reasoning needs, including image region cropping and zooming, object grounding and visual search, external text retrieval, webpage content extraction, and local Python code interpreter. Moreover, we construct a large-scale local visual corpus spanning categories such as person, animals, plants, cars, landmarks, and logos. We also build a new multimodal benchmark: MindWatcher-Evaluate Benchmark (MWE-Bench) for systematically evaluating agentic multimodal tool-use and reasoning capabilities.

At the system level, we design an RL training pipeline that supports asynchronous tool invocation, significantly improving learning efficiency. We also introduce a new GRPO-based agentic RL algorithm, which introduce step-wise normalization which ensure the optimization objective on individual action segments rather than the global token stream.

As shown in Figure 1, MindWatcher demonstrates strong generality and efficiency across a wide range of tasks on MWE-Bench. The 32B model achieves state-of-the-art (SOTA) performance in tool-augmented reasoning while maintaining robust general capabilities, and we distilled and open-sourced 2B, 3B, and 4B variants based on the MindWatcher, which also exhibit highly competitive results.

2 Method

2.1 Working Paradigm

To support flexible multimodal reasoning and autonomous tool use, MindWatcher models the TIR process as a Markov Decision Process (MDP). Given an initial user prompt s_0 , the agent interacts with the environment by generating an interleaved sequence of actions and tool-grounded observations:

$$Y = \{a_0, obs_0, a_1, obs_1, \dots, obs_{n-1}, a_n\}. \quad (1)$$

Each action a_j is executed against the environment—typically through a thinking process and a tool call—yielding an observation obs_j , which is appended to the context and becomes part of the next state. The agent iteratively continues this process until generating the final action a_n , which contains the concluding response.

Interleaved Thinking and Multimodal CoT MindWatcher implements this MDP through an autoregressive generation loop. At each step t , the Policy $\pi_\theta(a_t|s_t)$ (param-

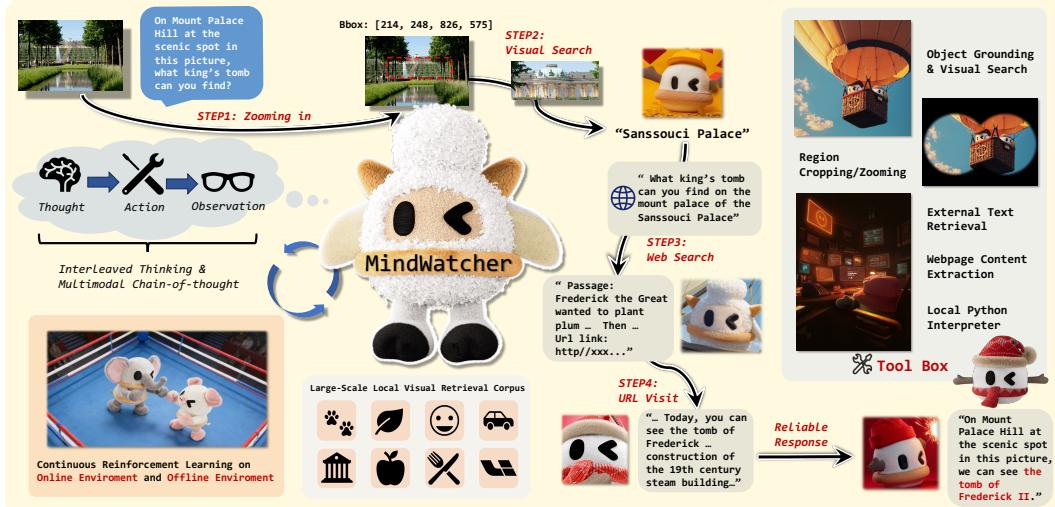


Figure 2: **The Working Paradigm of MindWatcher.** To address complex multimodal question answering tasks, we train our model using continuous RL to develop Multimodal CoT capabilities. By integrating interleaved thinking, the model is able to interact with the environment and autonomously invoke tools in the toolbox. Furthermore, to facilitate more accurate and lowcost visual search, MindWatcher have constructed a large-scale local retrieval corpus spanning eight major categories.

eterized by the LLM) conditions on the full history s_t . Distinct from traditional approaches where actions are strictly physical tool calls, we define a unified action space $\mathcal{A} = \mathcal{A}_{thought} \cup \mathcal{A}_{tool}$. In implementation, thoughts and tool calls are serialized through dedicated $<\text{think}> \dots <\backslash\text{think}>$ and $<\text{tool_call}> \dots <\backslash\text{tool_call}>$ tags, enabling the model to interleave reasoning and action generation within a single decoding sequence. MindWatcher further incorporates a multimodal CoT [38, 16] mechanism, which allows the agent to “think with images” by embedding image-dependent operations into its reasoning chain.

As shown in Figure 2, given an image and a complex query, the model enters an iterative planning and tool-call process. After each tool-call completes, the tool response of current stage is obtained. Subsequently, the next action is determined based on the tool response, ultimately yielding the query result.

2.2 Training Algorithm

While SFT remains the prevailing paradigm for training TIR agents, our empirical observations reveal significant limitations. We found that fine-tuning already robust instruction-following or thinking models on trajectory data often incurs a heavy “alignment tax”, severely degrading performance on general-purpose tasks. Furthermore, SFT tends to induce *tool abuse*, manifested as redundant invocations for trivial queries and excessive, ineffective looping in complex scenarios. Consequently, we adopt a pure RL approach to endow MindWatcher with genuine decision-making and self-correction capabilities.

2.2.1 Step-wise Normalized GRPO

We employ an enhanced version of Group Relative Policy Optimization (GRPO) [30] as our core learning algorithm. Standard GRPO typically normalizes advantages over a single dialogue turn or global sequence. However, in an agentic environment, observation tokens generated by the environment must be excluded from loss calculation. Let $\mathcal{O}_q = \{o_1, o_2, \dots, o_G\}$ be a group of trajectories generated from a user prompt q . For each trajectory o_i , we compute a sequence-level reward r_i . The advantage function is computed based on the distribution of rewards within the group:

$$\hat{A}_i = \frac{r_i - \mu_r}{\sigma_r}, \quad (2)$$

where μ_r and σ_r are the mean and standard deviation of the rewards for the G samples, respectively.

In a standard multi-turn agent setting, the expected objective function is typically formulated as summing over all action tokens:

$$J(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{\sum_{j=0}^n |a_j|} \sum_{j=0}^n \sum_{t=T_j}^{T_j+|a_j|} \min \left[\frac{\pi_\theta(t|s_t)}{\pi_{\theta_{old}}(t|s_t)} \cdot A_{i,t}, \text{clip} \left(\frac{\pi_\theta(t|s_t)}{\pi_{\theta_{old}}(t|s_t)}, 1 - \epsilon, 1 + \epsilon \right) \cdot A_{i,t} \right]. \quad (3)$$

However, in the context of Interleaved Thinking, a single trajectory comprises multiple “Think and Tool-call” cycles (episodes) with drastically varying action lengths. Simply summing gradients allows episodes to dominate optimization. To ensure balanced supervision across every reasoning step, we propose **Step-wise Normalization**. We define the optimization objective on individual *Action Segments* a_j rather than the global token stream. Assuming the i -th trajectory contains n_i action steps, and the j -th action segment a_j has a length of $|a_j|$, our optimized objective function $J(\theta)$ is formalized as:

$$J(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{n_i} \sum_{j=1}^{n_i} \frac{1}{|a_j|} \sum_{t \in a_j} \min \left[\frac{\pi_\theta(t|s_t)}{\pi_{\theta_{old}}(t|s_t)} \cdot \hat{A}_i, \text{clip} \left(\frac{\pi_\theta(t|s_t)}{\pi_{\theta_{old}}(t|s_t)}, 1 - \epsilon, 1 + \epsilon \right) \cdot \hat{A}_i \right]. \quad (4)$$

This formulation introduces a dual-normalization mechanism:

1. **Action-Step Normalization** ($\frac{1}{n_i}$): Weighs each trajectory equally regardless of the number of “Think and Tool-call” cycles.
2. **Token-Length Normalization** ($\frac{1}{|a_j|}$): Averages loss within each “Think and Tool-call” episode.

2.2.2 Hybrid Reward

To steer the model toward both syntactic correctness and factual accuracy, we design a hybrid reward function consisting of three components: Outcome Accuracy Reward, Format Reward, and Hallucination Tool-call Penalty.

1. Outcome Accuracy Reward (R_{acc}): This is a sparse reward computed only at termination. Given the complexity of open-ended multimodal QA, regular expressions are insufficient for verification. We employ a Model-based Judge to evaluate the factual consistency between the model output and the ground truth.

$$R_{acc} = \begin{cases} 1.0 & \text{if Judge returns "1" (Correct),} \\ 0.0 & \text{if Judge returns "0" (Incorrect).} \end{cases} \quad (5)$$

2. Format Reward (R_{fmt}): We implement a strict regex-based parser to enforce schema adherence. This includes:

- **Structural Integrity:** Verifying that tags such as `<think>`, `<tool_call>`, and `<answer>` appear in valid pairs and sequences.
- **Residue Penalty:** We strictly forbid “chitchat” outside of valid tags (e.g., outputting “I will now execute...” after a `<tool_call>` block). Any non-whitespace character outside tags incurs a penalty, as we observed that such residues often lead to output collapse during training.

$$R_{fmt} = \begin{cases} 0.5 & \text{if strictly follows schema,} \\ -0.5 - 0.01 \times \text{len(residue)} & \text{if format error or residue detected.} \end{cases} \quad (6)$$

3. Hallucination Tool-call Penalty (R_{halluc}): During experiments, *we observed a tendency for models to generate consecutive `<tool_call>` blocks without waiting*

for the environment feedback <tool_response>, effectively hallucinating execution results. To suppress this, we penalize the discrepancy between the number of model calls (N_{call}) and actual environmental responses (N_{resp}):

$$R_{halluc} = \min(0, (N_{resp} - N_{call}) \times 0.2). \quad (7)$$

This mechanism enforces a strict "Turn-taking" protocol, ensuring that only tool calls actually processed by the environment are considered valid behaviors.

The final reward is calculated as:

$$R_{total} = R_{acc} + \lambda_{fmt} \cdot R_{fmt} + \lambda_{halluc} \cdot R_{halluc}. \quad (8)$$

In this paper, we set $\lambda_{fmt} = 0.1$ and $\lambda_{halluc} = 0.05$.

2.3 Tool Platform Construction

2.3.1 Tool Functions

In this section, we present the comprehensive multimodal toolkit within MindWatcher, comprising the following five tools:

Region Cropping/Zooming: This tool encompasses diverse image processing operations designed to externalize visual reasoning and highlight critical regions to guide attention. It includes an image grounding tool for localizing and cropping target areas based on input boxes, thereby facilitating the 'thinking with images' reasoning paradigm.

Object Grounding and Visual Search: This tool accepts image interest regions and search categories, subsequently retrieving the corresponding knowledge from a large-scale local image retrieval database (described in Sec 2.3.2). By adaptively localizing query-relevant regions, our tool performs precise and targeted regional searches, effectively addressing complex visual search challenges.

External Text Retrieval: This tool leverages the search engines for information retrieval. It accepts textual queries as input and returns the top-10 ranked results, each comprising a title and an abstract.

Webpage Content Extraction: Taking a URL as input, this tool employs Jina [2] to retrieve the webpage content. The agent can read its full content, the content within the window it provides, or use an AI assistant to generate a structured summary based on the specific goal it provides.

Local Code Interpreter: This tool executes Python code within a sandbox environment isolated from external resources (e.g., files and the internet). It returns the execution results and supports the invocation of various Python libraries for diverse data computation tasks.

2.3.2 Local Multi-modal Retrieval Library

Conventional image search methods leverage massive Internet data. However, directly acquired Internet resources in fine-grained specialized domains contain erroneous knowledge. Additionally, the high cost of external visual search API calls can significantly increase training expenses during large-scale training. To alleviate the above issues, we construct a large-scale, high-quality image retrieval database. Based on the general taxonomy of world knowledge, we built our local search database through the following procedure: (1) We established knowledge entries that span eight major categories: Person, Car, Plant, Animal, Logo, Landmark, Fruit & Vegetable, and Dish. (2) Collect images corresponding to these knowledge entries from both Internet sources and professional museum databases. (3) We employ domain experts to conduct large-scale comparative filtering and knowledge categorization. Through rigorous identification and curation by domain experts, we ensure that the precision of our visual search image database exceeds 99%.

Ultimately, our constructed specialized image retrieval database, MindWatcher Multi-modal Retrieval Database (MWRD), encompasses eight major categories of knowledge images and associated information, covering a total of 50k retrieval entities. Each retrieval unit contains 3-10 high-quality images, amounting to over 300k images. To accommodate temporally dynamic data, we perform regular maintenance and updates on this image retrieval database.

3 Training Data and MWE-Bench

The RL training data of MindWatcher includes both online and offline environment training data. The online environment refers to real interactions with internet environments. In this paper, the online training data comprises three distinct sources: two types of data constructed based on automated pipelines, and data collected from open-source datasets.

3.1 Training Data Constructed from Private Images

To enable MindWatcher agents to master multimodal tools proficiently, we constructed a cross-modal question-answering (QA) dataset with progressively increasing difficulty. Unlike purely textual tasks, this task requires agents to jointly invoke visual perception and external search tools to solve problems. To ensure data robustness and training efficiency, we designed a Multimodal Knowledge-Augmented Pipeline comprising three core stages: source knowledge anchoring and generation, rigorous QA quality validation, and difficulty grading based on tool invocation.

3.1.1 Phase 1: Source Knowledge Annotation and Initial Generation

We first utilize a high-quality private multimodal database as seed data to construct a foundational multimodal dataset. To achieve deep alignment between visual signals and textual knowledge, we designed a generation mechanism comprising the following steps:

Fine-Grained Visual-Knowledge Mapping: We developed an integrated data processing pipeline combining “object localization” and “fine-grained retrieval.” This automated pipeline extracts bounding boxes and corresponding retrieval labels from source images, establishing precise image-text mappings.

Knowledge Graph Augmentation: Based on extracted visual labels, we utilize web search to construct dynamic knowledge graphs, acquiring relevant background knowledge and factual information. This external knowledge is then leveraged to generate initial question-answer pairs, ensuring questions rely not only on images but also integrate external world knowledge.

3.1.2 Phase 2: Timeliness and Uniqueness Verification

The accuracy of reward signals is critical in reinforcement learning training. We found that directly generated QA data often faces two major challenges, which may lead to misjudgments in reward models:

1. **Temporal Stability:** Search engine environments are dynamically changing. If a time gap exists between data production and actual training, updates to search results may cause answer drift.
2. **Answer Uniqueness and Non-openness:** Open-ended questions often have non-unique solutions. Even if an agent executes the correct search path, its generated answer may contain only partial correct information or be overly broad, making it difficult for the reward model to evaluate. In response to these limitations, we implemented a two-stage human-in-the-loop verification pipeline. This rigorous review ensures the final high-quality multimodal dataset maintains temporal consistency, with each question possessing a unique, unambiguous ground truth.

3.1.3 Phase 3: Difficulty Grading Based on Tool Invocation

Curriculum learning is an effective strategy for training agents, hinging on reasonable difficulty stratification. However, traditional difficulty assessments based on human subjective perception are often biased.

In tool-integrated scenarios, search engines can instantly resolve memory-based problems deemed “difficult” by humans, creating a disconnect between perceived difficulty and the actual challenges faced by agents. To achieve more precise difficulty screening, we designed a Tool-Invocation Screening Engine. This engine abandons subjective judgments, instead

defining sample difficulty through quantitative analysis of the “number of tool invocation rounds” required to solve problems and the “complexity of multi-tool combinations.” This approach constructs training data that truly aligns with the agent learning curve.

3.2 Training Data Constructed from Open-sourced News

Constructing a reliable reward signal for RL in open-ended web interactions is notoriously difficult due to the noisy nature of internet content. General web corpora often abound with subjective commentary, unverifiable rumors, and ambiguous “clickbait” titles, which can lead the reward model to provide incorrect optimization signals. Furthermore, factual information in niche domains is often buried in low-traffic sub-pages that are difficult for generic search engines to index instantly, causing agents to fail even when their reasoning path is correct.

To mitigate these challenges, we selected **Sports News** as the seed domain for our automated pipeline. Sports data possesses unique characteristics ideal for training TIR agents:

1. **Objective Verifiability:** Unlike social news, sports events have definitive outcomes (scores, winners, rankings) that constitute a unique ground truth.
2. **Resistance to Ambiguity:** Statistical facts in sports are less susceptible to the semantic pollution of opinions or fake news.
3. **Multimodal Richness:** Match reports are intrinsically multimodal, requiring the alignment of textual statistics with visual evidence (player jerseys, scoreboards, action shots).

We developed a robust Temporal-Aware Multimodal QA Pipeline to harvest and process this data, consisting of three sequential stages: Ingestion, Semantic Auditing, and Constraint-Aware Generation.

3.2.1 Domain-Specific Ingestion and Filtering

We deployed a focused crawler targeting authoritative sports portals to ensure information reliability. The raw stream captures article metadata, textual bodies, and associated image sets. A preliminary heuristic filter is applied to discard low-quality samples, retaining only articles with non-empty bodies and at least one relevant image. This creates a raw repository of event-centric multimedia documents.

3.2.2 Phase 1: LLM-Based Semantic Auditing

Quality control is paramount for RL training. We introduce a “Data Auditor” agent (powered by a strong LLM) to perform a feasibility check before generation. The auditor evaluates raw news based on a strict *Factuality Protocol*:

Retention Criteria: The content must describe a completed event with a clear timeline (e.g., match results, completed transactions). The text must provide key information (entities, actions) visually corresponding to the images.

Rejection Criteria: Purely subjective content, such as rumors, predictions of future games, gossip, or vague summaries without verifiable details, is discarded.

This phase filters out approximately 40% of the raw feed, ensuring that the downstream generation model operates only on solid factual ground.

3.2.3 Phase 2: Constraint-Aware QA Generation

The surviving samples are processed by a “Question Generator” agent. To prevent the model from learning shortcuts or hallucinating, we designed a Constraint-Aware Prompting Strategy that enforces strict rules on the generated QA pairs:

1. **Temporal Anchoring:** A critical challenge in time-sensitive QA is “Data Rot”—a question like “Who won the game yesterday?” becomes invalid over time. Our pipeline forces the generator to explicitly resolve relative time expressions into absolute timestamps

(e.g., converting "this season" to "the 2025 season") based on the publish time of the article. This ensures the question remains valid and unique indefinitely.

2. Visual-Textual Dependency: Questions are engineered to require information integration from both modalities. For instance, instead of explicitly naming a player, the question might refer to "the player in the No. 8 jersey on the right," compelling the agent to first identify the visual entity and then search for its identity using external knowledge.

3. De-referencing Context: To simulate real-world user queries, we strictly prohibit meta-references such as "According to the article." The agent receives only the standalone question and the image, forcing it to use search tools to retrieve the knowledge originally contained in the source article (which is hidden from the agent during training).

3.3 Open-sourced Training Data and Offline Training Data

Supplementing the autonomously constructed online training data detailed previously, we curated a focused collection of open-source datasets. These are strategically utilized to bolster the proficiency of MindWatcher in text-only search tasks and code-augmented mathematical reasoning.

Furthermore, distinct from the three aforementioned data categories designed for real-world environment training (manual, online-automatic, and open-source), MindWatcher incorporates a specialized offline training method for TIR. To facilitate this, we developed an automated pipeline to construct a substantial corpus of high-quality, multimodal QA pairs with stratified difficulty levels.

3.4 MWE-Benchmark

The MWE-Bench covers six primary categories: Car, Animal, Plant, Person, Landmark, and Sports. While these categories align with those in our automated data construction pipeline, we deliberately adopted a distinct construction methodology for the benchmark to ensure its integrity and prevent performance inflation caused by data-domain overlap.

Specifically, for data derived from private images, we utilized knowledge entries from our internal database that were strictly excluded from the training set. The construction process followed a multi-stage pipeline: we first expanded our knowledge base by collecting auxiliary web-based information to enrich the context. For each category, we then applied category-specific constraints and employed closed-source models to perform "uniqueness deconstruction"—extracting core factual statements that uniquely identify an entity. These statements formed the basis for constructing initial single-turn QA pairs, which were subsequently synthesized into more complex and challenging multi-step reasoning tasks. Finally, all generated samples underwent a two-tier verification process involving both automated model-based filtering and manual expert review to ensure quality and temporal accuracy.

For sports category data, based on the data construction method outlined in Section 3.2, we merged text and image corpora belonging to the same entity or event across news data, which are from entirely non-overlapping time points. We then employed a powerful LLM to extract atomic facts from all corpora. Subsequently, we constructed QA pairs with complex queries based on these atomic facts. Finally, data cleaning and filtering were performed following a process similar to that described in Section 3.2.3.

Following the aforementioned methodology, we have successfully constructed MWE-Bench. The dataset encompasses six categories: 373 car-related instances, 351 animal-related instances, 397 plant-related instances, 63 person-related instances, 90 landmark-related instances, and 142 sports-related instances.

4 Experiment

4.1 Application Details

The training data utilized in this study are segmented across online and offline environments. In the online training environment, we collected VQA data consisting of 1,639 samples based on private images and 2,949 samples derived from public news sources. The open-source domain data, primarily extracted from established benchmarks such as WebSailor [20], ToolStar [11], and SimpleDeepSearcher [35], totaled 5,000 samples. Furthermore, we leveraged approximately 20,000 samples within the offline RL training environment.

The RL process employed a curriculum learning strategy guided by data difficulty. Training was conducted on the Qwen2.5-VL-32B [4] model for one epoch. Our training framework features a synchronized rollout mechanism coupled with an asynchronous tool invocation logic (Details shown in the Appendix A.2). Specifically, within each step of the interleaved CoT trajectory, the presence of the `<im_end>` token triggers an immediate check for a `<tool_call>` token. If a `<tool_call>` is detected, it is dispatched instantly. Reward computation also utilizes an asynchronous model invocation method.

We used the fully trained MindWatcher-32B model to distill its multimodal reasoning and tool-use capabilities into smaller, cost-effective models. This process involved collecting an initial, diverse set of base datasets, including the VLAA SFT dataset [5] (126K samples), the text-only WebWalker silver dataset [40] (15K samples), and a self-built multimodal RAG QA dataset (30K samples). The MindWatcher-32B “teacher” model was then employed to roll out and generate 1-3 corresponding TIR trajectories for each sample. After a straightforward filtering process, the final distillation dataset comprised 124K samples, split into 100K multimodal and 24K pure text samples. By using Qwen3-VL-2B [3], Qwen2.5-VL-3B, and Qwen3-VL-4B as base models and training them for one epoch on the distillation dataset, we successfully produced the smaller distilled MindWatcher-2B, MindWatcher-3B, and MindWatcher-4B models, respectively.

To comprehensively validate the performance of MindWatcher, in addition to the MWE-Bench, we conducted comparative testing against the model performance on several other open-source benchmarks, including MMSearch (subset) [15], SimpleVQA (subset) [8], and WabWalkerQA [40]. All tests were conducted with a sampling temperature of 0.7 and a top- p setting of 0.95. The primary evaluation metric utilized was *pass@1*, with correctness assessed by employing the LLM-as-Judges methodology.

4.2 Main Results

Tabel 1 presents the detailed performance of different backbones on the MWE-Bench under both direct inference and React/Agent inference modes.

Disparity in Parametric Knowledge. Under the Direct Inference mode, we observe that the freshness of a model knowledge cutoff does not linearly correlate with its benchmark performance. Despite being the most recent release, the Qwen3-VL series achieves an average score of only 22.60. In contrast, Gemini 2.5 Pro—notwithstanding an older knowledge cutoff—attains a SOTA zero-shot score of 42.09. This discrepancy underscores a critical reality: when internal parameters alone are insufficient for handling long-tail or specialized world knowledge, the integration of external reasoning tools is necessary.

Tool-Augmented Performance Leap. Transitioning to the ReAct/Agent paradigm catalyzes a significant performance surge for models previously limited by their internal knowledge. For instance, the score of Qwen3-VL 32B nearly triples when equipped with tool-use capabilities. Similarly, GPT-5 mini exhibits a remarkable explosion in performance within the Sports domain, soaring from 13.38 to 80.28 upon gaining tool access.

MindWatcher Dominance. MindWatcher-32B achieves overall SOTA performance on MWE-Bench with a global score of 75.35, outperforming prominent closed-source commercial models such as Gemini2.5 Flash and GPT-5 mini. Notably, MindWatcher achieves the highest accuracy across four specific domains: Vehicle, Animal, Plant, and Person. Furthermore, the distilled variants, including MindWatcher-2B, 3B, and 4B, demonstrate per-

Table 1: Results on the MindWatcher-Evaluation Benchmark.

Method	MindWatcher-Evaluation Benchmark						Avg.
	Car	Animal	Plant	Person	Landmark	Sport	
<i>Direct inference</i>							
GPT-4o mini	14.48	26.21	24.43	9.52	25.56	4.23	19.63
GPT-4o	24.13	38.46	29.22	12.7	38.89	6.34	27.75
Qwen2.5-VL-7B	7.24	31.62	30.48	7.94	17.78	2.82	20.06
Qwen2.5-VL-32B	18.5	33.62	33.25	9.52	33.33	8.45	25.92
Qwen3-VL 8B Thinking	10.46	22.79	23.93	15.87	17.78	16.90	18.65
Qwen3-VL 30B-A3B Thinking	17.69	25.64	25.69	19.05	20	15.50	21.89
Qwen3-VL 32B Thinking	14.21	28.21	26.45	20.63	21.11	21.83	22.60
Openai o4-mini	21.98	35.61	31.49	11.11	47.78	6.34	27.61
Gemini 2.5 Pro	40.48	49.57	48.11	28.57	46.67	14.09	42.09
Gemini 2.5 Flash	35.12	45.14	43.58	22.22	40	17.60	37.96
GPT-5 mini	22.25	48.15	42.07	4.76	44.44	13.38	33.97
<i>ReAct/Agent</i>							
Gemini 2.5 Flash	68.17	71.76	77.26	57.14	46.51	37.32	66.65
GPT-5 mini	61.93	68.66	81.61	44.44	57.78	80.28	<u>69.91</u>
Doubao-Seed-1.6-vision	57.64	57.26	61.46	65.08	38.89	39.39	57.91
Qwen2.5-VL-32B	51.74	54.13	58.19	50.79	41.11	31.69	51.41
Qwen3-VL 32B Thinking	58.98	75.5	77.83	69.84	<u>48.89</u>	<u>46.48</u>	66.95
WebWatcher-7B	32.71	45.58	49.62	30.16	36.67	21.57	39.66
WebWatcher-32B	38.87	53.85	67.51	52.38	42.22	32.88	50.93
MindWatcher-2B	50.13	80.34	86.4	65.08	44.44	16.90	64.76
MindWatcher-3B	49.87	80.06	85.89	63.49	44.44	17.61	64.48
MindWatcher-4B	56.03	<u>84.62</u>	<u>87.66</u>	<u>68.25</u>	41.11	36.62	69.63
MindWatcher-32B	71.31	86.04	88.92	77.78	47.78	<u>46.48</u>	75.35

*Best results are in bold and the suboptimal results are in underline.

formance comparable to the Qwen3-VL 32B baseline. This empirically demonstrates that robust tool-call capabilities can effectively mitigate the knowledge gaps typically present in small-parameter models.

Table 2 presents the comparative performance of MindWatcher-32B against other models in identical environments across two filtered multimodal subsets (MMSearch and SimpleVQA) and one pure-text benchmark (WebWalkerQA). MindWatcher continues to deliver SOTA results on MMSearch among all open- and closed-source models evaluated. On the SimpleVQA subset, MindWatcher performance surpasses the next-generation Qwen3-

Table 2: Results on the Open-sourced Benchmarks.

Method	MMSearch	SimpleVQA			WebWalkerQA		
		Chinese	English	Avg.	Chinese	English	Avg.
Qwen2.5-VL-32B	51.13	21.61	21.86	21.99	42.51	17.55	26.62
Qwen3-VL 32B thinking	57.01	20.78	34.85	28.91	57.09	24.71	36.47
Gemini 2.5 Flash	57.92	30.03	<u>39.96</u>	35.93	48.18	22.17	31.62
GPT-5 mini	57.47	<u>32.96</u>	32.9	33.29	65.99	37.41	47.79
Doubao-Seed-1.6-vision	58.37	39.06	46.1	43.44	<u>63.56</u>	<u>28.87</u>	<u>41.47</u>
WebWatcher-7B	49.68	21.33	26.41	24.41	40.89	19.63	27.35
WebWatcher-32B	56.45	21.61	37.01	30.49	39.68	19.63	26.91
MindWatcher-32B	58.82	22.16	39.83	32.32	50.2	23.79	33.38

*Best results are in bold and the suboptimal results are in underline.

Table 3: Comparison Results of the Distilled Models and their Base Models.

Method	MindWatcher-Evaluation Benchmark						Avg.
	Car	Animal	Plant	Person	Landmark	Sport	
Qwen3-VL 2B Thinking	39.95	70.37	80.1	50.79	24.44	29.51	51.41
MindWatcher-2B	50.13	80.34	86.4	65.08	44.44	16.90	64.76
Qwen2.5-VL 3B Instruct	17.43	35.04	37.78	6.3	12.22	0.00	24.93
MindWatcher-3B	49.87	80.06	85.89	63.49	44.44	17.61	64.48
Qwen3-VL 4B Thinking	56.3	81.2	80.86	74.6	38.89	32.79	66.53
MindWatcher-4B	56.03	84.62	87.66	68.25	41.11	39.51	69.63

VL-32B base model. Importantly, on the pure-text WebWalkerQA benchmark, MindWatcher remains highly competitive. Compared to its base model, Qwen2.5-VL-32B, results indicate that our continuous multimodal agentic RL training has successfully enhanced agent capabilities without compromising its foundational text-based reasoning.

Figure 3 presents the win-tie-loss analysis comparing MindWatcher-32B against four representative models: Qwen3-VL 32B Thinking, WebWatcher-32B, Gemini 2.5 Flash, and GPT-5 mini, across both public open-source benchmarks and our MWE-Bench. The results indicate that MindWatcher-32B consistently outperforms its parameter equivalent 32B counterparts in both evaluation settings. Notably, on the MWE-Bench, MindWatcher-32B demonstrates superior performance even when compared to SOTA closed-source models, specifically Gemini 2.5 Flash and GPT-5 mini.

Table 3 further details the performance gains achieved by the three distilled small-scale models relative to their respective foundation models on the MWE-Bench. Among these, MindWatcher-3B (derived from Qwen2.5-VL-3B-Instruct) exhibits the most significant improvement, with its proficiency score surging from 24.93 to 64.48. This substantial leap underscores the effectiveness of our distilled training approach in empowering small-scale models with robust agentic capabilities.

4.3 Analysis

4.3.1 The Impact of the Tool Capacity

During experiments, we find that the proficiency of the integrated tools is a pivotal determinant of an agent’s final performance. This is particularly evident in external retrieval tasks, where the indexing and recall mechanisms of different search engines lead to highly heterogeneous outcomes for identical queries. Beyond direct downstream performance, we observed that the choice of search engine during RL training induces distinct tool-call behavioral adaptations and search patterns within the model.

To quantify this impact, we conducted experiments using sports-related datasets, subdivided into two domains (Football and Basketball) and two languages (Chinese and English). We

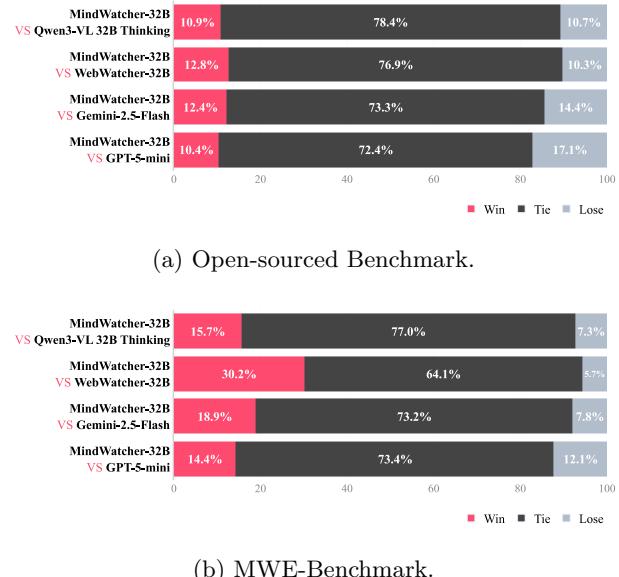


Figure 3: Benchmark Performance Comparison.

Table 4: Results on different search engines.

Search Engine	Basketball		Football		Avg.
	English	Chinese	English	Chinese	
Sogou Search	2.53	15.19	3.57	12.5	8.51
Bing search	<u>13.92</u>	<u>20.25</u>	<u>8.93</u>	<u>23.21</u>	<u>16.66</u>
Quark Search	20.25	39.24	28.57	55.36	34.81

* Best results are in bold and the suboptimal results are in underline.

evaluated the agent performance using three search engines—Sogou, Bing, and Quark—under a retrieval-only setting. The results, summarized in Table 4, demonstrate a substantial performance variance that frequently overshadows the variations attributed to algorithmic optimizations or foundation model scales. Specifically, in the most extreme case (football queries written by Chinese), the Quark search engine outperformed Sogou by a staggering 42.86%. However, these findings do not point to a universally "superior" search engine; rather, we found that the effectiveness of a search engine is highly volatile and contingent upon the specific domain and language of the query. This volatility highlights that the "capacity" of an agent is deeply coupled with its environment, suggesting that benchmark evaluations must account for tool-induced variance to ensure a fair assessment of a model's intrinsic reasoning abilities.

4.3.2 Genetic Inheritance in Agentic RL

We conduct a granular analysis of the relationship between tool-calling frequency and model accuracy. Specifically, we compare the behaviors and performance of MindWatcher against its own foundation model, Qwen2.5-VL-32B, and GPT-5 mini on the WME-Bench. The visualization of these results is presented in Figure 4.

Disparity in Decision Trigger Boundaries: A significant divergence is observed in the decision-making boundaries regarding tool invocation. As illustrated, GPT-5 mini opts to reason without any tool-calls (Round 0) in nearly one-sixth of the samples, yet achieves a mere 51.2% accuracy. This suggests a manifest "blind self-confidence" in GPT-5 mini; by relying on internal parameters for tasks requiring external knowledge, it forfeits substantial scores at the onset, leading to a lower overall performance compared to MindWatcher. Interestingly, when the number of tool-calls exceeds one, GPT-5 mini exhibits remarkable robustness in long-chain reasoning, with its accuracy showing negligible decay from Round 2 to Round 6.



(a) MindWatcher vs GPT-5 mini.

(b) MindWatcher vs Qwen2.5-VL-32B.

Figure 4: Comparison of Tool-use Behavior and Performance Distribution.

Case study 1: The Impact of Model World Knowledge on Agent Performance.



Question: In the image, a significant battle took place during the early 19th century in this northern South American city where a prominent female patriot played a crucial role by leading a company of soldiers and providing essential support to a key military leader's forces. This patriot was also known for her contributions to the independence movement of her country. In the main cathedral of this city, how many women are interred?

Ground-truth: 2

Qwen3-VL 32B:

<think>

Okay, let's try to figure this out. The user is asking about a northern South American city where a significant battle happened in the early 19th century. ...The flag has a red color with some emblem. Let me think. In South American independence, key figures include Simón Bolívar, and female figures like **Manuela Sáenz**. Wait, **Manuela Sáenz** was a prominent female patriot in the independence movement of Bolivia or Peru? ...<\think>

...

MindWatcher 32B:

<think>

The question provides enough context to proceed. The woman is described as a prominent female patriot who played a crucial role in a significant battle during the early 19th century in a northern South American city. She led a company of soldiers and provided essential support to a key military leader's forces. She is also known for her contributions to the independence movement of her country.

Given this information, **I can use the 'web_search' tool to search for prominent female patriots in northern South American cities who were involved in significant battles during the early 19th century**. Once I identify the correct patriot, I can determine the city she is associated with and then search for how many women are interred in the main cathedral of that city.

<\think>

<tool_call>

{ "name": "web_search", "arguments": { "query": prominent female patriot in northern South America early 19th century, } }

<\tool_call>

...

This phenomenon highlights that for high-capacity models, agentic performance can be bottlenecked by the decision trigger boundary rather than the executive action capability itself. Under autonomous settings, the model potential can be severely constrained by its initial failure to recognize the need for external tools.

Performance Shadowing and Genetic Inheritance in Agentic RL: While Mind-Watcher, trained via RL, significantly outperforms its foundation model (Qwen2.5-VL-32B), we observe a profound "Genetic Inheritance" in reasoning capacity. This is evidenced by the striking consistency in both accuracy trends and sample distribution across different tool-calling rounds.

As the required number of tool-calls increases, MindWatcher maintains a higher accuracy than Qwen2.5-VL-32B but fails to reverse the downward trend (identical decay slope) inherited from its foundation. Furthermore, MindWatcher's self-awareness—manifested in

its sample distribution across varying tool-call rounds—shows no significant deviation from Qwen2.5-VL-32B even after extensive RL training.

These observations suggest that while agentic RL can substantially refine tool-invocation and reasoning proficiency, it cannot fully breach the performance bottlenecks of the foundation model regarding long-range reasoning and multimodal processing. The foundation model imposes a fundamental performance constraint on the RL-derived agent; agentic RL serves as a strategy optimizer but remains fundamentally coupled with the base model’s capabilities. We term this phenomenon the **”Genetic Constraint” of the foundation model in agentic RL scenarios**. In Appendix A.3, we conducted further investigations into genetic inheritance in the agentic SFT scenario.

4.3.3 The Impact of Model World Knowledge on Agent Performance

Beyond numerical analysis, we conduct a case-based qualitative study to investigate how the world knowledge inherent in different foundation models affects downstream task performance. We observe that when the provided tools are insufficient for a ”low-knowledge” model to resolve a query, the model internal world knowledge becomes the decisive factor for downstream benchmark metrics.

Case Study 1 presents a visual comparison between MindWatcher (based on Qwen2.5-VL-32B) and the next-generation Qwen3-VL 32B Thinking on a specific case. In this example, neither model can correctly answer the question based solely on their internal capabilities. However, once provided with an external text-retrieval tool, a significant performance gap emerges: Qwen3-VL possesses the internal world knowledge to recognize the name of the person (Manuela Sáenz) in the artwork. This prior knowledge allows it to formulate precise search queries using the text-search tool, leading to a successful resolution. Conversely, MindWatcher (based on Qwen2.5-VL) lacks any prior information regarding this specific artwork. Without a starting point for inquiry and lacking auxiliary tools to bridge this knowledge gap, the model is unable to perform any viable reasoning or retrieval.

This case study demonstrates that in identical tool environments, performance metrics may not exclusively reflect TIR capabilities of a model. Current benchmarks contain a significant number of tasks that implicitly rely on the ”long-tail” knowledge of a model to catalyze the tool-use process. This coupling introduces substantial challenges in isolating and accurately evaluating the intrinsic TIR capacity of a model, as the benchmark results become confounded by the uneven distribution of world knowledge across different foundation models. For a given benchmark, when a vast number of queries cannot be adequately addressed by the provided tools, the evaluation of the agent functional capabilities essentially regresses into a test of the model’s internal world knowledge.

5 Related Work

5.1 TIR Agent

The landscape of TIR agents has witnessed a meteoric evolution over the past six months. By empowering models to autonomously select and invoke tools, the capability boundaries of LLMs—particularly those with smaller parameter counts—have been significantly expanded. However, a stark contrast persists: while contemporary LLMs exhibit reasoning capabilities comparable to human experts (often likened to ”PhD-level” cognition), their *action competence*—specifically, the precision and robustness of tool invocation—remains at a nascent, almost ”elementary” stage.

OpenAI o3 [23], as the first TIR agent deployed to a global user base, demonstrated astonishing proficiency. By actively manipulating images, performing complex calculations via code execution, and navigating file systems, o3 illuminated the vast potential of TIR agents to the research community. This paradigm shift catalyzed a surge of open-source initiatives inspired by o3, ranging from specialized code and search agents to DeepResearch systems. rStar2-Agent [28] leverages code execution as a verifier and solver to bolster mathematical reasoning. DeepEyes [44] introduces active visual tools, such as ”image zoom-in,” probing the ability of multimodal agents to resolve fine-grained visual details through iterative

manipulation. The Qwen DeepResearch [22, 13, 36] team has also made pivotal contributions to the open-source ecosystem by systematically diagnosing and addressing the multi-dimensional challenges inherent in long-horizon research tasks. Despite these strides, the chasm between an agent’s “thinking” and “acting” remains substantial. Critical challenges such as dynamic tool context management, long-term historical memory maintenance, and the attainment of training-free tool invocation capabilities represent significant hurdles that the field must address in the near future.

5.2 Training Method of TIR Agent

Unlike traditional LLM training, training TIR agents presents distinct challenges due to the necessity of interacting with external environments—specifically, executing tool calls and interpreting heterogeneous feedback during generation. Beyond mere planning, agents must learn to act adaptively within dynamic information contexts. Several works focus on the offline training stages of continual pre-training and SFT. [34] proposes an agent-specific continued pre-training method designed to endow base models with native action capabilities, thereby effectively supporting subsequent fine-tuning. WebDancer [39] and WebSailor [20] concentrate on methodologies for constructing high-quality TIR trajectories; while they incorporate elements of RL, they predominantly rely on SFT to shape agent behavior. The transition to online RL, which requires interaction with real-world environments, precipitates a steep rise in training complexity. Recent research has tailored algorithms specifically for this regime. ARPO [12] introduces an entropy balancing mechanism to prevent the training collapse often observed during TIR agent RL. LLDS [10] investigates the “lazy likelihood displacement” problem in agent RL, introducing likelihood preservation regularization to avert systemic stagnation in training. In summary, given the high cost of constructing premium trajectory data and the inherent instability of online RL, the path toward an optimal training strategy for TIR agents remains long and arduous.

6 Conclusion

In this paper, we introduced MindWatcher, a high-performance TIR agent. We developed robust visual question answering training data construction pipelines and proposed a specialized training methodology distinct from prior works. To empower MindWatcher, we established a comprehensive yet cost-effective multimodal toolbox and introduced MWE-Bench for rigorous performance evaluation. Our experiments demonstrate that MindWatcher, through its superior tool-invocation capabilities, can match or even exceed the performance of significantly larger or updated models. Beyond empirical results, this study reveals several critical experimental findings discovered during our development of the TIR agent. We hope our work provides unique insights and contributes to the future advancement of tool-augmented intelligence.

Acknowledgments and Disclosure of Funding

We extend our gratitude to Tongyi Qwen for their outstanding contributions to open-source LLMs; and we thank all colleagues at LiAuto Base Model for their support of the Mind-Watcher project.

Author List

Core Contributors (Equal contribution) Jiawei Chen Xintian Shen Lihao Zheng Zhenwei Shao

Contributors Hongyuan Zhang Pengfei Yu Xudong Rao Ning Mao Xiaobo Liu Lian Wen Chaoqun Du Feng Gu Wei He Qizhen Li Shanshan Li Zide Liu Jing Luo Lifu Mu Xuhao Pan Chang Ren Haoyi Sun Qian Wang Wei Wang Hongfu Yang Jiqing Zhan Chunpeng Zhou Zheng Zhou

Technique Leaders Hao Ma Tao Wei

Supervisors Pan Zhou Wei Chen

A Technical Appendices and Supplementary Material

A.1 Open-sourced Benchmark

In this work, rather than utilizing the full set or a naive random subset of open-source benchmarks, we implemented a rigorous data filtration pipeline. This decision stems from the observation that many existing benchmarks suffer from significant limitations, such as information lag due to insufficient temporal constraints. Furthermore, models released after benchmark publication may exhibit inflated performance due to inadvertent data leakage.

To address these issues, we established the Qwen3-VL 32B Thinking as a baseline for direct inference on the original benchmarks. All samples correctly answered by the model through direct inference were discarded. For the remaining samples, we conducted a meticulous manual review to filter out ambiguous questions or those with expired time-sensitivity. This process yielded a high-quality subset of open-source benchmarks, which we subsequently used to evaluate reasoning and tool-integrated capabilities under the ReAct/Agent paradigm.

Among these, the MMSearch subset contains 221 samples, while the simplevqa subset comprises 823 samples, including 361 Chinese samples and 462 English samples.

A.2 Infrastructure

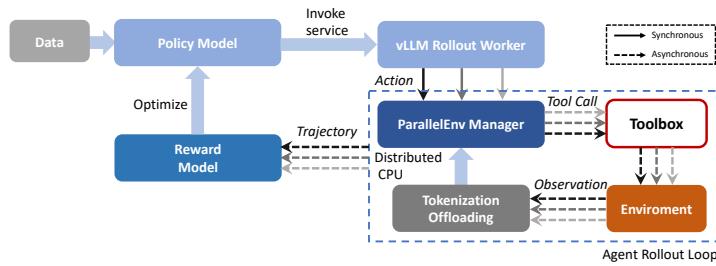


Figure 5: Step-wise Synchronous Sampling Framework of MindWatcher.

To facilitate efficient agentic reinforcement learning, we developed a step-wise synchronous sampling framework based on the Verl [32] to coordinate interactions between the agent and external environments as shown in Figure 5. In each rollout iteration, the vLLM engine [18] performs parallel batch inference to generate actions, followed by a synchronized barrier where the environment collects feedback. This design ensures trajectory consistency across massive batches while simplifying state management for the on-policy training process.

Empirical observations during training revealed that the primary bottleneck is not the trajectory generation itself, as the latency gap between synchronous and asynchronous sampling remains marginal. Instead, the dominant time expenditure arises from tool-calling latency. To mitigate this, we integrated an asynchronous tool invocation layer within the synchronous loop. By leveraging asyncio mechanisms and semaphore-based concurrency control, heterogeneous tools are dispatched and executed in parallel while strictly adhering to API QPS constraints. Furthermore, we implemented Tokenization Offloading, which offloads the computationally intensive task of processing environment observations' tokenization from the master node to distributed CPU workers. Additionally, the LLM-as-a-Judge reward model is invoked immediately upon the completion of each trajectory to minimize evaluation overhead. This hybrid architecture—synchroous in step control but asynchronous in tool execution—maximizes hardware utilization and significantly reduces the actual rollout time.

A.3 Genetic Inheritance in Agentic SFT

Building upon our analysis of genetic inheritance in the agentic RL paradigm (Section 4.3.2), we extend our investigation to the agentic SFT scenario using three distilled small-scale agent models. Figures 6, 7, and 8 illustrate the tool-use behavior and performance distribu-

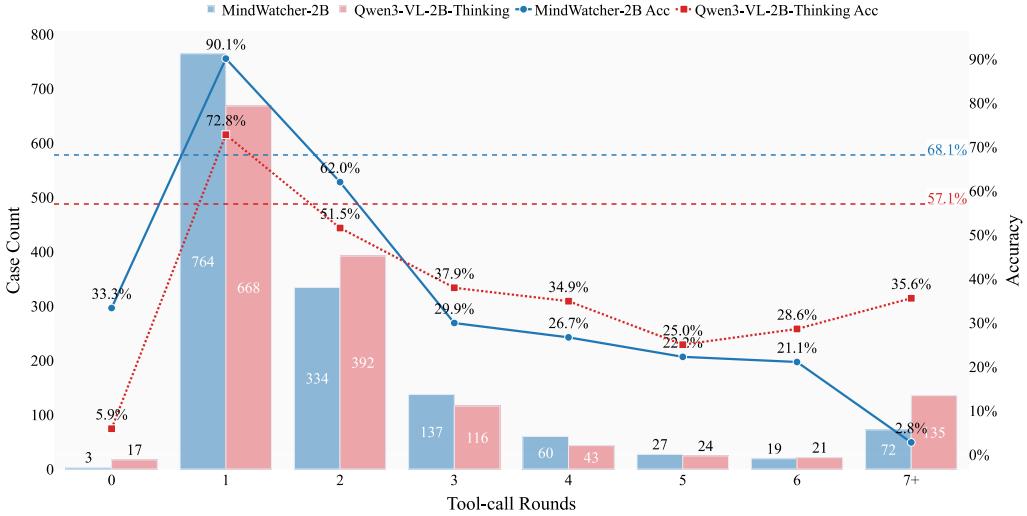


Figure 6: MindWatcher-2B vs Qwen3-VL 2B Thinking.

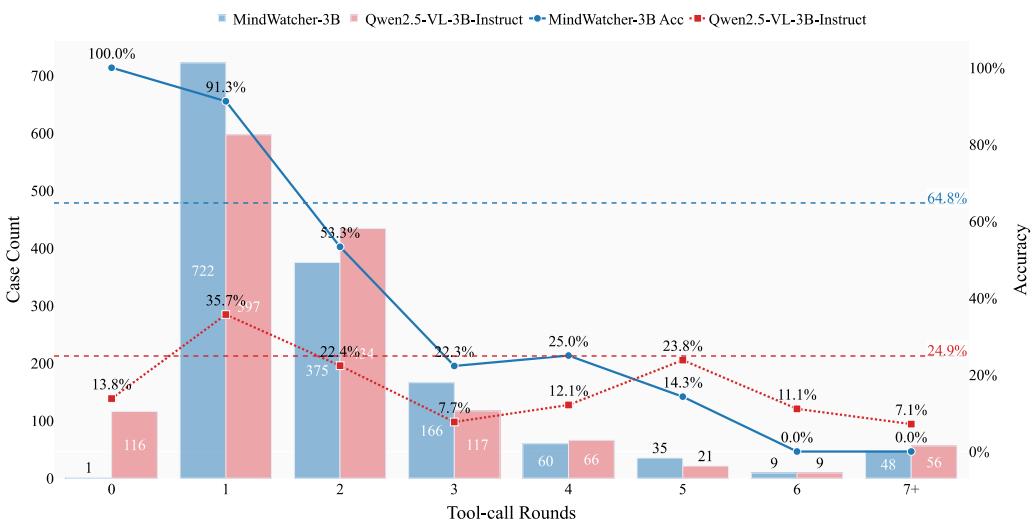


Figure 7: MindWatcher-3B vs Qwen2.5-VL-3B.

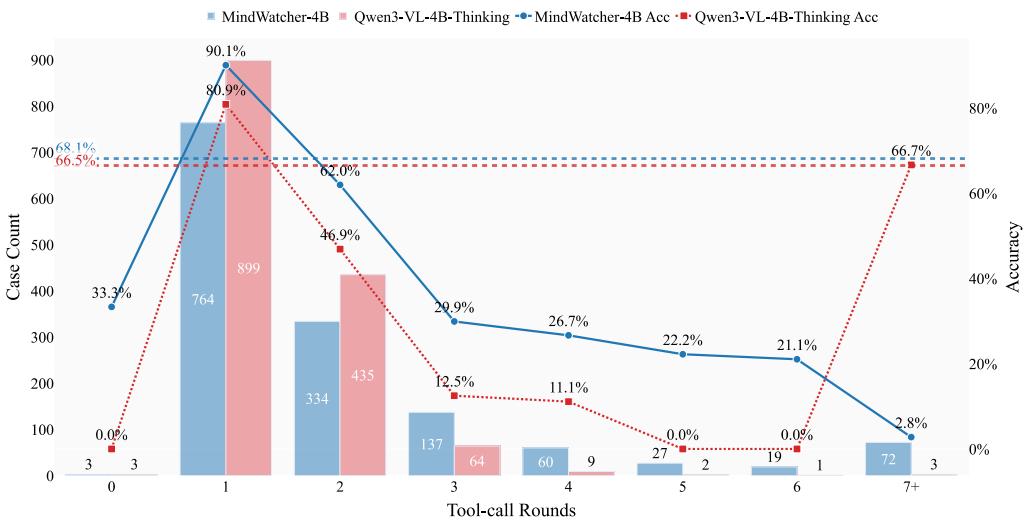


Figure 8: MindWatcher-4B vs Qwen3-VL 4B Thinking.

tions of MindWatcher-2B, 3B, and 4B alongside their respective foundations: Qwen3-VL-2B Thinking, Qwen2.5-VL-3B, and Qwen3-VL-4B Thinking.

Our observations indicate that, unlike the RL scenario, SFT-tuned models do not exhibit a consistent or predictable trend in tool-calling frequency relative to their base models. The decision trigger boundary in the SFT paradigm appears significantly less robust. For instance, after agentic SFT, the Qwen2.5-VL-3B model showed a dramatic shift, with Round 0 cases (no tool use) plummeting from 116 to just 1. Across the three distilled models, the distribution of tool-call rounds fluctuates inconsistently before and after SFT, lacking the stable behavioral alignment observed in the RL-tuned MindWatcher-32B.

Despite the behavioral volatility, the accuracy trends across different tool-call rounds reveal a phenomenon strikingly similar to that of agentic RL. As the complexity of the task increases (i.e., more tool-call rounds), both the SFT-tuned agents and their base models exhibit a synchronized downward trend in accuracy.

This reinforces the existence of "genetic inheritance" within the SFT paradigm: supervised fine-tuning is inherently limited by the base model capabilities in long-range reasoning and multimodal processing. Like RL, SFT serves as a method for policy alignment but fails to break through the fundamental cognitive "ceiling" established by the foundation model.

A key distinction arises in the "elegance" of the performance curves. In the agentic RL scenario, the performance of the agent and the base model decay at nearly identical slopes, showing a highly structured coupling. In contrast, the performance curves in the SFT are less congruent; while they share the same downward trajectory, the lack of a perfectly parallel slope suggests that SFT introduces more noise or less systematic optimization into the model's reasoning-tool-use integration compared to the more rigorous RL process.

A.4 Tool description for MindWatcher

Tool: Region Cropping/Zooming

Description: Zoom into a specific area of **the first input image** based on your provided bounding box.

input: the image and bounding box.

output:a new image.

Arguments:

- **bbox:** [x1, y1, x2, y2], # The bounding box you provided, where (x1, y1) is the top-left corner and (x2, y2) is the bottom-right corner.

Tool: Object Grounding & Visual Search

Description: Retrieve new similar images and their descriptions based on the provided bounding box area of **the first input image**.

input:image and bounding box.

output:only the most similar target's type name and the confidence score.

Arguments:

- **bbox_2d:**[x1, y1, x2, y2], # The bounding box you provided, where (x1, y1) is the top-left corner and (x2, y2) is the bottom-right corner.
- **category:**"the category" # The category of the image you want to search for, which can only be one of {plant, animal, car, person, landmark, vegetable, cuisine, logo}.

Tool: External Text Retrieval

Description: Retrieve external text information from the internet based on your provided text query.

input:only text query.

output:text.

Arguments:

- **query:** "the content"

Tool: Webpage Content Retrieval

Description: Visit a specified web page URL under 3 modes: you can read its full content, the content within the window you provided, or use an AI assistant to generate a structured summary based on the specific goal you provide.

input:A JSON object containing three arguments: url, window, goal.

output:A JSON object containing the visited URL and a structured result.

Arguments:

- **url:** "https://example.com/article", # The webpage url you want to visit.
- **window:** [a, b], # Select the content you want to read within [a, b] (Optional)
- **goal:** "". # What you want to get or find. (Optional)

Tool: Local Code Interpreter

Description: Execute Python code in a sandboxed environment. Useful for calculations, data manipulation, and other programmatic tasks. This tool is isolated and cannot access files or the internet.

input:A string of Python code. The tool call format can be a standard JSON object or a Python dictionary literal.

output:A JSON object containing stdout (from print statements), stderr (error messages), and the final result value.

Arguments:

- **code:** "the code such as: a = 5\nb = 10\nresult = a * b\nprint(result)" # The Python code to execute.

A.5 Trajectory Display

In this section, we visualize tool calling trajectory of MindWatcher to highlight the model multimodal chain-of-thought reasoning and interleaved thinking capabilities.

Case1



Question: 在 2025 年 11 月 18 日，图中这名球员与耐克续签代言后，他成为现役 NBA/WNBA 与耐克仍有签名鞋合同的第几位球星？截至 2025 年 11 月 18 日，图中这名来自底特律球队的球员的个人耐克签名鞋预计在哪个赛季首次亮相？在 2025 年 11 月 18 日这次续约时，图中球员所属的联赛共有多少位现役球员拥有与耐克的签名鞋合同？以 2025 年 11 月 18 日为准，与图中这名球员所属联赛不同的另一联赛共有多少位现役球员拥有与耐克的签名鞋合同？在 2025 年 11 月 18 日的情况下，图中这名身穿“DETROIT”球衣的球员与耐克续签的代言合同期限是几年？

Ground-truth: 第 9 位；2026-27 赛季；6 位；3 位；6 年

Step 1:

<think>

为了回答这个问题，我需要识别出图片中的球员，并获取以下信息：

1. 截至 2025 年 11 月 18 日，这位球员与耐克续签代言后，成为现役 NBA/WNBA 与耐克仍有签名鞋合同的第几位球星。
2. 截至 2025 年 11 月 18 日，这位球员的个人耐克签名鞋预计在哪个赛季首次亮相。
3. 在 2025 年 11 月 18 日这次续约时，这位球员所属的联赛共有多少位现役球员拥有与耐克的签名鞋合同。
4. 截至 2025 年 11 月 18 日，与这位球员所属联赛不同的另一联赛共有多少位现役球员拥有与耐克的签名鞋合同。
5. 截至 2025 年 11 月 18 日，这位球员与耐克续签的代言合同期限是几年。

首先，我需要识别出图片中的球员。可以通过调用工具来帮助我完成这个任务。

<\think>

<tool_call>

```
{ "name": "zoom_v_search", "arguments": { "bbox_2d": [194 649 807 676] } }
```

<\tool_call>

<tool_response>

名称：凯德·坎宁安，检索置信度：0.81

<\tool_response>

Step 2:

<think>

通过使用工具，我识别出图片中的球员是凯德·坎宁安（Cade Cunningham）。接下来，我需要查找与他相关的耐克代言信息，包括他在 2025 年 11 月 18 日与耐克续签代言后的排名、签名鞋首次亮相的赛季、续约时现役球员中有多少人拥有与耐克的签名鞋合同、WNBA 中拥有耐克签名鞋合同的现役球员数量以及他与耐克续约的合同期限。为了获取这些信息，我将使用网络搜索工具。

<\think>

<tool_call>

```
{ "name": "web_search", "arguments": { "query": "凯德·坎宁安 2025 年 11 月 18 日耐克续签" } }
```

<\tool_call>

Case1

```
<tool_response>
{
  'search_query': '凯德·坎宁安 2025 年 11 月 18 日耐克续签',
  'search_result': [
    {
      'id': 1,
      'title': '下一位超巨！坎宁安与耐克达成 6 年续约：签名鞋将于下赛季亮相',
      'content': '北京时间 11 月 18 日，底特律活塞队的后卫凯德·坎宁安将继续与耐克公司保持合作关系... 根据知名记者 Shams 的报道，坎宁安已经与耐克签署了一份新的为期六年的代言合同...',
      'url': 'https://m.sohu.com/a/956031720_121722346/?pvid=000115_3w_a',
      'date': '2025-11-19T10:43:00+08:00'
    },
    ...
    {
      'id': 8,
      'title': '耐克续约坎宁安，他将成为后詹杜库时代美国本土球员耐克的新门面',
      'content': '耐克把六年合同拍在坎宁安面前... 2027 年才能买到的第一双 CC 战靴，其实今天就已经在生产线排号...',
      'url': 'https://m.toutiao.com/a7574051184151937574/'
    }
  ]
}<\tool_response>
```

Step 3:

<think>

通过网络搜索，我获得了以下信息：

1. 凯德·坎宁安在 2025 年 11 月 18 日与耐克续签了一份为期六年的代言合同。
2. 他的签名鞋预计将在 2026-27 赛季的下半程正式发布。
3. 他是现役第六位拥有耐克签名鞋的 NBA 球员，与勒布朗·詹姆斯、凯文·杜兰特、扬尼斯·安特托昆博等超级巨星并列顶尖阵营。

根据这些信息，我可以回答问题：

1. 截至 2025 年 11 月 18 日，凯德·坎宁安与耐克续签代言后，成为现役 NBA/WNBA 与耐克仍有签名鞋合同的第 6 位球星。
2. 截至 2025 年 11 月 18 日，凯德·坎宁安的个人耐克签名鞋预计在 2026-27 赛季首次亮相。
3. 在 2025 年 11 月 18 日这次续约时，凯德·坎宁安所属的联赛（NBA）共有 6 位现役球员拥有与耐克的签名鞋合同。
4. 关于 WNBA 中拥有耐克签名鞋合同的现役球员数量，我没有找到具体的信息。为了回答这个问题，我将使用网络搜索工具。
5. 截至 2025 年 11 月 18 日，凯德·坎宁安与耐克续签的代言合同期限是 6 年。

接下来，我将使用网络搜索工具来查找 WNBA 中拥有耐克签名鞋合同的现役球员数量。

<\think>

```
<tool_call>
{ "name": "web_search", "arguments": { "query": "WNBA 现役耐克签名鞋合同数量" } }
<\tool_call>
```

<tool_response>

```
{
  'search_query': 'WNBA 现役耐克签名鞋合同数量',
  'search_result': [
    {
      'id': 1,
      'title': 'Nike 签约 6 位 NBA 及 3 位 WNBA 球员拥有签名鞋',
      'content': '目前拥有 Nike 签名鞋合同的 6 位 NBA 球员和 3 位 WNBA 球员。目前拥有 Nike 签名鞋合同的 6 位 NBA 球员和 3 位 WNBA 球员。...',
      'url': 'http://cj.sina.cn/articles/view/2432009827/p90f5866302701jglm',
      'date': '2025-11-18T12:00:00+08:00'
    },
    ...
}
```

Case1

```
{  
    'id': 8,  
    'title': '现役拥有 Nike 签名鞋合同 9 位球星！NBA6 人 WNBA3 人！',  
    'content': '11 月 18 日讯据 Shams Charania 报道，活塞球星康宁汉姆和耐克续签一份为 6 年的代言合同，并将拥有自己的签名鞋。球员的个人签名鞋预计将 26-27 赛季首次亮相。康宁汉姆也成为了 NBA/WNBA 现役球员中位与耐克还有签名鞋合同的球星...',  
    'url': 'http://17449.xdgk.com.cn/dianshiju/38952adc0ec9b8fNw6lg.html'  
    'date': '2025-12-20T01:37:34+08:00'  
}  
]  
}  
<\tool_response>
```

Step 4:

<think>

通过网络搜索，我获得了以下信息：

1. WNBA 中拥有耐克签名鞋合同的现役球员数量为 3 人。

现在我已经有了所有问题的答案：

1. 截至 2025 年 11 月 18 日，凯德·坎宁安与耐克续签代言后，成为现役 NBA/WNBA 与耐克仍有签名鞋合同的第 6 位球星。

2. 截至 2025 年 11 月 18 日，凯德·坎宁安的个人耐克签名鞋预计在 2026-27 赛季首次亮相。

3. 在 2025 年 11 月 18 日这次续约时，凯德·坎宁安所属的联赛（NBA）共有 6 位现役球员拥有与耐克的签名鞋合同。

4. 截至 2025 年 11 月 18 日，与凯德·坎宁安所属联赛不同的另一联赛（WNBA）共有 3 位现役球员拥有与耐克的签名鞋合同。

5. 截至 2025 年 11 月 18 日，凯德·坎宁安与耐克续签的代言合同期限是 6 年。

我将给出最终答案。

<\think>

<answer>

1. 截至 2025 年 11 月 18 日，凯德·坎宁安与耐克续签代言后，成为现役 NBA/WNBA 与耐克仍有签名鞋合同的第 6 位球星。

2. 截至 2025 年 11 月 18 日，凯德·坎宁安的个人耐克签名鞋预计在 2026-27 赛季首次亮相。

3. 在 2025 年 11 月 18 日这次续约时，凯德·坎宁安所属的联赛（NBA）共有 6 位现役球员拥有与耐克的签名鞋合同。

4. 截至 2025 年 11 月 18 日，与凯德·坎宁安所属联赛不同的另一联赛（WNBA）共有 3 位现役球员拥有与耐克的签名鞋合同。

5. 截至 2025 年 11 月 18 日，凯德·坎宁安与耐克续签的代言合同期限是 6 年。

<\answer>

Case2



Question: 这座山体雕刻上刻画了哪几位美国总统的面孔？其中哪一位人物在独立战争中发挥了关键作用？他哪一年写了信件，这封信授权美国首次正式开展间谍活动？

Ground-truth: 刻画了乔治·华盛顿、托马斯·杰斐逊、西奥多·罗斯福、亚伯拉罕·林肯；乔治·华盛顿在美国独立战争中发挥了关键作用；1777年

Step 1:

<think>

这张图片显示的是美国著名的地标——拉什莫尔山（Mount Rushmore）。根据图片中的特征，我可以识别出这座山体雕刻上刻画了几位美国总统的面孔。为了确认具体是哪几位总统，并找到与独立战争相关的关键人物及其写信授权美国首次正式开展间谍活动的年份，我需要调用工具来获取更多信息。

<\think>

<tool_call>

```
{ "name": "zoom_v_search", "arguments": { "bbox_2d": [193, 193, 806, 481], "category": "landmark" } }
```

<\tool_call>

<tool_response>

名称：拉什莫尔山国家纪念公园，检索置信度：0.97

<\tool_response>

Step 2:

<think>

根据工具返回的结果，确认这座山体雕刻是拉什莫尔山国家纪念公园。拉什莫尔山上雕刻了四位美国总统的面孔：乔治·华盛顿、托马斯·杰斐逊、西奥多·罗斯福和亚伯拉罕·林肯。现在，我需要找出其中哪一位在独立战争中发挥了关键作用，并确定他写信授权美国首次正式开展间谍活动的年份。为了获取这些详细信息，我将使用 web_search 工具来查找相关信息。

<\think>

<tool_call>

```
{ "name": "web_search", "arguments": { "query": "拉什莫尔山四位美国总统独立战争关键人物首次正式开展间谍活动" } }
```

<\tool_call>

<tool_response>

{

'search_query': '拉什莫尔山四位美国总统独立战争关键人物首次正式开展间谍活动',

'search_result': [

{

'id': 1,

'title': '美国南达科他州拉什莫尔山上的四个人头像是哪四位啊',

'content': '美国南达科他州拉什莫尔山上的四个人头像是哪四位啊美国南达科他州拉什莫尔山上的四座高 60 英尺的总统头像，象征着美国建国 150 年来的历史。这四位总统分别是乔治·华盛顿、托马斯·杰弗逊、西奥多·罗斯福和亚伯拉罕·林肯。华盛顿作为美...',

'url': 'https://zhidao.baidu.com/question/572951197473818884.html',

'date': '2024-10-23T00:00:00+08:00'

},

...

{

'id': 5,

'title': '独立战争期间美国反间谍第一人',

Case2

```
'content': '独立战争期间美国反间谍第一人，美国，威廉，杰伊，约翰，华盛顿，反间谍，克劳斯比，独立战争',
'url': 'https://m.163.com/dy/article/K30AB0GV0514R9P4.html'
'date': '2025-06-26T16:33:59+08:00'
}
]
}
<\tool_response>
```

Step 3:

<think>

通过工具返回的搜索结果，我了解到拉什莫尔山上雕刻的四位美国总统是乔治·华盛顿、托马斯·杰斐逊、西奥多·罗斯福和亚伯拉罕·林肯。为了找到在独立战争中发挥了关键作用的人物以及他写信授权美国首次正式开展间谍活动的年份，我需要进一步查看搜索结果中的相关内容。根据第三个搜索结果，提到独立战争期间美国反间谍第一人，并且涉及威廉·杰伊和约翰·华盛顿等人物。为了获取更多详细信息，我将访问相关网页以获取具体细节。

<\think>

<tool_call>

```
{ "name": "url_visit", "arguments": { "url": "https://m.163.com/dy/article/K30AB0GV0514R9P4.html", "goal": "找到在独立战争中发挥了关键作用的人物以及授权年份" } }
```

<\tool_call>

<tool_response>

{

```
"visited_url": "https://m.163.com/dy/article/K30AB0GV0514R9P4.html",
```

```
"mode": "summarize",
```

```
"result": {
```

"rational": "The user seeks to identify a key figure in the American Revolutionary War and the year they authorized formal espionage activities. The content highlights John Jay's pivotal role in establishing counterintelligence and mentions Washington's 1776 correspondence about British spies. The formation of Jay's committee in 1776 directly addresses the user's goal.",

"evidence": [

"华盛顿能够用谍报情，就知道反间工作同样必不可少。1776年3月24日，华盛顿在与部属交流的信中写道：‘我最担心的一件事，是他们的间谍，因此，我希望对它们保持最严密的监视。’于是，浩浩荡荡的清除英国间谍活动在大陆军中展开。”，”1776年6月，年仅30岁的约翰·杰伊受华盛顿的委派成立了‘纽约州侦查和挫败阴谋委员会’，负责情报搜集、逮捕英国间谍和调查可疑的英国同情者，拥有逮捕、审判、监禁、驱逐嫌疑犯的权力，约翰·杰伊出任该委员会主席至1777年2月中旬。”

],

"summary": "在独立战争中，约翰·杰伊(John Jay)发挥了关键作用。1776年3月24日，华盛顿在信中表达了对英国间谍的担忧，随后于同年6月授权杰伊成立‘纽约州侦查和挫败阴谋委员会’，这是美国首次正式开展反间谍活动的机构。该委员会赋予杰伊逮捕、审判和驱逐英国间谍及同情者的权力，标志着美国早期情报工作的制度化。”

}

<\tool_response>

Step 4:

<think>

根据工具返回的结果，我已经找到了问题的答案。拉什莫尔山上雕刻的四位美国总统是乔治·华盛顿、托马斯·杰斐逊、西奥多·罗斯福和亚伯拉罕·林肯。在独立战争中发挥关键作用的人物是约翰·杰伊(John Jay)。他在1776年被授权组建“纽约州侦查和挫败阴谋委员会”，这是美国首次正式开展反间谍活动的机构。因此，我可以给出最终答案。

<\think>

<answer>

拉什莫尔山上雕刻的四位美国总统是乔治·华盛顿、托马斯·杰斐逊、西奥多·罗斯福和亚伯拉罕·林肯。在独立战争中发挥关键作用的人物是约翰·杰伊(John Jay)。他在1776年被授权组建“纽约州侦查和挫败阴谋委员会”，这是美国首次正式开展反间谍活动的机构。

<\answer>

Case3



Question: What are the functions of these glasses?

Ground-truth: Camera; Open-back headphones; Vehicle control; AI assistant

Step 1:

<think>

The image shows a disassembled pair of glasses with various components labeled. To determine the functions of these glasses, I need to identify each component and understand their roles. However, the labels on the components are not clearly visible in the current image. I will zoom in on the image to get a better view of the labels and components.

<\think>

<tool_call>

{ "name": "im_zoom_in", "arguments": { "bbox_2d": [485, 385, 725, 580] } }

<\tool_call>

<tool_response>

<image>



<\tool_response>

Step 2:

<think>

The zoomed-in image provides a clearer view of the glasses and some of their components. However, the labels on the components are still not legible. To understand the functions of these glasses, I need to identify the components and their roles. Since the labels are not readable, I will perform a web search using the brand name "Livis" to find more information about these glasses and their functionalities.

<\think>

<tool_call>

{ "name": "web_search", "arguments": { "query": 'Livis glasses functions' } }

<\tool_call>

<tool_response>

{

 'search_query': 'Livis glasses functions',

 'search_result': [

 {

 'id': 1,
 'title': '视频评测 | 理想 Livis：首先是副好眼镜，这真的很重要'}

Case3

```
'content': '而为了让 Livis 作为一副眼镜更好用，理想也在镜片上花了不少功夫。理  
想 Livis 的镜片合作方为百年光学大厂蔡司...',  
    'url': 'https://new.qq.com/rain/a/LNK2025120406148900',  
    'date': '2025-12-04T12:05:00+08:00'  
},  
...  
{  
    'id': 5,  
    'title': '1699 元！今年最值得买的智能眼镜 Livis，居然来自一家车企',  
    'content': '依靠免唤醒和意图识别，你可以在任何地方直接说打开后备箱、查询胎  
压或者打开座椅加热，完全不需要掏出手机...',  
    'url': 'https://new.qq.com/rain/a/LNK2025120312421600"  
    'date': '2025-12-03T22:21:00+08:00'  
}  
]  
}  
<\tool_response>
```

Step 3:

```
<think>  
I am now ready to provide the final answer based on the information gathered.  
</think>
```

<answer>

The Livis glasses have multiple functions, including:

1. **Optical Function**: High-quality lenses made in collaboration with Carl Zeiss, offering a wide range of vision correction options.
2. **Sound System**: Sophisticated dual magnetic circuit sandwich speakers with intelligent frequency adjustment for leak-proof functionality and enhanced audio experience.
3. **Camera Function**: A 12-megapixel IMX681 sensor with a 105-degree ultra-wide-angle lens, supporting flexible multi-format shooting and live photo capture.
4. **Vehicle Control**: Voice-controlled vehicle operations, including opening the trunk, checking tire pressure, and adjusting seat heating.
5. **Voice Assistant**: Integrated smart voice assistant for hands-free communication and access to various services.
6. **Wireless Charging**: Dedicated charging case that supports wireless charging.
7. **AI Features**: Advanced AI technology for recognizing user intentions, enhancing audio quality, and optimizing user interactions.

These features make Livis glasses a multifunctional device that combines eyewear, audio, camera, vehicle control, and AI capabilities.

```
<\answer>
```

A.6 Prompt Design

In this section, we display the prompts utilized by policy model and evaluation judge model.

Prompt: Policy Model

你是一个 ReAct 范式的 agent，能够接受文本输入，回答用户问题对于一些复杂的问题，你可以选择调用网络搜索工具帮助你解决问题你可以调用的工具包括：

- (** 具体描述请参考工具描述部分 **)
1. Region Cropping/Zooming
 2. Object Grounding Visual Search
 3. External Text Retrieval
 4. Webpage Content Retrieval

对于每一个问题，你需要先思考，然后调用工具（如果需要），你会得到工具调用返回的结果，还可以根据工具的返回结果进行进一步的思考，最后给出答案你的思考过程，工具调用请求以及回答需要严格按照以下格式：

```
<\think>
你的思考过程
<\think>
<\tool_call>
"name": <function-name>, "arguments": <args-json-object>
<\tool_call> (如果需要调用工具，你的工具调用请求参考 usage 中的示例)
<\think>
你的思考过程
<\think> (如果需要进一步思考)
<\answer>
你的最终答案
<\answer>
```

请记住，你在每次调用工具之后，也就是输出 </tool_call> 之后，都需要结束本轮对话，等待工具调用的结果返回，再进行后续动作在输出回答之后，即在输出 <answer> 之后，你需要立即结束本轮对话，不要再输出任何内容你的思考次数和工具调用次数没有限制，但必须在最后给出你的答案对于任何问题，你不应该拒绝回答，而应该通过不断思考或调用工具，直到得到确信的结果。

Prompt: Evaluation Judge Model

你是一个公正的评测员，负责判断 【模型回答】 是否在事实和逻辑上符合 【标准答案】。

输入信息：

- ** 问题 **: question
- ** 模型回答 **: output
- ** 标准答案 (GT)**: ground_truth

请遵循以下 ** 通用等价性原则 ** 进行判断：

1. ** 核心事实一致性 **:
 - 忽略措辞、语序或详细程度的差异。只要大概语义一致，即为正确。
 - 2. ** 数学与单位等价性 **:
 - ** 自动换算 **: 如果模型回答的数值单位与 GT 或问题要求的单位不同，** 必须先进行数学换算 ** 再比较。
 - ** 精度容忍 **: 允许合理的精度误差，换算后的前两位有效数字正确就行。
 - ** 格式识别 **: 识别"45,610" (带逗号) 和"45610" 为同一个数。
 - 3. ** 跨语言与语义对齐 **:
 - 忽略中英文差异。
 - 识别同义指代 (例如"NYC" = "New York City"; "巴特勒码头" = "Butler's Wharf")。
 - 4. ** 非法问题处理 **:
 - 如果 GT 表示问题无解 (如"invalid question", "no answer")，只要模型 ** 否定了事件发生 ** 或 ** 拒绝回答 **，均判定为 ** 正确 **。
 - 只有当模型顺着错误前提编造事实时，才判错。
 - 5. ** 拒答处理 **:
 - 如果 GT 有明确答案，而模型回答“不知道”或“无法回答”，判定为错误。

Prompt: Evaluation Judge Model

Respond strictly in a valid JSON object matching this schema:

```
{  
    'extracted_final_answer': string,  
    'reasoning': string,  
    'result': "1" | "0",  
    'confidence': integer,  
    'strict': true  
}
```

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Jina AI. Jina, 2025. <https://jina.ai/>.
- [3] Shuai Bai, Yuxuan Cai, Ruizhe Chen, et al. Qwen3-vl technical report, 2025. URL <https://arxiv.org/abs/2511.21631>.
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [5] Hardy Chen, Haojin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. Sft or rl? an early investigation into training r1-like reasoning large vision-language models. *arXiv preprint arXiv:2504.11468*, 2025.
- [6] Jiawei Chen, Dingkang Yang, Yue Jiang, Mingcheng Li, Jinjie Wei, Xiaolu Hou, and Lihua Zhang. Efficiency in focus: Layernorm as a catalyst for fine-tuning medical visual language models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 3122–3130, 2024.
- [7] Jiawei Chen, Dingkang Yang, Tong Wu, Yue Jiang, Xiaolu Hou, Mingcheng Li, Shunli Wang, Dongling Xiao, Ke Li, and Lihua Zhang. Detecting and evaluating medical hallucinations in large vision language models. *arXiv preprint arXiv:2406.10185*, 2024.
- [8] Xianfu Cheng, Wei Zhang, Shiwei Zhang, Jian Yang, Xiangyuan Guan, Xianjie Wu, Xiang Li, Ge Zhang, Jiaheng Liu, Yuying Mai, et al. Simplevqa: Multimodal factuality evaluation for multimodal large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4637–4646, 2025.
- [9] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, 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.
- [10] Wenlong Deng, Yushu Li, Boying Gong, Yi Ren, Christos Thrampoulidis, and Xiaoxiao Li. On grpo collapse in search-r1: The lazy likelihood-displacement death spiral. *arXiv preprint arXiv:2512.04220*, 2025.
- [11] Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui Zhou, Zhicheng Dou, and Ji-Rong Wen. Tool-star: Empowering llm-brained multi-tool reasoner via reinforcement learning. *arXiv preprint arXiv:2505.16410*, 2025.
- [12] Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, et al. Agentic reinforced policy optimization. *arXiv preprint arXiv:2507.19849*, 2025.
- [13] Xinyu Geng, Peng Xia, Zhen Zhang, Xinyu Wang, Qiuchen Wang, Ruixue Ding, Chenxi Wang, Jialong Wu, Yida Zhao, Kuan Li, et al. Webwatcher: Breaking new frontier of vision-language deep research agent. *arXiv preprint arXiv:2508.05748*, 2025.
- [14] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [15] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanmin Wu, Jiayi Lei, Pengshuo Qiu, Pan Lu, Zehui Chen, Chaoyou Fu, Guanglu Song, et al. Mmsearch: Benchmarking the potential of large models as multi-modal search engines. *arXiv preprint arXiv:2409.12959*, 2024.

- [16] Yue Jiang, Jiawei Chen, Dingkang Yang, Mingcheng Li, Shunli Wang, Tong Wu, Ke Li, and Lihua Zhang. Comt: Chain-of-medical-thought reduces hallucination in medical report generation. In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2025. doi: 10.1109/ICASSP49660.2025.10887699.
- [17] Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive nlp. *arXiv preprint arXiv:2212.14024*, 2022.
- [18] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [19] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.
- [20] Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, et al. Websailor: Navigating super-human reasoning for web agent. *arXiv preprint arXiv:2507.02592*, 2025.
- [21] Mingcheng Li, Xiaolu Hou, Ziyang Liu, Dingkang Yang, Ziyun Qian, Jiawei Chen, Jinjie Wei, Yue Jiang, Qingyao Xu, and Lihua Zhang. Mccd: Multi-agent collaboration-based compositional diffusion for complex text-to-image generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13263–13272, 2025.
- [22] Zijian Li, Xin Guan, Bo Zhang, Shen Huang, Houquan Zhou, Shaopeng Lai, Ming Yan, Yong Jiang, Pengjun Xie, Fei Huang, et al. Webweaver: Structuring web-scale evidence with dynamic outlines for open-ended deep research. *arXiv preprint arXiv:2509.13312*, 2025.
- [23] OpenAI. Introducing openai o3 and o4-mini. <https://openai.com/index/introducing-o3-and-o4-mini/>, April 2025. Accessed: 2025-12-19.
- [24] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [25] MindGPT ov Team. Mindgpt-4ov: An enhanced mllm via a multi-stage post-training paradigm. *arXiv preprint arXiv:2512.02895*, 2025.
- [26] Zile Qiao, Guoxin Chen, Xuanzhong Chen, Donglei Yu, Wenbiao Yin, Xinyu Wang, Zhen Zhang, Baixuan Li, Huifeng Yin, Kuan Li, et al. Webresearcher: Unleashing unbounded reasoning capability in long-horizon agents. *arXiv preprint arXiv:2509.13309*, 2025.
- [27] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training, 2018.
- [28] Ning Shang, Yifei Liu, Yi Zhu, Li Lyna Zhang, Weijiang Xu, Xinyu Guan, Buze Zhang, Bingcheng Dong, Xudong Zhou, Bowen Zhang, et al. rstar2-agent: Agentic reasoning technical report. *arXiv preprint arXiv:2508.20722*, 2025.
- [29] Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. Prompting large language models with answer heuristics for knowledge-based visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14974–14983, June 2023.

- [30] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [31] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36:38154–38180, 2023.
- [32] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- [33] Yuchen Shi, Siqi Cai, Zihan Xu, Yuei Qin, Gang Li, Hang Shao, Jiawei Chen, Deqing Yang, Ke Li, and Xing Sun. Flowagent: Achieving compliance and flexibility for workflow agents. *arXiv preprint arXiv:2502.14345*, 2025.
- [34] Liangcai Su, Zhen Zhang, Guangyu Li, Zhuo Chen, Chenxi Wang, Maojia Song, Xinyu Wang, Kuan Li, Jialong Wu, Xuanzhong Chen, et al. Scaling agents via continual pre-training. *arXiv preprint arXiv:2509.13310*, 2025.
- [35] Shuang Sun, Huatong Song, Yuhao Wang, Ruiyang Ren, Jinhao Jiang, Junjie Zhang, Fei Bai, Jia Deng, Wayne Xin Zhao, Zheng Liu, et al. Simpledeepsearcher: Deep information seeking via web-powered reasoning trajectory synthesis. *arXiv preprint arXiv:2505.16834*, 2025.
- [36] Tongyi DeepResearch Team, Baixuan Li, Bo Zhang, Dingchu Zhang, Fei Huang, Guangyu Li, Guoxin Chen, Hufeng Yin, Jialong Wu, Jingren Zhou, et al. Tongyi deepresearch technical report. *arXiv preprint arXiv:2510.24701*, 2025.
- [37] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [38] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [39] Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, et al. Webdancer: Towards autonomous information seeking agency. *arXiv preprint arXiv:2505.22648*, 2025.
- [40] Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, et al. Webwalker: Benchmarking llms in web traversal. *arXiv preprint arXiv:2501.07572*, 2025.
- [41] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- [42] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [43] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*, 2022.
- [44] Ziwei Zheng, Michael Yang, Jack Hong, Chenxiao Zhao, Guohai Xu, Le Yang, Chao Shen, and Xing Yu. Deepeyes: Incentivizing” thinking with images” via reinforcement learning. *arXiv preprint arXiv:2505.14362*, 2025.