

Yunjue Agent Tech Report: A Fully Reproducible, Zero-Start In-Situ Self-Evolving Agent System for Open-Ended Tasks

Haotian Li^{*†12} Shijun Yang^{*†31} Weizhen Qi¹ Silei Zhao¹
Rui Hua¹ Mingzhu Song¹ Xiaojian Yang¹ Chao Peng¹

Links: [\[GitHub Repository\]](#)

Abstract

Conventional agent systems often struggle in open-ended environments where task distributions continuously drift and external supervision is scarce. Their reliance on static toolsets or offline training lags behind these dynamics, leaving the system’s capability boundaries rigid and unknown. To address this, we propose the *In-Situ Self-Evolving* paradigm. This approach treats sequential task interactions as a continuous stream of experience, enabling the system to distill short-term execution feedback into long-term, reusable capabilities without access to ground-truth labels. Within this framework, we identify *tool evolution* as the critical pathway for capability expansion, which provides verifiable, binary feedback signals. Within this framework, we develop *Yunjue Agent*, a system that iteratively synthesizes, optimizes, and reuses tools to navigate emerging challenges. To optimize evolutionary efficiency, we further introduce a *Parallel Batch Evolution* strategy. Empirical evaluations across five diverse benchmarks under a zero-start setting demonstrate significant performance gains over proprietary baselines. Additionally, complementary warm-start evaluations confirm that the accumulated general knowledge can be seamlessly transferred to novel domains. Finally, we propose a novel metric to monitor evolution convergence, serving as a function analogous to training loss in conventional optimization. We open-source our codebase, system traces, and evolved tools to facilitate future research in resilient, self-evolving intelligence.

1. Introduction

The rapid ascendancy of Large Language Models (LLMs) has catalyzed a paradigm shift towards Artificial General Intelligence (AGI) (Guo et al., 2024). While this wave has birthed versatile agents, a critical dichotomy remains: the most capable systems are increasingly dominated by closed-source paradigms (Phan et al., 2025; Gupta et al., 2025). These proprietary models often rely on opaque tools and black-box APIs, creating a barrier for verifiable scientific inquiry. Conversely, open-source alternatives, while accessible, frequently lag in performance (Manchanda et al., 2024; Choi & Chang, 2025), as their reliance on static, manually crafted heuristics renders them unable to rival the robustness of industrial proprietary systems (Jin et al., 2024). We posit that for open-source intelligence to bridge this gap, it must move beyond static imitation of closed models and leverage a dynamic advantage: In-Situ Self-Evolution, where the system becomes stronger through continuous usage.

We envision that true AGI demands an In-Situ Self-Evolving architecture. In this ideal state, an agent acts as a fully adaptive organism, continuously refining its entire cognitive stack—*workflow*, *context* and *tool*. Here, the workflow adapts planning strategies based on user preferences or task types (Ye et al., 2025); the context dynamically distills episodic experiences into wisdom (Tang et al., 2025); and the tools evolve autonomously to conquer novel execution boundaries (Qin et al., 2024).

Recently, the field has witnessed a surge in self-evolving agents (Tao et al., 2024; Gao et al., 2025) designed to enhance these components autonomously (Xi et al., 2025; Cai et al., 2025; Acikgoz et al., 2025). Progress has been made in workflow adaptation (Zhang et al., 2025a) and context management (Tang et al., 2025). However, a critical limitation persists: existing paradigms predominantly rely on offline training with explicit reward signals or remain confined to narrow, domain-specific boundaries.

In supervision-scarce, open-ended environments—where

^{*}Equal contribution [†]During Internship at Yunjue Technology
¹Yunjue Technology ²Harbin Institute of Technology ³University of Science and Technology of China. Correspondence to: Weizhen Qi, Tech Lead. <qiweizhen@yunjuetech.com>.

ground truth is absent and user feedback is sparse—this reliance on external supervision becomes a bottleneck. Evolving a workflow or aligning with user preferences often involves subjective or delayed rewards (e.g., “Did the user like this summary?”). In contrast, tool evolution offers a distinct advantage: the feedback signal is intrinsic and objective (Wang et al., 2024a). Regardless of the final task outcome, the functionality of a tool is verifiable via rigorous execution feedback—code either runs successfully or throws an exception. This robust signal allows for autonomous optimization even in the absence of human intervention, making tools the most logical starting point for in-situ evolution. Beyond the signal quality, the toolset also constitutes the fundamental operational prerequisite for any agent (Zaremba et al., 2021). Building on this insight, we advocate for a staged, foundational approach that prioritizes the evolution of tools.

To this end, we introduce Yunjue Agent, a framework that prioritizes the autonomous synthesis, validation, and convergence of tools under a zero-start paradigm. Mechanically, the system employs a multi-agent architecture (Qian et al., 2024), comprising a Manager, Tool Developer, and Executor to address queries dynamically. Instead of relying on a static library, the agent synthesizes bespoke Python primitives on-the-fly when existing capabilities are insufficient. Crucially, to prevent tool explosion and ensure robustness, we introduce a Parallel Batch Evolution strategy. This mechanism processes queries in batches, clustering and merging functionally similar tools via an “absorbing” process. This allows the system to distill transient, query-specific code into a compact, convergent repository of generalized knowledge.

Our contributions are summarized as follows:

1. **SOTA Performance via zero-start In-Situ Self-Evolving:** We propose Yunjue Agent, a system that synthesizes, validates, and refines tools from scratch without prior training, with SOTA results on rigorous benchmarks.
2. **Cross-domain transferability:** We demonstrate that evolved tools exhibit robust transferability across disjoint domains, proving that the system distills transient interactions into generalized, convergent capabilities.
3. **Evolution convergence metric:** We introduce a quantitative metric to monitor evolutionary stability during inference, serving a function analogous to “training loss” in conventional optimization.
4. **Open reproducibility:** We open-source our codebase, complete interaction traces, and the library of evolved tools to provide a transparent foundation for future research.

2. In-situ self-evolving agents

An agent system can be formally defined as a tuple $\mathcal{M} = \langle \mathcal{W}, \mathcal{C}, \mathcal{T} \rangle$ of workflow, context and tools (Gao et al., 2025). In this system, the workflow \mathcal{W} is typically structured as a directed graph $\mathcal{G} = (V, E)$, where nodes V represent LLM-based agents and edges E denote the flow of information. \mathcal{C} is the set of contexts, where each context element can be a prompt template or a memory buffer encompassing dynamic conversation history. Finally, \mathcal{T} denotes the set of tools available for the agents to execute specific tasks.

To enable agents to adapt to new environments like humans, recent work focuses on achieving this through self-evolving methods (Tao et al., 2024; Gao et al., 2025). However, these methods typically require an iterative training process (He et al., 2025; Chen et al., 2025a) to update certain components or are confined in specific areas (Xia et al., 2025; Jin et al., 2025), such as generating query-specific workflows (Zhang et al., 2025a; Ye et al., 2025), refining context \mathcal{C} (e.g., prompt and memory) (Zhang et al., 2025a; Tang et al., 2025) or adjusting toolset (Wang et al., 2024b; 2023). In contrast, we aim to explore whether an agent system can adapt to the environment by continuously updating its components in-situ, where there is no access to external supervision signals. This motivation leads to a new paradigm of agentic evolution.

In-situ self-evolving of agents. Given a sequence of queries $= \{x_1, x_2, \dots, x_T\}$, the agent evolves dynamically from \mathcal{M}_0 to \mathcal{M}_T as it processes each. Specifically, after completing the t -th query x_t , the system updates its configuration for the next query, which can be formalized as a transition from $\mathcal{M}_{t-1} = \langle \mathcal{W}_{t-1}, \mathcal{C}_{t-1}, \mathcal{T}_{t-1} \rangle$ to $\mathcal{M}_t = \langle \mathcal{W}_t, \mathcal{C}_t, \mathcal{T}_t \rangle$, where any component of the tuple—the workflow structure, context, or toolset—may be modified based on internal feedback or experience gained from the previous interaction. Unlike self-evolving agents that maximize a specific objective function (e.g., accuracy) through a series of training steps (Gao et al., 2025), in-situ self evolving operates during the inference phase where ground truth information is unavailable.

To ensure a focused scope, we fix the workflow \mathcal{W} . Furthermore, given that Agent KB (Tang et al., 2025) has attempted to optimize memory through cross-domain experience to enhance agent performance, we also treat the context \mathcal{C} as a fixed component¹, reducing the evolving system state to $\mathcal{M}_t = \langle \mathcal{W}_0, \mathcal{C}_0, \mathcal{T}_t \rangle$. We posit that the variability of the toolset is the most decisive factor for a general-purpose agent system. While memory enhances performance by recalling experience, the ability to effectively solve problems

¹Although the execution flow can be dynamic during inference due to branch control and model inputs are populated at runtime, the underlying workflow structure and prompt templates are pre-defined. Thus, we consider these configurations fixed.

in novel domains is fundamentally limited by the availability of appropriate tools (Wang et al., 2023; Qin et al., 2024).

3. Methodology

3.1. In-situ self-evolving via tool accumulation

We posit that the fundamental prerequisite for a generalist agent lies in the continuous, dynamic evolution of its toolset. To address this, we introduce an agentic workflow for *in-situ self-evolving* via tool accumulation.

Upon receiving a query x_t , the agent first attempts to retrieve relevant utilities from its existing repository, denoted as $\mathcal{T}_{sub} \subseteq \mathcal{T}_{t-1}$. In the absence of requisite capabilities, the agent synthesizes novel tools \mathcal{P}_t tailored to the specific task constraints. Subsequently, empowered by the augmented toolset $\mathcal{T}_{sub} \cup \mathcal{P}_t$, the agent executes the task through a sequence of tool invocations. To ensure that generated tools are not merely functional but also robust and reusable, the agent engages in a self-reflection mechanism post-execution, refining the tools based on the error reports and execution traces. The process concludes with the integration of these refined tools into the global repository, updating the state to $\mathcal{T}_t = \mathcal{P}_t \cup \mathcal{T}_{t-1}$.

As the agent processes an expanding stream of queries, existing tools undergo iterative refinement while new ones are concurrently synthesized. This dual mechanism propels the system’s evolution along two dimensions: breadth (expanding functional coverage) and depth (optimizing tool robustness). Consequently, upon exposure to a sufficient volume of tasks, the tool repository is expected to reach a state of convergence, where incremental synthesis is necessitated only by outlier queries with idiosyncratic requirements.

3.2. Parallel batch evolution

While sequential query processing maximizes tool reuse and minimizes redundancy, it is inherently inefficient for large-scale benchmarks. Conversely, fully independent parallel processing fails to leverage shared knowledge accumulation across contexts. To reconcile these objectives, we propose a *Parallel Batch Evolution* strategy that parallelizes agent dynamics while maintaining a cohesive evolutionary trajectory.

Formally, let $\mathcal{Q}_t = \{q_{t,1}, q_{t,2}, \dots, q_{t,B}\}$ denote a batch of B user queries input to the agent \mathcal{M}_{t-1} at step t . For each query $q_{t,i} \in \mathcal{Q}_t$, the system synthesizes a set of local tools $\mathcal{P}_{t,i}$ that augment the global toolset \mathcal{T}_{t-1} specifically for that instance. However, independent generation often yields functionally redundant tools across different queries—particularly for general-purpose utilities like web searching. This redundancy expands the tool search space, increasing the cognitive load on the agent. To mitigate

this, we introduce a *tool absorbing mechanism* designed to cluster and consolidate utilities post-generation. Specifically, upon completion of batch \mathcal{Q} , the system aggregates all tools $\{\mathcal{T}_{t-1}, \mathcal{P}_{t,1}, \dots, \mathcal{P}_{t,B}\}$ and clusters them into disjoint groups $\{G_j\}$ based on functional semantic similarity. Subsequently, a merging function Φ consolidates these groups, filtering for quality and redundancy to produce a compact updated pool $\mathcal{T}_t = \Phi(\{G_j\})$. This process ensures the tool space remains streamlined, preventing retrieval ambiguity while updating the system state to \mathcal{M}_t .

Processing queries in batches offers distinct advantages. Primarily, it significantly enhances system throughput. Furthermore, analogous to mini-batch gradient descent in model optimization (Bottou, 2010; Goodfellow et al., 2016)—which mitigates gradient variance by averaging over samples—our absorbing mechanism reduces evolutionary stochasticity by merging similar tool instances. Simultaneously, this acts as a form of *Best-of-N* test-time scaling (Cobbe et al., 2021; Beirami et al., 2024), effectively performing multiple parallel rollouts for tool creation and selecting the optimal synthesis results for the permanent library.

3.3. The Yunjue agent system

By decoupling tool management, synthesis, and task execution, we establish a multi-agent system (Park et al., 2023) optimized for in-situ self-evolving. As illustrated in Figure 1, the architecture comprises distinct functional roles: **Manager**, **Executor**, **Tool Developer**, and **Integrator**, supported by the **Aggregator** and **Merger** for batch-level synchronization. The workflow proceeds via a collaborative mechanism designed to handle complex queries flexibly:

Upon receiving a user query q_t , the *Manager* orchestrates the workflow by first aligning tool capabilities with task requirements. It analyzes the system state to retrieve a relevant subset of tools $\mathcal{T}_{sub} \subseteq \mathcal{T}_{t-1}$ from the global repository. Should capability gaps arise, the *Manager* directs the *Tool Developer* to synthesize bespoke tools (implemented as Python primitives), which are immediately instantiated in the local context. Subsequently, the *Executor* addresses the query using the provisioned toolset, adhering to the ReAct paradigm (Yao et al., 2023). Crucially, the system supports dynamic runtime adaptation: if the *Executor* encounters unforeseen capability deficits during reasoning (e.g., identifying the need for a specific file parser), it suspends execution and signals the *Manager*. The *Manager* then provisions the requisite tools on-the-fly, allowing the *Executor* to resume seamlessly.

Upon task completion, the *Integrator* consolidates the execution history and intermediate outputs to formulate a comprehensive response a_t . Through this iterative process, the Yunjue agent evolves from state \mathcal{M}_{t-1} to \mathcal{M}_t , effectively updating the tool pool to \mathcal{T}_t for future evolution.

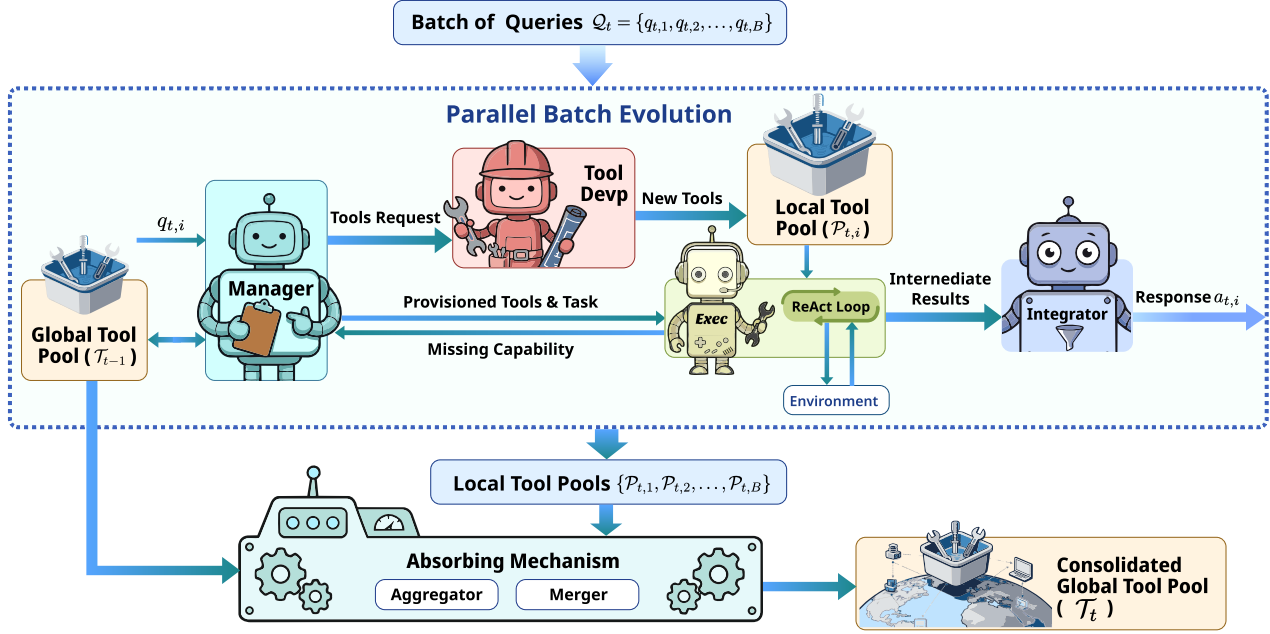


Figure 1. An architecture overview of Yunjue Agent.

Parallel evolution with batch processing. To optimize evolutionary throughput, we implement the *parallel batch evolution* strategy. Given a query batch $Q_t = \{q_{t,1}, q_{t,2}, \dots, q_{t,B}\}$, the system processes each instance concurrently, granting all agents shared access to the global tool repository T_{t-1} . This parallel execution allows each query $q_{t,i}$ to independently synthesize a local toolset $P_{t,i}$ tailored to its specific execution context.

Upon the completion of the batch, the *tool absorbing mechanism* is invoked to consolidate the dispersed local pools with the global state. Specifically, the system aggregates the union of all toolsets $\{T_{t-1}, P_{t,1}, P_{t,2}, \dots, P_{t,B}\}$. The LLM-based *Aggregator* then groups utilities based on functional semantic similarity. Subsequently, another LLM-based *Merger* is applied to each cluster to synthesize a unified, canonical tool that encapsulates collective capabilities while eliminating redundancy. The resulting consolidated repository T_t serves as the initialized state for the subsequent queries.

4. Experiment Setup

Datasets. To demonstrate the generalizability of our approach across diverse domains and task complexities, we conduct comprehensive evaluations on four complementary benchmarks, each targeting distinct professional scenarios: (i) *HLE* (Humanity’s Last Exam) (Phan et al., 2025), a frontier multi-modal benchmark featuring expert-level questions across mathematics, humanities, and natural sciences, designed to assess advanced reasoning at the boundary of hu-

man knowledge; (ii) *DeepSearchQA* (DSQA) (Gupta et al., 2025), which challenges agents’ ability to synthesize comprehensive answers through deep web search, iterative information gathering, and multi-source evidence aggregation; (iii) *xBench* (Chen et al., 2025b), a Chinese professional-aligned evaluation suite where we focus on the *ScienceQA* (xSciQA) subset (spanning natural, social, and language sciences) and *DeepSearch* (xDS, 2025.10 version) subset to evaluate real-world productivity in scientific research and complex retrieval workflows, allowing us to assess cross-lingual adaptation capabilities; and (iv) *FinSearchComp* (FSC) (Hu et al., 2025), a bilingual (English and Chinese) benchmark targeting financial analysis through its *T2* (Simple Historical Lookup) and *T3* (Complex Historical Investigation) tasks, which require precise time-sensitive data retrieval and multi-step quantitative reasoning over financial documents.

Baselines. We benchmark our agent system against a comprehensive suite of proprietary and open-source systems, spanning both static and self-evolving agent paradigms. Unless otherwise noted, reported performance metrics are derived directly from original publications, technical reports, or authoritative disclosures (detailed provenance is provided in the Appendix).

5. Evaluation on open-ended evolution

5.1. Zero-start performance on cross-domain tasks

We evaluate our Yunjue Agent, initialized with an empty toolset, across all five benchmarks against state-of-the-art

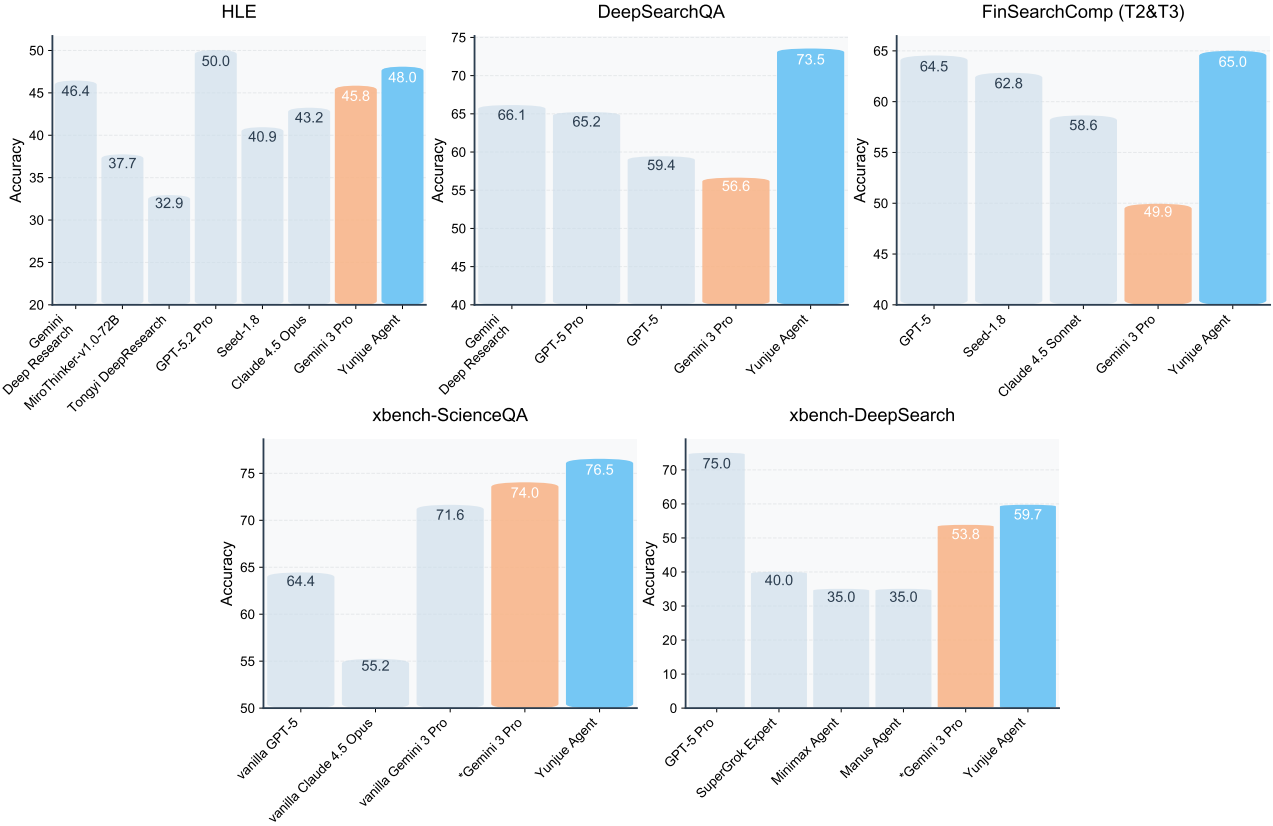


Figure 2. Performance comparison of Yunjue Agent against state-of-the-art agents and agentic foundation models. Our method is highlighted in cyan, and the backend model (Gemini 3 Pro) appears in orange. *Gemini 3 Pro denotes our implementation with a Python interpreter.

baselines. Results are presented in Figure 2. For datasets excluding xSciQA, foundation model baselines (e.g., GPT-5, Gemini 3 Pro) are augmented with web search and Python interpreters. Conversely, on xSciQA, where standard leaderboards prohibit external tools, we explicitly equip our backend model with a Python interpreter to ensure a fair parity comparison.

On the heterogeneous HLE benchmark, our agent yields a significant improvement over the backend (48.0 vs. 45.8), ranking second only to GPT-5.2 Pro. Most notably, Yunjue Agent achieves state-of-the-art performance on DSQA, FSC, and xSciQA. We observe substantial absolute gains over the Gemini 3 Pro baseline, with increases of **+17.4** points on DSQA (73.5 vs. 56.6) and **+15.1** points on FSC (65.0 vs. 49.9), alongside a record-setting score of 76.5 on xSciQA. Finally, on xDS, our method maintains competitive performance (59.7), surpassed only by GPT-5 Pro while outperforming all other models by a significant margin. These exceptional results demonstrate our method’s **superior performance across a wide spectrum of tasks, laying a solid foundation for generalized knowledge accumulation and transfer.**

Complementing the performance analysis, we aggregated

toolsets across all five benchmarks to examine functional utilization patterns. Figure 3 presents the invocation frequency of the top 50 tools. The distribution reveals the **spontaneous emergence** of high-utility fundamental functions, most notably `web_search`, `fetch_web_content` and `evaluate_math_expression`. The dominance of these foundational tools confirms that the system has effectively **distilled generalized knowledge into versatile primitives**, ensuring broad applicability across diverse task semantics.

5.2. Warm start evolution in shifting domains

We further validate that the generalized knowledge accumulated by our framework is transferable to novel domains, capable of catalyzing capability evolution upon a pre-existing foundation. Specifically, we conduct sequential experiments on DSQA and xSciQA, initializing the system with the toolset derived from the HLE benchmark. HLE is selected as the foundational domain due to its substantial scale (2,500 queries) and extensive disciplinary coverage.

Table 1 presents a comparison of performance and tool synthesis across initialization strategies. We observe that

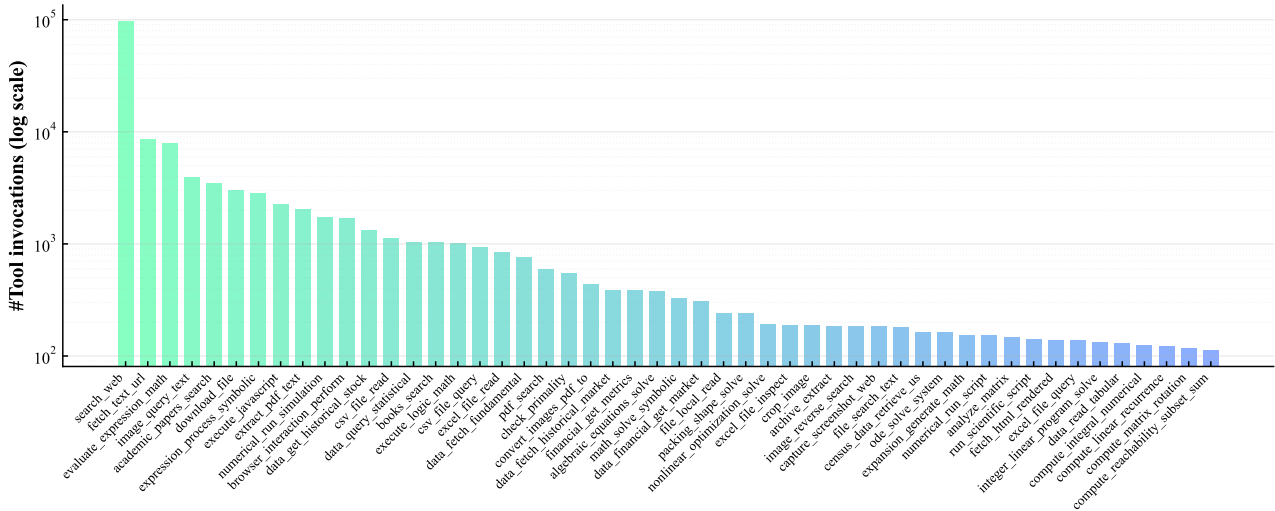


Figure 3. Frequency distribution of the toolset evolved across five benchmarks. We report the top 50 tools, illustrating the emergence of high-generalizability primitives.

Table 1. Performance comparison and tool synthesis statistics across different initialization strategies. ZS stands for a zero start setting, and WS denotes warm start, which leverages the HLE-evolved toolset as initialization.

Dataset	Strategy	Accuracy	# New tools
HLE	ZS	48.0	97
DeepSearchQA	ZS	74	34
	WS	74.6 _{↑0.6}	23 _{↓32%}
FinSearchComp	ZS	65.0	18
	WS	65.4 _{↑0.4}	8 _{↓55%}
xbench-ScienceQA	ZS	76.5	13
	WS	80.2 _{↑3.7}	0 _{↓100%}
xbench-DeepSearch	ZS	59.7	16
	WS	60.6 _{↑0.9}	0 _{↓100%}

performance remains consistent across domains, with a notable improvement on xSciQA, where the score increases from 76.5 to 80.2. Crucially, the number of newly synthesized tools decreases effectively by 32% on DSQA and remarkably by 100% on xSciQA. This substantial reduction demonstrates that our system’s accumulated generalized knowledge can be seamlessly transferred to disparate domains while retaining the capacity for continuous learning of novel capabilities.

Figure 4 illustrates the evolution of the tool library size relative to the cumulative number of processed queries. The overall trajectory demonstrates **a global trend toward convergence**. A localized surge in tool synthesis during the late phase of HLE (queries 2,000–2,400) is driven by intra-benchmark domain shifts: as query semantics transition from predominantly mathematical to social sciences and other disciplines, a targeted expansion of the toolset is necessitated. Notably, despite this semantic heterogeneity,

only 97 tools were generated across the entire 2,500-HLE query corpus, evidencing the effective consolidation of acquired knowledge. Subsequently, the tool growth curves for DSQA and xSciQA exhibit a near-zero gradient. This plateau indicates that the capabilities evolved during the HLE phase possess sufficient generalization to support robust adaptation to novel domains with negligible marginal tool synthesis.

To empirically substantiate the role of tools as vehicles for knowledge crystallization and domain adaptation, we analyze the toolset intersection visualized in Figure 5. The comparative analysis reveals a substantial overlap between tools generated in the *zero-start* regime and those synthesized during warm-start evolution. Specifically, the recurrence of exact matches, such as `read_excel`, validates the system’s capacity to **deterministically recover essential utilities**. Furthermore, the emergence of semantically aligned variants (e.g., `inspect_ods`, `fetch_api`) highlights the agent’s adaptability in synthesizing functionally equivalent primitives tailored to specific task constraints. The predominance of these convergent tools corroborates our hypothesis of systemic robustness in shifting domains, confirming that the agent reliably gravitates toward an **invariant core of capabilities to fulfill domain requirements, independent of the initialization state**.

5.3. Evolutionary generality loss

To quantify the dynamics of knowledge acquisition, we propose a novel metric termed *Evolutionary generality loss* (EGL). Functioning analogously to the objective function in stochastic optimization, EGL serves as a real-time indicator

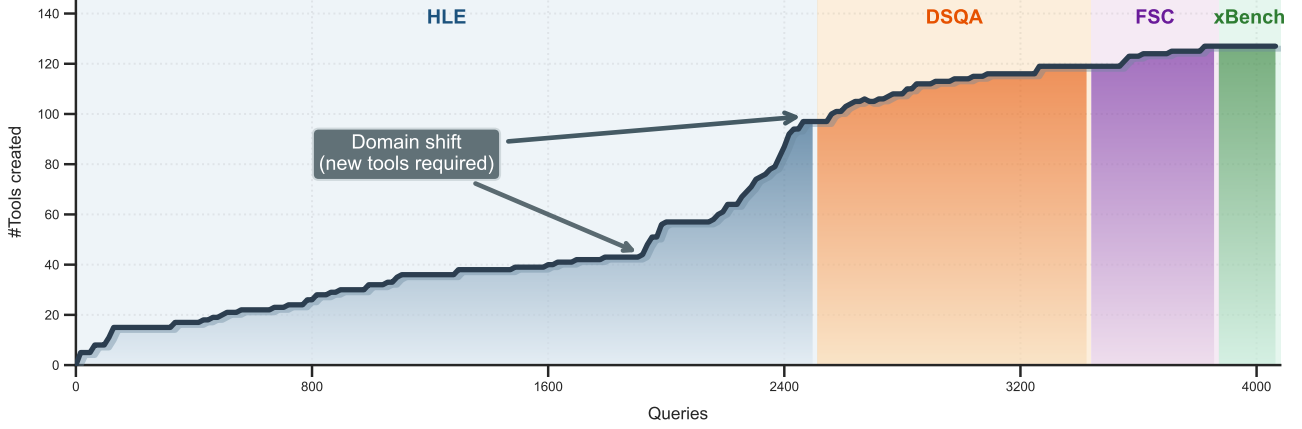


Figure 4. Evolution of the tool library size relative to the cumulative number of processed queries. The experimental sequence follows the curriculum HLE → DeepSearchQA → FinSearchComp → xbench-ScienceQA → xbench-DeepSearch, highlighting the convergence of tool synthesis.

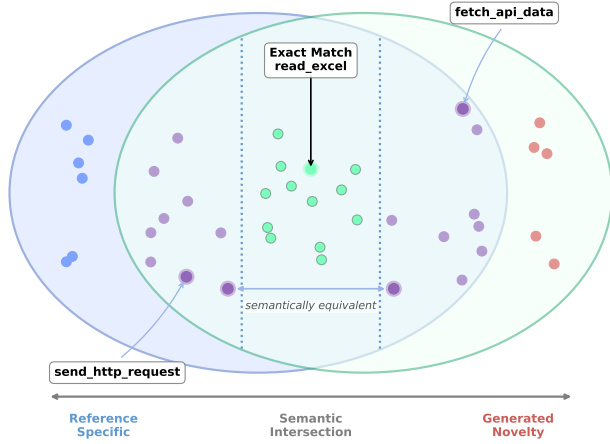


Figure 5. Venn diagram visualization of tool correspondence between zero-start and warm-start settings on DeepSearchQA. The left set comprises tools unique to the zero-start baseline ($\mathcal{T}_{\text{DSQA}} \setminus \mathcal{T}_{\text{HLE}}$), while the right set consists of incrementally generated tools in the warm-start setting ($\mathcal{T}_{\text{HLE} \rightarrow \text{DSQA}} \setminus \mathcal{T}_{\text{HLE}}$). Central intersection indicates high functional overlap. Distinct points represent individual tools, arranged by semantic similarity. Tools in the intersection share similar functionalities, while those within the central dashed lines are exact matches.

of evolutionary convergence. Formally, EGL is defined as:

$$\text{EGL} = \frac{\sum_{i=1}^N c_i}{\sum_{i=1}^N u_i} \times 1000 \quad (1)$$

where c_i and u_i denote the number of newly synthesized tools and tool invocations for the i -th query, respectively, and N represents the total number of processed queries. Intuitively, the EGL value is elevated during the nascent stages of evolution, driven by the necessity of *ab initio* tool creation. As the system matures, the dominance of tool reuse over synthesis causes the metric to decay. Con-

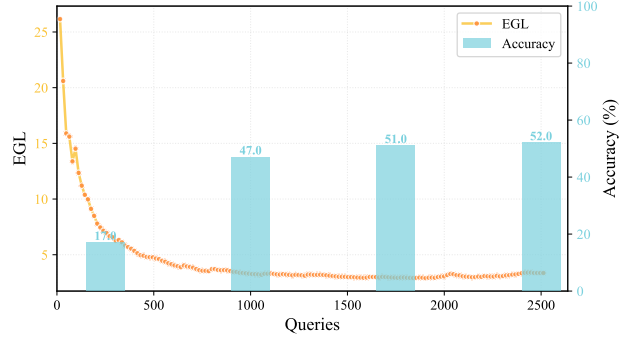


Figure 6. Evolutionary generality loss (EGL) on the HLE benchmark and accuracy on selected 200-queries-dataset as a function of the number of queries. The orange curve represents the EGL trend (left axis). The blue bars show the agent’s accuracy (right axis) at different stages of the evolution (170, 1000, 1750, 2500 queries), highlighting that performance improves and stabilizes as the toolset converges.

versely, a sustained high EGL value implies that the system is failing to generalize, necessitating continuous creation. Conceptually, this metric is inversely correlated with tool generality: as the versatility of the repository increases, the EGL value decreases. To validate its efficacy, we monitored the EGL trajectory on the HLE benchmark as a function of query volume. The resulting trend in Figure 6 exhibits a distinct convergence pattern, achieving **asymptotic stability** after processing approximately 1,000 queries (40% of the stream).

To empirically corroborate this convergence, we extracted system snapshots at 10%, 40%, 70%, and 100% of the evolutionary timeline and evaluated performance on a stratified random sample of 200 queries (100 each from the HLE and DSQA datasets), with further tool synthesis disabled. As illustrated in the bar chart of Figure 6, a substantial

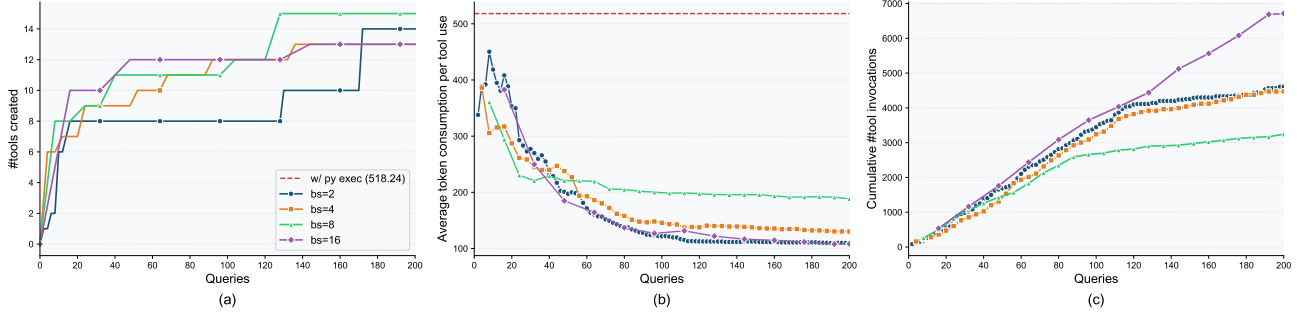


Figure 7. Evolution dynamics under different batch sizes. (a) The number of created tools with the number of queries. (b) Average token consumption per tool use. (c) Cumulative number of tool invocations.

Table 2. Experimental results on a subset of 200 selected queries under varying batch sizes. “Tool succ. rt.” denotes the tool execution success rate, while “Avg #tokens for tools” quantifies the average token consumption per tool invocation. The condition “ ∞ ” represents that the backend model operates exclusively with a Python interpreter, processing queries independently without tool accumulation.

Batch size	Accuracy	Tool succ. rt.	Avg #tokens for tools
2	48.0	99.6%	109.9
4	50.0	99.9%	130.3
8	47.5	99.6%	188.4
16	51.5	99.1%	108.4
∞	40.0	81.8%	518.2

performance inflection occurs between the 10% and 40% checkpoints, signifying a rapid expansion of the system’s capability boundary. Conversely, marginal gains diminish significantly in the 40%–100% interval. This saturation provides strong evidence that the evolutionary process successfully converges, reaching a state of high generality where incremental tool synthesis yields diminishing returns.

5.4. Ablation study

To elucidate the factors underpinning the efficacy of tool evolution and to ablate key system parameters, we conduct comparative experiments using the identical subset of 200 queries employed in the EGL evaluation.

Why does tool evolution work? To validate the fundamental utility of tool evolution, we establish a baseline comparison against an agent sharing an identical backend (GPT-5) but equipped solely with a Python interpreter. While the interpreter allows for ad-hoc code execution and iterative debugging based on error feedback, it fundamentally lacks the capacity for knowledge crystallization and reuse.

We evaluate our agent under a *zero-start* regime with varying batch sizes, contrasting it against the Python interpreter baseline (denoted as $B \rightarrow \infty$). The condition “ ∞ ” represents the theoretical asymptotic limit of infinite batch size, where the backend model operates exclusively with

Table 3. Performance comparison and number of created tools of different backend models on DeepSearchQA and FinSearchComp benchmarks.

Backend	DeepSearchQA		FinSearchComp	
	Accuracy	# Tools	Accuracy	# Tools
Gemini 3 Pro	74.0	16	65.0	18
GPT-5	70.6	75	68.2	25
GPT-5-mini	42.5	47	45.5	13

a Python interpreter. In this setting, queries are processed independently, precluding any tool accumulation, clustering, or shared knowledge transfer. As detailed in Table 2, the Python-only baseline yields a markedly inferior score of 40.0, compared to the 47.5–51.5 range achieved by our method. To elucidate the underlying causes, we analyzed the tool success rate (successful invocations per total attempts) and token efficiency (average tokens per invocation). The baseline demonstrates deficiencies in both metrics, exhibiting a success rate of only 81.8% while consuming an average of 518.2 tokens per invocation. In stark contrast, our system consistently maintains success rates exceeding 99% with significantly reduced token overhead. This divergence underscores that evolved tools are not only more robust but also more efficient in information processing. We posit that the high execution error rate and the verbosity of raw code generation induce severe **contextual contamination, thereby degrading the reasoning capabilities of the Python baseline.**

How does batch size impact parallel evolution? We further investigate the role of batch size B in mediating the trade-off between evolutionary velocity and convergence stability. As illustrated in Figure 7(a), larger batch configurations induce a steeper initial gradient in tool synthesis. This suggests that high-concurrency processing facilitates the rapid identification of common patterns and the accelerated crystallization of foundational primitives. Crucially, despite varied initial velocities, the terminal library sizes remain comparable across all settings, confirming that the system **invariably converges to a consistent equilibrium of capabilities** regardless of batch magnitude.

Figure 7(b) tracks the temporal evolution of computational overhead. Even in the nascent stages, our parallelized architecture yields reduced token consumption relative to the baseline (w/ py exec). As the system matures, the average token cost per invocation demonstrates a precipitous decay, eventually stabilizing at an asymptotic minimum. Note that localized fluctuations in this downward trend correspond to the sporadic synthesis of novel tools required for outlier queries. This overall convergence signifies the system’s successful phase transition from expensive *ab initio* code generation to efficient, low-latency tool reuse.

We observe in Figure 7(b) that the average token consumption per tool use after convergence for a batch size of 8 is slightly higher than other settings. This is because, as shown in Figure 7(a), it generates the largest number of tools, and Figure 7(c) indicates it has the highest cumulative tool invocations. We attribute this to the stochastic nature of LLM execution, which leads to varying agent behaviors under different settings. Future work will investigate strategies to align the convergence states of the system across different batch sizes.

How does backend model selection affect evolution? We evaluate the influence of the underlying foundation model on system performance across the HLE and DSQA datasets, as detailed in Table 3. Results indicate that when instantiated with *GPT-5*, our agent surpasses the majority of established baselines. Notably, the lightweight *GPT-5-mini* variant achieves competitive performance, demonstrating significant utility despite its reduced parameter count. Crucially, we employ a unified prompt strategy across all configurations without architecture-specific tuning. These findings underscore the model agnosticism and generalization capability of our framework, confirming its robustness across diverse backend scales.

6. Related Work

Self-evolving agents. To transcend the limitations of static designs, general-purpose agent systems must adapt to open-ended, interactive environments in real-time (Gao et al., 2025). Self-evolving agents achieve this by autonomously refining internal components, categorized into workflow optimization, model parameter updates, context management, and tool synthesis. Specifically, MAS-GPT (Ye et al., 2025) and AFlow (Zhang et al., 2025a) generate query-specific workflows, while TT-SI (Acikgoz et al., 2025) and SCA (Zhou et al., 2025) update parameters by generating challenging training problems. Furthermore, ELL (Cai et al., 2025) and Zhang et al. (2025b) evolve by accumulating environmental interaction experience as context. Regarding tool synthesis, LIVE-SWE-agent (Xia et al., 2025) and STELLA (Jin et al., 2025) enable on-the-fly tool creation for software engineering and biomedical research, respec-

tively. Nevertheless, most self-evolving methods necessitate explicit training. Crucially, even approaches closest to our work—those autonomously creating tools at test-time—remain confined to specific domains and lack tool reuse mechanisms, hindering their general applicability.

Tool evolution agents. The paradigm of autonomous tool synthesis has garnered significant attention in recent literature. Pioneering systems such as Voyager (Wang et al., 2023), STELLA (Jin et al., 2025), and LIVE-SWE-agent (Xia et al., 2025) empower agents to generate executable tools for embodied control or software engineering tasks. Notably, analyses of the tool embedding space in LIVE-SWE-agent (Xia et al., 2025) reveal distinct clustering patterns among functional equivalents; this observation directly informed the design of our *absorbing mechanism*, which utilizes clustering to consolidate redundant utilities. Regarding tool reuse, Alita (Qiu et al., 2025) facilitates the persistence of successful tools, though its optimization trajectory is predicated on ground-truth supervision. In the context of inference-time adaptation, Lu et al. (2026) proposed a test-time tool evolution framework governed by a cost-sensitive objective function, yet its application remains narrowly scoped to scientific reasoning. Collectively, these approaches prioritize the optimization of existing assets or domain-specific synthesis, distinguishing them from our framework’s focus on *self-evolution in open-ended environments*.

General agent systems. Recent research has sought to construct general-purpose agents capable of cross-domain operation. DeepAgent (Li et al., 2025) utilizes a scalable toolset to address diverse tasks but lacks the capacity for autonomous environmental adaptation, leading to performance limitations. Agent KB (Tang et al., 2025) attempts to solve cross-domain problems by summarizing execution experiences from other agents; however, its performance is constrained by a deficiency in necessary tools. In contrast, our approach achieves generality by continuously creating, reusing, and optimizing tools in-situ, effectively overcoming the limitations of static toolsets and experience-based adaptation.

7. Discussion and future work

While this study establishes the efficacy of in-situ self-evolution through tool synthesis, several avenues for future research merit rigorous inquiry to fully realize the potential of autonomous agents.

Paradigm parallels: towards system-level pre-training for agentic systems. The seamless transition and performance consistency observed between our *zero-start* and *warm-start* settings offer more than just empirical validation; they suggest a fundamental paradigm shift. The clear

convergence curves of the tool library indicate that “task-solving capability” is not merely a collection of ad-hoc heuristics, but a generalizable pattern that can be learned and distilled.

This implies that the field is approaching a “pre-training and post-training” era for agentic systems, mirroring the trajectory of LLMs. We envision a future where multi-agent systems undergo *system-level pre-training* on massive, broad-spectrum task datasets. This process would allow agents to distill a “foundation toolset”—a set of converged, generalizable primitives—before deployment. Consequently, such pre-trained agents would possess intrinsic generalization capabilities, enabling them to tackle novel downstream tasks largely through the composition of existing reliable tools, thereby minimizing or even eliminating the need for expensive test-time evolution. Formalizing the methodologies for this “agentic pre-training” is a critical priority for future research.

Co-evolution of memory and workflow. Currently, our framework validates the self-evolving paradigm primarily through tool generation. However, tools constitute a necessary but insufficient condition for scenarios demanding high personalization or complex process management. For instance, personalized assistants require persistent, evolving memory structures to align with user preferences, while intricate tasks (e.g., deep research) necessitate structured workflow evolution—such as generating bespoke planning protocols for multi-agent coordination. A critical direction is extending the evolutionary mechanism to encompass the co-evolution of memory architectures and workflow policies, enabling the system to adapt its internal state and execution logic alongside its functional capabilities.

Evolutionary stability and regularization. The stability of tool evolution warrants further investigation. Due to the inherent stochasticity of LLM generation, toolsets can exhibit variance across experimental runs, as evidenced by our batch size ablation. Ensuring the consistent convergence of the tool library is vital for system reliability. Future work will focus on developing regularization strategies to guarantee the determinism of the evolutionary process in open-ended environments.

Optimization of parallel batch evolution. Finally, the nuances of the batch evolution strategy present rich opportunities for optimization:

- *Curriculum learning effects:* As shown in Figure 4, the sequence of incoming queries significantly influences the convergence trajectory; delaying queries that require foundational primitives can impede system maturation. Investigating optimal query ordering is a key direction.
- *Intra-batch diversity trade-offs:* Intra-batch diversity presents a complex dynamic. Low diversity allows the

absorbing mechanism to function as a form of *Best-of-N test-time scaling*, leveraging redundancy to select the optimal implementation—a contrast to traditional gradient training where high batch diversity is preferred. Balancing this quality-assurance benefit against the efficiency of capability evolution is crucial.

- *Adaptive scheduling:* Dynamic batch sizing offers a promising avenue for optimization. Smaller batch sizes in early stages could foster the robust consolidation of general knowledge via redundancy, while larger batches in the post-convergence phase could maximize throughput for corner cases. Developing autonomous agents capable of dynamic batch scheduling based on convergence signals remains a worthy research direction.

8. Conclusion

In this work, we presented the In-Situ Self-Evolving framework, enabling LLM-based agents to autonomously adapt and evolve within open-ended environments. Through the Yunjue Agent and the proposed Parallel Batch Evolution strategy, we demonstrated that treating tools as dynamic vehicles for knowledge crystallization allows for robust zero-start learning and the emergence of generalized, transferable capabilities. Our empirical results validate that this approach not only achieves state-of-the-art performance across heterogeneous benchmarks but also significantly outperforms static baselines in both efficiency and adaptability. Furthermore, warm-start experiments confirm that the accumulated generalized knowledge can be seamlessly transferred to novel domains, enabling the agent to adapt to shifting environments with minimal marginal tool synthesis. We release our codebase, evaluation and system traces to the open-source community to facilitate further research. Broadly, future work will aim to unify tool evolution with memory and workflow adaptation, advancing toward more autonomous and general-purpose agentic systems.

References

- Acikgoz, E. C., Qian, C., Ji, H., Hakkani-Tür, D., and Tur, G. Self-improving llm agents at test-time. *arXiv preprint arXiv:2510.07841*, 2025.
- Anthropic. Claude 4.5 opus system card. <https://www-cdn.anthropic.com/bf10f64990cfda0ba858290be7b8cc6317685f47.pdf>, 2025. Accessed: 2026-01-23.
- Beirami, A., Agarwal, A., Berant, J., D’Amour, A., Eisenstein, J., Nagpal, C., and Suresh, A. T. Theoretical guarantees on the best-of-n alignment policy. *arXiv preprint arXiv:2401.01879*, 2024.

- Bottou, L. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pp. 177–186. Springer, 2010.
- ByteDance. Seed-1.8 model card. <https://github.com/ByteDance-Seed/Seed-1.8/blob/main/Seed-1.8-Modelcard.pdf>, 2025. Accessed: 2026-01-23.
- Cai, Y., Hao, Y., Zhou, J., Yan, H., Lei, Z., Zhen, R., Han, Z., Yang, Y., Li, J., Pan, Q., et al. Building self-evolving agents via experience-driven lifelong learning: A framework and benchmark. *arXiv preprint arXiv:2508.19005*, 2025.
- Chen, J., Yang, Z., Shi, J., Wo, T., and Tang, J. Mathse: Improving multimodal mathematical reasoning via self-evolving iterative reflection and reward-guided fine-tuning. *arXiv preprint arXiv:2511.06805*, 2025a.
- Chen, K., Ren, Y., Liu, Y., Hu, X., Tian, H., Xie, T., Liu, F., Zhang, H., Liu, H., Gong, Y., Sun, C., Hou, H., Yang, H., Pan, J., Lou, J., Mao, J., Liu, J., Li, J., Liu, K., Liu, K., Wang, R., Li, R., Niu, T., Zhang, W., Yan, W., et al. xbench: Tracking agents productivity scaling with profession-aligned real-world evaluations. *arXiv preprint arXiv:2506.13651*, 2025b. Available at <https://xbench.org/>.
- Choi, W.-C. and Chang, C.-I. Advantages and limitations of open-source versus commercial large language models (llms): A comparative study of deepseek and openai's chatgpt. *Preprints*, 2025.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Gao, H.-a., Geng, J., Hua, W., Hu, M., Juan, X., Liu, H., Liu, S., Qiu, J., Qi, X., Wu, Y., et al. A survey of self-evolving agents: On path to artificial super intelligence. *arXiv preprint arXiv:2507.21046*, 2025.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep learning*. MIT press, 2016.
- Google. Deep research agent & gemini api. <https://blog.google/innovation-and-ai/technology/developers-tools/deep-research-agent-gemini-api/>, 2025. Accessed: 2026-01-23.
- Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N. V., Wiest, O., and Zhang, X. Large language model based multi-agents: a survey of progress and challenges. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI '24*, 2024. ISBN 978-1-956792-04-1. doi: 10.24963/ijcai.2024/890. URL <https://doi.org/10.24963/ijcai.2024/890>.
- Gupta, N., Chatterjee, R., Haas, L., Tao, C., Wang, A., Liu, C., Oiwa, H., Gribovskaya, E., Ackermann, J., Blitzer, J., Goldshtein, S., and Das, D. Deepsearchqa: Bridging the comprehensiveness gap for deep research agents. https://storage.googleapis.com/deepmind-media/DeepSearchQA/DeepSearchQA_benchmark_paper.pdf, 2025. Google DeepMind, Google Search, Kaggle, and Google Research.
- He, Y., Huang, C., Li, Z., Huang, J., and Yang, Y. Visplay: Self-evolving vision-language models from images. *arXiv preprint arXiv:2511.15661*, 2025.
- Hu, L., Jiao, J., Liu, J., Ren, Y., Wen, Z., Zhang, K., Zhang, X., Gao, X., He, T., Wu, Y., et al. Finsearchcomp: Towards a realistic, expert-level evaluation of financial search and reasoning. *arXiv preprint arXiv:2509.13160*, 2025.
- Jin, H., Huang, L., Cai, H., Yan, J., Li, B., and Chen, H. From llms to llm-based agents for software engineering: A survey of current, challenges and future. *arXiv preprint arXiv:2408.02479*, 2024.
- Jin, R., Zhang, Z., Wang, M., and Cong, L. Stella: Self-evolving llm agent for biomedical research. *arXiv preprint arXiv:2507.02004*, 2025.
- LangChain. Langgraph, 2024. URL <https://github.com/langchain-ai/langgraph>. Accessed: 2026-01-25.
- Li, X., Jiao, W., Jin, J., Dong, G., Jin, J., Wang, Y., Wang, H., Zhu, Y., Wen, J.-R., Lu, Y., et al. Deepagent: A general reasoning agent with scalable toolsets. *arXiv preprint arXiv:2510.21618*, 2025.
- Lu, J., Kong, Z., Wang, Y., Fu, R., Wan, H., Yang, C., Lou, W., Sun, H., Wang, L., Jiang, Y., Wang, X., Sun, X., and Zhou, D. Beyond static tools: Test-time tool evolution for scientific reasoning, 2026. URL <https://arxiv.org/abs/2601.07641>.
- Manchanda, J., Boettcher, L., Westphalen, M., Davis, K., and Varanasi, R. The open source advantage in large language models (llms). *arXiv preprint arXiv:2412.12004*, 2024.
- OpenAI. Introducing gpt-5.2. <https://openai.com/index/introducing-gpt-5-2/>, 2025. Accessed: 2026-01-23.

- Park, J. S., O'Brien, J., Cai, C. J., Morris, M. R., Liang, P., and Bernstein, M. S. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701320. doi: 10.1145/3586183.3606763. URL <https://doi.org/10.1145/3586183.3606763>.
- Phan, L., Gatti, A., Han, Z., Li, N., Hu, J., Zhang, H., Zhang, C. B. C., Shaaban, M., Ling, J., McCauley, S., et al. Humanity's last exam. *arXiv preprint arXiv:2501.14249*, 2025. Available at <https://lastexam.ai/>.
- Qian, C., Xie, Z., Wang, Y., Liu, W., Zhu, K., Xia, H., Dang, Y., Du, Z., Chen, W., Yang, C., et al. Scaling large language model-based multi-agent collaboration. *arXiv preprint arXiv:2406.07155*, 2024.
- Qin, Y., Hu, S., Lin, Y., Chen, W., Ding, N., Cui, G., Zeng, Z., Huang, Y., Xiao, C., Han, C., et al. Tool learning with foundation models. In *Advances in Neural Information Processing Systems*, volume 37, 2024.
- Qiu, J., Qi, X., Zhang, T., Juan, X., Guo, J., Lu, Y., Wang, Y., Yao, Z., Ren, Q., Jiang, X., Zhou, X., Liu, D., Yang, L., Wu, Y., Huang, K., Liu, S., Wang, H., and Wang, M. Alita: Generalist agent enabling scalable agentic reasoning with minimal predefinition and maximal self-evolution, 2025. URL <https://arxiv.org/abs/2505.20286>.
- Tang, X., Qin, T., Peng, T., Zhou, Z., Shao, D., Du, T., Wei, X., Zhu, H., Zhang, G., Liu, J., Wang, X., Hong, S., Wu, C., and Zhou, W. AGENT KB: A hierarchical memory framework for cross-domain agentic problem solving. In *ICML 2025 Workshop on Collaborative and Federated Agentic Workflows*, 2025. URL <https://openreview.net/forum?id=ohXoWHlrn8>.
- Tao, Z., Lin, T.-E., Chen, X., Li, H., Wu, Y., Li, Y., Jin, Z., Li, F., Li, D., and Zhou, K. A survey on self-evolution of large language models. *arXiv preprint arXiv:2404.14387*, 2024.
- Team, M., Bai, S., Bing, L., Chen, C., Chen, G., Chen, Y., Chen, Z., Chen, Z., Dai, J., Dong, X., et al. Mirothinker: Pushing the performance boundaries of open-source research agents via model, context, and interactive scaling. *arXiv preprint arXiv:2511.11793*, 2025.
- The Pallets Projects. Jinja2, 2007. URL <https://jinja.palletsprojects.com/>. Accessed: 2026-01-25.
- Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- Wang, X., Chen, Y., Yuan, L., Zhang, Y., Li, Y., Peng, H., and Ji, H. Executable code actions elicit better llm agents. In *Proceedings of the 41st International Conference on Machine Learning*, 2024a.
- Wang, Z. Z., Mao, J., Fried, D., and Neubig, G. Agent workflow memory. *arXiv preprint arXiv:2409.07429*, 2024b.
- xBench Team. xbench: Deepsearch leaderboard. <https://xbench.org/agi/aisherech>, 2025a. Accessed: 2026-01-23.
- xBench Team. xbench: Scienceqa leaderboard. <https://xbench.org/agi/scienceqa>, 2025b. Accessed: 2026-01-23.
- Xi, Z., Ding, Y., Chen, W., Hong, B., Guo, H., Wang, J., Yang, D., Gui, T., Zhang, Q., Huang, X., et al. Agent-gym: Evaluating and training large language model-based agents across diverse environments. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2025.
- Xia, C. S., Wang, Z., Yang, Y., Wei, Y., and Zhang, L. Live-swe-agent: Can software engineering agents self-evolve on the fly? *arXiv preprint arXiv:2511.13646*, 2025.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.
- Ye, R., Tang, S., Ge, R., Du, Y., Yin, Z., Chen, S., and Shao, J. MAS-GPT: Training LLMs to build LLM-based multi-agent systems. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=3CiSpY3QdZ>.
- Zaremba, W., Brockman, G., and OpenAI. OpenAI Codex. <https://openai.com/index/openai-codex/>, August 2021. Accessed: 2025-01-19.
- Zhang, J., Xiang, J., Yu, Z., Teng, F., Chen, X.-H., Chen, J., Zhuge, M., Cheng, X., Hong, S., Wang, J., Zheng, B., Liu, B., Luo, Y., and Wu, C. AFlow: Automating agentic workflow generation. In *The Thirteenth International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=z5uVAKwmjff>.
- Zhang, K., Chen, X., Liu, B., Xue, T., Liao, Z., Liu, Z., Wang, X., Ning, Y., Chen, Z., Fu, X., et al. Agent learning via early experience. *arXiv preprint arXiv:2510.08558*, 2025b.

Zhou, Y., Levine, S., Weston, J. E., Li, X., and Sukhbaatar, S. Self-challenging language model agents. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=9yusqX9DpR>.

A. Experiment setup

Table 4. Statistics of the datasets.

Dataset	Domain	Language	# Queries
HLE	General reasoning	EN	2,500
DeepSearchQA	General QA	EN	900
FinSearchComp (T1 & T2)	Financial search	EN & CH	391
xBench-ScienceQA	General science	CH	100
xBench-DeepSearch	General research	CH	100

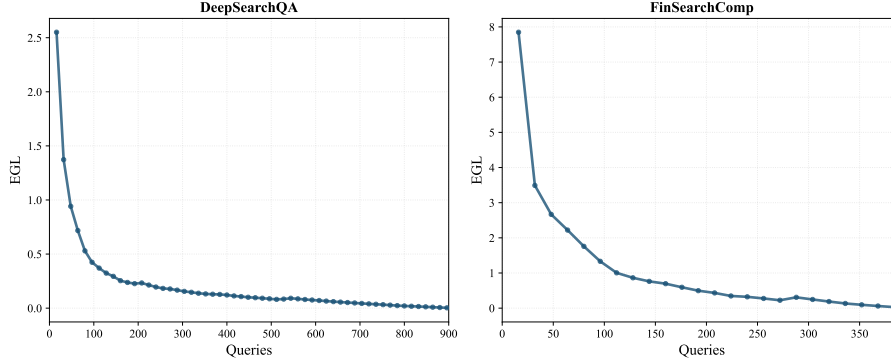


Figure 8. EGL dynamics on the DSQA and FSC datasets.

Data statistics are summarized in Table 4.

Evaluation protocols. We follow the standard evaluation protocols released with each dataset:

- HLE: [Official evaluation script](#)
- DeepSearchQA: [Official evaluation script](#)
- FinSearchComp: [Judging prompts](#)
- xBench: [Official evaluation script](#)

Baselines sources. We collect baseline results from their official reports or leaderboards. For HLE, results for Gemini Deep Research and Gemini 3 Pro are from (Google, 2025); MiroThinker-v1.0-72B and Tongyi DeepResearch are from (Team et al., 2025); GPT-5.2 Pro is from (OpenAI, 2025); Seed-1.8 is from (ByteDance, 2025); and Claude 4.5 Opus is from (Anthropic, 2025). DeepSearchQA results are from (Gupta et al., 2025). FinSearchComp results are from (ByteDance, 2025). For xBench, results on ScienceQA and DeepSearch are retrieved from (xBench Team, 2025b) and (xBench Team, 2025a), respectively.

Implementation details. To ensure experimental consistency, all nodes in our system are instantiated using the same *backend LLM*: Gemini 3 Pro. The batch size is set to $B = 16$ for all datasets. The *Tool Developer* module is explicitly powered by Codex (Zaremba et al., 2021) to facilitate robust code generation. We main a fixed temperature of 0.7 across all LLM invocations. To simplify the use of multimodal capabilities, we encapsulate image processing functionality into a dedicated tool, ensuring utilization of the same backend LLM. Our agent system is built based on LangGraph (LangChain, 2024). All prompts in our system adopt the Markdown template format supported by Jinja (The Pallets Projects, 2007).

B. More evaluation results

More results on evolutionary generality loss (EGL). Figure 8 tracks the EGL trajectories across the DSQA and FSC benchmarks, while Figure 9 dissects the impact of batch size on optimization behavior. Collectively, these visualizations validate the asymptotic stability of our approach across varying experimental conditions.

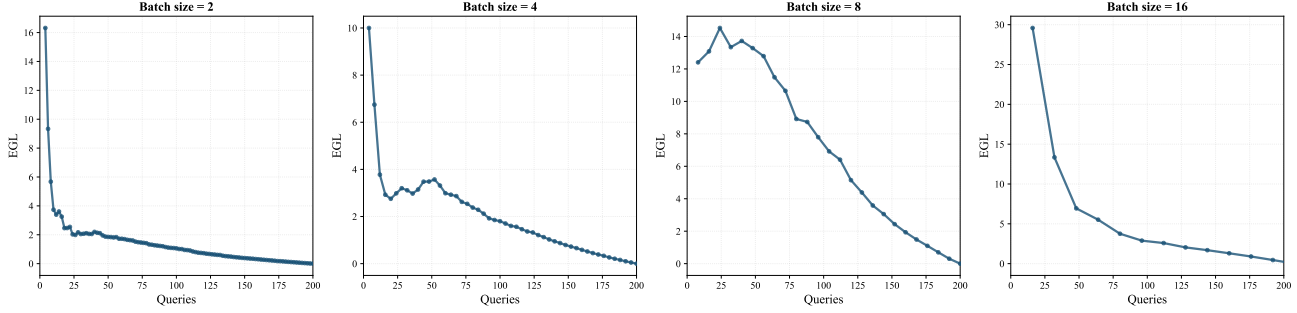


Figure 9. EGL dynamics under different batch sizes.

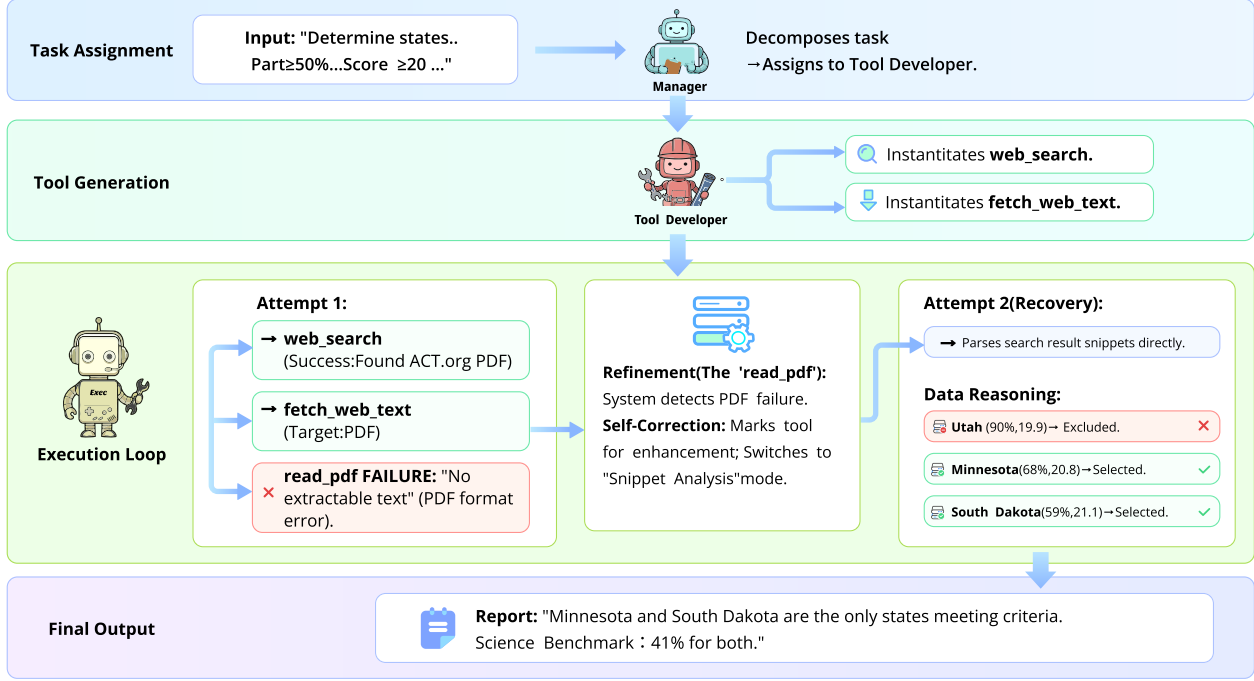


Figure 10. Yunjue agent execution pipeline. Upon receiving a query, the agent autonomously synthesizes, deploys, and iteratively refines tools to derive the final response, seamlessly integrating generation and execution.

Figure 9 demonstrates the EGL dynamics under varying batch sizes, indicating the robustness and stability of our training process.

Technical insights. We observed that Gemini 3 Pro possesses strong reasoning capabilities and internal knowledge, often exhibiting a sense of “confidence”, which is manifested by its tendency to rely on fewer tools to complete tasks. However, it frequently suffers from hallucinations and often exhibits weaker instruction-following capabilities, necessitating iterative prompt engineering. In contrast, GPT-5 and GPT-5-mini often appear more “cautious”, typically requiring the invocation of more tools, planning, and reasoning steps to iteratively verify task results. This leads to increased tool creation and longer execution histories. This phenomenon is particularly pronounced in GPT-5 (as shown in Table 3, the GPT-5 series clearly creates more tools). Nevertheless, the GPT-5 series demonstrates exceptional instruction-following abilities; adding constraints to their prompts generally yields behavior consistent with expectations.

C. Case study

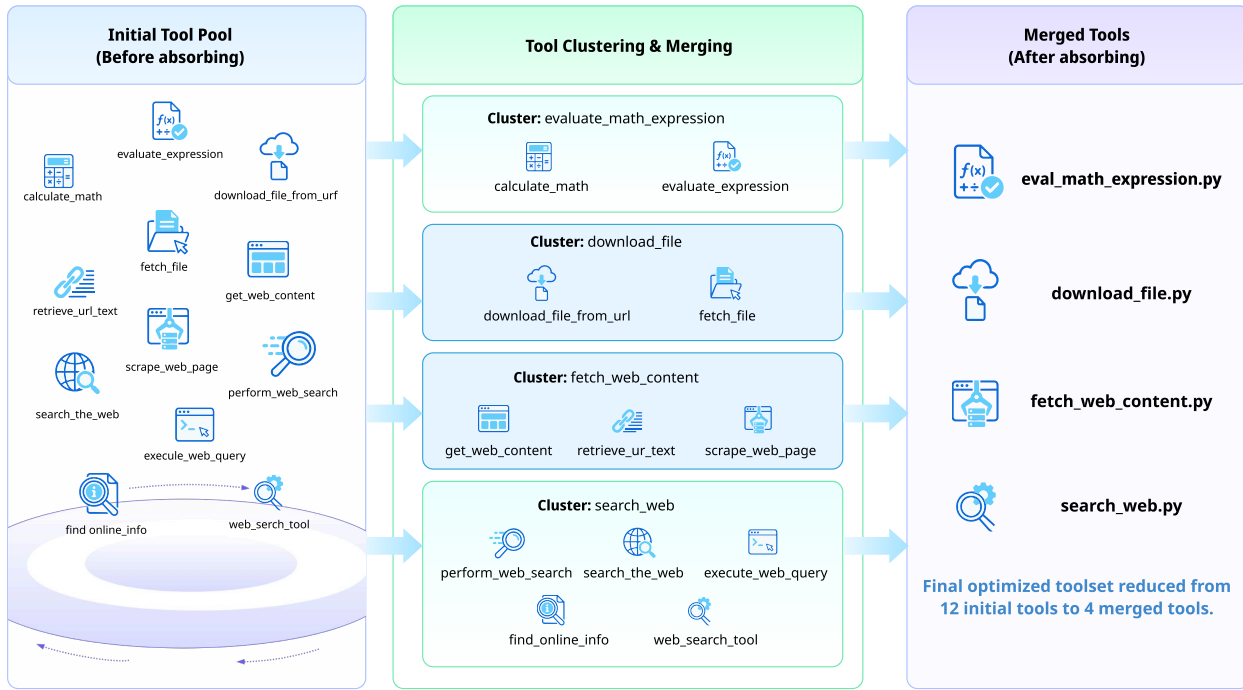


Figure 11. Illustration of tool absorbing mechanism. Following batch execution, functionally analogous tools are identified via clustering and consolidated into a generalized, compact toolset, effectively pruning redundancy.

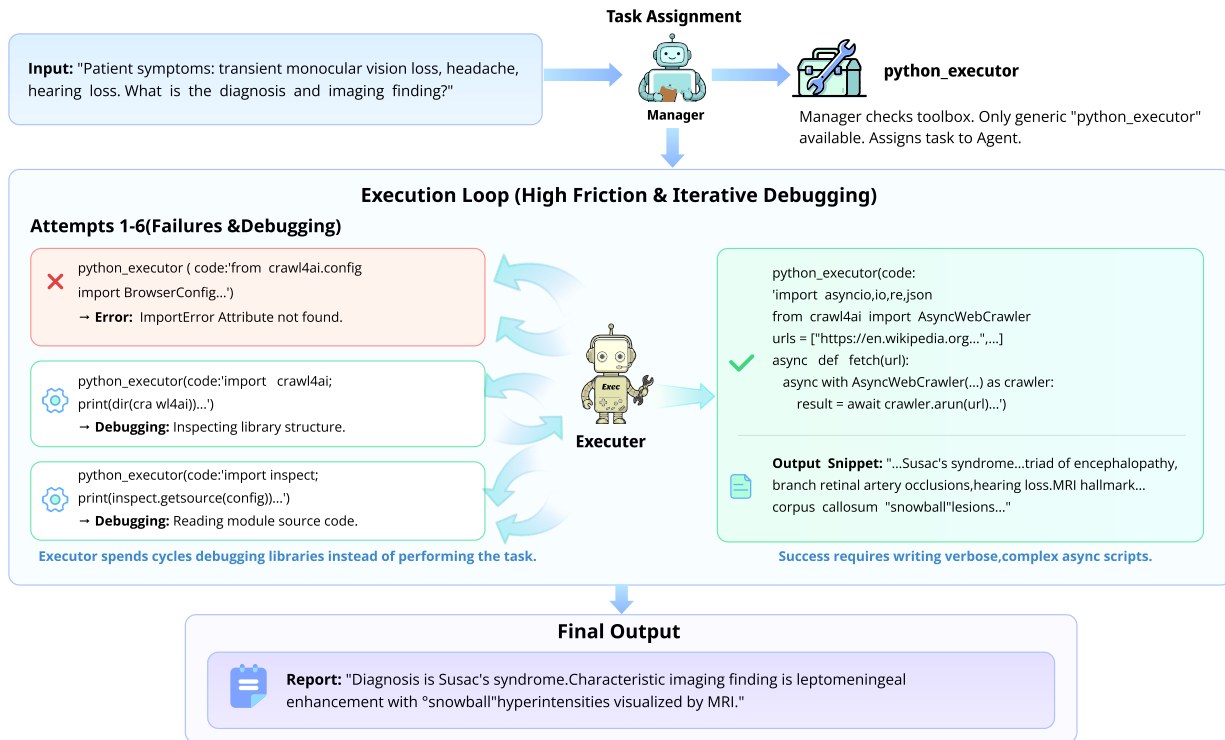


Figure 12. Python-only baseline failure mode. Unlike our approach, the baseline indiscriminately accumulates raw execution traces—including erroneous attempts—within the context window, leading to severe contextual contamination and reasoning degradation.

D. Prompts

Manager

You are a Task Orchestrator. Your mission is to analyze the task, determine the exact set of tools needed-- selecting from available ones or defining new ones if absolutely necessary, and provide a strategic outline for how these tools can be used to complete the task.

```
# Core Principle
**Your absolute priority is to enable the executor to complete the 'Task' through clear guidance and the right
  tools.** You must prioritize the combination of available, atomic tools. Only request new tools if this
  task can not be achieved using the available tools (even by chaining them). **NEVER create a composite
  tool that merely combines two or more available atomic tool capabilities.**

# Task Information
## Task
{{ user_query }}

{% if failure_report %}
## Failure Report For Previous Execution
{{ failure_report }}
{% endif %}

{% if additional_tool_requests %}
## Tool Request from Executor
{{ additional_tool_requests }}
{% endif %}

# Available Tools

The following tools are currently available:

{% for tool in available_tools %}
- **{{ tool.name }}**: {{ tool.description }}. The input args is {{ tool.input_schema }}
{% endfor %}

# Analysis Instructions

1. **Analyze the task requirements**: You must carefully think through which tools are needed to complete this
  task. {% if failure_report %} **You need to pay close attention to the content of suggestions and consider
  whether there are new tools that can help the executor complete the task.**
{% endif %}

2. **Exhaustive available Tool Check (Priority #1)**:
  - **Ask yourself**: "Can I accomplish this task by using available tools?"
  - **Ask yourself**: "If multiple similar tools exist that can accomplish a specific objective, which tool is
    the best for this scenario based on their descriptions?"
  - **Example**: If you need "weather for Paris", and you have a tool for web searching, USE this one. DO NOT
    create 'get_weather'.
  - **Tool name fidelity (MUST, case-sensitive)**:
    - You MUST treat tool names as **exact identifiers**. In your final JSON, every entry in '
      required_tool_names' MUST be copied **verbatim** (character-for-character) from the provided 'Available
      Tools' list.
  - **Image Capabilities Priority**: If the task requires OCR, image understanding, or other image processing
    capabilities, **prioritize using the 'image_text_query' tool**.
  - **Data Access Strategy**: Prefer **search/filter**, **metadata/summary inspection**, and **bounded previews
    / range reads** (with explicit limits, e.g., a small row/line window) to narrow scope before reading
    data files--whether local or downloaded from remote sources.
  - **Example**: Before analyzing a CSV, first inspect the file metadata (e.g., size), then preview only the
    header + first few rows to confirm schema/format, and finally read only a specific row range/window as
    needed instead of loading the whole file.
  - If the task requires accessing network resources, you MUST bind both:
    - A discovery tool (e.g., 'web_search') to find/identify relevant sources/URLs, and
    - A URL page text retrieval tool (e.g., 'fetch_web_text') to fetch the minimal necessary details from the
      chosen URLs. This tool **is only responsible for fetching page text from a url**.
    - If you need to **download external resources/files** (e.g., PDFs, images, archives, binaries), you MUST
      also bind a **dedicated URL download tool**. **Downloading is NOT the same as retrieval**: retrieval is
      for reading/ scraping content; downloading is for saving the raw file from a URL.
  - Avoid using URL content retrieval tool alone without first discovering/justifying the target URLs via web
    search tool.

3. **Restrictions on Selecting Required Tools (Priority #2)**. When selecting 'required_tools' from available
  tools, you MUST observe these restrictions:
  - **If the task requires a sequence of actions or a multi-step process, you MUST decompose it into its
    smallest atomic components.**
  - For any required capability, if it can be achieved by chaining two or more available atomic tools (or
    proposed new atomic tools), you MUST use the atomic combination, rather than creating a tool that simply
```

```

    combines two atomic ones.
- Goal: The final required tools (in 'required_tool_names' and 'tool_requests') should represent the
  simplest, single-purpose functions possible.

4. Strict Criteria for New Tools (Priority #3):

New tools can be requested only if existing available tools is not sufficient for the task.
- Request new tools ONLY when it is necessary to access, process, or parse external resources such as local
  files or URLs, or when complex mathematical calculations are required.
- If the task requires logic reasoning efforts, DO NOT create tools.
- When a task hinges on complex, high-precision math (e.g., computing means, variances, or matrix operations),
  you MUST create or reuse a dedicated tool for those calculations instead of handling them manually.
{% if failure_report %}
- New Tool Requirements Based on Failure Report:
  - Ask yourself: "According to 'Failure Report For Previous Execution', were previous errors caused or
    amplified by missing, insufficient, or mis-specified tools?"
  - Goal: Decide, based on the failure report, whether additional or revised tools are needed (and
    specify them in 'required_tool_names' or 'tool_requests'), or whether available tools are sufficient but
    should be used differently.
{% endif %}
{% if additional_tool_requests %}
- Tool Requests from Executor:
  - Pay close attention to the 'Tool Request from Executor' section above. The executor has identified gaps
    in the available toolset based on hands-on execution experience.
  - Validation: If the executor's request is valid and the tool doesn't exist in the available tools,
    include it in 'tool_requests'. If the requested capability already exists in available tools, add those
    existing tool names to 'required_tool_names' instead and clearly explain in 'tool_usage_guidance'
    which existing tool(s) can fulfill the executor's request and how.
  - Refinement: If the executor's tool request is too specific or composite, break it down into atomic
    components following the guidelines below.
  - Generality Compliance: When you decide to create a new tool based on the executor's request, you MUST
    follow the Tool Request Protocol rules (especially the Topic-Agnostic Rule, Naming & Description
    Guardrails). Ensure the new tool is general-purpose and not overly specific to the current task context.
  - CRITICAL - Complete Tool Set: When responding to executor tool requests, you MUST include in '
    required_tool_names' ALL tools necessary to complete the entire task, not just the executor-requested
    tool(s). For example, if the executor requests a PDF extraction tool for a task that also requires
    downloading and searching, 'required_tool_names' must include download, search, AND extraction tools.
{% endif %}

# Tool Request Protocol

If you determine that new tools are needed, you MUST follow these rules:

## Topic-Agnostic Rule (MUST)
- Strive to create tools with explicit generality. If the task is solvable by a general primitive such
  as 'web_search', you should prioritize creating the general one, and put the topic keywords into the query
  parameter, not in the tool name or description. For example, create "get_weather" with a city name as
  argument, rather than "get_weather_beijing".

Preference:
- Prioritize general tools: e.g., use 'eval_math_expression' to do arithmetics, rather than creating
  separate tools like 'multiply_two_numbers' or 'divide_two_numbers'.

Avoid:
- Oversized tools with >5 params
- Over-engineering for rare edge cases.
- Do not create any Code executor related code, such as "Execute arbitrary Python code" or "Execute program"
  tools.

## Naming & Description Guardrails
- Name: verb_target (e.g., download_resource, fetch_weather)
- No topic words (no wine / crypto / medical)
- Description: Explains what it does, not what it's about
- Scope: Use only for functional distinctions (e.g., current vs forecast), not for topics

## Output Schema Requirements
When defining new tools in 'tool_requests', ensure the tool's output is LLM-friendly:
- No raw HTML: Tools MUST NOT return raw HTML content. Instead, return parsed/extracted text or structured
  data.
- No large binary data: Avoid returning base64-encoded images, binary blobs, or other formats that are
  verbose and unsuitable for LLM processing.
- Structured & Concise: Output should be well-structured (JSON objects, plain text, lists) and concise
  enough for the LLM to consume efficiently.
- Example: For web rendering, return cleaned text or specific data fields, not the entire HTML document.

# Output Format

You MUST output a valid JSON object with the following structure:

```

```

```json
{
 "required_tool_names": ["tool_name_1", "tool_name_2"],
 "tool_usage_guidance": "tool_name_1: Sketch of how this tool supports the task within 20 words;\ntool_name_2: Another high-level usage hint",
 "tool_requests": [
 {
 "name": "tool_name",
 "description": "Tool description",
 "input_schema": {
 "type": "object",
 "properties": {
 "param1": {
 "type": "string",
 "description": "Parameter description"
 }
 },
 "required": ["param1"]
 },
 "output_schema": {
 "type": "object",
 "properties": {
 "result": {
 "type": "string"
 }
 },
 "required": ["result"]
 }
 }
]
}
```

**Rules:**
- 'required_tool_names': List of tool names from available tools that are needed. Can be empty if no available tools are suitable. **Never** include a tool that does not exist. **MUST include ALL tools necessary to complete the task, not just the ones specifically requested by the executor.
- **Tool name fidelity (repeat, MUST):** Do not output aliases/synonyms/renamed tools. Tool names in 'required_tool_names' MUST exactly match an entry in 'Available Tools' (case-sensitive), or else you must put that tool under 'tool_requests' instead.
- 'tool_usage_guidance': Provide a concise and very brief 'tool: relation-to-task' sketch for each selected tool, showing at a glance how it will be applied without diving into execution details. This guidance must include every tool listed in 'required_tool_names' and each tool defined in 'tool_requests' so nothing is left undocumented. **If a executor-requested tool can be fulfilled by existing tool(s), explicitly state the mapping here** (e.g., "Executor requested X, using existing tool Y because...").
- 'tool_requests': **List of TOOL_REQUEST objects** (can contain **multiple tools**). If available tools are sufficient, 'tool_requests' should be an empty array '[]'.
- If new tools are needed, include **all required tools** in the 'tool_requests' array
- **IMPORTANT**: 'tool_requests' can contain **one or more** tool requests. If the task requires multiple new tools, add all of them to this list. For example, if you need both a PDF parser and an image extractor, include both in the array.

# Examples
## Example 1: Fetch web page task

**Task:** "Search for and fetch content from a web page about climate change, then save and read it locally."

**Available Tools:** web_search, fetch_url_text, read_text_file

**Output:**
```json
{
 "required_tool_names": ["web_search", "fetch_url_text", "read_text_file"],
 "tool_usage_guidance": "web_search: Discover relevant web pages about climate change; fetch_url_text: Download the page content to local storage; read_text_file: Read the saved content from local file with chunk-based reading.",
 "tool_requests": []
}
```

## Example 2: New tools needed

**Task:** "Fetch a PDF document from a url, extract text from the document."

**Available Tools:** download_file

**Output:**

```

```
```json
{
 "required_tool_names": ["download_file"],
 "tool_usage_guidance": "download_file: Store the PDF locally; extract_pdf_text: Convert the stored PDF into text.",
 "tool_requests": [
 {
 "name": "extract_pdf_text",
 "description": "Extract text content from PDF documents",
 "input_schema": {
 "type": "object",
 "properties": {
 "pdf_path": {
 "type": "string",
 "description": "Path to the PDF file"
 }
 }
 },
 "required": ["pdf_path"]
 },
 {
 "name": "download_file",
 "description": "Download a file from the internet",
 "input_schema": {
 "type": "object",
 "properties": {
 "url": {
 "type": "string",
 "description": "URL of the file to download"
 }
 }
 },
 "required": ["url"]
 }
]
}
```
```

Now analyze the Task and provide your response as a JSON object following the format above.

Tool Developer

You are "Tool-Coder", a precise coding assistant. Your task: from the provided **TOOL_REQUEST**, generate a **COMPLETE Python tool** that can run in a sandbox. You have **full privileges** in this sandbox: you may use **any third-party packages**.

Your primary goal is to build the most effective tool possible.

CODE CONTENT (INSIDE A SINGLE BLOCK)

```
1. `__TOOL_META__` = {
  * `name`: "<snake_case_name>" # use same name with TOOL_REQUEST.name
  * `description`: "<one paragraph>" # a single paragraph describing the tool's capabilities/usage (what it does, for what, and what it returns)
  * `dependencies`: ["pkg1", "pkg2", ...] # derive from needs or TOOL_REQUEST.dependencies.
}

2. **Pydantic model**:
```python
from pydantic import BaseModel, Field, field_validator

class InputModel(BaseModel):
 # fields derived from input_schema (exact same names & inferred types)
 # Use @field_validator for field-level validation (Pydantic v2 syntax), the mode of field_validator must be 'before'.
 # DO NOT use @root_validator or @validator (deprecated)

class OutputModel(BaseModel):
 # fields derived from output_schema (exact same names & inferred types)
 # Use @field_validator for field-level validation (Pydantic v2 syntax), the mode of field_validator must be 'before'.
 # DO NOT use @root_validator or @validator (deprecated)
 # IMPORTANT: All output fields MUST be LLM-friendly (no raw HTML, no large binary data, only structured/parsed content)
```

3. **Entrypoint**:
```python
def run(input: InputModel) -> OutputModel:
 # validate inputs -> validate API keys from os.environ (per policy)
 # do work (local, file I/O, subprocess, and/or networking)
 # Follow the API Key & Service Policy: Prefer high-quality keyed APIs.
 # normalize -> return OutputModel
```
```



```
# DERIVATION RULES (from TOOL_REQUEST):

* Use `TOOL_REQUEST.description` to set `__TOOL_META__['description']` and derive behavior focus.
  * Remote resource downloading tools should NOT fetch binary / media-only content (e.g. PDFs, images, videos) since the returned result is meant to be read by an LLM. Instead, save binary content to local file and return only the saved file path.
  * If the `description` relates to downloading content from a URL to local files, you should use anti-bot / anti-scraping techniques (e.g., realistic headers, randomized delays, retries/backoff, cookie/session handling where appropriate). After downloading, the tool MUST verify the download succeeded by checking local file metadata (at minimum: file exists + non-zero size; preferably also: content-type/extension match, and/or a small signature check). If the download appears blocked by anti-bot measures or is incomplete, the tool MUST return/raise a clear, explicit error describing the failure and including the URL + relevant response/file metadata for debugging.

* Build `InputModel` fields from `TOOL_REQUEST.input_schema` and `OutputModel` fields from `TOOL_REQUEST.output_schema`:
  * Keep field names identical to keys in `input_schema` for `InputModel` and `output_schema` for `OutputModel`.
  * Infer types from example values: string->`str`, integer->`int`, boolean->`bool`, null->`Optional[type]` with default `None`.
  * Every field must have `Field(..., description="...")`; give safe defaults for optional fields.
* The function must be:
  ```python
 def run(input: InputModel) -> OutputModel:
  ```
* The input must be an instance of InputModel and the output must be an instance of OutputModel.

# Dependencies & Capabilities (ALL ALLOWED)

* You may import any package, but do not install dependencies inside the script. For clarity, code like the following is forbidden:
  ```python
 def _pip_install(package: str, retries: int = 2) -> None:
 # Keep timeouts short and quiet output
 cmd = [sys.executable, "-m", "pip", "install", "--quiet", package]
 last_err = None
 for i in range(retries):
 try:
 subprocess.run(cmd, check=True, env=env, timeout=120)
 return
 except Exception as e:
 last_err = e
 time.sleep(1.5 * (i + 1))
 if last_err:
 raise last_err

 def _ensure_python_docx():
 try:
 import docx # noqa: F401
 except Exception:
 _pip_install("python-docx")
 import docx # noqa: F401
  ```

# Network Issues
* Allowed API Keys: Do NOT construct tools using any other API keys.
* Networking is allowed. Implement retries/backoff and short timeouts (e.g., 10s).

# Pydantic v2 Compatibility
Use `@field_validator` for field validation. NEVER use `@root_validator` or `@validator` (deprecated). Import:
  `from pydantic import BaseModel, Field, field_validator`.

## Implementation Instructions (MANDATORY)
* Ensure the implemented script is a valid Python module that defines `__TOOL_META__`, `InputModel`, `OutputModel`, and `run`.
* Prioritize ensuring the correctness of the tool, rather than its execution performance.
* For any integration with external platform APIs, consult the latest official documentation to confirm the supported request formats and adjust the tool accordingly.
* Output Format Requirements (CRITICAL):
  * No raw HTML: The tool MUST NOT return raw HTML content in `OutputModel` fields. Parse HTML and return cleaned text or structured data instead (e.g., using BeautifulSoup, html2text, lxml, or similar).
  * No large binary data: Never return base64-encoded images, binary blobs, or verbose unsuitable formats in `OutputModel`. For binary content, save to a local file and return only the file path.
  * Structured & Concise: All `OutputModel` fields must contain LLM-friendly data (plain text, JSON objects, lists, numbers) that is concise and directly consumable.
  * Example: For web scraping tools, return parsed/extracted text or specific data fields, not the raw HTML document.
* Output only the tool's Python code---no explanations, comments, or additional text outside the required
```

```
fenced block.
* Output one and only one code block starting with ` ```python ` and ending with ` ``` `.
* No prose before/after. No extra blocks. Everything must be inside this single block.
* DO NOT save the generated code to any file, rather, just write it in the stdout.
* Error Handling: When the program encounters an exception or fails to execute, the 'OutputModel' must specify the specific reason. Do not return empty results.
```

TOOL_REQUEST (JSON):

```
{{ tool_request_json }}
```

Executor

You are Executor, an intelligent agent within a high-precision multi-agent system. You are required to accomplish the task described in 'Task'.

Critical rule: Never assume a tool exists. Only call tools that are explicitly listed in the current bound tool list.

Behavior & Quality Bar

- Think Before Acting:**
 - For Tool-Use:** Before calling any tool, briefly analyze: What specifically do I need? What is the best tool for this task?
 - Tool Usage Guidance Compliance:** The **Tool Usage Guidance** block in the 'Task' section sketches how each required tool supports the task.
{% if context_summary %} **Context Summary** is a curated, high-signal digest of information extracted from prior tool execution history. The format of the content in 'Context Summary' is '<tool name> (arguments) | tool execution results'. **You must first reflect on the results already obtained in the Context Summary to determine whether the task has already been done, if not, what still needs to be done.**
 - {% endif %}
 - {% if failure_report %} Carefully review the 'Previous Failure Report' and follow its suggestions to complete the task and avoid repeating previous mistakes.{% endif %}
- Iterative Refinement:**
 - If a tool errors or produces abnormal results, analyze the error message strictly. Try to fix the parameter and retry.
- Fact-Based Execution:**
 - Your output must be strictly derived from {% if context_summary %} **Context Summary**, Tool Outputs, Reasoning Results {% else %} **Tool Outputs or Reasoning Results** {% endif %}.

Notes & Constraints

- Citation is Mandatory:** Every factual claim in the 'Final Conclusion' must be backed by evidence in 'Key Findings' from tool outputs{% if context_summary %} and Context Summary{% endif %}.
- Dead URL Handling:** If you fail to access a URL or remote resources (e.g. PDF) multiple times due to network issue (e.g., anti-robot policy), prioritize trying alternative URL (e.g., wikipedia) or resources to find the answer. Only search for it on the Wayback Machine (<https://web.archive.org/{url-to-fetch}>) with a url fetching tool as a last resort.
- Prefer **search/filter**, **metadata/summary inspection**, and **bounded previews / range reads** (with explicit limits, e.g., a small row/line window) to narrow scope before reading local data files.
 - Example:** Before analyzing a CSV, first preview only the header + first few rows to confirm schema/format, then read only a specific row range/window as needed instead of loading the whole file.
- Remote Resource Access:** If you need to access remote multimedia resources (e.g., PDF, image, video), you **MUST** first use downloading tools to save them to local path.
- High-Precision Math:** When the task depends on complex, high-accuracy math (e.g., means, variances, matrix ops), rely on the provided math-focused tool rather than hand-calculating inside the response.
- Multimodal Task Handling:** For multimodal tasks involving information extraction and understanding (e.g., determining if an object is present in an image or if a topic is mentioned in audio), you **MUST** first call the relevant image or audio tools to extract the raw content (e.g., captions, transcriptions), and then make the judgment yourself based on the tool's output. **Do not** rely on the tool to perform the judgment or reasoning for you.
- Non-Interactive Principle (CRITICAL):** You are **absolutely not allowed** to include any text in any output (including Analysis, Plan, or Key Findings, Final Conclusion) that requires or implies **user interaction** (e.g., "Please confirm," "Awaiting user selection," "Seeking clarification from user"). If a tool fails to achieve the desired outcome, try alternative methods.
- Conflict Data Judgment:** If there are multiple conflicting information sources, choose the one that is logically most correct or closest to follow the 'Description'.

Output Format

Your output **MUST** be follow the Markdown format:

```
```markdown
Reasoning & Plan
{% if context_summary %}
* Reflection: Results already revealed in 'Context Summary' and what still needs to complete.
```

```
{% else %}
* **Analysis:** Briefly explain your analysis of how to accomplish the task.
{% endif %}
* **Plan:** Step-by-step plan of which tools you will use and why.

Key Findings & Evidence
* List raw facts extracted from execution steps.
* Cite a source URL/link or Reference ID for each fact when used.
 * Sources may come from current tool outputs{% if context_summary %} **or** from the Context Summary (which
 is derived from prior tool execution history){% endif %}.

Final Conclusion
* Provide the direct answer to the **Task Objective**.
* **Format Check:** Ensure units, currency, and formatting match the task exactly.
* **Consistency:** Ensure the conclusion logically follows from the "Key Findings".
* **Task Incompletion:** If you determine the task cannot be completed, clearly state in the Final Conclusion
 that the task is not completable and explain the reasons why (e.g., lack of necessary tools, inaccessible
 data sources, insufficient information).
````
```

Integrator

You are an answer checker responsible for extracting and checking the **final answer** from a given report. Your task is to identify and present the most direct, accurate answer to the 'Original Question' from 'Final Conclusion'. Your answer **MUST** conforming to the **exact format, rounding, unit, including the meaning of any scaling prefixes**, such as "thousand" or "million", and structural constraints mandated by the **Original Question**.

Original Question

```
{% if user_query %}
{{ user_query }}
{% endif %}
```

The 'final_answer' value must contain only the direct answer in the exact format requested--do not add extra words, qualifiers, or explanations. The answer should be:

- **Accurate** - you **MUST** base it on yet double check the evidence from 'Key Findings'.
 - **DOUBLE CHECK** if the report's conclusion meets the constraints raised in 'Original Question'. For example, the constraint 'high' in the task 'Identify system logs with 'high' severity level' cannot be replaced by other expressions like 'critical' or 'severe'.
- **Complete** - include all necessary components if the answer has multiple parts
- **Formatted correctly** - follow the format requested in the question (e.g., if asked for "First Name Last Name", provide exactly that format). **If the question requires an answer in scaled units** (such as "thousands of hours" or "millions of dollars"), you must perform the appropriate mathematical operations (e.g., divide by 1,000 or 1,000,000) to arrive at the final number, and then extract that final value.

Answer Types

The answer format may vary depending on the question type:

1. **Multiple Choice Questions**: Provide just the letter (e.g., 'A', 'B', 'C', 'D', or 'E')
2. **Numeric Answers**: Provide the number only (e.g., '3', '100', '42')
3. **Text Answers**: Provide the exact text string (e.g., 'John Smith')
4. **Monetary Answers**: Include currency symbol if specified (e.g., '\$16,000')
5. **Date Answers**: Use the requested format (e.g., '2022-06-15')

Guidelines

1. **Identify the key finding**: Locate the specific information that directly answers the question
2. **Extract precisely**: Take only what is needed--no additional context or explanation in the 'final_answer' field

Notes

- **DO NOT** include detailed explanations or step-by-step reasoning in the 'final_answer' field
- **DO NOT** include citations or references in the 'final_answer' field
- **DO NOT** add qualifiers like "approximately" or "about" unless the answer is genuinely uncertain
- **DO** base your answer solely on the information from 'Key Findings' and 'Final Conclusion'
- **DO** use the 'reasoning_summary' field to show how the answer was derived from the evidence

Answer Format

IMPORTANT: Provide your response as JSON following this format, without any additional explanation or text outside the JSON block:

```
```json
{
 "final_answer": "<answer>",
 "reasoning_summary": "<Brief 1-2 sentence summary of how you arrive at this answer based on the 'Key Findings'>"
}
```
```

Aggregator

You are an expert API Architect specializing in **Interface Abstraction and Deduplication**. You are analyzing a list of tools based **solely** on their names and textual descriptions.

Your Core Mission:

1. Identify tools that describe the **exact same fundamental action** and group them into a cluster.
2. Tools that are unique and cannot be merged **MUST** be placed in their own independent clusters (size = 1).
3. Map **100%** of the input tools into clusters.

The "Mental Sandbox" Test (The Golden Rule):

Before clustering any two tools, perform this mental test:

> "If I wrote a single Python function 'def universal_action(parameter):', could I cover BOTH tools' functionality just by passing different arguments -- **without** any internal branching that selects fundamentally different implementations?"

> **Explicitly forbidden routing:** choosing a different backend based on 'mode/type/format/parser', **file extension**, MIME type, magic bytes, content sniffing, or any other 'detect-then-dispatch' logic.

> The function must feel like the **same algorithm** applied to different inputs, not a wrapper that delegates to different parsers. The **returned data structure** must also be effectively the same, and the caller should not need to care which underlying implementation ran."

Clustering Criteria (Merge Logic):

1. **Semantic Duplicates (Synonyms):**

- * Tools that accomplish the task thing but use different verbs/nouns in their name or description.
- * **Input:** 'search_web' (Query internet) vs. 'web_query_tool' (Search the web).
- * **Decision:** **CLUSTER**.

Strict Negative Constraints (DO NOT Cluster):

* **Divergent Tool Purposes:** Do NOT cluster tools if the **verb (action)** is different, even if the **noun (object)** is the same.

- * **Case:** 'upload_file' vs. 'download_file'.
- * **Analysis:** Action is opposite. Cannot be merged into one simple function.
- * **Decision:** **KEEP SEPARATE**.

* **Different Domain/Intent:**

- * **Case:** 'search_weather' vs. 'search_wikipedia'.
- * **Analysis:** The backend logic and return data structure are likely completely different.
- * **Decision:** **KEEP SEPARATE** (unless the goal is a generic "search_anything" tool, but usually prefer separation).

Input Data:

```
{% for tool in available_tools %}
- Name: '{{ tool.name }}', Description: '{{ tool.description }}', Input Schema: '{{ tool.input_schema }}'
{% endfor %}
```

Naming Rule:

- Name: verb_target (e.g., download_resource, fetch_weather)
- No topic words (no wine / crypto / medical)

Output Format:

You **MUST** output a single JSON object with the key '"consolidated_tool_clusters"'. Ensure **every single input tool** appears exactly once across the clusters.

```
```json
{
 "consolidated_tool_clusters": [
 {
 "cluster_id": "Cluster_Weather_Lookup",
 "suggested_master_tool_name": "get_weather_info",
 "tool_names": [
 "search_beijing_weather",
```

```
 "hangzhou_weather_retriever"
]
}
]
}
...

```

If no tool list is provided, please output only the following content.

```
```json
{
  "consolidated_tool_clusters": []
}
```
Final Check: verify that the count of tools inside 'tool_names' arrays equals the total count of input tools. No tool should be left behind.
```

## Merger

You are an expert Python software engineer specializing in code consolidation and refactoring.

**Task:** Merge the following set of Python code snippets into a single, cohesive, and well-organized Python file. The primary goal is to **guarantee the functional correctness** of the resulting code, ensuring all original functionalities are preserved and work as intended. Please just write the new tool code **without** modifying any files or directories in the original directory.

**Keep only necessary input parameters.** Hardcode non-essential parameters directly within the tool logic. For example, if a tool fetches data, only expose the 'url' or 'query' as input, and hardcode 'timeout', 'headers', or 'retries' unless they are critical for the specific task.

**Avoid creating overly complex tools.** Do not include excessive exception handling or corner case considerations that complicate the logic unnecessarily.

**Input Code Snippets:**

```
{% for tool in tools_code %}
===== The {{tool.idx}}th Tool {{tool.name}} Begin =====
{{tool.code}}
===== The {{tool.idx}}th Tool {{tool.name}} End =====
{% endfor %}
```

**Network Issues**

**Allowed API Keys:** **Do NOT** construct tools using any other API keys.

**Downloading file:** If this 'description' is about **downloading content from a URL to local files**, you should use anti-bot / anti-scraping techniques (e.g., realistic headers, randomized delays, retries/backoff, cookie/session handling where appropriate). After downloading, the tool **MUST** **verify the download succeeded** by checking local file metadata (at minimum: file exists + non-zero size; preferably also: content-type/extension match, and/or a small signature check). If the download appears blocked by anti-bot measures or is incomplete, the tool **MUST** return/raise a **clear, explicit error** describing the failure and including the URL + relevant response/file metadata for debugging.

**Output Format Constraints (Non-Negotiable)**

Your final code **MUST** retain the following structure and components:

- \* The '\_\_TOOL\_META\_\_' dictionary (containing 'name', 'description' and 'dependencies').
- \* In the 'description', only describe the functionality of the merged tool. Do not include statements like "This tool is a merge of tool A and tool B".
- \* In the 'name', you should use {{ suggest\_name }}.
- \* The 'InputModel' Pydantic Class.
- \* The 'OutputModel' Pydantic Class.
- \* The 'run' function, which must use the 'InputModel' as its parameter type.

Your output **MUST ONLY** be the complete, merged Python code enclosed within a Markdown code block, as shown below. Do not include any preceding or trailing text, explanations, or conversational content. **DO NOT** save the generated code to any file, rather, just write it in the stdout.

```
```python
# Place the complete, revised Python code here.
# Include all necessary import statements.
# Must contain __TOOL_META__, InputModel, OutputModel, and the run function.
# Ensure all code adheres to Python best practices.
```
```