

Towards General Agentic Intelligence via Environment Scaling

Anonymous ACL submission

Abstract

Advanced agentic intelligence is a prerequisite for deploying Large Language Models in practical, real-world applications. Diverse real-world APIs demand precise, robust function-calling intelligence, which needs agents to develop these capabilities through interaction in varied *environments*. The breadth of function-calling competence is closely tied to the diversity of environments in which agents are trained. In this work, we scale up environments as a step towards advancing general agentic intelligence. This gives rise to two central challenges: (i) how to scale environments in a principled manner, and (ii) how to effectively train agentic capabilities from experiences derived through interactions with these environments. To address these, we design a scalable framework that automatically constructs heterogeneous environments that are fully simulated, broadening the space of function-calling scenarios. We further adapt a two-phase agent fine-tuning strategy: first endowing agents with fundamental agentic capabilities, then specializing them for domain-specific contexts. Extensive experiments on agentic benchmarks, τ -bench, τ^2 -Bench, and ACEBench, demonstrate that our trained model, **AgentScaler**, significantly enhances the models’ function-calling capability.

1 Introduction

Function calling empowers language agents to interface with the real world (Qin et al., 2023; Chen et al., 2024b; Qin et al., 2024; Schick et al., 2023; Su et al., 2025b). Yet, their progress is fundamentally constrained by the scarcity of agentic data¹, *i.e.*, trajectories generated by autonomous agents interacting with environments via explicit action executions, namely, tool calls (Zhou et al., 2023; Liu et al., 2024a). The community has gradually

¹In this paper, the terms “function-calling”, “tool”, “API”, “MCP” are used interchangeably; “agentic data” refers to trajectories involving such interactions.

transitioned from the era of raw corpora and human-curated data to the emerging *era of experience* (Silver and Sutton, 2025; Wu et al., 2025a; Li et al., 2025; Tao et al., 2025; Geng et al., 2025). Crucially, language agents must experience these interactions themselves in a predefined environment, which makes both data collection and reliable supervision highly challenging.

Several approaches have been attempted to generate synthetic agentic data. Broadly, previous methods fall into two categories. The first category follows a reverse paradigm, in which user queries are generated to match each assistant function call observed at every interaction turn (Yin et al., 2025), though the resulting trajectories may exhibit limited realism. The second category follows a forward paradigm, which we refer to as simulated agent–human interplay (Chen et al., 2024a; Liu et al., 2024b; Prabhakar et al., 2025a; Barres et al., 2025; Zeng et al., 2025). Such generated trajectories, however, may lack naturalness. In this category, a high-level user intent is first formulated to necessitate agent interaction. Agentic data is then constructed in a top-down manner based on this intent through human–agent interplay. Yet, the environment is not scalable: the absence of automated environment construction hinders large-scale deployment and inevitably entails some degree of manual intervention.

To address these challenges, we pursue the advancement of general agentic intelligence via systematic environment scaling. Our approach follows a principled two-stage pipeline: (i) **fully simulated environment construction and scaling**, responsible for establishing and expanding diverse agentic scenarios, and (ii) **agent experience learning**, which exploits these environments to foster generalizable intelligence.

In designing environment construction and scaling, we follow the principle that the core of an agent lies in its capacity for environment interaction, with

each environment instantiated as a read-write database (Barres et al., 2025; Zeng et al., 2025). Specifically, we collect a broad spectrum of APIs and organize them into domains using community detection, where each domain represents an environment aligned with a specific database structure. Then, we instantiate tools as executable code, thereby achieving programmatic materialization that enables direct operations on the underlying database structures. Finally, we sample from the domain-specific tool graph to generate parameters for the tool sequences and initialize the corresponding database state. We then integrate these components into an overall user intent, grounding tool executions directly on the database. This design enables verifiability at both the environment level and the tool-argument response level.

For learning from agent experience, our focus is on training the agent’s ability to perform tool calls and to respond effectively to users (Ye et al., 2025b; Su et al., 2025a). We begin by performing simulated human-agent interactions on the constructed agentic tasks (Prabhakar et al., 2025a), thereby collecting trajectories that serve as the agent’s experience and perform strict filtering.

To facilitate the acquisition of this capability, we adopt a two-stage agent experience learning framework: in stage 1, the agent acquires fundamental tool-calling skills across general domains; in stage 2, it is further trained within target vertical domains using domain-specific scenarios, enabling smoother and more context-aligned development of agentic capabilities.

Extensive experiments on agentic benchmarks, τ -bench (Yao et al., 2024), τ^2 -Bench (Barres et al., 2025), and ACEBench (Chen et al., 2025) show the effectiveness of our pipeline and trained models. Based on the above pipeline, we train our family of **AgentScaler** models (4B, 8B, 30B-A3B), built upon the Qwen-3 (Team, 2025b) series. At each comparable scale (4B, 8B), our models achieve state-of-the-art performance. Notably, AgentScaler-30-A3B sets a new state-of-the-art with significantly fewer parameters among models with comparable active parameter size, delivering results on par with existing 1T-parameter models and leading closed-source systems. We also provide a systematic analysis covering model generalization, stability, and the long-horizon tool-calling challenge, offering key insights into the development of general agentic intelligence.

2 Related Work

Tool-Use Environments The construction of tool-use environments primarily involves three approaches: real-world environments, LLM-simulated environment, and Simulated Environments based on a state config. Using real-world environments (Qin et al., 2023; Song et al., 2023; Mastouri et al., 2025; Wu et al., 2025b) to invoke actual tools yields the most authentic feedback and enhances the model’s robustness in practical applications. However, this requires frequent calls to MCP services, resulting in high costs and significant time overhead. Moreover, maintaining a highly available and stable MCP service is often difficult, posing major challenges for agentic data generation and online RL training of models. Many works use LLM-generated responses to simulate environments as a source of tool responses (Qin et al., 2024; Lu et al., 2024; Sun et al., 2025). By leveraging strong or fine-tuned LLMs, these approaches generate plausible responses given a tool call. However, such methods struggle with issues like hallucination and inconsistent response variability. To address the limitations of the above two approaches, some recent work (Ye et al., 2025b; Yao et al., 2024; Barres et al., 2025; Prabhakar et al., 2025b; Ye et al., 2025a) proposes building an offline tool execution environment for LLM training and evaluation. On one hand, offline environments avoid calling real tools, significantly reducing response generation cost and latency. On the other hand, mocked tool usage in such environments can still interact with real databases or state files through actual execution. However, these methods are more commonly applied in LLM evaluation rather than training, as constructing a reliable tool suite and a high-fidelity execution environment typically requires substantial manual effort. Furthermore, it is difficult to automatically validate the quality of such environments without human involvement, making scalability a significant challenge. Our approach enables domain scalability through sampling from a toolgraph, and eliminates the need for human intervention via a rigorous, rule-based validation pipeline. This makes scalable construction of tool execution environments feasible.

Tool Learning To enhance the agentic capabilities and tool-calling abilities of models, many works have attempted to improve tool utilization through various approaches. For instance, xLAMs (Prabhakar et al., 2025b; Zhang et al., 2024) and

ToolAce (Liu et al., 2024a) leverage large-scale agentic data synthesis pipelines to generate high-quality training data. DiaTool-DPO (Jung et al., 2025) employs DPO to enable models to learn from multi-turn positive and negative trajectories. Meanwhile, Tool-RL (Qian et al., 2025), Tool-N1 (Zhang et al., 2025) utilize reinforcement learning (RL) to enhance both the tool-calling proficiency and generalization ability of models, further pushing the performance boundaries beyond supervised fine-tuning. Overall, whether relying on agentic data synthesis or online interaction with environments via RL training, a reliable, and scalable execution environment is essential. Our method not only leverages accurately simulated tool environments to collect trajectories but also introduces verifiable environmental state changes, making each simulation response more reliable. Furthermore, we propose a state change based environment validation strategy, enabling a robust filtering mechanism for large-scale agentic data synthesis.

3 Environment Build and Scaling

Design Principal There already exist many real-world environments for agent interaction (Qin et al., 2023). However, these real environments suffer from several practical limitations that make them unsuitable for large-scale and stable training. First, many real online APIs are unstable, frequently affected by rate limits, fluctuating QPS capacity, and transient API failures. Such instability prevents agents from receiving consistent tool feedback, which is crucial for effective learning. This limitation has also been highlighted by multiple recent works (Guo et al., 2024, 2025), motivating the development of simulated environments to enable reliable and large-scale agent training. Second, simulated environments can provide deterministic, reproducible, and controllable tool feedback, which real-world environments cannot guarantee. This stability is especially important for training models that rely on precise multi-step tool interactions (Sun et al., 2025; Su et al., 2025b).

In essence, any function call can be interpreted as a read-write operation over an underlying environmental database \mathcal{D} (Guo et al., 2025). Specifically, each function $func$ can be assigned an operator type, $op(func) \in \{\text{read}, \text{write}\}$, where read-type function perform queries over \mathcal{D} (e.g., retrieval, inspection, monitoring), while write-type tools induce state transitions in \mathcal{D} (e.g., modifica-

tion, generation, actuation). Under this abstraction, a tool response is equivalent to evaluating the induced operator on \mathcal{D} , i.e., $API(func, \alpha) \equiv op(func)(\alpha; \mathcal{D})$, where the symbol α denotes the input arguments provided to a function call. Furthermore, let \mathcal{T}_d denote the set of tools within domain d . Tools in the same domain typically exhibit structurally similar read-write patterns, which can be captured by a common database schema \mathcal{S}_k . Consequently, the design problem reduces to defining a partition of the tool space into domains $\{\mathcal{T}_1 \dots, \mathcal{T}_M\}$, and assigning to each domain a database schema \mathcal{S}_k , where \mathcal{S}_k specifies the environment for that domain.

3.1 Environment Automatic Build

Building upon this design principle, we propose a systematic pipeline for leveraging a diverse set of tools as shown in Figure 1. We begin with **scenario collection**, which gathers a large corpus of real-world tools; proceed to **tool dependency graph modeling**, which induces well-structured domain partitions and distributions; and finally employ **function schema programmatic materialization**, which maps tool operations onto database interactions, thereby enabling the construction of the overall environment.

Scenario Collection We collected more than 30,000 APIs from ToolBench (Qin et al., 2023; Guo et al., 2024), API-Gen (Prabhakar et al., 2025b) and online tool repository. After applying rigorous filtering, including the removal of low-quality APIs and subsequent refinement, we rewrite some API descriptions to incorporate explicit input-output specifications (Fang et al., 2025). Building on this, we further constructed tool compositions by systematically exploiting the input-output relationships among APIs. This process ultimately resulted in API pools Θ_F whose size = N (over 30,000), providing a reliable foundation for subsequent experiments and analysis.

Tool Dependency Graph Modeling We construct a tool graph in which nodes are tools and edges encode compositional compatibility induced by function parameters. A tool $func$ consists of a description D_{func} and a list of parameters P_{func} . For a pair of tools, we can extract their respective parameter lists and convert them into vector representations ϕ to compute their cos-similarity. If the similarity exceeds a predefined threshold τ , we consider there to be a dependency relationship between the two tools. Accordingly, we insert an

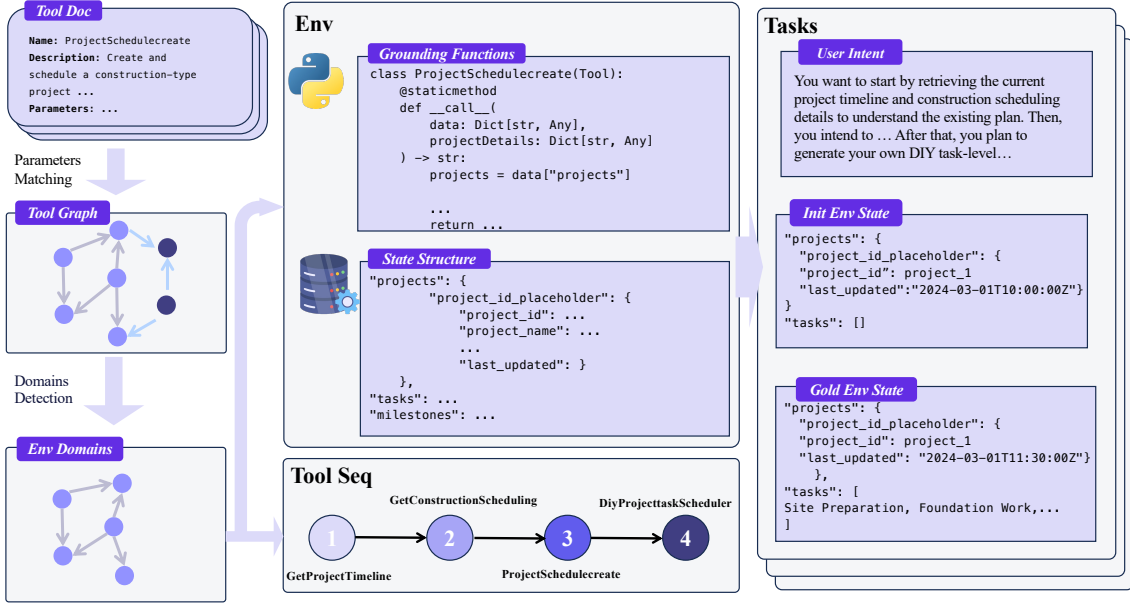


Figure 1: The overview of the environment automatic build, and agentic environment construction. (i) Tool schemata are matched by parameter similarity to form a tool graph; (ii) tools are clustered via community detection and sampled through random walks; (iii) executable functions and a state repository are constructed accordingly.

edge E between them in our graph.

$$E = \{(i, j) \mid \text{sim}(\phi(P_{func_i}), \phi(P_{func_j})) > \tau\} \quad (1)$$

Domain partitioning then reduces to a graph clustering problem. We employ Louvain community detection (Blondel et al., 2008) to identify coherent tool communities that serve as domains. For a segmented tool set, since parameter matching relies solely on vectorization and considers only individual parameter information, the overall inter-tool dependencies may be difficult to capture. Therefore, for tools within a given domain, we further employ an LLM to examine the dependencies between each pair of tools, thereby further improving the accuracy of edges in the tool graph. In total, we obtained M domains (exceeding 1,000).

Function Schema Programmatic Materialization We first leverage the parameters of all tools within a domain to generate a domain-specific database structure, which serves as the underlying state for subsequent tool operations. After obtaining the domain-specific tool set and the corresponding database schema in the previous stage, we can formalize each tool in python code, enabling it to perform read-write operations over the database schema. Interestingly, when generating database structures and formalizing code within specific domains of τ -bench, we observe through manual inspection that our outputs exhibit a high degree of

consistency with the official implementations provided by τ -bench (Yao et al., 2024).

3.2 Agentic Task Construction

We construct trajectories via forward simulated agent-human interplay, which allows us to fully simulate the environment, the user, and the agent. The critical step is to synthesize agentic tasks that elicit human tool usage while ensuring that the resulting trajectories remain **verifiable**. Concretely, we first initialize an environment state based on the domain-specific database schema, while encouraging as much diversity as possible in the initial state. Next, we sample logically coherent tool sequences from the domain's tool graph, specifically by constructing a directed dependency graph over APIs and traversing it to obtain valid sequences. Starting from a randomly selected initial node, we conduct a directed walk until either the maximum execution steps are reached or a node with no outgoing edges is encountered. This process yields a logically coherent tool sequence. For each step, we generate the corresponding arguments and perform the actual tool call, grounding the operations directly on the database and continuously tracking the evolving database state. This procedure enables verifiability at two complementary granularities: (i) database-level state consistency and (ii) exact matching of tool sequences.

4 Agent Experience Learning

We leverage user intent to drive interactions that yield agent experiences, and train the model through a two-phase process.

4.1 Experience Collection via Interplay

Interplay Motivated by (Yao et al., 2024), once we have constructed an agentic task, we proceed to perform human-agent interplay in the environment. Specifically, we instantiate a simulated user tasked with fulfilling a given overall intent. The agent then leverages domain-specific tools to address the user’s needs, continuing the interaction until the simulated user deems the task complete. This setup enables us to conduct **end-to-end simulation**, encompassing user simulation, agent, and environment, yielding a highly scalable framework. Each completed interaction trace constitutes an agent experience, which can subsequently be used for training. Importantly, since we possess both the gold tool sequences and arguments for the overall intent and the final environment state, we can apply these as supervision signals for experience filtering. All simulated trajectories are generated using the open-source model Qwen3-235B-A22B-Thinking.

Filtering We adopt a three-stage funnel-based trajectory filtering framework consisting of *validity control*, *environment state alignment*, and *function calling exact match*.

- *Validity control*, removes invalid trajectories to ensure well-formed alternating user assistant exchanges. Additionally, we apply an n -gram-based filtering procedure to eliminate severely repetitive reasoning segments. In such cases, we discard these data points.
- *Environment state alignment* retains only those trajectories whose final database state matches the golden state after the interplay, thereby validating the effectiveness of write operations. The filtering granularity at this stage is the **database/environment level**.
- *Function calling exact match* serves as the most stringent filtering stage, where the granularity is the **tool sequence**. Since a tool sequence consisting entirely of read operations without any write operations would cause state-based filtering to fail, we adopt a stricter exact match approach for filtering in such cases. A trajectory is preserved only if the

sequence of invoked tools and arguments exactly matches the overall intent, ensuring high-fidelity supervision.

It is worth noting that we do not filter out trajectories in which tool calls return errors. Thanks to the aforementioned filtering framework, such trajectories may still accomplish the intended goal despite intermediate failures. Retaining them in the training data helps improve the robustness of the model.

4.2 Agentic Experience Learning

Agentic Fine-tuning Given collected or filtered agent-human interplay experience trajectory $\mathcal{H} = (h_0, a_1, \dots, a_{n-1}, h_n, a_0)$, where each human instruction is denoted by h_t at t -round interaction, and each assistant turn a_t is decomposed as $a_t = (\tau_t, \rho_t, y_t)$. Here, τ_t represents the function call tokens, ρ_t the tool response tokens, and y_t the assistant response tokens. Our training objective is to optimize only the tool calls and assistant responses, while human instructions h_i and tool responses ρ_t are excluded from the loss. Formally, given an LLM $p_\theta(x_k | x_{<k})$, we define the loss as

$$\mathcal{L}(\theta) = -\frac{\sum_{k=1}^{|\mathcal{H}|} \mathbb{I}[x_k \in \mathcal{T}] \cdot \log \pi_\theta(x_k | x_{<k})}{\sum_{k=1}^{|\mathcal{H}|} \mathbb{I}[k_i \in \mathcal{T}]}, \quad (2)$$

where x_k denotes the k -th token in the trajectory, π_θ is the model distribution, $\mathbb{I}[\cdot]$ is the indicator function, \mathcal{T} is the set of tokens belonging to tool calls τ or assistant responses y . In practice, all tokens in ρ_i and h_i are masked out from supervision but remain visible in the context $x_{<k}$. This ensures that the model conditions on tool responses and human instructions, while gradients are only propagated through assistant-generated tool calls and natural-language responses.

Two-stage Experience Learning In the first phase, the agent is trained to acquire fundamental skills for tool usage and user interaction. We focus on general domains where a broad set of tools and tasks are available, allowing the agent to develop a robust understanding of when and how to invoke function calls, as well as how to integrate tool outputs into coherent user-facing responses. This stage emphasizes breadth and generality, ensuring that the agent builds a versatile foundation of agentic behaviors before domain-specific specialization. In the second phase, the agent undergoes fine-grained training in vertical domains, where tasks, tools, and

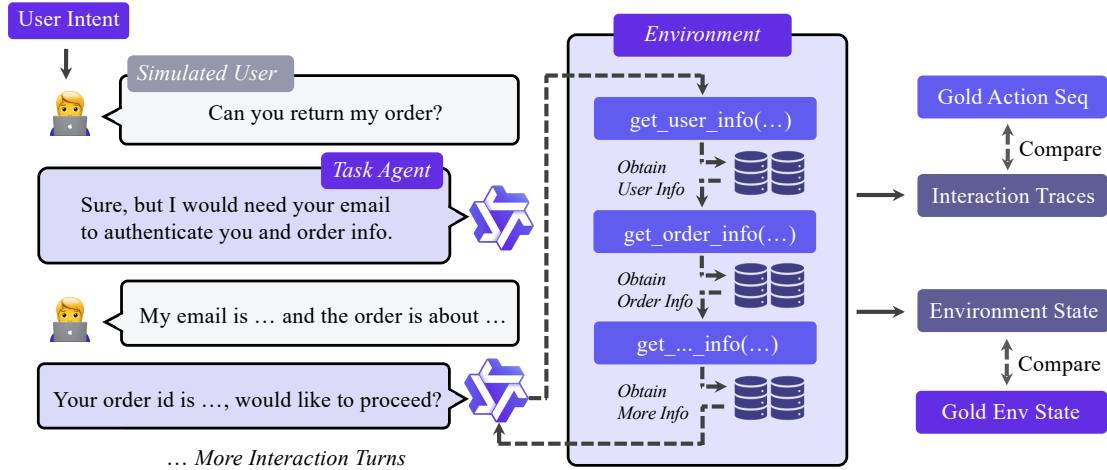


Figure 2: The agent interacts with the simulated user and changes the environment state through tool calling.

435 user intents exhibit domain-specific characteristics.
 436 In our setting, this stage primarily focuses on the
 437 τ -Bench and τ^2 -Bench. By grounding the learning
 438 process in realistic scenarios within a target domain,
 439 the agent refines its ability to select tools, param-
 440 eterize calls, and produce responses that are accu-
 441 rate, contextually appropriate, and aligned with
 442 domain-specific goals. This specialization ensures
 443 a smoother adaptation of agentic capabilities, en-
 444 abling the agent to operate effectively in real-world,
 445 task-oriented environments.

446 5 Experiments

447 5.1 Setup

448 **Benchmarks** We evaluate our methods on three es-
 449 tablished agentic benchmarks: τ -bench, τ^2 -Bench,
 450 and ACEBench-en. For τ -Bench and τ^2 -Bench,
 451 we adopt the pass¹ metric for evaluation and addi-
 452 tionally analyze the trend of pass^k, following the
 453 protocols in (Yao et al., 2024; Barres et al., 2025).
 454 For ACEBench-en, we use the accuracy metric.

455 **Baselines** We compare AgentScaler against:
 456 (i) *closed-sourced LLMs*, including Gemini-2.5-
 457 pro (Comanici et al., 2025), Claude-Sonnet-4 (An-
 458 thropic, 2025), GPT-o3 (OpenAI, 2025b), GPT-
 459 o4-mini (OpenAI, 2025b), and GPT-5-think (Ope-
 460 nAI, 2025a); (ii) *open-sourced LLMs*: GPT-OSS-
 461 120B-A5B (Agarwal et al., 2025), Deepseek-
 462 V3.1-671B-A37B (DeepSeek-AI, 2024), Kimi-K2-
 463 1T-A32B (Team et al., 2025), Qwen3-Thinking-
 464 235B-A22B (Team, 2025b), Seed-OSS-36B (Team,
 465 2025a), Qwen-Coder-30B-A3B (Hui et al., 2024),
 466 and xLAM-2 (Prabhakar et al., 2025a).

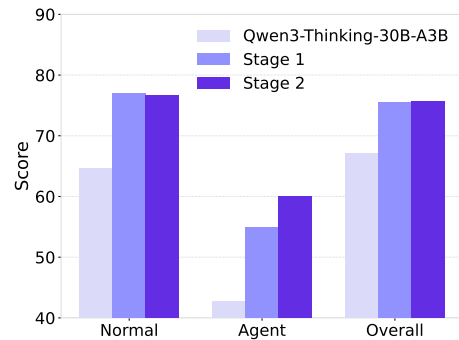


Figure 3: Performance comparison on the Normal, Agent, and Overall subsets of ACEBench-en for two-stage training models.

467 5.2 Experimental Results

468 **Main Results** From Table 1, we observe that
 469 closed-source large language models (LLMs) still
 470 maintain a clear performance advantage, consis-
 471 tently achieving the highest scores across most
 472 domains and benchmarks. This demonstrates
 473 the strength of industrial-scale training pipelines
 474 and proprietary optimization strategies. Neverthe-
 475 less, our proposed AgentScaler achieves a remark-
 476 able level of performance given its lightweight
 477 parameter scale. Specifically, it surpasses most
 478 open-source baselines with fewer than 1T param-
 479 eters, establishing a new state-of-the-art across τ -
 480 bench, τ^2 -Bench, and ACEBench-en. Notably,
 481 **AgentScaler-4B** achieves performance on par with
 482 30B-parameter models despite using the fewest
 483 parameters, highlighting the agentic potential of
 484 compact LLMs. Moreover, **AgentScaler-30B-A3B**
 485 delivers results that are comparable to trillion-
 486 parameter open-source models and, in several do-

Model	τ -bench		τ^2 -Bench			ACEBench-en			
	Retail	Airline	Retail	Airline	Telecom	Normal	Special	Agent	Overall
<i>Closed-Source Large Language Models</i>									
Gemini-2.5-pro	68.7	44.0	67.5	56.0	27.2	76.7	90.0	63.4	78.2
Claude-Sonnet-4	73.9	40.0	67.5	54.0	47.4	79.9	87.3	42.5	76.1
GPT-o3	70.4	52.0	80.2	64.8	58.2	78.3	86.7	63.3	78.2
GPT-o4-mini	70.4	46.0	70.2	56.0	46.5	79.9	84.0	60.0	77.9
GPT-5-think	78.3	44.0	81.1	62.6	96.7	76.7	85.3	32.5	72.2
<i>Open-Source Large Language Models</i>									
GPT-OSS-120B-A5B	67.8	49.2	57.0	38.0	45.6	79.1	84.0	50.8	76.0
Deepseek-V3.1-671B-A37B	66.1	40.0	64.9	46.0	38.5	80.3	62.0	40.8	69.3
Kimi-K2-1T-A32B	73.9	51.2	70.6	56.5	65.8	78.9	81.3	65.0	77.4
Qwen3-Thinking-235B-A22B	67.8	46.0	71.9	58.0	45.6	72.1	84.0	39.1	70.2
Seed-OSS-36B	70.4	46.0	68.4	52.0	41.2	79.1	82.0	58.4	76.7
Qwen-Coder-30B-A3B	68.7	48.0	60.5	42.0	30.7	74.0	41.3	24.1	57.5
xLAM-2-8B-fc-r	58.2	35.2	55.3	48.0	11.4	58.8	0.0	5.0	34.8
xLAM-2-32B-fc-r	64.3	45.0	55.3	52.0	16.7	69.2	24.7	13.4	52.5
xLAM-2-70B-fc-r	67.1	45.2	61.4	56.0	14.0	57.1	5.3	38.4	36.5
Qwen3-Thinking-4B	59.1	52.5	56.1	52.0	28.7	43.3	84.7	11.7	49.5
Qwen3-8B	45.2	25.0	41.2	30.5	23.5	71.4	75.3	29.1	65.9
Qwen3-14B	45.7	31.0	48.0	30.0	26.9	66.9	84.0	44.2	68.0
Qwen3-Thinking-30B-A3B	67.8	48.0	58.8	58.0	26.3	64.7	86.7	42.8	67.2
AgentScaler-4B	64.3	54.0	62.3	56.0	48.2	70.3	76.7	30.8	65.9
AgentScaler-8B	50.4	42.0	58.8	44.0	45.4	69.2	76.7	44.2	67.4
AgentScaler-30B-A3B	70.4	54.0	70.2	60.0	55.3	76.7	82.7	60.0	75.7

Table 1: **Main results** on τ -Bench, τ^2 -Bench, and ACEBench-en.

487 mains, approach those of closed-source counter-
488 parts. These findings highlight the efficiency of our
489 approach: agentic capabilities can be effectively
490 learned and deployed even in relatively compact
491 models, enabling competitive performance without
492 relying on massive parameter counts. This advantage
493 makes AgentScaler particularly well-suited
494 for practical deployment in resource-constrained or
495 latency-sensitive scenarios.

496 **Ablation Study** We further conduct an ablation
497 analysis to examine the effect of the proposed
498 two-stage agent experience learning framework on
499 ACEBench-en. As shown in Figure 3, both Stage 1
500 and Stage 2 training substantially improve performance
501 over the base model (Qwen3-Thinking-30B-A3B)
502 across all subsets. And through multi-steps
503 agent training in Stage 2, the model’s score on
504 the agent set has further improved, and the overall
505 score has also increased. These results validate
506 the design of the two-phase training pipeline: general
507 foundation learning is critical for establishing
508 tool-usage competence, and subsequent domain-
509 specialization further consolidates and contextual-
510 izes these capabilities.

Model	ACEBench-zh	
	Agent	Overall
Qwen3-Thinking-4B	6.7	43.9
AgentScaler-4B	38.4 ^{+31.7}	65.6 ^{+21.7}
Qwen3-8B	35.0	71.3
AgentScaler-8B	58.4 ^{+23.4}	73.7 ^{+2.4}
Qwen3-Thinking-30B-A3B	55.8	74.2
AgentScaler-30B-A3B	64.1 ^{+8.3}	81.5 ^{+7.3}

Table 2: The **results** on ACEBench-zh.

6 Analysis

511 **Our synthetic data approach enables efficient**
512 **knowledge transfer and strong robustness and**
513 **generalization.** We further evaluate our models
514 on ACEBench-zh, which represents an out-of-
515 distribution (OOD) scenario. As shown in Table 2,
516 the AgentScaler models consistently outperform
517 their Qwen baselines across all scales in terms of
518 overall score. In particular, AgentScaler-30B-A3B
519 achieves the best overall score of 81.5, demonstrat-
520 ing strong improvements in the Agent subsets. Not-
521 ably, the small Qwen3-4B model demonstrated
522 a remarkable improvement in agentic capabilities
523 after the two-stage training. Our evaluation setup
524

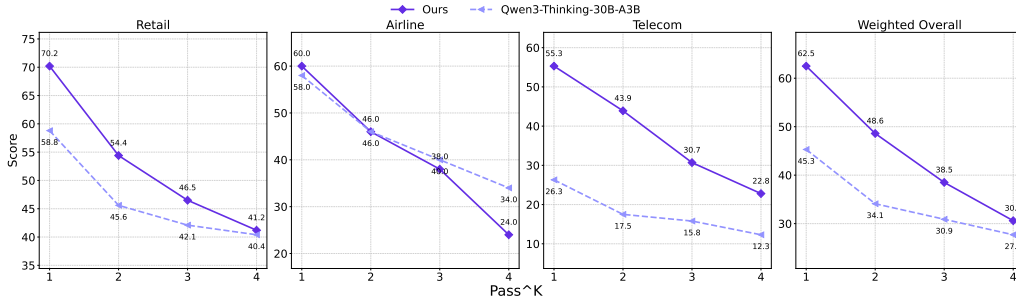


Figure 4: Pass^k metric results across all domains in the τ^2 -Bench.

525 does include domain- and format-level OOD generalization, not only cross-lingual robustness. **First**, ACEBench-en can be also seen an OOD evaluation for our system. The environments we construct use APIs sourced from ToolBench and API-Gen, which are not overlapping with the tool domains or schema structures in ACEBench-en. Therefore, ACEBench-en evaluates the model’s ability to generalize to unseen tool domains, rather than only testing in-domain performance. **Second**, the tool-calling format itself is out-of-distribution. Our training is performed in the Qwen3-Hermes tool-calling format, while ACEBench adopts its own custom parser format. Achieving strong performance under the ACEBench-en parsing and tool-calling rules demonstrates format-level generalization, showing that the model is not overfitted to a single schema or interaction protocol. **Third**, ACEBench-zh provides an additional cross-lingual generalization test, further validating robustness across languages.

546 **AgentScaler shows the strong consistency, stability.** To assess the stability of AgentScaler, Figure 4 reports the pass^k metric on the τ^2 -Bench, which denotes the accuracy achieved when the model correctly answers the same question in all k independent trials. According to the experimental results, the weighted overall score of AgentScaler-30B-A3B consistently surpasses that of Qwen3-Thinking-30B-A3B across all evaluated pass^k settings, indicating a substantial performance advantage of our model over Qwen3-Thinking-30B-A3B. Moreover, a clear downward trend in scores is observed as k increases, suggesting that the stability of existing LLMs remains a considerable challenge.

561 **Scaling Law of Environments.** As shown in Figure 5, we evaluate the model performance under varying environment scales on ACEBench-

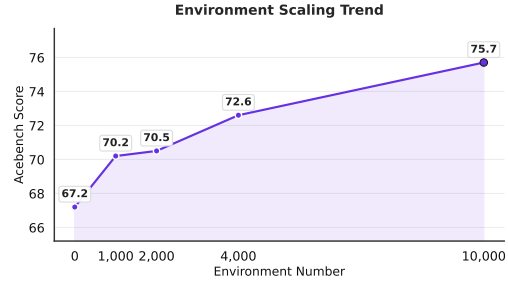


Figure 5: Performance under different environment scales on ACEBench-en.

564 en. The results demonstrate a clear and consistent improvement as the number of environments increases. The largest improvements occur in the low-to-medium scale regime, suggesting that increased environmental diversity substantially enhances tool-use generalization. This trend indicates that the proposed model effectively leverages larger and more diverse environments for learning, exhibiting strong scalability.

573 7 Conclusion

574 In this work, we presented a principled pipeline for advancing general agentic intelligence through systematic environment scaling and agent experience learning. By materializing tools as executable code and grounding them in database-structured environments, our approach enables large-scale construction of verifiable trajectories. Building on these environments, we introduced a two-stage agent experience learning framework that first equips agents with fundamental tool-usage capabilities and then specializes them for domain-specific contexts. Extensive experiments on three representative benchmarks, τ -bench, τ^2 -Bench, and ACEBench, demonstrate the effectiveness of our pipeline. Notably, our AgentScaler achieves state-of-the-art performance among open-source models under 1T parameters, and in several cases reaches parity with much larger or closed-source models.

592 Limitation

593 Although our proposed framework has demon-
594 strated promising results, several limitations re-
595 main, which point to ongoing efforts and potential
596 directions for future work.

597 **Reinforcement-Learning Integration** Although
598 the current system relies solely on two-stage su-
599 pervised fine-tuning, the simulator we have built
600 offers deterministic, low-latency feedback that is
601 ideal for reinforcement-learning optimization. In
602 future iterations we plan to add an RL stage using
603 policy gradient methods, to refine the agent’s long-
604 horizon decision-making and further improve its
605 emergent, agentic capabilities.

606 **Model Scale** Another limitation of our current
607 work lies in the model scale. Our method has so
608 far only been validated on a 30B-scale architec-
609 ture, without extension to larger models exceeding
610 200B or even trillion-parameter scales. While prior
611 work (Belcak et al., 2025) emphasizes that “small
612 language models are the future of agentic AI,” we
613 share the view that training agentic capabilities
614 in relatively smaller models is particularly mean-
615 ingful. Such models are easier to deploy on edge
616 devices, enable broader applicability across diverse
617 scenarios, and offer faster response times.

618 Ethical Considerations

619 This study strictly adheres to established ethi-
620 cal guidelines at every stage. During the tool-
621 collection phase, all code, models, and utilities
622 were obtained exclusively from publicly available,
623 open-source repositories or officially documented
624 APIs released under permissive licenses. No propri-
625 etary or restricted software was employed. Further-
626 more, every data point used in the experiments was
627 synthetically generated through algorithmic means.
628 Crucially, no personally identifiable information
629 was collected, accessed, or produced at any time.

630 References

631 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Alt-
632 man, Andy Applebaum, Edwin Arbus, Rahul K
633 Arora, Yu Bai, Bowen Baker, Haiming Bao, and 1
634 others. 2025. gpt-oss-120b & gpt-oss-20b model
635 card. *arXiv preprint arXiv:2508.10925*.

636 Anthropic. 2025. *System card: Claude opus 4 & claude*
637 *sonnet 4*.

Victor Barres, Honghua Dong, Soham Ray, Xujie Si,
and Karthik Narasimhan. 2025. tau2-bench: Evaluat-
ing conversational agents in a dual-control environ-
ment. *arXiv preprint arXiv:2506.07982*.

Peter Belcak, Greg Heinrich, Shizhe Diao, Yonggan
Fu, Xin Dong, Saurav Muralidharan, Yingyan Celine
Lin, and Pavlo Molchanov. 2025. Small language
models are the future of agentic ai. *arXiv preprint*
arXiv:2506.02153.

Vincent D Blondel, Jean-Loup Guillaume, Renaud
Lambiotte, and Etienne Lefebvre. 2008. Fast un-
folding of communities in large networks. *Jour-
nal of statistical mechanics: theory and experiment*,
2008(10):P10008.

Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang,
Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang,
Weinan Gan, Yuefeng Huang, and 1 others. 2025.
Acebench: Who wins the match point in tool usage?
arXiv preprint arXiv:2501.12851.

Mingyang Chen, Haoze Sun, Tianpeng Li, Fan Yang,
Hao Liang, Keer Lu, Bin Cui, Wentao Zhang, Zenan
Zhou, and Weipeng Chen. 2024a. Facilitating multi-
turn function calling for llms via compositional in-
struction tuning. *arXiv preprint arXiv:2410.12952*.

Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei
Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and
Feng Zhao. 2024b. Agent-flan: Designing data and
methods of effective agent tuning for large language
models. *arXiv preprint arXiv:2403.12881*.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann,
Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-
cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and
1 others. 2025. Gemini 2.5: Pushing the frontier with
advanced reasoning, multimodality, long context, and
next generation agentic capabilities. *arXiv preprint*
arXiv:2507.06261.

DeepSeek-AI. 2024. *Deepseek-v3 technical report*.
Preprint, arXiv:2412.19437.

Runnan Fang, Xiaobin Wang, Yuan Liang, Shuofei
Qiao, Jialong Wu, Zekun Xi, Ningyu Zhang, Yong
Jiang, Pengjun Xie, Fei Huang, and 1 others. 2025.
Synworld: Virtual scenario synthesis for agen-
tic action knowledge refinement. *arXiv preprint*
arXiv:2504.03561.

Xinyu Geng, Peng Xia, Zhen Zhang, Xinyu Wang,
Qiuchen Wang, Ruixue Ding, Chenxi Wang, Jia-
long Wu, Yida Zhao, Kuan Li, and 1 others. 2025.
Webwatcher: Breaking new frontiers of vision-
language deep research agent. *arXiv preprint*
arXiv:2508.05748.

Zhicheng Guo, Sijie Cheng, Yuchen Niu, Hao Wang,
Sicheng Zhou, Wenbing Huang, and Yang Liu. 2025.
Stabletoolbench-mirrorapi: Modeling tool environ-
ments as mirrors of 7,000+ real-world apis. *arXiv*
preprint arXiv:2503.20527.

804	ByteDance Seed Team. 2025a. Seed-oss open-source models. https://github.com/ByteDance-Seed/seed-oss .	Shaokun Zhang, Yi Dong, Jieyu Zhang, Jan Kautz, Bryan Catanzaro, Andrew Tao, Qingyun Wu, Zhiding Yu, and Guilin Liu. 2025. Nemotron-research-tool-1: Exploring tool-using language models with reinforced reasoning. <i>arXiv preprint arXiv:2505.00024</i> .	859
805			860
806			861
807	Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, and 1 others. 2025. Kimi k2: Open agentic intelligence. <i>arXiv preprint arXiv:2507.20534</i> .	Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, and 1 others. 2023. Agents: An open-source framework for autonomous language agents. <i>arXiv preprint arXiv:2309.07870</i> .	862
808			863
809			864
810			865
811			866
812	Qwen Team. 2025b. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388.		867
813			868
814	Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, and 1 others. 2025a. Web-dancer: Towards autonomous information seeking agency. <i>arXiv preprint arXiv:2505.22648</i> .		869
815			
816			
817			
818			
819	Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, and 1 others. 2025b. Web-walker: Benchmarking llms in web traversal. <i>arXiv preprint arXiv:2501.07572</i> .		
820			
821			
822			
823			
824	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. tau-bench: A benchmark for tool-agent-user interaction in real-world domains. <i>arXiv preprint arXiv:2406.12045</i> .		
825			
826			
827			
828	Junjie Ye, Zhengyin Du, Xuesong Yao, Weijian Lin, Yufei Xu, Zehui Chen, Zaiyuan Wang, Sining Zhu, Zhiheng Xi, Siyu Yuan, Tao Gui, Qi Zhang, Xuanjing Huang, and Jiecao Chen. 2025a. ToolHop: A query-driven benchmark for evaluating large language models in multi-hop tool use . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2995–3021, Vienna, Austria. Association for Computational Linguistics.		
829			
830			
831			
832			
833			
834			
835			
836			
837			
838	Junjie Ye, Changhao Jiang, Zhengyin Du, Yufei Xu, Xuesong Yao, Zhiheng Xi, Xiaoran Fan, Qi Zhang, Xuanjing Huang, and Jiecao Chen. 2025b. Feedback-driven tool-use improvements in large language models via automated build environments. <i>arXiv preprint arXiv:2508.08791</i> .		
839			
840			
841			
842			
843			
844	Fan Yin, Zifeng Wang, I Hsu, Jun Yan, Ke Jiang, Yanfei Chen, Jindong Gu, Long T Le, Kai-Wei Chang, Chen-Yu Lee, and 1 others. 2025. Magnet: Multi-turn tool-use data synthesis and distillation via graph translation. <i>arXiv preprint arXiv:2503.07826</i> .		
845			
846			
847			
848			
849	Yirong Zeng, Xiao Ding, Yuxian Wang, Weiwen Liu, Wu Ning, Yutai Hou, Xu Huang, Bing Qin, and Ting Liu. 2025. Boosting tool use of large language models via iterative reinforced fine-tuning. <i>arXiv e-prints</i> , pages arXiv–2501.		
850			
851			
852			
853			
854	Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, and 1 others. 2024. xlam: A family of large action models to empower ai agent systems. <i>arXiv preprint arXiv:2409.03215</i> .		
855			
856			
857			
858			

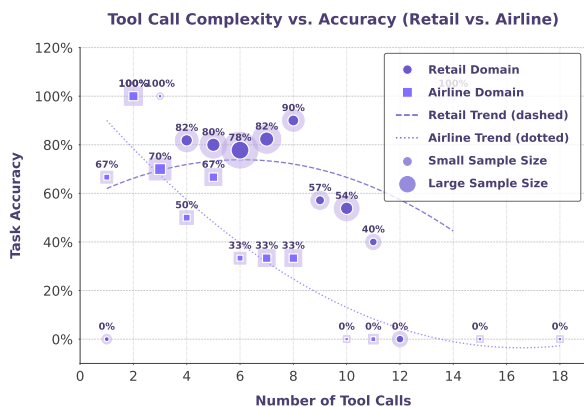


Figure 6: Accuracy by tool call count on τ -bench.

the top-p value to 0.95, and the top-k value to 20. Moreover, to speed up model inference, we deploy the model using vLLM (Kwon et al., 2023).

C Use of AI Assistant

We affirm that Large Language Models are employed solely as an assisted tool to refine wording and sentence structure during our paper writing process. Their use in the experiments is strictly for scientific research purposes, and all such usage has been explicitly documented in our Experimental Settings and Reproducibility Statement. No other reliance on LLMs is involved in this work.

A Further Analysis

Long-horizon tool calling remains a fundamental challenge for agentic models. To further analyze the model’s long-horizon tool-calling capability, we constructed a scatter plot on the τ -bench dataset showing the relationship between the number of tool calls in each trajectory and the corresponding trajectory accuracy, with a dashed line indicating the trend. As illustrated in Figure 6, there exists a clear negative correlation between the number of tool calls and task accuracy. Our AgentScaler models exhibit this trend as well, underscoring that handling extended tool-use chains is still an open problem that we plan to address in future work.

B Experiment Settings

Backbones We train the AgentScaler model series by training on Qwen3 models (Team, 2025b) of varying scales. Specifically, AgentScaler-4B and AgentScaler-30B-A3B are trained on Qwen3-Thinking-4B-2507 and Qwen3-Thinking-30B-A3B-2507, respectively, while AgentScaler-8B is trained on Qwen3-8B. During the agentic experience learning stage, all models in the AgentScaler series, including 4B, 8B, and 30B-A3B, are trained for three epochs with a batch size of 128. The model context length is set to 32768. Subsequently, we directly use the model from the final checkpoint. Additionally, the learning rate is set to $7e-6$, and a warm-up strategy was employed.

During the evaluation stage, we strictly follow all the official guidelines of the benchmarks. All baseline models are assessed using the benchmarks’ default configurations. For the AgentScaler series of models, we set the temperature parameter to 0.6,