TOUCAN: SYNTHESIZING 1.5M TOOL-AGENTIC DATA FROM REAL-WORLD MCP ENVIRONMENTS

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

033

034

037

038

040

041

042

043

044

046

047

048

049

051

052

Paper under double-blind review

ABSTRACT

Large Language Model (LLM) agents are rapidly emerging as powerful systems for automating tasks across domains. Yet progress in the open-source community is constrained by the lack of high quality permissively licensed tool-agentic training data. Existing datasets are often limited in diversity, realism, and complexity, particularly regarding multi-tool and multi-turn interactions. To address this gap, we introduce TOUCAN, the largest publicly available tool-agentic dataset to date, containing 1.5 million trajectories synthesized from nearly 500 real-world Model Context Protocols (MCPs). Unlike prior work, TOUCAN leverages authentic MCP environments to generate diverse, realistic, and challenging tasks with trajectories involving real tool execution. Our pipeline first produces a broad spectrum of tool-use queries using five distinct models, applies model-based quality filtering, and then generates agentic trajectories with three teacher models using two agentic frameworks. Rigorous rule-based and model-based validation ensures high-quality outputs. We also introduce three extension mechanisms to further diversify tasks and simulate multi-turn conversations. Models fine-tuned on TOU-CAN outperform larger closed-source counterparts on the BFCL V3 benchmark and establish a new Pareto optimum on MCP-Universe Bench.

1 Introduction

Large language models (LLMs) have become integral to AI applications, with LLM agents emerging as powerful systems for automating complex tasks across diverse domains Li et al. (2024). There is growing excitement about the potential of LLM agents to unlock new levels of automation across industries (Ferrag et al., 2025; Bousetouane, 2025). These agents handle multi-step workflows that require discovering the right tools from potentially large toolsets, calling them correctly with appropriate parameters, handle tool failures gracefully, and synthesizing results into accurate, context-aware responses Xu et al. (2025a). Recent advancements, such as the Model Context Protocol (MCP) (Anthropic, 2025), have streamlined tool integration by providing standardized interfaces, enabling seamless connections between LLMs and real-world environments and simplifying the process for LLM agents to discover, invoke, and execute external tools.

Despite these advancements, progress in the open-source community is constrained by the lack of high-quality, permissively licensed **tool-agentic data** for training more capable agentic LLMs. An instance of tool-agentic data comprises a task-trajectory pair, where trajectories capture sequences of planning, tool calls, tool responses, and the final model response. While previous efforts (Qin et al., 2023; Liu et al., 2024; 2025a; Prabhakar et al., 2025) have introduced datasets covering various tool-calling scenarios, they suffer from several limitations: restricted tool diversity, lack of authentic tool responses, focus on single-turn conversations between users and models, or insufficient scale, all of which constrain effective training of agentic capabilities. There is an urgent need for comprehensive, high-quality datasets that capture the full spectrum of tool-agentic interactions observed in production environments.

In this work, we bridge this gap by introducing TOUCAN, the largest publicly available tool-agentic dataset to date, comprising 1.5 million trajectories synthesized from nearly 500 real-world MCP servers. Unlike prior approaches that rely on simulated or limited toolsets, TOUCAN leverages authentic MCP environments with more than 2,000 tools to generate diverse, realistic, and challenging tasks spanning parallel and multi-step tool calls, as well as multi-turn conversations. Our pipeline

Table 1: TOUCAN comparison to open-source tool-agentic datasets. Comparison comprises total trajectories, tool calling scenarios ([S]ingle, [P]arallel, [M]ulti[S]tep) including no-tool-use edge case (irrelevance[IR]), number of multi-turn conversations, and other details about data generation. Note — indicates information not publicly available.

Dataset	Trajectories	Tool-Call Scenarios	Multi Turn	Tool Specs	Tool Response
APIGent-MT-5K (Prabhakar et al., 2025)	5,000	S P MS IR	5,000	From $ au$ -Bench	Executed
ToolACE (Liu et al., 2025a)	11,300	S P MS IR	509	Synthetic	Simulated
Hermes Function-Calling V1 (interstellarninja)	11,570	S P MS IR	1,890	Synthetic	Executed
Nemotron (Tools) (Nathawani et al., 2025)	310,051	S P MS –	199,610	-	-
TOUCAN (This Work)	1,527,259	S P MS IR	567,262	Real	Executed

begins by producing a broad spectrum of tool-use tasks using five distinct models with MCP server specifications, followed by model-based quality filtering to ensure relevance and difficulty. We then generate agentic trajectories with three teacher models, incorporating rigorous rule-based and model-based checks for high-quality outputs, including verification of tool execution and response accuracy. Our pipeline also integrates extensions to generate additional tasks targeting edge case scenarios, interactive conversations, and multi-turn dialogues.

Our experiments demonstrate the effectiveness of TOUCAN in enhancing LLM agentic capabilities. Models fine-tuned on TOUCAN surpass closed-source counterparts on the BFCL V3 benchmark (Patil et al., 2025), achieving superior performance in function calling accuracy across single-turn and multi-turn scenarios. Furthermore, they show substantial improvements on τ -Bench (Yao et al., 2024) and τ^2 -Bench (Barres et al., 2025), with gains in tool selection, execution fidelity, and multi-turn reasoning under dynamic user interactions. On the recent MCP-Universe benchmark (Luo et al., 2025), which evaluates LLMs on 231 realistic tasks using 11 real-world MCP servers, TOUCAN-tuned models achieve state-of-the-art performance within their parameter class, consistently outperforming leading models of comparable size. In summary, the contributions of our work are:

- TOUCAN Dataset. The largest open-source tool-agent training dataset, covering parallel and multi-step tool calls, multi-turn dialogues, and edge-case tool use. Recent reports on frontier LLM development, such as Kimi-K2 (Team et al., 2025b) and GLM-4.5 (Team et al., 2025a), highlight the value of large-scale trajectories with broad domain coverage, and TOUCAN provides an open-source alternative that bridges this gap.
- TOUCAN Pipeline. A pipeline that leverages any MCP specifications to generate diverse toolagent trajectories, supports tool execution through MCP servers, and can be seamlessly extended to new tools via the MCP standard.
- TOUCAN Checkpoints. Our experiments demonstrate that models fine-tuned on TOUCAN mixtures surpass closed-source counterparts on the BFCL V3 and MCP-Universe benchmarks.

2 Related Work

The past: Tool-calling datasets and benchmarks for LLMs. Early tool-calling datasets enabled LLMs to interact with tools like REST APIs and ML functions. The Gorilla project (Patil et al., 2023) demonstrated that fine-tuning on such data enhances tool-use over vanilla models, introducing the BFCL benchmark (Patil et al., 2025) as a standard. ToolAlpaca (Tang et al., 2023) offered cost-effective synthetic data with lower quality, while ToolLLM (Qin et al., 2023) expanded to 16,000+APIs across domains. API Pack (Guo et al., 2025a) added cross-language diversity (Python, Java, C++), and API Blend (Basu et al., 2024) optimized dataset mixtures for robustness, laying the foundation for tool-agent advancements. More recently, APIGen has focused on domain diversification, contributing a training dataset covering 21 domains Liu et al. (2024).

The present: Tool-calling benchmarks and datasets for LLM-agents. Recent research has shifted toward training LLM agents for effective tool use, exemplified by models like Kimi-K2 (Team et al., 2025b) and GLM-4.5 (Team et al., 2025a), with performance assessed via benchmarks such as BFCL (Patil et al., 2025), τ -Bench (Yao et al., 2024), and ACEBench (Chen et al., 2025). BFCL covers diverse scenarios including parallel, multi-step, and multi-turn tool use, while τ -Bench

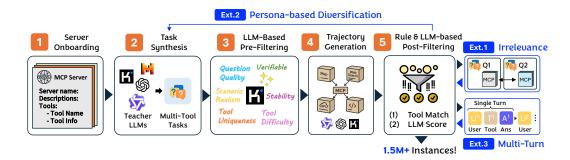


Figure 2: The TOUCAN construction pipeline: A systematic five-stage process from MCP server on-boarding through trajectory filtering, with three extensions for enhancing data diversity and realism.

focuses on realistic user-agent-tool interactions. ACEBench enhances evaluation by addressing edge cases and including a subset for tool-agent trajectories. Despite these advances, open-source training data for tool-agent trajectories remains limited. Existing datasets (interstellarninja; Liu et al., 2025a; Prabhakar et al., 2025; Nathawani et al., 2025) either lack dataset curation transparency, are small in size for SFT, or simulate tool responses via LLMs. Table 1 compares existing tool-agentic datasets with TOUCAN, which, at 1.5 million trajectories, offers the largest dataset, featuring extensive multi-turn dialogues, all tool-use scenarios, critical edge cases, and authentic tool responses from real-world environments.

The future: MCP benchmarks and datasets. As concurrent work, recent MCP benchmarks (Gao et al., 2025; Wang et al., 2025; Luo et al., 2025; Team, 2025a; Guo et al., 2025b; Yin et al., 2025; Liu et al., 2025b; Yan et al., 2025; Team, 2025b) aim to rigorously assess LLMs in tool-use settings beyond simple correctness. For instance, MCP-Radar (Gao et al., 2025) employs a five-dimensional evaluation including accuracy, tool selection efficiency, resource usage, parameter construction, and execution speed across software engineering, math, and problem-solving tasks with 300 queries and 42 MCP servers. Similarly, MCP-Bench (Wang et al., 2025) evaluates multi-step reasoning over 28 MCP servers and 250 tools, while MCP-Universe (Luo et al., 2025) focuses on execution-based metrics in six real-world domains. These advancements underscore the need for comprehensive training datasets to support the development of robust, open-source LLM agents.

3 TOUCAN: SCALING TOOL-AGENTIC DATA WITH REAL WORLD MCPS

3.1 TOUCAN GENERATION PIPELINE

TOUCAN is a comprehensive dataset comprising over 1.5 million tool-agent trajectories constructed using real-world tools from MCP servers. Each instance in our dataset contains a task description, a complete agent trajectory with its associated tools, quality and classification annotations, as well as comprehensive metadata. Appendix A provides a detailed schema description and demonstration samples. The construction of TOUCAN follows a systematic five-stage pipeline: MCP server onboarding, task synthesis, task filtering, trajectory generation, and trajectory filtering. Additionally, we implement three extension mechanisms to further enhance data diversity and realism. Figure 2 illustrates the complete construction pipeline. We detail each stage below.

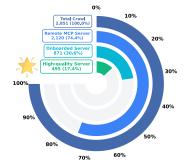


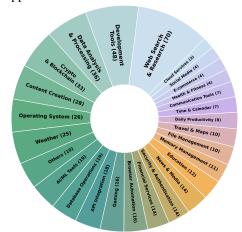
Figure 1: MCP servers filtering process

Stage 1: MCP Server Onboarding. To generate questions from diverse environments, the initial step involves onboarding as many high-quality MCP servers as possible. We sourced MCP server specification files from GitHub and Smithery ¹, a platform and registry for MCP servers that encapsulate modular execution environments. Each MCP server is accompanied by a structured JSON document detailing metadata about the server with a machine-readable definition of the tools it provides. From an initial crawl yielding approximately 2,800 MCP servers, we applied two key filtering

¹https://smithery.ai/

criteria: (1) retaining only remote MCP servers accessible via streamable HTTP to ensure compatibility with trajectory generation, and (2) excluding servers requiring third-party credentials (e.g., API keys) for tool invocation to maintain accessibility and reproducibility. This process reduced the dataset to 30.6% (871 servers). As a final step, we generated a small subset of test questions to evaluate each tool within the MCP servers, subsequently filtering out servers with problematic tools that returned error messages or failed to function correctly. This rigorous curation process resulted in a refined set of 495 high-quality MCP servers spanning diverse domains and functionalities. Figure 1 depicts the number of MCP servers retained at each filtering stage. Figure 3 demonstrates the domain distribution of the final server collection across diverse categories. The domain distribution is annotated by LLMs, where prompts can be found in Appendix D.1.

Stage 2: Task Synthesis. The next step involves synthesizing high-quality tasks from MCP servers, where each task comprises a question and the desired tool names from the MCP servers. The key challenge is ensuring that tasks are challenging, realistic, and cover edge cases. Therefore, we design diverse sampling strategies based on MCP server usage number from Smithery and server functionalities. To avoid potential bias from individual models, we utilized five open-source LLMs (Mistral-Small, DevStral-Small, GPT-OSS, Kimi-K2, and Qwen3-32B) as task generators to construct synthetic tasks (see the prompts in Appendix D.2). We apply the following three strategies to synthesize tasks, where the maximum number of tools is set to N=3 in our experiments:



Single Server: For a given MCP server, we synthesize tasks requiring the use of 1 to N tools, ensuring a balanced selection distribution guided by server usage statistics to reflect real-world applicability.

Figure 3: MCP servers distribution by domain, covering a wide range of categories. Values in parentheses indicate the number of servers belonging to each category.

Multi-Server: Leveraging LLM-based domain annotations derived from MCP metadata, we first sample N MCP servers from either the same or different categories. We then prompt LLMs to conduct a server analysis, outlining potential workflows that integrate tools across these servers, targeting two to N specific tools, and subsequently generating tasks that leverage functionalities from multiple servers.

Featured Server: Based on the original MCP file metadata, we manually selected 25 representative MCP servers from various domains, with the complete list available in Appendix B.1. In this approach, we provide all MCP server metadata within the context, specify an expected number of tools, and allow the LLM to freely explore combinations, devise realistic scenarios, select the necessary tools, and create comprehensive tasks.

Stage 3: Task Filtering. To ensure the quality of synthesized tasks, this stage involves annotating tasks across six dimensions and filter out suboptimal instances. We employed the Kimi-K2 model as the annotator, which was selected for its optimal balance between correlation with human annotations and cost efficiency. The correlation statistics are detailed in Appendix C.1, and the prompt template is provided in Appendix D.4. Each dimension is rated on a 1-5 Likert scale. The detailed evaluation metrics are as follows:

- Tool Selection Difficulty: Judges the difficulty of selecting the required tools from provided tools.
- *Tool Selection Uniqueness:* Assesses the uniqueness of the selected tool combination relative to the available tools, and whether viable alternatives could also solve the task.
- Question Quality: The task's overall quality, reflected by its clarity, specificity, and effectiveness.
- Scenario Realism: Evaluates the authenticity and realism of the task scenario.
- Verifiable: Evaluates how easily the final model answer can be verified given the question.
- Stability: Evaluates whether tool outputs remain consistent over time, across geolocation, and under stochastic variation.

Stage 4: Trajectory Generation. This step involves collecting trajectories including tool calls, tool responses, and reasoning steps in agentic environments given tasks synthesized and filtered from the previous steps. To ensure diversity, we employed three LLMs from different families (GPT-OSS-120B, Kimi-K2, and Qwen3-32B) in combination with two agent frameworks (Qwen-agent and OpenAI-agent) to produce high-quality agentic trajectories. The models are deployed remotely and accessed by the agent frameworks via streamable HTTP.

Stage 5: Rule&LLM-Based Post-Filtering. The trajectory filtering process combines rule-based verifiers with LLM-driven annotations to ensure high quality. Rule-based heuristics exclude trajectories that fail to start the agent or connect successfully with remote MCP servers, do not contain tool calls, exhibit failures in tool responses, or contain local file system paths. We also validate whether the trajectory uses the required tools specified by the task in the correct sequence, and report both the *desired tool use percentage* (coverage of required tools) and *order correctness* (adherence to expected sequence) metrics. We then employ GPT-OSS-120B as a judge to annotate each trajectory in terms of completeness and conciseness. The annotation prompt is provided in Appendix D.5, with metric definitions as follows:

- Completeness: Judges whether the assistant fulfills the user's request end-to-end.
- Conciseness: Judges whether the task is solved with the minimum necessary steps and verbosity.

This dual-stage filtering approach ensures that only high-quality, concise, and executable trajectories are retained in the final dataset.

3.2 TOUCAN EXTENSIONS

While the core pipeline generates high-quality trajectories, these are single-turn interactions between user and agent without follow-ups, which limits their practical applicability to real-world scenarios. In addition, since all available tools are contextually relevant, tool selection becomes trivial for LLMs, resulting in relatively low difficulty. To address these limitations and enhance the dataset's versatility, we apply three distinct procedures post-core pipeline (Steps 1 to 5) to generate new instances targeting specific objectives.

Ext.1: Irrelevance. To reduce hallucination, it is critical to train models to reject unanswerable queries or seek alternative solutions when desired tools are unavailable. To achieve this, we systematically generate queries unsolvable with the current toolset (Ext1 in Figure 2) by shuffling MCP server metadata across instances and repeating the task generation step.

Ext.2: Persona-based Diversification. We implement persona-based diversification (Ext2 in Figure 2) to create varied task versions. This involves two strategies: one enhances diversity by introducing new contexts and personas, while the other increases task complexity through additional constraints, all while utilizing the same target tools. This diversification process produces tasks similar yet distinct from those in the core pipeline. The prompts are detailed in Appendix D.3.

Ext.3: Multi-Turn. Recognizing that real-world user-agent-tool interactions seldom conform to single-turn conversations Yao et al. (2024), we introduce a self-simulation pipeline to generate multi-turn dialogues using the trajectory generation model. This is achieved through two methods: (1) splitting complex tasks requiring multi-tool coordination into sequential sub-questions, and (2) extending existing conversations by providing LLMs with context to formulate follow-up queries.

Finally, we repeat the core pipeline from steps 2 to 5 to build full trajectories with the new tasks. In the case of irrelevant tasks (Ext.1), we tighten trajectory filters to retain only instances with zero tool calls. Together, these data extensions yield a more realistic and robust TOUCAN dataset that covers all relevant tool-use scenarios and user question styles.

3.3 DATA ANALYSIS

This section analyzes the generated TOUCAN dataset from statistical analysis and LLM-based quality assessment.

Statistical Analysis of TOUCAN. We conduct comprehensive statistical analysis of MCP servers and data instances. The top MCP servers used in TOUCAN and tool statistics within each MCP servers are presented in Appendix B.2. Figure 4 provides a comprehensive analysis of the TOU-

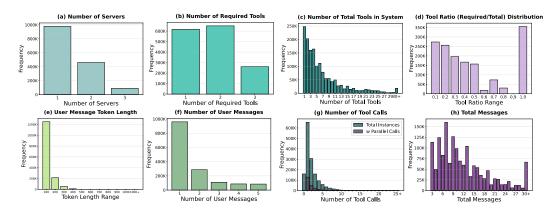


Figure 4: The figures above illustrate the TOUCAN dataset analysis. Subfigure (a) and (b) provide statistics on the number of servers and required tools per instance, highlighting TOUCAN 's comprehensive coverage of multi-server and multi-tool tasks. Subfigures (c) and (d) reveal that most tasks include more tools in the context than the targeted tools, underscoring the non-trivial tool selection challenges. Subfigure (e) displays the length of user messages in tokens. Subfigures (f) and (h) demonstrate the multi-turn nature of the tasks, characterized by extended and diverse interactions among users, agents, and tools. Subfigure (g) demonstrates that TOUCAN encompasses both single and parallel tool calls, which enhance the dataset's versatility in capturing diverse agent-tool interaction patterns.

CAN dataset. We observe that TOUCAN provides comprehensive coverage of multi-server and multi-tool tasks, and includes multi-turn conversations among users, agents, and tools. Additionally, most tasks contain more tools in the context than the required target tools, indicating non-trivial tool selection requirements. Figure 5 presents the subset statistics of TOUCAN across different trajectory generator LLMs and data partitions. We also provide embedding visualization of TOUCAN using UMAP projection in Appendix B.3, demonstrating the wide domain coverage of TOUCAN.

Quality Assessment of TOUCAN. Figure 6 presents a statistical analysis conducted by an LLM-as-a-judge on TOUCAN. From the task perspective (labels in ■), we observe that the majority of tasks exhibit exceptionally high question quality and scenario realism, indicating robust task design and alignment with real-world applications. Additionally, the dataset features a mixed difficulty range, encompassing both simple and challenging tasks. From the response perspective (label in ■), we find that trajectory quality is satisfactory, with most scores at or above 3 (medium) across both completeness and conciseness metrics.

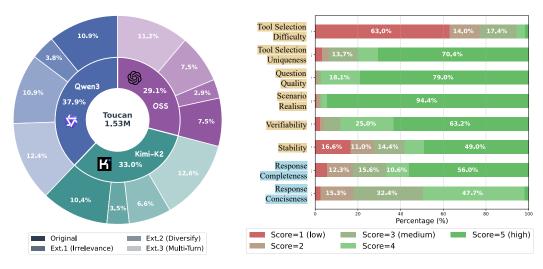


Figure 5: TOUCAN Subset Statistics

Figure 6: TOUCAN Quality Statistics

Table 2: This table compares the performance of TOUCAN -tuned models and baselines on the BFCL-V3 benchmark. We observe that TOUCAN remarkably improves baseline model performance through supervised fine-tuning (SFT) and enables smaller models to outperform larger models across different evaluation aspects.

Model	Overall	Single Turn		Multi Turn	Hallucination	
		Non-live (AST)	Live (AST)		Relevance	Irrelevance
DeepSeek-V3	64.71%	88.54%	77.34%	29.87%	83.33%	76.49%
Qwen2.5-72B-Instruct	64.37%	87.56%	78.68%	29.38%	72.22%	77.41%
Qwen3-235B-A22B	67.94%	87.90%	77.03%	40.12%	83.33%	76.32%
Qwen3-32B	69.25%	88.90%	77.83%	43.12%	72.22%	75.79%
o3-Mini	64.61%	86.15%	79.08%	28.75%	72.22%	82.96%
GPT-4.1	68.69%	85.42%	79.92%	40.50%	77.78%	85.95%
GPT-4.5-Preview	70.32%	86.12%	79.34%	45.38%	66.67%	83.64%
Qwen2.5-7B-Instruct	55.10%	84.19%	72.32%	12.88%	72.22%	67.93%
with TOUCAN	$58.26\%^{+3.16\%}$	78.52%	74.50%	22.62%	66.67%	75.18%
Qwen2.5-14B-Instruct	57.69%	83.38%	73.70%	19.75%	83.33%	68.46%
with TOUCAN	$65.09\%^{+7.40\%}$	85.42%	76.01%	35.25%	72.22%	75.96%
Qwen2.5-32B-Instruct	61.73%	85.58%	76.01%	26.38%	72.22%	72.68%
with TOUCAN	70.45% +8.72%	87.12%	78.90%	46.50%	77.78%	78.10%

Table 3: This table presents τ -Bench and τ^2 -Bench results for models fine-tuned on TOUCAN compared to their respective baselines. Improvements are observed across most evaluation scenarios.

Model	au-bench			$ au^2$ -bench			
	Avg.	Airline	Retail	Avg.	Airline	Retail	Telecom
Qwen2.5-7B-Instruct with TOUCAN	15.03%	8.75%	21.30%	16.08%	14.00%	17.54%	16.70%
	22.48% ^{+7.45} %	15.50%	29.46%	17.77% ^{+1.69%}	20.00%	22.80%	10.50%
Qwen2.5-14B-Instruct with TOUCAN	30.85%	17.25%	44.46%	24.46%	12.00%	41.20%	20.18%
	35.24% ^{+4.39%}	22.00%	48.48%	30.43% ^{+5.97%}	22.00%	49.10%	20.18%
Qwen2.5-32B-Instruct with TOUCAN	38.76%	26.00%	51.52%	29.40%	18.00%	49.10%	21.11%
	42.33% ^{+3.57} %	29.00%	55.65%	31.60% ^{+2.20%}	22.00%	52.60%	20.20%

4 EXPERIMENTS

In this section, we demonstrate the performance of TOUCAN by performing supervised fine-tuning (SFT) on baseline models of different sizes. We then compare the fine-tuned models' performance against existing model baselines across several widely used agentic tool-call benchmarks.

4.1 EXPERIMENT SETUP

Model and Baseline Setup. We perform supervised fine-tuning on <code>Qwen2.5-7B-Instruct</code>, <code>Qwen2.5-14B-Instruct</code>, and <code>Qwen2.5-32B-Instruct</code> (Team, 2024) to demonstrate the efficacy of TOUCAN across models of varying sizes. Detailed fine-tuning parameters are provided in Appendix C.2. We benchmark the performance of our fine-tuned models against models of comparable or larger scales, including <code>DeepSeek-V3</code> <code>DeepSeek-AI</code> et al. (2025), <code>Qwen2.5-72B-Instruct</code>, <code>Qwen3-235B-A22B</code>, <code>Qwen3-32B</code> Yang et al. (2025), and closed-source OpenAI models such as o3-mini, GPT-4.1, and GPT-4.5-Preview.

TOUCAN Setup. Given the large volume of the full dataset, we adopted a strategy similar to Xu et al. (2025b) by sampling from a high-quality subset of TOUCAN. This subset was selected based on the following criteria: question quality and scenario realism scores of 5, response completeness and conciseness scores of at least 4, and desired tool use percentage of 1.0 (indicating that trajectories fully utilize all required tools from the task). We performed necessary data re-balancing to ensure the dataset remains representative across different categories. The resulting SFT dataset comprises

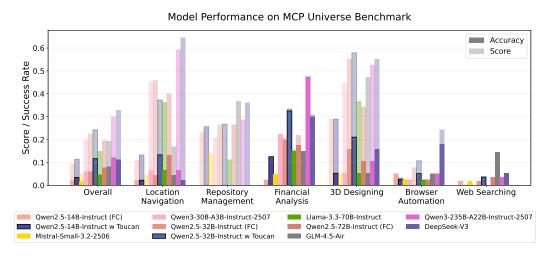


Figure 7: This figure compares the performance of TOUCAN-tuned models with other open-source models on MCP-Universe (Luo et al., 2025). Model sizes increase from left to right. Bars with darker colors represent task success rate (full task completion), while lighter colors represent average evaluation scores considering partial task completion. TOUCAN-tuned models are shown with black borders. TOUCAN-tuned models outperform other models of similar sizes across most tasks.

28.3K instances from the original pipeline, 40K instances from Ext.1 (Irrelevance), 15.8K instances from Ext.2 (Diversify), and 35.2K instances from Ext.3 (Multi-Turn), totaling 119.3K instances.

Benchmarks. We assess the performance of TOUCAN across several key tool-agentic benchmarks, including BFCL V3 Patil et al. (2025), τ -Bench Yao et al. (2024), τ^2 -Bench (Barres et al., 2025), and MCP-Universe Luo et al. (2025). All evaluations are conducted on an 8 \times H100 server. For BFCL-V3, we use the official evaluation setup. For τ -Bench and τ^2 -Bench, we employ GPT-40 as user simulators. For MCP-Universe, we configure the local evaluation environment as specified in the benchmark documentation.

4.2 EXPERIMENTAL RESULTS

TOUCAN Effectively Increases Agentic Tool-Calling Performance. Tables 2 and 3 present the experimental results of models fine-tuned on TOUCAN across BFCL V3, τ -Bench, and τ^2 -Bench, respectively. We make the following key observations: First, models fine-tuned with TOUCAN show performance improvements compared to baseline models without fine-tuning across almost all aspects of these three benchmarks, indicating that TOUCAN effectively enhances the agentic and tool-calling capabilities of models. Second, on BFCL V3, models fine-tuned on TOUCAN outperform larger production LLMs including DeepSeek-V3 and GPT-4.5-Preview in average scores and achieve top performance in the *multi-turn* subset. This demonstrates the effectiveness of TOUCAN and validates our dataset design.

TOUCAN Enhances Models' Performance on Using Real-World MCP Servers. Figure 7 demonstrates a performance comparison between TOUCAN-tuned models and other open-source models of similar or larger

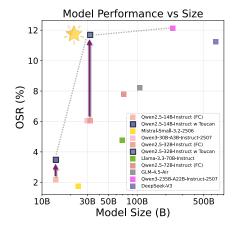


Figure 8: Model Performance vs Size on MCP-Universe Benchmark. We report overall task success rate (OSR). Our models achieve a new Pareto optimum.

sizes across six domains: Location Navigation, Repository Management, Financial Analysis, 3D Design, Browser Automation, and Web Search. We note that most servers in the benchmark require careful configurations and thus were not included in our data synthesis pipeline. Nevertheless, TOU-

CAN-tuned models show significant improvements on these challenging tasks compared to baselines, indicating that exposure to diverse tools enhances model performance on agentic tasks. Notably, our 32B model achieves the highest scores in 3D Design and strong performance in Financial Analysis, even outperforming much larger frontier open models like Llama-3.3-70B-Instruct, Qwen2.5-72B-Instruct, GLM-4.5-Air (106B), and DeepSeek-V3 (671B).

Figure 8 plots model performance versus model size on MCP-Universe benchmark. We observe that TOUCAN -tuned models establish a new Pareto optimum, indicating that TOUCAN can help models achieve superior performance-efficiency trade-offs in agentic tasks.

4.3 ABLATION ANALYSIS

To validate our extension designs, we perform ablation analysis on the Qwen2.5-14B-Instruct model, where we fine-tune on progressively extended versions of TOUCAN, allowing us to isolate the contributions of each extension described in Section 3.2. The experimental results are shown in Figrue 9. We observe that all components contribute to improved scores. Detailed benchmark scores for the BFCL ablation study are provided in Appendix C.3.

Figure 9: This table shows ablation analysis of TOUCAN extensions.

	BFCLv3	au-bench		
		Airline @1	Retail @ 1	
Qwen2.5-14B-Instruct	57.69%	17.25%	44.46%	
+ Single Turn	60.16%	15.50%	36.95%	
+ Irrelevance	64.74%	16.75%	41.63%	
+ Diversify	64.56%	17.25%	43.70%	
+ Multi-Turn	65.09%	22.00%	48.48%	

5 CONCLUSION AND FUTURE WORK

This paper introduces TOUCAN, a tool-agentic dataset containing 1.5M trajectories designed to train better agentic models. We propose a comprehensive pipeline for data generation and demonstrate that models fine-tuned on TOUCAN achieve superior performance on benchmarks including BFCL-V3 and MCP-Universe. TOUCAN represents the first step in a long-term effort to leverage tool use for building stronger LLM agents. Despite being a valuable contribution, we acknowledge our work exhibits certain limitations, which we plan to address through different initiatives.

Expanding to More MCP Servers. While our dataset is comprehensive, it was collected in June 2025, and new servers continue to emerge. We excluded MCP servers requiring special configurations (e.g., requires API keys or account setups), which simplifies the onboarding procedure but may overlook important servers and widely-used scenarios (e.g., Notion and GitHub). Manually onboarding more servers or developing automated onboarding agents could be valuable future work.

Expert models to simulate tool-responses. While real tool execution produces higher-quality results, it is often slow and costly, and therefore, not an option for everyone. To provide an alternative that also yields quality, we plan to develop an expert LLM capable of simulating tool execution. This artificial component will significantly reduce the cost of generating trajectory data involving tool use. Although the idea of tool-execution simulation is known within the community, it has most likely been implemented using off-the-shelf, closed-source LLMs.

MCP Benchmark for web search. As tool-use capabilities become central to both LLMs and LLM-agents, specific scenarios such as web search have gained prominence in the community as a means of synthesizing complex reasoning tasks. To advance this direction, we plan to develop an MCP benchmark focused on web search capabilities.

6 USE OF LARGE LANGUAGE MODELS (LLMS)

In our work, we used large language models (LLMs) to assist with improving the grammar, clarity, and overall readability of the manuscript, as well as to help generate the pipeline diagram included in the paper. All LLM-generated content was thoroughly verified by the authors as part of an iterative process to ensure accuracy, quality, and consistency with the scientific contributions of the work.

7 ETHICS STATEMENT

Developers planning to use Toucan for LLM fine-tuning should take into account certain considerations.

Data Ownership and Licensing. The MCP server specification files used to build TOUCAN were collected in June 2025 from https://smithery.ai/, a public platform hosting such specifications. These files were voluntarily published by their owners in accordance with the platform's privacy notice. Given the case a legitimate owner requests removal of their content from our dataset, we will honor that request through a take down process available via our GitHub repository.

Sensitive Information. The risk of exposing sensitive data in specification files is minimal, as they generally rely on placeholders rather than real information. However, human error may still lead to the inclusion of URLs, tokens, or email addresses. To mitigate this, we apply a pre-filtering stage with rule-based verifiers that detect common patterns of personally identifiable information (PII).

Data Evolution. Our data were collected in June 2025, so TOUCAN captures real-world tool-use scenarios available at that time. For example, responses from search MCP servers reflect information current through June 2025. To facilitate future updates and customization, we provide our modular data pipeline, allowing researchers and practitioners to expand domain coverage and tailor tool representations for their applications.

LLM Hallucinations. Only tasks and annotations in TOUCAN were generated with LLMs; trajectories were produced using LLMs in combination with agent frameworks and remote MCP servers. This integration ensures reliable tool call executions and responses, reducing the likelihood of code errors from hallucinations. Nevertheless, hallucinations remain a general risk when using LLMs, and outputs from models fine-tuned with TOUCAN should always be verified by humans.

8 REPRODUCIBILITY STATEMENT

We provide the code for our data generation pipeline, along with detailed instructions for executing the pipeline end-to-end, as well as sample dataset files in the supplementary materials. The main paper and appendix further document key implementation details, including prompt templates, hyperparameter configurations used during fine-tuning, and extensions of our data analysis and fine-tuning experiments. After publication, we plan to release the full codebase in a public GitHub repository and make our datasets publicly available on the HuggingFace platform.

REFERENCES

- Anthropic. Introducing the model context protocol. https://www.anthropic.com/news/model-context-protocol, 2025. Accessed: 2025-08-18.
- Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. τ²-bench: Evaluating conversational agents in a dual-control environment, 2025. URL https://arxiv.org/abs/2506.07982.
- Kinjal Basu, Ibrahim Abdelaziz, Subhajit Chaudhury, Soham Dan, Maxwell Crouse, Asim Munawar, Sadhana Kumaravel, Vinod Muthusamy, Pavan Kapanipathi, and Luis A. Lastras. Apiblend: A comprehensive corpora for training and benchmarking api llms, 2024. URL https://arxiv.org/abs/2402.15491.
- Fouad Bousetouane. Agentic systems: A guide to transforming industries with vertical ai agents, 2025. URL https://arxiv.org/abs/2501.00881.
- Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Yuefeng Huang, Wulong Liu, Xinzhi Wang, Defu Lian, Baoqun Yin, Yasheng Wang, and Wu Liu. ACEBench: Who Wins the Match Point in Tool Usage?, July 2025. URL http://arxiv.org/abs/2501.12851.arXiv:2501.12851 [cs].
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang,

541

543

544

546

547

548

549

550

551

552

553

554

556

558

559

560

561

563

564

565

566

567

568 569

570

571

572

573

574

575 576

577

578

579

580

581 582

583

584

585

586

587

588

589

590

592

Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL https://arxiv.org/abs/2412.19437.

- Mohamed Amine Ferrag, Norbert Tihanyi, and Merouane Debbah. From Ilm reasoning to autonomous ai agents: A comprehensive review, 2025. URL https://arxiv.org/abs/2504.19678.
- Xuanqi Gao, Siyi Xie, Juan Zhai, Shqing Ma, and Chao Shen. Mcp-radar: A multi-dimensional benchmark for evaluating tool use capabilities in large language models, 2025. URL https://arxiv.org/abs/2505.16700.
- Zhen Guo, Adriana Meza Soria, Wei Sun, Yikang Shen, and Rameswar Panda. Api pack: A massive multi-programming language dataset for api call generation, 2025a. URL https://arxiv.org/abs/2402.09615.
- Zikang Guo, Benfeng Xu, Chiwei Zhu, Wentao Hong, Xiaorui Wang, and Zhendong Mao. Mcpagentbench: Evaluating real-world language agent performance with mcp-mediated tools, 2025b. URL https://arxiv.org/abs/2509.09734.
- X. Li, S. Wang, S. Zeng, et al. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth*, 1:9, 2024. doi: 10.1007/s44336-024-00009-2. URL https://doi.org/10.1007/s44336-024-00009-2.
- Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, Zezhong Wang, Yuxian Wang, Wu Ning, Yutai Hou, Bin Wang, Chuhan Wu, Xinzhi Wang, Yong Liu, Yasheng Wang, Duyu Tang, Dandan Tu, Lifeng Shang, Xin Jiang, Ruiming Tang, Defu Lian, Qun Liu, and Enhong Chen. Toolace: Winning the points of llm function calling, 2025a. URL https://arxiv.org/abs/2409.00920.
- Zhiwei Liu, Jielin Qiu, Shiyu Wang, Jianguo Zhang, Zuxin Liu, Roshan Ram, Haolin Chen, Weiran Yao, Shelby Heinecke, Silvio Savarese, Huan Wang, and Caiming Xiong. Mcpeval: Automatic mcp-based deep evaluation for ai agent models, 2025b. URL https://arxiv.org/abs/2507.12806.

Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh Murthy, Liangwei Yang, Silvio Savarese, Juan Carlos Niebles, Huan Wang, Shelby Heinecke, and Caiming Xiong. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets, 2024. URL https://arxiv.org/abs/2406.18518.

- Ziyang Luo, Zhiqi Shen, Wenzhuo Yang, Zirui Zhao, Prathyusha Jwalapuram, Amrita Saha, Doyen Sahoo, Silvio Savarese, Caiming Xiong, and Junnan Li. Mcp-universe: Benchmarking large language models with real-world model context protocol servers, 2025. URL https://arxiv.org/abs/2508.14704.
- Leland McInnes and John Healy. Umap: Uniform manifold approximation and projection for dimension reduction. *ArXiv*, abs/1802.03426, 2018. URL https://api.semanticscholar.org/CorpusID:3641284.
- Dhruv Nathawani, Igor Gitman, Somshubra Majumdar, Evelina Bakhturina, Ameya Sunil Mahabaleshwarkar, , Jian Zhang, and Jane Polak Scowcroft. Nemotron-Post-Training-Dataset-v1, 2025. URL https://huggingface.co/datasets/nvidia/Nemotron-Post-Training-Dataset-v1.
- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis, 2023. URL https://arxiv.org/abs/2305.15334.
- Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.
- Akshara Prabhakar, Zuxin Liu, Ming Zhu, Jianguo Zhang, Tulika Awalgaonkar, Shiyu Wang, Zhiwei Liu, Haolin Chen, Thai Hoang, Juan Carlos Niebles, Shelby Heinecke, Weiran Yao, Huan Wang, Silvio Savarese, and Caiming Xiong. Apigen-mt: Agentic pipeline for multi-turn data generation via simulated agent-human interplay, 2025. URL https://arxiv.org/abs/2504.03601.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world apis, 2023. URL https://arxiv.org/abs/2307.16789.
- Donghao Ren, Fred Hohman, and Dominik Moritz. A scalable approach to clustering embedding projections, 2025. URL https://arxiv.org/abs/2504.07285.
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. ToolAlpaca: Generalized Tool Learning for Language Models with 3000 Simulated Cases, September 2023. URL http://arxiv.org/abs/2306.05301.arXiv:2306.05301 [cs].
- 5 Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang, Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei, Yean Cheng, Yifan An, Yilin Niu, Yuanhao Wen, Yushi Bai, Zhengxiao Du, Zihan Wang, Zilin Zhu, Bohan Zhang, Bosi Wen, Bowen Wu, Bowen Xu, Can Huang, Casey Zhao, Changpeng Cai, Chao Yu, Chen Li, Chendi Ge, Chenghua Huang, Chenhui Zhang, Chenxi Xu, Chenzheng Zhu, Chuang Li, Congfeng Yin, Daoyan Lin, Dayong Yang, Dazhi Jiang, Ding Ai, Erle Zhu, Fei Wang, Gengzheng Pan, Guo Wang, Hailong Sun, Haitao Li, Haiyang Li, Haiyi Hu, Hanyu Zhang, Hao Peng, Hao Tai, Haoke Zhang, Haoran Wang, Haoyu Yang, He Liu, He Zhao, Hongwei Liu, Hongxi Yan, Huan Liu, Huilong Chen, Ji Li, Jiajing Zhao, Jiamin Ren, Jian Jiao, Jiani Zhao, Jianyang Yan, Jiaqi Wang, Jiayi Gui, Jiayue Zhao, Jie Liu, Jijie Li, Jing Li, Jing Lu, Jingsen Wang, Jingwei Yuan, Jingxuan Li, Jingzhao Du, Jinhua Du, Jinxin Liu, Junkai Zhi, Junli Gao, Ke Wang, Lekang Yang, Liang Xu, Lin Fan, Lindong Wu, Lintao Ding, Lu Wang, Man Zhang, Minghao Li, Minghuan Xu, Mingming Zhao, Mingshu Zhai, Pengfan Du, Qian Dong, Shangde Lei, Shangqing Tu, Shangtong Yang, Shaoyou Lu, Shijie Li, Shuang Li, Shuang-Li, Shuxun Yang, Sibo Yi, Tianshu Yu, Wei Tian, Weihan Wang, Wenbo Yu, Weng Lam Tam, Wenjie Liang, Wentao Liu, Xiao Wang, Xiaohan Jia,

649

650

651

652

653

654

655

656 657

658

659

660

661

662

663

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682 683

684

685

686

687

688

689

690

691 692

693

694

696

697

698

699

700

Xiaotao Gu, Xiaoying Ling, Xin Wang, Xing Fan, Xingru Pan, Xinyuan Zhang, Xinze Zhang, Xiuqing Fu, Xunkai Zhang, Yabo Xu, Yandong Wu, Yida Lu, Yidong Wang, Yilin Zhou, Yiming Pan, Ying Zhang, Yingli Wang, Yingru Li, Yinpei Su, Yipeng Geng, Yitong Zhu, Yongkun Yang, Yuhang Li, Yuhao Wu, Yujiang Li, Yunan Liu, Yunqing Wang, Yuntao Li, Yuxuan Zhang, Zezhen Liu, Zhen Yang, Zhengda Zhou, Zhongpei Qiao, Zhuoer Feng, Zhuorui Liu, Zichen Zhang, Zihan Wang, Zijun Yao, Zikang Wang, Ziqiang Liu, Ziwei Chai, Zixuan Li, Zuodong Zhao, Wenguang Chen, Jidong Zhai, Bin Xu, Minlie Huang, Hongning Wang, Juanzi Li, Yuxiao Dong, and Jie Tang. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models, 2025a. URL https://arxiv.org/abs/2508.06471.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence, 2025b. URL https://arxiv.org/abs/2507.20534.

- Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.
- The MCPMark Team. Mcpmark: Stress-testing comprehensive mcp use. https://github.com/eval-sys/mcpmark, 2025a.
- The Scale Research Team. Actions, not words: Mcp-atlas raises the bar for agentic evaluation. https://scale.com/blog/mcp-atlas, September 2025b. Accessed: YYYY-MM-DD.
- Zhenting Wang, Qi Chang, Hemani Patel, Shashank Biju, Cheng-En Wu, Quan Liu, Aolin Ding, Alireza Rezazadeh, Ankit Shah, Yujia Bao, and Eugene Siow. Mcp-bench: Benchmarking toolusing llm agents with complex real-world tasks via mcp servers, 2025. URL https://arxiv.org/abs/2508.20453.
- Weikai Xu, Chengrui Huang, Shen Gao, and Shuo Shang. Llm-based agents for tool learning: A survey. *Data Science and Engineering*, 2025a. doi: 10.1007/s41019-025-00296-9. URL https://link.springer.com/article/10.1007/s41019-025-00296-9.
- Zhangchen Xu, Yang Liu, Yueqin Yin, Mingyuan Zhou, and Radha Poovendran. Kodcode: A diverse, challenging, and verifiable synthetic dataset for coding. *ArXiv*, abs/2503.02951, 2025b. URL https://api.semanticscholar.org/CorpusID:276782338.
- Yunhe Yan, Shihe Wang, Jiajun Du, Yexuan Yang, Yuxuan Shan, Qichen Qiu, Xianqing Jia, Xinge Wang, Xin Yuan, Xu Han, Mao Qin, Yinxiao Chen, Chen Peng, Shangguang Wang, and Mengwei Xu. Mcpworld: A unified benchmarking testbed for api, gui, and hybrid computer use agents, 2025. URL https://arxiv.org/abs/2506.07672.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.

Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for tool-agent-user interaction in real-world domains, 2024. URL https://arxiv.org/abs/2406.12045.

Ming Yin, Dinghan Shen, Silei Xu, Jianbing Han, Sixun Dong, Mian Zhang, Yebowen Hu, Shujian Liu, Simin Ma, Song Wang, et al. Livemcp-101: Stress testing and diagnosing mcp-enabled agents on challenging queries. *arXiv preprint arXiv:2508.15760*, 2025.

A DATASET SCHEMA AND EXAMPLES

756

758

759 760

761

762

764

765

766

768

769

770 771

772

773

774 775

776 777

778

779

780 781

782

An instance of TOUCAN contains the following columns:

- uuid: Unique sample identifier.
- **subset:** Annotation specifying which pipeline was used to generate the trajectory. Options: (1) *single-turn-original:* only the core processing (Stage 1 to 5) described in Section 3 are applied, (2) *irrelevant:* a server shuffle process applied on top of the *single-turn-original* pipeline, (3) *single-turn-diversify:* a question diversification process applied on top of the *single-turn-original* pipeline, and (4) *multi-turn:* a multi-turn extension of the *single-turn-original* and *single-turn-diversify* subsets.
- messages: The trajectory formatted with the chat template from the original LLM-agent used for generation. The system prompt includes the associated list of tools.
- question: The user task crafted to generate the trajectory.
- **target_tools:** The MCP tools used as seeds for question generation.
- question_quality_assessment: Task evaluation by an LLM-as-judge, covering quality, difficulty, realism, and uniqueness.
- response_quality_assessment: Response evaluation by an LLM-as-judge, covering completeness and conciseness.
- message_num_rounds: Total number of messages, including turns of all types.
- metadata: Original MCP server data collected and used as seed for generation, as well as respective LLM annotations.

This is the structure of an instance in TOUCAN:

```
783
             "uuid": "3ac8fdcc-b9b5-50d2-a840-947a42b558d2",
784
             "subset": "single-turn-original",
785
             "messages": "[{...long JSON string of messages...}]",
786
             "question": "I'm reviewing my company's policy document that needs
787
                quarterly updates for compliance purposes. I need to first
                 understand its current structure with all sections and
                subsections, then update every occurrence of 'Q3 2024' to 'Q1
789
              → 2025' throughout the document, and finally add explanatory
790
                footnotes to the last 3 sections that reference our new
791
                 regulatory guidance document published last month. Can you help
                 me accomplish these updates systematically?",
793
             "target_tools": "get_document_outline, search_and_replace,
                add_footnote_to_document"
794
             "question_quality_assessment_kimik2w4a16": {
               "tool_selection_difficulty": { "reasoning": "...", "score": 3 },
796
               "tool_selection_uniqueness": { "reasoning": "...", "score": 3 },
797
               "question_quality": { "reasoning": "...", "score": 4 },
798
               "scenario_realism": { "reasoning": "...", "score": 5 },
               "verifiable": { "reasoning": "...", "score": 4 },
799
               "stability": { "reasoning": "...", "score": 5 },
800
               "overall_score": 4.0
801
802
             "response_quality_assessment_oss120b": {
               "completeness": { "reasoning": "...", "score": 2 },
"conciseness": { "reasoning": "...", "score": 3 },
804
               "overall_score": 2.5,
805
               "desired_tools_used_percentage": 0.3333,
               "order_correctness": false
             },
808
             "messages_num_rounds": 9,
             "metadata": "{...long metadata JSON string...}"
```

811

812

813

Below, we present a TOUCAN trajectory, chat template tokens from the model generator and emojis have been removed, and some sections (e.g., tool list, tool output, assistant response) have been abbreviated to improve readability.

```
814
       [
815
816
           "role": "system",
817
           "content": "{... contains the list of tools ... }"
818
         },
819
           "role": "user",
820
           "content": "Can you explain what ENS (Ethereum Name Service)
821
           \hookrightarrow represents in the context of a transaction and why a name like
822
              'blockscout.eth' might appear instead of a raw address?"
823
         },
824
           "role": "assistant",
825
           "content": "I'll explain what ENS (Ethereum Name Service) represents,
826
           \hookrightarrow and show you how it resolves a name like 'blockscout.eth' into
827
           → its actual Ethereum address."
         },
828
829
           "role": "assistant",
830
           "content": ""
831
           "function_call": {
832
             "name": "blockscout-mcp-server-__unlock_blockchain_analysis___",
833
             "arguments": "{}"
834
         },
835
836
           "role": "function",
837
           "content": "{... long tool output ...}"
838
839
           "role": "assistant",
840
           "content": "Now, let me resolve 'blockscout.eth' to show the actual
841
           → address it represents."
842
         },
843
           "role": "assistant",
844
           "content": ""
845
           "function_call": {
846
             "name": "blockscout-mcp-server-get_address_by_ens_name",
847
             "arguments": "{\"name\": \"blockscout.eth\"}"
848
         },
849
850
           "role": "function",
851
           "name": "blockscout-mcp-server-get_address_by_ens_name",
852
           "content": "{\n \"data\": {\n \"resolved_address\":
               \"0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975\"\n },\n
853
               \"data_description\": null,\n \"notes\": null,\n
854
               \"instructions\": null,\n \"pagination\": null\n}"
855
         },
856
857
           "role": "assistant",
           "content": "{... parsed assistant response ...}"
858
         },
859
860
           "role": "user",
861
           "content": "Based on that explanation, could you look up the ENS name
862
               'blockscout.eth' and provide the Ethereum address it resolves
863

    to?"

         },
```

```
864
865
           "role": "assistant",
           "content": "I already resolved that for you in my previous response!
           → When I called the ENS lookup tool, `\"blockscout.eth\"` resolved
867

→ to the Ethereum

868
           → address:\n\n`0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975"
869
         }
870
       1
871
872
       In the previous trajectory example, we abbreviated the assistant's parsed response for readability.
873
       The full response is shown below.
874
875 1
       ## What ENS (Ethereum Name Service) Represents in Transactions
876
       ENS stands for **Ethereum Name Service**, and it functions much like "DNS
877
           for Ethereum addresses." It allows users to register human-readable
878
           names (like "blockscout.eth") that map to complex hexadecimal
879

→ blockchain addresses.

880 4
       ### Key Points About ENS:
881 5
882 7
       1. **Human-Readable Names for Complex Addresses**
          - Instead of displaying `0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975`
8 888
          \hookrightarrow (which is the raw address for blockscout.eth), interfaces show the \hookrightarrow friendly "blockscout.eth"
884
885
886 9
          - This reduces errors from copy/pasting long hex addresses
887 11
       2. **Reverse Resolution**
888 12
          - When you see "blockscout.eth" in transaction data, wallet labels, or

→ block explorers, it means:
889
890 13
            - ENS has been configured with this name
            - The system has performed a reverse lookup to find the name
891
            \hookrightarrow associated with the address
892 15
893 16
       3. **Self-Sovereign Identity**
          - ENS names are owned and controlled by users through smart contracts
894 17
          - They can be transferred, updated, or associated with additional
895 18
          → metadata (social profiles, avatars, etc.)
896
897 <sub>20</sub>
      ### Why "blockscout.eth" Appears Instead of Raw Address:
898 21
       - **User Experience**: "blockscout.eth" is immediately recognizable and
899 22
       \hookrightarrow memorable
901 23
      - **Trust Verification**: Users can verify they're interacting with the
       902 24
      - **Branding**: Organizations use ENS names for consistent branding
903

→ across DeFi protocols and wallet interfaces

904 25
      - **Error Prevention**: Reduces costly mistakes when users verify

→ transaction recipients

905
906 26
       Every time you see an ENS name like "blockscout.eth" in a transaction,
907
          the underlying system has actually resolved this to the precise
908
       \hookrightarrow blockchain address `0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975`-the
       \,\hookrightarrow\, ENS simply provides a human-friendly layer on top of the raw
909

→ cryptographic identity.

910
911
912
913
914
915
916
```

B MORE ON DATASET ANALYSIS

B.1 FEATURED SERVERS IN TOUCAN

Table 4 demonstrates featured servers of TOUCAN. We note that these featured servers span a wide range of domains, making combinations possible for complex tool use tasks.

Server Name	Repository Name	Domain	Tools Count
OKX Server	esshka/okx-mcp	Cryptocurrency & Blockchain	2
AI Research Assistant - Semantic Scholar	Access via Smithery ²	Web Search & Research	10
Book Search Server	Access via Smithery ³	Web Search & Research	1
PubMed MCP Server	JackKuo666/PubMed-MCP-Server	Web Search & Research	4
Flux ImageGen Server	falahgs/flux-imagegen-mcp-server	AI/ML Tools	3
Pokémcp	NaveenBandarage/poke-mcp	Data Analysis & Processing	4
Hotel Booking Server	jinkoso/jinko-mcp	E-commerce	6
Cloudflare Playwright	cloudflare/playwright-mcp	Browser Automation	24
Time MCP Server	yokingma/time-mcp	Time & Calendar	6
Exa Search	exa-labs/exa-mcp-server	Web Search & Research	8
Weather Forecast Server	iremaltunay55/deneme	Weather	5
Advanced Calculator Server	alan5543/calculator-mcp	Data Analysis & Processing	17
Dictionary Server	ceydasimsekk/dictionarymcp	Others	1
Airbnb Search and Listing Details Server	AkekaratP/mcp-server-airbnb	Web Search & Research	2
Code Runner MCP Server	formulahendry/mcp-server-code-runner	Development Tools	1
Movie Recommender	iremert/movie-recommender-mcp	Content Creation	1
United States Weather	smithery-ai/mcp-servers	Weather	6
Context7	upstash/context7-mcp	Development Tools	2
Think Tool Server	PhillipRt/think-mcp-server	Memory Management	1
OpenAPI MCP Server	janwilmake/openapi-mcp-server	API Integration	2
Film Information Server	zehranurugurr/film_mcp	Content Creation	1
Trends Hub	baranwang/mcp-trends-hub	News & Media	21
ClinicalTrials MCP Server	JackKuo666/ClinicalTrials-MCP-Server	Health & Fitness	7
Drawing Tool for AI Assistants	flrngel/mcp-painter	Content Creation	4
LeetCode	jinzcdev/leetcode-mcp-server	Development Tools	9

Table 4: Featured Server Information

B.2 More on MCP Server Analysis in Toucan

Figure 10 shows the distribution of the most frequently used MCP servers in our dataset, highlighting the diversity of servers and domains covered in TOUCAN. Figure 11 shows the distribution of tool counts across the 495 MCP servers employed by TOUCAN, revealing that most servers expose only a limited number of tools, with the majority containing fewer than 10 tools.

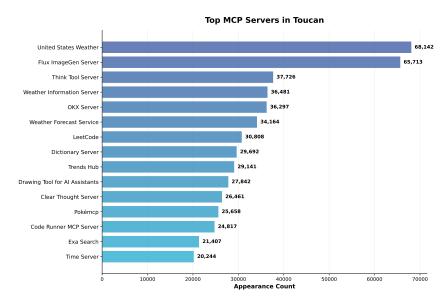


Figure 10: Distribution of the most frequently occurring MCP servers in the TOUCAN dataset.

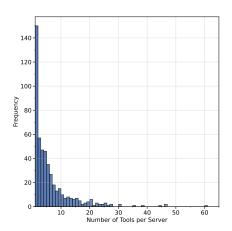


Figure 11: Tools Number distribution across MCP servers

B.3 EMEDDING VISUALIZATION

Figure 12 presents embedding visualization via Embedding Atlas (Ren et al., 2025) using the Xenova/multilingual-e5-small embedding model with UMAP projection McInnes & Healy (2018). The visualization demonstrates that TOUCAN covers a wide range of topics. In addition, the proposed TOUCAN extensions (e.g., diversification) effectively increase the overall dataset coverage.

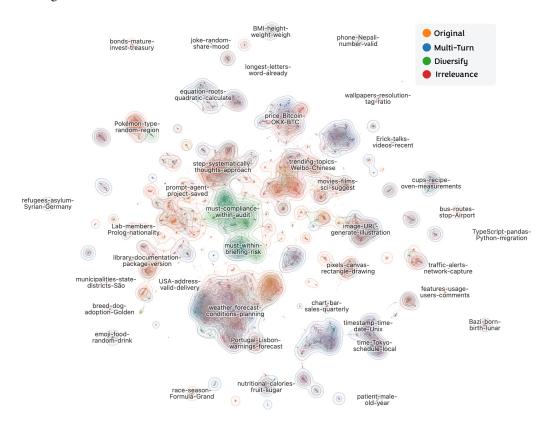


Figure 12: This figure is the visualization of 50K random-sampled TOUCAN instances via Embedding Atlas (Ren et al., 2025).

C MORE ON EXPERIMENTS

C.1 LLM ANNOTATION

 Figure 13 shows the Pearson correlation between human annotators and LLM-as-a-Judge evaluations across different models, based on 50 randomly sampled instances. The annotation prompt is available in Appendix D.4. We observe that ${\tt GPT-4.1}$ and ${\tt Kimi-K2}$ achieve the highest overall correlation with human judgments. Considering cost efficiency, we deploy ${\tt Kimi-K2}$ locally for our annotation pipeline.

C.2 FINE-TUNING HYPER-PARAMETERS

We fine-tune models with TOUCAN using a super computing cluster, which is outfitted with NVIDIA H100 GPUs. The fine-tuning hyper-parameters can be found in Table 5.

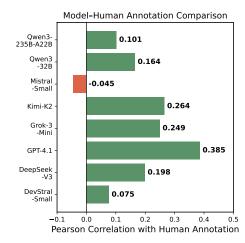


Figure 13: Pearson correlation between human annotators and LLM-as-a-Judge evaluations across different models.

Table 5: This table shows the hyper-parameters for supervised fine-tuning.

Hyper-parameter	Value
Tool-Call Template	Hermes
Learning Rate	2×10^{-5}
Number of Epochs	2
Number of Devices	8 or 64
Per-device Batch Size	1
Gradient Accumulation Steps	8 (8 GPUs) or 1 (64 GPUs)
Effective Batch Size	64
Optimizer	Adamw with $\beta s = (0.9, 0.999)$ and $\epsilon = 10^{-8}$
Deepspeed	zero3
Max Sequence Length	32768

C.3 MORE ON ABLATION STUDIES

Table 6 details the individual scores of the BFCL V3 benchmark for our ablation analysis. We observe that all extensions are meaningful in improving model performance.

Table 6: Ablation of TOUCAN Extensions on BFCL V3 Benchmark.

	Overall	Single Turn		Multi Turn	Hallucination	
		Non-live (AST)	Live (AST)		Relevance	Irrelevance
Qwen2.5-14B-Instruct	57.69%	83.38%	73.70%	19.75%	83.33%	68.46%
+ Single Turn	60.16%	87.50%	66.86%	34.38%	72.22%	46.88%
+ Irrelevance	64.74%	88.46%	77.25%	30.38%	72.22%	77.85%
+ Diversify	64.56%	86.06%	76.90%	32.50%	72.22%	75.45%
+ Multi-Turn	65.09%	85.42%	76.01%	35.25%	72.22%	75.96%

```
1080
      D PROMPTS
1081
1082
       D.1 MCP SERVER ANNOTATION PROMPT
1083
1084
       Below is the prompt for annotating MCP server categories.
1085
1086
     ## Task
1087<sub>2</sub> Generate **Server Labels** to categorize the provided MCP Server based on
       \rightarrow its description and available tools.
1089 3
      ## Objective
1090 <sup>4</sup>
      Analyze the provided MCP Server's description and available tools, then
1091 5
       → assign appropriate category labels that best describe its primary
1092
       \hookrightarrow functionality and use cases.
1093 6
10947 ## Guidelines
1095<sup>8</sup>
    9 ### Label Selection
1096
      - Analyze the MCP Server's core functionality and purpose
      - Consider the types of tools it provides and the problems it solves
109812 - Select labels that accurately represent the server's primary use cases
109913 - Choose from predefined categories when applicable, but also consider
       \rightarrow custom labels for unique functionality
1100
1101115
      ### Predefined Categories
1102<sub>16</sub>
      Choose from these established categories when appropriate:
110317
      - **Web Search & Research**: Tools for searching the web, gathering
       \hookrightarrow information, academic research
1104
      - **Browser Automation**: Web scraping, automated browsing, page
1105<sup>18</sup>

→ interaction

1106<sub>19</sub>
      - **Memory Management**: Data storage, retrieval, knowledge bases,
1107

→ note-taking

110820
      - **Operating System**: File operations, system commands, process

→ management

      - **Data Analysis & Processing**: Analytics, data transformation,
1110^{21}

→ statistical analysis

1111
      - **Cryptocurrency & Blockchain**: Trading, wallet management, DeFi,
1112
       \hookrightarrow blockchain interaction
111323
      - **Daily Productivity**: Task management, scheduling, personal

→ organization

1114
1115<sup>24</sup>
      - **File Management**: File operations, document handling, storage
       → management
1116
      - **Database Operations**: Data querying, database management, SQL
1117

→ operations

111826
      - **API Integration**: Third-party service integration, webhook handling
      - **Communication Tools**: Messaging, email, notifications, social
111927

→ interaction

      - **Development Tools**: Code analysis, debugging, version control, CI/CD
1121<sub>29</sub>
      - **Security & Authentication**: Password management, encryption, access
1122
       112330
      - **Cloud Services**: Cloud platform integration, serverless functions
      - **AI/ML Tools**: Machine learning, model interaction, AI-powered
112431
       1125
      - **Content Creation**: Writing, editing, media generation, publishing
1126 33
      - **Social Media**: Social platform integration, posting, analytics
1127<sub>34</sub>
      - **Financial Services**: Banking, payments, financial data, accounting
      - **E-commerce**: Shopping, product management, order processing
      - **Gaming**: Game-related tools, entertainment, interactive features
112936
1130<sup>37</sup>
      - **Education**: Learning tools, course management, educational content
      - **Health & Fitness**: Health monitoring, fitness tracking, medical
1131

→ tools

113239
      - **Travel & Maps**: Location services, travel planning, navigation
113340
      - **News & Media**: News aggregation, media consumption, journalism tools
```

- **Weather**: Weather data, forecasting, climate information

```
- **Time & Calendar**: Scheduling, time management, calendar integration
1135<sub>43</sub>
113644
      ### Custom Labels
113745 - If the server doesn't fit well into predefined categories, create a
       1138
1139<sup>46</sup>
      - Custom labels should be descriptive and specific to the server's unique
       - Use clear, concise terminology that would be useful for clustering and
1141
       → organization
114248
       ### Output Requirements
1143^{49}
       - **Primary Label**: The main category that best describes the server
1144<sup>50</sup>
       1145
      - **Secondary Labels**: Additional relevant categories (0-2 labels)
114652
      - **Custom Label**: A free-form descriptive label if the server has
1147
       \hookrightarrow unique functionality not covered by predefined categories
114853
1149<sup>54</sup>
       ## MCP Server Description
       {MCP_SERVER_NAME}: {MCP_SERVER_DESCRIPTION}
115056
1151<sub>57</sub> Available Tools:
115258
      {TOOL_LIST}
1153<sup>59</sup>
1154<sup>60</sup>
      ## Output
      Provide your response in the following XML format:
   61
1155
1156<sub>63</sub>
      <response>
115764
        <analysis>
          <!-- Briefly analyze the MCP Server's core functionality and the
115865
           → types of problems it solves based on its description and

→ available tools. -->

1160<sub>66</sub>
         </analysis>
1161<sub>67</sub>
        <reasoning>
           <!-- Brief explanation of why these labels were chosen and how they
116268

→ represent the server's functionality -->

         </reasoning>
1164<sup>69</sup>
        primary_label>
1165
        <!-- The main category that best describes this server's primary 

functionality -->
1166
116772 </primary_label>
1168<sup>73</sup>
      <secondary_labels>
1169<sup>74</sup>
          <!-- Additional relevant categories (0-2 labels), separated by commas

    if multiple -->

        </secondary_labels>
1171<sub>76</sub>
        <custom_label>
           <!-- A free-form descriptive label if the server has unique
117277
           \,\,\hookrightarrow\,\, functionality not covered by predefined categories. Leave empty
1173
           → if not needed. -->
         </custom_label>
1175<sub>79</sub>
      </response>
1176
1177
       D.2 TASK GENERATION PROMPT
1178
1179
       Below is an example of a task generation prompt for the single-server task synthesis. The prompt
1180
       generates a question targeting one tool.
1181
1182<sub>1</sub> ## Task
1183 2 Generate a **Tool Use Question** based on the provided MCP Server and its
       \rightarrow tool descriptions.
1184
1185<sup>3</sup>
      ## Objective
1186 5 Analyze the provided MCP Server and its available tools, then create a
1187
       → realistic user question that would naturally require the use of one
```

 \hookrightarrow of these tools to solve.

```
1188
1189 7
       ## Guidelines
1190<sub>8</sub>
11919 ### Ouestion Realism
119210 - Create questions that represent real-world scenarios where users would
       → need to interact with the MCP Server's tools
1193
      - The question should sound natural and authentic, as if asked by someone
       \rightarrow genuinely needing to accomplish a task
119512
      - Consider common use cases, problems, or workflows that would require
1196
       → the functionality provided by the MCP Server's tools
119713
      ### Tool Selection
1198<sup>14</sup>
      - Focus on **ONE specific tool** from the MCP Server that would be most
1199

→ appropriate to answer the question

1200<sub>16</sub>
      - Choose tools based on the core functionality they provide and how they
1201

→ would solve real user problems

      - Consider each tool's description and purpose when crafting the question
1202<sup>17</sup>
1203<sup>18</sup>
   19
      ### Question Complexity
120420
      - Create questions that are clear and specific enough to warrant tool
1205

→ usage

120621 - Avoid overly simple questions that could be answered without tools
_{1207^{22}} - Include relevant context or constraints that make the tool usage

→ necessary

1208
      - Do not contain the exact tool name in the question
1209_24
121025
       ### Output Format
121126
       Your response should include:
121227 1. **Tool Analysis**: Briefly analyze the MCP Server's available tools
       \hookrightarrow and their main functionalities.
1213
    28 2. **Target Tool**: The specific tool name from the MCP Server that
1214
       \hookrightarrow should be used to answer this question.
121529 3. **Question**: A clear, realistic user question that requires tool
1216

    usage.

1217<sup>30</sup>
1218<sup>31</sup>
       ## MCP Server Description
       {MCP_SERVER_NAME}: {MCP_SERVER_DESCRIPTION}
   32
1219
1220<sub>34</sub> Available Tools:
122135
      {TOOL_LIST}
122236
1223<sup>37</sup>
      ## Output
      Provide your response in the following XML format:
   38
1224 39
1225<sub>40</sub> <response>
122641
        <server_analysis>
          <!-- Briefly analyze the MCP Server's available tools and their main
122742

    functionalities. -->
1228
         </server_analysis>
122944
        <target_tool>
123045
           <!-- The specific tool name from the MCP Server that should be used
1231

→ to answer this question. -->

         </target_tool>
123246
1233<sup>47</sup>
         <question>
          <!-- A clear, realistic user question that requires tool usage. -->
123449
         </question>
1235_{50}
      </response>
1236
       Below is an example of a task generation prompt for the single-server task synthesis. The prompt
1237
       generates a question targeting multiple tools.
1238
1239 <sub>1</sub> ## Task
1240<sub>2</sub> Generate a **Tool Use Question** based on the provided MCP Server and its
1241

→ tool descriptions.
```

```
## Objective
1243 5
      Analyze the provided MCP Server and its available tools, then create a
1244
       \hookrightarrow realistic user question that would naturally require the use of
1245
       → **{NUM_TOOLS} tools** from this MCP Server to solve completely.
1246<sup>6</sup>
       ## Guidelines
1247 7
1248 9
       ### Question Realism
124910
       - Create questions that represent real-world scenarios where users would
       \rightarrow need to interact with the MCP Server's tools
125111
      - The question should sound natural and authentic, as if asked by someone
       \,\hookrightarrow\, genuinely needing to accomplish a task
1252
       - Consider common use cases, problems, or workflows that would require
1253
       \rightarrow the functionality provided by the MCP Server's tools
1254<sub>13</sub>
      ### Tool Selection
125514
      - Focus on **{NUM_TOOLS} tools** from the MCP Server that would work
1256<sup>15</sup>

→ together to answer the question

1257
      - The question should require a sequence or combination of tool calls to
1258
       \hookrightarrow solve completely
125917
      - Choose tools based on how they complement each other and create a
1260
       → logical workflow
      - Consider each tool's description and purpose when crafting the question
1261<sup>18</sup>
       \hookrightarrow that requires multiple steps
1262
1263<sub>20</sub>
      ### Ouestion Complexity
1264<sub>21</sub>
      - Create questions that are complex enough to warrant using {NUM TOOLS}
1265

→ tools

_{12662} - The question should have multiple components or require several steps

→ to solve

1267
      - Include relevant context or constraints that make the multi-tool usage
1268
1269_{24} - Do not contain the exact tool names in the question
127025 - Ensure the question cannot be reasonably answered with just a single
       → tool
1271
1272<sup>26</sup>
       ### Output Format
1273<sub>28</sub> Your response should include:
127429 1. **Tool Analysis**: Briefly analyze the MCP Server's available tools
       \rightarrow and their main functionalities.
1276<sup>30</sup>
      2. **Target Tools**: The specific tool names from the MCP Server that
       \hookrightarrow should be used together to answer this question, in the order they
1277
       \hookrightarrow would likely be called.
       3. **Question**: A clear, realistic user question that requires multiple
1279
       \hookrightarrow tool usage.
128032
      ## MCP Server Description
1281<sup>33</sup>
1282<sup>34</sup>
      {MCP_SERVER_NAME}: {MCP_SERVER_DESCRIPTION}
1283
      Available Tools:
128437
      {TOOL_LIST}
128538
1286<sup>39</sup>
      ## Output
1287<sup>40</sup>
       Ensure your question requires exactly {NUM_TOOLS} tools to solve
       → completely. Provide your response in the following XML format:
128841
128942
      <response>
129043
         <server_analysis>
129144
           <!-- Briefly analyze the MCP Server's available tools and their main

→ functionalities. -->

1292
         </server_analysis>
1293<sub>46</sub>
         <target_tools>
129447
           <!-- The specific tool names from the MCP Server that should be used
            \,\hookrightarrow\, together to answer this question, listed in order. e.g.,
1295

     <tool>create_twitter_post</tool> <tool>get_last_tweet</tool> -->
```

```
1296
_48
        </target_tools>
129749
        <question>
129850
          <!-- A clear, realistic user question that requires multiple tool

    usage. -->

        </question>
130051
1301<sup>52</sup>
      </response>
1302
      Below is an example of a task generation prompt for the multi-server task synthesis.
1303
13041 ## Task
13052 Generate a **Multi-Server Tool Use Question** based on the provided MCP
      → Servers and their tool descriptions.
1307 \frac{3}{4}
      ## Objective
1308 <sub>5</sub>
      Analyze the provided MCP Servers and their available tools, then create a
1309
      → realistic user question that would naturally require the use of
       \hookrightarrow **{NUM_TOOLS} tools from at least 2 different MCP servers** to solve
1310
      \hookrightarrow completely.
1311
1312 7
      ## Guidelines
1313<sub>8</sub>
13149 ### Question Realism
131510 - Create questions that represent real-world scenarios where users would
      1317 - The question should sound natural and authentic, as if asked by someone

→ genuinely needing to accomplish a complex task

1318_{12} - Consider workflows that span across different services/domains that
1319
          would require multiple servers
      - Think about how different MCP servers complement each other in
132013
       \hookrightarrow real-world use cases
1321
1322
      ### Server and Tool Selection
1323<sub>16</sub>
      - Use tools from **at least 2 different MCP servers** to answer the
1324

→ question

132517 - Select **{NUM_TOOLS} tools total** that work together across multiple

→ servers

1326
1327
      - The question should require a sequence or combination of tool calls

→ from different servers to solve completely

1328_{19} - Choose tools based on how they complement each other across different

→ services/domains

133020 - Consider each tool's description and purpose when crafting the

→ cross-server workflow

1331
1332 21
      - Ensure tools from different servers create a logical, interconnected
       → workflow
133322
133423
      ### Question Complexity
      - Create questions that are complex enough to warrant using {NUM_TOOLS}
1335^{24}

→ tools across multiple servers

The question should have multiple components or require several steps

→ that span different services

133826
      - Include relevant context or constraints that make the multi-server tool
1339

→ usage necessary

      - Do not contain the exact tool names or server names in the question
1340<sup>27</sup>
      - Ensure the question cannot be reasonably answered with tools from just

→ a single server

      - Create scenarios that naturally require different types of services
1343
       \hookrightarrow working together
134430
      ### Cross-Server Integration
134531
_{1346} - Think about how different servers' capabilities can be combined
33 - Consider data flow between different services (e.g., retrieving data

→ from one service to use in another)
1348_{34} - Create realistic scenarios where multiple services need to work
1349

→ together
```

- Focus on complementary functionalities across different domains

```
1350
36
1351
      ### Output Format
135238 Your response should include:
135339 1. **Server Analysis**: Briefly analyze all MCP Servers and their
       \hookrightarrow available tools, focusing on how they can work together.
1354
1355<sup>40</sup>
      2. **Cross-Server Workflow**: Describe the workflow showing how tools
       \,\,\hookrightarrow\,\, from different servers will be used together.
      3. **Target Tools**: The specific tool names from different MCP Servers
1357
          that should be used together, in the order they would likely be
1358
       \hookrightarrow called, with their server names.
       4. **Question**: A clear, realistic user question that requires
135942
       \hookrightarrow multi-server tool usage.
1360
136144
      ## Available MCP Servers
136245
      {SERVER DESCRIPTIONS}
136346
136447
1365<sup>48</sup>
      ## Output
      Ensure your question requires exactly {NUM_TOOLS} tools from at least 2
1366
       → different servers to solve completely. Provide your response in the
1367

→ following XML format:

136850
1369<sup>51</sup>
      <response>
1370<sup>52</sup>
        <server_analysis>
           <!-- Briefly analyze all MCP Servers and their available tools,
1371
           1372

→ domains/services. -->

137354
         </server_analysis>
137455
        <cross_server_workflow>
          <!-- Describe the workflow showing how tools from different servers
1375<sup>56</sup>
           → will be used together to solve the question. -->
1376<sub>57</sub>
        </cross_server_workflow>
137758
        <target_tools>
           <!-- The specific tool names from different MCP Servers that should
137859
           → be used together, listed in order with their server names. e.g.,
               <tool server="Server1">search_posts</tool> <tool</pre>
1380

    server="Server2">send_email</tool> -->

1381<sub>60</sub>
        </target_tools>
138261
         <question>
           <!-- A clear, realistic user question that requires multi-server tool
138362

→ usage spanning different services/domains. -->

1384
1385<sup>63</sup>
         </question>
1386
      </response>
1387
      Below is an example of a task generation prompt for the task synthesis for featured servers.
1388
1389 1 ## Task
1390 2 Generate a **Multi-Server Tool Use Question** based on featured MCP
       \rightarrow Servers and their tool descriptions.
1391
1392 4
      ## Objective
1393 5 Brainstorm a compelling real-world scenario, then analyze the provided
1394
      \hookrightarrow featured MCP Servers and their available tools to create a realistic
       → user question that would naturally require the use of **{NUM_TOOLS}
1395
       → tools from at least 2 different MCP servers** to solve completely.
1396
1397 7
      ## Guidelines
1398 <sub>8</sub>
1399 9
      ### Scenario Brainstorming
      - Think of realistic, specific scenarios where someone would need to use
140010
       → {NUM_TOOLS} different tools across multiple servers to accomplish a
1401

→ meaningful task

1402
      - Consider diverse real-world contexts such as:
1403<sub>12</sub>
        - Content creators managing their online presence across different
         → platforms
```

```
1404
         - Researchers gathering and analyzing information from multiple sources
1405
         - Developers building and deploying applications using different
1406

→ services

140715
         - Business professionals managing projects and communications across

→ platforms

1408
1409<sup>16</sup>
         - Students working on complex assignments requiring multiple tools
         - Entrepreneurs launching new ventures using various services
1410
      - The scenario should be detailed and authentic, representing genuine use
1411
       \hookrightarrow cases that span multiple services
141219
      ### Question Realism
1413^{20}
      - Create questions that represent real-world scenarios where users would
1414<sup>21</sup>

→ genuinely need tools from multiple MCP servers

1415,,
      - The question should sound natural and authentic, as if asked by someone
1416

→ with a specific goal

141723 - Include relevant context, constraints, and details that make the

→ question engaging

      - Consider workflows that require multiple complementary tools working
1419<sup>24</sup>
       \hookrightarrow together across different services
142025
      - Think about how different servers support each other in real-world use
1421

→ cases

142226
1423<sup>27</sup>
      ### Server and Tool Selection
      - Use tools from **at least 2 different MCP servers** to answer the
1424<sup>28</sup>
       \rightarrow question
      - Select **{NUM_TOOLS} tools total** that work together across multiple
1426

→ servers

142730 - The question should require a sequence or combination of tool calls
       \rightarrow from different servers to solve completely
1428
      - Choose tools based on how they complement each other across different
1429<sup>31</sup>

→ services/domains

143032
      - Consider each tool's description and purpose when crafting the
1431

→ cross-server workflow

      - Ensure tools from different servers create a logical, interconnected
143233
       → workflow
1433
1434<sup>34</sup>
      ### Question Complexity
1435<sub>36</sub>
      - Create questions that are complex enough to warrant using {NUM_TOOLS}
1436
       \hookrightarrow tools across multiple servers
143737 - The question should have multiple components or require several steps
       \hookrightarrow that span different services
1438
1439<sup>38</sup>
      - Include relevant context or constraints that make the multi-server tool

→ usage necessary

1440
      - Do not contain the exact tool names or server names in the question
144140
      - Ensure the question cannot be reasonably answered with tools from just
1442
       \hookrightarrow a single server
      - Create scenarios that naturally require different types of services
1443<sup>41</sup>
       → working together
1444
1445
      ### Cross-Server Integration
144644
      - Think about how different servers' capabilities can be combined
      - Consider data flow between different services (e.g., retrieving data
144745

→ from one service to use in another)

      - Create realistic scenarios where multiple services need to work
1449<sup>46</sup>

→ together

145047
      - Focus on complementary functionalities across different domains
1451<sub>48</sub>
      ### Output Format
145249
145350 Your response should include:
      1. **Server Analysis**: Briefly analyze the featured MCP Servers and
1454<sup>51</sup>
       \hookrightarrow their available tools, focusing on how they can work together.
      2. **Cross-Server Workflow**: Describe the workflow showing how tools
1456
       → from different servers will be used together.
1457
```

```
1458
      3. **Target Tools**: The specific tool names from different MCP Servers
1459
       \hookrightarrow that should be used together, in the order they would likely be
1460
       \hookrightarrow called, with their server names.
146154
      4. **Question**: A clear, realistic user question that requires
       → multi-server tool usage.
1462
1463<sup>55</sup>
       ## Available Featured MCP Servers
1464 57
146558
      {FEATURED_SERVER_DESCRIPTIONS}
146659
146760
      ## Output
      Ensure your question requires exactly {NUM_TOOLS} tools from at least 2
1468<sup>61</sup>
       \hookrightarrow different servers to solve completely. Provide your response in the
1469
          following XML format:
1470_{62}
147163
      <response>
        <server_analysis>
147264
           <!-- Briefly analyze the featured MCP Servers and their available
1473<sup>65</sup>
           \,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\,\, tools, focusing on how they can work together across different
1474

→ domains/services. -->

147566
      </server_analysis>
147667
      <cross_server_workflow>
           <!-- Describe the workflow showing how tools from different servers
147768
           → will be used together to solve the question. -->
1478
        </cross_server_workflow>
        <target_tools>
148071
           <!-- The specific tool names from different MCP Servers that should
           \hookrightarrow be used together, listed in order with their server names. e.g.,
1481
           1482
1483<sub>72</sub>

    server="Server2">send_email</tool> -->

        </target_tools>
1484
        <question>
148574
           <!-- A clear, realistic user question that requires multi-server tool
1486
           \rightarrow usage spanning different services/domains. -->
         </question>
1487<sup>75</sup>
      </response>
1488<sup>76</sup>
1489
      D.3 TASK DIVERSIFICATION PROMPT
1490
1491
      The following prompt aims to add diversity to the given task by introducing new contexts and per-
1492
      sonas.
1493
1494
    1 ## Task
1495 2 Generate **augmented variations** of a given question that maintain the
1496
      → same target tool(s) usage and complexity level but apply them across
1497
      \hookrightarrow different contexts and scenarios.
1498<sup>3</sup>
      ## Objective
1499 <sup>4</sup>
1500 5
      Take an existing question and its associated target tool(s), then create
       \hookrightarrow multiple variations that:
1501_6 - Use the same target tool(s) to achieve the core goal
15027 - Maintain the exact same tool usage order and final outcome
15038 - Apply the question to completely different contexts, scenarios, or

→ domains

1504
      - Keep the same level of complexity and constraints as the original
1505<sub>10</sub>
      - Demonstrate how the same tool usage pattern applies across diverse
1506

→ real-world scenarios

150711
      ## Guidelines
150812
      - Translate the question to distinctly different domains, user personas,
1509<sup>13</sup>
       → or situational contexts while preserving its original complexity
1510
       → level.
```

 1511 ₁₄ - Keep the tool usage sequence and final outcome identical across all

→ variations.

```
- Ensure each variation feels like a realistic scenario in its new
1513
       \hookrightarrow context and remains solvable with the same tool operations.
1514_{16} - Ensure the question does not contain any tool names or explicit
1515
       \rightarrow references to the target tools.
1516<sup>17</sup>
1517<sup>18</sup>
      ## Input Format
      **Original Question**: {ORIGINAL_QUESTION}
      **Target Tools**: {TARGET_TOOLS}
1519_{21}
      **Tool Descriptions**: {TOOL_DESCRIPTIONS}
152022
152123 ## Output Requirements
1522<sup>24</sup> Generate **{VARIATIONS_COUNT} augmented variations** of the original

→ question. Each variation should:

^{1523}_{25} 1. Maintain the same core goal that requires the target tool(s)
1524_{26} 2. Use the exact same tool(s) in the same order with the same final
1525

→ outcome

      3. Apply to a completely different context, scenario, or domain
1526<sup>27</sup>
1527<sup>28</sup>
       4. Keep the same complexity level and constraints as the original
      5. Feel like a natural, real-world scenario from a different setting
1528<sub>30</sub>
      6. Be meaningfully different from the original and other variations in
1529

→ terms of context only

153031 7. Avoid including any explicit mentions, hints, or references to the
       \rightarrow target tool names within the question text
1531
1532<sup>32</sup>
   33
      ## Output
1533
      Provide your response in the following XML format:
1534<sub>35</sub>
153536
      <response>
1536<sup>37</sup>
        <analysis>
           <!-- Briefly analyze the original question and target tool(s) to
1537<sup>38</sup>
            → understand the core goal, tool usage pattern, complexity level,
1538
           → and expected outcome, then identify how this can be applied
1539
            → across different domains while maintaining operational
1540

→ consistency -->

         </analysis>
1541<sup>39</sup>
1542<sup>40</sup>
         <variations>
          <!-- Generate {VARIATIONS_COUNT} variations, each with <variation_X>,
1543
           \hookrightarrow <context>, and <question> tags -->
154442
           <variation_1>
154543
             <cont.ext.>
1546<sup>44</sup>
               <!-- Brief description of the new domain/scenario introduced -->
1547<sup>45</sup>
             </context>
             <question>
154847
                <!-- The augmented question that maintains the same target
1549

→ tool(s) usage order, complexity, and outcome but in a
                \hookrightarrow different context -->
1550
             </question>
1551<sup>48</sup>
1552<sup>49</sup>
           </variation_1>
           <!-- Continue with variation_2, variation_3, etc. as needed based on
1553

→ number of variations -->

155451
         </variations>
155552
      </response>
1556
1557
       The prompt below is designed to enhance task complexity through the introduction of additional
1558
       constraints.
1559
1560<sub>1</sub> ## Task
1561 2 Generate **augmented variations** of a given question that maintain the
       → same target tool(s) usage and context but significantly increase the
1562
       → complexity and constraints required to solve the problem.
1563
1564
      ## Objective
1565<sub>5</sub>
      Take an existing question and its associated target tool(s), then create
       → multiple sophisticated variations that:
```

```
- Use the same target tool(s) to achieve the core goal while navigating
1567
       \hookrightarrow additional complexity layers
1568_{\,7} - Maintain the same general context and domain as the original question
1569 8 - Increase multi-dimensional complexity through realistic constraints,
       → competing requirements, stakeholder considerations, and
1570

→ interconnected dependencies

1571
      - Embed the tool usage within larger, more complex workflows that require
       \,\,\,\hookrightarrow\,\, strategic thinking and coordination
1573<sub>10</sub>
      - Demonstrate how the same core tool usage applies under vastly different
1574

→ complexity levels

157511
1576<sup>12</sup>
      ## Guidelines
1577
      - Introduce realistic constraints such as resource limits, compliance

→ requirements, tight timelines, or stakeholder conflicts

157814
      - Embed the same tool usage inside a broader workflow that requires
       \,\hookrightarrow\, coordination across teams or systems
1579
      - Escalate demands (performance, scalability, risk management) without
1580<sup>15</sup>

→ changing the original domain or context

1581
      - Ensure each variation targets a different primary complexity angle
1582
       → (organizational, technical, strategic) while preserving tool
1583

→ relevance

158417
      - Ensure the question does not contain any tool names or explicit
       \hookrightarrow references to the target tools.
1585
1586<sup>18</sup>
   19
       ## Input Format
1587
      **Original Question**: {ORIGINAL_QUESTION}
1588<sub>21</sub>
       **Target Tools**: {TARGET_TOOLS}
      **Tool Descriptions**: {TOOL_DESCRIPTIONS}
158922
1590^{23}
1591<sup>24</sup>
      ## Output Requirements
1592
      Generate **{VARIATIONS_COUNT} strategically augmented variations** of the
       → original question. Each variation should:
1593<sub>26</sub>
      1. Maintain the same core goal that requires the target tool(s) while
       \hookrightarrow adding multiple complexity layers
1594
       2. Keep the same general context and domain as the original question
1595<sup>27</sup>
       3. Introduce different but interconnected constraints and competing
1596<sup>28</sup>

→ requirements

1597,9
      4. Feel like natural, high-stakes, real-world scenarios that
1598

→ professionals encounter

      5. Be meaningfully different from the original and other variations in
159930

    → terms of complexity

1600
1601<sup>31</sup>
       6. Include specific details that make the constraints and requirements
       \hookrightarrow concrete and actionable
1602
       7. **Transform step-wise questions**: If the original question contains
1603

→ explicit steps, convert it to a goal-oriented format while

       \,\hookrightarrow\, maintaining the same tool usage requirements
1604
      8. Avoid including any explicit mentions, hints, or references to the
1605<sup>33</sup>
       \rightarrow target tool names within the question text
1606
1607
       ## Output
160836
      Provide your response in the following XML format:
160937
1610<sup>38</sup>
      <response>
1611<sup>39</sup>
         <analysis>
           <!-- Analyze the original question and target tool(s) to understand
1612
            \hookrightarrow the core goal, current complexity level, and identify multiple
1613
            → complexity dimensions that can be naturally introduced while
            → maintaining tool relevance and solution feasibility -->
1615<sup>41</sup>
         </analysis>
1616<sup>42</sup>
         <variations>
           <!-- Generate {VARIATIONS_COUNT} variations, each with <variation_X>,
1617

→ <constraints>, and <question> tags -->

161844
           <variation_1>
161945
             <constraints>
```

```
1620
               <!-- Specific organizational, stakeholder, or coordination
1621
               \hookrightarrow constraints that add realistic complexity -->
162247
             </constraints>
162348
             <question>
               <!-- The complex, organizationally-focused question that
162449
               → maintains the same target tool(s) usage within a more
1625

→ sophisticated workflow -->

             </guestion>
162751
           </variation_1>
162852
           <!-- Continue with variation_2, variation_3, etc. as needed based on
           \hookrightarrow number of variations -->
1629
1630<sup>53</sup>
        </variations>
      </response>
163155
1632
1633
1634
      D.4 TASK QUALITY ANNOTATION PROMPT
1635
1636
      ## Task
      Conduct a **Question Quality Assessment** of a tool use question across
1638
       \hookrightarrow six key dimensions to ensure it meets high standards for realistic
1639

→ tool usage scenarios.

1640<sup>3</sup>
1641 <sup>4</sup> ## Objective
5 Analyze the provided tool use question and assess its quality across six

→ primary dimensions:

1643 6 1. **Tool Selection Difficulty** - How challenging it is to determine
1644
       \rightarrow which tools to use giving all available tools
1645 7 2. **Tool Selection Uniqueness** - How unique and necessary the selected
       \hookrightarrow tools are for this specific task giving all available tools
1646
      3. **Question Quality** - Overall clarity, specificity, and effectiveness
1647 9
      4. **Scenario Realism** - How authentic and believable the scenario is
1648_{10} 5. **Verifiable** - How easy it is to verify the correctness of the final
       \hookrightarrow model answer
165011 6. **Stability** - How stable the answer will be when requested under
       \hookrightarrow different time and geolocation
1651
1652
      ## Assessment Criteria
1653<sub>14</sub>
165415
      ### 1. Tool Selection Difficulty
       **What to Evaluate**: How difficult it would be for a user to determine
1655<sup>16</sup>
       \rightarrow which specific tools are needed to solve this question.
1656
1657
      **Rating Guidelines**:
1658<sub>19</sub>
      - **very easy**: Question explicitly mentions tool names or makes tool
1659

→ selection obvious

      - **easy**: Tool selection is straightforward with clear indicators
1660^{20}
1661<sup>21</sup>
       - **medium**: Requires some reasoning but tool needs are fairly apparent
1662<sup>22</sup>
23
       - **hard**: Requires careful analysis to determine appropriate tools
       - **very hard**: Requires extensive expertise and deep reasoning to
1663
       \hookrightarrow identify the correct tools
166424
      ### 2. Tool Selection Uniqueness
166525
1666<sup>26</sup>
      **What to Evaluate**: How unique and necessary the selected tools are for

ightarrow accomplishing this specific task, and whether the task can only be
1667
       1668<sub>27</sub>
166928
      **Rating Guidelines**:
      - **not unique**: Many alternative tool combinations could accomplish the
1670<sup>29</sup>

→ same task equally well

      - **somewhat unique**: Some alternative approaches exist, but selected

→ tools offer advantages

1673<sub>31</sub>
      - **moderately unique**: Selected tools are well-suited, with limited

→ alternative approaches
```

```
1674
       - **quite unique**: Selected tools are particularly well-matched to the
1675

→ task requirements

167633
      - **highly unique**: Task can only be accomplished effectively with these
1677

→ specific tools in this sequence

167834
1679<sup>35</sup>
      ### 3. Question Quality
       **What to Evaluate**: Overall quality, clarity, and effectiveness of the
1680
       \rightarrow question as a realistic user query.
168137
168238
      **Rating Guidelines**:
      - **very poor**: Unclear, ambiguous, or poorly constructed question
168339
      - **poor**: Some clarity issues, missing important context
1684<sup>40</sup>
      - **average**: Clear and understandable, but could be more specific or
1685<sup>41</sup>

→ engaging

168642
      - **good**: Well-constructed, clear, specific, and realistic
168743
      - **excellent**: Exceptionally clear, detailed, engaging, and
       \hookrightarrow professionally written
1689<sup>44</sup>
   45
      ### 4. Scenario Realism
169046
      **What to Evaluate**: How authentic, believable, and true-to-life the
1691

→ described scenario is.

169247
1693<sup>48</sup>
      **Rating Guidelines**:
1694<sup>49</sup>
      - **unrealistic**: Artificial, contrived, or implausible scenario
      - **somewhat unrealistic**: Some realistic elements but feels forced or
1695

→ unlikely

1696<sub>51</sub>
      - **moderately realistic**: Believable scenario with minor authenticity
1697
       \hookrightarrow issues
      - **realistic**: Authentic scenario that represents genuine use cases
1698<sup>52</sup>
      - **highly realistic**: Completely natural, authentic scenario
1699<sup>53</sup>

→ indistinguishable from real user needs

170054
170155
      ### 5. Verifiable
170256
      **What to Evaluate**: How easy it is to verify the correctness of the

→ final model answer.

1703
1704<sup>57</sup>
      **Rating Guidelines**:
   58
1705<sub>59</sub>
      - **hard to verify**: Fully free-form answer that requires extensive
1706

→ human judgment

      - **somewhat hard**: Mostly subjective answer with some verifiable
170760

→ elements

1708
1709<sup>61</sup>
      - **moderately verifiable**: Short sentence that can be verified by LLM

→ comparison

1710<sub>62</sub>
      - **mostly verifiable**: Answer with clear, objective components and some
1711
       \hookrightarrow subjective elements
      - **easy to verify**: Answer can be verified by simple rules, exact
171263
       1713
1714<sup>64</sup>
      ### 6. Stability (1-5 Scale)
1715<sub>66</sub>
      **What to Evaluate**: How stable and consistent the answer will be when
1716
       → system contexts. Consider factors like temporal dependency,
1717
       \rightarrow geographical variations, operating system differences, network
1718
       \rightarrow environments, and software version variations.
1719
1720<sub>68</sub>
      **Rating Guidelines**:
1721<sub>69</sub>
      - **highly unstable**: Answer changes significantly across different

→ conditions (real-time data, location-specific, system-dependent)

1722
      - **somewhat unstable**: Answer may vary moderately based on
1723<sup>70</sup>
       \hookrightarrow environmental or system factors
1724
      - **moderately stable**: Answer mostly consistent with minor variations
1725
       \hookrightarrow due to context
172672 - **mostly stable**: Answer remains largely consistent across different
1727

→ conditions
```

```
1728
       - **highly stable**: Answer is completely independent of environmental
1729
       \hookrightarrow and system factors
173074
173175
       ## Question Analysis
1732<sup>76</sup>
       ### All Available Tools```
1733<sup>77</sup>
       {ALL_SERVER_AND_TOOL_INFORMATION}
1734
173580
173681
       ### Question Content
173782
1738<sup>83</sup>
       {QUESTION_CONTENT}
1739
1740<sub>86</sub>
      ### Intended Tool for This Question
174187
       {INTENDED_TOOL}
174288
1743<sup>89</sup>
1744
      ## Output Requirements
174592
174693 Provide analysis with detailed reasoning BEFORE scores for each of the
       \hookrightarrow six metrics.
1747
174894
       ## Output
1749<sub>96</sub>
      Provide your response in the following XML format:
1750<sub>97</sub>
175198
      <response>
        <tool_selection_difficulty>
175299
1753 100
           <reasoning>
              <!-- Detailed explanation including ambiguity level, domain
1754
              → knowledge required, and alternative solutions giving all
1755
              \hookrightarrow available tools -->
           </reasoning>
175@02
            <rating><!-- Rating: very easy, easy, medium, hard, very hard
175703
            1758
104
         </tool_selection_difficulty>
1759
105
176Q_{06}
        <tool_selection_uniqueness>
176107
            <reasoning>
              <!-- Detailed explanation of tool necessity, sequential
1762^{08}
              → dependencies, and alternative tool viability giving all
1763
              \hookrightarrow available tools -->
            </reasoning>
176 \bar{\phi}_{10}
            <rating><!-- Rating: not unique, somewhat unique, moderately unique,

→ quite unique, highly unique --></rating>

1766
         </tool_selection_uniqueness>
176711
1768^{12}
         <question_quality>
1769
114
            <reasoning>
177Q_{15}
              <!-- Detailed explanation covering linguistic quality, information
              \rightarrow architecture, and actionability -->
1771
            </reasoning>
177216
            <rating><!-- Rating: very poor, poor, average, good, excellent
177317

→ --></rating>

1774<sub>118</sub>
         </question_quality>
1775_{19}
         <scenario_realism>
177£20
            <reasoning>
177721
              <!-- Detailed explanation of industry authenticity, workflow
1778^{122}
              \rightarrow accuracy, and stakeholder behavior -->
1779_{123}
            </reasoning>
178 \varrho_{24}
            <rating><!-- Rating: unrealistic, somewhat unrealistic, moderately
1781
            → realistic, realistic, highly realistic --></rating>
         </scenario_realism>
  125
```

```
1782
126
178\frac{1}{127}
         <verifiable>
178428
           <reasoning>
178529
             <!-- Detailed explanation of answer format, objective criteria, and

→ ground truth availability -->

1786
1787 130
           </reasoning>
           <rating><!-- Rating: hard to verify, somewhat hard, moderately
1788
           → verifiable, mostly verifiable, easy to verify --></rating>
1789_{32}
         </verifiable>
179 133
        <stability>
179134
1792<sup>135</sup>
          <reasoning>
             <!-- Detailed explanation of temporal/geographical/system
1793
              \rightarrow dependencies and environmental factors -->
179437
           </reasoning>
179538
           <rating><!-- Rating: highly unstable, somewhat unstable, moderately

→ stable, mostly stable, highly stable --></rating>

1797 139
         </stability>
1797
1798
1798
       </response>
1799
1800
       D.5 TRAJECTORY ANNOTATION PROMPT
1801
1802
      ## Task
1803 <sup>1</sup>
      Conduct a **Response Quality Assessment** of a tool-use conversation
1804

ightarrow across two LLM-scored dimensions, with a third dimension computed
1805
       → automatically outside the LLM.
1806 3
1807 4
      ## Objective
1808 5 Analyze the provided conversation and assess its response quality across
       \hookrightarrow two primary dimensions scored by the LLM, while reserving an
1809
       → additional tool-call accuracy dimension for automated scoring:
^{1810}_{\,\,6} 1. Completeness - Whether the assistant fully accomplished the user's

→ request end-to-end

_{1812} 7 2. Conciseness - Whether the assistant solved the task using the minimum

→ necessary steps and verbosity

1813
1814 9
      ## Assessment Criteria
1815<sub>10</sub>
181611 ### 1. Completeness
      **What to Evaluate**: Did the assistant fully satisfy the user's goal
181712
       → given the conversation context? Consider whether the assistant:
1818
      - Executed all required steps end-to-end (including

→ saving/exporting/downloading where applicable)

^{1820}{}_{14} - Provided the final deliverable or a working alternative when blocked
1821

→ (e.g., tool failure with a usable fallback)
      - Included essential confirmations, paths, or instructions to achieve the
182215

→ out.come

1823
1824
      - Avoided missing key requirements or leaving the user with unresolved

→ gaps

1825<sub>17</sub>
      **Rating Guidelines**:
      - very incomplete: Major requirements missing; no usable outcome
182719
1828<sup>20</sup> - incomplete: Some key requirements missing; outcome is not directly

→ usable

1829
      - partially complete: Core steps attempted; outcome usable only with user
       \rightarrow effort or missing minor requirements
183122 - mostly complete: Meets most requirements; small omissions or minor

→ issues remain

1832
      - fully complete: All requirements met with a usable outcome delivered
1833<sup>23</sup>
1834 25
      ### 2. Conciseness
1835<sub>26</sub>
      **What to Evaluate**: Did the assistant achieve the goal with minimal
       → redundancy and steps? Consider whether the assistant:
```

```
- Avoided repetitive or unnecessary explanations/tool calls
1837<sub>28</sub>
      - Used the minimal set of steps/tools to complete the task
183829
       - Kept language concise while preserving clarity
183930
      **Rating Guidelines**:
1840<sup>31</sup>
1841.32
      - very redundant: Excessive repetition or unnecessary steps/tool calls
      - redundant: Noticeable verbosity or extra steps beyond what's needed
1842<sub>34</sub>
      - average: Reasonably concise with minor extraneous content
184335
       - concise: Efficient and to the point with minimal overhead
184436
       - very concise: Maximally efficient while clear and complete
1845<sup>37</sup>
1846<sup>38</sup>
       ## Response Analysis
184740
       ### Question Content
184841
184942
       {QUESTION_CONTENT}
1850<sup>43</sup>
1851<sub>45</sub>
       ### Intended Tool for This Question
1852<sub>46</sub>
185347
       {INTENDED_TOOL}
185448
1855<sup>49</sup>
1856<sup>50</sup>
       ### Conversation History
1857<sub>52</sub>
       {CONVERSATION_HISTORY}
185853
185954
1860<sup>55</sup>
       ## Output Requirements
       - Provide detailed reasoning BEFORE ratings for Completeness and
1861<sup>56</sup>

→ Conciseness

186257
      - Do NOT score Tool Call Accuracy; include placeholders only
186358
       ## Output
186459
       Provide your response in the following XML format:
1865<sup>60</sup>
1866<sup>61</sup>
      <response>
1867<sub>63</sub>
        <completeness>
186864
            <reasoning>
              <!-- Evaluate if the assistant delivered an end-to-end usable
186965
              \hookrightarrow outcome, addressed all requirements, handled tool failures with
1870
              → alternatives, and provided necessary confirmations/paths. -->
            </reasoning>
            <rating><!-- Rating: very incomplete, incomplete, partially complete,
1873
            → mostly complete, fully complete --></rating>
187468
         </completeness>
1875<sup>69</sup>
1876_{71}^{70}
         <conciseness>
1877<sub>72</sub>
              <!-- Evaluate if the assistant minimized redundant
1878
              \rightarrow messaging efficient while clear. -->
1879
            </reasoning>
1880<sup>73</sup>
1881<sup>74</sup>
            <rating><!-- Rating: very redundant, redundant, average, concise,

→ very concise --></rating>

1882<sub>75</sub>
         </conciseness>
188376
       </response>
1884
1885
1886
1887
1888
1889
```