

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

**Anonymous authors**

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 TOUCAN 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 M S IR	5,000	From $\tau$ -Bench	Executed
ToolACE (Liu et al., 2025a)	11,300	S P M S IR	509	Synthetic	Simulated
Hermes Function-Calling V1 (interstellarninja)	11,570	S P M S IR	1,890	Synthetic	Executed
Nemotron (Tools) (Nathawani et al., 2025)	310,051	S P M S –	199,610	–	–
TOUCAN (This Work)	1,527,259	S P M S 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  $\tau$ -Bench (Yao et al., 2024) and  $\tau^2$ -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 tool-agent 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),  $\tau$ -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  $\tau$ -Bench

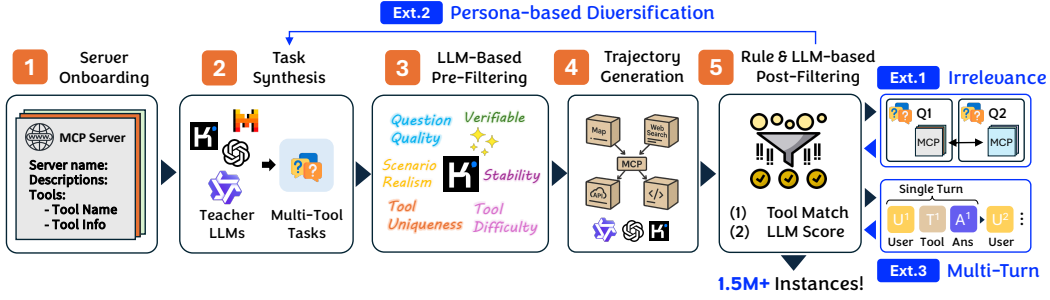


Figure 2: The TOUCAN construction pipeline: A systematic five-stage process from MCP server onboarding through trajectory filtering, with three extensions for enhancing data diversity and realism.

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, simulate tool responses via LLMs, or focus on VLMs rather than LLMs Gao et al. (2025b). Table 1 compares existing tool-agentic datasets for LLMs 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., 2025a; 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., 2025a) 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 meta-data. 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.

**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<sup>1</sup>, 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

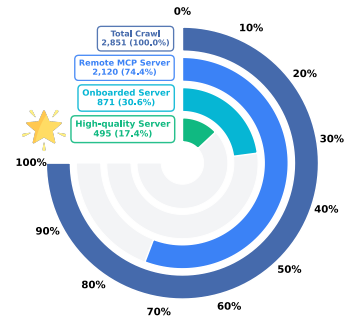


Figure 1: MCP servers filtering process

<sup>1</sup><https://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.

**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.

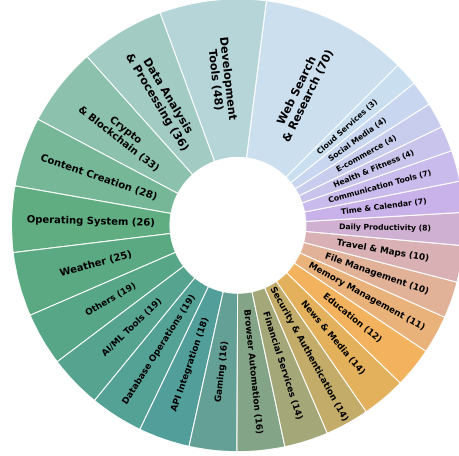


Figure 3: MCP servers distribution by domain, covering a wide range of categories. Values in parentheses indicate the number of servers belonging to each category.

**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 [all](#) 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-

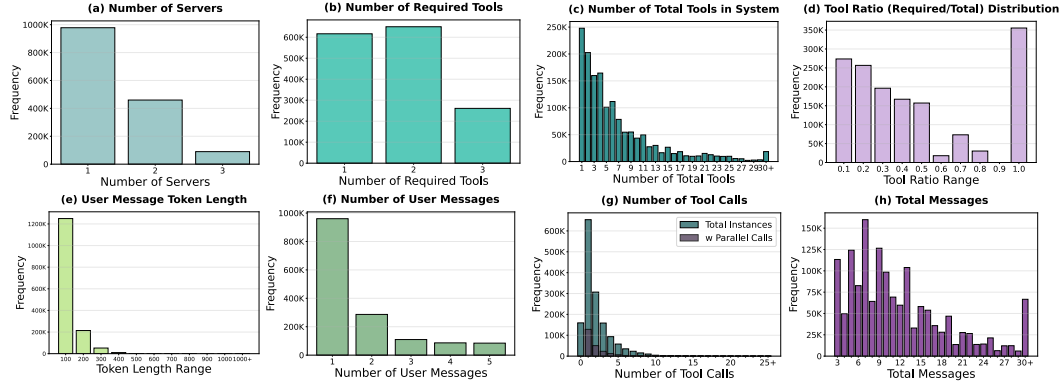


Figure 4: The figures above illustrate the TOUCAN dataset analysis. Subfigure (a) and (b) provide statistics on the number of servers and required tools per instance, highlighting TOUCAN’s comprehensive coverage of multi-server and multi-tool tasks. Subfigures (c) and (d) reveal that most tasks include 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.

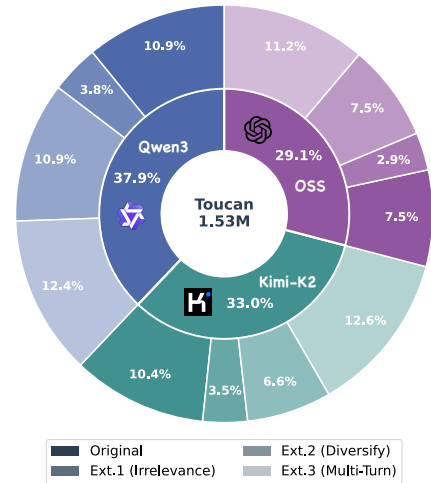


Figure 5: TOUCAN Subset Statistics

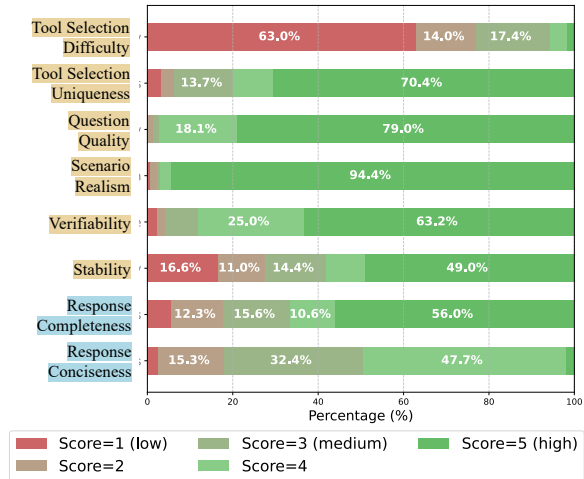


Figure 6: TOUCAN Quality Statistics



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%	<b>83.33%</b>	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%	<b>83.33%</b>	76.32%
Qwen3-32B	69.25%	<b>88.90%</b>	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%	<b>79.92%</b>	40.50%	77.78%	<b>85.95%</b>
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% <sup>+3.16%</sup>	78.52%	74.50%	22.62%	66.67%	75.18%
Qwen2.5-14B-Instruct	57.69%	83.38%	73.70%	19.75%	<b>83.33%</b>	68.46%
with TOUCAN	65.09% <sup>+7.40%</sup>	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	<b>70.45%</b> <sup>+8.72%</sup>	87.12%	78.90%	<b>46.50%</b>	77.78%	78.10%
Llama-3.1-8B-Instruct	26.23%	47.96%	33.63%	6.38%	94.44%	5.26%
with TOUCAN	58.46% <sup>+32.23%</sup>	83.44%	70.68%	24.88%	77.78%	64.85%
Llama-3.3-70B-Instruct	53.03%	85.23%	62.86%	16.38%	100.00%	48.50%
with TOUCAN	66.20% <sup>+13.17%</sup>	85.79%	73.48%	42.25%	77.78%	68.22%

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

Model	$\tau$ -bench			$\tau^2$ -bench			
	Avg.	Airline	Retail	Avg.	Airline	Retail	Telecom
Qwen2.5-7B-Instruct	15.03%	8.75%	21.30%	16.08%	14.00%	17.54%	16.70%
with TOUCAN	22.48% <sup>+7.45%</sup>	15.50%	29.46%	17.77% <sup>+1.69%</sup>	20.00%	22.80%	10.50%
Qwen2.5-14B-Instruct	30.85%	17.25%	44.46%	24.46%	12.00%	41.20%	20.18%
with TOUCAN	35.24% <sup>+4.39%</sup>	22.00%	48.48%	30.43% <sup>+5.97%</sup>	22.00%	49.10%	20.18%
Qwen2.5-32B-Instruct	38.76%	26.00%	51.52%	29.40%	18.00%	49.10%	21.11%
with TOUCAN	42.33% <sup>+3.57%</sup>	29.00%	55.65%	31.60% <sup>+2.20%</sup>	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 Qwen2.5-7B-Instruct, Qwen2.5-14B-Instruct, and Qwen2.5-32B-Instruct (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 DeepSeek-V3 DeepSeek-AI et al. (2025), Qwen2.5-72B-Instruct, Qwen3-235B-A22B, Qwen3-32B 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

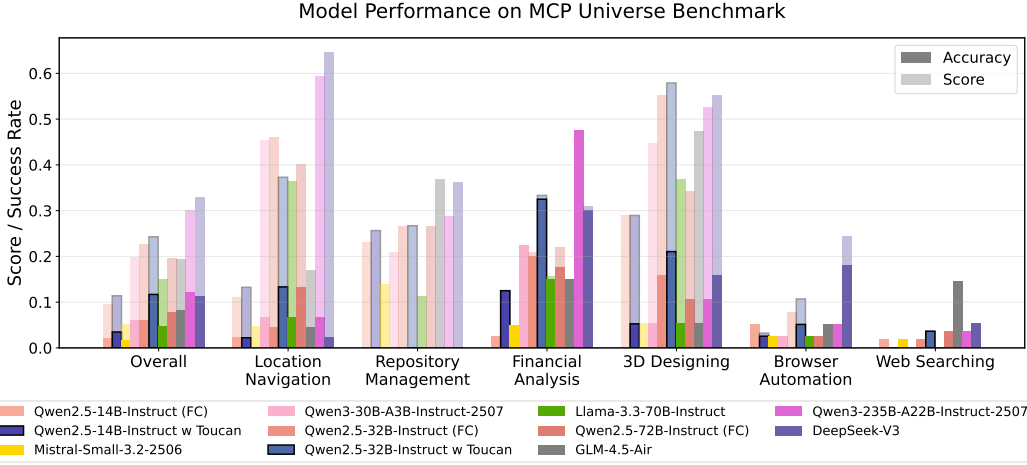


Figure 7: This figure compares the performance of TOUCAN-tuned models with other open-source models on MCP-Universe (Luo et al., 2025). Model sizes increase from left to right. Bars with darker colors represent task success rate (full task completion), while lighter colors represent average evaluation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

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 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),  $\tau$ -Bench Yao et al. (2024),  $\tau^2$ -Bench (Barres et al., 2025), and MCP-Universe Luo et al. (2025). All evaluations are conducted on an  $8 \times \text{H100}$  server. For BFCL-V3, we use the official evaluation setup. For  $\tau$ -Bench and  $\tau^2$ -Bench, we employ GPT-4o 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,  $\tau$ -Bench, and  $\tau^2$ -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 mod-

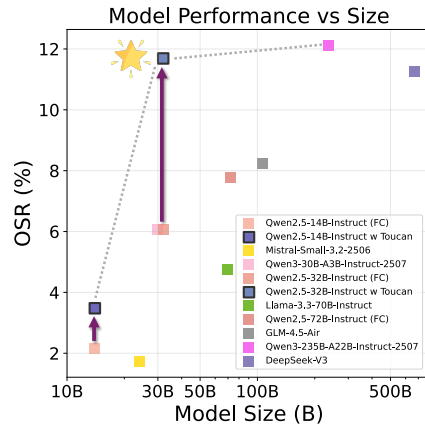


Figure 8: Model Performance vs Size on MCP-Universe Benchmark. We report overall task success rate (OSR). Our models achieve a new Pareto optimum.



els and other open-source models of similar or larger 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, TOUCAN-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 Figure 9. We observe that all components contribute to improved scores. Detailed benchmark scores for the BFCL ablation study are provided in Appendix C.3. [In addition, we include further ablations on tool scaling, dataset scaling, fine-tuning comparisons between TOUCAN and other baseline datasets, as well as trajectory annotation and filtering, in Appendix C.4–C.8.](#)

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

	BFCLv3	$\tau$ -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, extensions of our data analysis and fine-tuning experiments, [as well as compute requirements](#). 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.  $\tau^2$ -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. Api-blend: A comprehensive corpora for training and benchmarking api llms, 2024. URL <https://arxiv.org/abs/2402.15491>.

- Fouad Boussetouane. 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, Cheng-gang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, 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, Miaojuan 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, Shut-ing 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, Xi-aokang 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, Yudian 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 llm reasoning to au-tonomous 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, 2025a. URL <https://arxiv.org/abs/2505.16700>.
- Zhi Gao, Bofei Zhang, Pengxiang Li, Xiaojian Ma, Tao Yuan, Yue Fan, Yuwei Wu, Yunde Jia, Song-Chun Zhu, and Qing Li. Multi-modal agent tuning: Building a vlm-driven agent for efficient tool usage, 2025b. URL <https://arxiv.org/abs/2412.15606>.
- 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. Mcp-agentbench: Evaluating real-world language agent performance with mcp-mediated tools, 2025b. URL <https://arxiv.org/abs/2509.09734>.
- Teknium interstellarninja. Hermes function calling dataset v1. URL <https://huggingface.co/NousResearch/hermes-function-calling-v1>.

- 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 Awalgaoonkar, 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, Hui-long Chen, Ji Li, Jiajing Zhao, Jiamin Ren, Jian Jiao, Jiani Zhao, Jianshan 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, Sibao Yi, Tianshu Yu, Wei Tian, Weihang Wang, Wenbo Yu, Weng Lam Tam, Wenjie Liang, Wentao Liu, Xiao Wang, Xiaohan Jia, 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, Jiale Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Keli 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 tool-using 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.  $\tau$ -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

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 :

```
{
  "uuid": "3ac8fdcc-b9b5-50d2-a840-947a42b558d2",
  "subset": "single-turn-original",
  "messages": "[{...long JSON string of messages...}]",
  "question": "I'm reviewing my company's policy document that needs
  ↳ 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
  ↳ 2025' throughout the document, and finally add explanatory
  ↳ footnotes to the last 3 sections that reference our new
  ↳ regulatory guidance document published last month. Can you help
  ↳ me accomplish these updates systematically?",
  "target_tools": "get_document_outline, search_and_replace,
  ↳ add_footnote_to_document",
  "question_quality_assessment_kimik2w4a16": {
    "tool_selection_difficulty": { "reasoning": "...", "score": 3 },
    "tool_selection_uniqueness": { "reasoning": "...", "score": 3 },
    "question_quality": { "reasoning": "...", "score": 4 },
    "scenario_realism": { "reasoning": "...", "score": 5 },
    "verifiable": { "reasoning": "...", "score": 4 },
    "stability": { "reasoning": "...", "score": 5 },
    "overall_score": 4.0
  },
  "response_quality_assessment_oss120b": {
    "completeness": { "reasoning": "...", "score": 2 },
    "conciseness": { "reasoning": "...", "score": 3 },
    "overall_score": 2.5,
    "desired_tools_used_percentage": 0.3333,
    "order_correctness": false
  },
  "messages_num_rounds": 9,
  "metadata": "{...long metadata JSON string...}"
}
```

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.

```
[
  {
    "role": "system",
    "content": "{... contains the list of tools ... }"
  },
  {
    "role": "user",
    "content": "Can you explain what ENS (Ethereum Name Service)
↳ represents in the context of a transaction and why a name like
↳ 'blockscout.eth' might appear instead of a raw address?"
  },
  {
    "role": "assistant",
    "content": "I'll explain what ENS (Ethereum Name Service) represents,
↳ and show you how it resolves a name like 'blockscout.eth' into
↳ its actual Ethereum address."
  },
  {
    "role": "assistant",
    "content": "",
    "function_call": {
      "name": "blockscout-mcp-server-__unlock_blockchain_analysis__",
      "arguments": "{}"
    }
  },
  {
    "role": "function",
    "content": "{... long tool output ...}"
  },
  {
    "role": "assistant",
    "content": "Now, let me resolve 'blockscout.eth' to show the actual
↳ address it represents."
  },
  {
    "role": "assistant",
    "content": "",
    "function_call": {
      "name": "blockscout-mcp-server-get_address_by_ens_name",
      "arguments": "{\"name\": \"blockscout.eth\"}"
    }
  },
  {
    "role": "function",
    "name": "blockscout-mcp-server-get_address_by_ens_name",
    "content": "{\n  \"data\": {\n    \"resolved_address\":\n↳ \"0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975\"\n  },\n  \"data_description\": null,\n  \"notes\": null,\n  \"instructions\": null,\n  \"pagination\": null\n}"
  },
  {
    "role": "assistant",
    "content": "{... parsed assistant response ...}"
  },
  {
    "role": "user",
    "content": "Based on that explanation, could you look up the ENS name
↳ 'blockscout.eth' and provide the Ethereum address it resolves
↳ to?"
  },
  {
    "role": "assistant",
```

```

864     "content": "I already resolved that for you in my previous response!
865     ↳ When I called the ENS lookup tool, `\"blockscout.eth\"` resolved
866     ↳ to the Ethereum
867     ↳ address:\n\n`0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975`"
868   }
869 ]

```

In the previous trajectory example, we abbreviated the assistant's parsed response for readability. The full response is shown below.

```

873 1 ## What ENS (Ethereum Name Service) Represents in Transactions
874 2
875 3 ENS stands for **Ethereum Name Service**, and it functions much like "DNS
876 ↳ for Ethereum addresses." It allows users to register human-readable
877 ↳ names (like "blockscout.eth") that map to complex hexadecimal
878 ↳ blockchain addresses.
879 4
880 5 ### Key Points About ENS:
881 6
882 7 1. **Human-Readable Names for Complex Addresses**
883 ↳ Instead of displaying `0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975`
884 ↳ (which is the raw address for blockscout.eth), interfaces show the
885 ↳ friendly "blockscout.eth"
886 ↳ This reduces errors from copy/pasting long hex addresses
887 8
888 9 2. **Reverse Resolution**
889 ↳ When you see "blockscout.eth" in transaction data, wallet labels, or
890 ↳ block explorers, it means:
891 ↳ ENS has been configured with this name
892 ↳ The system has performed a reverse lookup to find the name
893 ↳ associated with the address
894 10
895 11 3. **Self-Sovereign Identity**
896 ↳ ENS names are owned and controlled by users through smart contracts
897 ↳ They can be transferred, updated, or associated with additional
898 ↳ metadata (social profiles, avatars, etc.)
899 12
900 13 ### Why "blockscout.eth" Appears Instead of Raw Address:
901 14
902 15 - **User Experience**: "blockscout.eth" is immediately recognizable and
903 ↳ memorable
904 16 - **Trust Verification**: Users can verify they're interacting with the
905 ↳ legitimate Blockscout entity rather than a random address
906 17 - **Branding**: Organizations use ENS names for consistent branding
907 ↳ across DeFi protocols and wallet interfaces
908 18 - **Error Prevention**: Reduces costly mistakes when users verify
909 ↳ transaction recipients
910 19
911 20 Every time you see an ENS name like "blockscout.eth" in a transaction,
912 ↳ the underlying system has actually resolved this to the precise
913 ↳ blockchain address `0x441e1e47a6Fa2DBfD3cd9b54291E9AB3a58D7975`-the
914 ↳ ENS simply provides a human-friendly layer on top of the raw
915 ↳ cryptographic identity.

```

Below example demonstrated an error recovery of TOUCAN. In this example, the assistant attempts to invoke a math evaluation tool multiple times with incorrect function syntax (*exp*, *e*), receives execution errors from the tool, and then recovers by reformulating the expression using a numeric approximation of *e* to successfully complete the calculation.

```

914 [
915   {
916     "role": "system",
917     "content": "{... contains the list of tools ... }"
918   },
919   {

```

```

918     "role": "user",
919     "content": "Calculate the company's revenue in year\u202f5 using the
920     ↪ exponential model Revenue\u202f=\u202f10,000\u202f\u00d7\u202fe^j
921     ↪ (0.15\u202f\u00d7\u202fyear). Show your work and give the
922     ↪ numerical result."
923 },
924 {
925     "role": "assistant",
926     "content": "I'll calculate the company's revenue in year 5 using the
927     ↪ exponential growth model you provided. Let me evaluate this
928     ↪ expression."
929 },
930 {
931     "role": "assistant",
932     "content": "",
933     "function_call": {
934         "name": "expression-evaluation-server-eval",
935         "arguments": "{\u201cexpression\u201c: \u201c10000 * exp(0.15 * 5)\u201c}"
936     }
937 },
938 {
939     "role": "function",
940     "content": "handler returned an error: unknown name exp (1:9)\n |
941     ↪ 10000 * exp(0.15 * 5)\n | .....^",
942     "name": "expression-evaluation-server-eval"
943 },
944 {
945     "role": "assistant",
946     "content": "Let me try with the proper function name:"
947 },
948 {
949     "role": "assistant",
950     "content": "",
951     "function_call": {
952         "name": "expression-evaluation-server-eval",
953         "arguments": "{\u201cexpression\u201c: \u201c10000 * e^(0.15 * 5)\u201c}"
954     }
955 },
956 {
957     "role": "function",
958     "content": "handler returned an error: unknown name e (1:9)\n | 10000
959     ↪ * e^(0.15 * 5)\n | .....^",
960     "name": "expression-evaluation-server-eval"
961 },
962 {
963     "role": "assistant",
964     "content": "Let me calculate this step by step, using the fact that e
965     ↪ \u2248 2.718281828:"
966 },
967 {
968     "role": "assistant",
969     "content": "",
970     "function_call": {
971         "name": "expression-evaluation-server-eval",
972         "arguments": "{\u201cexpression\u201c: \u201c10000 * 2.718281828^(0.15 *
973     ↪ 5)\u201c}"
974     }
975 },
976 {
977     "role": "function",
978     "content": "10000 * 2.718281828^(0.15 * 5) = 21170.00016344546",
979     "name": "expression-evaluation-server-eval"
980 },
981 {
982     "role": "assistant",

```

```

972     "content": "## Solution\n\n**Given Model:** Revenue = 10,000 \u00d7
973     ↪ e(0.15 \u00d7 year)\n\n**For Year 5:**\n\n1. **Calculate the
974     ↪ exponent:** 0.15 \u00d7 5 = 0.75\n\n2. **Calculate e0.75:0.75 \u2248 2.117\n\n3. **Calculate final revenue:** 10,000
976     ↪ \u00d7 2.117 = **$21,170.00**\n\n**Numerical Result:** The
977     ↪ company's revenue in year 5 is **$21,170** (rounded to the
978     ↪ nearest dollar).\n\nNote: This represents approximately 111.7%
979     ↪ growth over the starting revenue of $10,000, demonstrating the
980     ↪ power of exponential growth at 15% per year compounded
981     ↪ continuously."
981 },
982 ]
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

```

## 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 <sup>2</sup>	Web Search & Research	10
Book Search Server	Access via Smithery <sup>3</sup>	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
Pok��mcp	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	ceydasmsekk/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.

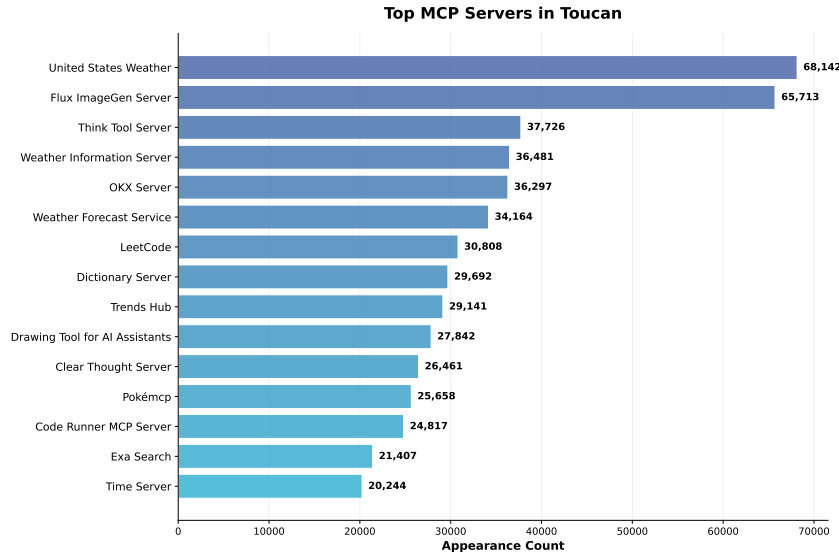


Figure 10: Distribution of the most frequently occurring MCP servers in the TOUCAN dataset.





#### B.4 DOMAIN COVERAGE COMPARISON BETWEEN TOUCAN AND MCP UNIVERSE

In what follows, we compare the MCP server names included in TOUCAN with those used for constructing MCP Universe. The results are summarized in Table 5. Our analysis shows that four MCP Universe domains are completely out-of-distribution (OOD) with respect to TOUCAN, indicating that our fine-tuned models demonstrate strong generalization performance on domains that were never seen during training.

Table 5: In-distribution (ID) and out-of-distribution (OOD) domain coverage of TOUCAN relative to MCP Universe.

Benchmark	Benchmark Domain	TOUCAN
MCP Universe	Location Navigation	OOD
	Repository Management	OOD
	Financial Analysis	ID
	3D Design	OOD
	Browser Automation	ID
	Web Searching	OOD

### C MORE ON EXPERIMENTS

#### C.1 LLM ANNOTATION

Figure 13 shows the Pearson correlation between human annotations and LLM-as-a-judge evaluations across different models on 50 randomly sampled instances. We observe that GPT-4.1 and Kimi-K2 achieve the highest correlation with human judgments. Notably, the Pearson correlation between the two human annotators is only 0.5028, indicating moderate inter-annotator agreement, especially on inherently subjective aspects such as tool selection uniqueness and scenario realism. This suggests that the relatively low model-human correlation is partly due to annotation subjectivity rather than model unreliability. Considering cost efficiency and maintaining an end-to-end open-source pipeline, we deploy Kimi-K2 locally as the annotator. Our annotation prompt is available in Appendix D.4.

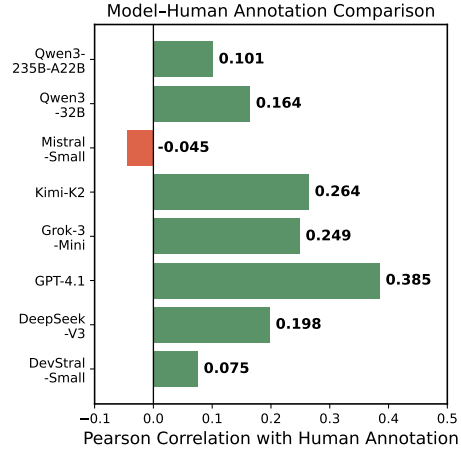


Figure 13: Pearson correlation between human annotator and LLM-as-a-Judge evaluations across different models.

#### 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 6.

Table 6: This table shows the hyper-parameters for supervised fine-tuning.

Hyper-parameter	Value
Tool-Call Template	Hermes
Learning Rate	$2 \times 10^{-5}$
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 ABLATION STUDIES ON DATA EXTENSIONS

Table 7 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 7: Ablation of TOUCAN Extensions on BFCL V3 Benchmark.

	Overall	Single Turn		Multi Turn	Hallucination	
		<i>Non-live (AST)</i>	<i>Live (AST)</i>		<i>Relevance</i>	<i>Irrelevance</i>
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%

### C.4 ABLATION STUDIES ON TOOL SCALING

In this experiment, we create five subsets of TOUCAN datasets, each doubling the number of tools relative to the previous one. The number of tools preserved in the dataset ranges from 100 to 1,600, while the total number of trajectories is kept constant at 20,000. We fine tuned Qwen2.5-32B-Instruct on these subsets, and evaluate the resulting models on BFCL V3 and  $\tau$ -Bench. Table 8 shows the results of this experiment. We observe a consistent upward trend in overall performance as tool diversity increases, indicating that a larger and more diverse tool set leads to better generalization rather than redundant learning.

Table 8: Ablation on Tool Diversity on BFCL V3 and  $\tau$ -Bench.

Model Variant	BFCL-V3				$\tau$ -Bench	
	Overall	Single Turn		Multi Turn	Hallucination	
		<i>Non-live (AST)</i>	<i>Live (AST)</i>		<i>Relevance</i>	<i>Irrelevance</i>
Qwen2.5-32B-Instruct-Toucan-100Tools-20K	60.38%	87.58%	64.82%	36.00%	<b>88.89%</b>	46.66%
Qwen2.5-32B-Instruct-Toucan-200Tools-20K	60.90%	86.56%	65.44%	37.00%	<b>88.89%</b>	49.50%
Qwen2.5-32B-Instruct-Toucan-400Tools-20K	61.99%	87.48%	65.08%	40.00%	83.33%	49.00%
Qwen2.5-32B-Instruct-Toucan-800Tools-20K	61.73%	<b>87.31%</b>	64.73%	<b>40.38%</b>	83.33%	45.98%
Qwen2.5-32B-Instruct-Toucan-1600Tools-20K	<b>62.26%</b>	86.27%	<b>67.57%</b>	39.38%	83.33%	<b>52.08%</b>

### C.5 ABLATION STUDIES ON LARGE DATASETS

We conduct an ablation experiment to investigate the relevance of TOUCAN’s large size for the research community. Specifically, we created two training datasets: TOUCAN -Full (including 1.5M trajectories), and TOUCAN -SFT (as detailed in Section 4.1). We fine-tuned Qwen2.5-14B-Instruct with each training dataset. Table 9 shows the evaluation results for BFCL V3. Overall, TOUCAN -Full slightly outperforms the SFT subset, and shows a remarkable improvement in the multi-turn setting. Our results also show that the model fine-tuned on the full dataset achieves a lower score on the Irrelevance setting, which suggests that the carefully rebalanced TOUCAN -SFT is more effective at reducing hallucinations.

Table 9: BFCL-V3 Results for TOUCAN Full and SFT datasets.

Dataset	Overall	Single Turn		Multi Turn	Hallucination	
		<i>Non-live (AST)</i>	<i>Live (AST)</i>		<i>Relevance</i>	<i>Irrelevance</i>
Qwen2.5-14B-Instruct	57.69%	83.38%	73.70%	19.75%	83.33%	68.46%
with Toucan-SFT	<b>65.09%</b>	<b>85.42%</b>	<b>76.01%</b>	35.25%	72.22%	<b>75.96%</b>
with Toucan-Full	<b>65.17%</b>	84.90%	74.63%	<b>39.13%</b>	<b>83.33%</b>	68.71%

### C.6 ABLATION STUDIES ON DATA SCALING

We perform a data-scale ablation by randomly sampling subsets of 20K, 40K, 60K, 80K, and 100K trajectories from TOUCAN -SFT (see Section 4.1). We then fine-tune Qwen2.5-32B-Instruct

on each subset and compare the results against training on the full dataset. We evaluated the results on the BFCL-V3 benchmark. Table 10 shows the results of this experiment. Overall, the models show a consistent performance gain as the data scale increases, with especially strong improvements in the multi-turn setting. We also observe diminishing returns and near-saturation beyond approximately 80K trajectories. This behavior mirrors common scaling trends in instruction tuning and suggests that our rebalanced subset already provides an effective cost-performance sweet spot.

Table 10: Ablation on SFT Data Scale on BFCL V3 and  $\tau$ -Bench.

Model Variant	BFCL-V3				$\tau$ -Bench			
	Overall	Single Turn		Multi Turn	Hallucination		Airline @1	Retail @1
		Non-live (AST)	Live (AST)		Relevance	Irrelevance		
Qwen2.5-32B-Instruct	61.73%	85.58%	76.01%	26.38%	72.22%	72.68%	26.00%	51.52%
Toucan-SFT-20K	68.21%	88.52%	74.99%	42.50%	83.33%	74.73%	27.75%	47.50%
Toucan-SFT-40K	68.82%	86.77%	77.30%	43.50%	77.78%	76.62%	28.50%	53.37%
Toucan-SFT-60K	68.55%	86.71%	77.08%	43.12%	83.33%	75.87%	28.25%	56.30%
Toucan-SFT-80K	69.62%	87.02%	77.65%	45.25%	77.78%	77.23%	30.00%	58.26%
Toucan-SFT-100K	69.83%	86.44%	78.76%	45.25%	77.78%	77.91%	28.00%	56.66%
Toucan-SFT-119K	<b>70.45%</b>	<b>87.12%</b>	<b>78.90%</b>	<b>46.50%</b>	77.78%	<b>78.10%</b>	29.00%	55.65%

### C.7 ABLATION STUDIES WITH COMPARABLE TOOL-CALLING DATASET

We perform a controlled experiment to compare TOUCAN with Nemotron-SFT (tool subset) under a similar data scale (see Table 1) and using the same model, Qwen2.5-14B-Instruct, as baseline. While the overall score of Nemotron-SFT achieves a comparable performance to TOUCAN on BFCL V3, our dataset shows substantially stronger performance on other tool-agentic benchmarks, especially  $\tau$ -Bench and  $\tau^2$ -bench. These benchmarks better reflect multi-tool reasoning and cross-domain generalization. Table 11 and Table 12 respectively report the results obtained.

Table 11: Comparison between Nemotron-SFT (tool subset) and TOUCAN on BFCL-V3.

Model	Overall	Single Turn		Multi Turn	Hallucination	
		Non-live (AST)	Live (AST)		Relevance	Irrelevance
Qwen2.5-14B with TOUCAN	65.09%	85.42%	76.01%	35.25%	72.22%	75.96%
Qwen2.5-14B with Nemotron-SFT(tools)	65.64%	85.02%	81.83%	30.00%	66.67%	85.45%

Table 12: Comparison between Nemotron-SFT (tool subset) and TOUCAN on  $\tau$ - and  $\tau^2$ -Bench.

Model	$\tau$ -bench			$\tau^2$ -bench			
	Avg.	Airline	Retail	Avg.	Airline	Retail	Telecom
Qwen2.5-14B with TOUCAN	35.24%	22.00%	48.48%	30.43%	22.00%	49.10%	20.18%
Qwen2.5-14B with Nemotron-SFT(tools)	24.38%	18.00%	30.76%	20.23%	16.00%	36.80%	7.90%

### C.8 ABLATION ON THE IMPORTANCE OF TRAJECTORY ANNOTATION AND FILTERING

We conducted an ablation study comparing datasets with and without the filtering step (Stage 5) of our generation pipeline. We report results in Table 13. We observed that the overall performance on BFCL V3, as well as the multi-turn and irrelevance setting benefit from filtering, which confirms the value of including an automated process to filter-out low-quality samples in data generation pipelines.

Table 13: Ablation on Stage 5 on BFCL-V3 Benchmark.

Model	Overall	Single Turn		Multi Turn	Hallucination	
		Non-live (AST)	Live (AST)		Relevance	Irrelevance
Qwen2.5-14B-Instruct (FC)	57.69%	83.38%	73.70%	19.75%	83.33%	68.46%
TOUCAN without Filtering (4 stages)	62.60%	<b>86.83%</b>	72.01%	32.25%	77.78%	67.01%
TOUCAN with Filtering (5 stages)	<b>65.09%</b>	85.42%	<b>76.01%</b>	<b>35.25%</b>	72.22%	<b>75.96%</b>

## D PROMPTS

### D.1 MCP SERVER ANNOTATION PROMPT

Below is the prompt for annotating MCP server categories.

```

1302 1  ## Task
1303 2  Generate **Server Labels** to categorize the provided MCP Server based on
1304 3  ↪ its description and available tools.
1305 4  ## Objective
1306 5  Analyze the provided MCP Server's description and available tools, then
1307 6  ↪ assign appropriate category labels that best describe its primary
1308 7  ↪ functionality and use cases.
1309 8  ## Guidelines
1310 9  ### Label Selection
1311 10 - Analyze the MCP Server's core functionality and purpose
1312 11 - Consider the types of tools it provides and the problems it solves
1313 12 - Select labels that accurately represent the server's primary use cases
1314 13 - Choose from predefined categories when applicable, but also consider
1315 14 ↪ custom labels for unique functionality
1316 15 ### Predefined Categories
1317 16 Choose from these established categories when appropriate:
1318 17 - **Web Search & Research**: Tools for searching the web, gathering
1319 18 ↪ information, academic research
1320 19 - **Browser Automation**: Web scraping, automated browsing, page
1321 20 ↪ interaction
1322 21 - **Memory Management**: Data storage, retrieval, knowledge bases,
1323 22 ↪ note-taking
1324 23 - **Operating System**: File operations, system commands, process
1325 24 ↪ management
1326 25 - **Data Analysis & Processing**: Analytics, data transformation,
1327 26 ↪ statistical analysis
1328 27 - **Cryptocurrency & Blockchain**: Trading, wallet management, DeFi,
1329 28 ↪ blockchain interaction
1330 29 - **Daily Productivity**: Task management, scheduling, personal
1331 30 ↪ organization
1332 31 - **File Management**: File operations, document handling, storage
1333 32 ↪ management
1334 33 - **Database Operations**: Data querying, database management, SQL
1335 34 ↪ operations
1336 35 - **API Integration**: Third-party service integration, webhook handling
1337 36 - **Communication Tools**: Messaging, email, notifications, social
1338 37 ↪ interaction
1339 38 - **Development Tools**: Code analysis, debugging, version control, CI/CD
1340 39 - **Security & Authentication**: Password management, encryption, access
1341 40 ↪ control
1342 41 - **Cloud Services**: Cloud platform integration, serverless functions
1343 42 - **AI/ML Tools**: Machine learning, model interaction, AI-powered
1344 43 ↪ features
1345 44 - **Content Creation**: Writing, editing, media generation, publishing
1346 45 - **Social Media**: Social platform integration, posting, analytics
1347 46 - **Financial Services**: Banking, payments, financial data, accounting
1348 47 - **E-commerce**: Shopping, product management, order processing
1349 48 - **Gaming**: Game-related tools, entertainment, interactive features
1350 49 - **Education**: Learning tools, course management, educational content
1351 50 - **Health & Fitness**: Health monitoring, fitness tracking, medical
1352 51 ↪ tools
1353 52 - **Travel & Maps**: Location services, travel planning, navigation
1354 53 - **News & Media**: News aggregation, media consumption, journalism tools
1355 54 - **Weather**: Weather data, forecasting, climate information
1356 55 - **Time & Calendar**: Scheduling, time management, calendar integration

```

```

1350
135143
135244 ### Custom Labels
135345 - If the server doesn't fit well into predefined categories, create a
135446   ↳ custom label
135547 - Custom labels should be descriptive and specific to the server's unique
135648   ↳ functionality
135749 - Use clear, concise terminology that would be useful for clustering and
135850   ↳ organization
135951
136052 ### Output Requirements
136153 - **Primary Label**: The main category that best describes the server
136254   ↳ (from predefined list or custom)
136355 - **Secondary Labels**: Additional relevant categories (0-2 labels)
136456 - **Custom Label**: A free-form descriptive label if the server has
136557   ↳ unique functionality not covered by predefined categories
136658
136759 ## MCP Server Description
136860 {MCP_SERVER_NAME}: {MCP_SERVER_DESCRIPTION}
136961
137062 Available Tools:
137163 {TOOL_LIST}
137264
137365 ## Output
137466 Provide your response in the following XML format:
137567
137668 <response>
137769   <analysis>
137870     <!-- Briefly analyze the MCP Server's core functionality and the
137971       ↳ types of problems it solves based on its description and
138072       ↳ available tools. -->
138173   </analysis>
138274   <reasoning>
138375     <!-- Brief explanation of why these labels were chosen and how they
138476       ↳ represent the server's functionality -->
138577   </reasoning>
138678   <primary_label>
138779     <!-- The main category that best describes this server's primary
138880       ↳ functionality -->
138981   </primary_label>
139082   <secondary_labels>
139183     <!-- Additional relevant categories (0-2 labels), separated by commas
139284       ↳ if multiple -->
139385   </secondary_labels>
139486   <custom_label>
139587     <!-- A free-form descriptive label if the server has unique
139688       ↳ functionality not covered by predefined categories. Leave empty
139789       ↳ if not needed. -->
139890   </custom_label>
139991 </response>

```

## 1391 D.2 TASK GENERATION PROMPT

1393 Below is an example of a task generation prompt for the single-server task synthesis. The prompt  
 1394 generates a question targeting **one tool**.  
 1395

```

13961 ## Task
13972 Generate a **Tool Use Question** based on the provided MCP Server and its
13983   ↳ tool descriptions.
13994
14005 ## Objective
14016 Analyze the provided MCP Server and its available tools, then create a
14027   ↳ realistic user question that would naturally require the use of one
14038   ↳ of these tools to solve.
14049
140510 ## Guidelines

```



```

1404 8
1405 9 ### Question Realism
1406 10 - Create questions that represent real-world scenarios where users would
1407 11   ↳ need to interact with the MCP Server's tools
1408 11 - The question should sound natural and authentic, as if asked by someone
1409 12   ↳ genuinely needing to accomplish a task
1410 12 - Consider common use cases, problems, or workflows that would require
1411 13   ↳ the functionality provided by the MCP Server's tools
1412 14 ### Tool Selection
1413 15 - Focus on **ONE specific tool** from the MCP Server that would be most
1414 16   ↳ appropriate to answer the question
1415 16 - Choose tools based on the core functionality they provide and how they
1416 17   ↳ would solve real user problems
1417 18 - Consider each tool's description and purpose when crafting the question
1418 19 ### Question Complexity
1419 20 - Create questions that are clear and specific enough to warrant tool
1420 21   ↳ usage
1421 21 - Avoid overly simple questions that could be answered without tools
1422 22 - Include relevant context or constraints that make the tool usage
1423 23   ↳ necessary
1424 24 - Do not contain the exact tool name in the question
1425 25 ### Output Format
1426 26 Your response should include:
1427 27 1. **Tool Analysis**: Briefly analyze the MCP Server's available tools
1428 28   ↳ and their main functionalities.
1429 28 2. **Target Tool**: The specific tool name from the MCP Server that
1430 29   ↳ should be used to answer this question.
1431 30 3. **Question**: A clear, realistic user question that requires tool
1432 31   ↳ usage.
1433 32 ## MCP Server Description
1434 33 {MCP_SERVER_NAME}: {MCP_SERVER_DESCRIPTION}
1435 34 Available Tools:
1436 35 {TOOL_LIST}
1437 36 ## Output
1438 37 Provide your response in the following XML format:
1439 38
1440 39 <response>
1441 40   <server_analysis>
1442 41     <!-- Briefly analyze the MCP Server's available tools and their main
1443 42      ↳ functionalities. -->
1444 43   </server_analysis>
1445 44   <target_tool>
1446 45     <!-- The specific tool name from the MCP Server that should be used
1447 46      ↳ to answer this question. -->
1448 47   </target_tool>
1449 48   <question>
1450 49     <!-- A clear, realistic user question that requires tool usage. -->
1451 50   </question>
1452 51 </response>
1453
1454 1 ## Task
1455 2 Generate a **Tool Use Question** based on the provided MCP Server and its
1456 3   ↳ tool descriptions.
1457 4 ## Objective

```

```

1458 5 Analyze the provided MCP Server and its available tools, then create a
1459 ↪ realistic user question that would naturally require the use of
1460 ↪ **{NUM_TOOLS} tools** from this MCP Server to solve completely.
1461 6
1462 7 ## Guidelines
1463 8
1464 9 ### Question Realism
1465 10 - Create questions that represent real-world scenarios where users would
1466 ↪ need to interact with the MCP Server's tools
1467 11 - The question should sound natural and authentic, as if asked by someone
1468 ↪ genuinely needing to accomplish a task
1469 12 - Consider common use cases, problems, or workflows that would require
1470 ↪ the functionality provided by the MCP Server's tools
1471 13
1472 14 ### Tool Selection
1473 15 - Focus on **{NUM_TOOLS} tools** from the MCP Server that would work
1474 ↪ together to answer the question
1475 16 - The question should require a sequence or combination of tool calls to
1476 ↪ solve completely
1477 17 - Choose tools based on how they complement each other and create a
1478 ↪ logical workflow
1479 18 - Consider each tool's description and purpose when crafting the question
1480 ↪ that requires multiple steps
1481 19
1482 20 ### Question Complexity
1483 21 - Create questions that are complex enough to warrant using {NUM_TOOLS}
1484 ↪ tools
1485 22 - The question should have multiple components or require several steps
1486 ↪ to solve
1487 23 - Include relevant context or constraints that make the multi-tool usage
1488 ↪ necessary
1489 24 - Do not contain the exact tool names in the question
1490 25 - Ensure the question cannot be reasonably answered with just a single
1491 ↪ tool
1492 26
1493 27 ### Output Format
1494 28 Your response should include:
1495 29 1. **Tool Analysis**: Briefly analyze the MCP Server's available tools
1496 ↪ and their main functionalities.
1497 30 2. **Target Tools**: The specific tool names from the MCP Server that
1498 ↪ should be used together to answer this question, in the order they
1499 ↪ would likely be called.
1500 31 3. **Question**: A clear, realistic user question that requires multiple
1501 ↪ tool usage.
1502 32
1503 33 ## MCP Server Description
1504 34 {MCP_SERVER_NAME}: {MCP_SERVER_DESCRIPTION}
1505 35
1506 36 Available Tools:
1507 37 {TOOL_LIST}
1508 38
1509 39 ## Output
1510 40 Ensure your question requires exactly {NUM_TOOLS} tools to solve
1511 ↪ completely. Provide your response in the following XML format:
1512 41
1513 42 <response>
1514 43   <server_analysis>
1515 44     <!-- Briefly analyze the MCP Server's available tools and their main
1516 ↪ functionalities. -->
1517 45   </server_analysis>
1518 46   <target_tools>
1519 47     <!-- The specific tool names from the MCP Server that should be used
1520 ↪ together to answer this question, listed in order. e.g.,
1521 ↪   <tool>create_twitter_post</tool> <tool>get_last_tweet</tool> -->
1522 48   </target_tools>
1523 49   <question>

```

```

1512      <!-- A clear, realistic user question that requires multiple tool
1513      ↪ usage. -->
1514    </question>
1515  </response>
1516
1517  Below is an example of a task generation prompt for the multi-server task synthesis.
1518
1519  ## Task
1520  Generate a **Multi-Server Tool Use Question** based on the provided MCP
1521  ↪ Servers and their tool descriptions.
1522
1523  ## Objective
1524  Analyze the provided MCP Servers and their available tools, then create a
1525  ↪ realistic user question that would naturally require the use of
1526  ↪ **{NUM_TOOLS} tools from at least 2 different MCP servers** to solve
1527  ↪ completely.
1528
1529  ## Guidelines
1530
1531  ### Question Realism
1532  - Create questions that represent real-world scenarios where users would
1533  ↪ need to interact with tools from multiple MCP Servers
1534  - The question should sound natural and authentic, as if asked by someone
1535  ↪ genuinely needing to accomplish a complex task
1536  - Consider workflows that span across different services/domains that
1537  ↪ would require multiple servers
1538  - Think about how different MCP servers complement each other in
1539  ↪ real-world use cases
1540
1541  ### Server and Tool Selection
1542  - Use tools from **at least 2 different MCP servers** to answer the
1543  ↪ question
1544  - Select **{NUM_TOOLS} tools total** that work together across multiple
1545  ↪ servers
1546  - The question should require a sequence or combination of tool calls
1547  ↪ from different servers to solve completely
1548  - Choose tools based on how they complement each other across different
1549  ↪ services/domains
1550  - Consider each tool's description and purpose when crafting the
1551  ↪ cross-server workflow
1552  - Ensure tools from different servers create a logical, interconnected
1553  ↪ workflow
1554
1555  ### Question Complexity
1556  - Create questions that are complex enough to warrant using {NUM_TOOLS}
1557  ↪ tools across multiple servers
1558  - The question should have multiple components or require several steps
1559  ↪ that span different services
1560  - Include relevant context or constraints that make the multi-server tool
1561  ↪ usage necessary
1562  - Do not contain the exact tool names or server names in the question
1563  - Ensure the question cannot be reasonably answered with tools from just
1564  ↪ a single server
1565  - Create scenarios that naturally require different types of services
1566  ↪ working together
1567
1568  ### Cross-Server Integration
1569  - Think about how different servers' capabilities can be combined
1570  - Consider data flow between different services (e.g., retrieving data
1571  ↪ from one service to use in another)
1572  - Create realistic scenarios where multiple services need to work
1573  ↪ together
1574  - Focus on complementary functionalities across different domains
1575
1576  ### Output Format
1577  Your response should include:

```

```

1566 1. **Server Analysis**: Briefly analyze all MCP Servers and their
1567 ↪ available tools, focusing on how they can work together.
1568 2. **Cross-Server Workflow**: Describe the workflow showing how tools
1569 ↪ from different servers will be used together.
1570 3. **Target Tools**: The specific tool names from different MCP Servers
1571 ↪ that should be used together, in the order they would likely be
1572 ↪ called, with their server names.
1573 4. **Question**: A clear, realistic user question that requires
1574 ↪ multi-server tool usage.
1575 ## Available MCP Servers
1576 {SERVER_DESCRIPTIONS}
1577 ## Output
1578 Ensure your question requires exactly {NUM_TOOLS} tools from at least 2
1579 ↪ different servers to solve completely. Provide your response in the
1580 ↪ following XML format:
1581 <response>
1582   <server_analysis>
1583     <!-- Briefly analyze all MCP Servers and their available tools,
1584     ↪ focusing on how they can work together across different
1585     ↪ domains/services. -->
1586   </server_analysis>
1587   <cross_server_workflow>
1588     <!-- Describe the workflow showing how tools from different servers
1589     ↪ will be used together to solve the question. -->
1590   </cross_server_workflow>
1591   <target_tools>
1592     <!-- The specific tool names from different MCP Servers that should
1593     ↪ be used together, listed in order with their server names. e.g.,
1594     ↪ <tool server="Server1">search_posts</tool> <tool
1595     ↪ server="Server2">send_email</tool> -->
1596   </target_tools>
1597   <question>
1598     <!-- A clear, realistic user question that requires multi-server tool
1599     ↪ usage spanning different services/domains. -->
1600   </question>
1601 </response>
1602
1603 Below is an example of a task generation prompt for the task synthesis for featured servers.
1604
1605 ## Task
1606 Generate a **Multi-Server Tool Use Question** based on featured MCP
1607 ↪ Servers and their tool descriptions.
1608
1609 ## Objective
1610 Brainstorm a compelling real-world scenario, then analyze the provided
1611 ↪ featured MCP Servers and their available tools to create a realistic
1612 ↪ user question that would naturally require the use of **{NUM_TOOLS}
1613 ↪ tools from at least 2 different MCP servers** to solve completely.
1614
1615 ## Guidelines
1616
1617 ### Scenario Brainstorming
1618 - Think of realistic, specific scenarios where someone would need to use
1619 ↪ {NUM_TOOLS} different tools across multiple servers to accomplish a
1620 ↪ meaningful task
1621 - Consider diverse real-world contexts such as:
1622   - Content creators managing their online presence across different
1623   ↪ platforms
1624   - Researchers gathering and analyzing information from multiple sources
1625   - Developers building and deploying applications using different
1626   ↪ services

```

```

162015 - Business professionals managing projects and communications across
162116   ↳ platforms
162217 - Students working on complex assignments requiring multiple tools
162318 - Entrepreneurs launching new ventures using various services
162419 - The scenario should be detailed and authentic, representing genuine use
162520   ↳ cases that span multiple services
162621
162722 ### Question Realism
162823 - Create questions that represent real-world scenarios where users would
162924   ↳ genuinely need tools from multiple MCP servers
163025 - The question should sound natural and authentic, as if asked by someone
163126   ↳ with a specific goal
163227 - Include relevant context, constraints, and details that make the
163328   ↳ question engaging
163429 - Consider workflows that require multiple complementary tools working
163530   ↳ together across different services
163631 - Think about how different servers support each other in real-world use
163732   ↳ cases
163833
163934 ### Server and Tool Selection
164035 - Use tools from **at least 2 different MCP servers** to answer the
164136   ↳ question
164237 - Select **{NUM_TOOLS} tools total** that work together across multiple
164338   ↳ servers
164439 - The question should require a sequence or combination of tool calls
164540   ↳ from different servers to solve completely
164641 - Choose tools based on how they complement each other across different
164742   ↳ services/domains
164843 - Consider each tool's description and purpose when crafting the
164944   ↳ cross-server workflow
165045 - Ensure tools from different servers create a logical, interconnected
165146   ↳ workflow
165247
165348 ### Question Complexity
165449 - Create questions that are complex enough to warrant using {NUM_TOOLS}
165550   ↳ tools across multiple servers
165651 - The question should have multiple components or require several steps
165752   ↳ that span different services
165853 - Include relevant context or constraints that make the multi-server tool
165954   ↳ usage necessary
166055 - Do not contain the exact tool names or server names in the question
166156 - Ensure the question cannot be reasonably answered with tools from just
166257   ↳ a single server
166358 - Create scenarios that naturally require different types of services
166459   ↳ working together
166560
166661 ### Cross-Server Integration
166762 - Think about how different servers' capabilities can be combined
166863 - Consider data flow between different services (e.g., retrieving data
166964   ↳ from one service to use in another)
167065 - Create realistic scenarios where multiple services need to work
167166   ↳ together
167267 - Focus on complementary functionalities across different domains
167368
167469 ### Output Format
167570 Your response should include:
167671 1. **Server Analysis**: Briefly analyze the featured MCP Servers and
167772   ↳ their available tools, focusing on how they can work together.
167873 2. **Cross-Server Workflow**: Describe the workflow showing how tools
167974   ↳ from different servers will be used together.
168075 3. **Target Tools**: The specific tool names from different MCP Servers
168176   ↳ that should be used together, in the order they would likely be
168277   ↳ called, with their server names.
168378 4. **Question**: A clear, realistic user question that requires
168479   ↳ multi-server tool usage.

```

```

167456 ## Available Featured MCP Servers
167557
167658 {FEATURED_SERVER_DESCRIPTIONS}
167759
167860 ## Output
167961 Ensure your question requires exactly {NUM_TOOLS} tools from at least 2
1680 ↪ different servers to solve completely. Provide your response in the
1681 ↪ following XML format:
168262
168363 <response>
168464   <server_analysis>
168565     <!-- Briefly analyze the featured MCP Servers and their available
1686 ↪ tools, focusing on how they can work together across different
1687 ↪ domains/services. -->
168866   </server_analysis>
168967   <cross_server_workflow>
169068     <!-- Describe the workflow showing how tools from different servers
1691 ↪ will be used together to solve the question. -->
169269   </cross_server_workflow>
169370   <target_tools>
169471     <!-- The specific tool names from different MCP Servers that should
1695 ↪ be used together, listed in order with their server names. e.g.,
1696 ↪ <tool server="Server1">search_posts</tool> <tool
1697 ↪ server="Server2">send_email</tool> -->
169872   </target_tools>
169973   <question>
170074     <!-- A clear, realistic user question that requires multi-server tool
1701 ↪ usage spanning different services/domains. -->
170275   </question>
170376 </response>
170477
170578
170679
170780
170881
170982
171083
171184
171285
171386
171487
171588
171689
171790
171891
171992
172093
172194
172295
172396
172497
172598
172699
1727100
1728101
1729102
1730103
1731104
1732105
1733106
1734107
1735108
1736109
1737110
1738111
1739112
1740113
1741114
1742115
1743116
1744117
1745118
1746119
1747120
1748121
1749122
1750123
1751124
1752125
1753126
1754127
1755128
1756129
1757130
1758131
1759132
1760133
1761134
1762135
1763136
1764137
1765138
1766139
1767140
1768141
1769142
1770143
1771144
1772145
1773146
1774147
1775148
1776149
1777150
1778151
1779152
1780153
1781154
1782155
1783156
1784157
1785158
1786159
1787160
1788161
1789162
1790163
1791164
1792165
1793166
1794167
1795168
1796169
1797170
1798171
1799172
1800173
1801174
1802175
1803176
1804177
1805178
1806179
1807180
1808181
1809182
1810183
1811184
1812185
1813186
1814187
1815188
1816189
1817190
1818191
1819192
1820193
1821194
1822195
1823196
1824197
1825198
1826199
1827200
1828201
1829202
1830203
1831204
1832205
1833206
1834207
1835208
1836209
1837210
1838211
1839212
1840213
1841214
1842215
1843216
1844217
1845218
1846219
1847220
1848221
1849222
1850223
1851224
1852225
1853226
1854227
1855228
1856229
1857230
1858231
1859232
1860233
1861234
1862235
1863236
1864237
1865238
1866239
1867240
1868241
1869242
1870243
1871244
1872245
1873246
1874247
1875248
1876249
1877250
1878251
1879252
1880253
1881254
1882255
1883256
1884257
1885258
1886259
1887260
1888261
1889262
1890263
1891264
1892265
1893266
1894267
1895268
1896269
1897270
1898271
1899272
1900273
1901274
1902275
1903276
1904277
1905278
1906279
1907280
1908281
1909282
1910283
1911284
1912285
1913286
1914287
1915288
1916289
1917290
1918291
1919292
1920293
1921294
1922295
1923296
1924297
1925298
1926299
1927300
1928301
1929302
1930303
1931304
1932305
1933306
1934307
1935308
1936309
1937310
1938311
1939312
1940313
1941314
1942315
1943316
1944317
1945318
1946319
1947320
1948321
1949322
1950323
1951324
1952325
1953326
1954327
1955328
1956329
1957330
1958331
1959332
1960333
1961334
1962335
1963336
1964337
1965338
1966339
1967340
1968341
1969342
1970343
1971344
1972345
1973346
1974347
1975348
1976349
1977350
1978351
1979352
1980353
1981354
1982355
1983356
1984357
1985358
1986359
1987360
1988361
1989362
1990363
1991364
1992365
1993366
1994367
1995368
1996369
1997370
1998371
1999372
2000373

```



```

1728 20 **Target Tools**: {TARGET_TOOLS}
1729 21 **Tool Descriptions**: {TOOL_DESCRIPTIONS}
1730 22
1731 23 ## Output Requirements
1732 24 Generate **{VARIATIONS_COUNT} augmented variations** of the original
1733 25 ↪ question. Each variation should:
1734 26 1. Maintain the same core goal that requires the target tool(s)
1735 27 2. Use the exact same tool(s) in the same order with the same final
1736 28 ↪ outcome
1737 29 3. Apply to a completely different context, scenario, or domain
1738 30 4. Keep the same complexity level and constraints as the original
1739 31 5. Feel like a natural, real-world scenario from a different setting
1740 32 6. Be meaningfully different from the original and other variations in
1741 33 ↪ terms of context only
1742 34 7. Avoid including any explicit mentions, hints, or references to the
1743 35 ↪ target tool names within the question text
1744 36
1745 37 ## Output
1746 38 Provide your response in the following XML format:
1747 39
1748 40 <response>
1749 41   <analysis>
1750 42     <!-- Briefly analyze the original question and target tool(s) to
1751 43     ↪ understand the core goal, tool usage pattern, complexity level,
1752 44     ↪ and expected outcome, then identify how this can be applied
1753 45     ↪ across different domains while maintaining operational
1754 46     ↪ consistency -->
1755 47   </analysis>
1756 48   <variations>
1757 49     <!-- Generate {VARIATIONS_COUNT} variations, each with <variation_X>,
1758 50     ↪ <context>, and <question> tags -->
1759 51     <variation_1>
1760 52       <context>
1761 53         <!-- Brief description of the new domain/scenario introduced -->
1762 54       </context>
1763 55       <question>
1764 56         <!-- The augmented question that maintains the same target
1765 57         ↪ tool(s) usage order, complexity, and outcome but in a
1766 58         ↪ different context -->
1767 59       </question>
1768 60     </variation_1>
1769 61     <!-- Continue with variation_2, variation_3, etc. as needed based on
1770 62     ↪ number of variations -->
1771 63   </variations>
1772 64 </response>
1773 65
1774 66 The prompt below is designed to enhance task complexity through the introduction of additional
1775 67 constraints.
1776 68
1777 69 ## Task
1778 70 Generate **augmented variations** of a given question that maintain the
1779 71 ↪ same target tool(s) usage and context but significantly increase the
1780 72 ↪ complexity and constraints required to solve the problem.
1781 73
1782 74 ## Objective
1783 75 Take an existing question and its associated target tool(s), then create
1784 76 ↪ multiple sophisticated variations that:
1785 77 - Use the same target tool(s) to achieve the core goal while navigating
1786 78 ↪ additional complexity layers
1787 79 - Maintain the same general context and domain as the original question
1788 80 - Increase multi-dimensional complexity through realistic constraints,
1789 81 ↪ competing requirements, stakeholder considerations, and
1790 82 ↪ interconnected dependencies
1791 83 - Embed the tool usage within larger, more complex workflows that require
1792 84 ↪ strategic thinking and coordination

```

```

1782 10 - Demonstrate how the same core tool usage applies under vastly different
1783      ↪ complexity levels
1784 11
1785 12 ## Guidelines
1786 13 - Introduce realistic constraints such as resource limits, compliance
1787      ↪ requirements, tight timelines, or stakeholder conflicts
1788 14 - Embed the same tool usage inside a broader workflow that requires
1789      ↪ coordination across teams or systems
1790 15 - Escalate demands (performance, scalability, risk management) without
1791      ↪ changing the original domain or context
1792 16 - Ensure each variation targets a different primary complexity angle
1793      ↪ (organizational, technical, strategic) while preserving tool
1794      ↪ relevance
1795 17 - Ensure the question does not contain any tool names or explicit
1796      ↪ references to the target tools.
1797 18
1798 19 ## Input Format
1799 20 **Original Question**: {ORIGINAL_QUESTION}
1800 21 **Target Tools**: {TARGET_TOOLS}
1801 22 **Tool Descriptions**: {TOOL_DESCRIPTIONS}
1802 23
1803 24 ## Output Requirements
1804 25 Generate **{VARIATIONS_COUNT} strategically augmented variations** of the
1805      ↪ original question. Each variation should:
1806 26 1. Maintain the same core goal that requires the target tool(s) while
1807      ↪ adding multiple complexity layers
1808 27 2. Keep the same general context and domain as the original question
1809 28 3. Introduce different but interconnected constraints and competing
1810      ↪ requirements
1811 29 4. Feel like natural, high-stakes, real-world scenarios that
1812      ↪ professionals encounter
1813 30 5. Be meaningfully different from the original and other variations in
1814      ↪ terms of complexity
1815 31 6. Include specific details that make the constraints and requirements
1816      ↪ concrete and actionable
1817 32 7. **Transform step-wise questions**: If the original question contains
1818      ↪ explicit steps, convert it to a goal-oriented format while
1819      ↪ maintaining the same tool usage requirements
1820 33 8. Avoid including any explicit mentions, hints, or references to the
1821      ↪ target tool names within the question text
1822 34
1823 35 ## Output
1824 36 Provide your response in the following XML format:
1825 37
1826 38 <response>
1827 39   <analysis>
1828 40     <!-- Analyze the original question and target tool(s) to understand
1829 41       ↪ the core goal, current complexity level, and identify multiple
1830 42       ↪ complexity dimensions that can be naturally introduced while
1831 43       ↪ maintaining tool relevance and solution feasibility -->
1832 44   </analysis>
1833 45   <variations>
1834 46     <!-- Generate {VARIATIONS_COUNT} variations, each with <variation_X>,
1835 47       ↪ <constraints>, and <question> tags -->
1836 48     <variation_1>
1837 49       <constraints>
1838 50         <!-- Specific organizational, stakeholder, or coordination
1839 51           ↪ constraints that add realistic complexity -->
1840 52       </constraints>
1841 53       <question>
1842 54         <!-- The complex, organizationally-focused question that
1843 55           ↪ maintains the same target tool(s) usage within a more
1844           ↪ sophisticated workflow -->
1845 56       </question>
1846 57     </variation_1>

```

```

1836      <!-- Continue with variation_2, variation_3, etc. as needed based on
1837      ↪ number of variations -->
1838    </variations>
1839  </response>
1840
1841

```

#### D.4 TASK QUALITY ANNOTATION PROMPT

```

1843
1844 1  ## Task
1845 2 Conduct a **Question Quality Assessment** of a tool use question across
1846 ↪ six key dimensions to ensure it meets high standards for realistic
1847 ↪ tool usage scenarios.
1848
1849 3 ## Objective
1850 4 Analyze the provided tool use question and assess its quality across six
1851 ↪ primary dimensions:
1852 5 1. **Tool Selection Difficulty** - How challenging it is to determine
1853 ↪ which tools to use giving all available tools
1854 6 2. **Tool Selection Uniqueness** - How unique and necessary the selected
1855 ↪ tools are for this specific task giving all available tools
1856 7 3. **Question Quality** - Overall clarity, specificity, and effectiveness
1857 8 4. **Scenario Realism** - How authentic and believable the scenario is
1858 9 5. **Verifiable** - How easy it is to verify the correctness of the final
1859 ↪ model answer
1860 10 6. **Stability** - How stable the answer will be when requested under
1861 ↪ different time and geolocation
1862
1863 11 ## Assessment Criteria
1864
1865 12 ### 1. Tool Selection Difficulty
1866 13 **What to Evaluate**: How difficult it would be for a user to determine
1867 ↪ which specific tools are needed to solve this question.
1868
1869 14 **Rating Guidelines**:
1870 15 - **very easy**: Question explicitly mentions tool names or makes tool
1871 ↪ selection obvious
1872 16 - **easy**: Tool selection is straightforward with clear indicators
1873 17 - **medium**: Requires some reasoning but tool needs are fairly apparent
1874 18 - **hard**: Requires careful analysis to determine appropriate tools
1875 19 - **very hard**: Requires extensive expertise and deep reasoning to
1876 ↪ identify the correct tools
1877
1878 20 ### 2. Tool Selection Uniqueness
1879 21 **What to Evaluate**: How unique and necessary the selected tools are for
1880 ↪ accomplishing this specific task, and whether the task can only be
1881 ↪ completed with these tools in the specified sequence.
1882
1883 22 **Rating Guidelines**:
1884 23 - **not unique**: Many alternative tool combinations could accomplish the
1885 ↪ same task equally well
1886 24 - **somewhat unique**: Some alternative approaches exist, but selected
1887 ↪ tools offer advantages
1888 25 - **moderately unique**: Selected tools are well-suited, with limited
1889 ↪ alternative approaches
1890 26 - **quite unique**: Selected tools are particularly well-matched to the
1891 ↪ task requirements
1892 27 - **highly unique**: Task can only be accomplished effectively with these
1893 ↪ specific tools in this sequence
1894
1895 28 ### 3. Question Quality
1896 29 **What to Evaluate**: Overall quality, clarity, and effectiveness of the
1897 ↪ question as a realistic user query.
1898
1899 30 **Rating Guidelines**:
1900 31 - **very poor**: Unclear, ambiguous, or poorly constructed question
1901

```

```

1890 40 - **poor**: Some clarity issues, missing important context
1891 41 - **average**: Clear and understandable, but could be more specific or
1892    ↳ engaging
1893 42 - **good**: Well-constructed, clear, specific, and realistic
1894 43 - **excellent**: Exceptionally clear, detailed, engaging, and
1895    ↳ professionally written
1896 44
1897 45 ### 4. Scenario Realism
1898 46 **What to Evaluate**: How authentic, believable, and true-to-life the
1899    ↳ described scenario is.
1900 47
1901 48 **Rating Guidelines**:
1902 49 - **unrealistic**: Artificial, contrived, or implausible scenario
1903 50 - **somewhat unrealistic**: Some realistic elements but feels forced or
1904    ↳ unlikely
1905 51 - **moderately realistic**: Believable scenario with minor authenticity
1906    ↳ issues
1907 52 - **realistic**: Authentic scenario that represents genuine use cases
1908 53 - **highly realistic**: Completely natural, authentic scenario
1909    ↳ indistinguishable from real user needs
1910 54
1911 55 ### 5. Verifiable
1912 56 **What to Evaluate**: How easy it is to verify the correctness of the
1913    ↳ final model answer.
1914 57
1915 58 **Rating Guidelines**:
1916 59 - **hard to verify**: Fully free-form answer that requires extensive
1917    ↳ human judgment
1918 60 - **somewhat hard**: Mostly subjective answer with some verifiable
1919    ↳ elements
1920 61 - **moderately verifiable**: Short sentence that can be verified by LLM
1921    ↳ comparison
1922 62 - **mostly verifiable**: Answer with clear, objective components and some
1923    ↳ subjective elements
1924 63 - **easy to verify**: Answer can be verified by simple rules, exact
1925    ↳ matches, or clear success criteria
1926 64
1927 65 ### 6. Stability (1-5 Scale)
1928 66 **What to Evaluate**: How stable and consistent the answer will be when
1929    ↳ the question is asked under different environmental conditions and
1930    ↳ system contexts. Consider factors like temporal dependency,
1931    ↳ geographical variations, operating system differences, network
1932    ↳ environments, and software version variations.
1933 67
1934 68 **Rating Guidelines**:
1935 69 - **highly unstable**: Answer changes significantly across different
1936    ↳ conditions (real-time data, location-specific, system-dependent)
1937 70 - **somewhat unstable**: Answer may vary moderately based on
1938    ↳ environmental or system factors
1939 71 - **moderately stable**: Answer mostly consistent with minor variations
1940    ↳ due to context
1941 72 - **mostly stable**: Answer remains largely consistent across different
1942    ↳ conditions
1943 73 - **highly stable**: Answer is completely independent of environmental
1944    ↳ and system factors
1945 74
1946 75 ## Question Analysis
1947 76
1948 77 ### All Available Tools```
1949 78 {ALL_SERVER_AND_TOOL_INFORMATION}
1950 79 ```
1951 80
1952 81 ### Question Content
1953 82 ```
1954 83 {QUESTION_CONTENT}
1955 84 ```

```

```

1944
194585
194686 ### Intended Tool for This Question
194787 ...
194888 {INTENDED_TOOL}
194989 ...
195090
195191 ## Output Requirements
195292
195393 Provide analysis with detailed reasoning BEFORE scores for each of the
195494 ↪ six metrics.
195595
195696 ## Output
195797 Provide your response in the following XML format:
195898
195999 <response>
1960100   <tool_selection_difficulty>
1961101     <reasoning>
1962102       <!-- Detailed explanation including ambiguity level, domain
1963103         ↪ knowledge required, and alternative solutions giving all
1964104         ↪ available tools -->
1965105     </reasoning>
1966106     <rating><!-- Rating: very easy, easy, medium, hard, very hard
1967107       ↪ --></rating>
1968108   </tool_selection_difficulty>
1969109
1970110   <tool_selection_uniqueness>
1971111     <reasoning>
1972112       <!-- Detailed explanation of tool necessity, sequential
1973113         ↪ dependencies, and alternative tool viability giving all
1974114         ↪ available tools -->
1975115     </reasoning>
1976116     <rating><!-- Rating: not unique, somewhat unique, moderately unique,
1977117       ↪ quite unique, highly unique --></rating>
1978118   </tool_selection_uniqueness>
1979119
1980120   <question_quality>
1981121     <reasoning>
1982122       <!-- Detailed explanation covering linguistic quality, information
1983123         ↪ architecture, and actionability -->
1984124     </reasoning>
1985125     <rating><!-- Rating: very poor, poor, average, good, excellent
1986126       ↪ --></rating>
1987127   </question_quality>
1988128
1989129   <scenario_realism>
1990130     <reasoning>
1991131       <!-- Detailed explanation of industry authenticity, workflow
1992132         ↪ accuracy, and stakeholder behavior -->
1993133     </reasoning>
1994134     <rating><!-- Rating: unrealistic, somewhat unrealistic, moderately
1995135       ↪ realistic, realistic, highly realistic --></rating>
1996136   </scenario_realism>
1997137
1998138   <verifiable>
1999139     <reasoning>
2000140       <!-- Detailed explanation of answer format, objective criteria, and
2001141         ↪ ground truth availability -->
2002142     </reasoning>
2003143     <rating><!-- Rating: hard to verify, somewhat hard, moderately
2004144       ↪ verifiable, mostly verifiable, easy to verify --></rating>
2005145   </verifiable>
2006146
2007147   <stability>
2008148     <reasoning>
2009149       <!-- Detailed explanation of temporal/geographical/system
2010150         ↪ dependencies and environmental factors -->

```

```

1998     </reasoning>
1999     <rating><!-- Rating: highly unstable, somewhat unstable, moderately
2000     ↪ stable, mostly stable, highly stable --></rating>
2001     </stability>
2002 </response>
2003
2004
2005 D.5 TRAJECTORY ANNOTATION PROMPT
2006
2007 ## Task
2008 Conduct a **Response Quality Assessment** of a tool-use conversation
2009 ↪ across two LLM-scored dimensions, with a third dimension computed
2010 ↪ automatically outside the LLM.
2011
2012 ## Objective
2013 Analyze the provided conversation and assess its response quality across
2014 ↪ two primary dimensions scored by the LLM, while reserving an
2015 ↪ additional tool-call accuracy dimension for automated scoring:
2016 1. Completeness - Whether the assistant fully accomplished the user's
2017 ↪ request end-to-end
2018 2. Conciseness - Whether the assistant solved the task using the minimum
2019 ↪ necessary steps and verbosity
2020
2021 ## Assessment Criteria
2022
2023 ### 1. Completeness
2024 **What to Evaluate**: Did the assistant fully satisfy the user's goal
2025 ↪ given the conversation context? Consider whether the assistant:
2026 - Executed all required steps end-to-end (including
2027 ↪ saving/exporting/downloading where applicable)
2028 - Provided the final deliverable or a working alternative when blocked
2029 ↪ (e.g., tool failure with a usable fallback)
2030 - Included essential confirmations, paths, or instructions to achieve the
2031 ↪ outcome
2032 - Avoided missing key requirements or leaving the user with unresolved
2033 ↪ gaps
2034
2035 **Rating Guidelines**:
2036 - very incomplete: Major requirements missing; no usable outcome
2037 - incomplete: Some key requirements missing; outcome is not directly
2038 ↪ usable
2039 - partially complete: Core steps attempted; outcome usable only with user
2040 ↪ effort or missing minor requirements
2041 - mostly complete: Meets most requirements; small omissions or minor
2042 ↪ issues remain
2043 - fully complete: All requirements met with a usable outcome delivered
2044
2045 ### 2. Conciseness
2046 **What to Evaluate**: Did the assistant achieve the goal with minimal
2047 ↪ redundancy and steps? Consider whether the assistant:
2048 - Avoided repetitive or unnecessary explanations/tool calls
2049 - Used the minimal set of steps/tools to complete the task
2050 - Kept language concise while preserving clarity
2051
2052 **Rating Guidelines**:
2053 - very redundant: Excessive repetition or unnecessary steps/tool calls
2054 - redundant: Noticeable verbosity or extra steps beyond what's needed
2055 - average: Reasonably concise with minor extraneous content
2056 - concise: Efficient and to the point with minimal overhead
2057 - very concise: Maximally efficient while clear and complete
2058
2059 ## Response Analysis
2060
2061 ### Question Content
2062 ...

```

```

2052  {QUESTION_CONTENT}
2053  ...
2054
2055  ### Intended Tool for This Question
2056  ...
2057  {INTENDED_TOOL}
2058  ...
2059
2060  ### Conversation History
2061  ...
2062  {CONVERSATION_HISTORY}
2063  ...
2064
2065  ## Output Requirements
2066  - Provide detailed reasoning BEFORE ratings for Completeness and
2067    ↳ Conciseness
2068  - Do NOT score Tool Call Accuracy; include placeholders only
2069
2070  ## Output
2071  Provide your response in the following XML format:
2072
2073  <response>
2074    <completeness>
2075      <reasoning>
2076        <!-- Evaluate if the assistant delivered an end-to-end usable
2077          ↳ outcome, addressed all requirements, handled tool failures with
2078          ↳ alternatives, and provided necessary confirmations/paths. -->
2079      </reasoning>
2080      <rating><!-- Rating: very incomplete, incomplete, partially complete,
2081        ↳ mostly complete, fully complete --></rating>
2082    </completeness>
2083
2084    <conciseness>
2085      <reasoning>
2086        <!-- Evaluate if the assistant minimized redundant
2087          ↳ steps/explanations, avoided unnecessary tool calls, and kept
2088          ↳ messaging efficient while clear. -->
2089      </reasoning>
2090      <rating><!-- Rating: very redundant, redundant, average, concise,
2091        ↳ very concise --></rating>
2092    </conciseness>
2093  </response>
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

```



## E COST AND ACCESSIBILITY

All models used to build TOUCAN (data generation and annotation) are open-source and can be deployed efficiently using vLLM servers, which substantially reduce inference cost. Our data generation pipeline demonstrates that producing high-quality, end-to-end synthetic data is feasible without relying on proprietary models. In this section, we provide additional guidance on reproducing our results and extending TOUCAN to new MCP servers, considering both standard and resource-constrained computational settings. Table 14 lists the models used to build TOUCAN and their corresponding GPU requirements, along with a set of open-source, resource-efficient alternative models that are fully compatible with the TOUCAN tool-trajectory generation pipeline. These alternative options, when combined with manual review and/or lightweight verification tools, could produce data of comparable quality and difficulty to TOUCAN.

Table 14: Approximate (H100) GPU requirements for the models used at each pipeline stage, as well as lightweight open-source alternatives. GPU requirements are provided for full precision inference, quantized versions of these models would further reduce resource requirements.

Stage	LLM Used	Approx. GPUs (vLLM, BF16)	Alternative (Smaller LLM)	Approx. GPUs (vLLM, BF16)	Notes
<b>Task Synthesis</b>	Mistral-Small-3.2-24B-Instruct-2506	1	N/A	N/A	Already efficient; runs on single GPU.
	Devstral-Small-2505	1	N/A	N/A	Already efficient; runs on single GPU.
	GPT-OSS-120B	4	GPT-OSS-20B	1	Suitable trade-off between performance and compute.
	Kimi-K2-Instruct	2	–	–	MoE with $\sim 32$ B active parameters ( $\approx 1$ T total).
	Qwen3-32B-Instruct	1	Qwen2.5-7B-Instruct	1	Lighter variant preserving coherence for synthesis.
<b>Task Filtering</b>	Kimi-K2-Instruct	32	GPT-OSS-20B	1	Best performing open model for filtering.
<b>Trajectory Generation</b>	GPT-OSS-120B	4	GPT-OSS-20B	1	Strong performance-compute balance in same agent framework.
	Kimi-K2-Instruct	32	–	–	Efficient, coherent model for synthesis tasks.
	Qwen3-32B-Instruct	1	Qwen2.5-7B-Instruct	1	Lighter model preserving coherence in trajectory generation.
<b>Trajectory Filtering</b>	GPT-OSS-120B	4	GPT-OSS-20B	1	Suitable performance-compute compromise for long trajectories.