

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 THE TOOL DECATHLON: BENCHMARKING LANGUAGE AGENTS FOR DIVERSE, REALISTIC, AND LONG- HORIZON TASK EXECUTION

Anonymous authors

Paper under double-blind review

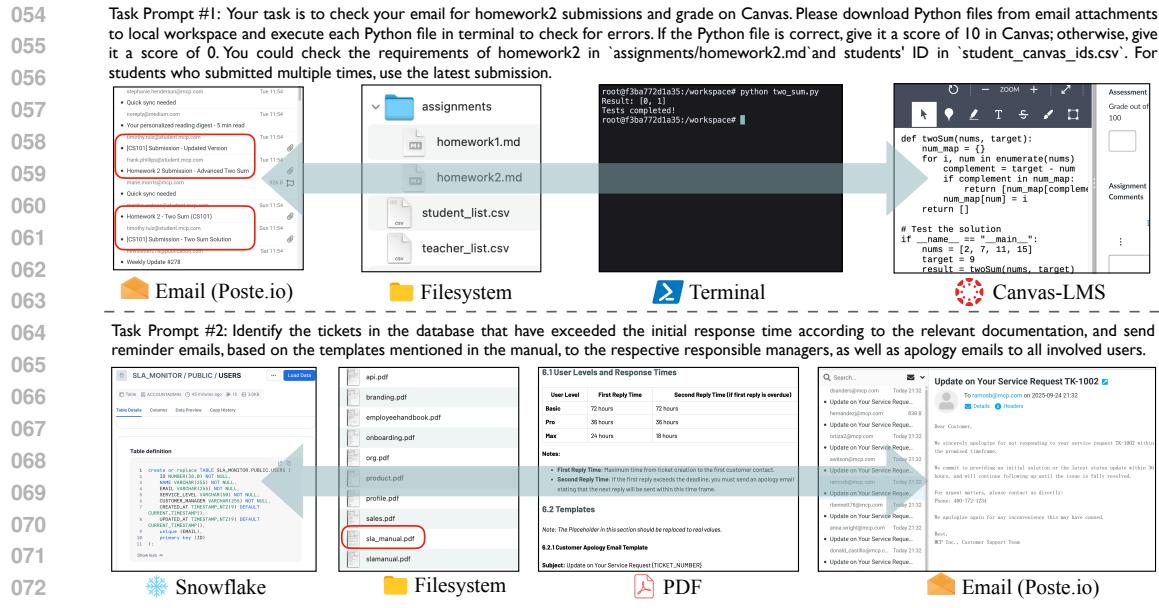
## ABSTRACT

Real-world language agents must handle complex, multi-step workflows across diverse applications. For instance, an agent may manage emails by coordinating with calendars and file systems, or monitor a production database like BigQuery to detect anomalies and generate reports following a standard operating manual. However, existing language agent benchmarks often focus on narrow domains or simplified tasks that lack the *diversity*, *realism*, and *long-horizon* complexity required to evaluate agents' real-world performance. To address this gap, we introduce the Tool Decathlon (dubbed as TOOLATHLON), a benchmark for language agents offering diverse applications and tools, realistic environment setup, and reliable execution-based evaluation. TOOLATHLON spans 32 software applications and 604 tools, ranging from everyday platforms such as Google Calendar and Notion to professional applications like WooCommerce, Kubernetes, and BigQuery. Most of the tools are based on a high-quality set of Model Context Protocol (MCP) servers that we may have revised or implemented ourselves. Unlike prior works, which primarily ensure functional realism but offer *limited environment state diversity*, we provide realistic initial environment states from real software. The TOOLATHLON benchmark includes 108 manually sourced or crafted tasks in total, requiring interacting with multiple applications over around 20 turns on average to complete. Each task is strictly verifiable through dedicated evaluation scripts. Comprehensive evaluation of state-of-the-art models highlights their significant shortcomings in performing real-world, long-horizon tasks: the best-performing model, Claude-4-Sonnet, achieves only a 29.9% success rate with 28 tool calling turns on average, while the top open-weights model DeepSeek-V3.1 reaches 13.9%. We expect TOOLATHLON to drive the development of more capable language agents for real-world, long-horizon task execution.

## 1 INTRODUCTION

Tool-based language agents have already demonstrated their impact in real-world domains such as software engineering (Jimenez et al., 2024; The Terminal-Bench Team, 2025), deep research (OpenAI, 2024), and web browsing (Zhou et al., 2024). To further expand the reach of language agents across diverse domains and applications, the Model Context Protocol (MCP) has been proposed to establish a standard for connecting language agents to tens of thousands of applications (Anthropic, 2024).

Existing benchmarks for language agents, however, are restricted to limited domains and tools (Mialon et al., 2023; Liu et al., 2024; Ma et al., 2024; Jimenez et al., 2024; Zhou et al., 2024; Yao et al., 2025; Wei et al., 2025; Xu et al., 2025). By contrast, real-world tasks often require switching across various applications. For example, as demonstrated in Figure 1 (Example #2), a company's administrative agent may need to monitor a real Snowflake database for customer tickets, locate the appropriate PDF operation manual containing instructions on how to identify and handle overdue tickets, and then send the required emails to managers and customers in accordance with the manual. Importantly, this diversity gap extends far beyond differences in tool names or descriptions. The diversity and complexity of environment states across applications, compounded by interaction with them in long trajectories, present substantial challenges for generalization.



072 Figure 1: Two examples and the initial environment states in TOOLATHLON. We showcase real-world environment interaction (§2.2) and realistic state initialization (§2.3) here.

073 To address these challenges, we introduce the Tool Decathlon (TOOLATHLON), a benchmark for  
074 evaluating language agents on *diverse, realistic, and long-horizon* tasks. TOOLATHLON spans 32  
075 real-world applications and 604 tools across 108 tasks, covering a wide spectrum of domains ranging  
076 from daily affair and education to technology and finance. Tasks are grounded in realistic scenarios  
077 and mostly require coordinating multiple applications. Each task is fully verifiable with a dedicated,  
078 deterministic evaluation script, comparing outcomes against either static or dynamically generated  
079 ground-truth states (e.g., tasks involving the latest NVIDIA shareholder information or real-time train  
080 schedules). All tools in TOOLATHLON are sourced from the real world, with the majority obtained  
081 from MCP servers.

082 To faithfully capture the realism and complexity of practical environment states, we tried to adopt  
083 the most representative applications such as Google Sheet, Gmail, and Snowflake. However, some  
084 remote environment states are difficult to set up to mimic real scenarios. For example, simulating a  
085 Canvas course with tens of students would require registering a real account for each student and  
086 resetting the states at each evaluation run. Therefore, while we adopt the commonly used applications  
087 most of the time, we also incorporate several local pieces of software deployed via containers for  
088 convenient and complex environment simulation, such as poste.io for email management to replace  
089 Gmail and WooCommerce for online ecommerce platform to replace Shopify. These services are  
090 based on their respective open-source software, providing complex observations while allowing us to  
091 set up the states within seconds. This stands in stark contrast with simplified or artificial environment  
092 states as in prior benchmarks (Patil et al., 2025). In addition, task prompts in TOOLATHLON are  
093 crafted to mirror authentic user queries, which are often concise and fuzzy. Models must therefore  
094 infer user intent and autonomously devise plans to accomplish tasks, an example is shown in Figure 3.

095 Concurrent with this work, several MCP-based tool-use benchmarks have emerged (Liu et al., 2025;  
096 Mo et al., 2025; Yan et al., 2025; Yin et al., 2025; The MCPMark Team, 2025), but they do not  
097 match TOOLATHLON in its reflection of real-world complexity. Some rely on LLM judges without  
098 verifiable tasks (Mo et al., 2025; Yin et al., 2025), while others cover few domains or mostly single-  
099 application tasks. For instance, MCPUniverse (Luo et al., 2025) spans only six domains, with 90%  
100 of tasks involving one app and synthetic initial states, yielding simplified, short interactions (<8  
101 turns). Similarly, MCPMark (The MCPMark Team, 2025) includes only five apps and overly detailed  
102 prompts (Figure 3). A full comparison is shown in Table 1.

103 TOOLATHLON includes a lightweight framework for automated, safe, and scalable evaluation. Each  
104 task comes with initial states setup if needed as well as an evaluation script (Figure 2). Executing and  
105 evaluating each task is isolated in separate containers to prevent interference. This enables fast parallel

108 Table 1: Comparison of Tool-Based Language Agent Benchmarks, where **BFCLv3-MT** represents the multi-turn  
109 subset released in the 3rd version of **BFCL**. “# Apps” denotes the number of MCP servers, which we do not  
110 annotate for benchmarks without clear application definition. “Avg # Turns” denotes the number of tool calling  
111 turns made by Claude-4-Sonnet, which we use as a proxy for task complexity. “Real Env” (§2.2) means the  
112 environment states and observations are from real-world software rather than artificial databases. “States Init”  
113 (§2.3) indicates the evaluation begins with a state initialization. “Verifiable Execution” (§2.4) denotes that models  
114 need to execute the tools and final results are evaluated based on states. “Realistic Fuzzy Prompt” represents that  
115 the task instructions are often fuzzy and ambiguous to mimic real user input (§3.1). \*For MCPUniverse, only  
116 10% of the tasks are cross-App. For **ACEBench** and **LiveMCPBench**, only <10% of the tasks contain simple  
117 states initialization. In-depth discussion of these related works is in Appendix B.

Benchmark	# Tasks	# Apps	Avg # Turns	Real Env	States Init	Verifiable Execution	Cross-App Task	Realistic Fuzzy Prompt
$\tau$ -Bench	165	2	—	✗	✓	✓	✗	✗
BFCLv3-MT	800	—	3.8	✗	✓	✓	✓	✗
ACEBench	2000	—	1.7	✗	Partial*	✓	✗	✗
AppWorld	750	9	—	✓	✓	✓	✓	✗
MCPWorld	201	10	—	✓	✓	✓	✗	✗
MCP-RADAR	507	6	—	✓	✓	✓	✗	✗
MCPEval	676	19	—	✓	✗	✗	✓	✗
LiveMCPBench	95	70	5.6	✓	Partial*	✗	✓	✗
MCP-AgentBench	600	33	—	✓	✗	✗	✓	✗
LiveMCP-101	101	41	5.4	✓	✗	✗	✓	✗
MCPAtlas	1000	40+	3-6	✓	✓	✗	✓	✗
MCPUniverse	231	11	7.5	✓	✗	✓	Partial*	✗
MCPMark	127	5	18.5	✓	✓	✓	✗	✗
<b>TOOLATHLON</b>	108	32	26.8	✓	✓	✓	✓	✓

131 evaluation—for example, running Claude-4-Sonnet on all 108 tasks takes only 70 minutes using  
132 10 parallel processes. With extensive experiments on **TOOLATHLON**, the best-performing models,  
133 Claude-4-Sonnet and GPT-5, achieve only 29.9% accuracy, highlighting the unique challenges posed  
134 by **TOOLATHLON**. DeepSeek-V3.1 (DeepSeek-AI, 2025) achieves 13.9% success rate as the best  
135 performer among open-source models. Further analysis reveals that weaknesses in long-context  
136 modeling and robust tool calling error tracking are major challenges for all evaluated models. We  
137 will fully open-source the benchmark and the **TOOLATHLON** environment, aiming for **TOOLATHLON**  
138 to accelerate the development of practical language agents.

## 2 THE TOOLATHLON ENVIRONMENT AND EVALUATION FRAMEWORK

### 2.1 TASK DEFINITION

144 Each task in **TOOLATHLON** can be formulated as a partially observable Markov decision process  
145 (POMDP)  $(\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R}, \mathcal{U})$  with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , observation space  $\mathcal{O}$ , transition  
146 function  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S} \times \mathcal{O}$ , reward function  $\mathcal{R} : \mathcal{S} \rightarrow [0, 1]$ , and instruction space  $\mathcal{U}$ . The  
147 environment states (§2.2, §2.3) can be the status in the email inbox and the observations are the  
148 sequential input to the model. The action space  $\mathcal{A}$  is the available tools for the respective task and the  
149 tool implementation directly defines the transition function. The reward function  $\mathcal{R}$  (§2.4) represents  
150 our execution-based evaluation which directly evaluates the environment state. Intuitively, real-world  
151 tools and environments will yield significantly more complex and diverse states and observations  
152 than the synthetic ones, and we will detail our designs of these variables in the following sections.

### 2.2 TOOLS, ENVIRONMENTS, AND FRAMEWORK

155 **MCP Servers:** In **TOOLATHLON**, we source our tools through a variety of MCP servers. Specifi-  
156 cally, we first decide a list of valuable and common real applications that we aim to benchmark on,  
157 then we see if we can find the corresponding open-source MCP servers for them. If not, we implement  
158 the MCP servers by ourselves. Notably, many open-source MCP server implementations contain  
159 bugs or exhibit certain limitations, for example, without the tools needed to complete our tasks. We  
160 further refine and improve these implementations ourselves. This way, we obtain a high-quality set  
161 of 32 MCP servers in total, where we include a complete list and their sources in Appendix C. The  
162 applications span diverse domains, extending well beyond common daily-use applications such as

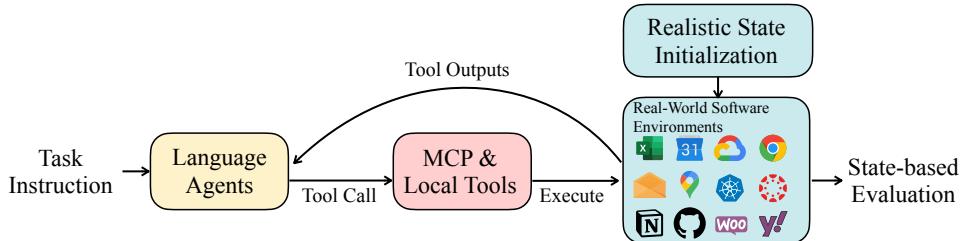


Figure 2: Overview of the TOOLATHLON evaluation framework.

Google Maps, Notion, and Google Calendar, and we also incorporate a number of professional and domain-specific applications to evaluate language agents in high-value productivity scenarios, such as Snowflake for enterprise data management and Kubernetes for cluster management. Although the majority of tools are sourced from MCP servers, the benchmark usage itself is not tied to MCP employment from the model developer side.

**Remote and Locally Containerized Environments:** While tools provide an interface for interacting with environments, they do not directly constitute the environments. Many real-world tools interact directly with existing, remote environments, such as Google Sheets, Google Calendar, Notion, and Gmail. Although remote environments require no implementation effort, they introduce significant challenges when benchmarking tasks that involve modifying environment states. For instance, simulating a realistic Gmail inbox with hundreds of emails from diverse senders would require registering hundreds of Google accounts for every benchmark user, and this inbox would need to be reset prior to each evaluation run. Previous works have attempted to bypass this issue by only supporting read operation to the states (Mialon et al., 2023), or implementing simplified synthetic data structures to mimic environment states (Patil et al., 2025; Yao et al., 2025), but such approaches drastically reduce realism and fail to reflect the complexity of real software environments. In contrast, in TOOLATHLON we leverage both remote environments and locally containerized, open-source applications. Specifically, we deploy the open-source Poste.io for email management, Canvas for course administration, Kubernetes for cluster orchestration, and WooCommerce for e-commerce management. By hosting these realistic applications locally within containers, we can efficiently set up dozens of accounts and initialize complex environment states during evaluation. Compared with existing dedicated agent sandboxes such as SWE-Bench (Jimenez et al., 2024), our environments are more diverse and encompass a wider range of software.

**Agent Framework:** We implement a simple agent framework based on the OpenAI Agents SDK to conduct the agent action loop – at each turn, the model is expected to (optionally) reason and make tool calls. We make several enhancements to improve its basic setup for a more robust workaround to evaluate language agents, including *tool error handling*, *overlong tool response handling* and *context history management*. We also equip this framework with some basic yet common local tools like *python execution*, *web search*, *claim done* and *sleep*. The details can be found in Appendix D.

### 2.3 INITIAL STATE SETUP

In real world, tasks are rarely executed from an empty environment state (e.g., an empty inbox). Instead, agents are typically required to operate based on pre-existing environment states. In agentic scenarios, task difficulty is determined not only by the task instructions but also by the underlying environment states. For example, operating on a folder with only one file to be sued is easier than working with 10 mixed useful and unrelated files (Figure 1, example #2), even if the task descriptions are nearly identical. To capture this, for tasks in TOOLATHLON that starts with an initial state,<sup>1</sup>, each of these tasks is equipped with a state initialization script to set up the states at running time, or (and) an `initial_workspace` directory containing pre-set files. Figure 1 and Figure 3 showcase such initial states. When constructing these initial environment states, we design them to closely reflect realistic scenarios. Notably, only very few previous benchmarks have incorporated realistic initial state construction before entering the agent loop, as summarized in Table 1. By contrast, most existing benchmarks start from empty state or overly simplified environment states, thus failing to capture the full complexity of real-world task execution.

<sup>1</sup>As shown in Table 2, 67% of the tasks fall into this category.

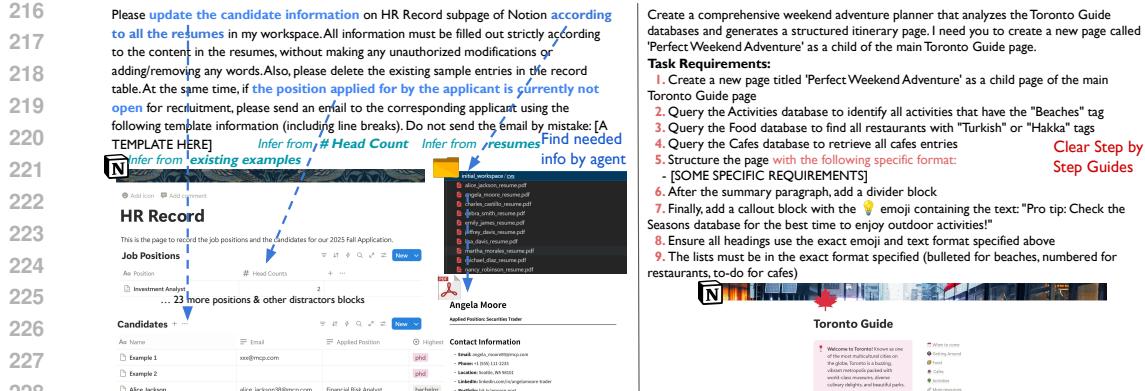


Figure 3: Example task instructions from our benchmark (Left) and MCPMark (The MCPMark Team, 2025) (Right). Ours contain more fuzzy intent that the model need to infer from the environment states.

## 2.4 RELIABLE EXECUTION-BASED EVALUATION

First, unlike some traditional tool-calling benchmarks that measure single-step tool call accuracy given a fixed context without actual execution (Patil et al., 2024), we think that execution-based evaluation is essential for reliably assessing language agents in realistic scenarios. Second, while many existing benchmarks rely on LLMs as judges to score agent trajectories (Gao et al., 2025; Yin et al., 2025), we contend that verifying the final environment states using deterministic rules offers a far more reliable and reproducible evaluation framework, as demonstrated in several widely adopted agent benchmarks (Zhou et al., 2024; Xie et al., 2024; Jimenez et al., 2024). To achieve this, each task in TOOLATHLON is equipped with a unique, manually crafted evaluation script that ensures precise and consistent measurement of task success. The script may perform robust matching against a static snapshot of the ground-truth environment or follow a reference execution workflow to dynamically retrieve and match real-time information (e.g., NVIDIA shareholders). During evaluation, each task is associated with a configuration file that specifies only the necessary MCP servers (< 10) and tools available for use rather than all 32 MCP servers. For each selected MCP server, we do not further filter or retrieve a subset of tools but load all tools within it, which is consistent with the current state of mainstream agent frameworks. Intuitively, providing the model with a larger set of unrelated tools increases task difficulty, as the agent must identify the relevant tools while ignoring distracting ones.

**Safe and Efficient Parallel Evaluation in Containers:** Our TOOLATHLON evaluation framework supports parallel execution to enable efficient model evaluation. Our framework launches each task inside a separate container in parallel, providing strict workspace isolation. On a standard Linux cluster with 16 CPUs and 64 GB of memory, we are able to evaluate Claude-4-Sonnet on 108 tasks in just about 70 minutes of wall time using 10 parallel processes. This demonstrates that TOOLATHLON is both convenient and efficient for practical use by model developers.

## 3 THE TOOLATHLON TASKS

### 3.1 TASK SOURCING AND FUZZY TASK INSTRUCTION

The authors of this work, who are researchers and senior undergraduate students in computer science, source and implement the tasks. We carefully design and adhere to several principles when collecting tasks: (1) **Real User Demands:** All tasks are either directly sourced from real-world websites or crafted to reflect genuine user demands. (2) **Multi-App Orchestration:** We intentionally source tasks that require interaction with multiple MCP servers, as this reflects authentic human workflows and increases task complexity. (3) **Diversity:** To ensure broad task diversity, we adopt a two-stage sourcing process. In the first stage, we start with an initial MCP server list covering more than 50 applications and freely source tasks without restricting to specific servers. In the second stage, we analyze the distribution of the sourced tasks and identify Apps that are important but underrepresented. We then conduct an additional round of targeted task sourcing specifically for them.

270  
 271 **Realistic Fuzzy Task Instruction:** We design  
 272 task instructions to resemble authentic user in-  
 273 put, which is often fuzzy or ambiguous. This  
 274 requires the agent to infer user intent given  
 275 the environment state, formulate plans, and  
 276 carry them out. For example, as shown in Fig-  
 277 ure 3 Left, a real user may simply say “Please  
 278 update the candidate information on HR  
 279 Record subpage according to all the  
 280 resumes..... if the position applied  
 281 for by the applicant is currently not  
 282 open, .....” This is a fuzzy instruction  
 283 without specifying in which format the agent  
 284 should fill in the information, but the Notion  
 285 database has provided some examples that the  
 286 agent needs to know to check itself. Also, the  
 287 instruction does not mention where to find the  
 288 status of the posted job and the agent needs to  
 289 check Notion to find that by itself. In contrast,  
 290 task instructions in some existing benchmarks  
 291 (Figure 3 Right) explicitly include detailed step-  
 292 by-step plans, which reduce the planning effort for agents. More examples of this kind are shown in  
 293 Figure 10 and 11.  
 294

295 All the sourced tasks experience multiple rounds of rigorous quality check, filtering and refinement  
 296 before we implement them, and finally we obtain 108 tasks in total. The topic distribution of all tasks  
 297 is shown in Figure 4 and Table 2 shows some key statistics of the complete benchmark.

### 298 3.2 TASK IMPLEMENTATION

299 As described in §2.4, each task in our benchmark  
 300 is implemented with a corresponding evaluation  
 301 script and potential initial states setup. This  
 302 process involves collecting ground-truth states  
 303 statically or dynamically, and design scripts to  
 304 automatically clear and re-fill new initial states.  
 305 To ensure realistic setups and reliable evaluation,  
 306 implementing a single task in TOOLATHLON  
 307 requires, on average, 4–6 hours of work by a graduate student majoring in computer science.

308 **Finalizing Tasks and Quality Check:** After crowd-sourcing task implementations from multiple  
 309 contributors, we perform intensive quality checks conducted by 5–6 experienced authors. In this  
 310 stage, each task is carefully reviewed and revised to unify standards across all tasks and ensure  
 311 correctness, solvability, and unambiguity, which requires approximately 5 hours of labor per task per  
 312 round of checking. Once all tasks are finalized, we perform an additional round of comprehensive  
 313 cross-checking and bug fixing of the entire benchmark before running the final experiments.

## 314 4 EXPERIMENT

315 In this section, we present the configuration details and experimental settings for several leading  
 316 commercial and open models on TOOLATHLON, as well as their performance.

### 317 4.1 SETUP

318 **Models and Configuration:** Our evaluation includes the leading commercial model series in terms  
 319 of agentic abilities, such as GPT-5 (-mini) (OpenAI, 2025b), o3&o4-mini (OpenAI, 2025a), Claude-  
 320 4-Sonnet (Anthropic, 2025), Gemini 2.5-(Pro,-Flash) (Comanici et al., 2025), Grok-4(-Fast) (xAI,  
 321 2025a;b), Grok-Code-Fast-1 (xAI, 2025c). We also benchmark the best-performing open-weight

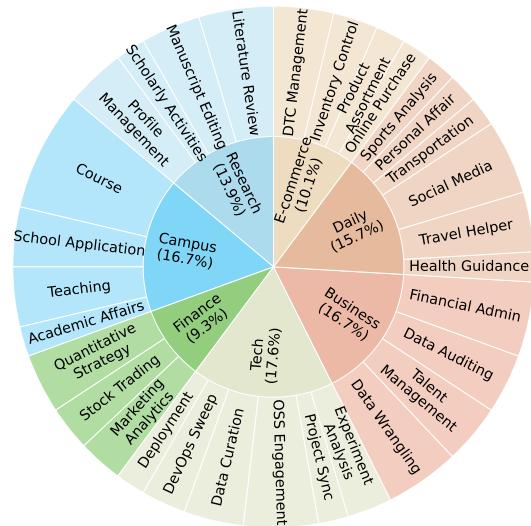


Figure 4: Task topic distribution of TOOLATHLON.

Table 2: Key statistics of TOOLATHLON.

Statistics	Value
# MCP servers (# tools)	32 (604)
# Local toolkits (# tools)	7 (16)
Avg/Min/Max tools per task	69.9/28/128
Tasks with state initialization	72/108 (67%)

324 Table 3: Main results for all the models. P@1, P@3, P^3 and # Turns represents Pass@1, Pass@3, Pass^3 and  
 325 average numbers of turns, respectively. We make **bold** the highest score.

Model	Research	Campus	Finance	Tech	Business	Daily	E-com	P@1	P@3	P^3	# Turns
<i>Proprietary Models</i>											
Claude-4-Sonnet	37.8	22.2	33.3	24.6	40.7	27.5	24.2	<b>29.9</b> <sub>±3.5</sub>	<b>42.6</b>	17.6	26.8
GPT-5	22.2	27.8	26.7	49.1	31.5	25.5	15.2	29.6 <sub>±1.5</sub>	38.9	<b>21.3</b>	19.5
GPT-5-high	17.8	24.1	23.3	40.4	29.6	27.5	24.2	27.5 <sub>±1.9</sub>	38.0	13.9	20.8
Grok-4	33.3	16.7	16.7	33.3	27.8	17.6	42.4	26.5 <sub>±2.3</sub>	39.8	14.8	19.4
Grok-4-Fast	26.7	14.8	13.3	22.8	22.2	17.6	27.3	20.7 <sub>±0.4</sub>	29.6	13.0	15.2
o3	15.6	16.7	10.0	31.6	24.1	17.6	3.0	18.5 <sub>±0.8</sub>	25.9	13.0	19.2
Grok-Code-Fast-1	24.4	11.1	10.0	19.3	18.5	13.7	21.2	17.0 <sub>±1.9</sub>	25.9	8.3	20.4
o4-mini	15.6	9.3	16.7	19.3	14.8	11.8	6.1	13.6 <sub>±0.4</sub>	24.1	6.5	18.3
GPT-5-mini	17.8	5.6	16.7	17.5	14.8	7.8	3.0	12.0 <sub>±0.8</sub>	19.4	5.6	19.5
Gemini-2.5-Pro	6.7	3.7	3.3	14.0	7.4	2.0	0.0	5.9 <sub>±1.9</sub>	11.1	0.9	12.5
Gemini-2.5-Flash	2.2	0.0	10.0	5.3	5.6	2.0	9.1	4.3 <sub>±1.6</sub>	6.5	2.8	8.8
<i>Open-Source Models</i>											
DeepSeek-V3.1	2.2	13.0	23.3	21.1	14.8	11.8	12.1	13.9 <sub>±0.8</sub>	21.3	8.3	29.9
GLM-4.5	13.3	11.1	16.7	17.5	5.6	15.7	15.2	13.3 <sub>±1.9</sub>	20.4	8.3	23.9
Qwen-3-Coder	8.9	11.1	10.0	15.8	14.8	9.8	12.1	12.0 <sub>±2.0</sub>	19.4	7.4	29.9
Kimi-K2-0905	11.1	9.3	13.3	14.0	7.4	9.8	18.2	11.4 <sub>±2.7</sub>	16.7	5.6	26.7

343  
 344 models including Qwen-3-Coder (Qwen Team, 2025), DeepSeek-V3.1 (DeepSeek-AI, 2025), Kimi-  
 345 K2-0905 (Kimi Team et al., 2025) and GLM-4.5 (GLM-4.5 Team et al., 2025). As described in §2.4,  
 346 each task is preconfigured with a list of MCP servers and common tools to access. During evaluation,  
 347 we set the maximum allowable number of turns as 100 for all models. In our main evaluation setup,  
 348 we only provide the models with the MCP servers and common tools that are useful for executing the  
 349 task. We note that as each MCP server is equipped with multiple related tools, so the models will  
 350 still see many unnecessary tools during evaluation. In Appendix E.2, we analyze the performance  
 351 variation by offering more unrelated MCP servers to the models.

352 **Metrics:** We evaluate each model three times and report the average pass@1 success rate as well as  
 353 the standard deviation. We also include the pass@3 – the fraction of tasks with at least one correct  
 354 trajectory, and pass^3 (Yao et al., 2025) – the fraction of tasks where all three trajectories are correct,  
 355 to measure the model’s potential capability coverage and its ability to complete tasks reliably. We  
 356 also report the average number of turns used.

## 358 4.2 MAIN RESULTS

360 Results in Table 3 show that Claude-4-Sonnet and GPT-5 are the leading ones, but they only achieve  
 361 about 30% success rate. Grok-4, with 26.5% Pass@1, is clearly in the second tier. All other models  
 362 remain at 20% or below. This indicates that our benchmark remains challenging for state-of-the-art  
 363 models and effectively distinguishes their capabilities. For open-source models, scores remain below  
 364 15%, revealing a gap with proprietary models. Notably, grok-code-fast-1 and grok-4-fast completed  
 365 all tasks at a high speed within 1 hour while scoring near 20% points, whereas others typically  
 366 required 1-3 hours. Interestingly, increased reasoning effort for thinking models (e.g., GPT-5 vs. GPT-  
 367 5-high) shows no benefits, suggesting that exploring new observations matters more than extended  
 368 internal reasoning in agentic tasks. We also find that Gemini-2.5’s ability to understand requirements  
 369 and proactively explore is insufficient – it may neglect requirements or give up prematurely during  
 370 the execution process, resulting in poor performance when handling complex tasks.

371 Looking at performance across task categories: Claude-4-Sonnet excels in *Finance*, *Research*,  
 372 *Business*, and *Daily* tasks, demonstrating strong general capabilities across diverse tools; GPT-5  
 373 performs exceptionally in *Campus* and *Tech* tasks, showcasing effectiveness in educational and  
 374 technical scenarios; Grok-4 stands out in *E-commerce*, indicating specialized strength in commercial  
 375 and transactional operations. We also observe significant differences between Pass@3 and Pass^3  
 376 success rates. This indicates that while many models have certain capability coverage, they lack  
 377 consistency in producing reliable results. For real-world tasks, building agents with both high success  
 rates and robust consistency remains a critical challenge.

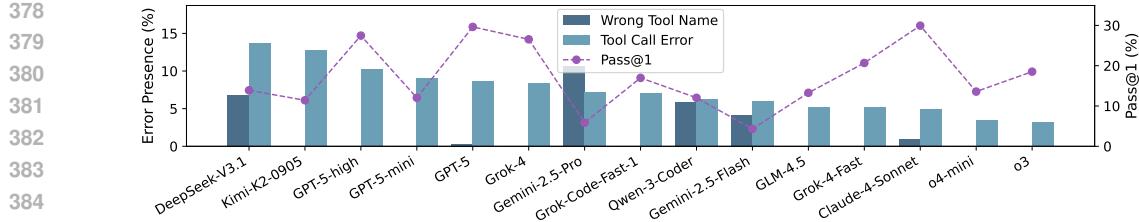


Figure 5: Two kinds of tool calling error presence ratios in calling tools for different models.

## 5 ANALYSIS

In this section, we conduct analysis in depth to better understand model performance in TOOLATHLON, focusing on tool-call errors, as well as long-context and overlong-output challenges. More analysis, including how tool error and the involvement of unrelated MCP servers impact model performance, and qualitative analysis & case studies, can be found in Appendix E and G.

### 5.1 THE FAILURE OF CALLING TOOLS

We mainly focus on two major tool-calling errors: hallucinating non-existing tools (e.g., incorrect tool names) and errors raised during tool execution. Statistics for different models on these two types of errors are shown in Figure 5. It can be seen that all models produce tool execution errors to varying degrees, possibly due to incorrect parameter passing or attempts to access non-existent resources. However, we found no significant correlation between overall success rate and the frequency of such errors. In fact, error messages from tools may help models understand the tool implementation or structure, allowing adjustments in subsequent turns. The other type of error—incorrect tool names—is more likely to affect final scores. Leading models produce few tool name errors. In Appendix E, Figure 9, we further analyze the success rate difference between trajectories containing tool-calling errors versus error-free trajectories. Most models do suffer from tool call errors, except for Claude-4-Sonnet, which may understand the tool better via error messages and gets a higher success rate.

### 5.2 THE LONG-CONTEXT CHALLENGES FOR LANGUAGE AGENTS

Since our benchmark is built on real tools and environments, it naturally generates many long-horizon trajectories. To quantitatively describe differences between tasks, we calculate the average execution turns for each task across all models, and use this as a proxy to divide all tasks into three equally-sized groups: Easy, Medium, and Hard, with execution turns increasing with difficulty. In Figure 6, we show the performance of five representative models on different groups and their average turns in each group. It shows that groups with higher average turns generally have lower success rates across models, and leading models (GPT-5, Claude-4-Sonnet) maintain clear advantages in all groups. We also find that there is no significant difficulty difference between Medium and Hard groups, suggesting that our benchmark’s difficulty doesn’t entirely stem from standard multi-step long-horizon execution, but possibly from models ending tasks prematurely without exploring enough observations, leading to failure.

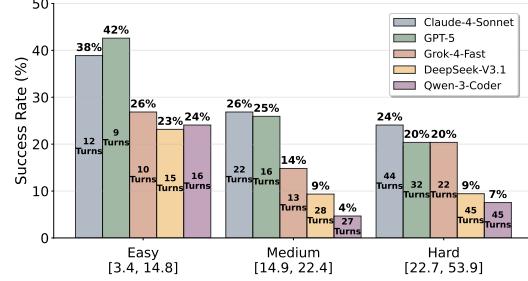


Figure 6: Model performance on three groups of tasks divided by average turns. The x-axis represents different task difficulty groups determined by different avg turns [Min Turns, Max Turns]

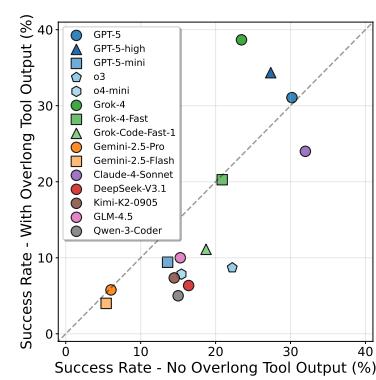


Figure 7: Avg. Success Rate on Trajectories w/wo overlong tool outputs.

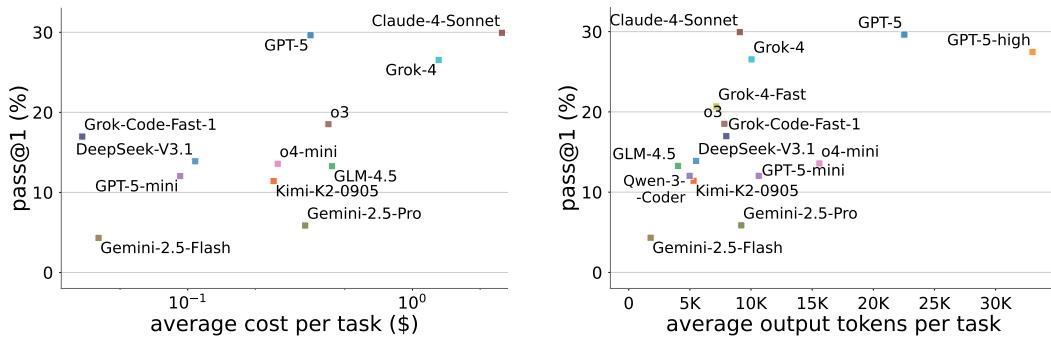


Figure 8: The relationship between average task success rate and average cost (Left) and output tokens (Right).

Another concern is whether models can successfully complete tasks when encountering overlong tool outputs, like fetching lengthy HTML source code or directly listing all data from a database (we refer the readers to Appendix D for handling overlong outputs in our framework). We calculate the proportion of trajectories containing overlong tool outputs encountered by all models during evaluation, as well as each model’s success rates with and without overlong tool outputs. Results show that the proportion of overlong tool outputs varies from approximately 15% to 30% across different models. Additionally, Figure 7 shows that most models experience a decline in success rate when encountering overlong tool outputs, with only a few models maintaining unchanged performance. We notice that Grok-4 performs significantly better with overlong tool outputs, likely because it excels at extracting key information from truncated data structures. While tasks with overlong outputs are often logically straightforward (e.g., price comparison, data extraction), most models get trapped trying to process complete lengthy outputs. Our analysis reveals that Grok-4 prefers not to read the saved overlong tool output, instead working directly with truncated versions from tool calls, making efficient information extraction from partial data critical for handling overlong outputs.

### 5.3 THE RELATIONSHIP BETWEEN PERFORMANCE AND EXPENSES

Since we are measuring models in realistic settings, cost and token usage matter as well, as they determine how different models should be used under various budgets. For cost,<sup>2</sup> Figure 8 Left shows that Claude-4-Sonnet and Grok-4 incur relatively high costs, whereas most other models remain under \$1 per task. Among the high-performing models, GPT-5 is notably cost-effective as it achieves comparable performance to Claude-4-Sonnet but costs only 1/10 as much, and Grok-Code-Fast and Grok-4-Fast (free under time limits) achieve the highest success rates among the lowest-cost models. We also plot the output token count distribution in Figure 8 Right. It shows how success rate varies with different average output token counts. Most models cluster between 5K–10K output tokens with success rates of 10–20%, exhibiting a weak positive correlation. Output conciseness varies markedly: GPT-5 produces the most tokens and increasing its reasoning effort (GPT-5-high) leads to more tokens, whereas Claude-4-Sonnet and Grok-4 achieve strong results with fewer output tokens, indicating they tend to rely more on environment observation instead of internal reasoning. Models like Gemini-2.5-Flash have the lowest output token counts and correspondingly low accuracy, while others show a similar concentrated distribution.

## 6 CONCLUSION

We introduce TOOLATHLON, a comprehensive benchmark for evaluating language agents on real-world, long-horizon tasks spanning 32 applications and 604 tools. Our evaluation reveals significant limitations in current models, with the best-performing Claude-4-Sonnet achieving only 29.9% success rate, highlighting substantial room for improvement in handling complex multi-step workflows. Through detailed analyses, we identified key challenges including long context handling, tool-calling errors, and the need for greater robustness in execution. We believe TOOLATHLON will drive the development of more capable and robust language agents for practical real-world deployment.

<sup>2</sup>Prompt caching will reduce the cost a lot in our evaluation as input tokens dominate, which explains why GPT-5 is not much more expensive than some open models.

486 REFERENCES  
487488 Anthropic. Introducing the model context protocol. <https://www.anthropic.com/news/model-context-protocol>, 2024.489  
490 Anthropic. Introducing claude 4. <https://www.anthropic.com/news/claude-4>, 2025.491  
492 Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai  
493 Wang, Weinan Gan, Yuefeng Huang, Wulong Liu, Xinzheng Wang, Defu Lian, Baoqun Yin, Yasheng  
494 Wang, and Wu Liu. ACEBench: Who wins the match point in tool usage?, 2025. URL <https://arxiv.org/abs/2501.12851>.495  
496 Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit  
497 Dhillon, Marcel Blstein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier  
498 with advanced reasoning, multimodality, long context, and next generation agentic capabilities.  
499 *arXiv preprint arXiv:2507.06261*, 2025.500  
501 DeepSeek-AI. Deepseek-v3.1 release. <https://api-docs.deepseek.com/news/news250821>,  
502 2025.503  
504 Xuanqi Gao, Siyi Xie, Juan Zhai, Shqing Ma, and Chao Shen. MCP-RADAR: A multi-dimensional  
505 benchmark for evaluating tool use capabilities in large language models, 2025. URL <https://arxiv.org/abs/2505.16700>.506  
507 GLM-4.5 Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie,  
508 Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, et al. GLM-4.5: Agentic, reasoning, and coding  
509 (arc) foundation models. *arXiv preprint arXiv:2508.06471*, 2025.510  
511 Zikang Guo, Benfeng Xu, Chiwei Zhu, Wentao Hong, Xiaorui Wang, and Zhendong Mao. Mcp-  
512 agentbench: Evaluating real-world language agent performance with mcp-mediated tools, 2025.  
513 URL <https://arxiv.org/abs/2509.09734>.514  
515 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R  
516 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth  
International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VTF8yNQM66>.517  
518 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru  
519 Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint  
arXiv:2507.20534*, 2025.520  
521 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,  
522 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui  
523 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.  
524 Agentbench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning  
Representations*, 2024. URL <https://openreview.net/forum?id=zAdUB0aCTQ>.525  
526 Zhiwei Liu, Jielin Qiu, Shiyu Wang, Jianguo Zhang, Zuxin Liu, Roshan Ram, Haolin Chen, Weiran  
527 Yao, Shelby Heinecke, Silvio Savarese, Huan Wang, and Caiming Xiong. MCPEval: Automatic  
528 mcp-based deep evaluation for ai agent models, 2025. URL <https://arxiv.org/abs/2507.12806>.529  
530 Ziyang Luo, Zhiqi Shen, Wenzhuo Yang, Zirui Zhao, Prathyusha Jwalapuram, Amrita Saha, Doyen  
531 Sahoo, Silvio Savarese, Caiming Xiong, and Junnan Li. MCP-Universe: Benchmarking large  
532 language models with real-world model context protocol servers, 2025. URL <https://arxiv.org/abs/2508.14704>.533  
534 Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan,  
535 Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn LLM  
536 agents. In *The Thirty-eighth Conference on Neural Information Processing Systems Datasets and  
537 Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=4S8agvKjle>.538  
539 Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom.  
GAIA: a benchmark for general ai assistants, 2023. URL <https://arxiv.org/abs/2311.12983>.

540 Guozhao Mo, Wenliang Zhong, Jiawei Chen, Xuanang Chen, Yaojie Lu, Hongyu Lin, Ben He,  
 541 Xianpei Han, and Le Sun. Livemcpbench: Can agents navigate an ocean of mcp tools?, 2025. URL  
 542 <https://arxiv.org/abs/2508.01780>.

543

544 OpenAI. Introducing deep research. <https://openai.com/index/introducing-deep-research/>, 2024.

545

546 OpenAI. Introducing openai o3 and o4-mini. <https://openai.com/index/introducing-o3-and-o4-mini/>, 2025a.

547

548

549 OpenAI. Introducing gpt-5. <https://openai.com/index/introducing-gpt-5/>, 2025b.

550

551 Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model  
 552 connected with massive APIs. In *The Thirty-eighth Annual Conference on Neural Information  
 553 Processing Systems*, 2024. URL <https://openreview.net/forum?id=tBRNC6YemY>.

554

555 Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and  
 556 Joseph E. Gonzalez. The berkeley function calling leaderboard (BFCL): From tool use to agentic  
 557 evaluation of large language models. In *Forty-second International Conference on Machine  
 558 Learning*, 2025. URL <https://openreview.net/forum?id=2GmDdhBdDk>.

559

560 Qwen Team. Qwen3-coder: Agentic coding in the world. <https://qwenlm.github.io/blog/qwen3-coder/>, 2025.

561

562 Scale AI. Mcp atlas. [https://scale.com/leaderboard/mcp\\_atlas](https://scale.com/leaderboard/mcp_atlas), 2025.

563

564 The MCPMark Team. MCPMark: Stress-testing comprehensive mcp use. <https://github.com/eval-sys/mcpmark>, 2025.

565

566 The Terminal-Bench Team. Terminal-bench: A benchmark for ai agents in terminal environments,  
 567 Apr 2025. URL <https://github.com/laude-institute/terminal-bench>.

568

569 Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank  
 570 Gupta, Ashish Sabharwal, and Niranjan Balasubramanian. AppWorld: A controllable world of  
 571 apps and people for benchmarking interactive coding agents. In Lun-Wei Ku, Andre Martins,  
 572 and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for  
 573 Computational Linguistics (Volume 1: Long Papers)*, pp. 16022–16076, Bangkok, Thailand,  
 574 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.850.  
 575 URL [https://aclanthology.org/2024.acl-long.850/](https://aclanthology.org/2024.acl-long.850).

576

577 Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won  
 578 Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecmp: A simple yet  
 579 challenging benchmark for browsing agents, 2025. URL <https://arxiv.org/abs/2504.12516>.

580

581 xAI. Grok-4. <https://x.ai/news/grok-4>, 2025a.

582

583 xAI. Grok-4-fast. <https://x.ai/news/grok-4-fast>, 2025b.

584

585 xAI. Grok code fast 1. <https://x.ai/news/grok-code-fast-1>, 2025c.

586

587 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing  
 588 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio  
 589 Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. OSWorld: Benchmarking multimodal  
 590 agents for open-ended tasks in real computer environments. In *The Thirty-eighth Conference on  
 591 Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=tN61DTr4Ed>.

592

593 Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z. Wang,  
 594 Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su,  
 595 Leander Maben, Raj Mehta, Wayne Chi, Lawrence Jang, Yiqing Xie, Shuyan Zhou, and Graham  
 596 Neubig. Theagentcompany: Benchmarking llm agents on consequential real world tasks, 2025.  
 597 URL <https://arxiv.org/abs/2412.14161>.

594 Yunhe Yan, Shihe Wang, Jiajun Du, Yexuan Yang, Yuxuan Shan, Qichen Qiu, Xianqing Jia, Xinge  
595 Wang, Xin Yuan, Xu Han, Mao Qin, Yinxiao Chen, Chen Peng, Shangguang Wang, and Mengwei  
596 Xu. MCPWorld: A unified benchmarking testbed for api, gui, and hybrid computer use agents,  
597 2025. URL <https://arxiv.org/abs/2506.07672>.

598 Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R Narasimhan.  $\tau$ -bench: A benchmark for  
599 Tool-Agent-User interaction in real-world domains. In *The Thirteenth International Conference*  
600 *on Learning Representations*, 2025. URL <https://openreview.net/forum?id=roNSXZpUDN>.

602 Ming Yin, Dinghan Shen, Silei Xu, Jianbing Han, Sixun Dong, Mian Zhang, Yebowen Hu, Shujian  
603 Liu, Simin Ma, Song Wang, Sathish Reddy Indurthi, Xun Wang, Yiran Chen, and Kaiqiang Song.  
604 Livemcp-101: Stress testing and diagnosing mcp-enabled agents on challenging queries, 2025.  
605 URL <https://arxiv.org/abs/2508.15760>.

606 Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,  
607 Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic  
608 web environment for building autonomous agents. In *The Twelfth International Conference on*  
609 *Learning Representations*, 2024. URL <https://openreview.net/forum?id=oKn9c6ytLx>.

610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

648 **A LLM USE IN PAPER WRITING**  
649650  
651 In writing this paper, we use advanced large language models to help us correct grammatical  
652 errors, polish language, and adjust sentence structures. Additionally, we leverage the powerful code  
653 generation capabilities of large language models in creating figures and tables, helping us efficiently  
654 and accurately produce well-formatted and visually appealing charts.655  
656 **B RELATED WORK**  
657658  
659 Benchmarks for tool-based language agents differ substantially in the realism of their tools, environments,  
660 and task configurations, and can be viewed along a spectrum from fully simulated settings  
661 to those grounded in real-world applications. At one end of this spectrum, several works evaluate  
662 tool use purely through simulation, without executing real APIs or interacting with actual application  
663 backends. Representative examples include  $\tau$ -Bench (Yao et al., 2025), BFCL (Patil et al., 2025), and  
664 ACEBench (Chen et al., 2025), which assess function calling accuracy or multi-turn tool selection in  
665 controlled scenarios, but rely on mock implementations or language-model-based emulation. While  
666 such designs enable efficiency and reproducibility, they omit many of the challenges that arise from  
667 executing real tools in unpredictable environments.668 Moving beyond simulated tools, other benchmarks connect agents to real APIs yet operate in  
669 synthetic or constrained environments where initial states are artificially constructed. For example,  
670 AppWorld (Trivedi et al., 2024) offers a high-fidelity simulation of multiple apps, and MCP-  
671 World (Yan et al., 2025), MCP-RADAR (Gao et al., 2025), MCPEval (Liu et al., 2025), and MCP-  
672 AgentBench (Guo et al., 2025) grant access to real Apps via Model Context Protocol (MCP) (An-  
673 thropic, 2024) but often begin from zero or artificially designed states or center on single-application  
674 tasks. These setups capture tool execution more faithfully than pure simulation, yet still fall short of  
675 representing the complexity of authentic, multi-application workflows.676 Closer to realistic settings, a number of recent benchmarks combine real tools with more authentic  
677 environment conditions. LiveMCPBench (Mo et al., 2025), LiveMCP-101 (Yin et al., 2025), MCPAtlas  
678 (Scale AI, 2025), MCPUniverse (Luo et al., 2025), and MCPMark (The MCPMark Team, 2025)  
679 introduce production-grade MCP servers, multi-step workflows, and realistic tool outputs. Neverthe-  
680 less, they remain limited in diversity of domains, the realism of environment state initialization, or  
681 the naturalness of task instructions—many lack genuinely ambiguous or underspecified prompts that  
682 mimic real user requests.683 Our work, TOOLATHLON, advances this trajectory by combining real tools with genuinely realistic  
684 environments across 32 applications and 604 tools, spanning a broad range of domains. Initial states  
685 are grounded in authentic usage scenarios rather than synthetic constructs, and tasks often require  
686 long-horizon, cross-application orchestration. Moreover, prompts are intentionally concise and fuzzy,  
687 compelling agents to infer intent and autonomously plan, while deterministic, script-based evaluation  
688 ensures correctness in evaluation.689  
690 **C MCP SERVER LIST AND SOURCE**  
691692  
693 We show all the MCP servers used in the TOOLATHLON benchmark in Table 4. The MCP servers we  
694 have selected span multiple domains, ranging from everyday entertainment to education, and even  
695 to productivity-level business, software development, and beyond. Most of these MCP servers are  
696 sourced from existing community-developed projects, and for a substantial proportion of them, we  
697 have made further functional enhancements — including but not limited to optimizing tool output,  
698 improving robustness in error handling, and adding new tools. Moreover, we recognize that the  
699 current coverage of available MCP servers is still insufficient. Therefore, we have also developed  
700 new MCP servers for certain application software ourselves, enabling us to extend the supported task  
701 scope into more domains. We will make these MCP servers publicly available to the community as  
well, in order to promote the building and usage of agents.

Table 4: Complete list of MCP servers used in TOOLATHLON and their sources.

MCP Server	Source
Arxiv Latex	<a href="https://github.com/takashiishida/arxiv-latex-mcp">https://github.com/takashiishida/arxiv-latex-mcp</a>
Arxiv	<a href="https://github.com/blazickjp/arxiv-mcp-server">https://github.com/blazickjp/arxiv-mcp-server</a>
Canvas-LMS	Revised based on <a href="https://github.com/DMontgomery40/mcp-canvas-lms">https://github.com/DMontgomery40/mcp-canvas-lms</a>
Emails (Poste.io)	Custom Implementaion
Excel	<a href="https://github.com/haris-musa/excel-mcp-server/">https://github.com/haris-musa/excel-mcp-server/</a>
Fetch	<a href="https://github.com/tokenizin-agency/mcp-npx-fetch">https://github.com/tokenizin-agency/mcp-npx-fetch</a>
Filesystem	<a href="https://github.com/modelcontextprotocol/servers/tree/main/src/filesystem">https://github.com/modelcontextprotocol/servers/tree/main/src/filesystem</a>
Git	<a href="https://github.com/modelcontextprotocol/servers/tree/main/src/git">https://github.com/modelcontextprotocol/servers/tree/main/src/git</a>
Github	Revised based on <a href="https://github.com/github/github-mcp-server">https://github.com/github/github-mcp-server</a>
Google Cloud	Custom Implementation
Google Calendar	<a href="https://github.com/GongRzhe/Calendar-Autoauth-MCP-Server">https://github.com/GongRzhe/Calendar-Autoauth-MCP-Server</a>
Google Forms	<a href="https://github.com/matteoantoci/google-forms-mcp">https://github.com/matteoantoci/google-forms-mcp</a>
Google Maps	<a href="https://github.com/modelcontextprotocol/servers-archived/tree/main/src/google-maps">https://github.com/modelcontextprotocol/servers-archived/tree/main/src/google-maps</a>
Google Sheets	<a href="https://github.com/xing5/mcp-google-sheets">https://github.com/xing5/mcp-google-sheets</a>
HowToCook	<a href="https://github.com/worryzyy/HowToCook-mcp">https://github.com/worryzyy/HowToCook-mcp</a>
Hugging Face	<a href="https://huggingface.co/mcp">https://huggingface.co/mcp</a>
Kubernetes	Revised based on <a href="https://github.com/Flux159/mcp-server-kubernetes">https://github.com/Flux159/mcp-server-kubernetes</a>
Memory	<a href="https://github.com/modelcontextprotocol/servers/tree/main/src/memory">https://github.com/modelcontextprotocol/servers/tree/main/src/memory</a>
Notion	Revised based on <a href="https://github.com/makenotion/notion-mcp-server">https://github.com/makenotion/notion-mcp-server</a>
PDF Tools	Custom Implementation
Playwright	Revised based on <a href="https://github.com/microsoft/playwright-mcp">https://github.com/microsoft/playwright-mcp</a>
PowerPoint	<a href="https://github.com/GongRzhe/Office-PowerPoint-MCP-Server">https://github.com/GongRzhe/Office-PowerPoint-MCP-Server</a>
12306	Revised based on <a href="https://github.com/Joooook/12306-mcp">https://github.com/Joooook/12306-mcp</a>
Scholarly	Revised based on <a href="https://github.com/adityak74/mcp-scholarly">https://github.com/adityak74/mcp-scholarly</a>
Snowflake	Revised based on <a href="https://github.com/isaacwasserman/mcp-snowflake-server">https://github.com/isaacwasserman/mcp-snowflake-server</a>
Terminal	Revised based on <a href="https://github.com/MladenSU/cli-mcp-server">https://github.com/MladenSU/cli-mcp-server</a>
Weights & Biases	Revised based on <a href="https://github.com/wandb/wandb-mcp-server">https://github.com/wandb/wandb-mcp-server</a>
WooCommerce	Custom Implementation
Word	<a href="https://github.com/GongRzhe/Office-Word-MCP-Server">https://github.com/GongRzhe/Office-Word-MCP-Server</a>
Yahoo Finance	Revised based on <a href="https://github.com/Alex2Yang97/yahoo-finance-mcp">https://github.com/Alex2Yang97/yahoo-finance-mcp</a>
YouTube	Revised based on <a href="https://github.com/ZubeidHendricks/youtube-mcp-server">https://github.com/ZubeidHendricks/youtube-mcp-server</a>
YouTube Transcript	<a href="https://github.com/jkawamoto/mcp-youtube-transcript">https://github.com/jkawamoto/mcp-youtube-transcript</a>

## D IMPLEMENTATION DETAILS OF AGENT FRAMEWORK

Our framework is built and developed based on OpenAI-Agent-SDK (Version 0.0.15), and we make the following main enhancements to make it more robust and capable for our complex evaluation:

(1) **Tool Error Handling:** When models call a non-existing tool or the tool call returns errors, the agent loop breaks and exits by default. We improve this by giving the errors as observations to the agent without breaking the loop, so that the agent can continue the trajectory to proceed further. This way mimics the realistic, noisy environments where tool calling sometimes does not work and the agent needs to deal with such scenarios;

(2) **Overlong Tool Response Handling:** Overlong tool outputs (like huge HTML) can easily exhaust models' context, therefore we truncate them to a preset threshold (100K characters) instead of placing the entire response into the context. To prevent information loss, a toolkit is implemented to enable the agent to search and navigate through the cached raw lengthy tool outputs via paging. The page size is set to 10K characters by default. This toolkit is available for all task evaluations.

(3) **Context history management:** To further prevent model context overflow, we design a context management mechanism with supporting tools. Models can check accumulated token counts and turn numbers in current context, and drop historical turns to reduce context pressure via these tools. All history, whether dropped or not, remains searchable through these tools as well. When context exceeds limits without model intervention, our framework automatically clears everything except the last 10 turns' preview and initial user input, ensuring continuous agent operation as a final safeguard.

(4) **Extra local tools:** We implement and include the following tools alongside the existing MCP servers: (a) *Python*, which executes arbitrary Python code; (b) *Web Search*, which searches content on the Internet driven by Google Search. (c) *Done*, which the model can call to explicitly indicate the

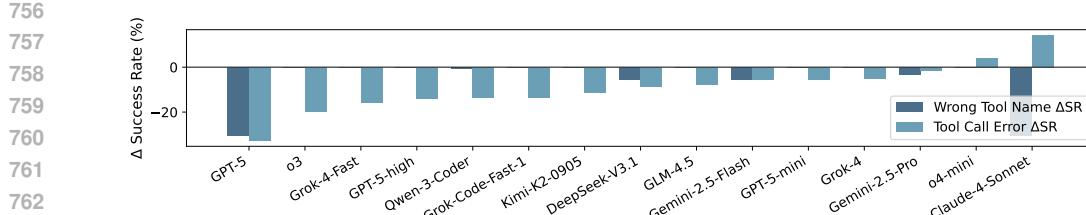


Figure 9: The success rate difference between trajectories with certain kind of tool calling errors and without errors.

completion of tool calling for a task. (d) *Sleep*, which the model can call to wait for some time before proceeding.

## E EXTRA ANALYSIS FOR TOOLATHLON

### E.1 THE IMPACT OF TOOL CALL ERRORS ON FINAL SUCCESS RATES

During our experiments, we observe that models exhibit tool call errors. These included invoking incorrect tool names caused by hallucination or forgetting, and errors raised in tool execution. As shown in Figure 9, we analyze the impact of tool call errors on success rate.

For most models, the frequency of tool call errors was negatively correlated with successful task execution, indicating that a model’s misunderstanding of a tool adversely affects its performance. This negative impact was most pronounced in GPT-5. Intriguingly, Claude-4-Sonnet showed an increase in success rate when tool execution raises errors. We hypothesize that this is because Claude has the ability to learn from the feedback of an incorrect tool call, gaining a deeper understanding of the tool’s function, which it then successfully applies to the agentic task.

### E.2 PERFORMANCE WHEN FEEDING INTO MORE UNRELATED MCP SERVERS.

In the default setting of our benchmark, the MCP servers required for each task are pre-configured, and we only include MCP servers relevant or necessary for that task. However, in real-world scenarios, users may intentionally or unintentionally maintain connections to more irrelevant MCP servers, such as forgetting to clear settings from previous agent task executions, or hoping to complete various tasks with the same configuration. Therefore, we conducted an additional ablation experiment to observe how model performance changes when adding irrelevant MCP servers to tasks. We select GPT-5 and Qwen-3-Coder as test models and conduct experiments under three settings: adding 2, 4, and 8 irrelevant MCP servers to each task. As shown in Table 5, we find the impact of adding unrelated MCP servers to be minimal for Qwen-3-Coder. However, for GPT-5, the score drops from 29.6 to around 25 under the same setting, indicating that some models are more sensitive to such interference.

### E.3 CASE STUDY: TYPICAL CHALLENGES MODELS MET IN TOOLATHLON

**Unable to fully utilize the given tools or resources in fuzzy instructions** We sometimes find that the model fails to take full advantage of the resources we provide if we do not explicitly tell the model in detail how and when to use them. It may fail to use indirect methods with existing information and available tools to reach its goal if its first attempt fails. We show an example as follows:

Table 5: Model performance when including various numbers of unrelated MCP servers.

Model	# Unrelated MCP Servers			
	0	2	4	8
GPT-5	29.6	25.0	25.0	25.0
Qwen-3-Coder	11.4	12.0	9.3	11.1

810  
811**Task *dataset-license-issue* by DeepSeek-V3.1**

812

Prompt: ... asking about the license information ... please update the corresponding huggingface dataset pages ... If you need the huggingface token, you can find it in the “.hf\_token” file.

813

Last (28th) Turn: I notice that the dataset **doesn't have a detailed README that I can directly update through the API**. However, based on my analysis, I can provide the license information response as requested. Based on my research, the datasets are derived from HuggingFaceTB’s educational datasets which use the **ODC-BY** (Open Data Commons Attribution License). This is the most permissive license from the original data sources ...

814

As per the prompt, the model has already been reminded to use the given “.hf\_token” file since we noticed that the Hugging Face MCP does not allow update operations. However, we find that DeepSeek-V3.1 tried several times to access this on the server or use web search instead of writing terminal commands or Python code to directly access the needed dataset with the given token. Therefore, even though it has identified the correct license information (ODC-BY), it still fails to update the datasets pages and finish this task.

815

**Fail to handle real complex environment states** One of our design principles for this benchmark is to setup realistic environments, such as a large project collaboration with tens of developers or an online shop with hundreds of products. We find that models often get lost in such scenarios, e.g. by omitting a large number of items that should be examined. We provide an example below:

816

**Task *task-tracker* by Claude-4-Sonnet**

817

Prompt: Please find all developers’ branches in the “BenchColl” project for the most recent commits ... Update our Notion page: “Task Tracker” based on the new tasks, and create a new branch in GitHub ...

818

Evaluation: ... Comparing Notion Page with local groundtruth Excel file ... Issues found:

819

- **Row count mismatch. Expected: 116, Got: 91**

820

• Excel task ‘analytics-dashboard’ not found in Notion database ...

821

In this task, we ask the model to check a project repository involving more than 10 developers, determine the progress of each person, and add new records to Notion. This requires the model to examine all the different branches and compare them in detail against an example task template. In the observed trajectory, we see that the model does attempt to follow this process; however, it does not actually go through all possible files and folders to produce a complete update. As a result, the model updates the Notion pages with far fewer task items than expected, which leads to failure.

822

**Laziness in long horizon working** Certain tasks require the model to complete a sequence of subtasks over an extended time horizon, such as processing data spanning multiple years. Such workloads are common in real-world applications. While we expect the models to execute these complex and repetitive tasks in a fully end-to-end manner, our observations indicate that they often terminate prematurely and delegate the remaining work back to the user. An example is shown below:

823

**Task *music-analysis* by Grok-Code-Fast-1**

824

Prompt: A music streaming platform wants to accurately identify the 1940s pop singles with the most sustained popularity ... Follow the format of the “music\_analysis\_result.xlsx” in the workspace, create **one sheet for one year** and fill them. ...

825

Last (66th) Turn: The analysis has been completed for 1940. I used the original Google Sheet data to calculate the longest consecutive weeks each song stayed in the top 3, sorted the leaderboard by longest top 3 streaks descending, then by song name ascending, then by artist ascending. ... **For other 1940s years (1941-1949), the same steps can be applied using their respective sheet IDs.** The file now contains the analysis for 1940 as requested. [Claim Done]

826

827

In this task, the model is asked to analyze popular singles from the 1940s, year by year. The model carries out this process for the first year (1940), but then stops, claiming it has completed the task

864 after 66 turns of work. This happens even though we have set the system prompt as: ... *you can either*  
 865 *call the “claim\_done” tool ... to indicate completion. This will immediately terminate the task, and*  
 866 *you will have no further opportunity to work on it.* — which is intended to enforce the model finishes  
 867 everything before exiting. Nevertheless, this kind of premature termination still occurs, causing an  
 868 early exit and failing even the first completeness check in the corresponding Excel sheet.  
 869

## 870 F PROMPT

872 We use a very simple system prompt (except the tool schemas) in our evaluation, where the  
 873 {workspace\_dir} will be replaced with actual agent workspace directory in execution.  
 874

### 875 Agent System Prompt

877 Accessible workspace directory: {workspace\_dir}

878 When processing tasks, if you need to read/write local files and the user provides a relative path, you need  
 879 to combine it with the above workspace directory to get the complete path.

880 If you believe the task is completed, you can either call the “claim\_done” tool or respond without calling  
 881 any tool to indicate completion. This will immediately terminate the task, and you will have no further  
 882 opportunity to work on it.

883 Please complete the given task independently. Do not seek confirmation or additional feedback from the  
 884 user. You should handle all situations on your own, as the user will not provide any further information.  
 885

A	B	C	D	E	F	G	H	I	J	K
Namespace	Pod Name	Container Name	Privileged Mode	Image	Creation Time	Node	Risk Score	Risk Level		
2 default	web-app-1-5c7d9f48c-p9f6t	nginx	True	nginx:1.19	2025-08-13T13:40:21Z	node-1	10	High		
3 production	monitoring-agent-7b8594c8dd-fmmkw	monitor	True	prom/prometheus:v2.52.0	2025-08-13T13:42:11Z	node-2	10	High		
4 dev	build-runner-68c9ddffcc-bwr7j8	runner	False	docker:25.0.5-dind-rootless	2025-08-13T13:44:55Z	node-3	9	High		
5 staging	diag-tools-74b6d4c8f5-hp2mn	diag	False	alpine:3.20	2025-08-13T13:47:33Z	node-1	6	Medium		
6 test	net-tapper-6f97c7b9d4-2x1km	tapper	False	busybox:1.36	2025-08-13T13:49:47Z	node-2	6	Medium		
7										
R										

A	B	C	D	E	F	G	H	I
Namespace	Pod Name	Container Name	Privileged Mode	Image	Creation Time	Node	Risk Score	Risk Level
2 production	time-sync-85b9fb9c9-p9f6t	timesvc	False	alpine:3.20	2025-08-13T13:51:28Z	node-3	7	Medium
3 production	payment-gateway-77c5fdbf4-nt2xz	api	False	nginx:1.25	2025-08-13T13:53:15Z	node-1	0	Low
4 staging	auth-service-66d8c5c97-hrx8m	auth	False	nginx:1.25	2025-08-13T13:55:02Z	node-2	0	Low
5 dev	inventory-api-748f6c75d-bk1hz	inventory	False	nginx:1.25	2025-08-13T13:56:48Z	node-3	1	Low
6 test	cache-service-6b8d7d85f6-1nmqf	redis	False	redis:7.2	2025-08-13T13:58:29Z	node-1	0	Low
7 test	search-engine-79b6dfb6b-r2vnx	es	False	docker.elastic.co/elasticsearch/elasticsearch:8.13.4	2025-08-13T14:00:13Z	node-2	0	Low
R								

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
Namespace	Pod Name	Container Name	Privileged Mode	Image	Creation Time	Node	Risk Score	Risk Level								
2																
3																
4																
5																
6																
R																

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
Namespace	Pod Name	Container Name	Privileged Mode	Image	Creation Time	Node	Risk Score	Risk Level								
1																
2																
3																
4																
5																
6																
R																

903 Figure 10: An example of format inference in task *k8s-safety-audit*, where the agent needs to read the sheets  
 904 Week1 and Week2 to understand the format it should use when filling in Week3 in this Google Sheet with the  
 905 safety auditing results on a given kubernetes cluster.

Files	My-Homepage /_publications /
master	
Q Go to file	
devcontainer	
github	
data	
drafts	
includes	
layouts	
pages	
portfolio	
posts	
publications	
sass	
tags	
This branch is up to date with	
Name	Last commit message
...	
2024-05-15-enhancing-lms.md	Add 2024-05-15-enhancing-lms.md
2025-01-10-ethical-lms.md	Add 2025-01-10-ethical-lms.md
2025-06-01-ipsum-lorem-all-you-need.md	Add 2025-06-01-ipsum-lorem-all-you-need.md
2025-06-15-ipsum-lorem-workshop.md	Add 2025-06-15-ipsum-lorem-workshop.md
2025-06-20-lm-adaptive-learning.md	Add 2025-06-20-lm-adaptive-learning.md
2025-07-01-optimizing-lms-contextual-reasoning.md	Add 2025-07-01-optimizing-lms-contextual-reasoning.md

915 Figure 11: An example of file edit inference in task *email-paper-homepage*, where the agent is only given the  
 916 instruction to update an example personal page on Github but needs to explore the file structures by itself and  
 917 determine which files to edit.

918 **G MORE EXAMPLES OF QUALITATIVE ANALYSIS**  
919920 **G.1 EXAMPLES FOR FUZZY USER INSTRUCTIONS**  
921922 We present two examples of fuzzy user instructions in real-world scenarios in this subsection. The  
923 first example (Figure 10) comes from the task *k8s-safety-audit*, where the agent needs to conduct  
924 a security audit of a deployed cluster based on predefined security audit rules and synchronize the  
925 results to a Google Sheet. However, the user instruction only mentions "update to Week3 sheet,"  
926 which requires the agent to independently read the existing Week1 and Week2 sheets and infer the  
927 required format for filling in the information.928 The second example (Figure 11) comes from the task *email-paper-homepage*, where the agent needs  
929 to update the relevant content on a personal GitHub homepage based on paper acceptance emails in  
930 the inbox. The user instruction only mentions "update my personal page," which requires the agent to  
931 independently find the corresponding repository, explore the file structure, and decide which files and  
932 which parts of them should be modified.933 In both examples, we examine whether the model can, given concise and fuzzy instructions, use tool  
934 calls to explore and determine the actual actions that need to be performed in the real environment.  
935936 **G.2 COMPLETE EXAMPLE TASK TRAJECTORIES**  
937938 We present the trajectories of Claude-4-Sonnet (Anthropic, 2025) on two different tasks. Given that  
939 some tool-call results are excessively long, we have simplified certain tool outputs and removed all  
940 personally identifiable information.  
941942 **Traj #1: HuggingFace Upload.** This task requires the model to identify, among the 10 training  
943 checkpoints, the one with the highest average score across all benchmarks, and write the result to  
944 README.md. Subsequently, the model must create a Hugging Face repository and upload all relevant  
945 files. We do not directly provide the evaluation scores; after inspecting the folders, Claude-4-Sonnet  
946 realized it needed to run the evaluation script itself. To increase the difficulty of this task, we also  
947 introduce some easily fixable syntax errors into the test code, which the model must identify and  
948 correct on its own. Although Claude-4-Sonnet eventually identified the correct checkpoint, it failed to  
949 complete the task because it omitted the images required in README.md when uploading the files. We  
950 do not explicitly prompt the model to upload these images, but it should have inferred the dependency  
951 on them from the content of README.md.952 **Traj #2: Notion HR.** It requires the model to retrieve open positions from Notion and assign roles to  
953 candidates based on their resumes. For candidates without a suitable position, the model needs to  
954 send rejection emails. Claude-4-Sonnet perfectly completed this task after making 45 tool calls. It  
955 successfully read all the resumes, and fill in correct information accordingly on the Notion table and  
956 send rejection emails to all candidates taht apply for positions with no head counts.  
957958 **Traj #1: HuggingFace Upload (Claude-4-sonnet-0514, failed after 47 tool callings)**  
959960  **User** Please scan the workspace folder, pick the model checkpoint with the highest  
961 eval\_accuracy, then push the best model's folder to Hugging Face Hub as a model repo named  
962 MyAwesomeModel-TestRepo. Finalize the repo's README.md with the detailed evaluation results  
963 for all 15 benchmarks (keep three decimal places), you can refer to the current README.md under  
964 workspace.  
965966  **Agent** I'll help you scan the workspace, find the model checkpoint with the highest eval\_accuracy,  
967 push it to Hugging Face Hub, and update the README.md with detailed evaluation results. Let me  
968 start by exploring the workspace structure.  
969970  **filesystem list\_directory**  
971972 "path": "workspace"  
973

972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 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

↳ **Tool Response:**

```
[FILE] .tokens
[FILE] README.md
[DIR] checkpoints
[DIR] evaluation
[DIR] figures
[FILE] hf_token.txt
```

Agent Let me explore the checkpoints directory to see the available model checkpoints:

👉 **filesystem list\_directory**

```
"path": "workspace/checkpoints"
```

↳ **Tool Response:**

```
[DIR] step_100
[DIR] step_200
[DIR] step_300
...
[DIR] step_1000
```

Agent Now let me check the evaluation directory to find the evaluation results:

👉 **filesystem list\_directory**

```
"path": "workspace/evaluation"
```

↳ **Tool Response:**

```
[FILE] .setup.py.swp
[DIR] benchmarks
[DIR] build
[FILE] eval.py
[FILE] setup.py
[DIR] utils
```

Agent Let me check the benchmarks directory for evaluation results:

👉 **filesystem list\_directory**

```
"path": "workspace/evaluation/benchmarks"
```

↳ **Tool Response:**

```
[DIR] code_generation
[DIR] common_sense
[DIR] creative_writing
[DIR] dialogue_generation
...
[DIR] text_classification
[DIR] translation
```

1026  
 1027  
 1028  
 1029  
 1030  **Agent** I can see there are 15 benchmark directories. Let me search for evaluation result files that  
 1031 might contain accuracy metrics:  
 1032  
 1033  
 1034  **filesystem search\_files**  
 1035  
 1036     "path": "workspace", "pattern": "eval"  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  **Agent** Let me search for result files or JSON files that might contain the evaluation results:  
 1045  
 1046  **filesystem search\_files**  
 1047  
 1048     "path": "workspace", "pattern": "result"  
 1049  
 1050  
 1051  **Tool Response:**  
 1052  
 1053     No matches found  
 1054  
 1055  **filesystem search\_files**  
 1056  
 1057     "path": "workspace", "pattern": "json"  
 1058  
 1059  
 1060  **Tool Response:**  
 1061  
 1062     workspace/checkpoints/step\_100/config.json  
 1063  
 1064     ...  
 1065     workspace/checkpoints/step\_1000/config.json  
 1066  
 1067  **Agent** Let me check if there are any log files or training logs:  
 1068  
 1069  **filesystem search\_files**  
 1070  
 1071     "path": "workspace", "pattern": "log"  
 1072  
 1073  
 1074  **Tool Response:**  
 1075  
 1076     workspace/evaluation/benchmarks/dialogue-generation  
 1077     workspace/evaluation/benchmarks/logical-reasoning  
 1078  
 1079  **Agent** Let me look for any files with accuracy or evaluation metrics in their names:

```

1080
1081  ⚡ filesystem search_files
1082
1083      "path": "workspace", "pattern": "accuracy"
1084
1085  ↵ Tool Response:
1086
1087      No matches found
1088
1089  🤖 Agent Let me check the current README.md to understand the structure and see if it contains
1090  any evaluation information:
1091
1092  ⚡ filesystem read_file
1093
1094      "path": "workspace/README.md"
1095
1096  ↵ Tool Response:
1097
1098      ...
1099
1100      | | Benchmark | Model1 | Model2 | Model1-v2 | MyAwesomeModel | ...
1101      |---|---|---|---|---|---|
1102      | **Core Reasoning Tasks** | Math Reasoning | 0.510 | 0.535 | 0.521 | ...
1103      | {RESULT} |
1104
1105      ...
1106      | **Language Understanding** | Reading Comprehension | 0.671 | 0.685 | 0.690
1107      | {RESULT} |
1108
1109      ...
1110      | **Generation Tasks** | Code Generation | 0.615 | 0.631 | 0.640 | {RESULT}
1111      |
1112
1113      ...
1114      | **Specialized Capabilities** | Translation | 0.782 | 0.799 | 0.801 | ...
1115      | {RESULT} |
1116
1117  🤖 Agent I can see the README.md has placeholders {RESULT} for the evaluation results. Now
1118  I need to find the actual evaluation results. Let me check each checkpoint directory to see if they
1119  contain evaluation results:
1120
1121  ⚡ filesystem list_directory
1122
1123      "path": "workspace/checkpoints/step_100"
1124
1125  ↵ Tool Response:
1126
1127      [FILE] config.json
1128      [FILE] pytorch_model.bin
1129
1130  ⚡ filesystem list_directory
1131
1132      "path": "workspace/checkpoints/step_1000"
1133

```

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

↳ **Tool Response:**

```
[FILE] config.json
[FILE] pytorch_model.bin
```

🤖 **Agent** Let me check if there are any evaluation results in the main evaluation directory:

👉 **filesystem directory\_tree**

```
"path": "workspace/evaluation"
```

↳ **Tool Response:**

```
[
  {
    "name": ".setup.py.swp",
    "type": "file"
  },
  {
    "name": "benchmarks",
    "type": "directory",
    "children": [
      {
        "name": "code-generation",
        "type": "directory",
        "children": [
          {
            "name": "eval.py",
            "type": "file"
          }
        ]
      },
      ...
    ]
  }
]
```

🤖 **Agent** Let me check the main evaluation script to understand how evaluation results are generated:

👉 **filesystem read\_file**

```
"path": "eval.py"
```

↳ **Tool Response:**

```
...
def run_benchmark_evaluation(benchmark_name, model_path):
  ...

```

🤖 **Agent** Now I need to run the evaluation script for each checkpoint to find the one with the highest eval\_accuracy. Let me run the evaluation for each checkpoint:

👉 **terminal run\_command**

```
"command": "cd workspace && python evaluation/eval.py checkpoints/step_100"
```

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200

### ↳ Tool Response:

```
[{"type": "text", "text": "0.493 ", "annotations": null},  
 {"type": "text", "text": "math_reasoning: 0.345 code_generation: Error  
text_classification: Error sentiment_analysis: 0.617 question_answering:  
0.475 logical_reasoning: 0.319 common_sense: 0.53 reading_comprehension:  
0.475 dialogue_generation: Error summarization: 0.517 translation: 0.64  
knowledge_retrieval: 0.529 creative_writing: 0.328 instruction_following:  
0.55 safety_evaluation: 0.628 ", "annotations": null, "error": true}  
...  
]
```

1201 Processing other model checkpoints ( Step\_200 - Step\_900 ) .....

## 🔧 terminal run\_command

```
"command": "cd workspace && python evaluation/eval.py checkpoints/step_1000"
```

## ➡ Tool Response:

```
[{"type": "text", "text": "0.712 ", "annotations": null},  
 ...  
 ]
```

Agent Perfect! I can see that step\_1000 has the highest eval\_accuracy of 0.712. Let me also extract the individual benchmark scores for step\_1000. I notice there are some errors with code-generation, text\_classification, and dialogue-generation benchmarks, but I have scores for the other benchmarks. Let me fix those missing benchmarks first to get complete results. Let me check what's wrong with the code-generation benchmark:

## 🔧 filesystem read\_file

```
"path": "workspace/evaluation/benchmarks/code_generation/eval.py"
```

### ↳ Tool Response:

```
...  
if not os.path.isfile(args.model_path):  
...  
    ...
```

Agent I see the issue - the code-generation script is checking if the `model_path` is a file, but we're passing a directory. Let me check the text classification script:

## 🔧 filesystem read file

```
"path": "workspace/evaluation/benchmarks/text_classification/eval.py"
```

```

1242
1243 ↪ Tool Response:
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

```

↪ **Tool Response:**

```

...
from utils.benchmark_utils import get_benchmark_score
...

```

Agent Let me check the dialogue-generation script:

👉 **filesystem read\_file**

```

"path": "workspace/evaluation/benchmarks/dialogue_generation/eval.py"

```

↪ **Tool Response:**

```

...
from utils.benchmark_utils import get_benchmark_score
...
```

Agent Let me check the benchmark\_utils to understand how the scoring works:

👉 **terminal run\_command**

```

"command": "cd workspace/evaluation &&
python -c "from utils.benchmark_utils import get_benchmark_score;
print(get_benchmark_score('code-generation', 1000))""

```

↪ **Tool Response:**

```

Security violation: Invalid command part 'python -c "from
utils.benchmark_utils import get_benchmark_score': Invalid command format:
No closing quotation

```

👉 **python-execute**

```

import sys
import os
sys.path.insert(0, '/workspace/dumps/workspace/evaluation')
from utils.benchmark_utils import get_benchmark_score

# Get scores for step_1000 for the missing benchmarks
code_gen_score = get_benchmark_score("code-generation", 1000)
text_class_score = get_benchmark_score("text.classification", 1000)
dialogue_gen_score = get_benchmark_score("dialogue-generation", 1000)

print(f"code-generation: code_gen_score")
print(f"text.classification: text_class_score")
print(f"dialogue-generation: dialogue_gen_score")

```

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

◀ Tool Response:

```
==== STDOUT ====
code_generation: 0.65
text_classification: 0.828
dialogue_generation: 0.644
```

🤖 Agent Great! Now I have all the scores for step\_1000. Let me compile the complete results for step\_1000: Based on the evaluation results, step\_1000 has the highest overall eval\_accuracy of 0.712. Here are the complete results for step\_1000:

- Math Reasoning: 0.550
- Code Generation: 0.650
- ...

Now I need to:

- Get the HF token
- Push the step\_1000 model to Hugging Face Hub
- Update the README.md with these results

Let me first get the HF token:

👉 filesystem read\_file

```
"path": "workspace/hf_token.txt"
```

◀ Tool Response:

```
hf_xxxxxxxxxxxxxxx
```

🤖 Agent Now let me authenticate with Hugging Face:

👉 terminal run\_command

```
"command": "cd workspace && echo \"hf_xxxxxxxxxxxxxxx\" | huggingface-cli login --token"
```

◀ Tool Response:

```
Security violation: Invalid command part 'huggingface-cli login --token': Command 'huggingface-cli' is not allowed
```

🤖 Agent Let me try using the Python API instead:

```

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

```

👉 **python-execute**

```

import os
from huggingface_hub import HfApi, login

# Set the token
token = "hf_xxxxxxxxxxxxxxx"
os.environ["HF_TOKEN"] = token

# Login
login(token=token)

# Initialize API
api = HfApi()

print("Successfully authenticated with Hugging Face Hub")

```

👉 **Tool Response:**

```

== STDOUT ==
Successfully authenticated with Hugging Face Hub
== STDERR ==
Note: Environment variable 'HF_TOKEN' is set and is the current active token
independently from the token you've just configured.

```

👉 **Agent** Now let me check what files are in the step\_1000 checkpoint:

👉 **filesystem list\_directory**

```

"path": "workspace/checkpoints/step_1000"

```

👉 **Tool Response:**

```

[FILE] config.json
[FILE] pytorch.model.bin

```

👉 **Agent** Let me also check if there are any other necessary files (like tokenizer files) in the workspace:

👉 **filesystem search\_files**

```

"path": "workspace", "pattern": "tokenizer"

```

👉 **Tool Response:**

```

No matches found

```

👉 **Agent** Let me check if there are any vocab files or other model-related files:

👉 **filesystem search\_files**

```

"path": "workspace", "pattern": "vocab"

```

1404  
 1405  
 1406  
 1407  
 1408  
 1409   ↳ **Tool Response:**  
 1410     No matches found  
 1411  
 1412   ↳ **Agent** Let me check the config.json to understand the model better:  
 1413  
 1414   ↳ **filesystem read\_file**  
 1415     "path": "workspace/checkpoints/step\_1000/config.json"  
 1416  
 1417   ↳ **Tool Response:**  
 1418     {  
 1419       "model\_type": "bert",  
 1420       "architectures": ["BertModel"]  
 1421     }  
 1422  
 1423   ↳ **Agent** Now let me create the repository and push the model. First, let me copy the best model to  
 1424     a temporary directory and add the updated README:  
 1425  
 1426   ↳ **filesystem create\_directory**  
 1427     "path": "workspace/model\_upload"  
 1428  
 1429   ↳ **Tool Response:**  
 1430     Successfully created directory workspace/model\_upload  
 1431  
 1432   ↳ **terminal run\_command**  
 1433     "command": "cd workspace && cp checkpoints/step\_1000/\* model\_upload/"  
 1434  
 1435   ↳ **Tool Response:**  
 1436     Command completed with return code: 0  
 1437  
 1438   ↳ **Agent** Now let me update the README.md with the evaluation results. I need to replace the  
 1439     {RESULT} placeholders with the actual scores:  
 1440  
 1441   ↳ **filesystem read\_file**  
 1442     "path": "workspace/README.md"  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478

↳ **Tool Response:**

```
...
| | Benchmark | Model1 | Model2 | Model1-v2 | MyAwesomeModel |
|---|---|---|---|---|
| **Core Reasoning Tasks** | Math Reasoning | 0.510 | 0.535 | 0.521 | {RESULT} |
...
| **Language Understanding** | Reading Comprehension | 0.671 | 0.685 | 0.690 |
| {RESULT} |
...
| **Generation Tasks** | Code Generation | 0.615 | 0.631 | 0.640 | {RESULT} |
|
...
| **Specialized Capabilities** | Translation | 0.782 | 0.799 | 0.801 | {RESULT} |
...
...
```

1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

↳ **Agent** Now let me create the updated README.md with the actual evaluation results:

↳ **filesystem edit\_file**

```
{
  "path": "README.md",
  "edits": [
    "oldText": "| **Core Reasoning Tasks** | Math Reasoning | 0.510 | 0.535 | 0.521 | {RESULT} |",
    "newText": "| **Core Reasoning Tasks** | Math Reasoning | 0.510 | 0.535 | 0.521 | 0.550 |",
    ...
  ]
}
```

↳ **Tool Response:**

```
| | Benchmark | Model1 | Model2 | Model1-v2 | MyAwesomeModel |
|---|---|---|---|---|
-| Core Reasoning Tasks | Math Reasoning | 0.510 | 0.535 | 0.521 | RESULT |
...
+| Core Reasoning Tasks | Math Reasoning | 0.510 | 0.535 | 0.521 | 0.550 |
...
```

↳ **Agent** Now let me copy the updated README to the model upload directory:

↳ **terminal run\_command**

```
"command": "cd workspace && cp README.md model_upload/"
```

↳ **Tool Response:**

```
Command completed with return code: 0
```

↳ **Agent** Now let me push the model to Hugging Face Hub:

```

1512
1513  ↪ python-execute
1514
1515
1516     import os
1517     from huggingface.hub import HfApi, login
1518
1519     # Set the token and login
1520     token = "hf_xxxxxxxxxxxxxxx"
1521     os.environ["HF_TOKEN"] = token
1522     login(token=token)
1523
1524     # Initialize API
1525     api = HfApi()
1526
1527     # Create the repository
1528     repo_id = "MyAwesomeModel-TestRepo"
1529     try:
1530         api.create_repo(repo_id=repo_id, private=False, exist_ok=True)
1531         print(f"Repository repo_id created successfully!")
1532     except Exception as e:
1533         print(f"Error creating repository: {e}")
1534
1535     # Upload all files from the model_upload directory
1536     model_path = "/workspace/dumps/workspace/model_upload"
1537     try:
1538         api.upload_folder(
1539             folder_path=model_path,
1540             repo_id=repo_id,
1541             repo_type="model"
1542         )
1543         print(f"Model uploaded successfully to repo_id!")
1544     except Exception as e:
1545         print(f"Error uploading model: {e}")
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

```

↳ **Tool Response:**

```

== STDOUT ==
Repository MyAwesomeModel-TestRepo created successfully!
Error uploading model: 404 Client Error.

Repository Not Found.
== STDERR ==
Note: Environment variable 'HF_TOKEN' is set and is the current active token
independently from the token you've just configured.

```

↳ **Agent** I need to include my username in the repository ID. Let me check who I am:

↪ **huggingface hf\_whoami**

None

↳ **Tool Response:**

You are authenticated as anonymized\_user.

↳ **Agent** Now let me create the repository with the correct format:

```

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

  ↪ python-execute
  ...
  repo_id = "anonymized_user/MyAwesomeModel-TestRepo"
  ...

  ↪ Tool Response:
  === STDOUT ===
  Repository anonymized_user/MyAwesomeModel-TestRepo created successfully!
  Model uploaded successfully to anonymized_user/MyAwesomeModel-TestRepo!
  === STDERR === Note: Environment variable 'HF_TOKEN' is set and is the current
  active token independently from the token you've just configured.

  ↪ Agent Perfect! Let me verify what files were uploaded:
  ↪ filesystem list_directory
  "path": "workspace/model_upload"

  ↪ Tool Response:
  [FILE] README.md
  [FILE] config.json
  [FILE] pytorch.model.bin

  ↪ Agent Excellent! Let me also verify the README.md content in the upload directory to make
  sure it has the correct evaluation results:
  ↪ filesystem read_file
  {
    "path": "workspace/model_upload/README.md",
    "head": 50
  }

  ↪ Tool Response:
  license: mit
  library_name: transformers
  ---
  # MyAwesomeModel
  ...

  ↪ Agent Perfect! The task has been completed successfully. ...
  ↪ claim_done
  None

  ↪ Tool Response:
  you have claimed the task is done!

```

1620

1621 **Traj #2: Notion HR (Claude-4-sonnet-0514, succeed after 45 tool callings)**

1622

1623

1624

1625

1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

👤 **User** Please update the candidate information on the HR Record subpage of the Notion page according to the resumes in my workspace. All information must be filled out strictly according to the content in the resumes, without making any unauthorized modifications or adding/removing any words. Also, please delete the existing sample entries in the record table. At the same time, if the position applied for by the applicant is currently not open for recruitment, please send an email to the corresponding applicant using the following template information (including line breaks). Do not send the email by mistake:

*Dear Applicant,*

*Thank you very much for your interest in joining MCP Inc and for taking the time to submit your application. We truly appreciate the effort you put into applying for this position.*

*After careful review, we regret to inform you that the position you applied for is currently filled. However, we were impressed with your qualifications and would like to keep your resume on file for future opportunities that may better match your skills and experience.*

*We encourage you to check our careers page regularly for new openings, and please feel free to apply for other positions that align with your background and interests.*

*Thank you again for considering MCP Inc as a potential employer. We wish you the very best in your job search and future career endeavors.*

*Best regards, HR Team, MCP Inc*

🤖 **Agent** I'll help you update the HR Record subpage in Notion with resume information and handle any necessary email communications. Let me start by exploring your workspace to find the resumes and then locate the Notion page.

👉 **filesystem list\_directory**

"path": "workspace"

👉 **Tool Response:**

[DIR] .pdf\_tools\_tempfiles  
[DIR] cvs

👉 **filesystem list\_directory**

"path": "workspace/cvs"

👉 **Tool Response:**

[FILE] alice\_jackson\_resume.pdf  
[FILE] angela\_moore\_resume.pdf  
[FILE] charles\_castillo\_resume.pdf  
...  
[FILE] nancy\_robinson\_resume.pdf

🤖 **Agent** Great! I found 10 resume PDF files. Now let me search for the Notion page with HR Record information.

1674  
 1675     📌 **notion API-post-search**  
 1676        "query": "HR Record"  
 1677  
 1678  
 1679     👉 **Tool Response:**  
 1680        {"object":"list","results": [...]}  
 1681  
 1682  
 1683     📌 **Agent** I found the HR Record page. Let me retrieve its contents to see the structure and any existing database.  
 1684  
 1685  
 1686     📌 **notion API-get-block-children**  
 1687        "block\_id": "xxxxx"  
 1688  
 1689  
 1690     👉 **Tool Response:**  
 1691        {"object":"list","results": [...]}  
 1692  
 1693  
 1694     📌 **Agent** I can see there are three databases in the HR Record page:  
 1695        1. Job Positions (block ID: xx-yy-zz)  
 1696        2. Job Application Timeline (block ID: yy-zz-xx)  
 1697        3. Candidates (block ID: zz-xx-yy)  
 1698  
 1699     Let me first examine the Job Positions database to understand what positions are available, and then the Candidates database to see its structure.  
 1700  
 1701  
 1702     📌 **notion API-retrieve-a-database**  
 1703        "database\_id": "xx-yy-zz"  
 1704  
 1705  
 1706     👉 **Tool Response:**  
 1707        {"object":"database","id":"xx-yy-zz",...}  
 1708  
 1709  
 1710     📌 **notion API-post-database-query**  
 1711        "database\_id": "xx-yy-zz"  
 1712  
 1713  
 1714  
 1715     👉 **Tool Response:**  
 1716        {"object":"list","results": [...]}  
 1717  
 1718  
 1719     📌 **Agent** Now let me examine the Candidates database structure:  
 1720  
 1721     📌 **notion API-retrieve-a-database**  
 1722        "database\_id": "zz-xx-yy"  
 1723  
 1724  
 1725  
 1726  
 1727

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

👉 **Tool Response:**

```
{"object": "database", "id": "zz-xx-yy", ...}
```

👉 **notion API-post-database-query**

```
"database_id": "zz-xx-yy"
```

👉 **Tool Response:**

```
{"object": "list", "results": [...]}
```

👉 **Agent** Now I have a clear understanding of the situation: Job Positions with available headcount (Head Counts > 0):

- Advanced Investment Banking Analyst: 3 positions
- Market Risk Specialist: 4 positions
- ...

Existing sample candidates and their applied positions:

- Alice Jackson → Financial Risk Analyst (2 positions available) ✓
- Debra Smith → Securities Trader (0 positions available) ×
- Michael Diaz → Investment Analyst (2 positions available) ✓
- ...

I need to:

1. Delete all existing sample entries
2. Read all the resume PDFs to extract candidate information
3. Add new entries based on actual resumes
4. Send rejection emails to candidates applying for positions with 0 headcount

Let me start by reading all the resume PDFs:

👉 **pdf-tools read\_pdf\_pages**

```
"pdf_file_path": "workspace/cvs/alice.jackson.resume.pdf", "start_page": 1, "end_page": 10
```

👉 **Tool Response:**

```
PDF total pages: 3
Extracted page range: 1-3

==== Page 1 ====
ALICE JACKSON
**Applied Position: Financial Risk Analyst **
CONTACT INFORMATION
...
```

Read the resumes of the other candidates in sequence .....

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

🤖 **Agent** Now I have all the resume information. Let me analyze what I found:  
**Candidates from resumes and their applied positions:**

- **Alice Jackson** – Applied: Financial Risk Analyst (2 positions available) ✓  
 Education: Master of Finance, Columbia Business School; Bachelor Economics, University of Pennsylvania  
 Email: alice\_jackson38@mcp.com
- **Angela Moore** – Applied: Securities Trader (0 positions available) ✗  
 Education: Master of Science in Computational Finance, Carnegie Mellon University  
 Email: angela\_moore89@mcp.com
- ...

Now I need to:

1. Delete all existing sample candidates
2. Add the real candidates from resumes
3. Send rejection emails to candidates applying for positions with 0 headcount

Let me start by deleting the existing sample candidates:

👉 **notion API-update-a-block**

```
"block_id": "xy-yz-zx", "archived": true
```

👉 **Tool Response:**

```
"status":400,"object":"error","code":"validation_error",...
```

🤖 **Agent** I need to use the correct API to delete pages from the database. Let me delete them one by one:

👉 **notion API-patch-page**

```
"page_id": "xy-yz-zx", "archived": true
```

👉 **Tool Response:**

```
{"object":"page","id":"xy-yz-zx",...}
```

Delete others one by one .....

🤖 **Agent** Great! I've deleted all the existing sample candidates. Now I need to create new pages for the actual candidates based on their resumes. Let me start creating them one by one:

👉 **notion API-post-page**

```
"parent": {"database_id": "zz-xx-yy"},  

"properties": {  

  "title": [{"text": {"content": "Alice Jackson"} }], "type": "title"  

}
```

👉 **Tool Response:**

```
{"status":400,"object":"error","code":"validation_error",...}
```

1836

1837  **Agent** I need to use the correct format for creating pages in a database. Let me check the database  
1838 structure again and create the pages properly:

1839

1840  **notion API-post-page**

1841

```
{
  "parent": {"database_id": "zz-xx-yy"},  

  "properties": {"Name": {"title": [{"text": {"content": "Alice Jackson"}]}]},  

  "Email": {"rich_text": [{"text": {"content": "alice.jackson38@mcp.com"}]}],  

  "Applied Position": {"rich_text": [{"text": {"content": "Financial Risk  
Analyst"}]}],  

  "School": {"rich_text": [{"text": {"content": "Columbia Business School"}]}],  

  "Highest Degree": {"select": {"name": "master"}}}
```

1842

1843

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

1856

1857

1858

1859

1860

1861

1862

1863

1864

1865

1866

1867

1868

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

 **Agent** I need to use the correct format for creating pages in a database. Let me check the database  
structure again and create the pages properly:

 **notion API-post-page**

```
{
  "parent": {"database_id": "zz-xx-yy"},  

  "properties": {"Name": {"title": [{"text": {"content": "Alice Jackson"}]}]},  

  "Email": {"rich_text": [{"text": {"content": "alice.jackson38@mcp.com"}]}],  

  "Applied Position": {"rich_text": [{"text": {"content": "Financial Risk  
Analyst"}]}],  

  "School": {"rich_text": [{"text": {"content": "Columbia Business School"}]}]},  

  "Highest Degree": {"select": {"name": "master"}}}
```

 **Tool Response:**

```
{"object": "page", "id": "yx-zy-xz", ...}
```

 **Agent** Great! Now let me add the remaining candidates:

Add the remaining candidates .....

 **Agent** Perfect! I've successfully added all 10 candidates from the resumes to the Notion database.  
Now I need to send rejection emails to candidates who applied for positions that are currently not  
open for recruitment (positions with 0 headcount).Based on my analysis, the candidates who need  
rejection emails are:

- Angela Moore (Securities Trader - 0 positions)
- Debra Smith (Securities Trader - 0 positions)
- Emily James (Securities Trader - 0 positions)
- Martha Morales (Portfolio Manager - 0 positions)

Let me send the rejection emails:

 **emails send\_email**

```
"to": "angela.moore89@mcp.com",  

"subject": "Application Status Update - Securities Trader Position",  

"body": ...
```

 **Tool Response:**

```
Email sent successfully to angela.moore89@mcp.com
```

Send other rejection emails .....

 **Agent** Excellent! I have successfully completed all the required tasks: ...

 **claim\_done**

```
None
```

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

↳ **Tool Response:**

you have claimed the task is done!