

---

# Eta Given Delta: Defining LLM Tool Efficiency With Marginal Tool Utility

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper introduces **tool efficiency**, a new quantitative metric to evaluate the rate  
2 of *useful* tool calls in an LLM agent trajectory. To ensure that tool efficiency is well-  
3 defined, we also introduce **marginal tool utility**, a new quantitative metric defined  
4 per tool call indicating whether a tool is useful or whether it can be safely removed  
5 from the tool suite without affecting accuracy while increasing tool efficiency;  
6 in this paper, we determine the sign of marginal tool utility for each tool call in  
7 a trajectory using LLM-as-a-Judge. While much prior work has been done to  
8 develop techniques that improve tool use by LLMs and design evaluation methods  
9 measuring efficiency indirectly using accuracy as a proxy, our work is centered on  
10 measuring efficiency *directly* via the quantitative metric proposed in this paper in  
11 post hoc trajectory analyses. It is our intention that this work contributes to the  
12 frontier of LLM evaluation research as a springboard for future benchmark designs  
13 and agent harness engineering (specifically with regards to creating lean tool suites)  
14 that optimize for metrics that complement but are distinct from accuracy.

## 15 1 Introduction

16 Large Language Models (LLMs) are increasingly used in real-world tasks such as software engineer-  
17 ing (Chen et al., 2021; Athiwaratkun et al., 2023; Austin et al., 2021), mathematical reasoning and  
18 autoformalization (Ahn et al., 2024; Guo et al., 2025; Manem et al., 2025), and scientific research  
19 (Zhuang et al., 2025; Chaturvedi et al., 2026). However, should LLMs remain simply as language  
20 modeling artifacts (Brown et al., 2020) without access to external up-to-date knowledge (Cheng  
21 et al., 2024) and no way to act or modify its environment (Yao et al., 2020; Huang et al., 2022; Ahn  
22 et al., 2022), they can hardly be called agents (Russell and Norvig, 2020). To address this, LLMs  
23 are equipped with tools (Schick et al., 2023), which have become an integral part of LLM agentic  
24 workflows that we see today (Yao et al., 2022).

25 While tools are generally useful, not all tools are equally useful as others (Huang et al., 2024; Yu  
26 et al., 2025), this finding acknowledged even by frontier industry labs like Anthropic (2025). To that  
27 end, much prior work focuses on developing techniques that improve tool use (Sivakumaran et al.,  
28 2026) and on creating benchmarks that such techniques can evaluate against and optimize for (Ning  
29 et al., 2024). A common ground shared by existing work is that the evaluation metric used is almost  
30 always outcome accuracy (Jin et al., 2025; Li et al., 2025b), even if accuracy is not the primary metric  
31 of interest, rather it being used as a proxy to measure the efficiency of tool use (Ning et al., 2024).

32 Given this, we notice a gap in the literature: there does not yet exist a direct way to measure tool  
33 efficiency, which we define as the rate of useful tool calls, given an LLM agent trajectory. While  
34 recent, even contemporaneous works like TRM by Ma et al. (2026) take the first steps towards defining  
35 reward functions centered on intermediate tool invocations for reinforcement learning to improve tool  
36 call quality, we take this idea one step further by enabling post hoc analyses on agent trajectories with

37 an aggregate tool efficiency metric, and keeping to solely inference-time methods easily-adoptable by  
38 application-layer engineering teams to assist them in designing leaner and easier-to-maintain agent  
39 tool suites.

40 In attempting to define tool efficiency, we realize that this necessarily includes defining what use-  
41 fulness means with respect to a tool call. We arrive at the idea of marginal tool utility, which is a  
42 metric relevant to each particular tool call instance in an agent trajectory that directly determines the  
43 usefulness of these tool calls, and whose aggregate per tool determines the usefulness of the tool itself.  
44 We also realize that this definition must be supported by empirical results, specifically those obtained  
45 from tool ablations and observing the change in accuracy (or lack thereof) depending on which tool  
46 was ablated. To that end, we used the APEX-SWE Observability benchmark (Kottamasu et al., 2026),  
47 a benchmark in which today’s frontier models still struggle to achieve satisfactory accuracy levels,  
48 and a benchmark that by default includes tools that we hypothesize (and later prove) are necessary  
49 (useful) as well as unnecessary (safely disposable).

50 As such, **our contributions** in this paper are four-fold:

- 51 1. Introduce marginal tool utility  $\Delta_\alpha$ , specifically for pass/fail tasks.
- 52 2. Introduce tool efficiency  $\eta_\alpha$ , where a useful tool call  $\alpha_i$  is one with  $\Delta_\alpha(\alpha_i) > 0$ .
- 53 3. Implement a simple and cost-efficient method for calculating  $\eta_\alpha$  for pass/fail tasks by using  
54 LLM-as-a-Judge to determine the sign of each  $\Delta_\alpha(\alpha_i)$ .
- 55 4. Empirically support the correctness of our definitions and LLM-as-a-Judge evaluator through  
56 tool ablation results with the APEX-SWE benchmark; critically, adding tools with generally  
57  $\Delta_\alpha \leq 0$  does not increase the accuracy of the agent, while adding tools with generally  
58  $\Delta_\alpha > 0$  does increase the accuracy of the agent, all as expected.

## 59 2 Related Works

60 **LLM evaluations.** There exists a plethora of benchmarks to assess an LLM agent’s accuracy in  
61 completing various tasks: SWE-bench (Jimenez et al., 2024), Terminal-Bench (Merrill et al., 2026),  
62 BrowseComp (Wei et al., 2025), OSWorld-Verified (Xie et al., 2024). Benchmarks specifically  
63 designed to evaluate an agent’s ability to use tools to correctly complete tasks include MCP-Atlas  
64 (Bandi et al., 2026), ToolBench (Qin et al., 2023), StableToolBench (Guo et al., 2024), and Toolathlon  
65 (Li et al., 2025a). The aforementioned benchmarks focus on evaluating the accuracy of an agent  
66 in completing the included tasks, this dependence on accuracy pervading even benchmarks like  
67 WTU-EVAL by Ning et al. (2024) that use accuracy as a proxy to determine whether tool usage is  
68 done correctly. In this paper, we seek to bypass accuracy altogether, introducing a new quantifiable  
69 metric in tool efficiency that we hope future benchmarks can use as a complementary metric to  
70 accuracy.

71 **LLM tool call optimizations.** The literature on optimizing LLM tool usage include reward models  
72 (Agarwal et al., 2026; Ma et al., 2026), self-play reinforcement learning (Acikgoz et al., 2026),  
73 and tool retrieval (Moura, 2025; Erdogan et al., 2024). Other interesting related work include ART  
74 (Paranjape et al., 2023) that approach tool call optimization not necessarily from the lens of improving  
75 tool usage, but using tool outputs to improve generation dynamically. Several works like Yu et al.  
76 (2025) concluded that some tasks benefit from tool usage while others do not, suggesting that tool  
77 usefulness is task-dependent and is not equal across all tools. Similar to Ma et al. (2026), we extend  
78 this idea and hypothesize that usefulness is not equal across all tool *calls*, i.e., particular instances of  
79 tool invocations dotted throughout an agent trajectory. Our work does not directly contribute to tool  
80 call optimizations, instead its contributions aim to enable future optimizations that directly index on  
81 tool efficiency.

82 **Methods improving agentic efficiency.** To the best of our knowledge, tool efficiency is not a well-  
83 defined/quantitative metric prior to this paper. However, this does not mean that there has been no  
84 prior work that is focused on improving efficiency in LLM agentic workflows. For instance, Recursive  
85 Language Models (Zhang et al., 2026) and FrugalGPT (Chen et al., 2024) focus on minimizing API  
86 cost, LLMingua (Jiang et al., 2023) on minimizing inference latency, and an instruction-refinement  
87 framework by Wu et al. (2025) on maximizing Cost-Aware Pass Rate (which is a metric introduced  
88 in the same paper). We consider the metrics these methods aim to optimize *efficiency-adjacent*,

89 reserving the term *efficiency* for a metric that calculates the rate of *useful* items as is accepted in  
 90 physics.

### 91 3 Methodology

92 To enable direct and explicit measurements of efficiency, we first introduce **marginal tool utility**,  
 93 particularly in the case of pass/fail tasks. We then introduce **tool efficiency**, where the usefulness of a  
 94 tool call is conditioned on the sign of marginal tool utility of said tool call.

#### 95 3.1 Problem Formulation

96 We formalize a multi-turn LLM agent trajectory as an ordered sequence of token sequences. An LLM  
 97  $\pi$  alternates between reasoning tokens and tool calls as implemented according to the ReAct paradigm  
 98 (Yao et al., 2022), which is the paradigm used for the agent harness in APEX-SWE Observability  
 99 (Kottamasu et al., 2026).

100 Formally, consider LLM  $\pi$  as well as system message  $s$ , and user message  $u$ . Given prompt  $p = (s, u)$ ,  
 101 LLM  $\pi$  executes multiple tool calls sequentially and outputs reasoning tokens in between tool calls to  
 102 synthesize the information it has access to at each timestep. This iterative process terminates once  
 103 LLM  $\pi$  reaches the maximum allowable step count or calls a terminal tool to produce final answer  $y$ .  
 104 Hence, we define agent trajectory  $\tau$  as

$$\tau \triangleq (p, t_1, r_1, o_1, \dots, t_N, r_N, o_N, t_{N+1}, y) \quad (1)$$

105 where each  $t_i (1 \leq i \leq N + 1)$  denotes a sequence of reasoning tokens, each  $r_i (1 \leq i \leq N)$  denotes  
 106 a tool request, each  $o_i (1 \leq i \leq N)$  denotes the tool output corresponding to  $r_i$ , and  $N \in \mathbb{N}$  denotes  
 107 the total number of steps taken by the agent that is no greater than the maximum step count set as a  
 108 hyperparameter. We thus define a tool call  $\alpha_i = (r_i, o_i)$ .

109 To validate our paper’s proposed definitions, we must define the expected correctness of final answer  
 110  $y$ . In APEX-SWE Observability, final answer  $y$  is the final code patch generated by the agent to solve  
 111 its automatic program repair task (de Souza et al., 2017). Correctness is determined by unit tests  
 112 following the FAIL\_TO\_PASS/PASS\_TO\_PASS methodology inspired by SWE-bench (Jimenez  
 113 et al., 2024). We define each unit test as function  $f_j: \mathbb{N}^n \rightarrow \{0, 1\}$ , where  $f_j(y) = 0$  when final  
 114 answer  $y$  causes unit test  $j$  to fail and  $f_j(y) = 1$  when final answer  $y$  causes unit test  $j$  to pass. The  
 115 expected correctness of final answer  $y$  over trajectories generated by LLM  $\pi$  is thus

$$\mathbb{E}_{\tau \sim \pi} [\mathbb{I}(f_1(y) \cdot \dots \cdot f_M(y) = 1)] \quad (2)$$

116 where  $M \in \mathbb{N}$  denotes the total number unit tests such that  $1 \leq j \leq M$ , and  $\mathbb{I}(\cdot)$  is the indicator  
 117 function.

#### 118 3.2 Defining Marginal Tool Utility

119 Given a trajectory  $\tau$  of an agent tasked to solve a pass/fail task, the marginal tool utility of the  $i^{\text{th}}$  tool  
 120 call  $\Delta_\alpha(\alpha_i)$  is defined as the difference between the likelihood that the task is solved correctly given  
 121 tool calls up to and *including*  $\alpha_i$  and the likelihood that the task is solved correctly given tool calls up  
 122 to and *excluding*  $\alpha_i$ . Among other methods, this can be determined using repeated policy  $\pi$  rollouts  
 123 or using LLM-as-a-Judge (Zheng et al., 2023) comparing the before and after trajectories. The latter  
 124 is more computationally efficient and practical than the former, hence we opt to use LLM-as-a-Judge  
 125 in this paper, specifically to determine  $\text{sgn}(\Delta_\alpha(\alpha_i))$ , where  $\text{sgn}(\cdot)$  is the signum function. We argue  
 126 that this is an acceptable method as a tool call  $\alpha_i$  is *useful* if and only if upon its execution the  
 127 likelihood of correctness increases ( $\Delta_\alpha(\alpha_i) > 0$ ). Hence, the additional information afforded by  
 128 repeated policy  $\pi$  rollouts, that being the exact likelihood difference before and after the tool call, is  
 129 unnecessary for our purposes.

130 Formally, consider judge LLM  $\pi^{(J)}$  as well as system message  $s^{(J)}$  and user message  $u^{(J)}(\alpha_i, \tau)$ ;  
 131 the specific  $s^{(J)}$  and  $u^{(J)}$  are outlined in Appendix A, and the particular model used is GPT-5.4 on  
 132 Microsoft Azure (OpenAI, 2026b). Given prompt  $p^{(J)} = (s^{(J)}, u^{(J)}(\alpha_i, \tau))$ , judge  $\pi^{(J)}$  generates

133 final answer  $y^{(J)}$ , which is a structured output that can be represented as a triplet  $(L, C, R)$ , where  
 134  $L = \{0, 1\}$  denoting two possible marginal tool utility labels (positive or non\_positive),  $C =$   
 135  $[0, 1]$  denoting a float confidence score, and  $R = \mathbb{N}^n$  denoting tokens that read as the judge’s rationale  
 136 for its classification.

137 The judge trajectory  $\tau^{(J)}$  is thus simply

$$\tau^{(J)} \triangleq (p^{(J)}, y^{(J)}) \quad (3)$$

138 Given judge  $\pi^{(J)}$ , the  $i^{\text{th}}$  tool call  $\alpha_i$ , the trajectory up to and excluding the  $i^{\text{th}}$  tool call  $\tau_{1,\dots,i-1}$  (i.e.,  
 139 before trajectory), and the trajectory up to and including the  $i^{\text{th}}$  tool call  $\tau_{1,\dots,i}$  (i.e., after trajectory),  
 140 we get  $y_i^{(J)}$ . In particular,

$$\begin{cases} L_i = 1 & \iff \Delta_\alpha(\alpha_i) > 0 & \text{tool call is useful} \\ L_i = 0 & \iff \Delta_\alpha(\alpha_i) \leq 0 & \text{otherwise} \end{cases} \quad (4)$$

141 Notice that marginal tool utility is defined per tool *call*. For a particular trajectory  $\tau$ , a *tool* is useful  
 142 if and only if it has more useful ( $\Delta_\alpha(\alpha_i) > 0$ ) calls than it has non-useful ( $\Delta_\alpha(\alpha_i) \leq 0$ ) calls.  
 143 Equivalently, for all tool calls  $\alpha_i$  of a particular tool,

$$\sum \Delta_\alpha(\alpha_i) \begin{cases} > 0 & \text{tool is useful} \\ \leq 0 & \text{otherwise} \end{cases} \quad (5)$$

144 We can also define this aggregate for a set of independent trajectories. This is only advisable in  
 145 practical settings if the independent trajectories share the same tool suite and each corresponds to a  
 146 task comparable to each other (e.g., each corresponds to an observability task using the same agent  
 147 harness within the same benchmark).

### 148 3.3 Defining Tool Efficiency

149 Having determined  $\text{sgn}(\Delta_\alpha(\alpha_i)) \forall i$ , we can define tool efficiency as the ratio of the number of useful  
 150 tool calls to the total number of tool calls in a trajectory. Formally, for trajectory  $\tau$  with  $N$  tool calls  
 151 such that  $1 \leq i \leq N$ ,

$$\eta_\alpha(\tau) \triangleq \frac{|\{\alpha_i | \Delta_\alpha(\alpha_i) > 0\}|}{N} \quad (6)$$

152 We can thus also define the mean tool efficiency of a set of  $L$  independent trajectories

$$\bar{\eta}_\alpha(\tau^{(1)}, \dots, \tau^{(L)}) = \frac{\eta_\alpha(\tau^{(1)}) + \dots + \eta_\alpha(\tau^{(L)})}{L} \quad (7)$$

153 We hypothesize, and later empirically show, that if a non-useful tool is removed from the tool suite,  
 154 we expect to see tool efficiency to increase from prior to the removal of said tool. This is not obviated  
 155 by the definition of tool efficiency, most significantly because the total number of tool calls  $N$  can  
 156 vary given a change in the tool suite.

### 157 3.4 Experimental Setup

158 We conduct our empirical experiments to show that the concept of marginal tool utility aligns with  
 159 what we expect to see in practical settings when performing tool ablations. Particularly, if a tool has  
 160 positive aggregate tool utility, the expected correctness of the agent’s final answer increases with  
 161 the inclusion of the tool in the agent’s tool suite; if a tool has non-positive aggregate tool utility, the  
 162 expected correctness of the agent’s final answer remains constant or decreases with the inclusion  
 163 of the tool in the agent’s tool suite. We show that these hypotheses are supported by the empirical  
 164 correctness results.

165 We evaluate frontier LLMs against the 25 public observability tasks of APEX-SWE (Kottamasu  
166 et al., 2026) available on HuggingFace. (Kottamasu et al., 2026) noted that this public subset is  
167 representative of the whole set of 100 observability tasks, of which 75 are private.

168 In APEX-SWE Observability, the default agent harness includes native agent tools (`bash`,  
169 `search_files`, `read_file`, `apply_patch`, `update_plan`) and three read-only MCP tools ex-  
170 posed via `bash`: Grafana/Loki for logs, Mattermost for discussions between human developers,  
171 and Plane for software issues/tickets/specifications. Being familiar with how observability tasks  
172 for production software are often solved, we hypothesize that Grafana/Loki is a useful MCP tool  
173 with positive aggregate tool utility, while Mattermost and Plane are non-useful MCP tools with  
174 non-positive aggregate tool utility.

175 To that end, we conducted tool ablations to obtain three variants of the agent harness:

- 176 • `default`: Full suite of MCP tools (Grafana/Loki, Mattermost, Plane).
- 177 • `grafana`: Only Grafana/Loki MCP is available.
- 178 • `no-mcp`: Absolutely no MCP tools are available.

179 Through these ablations, we are able to verify whether (1) the inclusion/omission of Mattermost and  
180 Plane affect accuracy, and (2) the inclusion/omission of Grafana/Loki affect accuracy. In Section  
181 4, we present the task accuracy results we obtained and cross-reference them with the aggregate  
182 tool utility results of the three MCP tools. In our experiments, we only consider MCP tools as tools  
183 for simplicity of analysis, and because the native agent tools are either necessary for the task to be  
184 completed (e.g., `apply_patch`) or are necessary for the MCP tools to be executed (e.g., `bash`).

185 **Agent specifications.** We used GPT-5.3-Codex on Microsoft Azure (OpenAI, 2026a) and Gemini  
186 3.1 Pro on Google Agent Platform (formerly Vertex AI) (Google, 2026) to create two independent  
187 agent instances. Therefore, we obtained a total of  $25 \cdot 3 \cdot 2 = 150$  agent trajectories.

188 **Compute specifications.** The proprietary LLMs were accessed via their respective provider APIs  
189 and thus run on their respective providers’ compute clusters. The experiments were otherwise run on  
190 AWS Batch using EC2 CPU workers, each as an isolated containerized job with Docker-in-Docker  
191 enabled. Jobs used 16 vCPUs, 110 GB RAM, and a 600 GB gp3 root volume. The AWS Batch  
192 compute environment used the `r7i.4xlarge`, `r7i.8xlarge`, and `m7i.8xlarge` instance types.

## 193 4 Results

194 In this section, we report the results we obtained from all 150 agent trajectories, namely task accuracy,  
195 marginal and aggregate tool utility, and mean tool efficiency. We show that task accuracy increases  
196 with the inclusion of a tool with positive aggregate tool utility, and tool efficiency increases with the  
197 removal of non-useful tools.

### 198 4.1 Task Accuracy

199 We find that task accuracy remains approximately constant when both Mattermost and Plane are  
200 removed (`default` vs `grafana`) and that task accuracy decreases when Grafana/Loki is removed  
201 (`grafana` vs `no-mcp`), as seen in Tables 1 and 2.

Table 1: Accuracy across 25 trajectories per variant for GPT-5.3-Codex.

Variant	Pass	Fail	Accuracy
<code>default</code>	8	17	0.32
<code>grafana</code>	9	16	0.36
<code>no-mcp</code>	6	19	0.24

Table 2: Accuracy across 25 trajectories per variant for Gemini 3.1 Pro.

Variant	Pass	Fail	Accuracy
default	9	16	0.36
grafana	9	16	0.36
no-mcp	5	20	0.20

202 A breakdown of which specific tasks passed/failed can be found in Appendix B. Comparing `default`  
 203 and `grafana`, when a task passes for one it also generally passes for the other, ditto for fails.

204 From these results, note that we can qualitatively claim that Grafana/Loki is a useful tool, while  
 205 Mattermost and Plane are not. This claim aligns with the nature of the tasks: an agent completing an  
 206 observability task benefits more from reading logs (primary source of information) than from reading  
 207 discussions/specifications (secondary sources of information) which may actually mislead the agent  
 208 towards investigating parts of the codebase that are actually not broken due to directionally-wrong  
 209 bugfix discussions between human developers and/or outdated specifications.

## 210 4.2 Marginal Tool Utility

211 Using the judge LLM, we first determine the sign of marginal tool utility of each tool call in each  
 212 trajectory. We also analyze the rationales outputted by the judge LLM for its classifications. Then, we  
 213 calculate the aggregate tool utility of each tool across all trajectories separated by model and variant.

214 In Section 4.1, we see that Grafana/Loki affects accuracy when added/removed from the tool suite,  
 215 and that neither Mattermost nor Plane do. We thus expect the following aggregate tool utilities:

$$\begin{cases} \sum_{\text{grafana}} \Delta_{\alpha}(\alpha_i) > 0 \\ \sum_{\text{mattermost}} \Delta_{\alpha}(\alpha_i) \leq 0 \\ \sum_{\text{plane}} \Delta_{\alpha}(\alpha_i) \leq 0 \end{cases} \quad (8)$$

216 Investigating the 50 trajectories belonging to the `default` variant, the marginal tool utility results  
 217 in Tables 4 and 3 are consistent with that expectation, supporting our proposed definition and  
 218 implementation. For `default` with GPT-5.3-Codex, the aggregate tool utility of Grafana/Loki is  
 219 +25, Mattermost is -35, and Plane is -30. For `default` with Gemini 3.1 Pro, the aggregate tool  
 220 utility of Grafana/Loki is +5, Mattermost is -17, and Plane is -28. For `grafana`, the aggregate tool  
 221 utility of Grafana/Loki is  $58 - 29 = +29$  with GPT-5.3-Codex and  $23 - 20 = +3$  with Gemini 3.1  
 222 Pro.

223 In determining the signs of marginal tool utilities, the judge was prompted to output its rationale. The  
 224 most common rationales to justify  $\Delta_{\alpha} > 0$  classifications are as follows: Grafana/Loki logs exposed  
 225 concrete failures (exact error logs that hint at the likely fix); more targeted queries for Grafana/Loki  
 226 logs at later steps narrowed the investigation and provided more focused signals; and Mattermost  
 227 and Plane context helped orient the agent towards useful conversation and/or discussion context.  
 228 To justify  $\Delta_{\alpha} \leq 0$  classifications: generic, irrelevant, or noisy context from logs, discussions, and  
 229 specifications distracted the agent; malformed queries by the agent wasted steps on retries; and valid  
 230 queries but those resulting in empty outputs wasted steps on pivoting to valid queries that actually  
 231 gave non-empty outputs.

Table 3: Marginal tool utility values of each tool in the `default` harness for GPT-5.3-Codex.

Tool	$\Delta_{\alpha} > 0$	$\Delta_{\alpha} \leq 0$	Mean Confidence ( $\Delta_{\alpha} > 0$ )	Mean Confidence ( $\Delta_{\alpha} \leq 0$ )
Grafana/Loki	52	27	0.819	0.890
Mattermost	7	42	0.734	0.927
Plane	23	53	0.775	0.928

Table 4: Marginal tool utility values of each tool in the default harness for Gemini 3.1 Pro.

Tool	$\Delta_\alpha > 0$	$\Delta_\alpha \leq 0$	Mean Confidence ( $\Delta_\alpha > 0$ )	Mean Confidence ( $\Delta_\alpha \leq 0$ )
Grafana/Loki	26	21	0.825	0.870
Mattermost	3	20	0.673	0.930
Plane	3	31	0.607	0.917

232 Analyzing agent trajectories, we notice that tool calls with  $\Delta_\alpha > 0$  occur more often near the  
 233 beginning of trajectories, followed by near the end of trajectories; tool calls with  $\Delta_\alpha \leq 0$  are found  
 234 more often in the middle of trajectories. Most early calls are high-signal Grafana/Loki calls (so are  
 235 most early positive marginal tool utility calls); middle calls are mostly spent querying noisy sources  
 236 and performing broad conversation/specification discovery; later calls generally surface more specific  
 237 evidence to support the agent’s code patch generation.

### 238 4.3 Tool Efficiency

239 Given the signs of marginal tool utility determined by the judge LLM for each tool call, we calculate  
 240 tool efficiency for each trajectory in default and grafana. Then, we calculate the mean tool  
 241 efficiency separated by model and variant. Since the no-mcp variant has no MCP tools, tool efficiency  
 242 is undefined for that particular variant. Our results are found in Tables 5 and 6.

Table 5: Mean tool efficiency for GPT-5.3-Codex out of 25 agent trajectories per variant. Included are the number of useful tool calls and total number of tool calls from all trajectories combined; the quotient of these two values is not equivalent to mean tool efficiency as agent trajectories vary in length.

Variant	Useful Call Count (All)	Total Call Count (All)	Mean Tool Efficiency
default	82	204	0.359
grafana	58	87	0.720

Table 6: Mean tool efficiency for Gemini 3.1 Pro out of 25 agent trajectories per variant. Included are the number of useful tool calls and total number of tool calls from all trajectories combined; the quotient of these two values is not equivalent to mean tool efficiency as agent trajectories vary in length.

Variant	Useful Call Count (All)	Total Call Count (All)	Mean Tool Efficiency
default	32	104	0.367
grafana	23	43	0.593

243 As we can see, tool efficiency increases when Mattermost and Plane are removed. This finding is also  
 244 consistent with the conclusion that these two MCP tools are unnecessarily included in the agent’s  
 245 tool suite and more often than not act as distractors to the agent in attempting to solve its given task.

## 246 5 Discussion

247 Integrating all our results together, we find several promising directions for future work the community  
 248 can engage in that leverage marginal tool utility and tool efficiency.

249 **Minimizing sub-optimal middle calls.** The trend that middle calls are often sub-optimal (see Section  
 250 4.2) is a useful finding that may inspire future LLM reinforcement learning techniques, similar to  
 251 those by Ma et al. (2026); Shao et al. (2024); Ouyang et al. (2022); Uesato et al. (2022), integrating  
 252 marginal tool utility as a reward signal. Alternatively, marginal tool utility can be used for inference-time  
 253 optimizations, like designing agents that can backtrack (Yang et al., 2025) given  $N$  consecutive tool  
 254 calls with non-positive marginal tool utility.

255 **Orthogonal tool suites.** As seen through our tool ablations of APEX-SWE Observability, tool  
256 suites – even those in publicly available benchmarks – can be more bloated than necessary, and other  
257 works like those by Liu et al. (2025) agree, albeit from a different perspective. The tool efficiency  
258 metric can be directly used to assess the overall health of a tool suite such that it remains as lean and  
259 maintainable as possible, and the aggregate tool utility of each tool to decide whether to preserve the  
260 tool in improved versions of the harness. This enables what we call *orthogonal* tool suites, tool suites  
261 that contain tools that serve non-redundant and/or mutually exclusive functions, which may not be  
262 obviously orthogonal given differences in their providers and/or descriptions (e.g., Mattermost vs  
263 Plane vs Grafana/Loki), yet actually are due to the objective of the task the agent must solve.

264 **More informed harness engineering.** We observe in Section 4 that the same task accuracy can be  
265 reached by agents using different backbone models despite having different tool efficiencies. We  
266 posit that this is primarily a product of different model training methodologies as the agent harness  
267 remains a constant variable, but an important finding to note nonetheless as this supports the view that  
268 harness engineering must be informed by each agent’s particular backbone model. In other words,  
269 a one-size-fits-all approach to harness engineering may be sub-optimal, hence workflows that are  
270 designed to automate harness engineering Lee et al. (2026) can benefit from our proposed metrics.

271 **Self-improving agents indexing on tool efficiency.** Taking automated harness engineering one  
272 step further, the agents themselves can recursively self-improve (Xia et al., 2025; Xiong et al.,  
273 2026) without relying on another agent that optimizes the harness offline. Not only can this online  
274 improvement be based on the usual metrics like task accuracy, tool efficiency and marginal tool utility  
275 can be leveraged mid-execution to provide rich reward signals to optimize context (including tool  
276 descriptions, input schemas, and output data shapes) or update model weights.

## 277 6 Limitations

278 While confidence scores outputted by an LLM as decoded tokens are not the most rigorous (Xiong  
279 et al., 2024; Yang et al., 2026), the judge’s confidence score per classification in Tables 4 and 3 indicate  
280 that the judge is generally more confident on non-positive classifications. We did not investigate why  
281 this is so, nor did we attempt to balance the confidence scores, as we view this work as outside the  
282 scope of this paper. Alternatively, a different LLM-as-a-Judge implementation could have been used,  
283 such as replacing the language modeling head of the judge with a binary classification head (Ma  
284 et al., 2026), though this requires using an open-weight LLM.

285 Most notably, we only ran tool ablations for read-only tools (the three MCPs are all read-only). We  
286 cannot remove a tool with write privileges (e.g., `apply_patch`) from the agent’s tool suite as that  
287 would have prevented the LLM from completing the task at all. Having said that, we seek to find  
288 ways in which we can assess the marginal tool utility of write tool calls, either directly or indirectly  
289 (the latter case by investigating how read-only tool calls are affected by what was written using the  
290 write tool calls), in future experiments.

## 291 7 Conclusion

292 We introduce tool efficiency and marginal tool utility, both being new quantitative metrics that are  
293 relevant in determining the usefulness of tools and thus in constructing the leanest possible tool suite  
294 for an LLM agent given a particular task. Our findings suggest that intuitive judgments regarding a  
295 tool’s usefulness as well as empirical tool ablation results using accuracy as a proxy in discerning  
296 usefulness/necessity of different tools agree with conclusions drawn from analyzing marginal tool  
297 utility and tool efficiency values.

298 All in all, we look forward to what our new definitions enable in the development of LLMs and LLM  
299 agents. We hope that the literature matures not only in terms of new methods, architectures, and  
300 workflows, but also in new evaluation dimensions against which we can measure the efficacy of these  
301 new artifacts. As well-defined quantitative metrics are more tractable to optimize for, we believe  
302 that evaluation/metrics research will accelerate progress in LLM research and engineering further,  
303 providing necessary direction and focus when iterating on novel methods.

304 **Societal Impacts.** As large language models and agents built on them become more capable, many  
305 have raised concerns regarding our over-reliance, even dependence, on them. While this work does

306 not introduce a new artifact directly relevant to such concerns, this work does enable new artifacts to  
307 be more capable than their existing versions today. We want to acknowledge that, and it is our hope  
308 and belief that agents and humans are able to work in tandem, each tackling a dimension of work the  
309 other is less-suited for. The onus is on us, those at the frontier of the development of this technology,  
310 to understand and thus educate the public on what those dimensions are.

## 311 **References**

- 312 Acikgoz, E. C., Qian, C., Hübotter, J., Ji, H., Hakkani-Tür, D., and Tur, G. (2026). Tool-r0: Self-  
313 evolving llm agents for tool-learning from zero data.
- 314 Agarwal, M., Abdelaziz, I., Basu, K., Unuvar, M., Lastras, L. A., Rizk, Y., and Kapanipathi, P. (2026).  
315 Toolrm: Outcome reward models for tool-calling large language models.
- 316 Ahn, J., Verma, R., Lou, R., Liu, D., Zhang, R., and Yin, W. (2024). Large language models for  
317 mathematical reasoning: Progresses and challenges.
- 318 Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan,  
319 K., Hausman, K., Herzog, A., Ho, D., Hsu, J., Ibarz, J., Ichter, B., Irpan, A., Jang, E., Ruano, R. J.,  
320 Jeffrey, K., Jesmonth, S., Joshi, N. J., Julian, R., Kalashnikov, D., Kuang, Y., Lee, K.-H., Levine,  
321 S., Lu, Y., Luu, L., Parada, C., Pastor, P., Quiambao, J., Rao, K., Rettinghouse, J., Reyes, D.,  
322 Sermanet, P., Sievers, N., Tan, C., Toshev, A., Vanhoucke, V., Xia, F., Xiao, T., Xu, P., Xu, S., Yan,  
323 M., and Zeng, A. (2022). Do as i can, not as i say: Grounding language in robotic affordances.
- 324 Anthropic (2025). Introducing advanced tool use on the claude developer platform.
- 325 Athiwaratkun, B., Gouda, S. K., Wang, Z., Li, X., Tian, Y., Tan, M., Ahmad, W. U., Wang, S., Sun, Q.,  
326 Shang, M., Gonugondla, S. K., Ding, H., Kumar, V., Fulton, N., Farahani, A., Jain, S., Giaquinto,  
327 R., Qian, H., Ramanathan, M. K., Nallapati, R., Ray, B., Bhatia, P., Sengupta, S., Roth, D., and  
328 Xiang, B. (2023). Multi-lingual evaluation of code generation models.
- 329 Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M.,  
330 Le, Q., and Sutton, C. (2021). Program synthesis with large language models.
- 331 Bandi, C., Hertzberg, B., Boo, G., Polakam, T., Da, J., Hassaan, S., Sharma, M., Park, A., Hernandez,  
332 E., Rambado, D., Salazar, I., Cruz, R., Rane, C., Levin, B., Kenstler, B., and Liu, B. (2026).  
333 Mcp-atlas: A large-scale benchmark for tool-use competency with real mcp servers.
- 334 Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam,  
335 P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R.,  
336 Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray,  
337 S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D.  
338 (2020). Language models are few-shot learners.
- 339 Chaturvedi, S. S., Bergerson, J., and Mallick, T. (2026). Toward reliable, safe, and secure llms for  
340 scientific applications.
- 341 Chen, L., Zaharia, M., and Zou, J. (2024). FrugalGPT: How to use large language models while  
342 reducing cost and improving performance. *Transactions on Machine Learning Research*. Featured  
343 Certification.
- 344 Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J., Edwards, H., Burda,  
345 Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G.,  
346 Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter,  
347 C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss,  
348 A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S.,  
349 Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford,  
350 A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D.,  
351 McCandlish, S., Sutskever, I., and Zaremba, W. (2021). Evaluating large language models trained  
352 on code.
- 353 Cheng, J., Marone, M., Weller, O., Lawrie, D., Khashabi, D., and Durme, B. V. (2024). Dated data:  
354 Tracing knowledge cutoffs in large language models. In *First Conference on Language Modeling*.

- 355 de Souza, H. A., Chaim, M. L., and Kon, F. (2017). Spectrum-based software fault localization: A  
356 survey of techniques, advances, and challenges.
- 357 Erdogan, L. E., Lee, N., Jha, S., Kim, S., Tabrizi, R., Moon, S., Hooper, C., Anumanchipalli, G.,  
358 Keutzer, K., and Gholami, A. (2024). Tinyagent: Function calling at the edge.
- 359 Google (2026). Gemini 3.1 pro: A smarter model for your most complex tasks.
- 360 Guo, C., Patel, M., Hartmann, B., Zamfirescu-Pereira, J., Chasins, S., and Ranade, G. (2025).  
361 Unspoken logic: Understanding and bridging the gap between free-form and LLM-interpretable  
362 natural language mathematical proofs. In *The 5th Workshop on Mathematical Reasoning and AI at  
363 NeurIPS 2025*.
- 364 Guo, Z., Cheng, S., Wang, H., Liang, S., Qin, Y., Li, P., Liu, Z., Sun, M., and Liu, Y. (2024).  
365 Stabletoolbench: Towards stable large-scale benchmarking on tool learning of large language  
366 models.
- 367 Huang, W., Abbeel, P., Pathak, D., and Mordatch, I. (2022). Language models as zero-shot planners:  
368 Extracting actionable knowledge for embodied agents.
- 369 Huang, Y., Shi, J., Li, Y., Fan, C., Wu, S., Zhang, Q., Liu, Y., Zhou, P., Wan, Y., Gong, N. Z., and  
370 Sun, L. (2024). Metatool benchmark for large language models: Deciding whether to use tools  
371 and which to use.
- 372 Jiang, H., Wu, Q., Lin, C.-Y., Yang, Y., and Qiu, L. (2023). LLMingua: Compressing prompts for  
373 accelerated inference of large language models. In Bouamor, H., Pino, J., and Bali, K., editors,  
374 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages  
375 13358–13376, Singapore. Association for Computational Linguistics.
- 376 Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. (2024). Swe-  
377 bench: Can language models resolve real-world github issues?
- 378 Jin, B., Zeng, H., Yue, Z., Yoon, J., Arik, S. O., Wang, D., Zamani, H., and Han, J. (2025). Search-r1:  
379 Training LLMs to reason and leverage search engines with reinforcement learning. In *Second  
380 Conference on Language Modeling*.
- 381 Kottamasu, A., Mahapatra, C., Lee, S., Pan, B., Barthwal, A., Datta, A., Gupta, A., Mehta, P., Arun,  
382 A., Alberti, S., Hiremath, A., Foody, B., and Vidgen, B. (2026). Apex-swe.
- 383 Lee, Y., Nair, R., Zhang, Q., Lee, K., Khattab, O., and Finn, C. (2026). Meta-harness: End-to-end  
384 optimization of model harnesses.
- 385 Li, J., Zhao, W., Zhao, J., Zeng, W., Wu, H., Wang, X., Ge, R., Cao, Y., Huang, Y., Liu, W., Liu,  
386 J., Su, Z., Guo, Y., Zhou, F., Zhang, L., Michelini, J., Wang, X., Yue, X., Zhou, S., Neubig, G.,  
387 and He, J. (2025a). The tool decathlon: Benchmarking language agents for diverse, realistic, and  
388 long-horizon task execution.
- 389 Li, X., Zou, H., and Liu, P. (2025b). Torl: Scaling tool-integrated rl.
- 390 Liu, M. M., Garcia, D., Parllaku, F., Upadhyay, V., Shah, S. F. A., and Roth, D. (2025). Toolscope:  
391 Enhancing llm agent tool use through tool merging and context-aware filtering.
- 392 Ma, D., Yang, Z., Xu, H., Fang, H., Yu, K., and Chen, L. (2026). Empowering LLM tool invo-  
393 cation with tool-call reward model. In *The Fourteenth International Conference on Learning  
394 Representations*.
- 395 Manem, C., Brahma, P. P., Mishra, P., Liu, Z., and Barsoum, E. (2025). SAND-math: Using LLMs to  
396 generate novel, difficult and useful mathematics questions and answers. In *The 5th Workshop on  
397 Mathematical Reasoning and AI at NeurIPS 2025*.
- 398 Merrill, M. A., Shaw, A. G., Carlini, N., Li, B., Raj, H., Bercovich, I., Shi, L., Shin, J. Y., Walshe,  
399 T., Buchanan, E. K., Shen, J., Ye, G., Lin, H., Poulos, J., Wang, M., Nezhurina, M., Jitsev, J.,  
400 Lu, D., Mastromichalakis, O. M., Xu, Z., Chen, Z., Liu, Y., Zhang, R., Chen, L. L., Kashyap,  
401 A., Uslu, J.-L., Li, J., Wu, J., Yan, M., Bian, S., Sharma, V., Sun, K., Dillmann, S., Anand, A.,

402 Lanpouthakoun, A., Koopah, B., Hu, C., Guha, E., Dreiman, G. H. S., Zhu, J., Krauth, K., Zhong,  
403 L., Muennighoff, N., Amanfu, R., Tan, S., Pimpalgaonkar, S., Aggarwal, T., Lin, X., Lan, X., Zhao,  
404 X., Liang, Y., Wang, Y., Wang, Z., Zhou, C., Heineman, D., Liu, H., Trivedi, H., Yang, J., Lin, J.,  
405 Shetty, M., Yang, M., Omi, N., Raoof, N., Li, S., Zhuo, T. Y., Lin, W., Dai, Y., Wang, Y., Chai, W.,  
406 Zhou, S., Wahdany, D., She, Z., Hu, J., Dong, Z., Zhu, Y., Cui, S., Saiyed, A., Kolbeinsson, A.,  
407 Hu, J., Rytting, C. M., Marten, R., Wang, Y., Dimakis, A., Konwinski, A., and Schmidt, L. (2026).  
408 Terminal-bench: Benchmarking agents on hard, realistic tasks in command line interfaces.

409 Moura, A. (2025). Antl3x/toolrag: Unlimited llm tools, zero context penalties - toolrag serves exactly  
410 the llm tools your user-query demands.

411 Ning, K., Su, Y., Lv, X., Zhang, Y., Liu, J., Liu, K., and Xu, J. (2024). Wtu-eval: A whether-or-not  
412 tool usage evaluation benchmark for large language models.

413 OpenAI (2026a). Introducing gpt-5.3-codex | openai.

414 OpenAI (2026b). Introducing gpt-5.4 | openai.

415 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal,  
416 S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A.,  
417 Welinder, P., Christiano, P., Leike, J., and Lowe, R. (2022). Training language models to follow  
418 instructions with human feedback.

419 Paranjape, B., Lundberg, S., Singh, S., Hajishirzi, H., Zettlemoyer, L., and Ribeiro, M. T. (2023).  
420 Art: Automatic multi-step reasoning and tool-use for large language models.

421 Qin, Y., Liang, S., Ye, Y., Zhu, K., Yan, L., Lu, Y., Lin, Y., Cong, X., Tang, X., Qian, B., Zhao, S.,  
422 Tian, R., Xie, R., Zhou, J., Gerstein, M., Li, D., Liu, Z., and Sun, M. (2023). Toolllm: Facilitating  
423 large language models to master 16000+ real-world apis.

424 Russell, S. J. and Norvig, P. (2020). *Artificial Intelligence: A Modern Approach (4th Edition)*.  
425 Pearson.

426 Schick, T., Dwivedi-Yu, J., Dessi, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., Cancedda, N., and  
427 Scialom, T. (2023). Toolformer: Language models can teach themselves to use tools.

428 Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and  
429 Guo, D. (2024). Deepseekmath: Pushing the limits of mathematical reasoning in open language  
430 models.

431 Sivakumaran, N., Chen, J., Wan, D., Zhang, Y., Yoon, J., Stengel-Eskin, E., and Bansal, M. (2026).  
432 DART: Leveraging multi-agent disagreement for tool recruitment in multimodal reasoning. In  
433 Demberg, V., Inui, K., and Marquez, L., editors, *Proceedings of the 19th Conference of the  
434 European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
435 pages 5445–5464, Rabat, Morocco. Association for Computational Linguistics.

436 Uesato, J., Kushman, N., Kumar, R., Song, F., Siegel, N., Wang, L., Creswell, A., Irving, G., and  
437 Higgins, I. (2022). Solving math word problems with process- and outcome-based feedback.

438 Wei, J., Sun, Z., Papay, S., McKinney, S., Han, J., Fulford, I., Chung, H. W., Passos, A. T., Fedus, W.,  
439 and Glaese, A. (2025). Browsecomp: A simple yet challenging benchmark for browsing agents.

440 Wu, B., Meij, E., and Yilmaz, E. (2025). A joint optimization framework for enhancing efficiency of  
441 tool utilization in LLM agents. In Che, W., Nabende, J., Shutova, E., and Pilehvar, M. T., editors,  
442 *Findings of the Association for Computational Linguistics: ACL 2025*, pages 22361–22373, Vienna,  
443 Austria. Association for Computational Linguistics.

444 Xia, P., Zeng, K., Liu, J., Qin, C., Wu, F., Zhou, Y., Xiong, C., and Yao, H. (2025). Agent0:  
445 Unleashing self-evolving agents from zero data via tool-integrated reasoning.

446 Xie, T., Zhang, D., Chen, J., Li, X., Zhao, S., Cao, R., Hua, T. J., Cheng, Z., Shin, D., Lei, F., Liu, Y.,  
447 Xu, Y., Zhou, S., Savarese, S., Xiong, C., Zhong, V., and Yu, T. (2024). Osworld: Benchmarking  
448 multimodal agents for open-ended tasks in real computer environments.

449 Xiong, M., Hu, Z., Lu, X., Li, Y., Fu, J., He, J., and Hooi, B. (2024). Can llms express their  
450 uncertainty? an empirical evaluation of confidence elicitation in llms.

451 Xiong, Y., Hu, S., and Clune, J. (2026). Learning to continually learn via meta-learning agentic  
452 memory designs.

453 Yang, S.-H., Wu, C.-K., Lin, C.-Y., Chen, Y.-N., yi Lee, H., and Sun, S.-H. (2026). On calibration of  
454 large language models: From response to capability.

455 Yang, X.-W., Zhu, X.-Y., Wei, W.-D., Zhang, D.-C., Shao, J.-J., Zhou, Z., Guo, L.-Z., and Li, Y.-F.  
456 (2025). Step back to leap forward: Self-backtracking for boosting reasoning of language models.

457 Yao, S., Rao, R., Hausknecht, M., and Narasimhan, K. (2020). Keep calm and explore: Language  
458 models for action generation in text-based games.

459 Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. (2022). React: Synergizing  
460 reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.

461 Yu, B., Baker, F. N., Chen, Z., Herb, G., Gou, B., Adu-Ampratwum, D., Ning, X., and Sun, H. (2025).  
462 Chemtoolagent: The impact of tools on language agents for chemistry problem solving.

463 Zhang, A. L., Kraska, T., and Khattab, O. (2026). Recursive language models.

464 Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E.,  
465 Zhang, H., Gonzalez, J. E., and Stoica, I. (2023). Judging LLM-as-a-judge with MT-bench and  
466 chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets  
467 and Benchmarks Track*.

468 Zhuang, Z., Chen, J., Xu, H., Jiang, Y., and Lin, J. (2025). Large language models for automated  
469 scholarly paper review: A survey. *Information Fusion*, 124:103332.

## 470 A Prompt for LLM-As-A-Judge for Marginal Tool Utility Sign

471  
472

### System Prompt

You are an expert evaluator of agent trajectories completing a SWE bugfix observability task.

Your task is to determine whether a tool call increases the likelihood the bugfix task is solved correctly. This is measured by a metric called marginal tool utility (MTU).

Concretely, you are tasked to determine the sign of marginal tool utility (MTU) for one tool call.

Definition:

$$MTU(i) = p(\text{task solved correctly} \mid \text{tool calls } 0..i) - p(\text{task solved correctly} \mid \text{tool calls } 0..i-1)$$

Equivalently:

MTU(i) is positive if the likelihood of the bugfix task being solved correctly strictly increases after tool call i is added into the trajectory. MTU(i) is non-positive otherwise.

You are given:

- BEFORE context: state of trajectory before tool call i executes.
- AFTER context: state of trajectory after tool call i completes.

473

- The specific tool call and tool result.

Output strict JSON only:

```
{
  "label": "positive" | "non_positive",
  "confidence": float between 0 and 1,
  "rationale": "brief explanation"
}
```

Use "positive" only when the call clearly increases solve likelihood.

Use "non\_positive" when likelihood is unchanged or reduced.

474  
475  
476

### User Prompt

Decide MTU sign for this tool call.

Target tool call:

- tool\_call\_id: {TOOL\_CALL\_ID}
- tool\_name: {TOOL\_NAME}
- arguments: {ARGUMENTS\_TRUNCATED}
- tool\_result\_excerpt: {TOOL\_RESULT}

BEFORE context (through tool calls 0..i-1 completion):

=== BEFORE START ===

{BEFORE\_CONTEXT}

=== BEFORE END ===

AFTER context (through tool calls 0..i completion):

=== AFTER START ===

{AFTER\_CONTEXT}

=== AFTER END ===

Return strict JSON only.

477

## 478 B APEX-SWE Observability Per-Task Correctness Breakdown

479 A score of 1.0 indicates a pass, while a score of 0.0 indicates a fail.

### 480 1. GPT-5.3-Codex

#### 481 (a) default

- 482 i. 0xpolygon-bor-1710-observability: 0.0
- 483 ii. 0xpolygon-bor-1728-1748-observability: 0.0
- 484 iii. 0xpolygon-bor-1743-observability: 1.0
- 485 iv. chainsafe-gossamer-4286-4720-observability: 1.0
- 486 v. chainsafe-gossamer-4489-4640-observability: 0.0
- 487 vi. chainsafe-gossamer-4573-4683-observability: 0.0
- 488 vii. containers-podman-compose-1221-1231-observability: 1.0
- 489 viii. containers-podman-compose-1226-1243-observability: 0.0

- 490 ix. containers-podman-compose-2-1238-observability: 1.0
- 491 x. containers-podman-compose-23-1214-observability: 0.0
- 492 xi. ethereum-optimism-op-geth-553-558-observability: 0.0
- 493 xii. ethereum-optimism-op-geth-655-observability: 0.0
- 494 xiii. ethereum-optimism-op-geth-666-observability: 0.0
- 495 xiv. ethereum-optimism-op-geth-675-observability: 0.0
- 496 xv. ethereum-optimism-op-geth-679-680-observability: 0.0
- 497 xvi. git-bug-git-bug-132-449-observability: 0.0
- 498 xvii. git-bug-git-bug-1367-1370-observability: 0.0
- 499 xviii. git-bug-git-bug-264-274-observability: 1.0
- 500 xix. git-bug-git-bug-338-341-observability: 0.0
- 501 xx. paperless-ngx-paperless-ngx-10195-10196-observability: 1.0
- 502 XXI. paperless-ngx-paperless-ngx-10555-observability: 0.0
- 503 xxii. paperless-ngx-paperless-ngx-5228-10559-observability: 0.0
- 504 xxiii. paperless-ngx-paperless-ngx-6341-9305-observability: 1.0
- 505 xxiv. paperless-ngx-paperless-ngx-8910-8912-observability: 1.0
- 506 xxv. paperless-ngx-paperless-ngx-9784-observability: 0.0

507 (b) grafana

- 508 i. 0xpolygon-bor-1710-observability: 0.0
- 509 ii. 0xpolygon-bor-1728-1748-observability: 0.0
- 510 iii. 0xpolygon-bor-1743-observability: 1.0
- 511 iv. chainsafe-gossamer-4286-4720-observability: 1.0
- 512 v. chainsafe-gossamer-4489-4640-observability: 0.0
- 513 vi. chainsafe-gossamer-4573-4683-observability: 0.0
- 514 vii. containers-podman-compose-1221-1231-observability: 1.0
- 515 viii. containers-podman-compose-1226-1243-observability: 0.0
- 516 ix. containers-podman-compose-2-1238-observability: 1.0
- 517 x. containers-podman-compose-23-1214-observability: 0.0
- 518 xi. ethereum-optimism-op-geth-553-558-observability: 0.0
- 519 xii. ethereum-optimism-op-geth-655-observability: 0.0
- 520 xiii. ethereum-optimism-op-geth-666-observability: 0.0
- 521 xiv. ethereum-optimism-op-geth-675-observability: 0.0
- 522 xv. ethereum-optimism-op-geth-679-680-observability: 0.0
- 523 xvi. git-bug-git-bug-132-449-observability: 0.0
- 524 xvii. git-bug-git-bug-1367-1370-observability: 0.0
- 525 xviii. git-bug-git-bug-264-274-observability: 1.0
- 526 xix. git-bug-git-bug-338-341-observability: 1.0
- 527 xx. paperless-ngx-paperless-ngx-10195-10196-observability: 1.0
- 528 XXI. paperless-ngx-paperless-ngx-10555-observability: 0.0
- 529 xxii. paperless-ngx-paperless-ngx-5228-10559-observability: 0.0
- 530 xxiii. paperless-ngx-paperless-ngx-6341-9305-observability: 1.0
- 531 xxiv. paperless-ngx-paperless-ngx-8910-8912-observability: 1.0
- 532 xxv. paperless-ngx-paperless-ngx-9784-observability: 0.0

533 (c) no-mcp

- 534 i. 0xpolygon-bor-1710-observability: 0.0
- 535 ii. 0xpolygon-bor-1728-1748-observability: 0.0
- 536 iii. 0xpolygon-bor-1743-observability: 0.0
- 537 iv. chainsafe-gossamer-4286-4720-observability: 0.0
- 538 v. chainsafe-gossamer-4489-4640-observability: 0.0
- 539 vi. chainsafe-gossamer-4573-4683-observability: 0.0
- 540 vii. containers-podman-compose-1221-1231-observability: 1.0
- 541 viii. containers-podman-compose-1226-1243-observability: 0.0
- 542 ix. containers-podman-compose-2-1238-observability: 1.0
- 543 x. containers-podman-compose-23-1214-observability: 1.0

- 544 xi. ethereum-optimism-op-geth-553-558-observability: 0.0
- 545 xii. ethereum-optimism-op-geth-655-observability: 0.0
- 546 xiii. ethereum-optimism-op-geth-666-observability: 0.0
- 547 xiv. ethereum-optimism-op-geth-675-observability: 0.0
- 548 xv. ethereum-optimism-op-geth-679-680-observability: 0.0
- 549 xvi. git-bug-git-bug-132-449-observability: 0.0
- 550 xvii. git-bug-git-bug-1367-1370-observability: 0.0
- 551 xviii. git-bug-git-bug-264-274-observability: 0.0
- 552 xix. git-bug-git-bug-338-341-observability: 0.0
- 553 xx. paperless-ngx-paperless-ngx-10195-10196-observability: 1.0
- 554 xxi. paperless-ngx-paperless-ngx-10555-observability: 0.0
- 555 xxii. paperless-ngx-paperless-ngx-5228-10559-observability: 0.0
- 556 xxiii. paperless-ngx-paperless-ngx-6341-9305-observability: 1.0
- 557 xxiv. paperless-ngx-paperless-ngx-8910-8912-observability: 1.0
- 558 xxv. paperless-ngx-paperless-ngx-9784-observability: 0.0

## 2. Gemini 3.1 Pro

### (a) default

- 560 i. Oxpolygon-bor-1710-observability: 0.0
- 562 ii. Oxpolygon-bor-1728-1748-observability: 0.0
- 563 iii. Oxpolygon-bor-1743-observability: 1.0
- 564 iv. chainsafe-gossamer-4286-4720-observability: 1.0
- 565 v. chainsafe-gossamer-4489-4640-observability: 0.0
- 566 vi. chainsafe-gossamer-4573-4683-observability: 0.0
- 567 vii. containers-podman-compose-1221-1231-observability: 1.0
- 568 viii. containers-podman-compose-1226-1243-observability: 0.0
- 569 ix. containers-podman-compose-2-1238-observability: 1.0
- 570 x. containers-podman-compose-23-1214-observability: 1.0
- 571 xi. ethereum-optimism-op-geth-553-558-observability: 0.0
- 572 xii. ethereum-optimism-op-geth-655-observability: 0.0
- 573 xiii. ethereum-optimism-op-geth-666-observability: 0.0
- 574 xiv. ethereum-optimism-op-geth-675-observability: 0.0
- 575 xv. ethereum-optimism-op-geth-679-680-observability: 0.0
- 576 xvi. git-bug-git-bug-132-449-observability: 0.0
- 577 xvii. git-bug-git-bug-1367-1370-observability: 1.0
- 578 xviii. git-bug-git-bug-264-274-observability: 1.0
- 579 xix. git-bug-git-bug-338-341-observability: 0.0
- 580 xx. paperless-ngx-paperless-ngx-10195-10196-observability: 1.0
- 581 xxi. paperless-ngx-paperless-ngx-10555-observability: 0.0
- 582 xxii. paperless-ngx-paperless-ngx-5228-10559-observability: 0.0
- 583 xxiii. paperless-ngx-paperless-ngx-6341-9305-observability: 1.0
- 584 xxiv. paperless-ngx-paperless-ngx-8910-8912-observability: 0.0
- 585 xxv. paperless-ngx-paperless-ngx-9784-observability: 0.0

### (b) grafana

- 587 i. Oxpolygon-bor-1710-observability: 0.0
- 588 ii. Oxpolygon-bor-1728-1748-observability: 0.0
- 589 iii. Oxpolygon-bor-1743-observability: 1.0
- 590 iv. chainsafe-gossamer-4286-4720-observability: 1.0
- 591 v. chainsafe-gossamer-4489-4640-observability: 1.0
- 592 vi. chainsafe-gossamer-4573-4683-observability: 0.0
- 593 vii. containers-podman-compose-1221-1231-observability: 1.0
- 594 viii. containers-podman-compose-1226-1243-observability: 0.0
- 595 ix. containers-podman-compose-2-1238-observability: 1.0
- 596 x. containers-podman-compose-23-1214-observability: 0.0
- 597 xi. ethereum-optimism-op-geth-553-558-observability: 0.0

598           xii. ethereum-optimism-op-geth-655-observability: 0.0  
 599           xiii. ethereum-optimism-op-geth-666-observability: 0.0  
 600           xiv. ethereum-optimism-op-geth-675-observability: 0.0  
 601           xv. ethereum-optimism-op-geth-679-680-observability: 0.0  
 602           xvi. git-bug-git-bug-132-449-observability: 0.0  
 603           xvii. git-bug-git-bug-1367-1370-observability: 0.0  
 604           xviii. git-bug-git-bug-264-274-observability: 1.0  
 605           xix. git-bug-git-bug-338-341-observability: 1.0  
 606           xx. paperless-ngx-paperless-ngx-10195-10196-observability: 1.0  
 607           xxi. paperless-ngx-paperless-ngx-10555-observability: 0.0  
 608           xxii. paperless-ngx-paperless-ngx-5228-10559-observability: 0.0  
 609           xxiii. paperless-ngx-paperless-ngx-6341-9305-observability: 1.0  
 610           xxiv. paperless-ngx-paperless-ngx-8910-8912-observability: 0.0  
 611           xxv. paperless-ngx-paperless-ngx-9784-observability: 0.0  
 612       (c) no-mcp  
 613           i. 0xpolygon-bor-1710-observability: 0.0  
 614           ii. 0xpolygon-bor-1728-1748-observability: 0.0  
 615           iii. 0xpolygon-bor-1743-observability: 1.0  
 616           iv. chainsafe-gossamer-4286-4720-observability: 1.0  
 617           v. chainsafe-gossamer-4489-4640-observability: 0.0  
 618           vi. chainsafe-gossamer-4573-4683-observability: 0.0  
 619           vii. containers-podman-compose-1221-1231-observability: 1.0  
 620           viii. containers-podman-compose-1226-1243-observability: 0.0  
 621           ix. containers-podman-compose-2-1238-observability: 0.0  
 622           x. containers-podman-compose-23-1214-observability: 0.0  
 623           xi. ethereum-optimism-op-geth-553-558-observability: 0.0  
 624           xii. ethereum-optimism-op-geth-655-observability: 0.0  
 625           xiii. ethereum-optimism-op-geth-666-observability: 0.0  
 626           xiv. ethereum-optimism-op-geth-675-observability: 0.0  
 627           xv. ethereum-optimism-op-geth-679-680-observability: 0.0  
 628           xvi. git-bug-git-bug-132-449-observability: 0.0  
 629           xvii. git-bug-git-bug-1367-1370-observability: 0.0  
 630           xviii. git-bug-git-bug-264-274-observability: 0.0  
 631           xix. git-bug-git-bug-338-341-observability: 0.0  
 632           xx. paperless-ngx-paperless-ngx-10195-10196-observability: 1.0  
 633           xxi. paperless-ngx-paperless-ngx-10555-observability: 0.0  
 634           xxii. paperless-ngx-paperless-ngx-5228-10559-observability: 0.0  
 635           xxiii. paperless-ngx-paperless-ngx-6341-9305-observability: 1.0  
 636           xxiv. paperless-ngx-paperless-ngx-8910-8912-observability: 0.0  
 637           xxv. paperless-ngx-paperless-ngx-9784-observability: 0.0

638 **NeurIPS Paper Checklist**

639 **1. Claims**

640 Question: Do the main claims made in the abstract and introduction accurately reflect the  
641 paper’s contributions and scope?

642 Answer: [Yes]

643 Justification: In the abstract and introduction, we claim to define tool efficiency and marginal  
644 tool utility, and implement an LLM-as-a-Judge evaluator to determine the usefulness of a  
645 tool. These claims accurately reflect the paper’s contributions and scope.

646 **2. Limitations**

647 Question: Does the paper discuss the limitations of the work performed by the authors?

648 Answer: [Yes]

649 Justification: The paper discusses the limitations of the work performed by the authors in a  
650 dedicated section. We self-identified limitations in our marginal tool utility calculation and  
651 types of tools we were able to assess.

652 **3. Theory assumptions and proofs**

653 Question: For each theoretical result, does the paper provide the full set of assumptions and  
654 a complete (and correct) proof?

655 Answer: [N/A]

656 Justification: This paper does not include theoretical results. We formalized definitions that  
657 need not be proven, rather simply calculated and supported by empirical results.

658 **4. Experimental result reproducibility**

659 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
660 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
661 of the paper (regardless of whether the code and data are provided or not)?

662 Answer: [Yes]

663 Justification: We have provided the system and user prompts given to the LLM-as-a-Judge  
664 and disclosed its exact provider. The experiments were ran against a publicly available  
665 benchmark and used proprietary LLMs for which we have disclosed the exact providers.

666 **5. Open access to data and code**

667 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
668 tions to faithfully reproduce the main experimental results, as described in supplemental  
669 material?

670 Answer: [Yes]

671 Justification: The code to run the benchmarks are all publicly available. We have also  
672 provided open access to our original code with sufficient instructions to faithfully reproduce  
673 the main experimental results.

674 **6. Experimental setting/details**

675 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
676 rameters, how they were chosen, type of optimizer) necessary to understand the results?

677 Answer: [Yes]

678 Justification: Our paper does not involve training/fine-tuning, but we have disclosed all  
679 relevant hyperparameters/information to reproduce our results.

680 **7. Experiment statistical significance**

681 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
682 information about the statistical significance of the experiments?

683 Answer: [No]

684 Justification: Unfortunately, each experiment run is very time-intensive and computationally-  
685 expensive, hence we are only able to run a handful of experiments insufficient to construct  
686 error bars. However, in the limited samples we have run, the results remain consistent.

687 **8. Experiments compute resources**  
688 Question: For each experiment, does the paper provide sufficient information on the com-  
689 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
690 the experiments?  
691 Answer: [Yes]  
692 Justification: We have provided sufficient information on the computer resources needed to  
693 reproduce the experiments. Further details can be found in the shared code repository.

694 **9. Code of ethics**  
695 Question: Does the research conducted in the paper conform, in every respect, with the  
696 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)  
697 Answer: [Yes]  
698 Justification: We have reviewed the NeurIPS Code of Ethics and we confirm that our research  
699 in this paper conforms in every respect with it.

700 **10. Broader impacts**  
701 Question: Does the paper discuss both potential positive societal impacts and negative  
702 societal impacts of the work performed?  
703 Answer: [Yes]  
704 Justification: The paper discusses both potential positive societal impacts and negative  
705 societal impacts of the work performed.

706 **11. Safeguards**  
707 Question: Does the paper describe safeguards that have been put in place for responsible  
708 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
709 image generators, or scraped datasets)?  
710 Answer: [N/A]  
711 Justification: Our work in this paper does not involve the aforementioned class of artifacts  
712 and therefore our paper does not pose such risks.

713 **12. Licenses for existing assets**  
714 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
715 the paper, properly credited and are the license and terms of use explicitly mentioned and  
716 properly respected?  
717 Answer: [Yes]  
718 Justification: We have properly respected the CC-BY-4.0 license of the APEX-SWE bench-  
719 mark.

720 **13. New assets**  
721 Question: Are new assets introduced in the paper well documented and is the documentation  
722 provided alongside the assets?  
723 Answer: [Yes]  
724 Justification: We have written documentation and provided said documentation with the  
725 assets.

726 **14. Crowdsourcing and research with human subjects**  
727 Question: For crowdsourcing experiments and research with human subjects, does the paper  
728 include the full text of instructions given to participants and screenshots, if applicable, as  
729 well as details about compensation (if any)?  
730 Answer: [N/A]  
731 Justification: The paper does not involve crowdsourcing nor research with human subjects.

732 **15. Institutional review board (IRB) approvals or equivalent for research with human**  
733 **subjects**

734 Question: Does the paper describe potential risks incurred by study participants, whether  
735 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
736 approvals (or an equivalent approval/review based on the requirements of your country or  
737 institution) were obtained?

738 Answer: [N/A]

739 Justification: The paper does not involve crowdsourcing nor research with human subjects.

740 **16. Declaration of LLM usage**

741 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
742 non-standard component of the core methods in this research? Note that if the LLM is used  
743 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
744 scientific rigor, or originality of the research, declaration is not required.

745 Answer: [N/A]

746 Justification: The core method development in this research does not involve LLMs as any  
747 important, original, or non-standard components.