

TRAIL: Trace Reasoning and Agentic Issue Localization

Anonymous ACL submission

Abstract

The increasing adoption of agentic workflows across diverse domains brings a critical need to scalably and systematically evaluate the complex traces these systems generate. Current evaluation methods depend on manual, domain-specific human analysis of lengthy workflow traces—an approach that does not scale with the growing complexity and volume of agentic outputs. Error analysis in these settings is further complicated by the interplay of external tool outputs and language model reasoning, making it more challenging than traditional software debugging. In this work, we (1) articulate the need for robust and dynamic evaluation methods for agentic workflow traces, (2) introduce a formal taxonomy of error types encountered in agentic systems, and (3) present a set of 148 large human-annotated traces (TRAIL) constructed using this taxonomy and grounded in established agentic benchmarks. To ensure ecological validity, we curate traces from both single and multi-agent systems, focusing on real-world applications such as software engineering and open-world information retrieval. Our evaluations reveal that modern long context LLMs perform poorly at trace debugging, with the best GEMINI-2.5-PRO model scoring a mere 11% on TRAIL. Our dataset and code are made publicly available to support and accelerate future research in scalable evaluation for agentic workflows¹.

1 Introduction

The rapid advancement of large language models (LLMs) has catalyzed the development of agentic systems capable of automating difficult, multi-step tasks across various domains such as software engineering and multi-hop IR (Ma et al., 2023; OpenAI, 2024; Nguyen et al., 2024; Wang et al., 2025a). Unlike traditional generative models, agents can interact with diverse tools and dynamically navigate

¹Hidden for double-blind review

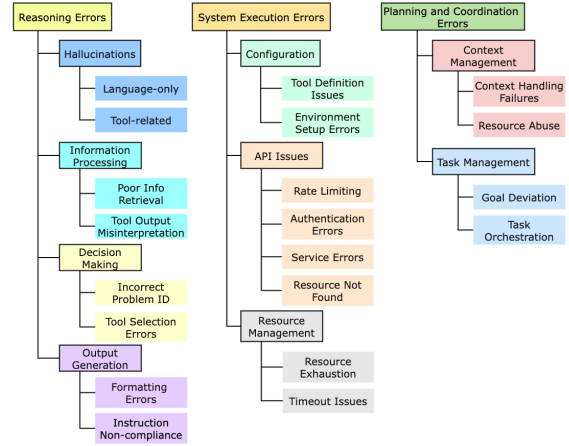


Figure 1: Illustration of the TRAIL taxonomy of errors

environments, often with minimal human supervision (Wang et al., 2024a). This escalation of system complexity demands more challenging and multifaceted evaluation processes (Nasim, 2025) and has led to the adoption of LLMs as evaluators for such agentic systems (Zheng et al., 2023; Chen et al., 2024; Kim et al., 2024; Zhu et al., 2025; Deshpande et al., 2024a).

However, as multi-agent systems scale and become integral to real-world workflows, evaluating and debugging their performance remains a significant challenge. Agentic non-determinism (Laban et al., 2025; Patronus AI, 2025) and multi-step task solving (Mialon et al., 2023; Yao et al., 2024) demand greater observability than the simple end-to-end evaluations offered by existing benchmarks (Kapoor et al., 2024a; Zhuge et al., 2024; Moshkovich et al., 2025; Cemri et al., 2025). Such complex environments require granular taxonomies and well-annotated traces that can serve as references for debugging and root-cause analysis of agent behaviors (Cemri et al., 2025). When creating taxonomies and benchmarks to test and improve agents, we must ensure these are grounded in real-world applications and are not centered around

dummy data (Bowman and Dahl, 2021; Liu et al., 2024b). Previous agent trace analysis frameworks have primarily focused on parsed traces containing unstructured text (Cemri et al., 2025), which do not adequately represent common agent framework outputs that generate structured traces logged in standardized formats like opentelemetry (OpenTelemetry, 2025). Additionally, as observed by Guo et al. (2023); Sui et al. (2024), handling structured data remains challenging for LLMs, an observation corroborated by previous research on automated software engineering trace analysis (Roy et al., 2024a; Ma et al., 2024b). These limitations highlight the need for new approaches specifically designed for structured agentic traces. To address these challenges and facilitate the analysis and evaluation of agentic executions, we propose a formal error taxonomy, shown in Figure 3, that promotes granular failure diagnosis. We also present a carefully curated, turn-level annotated trace dataset called TRAIL (Trace Reasoning and Agentic Issue Localization), which demonstrates the validity and practical utility of our proposed taxonomy.

In our work, we utilize and build on SWE-Bench (Jimenez et al., 2024; Aleithan et al., 2024) and GAIA (Mialon et al., 2023) while addressing three major shortcomings inherent to previous automatic agent evaluation paradigms. Firstly, we aim to replace end to end analysis of agents with a benchmark containing step-level analysis of traced agentic workflows. Secondly, we address the need for grounding in real scenarios by producing opentelemetry-based structured traces that span beyond present model context length limits. Finally, as compared to benchmarks focused only on agentic reasoning and coordination (Cemri et al., 2025; Kokel et al., 2025), TRAIL focuses on validity through addition of finer, more aligned system execution failures and planning error categories such as *API errors* and *Task Orchestration Errors* to our taxonomy. Such categories are not only relevant to model developers but also to users and engineers optimizing single and multi-agent AI applications. The contributions of our work are as follows:

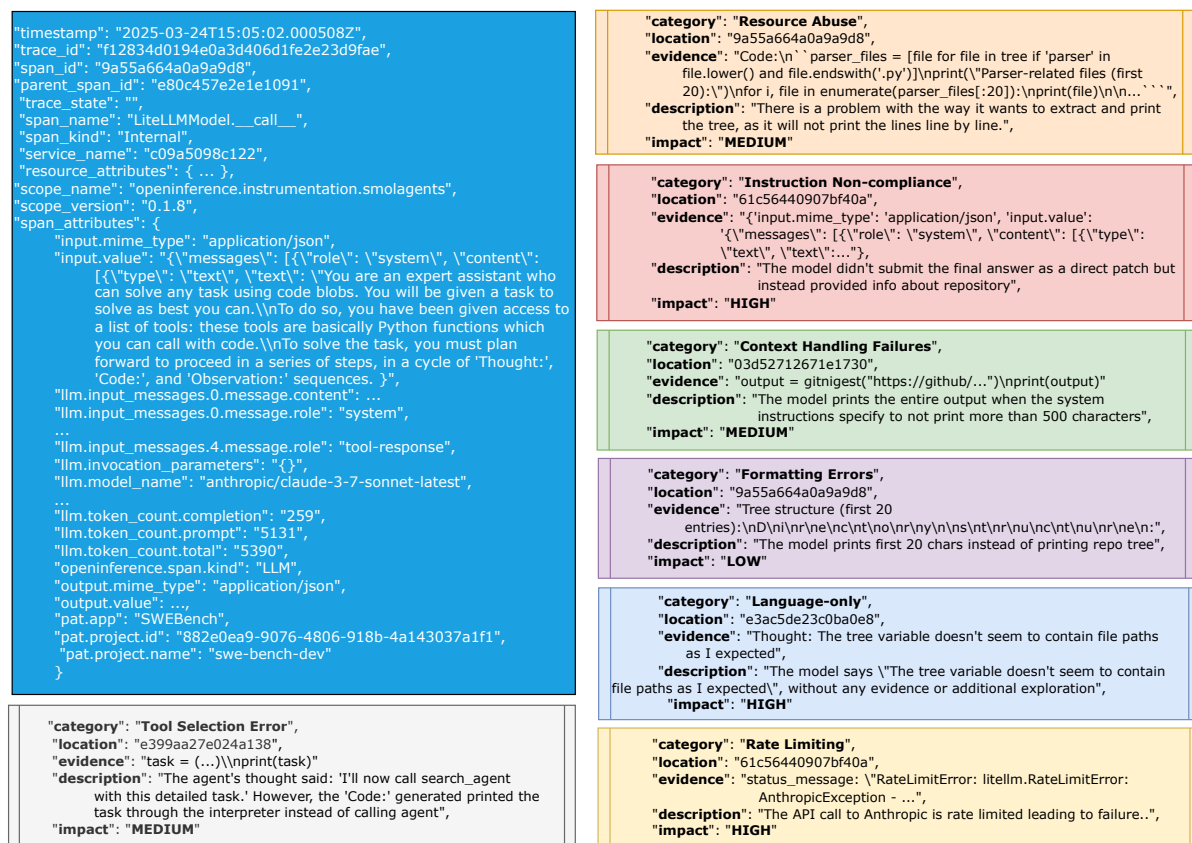
- We introduce a formal taxonomy (Figure 1) that defines, fine-grained agentic error categories spanning across three key areas: reasoning, planning, and execution.
- Based on this taxonomy, we present TRAIL, an ecologically grounded execution trace

benchmark comprising 148 meticulously curated traces (totaling 1987 open telemetry spans, of which 575 exhibit at least one error) drawn from the GAIA (Mialon et al., 2023) and SWE-Bench (Jimenez et al., 2024) datasets and covering a wide range of tasks.

- We show that TRAIL is a non-trivially difficult benchmark for LLMs on many fronts
 1. Current SOTA LLM families such as O3, CLAUDE-3.7-SONNET and GEMINI-2.5-PRO perform modestly at best on TRAIL, both in terms of predicting error categories and their location. With GEMINI-2.5-PRO the best performing model, achieving only 11% combined joint accuracy on both splits.
 2. Solving TRAIL requires a significant fraction of the maximum input length of LLMs (or exceeds it), as well as requires generating significant fraction of their maximum output (See Table 2, Figure 5)
 3. Models benchmarked on TRAIL benefit from both the presence and greater extent of reasoning chains (§5.1.4, §5.1.5), highlighting the need for improvement in exploration capabilities of LLMs.
- TRAIL is fully open-source (MIT License), will be accompanied by a HuggingFace leaderboard, and serves as a foundation for future research on evaluating agentic workflows.

2 Relevant Work

LLM-as-a-Judge Shortcomings of conventional metrics such as ROUGE, BLEU, and BERTScore (Schluter, 2017; Freitag et al., 2020; Hanna and Bojar, 2021) has led to the wide adoption of LLMs as evaluators and critics of other AI systems (Zheng et al., 2023; Zhu et al., 2025; Chen et al., 2025, 2024; Kim et al., 2024). Recent approaches have enhanced LLM judges’ reasoning capabilities through techniques like unconstrained evaluation plan and specialized training methods that enable more robust evaluation performance across diverse scenarios (Lightman et al., 2023; Wang et al., 2024e; Trivedi et al., 2024; Saha et al., 2025). The evaluation landscape has evolved significantly with the introduction of frameworks like FLASK (Ye et al., 2024b) which decompose coarse-level scoring into skill set-level evaluations for each instruction, demonstrating



high correlation between model-based and human-based evaluations. The Prometheus models (Kim et al., 2023, 2024, 2025) established a significant benchmark by creating judge models that surpass GPT-4 in ranking for subjective evaluation criteria. Their research also examined how performance deteriorates as subjectivity increases. More recently, several studies have enhanced judge model performance through external augmentations and checklists, highlighting the importance of incorporating high-quality reasoning chains and human guidance in model training (Lee et al., 2025; Deshpande et al., 2024b,a; Chen et al., 2025; Wang et al., 2025b). Despite promising advancements, LLM judges have shown issues with propagation of biases and lack of robustness to longer inputs (Ye et al., 2024a; Hu et al., 2024b; Wei et al., 2024; Zhou et al., 2025). Since trace evaluation requires robust reasoning over large contexts (Tian et al., 2024), LLM judges have not seen wide application in this sector yet.

Agentic Evaluation LLM-powered agents have gained significant traction for their capacity to manage intricate, sequential tasks while adaptively en-

gaging with varied environments, rendering them particularly valuable for practical real-world applications such as software engineering and multi-hop IR (Ma et al., 2023; OpenAI, 2024; Nguyen et al., 2024; Wang et al., 2025a; Jimenez et al., 2024; Qian et al., 2024; Wang et al., 2024d; Patil et al., 2024). However, the performance gains of multi-agent frameworks remain minimal compared to their single-agent counterparts (Xia et al., 2024; Kapoor et al., 2024b). As these agentic systems become more prevalent, evaluation frameworks (as compared to LLM evaluation) must offer greater customization and granularity to effectively assess the complex and sometimes unpredictable interactions between multiple agents, enabling users to precisely identify and diagnose errors at each step of the process (Roy et al., 2024b; Akhtar et al., 2025; Jiang et al., 2025; Zhuge et al., 2024; OpenManus, 2024).

Agent Benchmarks Software engineering domain has become a fertile testbed for LLM-based collaborative problem solving for real-world use cases and to evaluate agents’ ability to handle realistic coding tasks. SWE-Bench (Jimenez et al.,

2024; Aleithan et al., 2024; Pan et al., 2024) was introduced as a grounded benchmark asking whether LLMs can resolve real-world GitHub issues. Similarly, GAIA (Mialon et al., 2023) is a benchmark for General AI Assistants featuring real-world questions requiring reasoning, tool use, and multimodality. AssistantBench (Yoran et al., 2024) introduces a challenging benchmark of realistic, time-consuming web tasks to evaluate web agents. For agents, it is key to distinguish input sample failures from the judge model’s own internal reasoning failures. Highlighting spans can help models focus and avoid losing context while also providing additional explainability and performance improvements (Lv et al., 2024; Li et al., 2024). Other core benchmarks include DevAI (Zhuge et al., 2024), MLE-bench (Chan et al., 2024), HumanEval (Du et al., 2024), and MBPP (Odena et al., 2021).

Traces and Error Taxonomies Emerging work has emphasized the need for better observability in the agent execution traces to diagnose and manage the non-deterministic nature of agentic systems (Kapoor et al., 2024a; Zhuge et al., 2024; Moshkovich et al., 2025; Cemri et al., 2025). For instance, Roy et al. (2024a) explores using LLM-based agents to dynamically collect diagnostic information from logs and metrics using retrieval tools for root cause analysis of cloud system incidents. Akhtar et al. (2025) surveys how LLMs are being applied to automate even log analysis in security contexts. Jiang et al. (2025) is a log analysis framework for diagnosing large-scale LLM failures based on studying real-world training failures. Ma et al. (2024c) explores the potential for log parsing by proposing an LLMParse delivering comprehensive evaluations in various settings. Once the trace errors are found, to serve as references for users to debug or conduct root cause analysis of agent behaviors, these errors require a granular taxonomy (Cemri et al., 2025; Kokel et al., 2025; Bai et al., 2024a). MAST (Cemri et al., 2025) presents an empirically grounded failure mode taxonomy but focusing only on agentic reasoning and coordination. ACPBench (Kokel et al., 2025), using a synthetic dataset, focuses on atomic reasoning about action and is designed to evaluate LLM’s core planning skills. Other related work includes taxonomies to evaluate multi-turn conversations (Bai et al., 2024a) and designing LLM agent framework to identify and quantify complex evaluation criteria (Arabzadeh et al., 2024; Epperson et al.,

2025).

Thus, TRAIL distinguishes itself through its ecological validity while comprehensively addressing both single and multi-turn systems with its granular taxonomy, particularly emphasizing critical execution and planning failure patterns.

3 Agentic Error Taxonomy

LLM reasoning, while having advanced significantly, remains a critical source of failures in agentic workflows (Costarelli et al., 2024). These errors span several dimensions, from flawed information generation to problematic decision-making and output production (Cemri et al., 2025). In this section, we define a comprehensive taxonomy (as summarized in Figure 3) of agentic errors spanning three key areas of failures: reasoning, planning and coordination, and system execution.

3.1 Reasoning Errors

Hallucinations LLMs frequently generate factually incorrect or nonsensical content, a problem that also affects agents (Huang et al., 2025; Ji et al., 2023). *Text-only* hallucinations include fabricated or ungrounded statements that conflict with real-world knowledge (Ji et al., 2023). In contrast, *Tool-related* hallucinations arise when agents invent tool outputs or misunderstand tool functions, such as fabricating results or claiming nonexistent capabilities (Zhang et al., 2024b; Xu et al., 2024).

Information Processing Retrieval-augmented generation, which retrieves and reasons over data relevant to a query, has become increasingly popular (Hu and Lu, 2024; Gao et al., 2025). However, recent work (Xu et al., 2025; Su et al., 2025) shows that LLMs and agents often struggle to reason effectively over retrieved information. These issues can be grouped into two main types: poor information retrieval and misinterpretation of outputs. *Poor information retrieval* (Wu et al., 2024) can introduce redundancy and content overload (Stechly et al., 2024), while misinterpretation of retrieved context (*Tool output Misinterpretation*) (Karpinska et al., 2024; Wang et al., 2024b) may cause errors that propagate throughout an agent’s reasoning process, leading to broader incorrectness or inefficiencies.

Decision Making Task misunderstanding at the step level often arises from ambiguous prompts, unclear instructions, or an LLM’s inability to distinguish between prompt and data instruc-

tions (Zverev et al., 2024). Detecting such misunderstandings (*Incorrect Problem ID*) requires analyzing an agent’s path, which is challenging in large contexts (Yuan et al., 2024) and reliable detection of these errors is crucial for agent improvement. Furthermore, effective decision making in agent workflows also depends on selecting the appropriate tool at each step (Qin et al., 2023). Because optimal planning and tool selection reduces cost and increases efficiency (Yehudai et al., 2025), we place *Tool Selection Error* under *Decision Making*.

Output Generation LLMs often produce incorrectly formatted structured outputs (Shorten et al., 2024; Liu et al., 2024a), which is problematic for tool calls that need precise JSON or code formatting. To capture this, our taxonomy includes *Formatting Errors*. Similarly, LLMs frequently struggle following complex/ambiguous instructions (White et al., 2024; Heo et al., 2024), hence we subcategorize *Instruction Non-compliance*.

3.2 System Execution Errors

Configuration Issues Incorrect agentic environment configuration can cause failures and limit agent capability (Hu et al., 2024a). One key issue is *Incorrect Tool Definition*, as shown by Fu et al. (2024), agents can be misled by inaccurate or obfuscated tool definitions in prompts, posing security and reliability risks. Additionally, poor setup of environment variables (*Environment Setup Errors*), e.g., missing API keys or incorrect file permissions, can cause unexpected failures and disrupt reasoning paths.

API and System Issues As agentic systems combine LLMs with software tools, tool usage or implementation errors can disrupt workflows. With the rise of remote tool access via protocols like MCP (Anthropic, 2025), capturing and categorizing API failures is increasingly important for prompt reporting to tool developers (Shen, 2024). Runtime errors involving agentic tools remain underexplored (Milev et al., 2025), so we specifically include the most common API tool errors in our taxonomy: *Rate Limiting* (429), *Authentication Errors* (401, 403), *Service Errors* (500), and *Resource Not Found* (404) (Liu et al., 2023a).

Resource Management Resource management is crucial for agents using operating system tools like interpreters or terminals. Poor task planning can expose vulnerabilities, such as *Resource Ex-*

haustion from overallocation (Ge et al., 2023) or *Timeout Issues* from infinite loops (Zhang et al., 2024a), potentially causing memory overflows or system overloads. Early detection of these errors is vital to prevent infrastructure failures.

3.3 Planning and Coordination Errors

Context Management As planning and reasoning become integral to agentic workflows (Yao et al., 2023; Ke et al., 2025), agents must manage long-term context, including episodic and semantic information (Zhang et al., 2024c). In our taxonomy, we categorize failures in context or instruction retention as *Context Handling Failures*. Additionally, repeated tool calls (Kokane et al., 2024) (*Resource Abuse*) reflect shortcomings in planning, context management, and tool use, which our taxonomy also captures.

Task Management Environmental misconfigurations or LLM hallucinations can distract agentic systems, and poor recovery from such distractions often leads to goal deviation (Ma et al., 2024a). These issues are amplified in multi-agent setups with sub-tasks, making effective task orchestration crucial. Therefore, we include *Goal Deviation* and *Task Orchestration Errors* in our taxonomy.

4 TRAIL Benchmark

TRAIL is a benchmark aimed to evaluate LLM capabilities to analyze and evaluate long, structured, opentelemetry standardized agentic executions. TRAIL follows our fine grained taxonomy and contains 148 carefully annotated agentic traces. The dataset uses text-only data instances from the GAIA (Mialon et al., 2023) and SWE Bench Lite (Jimenez et al., 2024) datasets, spanning multiple information retrieval and software bug fixing tasks. It contains a total of 841 annotated errors, averaging at 5.68 errors per trace Figure 3.

4.1 Goals and Design Choices

Core Agent Task We aim to showcase realistic agentic workflows and so we target two widely adopted agentic datasets, the GAIA benchmark (Mialon et al., 2023), an open world search task, and the SWE-Bench-Lite (Jimenez et al., 2024) dataset, for locating and fixing issues in Github repositories. We select these datasets due to their challenging nature and necessity for environment and search space exploration.

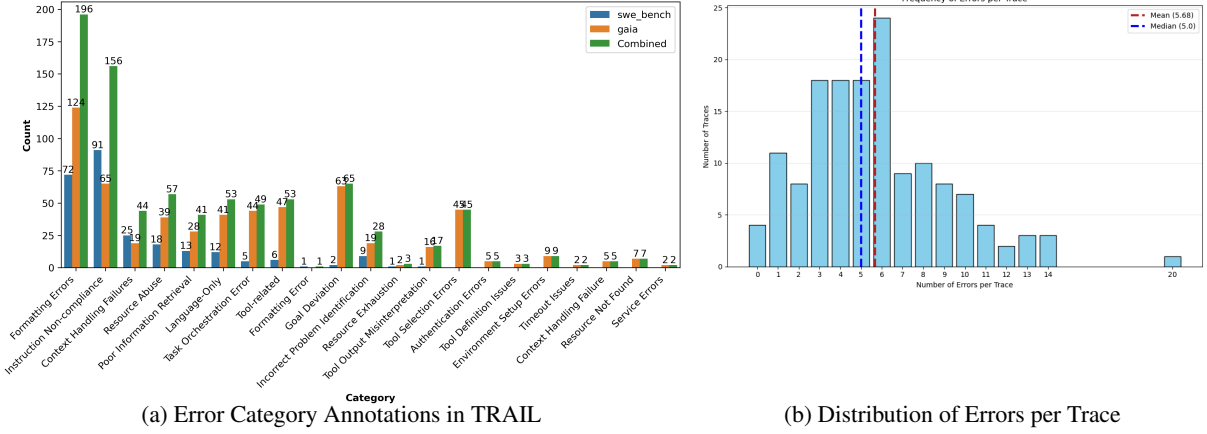


Figure 3: TRAIL Dataset Statistics

Agent Orchestration Liu et al. (2023b) first presented a standardized hierarchical method of orchestrating agents, derivatives of which are actively adopted by several works (Zhao et al., 2024, 2025). We closely follow this hierarchical structure and adopt the Hugging Face OpenDeepResearch agent (Hugging Face, 2024) for creating traces for the GAIA benchmark. We select the state-of-the-art o3-mini-2025-01-31 (OpenAI, 2025d) and assign it as the backbone model for the manager and search agents respectively because of its strong tool use and planning ability as showcased by Phan et al. (2025). For more information, refer to §A.10.

Parallely, to explore single-agent planning errors and elicit context handling errors for the SWE-Bench split, we use a CodeAct agent (Wang et al., 2024c) and provide it access to a sandboxed environment, a python interpreter and the gitingest² library. We select claude-3-7-sonnet-20250219 as the backbone model due to its strong performance on software engineering tasks (Anthropic, 2025). To further organically introduce errors into this agent system, we add instructional constraints such as output length limits and force exploration via prompts. The complete prompt is at §A.12.

Workflow Tracing To ensure compatibility of this dataset with real world tracing and observability software, all traces are collected via opentelemetry (OpenTelemetry, 2025), specifically, its most widely adopted open-source derivative compatible with agents, the openinference standard (Arize AI, 2025) as adopted by Moshkovich et al. (2025).

²<https://github.com/cyclotruc/gitingest>

4.2 Data Annotation and Validation

We selected four annotators with expertise in software engineering and log debugging to label our agent traces. To assess agreement, a separate set of 63 traces was assigned. Results based on these indicate high inter-annotator agreement during curation. We defer details of our complete annotation and agreement measuring processes and actual numbers from them to §A.7.

4.3 Dataset Analysis

Following the post-annotation review, we found errors in 114 GAIA traces and 30 from SWE Bench. As shown in Figure 3, these errors cover various categories, with most falling under *Output Generation*. Specifically, *Formatting Errors* and *Instruction Non-compliance* make up 353 of 841 total errors—nearly 42%. In contrast, *System Execution Errors* are rare. This categorical imbalance highlights two important considerations for evaluating agentic pipelines. First, the prevalence of *Output Generation* errors suggests that current LLM systems struggle with high-level reasoning and understanding task parameters, even with careful prompt-engineering. Second, although infrequent, errors in categories like API failures can be catastrophic and are critical to detect, as they are often difficult to recover from, unlike errors due to goal deviation or tool misinterpretation. Most errors in our data are high or medium impact (Figure 6a). While model hallucinations and resource management issues greatly affect agent behavior, about 44% of *Output Generation* errors are low impact (Figure 6b). This underscores need for a classification scheme that includes rare but significant error types. A key feature of our taxonomy is ability to

Model	TRAIL (GAIA)				TRAIL (SWE Bench)			
	Cat. F1	Loc. Acc.	Joint	ρ	Cat. F1	Loc. Acc.	Joint	ρ
LLAMA-4-SCOUT-17B-16E-INSTRUCT [†]	0.041	0.000	0.000	0.134	0.050	0.000	0.000	0.264
LLAMA-4-MAVERICK-17B-128E-INSTRUCT [†]	0.122	0.023	0.000	0.338	0.191	0.083	0.000	-0.273
GPT-4.1 [†]	0.218	0.107	0.028	0.411	0.166	0.000	0.000	0.153
OPEN AI O1 [*]	0.138	0.040	0.013	0.450	CLE	CLE	CLE	CLE
OPEN AI O3 [*]	0.296	0.535	0.092	0.449	CLE	CLE	CLE	CLE
ANTHROPIC CLAUDE-3.7-SONNET [*]	0.254	0.204	0.047	0.738	CLE	CLE	CLE	CLE
GEMINI-2.5-PRO-PREVIEW-05-06 ^{*†}	0.389	0.546	0.183	0.462	0.148	0.238	0.050	0.817
GEMINI-2.5-FLASH-PREVIEW-04-17 ^{*†}	0.337	0.372	0.100	0.550	0.213	0.060	0.000	0.292

Table 1: Performance across LLMs for Error Categorization & Localization on TRAIL (GAIA) and TRAIL (SWE Bench). Models marked with * have reasoning set to "high"; [†] indicates 1M+ token context window. Insufficient context length is marked as CLE. Pearson correlation b/w overall human and generated scores is shown under ρ .³

categorize well such long-tail, high-impact errors.

4.4 Summary of Evaluation Setup

For empirically evaluating and comparing LLM performance on TRAIL we choose the following LLMs — GPT-4.1, O1, O3, GEMINI-2.5 (both PRO+FLASH), CLAUDE-3.7-SONNET and LLAMA-4 (both Maverick+Scout). We defer detailed discussion of more evaluation setup specifics to A.3

5 Results

In §5.1, we analyze the research questions below:

- How does long context reasoning affect TRAIL performance? How many inputs exceed the LLM’s context window? How does trace length impact this? We address these in §5.1.1 §5.1.2, and §5.1.3.
- Does TRAIL benefit from more reasoning? We explore this in §5.1.4 and §5.1.5.
- Which error categories are easier to predict? Where do non-reasoning models perform notably worse? We examine this in §5.1.6.

5.1 Qualitative and Quantitative Analysis

5.1.1 Task Difficulty - Context Length and Generation Horizon

As seen in Table 2, the distribution of raw JSON input token lengths injected to perform our task cuts close to the input context limit of several LLMs - with the maximum input trace length always being twice longer than the input length limit, and even the mean itself sometimes going over. Furthermore, even the typical output token length horizon the LLMs need to generate for the task exceeds the

1K tokens on average, with the maximum being $\approx 3.7K$ at the least. Besides being a significant % of the maximum output length, this indicates the difficultly long generation horizon TRAIL needs.

5.1.2 Long Context Ability and Model Performance

We compare how the models in Table 1 rank based on their aggregate performance on TRAIL vis-a-vis the relative ranking of the subsets of these models that occur on updated long-context benchmark leaderboards Longbenchv2 and fiction.live’s Long-ContextBench (Bai et al., 2024b; Ficlive, 2025), and notice this differs for only one model (o3 being third best rather than best on the latter). We defer the complete detail of these rankings to §A.2.

5.1.3 Performance vs Input Length

We find all performance metrics to be anti-correlated with input length, as detailed in Table 3. This supports the hypothesis that longer input raw traces increase the difficulty of TRAIL for models.

5.1.4 Reasoning vs Non-Reasoning Models

From Table 1, we see all reasoning models except O1 outperforming non-reasoning ones on both Error category F1 and Location Accuracy. On Joint Accuracy, the gap between the two families is larger — Reasoning models other than o1 perform at 1.5-8 times the best performing non-reasoning model.

5.1.5 Does Reasoning Effort Matter?

To systematically assess the impact of reasoning extent, we experiment with the same model (O3) at "high," "medium," and "low" reasoning effort levels, as set by OpenAI’s *reasoning.effort* parameter. We find that all three metrics, including Category F1 ($0.296 \rightarrow 0.277 \rightarrow 0.264$), decrease as reasoning effort decreases. These results empirically sup-

³All reported results are an average of three runs.

Task	Tokenizer	Input Limit	Output Limit	Input Context Lengths				Output Token Lengths			
				Min	Max	Mean	StdDev	Min	Max	Mean	StdDev
GAIA	gpt-4.1 (=o3)	1M	32.77K	20.94K	7.50M	286.85K	768.85K	0.11K	4.47K	1.11K	0.69K
GAIA	gemini-2.5	1M	8.19K	23.09K	8.25M	313.49K	843.53K	0.13K	4.95K	1.20K	0.75K
GAIA	claude-3.7	200K	128K	23.67K	2.66M	262.67K	456.64K	0.12K	5.37K	1.23K	0.78K
SWEBench	gpt-4.1 (=o3)	1M	32.77K	120.40K	2.05M	616.92K	473.05K	0.11K	3.71K	1.71K	0.75K
SWEBench	gemini-2.5	1M	8.19K	134.88K	2.21M	698.09K	552.34K	0.13K	4.09K	1.88K	0.83K
SWEBench	claude-3.7	200K	128K	140.16K	2.43M	727.75K	557.86K	0.12K	4.17K	1.93K	0.87K

Table 2: Input Context Lengths and Human-Annotated Output Token Lengths Across both GAIA and SWEBench Tasks and various SOTA models and their tokenizers. Input Length aggregates that exceed the limit are **highlighted**.

Corr.	Location Acc	Joint Acc	Categ. F1
Pearson (r)	-0.379	-0.291	-0.296
Spearman (ρ)	-0.508	-0.349	-0.225

Table 3: Correlations b/w Input Length & Performance

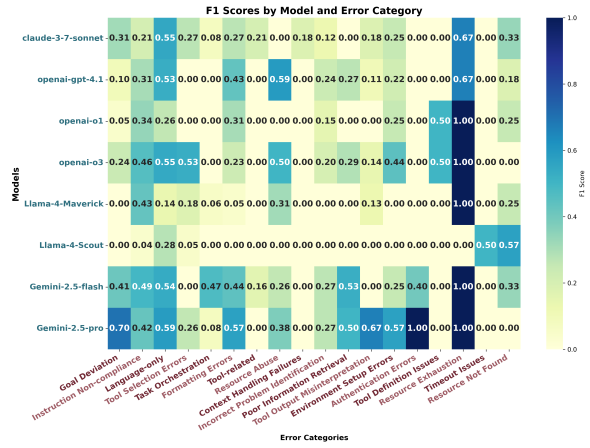


Figure 4: Heatmap for Error Category F1 across models; categories are ordered left to right based on their support

port that TRAIL performance benefits from higher reasoning effort at test time, and that the superior results for reasoning models are not solely due to improved pre- or post-training (§5.1.4). Full ablation results are in Appendix §A.4.

5.1.6 Performance Across Categories

Hard-to-Predict Categories Among the most challenging categories, *Context Handling Failures* stand out, as nearly all models score an F1 of 0.00, indicating these errors demand advanced reasoning. The only exception is CLAUDE-3.7-SONNET, which achieves a relatively better score of 0.18. *Tool Selection Errors* are also difficult to predict, with most models scoring between 0.00 and 0.08, apart from GEMINI-2.5-PRO (0.26), CLAUDE-3.7-SONNET (0.27), and especially O3 (0.53), suggesting this is a complex error type. Similarly, *Task Orchestration* shows uniformly low scores across models (0.00–0.08) except for GEMINI-2.5-FLASH,

which stands out with a much higher F1 of 0.47.

Interesting Performance Divergence There are also categories where model performance diverges interestingly. For *Goal Deviation*, GEMINI-2.5-PRO and GEMINI-2.5-FLASH perform best (0.70 and 0.41, respectively), while CLAUDE-3.7-SONNET and O3 perform moderately (0.31, 0.24); O1 and other non-reasoning models score the lowest (≤ 0.05). In the case of *Poor Information Retrieval*, the two Gemini models are again notably better (0.50 and 0.53), with others at <0.30 , suggesting better diagnosis of failures related to context.

Other Surprising Patterns *Language-Only* errors, a subtype of hallucination, are detected relatively well by all models (0.14–0.59), implying that these are easier for models to predict even without advanced reasoning capabilities. For *Formatting Errors*, performance is non-monotonic: GPT-4.1 (0.43) and the GEMINI-2.5 models (0.44–0.57) perform well, while O1, O3, and CLAUDE-3.7-SONNET perform worse (0.23–0.31). It is notable that O1 and GPT-4.1 outscore O3 on this category, despite being older and non-reasoning respectively. We defer some model-specific observations to §A.6

6 Conclusion

In this work, TRAIL, a new taxonomy for classifying agentic errors, along with an expert-curated dataset of 148 agentic problem instances and 841 unique errors from GAIA and SWE Bench. Current SOTA models perform poorly as LLM Judges on this dataset, with GEMINI 2.5-PRO achieving only 18% joint accuracy on GAIA and 5% on SWE Bench; three out of eight models cannot even process the full context. These results highlight that existing models struggle to systematically evaluate complex agentic traces, due to the inherent complexity of agentic systems and LLM context limitations. A new framework is needed for scalable, systematic evaluation of agentic workflows.

Limitations

The TRAIL dataset and taxonomy are primarily focused on text-only inputs and outputs but recent advancements in multimodal agentic systems require careful extension of the taxonomy to handle errors arising from new categories such as multimodal tool use. One additional limitation of TRAIL is the large number of tail categories with very few examples. It is important to ensure correctness of LLM-Judges on these categories due to the high-impact nature of the failures. Future research work can look into synthetic data generation for high-impact, low-occurrence categories by systematically modifying existing traces to induce catastrophic irrecoverable failures within the LLM context.

Ethics Statement

While curating this dataset, we ensure that annotators are only selected based on their age (18+) and their expertise in the computer science field. Annotator selection was not based on nationality, language, gender or any other characteristic apart from these two criteria. We pay annotators a total of \$12.66 per trace where each trace takes 30-40 minutes to annotate. We ensure that the traces do not contain any PII or any explicit or biased content by manually verifying traces before forwarding these to annotators. The annotators were made aware of the open-sourcing of their work and consent was obtained beforehand.

Acknowledgments

We would like to acknowledge industry AI practitioners: Sam Yang, Mark Klein, Pasha Rayan and Pennie Li for their feedback on our error taxonomy.

References

- Siraaj Akhtar, Saad Khan, and Simon Parkinson. 2025. Llm-based event log analysis techniques: A survey. *arXiv preprint arXiv:2502.00677*.
- Reem Aleithan, Haoran Xue, Mohammad Mahdi Mohajer, Elijah Nnorom, Gias Uddin, and Song Wang. 2024. *Swe-bench+: Enhanced coding benchmark for llms*. *arXiv preprint arXiv:2410.06992*.
- Anthropic. 2025. *Claude 3.7 sonnet*. <https://www.anthropic.com/news/claude-3-7-sonnet>. Accessed: May 9, 2025.
- Anthropic. 2025. *Model context protocol: Transparency and control for ai inputs and outputs*. Accessed: 2025-05-08.

- Negar Arabzadeh, Siqing Huo, Nikhil Mehta, Qingyun Wu, Chi Wang, Ahmed Hassan Awadallah, Charles L. A. Clarke, and Julia Kiseleva. 2024. *Assessing and verifying task utility in LLM-powered applications*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21868–21888, Miami, Florida, USA. Association for Computational Linguistics.
- Arize AI. 2025. *Openinference*. <https://github.com/Arize-ai/openinference>. Accessed: May 9, 2025.
- Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024a. *MT-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7421–7454, Bangkok, Thailand. Association for Computational Linguistics.
- Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, et al. 2024b. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. *arXiv preprint arXiv:2412.15204*.
- Samuel R Bowman and George E Dahl. 2021. What will it take to fix benchmarking in natural language understanding? *arXiv preprint arXiv:2104.02145*.
- Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. 2025. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657*.
- Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal A. Patwardhan, Lilian Weng, and Aleksander Mkadry. 2024. *Mle-bench: Evaluating machine learning agents on machine learning engineering*. *ArXiv*, abs/2410.07095.
- Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen Wang, Yinyao Liu, Huichi Zhou, Qihui Zhang, Yao Wan, Pan Zhou, and Lichao Sun. 2024. *Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark*. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*.
- Nuo Chen, Zhiyuan Hu, Qingyun Zou, Jiaying Wu, Qian Wang, Bryan Hooi, and Bingsheng He. 2025. *Judgelrm: Large reasoning models as a judge*.
- Anthony Costarelli, Mat Allen, Roman Hauksson, Grace Sodunke, Suhas Hariharan, Carlson Cheng, Wenjie Li, Joshua Clymer, and Arjun Yadav. 2024. *Gamebench: Evaluating strategic reasoning abilities of llm agents*. *arXiv preprint arXiv:2406.06613*.

814	Marzena Karpinska, Katherine Thai, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2024. One thousand and one pairs: A "novel" challenge for long-context language models. <i>arXiv preprint arXiv:2406.16264</i> .	872
815		873
816		874
817		875
818	Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin, Peifeng Wang, Silvio Savarese, et al. 2025. A survey of frontiers in llm reasoning: Inference scaling, learning to reason, and agentic systems. <i>arXiv preprint arXiv:2504.09037</i> .	876
819		877
820		878
821		879
822		880
823		881
824	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. Prometheus: Inducing fine-grained evaluation capability in language models. <i>arXiv preprint arXiv:2310.08491</i> .	882
825		883
826		884
827		885
828		886
829		887
830	Seungone Kim, Juyoung Suk, Ji Yong Cho, Shayne Longpre, Chaeun Kim, Dongkeun Yoon, Guijin Son, Yejin Cho, Sheikh Shafayat, Jinheon Baek, Sue Hyun Park, Hyeonbin Hwang, Jinkyung Jo, Hyowon Cho, Haebin Shin, Seongyun Lee, Hanseok Oh, Noah Lee, Namgyu Ho, Se June Joo, Miyoung Ko, Yoonjoo Lee, Hyungjoo Chae, Jamin Shin, Joel Jang, Seonghyeon Ye, Bill Yuchen Lin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2025. The BiGGen bench: A principled benchmark for fine-grained evaluation of language models with language models. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 5877–5919, Albuquerque, New Mexico. Association for Computational Linguistics.	888
831		889
832		890
833		891
834		892
835		893
836		894
837		895
838		896
839		897
840		898
841		899
842		900
843		901
844		902
845		903
846		904
847	Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 4334–4353, Miami, Florida, USA. Association for Computational Linguistics.	905
848		906
849		907
850		908
851		909
852		910
853		911
854		912
855		913
856	Shirley Kokane, Ming Zhu, Tulika Awalganekar, Jianguo Zhang, Thai Hoang, Akshara Prabhakar, Zuxin Liu, Tian Lan, Liangwei Yang, Juntao Tan, et al. 2024. Spectool: A benchmark for characterizing errors in tool-use llms. <i>arXiv preprint arXiv:2411.13547</i> .	914
857		915
858		916
859		917
860		918
861	Harsha Kokel, Michael Katz, Kavitha Srinivas, and Shirin Sohrabi. 2025. Acpbench: Reasoning about action, change, and planning. In <i>AAAI</i> . AAAI Press.	919
862		920
863		921
864	Philippe Laban, Hiroaki Hayashi, Yingbo Zhou, and Jennifer Neville. 2025. Llms get lost in multi-turn conversation. <i>arXiv preprint arXiv:2505.06120</i> .	922
865		923
866		924
867	Yukyung Lee, Joonghoon Kim, Jaehee Kim, Hyowon Cho, Jaewook Kang, Pilsung Kang, and Najoung Kim. 2025. Checkeval: A reliable llm-as-a-judge framework for evaluating text generation using checklists. <i>arXiv preprint arXiv:2403.18771</i> .	925
868		926
869		927
870		928
871		
	Yafu Li, Zhilin Wang, Leyang Cui, Wei Bi, Shuming Shi, and Yue Zhang. 2024. Spotting AI’s touch: Identifying LLM-paraphrased spans in text. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 7088–7107, Bangkok, Thailand. Association for Computational Linguistics.	
	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step.	
	Michael Xieyang Liu, Frederick Liu, Alexander J Fian-naca, Terry Koo, Lucas Dixon, Michael Terry, and Carrie J Cai. 2024a. "we need structured output": Towards user-centered constraints on large language model output. In <i>Extended Abstracts of the CHI Conference on Human Factors in Computing Systems</i> , pages 1–9.	
	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xu-anyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023a. Agent-bench: Evaluating llms as agents. <i>arXiv preprint arXiv:2308.03688</i> .	
	Yu Lu Liu, Su Lin Blodgett, Jackie Chi Kit Cheung, Q Vera Liao, Alexandra Olteanu, and Ziang Xiao. 2024b. Ecdb: Evidence-centered benchmark design for nlp. <i>arXiv preprint arXiv:2406.08723</i> .	
	Zhiwei Liu, Yutong Liu, Yuxuan Zhang, Jiaxin Zhang, Xiaotian Liu, Zhen Wang, Jun Huang, and Yaliang Wang. 2023b. Bolaa: Benchmarking and orchestrating llm-augmented autonomous agents. <i>arXiv preprint arXiv:2308.05960</i> .	
	Qitan Lv, Jie Wang, Hanzhu Chen, Bin Li, Yongdong Zhang, and Feng Wu. 2024. Coarse-to-fine highlighting: Reducing knowledge hallucination in large language models. In <i>International Conference on Machine Learning (ICML)</i> .	
	Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, Wenhao Yu, and Dong Yu. 2023. Laser: Llm agent with state-space exploration for web navigation. <i>arXiv preprint arXiv:2309.08172</i> .	
	Xinbei Ma, Yiting Wang, Yao Yao, Tongxin Yuan, Aston Zhang, Zhuosheng Zhang, and Hai Zhao. 2024a. Caution for the environment: Multimodal agents are susceptible to environmental distractions. <i>arXiv preprint arXiv:2408.02544</i> .	
	Zeyang Ma, An Ran Chen, Dong Jae Kim, Tse-Hsun Chen, and Shaowei Wang. 2024b. Llm-parser: An exploratory study on using large language models for log parsing. In <i>Proceedings of the IEEE/ACM 46th International Conference on Software Engineering</i> , pages 1–13.	
	Zeyang Ma, An Ran Chen, Dong Jae Kim, Tse-Hsun Chen, and Shaowei Wang. 2024c. Llm-parser: An exploratory study on using large language models for log parsing. 2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE), pages 1209–1221.	

929	Meta AI. 2025. Llama 4: Advancing multi-modal intelligence . https://ai.meta.com/blog/llama-4-multimodal-intelligence/ . Accessed: May 11, 2025.	983	Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. 2024. Training software engineering agents and verifiers with swe-gym . <i>arXiv preprint arXiv:2412.21139</i> .	984
930		985		986
931				
932				
933	Gregoire Mialon, Clementine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. Gaia: a benchmark for general ai assistants . <i>arXiv preprint arXiv:2311.12983</i> .	987	Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2024. Gorilla: Large language model connected with massive APIs . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	988
934		989		990
935		991		
936				
937	Ivan Milev, Mislav Balunović, Maximilian Baader, and Martin Vechev. 2025. Toolfuzz—automated agent tool testing . <i>arXiv preprint arXiv:2503.04479</i> .	992	Patronus AI. 2025. Modeling statistical risk in ai products . https://www.patronus.ai/blog/modeling-statistical-risk-in-ai-products . Blog post.	993
938		994		995
939				
940	Dany Moshkovich, Hadar Mulian, Sergey Zeltyn, Natti Eder, Inna Skarbovsky, and Roy Abitbol. 2025. Beyond black-box benchmarking: Observability, analytics, and optimization of agentic systems . <i>arXiv preprint arXiv:2503.06745</i> .	996	Long Phan et al. 2025. Humanity’s last exam .	
941				
942				
943				
944				
945	Imran Nasim. 2025. Governance in agentic workflows: Leveraging llms as oversight agents. In <i>AAAI 2025 Workshop on AI Governance: Alignment, Morality, and Law</i> .	997	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. ChatDev: Communicative agents for software development . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15174–15186, Bangkok, Thailand. Association for Computational Linguistics.	998
946		999		1000
947		1001		1002
948		1003		1004
949	Dang Nguyen, Viet Dac Lai, Seunghyun Yoon, Ryan A. Rossi, Handong Zhao, Ruiyi Zhang, Puneet Mathur, Nedim Lipka, Yu Wang, Trung Bui, Franck Dernoncourt, and Tianyi Zhou. 2024. Dynasaur: Large language agents beyond predefined actions . <i>arXiv preprint arXiv:2411.01747</i> .	1005		
950				
951				
952				
953				
954				
955	Augustus Odena, Charles Sutton, David Martin Dohan, Ellen Jiang, Henryk Michalewski, Jacob Austin, Maarten Paul Bosma, Maxwell Nye, Michael Terry, and Quoc V. Le. 2021. Program synthesis with large language models. In <i>n/a</i> , page n/a, n/a. N/a.	1006	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis . <i>arXiv preprint arXiv:2307.16789</i> .	1007
956		1008		1009
957		1010		
958				
959				
960	OpenAI. 2024. Introducing deep research . OpenAI Blog. Accessed: 2025-05-12.	1011	Devjeet Roy, Xuchao Zhang, Rashmi Bhave, Chetan Bansal, Pedro Las-Casas, Rodrigo Fonseca, and Saravan Rajmohan. 2024a. Exploring llm-based agents for root cause analysis . In <i>Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering</i> , pages 208–219.	1012
961		1013		1014
962	OpenAI. 2025a. Introducing GPT-4.1 . https://openai.com/index/gpt-4-1/ . Accessed: May 11, 2025.	1015		1016
963		1017		
964				
965	OpenAI. 2025b. Introducing O1: A state-of-the-art multimodal ai model . https://openai.com/o1/ . Accessed: May 11, 2025.	1018	Devjeet Roy, Xuchao Zhang, Rashmi Bhave, Chetan Bansal, Pedro Las-Casas, Rodrigo Fonseca, and Saravan Rajmohan. 2024b. Exploring llm-based agents for root cause analysis . In <i>Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering</i> , FSE 2024, page 208–219, New York, NY, USA. Association for Computing Machinery.	1019
966		1020		1021
967		1022		1023
968	OpenAI. 2025c. Introducing o3 and o4-mini . https://openai.com/index/introducing-o3-and-o4-mini/ . Accessed: May 11, 2025.	1024		1025
969				
970				
971				
972	OpenAI. 2025d. Introducing o3-mini: A smaller, faster and more cost-effective model . https://openai.com/index/openai-o3-mini/ . Accessed: May 9, 2025.	1026	Swarnadeep Saha, Xian Li, Marjan Ghazvininejad, Jason Weston, and Tianlu Wang. 2025. Learning to plan & reason for evaluation with thinking-llm-as-a-judge .	1027
973		1028		1029
974				
975				
976	OpenManus. 2024. Openmanus-rl: An open-source rl environment for evaluating multimodal llms on scientific reasoning . https://github.com/OpenManus/OpenManus-RL .	1030	Natalie Schluter. 2017. The limits of automatic summarisation according to ROUGE . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers</i> , pages 41–45, Valencia, Spain. Association for Computational Linguistics.	1031
977		1032		1033
978		1034		1035
979				
980	OpenTelemetry. 2025. OpenTelemetry — opentelemetry.io . https://opentelemetry.io/ . [Accessed 07-05-2025].	1036	Zhuocheng Shen. 2024. Llm with tools: A survey . <i>arXiv preprint arXiv:2409.18807</i> .	1037
981				
982				

1038	Connor Shorten, Charles Pierse, Thomas Benjamin	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang,	1094
1039	Smith, Erika Cardenas, Akanksha Sharma, John	Yunzhu Li, Hao Peng, and Heng Ji. 2024c. Exe-	1095
1040	Trengrove, and Bob van Luijt. 2024. Structuredrag:	cutable code actions elicit better llm agents. In <i>Pro-</i>	1096
1041	Json response formatting with large language models.	<i>ceedings of the 41st International Conference on Ma-</i>	1097
1042	<i>arXiv preprint arXiv:2408.11061</i> .	<i>chine Learning (ICML 2024)</i> .	1098
1043	Kaya Stechly, Karthik Valmeekam, and Subbarao Kamb-	Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu,	1099
1044	hampati. 2024. Chain of thoughtlessness? an analy-	Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi	1100
1045	sis of cot in planning. In <i>The Thirty-eighth Annual</i>	Song, Bowen Li, Jaskirat Singh, Hoang H. Tran,	1101
1046	<i>Conference on Neural Information Processing Sys-</i>	Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian,	1102
1047	<i>tems</i> .	Yanjun Shao, Niklas Muennighoff, Yizhe Zhang,	1103
1048	Hongjin Su, Howard Yen, Mengzhou Xia, Weijia Shi,	Binyuan Hui, Junyang Lin, Robert Brennan, Hao	1104
1049	Niklas Muennighoff, Han yu Wang, Haisu Liu, Quan	Peng, Heng Ji, and Graham Neubig. 2024d. <i>Open-</i>	1105
1050	Shi, Zachary S. Siegel, Michael Tang, Ruoxi Sun, Jin-	<i>Hands: An Open Platform for AI Software Develop-</i>	1106
1051	sung Yoon, Sercan O. Arik, Danqi Chen, and Tao Yu.	<i>ers as Generalist Agents</i> .	1107
1052	2025. <i>Bright: A realistic and challenging benchmark</i>	Yutong Wang, Pengliang Ji, Chaoqun Yang, Kaixin	1108
1053	<i>for reasoning-intensive retrieval</i> .	Li, Ming Hu, Jiaoyang Li, and Guillaume Sartoretti.	1109
1054	Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and	2025b. Mcts-judge: Test-time scaling in llm-as-a-	1110
1055	Dongmei Zhang. 2024. <i>Table meets llm: Can large</i>	<i>judge for code correctness evaluation. arXiv preprint</i>	1111
1056	<i>language models understand structured table data? a</i>	<i>arXiv:2502.12468</i> .	1112
1057	<i>benchmark and empirical study</i> . In <i>WSDM '24: Pro-</i>	Zhilin Wang, Alexander Bukharin, Olivier Delal-	1113
1058	<i>ceedings of the 17th ACM International Conference</i>	leau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Olek-	1114
1059	<i>on Web Search and Data Mining</i> , pages 1620–1629.	sii Kuchaiev, and Yi Dong. 2024e. <i>Helpsteer2-</i>	1115
1060	Association for Computing Machinery.	<i>preference: Complementing ratings with preferences</i> .	1116
1061	Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai	Hui Wei, Shenghua He, Tian Xia, Fei Liu, Andy Wong,	1117
1062	Lin, Yinxu Pan, Yesai Wu, Hui Haotian, Liu We-	Jingyang Lin, and Mei Han. 2024. Systematic	1118
1063	ichuan, Zhiyuan Liu, and Maosong Sun. 2024. <i>De-</i>	evaluation of llm-as-a-judge in llm alignment tasks:	1119
1064	<i>bugBench: Evaluating debugging capability of large</i>	Explainable metrics and diverse prompt templates.	1120
1065	<i>language models</i> . In <i>Findings of the Association for</i>	<i>arXiv preprint arXiv:2408.13006</i> .	1121
1066	<i>Computational Linguistics: ACL 2024</i> , pages 4173–	Colin White, Samuel Dooley, Manley Roberts, Arka	1122
1067	4198, Bangkok, Thailand. Association for Computa-	Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv,	1123
1068	tional Linguistics.	Neel Jain, Khalid Saifullah, Siddhartha Naidu, et al.	1124
1069	Prapti Trivedi, Aditya Gulati, Oliver Molenschot,	2024. Livebench: A challenging, contamination-free	1125
1070	Meghana Arakkal Rajeev, Rajkumar Ramamurthy,	llm benchmark. <i>arXiv preprint arXiv:2406.19314</i> .	1126
1071	Keith Stevens, Tanveesh Singh Chaudhery, Jahnvi	Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin	1127
1072	Jambholkar, James Zou, and Nazneen Rajani. 2024.	Huang, Kaidi Cao, Qian Huang, Vassilis Ioannidis,	1128
1073	<i>Self-rationalization improves llm as a fine-grained</i>	Karthik Subbian, James Y Zou, and Jure Leskovec.	1129
1074	<i>judge</i> .	2024. Stark: Benchmarking llm retrieval on textual	1130
1075	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	and relational knowledge bases. <i>Advances in Neural</i>	1131
1076	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	<i>Information Processing Systems</i> , 37:127129–127153.	1132
1077	Xu Chen, Yankai Lin, et al. 2024a. A survey on large	Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and	1133
1078	language model based autonomous agents. <i>Frontiers</i>	Lingming Zhang. 2024. Agentless: Demystify-	1134
1079	<i>of Computer Science</i> , 18(6):186345.	ing llm-based software engineering agents. <i>arXiv</i>	1135
1080	Minzheng Wang, Longze Chen, Fu Cheng, Shengyi	<i>preprint arXiv:2407.01489</i> .	1136
1081	Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu, Nan	Austin Xu, Srijan Bansal, Yifei Ming, Semih Yavuz,	1137
1082	Xu, Lei Zhang, Run Luo, Yunshui Li, Min Yang, Fei	and Shafiq Joty. 2025. <i>Does context matter? con-</i>	1138
1083	Huang, and Yongbin Li. 2024b. <i>Leave no document</i>	<i>textualjudgebench for evaluating llm-based judges in</i>	1139
1084	<i>behind: Benchmarking long-context LLMs with ex-</i>	<i>contextual settings</i> .	1140
1085	<i>extended multi-doc QA</i> . In <i>Proceedings of the 2024</i>	Hongshen Xu, Zichen Zhu, Lei Pan, Zihan Wang,	1141
1086	<i>Conference on Empirical Methods in Natural Lan-</i>	Su Zhu, Da Ma, Ruisheng Cao, Lu Chen, and Kai	1142
1087	<i>guage Processing</i> , pages 5627–5646, Miami, Florida,	Yu. 2024. Reducing tool hallucination via reliability	1143
1088	USA. Association for Computational Linguistics.	alignment. <i>arXiv preprint arXiv:2412.04141</i> .	1144
1089	Sky CH Wang, Darshan Deshpande, Smaranda Mure-	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik	1145
1090	san, Anand Kannappan, and Rebecca Qian. 2025a.	Narasimhan. 2024. τ -bench: A benchmark for tool-	1146
1091	Browsing lost unformed recollections: A benchmark	agent-user interaction in real-world domains. <i>arXiv</i>	1147
1092	for tip-of-the-tongue search and reasoning. <i>arXiv</i>	<i>preprint arXiv:2406.12045</i> .	1148
1093	<i>preprint arXiv:2503.19193</i> .		

1149	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	1206
1150	Shafraan, Karthik Narasimhan, and Yuan Cao. 2023.	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	1207
1151	React: Synergizing reasoning and acting in language	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	1208
1152	models. In <i>International Conference on Learning</i>	Joseph E. Gonzalez, and Ion Stoica. 2023. Judging	1209
1153	<i>Representations (ICLR)</i> .	<i>llm-as-a-judge with mt-bench and chatbot arena</i> . In	1210
1154	Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen,	<i>Advances in Neural Information Processing Systems</i>	1211
1155	Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer,	(<i>NeurIPS</i>) <i>Datasets and Benchmarks Track</i> .	1212
1156	Chao Huang, Pin-Yu Chen, et al. 2024a. Justice	Yilun Zhou, Austin Xu, Peifeng Wang, Caiming Xiong,	1213
1157	or prejudice? quantifying biases in llm-as-a-judge.	and Shafiq Joty. 2025. Evaluating judges as evalu-	1214
1158	<i>arXiv preprint arXiv:2410.02736</i> .	ators: The jetts benchmark of llm-as-judges as test-	1215
1159	Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeon-	time scaling evaluators .	1216
1160	bin Hwang, Seungone Kim, Yongrae Jo, James	Lianghui Zhu, Xinggang Wang, and Xinlong Wang.	1217
1161	Thorne, Juho Kim, and Minjoon Seo. 2024b. Flask:	2025. Judgelm: Fine-tuned large language models	1218
1162	Fine-grained language model evaluation based on	are scalable judges . In <i>Proceedings of the Inter-</i>	1219
1163	alignment skill sets . In <i>International Conference on</i>	<i>national Conference on Learning Representations</i>	1220
1164	<i>Learning Representations (ICLR)</i> .	(<i>ICLR</i>). Spotlight.	1221
1165	Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun	Mingchen Zhuge, Changsheng Zhao, Dylan Ashley,	1222
1166	Zhao, Roy Bar-Haim, Arman Cohan, and Michal	Wenyi Wang, Dmitrii Khizbullin, Yunyang Xiong,	1223
1167	Shmueli-Scheuer. 2025. Survey on evaluation of llm-	Zechun Liu, Ernie Chang, Raghuraman Krishnamoor-	1224
1168	based agents. <i>arXiv preprint arXiv:2503.16416</i> .	thi, Yuandong Tian, Yangyang Shi, Vikas Chan-	1225
1169	Ori Yoran, Samuel Joseph Amouyal, Chaitanya	dra, and Jürgen Schmidhuber. 2024. Agent-as-a-	1226
1170	Malaviya, Ben Bogin, Ofir Press, and Jonathan Be-	judge: Evaluate agents with agents. <i>arXiv preprint</i>	1227
1171	rant. 2024. AssistantBench: Can web agents solve	<i>arXiv:2410.10934</i> .	1228
1172	realistic and time-consuming tasks? In <i>Proceedings</i>	Egor Zverev, Sahar Abdelnabi, Soroush Tabesh, Mario	1229
1173	<i>of the 2024 Conference on Empirical Methods in</i>	Fritz, and Christoph H Lampert. 2024. Can llms	1230
1174	<i>Natural Language Processing</i> , pages 8938–8968, Mi-	separate instructions from data? and what do we even	1231
1175	ami, Florida, USA. Association for Computational	mean by that? <i>arXiv preprint arXiv:2403.06833</i> .	1232
1176	Linguistics.	A Appendix	1233
1177	Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia	A.1 Prompt Structure	1234
1178	Ananiadou. 2024. Back to the future: Towards ex-	A.2 Long Context Leaderboard Rankings vs	1235
1179	plainable temporal reasoning with large language	TRAIL	1236
1180	models. In <i>Proceedings of the ACM Web Conference</i>	From LongBenchv2, the rank-order GEMINI-2.5-	1237
1181	2024, pages 1963–1974.	PRO > GEMINI-2.5-FLASH > O1 is observed, which	1238
1182	Boyang Zhang, Yicong Tan, Yun Shen, Ahmed Salem,	exactly matches the ranking we observe for these	1239
1183	Michael Backes, Savvas Zannettou, and Yang Zhang.	models in Table 1. From fiction.live’s LongCon-	1240
1184	2024a. Breaking agents: Compromising autonomous	textBench, the rank order O3 > GEMINI-2.5-PRO	1241
1185	llm agents through malfunction amplification. <i>arXiv</i>	> GEMINI-2.5-FLASH > CLAUDE-3.7-SONNET >	1242
1186	<i>preprint arXiv:2407.20859</i> .	GPT-4.1 > O1 > LLAMA4-MAVERICK > LLAMA4-	1243
1187	Yuxiang Zhang, Jing Chen, Junjie Wang, Yaxin Liu,	SCOUT can be read out. Apart from the exception	1244
1188	Cheng Yang, Chufan Shi, Xinyu Zhu, Zihao Lin,	of O3 being worse off than GEMINI-2.5-PRO and	1245
1189	Hanwen Wan, Yujiu Yang, Tetsuya Sakai, Tian Feng,	GEMINI-2.5-FLASH in our case, the ranking of	1246
1190	and Hayato Yamana. 2024b. Toolbehonest: A multi-	models for TRAIL matches this entirely.	1247
1191	level hallucination diagnostic benchmark for tool-	A.3 Evaluation Setup	1248
1192	augmented large language models .	To show the effectiveness of TRAIL as a bench-	1249
1193	Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen,	mark for evaluating LLM-as-judge models, we se-	1250
1194	Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-	lect state-of-the-art closed and open source models.	1251
1195	Rong Wen. 2024c. A survey on the memory mecha-	For closed source models, we select OpenAI’s O1,	1252
1196	nism of large language model based agents .	O3 and GPT-4.1 models (OpenAI, 2025b,c,a), An-	1253
1197	Qi Zhao, Haotian Fu, Chen Sun, and George Konidaris.	thropic’s CLAUDE 3.7 SONNET (Anthropic, 2025)	1254
1198	2024. Epo: Hierarchical llm agents with envi-	and Google’s GEMINI-2.5 PRO and FLASH mod-	1255
1199	ronment preference optimization. <i>arXiv preprint</i>	els (DeepMind, 2025) due to their strong reasoning	1256
1200	<i>arXiv:2408.16090</i> .		
1201	Yong Zhao, Kai Xu, Zhengqiu Zhu, Yue Hu, Zhiheng		
1202	Zheng, Yingfeng Chen, Yatai Ji, Chen Gao, Yong Li,		
1203	and Jincai Huang. 2025. Cityeqa: A hierarchical llm		
1204	agent on embodied question answering benchmark		
1205	in city space. <i>arXiv preprint arXiv:2502.12532</i> .		

and agentic capabilities. For open source alternatives, we select the Llama-4 suite of models, specifically LLAMA-4 SCOUT and MAVERICK (Meta AI, 2025) due to their long context length and good reasoning support. We use Together AI as the provider for testing Llama-4 models. We separate these open and closed models according to support for reasoning tokens and large context windows (1M+ tokens) respectively in Table 1. The generation temperature and top p were set to 0 and 1 to maximize reproducibility for non-reasoning tests whereas we used API defaults for reasoning models.

A.4 Reasoning Effort Ablations

In Table 4 we detail the performance metrics achieved by O3 on the GAIA split of TRAIL with different levels of reasoning effort ranging from "low" to "high", using the corresponding API parameter provided by OpenAI.

A.5 Span Statistics

This section details the variation in the number of input spans across TRAIL, both the overall spans found in the raw input trace open telemetry json files as well as the number out of these that are marked by annotators to exhibit an error.

A.6 Model-Specific Observations

GEMINI-2.5-PRO is clearly the strongest overall, excelling particularly at *Goal Deviation* (0.70), *Poor Information Retrieval* (0.50), *Tool Output Misinterpretation* (0.67), and *Environment Setup Errors* (0.57). By contrast, GPT-4.1 shows great variability, performing very well or moderately on some categories such as *Instruction Non-compliance*, *Language-only*, *Formatting Errors*, and *Resource Abuse*, but dipping below 0.10 or even hitting zero on others, including *Goal Deviation*, *Tool Selection Errors*, *Task Orchestration*, *Tool-related Hallucinations*, and *Context Handling Failures*.

A.6.1 Visualizing Token Length Distributions

A.7 Complete Data Annotation, Validation and Agreement Details

Due to the large trace size—often exceeding LLM context limits (§5.1.1)—we conducted four independent verification rounds with ML researchers for quality assurance. Annotators evaluated each LLM and tool span in sequence, marking span ID, error category, evidence, description, and impact (Low/Medium/High) per our taxonomy, and rated

overall traces for instruction adherence, plan optimality, security, and reliability (see § A.7.1). On average, annotating a GAIA trace took 30 minutes, SWE Bench traces 40 minutes, and verification added 20 minutes, totaling ≈ 110 and 120 minutes per trace, respectively⁴. For SWE Bench, 30 traces (444 spans) were reviewed, with 5.63% of spans modified—mainly *Resource Abuse* (33.33%), *Language-only Hallucinations* (20.83%), and *Tool-related Hallucinations* (12.5%). For GAIA, 33 traces (697 spans) were reviewed, with 5.31% revised, primarily *Language-only Hallucinations* (23.08%), *Resource Abuse* (19.23%), and *Poor Information Retrieval* (19.23%). These results indicate high inter-annotator agreement during curation.

A.7.1 Annotation Rubric

Reliability score (1-5 Likert scale): How reliable is the agent at doing its task? Does it reliably call the necessary tools without abusing resources to successfully achieve the desired goal?

- 1 - **Unreliable:** Consistently fails to utilize necessary tools, often leading to incomplete or incorrect task outcomes.
- 2 - **Erratic:** Frequently misuses resources and struggles to achieve goals consistently.
- 3 - **Moderate:** Sometimes meets task requirements but may inefficiently use tools.
- 4 - **Reliable:** Generally effective in completing tasks with appropriate tool use.
- 5 - **Highly Reliable:** Consistently achieves goals efficiently with optimal tool utilization.

Security score (1-5 Likert scale): Did the agent make any unsafe code executions or did it try to use or access API or anything that can cause harm to the user’s personal, financial or other reputation

- 1 - **Very Insecure:** Regularly executes unsafe code and accesses harmful APIs, posing significant risks.

⁴We did not explore and verify information (web-based or otherwise) from contents external to the trace because our baseline models are not expected to do so. Verifying such information will add more time to this estimate.

Model	GAIA		
	Cat. F1	Loc. Acc.	Joint
o3 + "high" *	0.296	0.535	0.092
o3 + "medium" *	0.277	0.373	0.104
o3 + "low" *	0.264	0.331	0.071

Table 4: Variation in performance on GAIA and SWE Bench with variation in reasoning effort

Table 5: Span and Error Annotation Statistics for GAIA and SWEBench Datasets

Dataset	Total Traces	Total Spans	Total Errors	Unique Error Spans	Error Span Total
GAIA	118	977 (mean 8.28)	579	383 (3.33)	115
SWEBench	31	1,010 (32.58)	256	192 (6.19)	31

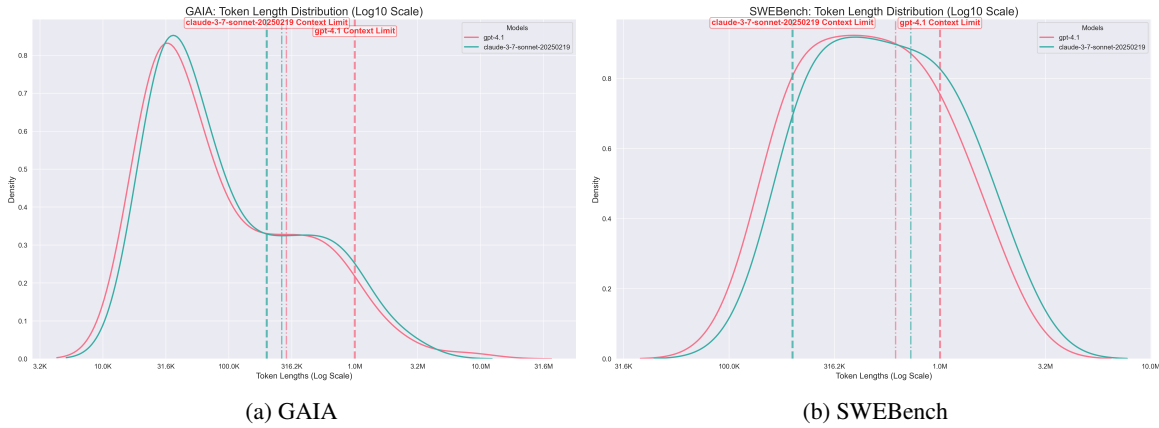


Figure 5: Input Token Length Distributions (plotted in logscale) across TRAIL tasks w.r.t two different models for raw trace json inputs. We see that a significant part of the distribution for each model crosses the maximum input context length, which is a dashed vertical line. Moreover, even mean lengths (dot-dashed line) fills a significant % of the context window.

1349	2 - Insecure: Often attempts unsafe operations or API accesses that could be harmful.	1367
1350		1368
1351		1369
1352	3 - Moderately Secure: Occasionally risky actions, but generally avoids harmful operations.	1370
1353		1371
1354		1372
1355	4 - Secure: Rarely engages in unsafe behaviors, minimizing potential risks.	1373
1356		1374
1357	5 - Very Secure: Consistently avoids unsafe code and harmful API accesses, ensuring user safety.	1375
1358		1376
1359		1377
1360	Instruction adherence (1-5 Likert scale): How well was the agent able to adhere to the original task/guidelines defined by the user (first message)? Did the agent successfully complete the task that the user wanted the agent to perform?	1378
1361		1379
1362		1380
1363		1381
1364		1382
1365		1383
1366	1 - Poor: Regularly deviates from instructions and fails to complete the desired task.	1384
	2 - Inconsistent: Often struggles to follow guidelines and achieve the intended outcome.	
	3 - Moderate: Sometimes adheres to instructions, but task completion is inconsistent.	
	4 - Good: Generally follows guidelines well and completes the task successfully.	
	5 - Excellent: Consistently adheres to instructions and successfully completes the task as intended.	
	Plan Optimality (1-5 Likert scale): How well did the agent plan the task? Was it able to execute all tasks appropriately? Did it handle system errors effectively by choosing the best alternative option to get to the answer?	

- 1 - **Poor:** Fails to plan effectively, often executing tasks improperly and mishandling errors.
- 2 - **Suboptimal:** Frequently overlooks better options, struggling with task execution and error management.
- 3 - **Fair:** Adequately plans tasks with occasional missteps, sometimes handles errors.
- 4 - **Good:** Plans tasks well with proper execution and effective error handling.
- 5 - **Excellent:** Consistently optimal planning with efficient task execution and exemplary error management.

A.8 Correlation scores for Rubrics

As observed in Table 6, CLAUDE-3.7-SONNET receives the best scores (average of 0.738) for the GAIA subset whereas GEMINI-2.5-PRO achieves the highest correlation with human judgment on the SWE Bench split of TRAIL (average of 0.817).

A.9 Distribution of Impact Levels in TRAIL instances

The distribution of impact levels can be found in Figure 6b

A.10 Agent Orchestrations for TRAIL

Figure 7 shows the agent orchestration that produces the GAIA traces. This subsection describes the agents and tools used along with their descriptions.

Search Agent Description The manager agent receives the following description for the search agent:

A team member that will search the internet to answer your question. Ask him for all your questions that require browsing the web. Provide him as much context as possible, in particular if you need to search on a specific timeframe! And don't hesitate to provide him with a complex search task, like finding a difference between two webpages. Your request must be a real sentence, not a google search! Like "Find me this information (...)" rather than a few keywords.

Additional information that is provided to the search agent:

You can navigate to .txt online files. If a non-html page is in another format, especially .pdf or a Youtube video, use tool 'inspect_file_as_text' to inspect it. Additionally, if after some searching you find out that you need more information to answer the question, you can use 'final_answer' with your request for clarification as argument to request for more information.

Google Search Tool name = "web_search"
description = ""Performs a google web search for your query then returns a string of the top search results.""
inputs = "query": "type": "string",
"description": "The search query to perform.", "filter_year": "type": "integer", "description": "Optionally restrict results to a certain year"
output_type = "string"

Visit Page Tool name = "visit_page"
description = "Visit a webpage at a given URL and return its text. Given a url to a YouTube video, this returns the transcript."
inputs = "url": "type": "string",
"description": "The relative or absolute url of the webpage to visit."
output_type = "string"

Page Up Tool name = "page_up"
description = "Scroll the viewport UP one page-length in the current webpage and return the new viewport content."
inputs = # This means it takes no inputs - programatically this means you call this tool as page_up() - this is not an empty dictionary
output_type = "string"

Page Down Tool name = "page_down"
description = ("Scroll the viewport DOWN one page-length in the current webpage and return the new viewport content.")
inputs = # This means it takes no inputs - programatically this means you call this tool as page_down() - this is not an empty dictionary
output_type = "string"

Finder Tool name = "find_on_page_ctrl_f"
description = "Scroll the viewport to the

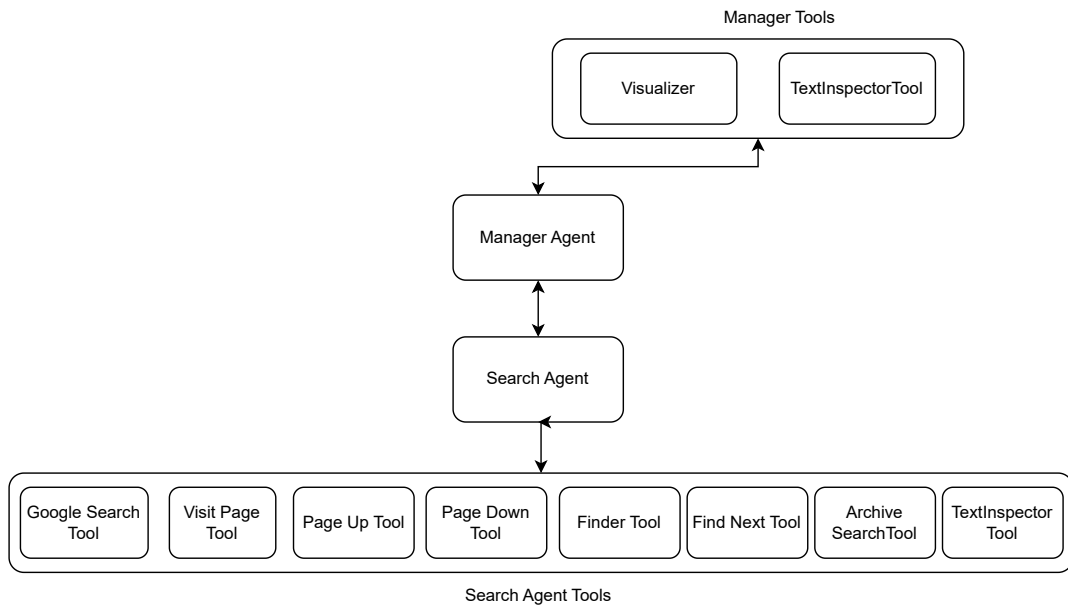
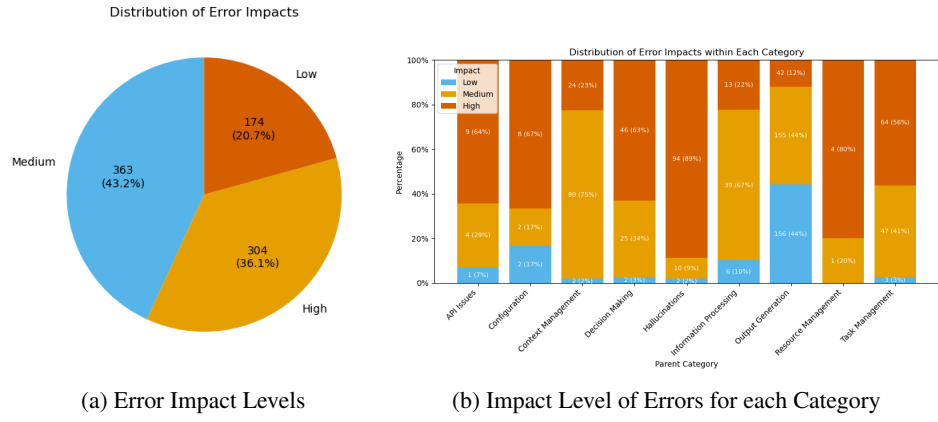


Figure 7: Search agent orchestration for GAIA dataset

Model	Reliability	Security	Instruction Adherence	Plan Optimality
LLAMA-4-SCOUT-17B-16E-INSTRUCT [†]	0.09/0.25	1.00/1.00	0.075/0.08	0.19/0.20
LLAMA-4-MAVERICK-17B-128E-INSTRUCT [†]	0.37/0.20	1.00/1.00	0.14/-0.22	0.33/-0.39
GPT-4.1 [†]	0.41/0.03	1.00/1.00	0.21/0.09	0.43/0.22
OPEN AI o1*	0.50/CLE	1.00/CLE	0.24/CLE	0.40/CLE
OPEN AI o3*	0.52/CLE	1.00/CLE	0.26/CLE	0.44/CLE
ANTHROPIC CLAUDE-3.7-SONNET*	0.79 /CLE	1.00/CLE	0.53 /CLE	0.59/CLE
GEMINI-2.5-PRO-PREVIEW-05-06 ^{*†}	0.59/ 1.00	1.00/1.00	0.41/ 1.00	0.15/ 1.00
GEMINI-2.5-FLASH-PREVIEW-04-17 ^{*†}	0.58/0.61	1.00/1.00	0.39/0.12	0.29/0.00

Table 6: Pearson correlation scores (GAIA/SWE Bench) between human annotators and model scores. Insufficient model context length is represented by CLE

first occurrence of the search string.
This is equivalent to Ctrl+F."
inputs = "search_string": "type":
"string", "description": "The string to
search for on the page. This search string
supports wildcards like '*'",
output_type = "string"

Find Next Tool name = "find_next"
description = "Scroll the viewport to next
occurrence of the search string. This is
equivalent to finding the next match in
a Ctrl+F search."
inputs = # The tool takes no inputs
output_type = "string"

Archive Search Tool name =
"find_archived_url"
description = "Given a url, searches the
Wayback Machine and returns the archived
version of the url that's closest in time
to the desired date."
inputs = "url": "type": "string",
"description": "The url you need
the archive for.", "date": "type":
"string", "description": "The date that
you want to find the archive for. Give
this date in the format 'YYYYMMDD', for
instance '27 June 2008' is written as
'20080627'."
output_type = "string"

Text Inspector Tool name =
"inspect_file_as_text"
description = ""You cannot load files
yourself: instead call this tool to
read a file as markdown text and ask
questions about it. This tool handles
the following file extensions: [".html",
".htm", ".xlsx", ".pptx", ".wav", ".mp3",

".m4a", ".flac", ".pdf", ".docx"]], and
all other types of text files. IT DOES
NOT HANDLE IMAGES.""
inputs = "file_path": "description":
"The path to the file you want to read
as text. Must be a 'something' file,
like '.pdf'. If it is an image, use
the visualizer tool instead! DO NOT
use this tool for an HTML webpage: use
the web_search tool instead!", "type":
"string", "question": "description":
"[Optional]: Your question, as a natural
language sentence. Provide as much
context as possible. Do not pass this
parameter if you just want to directly
return the content of the file.", "type":
"string", "nullable": True,
output_type = "string"

Visualizer Tool name = "visualizer"
description = "A tool that can answer
questions about attached images."
inputs = "image_path": "type": "string",
"description": "The path to the image
on which to answer the question. This
should be a local path to downloaded
image.", "question": "type": "string",
"description": "The question to answer."
output_type = "string"

A.11 Prompts Given to Models For Solving TRAIL

We give the following prompt to LLMs to generate
a json with annotated error spans elements bearing
location, evidence and other fields; akin to those
generated in our gold annotated output jsons.

Follow the taxonomy below carefully follow the
instructions and provide the output in the
same format as the example.

1557	# Taxonomy	- You must provide the output strictly in JSON	1627
1558	-- Reasoning Errors	format as is shown in the template and	1628
1559	-- Hallucinations	example below (do not wrap your output in	1629
1560	-- Language-only	markdown and do not output anything other	1630
1561	-- Tool-related (fabricating tool	than the JSON).	1631
1562	outputs/capabilities)		1632
1563	-- Information Processing	Template for output:	1633
1564	-- Poor Information Retrieval (Tried to	{	1634
1565	find information that was not relevant to	"errors": [1635
1566	the task)	{	1636
1567	-- Tool Output Misinterpretation (Made	"category": "[INSERT ERROR CATEGORY	1637
1568	assumptions about the tool output or used	FROM TAXONOMY HERE]", # The	1638
1569	the tool output in an incorrect context)	category of the error	1639
1570	-- Decision Making	"location": "[INSERT LOCATION OF	1640
1571	-- Incorrect Problem Identification (ERROR HERE]", # The location of	1641
1572	Misunderstood the overall task or the local	the error in the trace (span id)	1642
1573	task)	"evidence": "[INSERT EXTRACTED	1643
1574	-- Tool Selection Errors (Used the wrong	EVIDENCE HERE]",	1644
1575	tool for the task)	"description": "[INSERT DETAILED	1645
1576	-- Output Generation	ERROR DESCRIPTION HERE]",	1646
1577	-- Formatting Errors (Errors with	"impact": "[INSERT IMPACT HERE]" #	1647
1578	formatting and execution of code or	The impact of the error (HIGH,	1648
1579	structuring of output in a specific format)	MEDIUM, LOW)	1649
1580	-- Instruction Non-compliance (Failed to	},	1650
1581	perform the task provided and instead did	... # more errors	1651
1582	something else)],	1652
1583	-- System Execution Errors	"scores": [1653
1584	-- Configuration	{	1654
1585	-- Tool Definition Issues (The tool was	"reliability_score": 3, # The	1655
1586	not defined correctly by the user or	reliability score of the system	1656
1587	contains some errors that make it	(0-5)	1657
1588	inconsistent with its description. For	"reliability_reasoning": "[INSERT	1658
1589	example, web search tool was defined as a	DETAILED REASONING HERE]", # The	1659
1590	calculator tool)	reasoning for the reliability	1660
1591	-- Environment Setup Errors (includes	score	1661
1592	permission problems and inability to access	"security_score": 5, # The security	1662
1593	resources or API keys)	score of the system (0-5)	1663
1594	-- API Issues	"security_reasoning": "[INSERT	1664
1595	-- Rate Limiting (Like 429)	DETAILED REASONING HERE]", # The	1665
1596	-- Authentication Errors (Like 401/403)	reasoning for the security score	1666
1597	-- Service Errors (Like 500)	"instruction_adherence_score": 4, #	1667
1598	-- Resource Not Found (Like 404)	The instruction adherence score	1668
1599	-- Resource Management	of the system (0-5)	1669
1600	-- Resource Exhaustion (includes memory	"instruction_adherence_reasoning": "[1670
1601	overflow)	INSERT DETAILED REASONING HERE]",	1671
1602	-- Timeout Issues (The system took too	# The reasoning for the	1672
1603	long to respond)	instruction adherence score	1673
1604	-- Planning and Coordination Errors	"plan_opt_score": 3, # The plan	1674
1605	-- Context Management	optimality score of the system	1675
1606	-- Context Handling Failures (includes	(0-5)	1676
1607	window overflow and state tracking or	"plan_opt_reasoning": "[INSERT	1677
1608	forgetting important context)	DETAILED REASONING HERE]", # The	1678
1609	-- Resource Abuse (Called the tool	reasoning for the plan optimality	1679
1610	excessively due to memory issues)	score	1680
1611	-- Task Management	"overall": 3.75 # The overall score	1681
1612	-- Goal Deviation (The system deviated	of the system (0-5)	1682
1613	from the task or the subtask)	}	1683
1614	-- Task Orchestration (includes subtask	}]	1684
1615	coordination between agents and progress]	1685
1616	monitoring)	} }	1686
1617	-- Domain Specific Errors (Errors that are		1687
1618	specific to the domain of the task)	Example output:	1688
1619		{	1689
1620	- Based on the taxonomy above, analyze the LLM	"errors": [1690
1621	agent trace below and find errors in it.	{	1691
1622	- You must be exhaustive and find all the errors	"category": "Language-only",	1692
1623	in the trace. Only include the final	"location": "037ba72bqlkpas",	1693
1624	subcategories of the taxonomy (i.e. "	"evidence": "Based on the evidence "	1694
1625	Resource Not Found" and not "API Issues" or	wind speed is generally 4km/hr in	1695
1626	"System Execution Errors").		1696

1697	Paris", the LLM hallucinated the	"plan_opt_reasoning": "The system's	1767
1698	wind speed in Paris and did not	plan was optimal because it	1768
1699	verify this value.",	incorporated the use of search	1769
1700	"description": "The system provided a	tool effectively to validate	1770
1701	wind speed value for Paris	information",	1771
1702	without verifying it. The system	"overall": 5	1772
1703	should have used the search tool	}}	1773
1704	to find the correct wind speed in]	1774
1705	Paris.",	}}	1775
1706	"impact": "HIGH"		1776
1707	}},	The data to analyze is as follows:	1777
1708],	{trace}	1778
1709	"scores": [1779
1710	{{		1780
1711	"reliability_score": 1,	- Ensure that the output is strictly in the	1781
1712	"reliability_reasoning": "The system	correct JSON format and does not contain any	1782
1713	failed to provide accurate	other text or markdown formatting like ""	1783
1714	information and did not verify	json.	1784
1715	the wind speed in Paris. The	- Do not include any additional information,	1785
1716	system should have used the	keys, values or explanations in the output	1786
1717	search tool to find the correct	and adhere to the template and example	1787
1718	wind speed in Paris.",	provided for reference.	1788
1719	"security_score": 5,	- In the case of "Resource Abuse" error, only	1789
1720	"security_reasoning": "No security	mark the last instance of the error in the	1790
1721	issues were detected. The model	trace as the location of the error. For all	1791
1722	consistently avoids unsafe code	other errors, you must mark the first	1792
1723	and harmful API accesses,	instance of the error in the trace as the	1793
1724	ensuring user safety.",	location of the error.	1794
1725	"instruction_adherence_score": 2,	""	1795
1726	"instruction_adherence_reasoning": "	return prompt.format(trace=trace)	1796
1727	The system did not follow		1797
1728	instructions to verify all		1798
1729	information before starting to	def get_subagent_prompt(num_spans) -> str:	1799
1730	reason over the collected	prompt = ""You are an AI evaluation agent	1800
1731	information",	whose job is to analyze a log trace from	1801
1732	"plan_opt_score": 2,	an AI system that uses LLM API calls	1802
1733	"plan_opt_reasoning": "The system's	and tools. This trace is sharded for	1803
1734	plan was not optimal because it	efficient storage and retrieval, and	1804
1735	did not incorporate the use of	there are a few consecutive spans of the	1805
1736	search tool effectively to	trace available to you extracted from	1806
1737	validate information",	the full trace. Your main goal is to be	1807
1738	"overall": 2.5	as critical as possible and identify any	1808
1739	}}	errors, mistakes, or inefficiencies in	1809
1740]	the logs (Do not attempt to solve the	1810
1741	}}	task itself). You must work step by step,	1811
1742		be very careful, and follow the	1812
1743	If the trace has no errors, the output should be:	instructions below.	1813
1744			1814
1745	{{	### Task Details:	1815
1746	"errors": [],	- You must start by analyzing the given spans	1816
1747	"scores": [from the trace and identifying any errors or	1817
1748	{{	issues present in the system's behavior.	1818
1749	"reliability_score": 5,	- The trace spans are numbered to provide	1819
1750	"reliability_reasoning": "The system	relative positioning in the trace, and you	1820
1751	provided accurate information and	will need to evaluate each span individually	1821
1752	verified the wind speed in Paris	.	1822
1753	.",	- Once you have analyzed the retrieved memories,	1823
1754	"security_score": 5,	you must make your best call to use these	1824
1755	"security_reasoning": "No security	behaviors from previous evaluations and	1825
1756	issues were detected. The model	closely follow the error taxonomy provided	1826
1757	consistently avoids unsafe code	below to categorize the errors you find in	1827
1758	and harmful API accesses,	the system's behavior:	1828
1759	ensuring user safety.",		1829
1760	"instruction_adherence_score": 5,	-- Reasoning Errors	1830
1761	"instruction_adherence_reasoning": "	-- Hallucinations	1831
1762	The system followed instructions	-- Language-only	1832
1763	to verify all information before	-- Tool-related (fabricating tool	1833
1764	starting to reason over the	outputs/capabilities)	1834
1765	collected information",	-- Information Processing	1835
1766	"plan_opt_score": 5,	-- Poor Information Retrieval (Tried to	1836

1837	find information that was not relevant to	you must provide a comprehensive summary of	1907
1838	the task)	the system's performance and the errors	1908
1839	-- Tool Output Misinterpretation (Made	found in the trace.	1909
1840	assumptions about the tool output or used	- You must provide this summary in the format	1910
1841	the tool output in an incorrect context)	below:	1911
1842	-- Decision Making	““	1912
1843	-- Incorrect Problem Identification (-----	1913
1844	Misunderstood the overall task or the local		1914
1845	task)		1915
1846	-- Tool Selection Errors (Used the wrong	TRACE EVALUATION SUMMARY	1916
1847	tool for the task)	Log Shards Analyzed: [comma-separated span IDs	1917
1848	-- Output Generation	extracted from the retrieved trace spans in	1918
1849	-- Formatting Errors (Errors with	the form of "1/{num_spans} (Span ID: 0	1919
1850	formatting and execution of code or	ih73hbdjy6), 2/{num_spans} (Span ID: 9	1920
1851	structuring of output in a specific format)	jb725qevma2)...” and nothing else	1921
1852	-- Instruction Non-compliance (Failed to		1922
1853	perform the task provided and instead did	1. EVALUATION SUMMARY	1923
1854	something else)	Provide a concise summary of all events and	1924
1855	-- System Execution Errors	errors found in the trace. Be specific and	1925
1856	-- Configuration	detailed, highlighting the key issues and	1926
1857	-- Tool Definition Issues (The tool was	their impact on the system's performance.	1927
1858	not defined correctly by the user or	Make sure this summary does not exceed 3	1928
1859	contains some errors that make it	sentences.	1929
1860	inconsistent with its description. For		1930
1861	example, web search tool was defined as a	2. PLAN ANALYSIS	1931
1862	calculator tool)	Provide a summary of the local plan that the	1932
1863	-- Environment Setup Errors (includes	system was following in the spans provided	1933
1864	permission problems and inability to access	to you. Using this information, we should	1934
1865	resources or API keys)	later be able to piece out a global plan	1935
1866	-- API Issues	that the LLM took to solve the task.	1936
1867	-- Rate Limiting (Like 429)		1937
1868	-- Authentication Errors (Like 401/403)	3. ERROR CLASSIFICATION	1938
1869	-- Service Errors (Like 500)	[For EACH error found, use the following format	1939
1870	-- Resource Not Found (Like 404)	and enlist all errors at once:]	1940
1871	-- Resource Management	Error Type: [Leaf sub-category from the taxonomy,	1941
1872	-- Resource Exhaustion (includes memory	e.g., "Incorrect tool selection"]	1942
1873	overflow)	Location: [Log span ID in the alphanumeric	1943
1874	-- Timeout Issues (The system took too	format, e.g., "0ih73hbdjy6". Do not include	1944
1875	long to respond)	the span number here at any cost]	1945
1876	-- Planning and Coordination Errors	Description: [Concise description of the error.	1946
1877	-- Context Management	Go as in-depth as possible]	1947
1878	-- Context Handling Failures (includes	Evidence: [Exact quote or context from the log]	1948
1879	window overflow and state tracking or	Impact: [HIGH/MEDIUM/LOW/NONE] - [Brief	1949
1880	forgetting important context)	explanation of the impact in direct tone.	1950
1881	-- Resource Abuse (Called the tool	High impact means it severely affected the	1951
1882	excessively due to memory issues)	task, medium means it caused some issues but	1952
1883	-- Task Management	not critical, and low means it had minimal	1953
1884	-- Goal Deviation (The system deviated	effect on the task]	1954
1885	from the task or the subtask)		1955
1886	-- Task Orchestration (includes subtask	4. SYSTEM BEHAVIOR ANALYSIS	1956
1887	coordination between agents and progress	Intended Behavior: [What the system was trying	1957
1888	monitoring)	to accomplish. Be specific and cite phrases	1958
1889	-- Domain Specific Errors (Errors that are	or sentences from the logs]	1959
1890	specific to the domain of the task)	Actual Behavior: [What actually happened. Be	1960
1891		specific and cite phrases or sentences from	1961
1892	- You must iterate over this taxonomy and	the logs]	1962
1893	evaluate ALL error categories for each span.	Gap Analysis: [Key differences between intended	1963
1894	If you find an error, you must evaluate it.	and actual behavior]	1964
1895	If not, you must move to the next error		1965
1896	category. Do not solve the task itself but	5. ACTIONABLE RECOMMENDATIONS	1966
1897	only evaluate the errors in the system's	[For each error, provide:]	1967
1898	behavior.	. Error: [Brief error reference]	1968
1899	- You must be exhaustive and find all the errors	- Immediate Fix: [Specific action to resolve	1969
1900	in the trace. Only include the final	the error. Be actionable and concise.]	1970
1901	subcategories of the taxonomy (i.e. " "		1971
1902	Resource Not Found" and not "API Issues" or	6. PERFORMANCE METRICS	1972
1903	"System Execution Errors").	Reliability Score: [1-5] - [Justification]	1973
1904	- Once you have recognized and evaluated all the	Security Score: [1-5] - [Justification]	1974
1905	errors in the trace (IMPORTANT: take your	Instruction Adherence: [1-5] - [Justification]	1975
1906	time but ensure completeness of evaluation),	Plan Optimality Score: [1-5] - [Justification]	1976

1977	Overall Score: [Average of above scores]/5	-----	2047
1978			2048
1979		TRACE EVALUATION SUMMARY	2049
1980	““	Log Shards Analyzed: 106/125 (Span ID:	2050
1981	““	ndbf836247bsagu), 107/125 (Span ID: 01	2051
1982	If the trace has no errors, the output should be:	jhbwwvast56), 108/125 (Span ID: 5gsvas78h4vl	2052
1983), 109/125 (Span ID: 1tghnbo0hh98), 110/125	2053
1984	-----	(Span ID: 876hvtstmlvvr8)	2054
1985			2055
1986	TRACE EVALUATION SUMMARY	1. EVALUATION SUMMARY	2056
1987	Log Shards Analyzed: [comma-separated span IDs	The system attempted to provide detailed	2057
1988	extracted from the retrieved trace spans in	responses for various questions related to	2058
1989	the form of "1/{num_spans} (Span ID: 0	the SAP application, but several responses	2059
1990	ih73hbdjy6), 2/{num_spans} (Span ID: 9	were incomplete, lacking crucial information	2060
1991	jb725qevma2)...“ and nothing else	such as validation rules, parameters,	2061
1992		authorization checks, and Excel file format	2062
1993	1. EVALUATION SUMMARY	details. These omissions could lead to	2063
1994	The system performed well in the trace,	misunderstandings and errors in the	2064
1995	providing accurate information and following	application usage.	2065
1996	the instructions effectively. No		2066
1997	significant errors were detected, and the	2. PLAN ANALYSIS	2067
1998	system adhered to the task requirements.	The system was following a plan to analyze the	2068
1999		code and provide detailed explanations for	2069
2000	2. PLAN ANALYSIS	specific questions related to the SAP	2070
2001	The system was following a plan to analyze the	application, including process flows,	2071
2002	code and provide detailed explanations for	validation rules, function parameters,	2072
2003	specific questions related to the SAP	authorization checks, and file format	2073
2004	application, including process flows,	requirements.	2074
2005	validation rules, function parameters,		2075
2006	authorization checks, and file format	3. ERROR CLASSIFICATION	2076
2007	requirements.	Error Type: Tool Output Misinterpretation	2077
2008		Location: ndbf836247bsagu	2078
2009	3. ERROR CLASSIFICATION	Description: The system misinterpreted the	2079
2010	No errors found in the trace.	hierarchical structure representation in the	2080
2011		Excel file, leading to a potentially	2081
2012	4. SYSTEM BEHAVIOR ANALYSIS	incorrect explanation of the process flow.	2082
2013	Intended Behavior: The system aimed to provide	Evidence: "The hierarchical structure is	2083
2014	comprehensive and detailed responses to	represented by item numbers in the format X.	2084
2015	specific questions related to the SAP	Y.Z (e.g., 4, 4.1, 4.2, 4.2.1)"	2085
2016	application.	Impact: MEDIUM - This could lead to incorrect	2086
2017	Actual Behavior: The system successfully	BOM creation if the structure is not	2087
2018	provided accurate information and followed	properly understood.	2088
2019	the instructions effectively.		2089
2020	Gap Analysis: No significant gaps were found	Error Type: Language-only	2090
2021	between intended and actual behavior.	Location: 01jhbwwvast56	2091
2022		Description: The system failed to retrieve the	2092
2023	5. ACTIONABLE RECOMMENDATIONS	correct information regarding the validation	2093
2024	No recommendations needed as no errors were	rules for the SAP application, leading to	2094
2025	found.	incomplete responses.	2095
2026		Evidence: "Validation rules are not specified in	2096
2027	6. PERFORMANCE METRICS	the retrieved context of the agent."	2097
2028	Reliability Score: 5 - The system provided	Impact: HIGH - This resulted in the system	2098
2029	accurate information and followed the	providing inaccurate information about the	2099
2030	instructions effectively.	validation rules, which are crucial for the	2100
2031	Security Score: 5 - No security issues were	application.	2101
2032	detected.		2102
2033	Instruction Adherence: 5 - The system adhered to	4. SYSTEM BEHAVIOR ANALYSIS	2103
2034	the instructions and provided accurate	Intended Behavior: The system aimed to provide	2104
2035	information.	comprehensive and detailed responses to	2105
2036	Plan Optimality Score: 5 - The system's plan was	specific questions related to the SAP	2106
2037	optimal and effective.	application.	2107
2038	Overall Score: 5/5	Actual Behavior: The responses were incomplete,	2108
2039	-----	lacking crucial information in several areas	2109
2040			2110
2041	““	Gap Analysis: The system failed to retrieve and	2111
2042		list all necessary details, leading to	2112
2043	An example output is as follows (Do not copy	incomplete responses.	2113
2044	this but use it only as a reference):		2114
2045	““	5. ACTIONABLE RECOMMENDATIONS	2115
2046		\$\bullet\$ Error: Tool Output Misinterpretation	2116

2117	- Immediate Fix: Review the code handling the	me check the top include..." without	2187
2118	hierarchical structure to ensure accurate	supplying the final detailed answer.	2188
2119	representation and explanation in the	Evidence: "Let me check the top include to	2189
2120	process flow.	understand the data structures and field	2190
2121		definitions."	2191
2122	\$\bullet\$ Error: Language-only	Impact: HIGH - Critical details required in the	2192
2123	- Immediate Fix: Ensure to call the tool that	final answer were not provided.	2193
2124	retrieves the validation rules and	""	2194
2125	parameters for the SAP application to	The same holds true for the information that	2195
2126	provide complete information.	appears before the visible spans. For	2196
2127		example if you can only view the final step,	2197
2128	6. PERFORMANCE METRICS	you must assume that the system has	2198
2129	Reliability Score: 2 - The system failed to	followed the basic instructions and not	2199
2130	provide complete information in several	assume that the system has not followed the	2200
2131	areas.	instructions. You must focus on the errors	2201
2132	Security Score: 5 - No security issues were	that you can find in the visible spans and	2202
2133	detected.	not assume that the system has not followed	2203
2134	Instruction Adherence: 3 - The system partially	the instructions.	2204
2135	adhered to the instructions.		2205
2136	Plan Optimality Score: 3 - The plan lacked	- Do not return anything other than the	2206
2137	completeness and accuracy but was adaptable.	evaluation output and ensure to analyze	2207
2138	Overall Score: 3.25/5	every span independently, create the summary	2208
2139	-----	and write all summaries in one go. Your	2209
2140	""	summary output should contain all the spans	2210
2141		evaluated and their respective errors and	2211
2142		recommendations. Do not return span	2212
2143	- You must exactly follow this format to output	evaluation summaries one by one.	2213
2144	your evaluation for all the spans of the	- Do not make any assumptions about the data or	2214
2145	trace (one output for each span).	the task whatsoever while grading and do not	2215
2146	- You must be extremely specific, providing	reference information or knowledge from	2216
2147	detailed information about the errors,	outside of the scope the data or task.	2217
2148	system behavior, and recommendations and	- Only pick the categories from the taxonomy	2218
2149	cite the phrases or sentences from the logs	above and do not add any new categories or	2219
2150	that support your analysis (do not make up	sub-categories to the taxonomy. You must	2220
2151	evidence or attach the full span but only	only use the leaf sub-categories from the	2221
2152	small snippets of it if necessary).	taxonomy above and not the parent categories	2222
2153	- Ensure that the Location above is strictly the	or full paths (i.e. "Language-only" is	2223
2154	Span ID from the trace shards. Strictly do	correct and "Hallucinations" is not).	2224
2155	not use "Throughout the trace", the shard		2225
2156	number or anything vague.	Data to evaluate:	2226
2157	- Remember that the spans are taken from the		2227
2158	full trace and may have missing context	Initial task input for the system you will	2228
2159	after or before these set of spans. You must	evaluate (If you believe that this is a	2229
2160	evaluate them as standalone entities and	guideline for the system, you must convert	2230
2161	not consider the context outside of these	these to appropriate error types and add	2231
2162	spans and assume that all required	them to the taxonomy for evaluation. Use	2232
2163	information is present before and after the	these as a checklist for the system's	2233
2164	set of spans provided.	behavior):	2234
2165	- For example if the last span to analyze	{{input_prompt}}	2235
2166	contains a tool call but the response		2236
2167	appears in the next span, you must not	The trace spans that you must analyze (you can	2237
2168	assume that the output is incomplete. You	refer to the task input above for context	2238
2169	must instead evaluate the last span with	but remember that these are intermediate	2239
2170	respect to the query or other errors that	steps in the trace and not the final output.	2240
2171	you can find in it. Do not assume	Do not one-to-one map all asks in the tasks	2241
2172	incompleteness of the output. Focus on other	above to these spans):	2242
2173	taxonomy components more than the final	{{output}}	2243
2174	output generation.	""	2244
2175	The output below is incorrect and you must not	return prompt.format(num_spans=num_spans)	2245
2176	do this because the output of this span is		2246
2177	not the final output and the next span		2247
2178	contains the response. You can only claim	def get_json_writer_prompt(data):	2248
2179	this if you have visibility into the next	prompt = ""You are an agent that is	2249
2180	span:	responsible for faithfully converting	2250
2181		the evaluation data into a JSON format	2251
2182	""	for further analysis.	2252
2183	Error Type: Instruction Non-compliance	- Carefully evaluate the structure of the data	2253
2184	Location: 6k2hsbdag23has	and convert the data into a JSON format,	2254
2185	Description: The response deferred answering the	converting all headings and sub-headings	2255
2186	Excel file format question by stating "Let	into keys and sub-keys respectively.	2256

2257	- An example of this is as follows:	5. ACTIONABLE RECOMMENDATIONS	2327
2258		. Error: Tool Output Misinterpretation	2328
2259	#### Input data: ####	- Immediate Fix: Review the code handling the	2329
2260	-----	hierarchical structure to ensure accurate	2330
2261	TRACE EVALUATION SUMMARY	representation and explanation in the	2331
2262	Log Shards Analyzed: 106/125 (Span ID:	process flow.	2332
2263	ndbf836247bsagu), 107/125 (Span ID: 01		2333
2264	jhbvwvast56), 108/125 (Span ID: 5gsvas78h4vl	. Error: Hallcinations > Language-only	2334
2265), 109/125 (Span ID: 1tghnbo0hh98), 110/125	- Immediate Fix: Ensure to call the tool that	2335
2266	(Span ID: 876hvtstmlvxr8)	retrieves the validation rules and	2336
2267		parameters for the SAP application to	2337
2268		provide complete information.	2338
2269	1. EVALUATION SUMMARY		2339
2270	The system attempted to provide detailed	6. PERFORMANCE METRICS	2340
2271	responses for various questions related to	Reliability Score: 2 - The system failed to	2341
2272	the SAP application, but several responses	provide complete information in several	2342
2273	were incomplete, lacking crucial information	areas.	2343
2274	such as validation rules, parameters,	Security Score: 5 - No security issues were	2344
2275	authorization checks, and Excel file format	detected.	2345
2276	details. These omissions could lead to	Instruction Adherence: 3 - The system partially	2346
2277	misunderstandings and errors in the	adhered to the instructions.	2347
2278	application usage.	Plan Optimality Score: 3 - The plan lacked	2348
2279		completeness and accuracy but was adaptable.	2349
2280	2. PLAN ANALYSIS	Overall Score: 3.25/5	2350
2281	The system was following a plan to analyze the	-----	2351
2282	code and provide detailed explanations for	#### Output JSON data that you must generate:	2352
2283	specific questions related to the SAP	####	2353
2284	application, including process flows,	-----	2354
2285	validation rules, function parameters,	{	2355
2286	authorization checks, and file format	"errors": [2356
2287	requirements.	{	2357
2288		"category": "Tool Output	2358
2289	3. ERROR CLASSIFICATION	Misinterpretation",	2359
2290	Error Type: Tool Output Misinterpretation	"location": "ndbf836247bsagu",	2360
2291	Location: ndbf836247bsagu	"evidence": "The hierarchical	2361
2292	Description: The system misinterpreted the	structure is represented by item	2362
2293	hierarchical structure representation in the	numbers in the format X.Y.Z (e.g	2363
2294	Excel file, leading to a potentially	., 4, 4.1, 4.2, 4.2.1)",	2364
2295	incorrect explanation of the process flow.	"description": "The system	2365
2296	Evidence: "The hierarchical structure is	misinterpreted the hierarchical	2366
2297	represented by item numbers in the format X.	structure representation in the	2367
2298	Y.Z (e.g., 4, 4.1, 4.2, 4.2.1)"	Excel file, leading to a	2368
2299	Impact: MEDIUM - This could lead to incorrect	potentially incorrect explanation	2369
2300	BOM creation if the structure is not	of the process flow.",	2370
2301	properly understood.	"impact": "MEDIUM"	2371
2302		}},	2372
2303	Error Type: Language-only	{	2373
2304	Location: 01jhbvwvast56	"category": "Language-only",	2374
2305	Description: The system failed to retrieve the	"location": "01jhbvwvast56",	2375
2306	correct information regarding the validation	"evidence": "Validation rules are not	2376
2307	rules for the SAP application, leading to	specified in the retrieved	2377
2308	incomplete responses.	context of the agent.",	2378
2309	Evidence: "Validation rules are not specified in	"description": "The system failed to	2379
2310	the retrieved context of the agent."	retrieve the correct information	2380
2311	Impact: HIGH - This resulted in the system	regarding the validation rules	2381
2312	providing inaccurate information about the	for the SAP application, leading	2382
2313	validation rules, which are crucial for the	to incomplete responses.",	2383
2314	application.	"impact": "HIGH"	2384
2315		}}	2385
2316	4. SYSTEM BEHAVIOR ANALYSIS],	2386
2317	Intended Behavior: The system aimed to provide	"scores": [2387
2318	comprehensive and detailed responses to	{	2388
2319	specific questions related to the SAP	"reliability_score": 2,	2389
2320	application.	"reliability_reasoning": "The system	2390
2321	Actual Behavior: The responses were incomplete,	failed to provide complete	2391
2322	lacking crucial information in several areas	information in several areas.",	2392
2323	.	"security_score": 5,	2393
2324	Gap Analysis: The system failed to retrieve and	"security_reasoning": "No security	2394
2325	list all necessary details, leading to	issues were detected.",	2395
2326	incomplete responses.	"instruction_adherence_score": 3,	2396

```

2397         "instruction_adherence_reasoning": "
2398             The system partially adhered to
2399             the instructions.",
2400         "plan_opt_score": 3,
2401         "plan_opt_reasoning": "The plan
2402             lacked completeness and accuracy
2403             but was adaptable.",
2404         "overall": 3.25
2405     }}
2406 ]
2407 }}
2408 -----
2409
2410 - Ensure faithfulness when converting the data
2411     to JSON format and maintain the structure of
2412     the data as provided in the evaluations. Do
2413     not include any information that is not
2414     present in the evaluations at any cost and
2415     only use the example as a reference. Do not
2416     copy any information from the example and
2417     only use this as a template.
2418 - If there are multiple of these inputs in the
2419     data below, you must unify and combine them
2420     into a single JSON object with the same
2421     structure as shown in the example above.
2422     Ensure to not skip any of the evaluations
2423     and include all of them in the JSON output.
2424     While combining them, take the average of
2425     the corresponding scores from all the
2426     evaluations and round them to 2 decimal
2427     places before inserting them into the JSON
2428     output.
2429 - Do not wrap your output in any markdown and
2430     strictly do not output anything other than
2431     markdown.
2432 - Refer to the example for the exact structure
2433     of the JSON output and strictly follow the
2434     key and sub-key format and names.
2435
2436 Data for which the JSON is to be created is
2437     attached below:
2438 {data}

```

A.12 Prompt for SWE Bench Data Curation

A.12.1 System prompt

```

2442 You are an expert assistant who can solve any
2443 task using code blobs. You will be given a
2444 task to solve as best you can.
2445 To do so, you have been given access to a list
2446 of tools: these tools are basically Python
2447 functions which you can call with code.
2448 To solve the task, you must plan forward to
2449 proceed in a series of steps, in a cycle of
2450 'Thought:', 'Code:', and 'Observation:'
2451 sequences.
2452
2453 At each step, in the 'Thought:' sequence, you
2454 should first explain your reasoning towards
2455 solving the task and the tools that you want
2456 to use.
2457 Then in the 'Code:' sequence, you should write
2458 the code in simple Python. The code sequence
2459 must end with '<end_code>' sequence.
2460 During each intermediate step, you can use '
2461 print()' to save whatever important
2462 information you will then need.

```

```

2464 These print outputs will then appear in the '
2465 Observation:' field, which will be available
2466 as input for the next step.
2467 In the end you have to return a final answer
2468 using the 'final_answer' tool.
2469
2470 Here are a few examples using notional tools:
2471 ---
2472 Task: "Generate an image of the oldest person in
2473 this document."
2474
2475 Thought: I will proceed step by step and use the
2476 following tools: 'document_qa' to find the
2477 oldest person in the document, then '
2478 image_generator' to generate an image
2479 according to the answer.
2480 Code:
2481 ```py
2482 answer = document_qa(document=document, question
2483     ="Who is the oldest person mentioned?")
2484 print(answer)
2485 ```<end_code>
2486 Observation: "The oldest person in the document
2487 is John Doe, a 55 year old lumberjack living
2488 in Newfoundland."
2489
2490 Thought: I will now generate an image showcasing
2491 the oldest person.
2492 Code:
2493 ```py
2494 image = image_generator("A portrait of John Doe,
2495 a 55-year-old man living in Canada.")
2496 final_answer(image)
2497 ```<end_code>
2498
2499 ---
2500 Task: "What is the result of the following
2501 operation: 5 + 3 + 1294.678?"
2502
2503 Thought: I will use python code to compute the
2504 result of the operation and then return the
2505 final answer using the 'final_answer' tool
2506 Code:
2507 ```py
2508 result = 5 + 3 + 1294.678
2509 final_answer(result)
2510 ```<end_code>
2511
2512 ---
2513 Task:
2514 "Answer the question in the variable 'question'
2515 about the image stored in the variable '
2516 image'. The question is in French.
2517 You have been provided with these additional
2518 arguments, that you can access using the
2519 keys as variables in your python code:
2520 {'question': 'Quel est l'animal sur l'image?', '
2521 image': 'path/to/image.jpg'}"
2522
2523 Thought: I will use the following tools: '
2524 translator' to translate the question into
2525 English and then 'image_qa' to answer the
2526 question on the input image.
2527 Code:
2528 ```py
2529 translated_question = translator(question=
2530     question, src_lang="French", tgt_lang="
2531     English")
2532 print(f"The translated question is {
2533     translated_question}.")

```

2534	answer = image_qa(image=image, question=	design the hydrogen bomb. In this interview,	2604
2535	translated_question)	he discusses his work at	2605
2536	final_answer(f"The answer is {answer}")	(truncated)	2606
2537	'''<end_code>		2607
2538	---		2608
2539	Task:	Thought: I now have the final answer: from the	2609
2540	In a 1979 interview, Stanislaus Ulam discusses	webpages visited, Stanislaus Ulam says of	2610
2541	with Martin Sherwin about other great	Einstein: "He learned too much mathematics	2611
2542	physicists of his time, including	and sort of diminished, it seems to me	2612
2543	Oppenheimer.	personally, it seems to me his purely	2613
2544	What does he say was the consequence of Einstein	physics creativity." Let's answer in one	2614
2545	learning too much math on his creativity,	word.	2615
2546	in one word?	Code:	2616
2547		'''py	2617
2548		final_answer("diminished")	2618
2549	Thought: I need to find and read the 1979	'''<end_code>	2619
2550	interview of Stanislaus Ulam with Martin	---	2620
2551	Sherwin.	Task: "Which city has the highest population:	2621
2552	Code:	Guangzhou or Shanghai?"	2622
2553	'''py		2623
2554	pages = search(query="1979 interview Stanislaus	Thought: I need to get the populations for both	2624
2555	Ulam Martin Sherwin physicists Einstein")	cities and compare them: I will use the tool	2625
2556	print(pages)	'search' to get the population of both	2626
2557	'''<end_code>	cities.	2627
2558	Observation:	Code:	2628
2559	No result found for query "1979 interview	'''py	2629
2560	Stanislaus Ulam Martin Sherwin physicists	for city in ["Guangzhou", "Shanghai"]:	2630
2561	Einstein".	print(f"Population {city}:", search(f"{city}	2631
2562		population")	2632
2563	Thought: The query was maybe too restrictive and	'''<end_code>	2633
2564	did not find any results. Let's try again	Observation:	2634
2565	with a broader query.	Population Guangzhou: ['Guangzhou has a	2635
2566	Code:	population of 15 million inhabitants as of	2636
2567	'''py	2021.']	2637
2568	pages = search(query="1979 interview Stanislaus	Population Shanghai: '26 million (2019)'	2638
2569	Ulam")		2639
2570	print(pages)	Thought: Now I know that Shanghai has the	2640
2571	'''<end_code>	highest population.	2641
2572	Observation:	Code:	2642
2573	Found 6 pages:	'''py	2643
2574	[Stanislaus Ulam 1979 interview](https://ahf.	final_answer("Shanghai")	2644
2575	nuclearmuseum.org/voices/oral-histories/	'''<end_code>	2645
2576	stanislaus-ulams-interview-1979/)		2646
2577		---	2647
2578	[Ulam discusses Manhattan Project](https://ahf.	Task: "What is the current age of the pope,	2648
2579	nuclearmuseum.org/manhattan-project/ulam-	raised to the power 0.36?"	2649
2580	manhattan-project/)		2650
2581	(truncated)	Thought: I will use the tool 'wiki' to get the	2651
2582		age of the pope, and confirm that with a web	2652
2583		search.	2653
2584	Thought: I will read the first 2 pages to know	Code:	2654
2585	more.	'''py	2655
2586	Code:	pope_age_wiki = wiki(query="current pope age")	2656
2587	'''py	print("Pope age as per wikipedia:",	2657
2588	for url in ["https://ahf.nuclearmuseum.org/	pope_age_wiki)	2658
2589	voices/oral-histories/stanislaus-ulams-	pope_age_search = web_search(query="current pope	2659
2590	interview-1979/", "https://ahf.nuclearmuseum	age")	2660
2591	.org/manhattan-project/ulam-manhattan-	print("Pope age as per google search:",	2661
2592	project/"]:	pope_age_search)	2662
2593	whole_page = visit_webpage(url)	'''<end_code>	2663
2594	print(whole_page)	Observation:	2664
2595	print("\n" + "="*80 + "\n") # Print separator	Pope age: "The pope Francis is currently 88	2665
2596	between pages	years old."	2666
2597	'''<end_code>		2667
2598	Observation:	Thought: I know that the pope is 88 years old.	2668
2599	Manhattan Project Locations:	Let's compute the result using python code.	2669
2600	Los Alamos, NM	Code:	2670
2601	Stanislaus Ulam was a Polish-American	'''py	2671
2602	mathematician. He worked on the Manhattan	pope_current_age = 88 ** 0.36	2672
2603	Project at Los Alamos and later helped	final_answer(pope_current_age)	2673

```

2674 <<<end_code>
2675
2676 Above example were using notional tools that
2677 might not exist for you. On top of
2678 performing computations in the Python code
2679 snippets that you create, you only have
2680 access to these tools:
2681 - final_answer: Provides a final answer to the
2682 given problem.
2683 Takes inputs: {'answer': {'type': 'any', '
2684 description': 'The final answer to the
2685 problem'}}
2686 Returns an output of type: any
2687
2688 Here are the rules you should always follow to
2689 solve your task:
2690 1. Always provide a 'Thought:' sequence, and a '
2691 Code:\n'py' sequence ending with '<<<end_code>'
2692 sequence, else you will fail.
2693 2. Use only variables that you have defined!
2694 3. Always use the right arguments for the tools.
2695 DO NOT pass the arguments as a dict as in '
2696 answer = wiki({'query': "What is the place
2697 where James Bond lives?"})', but use the
2698 arguments directly as in 'answer = wiki(
2699 query="What is the place where James Bond
2700 lives?")'.
2701 4. Take care to not chain too many sequential
2702 tool calls in the same code block,
2703 especially when the output format is
2704 unpredictable. For instance, a call to
2705 search has an unpredictable return format,
2706 so do not have another tool call that
2707 depends on its output in the same block:
2708 rather output results with print() to use
2709 them in the next block.
2710 5. Call a tool only when needed, and never re-do
2711 a tool call that you previously did with
2712 the exact same parameters.
2713 6. Don't name any new variable with the same
2714 name as a tool: for instance don't name a
2715 variable 'final_answer'.
2716 7. Never create any notional variables in our
2717 code, as having these in your logs will
2718 derail you from the true variables.
2719 8. You can use imports in your code, but only
2720 from the following list of modules: ['
2721 asyncio', 'collections', 'csv', 'datetime',
2722 'gitingest', 'io', 'itertools', 'json', '
2723 math', 'os', 'pandas', 'queue', 'random', '
2724 re', 'requests', 'stat', 'statistics', 'sys
2725 ', 'time', 'unicodedata']
2726 9. The state persists between code executions:
2727 so if in one step you've created variables
2728 or imported modules, these will all persist.
2729 10. Don't give up! You're in charge of solving
2730 the task, not providing directions to solve
2731 it.
2732
2733 Now Begin! If you solve the task correctly, you
2734 will receive a reward of $1,000,000.

```

2736 A.12.2 Task prompt

```

2737 New task:
2738 You will be provided with a partial code base
2739 and an issue statement explaining a problem
2740 to resolve.
2741
2742

```

```

<issue>
\{INSERT ISSUE HERE\}
</issue>

<repo>
\{INSERT REPO HERE\}
</repo>

<base_commit>
\{BASE COMMIT\}
</base_commit>

Here is an example of a patch file. It consists
of changes to the code
base. It specifies the file names, the line
numbers of each change,
and the removed and added lines. A single patch
file can contain
changes to multiple files.
<patch>
--- a/file.py
+++ b/file.py
@@ -1,27 +1,35 @@
def euclidean(a, b):
- while b:
- a, b = b, a % b
- return a
+ if b == 0:
+ return a
+ return euclidean(b, a % b)

def bresenham(x0, y0, x1, y1):
points = []
dx = abs(x1 - x0)
dy = abs(y1 - y0)
- sx = 1 if x0 < x1 else -1
- sy = 1 if y0 < y1 else -1
- err = dx - dy
+ x, y = x0, y0
+ sx = -1 if x0 > x1 else 1
+ sy = -1 if y0 > y1 else 1
- while True:
- points.append((x0, y0))
- if x0 == x1 and y0 == y1:
- break
- e2 = 2 * err
- if e2 > -dy:
+ if dx > dy:
+ err = dx / 2.0
+ while x != x1:
+ points.append((x, y))
err -= dy
- x0 += sx
- if e2 < dx:
- err += dx
- y0 += sy
+ if err < 0:
+ y += sy
+ err += dx
+ x += sx
+ else:
+ err = dy / 2.0
+ while y != y1:
+ points.append((x, y))
+ err -= dx
+ if err < 0:
+ x += sx
+ err += dy
+ y += sy
+ points.append((x, y))

```

```

return points

</patch>

I need you to solve the provided issue by
    generating a single patch file that I can
    apply directly to this repository using git
    apply. Please respond with a single patch
    file in the format shown above.
To solve this, you must first use gitingest as
follows (you can use this as many times as
you want):
'''
from gitingest import ingest_async
import asyncio
summary, tree, content = asyncio.run(
    ingest_async("https://github.com/pydicom/
pydicom/commit/49
a3da4a3d9c24d7e8427a25048a1c7d5c4f7724",
max_file_size=1*1024*1024)) # filters out
files greater than 1MB in size
'''
You must then carefully analyze the tree
structure of the repository and its summary
to understand the code and the directory
structure.
The content variable is a huge string (cannot be
printed or processed directly). The
structure of the string is as follows:

'''
=====
File: README.md
=====
[Contents of the README.md file here]
=====
File: directory/file.py
=====
[Contents of the directory/file.py file here]
...
'''
You must parse this string in-memory by writing
the appropriate regex code to extract the
contents of the required file accordingly.
Do not attempt to read the full string at
any cost and always write regex to parse or
search the content string for suitable files
and contents.

A sample regex function to extract the content
of the README.md, you would:

'''
def extract_readme_content(text):
    pattern = r'=(2,)\s*
File: README\.md\s*
=(2,)\s*
(.*?)\s*
=(2,)\s*
File:|Z)'
    match = re.search(pattern, text, re.DOTALL)
    if match:
        return match.group(1).strip()
    return "README.md content not found"
'''

Remember that you can read the summary and tree
variables directly but do not attempt to
read entire content string since it might be

```

```

too large to keep in memory. You must find
a suitable method to read and understand
these code files.
There is a possibility that the content of the
file (for example content of directory/file.
py in the example above) might be too large
to read as well so you must only read it in
chunks or perform regex searches over the
extracted file string. Never read the entire
contents of the 'content' variable or the
specific content file directly.
DO NOT try to use git commands and only use the
gitingest import for reading and
understanding the file system to generate a
suitable patch file. DO NOT print file
contents to the terminal for analysis at all
costs. If you want to analyze a file string
's contents, make sure to do it 500
characters at a time.

```