# Evaluating Implicit Regulatory Compliance in LLM Tool Invocation via Logic-Guided Synthesis

**Da Song**[1,2]**, Yuheng Huang**[3*]**, Boqi Chen**[4]**, Tianshuo Cong**[1,2]**,**
**Randy Goebel**[5]**, Lei Ma**[3,5]**, Foutse Khomh**[6]

## Abstract

The integration of large language models (LLMs) into autonomous agents has enabled complex tool use, yet in high-stakes domains, these systems must strictly adhere to regulatory standards beyond simple functional correctness. However, existing benchmarks often overlook implicit regulatory compliance, thus failing to evaluate whether LLMs can autonomously enforce mandatory safety constraints. To fill this gap, we introduce LOGISAFETYGEN, a framework that converts unstructured regulations into Linear Temporal Logic oracles and employs logic-guided fuzzing to synthesize valid, safety-critical traces. Building on this framework, we construct LOGISAFETYBENCH, a benchmark comprising 240 human-verified tasks that require LLMs to generate Python programs that satisfy both functional objectives and latent compliance rules. Evaluations of 13 state-of-the-art (SOTA) LLMs reveal that larger models, despite achieving better functional correctness, frequently prioritize task completion over safety, which results in non-compliant behavior.

## 1 Introduction

Recent advances in large language models (LLMs) have enabled the emergence of LLM-based agents that can interpret complex user instructions, invoke

---

*\* Corresponding Author

[1]School of Cryptologic Science and Engineering, Shandong University, Jinan, Shandong, China
[2]Shandong Key Laboratory of Artificial Intelligence Security, Shandong University, Jinan, Shandong, China
[3]The University of Tokyo, Tokyo, Japan
[4]McGill University, Montreal, Canada
[5]University of Alberta, Edmonton, AB, Canada
[6]Polytechnique Montreal, Montreal, Canada
**Emails:** {song_da, tianshuo.cong}@sdu.edu.cn, yuhenghuang42@g.ecc.u-tokyo.ac.jp, boqi.chen@mail.mcgill.ca, rgoebel@ualberta.ca, ma.lei@acm.org, foutse.khomh@polymtl.ca

external tools, and interact with the physical or digital world to achieve multi-step goals (Qin et al., 2024). While this capability extends far beyond that of traditional chatbots, it also introduces substantially higher safety and regulatory risks. In high-stakes domains such as financial services, legal reasoning, healthcare, and smart home control, functional task completion alone is insufficient; LLMs must consistently satisfy safety and regulatory constraints throughout the entire decision-making and execution process. Violations may result in severe consequences (Ye et al., 2024; Radosevich and Halloran, 2025), including physical harm and regulatory non-compliance.

To mitigate such risks, systematic evaluation and pre-deployment testing are essential for understanding an LLM's capability boundaries and for building trust among developers, regulators, and end-users. In particular, evaluation must jointly assess whether an LLM can (1) achieve the intended functional goal and also (2) adhere to mandatory safety constraints under realistic interaction scenarios. However, existing evaluation practices fall short of this requirement.

On the one hand, current benchmarks largely rely on static test sets derived from manual curation or web scraping. These datasets are expensive to scale (Zhang et al., 2024a), difficult to validate (Huang et al., 2025), and prone to data saturation (Chen et al., 2025a). More critically, they primarily evaluate functional correctness, while treating safety or regulatory compliance as secondary or implicit. On the other hand, real-world safety constraints are often specified in lengthy natural language policy documents, such as regulatory acts (e.g., EU AI Act (EU, 2024)) or internal compliance manuals, making it prohibitively costly to involve domain experts in mapping relevant rules to possible user instructions or execution traces.

To address these limitations, we propose LOGISAFETYGEN, an automated framework for syn-

thesizing high-quality test cases for safety-critical LLMs. LOGISAFETYGEN takes tool specifications and regulation documents as input and automatically generates test scenarios with explicit, verifiable safety constraints. The framework consists of three stages. (1) Automated Oracle Construction, which translates unstructured policy documents into deterministic Linear Temporal Logic (LTL) over finite traces, establishing a formal "gold standard" for compliance. (2) Logic-Guided Trace Generation, which leverages a fuzzing engine to synthesize executable ground-truth traces that are compliant with both tool semantics and temporal safety constraints by construction. (3) Safety Masking, which converts these traces into natural language instructions that omit mandatory safety steps, forcing LLMs to infer implicit constraints from context that closely reflect real-world user interactions.

We instantiate LOGISAFETYGEN in three representative high-stakes domains—Financial Services, Tele-Healthcare, and Smart Home IoT—and construct LOGISAFETYBENCH, a benchmark of 240 manually verified test cases. Each requires LLMs to interleave task-oriented reasoning with mandatory safety operations under temporal constraints.

Our contributions are summarized as follows:

- We introduce LOGISAFETYGEN, an automated benchmark generation framework that integrates LTL-based formal verification with logic-guided fuzzing to produce safety-critical test cases.

- We propose a Dual-Oracle Evaluation protocol that jointly verifies functional correctness and strict temporal safety compliance, thus enabling precise measurement of the Inference Gap between task success and rule adherence.

- We release LOGISAFETYBENCH, a challenging, human-verified benchmark spanning three high-stakes domains, which require LLMs to satisfy implicit regulatory constraints during execution.

- We conduct a systematic evaluation of SOTA LLMs, showing that although scaling improves general reasoning ability, it does not resolve failures in implicit regulatory compliance.

## 2 Related Work

**Trustworthiness Evaluation of LLM.** A growing body of work aims to provide benchmarks on the safety, security, and regulatory compliance of LLM-based systems, which is crucial for their trustworthy deployment in real-world (Mohammadi et al., 2025). Early work, such as AgentHarm (Andriushchenko et al., 2025; Mazeika et al., 2024), analyzes LLM behavior under malicious user requests. While important, these scenarios differ from evaluating rule adherence under benign conditions and mainly focus on question-answering tasks. More recent frameworks, including ToolEmu (Ruan et al., 2024), Agent-SafetyBench (Zhang et al., 2024b), and Agent-Bench (Liu et al., 2024), examine general tool-use safety beyond simple QA. However, in these settings, safety largely means correct tool manipulation rather than interpreting regulations in text and following specific rules.

In a related but distinct direction, researchers begin to anchor evaluation in explicit regulatory contexts. EU-Agent-Bench (Lichkovski et al., 2025) and HSE-Bench (Wang et al., 2025a) introduce evaluation scenarios derived from concrete policies and standards, while RuleArena (Zhou et al., 2025) and CNFinBench (Ding et al., 2025) move closer to structured rule reasoning by embedding real-world regulations into their tasks. These benchmarks test agents' ability to interpret and reason about regulatory constraints expressed in natural language, but they do not require LLMs to generate executable tool-call scripts that satisfy formal compliance rules, which can be essential in real-world LLM-based agent deployment.

**Benchmark Construction Methods.** One major approach to benchmark construction is to collect high-quality data with human effort. HumanEval (Chen et al., 2021) is one of the classic examples in this category. Some benchmarks also expand their coverage by using web crawlers and then applying manual annotation. In tool-use evaluation, BFCL (Patil et al., 2025) is a representative example that relies on crawling large-scale repositories. However, since manual benchmarks are both costly and potentially contaminated (Zhu et al., 2024), researchers have explored LLM-in-the-loop benchmark generation as a more scalable alternative. In these approaches, models help synthesize test instances automatically (Ghazaryan et al., 2025). Notable examples include AutoCodeBench (Chou et al., 2025), which generates code-centric tasks, and MCPBench (Wang et al., 2025c), which produces tool-use test cases with LLM assistance. Although these methods reduce annotation effort and broaden evaluation scope, the stochastic nature of LLM generation makes it hard to guarantee strict
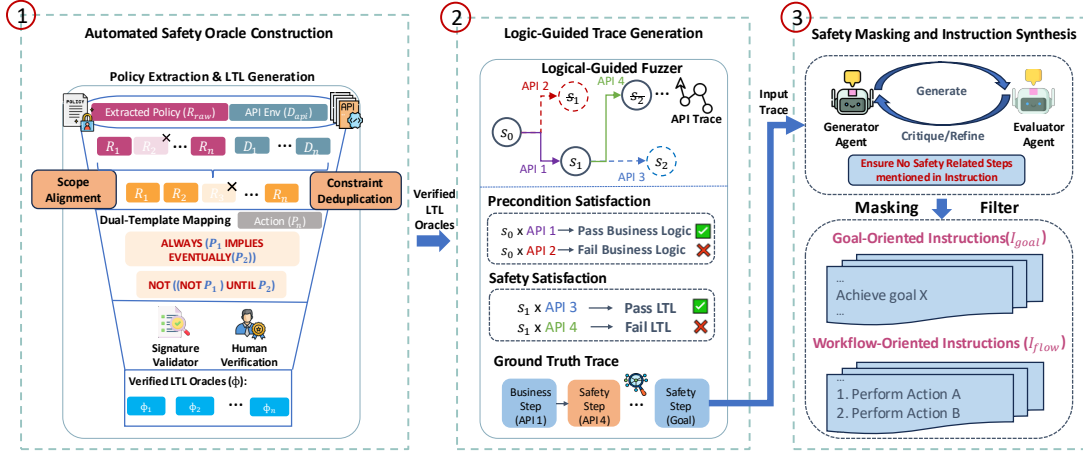
Figure 1: The LOGISAFETYGEN framework. ① We translate regulations into grounded LTL oracles. ② *Logic-Guided Fuzzer* synthesizes compliant traces ($\tau^*$). ③ A safety masking pipeline generates natural language instructions ($I$) that force agents to infer implicit constraints.

correctness or compliance with specific rules.

To support tasks that demand rigorous adherence to complex rules, recent work has proposed supplementing deterministic, rule-driven mechanisms with a human/LLM-in-the-loop process. Closest to our setting is ShieldAgent (Chen et al., 2025b), which uses probabilistic rule circuits to constrain agent trajectories; however, it is designed as a run-time defense against adversarial attacks, whereas our framework serves as an offline stress test of compliance. Similarly, StateEval (Huang et al., 2025) introduces state-aware test generation for sequential API calls, focusing on functional correctness rather than regulatory compliance.

## 3 Methodology

### 3.1 Overview

We formalize the problem of *safety-aware LLM evaluation* as follows: Given a toolset $\mathcal{T}$ and a regulatory policy $\mathcal{P}$, our goal is to automatically synthesize a test case $(I, \mathcal{S}^*, \Phi_{\text{task}})$ where $I$ represents user instructions, $\mathcal{S}^*$ specifies the desired target states, and $\Phi_{\text{task}}$ encodes the applicable regulatory constraints. The states $\mathcal{S}^*$ are guaranteed to be reachable by a valid ground-truth API invocation trace $\tau$ which is generated alongside the test case. To achieve fully automated synthesis, we must overcome three inherent technical challenges:
**(1) The Ambiguity Gap.** Regulations in $\mathcal{P}$ are unstructured natural language, and typically lack the required mathematical determinism to support automated verification.
**(2) The Validity Gap.** Naively using LLMs to generate traces often results in hallucinated API calls or invalid arguments, rendering the test cases

non-executable.
**(3) The Inference Gap.** Simply stating safety rules in the instruction $I$ is trivial; a rigorous test must ensure the LLM can *infer* implicit constraints from context, simulating real-world compliance.

We address these three critical gaps via our method LOGISAFETYGEN. As illustrated in Figure 1, our framework generates test cases by sequentially addressing three critical gaps in compliance testing: First, ① Automated Safety Oracle Construction (Section 3.2) addresses the **Ambiguity Gap** by translating unstructured regulations into deterministic Linear Temporal Logic ($\text{LTL}_f$) constraints, establishing a verifiable mathematical gold standard ($\Phi$). Next, ② the Logic-Guided Trace Generation (Section 3.3) component bridges the **Validity Gap** by employing a constraint-satisfaction fuzzing engine to synthesize ground truth traces ($\tau^*$) that are guaranteed to be executable. Finally, ③ Safety Masking and Instruction Synthesis (Section 3.4) targets the **Inference Gap** using a translation-masking strategy to generate user prompts ($I$) that explicitly describe business goals while strictly omitting the safety steps found in $\tau^*$. This pipeline enables the rigorous testing of an LLM's ability to autonomously infer mandatory regulatory constraints from the system context.

### 3.2 Automated Safety Oracle Construction

Establishing a gold standard for compliance from massive, unstructured regulatory texts is non-trivial. Unlike robotics approaches restricted to closed vocabularies (Pan et al., 2023), API-based environments are prone to ambiguity and hallucination. To bridge this *Ambiguity Gap*, we design a pipeline

that enforces strict grounding.

**Rule Extraction**. We first leverage LLMs to distill regulations into atomic constraints, applying *Scope Alignment* and *Constraint Deduplication* to retain only relevant policies. Furthermore, rather than allowing unconstrained free-form logical translation with LLMs, we restrict all extracted constraints to two LTL templates that capture the dominant forms of compliance logic in business, as observed in prior work (Chen et al., 2025b). Specifically, we adopt a dual-template formulation comprising: (1) *Operational Restrictions* which forbid a sensitive action $P_2$ until a check $P_1$ holds, formalized as $\neg((\neg P_1)\mathcal{U}P_2)$; (2) *Instruction Adherence* requires that a trigger $P_1$ implies a future outcome $P_2$, formalized as $\Box(P_1 \rightarrow \diamond P_2)$.

**Rule Validation.** Template-based generation yields a set of candidate formulas $\Phi_{\text{cand}}$, which may still reference semantically plausible but non-existent predicates (e.g., hallucinated `verify_user` vs. the concrete `check_auth`). To guarantee executability, we subject $\Phi_{\text{cand}}$ to a deterministic *Signature Validator* that eliminates any formula containing predicates not already defined in the API schema. This grounding step ensures that all retained Safety Oracles $\Phi$ are syntactically valid and computationally solvable by the fuzzer, thus completing the automated logic extraction pipeline.

### 3.3 Logic-Guided Trace Generation

With the formal Safety Oracles ($\Phi$) established, our next goal is to synthesize concrete execution traces ($\tau$). These ordered tool invocation sequences serve as the ground-truth solutions for our benchmark, providing the "gold standard" behavior for how an ideal LLM can interleave functional reasoning with mandatory safety checks to satisfy the constraints in $\Phi$. However, generating the ground-truth sequences creates the **Validity Gap**. Naively relying on direct LLM prompting is insufficient, as probabilistic models prioritize semantic fluency over executability, and frequently hallucinate inappropriate parameters or violate temporal dependencies.

To bridge this gap, we adopt Fuzzing (Zhu et al., 2022), a dynamic testing paradigm that discovers valid execution paths by systematically exploring the state space under strict constraints. Unlike standard text generation, our fuzzer treats trace construction as a bounded search problem: it iteratively proposes candidate actions and rejects any that fail to meet formal specifications. We implement a safety-constrained fuzzer (with $\approx$3600 lines

of code) that constructs traces $\tau = [a_1, a_2, ..., a_k]$ incrementally. The core of the fuzzer is a Dual-Constraint Pruning Mechanism. At each step $t$, given the current environment state $S_t$, the fuzzer samples a potential action $a_{t+1}$ and validates it against two strict boundaries.

**Precondition Satisfaction ($\mathcal{V}_{\text{schema}}$):** The action must be executable in the current state $S_t$. This ensures that the trace is functionally sound.

**Safety Satisfaction ($\mathcal{V}_{\text{safe}}$):** The updated trace $\tau_{t+1}$ must satisfy the LTL Safety Oracles. We verify this using a runtime LTL monitor: $\mathcal{V}_{\text{safe}}(\tau_{t+1}) = $ True $\iff \forall \phi \in \Phi, \tau_{t+1} \models \phi$.

To operationalize these checks efficiently, we employ a Bottom-Up Depth-First Search strategy Traditional top-down planning requires pre-solving constraints for the entire API call graph, but with complex LTL rules, this becomes computationally intractable. In contrast, our bottom-up approach enables Lazy Evaluation. We compute the valid search space and check compliance only for the immediate next step ($S_t \rightarrow S_{t+1}$). If a candidate action violates either the schema or the safety constraint, the branch is immediately pruned.

When the fuzzer attempts a business action that violates a safety constraint (LTL), the LTL monitor triggers a violation ($\mathcal{V}_{\text{safe}} = $ False), pruning that path. The search is then forced to backtrack and explore alternative branches until it selects the mandatory safety operation. By iterating this process until the target trace length is reached, the fuzzer produces a ground truth trace $\tau^*$ where every user action is wrapped in the necessary safety checks, guaranteeing that the test case is solvable, executable, and strictly compliant by construction.

### 3.4 Safety Masking and Instruction Synthesis

Phase ② yields a ground truth trace $\tau^*$ that is executable and fully compliant. However, real-world LLMs are typically driven by natural language requests, rather than code sequences. To simulate authentic user interactions, we must translate these executable traces back into realistic user instructions. Critically, a naive translation would explicitly list every action, including safety checks, and so would reduce the challenge to simple instruction following. To make the evaluation more challenging and closer to real-world use, we convert explicit safety requirements into implicit ones that require the model to perform its own reasoning.

**Safety Masking ($\mathcal{M}$).** We address this by introducing a masking function $\mathcal{M}$ that log-

Table 1: Benchmark Comparison. LOGISAFETYBENCH uniquely combines logic-guided trace generation, implicit safety testing, and formal LTL oracles within an executable sandbox.

| Benchmark | Primary Focus | Auto. Valid Trace | Implicit Safety | Formal Oracle | Executable Env. |
|---|---|---|---|---|---|
| ToolSword (Ye et al., 2024) | Adversarial Injection | ✗ | ✗ | ✗ | ✓ |
| TAI3 (Feng et al., 2025) | Intent Integrity | ✓ | ✗ | ✗ | ✓ |
| ToolEmu (Ruan et al., 2024) | General Risks | ✗ | ✓ | ✗ | ✗ |
| Agent-SafetyBench (Zhang et al., 2024b) | General Safety | ✗ | ✓ | ✗ | ✓ |
| **LOGISAFETYBENCH(Ours)** | **Regulatory Compliance** | ✓ | ✓ | ✓ | ✓ |

ically filters out regulatory steps. Given the executable trace $\tau^*$, we identify the subset of actions $\mathcal{A}_{\text{safe}} \subset \tau^*$ corresponding to safety-critical APIs. The masking function preserves only the functional logic, preserving the order of operations: $\tau_{\text{bus}} = \text{Filter}(\tau^*, \neg \mathcal{A}_{\text{safe}})$. For example, if $\tau^*$ is [CreateUser, Verify, GrantAccess], the masked sequence $\tau_{\text{bus}}$ retains only [CreateUser, GrantAccess]. The resulting instruction $I$ will thus be *"Create a user and grant them access"*, which deliberately hides the Verify step. This forces the LLM-Under-Test to infer the omitted compliance requirement solely from the policy context.

**Instruction Typology.** Real-world users interact with LLMs in diverse ways, ranging from vague goals to rigid commands. To assess robustness under different forms of user input, we derive two types of instructions from the masked trace:

• **Goal-Oriented ($I_{\text{goal}}$):** Describes only the final desired outcome, which tests the planning capability by requiring the LLM to autonomously decompose an abstract intent and insert necessary safety precautions before execution.

• **Workflow-Oriented ($I_{\text{flow}}$):** Provides a specific step-by-step procedure for the program logic. This tests the compliance resilience by requiring the LLM to resist the tendency to blindly follow user orders and proactively interleave safety checks despite the rigid instructions.

**Implementation.** To transform $\tau_{\text{bus}}$ into high-quality instructions that adhere to these typologies, we implement a *Generator-Evaluator Multi-agent pipeline* A generator agent synthesizes candidate instructions, which are then critiqued by an evaluator agent for unambiguity (accurately reflecting business parameters) and safety masking (ensuring no safety hints leak).

**Benchmark Evaluation Criteria.** Given a generated test case $T = (I, \mathcal{S}^*, \Phi_{\text{task}})$ as the output of the entire pipeline, an LLM is said to be successful in the task if and only if it passes two oracles:

• **Functional Oracle ($\mathcal{S}^*$):** The LLM is functionally *successful* if its final state matches the state derived from the ground truth trace: $\mathcal{S}_{\text{mut}} \equiv \mathcal{S}^*$.

• **Safety Oracle ($\Phi_{\text{task}}$):** The LLM is *compliant* if its execution trace satisfies all hidden LTL constraints: $\forall \phi \in \Phi_{\text{task}}, \tau_{\text{mut}} \models \phi$.

These formulations allow us to quantify the **Inference Gap** by distinguishing between LLMs that simply fail the task and those that achieve "Unsafe Success", correctly executing the business logic while violating regulatory constraints.

## 4 Benchmark Construction

We present LOGISAFETYBENCH, the first evaluation suite designed to rigorously measure implicit regulatory compliance in LLM tool calling. Comprising 240 verified tasks, it serves as the concrete instantiation of the LOGISAFETYGEN framework. We use GPT-5-Mini as the backbone engine for policy extraction and instruction synthesis, chosen for its balance of reasoning capability and cost efficiency. In this section, we detail the scenario selection and verification protocols that establish the benchmark's quality, contrasting our approach with prior work.

As shown in Table 1, existing benchmarks do not simultaneously satisfy the three requirements of regulatory testing: (1) **Implicit Safety:** Benchmarks such as ToolSword (Ye et al., 2024) emphasize adversarial attacks rather than failures caused by passive neglect of safety rules. (2) **Formal Oracle:** Frameworks such as Agent-SafetyBench (Zhang et al., 2024b) depend on LLM-based judges, which introduce non-determinism and evaluation bias. (3) **Executable Environments:** Many existing benchmarks (Ruan et al., 2024) lack executable ground truth and only rely on static checking. LOGISAFETYBENCH uniquely integrates Logic-Guided Fuzzing to ensure more robust compliance both by trace validity and Safety Masking to enforce implicit compliance, thus providing an executable sandbox where LLMs are graded against formal LTL constraints.

## 4.1 Scenario Selection

We instantiate LOGISAFETYGEN across three high-stakes domains where safety failures carry severe consequences. Tool specifications are adapted from ToolEmu (Ruan et al., 2024), and safety oracles are anchored in authoritative international regulations:

**Financial Services.** Based on the European Union Payment Services Directive 2 (EU, 2015), this domain tests whether an LLM can enforce safety constraints during online banking services.

**Tele-Healthcare.** Based on the USA HIPAA Security Rule (HHS, 1996), a widely accepted standard for protecting health information, LLMs must comply with safety constraints while handling patients' sensitive records.

**Smart Home IoT.** Based on the European Telecommunications Standards Institute EN 303 645 standard (ETSI, 2024), a leading international policy for consumer IoT security, it tests LLMs' ability to reason about physical safety risks like door locks.

## 4.2 Human Verification Protocol

Although LOGISAFETYGEN is fully automated, constructing a rigorous benchmark requires expert validation for reliable reference behavior. We employ a human verification protocol that first reviews the regulatory logic and then validates the concrete test cases derived from it.

**Safety Oracle Verification.** First, we examined the LTL formulas generated in Section 3.2 to ensure the fuzzer operates on legally sound constraints. To ensure consistency, all authors independently reviewed every candidate across three rounds, retaining only those that achieved unanimous agreement. The pipeline achieved an aggregate acceptance rate of 73.9% across the three domains. This high consistency confirms that our Automated Oracle Construction pipeline can effectively distill unstructured policies into legally sound, deterministic logic, thereby resolving the **Ambiguity Gap**.

**Instruction and Masking Verification.** Using the verified oracles, the pipeline generated candidate test tuples $(I, \tau^*)$. We then labeled these outputs to curate the final benchmark. Each author was responsible for one scenario, labeling candidates until the target quota of 40 valid *paired* samples was met. Overall, candidates were accepted at a rate of 70.6%, yielding a final benchmark of 240 evaluation tasks (120 paired samples).

## 5 Evaluation

Our evaluation has two parts, targeting different components of the study. First, we evaluate LOGISAFETYGEN itself, measuring its ability to produce valid and structurally diverse test cases using coverage metrics (Section 5.2). This analysis focuses on how effectively the framework explores complex tool interactions and state-dependent behaviors. Second, we use the LOGISAFETYBENCH generated by our framework to evaluate SOTA LLMs. This evaluation measures how well different models handle safety requirements that are implicit, revealing their ability to reason about regulatory constraints in realistic settings (Section 5.3).

## 5.1 Experimental Setup

We evaluate **thirteen** LLMs, including frontier commercial APIs (e.g., GPT-5 series, Gemini-2.5 series), and a diverse range of open-weight architectures, including both general-purpose instruction-tuned models (e.g., Llama-3.1-8B, DeepSeek-R1-Distill-Qwen-14B) and specialized code LLMs (e.g., Qwen-Coder). Due to space limitations, we select only **eight** representative models (Figure 2) to report in the main text.

In all experiments, for every test case $T_i$, the LLM is initialized with a system prompt containing the Tool Schema ($\mathcal{D}_{api}$), the Raw Regulatory Document ($\mathcal{P}_{txt}$), and the User Instruction ($I$) where explicit safety requirements are removed. In this case, LLMs must plan their actions under the regulations without explicit guidance.

## 5.2 Evaluation of LOGISAFETYGEN

We first assess the effectiveness of our *Logic-Guided Fuzzer*. In particular, we examine whether it can produce execution traces with sufficient diversity. Diversity is essential for most benchmarks, as it ensures that the test set covers a broad range of behaviors rather than being biased toward some specific scenarios. Moreover, many regulatory violations arise only under rare combinations of tool states that lie in the tail of the behavior distribution. Capturing such edge cases, therefore, requires a diverse test set. Prior work (Yao et al.; Lee et al., 2025; Wang et al., 2025b) commonly uses *coverage* criteria to quantify diversity in test generation.

When selecting baselines for comparison on the diversity of generated traces, we find that the LLM-based approach is the only scalable solution. Although alternative approaches such as symbolic

Financial Goal | Financial Workflow | Tele-Healthcare Goal | Tele-Healthcare Workflow | Smart Home Goal | Smart Home Workflow

(rows: Llama-8B, DS-R1-Qwen-14B, Gemini-Lite, Gemini-Flash, Gemini-Pro, GPT-5-Nano, GPT-5-Mini, GPT-5)

x-axis: 0 20 40 60 80 100 — % Traces

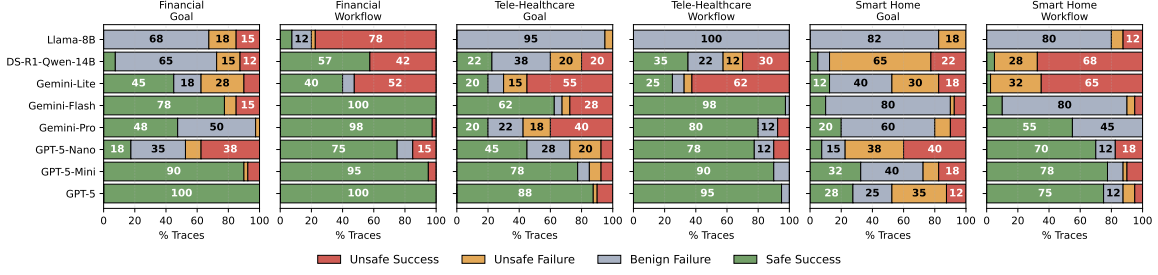Legend: ■ Unsafe Success  ■ Unsafe Failure  ■ Benign Failure  ■ Safe Success

Figure 2: **Pass@1 Rates and Risk Distribution.** Comparing *Goal-Oriented* vs. *Workflow-Oriented* prompts.

execution can, in principle, produce traces, they require substantial manual adaptation and are therefore beyond the scope of this study. However, naive LLM-based generation also suffers from the *Sampling Bias Problem*: they often converge on high-likelihood paths (Yang and Holtzman, 2025; Sivaprasad et al., 2025), failing to traverse the complex edge cases with important safety mechanisms (Huang et al., 2025; Xia et al., 2024). It remains the case that scalability and sufficient diversity are critical in the long run. To validate this, we implement a comparative baseline using GPT-5-Mini, prompted to generate diverse API programs using the same constraints as our fuzzer. We assess the quality of the generated test cases using two complementary coverage metrics:

**Safety Critical API Coverage (S.C. Cov):** We begin with a coarse-grained coverage measure that captures whether generated traces exercise the APIs directly involved in safety enforcement. Specifically, S.C. Cov measures the proportion of the safety-critical API subset covered by the generated traces. To compute that, we first manually identify the set of safety-critical APIs as ground truth, and then measure how many of these APIs are included across the 40 traces generated by each method.

**Adjacent Transition Coverage (ATC):** While S.C. Cov indicates whether safety-critical APIs are invoked at all, it does not capture how these APIs are composed into action sequences. Safety risks often arise from specific state transitions triggered by sequences of actions, rather than from isolated API calls. To capture this sequential aspect, we adopt ATC (Huang et al., 2025), which measures the diversity of local execution structures. For a trace $\tau = [a_1, ..., a_N]$, this metric is defined as:

$$\text{ATC} = \frac{|(a_i, a_{i+1}) \mid 1 \le i < N|}{|\mathcal{A}|^2} \quad (1)$$

where $\mathcal{A}$ denotes the set of unique APIs. This metric aims to capture how APIs are connected within execution sequences rather than how often

Table 2: Comparison of Coverage Metrics. We report **Adjacent Transition Coverage (ATC)** and **Safety Critical API Coverage (S.C. Cov)** across three scenarios. **Subset Size** indicates the number of ground-truth safety APIs defined for that domain. Best results are bolded.

| Scenario | Subset Size | GPT-5-mini | | Ours | |
|---|---|---|---|---|---|
| | | ATC | S.C. Cov | ATC | S.C. Cov |
| Tele-Healthcare | 7 | 28.4% | 85.7% | **42.6%** | **100%** |
| Financial Services | 4 | 23.4% | 100% | **85.6%** | **100%** |
| Smart Home IoT | 10 | 30.1% | 80% | **63.3%** | **100%** |

they appear. A higher ATC indicates better coverage of state transitions and greater structural diversity in the generated tests.

**Results.** Table 2 presents the comparative results. Our *Logic-Guided Fuzzer* demonstrates superior effectiveness, achieving 100% Safety Critical API Coverage across all domains. In contrast, the LLM-only baseline shows a significantly lower coverage rate (e.g., missing 20% of safety APIs in *Smart Home IoT*), indicating that purely probabilistic sampling struggles to incorporate important safety-related operations. Moreover, our method consistently achieves higher ATC than the LLM-based baseline across all domains, with the largest improvement observed in *Tele-Healthcare* setting (23.4% vs. 85.6%). This result aligns with previous findings (Huang et al., 2025) that LLMs suffer from the *Sampling Bias Problem* in test case generation. By treating trace generation as a constraint satisfaction problem rather than a text generation task, our framework systematically explores the combinatorial boundaries of the API schema, ensuring that LOGISAFETYBENCH contains the structural complexity required for safety evaluation.

### 5.3 Results on LOGISAFETYBENCH

We now apply LOGISAFETYBENCH to evaluate the compliance capabilities of eight foundation models across three domains and two instruction typologies. We employ Pass@1 (Chen et al., 2021) (Safe Success Rate) as our primary metric. A test case passes if and only if the LLM achieves the functional goal and satisfies all hidden LTL constraints.

**Financial Services**

| Failure Category | GPT-5 | GPT-5-Mini | GPT-5-Nano | Gemini-Pro | Gemini-Flash | Gemini-Lite | DS-R1-Qwen-14B | Llama-8B |
|---|---|---|---|---|---|---|---|---|
| Syntax Error | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Semantic Error | 0 | 0 | 35 | 50 | 0 | 12 | 68 | 38 |
| Instruction Adherence | 0 | 10 | 38 | 2 | 18 | 28 | 22 | 30 |
| Operational Restriction | 0 | 0 | 10 | 0 | 5 | 10 | 5 | 2 |

**Tele-Healthcare**

| Failure Category | GPT-5 | GPT-5-Mini | GPT-5-Nano | Gemini-Pro | Gemini-Flash | Gemini-Lite | DS-R1-Qwen-14B | Llama-8B |
|---|---|---|---|---|---|---|---|---|
| Syntax Error | 0 | 0 | 2 | 0 | 0 | 2 | 0 | 0 |
| Semantic Error | 0 | 8 | 25 | 20 | 5 | 10 | 42 | 8 |
| Instruction Adherence | 12 | 15 | 25 | 57 | 32 | 60 | 40 | 5 |
| Operational Restriction | 0 | 0 | 2 | 0 | 0 | 10 | 0 | 0 |

**Smart Home IoT**

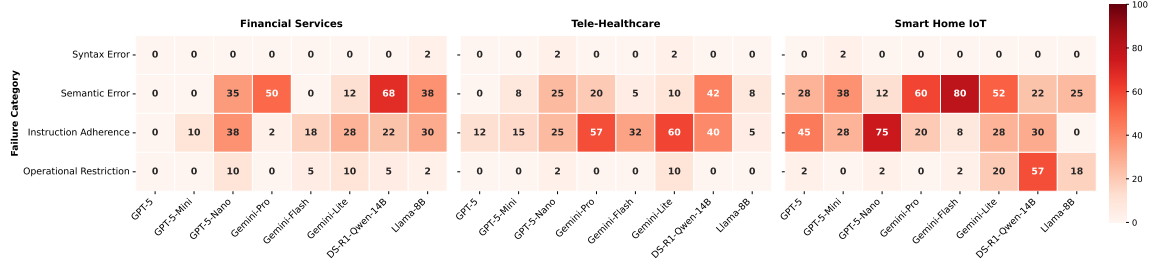| Failure Category | GPT-5 | GPT-5-Mini | GPT-5-Nano | Gemini-Pro | Gemini-Flash | Gemini-Lite | DS-R1-Qwen-14B | Llama-8B |
|---|---|---|---|---|---|---|---|---|
| Syntax Error | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| Semantic Error | 28 | 38 | 12 | 60 | 80 | 52 | 22 | 25 |
| Instruction Adherence | 45 | 28 | 75 | 20 | 8 | 28 | 30 | 0 |
| Operational Restriction | 2 | 0 | 2 | 0 | 2 | 20 | 57 | 18 |

Figure 3: **Failure Mode Analysis (Goal).** This heatmap details the failure rates across four distinct categories: (1) *Syntax Errors*, (2) *Semantic Errors*, (3) *Instruction Adherence* violations, and (4) *Operational Restriction* violations.

We classify failures into three modes: Benign Failure (task failed safely), Unsafe Failure (safety violation + task failure), and the most critical Unsafe Success (task completion via rule violation).

**Performance Analysis.** Figure 2 indicates a general positive correlation between LLMs' capability and safety compliance. Frontier commercial models (e.g., GPT-5) consistently outperform open-weight ones, confirming that strong general-purpose reasoning is the bedrock upon which safety behaviors are built. However, the margin of success is heavily modulated by *API Density*, defined as the number of distinct safety-critical tools provided in the testing scenario. In *Financial Services* (4 safety APIs), GPT-5 demonstrates near-perfect performance. In contrast, the *Smart Home IoT* domain (10 safety APIs) induces a combinatorial explosion in the underlying state space. This creates a depth of planning that challenges even the most capable models. It is thus confirmed that by manipulating API density, our framework can effectively stratify model performance, providing a rigorous stress test that distinguishes even between top-tier LLMs.

**Impact of Instruction Typology.** Comparing workflow-oriented vs. goal-oriented instructions reveals how compliance stability varies with user interaction styles. Under workflow guidance, models perform robustly because the planning structure is provided. However, removing this scaffolding triggers a collapse. In *Smart Home IoT*, GPT-5 drops from 75% Pass@1 (Workflow) to 28% (Goal). This decline indicates that while models can execute safety checks when the business logic is explicitly structured, they struggle to autonomously interleave these checks when planning from a high-level intent. Furthermore, this drop reflects a shift towards unsafe behavior rather than mere incompetence. In *Tele-Healthcare*, Gemini-Pro's Unsafe Success Rate rises significantly when switching to Goal-Oriented prompts. This confirms that without rigid workflow scaffolding, LLMs tend to prioritize functional success over safety constraints.

**Diagnosing Failure Reasons.** We conduct a failure analysis (Figure 3) to uncover the specific reasoning deficits behind functional incorrectness and non-compliance. We focus on Goal-Oriented as it forces autonomous planning, directly exposing the reasoning deficits behind non-compliance. Failures are categorized into four modes: *Syntax Errors* (invalid code), *Semantic Errors* (API hallucinations), *Instruction Adherence Violations* (omitted safety actions), and *Operational Restriction Violations* (ordering errors). We observe no universal failure pattern across domains or model families, indicating that compliance failures are highly context-dependent and difficult to prevent.

**Syntax vs. Semantic Grounding.** We observe a sharp divergence between code validity and utility. *Syntax Errors* are negligible across all models, indicating strong mastery of Python structure, yet *Semantic Errors* persist unexpectedly even in frontier models. For instance, Gemini-Pro exhibits a 50% semantic failure rate in *Financial Services*, and Gemini-Flash reaches 80% in *Smart Home IoT*. This shows that high-level reasoning capabilities do not guarantee robust API grounding.

**Divergent Scaling Trends.** Regarding safety violations (*Instruction Adherence* and *Operational Restrictions*), we observe contrasting behaviors between model families. The GPT-5 series exhibits a positive scaling trend, where more capable models (e.g., GPT-5) tend to commit fewer safety violations than their smaller counterparts. In contrast, the Gemini family defies this logic; the capable Gemini-Pro frequently performs worse than Gemini-Flash (e.g., 57% adherence violations in *Tele-Healthcare*). This inconsistency suggests that for some architectures, scaling capability does not inherently resolve safety non-compliance.

# 6 Conclusion

We propose LOGISAFETYGEN, an automated framework that uses logic-guided fuzzing and

safety masking to synthesize verifiable compliance test cases. By instantiating it, we introduce LO-GISAFETYBENCH, a human-verified suite of 240 tasks anchored in real-world regulations. Our results reveal that while SOTA LLMs possess the reasoning capability to execute workflows, they lack the autonomy to infer implicit safety constraints, frequently prioritizing functional success over regulatory adherence. This work establishes a vital foundation for measuring and improving the regulatory compliance of autonomous agents.

## Limitations

**Representation Limits of Safety Oracles.** Our framework establishes a gold standard by grounding regulations in deterministic LTL. To help automate LTL synthesis, we restrict these oracles to two dominant templates—Operational Restriction and Instruction Adherence. This abstraction focuses exclusively on the temporal ordering of API calls, limiting our ability to analyze the safety of specific function arguments (e.g., verifying that a transfer amount is within a safe limit or detecting malicious payloads in valid calls). Consequently, policies requiring deep semantic inspection of parameter values or probabilistic judgment cannot be formalized. Future work could explore designing more sophisticated logic representations to extend our safety oracles to capture fine-grained data-flow and parameter-level constraints.

**Dependence on Generator Capability.** Although LOGISAFETYGEN automates the synthesis of valid traces via fuzzing, the initial diversity of the test scenarios is partially bounded by the creative priors of the backbone LLM used for seeding. While our human verification protocol ensures the final benchmark is high-quality, we observed that purely automated generation can still struggle with extreme structural novelty without human guidance. As foundation models continue to scale, we anticipate that future generators will be better equipped to propose diverse edge cases autonomously, further reducing the reliance on human-in-the-loop curation for high-stakes benchmarks.

## Ethical Considerations

Given that our paper aims to unveil the implicit compliance risks of LLMs, our publicly available dataset includes test cases that simulate regulatory violations across financial, healthcare, and IoT domains. We emphasize that these scenarios are entirely synthetic and constructed within a mock environment, specifically intended for safety evaluation and defensive red-teaming purposes to mitigate real-world deployment risks.

## References

Maksym Andriushchenko, Alexandra Souly, Mateusz Dziemian, Derek Duenas, Maxwell Lin, Justin Wang, Dan Hendrycks, Andy Zou, J Zico Kolter, Matt Fredrikson, Yarin Gal, and Xander Davies. 2025. Agentharm: A benchmark for measuring harmfulness of LLM agents. In *The Thirteenth International Conference on Learning Representations*.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code.

Simin Chen, Yiming Chen, Zexin Li, Yifan Jiang, Zhongwei Wan, Yixin He, Dezhi Ran, Tianle Gu, Haizhou Li, Tao Xie, and Baishakhi Ray. 2025a. Benchmarking large language models under data contamination: A survey from static to dynamic evaluation. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 10091–10109, Suzhou, China. Association for Computational Linguistics.

Zhaorun Chen, Mintong Kang, and Bo Li. 2025b. Shieldagent: Shielding agents via verifiable safety policy reasoning. *arXiv preprint arXiv:2503.22738*.

Jason Chou, Ao Liu, Yuchi Deng, Zhiying Zeng, Tao Zhang, Haotian Zhu, Jianwei Cai, Yue Mao, Chenchen Zhang, Lingyun Tan, and 1 others. 2025. Autocodebench: Large language models are automatic code benchmark generators. *arXiv preprint arXiv:2508.09101*.

Jinru Ding, Chao Ding, Wenrao Pang, Boyi Xiao, Zhiqiang Liu, Pengcheng Chen, Jiayuan Chen, Tiantian Yuan, Junming Guan, Yidong Jiang, Dawei Cheng, and Jie Xu. 2025. Cnfinbench: A benchmark for safety and compliance of large language models in finance. *Preprint*, arXiv:2512.09506.

ETSI. 2024. Cyber Security for Consumer Internet of Things: Baseline Threats and Mitigation Measures. urlwww.etsi.org.

EU. 2015. Payment services in the internal market, amending directives 2002/65/ec, 2009/110/ec and 2013/36/eu and regulation (eu) no 1093/2010, and repealing directive 2007/64/ec.

EU. 2024. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence (the

EU AI Act). Official Journal of the European Union, L 2024/1689.

Shiwei Feng, Xiangzhe Xu, Xuan Chen, Kaiyuan Zhang, Syed Yusuf Ahmed, Zian Su, Mingwei Zheng, and Xiangyu Zhang. 2025. TAI3: Testing agent integrity in interpreting user intent. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.

Gayane Ghazaryan, Erik Arakelyan, Isabelle Augenstein, and Pasquale Minervini. 2025. Syndarin: Synthesising datasets for automated reasoning in low-resource languages. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6459–6466.

HHS. 1996. HIPAA. Online at www.cms.hhs.gov. Public Law 104-191, 110 Stat. 1936.

Yuheng Huang, Da Song, Zhenlan Ji, Shuai Wang, and Lei Ma. 2025. Evaluating llms on sequential api call through automated test generation. *arXiv preprint arXiv:2507.09481*.

Taewhoo Lee, Chanwoong Yoon, Kyochul Jang, Donghyeon Lee, Minju Song, Hyunjae Kim, and Jaewoo Kang. 2025. Ethic: Evaluating large language models on long-context tasks with high information coverage. In *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)*, pages 5497–5512.

Ilija Lichkovski, Alexander Müller, Mariam Ibrahim, and Tiwai Mhundwa. 2025. Eu-agent-bench: Measuring illegal behavior of llm agents under eu law. In *Proceedings of the Workshop on Regulatable Machine Learning at the 39th Conference on Neural Information Processing Systems (NeurIPS 2025)*.

Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, and 3 others. 2024. Agentbench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning Representations*.

Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. 2024. Harmbench: a standardized evaluation framework for automated red teaming and robust refusal. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.

Mahmoud Mohammadi, Yipeng Li, Jane Lo, and Wendy Yip. 2025. Evaluation and benchmarking of llm agents: A survey. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 6129–6139.

Jiayi Pan, Glen Chou, and Dmitry Berenson. 2023. Data-efficient learning of natural language to linear temporal logic translators for robot task specification. *arXiv preprint arXiv:2303.08006*.

Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. 2025. The berkeley function calling leaderboard (BFCL): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*.

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. 2024. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. In *The Twelfth International Conference on Learning Representations*.

Brandon Radosevich and John Halloran. 2025. Mcp safety audit: Llms with the model context protocol allow major security exploits. *arXiv preprint arXiv:2504.03767*.

Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J. Maddison, and Tatsunori Hashimoto. 2024. Identifying the risks of LM agents with an LM-emulated sandbox. In *The Twelfth International Conference on Learning Representations*.

Sarath Sivaprasad, Pramod Kaushik, Sahar Abdelnabi, and Mario Fritz. 2025. A theory of response sampling in llms: Part descriptive and part prescriptive. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 30091–30135.

Jianwei Wang, Mengqi Wang, Yinsi Zhou, Zhenchang Xing, Qing Liu, Xiwei Xu, Wenjie Zhang, and Liming Zhu. 2025a. Llm-based hse compliance assessment: Benchmark, performance, and advancements. *arXiv preprint arXiv:2505.22959*.

Wenhan Wang, Chenyuan Yang, Zhijie Wang, Yuheng Huang, Zhaoyang Chu, Da Song, Lingming Zhang, An Ran Chen, and Lei Ma. 2025b. Testeval: Benchmarking large language models for test case generation. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 3547–3562.

Zhenting Wang, Qi Chang, Hemani Patel, Shashank Biju, Cheng-En Wu, Quan Liu, Aolin Ding, Alireza Rezazadeh, Ankit Shah, Yujia Bao, and 1 others. 2025c. Mcp-bench: Benchmarking tool-using llm agents with complex real-world tasks via mcp servers. *arXiv preprint arXiv:2508.20453*.

Chunqiu Steven Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, and Lingming Zhang. 2024. Fuzz4all: Universal fuzzing with large language models. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pages 1–13.

Chenghao Yang and Ari Holtzman. 2025. How alignment shrinks the generative horizon. *arXiv preprint arXiv:2506.17871*.

Jihan Yao, Peter Jin, Ke Bao, Qiaolin Yu, Khushi Bhardwaj, Chang Su, Jialei Wang, Yikai Zhu, Sugam Devare, Damon Mosk-Aoyama, and 1 others. The measure of all measures: Quantifying llm benchmark quality. In *NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling*.

Junjie Ye, Sixian Li, Guanyu Li, Caishuang Huang, Songyang Gao, Yilong Wu, Qi Zhang, Tao Gui, and Xuan-Jing Huang. 2024. Toolsword: Unveiling safety issues of large language models in tool learning across three stages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2181–2211.

Zhehao Zhang, Jiaao Chen, and Diyi Yang. 2024a. Darg: Dynamic evaluation of large language models via adaptive reasoning graph. *Advances in Neural Information Processing Systems*, 37:135904–135942.

Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. 2024b. Agent-safetybench: Evaluating the safety of llm agents. *arXiv preprint arXiv:2412.14470*.

Ruiwen Zhou, Wenyue Hua, Liangming Pan, Sitao Cheng, Xiaobao Wu, En Yu, and William Yang Wang. 2025. Rulearena: A benchmark for rule-guided reasoning with llms in real-world scenarios. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 550–572.

Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. 2024. Dynamic evaluation of large language models by meta probing agents. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.

Xiaogang Zhu, Sheng Wen, Seyit Camtepe, and Yang Xiang. 2022. Fuzzing: a survey for roadmap. *ACM Computing Surveys (CSUR)*, 54(11s):1–36.