

# Beyond Superficial Tests: Adversarial Refinement for Reliable Property-Based Testing

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have demonstrated remarkable proficiency in code generation, yet their application to Property-Based Testing (PBT) remains fraught with a “superficiality gap”. While LLMs can readily generate syntactically correct tests, they often struggle to bridge the semantic gap between code implementation and its intended invariant logic, resulting in weak properties that provide a false sense of security. To address this, we introduce PROBE, an agentic framework that hardens software properties through *Adversarial Refinement*. Unlike traditional generation approaches, PROBE treats test generation as a game of “semantic asymmetry”: it employs a Validator agent to actively generate counter-implementations, which are semantically incorrect codes that satisfy the generated property, to expose loopholes in the specification. Furthermore, PROBE constructs a cross-functional semantic graph to capture deep dependencies often missed by local analysis. Extensive evaluation reveals that PROBE increases mutation scores by 9.79% over baselines. In real-world deployment, PROBE identified 45 previously unknown bugs in top-tier libraries that have been confirmed by developers, demonstrating its ability to uncover deep semantic defects.

## 1 Introduction

Property-Based Testing (PBT) has emerged as a rigorous paradigm for ensuring software correctness (Fink and Bishop, 1997). Unlike Example-Based Testing (EBT) which verifies discrete input-output pairs (Daka and Fraser, 2014), PBT validates programs against executable *properties*—universal invariants that must hold across the entire input domain. Since its inception with QuickCheck (Claessen and Hughes, 2000), this approach has proven uniquely capable of uncovering deep corner-case defects. However, despite its theoretical superiority, PBT remains underuti-

lized in practice due to a substantial **cognitive barrier**: defining non-trivial properties and implementing effective input generators are notoriously difficult and time-consuming tasks for human developers (Goldstein et al., 2024).

The advent of Large Language Models (LLMs) promised to automate this labor-intensive process, given their proficiency in code synthesis (Islam et al., 2024; Liu et al., 2024). Yet, current approaches to automated PBT generation face a critical “**superficiality gap**”. While LLMs can readily generate syntactically correct test scripts, they struggle to bridge the semantic distance between implementation details and high-level behavioral contracts. Recent studies reveal two distinct failure modes: (1) **Low Executability**: Even with optimized prompting, over 50% of generated tests fail to compile or execute due to hallucinations (Vikram et al., 2024). (2) **Semantic Triviality**: More insidiously, even valid tests often enforce weak properties (e.g., generic type checks) that provide a false sense of security without exercising the target function’s core logic. This limitation stems from the lack of **semantic grounding**—general-purpose models treat PBT as a localized translation task, ignoring the broader logical dependencies (e.g., cross-function contracts) required for rigorous specification.

To bridge this gap, we propose PROBE, the first open-source agentic framework explicitly designed to **harden** properties via **Adversarial Refinement**. Unlike traditional methods that treat PBT as a one-shot translation, PROBE frames it as a game of *semantic asymmetry*. We leverage the insight that **construction is harder than falsification**: while synthesizing a comprehensive property requires complex *logic abstraction* (a high-difficulty task), identifying a loophole via a specific *counter-implementation* is a much simpler *concrete validation* task. PROBE exploits this asymmetry by employing a Validator agent that actively seeks to

“break” the generated properties by constructing degenerate implementations—semantically incorrect code but satisfies the generated property—to expose loopholes in the generated weak property. These discovered *weaknesses* are then fed back to a Generator agent to iteratively refine and strengthen the property.

Furthermore, PROBE addresses the limitations of context-isolation through two key innovations: (1) **Semantic Planning:** It constructs a *Cross-function Semantic Graph* to retrieve logical dependencies often missed by local analysis, ensuring properties adhere to global function contracts (e.g., round-trip consistency between encode/decode). (2) **Contextual Grounding:** It utilizes AST-based static analysis to derive physical input constraints, providing a “reasoning scaffolding” that prevents the agent from exploring invalid search domains.

Extensive evaluations demonstrate that PROBE consistently outperforms state-of-the-art baselines, improving mutation scores (Just et al., 2014) by up to **9.79%** while maintaining robust performance gains across diverse LLM backbones. Manual analysis reveals that PROBE not only bridges the semantic gap with a **95% property correctness rate** (compared to 65% for standalone models) but also exhibits superior discriminative power, uniquely uncovering subtle logic flaws overlooked by other approaches. Most notably, in a large-scale deployment on 21 high-impact repositories, PROBE successfully identified **45 previously unknown bugs** in foundation-level libraries such as CPython, scipy, and cryptography, with 32 already fixed by developers. To the best of our knowledge, PROBE is the first open-source framework that treats property strength as an explicit optimization objective through adversarial agentic games.

Our contributions are summarized as follows:

- We present PROBE, a novel agentic framework to automate whole PBT lifecycle, from constraint grounding to adversarial property hardening.
- We introduce *Adversarial Refinement* and *Semantic Planning*, mechanisms that exploit semantic asymmetry and cross-function dependencies to synthesize high-strength properties.
- We provide an extensive empirical evaluation and demonstrate practical impact by detecting 45 previously unknown bugs (with 32 merged fixes) in widely-used Python libraries.

## 2 Background and Related Work

This section outlines the theoretical foundations of Property-Based Testing and surveys recent advancements in Large Language Model (LLM) driven test generation.

### 2.1 Property-Based Testing

Property-Based Testing (PBT) is a dynamic verification framework that validates programs against executable specifications, known as *properties* (Claessen and Hughes, 2000; Goldstein et al., 2024). Unlike example-based testing, which verifies specific input-output pairs, PBT employs strategy-based generators to sample diverse inputs  $x$  from a domain  $\mathcal{X}$ . A property  $\varphi$  is a predicate over the execution trace:

$$\varphi : (x, f(x)) \rightarrow \{\text{true}, \text{false}\} \quad (1)$$

We say a function  $f$  satisfies  $\varphi$  if :

$$\forall x \in \mathcal{X} : \varphi(x, f(x)) = \text{true} \quad (2)$$

When a violation is detected, PBT frameworks typically perform *shrinking* to minimize the failing input into a concise counterexample (MacIver and Donaldson, 2020).

A property represents a *necessary condition* for correctness. For example, a correct `sort` function must satisfy both monotonicity ( $\varphi_{\text{ord}}$ ) and permutation ( $\varphi_{\text{perm}}$ ). As shown in Figure 2, each property alone is insufficient. A function returning a constant list satisfies  $\varphi_{\text{ord}}$  but violates  $\varphi_{\text{perm}}$ , while reversing the input satisfies  $\varphi_{\text{perm}}$  but violates  $\varphi_{\text{ord}}$ . Only their conjunction  $\varphi_{\text{ord}} \wedge \varphi_{\text{perm}}$  adequately constrains the behavior.

To formalize property effectiveness, we introduce the notion of *behavior sets*. A *behavior* is an input-output pair  $(x, f(x))$  representing one possible execution. Let  $\mathcal{S}$  denote the set of correct behaviors and  $\mathcal{A}_\varphi$  denote the behaviors accepted by  $\varphi$ . A valid property requires  $\mathcal{S} \subseteq \mathcal{A}_\varphi$ . Our framework leverages the principle that conjunction strengthens properties:  $\mathcal{A}_{\varphi_1 \wedge \varphi_2} = \mathcal{A}_{\varphi_1} \cap \mathcal{A}_{\varphi_2}$ . By synthesizing and combining multiple properties, we progressively narrow the accepted behavior set  $\mathcal{A}_\varphi$  towards  $\mathcal{S}$ , leaving minimal room for incorrect implementations to pass. In this work, we utilize Hypothesis (MacIver et al., 2019), a mature Python PBT framework, as our execution backend. Figure 3 illustrates how to test `sort` function by using Hypothesis.

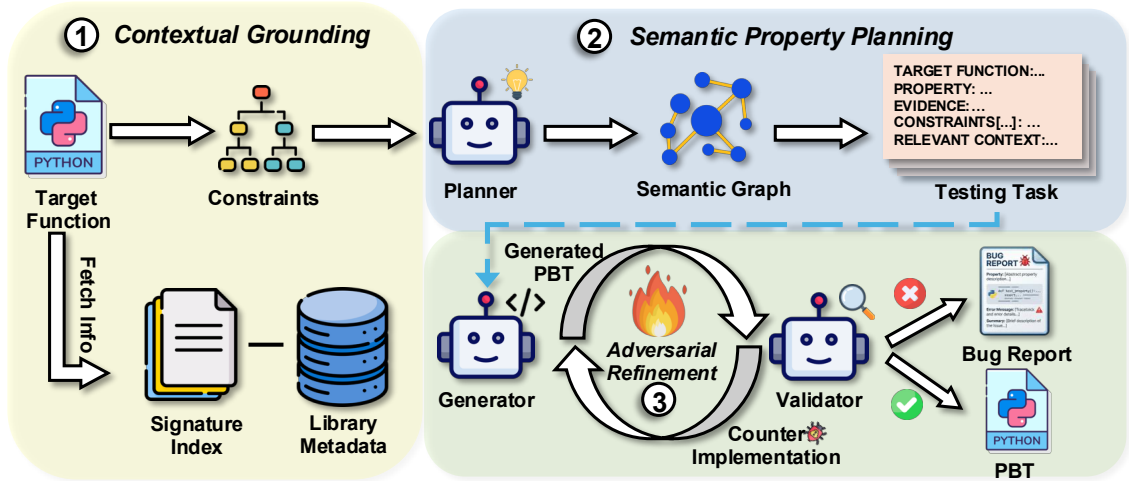


Figure 1: Overall structure of PROBE.

```
def sort(x: list) -> list:
    return [i for i in range(len(x))]

def sort(x: list) -> list:
    return x[::-1]
```

Figure 2: Incorrect implementations of sort.

```
from collections import Counter
from hypothesis import given, strategies as st
@given(st.lists(st.integers()))
def test_sort(arr):
    sorted_arr = sorted(arr)
    for i in range(1, len(arr)):
        if sorted_arr[i-1] > sorted_arr[i]:
            assert False
    assert Counter(sorted_arr) == Counter(arr)
```

Figure 3: An example of a property-based test using Hypothesis framework.

## 2.2 LLM-based Test Generation

Recent research has explored leveraging LLMs for automated software testing (Chen et al., 2024; Yang et al., 2024; Zhang et al., 2025). Approaches such as MuTAP (Dakhel et al., 2024) utilize mutation testing to guide LLM generation, while Test4Py (Liu et al., 2025) infers parameter types to construct valid inputs. Frameworks like AEGIS (Wang et al., 2024) further employ agent-based architectures with multi-feedback optimization to generate reproduction scripts.

Despite these successes in unit testing, automating property-based testing remains challenging due to the complexity of specifying semantic invariants. Heuristic tools like Hypothesis ghostwriter (Hypothesis team) provide syntactic skeletons but still necessitate significant manual effort to fill in se-

semantic logic. Recent works have attempted to leverage LLMs to reduce this burden: Vikram et al. (2024) explored directly prompting LLMs to translate documentation into PBTs, while Maaz et al. (2025) applied general-purpose commercial agents to existing frameworks. However, these approaches largely rely on the intrinsic coding capabilities of LLMs without task-specific optimization. In contrast, PROBE introduces an approach that combines global context with adversarial refinement to generate more rigorous PBTs.

## 3 Methodology

In this section, we present PROBE, an agentic framework for automatically generating high quality PBTs. Unlike single-pass generation methods, PROBE orchestrates a collaborative workflow among specialized agents to produce PBTs that are both syntactically valid and semantically robust.

As illustrated in Figure 1, the PROBE comprises three stages: (1) **Contextual Grounding**, employing AST-based static analysis to derive physical constraints defining valid input spaces; (2) **Semantic Property Planning**, which contextualizes the target function within a repository-wide semantic graph to infer implicit properties; and (3) **Adversarial Refinement**, where the Generator and Validator engage in a minimax game to iteratively harden PBT strength by synthesizing and falsifying counter-implementations. This workflow produces either a semantically robust PBT or a bug report grounded in rigorous evidence.

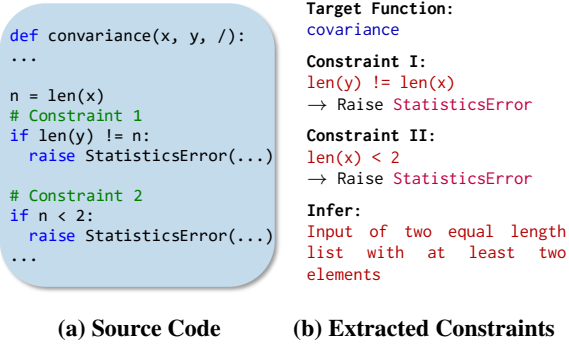


Figure 4: Static constraints extraction of covariance from statistics package.

### 3.1 Contextual Grounding

To avoid trivial inputs that merely trigger intentional error handlers (e.g., `ValueError`), PROBE implements **Contextual Grounding** to derive physical constraints directly from the source code. The system first analyzes the Abstract Syntax Tree (AST) to identify *guard clauses* (e.g., `raise` or `assert`). Instead of treating these as isolated code snippets, PROBE resolves them into **Physical Constraints** through backward data-flow analysis. As illustrated in Figure 4, it maps local variables back to inputs (e.g., converting `if n < 2` to `len(x) >= 2`). These constraints prune the input search space, establishing a “logical safety zone” that allows the agent to focus on inferring implicit semantic guarantees rather than satisfying basic preconditions.

### 3.2 Semantic Property Planning

The primary challenge in automated PBT generation lies in the *semantic isolation* of LLMs, which often struggle to infer invariant logic by analyzing a single function in a vacuum. To bridge this gap, PROBE treats PBT generation not as a localized coding task, but as a systematic planning process over a **Cross-function Semantic Graph**  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . Here, nodes  $\mathcal{V}$  represent functions within the repository, and edges  $\mathcal{E}$  encode implicit logical contracts and behavioral dependencies.

The Planner agent orchestrates a three-stage workflow to construct this graph and synthesize high-strength properties:

#### 1 Relational Neighborhood Identification

Rather than ingesting the entire repository, which introduces excessive noise, the Planner first isolates the *logical neighborhood* of the target function  $f_{target}$ . It traverses the import dependencies and class hierarchies to isolate a subset of modules  $\mathcal{M}_{rel}$ . This minimizes noise and establishes strict

boundaries for relevant context retrieval.

**2 Contractual Edge Discovery** The core innovation of the Planner lies in its ability to identify **Implicit Contracts** — behavioral consistencies that span multiple functions. Using the signature index  $\mathcal{I}$ , the Planner iteratively scans  $\mathcal{M}_{rel}$  to identify functions that share relational signatures with  $f_{target}$ . Specifically, it searches for three types of semantic edges:

- **Inverse Operations:** Identifying pairs such as encode/decode, implying the round-trip consistency ( $f_{inv}(f_{target}(x)) = x$ ).
- **State Invariants:** Identifying operations that should preserve a global state, such as push/pop in a stack or add/remove in a container.
- **Relational Identities:** Identifying logical properties between distinct functions that share underlying semantics. For example, the standard deviation function should yield a result consistent with the square root of the variance function.

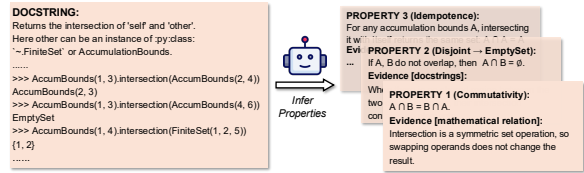


Figure 5: Example of evidence-grounded property.

To mitigate hallucination, PROBE enforces **Evidence-Grounded Inference**. As demonstrated in Figure 5, the Planner is strictly required to map every synthesized property  $\varphi$  to explicit evidence  $e$  (e.g., docstrings or algebraic identities) extracted during the discovery phase.

The result is consolidated into a formal **Testing Task**  $t = \langle f_{target}, ic, \varphi, c_\varphi, e \rangle$ , serving as a traceable blueprint for Algorithm 1. Here,  $ic$  denotes the physical input constraints derived in Section § 3.1 and  $c_\varphi$  the enriched context.

### 3.3 Adversarial Refinement

Given a testing task  $t$ , the Generator synthesizes a PBT that encodes the specified property  $\varphi$ . Initially, the generated PBT is executed in an isolated environment. If the execution fails, the Validator performs root cause analysis based on the execution logs and  $t$ , classifying failures into three categories:

- **Test Code Defects:** Syntax errors or incorrect function usage.
- **Property Design Defects:** semantic misalignment between the property and evidence, or discrepancies between the evidence and docstrings.

---

**Algorithm 1** Generator-Validator Workflow

---

**Require:** Testing Task  $t = \langle f_{target}, ic, \varphi, c_\varphi, e \rangle$   
**Ensure:** Hardened PBT or Bug Report

```
1:  $PBT \leftarrow \text{GENERATOR.GENPBT}(t)$ 
2:  $Env \leftarrow \text{ENV.CREATE}(t)$ 
3:  $cnt_{fix}, cnt_{ref} \leftarrow 0, 0$   $\triangleright$  Counters for  $T_{fix}$  and  $T_{ref}$ 
4: while  $cnt_{fix} \leq T_{fix}$  do
5:    $res \leftarrow \text{RUN}(PBT, f_{target}, Env)$ 
6:   if  $res.status == \text{PASS}$  then
7:     if  $cnt_{ref} == T_{ref}$  then
8:       return  $PBT$ 
9:      $f' \leftarrow \text{VALIDATOR.GENADVERSARY}(PBT, t)$ 
10:    if  $f' \neq \emptyset$  then
11:       $res' \leftarrow \text{RUN}(PBT, f', Env)$ 
12:      if  $res'.status == \text{PASS}$  then
13:         $t \leftarrow \text{PLANNER.REFINE}(PBT, f', t)$ 
14:         $PBT \leftarrow \text{GENERATOR.GENPBT}(t)$ 
15:         $cnt_{ref} \leftarrow cnt_{ref} + 1; cnt_{fix} \leftarrow 0$ 
16:      else
17:        return  $PBT$ 
18:      else
19:        return  $PBT$ 
20:    else
21:       $cause, info \leftarrow$ 
22:       $\text{VALIDATOR.DIAGNOSE}(res, PBT, t)$ 
23:      if  $cause == \text{CODE\_DEFECT}$  then
24:         $PBT \leftarrow \text{FIXCODE}(PBT, info)$ 
25:      else if  $cause == \text{LIB\_DEFECT}$  then
26:        return  $\text{CREATEBUGREPORT}(t, info)$ 
27:      else if  $cause == \text{PROP\_DEFECT}$  then
28:         $t \leftarrow \text{PLANNER.REPLAN}(t, info)$ 
29:         $PBT \leftarrow \text{GENERATOR.GENPBT}(t)$ 
30:       $cnt_{fix} \leftarrow cnt_{fix} + 1$ 
31: return  $PBT$ 
```

---

- *Library Defects:* The property is valid and the test code is correctly implemented, yet execution fails. This triggers an immediate bug report.

The diagnostic loop is bounded by a limit  $T_{fix}$ . However, achieving a passing state in this loop merely guarantees executability. It does not ensure the property’s rigor.

Transitioning from correctness to quality, we address a fundamental limitation in PBT generation known as the *propensity for triviality*: LLMs often converge on weak properties that are syntactically valid but semantically vacuous. To address this, PROBE frames property refinement as a min-max game predicated on **Difficulty Asymmetry** — the principle that verifying or falsifying a solution is often significantly easier than generating it (Wei, 2024). We posit that while synthesizing a globally optimal property is a high-entropy inductive challenge, identifying a local loophole via a *counter-implementation* is a significantly more tractable deductive task.

**The Adversarial Objective** Upon receiving a passing PBT, the Validator assumes the role of a *Se-*

*mantic Adversary* to challenge the sufficiency of  $\varphi$ . Instead of verifying  $\varphi$ , it actively seeks to construct a *counter-implementation*  $f'$ —a code variant that is *syntactically plausible* and satisfies the current property  $\varphi$ , yet is *semantically incorrect* (e.g., dropping the last element of the output list for a sort function). This process exploits the inherent asymmetry between *construction* and *breaking*. Constructing  $f'$  provides the system with a *concrete falsification witness*. When  $f'$  successfully “cheats” the current  $\varphi$ , it yields a high-signal feedback message for the Generator: “The current property is too weak to distinguish intended behavior from this specific degenerate  $f'$ .” This evidence-driven feedback is more effective as it anchors the LLM’s reasoning in a tangible failure case.

As illustrated in Algorithm 1, this adversarial process is iterative. Each successful  $f'$  acts as a high-signal feedback, forcing the Generator to tighten the PBT with more rigorous properties (e.g., moving from “output is sorted” to “output is a permutation of the input and is sorted”). The loop continues until the Validator fails to synthesize a property-passing  $f'$  within  $T_{ref}$  attempts, effectively collapsing the accepted behavior set toward the ground truth intended behaviors.

## 4 Experiment

To comprehensively evaluate the performance and practical value of PROBE, we design experiments to answer the following research questions:

- **RQ-1 Effectiveness:** How effectively can PROBE generate non-trivial PBTs that distinguish correct implementations from erroneous ones compared to baselines.
- **RQ-2 Validity & Distinctiveness:** What proportion of generated PBTs are syntactically and semantically correct, and how distinctive are the properties discovered by each tool?
- **RQ-3 Bug Finding:** Can PROBE discover previously unknown bugs in widely-used libraries?
- **RQ-4 Ablation Study:** How does each component contribute to the overall performance of PROBE?

### 4.1 Experiment Setup

**Baseline Methods** To our knowledge, current methods for automated PBT generation primarily utilize LLM through two main paradigms: employing an LLM via direct prompting (Vikram et al., 2024) or leveraging a commercial agent (Maaz

et al., 2025). To thoroughly evaluate PROBE, we categorize our baselines into three configurations:

- **LLM:** We select several leading LLMs to serve as our primary baseline, including DeepSeek-V3.2 (DeepSeek-AI et al., 2025), GPT-5 (OpenAI, 2025) and Qwen3-Next-80B-A3B-Thinking (QwenTeam, 2025a). This category assesses the inherent reasoning and code generation capabilities of LLMs.
- **Retrieval-Augmented LLM (LLM+RAG):** We implement the second baseline by equipping the LLM with Retrieval-Augmented Generation (RAG). This setup provides the model with a more precise context of the repository.
- **LLM-Based Agent:** Finally, we include an LLM-based agent to represent a method optimized for complex software tasks. Claude Code (Anthropic, 2025a) has been explored to generate PBTs. We select it with claude sonnet 4.5 (Anthropic, 2025b) as our third baseline, which features an state-of-art agent designed for code generation beyond simple prompting.

Baselines details are provided in Appendix B.4.

**Datasets** We construct two datasets for evaluation. *Dataset I* serves as the primary testbed for RQ1, RQ2, and RQ4. We derived 256 non-trivial functions from four diverse PyPI packages: `sympy` (mathematics), `sortedcontainers` (data structures), `more-itertools` (utilities), and `simplejson` (serialization). These libraries span various domains, reflecting the diversity of properties encountered in practice. To ensure meaningful evaluation, we apply a systematic filtering pipeline based on code complexity metrics (detailed in Appendix B.2). *Dataset II* is curated for RQ3, which encompasses 21 high-impact repositories, including 8 Python Standard Library modules and 13 prominent PyPI packages. Within each package, we manually identified functions characterized by intricate semantic logic.

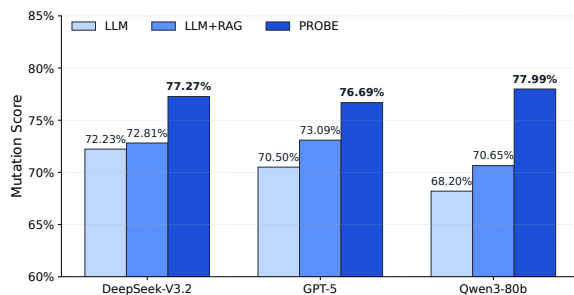


Figure 6: Performance comparison of PROBE against baselines across different backbone models.

## 4.2 RQ-1: Effectiveness

To quantify PBT strength, we employ mutation testing (DeMillo et al., 1978), a standard method for evaluating fault-detection capability of a test suite. The core principle of mutation testing is to systematically inject artificial faults into the original program using predefined mutation operators, producing modified variants known as *mutants*. Previous studies have demonstrated a statistically significant correlation between mutant detection and real fault detection (Just et al., 2014; Ravi and Coblenz, 2025). A mutant is considered *killed* if the PBT passes on the original implementation but fails on the mutated version. Higher mutation scores indicate stronger discriminative power of tests.

To ensure evaluation independence, we utilized the MUTMUT tool (Mutmut contributors, 2025) to generate up to 10 mutants per function. We compute mutation score exclusively on the *intersection* of functions where all evaluated tools produce PBTs that pass on the original implementation. This intersection-based protocol eliminates potential bias arising from tools that generate valid PBTs on different subsets of functions. For completeness, we report the mutation scores per-tool on the full dataset in Appendix B.1.

**Results.** Table 1 presents the mutation scores on a set of 82 functions (695 mutants). The best performance in each category is highlighted in bold. Among all evaluated configurations, PROBE (Qwen3-80b) achieves the best overall performance, attaining the highest mutation score of **77.99%** by killing 542 out of 695 mutants. This represents a 9.79 percentage point improvement over the vanilla Qwen3-80b baseline.

This performance advantage is consistent across architectures. As illustrated in Figure 6, when applied to DeepSeek-V3.2 and GPT-5, PROBE boosts mutation scores to **77.27%** and **76.69%**, respectively. Crucially, PROBE outperforms Claude Code (Sonnet 4.5) by margins of 5.90~7.20 percentage points. These results suggests that our domain-specific structured approach is more effective than general-purpose methods for PBT generation.

## 4.3 RQ-2: Validity & Distinctiveness

To evaluate the semantic validity of generated PBTs, we manually analyzed PBTs in representative top-performers in RQ-1.

We conduct a systematic manual analysis to evaluate the **validity** of PBTs. Specifically, three senior



Table 3: Summary of defects discovered by PROBE. *Rep.* denotes the total issues reported. *Conf.* indicates bugs confirmed by developers. *Fixed* denotes issues resolved with merged patches. *Doc.* refers to documentation updates instead of code changes.

Repository (Star Count)	Rep.	Conf.	Fixed	Doc.
<b>Python Standard Library (70.5k)</b>				
json / html / re	5	4	3	1
functools	1	1	1	1
collections	2	1	1	0
statistics	2	1	1	0
shlex	1	1	0	0
types	1	1	0	0
<b>PyPI Packages</b>				
boto / boto3 (9.6k)	1	0	0	0
marshmallow (7.2k)	3	2	1	0
scipy / scipy (14.3k)	4	4	2	0
sympy / sympy (14.2k)	13	13	9	0
pyproj4 / pyproj (1.2k)	2	1	0	1
tobgu / pyrsistent (2.2k)	1	0	0	0
cpburnz / pathspec (205)	3	1	1	0
pyca / cryptography (7.4k)	2	2	2	0
scikit-hep / awkward (929)	8	6	4	0
networkx / networkx (16.4k)	1	1	1	0
aws-lambda-powertools (3.2k)	3	3	3	2
google-deepmind / optax (2.1k)	2	2	2	0
huggingface / tokenizers (10.3k)	1	1	1	0
<b>Total</b>	<b>56</b>	<b>45</b>	<b>32</b>	<b>5</b>

**Results.** As summarized in Table 3, we have submitted 56 issues, of which 45 were confirmed by the developers with 32 fixes. Notably, all confirmed defects were previously unknown, demonstrating PROBE’s ability to detect subtle semantic bugs. Among the confirmed issues, 5 were resolved through documentation amendments rather than code changes, as developers attempted to preserve backward compatibility for existing users. Details of all reported defects are provided in Appendix C. The remaining raw reports, which involve highly specialized domain logic beyond our immediate expertise, will be released on our repository to invite community adjudication.

Table 4: *w/o SeGra.* and *w/o Adv.* exclude the Semantic Graph and Adversarial Refinement, respectively. Red percentages denote performance drops.

Backbone Model	Full System	w/o SeGra.	w/o Adv.
<b>DeepSeek-V3.2</b> (680B)	<b>83.57</b>	73.43 (-10.1%)	76.78 (-6.8%)
<b>Qwen3-80B</b> (80B)	<b>81.93</b>	70.36 (-11.6%)	77.69 (-4.2%)
<b>GPT-5</b> (Closed)	<b>82.23</b>	74.78 (-7.5%)	77.01 (-5.2%)

Table 5: Performance of PROBE across different LLM backbones. ( $\Delta$ ) denotes the improvement compared to the corresponding vanilla LLM.

Backbone LLM	Size	Mut. Score	$\Delta$
Qwen3-80b	80B	77.99%	+9.79%
DeepSeek-V3.2	685B	77.27%	+5.04%
GPT-5	-	76.69%	+6.19%

#### 4.5 RQ-4: Ablation Study

To quantify the impact of PROBE’s core components, we performed an ablation study by disabling semantic graph and adversarial refinement.

**Results.** Results in Table 4 are reported on the intersection of functions where all variants successfully generated passing PBTs. Semantic graph proves to be a critical module. Its removal precipitates the performance decline with drops ranging from 7.5% to 11.6% across models. This underscores that retrieving cross-function dependencies is essential for strengthening PBTs. Adversarial refinement also plays an indispensable role in hardening properties. Disabling it reduces mutation scores by 4.2% to 6.8%, which confirms that the adversarial loop effectively filters out trivial properties that fail to detect subtle bugs.

Crucially, the relative impact of these components remains consistent across different LLMs. As shown in Table 5, the full PROBE system consistently outperforms its vanilla LLM regardless of the backbone LLMs. This suggests that the performance gains are intrinsic to our structured framework design rather than being an artifact of a specific underlying LLM’s capability.

## 5 Conclusion

In this paper, we introduce PROBE, the first open-source agentic framework designed to automate the lifecycle of Property-Based Testing. PROBE addresses the critical limitations of existing methods. Our approach ensures that generated PBTs are not only executable but also able to capture complex semantic logic with cross-function properties. Extensive evaluations demonstrate that PROBE significantly outperforms the baselines in mutation scores, which achieve up to 9.79%. Notably, PROBE discovered 45 previously unknown bugs from top-tier Python libraries, with 32 already fixed and merged. In addition, the ablation study shows that each component of PROBE has contribution to the overall performance.

## 572 Limitations

573 We acknowledge several limitations inherent to  
574 PROBE. First, agentic architectures that priori-  
575 tize test quality through iterative refinement inher-  
576 ently involve more computation than pure LLM  
577 generation methods, which is a trade-off between  
578 test robustness and inference efficiency common  
579 to all refinement-based approaches. We mitigate  
580 this overhead through static analysis to constrain  
581 the search space and modular agent design to min-  
582 imize context sizes, though further optimization  
583 remains an avenue for future work. Second, the  
584 effectiveness of LLM-based testing frameworks is  
585 partially bounded by the reasoning capabilities of  
586 the underlying language models. We address it  
587 by designing targeted prompts and decomposing  
588 complex tasks across specialized agents to reduce  
589 the burden on any single inference call. Finally,  
590 a fundamental challenge in automated PBT gen-  
591 eration—shared across existing methods—is dis-  
592 tinguishing genuine specification violations from  
593 intentional design. Properties derived from doc-  
594 strings may not always align with developer’s in-  
595 tent, occasionally yielding false positives. Fully  
596 resolving this challenge remains an open problem.

## 597 References

- 598 Anthropic. 2025a. Claude code. <https://www.anthropic.com/claude-code>. Powered by Claude  
599 Sonnet 4.5.  
600
- 601 Anthropic. 2025b. Introducing claude sonnet  
602 4.5. <https://www.anthropic.com/news/claude-sonnet4-5>.  
603
- 604 G. Ann Campbell. 2018. [Cognitive complexity: an  
605 overview and evaluation](#). In *Proceedings of the  
606 2018 International Conference on Technical Debt*,  
607 TechDebt ’18, page 57–58, New York, NY, USA.  
608 Association for Computing Machinery.
- 609 Yinghao Chen, Zehao Hu, Chen Zhi, Junxiao Han,  
610 Shuiguang Deng, and Jianwei Yin. 2024. Chatu-  
611 nitest: A framework for llm-based test generation.  
612 In *Companion Proceedings of the 32nd ACM Inter-  
613 national Conference on the Foundations of Software  
614 Engineering*, pages 572–576.
- 615 Koen Claessen and John Hughes. 2000. [Quickcheck:  
616 a lightweight tool for random testing of haskell pro-  
617 grams](#). In *Proceedings of the Fifth ACM SIGPLAN  
618 International Conference on Functional Program-  
619 ming*, ICFP ’00, page 268–279, New York, NY, USA.  
620 Association for Computing Machinery.
- 621 Ermira Daka and Gordon Fraser. 2014. [A survey on  
622 unit testing practices and problems](#). In *2014 IEEE*

*25th International Symposium on Software Reliabil-  
ity Engineering*, pages 201–211. 623  
624

- 625 Arghavan Moradi Dakhel, Amin Nikanjam, Vahid Ma-  
626 jdinassab, Foutse Khomh, and Michel C Desmarais.  
627 2024. Effective test generation using pre-trained  
628 large language models and mutation testing. *Infor-  
629 mation and Software Technology*, 171:107468.
- 630 DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin,  
631 Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao  
632 Wu, Bowei Zhang, Chaofan Lin, Chen Dong,  
633 Chengda Lu, Chenggang Zhao, Chengqi Deng, Chen-  
634 hao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian  
635 Yang, and 245 others. 2025. [Deepseek-v3.2: Pushing  
636 the frontier of open large language models](#). *Preprint*,  
637 arXiv:2512.02556.
- 638 R. A. DeMillo, R. J. Lipton, and F. G. Sayward. 1978.  
639 [Hints on test data selection: Help for the practicing  
640 programmer](#). *Computer*, 11(4):34–41.
- 641 George Fink and Matt Bishop. 1997. [Property-based  
642 testing: a new approach to testing for assurance](#). *SIG-  
643 SOFT Softw. Eng. Notes*, 22(4):74–80.
- 644 Harrison Goldstein, Joseph W. Cutler, Daniel Dick-  
645 stein, Benjamin C. Pierce, and Andrew Head. 2024.  
646 [Property-based testing in practice](#). In *Proceedings  
647 of the IEEE/ACM 46th International Conference on  
648 Software Engineering*, ICSE ’24, New York, NY,  
649 USA. Association for Computing Machinery.
- 650 Hypothesis team. Integrations Reference: Ghostwriter  
651 (Hypothesis 6.148.9 Documentation). [https://hypothesis.readthedocs.io/en/latest/  
652 reference/integrations.html#ghostwriter](https://hypothesis.readthedocs.io/en/latest/reference/integrations.html#ghostwriter).  
653 Accessed: 2025-11-01. 654
- 655 Md. Ashraf Islam, Mohammed Eunus Ali, and  
656 Md Rizwan Parvez. 2024. [MapCoder: Multi-agent  
657 code generation for competitive problem solving](#). In  
658 *Proceedings of the 62nd Annual Meeting of the As-  
659 sociation for Computational Linguistics (Volume 1:  
660 Long Papers)*, pages 4912–4944, Bangkok, Thailand.  
661 Association for Computational Linguistics.
- 662 René Just, Darioush Jalali, Laura Inozemtseva,  
663 Michael D. Ernst, Reid Holmes, and Gordon Fraser.  
664 2014. [Are mutants a valid substitute for real faults in  
665 software testing?](#) In *Proceedings of the 22nd ACM  
666 SIGSOFT International Symposium on Foundations  
667 of Software Engineering*, FSE 2014, page 654–665,  
668 New York, NY, USA. Association for Computing  
669 Machinery.
- 670 Runlin Liu, Yuhang Lin, Yunge Hu, Zhe Zhang, and  
671 Xiang Gao. 2024. [Llm-based java concurrent pro-  
672 gram to arks converter](#). In *Proceedings of the 39th  
673 IEEE/ACM International Conference on Automated  
674 Software Engineering*, ASE ’24, page 2403–2406,  
675 New York, NY, USA. Association for Computing  
676 Machinery.
- 677 Runlin Liu, Zhe Zhang, Yunge Hu, Yuhang Lin, Xiang  
678 Gao, and Hailong Sun. 2025. Llm-based unit test

679	generation for dynamically-typed programs. <i>arXiv preprint arXiv:2503.14000</i> .	Zhe Zhang, Xingyu Liu, Yuanzhang Lin, Xiang Gao, Hailong Sun, and Yuan Yuan. 2025. <a href="#">Reference-based retrieval-augmented unit test generation</a> . <i>ACM Trans. Softw. Eng. Methodol.</i> Just Accepted.	732
680			733
681	Muhammad Maaz, Liam DeVoe, Zac Hatfield-Dodds, and Nicholas Carlini. 2025. <a href="#">Agentic property-based testing: Finding bugs across the python ecosystem</a> . In <i>NeurIPS 2025 Fourth Workshop on Deep Learning for Code</i> .		734
682			735
683			
684			
685			
686	David MacIver, Zac Hatfield-Dodds, and Many Contributors. 2019. <a href="#">Hypothesis: A new approach to property-based testing</a> . <i>Journal of Open Source Software</i> , 4:1891.		
687			
688			
689			
690	David R. MacIver and Alastair F. Donaldson. 2020. <a href="#">Test-Case Reduction via Test-Case Generation: Insights from the Hypothesis Reducer</a> . In <i>34th European Conference on Object-Oriented Programming (ECOOP 2020)</i> , volume 166 of <i>Leibniz International Proceedings in Informatics (LIPIcs)</i> , pages 13:1–13:27, Dagstuhl, Germany. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.		
691			
692			
693			
694			
695			
696			
697			
698	Thomas J. McCabe. 1976. A complexity measure. In <i>Proceedings of the 2nd International Conference on Software Engineering, ICSE '76</i> , page 407, Washington, DC, USA. IEEE Computer Society Press.		
699			
700			
701			
702	Mutmut contributors. 2025. <a href="#">mutmut: Mutation testing for python</a> . Accessed: 2025-12-03.		
703			
704	OpenAI. 2025. <a href="#">Gpt-5 system card</a> . Accessed: 2025-11-29.		
705			
706	QwenTeam. 2025a. <a href="#">Qwen3-next: Towards ultimate training &amp; inference efficiency</a> . Accessed: 2025-11-29.		
707			
708			
709	QwenTeam. 2025b. <a href="#">Qwen3: Think deeper, act faster</a> . Accessed: 2025-11-29.		
710			
711	Savitha Ravi and Michael Coblenz. 2025. <a href="#">An empirical evaluation of property-based testing in python</a> . <i>Proc. ACM Program. Lang.</i> , 9(OOPSLA2).		
712			
713			
714	Vasudev Vikram, Caroline Lemieux, Joshua Sunshine, and Rohan Padhye. 2024. <a href="#">Can large language models write good property-based tests?</a> <i>Preprint</i> , arXiv:2307.04346.		
715			
716			
717			
718	Xinchen Wang, Pengfei Gao, Xiangxin Meng, Chao Peng, Ruida Hu, Yun Lin, and Cuiyun Gao. 2024. <a href="#">Aegis: An agent-based framework for general bug reproduction from issue descriptions</a> . <i>arXiv preprint arXiv:2411.18015</i> .		
719			
720			
721			
722			
723	Jason Wei. 2024. <a href="#">Asymmetry of verification and verifier's rule</a> . Accessed: 2025-12-03.		
724			
725	Lin Yang, Chen Yang, Shutao Gao, Weijing Wang, Bo Wang, Qihao Zhu, Xiao Chu, Jianyi Zhou, Guangtai Liang, Qianxiang Wang, and 1 others. 2024. <a href="#">On the evaluation of large language models in unit test generation</a> . In <i>Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering</i> , pages 1607–1619.		
726			
727			
728			
729			
730			
731			
		<b>A Case Analysis</b>	736
		<b>A.1 Failed Case Analysis</b>	737
		To better understand the validity results in RQ-2, we analyze the single failed case generated by PROBE. The target function is <code>sympy.series.gruntz.compare(a, b, x)</code> , a core utility within the Gruntz limit algorithm. The docstring for <code>compare</code> succinctly states: <i>"returns '&lt;' if a &lt; b, '=' if a == b, and '&gt;' if a &gt; b."</i> In isolation, this wording strongly resembles a conventional total order, which can lead a general-purpose LLM to interpret <code>compare</code> as comparing numeric values. Guided by this interpretation, PROBE synthesized the following PBT, asserting that distinct powers of $x$ should be strictly ordered:	738
			739
			740
			741
			742
			743
			744
			745
			746
			747
			748
			749
			750
		$\forall m, n \in \mathbb{Z}^+, \quad m < n \implies$ $\text{compare}(x^m, x^n, x) = "<"$	751
		This property is <i>invalid</i> under the intended semantics of <code>compare</code> . In the Gruntz framework, the comparison operators correspond to a coarse asymptotic classification used for limit computation, rather than a fine-grained ordering over polynomial degrees. In particular, polynomial powers typically fall into the same comparability class (e.g., $x^m \asymp x^n$ for $m, n > 0$ ), so <code>compare</code> may legitimately return <code>"="</code> even when $m < n$ . Consequently, the generated PBT fails not because of a defect in SymPy, but due to a <i>semantic gap</i> between an underspecified docstring and a domain-specific interpretation.	752
			753
			754
			755
			756
			757
			758
			759
			760
			761
			762
			763
			764
			765
			766
			767
			768
			769
			770
			771
			772
		<b>A.2 Bug Cases Analysis</b>	773
		To provide qualitative insights into PROBE's efficacy, we present a comparative analysis of PBT generation in Figure 8 and a showcase of discovered real-world defects in Figure 9.	774
			775
			776
			777
		Figure 8 illustrates that using the <code>mark_ends</code> function, where a logic bug was injected. The LLM	778
			779

```

def mark_ends(iterable):
    """
    Yield 3-tuples of the form ``(is_first, is_last,
    item)``.
    ...
    Use this when looping over an iterable to take
    special action on its first and/or last items
    ...
    """
    it = iter(iterable)
    for a in it:
        first = True
        for b in it:
            yield first, False, a
            yield first, True, a
            a = b
            first = False
        yield first, True, a

```

**Bug Injection**

```

@settings(...)
@given(...)
def test_properties(lst):
    ...
    result = list(mark_ends(lst))
    assert result[0][0] is True
    result = list(mark_ends(lst))
    assert len(result) == len(lst)
    ...

```

**Failed to catch it**

```

@settings(...)
@given(...)
def test_properties(iterable):
    ...
    if not iterable:
        assert not result
    else:
        first = sum(1 for is_first, _, _ in result if is_first)
        last = sum(1 for _, is_last, _ in result if is_last)
        assert first + last == 2
    ...

```

**Catch the bug**

Figure 8: Comparative analysis of generated PBTs between PROBE and LLM.

```

from sympy.ntheory.elliptic_curve import EllipticCurve,
EllipticCurvePoint
# Any elliptic curve
curve = EllipticCurve(0, 0, 0, 1, 1, 0)
P = curve(0, 0, 1)
identity = EllipticCurvePoint.point_at_infinity(curve)
# Compute P + (-P) - should return identity
assert P + (-P) == identity # Fail

```

**Sympy**

(a) Sympy: Violation of group property of elliptic curves.

```

import awkward as ak
# Any random float number
value = -943305.0469559873
arr = ak.Array([[value]])
json_str = ak.to_json(arr)
back = ak.from_json(json_str)
assert ak.array_equal(arr, back, equal_nan=True) # Fail

```

**Awkward Array**

(b) Awkward: Data loss through the round-trip.

```

from shlex import shlex
lexer = shlex('a')
print(lexer.get_token()) # a
print(lexer.get_token()) # EOF -> state is None
lexer.push_source('b', None)
print(lexer.get_token()) # actual '' (EOF), expected: 'b'

```

**Python**

(c) CPython: Shlex failed to reset the lexer state.

```

from cryptography.x509 import Name
from cryptography.x509.name import NameAttribute
from cryptography.x509.oid import NameOID
name = Name([NameAttribute(NameOID.COMMON_NAME, "0"),
NameAttribute(NameOID.ORGANIZATION_NAME, " ")])
s = name.rfc4514_string() # yields "CN=0,0=\\ "
back = Name.from_rfc4514_string(s) # raises ValueError
assert back == name # Fail

```

**Cryptography**

(d) Cryptography: Non-parseable output for trailing spaces.

Figure 9: Exemplary round-trip like bugs from reported issues.

780 generates superficial assertions that pass despite  
781 the bug, exemplifying the weak property problem  
782 inherent in standard prompting. In contrast, PROBE  
783 synthesizes a precise invariant verifying that the  
784 sum of `is_first` and `is_last` flags must equal 2 for the  
785 sequence, successfully catching the bug.

786 Figure 9 shows 4 minimal reproduced codes for  
787 exemplary round-trip-like bugs from our reported  
788 issues. In Sympy, PROBE identified the fundamen-  
789 tal group property  $P + (-P) = O$  (Identity) viola-  
790 tion. The library computed an incorrect sum for a  
791 point and its inverse on a specific elliptic curve, vi-  
792 olating the mathematical axioms the library claims  
793 to support. In package Awkward, PROBE detected  
794 a precision loss when converting floating-point ar-

795 rays to JSON and back. The deserialized value  
796 differed from the original input, indicating a failure  
797 in preserving data fidelity during the serialization  
798 process. In CPython, the shlex lexer failed to cor-  
799 rectly reset its internal EOF state when a new input  
800 source was pushed after the previous one was ex-  
801 hausted. This caused the lexer to prematurely stop  
802 processing valid subsequent inputs. Finally, PROBE  
803 detected that the Cryptography library successfully  
804 serialized but failed to parse that the same string  
805 back into an object, raising a `ValueError`. These  
806 bugs are detected through cross-function properties,  
807 which are typically unreachable by trivial proper-  
808 ties.

## B Experiment Details

All experiments were conducted on a Ubuntu 20.04 server with dual Intel(R) Xeon(R) Gold 5218 @2.30GHz CPU with 128 GB of RAM, 2 NVIDIA A100-SXM4-80GB GPUs. We set both the diagnostic retry limit  $T_{fix}$  and the adversarial refinement limit  $T_{ref}$  to 3 across all experiments.

### B.1 Mutation Score with full dataset

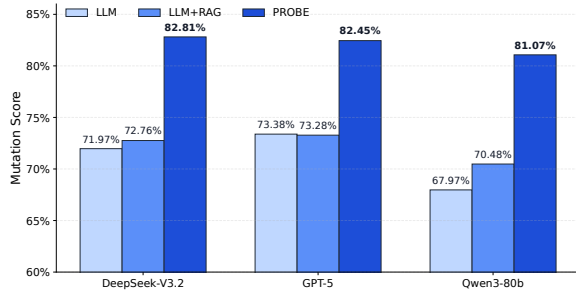


Figure 10: Performance comparison of PROBE against baselines across different backbone models (full dataset).

Method	Pass%	Killed / Total	Score
<b>LLM</b>			
Qwen3-80b	73.04%	1193 / 1608	67.97%
DeepSeek-V3.2	58.99%	937 / 1302	71.97%
GPT-5	84.77%	1384 / 1886	73.38%
<b>LLM + RAG</b>			
Qwen3-80b + RAG	52.73%	814 / 1155	70.48%
DeepSeek-V3.2 + RAG	47.27%	708 / 973	72.76%
GPT-5 + RAG	77.34%	1245 / 1699	73.28%
<b>LLM-Based Agent</b>			
Claude-Code (sonnet 4.5)	<b>97.66%</b>	1521 / 2186	69.58%
<b>Ours</b>			
PROBE (Qwen3-80b)	76.20%	1366 / 1685	81.07%
PROBE (DeepSeek-V3.2)	74.22%	<b>1363 / 1646</b>	<b>82.81%</b>
PROBE (GPT-5)	90.60%	1668 / 2029	82.45%

Table 6: *Pass%* denotes the percentage of functions for which the tool generated a passing PBT. *Killed / Total* denotes the number of killed mutants / total mutants. *Score* denotes the mutation score.

Table 6 and Figure 10 present mutation testing results on the complete Dataset I (256 functions), complementing the intersection-based analysis in Section 4.2. Unlike the main evaluation, which controls for function coverage by restricting to the intersection where all tools produce passing PBTs, this analysis evaluates each tool on all functions for which it successfully generates a passing PBT.

PROBE achieves substantially higher mutation scores compared to all baselines across the full

dataset. PROBE (DeepSeek-V3.2) attains the highest mutation score of 82.81%, followed by PROBE (GPT-5) at 82.45% and PROBE (Qwen3-80b) at 81.07%. These scores represent improvements of 9~14 percentage points over the corresponding LLM baselines and 7~10 percentage points over LLM+RAG configurations. Notably, the performance gap between PROBE and baselines is even more pronounced on the full dataset than on the intersection (Table 1), suggesting that PROBE’s advantages are particularly evident on challenging functions where baselines struggle to generate effective PBTs.

Pass rate measures the proportion of functions for which a tool generates a syntactically valid PBT that passes on the original implementation. Claude Code achieves the highest pass rate (97.66%), followed by PROBE (GPT-5) at 90.60%. However, high pass rate does not necessarily indicate effective PBT generation. Claude Code, despite its superior pass rate, attains a mutation score of only 69.58%. This confirms our observation that existing methods tend to generate conservative properties that pass easily but lack discriminative power.

### B.2 Dataset Construction Details

To construct a representative dataset, we applied a systematic filtering process to functions collected from `sympy`, `more-itertools`, `simplejson`, and `sortedcontainers` packages. Starting from an initial corpus of 21,319 candidate functions, we applied a systematic filtering process based on code complexity metrics. Our selection criteria established both minimum and maximum thresholds to exclude trivial implementations and overly complex functions that would confound evaluation. Minimum thresholds required at least 5 lines of code (LOC), 3 source lines of code (SLOC), and a McCabe’s Cyclomatic Complexity (CC) (McCabe, 1976) of at least 2 to ensure meaningful implementations containing at least one decision point. Maximum thresholds constrained  $LOC \leq 200$ ,  $CC \leq 15$ , Cognitive Complexity (Campbell, 2018)  $\leq 25$ , maximum nesting depth  $\leq 5$ , and number of parameters  $\leq 7$ . Additionally, we required all functions to include docstrings and available Python source code. To ensure library diversity and prevent overrepresentation, we sampled at most 100 functions per package.

After applying all criteria, 256 functions were retained for the final dataset, comprising 100 functions from `sympy`, 100 from `more-itertools`, 48

Metric	Min	Max	Mean	Median
Cyclomatic Complexity	2	15	4.82	4
Cognitive Complexity	0	21	5.20	3
Lines of Code	5	194	35.80	32

Table 7: Complexity statistics of selected functions

from sortedcontainers, and 8 from simplejson. Table 7 presents the complexity distribution of the final dataset.

### B.3 LLMs Configuration

We access all language models through official APIs or local deployment. For local deployment, we host Qwen3-30B-A3B-Thinking using **SGLang**<sup>1</sup>. For API-based models, we access DeepSeek-V3.2 through the official DeepSeek API<sup>2</sup>, GPT-5 through OpenRouter<sup>3</sup>, and Qwen3-Next-80B-A3B-Thinking through Alibaba Cloud<sup>4</sup>. For all models, we retain the default temperature and sampling configurations provided by each platform without modification.

### B.4 Baseline Details

#### B.4.1 LLM

We provide the LLM with the target function’s metadata, and prompt it to identify properties and generate corresponding PBTs. The generated PBTs are then executed against the target function. Since we found that single-pass generation yields prohibitively low pass rates due to the syntax or logic error, we incorporate a self-repair mechanism: if execution fails, we feed the error message back to the LLM for repair, allowing a maximum of 3 fix attempts. This provides a more fair comparison.

#### B.4.2 Retrieval-Augmented LLM (LLM+RAG)

We adopt LLM+RAG framework as another baseline. Given a target project, we decompose the codebase into function-level units, where each unit corresponds to a single function definition. This granularity preserves the semantic coherence of each functional unit and is well aligned with property-based test generation, as it facilitates the identification of properties that recur across functions.

<sup>1</sup><https://github.com/sgl-project/sglang>

<sup>2</sup><https://api-docs.deepseek.com/>

<sup>3</sup><https://openrouter.ai/>

<sup>4</sup><https://bailian.console.aliyun.com/>

All function units are embedded using **bge-m3**<sup>5</sup> and indexed in the **Chroma**<sup>6</sup> vector store. During test generation, the language model first formulates a set of abstract property descriptions for the target function. For each property, we retrieve the top-k semantically related functions from the codebase based on embedding similarity, followed by a lightweight cosine-similarity re-ranking step to refine relevance.

The retrieved functions are then supplied to the model as auxiliary context, together with the target function and the corresponding property description, to guide the generation of executable PBTs. All generated tests are subsequently executed, and their outcomes are analyzed by the language model to assess compliance with the intended properties. When violations are observed, the model iteratively refines the generated tests to resolve inconsistencies between the specified properties and observed execution behavior.

#### B.4.3 Claude Code

We employ the official Claude CLI ([Anthropic, 2025a](#)) with Claude Sonnet 4.5 as the underlying model. Following the prompt design from [Maaz et al. \(2025\)](#)’s work, we instruct the agent to analyze the target function and generate PBTs. Claude Code operates as a fully autonomous agent capable of reading files, executing code, and iteratively fixing its outputs without explicit prompts. This baseline represents the current state-of-the-art in commercial LLM-based coding agents.

## C Summary of Reported Issues

This section provides detailed information on the bugs discovered by PROBE during our evaluation in RQ-3 (Section § 4.4). We only reported issues that we assessed as genuine bugs with potential impact on real-world users. Table 8 summarizes the detailed information. For each reported bug shown in table, we provide the affected repository with star count, module, the issue tracking ID (anonymized with \*), current resolution status, fix type, and a brief description of the defect. Complete reports are available on our GitHub repository.<sup>7</sup>

## D Prompt Details

<sup>5</sup><https://huggingface.co/BAAI/bge-m3>

<sup>6</sup><https://github.com/chroma-core/chroma>

<sup>7</sup><https://anonymous.4open.science/r/PROBE-C7A5>

Table 8: Summary of Bug Findings.

Repo Name (Star)	Module	Issue ID	Issue State	Fix Type	Description
cpython (70.5k)	re	*	FIXED	CODE FIX	Crash when using multiple capturing groups in re.Scanner.
	html	*	FIXED	CODE FIX	html.parser() silently drops data.
	html	*	CONFIRMED	-	HTMLParser.unknown_decl receives corrupted data in corner cases.
	html	*	REJECTED	-	check_for_whole_start_tag crashes with AssertionError on empty string .
	functools	*	FIXED	DOC FIX	Inconsistent behavior with docstring.
	collections	*	FIXED	CODE FIX	collections.UserString.rindex() fails to accept UserString as a sub argument.
	collections	*	DUPLICATED	-	collections.namedtuple fails to reconstruct from _asdict().
	statistics	*	FIXED	CODE FIX	statistics.stdev() and statistics.variance do not satisfy the relationship in corner cases.
	statistics	*	NOT PLANNED	-	NormalDist.overlap() returns incorrect values due to overflow.
	shlex	*	CONFIRMED	-	shlex.push_source() fails to reset lexer state, causing new stream to be ignored if called after EOF.
	types	*	FIXED	CODE FIX	DynamicClassAttribute drops explicitly provided empty docstrings.
json	*	FIXED	DOC FIX	JSONEncoder incorrectly escapes DEL character.	
tokenizer (10.3k)	tokenizers	*	FIXED	CODE FIX	BaseTokenizer.normalize method calls non-existent method.
sympy (14.2k)	ntheory	*	CONFIRMED	CODE FIX	Point addition $P + (-P)$ fails to return point at infinity.
	ntheory	*	FIXED	CODE FIX	EllipticCurvePoint.neg returns non-canonical infinity point.
	algebra	*	FIXED	CODE FIX	Quaternion.log return nan in corner case.
	calculus	*	CONFIRMED	CODE FIX	is_monotonic misclassifies monotonic functions with stationary points.
	calculus	*	FIXED	CODE FIX	AccumBounds.intersection violates the commutativity.
	calculus	*	CONFIRMED	CODE FIX	is_decreasing crashes on single-point intervals due to ConditionSet substitution error.
	calculus	*	FIXED	CODE FIX	is_monotonic ignores discontinuities inside the interval.
	geometry	*	FIXED	CODE FIX	Plane.intersection fails to find point in plane.
	combinatorics	*	FIXED	CODE FIX	AbelianGroup.random() crashes on trivial group.
	combinatorics	*	FIXED	CODE FIX	AlternatingGroup(2) has wrong degree.
	combinatorics	*	FIXED	CODE FIX	collector.induced_pcggs drops generators at final depth.
	combinatorics	*	FIXED	CODE FIX	Cycle().list(size=0) returns [0] instead of the empty list.
	combinatorics	*	CONFIRMED	CODE FIX	PermutationGroup.coset_rank mishandle trivial groups.
pathspec (205)	pathspec	*	NOT PLANNED	-	CheckResult.include returns None instead of bool.
	pathspec	*	NOT PLANNED	-	check_tree_files and match_files return inconsistent results.
	pathspec	*	FIXED	CODE FIX	UnboundLocalError in RegexPattern.

Repo Name (Star)	Module	Issue ID	Issue State	Fix Type	Description
cryptography (7.4k)	cryptography	*	FIXED	CODE FIX	IssuingDistributionPoint.hash fails for valid inputs.
	cryptography	*	FIXED	CODE FIX	Name.rfc4514_string emits non-parseable output for trailing spaces.
pyproj (1.2k)	pyproj	*	PENDING	-	Hash/Equality contract violation.
	pyproj	*	CONFIRMED	DOC FIX	Transformer.to_json() returns None for valid CRS pair.
awkward (929)	awkward	*	CONFIRMED	CODE FIX	Crash ak.to_json loses float64 precision on round-trip.
	awkward	*	FIXED	CODE FIX	ak.moment applies weights incorrectly for $n \geq 2$ (weights exponentiated).
	awkward	*	FIXED	CODE FIX	ak.Record rejects nested dicts with scalar leaves.
	awkward	*	CONFIRMED	-	ak.ArrayBuilder.show passes Formatter to ak.Array.show, causing TypeError.
	awkward	*	PENDING	-	ByteMaskedArray.mask_as_bool returns inverted validity for valid_when=False.
	awkward	*	FIXED	CODE FIX	IndexedOptionArray.to_IndexedOptionArray64 fails when index dtype is int32.
	awkward	*	PENDING	-	Silent data corruption using tolist().
	awkward	*	FIXED	CODE FIX	ArrayBuilder.__bool__ raises TypeError when builder length is 1 .
pws-python (3.2k)	logging	*	FIXED	DOC FIX	append_context_keys drops previously appended keys.
	event_handler	*	FIXED	DOC FIX	BedrockResponse.is_json always returns True regardless of ContentType.
	event_handler	*	FIXED	DOC FIX	UnauthorizedException raises TypeError when using documented keyword arguments.
optax (2.1k)	optax	*	FIXED	CODE FIX	scale_by_distance_over_gradients returns nan.
	optax	*	FIXED	CODE FIX	contrib.momo crashes when loss value is a Python float.
pyrsistent (2.2k)	pyrsistent	*	PENDING	-	pdeque ignores maxlen=0 during construction.
marshmallow (7.2k)	marshmallow	*	CONFIRMED	CODE FIX	fields.Constant rejects None constant values during load.
	marshmallow	*	PENDING	-	OneOf.options() emits extra pairs when labels outnumber choices.
	marshmallow	*	FIXED	CODE FIX	URL validator treats custom schemes case-sensitively.
scipy (14.3k)	interpolate	*	CONFIRMED	CODE FIX	PPoly.roots crashes with Internal error in root finding for valid cubic spline.
	cluster	*	FIXED	CODE FIX	scipy.cluster.vq.kmeans returns a distortion value from a different iteration than the final centroids.
	cluster	*	CONFIRMED	-	io.hb_write crashes on sparse matrices with zero stored entries.
	interpolate	*	FIXED	CODE FIX	RectBivariateSpline.partial_derivative rejects valid x-derivative orders when $k_x > k_y$ .
boto3 (9.6k)	resources	*	NOT PLANNED	-	ResourceCollection.limit(0) yields one item.
networkx (16.4k)	networkx	*	FIXED	DOC FIX	chordal_cycle_graph returns a 6-regular multigraph with self-loops.

## Properties Planning Prompt

You are an expert in **Property-Based Testing (PBT)** for Python functions.

Your Goal:

Generate verifiable properties description for a target function at a time **WITH DIRECT EVIDENCE**.

**CRITICAL:** Only generate properties that are **DIRECTLY** stated in the evidence!

Process:

1. Analyze the docstrings and code.
2. Find statements that **DIRECTLY** and **EXPLICITLY** describe behavior.
3. Generate a property that is a **DIRECT TRANSLATION** of that statement.

Guidelines:

1. Property must be a direct translation of evidence, **NOT** an assumption.
2. Evidence should be exact quotes, not paraphrases.
3. Don't generalize from examples to "always" behavior.

Output strict JSON format:

```
{
  "property": "The property description",
  "evidence": "evidence from docstring/comments that DIRECTLY states this property",
  "evidence_type": "docstring" | "code comments" | "mathematical_convention" | "type_hint",
  "confidence": "high" (if directly stated) | "medium" (if strongly implied) | "low" (if inferred)
}
```

Here is the Target Function information:

```
{func_info}
```

Here are the current properties that have been identified for this function:

```
{properties}
```

Here are some Reference Property Heuristics (Use these as inspiration):

```
{exemplary_properties}
```

Based on the heuristics and the function info, generate new non-trivial properties description that is logically different from the **Existing Properties**.

## Context Fetching Prompt

You are an expert in selecting relevant functions for **Property-Based Testing**.

Given a target function with a specific property, and a list of all available functions in the library, your task is to select functions that would be useful for testing the property.

Consider:

1. Complementary functions (e.g., encode/decode for round-trip)
2. Related operations that could be composed
3. Helper functions that create test inputs
4. Validation functions

Respond in JSON format:

```
{ "selected_funcs": [func1, "func2", ...] }
```

**Note:** You should only return the name of the function from the provided list, no extra information.

Target Function:

```
{func_info}
```

Property to test:

```
{property}
```

Available functions in the library:

```
{signature_index}
```

Select several most relevant functions from the available functions in the library for testing this property of target function.

## PBT Generation Prompt

You are an expert Python QA engineer specializing in **Property-Based Testing** using the hypothesis library and pytest. Your task is to write a single, executable Python test function that verifies a specific property for a given target function.

### Strict Output Requirements:

1. Output one Python code block.
2. The code must be clean and executable.
3. Do not generate any text outside the code block.
4. Include necessary imports (e.g., from hypothesis import given, settings, strategies as st) at the beginning of the code.
5. Infer the appropriate hypothesis.strategies based on the target function's argument types and the property description.
6. The test function name should start with test\_.

Target Function:

{func\_info}

Relevant Context:

{relevant\_context}

Hypothesis guidelines:

{hypothesis\_doc}

Property to test:

{property}

### CRITICAL Hypothesis Guidelines:

- \* Always use @settings(max\_examples=1000) decorator on **EVERY** test function for thorough testing.
- \* Use math.isclose() or pytest.approx() for float comparisons
- \* Focus on properties that reveal genuine bugs when violated
- \* For recursive data generation, use st.recursive().
- \* Keep PBT code simple. Prefer using basic strategies like st.text(), st.integers(), st.lists() directly rather than complex custom composite strategies when possible.

### Example Structure:

{Some PBT examples}

## Bug determination Prompt

You are an expert software testing engineer analyzing Property-Based Testing (PBT) execution results.

Goal: Determine the type of failure from three categories:

1. **Code Defect**: malformed strategies, invalid dimension assumptions, missing mocks, overly strict tolerances, nondeterministic seeding, import errors, syntax errors, or any mistake that makes the property check unreliable or impossible.
2. **Property Defect**: the property itself is incorrectly designed based on the evidence, not matching what the docstrings promises.
3. **Library Defect\*\***: the property and evidence correctly reflects what the docstrings promises, but the actual code implementation violates it.

**CRITICAL: Ground Truth is NOT CODE IMPLEMENTATION!** - The docstring is the specification. The code is the implementation.

PBT Code: {code}

Execution Result: {msg}

Property to test: {property}

Property Evidence: {evidence}

Evidence Type: {evidence\_type} Evidence Confidence: {confidence}

Target Function: {func\_info}

Auxiliary Functions Information: {full\_context}

Output strict JSON { "bug\_type": "code\_defect" | "property\_defect" | "library\_defect", "reasoning": "Succinct justification quoting the most relevant evidence and explaining your reasoning.", "fix\_suggestion": "If code\_defect, provide a specific fix suggestion. If property\_defect, explain what's wrong with the property and how to redesign it. If library\_defect, leave empty." }

## Counter-Impl Generation Prompt

You are an expert software engineer performing **ADVERSARIAL REFINEMENT** on Property-Based Tests. Your task is to generate **ONE COUNTER-IMPLEMENTATION**: a code modification of the target function that preserves syntax but alters behavior in ways the current PBT **SHOULD** detect but **DOESN'T**.

A good counter-implementation  $f'$ :

1. Is syntactically valid Python code.
2. Violates the **INTENDED** semantics (what the function should do).
3. Would **PASS** the current PBT (exposing PBT weakness).

Examples of counter-implementations:

- For a sorting function: return the list with the last element dropped (passes "ordered" check but wrong)
- For an addition function: return 0 for certain inputs (passes some tests but wrong for edge cases)
- For a contains function: always return True (passes positive tests but wrong for negatives)

Generate **EXACTLY ONE** counter-implementation with **ACTUAL EXECUTABLE CODE**.

Output strict JSON:

```
{
  "has_counter_impl": true | false, "counter_implementation": { "description": "Brief description of what this counter-implementation does wrong", "code": "The complete modified function code (executable Python)", "what_it_violates": "What semantic property this violates" }, "reasoning": "Why this counter-implementation would pass the current PBT"
}
```

If the PBT is already strong enough (no counter-implementation can pass), output:

```
{
  "has_counter_impl": false, "counter_implementation": null, "reasoning": "explanation"
}
```

## Strengthen Property Prompt

You are an expert in Property-Based Testing (PBT) for Python functions.

Your task is to design a **STRONGER** property based on **SURVIVING COUNTER-IMPLEMENTATIONS**.

Context: we generated counter-implementations  $f'$  (modified versions of the target function).

These counter-implementations **PASSED** the current PBT, meaning the PBT is too weak to detect them.

You must design a property that would **FAIL** on these counter-implementations.

The new property should:

1. Specifically **FAIL** on the surviving counter-implementations (catch their incorrect behavior).
2. Be verifiable through PBT.
3. Not overlap with existing properties.

Target Function: {func\_info}

Existing Properties: {existing\_properties}

Current PBT: {PBT}

SURVIVING Counter-Implementations: {surviving\_counter\_impls}

Analyze each counter-implementation:

- What incorrect behavior does it exhibit?
- What property would catch this behavior?

Output strict JSON:

```
{
  "property": "The stronger property description that catches the counter-implementations", "evidence": "Reference to which counter-implementation behavior this catches and why", "evidence_type": "counter_implementation"
}
```