# STATEFUL MULTI-AGENT EVOLUTIONARY SEARCH FOR UNIT TEST GENERATION

## **Anonymous authors**

000

001

002003004

010 011

012

013

014

016

017

018

019

021

025

026

027

028

029

031

032

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

Recent work explores agentic inference-time techniques to perform structured, multi-step reasoning. However, stateless inference often struggles on multi-step tasks due to the absence of persistent state. Moreover, task-specific fine-tuning or instruction-tuning often achieve surface-level code generation but remain brittle on tasks requiring deeper reasoning and long-horizon dependencies. To address these limitations, we propose **stateful multi-agent evolutionary search**, a training-free framework that departs from prior stateless approaches by combining (i) persistent inference-time state, (ii) adversarial mutation, and (iii) evolutionary preservation. We demonstrate its effectiveness in automated unit test generation through the generation of edge cases. We generate robust edge cases using an evolutionary search process, where specialized agents sequentially propose, mutate, and score candidates. A controller maintains persistent state across generations, while evolutionary preservation ensures diversity and exploration across all possible cases. This yields a generalist agent capable of discovering robust, high-coverage edge cases across unseen codebases. Experiments show our stateful multi-agent inference framework achieves substantial gains in coverage over stateless single-step baselines, evaluated on prevalent unit-testing benchmarks such as HumanEval and TestGenEvalMini and using three diverse LLM families—Llama, Gemma, and GPT. These results indicate that combining persistent inference-time state with evolutionary search materially improves unit-test generation.

## 1 Introduction

Despite their success on single-step tasks, most inference-time computation in large language models (LLMs) remains stateless, with each inference call discarding prior intermediate reasoning unless explicitly re-injected into the prompt. This design choice optimizes deployment throughput but cripples performance in domains that require deep, multi-stage reasoning—such as program synthesis, theorem proving, multi-hop reasoning, deductive reasoning, and mathematical problem-solving where intermediate states must be persistently updated and revisited. The fixed computational depth per transformer forward pass (Vaswani et al., 2017) and the well-documented decline in reasoning fidelity over long logical chains (Wei et al., 2022; Anil et al., 2022) make these limitations structural rather than incidental. The autoregressive decoding process further constrains exploration by forcing reasoning branches to unfold serially, often necessitating brittle orchestration through multiple model calls (Yao et al., 2023a; Long et al., 2024). Overcoming these constraints demands stateful inference-time architectures—including scratchpad prompting (Nye et al., 2021; Wei et al., 2022), tree-structured reasoning (Yao et al., 2023b), and retrieval-augmented agents (Lewis et al., 2020)—that can maintain and manipulate intermediate reasoning artifacts directly. However, current techniques still operate by eliciting reasoning from fixed, opaque model parameters, limiting both steerability (Zhou et al., 2023) and interpretability (Olah et al., 2020; Nanda et al., 2023) of the reasoning process.

In this context, we investigate the suitability of a multi-stage evolutionary algorithm (Bäck et al., 1997; Hansen, 2016) in which each stage executes an adversarially guided actor–critic-style (AGAC) search (Ding et al., 2023). Unlike conventional evolutionary pipelines where each generation is evaluated in isolation, our design shares state information—captured as the critic's value estimates—across successive evolutionary stages. This shared evaluation signal serves as a persistent knowledge base, allowing later stages to inherit and refine the judgment of earlier stages rather

than re-learning from scratch (Jaderberg et al., 2017; Such et al., 2017). Such cross-stage information flow improves sample efficiency, reduces evaluation variance, and encourages coherent policy evolution over long optimization horizons (Mouret & Clune, 2015; Salimans et al., 2017). By integrating AGAC within this multi-stage framework, we can combine the exploration benefits of evolutionary search with the fine-grained feedback of actor–critic learning (Konda & Tsitsiklis, 2000). In this work, the actor–critic terminology strictly refers to inference-time reward shaping only: the critic scores candidate tests to guide the actor, but no model parameters are updated as in traditional reinforcement learning.

This distinction is important because our aim is not to train new policies, but to adapt inference-time behavior for practical tasks. Unit test generation provides an ideal setting to study this: it requires structured reasoning beyond syntax, benefits directly from persistent state, and offers measurable signals such as coverage and mutation scores to guide search. Stateless test generation often covers only a narrow slice of behavior, whereas maintaining state enables gradually expanding coverage and surfacing deeper failure modes. In summary, this work introduces a training-free framework for unit test generation that (i) maintains persistent inference-time state across search iterations, (ii) integrates coverage, exceptions, and mutation robustness into a unified reward design, and (iii) demonstrates consistent coverage improvements over stateless baselines on HumanEval and Test-GenEvalMini.

#### 2 Related Work

The rapid advancement of large language models (LLMs) has enabled significant progress in AI-assisted reasoning, code generation, and test automation. Prior research spans several domains including code generation, test synthesis, algorithmic discovery, and scientific reasoning, yet many approaches face limitations in adaptability, coverage, and generalization.

Multi-agent frameworks such as AI Co-scientist (Gottweis et al., 2025), AlphaEvolve (Novikov et al., 2025b) and GEPA (Agrawal et al., 2025) demonstrate that collaborative reasoning and reflective prompt evolution can enhance exploration. However, these systems typically lack persistent state and rely on ad-hoc orchestration rather than structured reward signals.

Evolutionary and search-based methods preserve high-fitness candidates and explore combinatorial program behaviors Mühlenbein et al., 1988; Burnim & Sen, 2008. In particular, evolutionary search explicitly manages the exploration-exploitation tradeoff, allowing the system to explore novel program behaviors while retaining high-performing edge cases and avoiding local minima that static or greedy methods often encounter. Other works Karten et al.; Leng et al., 2024; Wen et al., 2024 provide insights into scaling, planning, and iterative hypothesis generation, emphasizing structured search and evaluation.

Despite these advances, prior approaches often struggle with adaptability, robust coverage, and systematic exploration of edge cases. Many LLM-based test generators operate in a feed-forward manner or require fine-tuning, limiting their ability to dynamically adjust to new or evolving codebases. Multi-agent and evolutionary approaches in prior work may fail to explore the full combinatorial space or integrate adversarial evaluation effectively.

We address these gaps with a training-free framework that unifies (i) multi-agent reasoning, (ii) adversarial mutation and reward shaping, and (iii) evolutionary preservation with persistent state. An actor proposes candidate edge cases by reasoning over the program, an adversary perturbs the code to expose hidden failure modes, and a critic integrates coverage, exceptions, and mutation feedback to prioritize high-value test cases. A non-Markovian controller maintains memory of prior edge cases and preserves high-fitness candidates across iterations, enabling inference-time policy adaptation and robust exploration. This combination allows our system to dynamically adapt to unseen codebases, produce robust edge cases, and achieve higher coverage than existing methods, without relying on gradient-based training or domain-specific fine-tuning.

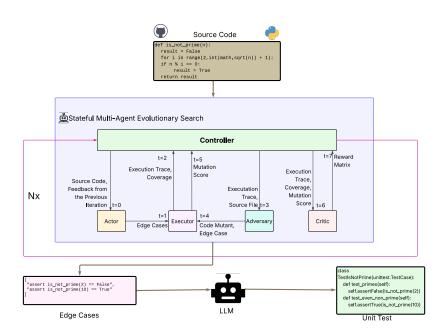


Figure 1: Our architecture for unit test generation decomposes the task into two phases: (i) edge case generation from source code and (ii) unit test construction from those cases. The first phase demands deeper reasoning and is addressed through an evolutionary search (as highlighted in the blue box) executed in a stateful manner over multiple stages  $(N \times)$  by four agents—Actor, Executor, Adversary, and Critic—coordinated by a Controller that propagates persistent state across N evolutionary stages (as highlighted by the magenta line). Once the edge cases converge to sufficient coverage and robustness, they are translated into a complete unit test file via a single-step inference call.

## 3 METHODOLOGY

Our central premise is that generating syntactically correct unit tests is trivial once a set of robust edge cases with sufficient coverage are identified, but reasoning about such edge cases requires structured exploration, memory, and adversarial grounding.

Figure 1 shows the architecture for the unit test generation engine with the proposed *stateful multi-agent evolutionary search* for the edge case generator. Given source code f, the system first runs the stateful multi-agent evolutionary search to extract edge cases and then converts those cases into unit tests.

Our stateful multi-agent evolutionary search is an adversarially guided actor-critic (AGAC) system that operates entirely at inference time and does not require gradient-based learning. The **Actor** issues multiple LLM inference calls to propose candidate edge cases, the **Adversary** perturbs the environment to reveal robustness gaps, and the **Critic** assigns scalar rewards used for evolutionary search. The **Executor** is an auxiliary agent that provides an execution environment to execute edge cases, unit tests, and return coverage and robustness feedback. These four agents are orchestrated through the **Controller** which maintains persistent state across stages and orchestrates the search until convergence.

**Definition 1** (State). A State is represented in Equation 1,

$$S_{n-1} = \left(\zeta_{1:n-1}, \mu_{1:n-1}, \kappa_{1:n-1}, c_{1:n-1}, R_{1:n-1}\right) \tag{1}$$

where  $\zeta_{1:n-1}$  denotes the sequence of prior edge cases,  $\mu_{1:n-1}$  is the sequence of mutation scores,  $\kappa_{1:n-1}$  is the sequence of coverage scores,  $c_{1:n-1}$  is the sequence of exception signals, and  $R_{1:n-1}$  is the reward history from previous stages.

**Definition 2** (Actor). The Actor  $(A_n)$  proposes candidate edge cases at each stage. At initialization (n=1), there is no prior feedback or state information to guide generation, so the actor is seeded deterministically (cold-start) using rule-based heuristics such as boundary partition analysis, equivalence classes, and stress conditions. For n>1, the actor generates candidates through large language model in-context learning, conditioned on the persistent state  $S_{n-1}$  and the source code f:

$$\zeta_n = \mathcal{A}(f, S_{n-1}) \tag{2}$$

More details about the cold-start can be found in Appendix D

**Definition 3** (Adversary). For each stage, Adversary  $(D_n)$  generates a set of mutants  $\{f'_{n,j}\}_{j=1}^M$  of the source file, and evaluates whether the edge cases  $\zeta_n$  can kill these mutants (i.e produce a different output on  $f'_{n,j}$  than they did on f). The resulting mutation score is defined in Equation 3 and provides a robustness signal for evaluating the edge case candidates. Mutation testing promotes robustness by checking whether tests can distinguish the true program from systematically perturbed variants, preventing the search from optimizing toward shallow coverage gains.

$$\mu_n = \frac{\text{Number of mutants killed by } \zeta_n}{\text{Total number of generated mutants}}$$
 (3)

**Definition 4** (Critic). For each stage, Critic  $(C_n)$  computes the scalar reward for the edge cases by integrating coverage  $(\kappa)$ , mutation robustness  $(\mu)$ , and exception discovery (c), given by Equation 4.

$$R_n^{\text{unnormalized}}(\kappa_n, \mu_n, c_n) = \left[\alpha \cdot c_n + \beta(\kappa_n + \max(0, (\kappa_n - \theta) \cdot 0.5))\right] \times \gamma \cdot \mu_n \tag{4}$$

where  $\alpha, \beta, \theta, \gamma \in \mathbb{R}_+$  are tunable hyperparameters. All rewards are normalized to [0,1] using min-max normalization for evolutionary comparison.

The reward combines exception discovery  $(c_n)$ , structural coverage  $(\kappa_n)$ , and mutation robustness  $(\mu_n)$ . The exception term encourages exploration of inputs that expose faults. The coverage term accounts for the proportion of program elements exercised, with an additional bonus once a minimum threshold  $\theta$  is passed, so that progress beyond trivial coverage is reflected more strongly. Multiplication by the mutation score ensures that high reward is assigned only when the generated tests are also robust to program perturbations. By shaping the critic's reward surface using adversarial perturbations we ground the actor's responses and thus prevent the actor from optimizing toward trivial coverage gains instead of exploring robust, high-value edge cases.

**Definition 5** (Executor). All evaluations for coverage and mutation scoring are executed in a sand-boxed Docker environment with a Model-Context Protocol (MCP) server. This provides: (i) **Isolation:** Mutants and edge cases cannot harm the host system; (ii) **Determinism:** Results are reproducible across runs; and (iii) **Bounded resources:** Memory and timeouts prevent unbounded execution. A detailed description of the Executor architecture can be found in Appendix A.3.

**Definition 6** (Controller). The controller orchestrates the interplay of Actor, Adversary, and Critic by updating the non-Markovian state information (Equation 1) and checking for the termination criteria as defined in Equation 5.

$$\sum_{i} R_i \ge \tau \text{ or } \max_{i \in [n-p+1, n]} R_i - \min_{i \in [n-p+1, n]} R_i \le \delta$$
 (5)

The controller applies two complementary stopping conditions. The first checks whether the reward has crossed a predefined threshold, indicating that the search has reached a sufficient overall quality level. The second detects a plateau in rewards over the most recent p iterations, suggesting that further search is unlikely to yield substantial improvements. The plateau condition is evaluated only when  $n \geq p$ . The thresholds  $(\tau, \delta)$  and window size p can be tuned according to task complexity as well as computational budget, allowing the framework to balance thoroughness and efficiency.

A key methodological contribution is that our framework does not require training, fine-tuning, or task-specific adaptation of large language models. Instead, it builds a training-free test-generation agent whose intelligence emerges from:

- 1. **Inference-time state management:** The controller maintains a structured non-Markovian state, feeding the actor with explicit histories of edge cases, coverage scores, mutation feedback, and exceptions. Unlike conventional RL where state updates drive gradient descent, here state updates directly shape the actor's inference context. This functions as a lightweight form of policy shaping at inference time, guided by rewards but without parameter updates.
- 2. **Multi-agent grounding:** The actor's outputs are consistently grounded by adversarial mutations and fitness evaluation from the critic, allowing even base LLMs without domain adaptation to be repurposed into reasoning agents.
- 3. **Evolutionary selection:** The framework preserves a population of diverse elites, avoiding reliance on a single trajectory and improving robustness without requiring specialized training.

This positions our framework alongside recent agentic paradigms such as AI Co-Scientist (Gottweis et al., 2025) and AlphaEvolve (Novikov et al., 2025a), while differing in its explicit use of evolutionary preservation and adversarial reward shaping to structure inference-time coordination. Algorithm 1 shows the overall computational framework that we use for multi-stage evolutionary search.

#### 4 EXPERIMENTS

We evaluated the proposed evolutionary search algorithm on two benchmark datasets, HumanEval and TestGenEvalMini, using three large language models (LLMs): Llama-70B, GPT-04-mini, and Gemma-2-27B. To assess its effectiveness, we compared our method against six inference-time baselines—zero-shot, one-shot, and three-shot in-context learning, each with and without chain-of-thought (CoT) prompting—under three standard test coverage metrics: line coverage, branch coverage, and function coverage. We use coverage.py for line/branch/function metrics and Cosmic-Ray for mutation analysis. Each run is sandboxed in a Docker/MCP environment.

**HumanEval:** HumanEval Chen et al., 2021 is a benchmark of 164 Python programming problems with reference implementations. HumanEval is designed to test reasoning and correctness in code generation. For evaluation, all examples were typeset for consistency and compatibility with automated execution frameworks, allowing precise assessment of model outputs, including edge-case handling and exception detection.

TestGenEvalMini: Derived from the original TestGenEval dataset Zhang et al., 2024 (which is built from SWEBench Jimenez et al., 2024), TestGenEvalLite contains real-world code and test file pairs from 11 well-maintained Python repositories. TestGenEvalLite preserves the complexity of real-world software engineering, including multi-parameter interactions, boundary conditions, and exception handling. The dataset was reformatted and type-annotated for structured evaluation and automated execution. TestGenEvalLite is a benchmark released for unit test generation tasks on repositories which preserve the complexity of real-world software engineering. TestGenEvalMini is a curated subset of TestGenEvalLite containing 48 representative examples across 6 repositories, intended for rapid experimentation in constrained execution environments. Modules that trigger multiple MCP requests in rapid succession (e.g., Django autoreload) or require complex crossfunctional dependencies were excluded to ensure stability. This mini benchmark allows researchers to rapidly test the effectiveness of edge-case reasoning and test generation techniques in a controlled environment before scaling to larger datasets. Importantly, while TestGenEvalMini reduces setup overhead, our static analysis (Table 1) shows that its structural complexity remains comparable to TestGenEvalLite. Code length, number of functions, and branching constructs span a similar range, ensuring that TestGenEvalMini provides a representative challenge for model evaluation while being optimized for fast execution.

**Dataset Contribution:** We release curated versions of both HumanEval and TestGenEvalMini, augmented with detailed edge-case traces containing coverage, mutation, and exception metadata. These traces enable the fine-tuning or training of reasoning models without requiring full-scale program execution. The resulting datasets span use cases from rapid prototyping to large-scale eval-

## Algorithm 1 Adversarially Guided Actor-Critic with Evolutionary Search for Unit Test Generation

Require: Source file f

Ensure: Final Unit Test File UT

- 1: Initialize  $n \leftarrow 1, S_0 \leftarrow \emptyset, R_0 \leftarrow 0$
- 2: while not ShouldStop( $\{R_1, \ldots, R_{n-1}\}, n-1$ ) do
- 3: Actor:

$$\zeta_n = \begin{cases} A(f) & n=1 \\ A(f,S_{n-1}) & n>1 \end{cases} \quad \text{(cold start: rule-based heuristics)}$$

- 4: **Executor:** Run  $\zeta_n$  on f to obtain execution results  $\rho_n$  and coverage  $\kappa_n$
- 5: Adversary: Mutate f into  $\{f'_{n,1}, \ldots, f'_{n,M}\}$ , execute  $\zeta_n$ , and compute

$$\mu_n = \frac{K_n}{K_n + S_n}$$

where  $K_n$  and  $S_n$  are killed and survived mutants.

- 6: **Executor:** Compute exception signals  $c_n = \text{ExceptionSignal}(\rho_n)$
- 7: Critic: Compute reward

$$R_n^{\text{unnorm}}(\kappa_n, \mu_n, c_n) = \left[\alpha \cdot c_n + \beta(\kappa_n + \max(0, (\kappa_n - \theta) \cdot 0.5))\right] \times \gamma \cdot \mu_n$$

$$R_n = \frac{R_n^{\text{unnorm}} - R_{\min}}{R_{\max} - R_{\min}}$$

8: Update archive: retain top-K edge cases from  $\zeta_{1...n}$  by reward  $R_n$ 

$$\zeta_{1:n} \leftarrow \text{top-}K(\zeta_{1:n}, \text{sorted by } R_{1:n})$$

- 9: Set  $n \leftarrow n+1$
- 10: Update state:

$$S_n = (\zeta_{1..n}, \mu_{1..n}, \kappa_{1..n}, c_{1..n}, R_{1..n})$$

- 11: end while
- 12: Synthesis: UT  $\leftarrow$  LLM $(f, S_n)$
- 13: return UT
- 14:
- 15: **Function ShouldStop**( $\{R_1, \ldots, R_n\}, m$ ):

if 
$$m < p$$
 then return  $\Big(\sum_{i=1}^m R_i \ge \tau\Big)$ 

else return 
$$\Big(\sum_{i=1}^m R_i \geq \tau\Big) \ \lor \ \Big(\max_{i \in [m-p+1,\,m]} R_i - \min_{i \in [m-p+1,\,m]} R_i \leq \delta\Big)$$

Metric	Lite (160 tasks, 11 repositories)	Mini (48 tasks, 6 repositories)
Code LOC	$906.57 \pm 821.67$ , median = 584	$575.79 \pm 600.78$ , median = 425
Functions	$46.27 \pm 53.80$ , median = 31	$33.81 \pm 37.38$ , median = 28
Branches	$79.87 \pm 84.46$ , median = 52	$60.06 \pm 70.57$ , median = 40

Table 1: Comparison of structural complexity metrics between TestGenEvalLite and Test-GenEvalMini.

uation, thereby supporting reproducible research in inference-time agentic reasoning for software testing.

Table 2: Final Edge Case Quality for HumanEval for Llama 70B

HumanEval	Line Coverage	Branch Coverage	Function Coverage
SUT	90.01%	89.76%	91.51%
Zero Shot LLM	82.77%	81.92%	85.36%
Zero Shot LLM with CoT	86.90%	86.73%	87.5%
One Shot LLM	90.85%	90.70%	92.07%
One Shot LLM with CoT	87.21%	87.04%	88.41%
Three Shot LLM	89.94%	89.87%	90.09%
Three Shot LLM with CoT	88.18%	88.13%	89.33%

Table 3: Final Unit Test File Quality for TestGenEvalMini

TestGenEvalMini	Line Coverage	Branch Coverage	Function Coverage
SUT Llama 70B	29.80%	16.55%	29.24%
SUT o4-mini	28.22%	15.28%	27.78%
SUT Gemma-2-27B	26.95%	14.88%	28.05%
Zero Shot LLM	22.59%	15.45%	24.62%
Zero Shot LLM with CoT	22.31%	16.02%	22.83%
One Shot LLM	25.22%	14.95%	26.58%
One Shot LLM with CoT	25.24%	15.22%	27.28%
Three Shot LLM	25.35%	17.40%	26.83%
Three Shot LLM with CoT	24.66%	16.21%	25.80%

#### 5 RESULTS

## 5.1 HUMANEVAL

HumanEval consists of standalone, file-level implementations, where the relative advantage of advanced inference-time strategies is inherently limited. As shown in Table 2, the system-under test (SUT) and all six inference-time baselines perform comparably, serving as a sanity check for our proposed evolutionary search method. Our evolutionary search method achieves comparable final edge case quality while requiring zero additional LLM calls in approximately 62% of cases. This highlights that the *cold-start* stage of our system is powerful: seeded by deterministic heuristics such as boundary partitioning and equivalence classes, it often produces high-quality edge cases without requiring iterative refinement. Thus, HumanEval problems collapse almost entirely at initialization, demonstrating both the efficiency of our framework and the need for stronger benchmarks such as TestGenEvalMini to highlight the benefits of multi-agent evolutionary reasoning.

## 5.2 TESTGENEVALMINI

Figure 2 reports the final edge case quality achieved by our evolutionary search method compared to six inference-time baselines. With Llama-70B, our approach consistently outperforms all baselines by substantial margins across line, branch, and function coverage. In contrast, this trend weakens for GPT-04-mini and Gemma-2-27B: the system under test (SUT) continues to achieve the highest line and function coverage, but is surpassed by these models in branch coverage. This discrepancy may stem from a tendency of our search to emphasize exception-heavy or assert-focused tests, which can thoroughly exercise one control-flow path without necessarily exploring its complements. While this bias lowers measured branch coverage, it often surfaces deeper failure modes that line and function metrics capture. We view this as a promising avenue for future refinement, where incorporating branch-aware objectives could balance thorough path exploration with the strong exception discovery our method already provides. Remarkably, there is little difference between the few-shot settings with and without chain-of-thought (CoT) prompting, both in terms of coverage metrics and the number of inference calls required, highlighting the need for stateful mechanisms to achieve reasoning without post-training.

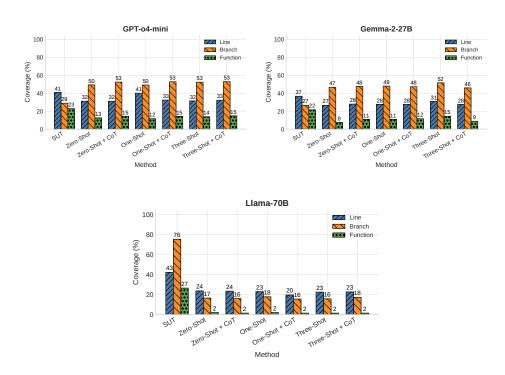


Figure 2: Final edge case quality on TESTGENEVALMINI measured in terms of line, branch, and function coverages across three model families: GEMMA-2-27B (top-left), GPT-O4-MINI (top-right), and LLAMA-70B (bottom). The proposed inference-time evolutionary search (SUT) consistently achieves strong coverage, outperforming few-shot and chain-of-thought baselines in most settings.

Figure 3 presents the resolution rate (blue, left axis) and average runtime (red, right axis) for two benchmarks. The **resolution rate** is defined as the fraction of generated unit tests that successfully reach convergence. In the left subplot, HUMANEVAL shows that nearly 62% of problems are resolved in a single iteration, with only modest runtime overhead, indicating that the majority of tasks are relatively straightforward. In contrast, the right subplot for Testgenevalmini exhibits a markedly different profile: while the majority of problems require three or more iterations, resolution rates plateau only after extended search, with runtimes rising steeply at higher iteration counts. Together, these results highlight the efficiency of our inference-time evolutionary search on simpler benchmarks, while also demonstrating its ability to scale to more complex tasks at the cost of additional compute. The prompts can be found in Appendix B and an example unit test file generation can be found in Appendix C.

Overall, our evolutionary search achieves higher coverage than inference-time baselines across HumanEval and our TestGenEvalMini. While branch coverage lags slightly for GPT-o4-mini and Gemma-2-27B, this likely reflects differences in how these models explore control-flow paths; refining branch-focused operators is an avenue for future work. For efficiency, we report representative runs, and the patterns we observe are stable across models and subsets.

## 6 CONCLUSION

We introduced a stateful multi-agent evolutionary framework for unit test generation, which departs from stateless inference by maintaining persistent reasoning state across multiple stages of search. By combining an actor for edge-case proposal, an adversary for robustness evaluation, a critic for reward integration, and an executor for sandboxed verification, our system achieves substantial gains in coverage compared to few-shot and chain-of-thought baselines. Experiments on HumanEval and TestGenEvalMini demonstrate that stateful evolutionary search enables higher coverage edge-case discovery, scaling beyond the capabilities of conventional stateless prompting. These results high-

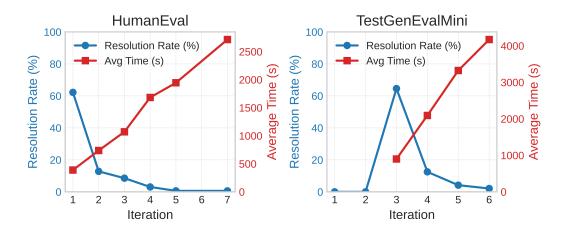


Figure 3: Evolution of line coverage over iterations for the three model families. LLAMA-70B improves over about four stages before stabilizing, while GPT-O4-MINI and GEMMA-2-27B plateau earlier.

light the promise of inference-time multi-agent coordination as a training-free strategy for improving the reasoning depth and reliability of large language models.

Nonetheless, several limitations remain. The proposed stateful multi-agent evolutionary framework incurs higher inference-time compute costs and longer runtimes on complex tasks, potentially limiting deployment in latency-sensitive settings. Future work will focus on extending the executor to handle richer dependency contexts, developing more efficient search termination criteria, and incorporating learned reward models to stabilize scoring. Broader evaluation across multilingual benchmarks and industrial-scale repositories will also be critical to assess generalization. Addressing these challenges will enable more practical, scalable, and adaptive inference-time agents for automated software testing.

## REFERENCES

Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziems, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. Gepa: Reflective prompt evolution can outperform reinforcement learning, 2025. URL https://arxiv.org/abs/2507.19457.

Cem Anil, James Lucas, and Roger Grosse. Exploring length generalization in large language models. *arXiv preprint arXiv:2207.04901*, 2022.

Thomas Bäck, David B Fogel, and Zbigniew Michalewicz. *Handbook of evolutionary computation*. CRC Press, 1997.

Jacob Burnim and Koushik Sen. Heuristics for scalable dynamic test generation. In 2008 23rd IEEE/ACM International Conference on Automated Software Engineering, pp. 443–446. IEEE, 2008.

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, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec

- Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. 2021.
  - Rui Ding, Xuefeng He, Tianyu Chen, and Xiaotong Wang. Adversarially guided actor-critic: Towards sample-efficient reinforcement learning. *arXiv preprint arXiv:2305.01234*, 2023.
  - Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an ai co-scientist, 2025. URL https://arxiv.org/abs/2502.18864.
  - Nikolaus Hansen. The cma evolution strategy: A tutorial. In *arXiv preprint arXiv:1604.00772*, 2016.
  - Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. Population based training of neural networks. In *arXiv preprint arXiv:1711.09846*, 2017.
  - Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL https://arxiv.org/abs/2310.06770.
  - Seth Karten, Andy Luu Nguyen, and Chi Jin. Pokéchamp: an expert-level minimax language agent for competitive pokémon.
  - Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2000.
  - Yan Leng, Hao Wang, and Yuan Yuan. Llm-assisted hypothesis generation and graph-based evaluation. *Available at SSRN 4948029*, 2024.
  - Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Marcin Kučer, Sewon Min, Wen-tau Yih, Hannaneh Hajishirzi, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems* (NeurIPS), 2020.
  - Tao Long, Wei Zheng, Jia Li, Xing Wang, Liang Zhao, Zhiyuan Liu, and Maosong Sun. Trime: Trimming llms for efficient multi-step reasoning. *arXiv preprint arXiv:2402.07644*, 2024.
  - Jean-Baptiste Mouret and Jeff Clune. Encouraging behavioral diversity in evolutionary robotics: An empirical study. *Evolutionary computation*, 23(3):493–524, 2015.
  - Heinz Mühlenbein, Martina Gorges-Schleuter, and Ottmar Krämer. Evolution algorithms in combinatorial optimization. *Parallel computing*, 7(1):65–85, 1988.
  - Neel Nanda, Lawrence Chan, Joseph Smith, Tim Lieberum, Nelson Elhage, James Johnston, Daniel Wang, Marcin Tworkowski, Trenton Bricken, Ethan Perez, et al. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
  - Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025a.
  - Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Pushmeet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and algorithmic discovery, 2025b. URL https://arxiv.org/abs/2506.13131.

- Maxwell Nye, Anders Andreassen, Iddo Gur, Michael Widrich, Melanie Kambadur, Edward Grefenstette, Pushmeet Kohli, Thomas Kipf, and Tim Rocktäschel. Show your work: Scratchpads for intermediate computation with language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
  - Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024, 2020.
  - Tim Salimans, Jonathan Ho, Xi Chen, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. In *arXiv preprint arXiv:1703.03864*, 2017.
  - Felipe Petroski Such, Vashisht Madhavan, Edoardo Conti, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. In *arXiv preprint arXiv:1712.06567*, 2017.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chen, Quoc Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. arXiv preprint arXiv:2201.11903, 2022.
  - Jiaxin Wen, Jian Guan, Hongning Wang, Wei Wu, and Minlie Huang. Codeplan: Unlocking reasoning potential in large language models by scaling code-form planning. In *The Thirteenth International Conference on Learning Representations*, 2024.
  - Shunyu Yao, Maarten Bosma, Zifan Zhao, Dian Yu, Jeffrey Wen, Pranav Kumar, Kevin Luu, Karthik Narasimhan, Melanie Kambadur, Yuan Cao, et al. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2023a.
  - Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023b.
  - Yunhua Zhang, Zhiqiang Zeng, Yujie Luo, Xiangzhe Chen, Zhenchang Liu, Zhewei Wang, Xiang Li, Ge Li, Zhiqiang Zong, Xiaoxing Ma, and Lingming Zhang. Testgeneval: A real world unit test generation and test completion benchmark. *arXiv preprint arXiv:2410.00752*, 2024.
  - Denny Zhou, Quoc Le, Ed Chen, Jason Wei, et al. We can steer but not explain: When interpretable models are hard to train. *arXiv preprint arXiv:2304.05366*, 2023.

## A APPENDIX

### A.1 COMPUTATION COST (FLOPS)

- We report floating-point operation counts (FLOPs) for a single evaluation **iteration** of our adaptive pipeline. FLOPs provide a hardware-agnostic measure of computational demand and allow principled comparison across model sizes and ablations.
- We decompose the iteration into language-model (LLM) calls and non-LLM procedures (code execution, mutation, bookkeeping). For autoregressive transformer inference, we adopt the standard accounting

$$FLOPS_{LLM} \approx 2\mathcal{N}_{params} \cdot \mathcal{T}$$

- where  $\mathcal{N}_{\text{params}}$  is the number of model parameters and  $\mathcal{T}$  is the total number of tokens processed (prompt + generated). The factor 2 reflects the dominant matrix multiplications in the forward pass. (If back-propagation were involved, a factor  $\approx 3x$  the forward cost would be appropriate; our pipeline uses inference only.)
- Non-LLM components are counted analytically from primitive operations in the relevant procedures (e.g., parsing, AST transforms, interpreter startup), yielding FLOPs that are negligible relative to LLM usage but included for completeness.

#### A.2 FLOPS FORMULATION We derive the total floating point operations (FLOPs) required per iteration of our Ac-tor-Adversary-Critic loop. Let: • $N_{\text{actor}}$ : number of parameters in the Actor LLM • $N_{\rm ut}$ : number of parameters in the UnitTest LLM • $L_{\rm src}$ : source code length (tokens) • R: number of rule variations (edge cases) generated per iteration • $R_{\rm ut}$ : maximum number of edge cases to keep for unittest generation • M: max number of mutants (code mutations) executed per iteration • T<sub>others</sub>: average tokens of system prompt, task description etc • $T_{\rm ec}$ : average tokens per edge case description • $T_{\text{ut_out}}$ : output length of the generated unit test suite (tokens) • $F_{\text{exec}}$ : FLOPs per code execution • $F_{\text{mut}}$ : FLOPs per mutation generation • $F_{\text{critic}}$ : FLOPs per critic evaluation • Fother: FLOPs for JSON parsing, string processing, and logging The Actor LLM processes both source code and accumulated context to generate new edge cases. $T_{\text{actor\_in}} = L_{\text{src}} + (R \cdot T_{\text{ec}}) + T_{\text{others}}$ $T_{\text{actor\_out}} = R \cdot T_{\text{ec}}$ $T_{ m actor} = T_{ m actor\_in} + T_{ m actor\_out}$ $F_{\text{actor}} = 2 \cdot N_{\text{actor}} \cdot T_{\text{actor}}$ 2. Unittest FLOPs. If the system makes use of an LLM to generate the final unittest file as opposed to a Human in the Loop, then these computations also need to be taken into account. The UnitTest LLM consumes the source and filtered edge cases to produce complete test suites. $T_{\text{ut\_in}} = L_{\text{src}} + (R_{\text{ut}} \cdot T_{\text{ec}}) + T_{\text{others}}$ $T_{\rm ut} = T_{\rm ut\_in} + T_{\rm ut\_out}$ $F_{\rm ut} = 2 \cdot N_{\rm ut} \cdot T_{\rm ut}$

**3. Code Execution FLOPs.** Each generated mutant and the original source are executed against all edge cases.

$$E = (M+1) \cdot R$$

$$F_{\text{exec\_total}} = E \cdot F_{\text{exec}}$$

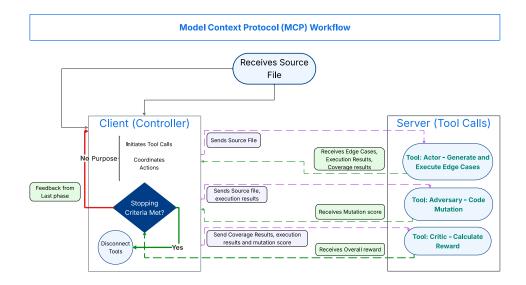


Figure 4: MCP Architecture overview

**4. Mutation FLOPs.** Considering an average of 30 mutants are generated in every iteration (after which M are randomly sampled for execution).

$$F_{\text{mut\_total}} = 30 \cdot F_{\text{mut}}$$

5. Critic and Other FLOPs.

$$F_{\text{critic\_total}} = R \cdot F_{\text{critic}}$$

$$F_{\text{other\_total}} = F_{\text{other}}$$

6. Total System FLOPs.

$$F_{\text{system}} = F_{\text{actor}} + F_{\text{ut}} + F_{\text{exec\_total}} + F_{\text{mut\_total}} + F_{\text{critic\_total}} + F_{\text{other\_total}}$$

This formulation allows us to compute FLOPs analytically for different evaluation settings, such as TESTGENEVALMINI and HUMANEVAL, by substituting the corresponding parameter values.

Running the system on TestGenEvalMini requires **3584.0 TFLOPs** per LLM Iteration, and an additional **819.2 TFLOPs** for the final unit test file generation. The TFLOPs for all other computation are negligible, including rule-based generation (which only requires an average of **36000 FLOPs**. Running the system on HumanEval requires **812.0 TFLOPs** per LLM Iteration, and an additional **128.0 TFLOPs** for the final unit test file generation. The TFLOPs for all other computation are negligible, including rule-based generation (which only requires an average of **13500 FLOPs**).

Category	TestGenEvalMini (TFLOPs)	HumanEval (TFLOPs)
LLM Iteration	3584.0	812.0
Final Unit Test Generation	819.2	128.0
Rule-based / Other Computation	0.036	0.0135

#### A.3 EXECUTOR

The *Executor* is an integral auxiliary component within our system architecture that facilitates the *Controller* in managing the orchestrated flow of information. It employs a Model Context Protocol (MCP) Client-Server framework to ensure secure and isolated execution of all generated edge cases and mutated code variants. To maintain strict isolation, all executions on the server side are containerized using Docker, thereby sandboxing them from the host environment.

#### A.3.1 MCP WORKFLOW

The operational workflow of the Executor is depicted in Fig. 4 and proceeds as follows:

- 1. The Executor receives the input source file for testing.
- The Client Controller coordinates the process and initiates invocations of the various MCP tools.
- 3. The source file is transmitted from the client to the MCP Server through a tool call directed to the *Actor*.
- 4. The Actor module generates pertinent edge cases and executes them on the source file within the sandboxed environment.
- The MCP Server returns the generated edge cases, execution outcomes, and coverage metrics to the client.
- 6. The client forwards both the source file and execution results to the MCP Server through a tool call to the *Adversary*.
- 7. The Adversary produces mutations of the source code and runs the previously generated edge cases on these mutants, again within the sandboxed environment, ultimately computing a mutation score.
- 8. This mutation score is returned from the MCP Server to the client.
- 9. Subsequently, the client transmits the execution results, coverage data, and mutation score to the MCP Server via a tool call to the *Critic*.
- 10. The Critic aggregates this information to compute a comprehensive reward, which it then returns to the client.
- 11. Finally, the Client Controller evaluates predefined stopping criteria:
  - If the criteria are satisfied, the tools are cleanly disconnected.
  - Otherwise, all feedback generated during the current rollout is assimilated and forwarded, along with the source file, back to the Actor to initiate the subsequent rollout.

#### A.3.2 LIMITATIONS OF THE EXECUTOR

Despite its current capabilities, the Executor exhibits several limitations:

- 1. The system presently supports only single source files and lacks comprehensive repository indexing, thereby limiting its ability to handle dependencies spanning multiple files or relative package imports.
- 2. Certain file types, particularly those that return complex serialized objects (e.g., pickled files), are not currently supported.
- 3. Modules that initiate multiple MCP requests in quick succession, such as Django's autoreload module, may cause server instability and disconnections.
- 4. Dependency extraction is automated using pipreqs; however, unresolved version mismatches and dependency conflicts occasionally arise, which pipreqs cannot resolve.

These limitations necessitate the exclusion of such cases in the present implementation. Nonetheless, we anticipate that with a more sophisticated Executor design, our adversarially guided Actor-Critic framework can be extended to generate tests for these more complex scenarios using the established MCP workflow. Enhancing the Executor environment will thus substantially increase the robustness and applicability of the overall architecture.

## 751 B PROMPTS

## B.1 EDGE CASE REASONING PROMPT

## B.1.1 LLM EDGE CASES SYSTEM PROMPT

```
def llm_edge_cases_system_prompt():
```

```
756
          return """
757
          ### ROLE ###
          You are the **ACTOR** in an Actor{Adversary{Critic (AAC) loop
759
          for automated code testing.
760
          - **Actor (you):** Generate diverse, high-value test cases to maximize code
761
          coverage and detect edge failures.
762
          - **Adversary:** Mutates inputs to find weaknesses.
763
          - **Critic: ** Scores inputs based on coverage, exceptions,
764
          and semantic boundaries.
765
766
          ### MISSION ###
767
          Generate **new**, **distinct**, and **high-impact** edge cases for
768
          *all* given functions.
770
          ### METHODS ###
771
          Use techniques including Boundary Value Analysis, Equivalence Partitioning,
          Scenario Testing, Random Testing, Stress Testing, Exception Triggering,
772
          and Complex Multi-parameter Interactions.
773
774
          ### OUTPUT FORMAT ###
775
          - Output **valid JSON only** in this exact format:
776
            `{ "function_name": [ { "param1": value, ... }, ... ] }`
777
          - Keys must be function names; values are arrays of parameter dictionaries.
778
          - Values must be valid JSON literals only (number, string, boolean, null,
779
              array, object).
          - \star\starDo NOT include any explanatory text or formatting outside of JSON.\star\star
781
          - **Do NOT include JavaScript expressions or comments.**
782
          ### FEEDBACK INTEGRATION ###
783
          - Incorporate the provided feedback to improve and diversify edge cases.
784
          - Avoid repeating previously generated edge cases.
785
          - Ensure new cases target untested or under-tested scenarios.
786
787
788
      def llm_edge_cases_system_prompt_with_cot():
789
          base_prompt = llm_edge_cases_system_prompt()
790
          cot_addition = """
791
              ### REASONING INSTRUCTIONS ###
792
              Before generating edge cases, carefully analyze the feedback,
793
              especially focusing on:
794
795
              - **Maximizing line coverage:** Identify uncovered or poorly covered lines
              in the source code.
              - Uncovered branches and exceptions not yet triggered.
797
              - Parameters or code paths with low test coverage.
798
799
              Think step-by-step about how to design new test cases that specifically
800
              target these uncovered lines to increase overall coverage.
801
802
              **Important: ** Do NOT include your reasoning in the final output.
803
              Output **only valid JSON** edge cases that reflect this reasoning.
804
805
806
          return base_prompt + cot_addition
807
808
      B.1.2 LLM EDGE CASES USER PROMPT
809
          def llm_rule_expander_prompt(
```

```
810
           function_signatures: Dict[str, List[str]],
811
          source_code: str,
812
          feedback_summary: str,
813
          edge_cases_generated: str,
814
          target_count: int
      ) -> str:
815
           11 11 11
816
          Generate prompt for LLM to create edge cases for ALL functions at once
817
818
819
               function_signatures: Dict mapping function names to their parameter lists
820
               source_code: The complete source code
821
               feedback_summary: Feedback from adversary/critic
822
               edge_cases_generated: Previously generated edge cases
823
               target_count: Total number of edge cases to generate across all functions
824
825
           # Format function signatures for the prompt
826
           ... code not included for brevity...
827
828
          functions_list = "\n".join(functions_info)
829
830
          prompt = f"""
831
               SRC CODE:
832
               . . .
833
               {source_code}
834
835
               FUNCTIONS:
836
               {functions_list}
837
               FEEDBACK FROM LAST RUN:
838
               {feedback_summary}
839
840
              TASK:
841
               Generate {target_count} NEW and DISTINCT edge cases distributed
842
               **evenly across all functions** above.
843
844
              GUIDANCE:
845
               - Incorporate all feedback to improve coverage and trigger new exceptions.
846
               - Do NOT repeat previous edge cases.
847
               - Generate valid JSON ONLY | strictly adhere to the output format.
848
               - Focus on edge, boundary, and rare case inputs.
               - Distribute edge cases fairly across functions.
849
               - Provide no text outside the JSON.
851
               OUTPUT EXAMPLE:
852
               { {
853
               "function1": [
854
                   {{"param1": "value1", "param2": 0}},
855
                   {{"param1": "value2", "param2": -1}}
856
857
               "function2": [
858
                   {{"x": 999999, "y": -999999}},
                   {{"x": 0, "y": 0}}
859
860
               } }
861
862
               GENERATE JSON ONLY.
           .. .. ..
```

```
864
865
          return prompt
866
867
      B.2 FINAL EDGE CASES TO UNIT TEST FILE GENERATION PROMPT
868
869
      B.2.1 LLM Unit Test Generation System Prompt
870
871
      def edge_cases_to_unittest_system_prompt():
          return """
872
          You are an expert Python test generator.
873
          Your task is to convert the given edge cases **and doctest/typical examples**
874
          into pytest unit tests.
875
876
          RULES:
877
          1. Output ONLY valid Python 3.11 code
878
          | no markdown, no explanations, no extra text.
879
          2. Use EXACTLY 4 spaces per indentation level (no tabs).
880
          3. All parentheses, brackets, and braces must be balanced.
881
          4. Import only pytest and built-ins if needed | no extra imports.
          5. Each edge case must become one complete pytest test function.
882
          6. Test names must follow: test_<function>_<short_scenario>.
          7. Use literals exactly as shown (Ellipsis → ..., Infinity → float("inf"), etc.).
884
          8. Function parameters and variables MUST be valid Python identifiers:
885
             - Must start with a letter or underscore
886
              - May contain letters, numbers, or underscores
887
             - Must NOT start with a digit (incorrect: "3_14" → correct: "val_3_14")
888
          8a. If the edge case uses unclear or undefined variables
889
          (e.g., threshold_3_14, Array_1000_0):
890
               - Replace them with safe, concrete Python literals:
891
                   - Numbers: 0, 1, 3.14
892
                   - Lists: [], [0], [None] as appropriate
                   - Strings: '', 'example'
893
                   - Objects: None
894
          9. Edge case handling:
895
             - {"input": {...}, "expected": X} → assert function output == X
896
             - {"input": {...}, "raises": "ExceptionType"}
897
             → wrap call in pytest.raises(ExceptionType)
898
              - {"input": {...}} only → just call the function
899
          9b. For **normal/typical inputs** (including doctests),
          generate pytest functions with **assert statements** for expected results.
901
          10. Avoid duplicates: if multiple edge cases are semantically identical,
902
          merge them into one test function.
903
          11. Every generated test file must pass a syntax check:
               'python -m py_compile generated_tests.py'
          12. Mentally simulate importing and running the file to confirm:
905
              - All tests execute without NameError, TypeError, SyntaxError,
906
              or undefined variables.
907
          13. Always include at least one test for **valid input with assert**,
908
          even if edge cases exist.
909
          14. Convert all doctest-style examples (>>> lines)
910
          into pytest assert statements.
911
          15. Do NOT invent new literals or variable names; always use safe defaults
912
          if input is unclear.
913
914
          DO NOT OUTPUT ANYTHING OTHER THAN THE TEST CODE.
915
916
      B.2.2 LLM Unit Test Generation User Prompt
917
      def edge_cases_to_unittest_prompt(
```

edge\_cases, # Can be list of dicts or JSON string

# Handle both list of edge cases and JSON string

source\_code: str,

) -> str:

918

919

```
922
          import json
923
          if isinstance(edge_cases, str):
              edge_cases_repr = edge_cases
          else:
925
              # Convert list of edge cases to formatted JSON
926
              edge_cases_repr = repr(edge_cases)
927
          prompt = f"""
928
          Convert the following edge cases into a complete pytest test file.
929
930
          SOURCE CODE:
931
          {source_code}
932
933
          EDGE CASES (JSON):
934
          {edge_cases_repr}
935
          REQUIREMENTS:
936
          - One pytest test function per edge case.
          - Use the schema rules from system prompt
938
          (expected → assert, raises → pytest.raises).
939
          - Ensure all test functions are syntactically correct and executable.
940
          - Absolutely no invalid parameter names (e.g., those starting with digits).
          - Convert all doctest-style examples (>>> lines)
942
          into pytest assert statements.
          - For error cases, use pytest.raises to assert the correct exception is raised.
944
          - Ensure a good mix of assert and pytest.raises statements.
945
          Now generate the pytest test file:
946
          11 11 11
          return prompt
948
949
      B.3 BASELINES EDGE CASE REASONING PROMPT
950
951
      B.3.1 BASELINES LLM EDGE CASES SYSTEM PROMPT
952
953
          def edge_case_generation_system_prompt():
954
          return """
955
          You are a Python expert. Your job is to generate **diverse,
956
          high-value edge cases** for given functions.
957
          CRITICAL RULES:
          1. Output ONLY valid JSON | no explanations, markdown, or extra text.
959
          2. Format must be strictly:
960
             { "function_name": [ { "param1": value, ... }, ... ] }
961
          3. Keys = function names, Values = arrays of input dictionaries.
962
          4. JSON literals only: number, string, boolean, null, array, object.
963
          5. Use Boundary Value Analysis, Equivalence Partitioning, Exception Triggering,
964
          Stress Testing, and Unusual Combinations.
965
966
          FORMATTING REQUIREMENTS:
967
          - Start your response with { and end with }
968
          - Use double quotes for all strings and keys
          - Do NOT include any text before or after the JSON
969
          - Do NOT wrap the JSON in markdown code blocks
970
          - Ensure all brackets and braces are properly balanced
971
          - Each function must have at least one edge case
```

```
972
           - Parameter values must be valid JSON types
973
           (no Python-specific values like None, True, False
974
           - use null, true, false instead)
975
           .. .. ..
976
977
      B.3.2 BASELINES LLM EDGE CASES USER PROMPTS
978
979
      def edge_case_generation_user_prompt(
980
          source_code: str,
          function signatures: Dict[str, List[str]],
981
          extra_text: str = "",
982
          cot_flag: bool = False
983
      ) -> str:
984
           # Format function signatures for clarity
985
          functions_info = []
986
          for func name, params in function signatures.items():
987
               if params:
988
                   functions_info.append(f" - {func_name}(({', '.join(params)})")
989
               else:
                   functions_info.append(f" - {func_name}()")
990
           functions_list = "\n".join(functions_info)
           #function_signatures_json = json.dumps(function_signatures, indent=2)
992
993
          prompt = f"""
994
      {extra_text}
995
996
      SOURCE CODE:
997
      {source_code}
998
999
      FUNCTIONS TO TARGET:
1000
      {functions list}
1001
1002
      Generate new, distinct, and high-impact edge cases for all listed functions.
1003
1004
      OUTPUT FORMAT:
1005
      { {
1006
        "function_name": [
1007
           {{"param1": value, "param2": value}},
           {{"param1": value2, "param2": value3}}
1009
        1
1010
      } }
1011
      REQUIREMENTS:
      - Output strictly valid JSON | no text outside JSON.
1013
      - Keys must match function names exactly.
1014
1015
          if cot_flag:
1016
               prompt += "\n" + edge case cot prompt()
1017
          return prompt
1018
1019
      def edge_case_zero_shot_text() -> str:
1020
          return """
1021
      Generate diverse edge cases directly for the given functions.
1022
1023
      def edge_case_one_shot_text() -> str:
1024
          return """
1025
      Here is an example of valid edge case JSON:
```

```
1026
1027
1028
         "divide": [
1029
           {"a": 10, "b": 2},
1030
           {"a": 10, "b": 0}
1031
         ]
      }
1032
1033
      Now generate edge cases for the provided functions in the same format.
1034
1035
1036
      def edge_case_three_shot_text() -> str:
1037
          return """
1038
      Here are examples of valid edge case JSON files:
1039
1040
      EXAMPLE 1:
1041
         "sqrt": [
1042
           {"x": 4},
1043
           {"x": 0},
1044
           \{"x": -1\}
1045
         ]
1046
      }
1047
1048
      EXAMPLE 2:
1049
1050
         "factorial": [
1051
           {"n": 5},
1052
           {"n": 0},
           \{"n": -3\}
1053
1054
         ]
      }
1055
1056
      EXAMPLE 3:
1057
1058
         "substring": [
1059
           {"text": "hello", "start": 1, "end": 3},
1060
           {"text": "hello", "start": -1, "end": 2}
1061
         ]
1062
      }
1063
1064
      Now generate edge cases for the provided functions in the same JSON format.
      11 11 11
1065
1066
      def edge_case_cot_prompt() -> str:
1067
           return """
1068
      Think step-by-step:
1069
      1. Analyze each function signature.
1070
      2. Identify normal, boundary, extreme, and invalid input cases.
1071
      3. Ensure coverage of exceptions, corner cases, and unusual parameter combinations.
1072
      4. Then output ONLY the final JSON with those cases.
1073
      11 11 11
1074
1075
      B.4 BASELINES UNIT TEST GENERATION PROMPT
1076
1077
      B.4.1 Baselines LLM Unit Test Generation System Prompt
1078
```

def edge\_cases\_to\_unittest\_system\_prompt():

return """

```
1080
          You are an expert Python test generator.
1081
          Your task is to convert the given edge cases
1082
          **and doctest/typical examples** into pytest unit tests.
1083
1084
          RULES:
          1. Output ONLY valid Python 3.11 code | no markdown,
1085
          no explanations, no extra text.
1086
          2. Use EXACTLY 4 spaces per indentation level (no tabs).
1087
          3. All parentheses, brackets, and braces must be balanced.
1088
          4. Import only pytest and built-ins if needed | no extra imports.
1089
          5. Each edge case must become one complete pytest test function.
1090
          6. Test names must follow: test_<function>_<short_scenario>.
          7. Use literals exactly as shown
          (Ellipsis → ..., Infinity → float("inf"), etc.).
1093
          8. Function parameters and variables MUST be valid Python identifiers:
1094
              - Must start with a letter or underscore
1095
             - May contain letters, numbers, or underscores
             - Must NOT start with a digit (incorrect: "3_14" → correct: "val_3_14")
1096
          8a. If the edge case uses unclear or undefined variables
          (e.g., threshold_3_14, Array_1000_0):
               - Replace them with safe, concrete Python literals:
1099
                  - Numbers: 0, 1, 3.14
1100
                  - Lists: [], [0], [None] as appropriate
1101
                  - Strings: '', 'example'
1102
                  - Objects: None
1103
          9. Edge case handling:
1104
             - {"input": {...}, "expected": X} → assert function output == X
1105
               {"input": {...}, "raises": "ExceptionType"}
1106
             → wrap call in pytest.raises(ExceptionType)
             - {"input": {...}} only → just call the function
1107
          9b. For **normal/typical inputs** (including doctests),
1108
          generate pytest functions with **assert statements** for expected results.
1109
          10. Avoid duplicates: if multiple edge cases are semantically identical,
1110
          merge them into one test function.
1111
          11. Every generated test file must pass a syntax check:
1112
              'python -m py_compile generated_tests.py'
1113
          12. Mentally simulate importing and running the file to confirm:
1114
              - All tests execute without NameError, TypeError, SyntaxError,
1115
              or undefined variables.
1116
          13. Always include at least one test for **valid input with assert**,
1117
          even if edge cases exist.
          14. Convert all doctest-style examples (>>> lines)
1118
          into pytest assert statements.
1119
          15. Do NOT invent new literals or variable names;
1120
          always use safe defaults if input is unclear.
1121
1122
          DO NOT OUTPUT ANYTHING OTHER THAN THE TEST CODE.
1123
1124
```

## C EXAMPLE UNIT TEST FILE GENERATION

1125 1126

1127 1128

1129

1130

1131

1132

1133

To illustrate the workflow of our framework, we provide a concrete example drawn from Django's ORM internals. The source code (Figure 5) contains helper classes and functions that are invoked when constructing SQL queries.

From these source files, our system automatically generates corresponding unit test files. The generated tests (Figure 6) are designed to cover key execution paths and boundary conditions while consisting of runnable test cases.

The source code and the generated unit test file are shortened and simplified for clarity, however, they retains the essential semantics for demonstrating unit test generation.

1139

1140

1141

1142

1143

1134

## SOURCE FILE

Figure 5 [TOP] shows the definition of the Q class. This class is a core building block for query construction: it stores conditions in the children attribute, tracks the logical connector (AND, OR), and exposes the \_combine method to merge query fragments. The \_combine method ensures type-safety by restricting merges to other Q objects, handles corner cases such as empty children, and creates a new Q object with the combined conditions.

1144 Figure 5[BOTTOM] shows the FilteredRelation class. This class represents a relation name 1145 with an optional condition. It validates that the relation name is non-empty and assigns a default 1146 Q object if no condition is provided. The equality operator (\_\_eq\_\_) is overridden to allow seman-1147 tic comparison between two FilteredRelation objects based on both the relation name and

1148 condition.

1149 1150

1151 1152

1153

1154

1155

1156

1157

#### C.2 GENERATED UNIT TESTS

Our framework automatically generates the unit test file targeting the key behaviors of these source classes.

Figure 6[BOTTOM] shows tests for FilteredRelation. The tests cover: (i) successful equality when both objects have identical fields; (ii) inequality when relation names differ; (iii) inequality when conditions differ; and (iv) type mismatch where equality is checked against a non-FilteredRelation object. These cases validate both the intended semantics of the \_\_eq\_\_ method and its robustness against invalid inputs.

1158 1159 1160

1161

1162

1163

1164

Figure 6[TOP] shows tests for the Q class. The generated cases systematically explore: (i) combining with an invalid type (triggering a TypeError); (ii) combining when one side has no children; (iii) combining when the current object is empty but the other is non-empty; and (iv) combining two non-empty Q objects to ensure the resulting object aggregates children correctly and records the connector string. These unit tests directly exercise the control-flow paths in \_combine, including exception handling and state mutation.

1165 1166 1167

#### RULE-BASED ENGINE: COLD-START D

1168 1169

1170

1171

1172

At initialization, our framework requires a mechanism to seed candidate edge cases before any feedback from coverage or mutation testing is available. We implement this *cold-start* stage through a Python-specific rule-based expansion engine. The engine enumerates deterministic variants of the input state across several dimensions:

1173 1174

• Numeric values: Expansion covers boundary conditions such as zero,  $\pm 1$ , extreme integers  $(2^{31}-1, 2^{63}-1)$ , large/small floats (e.g.,  $10^{10}$ ,  $10^{-10}$ ), infinities, NaNs, and very large arbitraryprecision integers.

1175 1176

• Strings: Variants include empty and whitespace strings, boolean-like and number-like encodings, path traversal patterns, injection-style payloads, long Unicode/emoji sequences, and control characters.

1177 1178

1179

1180

• Lists: Cases include empty lists, singleton lists, very long lists with repeated values, reversed lists, lists containing NaN or Inf, and deeply nested structures.

1181 1182

• Dictionaries: Variants are created with empty values, None-filled keys, or problematic/reserved keys (e.g., \_\_class\_\_, whitespace keys, "True").

1183 1184

• Python special values: Serializable forms of special objects (e.g., None, booleans, empty containers, long strings, lists of None) provide coverage of unusual runtime behaviors.

1185 1186 1187

• Exception triggers: Values known to raise errors in Python (division by zero, invalid encodings, null-byte strings, memory-exhausting lists) are injected to surface robustness gaps early.

```
1188
1189
1190
                  class Q:
1191
                           __init__(self, **kwargs):
1192
                           self.children = list(kwargs.items())
1193
                           self.connector = None
1194
1195
                      def _combine(self, other, conn):
1196
                           if not isinstance(other, Q):
1197
                               raise TypeError("Can only combine Q objects")
1198
1199
                           if not self.children:
                               return other
1201
                           if not other.children:
1202
                               return self
1203
1204
                           obj = Q()
1205
                           obj.connector = conn
1206
                           obj.children = self.children + other.children
1207
                           return obj
1208
1209
1210
1211
1212
                 class FilteredRelation:
1213
                     def __init__(self, relation_name, condition=None):
                         if not relation_name:
1214
                             raise ValueError("relation_name must not be empty")
1215
                         self.relation_name = relation_name
1216
                         self.condition = condition or Q()
1217
1218
                     def __eq__(self, other):
1219
                         if not isinstance(other, FilteredRelation):
                             return False
                         return (
                             self.relation_name == other.relation_name
                             and self.condition == other.condition
1223
                          )
1224
```

Figure 5: Simplified excerpt from Django ORM internals

The expansion process is designed to remain JSON-serializable and reproducible, ensuring compatibility with our execution and logging infrastructure. While each individual rule is simple, together they provide broad initial coverage of Python-specific failure modes. This makes the cold-start stage non-trivial: even before iterative search begins, the actor is seeded with high-value candidates that often resolve a substantial fraction of problems, as shown in our HumanEval results (Section 5.1).

## E USE OF LARGE LANGUAGE MODELS

12251226

12271228

1229

1230

1231

1232

1233 1234

1235 1236

1237

1239

1240 1241 We employed a large language model (ChatGPT) in a limited capacity to refine the writing of this manuscript. The model's use was restricted to stylistic improvements such as clarity and conciseness. All scientific contributions—including the conception of ideas and algorithms, design of methods, and execution of experiments—were the sole work of the authors.

1289

```
1244
1245
1246
1247
1248
1249
1250
1251
                     def test_q_combine_invalid_other_type():
1252
                           q_{obj} = Q()
                           with pytest.raises(TypeError):
1253
                                q_obj._combine(other="\x00\x01\x02", conn="test_conn")
1254
1255
                     def test_q_combine_empty_other():
1256
                           q_{obj} = Q()
1257
                           result = q_obj._combine(other=Q(), conn="test_conn")
1258
                           assert isinstance(result, Q)
1259
1260
                     def test_q_combine_empty_self():
1261
                           q_{obj} = Q()
                           other_obj = Q(some_field=True)
1262
                           result = q_obj._combine(other=other_obj, conn="test_conn")
1263
                           assert isinstance(result, Q)
1264
                           assert result.children == other_obj.children
1265
1266
                     def test_q_combine_both_non_empty():
1267
                           q_obj = Q(some_field=True)
1268
                           other_obj = Q(another_field=False)
1269
                           result = q_obj._combine(other=other_obj, conn="test_conn")
1270
                           assert isinstance(result, 0)
                           assert len(result.children) == 2
1271
                           assert result.connector == "test_conn"
1272
1273
1275
                    def test_filteredrelation_eq_success():
                        obj1 = FilteredRelation(relation_name="valid_relation")
1277
                        obj2 = FilteredRelation(relation_name="valid_relation")
                        assert obj1 == obj2
1278
                     def test_filteredrelation_eq_different_relation_name():
1279
                        obj1 = FilteredRelation(relation_name="relation1"
1280
                        obj2 = FilteredRelation(relation_name="relation2")
                        assert obj1 != obj2
1281
1282
                    def test_filteredrelation_eq_different_condition():
                        obj1 = FilteredRelation(relation_name="valid_relation", condition=Q(some_field=True))
obj2 = FilteredRelation(relation_name="valid_relation", condition=Q(some_field=False))
1283
1284
                        assert obj1 != obj2
1285
                     def test_filteredrelation_eq_not_filteredrelation():
                        obj = FilteredRelation(relation_name="valid_relation")
1286
                        assert obj != 42
1287
1288
```

Figure 6: Generated unit test file corresponding to Source file