

Human-in-the-Loop Automated Software Testing Enhancing Coverage and Reducing False Positives Through Interactive Machine Learning

Neil Javier,

Computational Biology Scientist, USA.

Citation: Javier, N. (2025). Human-in-the-Loop Automated Software Testing Enhancing Coverage and Reducing False Positives Through Interactive Machine Learning. *International Journal of Scientific Research in Computer Science and Information Technology (IJSRCSIT)*, 6(2), 1-7.

Abstract

Software testing is a critical component of software development, yet traditional automation often struggles with low coverage and high false positives. This paper introduces a Human-in-the-Loop (HITL) automated software testing framework that leverages interactive machine learning (IML) to enhance test coverage and reduce false positives. By incorporating domain experts into the machine learning loop, the system dynamically refines test case generation and bug classification in real time. We propose a hybrid architecture integrating human feedback loops with test synthesis and prioritization modules. Experimental validation on multiple open-source projects demonstrates improved precision and recall, especially in anomaly detection and regression testing. This work advances the field of intelligent software testing by aligning human insight with algorithmic rigor, enabling more reliable and scalable QA workflows.

Keywords: Software Testing, Human-in-the-Loop, Interactive Machine Learning, Test Coverage, False Positives, Program Analysis, Automation, QA Engineering

1. Introduction

1.1 Improving Test Case Relevance

Automated software testing has made significant strides with the integration of static and dynamic analysis, symbolic execution, and machine learning. However, automated systems still suffer from **limited contextual understanding**, leading to insufficient test coverage and high false alarm rates.

In contrast, **Human-in-the-Loop (HITL)** methods, particularly through **Interactive Machine Learning (IML)**, allow systems to learn from human feedback, rapidly adapting to edge cases and domain-specific behaviors. This paper explores how HITL methodologies can enhance test automation. In conventional automated software testing, test cases are often generated by heuristics, symbolic execution, or random input generation. While efficient, these methods frequently produce **irrelevant or redundant test cases** that do not align with real-world usage or critical functional paths. This results in **low test efficiency** and the potential overlooking of business-critical issues.

The incorporation of human feedback into the test generation loop significantly **increases the semantic relevance** of test cases. By allowing domain experts or developers to provide guidance on what constitutes a meaningful input, behavior, or path, the system can:

- Prioritize functions or components known to be sensitive or business-critical.
- Avoid generating tests for deprecated or non-functional features.
- Tailor tests to real-world user behavior, especially in complex GUI or API environments.

Over time, interactive learning enables the test generation model to learn from human feedback, effectively **re-weighting its priorities** to align with software goals and user expectations.

1.2 Dynamically Reducing Redundant or Misleading Alerts:

A major limitation of many automated testing tools is their **high false positive rate**. Automated test suites, especially those incorporating static analysis or property-based fuzzing, may flag harmless anomalies, configuration-specific behaviors, or even legacy artifacts as potential bugs.

1.3 Enabling Faster Adaptation to Novel Software Behaviors

Modern software evolves rapidly, with frequent updates, new modules, and changing external dependencies. Traditional static test generation struggles to **keep pace with dynamic changes**, often requiring complete re-generation of tests or manual reconfiguration.

In contrast, HITL systems can **rapidly incorporate knowledge about new or unexpected behaviors**:

- Developers can label anomalous outputs from new code segments.
- Human feedback provides contextual insights that models cannot infer from code alone.
- The system continuously updates its models to understand new workflows, input types, or execution patterns.

This results in a test system that is not only **reactive but adaptive**, capable of updating itself on the fly without the need for exhaustive reconfiguration. This is particularly valuable in **agile, DevOps, or CI/CD environments** where fast iteration and real-time QA feedback are essential.

Our framework integrates continuous human input into key points of the automated pipeline, resulting in a more **resilient, interpretable, and accurate testing system**.

2. Literature Review

Recent advances demonstrate the growing importance of HITL in machine learning and software systems. Chai and Li (2020) reviewed the taxonomy of HITL approaches in ML pipelines, emphasizing their role in refining model decisions. Zhang et al. (2020) outlined testing challenges in machine learning systems, especially the lack of robustness and oracle problems.

Gil et al. (2020) presented HITL architectures in cyber-physical systems where human feedback influenced system state monitoring. Meanwhile, Böhme et al. (2020) explored **HITL program repair**, where humans guided automated patches for faulty code.

Mantovani (2022) and Monarch (2021) discussed how HITL can be applied to **binary analysis and annotation**, enabling precise feature selection and defect isolation. Other notable contributions include:

- Kothari et al. (2014) on resolving complex software anomalies through human expertise.
- Kumar et al. (2024) and Retzlaff et al. (2024) on IML-enhanced model debugging and label correction.

These studies validate that **combining automation with human oversight** leads to smarter systems, especially in uncertain or ambiguous problem spaces like software testing.

3. Framework Overview

Our proposed system consists of four integrated modules:

1. Test Generator:

Uses static analysis and coverage criteria (e.g., path-based, branch-based) to synthesize initial test suites.

2. Interactive ML Classifier:

Applies supervised learning to classify test results as “bug”, “false positive”, or “inconclusive” based on labeled data, updating its model incrementally with human input.

3. Human Review Interface:

Displays ambiguous test outputs to developers for labeling and confirmation, feeding back corrections and confidence scores.

4. Feedback Loop Engine:

Adjusts test priorities and modifies classifiers dynamically, giving preference to test cases with lower confidence or higher developer disagreement.

4. Implementation and Experimental Setup

To validate the proposed Human-in-the-Loop (HITL) framework for automated software testing, we developed a full-stack implementation comprising test generation, interactive learning modules, and real-time human feedback integration. The following technologies were employed to construct and evaluate the system:

4.1 Toolchain and Technologies

- Python (TensorFlow, Scikit-learn)

Used for developing the Interactive Machine Learning (IML) components, including the test classification model and retraining pipeline. TensorFlow enabled real-time LSTM retraining, while Scikit-learn supported preprocessing and basic ensemble learners for anomaly detection.

- Pynguin and EvoSuite

These test generation tools were integrated to automatically synthesize unit and functional test cases. Pynguin was used primarily for Python-based projects, and EvoSuite was adopted for Java environments. Both tools provided path coverage strategies that served as initial input seeds to the IML loop.

- Web-Based Annotation Tool

A lightweight human interface was developed for presenting uncertain or ambiguous test results to domain experts. The interface allowed for quick annotation—labeling outputs as true bugs, false positives, or environment-specific variations—and dynamically fed this data back into the classifier for real-time learning.

- GitHub Actions for CI/CD Integration

To simulate deployment in real-world environments, the framework was integrated with GitHub Actions. This allowed automated test runs and retraining cycles to be triggered as part of each code commit, supporting real-time HITL evaluation in a continuous delivery pipeline.

5. Results and Analysis

Model	Precision	Recall	F1-score	HER ↓	TSGR ↑
Baseline AutoML	0.72	0.65	0.68	—	+10%
HITL-IML	0.87	0.81	0.84	0.18	+36%

Key observations:

- False positive rate dropped by **32%**
- Human reviewers were required in only **18% of test runs**
- HITL feedback led to improved anomaly classification and regression catch rates

6. Discussion

The results affirm that **selective human intervention** in testing workflows significantly improves performance without overburdening engineers. The adaptive learning process ensures that:

- Common false positives are quickly deprioritized
- Rare or complex bugs become **learning milestones** for the system
- Developers retain **control and interpretability** in critical decision points

Challenges include balancing **review fatigue** with automation speed and designing intuitive **user feedback interfaces**.

7. Future Directions

Several enhancements are proposed:

As the complexity and velocity of software development continue to rise, the demand for more intelligent, context-aware, and adaptive testing mechanisms is growing. Human-in-the-Loop (HITL) systems have demonstrated significant promise in addressing many of the current limitations in automated testing. Looking forward, several high-potential directions could further enhance the effectiveness and scope of these systems:

7.1 Incorporating Reinforcement Learning with Human Preferences

Future HITL testing frameworks can benefit from **reinforcement learning (RL)**, particularly by modeling human feedback as a reward signal. Unlike supervised learning, which relies on labeled data, RL can continuously improve policies for test generation and bug triaging by:

- Learning which test behaviors (coverage paths, assertion types) are preferred by human reviewers.
- Automatically adjusting exploration vs. exploitation strategies during test generation.
- Reducing dependence on exhaustive human annotation by learning generalized preference patterns.

This hybrid setup would enable testing agents to **evolve policies** over time that align more closely with **developer expectations and quality assurance standards**.

7.2 Expanding to UI and Mobile App Testing via Event Trace Models

While most automated testing focuses on backend systems or APIs, **UI testing**, especially in mobile apps, remains a manual and error-prone process. Integrating HITL with **event trace modeling** offers a promising avenue for enhancing front-end test coverage.

Key opportunities include:

- Capturing user interaction flows as **event graphs**.
- Using human feedback to label or correct unexpected UI transitions.

- Automating gesture-based testing across platforms and devices using learned user behavior models.

By training models to understand **intent and layout semantics**, testers can simulate complex usage scenarios and adapt tests to UI changes more gracefully.

7.3 Combining HITL with Fuzz Testing and Mutation-Based Exploration

Fuzz testing and mutation-based techniques are effective in identifying corner-case failures, but they often generate **uninformed or redundant inputs**. HITL frameworks can complement these methods by:

- Steering fuzzers toward unexplored but meaningful code regions based on human insight.
- Filtering out inputs that are semantically irrelevant or lead to known false positives.
- Prioritizing mutations that affect security-sensitive or frequently modified code blocks.

This integration enables **guided exploration** rather than blind input generation, improving both **depth and precision** in vulnerability detection.

Real-time HITL testing in DevOps pipelines remains a promising direction for **adaptive, scalable QA**.

8. Conclusion

This paper presents a **HITL-driven automated testing framework** that significantly improves test reliability and precision through human interaction. Our approach bridges the gap between brute-force automation and human intuition, making software testing **smarter, faster, and more trustworthy**.

References

- [1] Mantovani, A. (2022). An analysis of human-in-the-loop approaches for binary analysis automation.
- [2] Chai, C., & Li, G. (2020). Human-in-the-loop techniques in machine learning.
- [3] Liu, L., Perez-Concha, O., & Nguyen, A. (2023). Web-based HITL deep learning for deidentifying medical text.
- [4] Zhang, J. M., Harman, M., & Ma, L. (2020). Machine learning testing: Survey, landscapes and horizons.
- [5] Fontes, A. (2025). Context-infused automated software test generation.
- [6] Wu, X., Xiao, L., Sun, Y., Zhang, J., & He, L. (2022). A survey of human-in-the-loop for machine learning.
- [7] Maity, A., Banerjee, A., & Gupta, S. K. S. (2025). Detection of unknown-unknowns in HITL safety systems.
- [8] Kothari, S., Deepak, A., & Tamrawi, A. (2014). A human-in-the-loop approach for resolving complex software anomalies.

- [9] Singh, H. P. (2020). Non-intrusive monitoring of HMI states for detecting human-in-the-loop error precursors.
- [10] Böhme, M., Geethal, C., & Pham, V. T. (2020). Human-in-the-loop automatic program repair.
- [11] Gil, M., Albert, M., & Pelechano, V. (2020). Human-in-the-loop interactions in cyber-physical systems.
- [12] Mantovani, A. (2022). Reverse engineering automation using HITL machine learning.
- [13] Kumar, S., Datta, S., & Singh, V. (2024). Interactive ML solutions for oracle label accuracy.
- [14] Retzlaff, C. O., Das, S., & Wayllace, C. (2024). Human-in-the-loop reinforcement learning.
- [15] da Cunha, F. M. T. (2023). Keeping the human in the loop: IML platform engineering.
- [16] Khuat, T. T., Kedziora, D. J., & Gabrys, B. (2022). Roles and modes of HITL in AutoML systems.
- [17] Kumar, S., Datta, S., & Singh, V. (2024). Future directions in human-in-the-loop learning.
- [18] Monarch, R. M. (2021). Human-in-the-loop machine learning: Active learning and annotation.