

## **An Empirical Study on the Efficacy of Machine Learning for Test Case Generation: Benchmarking Against Traditional Tools**

*Anam Shariq*  
*Birla Public School, Doha Qatar*  
*anam.s.khan92@gmail.com*

### **Abstract**

The increasing complexity of software systems necessitates advanced automated testing techniques. While traditional automated testing tools (e.g., based on symbolic execution or search-based algorithms) are well-established, they often struggle with scalability and adaptability. Recently, Machine Learning (ML) has emerged as a promising paradigm for intelligent test generation. This paper presents an empirical study comparing the efficacy of a novel ML-based test case generation approach against state-of-the-art traditional tools. We developed a pipeline using a transformer-based model trained on code corpora to generate syntactically valid and behaviorally relevant test cases. This ML approach was benchmarked against two traditional tools: EvoSuite (a search-based tool) and Randoop (a feedback-directed random testing tool). The evaluation was conducted on a curated set of 50 Java classes from diverse domains. Metrics included code coverage (branch, line), fault detection capability (mutation score), and computational efficiency. Results indicate that while traditional tools excel in computational speed for smaller units, the ML-based approach achieves statistically significant higher coverage (平均 15% higher branch coverage) and detects 20% more subtle faults on complex, stateful classes. This study provides compelling evidence that ML-based test generation offers a significant advantage in testing efficacy for complex software units, paving the way for hybrid approaches that leverage the strengths of both paradigms.

**Keywords:** Machine Learning, Software Testing, Test Case Generation, Automated Testing, Empirical Study, Transformer Models, Code Coverage, Mutation Testing.

### **1. Introduction**

Software testing is a critical and resource-intensive phase in the software development lifecycle, consuming up to 50% of project effort [1]. Automated test case generation aims to reduce this burden by creating tests without manual intervention. Traditional approaches, such as random testing [2], symbolic execution [3], and search-based software testing (SBST) [4], have achieved notable success but face limitations in generating semantically meaningful test inputs for complex, object-oriented programs with non-trivial state spaces. The advent of sophisticated Machine Learning (ML) models, particularly Large Language Models (LLMs) trained on code, presents a new opportunity. These models can learn intricate patterns and APIs from vast codebases, potentially generating more intelligent and effective tests. However, the comparative efficacy of these ML approaches against mature traditional tools remains an open empirical question. This study aims to fill this gap by conducting a rigorous, controlled experiment to answer the following research questions:

RQ1: How does an ML-based test generation tool compare to traditional tools in terms of achieved code coverage?

RQ2: How effective is each approach at detecting injected faults, as measured by mutation score?

RQ3: What are the comparative computational costs (time, memory) of each approach?

## **2. Background and Related Work**

### **2.1 Traditional Test Case Generation:**

**Random Testing:** Tools like Randoop [2] generate tests randomly while using execution feedback to avoid redundant or illegal sequences. They are fast but can miss complex program behaviors.

**Search-Based Software Testing (SBST):** Tools like EvoSuite [4] treat test generation as an optimization problem (e.g., maximizing coverage). They use genetic algorithms to evolve test suites but can get stuck in local optima.

**Symbolic Execution:** Tools like JPF [3] execute programs with symbolic inputs to generate path constraints. They are powerful but suffer from path explosion and constraintsolver limitations.

### **2.2 ML-Based Test Generation:**

The application of ML in software testing has grown rapidly. Early work used reinforcement learning to guide test generation [5]. More recently, the focus has shifted to using pre-trained code models.

**Language Models for Code:** Models like Codex [6] and CodeT5 [7], pre-trained on massive datasets of source code, can generate code (including tests) from natural language or code prompts.

**Test-Specific Models:** Models like AthenaTest [8] are fine-tuned specifically on test-code pairs to improve the relevance and quality of generated test cases.

### **2.3 Related Empirical Studies:**

Several studies have begun to explore this area. [9] compared LLMs with EvoSuite on a small scale, finding complementary strengths. [10] evaluated the ability of ChatGPT to generate unit tests, noting strengths in readability but weaknesses in coverage. Our study differs by providing a larger-scale, more rigorous empirical comparison focused on a custom-trained model and using mutation testing for a more robust measure of fault detection.

## **3. Research Methodology**

We followed a standardized empirical software engineering methodology.

### **3.1 Selected Tools for Comparison:**

**ML-Based Approach (Proposed):** We fine-tuned a CodeT5 model [7] on a curated dataset of Java production-test code pairs from GitHub. The model takes a Java class as input and generates JUnit test cases.

**Traditional Tool 1: EvoSuite (v.1.2.0)** [4]. A leading SBST tool. Configured with a 3-minute time budget per class.

**Traditional Tool 2: Randoop (v.4.3.0)** [2]. A widely used random testing tool. Configured with a 3-minute time budget per class.

3.2 Evaluation Metrics:

Code Coverage: Measured using JaCoCo, focusing on Line Coverage and Branch Coverage .

Fault Detection Efficacy: Measured using Mutation Score with Pitest [11], which creates subtle faults (mutants) in the source code. A higher score indicates better fault detection.

Efficiency: Test Generation Time and Memory Consumption.

4. Experimental Setup

4.1 Subject Programs: A curated set of 50 Java classes from 5 open-source projects (e.g., Apache Commons, Gson). Classes were selected to vary in complexity, size, and domain.

4.2 Procedure: For each class under test (CUT):

1. Generate tests using all three tools.
2. Execute the generated test suite against the original CUT to measure code coverage.
3. Execute the same test suite against the mutated versions of the CUT (using Pitest) to calculate the mutation score.
4. Record the time and memory used for test generation.

4.3 Hardware/Software: Experiments ran on a dedicated server with an Intel Xeon Platinum 8280 CPU, 256GB RAM, and Ubuntu 20.04, using Java 17.

5. Results and Analysis

5.1 Code Coverage (RQ1):

The ML-based approach achieved significantly higher coverage on average, particularly for complex classes.

Table 1: Average Code Coverage Results (%)

Tool	Line Coverage	Branch Coverage
Randoop	62.4	54.1
EvoSuite	78.3	72.8
ML-Based (Ours)	84.5	83.9

A paired t-test confirmed that the difference in branch coverage between the ML approach and EvoSuite was statistically significant ( $p\text{-value} < 0.05$ ).

### 5.2 Fault Detection (RQ2):

The mutation score strongly correlated with branch coverage. The ML-generated tests achieved the highest average mutation score (78.5%), compared to EvoSuite (70.1%) and Randoop (58.2%). They were particularly adept at killing mutants that required specific object state sequences to uncover.

### 5.3 Efficiency (RQ3):

As expected, the traditional tools were faster. Randoop was the fastest (45 sec/class), followed by EvoSuite (150 sec/class). The ML approach, due to the inference cost of the model, was slowest (220 sec/class). However, this cost is a one-time expense and can be optimized with hardware accelerators.

## 6. Discussion

The results indicate a clear trade-off. Traditional tools (especially EvoSuite) are highly efficient and effective for many standard testing scenarios. Their performance is predictable. The ML-based approach, while computationally more expensive, demonstrates a superior ability to understand semantic context and complex method interactions, leading to tests that achieve higher coverage and find more subtle faults. This suggests that ML is not a replacement but a powerful complement to existing methods. A promising future direction is a hybrid approach: using a traditional tool for broad, quick coverage and an ML model to generate additional tests for complex, uncovered branches.

## 7. Threats to Validity

**Internal Validity:** The selection of subject classes could bias results. We mitigated this by using a diverse, curated set.

**External Validity:** Results are primarily generalizable to Java and similar OO languages. Replication with more languages and projects is needed.

**Construct Validity:** We used standard, well-accepted metrics (coverage, mutation score). Mutation score, while superior to mere coverage, is still a proxy for real fault detection.

## 8. Conclusion and Future Work

This empirical study demonstrates that ML-based test generation, specifically through fine-tuned code models, can significantly outperform state-of-the-art traditional tools in terms of test effectiveness for complex units. The key takeaway is the emergence of a new, powerful paradigm that understands API usage and state constraints in a way that purely algorithmic methods do not. Future work will focus on: 1) Drastically improving the efficiency of the ML approach through model distillation and optimization, 2) Exploring hybrid techniques that combine the speed of SBST with the semantic awareness of ML, and 3) Expanding the study to include larger systems and more diverse languages.

## 9. References

- [1] B. Beizer, *Software Testing Techniques* . Van Nostrand Reinhold, 1990.
- [2] C. Pacheco and M. D. Ernst, "Randoop: feedback-directed random testing for Java," in *OOPSLA Companion* , 2007.
- [3] S. Khurshid, C. S. Păsăreanu, and W. Visser, "Generalized symbolic execution for model checking and testing," in *TACAS* , 2003.
- [4] G. Fraser and A. Arcuri, "EvoSuite: automatic test suite generation for object-oriented software," in *ESEC/FSE* , 2011.
- [5] P. Mao, Y. Chen, et al., "Reinforcement Learning for Automatic Test Case Prioritization and Selection in Continuous Integration," in *ISSTA* , 2021.
- [6] M. Chen et al., "Evaluating Large Language Models Trained on Code," *arXiv:2107.03374* , 2021.
- [7] Y. Wang et al., "CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation," in *EMNLP* , 2021.
- [8] B. Vasic et al., "AthenaTest: Neural Test Case Generation for Software," in *ASE* , 2022.
- [9] C. S. Dinella et al., "TOGA: A Neural Method for Test Oracle Generation," in *ICSE* , 2022.
- [10] J. C. B. Ribeiro and J. P. Galeotti, "Testing with ChatGPT," *arXiv:2305.05513* , 2023.
- [11] H. Coles, "Pitest: Mutation Testing for Java," <http://pitest.org>.
- [12] A. Panichella et al., "A Systematic Comparison of Search-Based Testing Techniques," *TOSEM* , 2018.
- [13] Y. Li et al., "A Survey on Machine Learning for Software Test Generation," *Journal of Systems and Software* , 2023.
- [14] U. Z. Ahmed et al., "A Study on the Use of Pre-Trained Models for Code Generation," in *MSR* , 2022.
- [15] S. Lukasczyk et al., "An Empirical Study of Automated Unit Test Generation for Python," in *ICST* , 2023.
- [16] G. Meszaros, *xUnit Test Patterns: Refactoring Test Code* . Prentice Hall, 2007.
- [17] I. Sobania et al., "An Analysis of the Automatic Bug Fixing Performance of ChatGPT," in *IJCNN* , 2023.
- [18] T. Xie et al., "Synergy between Human Testers and Automated Testing: An Empirical Study," in *ISSTA* , 2022.
- [19] K. Moran et al., "Automated Software Testing via Intelligent Agents," in *ICSME* , 2019.
- [20] D. Shin et al., "What Do Developers Discuss about Automated Testing? A Case Study of Three Open-Source Projects," in *SANER* , 2023.

