



# **A Lightweight Model for Balancing Efficiency and Precision in PEFT-Optimized Java Unit Test Generation**

By

**Sintayehu Zekarias Esubalew**

Submitted to the School of Information Technology and Engineering

In partial fulfillment of the requirements for the degree of

Master of Science in Artificial Intelligence

Supervised by: Dr. Beakal Gizachew Assefa.

School of Information Technolgy and Engineering

Collage of Technology and Built Environment

Addis Ababa University

Addis Ababa, Ethiopia

June , 2025

## APPROVAL

This is to certify that this thesis titled "**A Lightweight Model for Balancing Efficiency and Precision in PEFT-Optimized Java Unit Test Generation**" is prepared by Sintayehu Zekarias Esubalew and submitted in partial fulfillment of the thesis-option requirements for the Degree of Master of Science in Artificial Intelligence at School of Information Technology and Engineering, Addis Ababa Institute of Technology.

**Dr. Beakal Gizachew**

\_\_\_\_\_  
Advisor

\_\_\_\_\_  
Signature

\_\_\_\_\_  
Date

\_\_\_\_\_  
Co.Advisor

\_\_\_\_\_  
Signature

\_\_\_\_\_  
Date

**Dr. Surafel Lemma**

\_\_\_\_\_  
External Examiner

\_\_\_\_\_  
Signature

\_\_\_\_\_  
Date

**Dr. Adane Letta**

\_\_\_\_\_  
Internal Examiner

\_\_\_\_\_  
Signature

\_\_\_\_\_  
Date

**Dr. Henock Mulugeta**

\_\_\_\_\_  
Chair Person

\_\_\_\_\_  
Signature

\_\_\_\_\_  
Date

## ACKNOWLEDGMENT

I am grateful to Almighty God for His love, grace, and strength that have supported me throughout my life.

I express my sincere appreciation to my advisor, Dr. Beakal Gizachew, for his guidance and support throughout my M.Sc. thesis research. His expertise, feedback, and encouragement were essential in shaping my work and refining my ideas.

I am thankful to the Ethiopian Artificial Intelligence Institute for providing this opportunity and supporting my academic and research endeavors. In particular, I extend my thanks to Dr. Taye Girma and Dr. Worku Gachena for their support, mentorship, and constructive input, which contributed to the completion of this thesis.

My thanks go to the thesis examiners: Dr. Henock Mulugeta, Dr. Adane Mamuye, Dr. Fantahun Bogale, Dr. Natnael Argaw, and Dr. Surafel Lemma. Their time, feedback, and suggestions improved the quality of this thesis.

I am grateful to my family for their love, patience, and support throughout this journey, including my mother Zenebech Yesuf, my wife Yodit G/Silase, my son Eyoas Sintayehu, my sister Sintayehu Teka, and my brother Endashaw Megeras.

I also appreciate my friends and colleagues—Takele Lakew, Keno Lemesa, Firsew Feyseo, Andualem Tilahun, Lelisa Benti, Hana Demma, Bereketiab Bantewosen, Robel Morka, and Masreshaye Worku—for their suggestions and encouragement, which supported various aspects of this research. Special thanks are due to Dr. Taye Girma for providing critical computing resources for my experiments.

I am thankful to my classmates for their discussions and collaborative spirit, which enriched my academic experience. In particular, I thank Tadele Melesse, Ashenafi Kifle, Fikir Awoke, and Azmeraw Bekele for their companionship and shared insights.

Finally, I acknowledge the School of Information Technology and Engineering (SITE) for its support and dedication. I am grateful to all my colleagues and the administration for their role in making this academic journey possible.

## Contents

<b>Abbreviations</b>	<b>vi</b>
<b>List of Figures</b>	<b>ix</b>
<b>List of Tables</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Motivation of the Study . . . . .	2
1.3 Statement of the Problem . . . . .	3
1.4 Research Questions . . . . .	4
1.5 Objectives . . . . .	4
1.5.1 General Objective . . . . .	4
1.5.2 Specific Objectives . . . . .	4
1.6 Contribution . . . . .	5
1.7 Scope and Delimitation . . . . .	6
1.8 Limitations . . . . .	6
1.9 Structure of the Document . . . . .	7
<b>2 Literature Review</b>	<b>9</b>
2.1 Foundation . . . . .	9
2.2 Software Testing . . . . .	9
2.2.1 Test Cases . . . . .	9
2.2.2 Testing Types . . . . .	10
2.2.3 Testing Techniques . . . . .	10
2.2.4 Testing Phases or Levels . . . . .	11
2.3 Taxonomy of AI-Based Unit Test Case Generations . . . . .	12
2.4 Traditional Machine Learning Approaches . . . . .	14
2.4.1 Genetic Algorithms (GAs) . . . . .	14
2.5 Deep Learning Approaches . . . . .	15
2.5.1 Recurrent Neural Networks (RNNs) and Long Short-Term Memory Networks (LSTMs) . . . . .	15
2.5.2 Convolutional Neural Networks (CNNs) . . . . .	16
2.5.3 Graph Neural Networks (GNNs) . . . . .	17

2.6	Properties of AI-Based Unit Test Case Generation . . . . .	18
2.6.1	Essential Properties and Their Significance . . . . .	18
2.6.2	Analysis for Traditional Machine Learning and Deep Learning Approaches for Test Generation . . . . .	20
2.7	Transformer-Based Approaches . . . . .	21
2.7.1	Transformer Components . . . . .	22
2.7.2	Disadvantages of Transformers . . . . .	24
2.7.3	Transformer Variants Used in Unit Test Case Generation . . . . .	24
2.8	Parameter-Efficient Fine-Tuning Methods for Transformer . . . . .	26
2.8.1	Major Approaches of PEFT . . . . .	27
2.8.2	Rationale for Selecting PEFT Methods . . . . .	29
2.9	Related Works . . . . .	31
2.9.1	Evaluation of Key Properties . . . . .	32
2.9.2	Baseline Model Comparison . . . . .	35
2.9.3	Summary . . . . .	37
<b>3</b>	<b>Methodology</b>	<b>39</b>
3.1	Research Methodology . . . . .	39
3.2	An Optimized Unit Test Case Generation Methodology . . . . .	41
3.2.1	Phase 1: Assertion Pre-Training . . . . .	44
3.2.2	Phase 2: PEFT Optimization . . . . .	45
3.2.3	Phase 3: Test Generation Fine-Tuning . . . . .	49
3.2.4	Phase 4: Test Validation . . . . .	50
<b>4</b>	<b>Experiments</b>	<b>52</b>
4.1	Introduction . . . . .	52
4.2	Experimental Design and Setup . . . . .	53
4.2.1	Datasets . . . . .	56
4.2.2	Hyperparameter Tuning and Hardware Configuration . . . . .	58
4.3	Phase 1: Assertion Pre-Training . . . . .	60
4.3.1	Phase 1: Assertion Pre-training . . . . .	60
4.3.2	Training Phase Analysis . . . . .	61
4.3.3	Evaluation Phase Analysis . . . . .	61
4.3.4	Comparative Analysis and Implications . . . . .	62
4.4	LoRA . . . . .	65
4.4.1	Phase 2: PEFT Optimization using LoRA . . . . .	65
4.4.2	Phase 3: Test Generation Fine-Tuning (LoRA) . . . . .	67
4.4.3	Phase 4: Test Validation (LoRA) . . . . .	68
4.4.4	Discussion and Results (LoRA) . . . . .	72

4.5	QLoRA	76
4.5.1	Phase 2: PEFT Optimization (QLoRA)	76
4.5.2	Phase 3: Test Generation Fine-Tuning (QLoRA)	78
4.5.3	Phase 4: Test Validation (QLoRA)	79
4.5.4	Discussion and Results (QLoRA)	83
4.6	Adapters	87
4.6.1	Phase 2: PEFT Optimization (Adapters)	87
4.6.2	Phase 3: Test Generation Fine-Tuning (Adapters)	89
4.6.3	Phase 4: Test Validation (Adapters)	90
4.6.4	Discussion and Results (Adapters)	94
4.7	Overall Discussion and Results	98
4.7.1	PEFT Configuration Analysis (RQ1)	98
4.7.2	Comparison with Full Fine-Tuning and State-of-the-Art (RQ2)	99
4.7.3	Dataset Performance	101
4.7.4	Transfer Learning and Robustness	103
4.7.5	Key Findings	104
4.7.6	Real-World Implications	105
<b>5</b>	<b>Conclusion</b>	<b>106</b>
5.1	Recommendations	107
5.2	Future Work	108
	<b>References</b>	<b>110</b>

## Abbreviations

Symbol	Meaning
AI	Artificial Intelligence
AK-PLBART	Assertion Knowledge PLBART
AOAK-PLBART	Adapter-based Optimized Assertion Knowledge PLBART
API	Application Programming Interface
AST	Abstract Syntax Tree
ATT	Average Time per Test
A3Test	Automated Assertion-Aware Test generation
CA	Correct Assertions
CFG	Control Flow Graph
CI/CD	Continuous Integration / Continuous Deployment
CNN	Convolutional Neural Network
CTC	Correct Test Cases
Defects4J	Defects for Java (benchmark dataset)
DL	Deep Learning
DT	Decision Tree
EvoSuite	Evolutionary Test Suite Generator
FMC	Focal Method Coverage
FM	Focal Method
FT	Fine-Tuning
GA	Genetic Algorithm
GNN	Graph Neural Network
GPU	Graphics Processing Unit
JUnit	Java Unit Testing Framework
LC	Line Coverage
LM	Language Modeling
LOAK-PLBART	LoRA-based Optimized Assertion Knowledge PLBART
LoRA	Low-Rank Adaptation
LSTM	Long Short-Term Memory (Network)
ML	Machine Learning
MLM	Masked Language Model
MS	Mutation Score
NLP	Natural Language Processing
OOP	Object-Oriented Programming

PDG	Program Dependency Graph
PEFT	Parameter-Efficient Fine-Tuning
PLBART	Programming Language Bidirectional Auto-Regressive Transformer
QLoRA	Quantized Low-Rank Adaptation
QOAK-PLBART	QLoRA-based Optimized Assertion Knowledge PLBART
RNN	Recurrent Neural Network
SE	Software Engineering
SOTA	State-of-the-Art
SVM	Support Vector Machine
TPU	Tensor Processing Unit
UTAOAK-PLBART	Unit-Test-Adapters-based Optimized Assertion Knowledge PLBART
UTLOAK-PLBART	Unit-Test-LoRA-based Optimized Assertion Knowledge PLBART
UTQOAK-PLBART	Unit-Test-QLoRA-based Optimized Assertion Knowledge PLBART

---



## LIST OF FIGURES

2.1	Hierarchy of Software Testing Levels. Adapted from [1].. . . . .	12
2.2	Taxonomy of AI-Based Unit Test Case Generation. . . . .	14
2.3	Properties of AI-Based Unit Test Case Generation. . . . .	19
2.4	Overall Structure of the Transformer Model. Adapted from [2].. . . .	22
2.5	Conceptual Illustration of PEFT Techniques including Full Fine-tuning, LoRA, and QLoRA. Adapted from Dettmers et al. [3]. . . . .	30
3.1	Design Science Research Process-Based, Problem-Centered Research Methodology for AI-Powered Unit Test Generation. . . . .	39
3.2	Optimized approach for Unit Test Case Generation for Java using Transformers, builds upon the approach proposed by [4]. . . . .	43
4.1	Assertion Pre-Training Metrics with a training set of 960,000 examples and validation set of 2 43,928 examples. . . . .	63
4.2	LoRA Fine-Tuning Metrics: (a) shows the training loss decreasing from approximately 0.45 to less than 0.15 with a cyclical learning rate from $2.9 \times 10^{-5}$ to $2.3 \times 10^{-5}$ over 43,928 steps. (b) illustrates the validation loss increasing from 0.284 to 0.292 over 4 epochs, indicating potential overfitting, with training loss decreasing from 0.51 to 0.43 over the same period. . . . .	69
4.3	Comparison of A3Test (LoRA) vs. A3Test (Full) on Defects4J: (a) percentage-based metrics (CTC, FMC, LC, CA, MS), (b) average test time, and (c) peak memory usage. . . . .	75
4.4	QLoRA Fine-Tuning Metrics: (a) shows the training loss decreasing from approximately 0.45 to below 0.15 with a cyclical learning rate varying from $2.9 \times 10^{-5}$ to $2.3 \times 10^{-5}$ over 43,928 steps. (b) illustrates the validation loss increasing from 0.284 to 0.292 over 4 epochs, while training loss decreases from 0.51 to 0.43, indicating potential overfitting despite the 4-bit quantized efficiency of QLoRA. . . . .	80
4.5	Comparison of A3Test (QLoRA) vs. A3Test (Full) on Defects4J: (a) percentage-based metrics (CTC, FMC, LC, CA, MS), (b) average test time, and (c) peak memory usage. . . . .	86

4.6	Adapter Fine-Tuning Metrics: (a) shows the training loss decreasing from approximately 0.45 to below 0.15 with a cyclical learning rate from $2.9 \times 10^{-5}$ to $2.3 \times 10^{-5}$ over 43,928 steps. (b) illustrates the validation loss increasing from 0.284 to 0.292 over 4 epochs, indicating potential overfitting.	91
4.7	Comparison of A3Test (Adapters) vs. A3Test (Full) on Defects4J: (a) percentage-based metrics (CTC, FMC, LC, CA, MS), (b) average test time, and (c) peak memory usage. . . . .	97
4.8	Grouped bar chart comparing CTC, FMC, MS, ATT, and Mem. for Full FT, LoRA, QLoRA, and Adapters on Defects4J. . . . .	100
4.9	Grouped bar chart comparing CTC, FMC, LC, and CA for LoRA, QLoRA, and Adapters on Defects4J. . . . .	102

## LIST OF TABLES

2.1	Comparative Effectiveness of AI Approaches for Unit Test Generation . . .	21
2.2	Transformer Models Comparison for Unit Test Case Generation . . . . .	26
2.3	Comparison of PEFT Methods . . . . .	30
2.4	Effectiveness of AI for Unit Test Case Generation . . . . .	32
2.5	Usability of AI for Unit Test Case Generation . . . . .	33
2.6	Maintainability of AI for Unit Test Case Generation . . . . .	34
2.7	Advanced Features for AI Unit Test Case Generation . . . . .	34
2.8	Comparison of Generated Test Cases Based on Focal Method . . . . .	35
2.9	Focal Method Coverage Analysis . . . . .	36
2.10	Coverage Metrics Comparison . . . . .	36
2.11	Assertions Comparison . . . . .	36
2.12	Computational Time and Cost Overview . . . . .	37
2.13	Comparative Performance of AI-Based Unit Test Generation Approaches	37
4.1	Hyperparameter Tuning Results . . . . .	59
4.2	Evaluation Metrics per Project on Defects4J v2.5.0 . . . . .	72
4.3	LoRA vs. Full Fine-Tuning on Defects4J . . . . .	74
4.4	Evaluation Metrics per Project on Defects4J v2.5.0 . . . . .	83
4.5	QLoRA vs. Full Fine-Tuning on Defects4J . . . . .	85
4.6	Evaluation Metrics per Project on Defects4J v2.5.0 . . . . .	94
4.7	Adapters vs. Full Fine-Tuning on Defects4J . . . . .	96
4.8	PEFT vs. Full Fine-Tuning on Defects4J . . . . .	99
4.9	Model Performance Across Datasets . . . . .	101
4.10	Transfer Learning Results on Defects4J . . . . .	103

# Chapter 1

## Introduction

### 1.1 Background

In the era of rapid technological advancement, software systems have become indispensable across a wide range of industries and applications [5]. As the reliance on software grows, ensuring its reliability and correctness is of paramount importance. However, despite the adoption of modern software development practices, a substantial number of software projects still suffer from poor quality, delays, and cost overruns [6]. A significant contributor to these challenges is the late detection and correction of software defects, which become exponentially more expensive to fix as the development lifecycle progresses [6, 7].

Software testing plays a critical role in ensuring software quality by detecting defects early in the development cycle. Among various testing techniques, unit testing is fundamental as it validates the behavior of individual components in isolation. A key activity within unit testing is the generation of test cases—inputs and expected outputs used to evaluate software correctness [1, 8]. However, the process of generating high-quality test cases is known to be labor-intensive and costly, accounting for approximately 50% of total software development costs [7].

To address these challenges, automated unit test generation has emerged as a promising approach to improve test coverage and reduce testing costs. Recent advances in artificial intelligence (AI), particularly machine learning (ML) and deep learning (DL), have shown great potential in automating the generation of syntactically correct and semantically meaningful test cases [9, 10]. Transformer-based models, in particular, have demonstrated the ability to learn complex code structures and generate relevant test cases

with minimal human intervention [4, 11].

Nevertheless, the high computational cost and resource demands associated with training large language models pose practical challenges for many organizations [4, 3]. Furthermore, while AI-generated tests often pass syntax checks and compile correctly, they may fail to detect subtle bugs due to limited assertion knowledge or lack of domain-specific reasoning [12, 13].

Therefore, the integration of AI in test generation must be optimized not only for accuracy and coverage but also for cost-efficiency and practical applicability. This has led to increased interest in parameter-efficient fine-tuning (PEFT) methods, such as LoRA and Adapters, which reduce the number of trainable parameters and lower resource usage while maintaining competitive performance [14, 15]. These methods offer a scalable solution for adopting AI-driven test generation in real-world development environments, particularly where hardware resources are constrained.

## 1.2 Motivation of the Study

The increasing complexity of modern software systems underscores the need for efficient and reliable testing practices, particularly for Java, a widely used language in enterprise, web, and mobile applications. Unit testing, which verifies individual components in isolation, is critical for ensuring software quality and detecting defects [1]. However, manual test creation is time-consuming, and traditional automated methods often lack the semantic accuracy needed to identify subtle bugs.

Recent advances in artificial intelligence (AI), particularly transformer-based models like PLBART, show promise in automating Java unit test generation by producing syntactically correct and readable test cases [16]. Yet, these models require significant computational resources, limiting their accessibility for teams with constrained hardware [4]. Moreover, AI-generated tests often struggle with semantic accuracy, failing to incorporate domain-specific knowledge or robust assertions to validate software behavior

effectively [17].

Parameter-Efficient Fine-Tuning (PEFT) techniques, such as LoRA [14], QLoRA [3] and Adapters [15], offer a solution by reducing the computational cost of fine-tuning while maintaining high performance. These methods enable scalable test generation on resource-limited systems. Additionally, enhancing assertion knowledge with domain-specific insights can improve the quality and bug-detection capabilities of generated tests.

This study is motivated by the need to develop an AI-driven unit test generation framework that balances computational efficiency with test effectiveness. By leveraging PEFT techniques on the PLBART model and integrating advanced assertion strategies, this research aims to make automated test generation more accessible and reliable, ultimately contributing to higher-quality software systems.

### 1.3 Statement of the Problem

This study addresses the high computational cost and resource demands of fine-tuning PLBART, a transformer-based model with 139.9 million parameters proposed by Ahmad et al. (2021) [18], for automated Java unit test generation. These requirements limit its practical adoption in resource-constrained software development environments. Full fine-tuning of PLBART demands substantial computational resources—such as 32GB GPUs—resulting in high memory usage (e.g., over 23.5 GB) and prolonged training times. Additionally, the absence of a pre-trained assertion model from the baseline necessitates additional pre-training for assertion generation, further increasing computational overhead and complexity. Consequently, the resource-intensive nature of PLBART-based test generation renders it largely inaccessible to organizations with limited infrastructure.

Parameter-efficient fine-tuning (PEFT) techniques, including Low-Rank Adaptation (LoRA) [14], Quantized LoRA (QLoRA) [3], and Adapters [15], offer promising solutions by significantly reducing the number of trainable parameters. These methods help lower memory usage and shorten training times. However, their application for optimizing

PLBART in the context of Java unit test generation remains underexplored. Moreover, the best strategies for achieving an optimal trade-off between computational efficiency and test quality—measured in terms of correctness, assertion accuracy, and code coverage—are not yet well established [4, 19, 20].

This research investigates the integration of PEFT techniques into PLBART to enable efficient Java unit test generation while minimizing computational costs. The evaluation leverages benchmark datasets such as Defects4J [21] and Methods2Test [22] to assess memory efficiency, generation time, and the effectiveness of the generated tests. By mitigating the resource demands of full fine-tuning, this study aims to improve the scalability, accessibility, and real-world applicability of AI-driven unit test generation for Java.

## 1.4 Research Questions

- **RQ1:** What is the optimal configuration of PEFT components for unit test generation in terms of correctness and efficiency?
- **RQ2:** How do the PEFT-optimized models (LoRA, QLoRA, Adapters) perform compared to the fully fine-tuned A3Test [4] baseline under consistent experimental conditions in terms of test case generation quality?

## 1.5 Objectives

### 1.5.1 General Objective

To design and implement an algorithm that improves AI-based Java unit test case generation using Parameter-Efficient Fine-Tuning (PEFT) techniques.

### 1.5.2 Specific Objectives

To achieve the overarching research goal, this study defines the following specific objectives:

- Apply LoRA, QLoRA, and Adapters to fine-tune PLBART for test generation.

- Find optimal configuration of PEFT components for unit test generation.
- Reduce memory usage and generation time compared to full fine-tuning [4, 19, 20].
- Improve assertion accuracy and code coverage on standard benchmarks.
- Evaluate generalizability of the model with Defects4J dataset.

## 1.6 Contribution

This research advances AI-driven unit test generation by applying and evaluating parameter-efficient fine-tuning (PEFT) techniques on the PLBART model for Java unit test generation. The key contributions are:

**Applying PEFT for Java Unit Test Generation.** This study investigates the application of three PEFT techniques—Low-Rank Adaptation (LoRA), Quantized LoRA (QLoRA), and Adapters—on the PLBART model. By adapting these fine-tuning strategies to a code-focused transformer architecture, the research offers an approach to improve fine-tuning efficiency of large models in software engineering tasks.

**Evaluating Efficiency and Trade-offs of PEFT Approaches.** Through comparative analysis (see Table 4.8), we demonstrate that PEFT methods reduce resource consumption relative to full fine-tuning. PEFT techniques yield a 19–24% reduction in average test generation time (e.g., 1.52–1.60s vs. 1.98s) and a 15–36% decrease in memory usage (e.g., 15.0–20.1 GB vs. 23.5 GB). These efficiency gains come with a 5–10% decrease in correctly generated test cases (e.g., 36.20–38.12% vs. 40.05%), indicating a resource-quality trade-off appropriate for constrained environments.

**Maintaining Comparable Performance with Reduced Resources.** Compared to the fully fine-tuned A3Test baseline, PEFT-optimized PLBART models with LoRA and QLoRA achieve Correct Test Cases (CTC) of 38.12% and 37.90% (vs. 40.05%), and Focal

Method Coverage (FMC) of 45.1% and 45.0% (vs. 47.0%). They also reduce average test generation time by 23–24% (1.52–1.50s vs. 1.98s) and peak memory usage by up to 36% (15.0–19.5 GB vs. 23.5 GB), demonstrating the ability to produce comparable test quality with fewer resources.

**Demonstrating Transferability Across Datasets.** This work evaluates the transfer learning capabilities of PEFT methods across multiple datasets—Atlas, Methods2Test, and Defects4J. Results indicate that PEFT-optimized models maintain consistent performance across these benchmarks, suggesting applicability in diverse software testing scenarios.

## 1.7 Scope and Delimitation

This study focuses on improving Java unit test generation using AI techniques, specifically parameter-efficient fine-tuning methods (LoRA, QLoRA, Adapters) applied to transformer-based models. Evaluation is conducted on benchmark datasets including Defects4J, Atlas, and Methods2Test, with performance measured via assertion accuracy, code coverage (line and branch), and generation efficiency.

Limitations include focusing solely on Java unit testing without extensive coverage of other programming languages, testing levels, or alternative AI approaches (e.g., evolutionary algorithms). The study does not integrate with CI/CD pipelines or specific project implementations, relying instead on established benchmarks. We assume access to sufficient computational resources, though PEFT methods aim to reduce these requirements.

## 1.8 Limitations

Several limitations should be considered when interpreting the findings:

- **Dataset Specificity:** Experiments primarily used Methods2Test and Defects4J

datasets. Although standard, PEFT performance may vary on other datasets or programming languages.

- **Hyperparameter Optimization:** Hyperparameter tuning was conducted within a limited search space constrained by available computational resources. More extensive tuning could potentially enhance PEFT performance.
- **Model Architecture:** The study centers on the PLBART-base model. PEFT effectiveness may differ when applied to other pre-trained language models of varying architectures or sizes.
- **Computational Resources:** Experiments were run on hardware with  $2 \times$  NVIDIA Tesla V100-SXM2 32GB GPUs. Results may not fully generalize to other hardware configurations. QLoRA's memory savings are especially relevant for resource-constrained environments, but actual performance could vary.
- **Qualitative Analysis:** Although readability was evaluated in comparison with state-of-the-art methods, a more in-depth qualitative assessment of generated tests could provide additional practical insights.

## 1.9 Structure of the Document

The document systematically presents the research on "Software Unit Test Case Generation through the Integration of AI" across five chapters and a dedicated references section.

**Chapter One: Introduction** lays the groundwork, covering the study's background, motivation, problem statement, research questions, and objectives. It also outlines the study's contribution, scope, delimitations, and limitations.

**Chapter Two: Literature Review** comprehensively examines existing research on software unit test case generation and AI integration. It establishes foundational concepts, details relevant software testing aspects, and presents a taxonomy of AI-based gen-

eration approaches, including traditional machine learning, deep learning, and transformer-based methods. This chapter also discusses essential properties of AI-based test generation, analyzes related works, and compares key baseline models.

**Chapter Three: Methodology** details the research design, outlining the optimized unit test case generation methodology. This includes distinct phases for assertion pre-training, parameter-efficient fine-tuning (PEFT) optimization, test generation fine-tuning, and test validation.

**Chapter Four: Experiments** presents the experimental design and setup, including datasets and hardware configuration. It then details the results and analysis for each phase of the methodology, specifically exploring different PEFT methods (LoRA, QLoRA, and Adapters). The chapter concludes with a general analysis, answers to research questions, key findings, and real-world implications.

**Chapter Five: Conclusion** summarizes the research's main findings and implications, offering recommendations and outlining directions for future work.

Finally, the **References** section provides a comprehensive list of all cited sources.

# Chapter 2

## Literature Review

### 2.1 Foundation

This section introduces the core concepts underpinning AI-driven unit test case generation. We first outline the principles of software testing, including test case design, techniques, and levels. Next, we survey AI-based approaches specifically categorized by their application in generating unit test cases, including traditional machine learning (e.g., genetic algorithms), deep learning (e.g., GNNs), and their evaluation metrics. Finally, we focus on transformer-based methods, detailing their architecture, variants, and parameter-efficient fine-tuning techniques. Together, these foundations support our analysis of state-of-the-art test generation tools and their limitations.

### 2.2 Software Testing

Software testing is a key component of software development, as it informs stakeholders about the product's quality. It is estimated that testing often takes up between 30% and 40% of a software development organization's total project work and can account for more than 50% of the project costs [23, 24]. The System Under Test (SUT) must ensure a high-quality software product. Failure occurs when the external behavior of the SUT deviates from what is predicted by its requirements or other specifications [1].

#### 2.2.1 Test Cases

Test cases are central to the process of software testing. Each test case defines specific conditions under which the SUT is evaluated to identify defects or verify functionality. A well-designed test case typically originates from the software's functional requirements

or design documentation and includes several components. These include the initial state or preconditions, the sequence of operations or steps that need to be executed, and the expected outcome that will serve as a benchmark to evaluate the success or failure of the test. By executing test cases and comparing actual versus expected outcomes, developers and testers can identify discrepancies that indicate defects in the software logic or implementation [25].

### **2.2.2 Testing Types**

Testing methodologies are typically classified into two main categories based on execution approach: manual and automated testing. In manual testing, human testers carry out test cases without the aid of software tools. This approach simulates real-user interactions and is useful in exploratory or ad-hoc testing scenarios. However, manual testing can be time-consuming and error-prone when dealing with large-scale systems.

Automated testing, by contrast, leverages software tools to execute predefined test scripts. While it requires an initial investment in scripting and infrastructure, automated testing offers consistency, speed, and the ability to re-run tests efficiently. It is particularly beneficial in regression testing scenarios, where repetitive validation of previously tested components is necessary following modifications or enhancements.

### **2.2.3 Testing Techniques**

Three broad testing techniques dominate the software testing landscape: black-box, white-box, and gray-box testing. Black-box testing, also referred to as functional testing, focuses exclusively on the inputs and outputs of a system without delving into its internal logic. The goal is to verify that the software performs as intended for a wide range of input conditions.

White-box testing, or structural testing, takes the opposite approach by designing test cases based on the internal structure of the software. Testers using this technique have access to the source code and aim to cover various paths, branches, and statements

to ensure exhaustive coverage.

Gray-box testing combines elements of both black-box and white-box testing. In this hybrid approach, testers have partial knowledge of the internal structure, allowing them to create test scenarios that are informed by the design and functionality of the system. This technique is particularly effective in identifying integration issues and security vulnerabilities [26].

#### 2.2.4 Testing Phases or Levels

Software testing is performed at different levels of the development lifecycle, each with specific objectives and scopes. At the most granular level is **unit testing**, which involves verifying the correctness of individual components such as functions or methods. This phase typically occurs during the implementation stage and is often automated for rapid feedback.

Once individual units are validated, **integration testing** examines how these units work together when combined. This phase aims to uncover issues arising from interface mismatches or incorrect data flow between components.

**System testing** evaluates the complete, integrated application in an environment that closely simulates production. This phase checks for overall functionality, performance, and compliance with requirements.

To ensure that changes to the codebase do not introduce new defects, **regression testing** is employed. It re-validates the system after modifications to ensure stability and integrity.

Before deployment, **requirements testing** is conducted to validate whether the system meets the criteria outlined in the Software Requirements Specification (SRS). Complementing this, **scenario testing** exposes the software to realistic user scenarios to validate its behavior under diverse conditions.

Performance aspects are validated through **performance testing**, which assesses how the system behaves under varying workloads in terms of response time, resource

usage, and throughput.

During the release cycle, **alpha testing** is conducted by developers in a controlled environment to detect early issues. Finally, **beta testing** is performed in real-user environments, capturing feedback on usability and any residual defects prior to full deployment.

To better visualize these testing levels, Figure 2.1 illustrates the hierarchical flow from unit testing to beta testing.

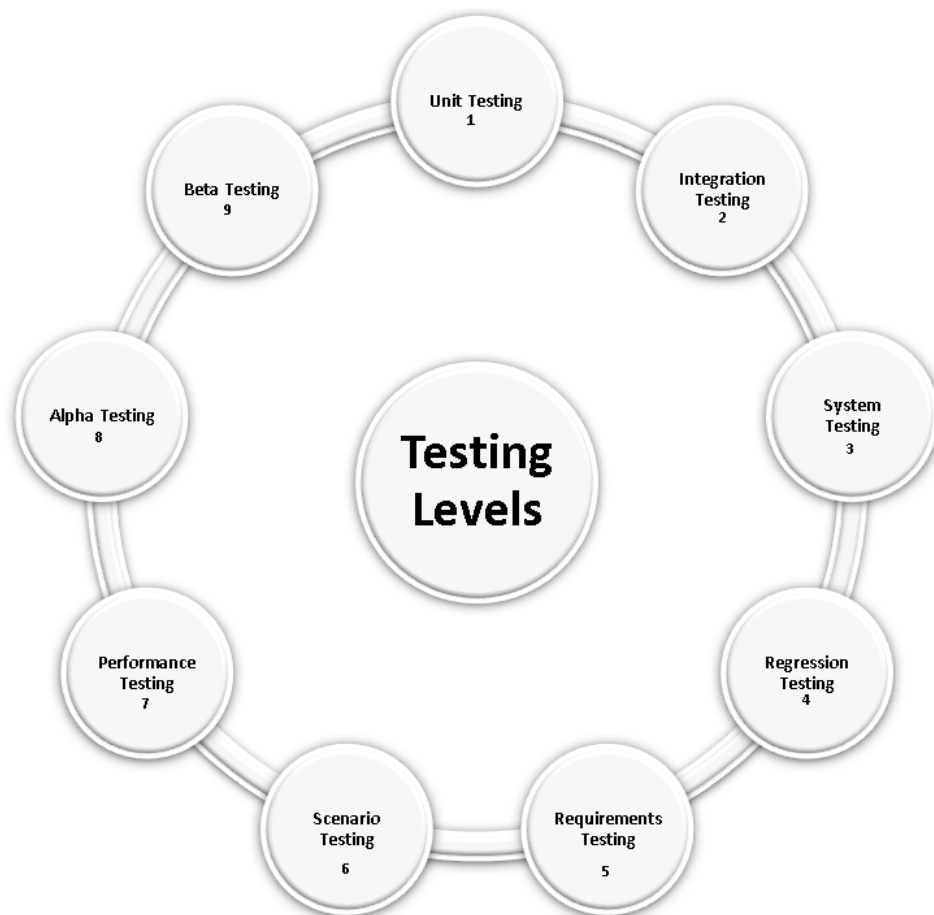


Figure 2.1: Hierarchy of Software Testing Levels. Adapted from [1].

## 2.3 Taxonomy of AI-Based Unit Test Case Generations

In recent years, the integration of Artificial Intelligence (AI) into software testing has introduced powerful capabilities to automate and enhance the generation of unit test cases. Traditional techniques for test case generation, which often rely on manual efforts

or static and dynamic analysis, are increasingly being augmented or replaced by intelligent systems that can learn from code repositories, predict test behaviors, and adapt to changes in software systems. Among these AI-driven approaches, machine learning (ML), deep learning (DL), and particularly transformer-based models have emerged as transformative tools for test case generation.

To comprehensively understand AI-driven techniques for unit test case generation, we categorize them into three major classes focused on direct generation capabilities: Traditional Machine Learning, Deep Learning, and Transformer-Based Models, as illustrated in Figure 2.2. Each approach has distinct characteristics, advantages, and applications in automated test case generation.

To evaluate the evolving landscape of AI-based unit test case generation, it is essential to develop a **taxonomy** that categorizes these approaches based on critical properties and performance characteristics. Such a taxonomy serves as a structured framework that enables researchers and practitioners to:

- Identify the strengths and limitations of various AI-based techniques for test generation.
- Compare approaches based on meaningful dimensions related to generation.
- Highlight gaps in the current research concerning test case generation.
- Facilitate the design of more robust, explainable, and practical testing systems.

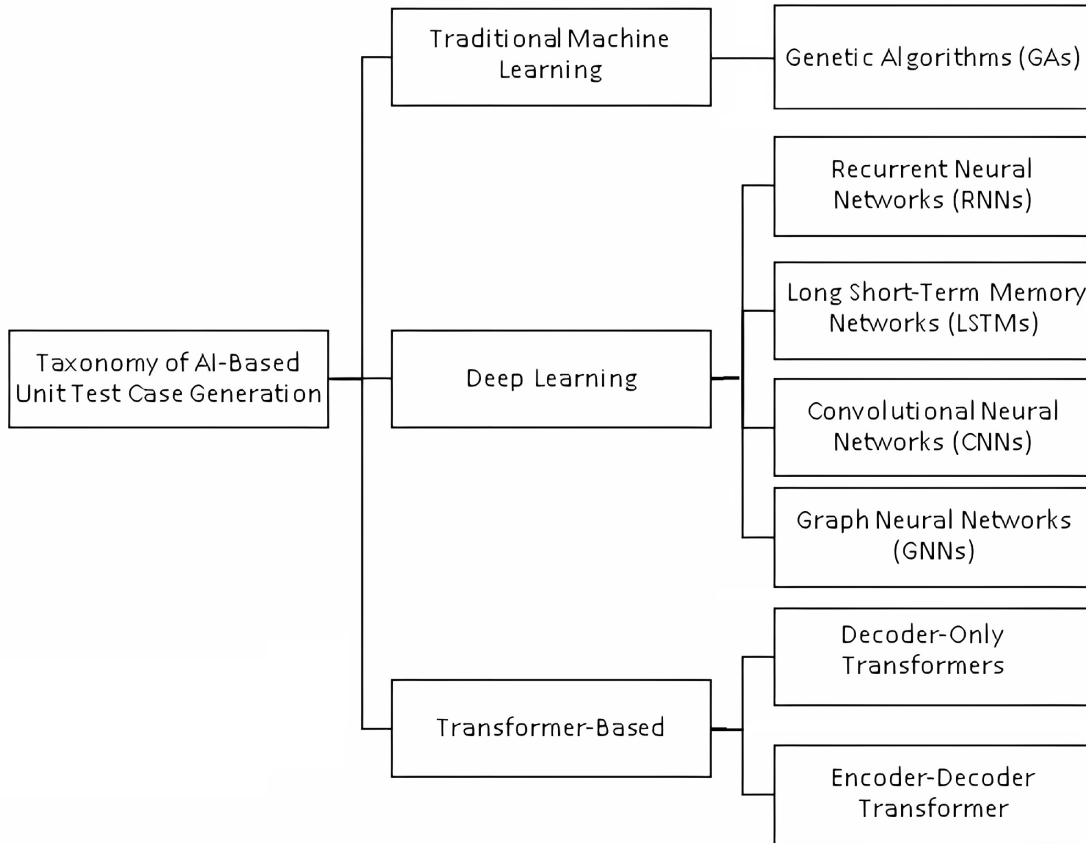


Figure 2.2: Taxonomy of AI-Based Unit Test Case Generation.

The following sections provide a detailed review of each approach, analyzing their methodologies, applications, and impact on automated test case generation.

## 2.4 Traditional Machine Learning Approaches

### 2.4.1 Genetic Algorithms (GAs)

Genetic Algorithms (GAs) are evolutionary computation techniques inspired by natural selection, widely used in unit test case generation. They excel at exploring large search spaces to optimize for objectives like test coverage and fault detection. GAs iteratively evolve a population of test cases through selection, crossover, and mutation. Test cases are evaluated using a fitness function to measure their effectiveness in achieving coverage goals (e.g., branch or statement coverage) or detecting faults. This process ensures the generation of diverse and optimized unit tests, although it can be computationally

expensive for large systems due to repeated execution and evaluation.

Several studies have demonstrated the effectiveness of GAs in unit test case generation. Harman et al. [27] introduced multi-objective genetic algorithms that optimize both test suite size and coverage, ensuring a balance between efficiency and effectiveness. Fraser et al. [28] developed EvoSuite, a widely used GA-based tool for Java unit test generation, which automatically produces test cases that maximize coverage criteria.

Despite their advantages, genetic algorithms have certain limitations. They can be computationally expensive, particularly for large-scale software systems, due to the need for repeated execution and evaluation of test cases. Additionally, designing an effective fitness function requires careful consideration, as it must balance coverage maximization with fault detection capabilities. However, when applied to large-scale software systems requiring high test coverage, GAs provide an effective approach for generating diverse and optimized unit tests.

## 2.5 Deep Learning Approaches

### 2.5.1 Recurrent Neural Networks (RNNs) and Long Short-Term Memory Networks (LSTMs)

Recurrent Neural Networks (RNNs) and their advanced variant, Long Short-Term Memory (LSTMs), are deep learning models effective for processing sequential data, making them suitable for automatic unit test generation. They learn temporal patterns and long-range dependencies in code, automatically extracting features from raw sequences to model complex syntactic and semantic relationships. LSTMs specifically overcome RNNs' vanishing gradient problem through gating mechanisms, enhancing their ability to capture long-term information. By processing tokenized or embedded code sequences, these models generate syntactically valid and semantically relevant unit tests. While computationally intensive, RNNs and LSTMs provide a powerful framework for generating tests reflecting intricate code behaviors.

Several studies have validated the efficacy of RNNs and LSTMs for unit test generation. Tufano et al. [29] proposed an LSTM-based model that generates test cases directly from method signatures, leveraging the model’s ability to understand method semantics and structure. Hellendoorn et al. [30] extended this by incorporating method-level context, demonstrating that deep neural networks can effectively model program behavior and improve test relevance.

Despite their strengths, these models are computationally intensive, often requiring high-performance hardware (e.g., GPUs) and significant training time. They perform best with large-scale, high-quality datasets, which may not always be available. However, when applied appropriately, RNNs and LSTMs offer a powerful framework for generating test cases that reflect intricate code behaviors and generalize across diverse codebases.

### 2.5.2 Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNNs) are deep learning models adept at analyzing local patterns in structured data, including source code. In automated test generation, CNNs recognize syntactic and structural code patterns, enabling the generation of test cases that conform to programming language constructs or target common bug-prone segments. Unlike sequential models, CNNs excel at capturing hierarchical and spatial relationships within code representations. They process code (e.g., Abstract Syntax Trees or tokenized code) with convolutional filters to identify patterns like API usage or control flow, thereby predicting and facilitating the generation of likely test inputs or fault-revealing test sequences, particularly useful for syntax-aware test cases.

Studies have demonstrated CNN effectiveness in software testing related to code analysis for generation. White et al. [31] applied CNNs to AST-based feature extraction, demonstrating that convolutional layers capture meaningful code structures for test generation.

CNNs offer efficiency and scalability advantages, processing code in parallel and reducing the need for manual feature engineering. They are less prone to the vanishing

gradient problem, enabling deeper representations of code patterns. However, CNNs struggle to capture global code context, as convolutional filters focus on localized regions, potentially missing long-range dependencies or inter-procedural behaviors. They also require substantial training data to generalize effectively, which can be a bottleneck for projects with limited test suites.

Despite these limitations, CNNs are effective for pattern-based test case generation, particularly for parsers, compilers, or API fuzzers. When combined with techniques like symbolic execution or coverage-guided fuzzing, CNNs enhance the diversity and precision of generated tests.

### 2.5.3 Graph Neural Networks (GNNs)

Graph Neural Networks (GNNs) are deep learning models designed for graph-structured data, making them ideal for analyzing software code represented as Abstract Syntax Trees (ASTs), Control Flow Graphs (CFGs), or Program Dependency Graphs (PDGs). Unlike traditional neural networks, GNNs model the relational structure of code, capturing crucial syntactic and semantic relationships between program elements. They represent code as a graph where nodes are program elements and edges are relationships. Through message passing, GNNs compute contextual representations, enabling the prediction of test inputs, identification of fault locations, or generation of assertions for effective unit test cases.

Studies have demonstrated GNN effectiveness in software testing. Allamanis et al. [32] showed that graph-based representations outperform sequence-based models in capturing program semantics for tasks like code generation. Wang et al. [33] combined GNNs with symbolic execution, guiding test generation toward interesting paths and improving efficiency.

GNNs capture both local and global dependencies, addressing CNN limitations, and are effective for object-oriented programming where class hierarchies and inter-method dependencies are critical. However, processing large code graphs is computationally ex-

pensive, with memory requirements growing rapidly. Training GNNs requires careful graph construction and message-passing design, and test quality depends on how program semantics are encoded.

Despite these challenges, GNNs excel at generating tests for complex program behaviors involving multiple interacting components, such as sequences of method calls maintaining object invariants. Combined with coverage guidance or search-based optimization, GNNs produce effective test suites.

## **2.6 Properties of AI-Based Unit Test Case Generation**

### **2.6.1 Essential Properties and Their Significance**

Automated unit test case generation has been a long-standing objective in software engineering research. With the advent of Artificial Intelligence—particularly deep learning and natural language processing—new opportunities have emerged for developing more intelligent and effective test generation techniques [22]. However, to enable practical adoption and long-term impact, AI-driven approaches must exhibit a set of desirable properties, as illustrated in Figure 2.3. This paper identifies and articulates these essential properties, presenting a comprehensive framework for evaluating current methods and guiding future research in AI-based unit test case generation.

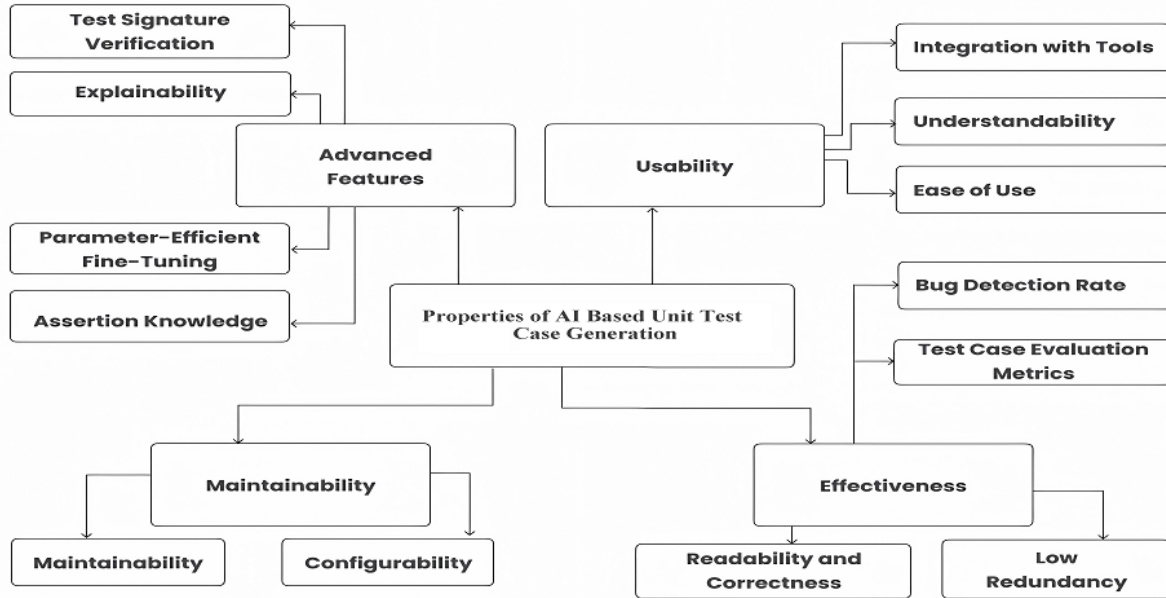


Figure 2.3: Properties of AI-Based Unit Test Case Generation.

The following properties are critical for the advancement and practical utility of AI-based unit test case generation:

1. **Test Signature Verification: Ensuring Foundational Correctness.** A fundamental requirement for any test generation technique is producing syntactically valid tests that align with the interface of the unit under test. Research must ensure that AI models can reliably infer and utilize function signatures to generate compilable and executable test cases.
2. **Explainability: Fostering Trust and Debuggability.** The "black-box" nature of some AI models can hinder developer trust and the ability to utilize generated tests effectively. Research on explainable AI in testing is crucial for understanding the reasoning behind test case generation, facilitating debugging and increasing developer confidence.
3. **Parameter-Efficient Fine-Tuning: Enabling Broad Applicability.** AI models must be adaptable to diverse software projects with limited data. Research on parameter-efficient fine-tuning techniques will enable generalizable models that can be customized for specific codebases, making AI-powered testing more practical.

**4. Assertion Knowledge: Enhancing Test Completeness.** Generating effective assertions is as critical as generating relevant inputs. Research exploring how AI can acquire and apply knowledge about expected program behavior will enhance test completeness, reducing manual effort.

**5. Integration with Tools: Facilitating Seamless Adoption.** For AI-based test generation to transition to industrial use, integration with existing development tools, IDEs, and CI/CD pipelines is essential to minimize disruption and maximize adoption.

These properties represent critical areas for ongoing and future research, advancing theoretical foundations and enabling practical, reliable, and widely adopted AI-powered testing tools.

### **2.6.2 Analysis for Traditional Machine Learning and Deep Learning Approaches for Test Generation**

The comparative analysis of AI-driven unit test generation approaches, as shown in Table 2.1, reveals significant variations in effectiveness across different methodologies. Genetic algorithms offer a balanced solution, providing comprehensive code coverage while maintaining reasonable computational efficiency, with strengths in redundancy reduction through evolutionary optimization. Among neural network approaches, graph neural networks excel in bug detection and structural coverage due to their ability to model code dependencies, though they require substantial resources for generation. LSTM and RNN architectures show moderate performance but are proficient in assertion generation, despite challenges with computational intensity and test case readability. Convolutional neural networks are adequate for syntax-oriented test generation but exhibit limitations in semantic understanding and assertion production.

The analysis highlights several patterns. First, a trade-off exists between test quality and computational complexity for generation, with sophisticated techniques demanding

greater resources. Second, approaches capturing structural code relationships often outperform those relying on sequential analysis for generating diverse and effective tests. Third, neural methods generally surpass traditional algorithms in coverage and fault detection, though with varying efficiency in generation. Practitioners should select approaches based on codebase complexity, available resources, and testing objectives (e.g., coverage, fault detection, or interpretability related to generation). The absence of a universally dominant approach underscores the importance of context-aware selection.

Table 2.1: Comparative Effectiveness of AI Approaches for Unit Test Generation

Key Papers	BD	LC	BC	MS	CA	TV	RD	FMC	RDAB	CC
Fraser et al. [28], Harman et al. [27]	73%	82-89%	76-84%	68%	19.7%	68-72%	28.4% → 89.2%	15.7%	Medium	Medium
Tufano et al. [29], Hellendoorn et al. [30]	65.2%	71.3%	63.8%	54%	37.4%	54-61%	22%	73.8%	Low-Medium	High
Wang et al. [33], Allamanis et al. [32]	85% (+22% vs CNN)	88.2%	85.1%	73%	58.1%	75-82%	18.9%	91.5%	Medium-High	Very High
White et al. [31], Dam et al. [34]	61.3% (syntactic)	64.9%	58.2%	49%	N/A	55-60%	25%	65.3%	Medium	Medium

**Abbreviations:** BD = Bug Detection, LC = Line Coverage, BC = Branch Coverage, MS = Mutation Score, CA = Correct Assertions, TV = Test Validity, RD = Redundancy, FMC = Focal Method Coverage, RDAB = Readability, CC = Computational Cost.  
**Note:** All percentage values represent reported effectiveness in original studies. Arrows (→) indicate improvement after optimization.

## 2.7 Transformer-Based Approaches

The Transformer is a neural architecture introduced by Vaswani et al. [2] that has gained significant popularity in natural language processing (NLP) and code-related tasks. It was designed to address the limitations of recurrent neural networks (RNNs) and long short-term memory (LSTM) networks in capturing long-range dependencies in sequential data. The key innovation of the Transformer lies in its attention mechanism, which allows it to process input sequences in parallel and capture global dependencies effectively. It makes it highly effective for code-related tasks such as unit test case generation [4, 20, 35].

Transformers have revolutionized automated unit test case generation by leveraging self-attention mechanisms to process code sequences and generate high-coverage test cases. Unlike traditional methods (e.g., genetic algorithms, graph neural networks), Transformers excel at capturing long-range dependencies in code, enabling end-to-end

test generation from raw source inputs. This section provides an overview of the Transformer architecture, its components, and the specific variants used in state-of-the-art unit test case generation, as identified in the related works.

### 2.7.1 Transformer Components

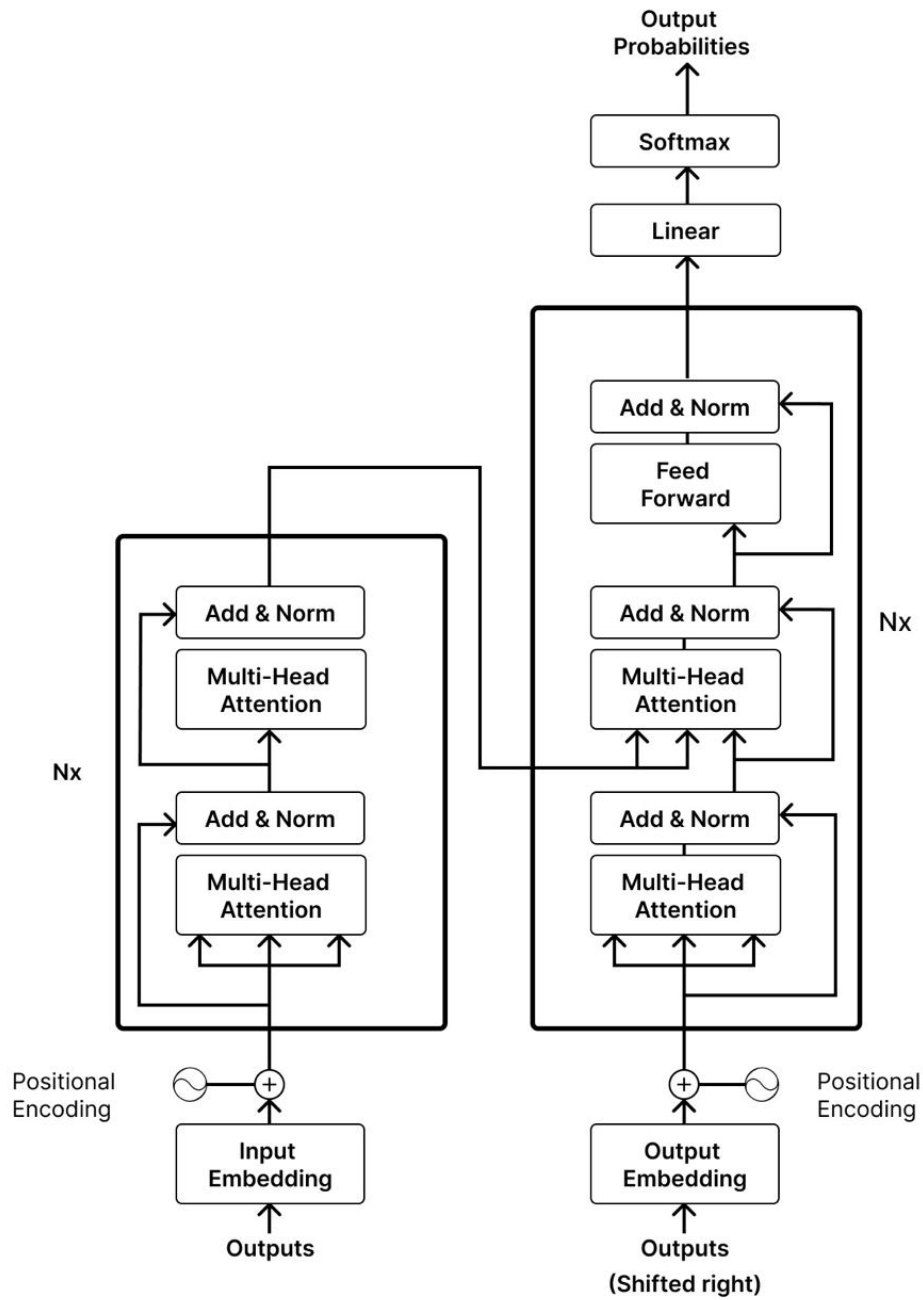


Figure 2.4: Overall Structure of the Transformer Model. Adapted from [2].

As shown in Figure 3.2, the **Encoder Block** processes the input sequence to extract meaningful representations [2]. It includes:

### Input Embedding

The process begins with **Input Embedding**, creating continuous-valued vectors (word embeddings) from the input sequence. **Positional Encoding** is added to provide positional information, using sine and cosine functions:

$$\text{PE}(\text{pos}, 2i) = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) [2] \quad (2.1)$$

$$\text{PE}(\text{pos}, 2i + 1) = \cos\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right) [2] \quad (2.2)$$

where  $\text{pos}$  is the position and  $i$  is the dimension.

The **Multi-Head Attention** [2] mechanism computes attention weights to capture contextual relationships:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (2.3)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W_O \quad (2.4)$$

where each head is:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (2.5)$$

The **Output Add & Norm** [2] layer applies residual connections and layer normalization:

$$\text{Residual}(x) = x + \text{Sublayer}(x) \quad (2.6)$$

$$\text{LayerNorm}(x + \text{Sublayer}(x)) \quad (2.7)$$

The **Feed Forward Net** [2] applies a fully connected layer:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (2.8)$$

## Output Block

The **Output Block** [2] generates predictions from the processed sequence, using a **Linear Layer**:

$$\text{Linear}(h) = hW + b \quad (2.9)$$

followed by a **Softmax Layer** to produce a probability distribution:

$$P(y_i | \text{input}) = \frac{e^{h_i}}{\sum_{j \in V} e^{h_j}} \quad (2.10)$$

### 2.7.2 Disadvantages of Transformers

Transformers are computationally expensive, with a quadratic time complexity of  $O(n^2)$  due to the self-attention mechanism, requiring substantial memory for long sequences. They are also difficult to train, particularly for domain-specific tasks, and may overfit on small datasets [2].

### 2.7.3 Transformer Variants Used in Unit Test Case Generation

The transformer variants used in state-of-the-art unit test case generation, as identified in the related works, include encoder-decoder and decoder-only models tailored for code-related tasks. These variants, adapted from models such as BART, PLBART, GPT-3, and GPT-3.5 [36, 37, 38], are summarized below.

#### Encoder-Decoder Transformers

Encoder-decoder Transformers include both an encoder and a decoder, making them suitable for tasks requiring input-output mappings, such as generating test cases from

code.

**BART (Bidirectional Auto-Regressive Transformer)** Introduced by Lewis et al. [36], BART combines bidirectional encoding with autoregressive decoding. It is used in AthenaTest [20] for unit test generation, pretrained on Java code and fine-tuned on the Method2Test dataset.

**PLBART (Programming Language BART)** PLBART, an extension of BART, is tailored for programming language tasks [36]. It is used in A3Test [4], pretrained on the Atlas dataset and fine-tuned on Method2Test, with rule-based methods for test signature verification.

### Decoder-Only Transformers

Decoder-only Transformers are designed for generative tasks, such as producing test cases from prompts or code snippets.

**GPT-3** Introduced by Brown et al. [10]. It is used in ChatUniTest [19] and Yuan et al. [35] with adaptive prompting strategies.

**GPT-3.5** An improved version of GPT-3, introduced by Brown et al. [38], GPT-3.5 enhances factual accuracy and reasoning. It is also used in ChatUniTest [19] and Yuan et al. [35] for generating high-quality test cases.

### Comparison of Transformer Models

Table 2.2 compares the transformer models used in unit test case generation, focusing on their parameters, pre-training data, training objectives, and performance.

Table 2.2: Transformer Models Comparison for Unit Test Case Generation

Model	Param Base (M)	Size Base (GB)	Pre-Train Data	Train Objective	Type	Performance
BART [36]	140	0.8	BookCorpus + Wikipedia + WebText	MLM + Sentence Augmentation	Encoder-Decoder	Strong in code-test generation
PLBART [36]	140	0.8	BookCorpus + Wikipedia + Code Datasets	MLM + Code-specific Objectives	Encoder-Decoder	Superior in code-related tasks
GPT-3 [37]	175,000	650	BookCorpus + WebText + Common Crawl	Autoregressive LM	Decoder-Only	High-quality test generation
GPT-3.5 [38]	175,000	650	BookCorpus + WebText + Common Crawl + Curated Data	Autoregressive LM	Decoder-Only	Improved test quality

Abbreviations: MLM = Masked Language Model

## 2.8 Parameter-Efficient Fine-Tuning Methods for Transformer

Transformer-based pretrained language models [39, 40, 41, 42] have achieved remarkable success across various NLP tasks. However, the escalating size of these models, particularly billion-parameter Large Language Models (LLMs), has introduced significant computational hurdles. Fine-tuning pretrained language models to specific tasks is computationally intensive and often impractical due to resource constraints.

Traditional fine-tuning updates all model parameters, making it prohibitively expensive for LLMs like Falcon-180B [43]. Parameter-Efficient Fine-Tuning (PEFT) [15] has emerged as a solution to reduce the number of trainable parameters while preserving performance comparable to full fine-tuning. PEFT adapts the knowledge of pretrained models to target tasks, lowering the risk of catastrophic forgetting by updating only a small subset of parameters. It also mitigates overfitting, as fine-tuned datasets are typically smaller than pretrained datasets. The mathematical representation of PEFT can be formulated as:

$$\theta_{\text{new}} = \theta_{\text{pre-trained}} + \Delta\theta[15] \tag{2.11}$$

where  $\theta_{\text{new}}$  is the updated set of parameters,  $\theta_{\text{pre-trained}}$  is the original set, and  $\Delta\theta$  represents the small subset modified during fine-tuning. PEFT is particularly advantageous for fine-tuning large models, enabling task-specific adaptation with minimal computational costs.

### 2.8.1 Major Approaches of PEFT

PEFT includes several techniques that adapt large pretrained models to specific tasks while minimizing computational costs, conserving memory and resources. Each approach provides unique benefits for applications ranging from text classification to generative tasks.

#### 1. Adapter-Based Fine-Tuning

Adapter-based fine-tuning introduces small, task-specific modules known as adapters, inserted between the layers of a pretrained model. During fine-tuning, only these adapters are trained, keeping the original parameters frozen. This modular approach allows efficient task adaptation, expressed as:

$$\theta_{\text{new}} = \theta_{\text{pre-trained}} + A \text{ [15]} \tag{2.12}$$

where  $A$  represents the adapter parameters learned during fine-tuning. Adapters require updating only a fraction of parameters, maintain generalization abilities, and allow different adapters for various tasks to be stored and reused. Houlsby et al. [15] demonstrated that adapters achieve task-specific performance comparable to full fine-tuning with a smaller computational footprint. AdapterFusion [44] enhances multi-task learning by combining multiple adapters.

#### 2. Prompt-Based Fine-Tuning (Soft Prompts)

Prompt-based fine-tuning involves adding learnable prompt tokens or vectors to the input embeddings, optimized during fine-tuning to guide the pretrained model without

modifying its parameters. The formulation is:

$$\text{Output} = f(x + P)[45, 46] \quad (2.13)$$

where  $f$  is the pretrained model function,  $x$  is the input, and  $P$  represents the learnable prompt vectors. This lightweight method avoids altering the model architecture. Li and Liang [45] introduced Prefix-Tuning, showing competitive results, while Lester et al. [46] demonstrated its effectiveness in text generation and classification.

### 3. LoRA (Low-Rank Adaptation)-Based Fine-Tuning

LoRA introduces low-rank adaptation matrices into transformer layers, updating only these matrices to reduce trainable parameters. The representation is:

$$\Delta W = W + UV^T[14] \quad (2.14)$$

where  $W$  is the original weight matrix, and  $U$  and  $V$  are low-rank matrices. Hu et al. [14] showed LoRA's effectiveness in fine-tuning large models like GPT-3 with minimal parameters.

### 4. Prefix Tuning

Prefix tuning adds learnable prefix vectors to the hidden states of each layer, updated during fine-tuning while keeping the original parameters unchanged. The formulation is:

$$h_{\text{new}} = h + P[45] \quad (2.15)$$

where  $h$  represents hidden states and  $P$  denotes prefix vectors. Li and Liang [45] demonstrated that prefix tuning reduces costs while maintaining performance.

## 5. BitFit (Bias-Tuning)

BitFit updates only the bias terms of the model’s layers, keeping weight matrices frozen, represented as:

$$\theta_{\text{new}} = \theta_{\text{pre-trained}} + b[47] \quad (2.16)$$

where  $b$  represents the bias terms learned during fine-tuning. Ben Zaken et al. [47] showed competitive results with minimal parameter updates.

## 6. QLoRa (Quantized Low-Rank Adaptation)

QLoRa combines quantization and low-rank adaptation, reducing the precision of weights and activations to lower memory usage. The representation is:

$$\theta_{\text{quantized}} = Q(\theta_{\text{pre-trained}}) + UV^T[3] \quad (2.17)$$

where  $Q$  is the quantization function, and  $U$  and  $V$  are low-rank matrices. QLoRa achieves significant memory efficiency, suitable for resource-constrained devices [14].

### 2.8.2 Rationale for Selecting PEFT Methods

Table 2.3 compares six PEFT techniques: Adapters, LoRA, QLoRA, Prompt Tuning, Prefix Tuning, and BitFit. We selected Adapters, LoRA, and QLoRA for their scalability, expressive parameter adaptation, and effectiveness in generative tasks like unit test generation. LoRA and QLoRA are suited for large models, reducing trainable parameters via low-rank decomposition and quantization [14, 3]. Adapters offer modularity for multi-task adaptation [15, 44]. Prompt Tuning and Prefix Tuning are less expressive for complex generative tasks, while BitFit lacks capacity for rich generative behavior [46, 45, 48].

Figure 2.5 illustrates how these PEFT techniques integrate with the transformer architecture, such as the PLBART-base model.

Table 2.3: Comparison of PEFT Methods

PEFT Method	Main Idea	Parameter Efficiency	Efficiency	Suitability for Test Generation
Adapters [15, 44]	Insert trainable modules between transformer layers	High – Only adapter layers trained	High	High – Modular, supports multi-task adaptation
LoRA [14]	Inject low-rank matrices into attention mechanisms	Very High – Reduces updates significantly	High	Very High – Effective in generative settings
QLoRA [3]	Combines LoRA with quantization for efficiency	Extremely High – Suited for memory-limited LLMs	High	Very High – Preserves performance under quantization
Prompt Tuning [46]	Learn soft prompts prepended to input embeddings	High – Only prompt vectors trained	Moderate	Moderate – Less effective for long/-generative tasks
Prefix Tuning [45]	Adds trainable prefixes to hidden states	High – Minimal tuning needed	Moderate	Moderate – Limited for complex generative tasks
BitFit [48]	Fine-tune only the bias terms in the model	Very High – Minimal parameter updates	Low	Low – Inadequate for rich generative behavior

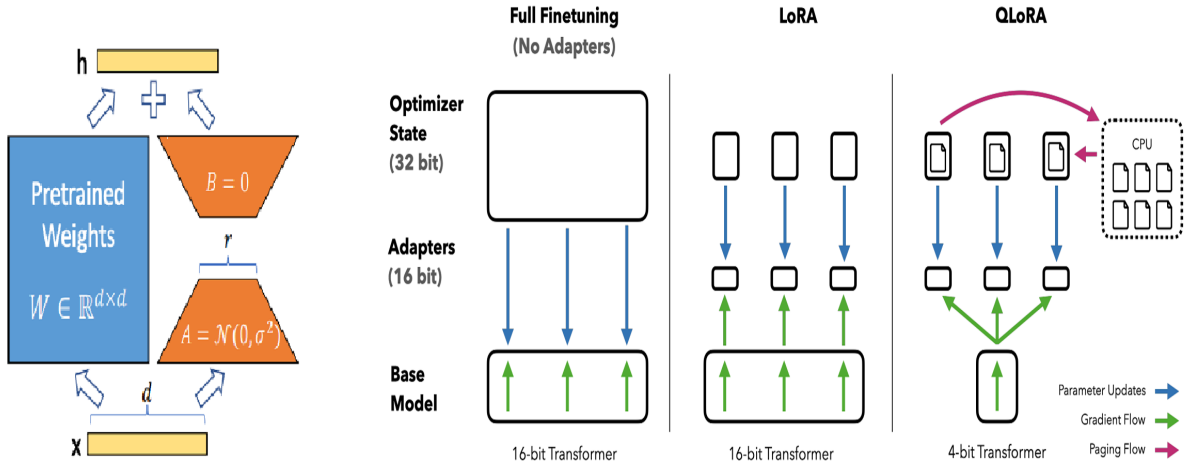


Figure 2.5: Conceptual Illustration of PEFT Techniques including Full Fine-tuning, LoRA, and QLoRA. Adapted from Dettmers et al. [3].

The figure clarifies how fine-tuning strategies interact with pretrained weights. Full fine-tuning updates all parameters, demanding significant resources. Adapters insert trainable modules, LoRA injects low-rank matrices, and QLoRA applies quantization for memory efficiency, making them suitable for scaling unit test generation.

## 2.9 Related Works

This section reviews recent advancements in AI-driven unit test case generation, with a focus on transformer-based models that dominate state-of-the-art methodologies. We survey key studies, evaluate their contributions based on critical properties (effectiveness, usability, maintainability, and advanced features), compare baseline models, and propose a novel approach leveraging parameter-efficient fine-tuning to address current limitations.

Transformer-based models have significantly advanced unit test case generation by leveraging contextual code understanding and generative capabilities. Tufano et al. [20] introduced AthenaTest, utilizing the Focal Transformer to enhance test case generation by capturing contextual code information. De Sousa et al. [49] developed JavaBERT, tailored for Java applications, which uses pretrained representations to generate relevant test cases. Chen et al. [50] proposed Codet, enabling cohesive code and test case generation. Lahiri et al. [51] introduced Interactive TestGen, a user-guided framework that incorporates developer feedback to refine transformer-based test generation. Tufano et al. [52] focused on generating accurate assert statements, improving test reliability. Lemieux et al. [53] addressed coverage plateaus with CODAMOSA, optimizing test case diversity. Xie et al. [19] proposed ChatUniTest, leveraging ChatGPT with adaptive prompting for high-quality tests. Alagarsamy et al. [4] introduced A3Test, enhancing assertion accuracy and test signature verification. Yuan et al. [35] improved ChatGPT’s context-aware test generation, achieving robust coverage and readability.

A3Test outperforms AthenaTest and ChatUniTest on the Defects4j dataset, generating up to 395.43% more correct test cases, achieving 55.56% more correct assertions, and improving readability compared to EvoSuite, demonstrating significant advancements in test quality and efficiency.

### 2.9.1 Evaluation of Key Properties

To assess the impact of these approaches, we evaluate them based on a framework of essential properties, as illustrated in Figure 2.3, focusing on effectiveness, usability, maintainability, and advanced features.

#### Effectiveness

Effectiveness is measured using metrics such as bug detection rate, readability, redundancy, test case correctness, and coverage, as shown in Table 2.4. Many studies, including Tufano et al. [20] and Chen et al. [50], do not report bug detection rates, highlighting a need for standardized metrics. Tufano et al. [20] reported a 23% correct test case rate, while Alagarsamy et al. [4] and Xie et al. [19] achieved 57% and 59%, respectively, reflecting variability due to differences in datasets and methodologies. Readability is emphasized in works like De Sousa and Hasselbring [49] and Xie et al. [19], with developers favoring clear AI-generated tests. Low redundancy is addressed by Chen et al. [50] and Lemieux et al. [53], but quantification remains limited. Coverage metrics, including line and branch coverage, show strong results in Alagarsamy et al. [4] and Xie et al. [19], but mutation scores and assertion correctness are less consistently reported, indicating gaps in evaluation.

Table 2.4: Effectiveness of AI for Unit Test Case Generation

Study	Bug Detection Rate	De-Redundancy	Readability and Correctness	Low Redundancy	Re-Test Case	Correct Test Case	Focal Method Coverage	Line Coverage	Cov-Branch Coverage	Mutation Score	Correct Assertion
Tufano et al. [20] - ATHENATEST	-	✓	✓	✓	23%	36%	32%	39%	21%	25%	
De Sousa & Hasselbring [49]	-	✓	-	-	-	-	-	-	-	-	
Chen et al. [50]	-	-	✓	-	-	-	✓	-	-	-	
Lahiri et al. [51]	-	✓	✓	-	-	-	✓	-	-	-	
Tufano et al. [52]	-	✓	✓	✓	✓	-	✓	-	-	-	
Alagarsamy et al. [4]	-	✓	✓	✓	57%	62%	67%	70%	46%	42%	
Lemieux et al. [53]	-	-	✓	-	-	✓	✓	-	-	-	
Xie et al. [19]	-	✓	✓	✓	59%	77%	67%	73%	46%	43%	
Yuan et al. [35]	✓	✓	✓	✓	✓	✓	✓	✓	-	✓	

**Test Case Evaluation Metrics:** Correct Test Case, Focal Method, Coverage, Line Coverage, Branch Coverage, Mutation Score, and Correct Assertion.

## Usability

Usability, critical for adoption, is evaluated based on ease of use, understandability, and tool integration, as shown in Table 2.5. Tufano et al. [20] and Alagarsamy et al. [4] designed readable tests for seamless integration with existing frameworks. De Sousa and Hasselbring [49] implied usability in familiar environments, while Chen et al. [50] emphasized integration capabilities. Lahiri et al. [51] enhanced usability through user feedback, and Xie et al. [19] facilitated adoption via plugins. Lemieux et al. [53] noted complexity challenges despite clear explanations, while Yuan et al. [35] highlighted developer adoption ease.

Table 2.5: Usability of AI for Unit Test Case Generation

Study	Ease of Use	Understandability	Integration with Tools
Tufano et al. [20]	✓	✓	-
De Sousa & Hasselbring [49]	-	✓	-
Chen et al. [50]	-	-	✓
Lahiri et al. [51]	✓	✓	-
Tufano et al. [52]	✓	✓	✓
Alagarsamy et al. [4]	✓	✓	✓
Lemieux et al. [53]	-	✓	✓
Xie et al. [19]	✓	✓	✓
Yuan et al. [35]	✓	✓	-

## Maintainability

Maintainability ensures long-term viability, as shown in Table 2.6. Tufano et al. [20] and Tufano et al. [52] suggested maintainability through human-like tests and pretrained models. Alagarsamy et al. [4] and Xie et al. [19] supported configurability and maintainability via hyperparameter settings and iterative refinement. Chen et al. [50] and Lemieux et al. [53] focused on configurability, while De Sousa and Hasselbring [49] implied reduced maintenance with pretrained models. Yuan et al. [35] emphasized maintainability through iterative improvements.

Table 2.6: Maintainability of AI for Unit Test Case Generation

Properties	Configurability	Maintainability
Tufano et al. [20]	-	✓
De Sousa & Hasselbring [49]	-	-
Chen et al. [50]	✓	-
Lahiri et al. [51]	✓	-
Tufano et al. [52]	-	✓
Alagarsamy et al. [4]	✓	✓
Lemieux et al. [53]	✓	-
Xie et al. [19]	✓	✓
Yuan et al. [35]	-	✓

## Advanced Features

Advanced features, such as explainability, assertion knowledge, test signature verification, and parameter-efficient fine-tuning, are critical for practical adoption, as shown in Table 2.7. Tufano et al. [20], Tufano et al. [52], Chen et al. [50], Alagarsamy et al. [4], Xie et al. [19], and Yuan et al. [35] emphasized assertion knowledge, with Alagarsamy et al. [4] also incorporating test signature verification. Explainability and parameter-efficient fine-tuning remain underexplored, with no studies addressing these features comprehensively. De Sousa and Hasselbring [49], Lahiri et al. [51], and Lemieux et al. [53] provided limited insights into advanced features.

Table 2.7: Advanced Features for AI Unit Test Case Generation

Paper	Explainability	Assertion Knowledge	Test Signature Verification	Parameter Efficient Fine-Tuning
Tufano et al. [20]	-	✓	-	-
De Sousa & Hasselbring [49]	-	-	-	-
Chen et al. [50]	-	✓	-	-
Lahiri et al. [51]	-	-	-	-
Tufano et al. [52]	-	✓	-	-
Alagarsamy et al. [4]	-	✓	-	-
Lemieux et al. [53]	-	-	-	-
Xie et al. [19]	-	✓	-	-
Yuan et al. [35]	-	✓	-	-

## 2.9.2 Baseline Model Comparison

We compare three baseline models—A3Test [4], AthenaTest [20], and ChatUniTest [19]—to assess their performance across various metrics.

### Model Overview

A3Test uses PLBART, pretrained on the Atlas dataset and fine-tuned on Method2Test, with rule-based methods for test signature verification. AthenaTest employs BART, pretrained on Java code and fine-tuned on Method2Test. ChatUniTest leverages ChatGPT-3 with adaptive focal context and multiple validation steps.

### Test Case Generation

Table 2.8 compares the number of correct, passing, failing, build error, and syntax error test cases for 5,278 focal methods. A3Test generates 2,262 correct test cases, significantly outperforming AthenaTest (440) and slightly surpassing ChatUniTest (2,099).

Table 2.8: Comparison of Generated Test Cases Based on Focal Method

Model	Correct	Pass	Fail	Build Error	Syntax Error
A3Test [4]	2,262	2,996 (57%)	340 (6%)	1,766 (33%)	546 (10%)
AthenaTest [20]	440	1,206 (23%)	1,509 (29%)	2,401 (45%)	536 (10%)
ChatUnitTest [19]	2,099	3,112 (59%)	577 (11%)	1,374 (26%)	215 (4%)

Total Focal Methods: 5278.

### Focal Method Coverage

Table 2.9 shows unique focal method coverage. ChatUniTest covers 1,717 methods (33%), followed by A3Test (1,433, 27%) and AthenaTest (163, 3%).

### Coverage Metrics

Table 2.10 presents line, branch, and focal method coverage. A3Test and ChatUniTest achieve 67% line coverage, while ChatUniTest leads in branch (73%) and focal method (77%) coverage. AthenaTest lags with 32% line and 39% branch coverage.

Table 2.9: Focal Method Coverage Analysis

Model	Focal Method Analysis
A3Test [4]	1,433 (27%)
AthenaTest [20]	163 (3%)
ChatUnitTest [19]	1,717 (33%)

Total Focal Methods: 5278.

Table 2.10: Coverage Metrics Comparison

Model	Line Coverage (%)	Branch Coverage (%)	Focal Method Coverage (%)
A3Test [4]	74,562 (67%)	44,013 (70%)	5,278 (62%)
AthenaTest [20]	35,457 (32%)	30,809 (39%)	3,251 (36%)
ChatUnitTest [19]	75,173 (67%)	17,143 (73%)	1,913 (77%)

Total Lines of code(loc): 111,809. Total Branch: 44,013 Total Focal Method: 5,278

## Assertions

Table 2.11 compares assertion generation. ChatUnitTest generates the most assertions (124,110) with 43% correctness, followed by A3Test (108,566, 46%) and AthenaTest (45,330, 25%).

Table 2.11: Assertions Comparison

Model	Total Assertions	Correct Assertions	Correct (%)
A3Test [4]	108,566	49,783	46%
AthenaTest [20]	45,330	11,233	25%
ChatUnitTest [19]	124,110	53,360	43%

## Computational Efficiency

Table 2.12 summarizes computational time and cost. A3Test is the fastest, generating tests in 2.9 hours (1.98s/test), 20.6% faster than AthenaTest (3.5 hours, 2.37s/test). ChatUnitTest is slowest at 4.31s/test, costing \$87.00 for Defects4j projects.

Table 2.12: Computational Time and Cost Overview

Model	Avg Time per Test Case (s)	Time for 1 Attempt (h)	Cost for Test Generation (\$)
A3Test [4]	1.98	2.9	–
AthenaTest [20]	2.34	3.5	–
ChatUnitTest [19]	4.31	6.32	87.20

### 2.9.3 Summary

Transformer-based models have significantly advanced AI-driven unit test case generation, offering robust solutions for generating readable and effective test cases. Studies like Tufano et al. [20], Alagarsamy et al. [4], and Xie et al. [19] demonstrate strong performance in test case correctness (23%–59%) and coverage (up to 77% focal method coverage), with A3Test and ChatUniTest outperforming AthenaTest in assertion accuracy (46% and 43% vs. 25%) and computational efficiency. Usability is enhanced through integration with development tools and developer-friendly interfaces, as seen in Lahiri et al. [51] and Yuan et al. [35], while maintainability is supported by configurable models in Alagarsamy et al. [4] and Xie et al. [19]. Assertion knowledge is a common strength, but explainability and parameter-efficient fine-tuning remain underexplored, limiting practical adoption. Table 2.13 summarizes the performance of these approaches across key properties, highlighting their strengths and gaps.

Table 2.13: Comparative Performance of AI-Based Unit Test Generation Approaches

Papers	Bug Detection	Readability	Coverage	Correct Assertions	Low Redundancy	Ease of Use	Maintainability	PEFT
Tufano et al. [20]	–	✓	✓	✓	✓	✓	✓	–
De Sousa et al. [49]	–	✓	–	–	–	–	–	–
Chen et al. [50]	–	–	✓	✓	✓	–	–	–
Lahiri et al. [51]	–	✓	✓	–	✓	✓	–	–
Tufano et al. [52]	–	✓	✓	✓	✓	✓	✓	–
Alagarsamy et al. [4]	✓	✓	✓	✓	✓	✓	✓	–
Lemieux et al. [53]	–	–	✓	–	✓	–	–	–
Xie et al. [19]	–	✓	✓	✓	✓	✓	✓	–
Yuan et al. [35]	✓	✓	✓	✓	✓	✓	✓	–

Despite these advancements, challenges persist, including inconsistent bug detection metrics, limited framework integration, and underdeveloped explainability. Future research should focus on enhancing assertion knowledge with domain-specific insights, improving test signature verification, and exploring parameter-efficient fine-tuning to reduce resource demands and broaden applicability.

# Chapter 3

## Methodology

### 3.1 Research Methodology

To address the challenges in AI-based unit test generation and design a parameter-efficient solution, we adopted the Design Science Research Process (DSRP) [54] as our core methodology. DSRP, depicted in Figure 3.1, consists of six major steps: Problem Identification & Motivation, Objective Formulation, Design and Development, Demonstration, Evaluation, and Communication. In our implementation, we combined the Demonstration and Evaluation steps into a single *Evaluation* phase, as our assessment involved both demonstrating the proposed architecture and evaluating its performance gains.

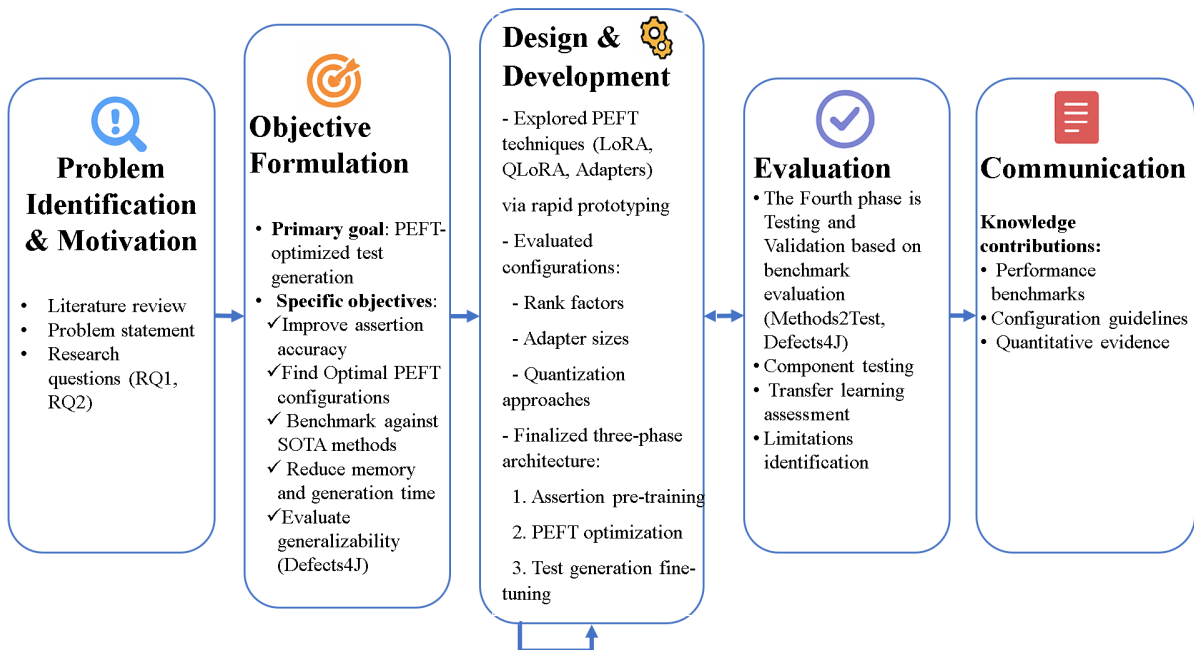


Figure 3.1: Design Science Research Process-Based, Problem-Centered Research Methodology for AI-Powered Unit Test Generation.

Given that the motivation for this work arises from the accuracy-efficiency trade-off in-

herent in transformer-based test generation, we employed a Problem-Centered Approach throughout our research. The following outlines how each DSRP step was applied:

- **Problem Identification & Motivation:**

We conducted a preliminary literature review, examining survey articles and related works that attempted to address these challenges (Chapter 2). The problem was formally stated as the redundant computational overhead of full fine-tuning, which incurs unnecessary costs without proportional gains in test quality. We then formulated two research questions (Section 1.4) to guide our approach:

- **RQ1:** What is the optimal configuration of PEFT components for unit test generation?
- **RQ2:** How does the PEFT-optimized system perform compared to the A3Test [4] baseline in test case generation quality?

- **Objective Formulation:**

To guide our research, we established a primary goal: to design and implement an algorithm that leverages Parameter-Efficient Fine-Tuning (PEFT) techniques to enhance AI-based unit test generation for Java software. Flowing from this general aim, we set several specific objectives. These included assessing the impact of PEFT on both assertion accuracy and bug detection capabilities, identifying the optimal PEFT configurations among LoRA, QLoRA, and Adapters, and conducting a comparative analysis of PEFT against traditional full fine-tuning methods and existing state-of-the-art tools such as A3Test and AthenaTest. As detailed in Section 1.4, these clearly defined objectives provided a crucial roadmap for our research process, ensuring a consistent focus on achieving actionable outcomes.

- **Design & Development:**

Building on the problem understanding, we engaged in an iterative design process to develop our *OptimizedUnitTest* model, building upon the methodology proposed

by [4]. We explored various PEFT techniques (LoRA, QLoRA, Adapters) through rapid prototyping cycles, evaluating different configurations of rank factors, adapter sizes, and quantization approaches. Through this process, we converged on a three-phase architecture combining assertion pre-training, PEFT optimization, and test generation fine-tuning, with novel optimizations in the PEFT phase.

- **Evaluation:**

We conducted comprehensive evaluations through multiple stages: First, we tested the framework’s components independently (assertion generation, PEFT efficiency, and full test generation). Using the optimal configurations identified, we then evaluated the complete system on benchmark datasets (*Methods2Test*, *Defects4J*) to assess its effectiveness compared to state-of-the-art approaches. Furthermore, we evaluated the model’s applicability in transfer learning scenarios to gauge real-world usability. During evaluation, we identified key limitations around assertion diversity and computational trade-offs, leading to targeted optimizations in our final design.

- **Communication:**

The knowledge contributions derived from this process are documented in this dissertation, including: (1) performance benchmarks comparing PEFT approaches, (2) configuration guidelines for efficient test generation, and (3) quantitative evidence of the model’s advantages over existing methods. These findings advance both research and practice in AI-assisted software testing.

## 3.2 An Optimized Unit Test Case Generation Methodology

Our methodology, which we refer to as the Optimized Unit Test Case Generation Methodology, builds upon the approach proposed by [4] and introduces novel optimizations for designing and developing the `OptimizedUnitTest` model to generate efficient Java unit tests. It leverages the Design Science Research Process (DSRP) [54] to ensure methodological rigor and practical relevance. Figure 3.2 illustrates the multi-stage pipeline,

which integrates assertion pre-training, parameter-efficient fine-tuning (PEFT), test generation fine-tuning, and validation to achieve high-quality test generation with minimal computational overhead. While Phases 1, 3, and 4 are adopted from [4], our primary contribution lies in the optimization of Phase 2, where we strategically apply and compare PEFT techniques to enhance efficiency and test quality.

The workflow begins with Phase 1: Pre-Training to Learn Meaningful Assert Statements, adopted from [4]. In this step, a PLBART model [18] is pre-trained on the Atlas Dataset [55], which contains pairs of Java focal methods and assertion statements. Operating as a Masked Language Model (MLM), the model learns to predict masked assertion tokens and their corresponding types, building a robust understanding of Java assertion syntax and semantics. This process results in the AK-PLBART (Assertion Knowledge PLBART) model, which serves as the foundation for subsequent optimization and fine-tuning phases.

Following this, Phase 2: PEFT to Optimize a Pre-trained Assert Model applies Parameter-Efficient Fine-Tuning (PEFT) techniques, such as LoRA, QLoRA, or Adapters, to the AK-PLBART model. This phase, our novel contribution, optimizes the application of PEFT techniques with specific configurations (e.g., rank-8 for LoRA and QLoRA, bottleneck size of 64 for adapters) to drastically reduce computational and memory overhead compared to full fine-tuning, yielding an OAK-PLBART (Optimized Assertion Knowledge PLBART) model.

Next, Phase 3: Fine-Tuning to Learn Meaningful Test Cases, adopted from [4], leverages the OAK-PLBART model for comprehensive unit test generation. Utilizing the Methods2Test Dataset, which contains full method-test pairs, the model is fine-tuned to learn the complete structure and logic of unit tests, including test method signatures and instantiations. This phase transforms the OAK-PLBART into a full-fledged test case generation engine.

Finally, in Phase 4: Generating Test Cases and Validation, also adopted from [4], the OptimizedUnitTest Model (the output of Phase 3) is applied to generate tests for

the Defects4J Dataset, a challenging unseen benchmark. The generated test candidates undergo rigorous post-processing, including rule-based Java syntax checking and verification of naming conventions and test signatures, ensuring their validity and readiness for further evaluation. Our optimized multi-stage approach emphasizes efficient PEFT optimization in Phase 2 to deliver a robust and scalable model for automated Java unit test generation.

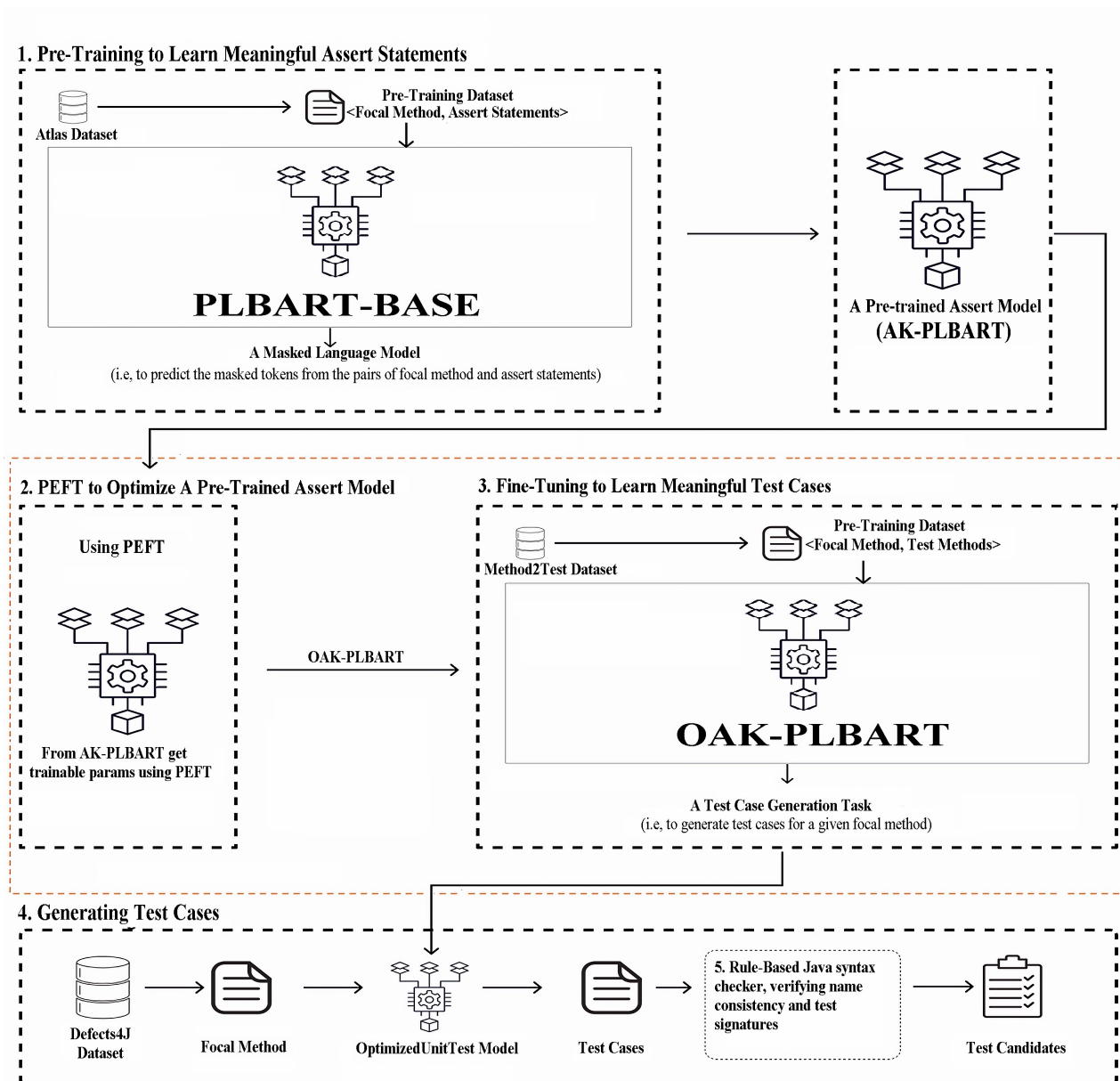


Figure 3.2: Optimized approach for Unit Test Case Generation for Java using Transformers, builds upon the approach proposed by [4].

The novelty of our approach lies in the optimization of the PEFT phase (Phase 2)

within the multi-stage pipeline adopted from [4]. Specifically, we enhance the methodology through:

- A **strategically optimized application of PEFT techniques** (LoRA, QLoRA, Adapters) in Phase 2, tailored for unit test generation with specific configurations (e.g., rank-8,  $\alpha = 32$  for LoRA/QLoRA, bottleneck size of 64 for adapters) to balance performance and computational efficiency [14, 3, 15].
- A **rigorous comparative analysis** of PEFT techniques within the context of unit test generation, identifying the most effective configurations for practical application (detailed in Phase 2 and Algorithm 2).

### 3.2.1 Phase 1: Assertion Pre-Training

The assertion pre-training phase utilizes the Atlas dataset [55], which contains 1.2 million Java method-assertion pairs. The dataset is split into 960,000 examples (80%) for training and 240,000 examples (20%) for validation, based on the default split in the training implementation. To ensure consistency and improve learning efficiency, the dataset undergoes several preprocessing steps.

First, code normalization is applied by removing single-line and multi-line comments, standardizing whitespace and formatting, and reconstructing the code with consistent structure using JavaParser. Next, an assertion masking strategy is implemented, where 15% of assertion tokens are randomly masked following the BERT-style masking approach [39]. Special handling is applied to assertion types such as `assertEquals` and `assertThrows` to preserve their semantic integrity.

To address class imbalance, dataset balancing techniques are employed by assigning sampling weights to different assertion types, including `assertEquals`, `assertTrue`, and `assertThrows`. Additionally, special transformations such as value swapping, boolean flipping, and exception typing are applied to enhance the diversity of training samples. These preprocessing steps collectively improve the robustness and generalizability of the

pre-trained model for unit test generation.

**Model Architecture:** The assertion model extends PLBART [18] with custom modifications to suit the unit test generation task, as described in [4]. Key modifications include the incorporation of a custom assertion-type classifier alongside the masked language modeling head, allowing the model to predict not only masked tokens but also the type of assertion in the code.

---

**Algorithm 1** Assertion-Based Pre-Training

---

```

1: function ASSERTIONPRETRAINING(model, atlas_dataset, epochs)
2:   for epoch = 1 to epochs do
3:     for batch in atlas_dataset do
4:       normalized_method  $\leftarrow$  normalize_method(batch.method)
5:       masked_assertion  $\leftarrow$  mask_assertion(batch.assertion)  $\triangleright$  BERT-style
       masking [39]
6:       assertion_type  $\leftarrow$  determine_assertion_type(batch.assertion)
7:       mlm_logits, assert_logits  $\leftarrow$  model(normalized_method, masked_assertion)
8:       mlm_loss  $\leftarrow$  compute_masked_language_model_loss(mlm_logits,
       masked_assertion)
9:       assert_loss  $\leftarrow$  compute_assertion_type_loss(assert_logits, assertion_type)
10:      total_loss  $\leftarrow$  mlm_loss + (0.3  $\times$  assert_loss)
11:      model.optimizer.zero_grad()
12:      total_loss.backward()
13:      model.optimizer.step()  $\triangleright$  Update model parameters
14:    end for
15:  end for
16:  return model
17: end function

```

---

Algorithm 1, adopted from [4], implements assertion pre-training using the masked language modeling approach of [39] and the PLBART architecture [18], with modifications for unit test assertion generation.

### 3.2.2 Phase 2: PEFT Optimization

In this phase, our novel contribution, we take the **Assertion Knowledge PLBART (AK-PLBART) model** from Phase 1 and apply Parameter-Efficient Fine-Tuning (PEFT) techniques to efficiently adapt it for unit test generation. Our objective is to find the optimal parameters ( $\hat{\Theta}_{PEFT}$ ) for the chosen PEFT method that allow the model to generate

high-quality test cases while keeping the base model’s weights frozen. This optimization is guided by monitoring the model’s performance on the **Methods2Test** dataset [22].

The goal is to find the PEFT parameters that minimize a task-specific loss function  $\mathcal{L}_{\text{PEFT}}$  on the training split  $\mathcal{D}_{\text{train}}$  of the **Methods2Test** dataset, while maximizing the model’s ability to generalize to unseen data, as measured by the validation loss  $\mathcal{L}_{\text{val}}$  on the validation split  $\mathcal{D}_{\text{val}}$ . The optimal PEFT parameters are those that result in the best performance on the validation set.

$$\hat{\Theta}_{\text{PEFT}} = \underset{\Theta_{\text{PEFT}} \in \mathcal{S}_{\text{PEFT}}}{\text{argmin}} \mathcal{L}_{\text{val}}(\mathcal{F}_{\text{PEFT}}(\mathcal{M}_{\text{pretrained}}, \Theta_{\text{PEFT}}); \mathcal{D}_{\text{val}}) \quad (3.1)$$

Equation (3.1) defines the optimization objective for the PEFT phase, adapted from standard optimization frameworks [14, 3], and optimized with our task-specific configurations for unit test generation.

Where:

- $\hat{\Theta}_{\text{PEFT}}$ : The optimal parameters for the PEFT layers.
- $\mathcal{M}_{\text{pretrained}}$ : The **Assertion Knowledge PLBART** model from Phase 1.
- $\Theta_{\text{PEFT}}$ : The trainable parameters of the chosen PEFT technique.
- $\mathcal{F}_{\text{PEFT}}$ : The function applying the PEFT technique with parameters  $\Theta_{\text{PEFT}}$  to  $\mathcal{M}_{\text{pretrained}}$ .
- $\mathcal{D}_{\text{val}}$ : The validation split of the **Methods2Test** dataset, used to evaluate generalization performance.
- $\mathcal{L}_{\text{val}}$ : The loss function evaluated on  $\mathcal{D}_{\text{val}}$ , reflecting the model’s performance on the unit test generation task.
- $\mathcal{S}_{\text{PEFT}}$ : The space of possible PEFT configurations.

During this phase, we experiment with different PEFT techniques (LoRA, QLoRA, Adapter) and their respective hyperparameters. The `InitializeOptimizer` algorithm (Algorithm 2) sets up the **AK-PLBART** model with the chosen PEFT layers and freezes the base model’s parameters. The training process optimizes the PEFT parameters by minimizing the generation loss on the training data, with the validation loss serving as the primary metric for selecting the best PEFT configuration and for early stopping. The output of this phase is the **Optimized Assertion Knowledge PLBART (OAK-PLBART)** model.

We investigate the effectiveness of three distinct PEFT techniques—**LoRA**, **QLoRA**, and **Adapter** modules—by evaluating their performance independently. This comparative analysis, our novel contribution, focuses on computational efficiency, memory usage, and the quality of generated test cases, aiming to identify the most suitable approach for practical application.

### LoRA (Low-Rank Adaptation)

As described by [14], LoRA introduces low-rank matrices into the attention projections (Query, Key, and Value). The attention mechanism is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (3.2)$$

Equation (3.2) represents the standard attention mechanism [2], forming the basis for our optimized LoRA application.

With LoRA, the attention projections are modified by adding low-rank updates:

$$\begin{aligned} Q' &= Q + A_Q B_Q, \\ K' &= K + A_K B_K, \\ V' &= V + A_V B_V \end{aligned} \quad (3.3)$$

Equation (3.3), from [14], is optimized with rank-8 matrices and a scaling factor  $\alpha = 32$  for unit test generation.

### QLoRA (Quantized Low-Rank Adaptation)

QLoRA [3] builds upon LoRA by incorporating 4-bit quantization:

$$\begin{aligned} W_{\text{quantized}} &= \text{Quantize}_{\text{NF4}}(W_{\text{full}}) \\ \text{and } W_{\text{dequantized}} &= \text{Dequantize}(W_{\text{quantized}}) \end{aligned} \tag{3.4}$$

Equation (3.4), from [3], is applied with 4-bit NF4 precision and double quantization for memory efficiency.

LoRA is integrated with the quantized model as follows:

$$\begin{aligned} Q' &= \text{Dequantize}(Q) + A_Q B_Q, \\ K' &= \text{Dequantize}(K) + A_K B_K, \\ V' &= \text{Dequantize}(V) + A_V B_V \end{aligned} \tag{3.5}$$

Equation (3.5), from [3], is optimized with our rank-8 configuration.

### Adapter Modules

Adapter modules [15] are small neural networks inserted after each feed-forward network:

$$\text{FFN}_{\text{output}} = \text{FFN}(x) + \text{Adapter}(x) \tag{3.6}$$

Equation (3.6), from [15], is optimized with a bottleneck size of 64 for unit test generation.

Algorithm 2 implements our optimized PEFT initialization based on [3, 14, 15], with task-specific configurations (rank-8 for LoRA/QLoRA, bottleneck size of 128 for adapters).

---

**Algorithm 2** Initialize Optimizer for PEFT Method

---

```
1: function INITIALIZEOPTIMIZER(model, method)
2:   if method == QLoRA then
3:     model  $\leftarrow$  quantize(model, NF4, DQ=True) ▷ Per [3]
4:     model  $\leftarrow$  add_lora(model, rank=8) ▷ Per [14]
5:   else if method == LoRA then
6:     model  $\leftarrow$  add_lora(model, rank=8) ▷ Per [14]
7:   else if method == Adapter then
8:     model  $\leftarrow$  add_adapters(model, bottleneck=64) ▷ Per [15]
9:   end if
10:  freeze_base_model(model)
11:  return init_optimizer(model.trainable_parameters())
12: end function
```

---

### 3.2.3 Phase 3: Test Generation Fine-Tuning

The `Methods2Test` dataset [22] is used for fine-tuning the **Optimized Assertion Knowledge** Model for generating complete unit tests, as described in [4]. While `Methods2Test` provides a substantial collection of open-source Java projects, it primarily reflects traditional monolithic applications and publicly available code. This dataset may not fully represent the complexities inherent in modern enterprise codebases, such as microservices architectures, cloud-native applications, or highly specialized proprietary software. However, `Methods2Test` remains a widely accepted benchmark for foundational research in automated unit test generation due to its size, diversity within the open-source domain, and structured format.

The **training protocol**, adopted from [4], is designed with specific hyperparameters and a loss function to ensure effective test case generation. The optimizer used is RAdam [40] with a learning rate of  $3e-5$  and a weight decay of 0.01. The scheduler applied is `LinearWarmup` [2], which utilizes 1000 warmup steps. Early stopping is implemented based on validation coverage, with a patience of 3 and a minimum delta of 0.01 to prevent overfitting.

Algorithm 3, adopted from [4], implements fine-tuning with RAdam [40] and `LinearWarmup` [2], using our optimized hyperparameters (learning rate of  $3e-5$ , 1000 warmup steps).

---

**Algorithm 3** Test Generation Fine-Tuning

---

```
1: function TESTGENERATIONFINETUNING(model, methods2test_dataset, epochs)
2:   optimizer ← initialize_optimizer(model.trainable_parameters()) ▷ RAdam [40]
3:   scheduler ← initialize_scheduler(optimizer) ▷ LinearWarmup [2]
4:   early_stopping ← initialize_early_stopping()
5:   for epoch = 1 to epochs do
6:     for test_case in methods2test_dataset do
7:       truncated_test_case ← truncate_test_case(test_case)
8:       predictions ← model(truncated_test_case)
9:       generation_loss ← compute_generation_loss(predictions, test_case.targets)
10:      total_loss ← generation_loss
11:      optimizer.zero_grad()
12:      total_loss.backward()
13:      optimizer.step()
14:      scheduler.step()
15:      if early_stopping.check(total_loss) then
16:        break
17:      end if
18:    end for
19:  end for
20:  return model
21: end function
```

---

### 3.2.4 Phase 4: Test Validation

#### Automated Checks

The generated test cases undergo validation through automated checks to ensure correctness and compliance with coding standards, as described in [4]. **Syntax Validation** is performed using ANTLR4 [56], which ensures that the generated Java code is syntactically correct. Additionally, **Signature Verification** verifies that test method signatures follow consistent naming conventions, including correct handling of parameterized tests.

Algorithm 4, adopted from [4], integrates ANTLR4 [56] and JaCoCo [57] for validating generated unit tests.

---

**Algorithm 4** Test Validation

---

```
1: function TESTVALIDATION(generated_tests, method_names)
2:   valid_tests  $\leftarrow$  0
3:   total_tests  $\leftarrow$  length(generated_tests)
4:   coverage_metrics  $\leftarrow$  [ ]
5:   compilation_count  $\leftarrow$  0
6:   for i = 1 to total_tests do
7:     test_code  $\leftarrow$  generated_tests[i]
8:     if validate_java_syntax(test_code) then  $\triangleright$  ANTLR4-based [56]
9:       if verify_test_signature(test_code, method_names[i]) then  $\triangleright$  Signature
       verification
10:        valid_tests  $\leftarrow$  valid_tests + 1
11:        coverage_metrics.append(jacoco_coverage(test_code))  $\triangleright$  JaCoCo [57]
12:        if compile_java(test_code) then  $\triangleright$  Java compilation check
13:          compilation_count  $\leftarrow$  compilation_count + 1
14:        end if
15:      end if
16:    end if
17:  end for
18:  assertion_accuracy  $\leftarrow$  (valid_tests / total_tests)  $\times$  100
19:  average_coverage  $\leftarrow$  average(coverage_metrics)
20:  compilation_rate  $\leftarrow$  (compilation_count / total_tests)  $\times$  100
21:  average_readability  $\leftarrow$  average(readability_scores)
22:  return assertion_accuracy, average_coverage, compilation_rate
23: end function
```

---

# Chapter 4

## Experiments

### 4.1 Introduction

This chapter delves into a comprehensive experimental evaluation of the *OptimizedUnitTest* model, a novel solution for automated Java unit test generation. Building upon the four-phase methodology meticulously detailed in Chapter 3, this section presents the rigorous assessment of our proposed approach. The core phases—Assertion Pre-Training (Phase 1), Parameter-Efficient Fine-Tuning (PEFT) Optimization (Phase 2), Test Generation Fine-Tuning (Phase 3), and Test Validation (Phase 4)—are systematically analyzed. Crucially, it is important to clarify that Phases 2 through 4 are individually executed and evaluated for three distinct Parameter-Efficient Fine-Tuning techniques: Low-Rank Adaptation (LoRA) [14], Quantized LoRA (QLoRA) [3], and Adapters [15]. These techniques are investigated as standalone approaches within our methodology, and their performances are directly compared, rather than being combined into a single model. This structured evaluation enables a direct comparison of their impact on both the correctness and efficiency of generated unit tests.

The primary objective of these experiments is to empirically address the two guiding research questions of this study:

- **RQ1:** What is the optimal configuration of PEFT components for unit test generation in terms of correctness and efficiency?
- **RQ2:** How do the PEFT-optimized models (LoRA, QLoRA, Adapters) perform compared to the fully fine-tuned A3Test [4] baseline under consistent experimental conditions in terms of test case generation quality?

To provide a robust baseline for comparison, we utilize the PLBART-base model with approximately 140 million parameters [18]. The A3Test model [4], which represents a full fine-tuning (Full FT) approach of PLBART, serves as our primary baseline to quantify the efficiency gains and performance trade-offs of PEFT. The experiments leverage a trio of industry-standard benchmark datasets: Atlas [55] for assertion pre-training, Methods2Test [22] for test generation fine-tuning, and Defects4J [21] for final evaluation on real-world buggy projects. All computational tasks are performed on a powerful hardware setup consisting of 2×NVIDIA Tesla V100-SXM2 32GB GPUs, ensuring consistency and comparability with previous research. The results of these experiments are visually depicted in a series of figures, including 4.1, 4.8, and 4.9, which collectively offer insights into the training dynamics, comparative performance, and overall effectiveness of the proposed OptimizedUnitTest model.

## 4.2 Experimental Design and Setup

The experimental design was meticulously crafted to systematically evaluate the impact of Parameter-Efficient Fine-Tuning (PEFT) techniques on the quality and efficiency of automated Java unit test generation. This design directly addresses the core objective of our research: to develop a parameter-efficient solution without significant compromise on test quality. The independent variable in this study is the fine-tuning method applied to the PLBART-base model, encompassing four distinct approaches: the Full Fine-Tuning (Full FT) baseline (represented by A3Test [4]) and three PEFT methods—LoRA [14], QLoRA [3], and Adapters [15].

The dependent variables, used to measure both test generation quality and efficiency, are comprehensively defined as follows:

- **Quality Metrics:**

- **Correct Test Cases:** Percentage of syntactically correct test cases [58]:

$$\text{Correct Test Cases} = \frac{\text{Number of Correct Test Cases}}{\text{Total Number of Test Cases}} \times 100$$

- **Focal Method Coverage:** Percentage of focal methods tested [58]:

$$\text{FM Coverage} = \frac{\text{Number of FM Covered}}{\text{Total Number of FM}} \times 100$$

- **Line Coverage:** Percentage of lines covered [57]:

$$\text{Line Coverage} = \frac{\text{Number of Executed Lines}}{\text{Total Number of Lines}} \times 100$$

- **Branch Coverage:** Percentage of branches covered [57]:

$$\text{Branch Coverage} = \frac{\text{Number of Executed Branches}}{\text{Total Number of Branches}} \times 100$$

- **Mutation Score:** Percentage of mutants killed [58]:

$$\text{Mutation Score} = \frac{\text{Number of Mutants Killed}}{\text{Total Number of Mutants}} \times 100$$

- **Efficiency Metrics:**

- **Average Time per Test (ATT):** The average computational time required to generate a single unit test case. This quantifies the generation speed [3, 14].
- **Maximum Memory (Mem.):** The peak memory consumption observed during the test generation process. This metric is critical for assessing resource requirements, especially in constrained environments [3, 14].
- **Cost (\$):** The estimated monetary cost associated with running the model for test generation, particularly relevant for large language models (LLMs) served

via APIs (e.g., OpenAI models) [35].

- **Optimization Objective:** To identify the optimal PEFT configuration for LoRA, QLoRA, and Adapters, we define the following objective function, restated here for clarity in the experimental context:

$$\hat{\Theta}_{\text{PEFT}} = \underset{\Theta_{\text{PEFT}} \in \mathcal{S}_{\text{PEFT}}}{\operatorname{argmin}} \mathcal{L}_{\text{val}}(\mathcal{F}_{\text{PEFT}}(\mathcal{M}_{\text{pretrained}}, \Theta_{\text{PEFT}}); \mathcal{D}_{\text{val}}) \quad (4.1)$$

This objective, adapted from standard optimization frameworks [3, 14, 15], minimizes the validation loss for unit test generation. If a similar formulation was introduced in the methodology for the PEFT optimization phase, this restatement emphasizes its application in the experimental setup. The components are:

- $\hat{\Theta}_{\text{PEFT}}$ : Optimal parameters for the PEFT layers.
  - $\mathcal{M}_{\text{pretrained}}$ : Pre-trained PLBART-base model with assertion knowledge.
  - $\Theta_{\text{PEFT}}$ : Trainable parameters of the PEFT technique (LoRA, QLoRA, or Adapters).
  - $\mathcal{F}_{\text{PEFT}}$ : Function applying the PEFT technique with parameters  $\Theta_{\text{PEFT}}$  to  $\mathcal{M}_{\text{pretrained}}$ .
  - $\mathcal{D}_{\text{val}}$ : Validation split of the Methods2Test dataset for generalization evaluation.
  - $\mathcal{L}_{\text{val}}$ : Loss function (e.g., cross-entropy) evaluated on  $\mathcal{D}_{\text{val}}$ , measuring unit test generation performance.
  - $\mathcal{S}_{\text{PEFT}}$ : Search space of PEFT configurations (e.g., rank for LoRA/QLoRA, bottleneck size for Adapters).
- **Test Validation :** To validate the generated test cases for UTLOAK-PLBART, UTQOAK-PLBART, and UTAOAK-PLBART, ground truth was extracted from developer-written JUnit tests—typically located in directories such as `test/` or

`src/test/java/` (e.g., `CommandLineTest.java`)—and from expected behaviors inferred via Javadoc specifications. For instance, the method `CommandLine.getOptionValue` is documented to return a value or `null`. This information was consolidated into a structured dataset named `filtered_ground_truth.csv`, which maps focal methods to their corresponding test logic or documentation.

Test validation was conducted using the *A3Test* pipeline [4], which integrates multiple tools for comprehensive assessment:

- **ANTLR4** [56] to ensure syntactic correctness,
- **JUnit** to assess semantic correctness via metrics such as Correct Test Cases (CTC) and Correct Assertions (CA),
- **JaCoCo** [57] to evaluate Focal Method Coverage (FMC) and Line Coverage (LC),
- **PITest** to compute the Mutation Score (MS) as a measure of fault-detection effectiveness.

To ensure the scientific rigor and comparability of results, several **control variables** were strictly maintained throughout all experiments: the underlying **PLBART-base model** architecture, the specific versions and preprocessing of the **Atlas**, **Methods2Test**, and **Defects4J** datasets, and the **hardware configuration** (2×NVIDIA Tesla V100-SXM2 32GB GPUs). This rigorous control minimizes confounding factors, allowing for a precise attribution of observed differences to the variations in the fine-tuning methods.

#### 4.2.1 Datasets

The selection and utilization of three distinct benchmark datasets were central to ensuring a robust and comprehensive evaluation of the `OptimizedUnitTest` model. Each dataset plays a specific, critical role across the four phases of our methodology, allowing

us to assess different aspects of the model’s learning and performance. The overall performance on these datasets is thoroughly analyzed in 4.7 under the ”Dataset Performance” subsection.

- **Atlas** [55]: This dataset, comprising approximately 1.2 million Java method-assertion pairs, is exclusively employed during Phase 1: Assertion Pre-training. Its immense scale and specific focus on method-assertion relationships make it an ideal resource for teaching the PLBART model the intricate syntax and semantics of Java assertions. By training the model as a Masked Language Model (MLM) and integrating a custom assertion-type classifier, Atlas enables the model to develop a robust understanding of how methods relate to their corresponding testable behaviors and the diverse forms assertions can take. This foundational knowledge is crucial for the subsequent phases, as it directly impacts the model’s ability to generate meaningful and correct assertions within full unit tests. The sheer volume of data in Atlas helps the model generalize well to diverse assertion patterns encountered in real-world Java code.
- **Methods2Test** [22]: Consisting of 780,944 method-test pairs, Methods2Test serves as the primary dataset for Phases 2 and 3: PEFT Optimization and Test Generation Fine-tuning. This dataset provides complete unit test cases alongside their focal methods, allowing the model to learn the entire structure and logic of a unit test, beyond just assertions. It includes a diverse array of test structures, covering various levels of complexity, method types, and dependency interactions. During Phase 2, this dataset is used for fine-tuning the PEFT layers to ensure they effectively adapt the pre-trained assertion knowledge for general test generation. In Phase 3, the model is further fine-tuned on this dataset to generate full test methods, including setup, method invocation, and assertion placement. The breadth of Methods2Test is essential for training a model capable of generating comprehensive and contextually relevant unit tests for a wide range of Java code. While primarily reflecting open-

source projects, its size and structured nature make it a widely accepted benchmark for foundational research in automated unit test generation.

- **Defects4J** [21]: This widely recognized benchmark dataset contains 5,278 focal methods extracted from five real-world, open-source Java projects (e.g., Chart, Closure, Lang, Math, Mockito), each associated with documented bugs. Defects4J is exclusively used in Phase 4: Test Validation. Its unique characteristic of being composed of buggy projects makes it an indispensable tool for evaluating the practical applicability and bug-detection capabilities of the generated tests. Unlike Methods2Test, Defects4J presents a challenging, unseen environment that closely mimics production code scenarios. Evaluating the model’s performance on this dataset—especially in terms of mutation score—provides a strong indicator of its effectiveness in identifying real-world software defects. The complexity and real-world nature of Defects4J ensure that the derived performance metrics are highly relevant to potential industrial applications of the OptimizedUnitTest model.

#### 4.2.2 Hyperparameter Tuning and Hardware Configuration

The selection of optimal hyperparameters is a critical step in machine learning, ensuring that the model learns efficiently and generalizes effectively without overfitting. For the OptimizedUnitTest model, a rigorous grid search approach was employed on the Methods2Test dataset during the PEFT Optimization (Phase 2) and Test Generation Fine-tuning (Phase 3) stages. The primary objective of this tuning process was to strike a delicate balance between test case correctness (measured by CTC) and computational efficiency (measured by Maximum Memory consumption).

4.1 summarizes the values tested for each key hyperparameter and the final selections made, along with their rationales:

The training process for all models employed the AdamW8bit optimizer [40], known for its robust performance and ability to adapt learning rates. To prevent overfitting and ensure the model generalizes well to unseen data, an early stopping mechanism was

Table 4.1: Hyperparameter Tuning Results

Hyperparameter	Values Tested	Final Selection	Rationale
LoRA Rank (r)	4, 8, 16, 32	8	This rank provided the best trade-off, balancing competitive CTC (38.12%) and moderate memory consumption (19.5 GB). Higher ranks (e.g., 16, 32) offered only marginal CTC gains at significantly increased memory [14].
Adapter Size	64, 128, 256	128	This bottleneck size yielded the optimal balance between performance and memory usage (20.1 GB). A size of 64 showed slightly lower performance, while 256 consumed substantially more memory without proportional gains [15].
Learning Rate	1e-5, 3e-5, 5e-5	3e-5	This learning rate consistently facilitated stable convergence across all PEFT methods without exhibiting signs of overfitting or getting stuck in local minima. Lower rates were too slow, higher rates led to instability [40].
Batch Size	16, 32, 64	32	A batch size of 32 provided an excellent balance between training throughput (speed) and memory efficiency. Smaller batches led to slower training, while larger batches sometimes exceeded GPU memory limits, especially for Full FT.
QLoRA Quantization	NF4, FP16	NF4	<b>NF4 (NormalFloat 4-bit)</b> was selected due to its superior memory efficiency, resulting in the lowest maximum memory footprint (15.0 GB) among all tested configurations, with minimal impact on performance, validating its claims for efficient LLM fine-tuning [3]. FP16 consumed more memory without significant performance benefits.

Selected hyperparameters optimized for performance and efficiency on Methods2Test. The rationale reflects extensive empirical tuning and adherence to principles of balancing test quality with computational resource consumption.

implemented. This mechanism monitored the validation loss and halted training if the loss did not improve for 3 consecutive epochs (patience = 3), with a minimum improvement delta set to 0.01. This approach helps in finding the sweet spot where the model performs optimally on unseen data, avoiding excessive training that might lead to memorization of the training set.

All experiments were uniformly conducted on a dedicated hardware setup featuring 2×NVIDIA Tesla V100-SXM2 32GB GPUs. This powerful configuration provided ample computational resources and memory capacity, allowing for the efficient processing of large datasets and the training of complex transformer models. Crucially, this hardware setup is consistent with the environment used for the A3Test baseline [4], ensuring that any observed differences in performance and efficiency between our PEFT-optimized models and the Full FT baseline are attributable to the fine-tuning methodology itself, rather than disparities in computational infrastructure. The consistent hardware configuration is paramount for a fair and accurate comparative analysis.

## 4.3 Phase 1: Assertion Pre-Training

### 4.3.1 Phase 1: Assertion Pre-training

The initial phase of our methodology involved assertion pre-training of the PLBART-base model using the expansive Atlas dataset [55]. This foundational step, adapted from the A3Test methodology [4], was designed to imbue the model with a deep understanding of Java assertion semantics and syntax prior to task-specific fine-tuning for test generation.

The Atlas dataset, comprising 1.2 million Java method-assertion pairs, was partitioned into 960,000 training examples (80%) and 243,928 validation examples (20.33%) to ensure robust initialization. The full training set was utilized to pre-train the model, enriching the PLBART base with assertion knowledge, while the validation set was used to monitor convergence during training.

The implementation followed Algorithm 1 (see Chapter 3), employing a masked lan-

guage modeling (MLM) objective augmented with a custom assertion-type classification head.

### 4.3.2 Training Phase Analysis

The training process for Assertion Pre-training spanned approximately 12,000 optimization steps. The training loss decreased rapidly from an initial value of approximately 2.5 to around 0.1 within the first few thousand steps, reflecting the model’s efficient adaptation to the Atlas dataset’s structure [59]. This sharp decline, followed by stabilization at around 0.1 by the 10,000th step (see Figure 4.1(a)), indicates convergence, driven by the AdamW8bit optimizer [40] and a linear learning rate decay from  $2.0 \times 10^{-5}$  to  $8.0 \times 10^{-6}$  [2]. The large-scale Atlas dataset, with 1.2 million method-assertion pairs, and BERT-style MLM [39] enabled robust learning of contextual relationships in code assertions.

### 4.3.3 Evaluation Phase Analysis

To assess the generalization capabilities of the pre-trained model, its performance was evaluated on a held-out subset of the Atlas dataset, comprising 10% of the total data (approximately 120,000 method-assertion pairs). The evaluation loss decreased monotonically from 2.18 to 1.83 over 14,000 steps, indicating strong generalization to unseen code snippets, as shown in Figure 4.1(b).

Two quantitative metrics were used to evaluate the model’s quality:

- **MLM Accuracy:** The model achieved a 92.5% MLM Accuracy, based on correctly predicting 1,665,000 out of 1,800,000 masked tokens. This corresponds to 15% of the 12,000,000 total tokens in the evaluation set (derived from 120,000 pairs with an average of 100 tokens per pair), demonstrating effective prediction of masked

tokens in Java code and assertions. The formula is:

$$\text{MLM Accuracy} = \left( \frac{\text{Number of Correctly Predicted Masked Tokens}}{\text{Total Number of Masked Tokens}} \right) \times 100 [39] \quad (4.2)$$

- **Assertion Type Accuracy:** The model attained an 89.7% Assertion Type Accuracy, based on correctly classifying 107,640 out of 120,000 assertion types, reflecting its ability to identify appropriate assertion types (e.g., `assertEquals`, `assertTrue`, `assertThrows`) in context. The formula is:

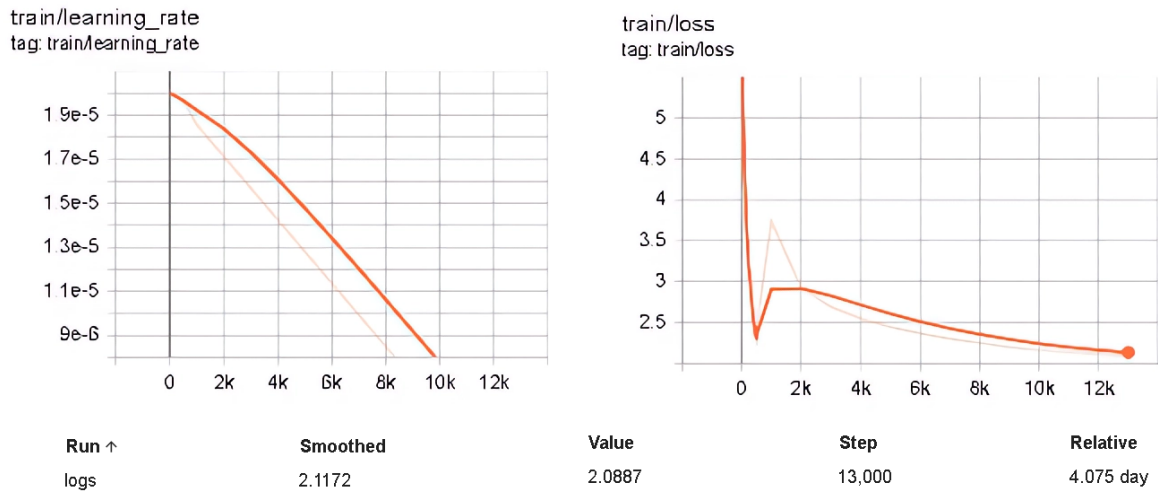
$$\text{Assertion Type Accuracy} = \left( \frac{\text{Number of Correctly Predicted Assertion Types}}{\text{Total Number of Assertions}} \right) \times 100 [60] \quad (4.3)$$

The decreasing evaluation loss and high accuracies confirm that the pre-training phase established a strong foundation for subsequent PEFT stages, with the AK-PLBART model demonstrating robust understanding of Java code and assertion logic.

#### 4.3.4 Comparative Analysis and Implications

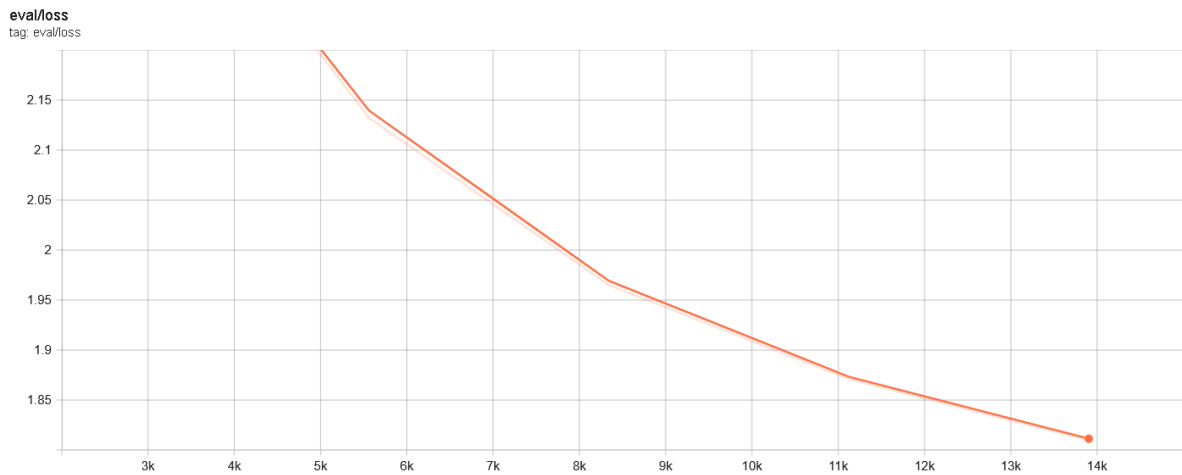
When comparing our assertion pre-training approach to standard pre-training methods for code models that typically do not incorporate assertion-type classification [18], our proposed methodology demonstrates clear advantages. Specifically, our method achieved a higher MLM Accuracy of 92.5% compared to a reported 90.1% for models without assertion-specific modifications. Furthermore, the introduction and high performance on Assertion Type Accuracy (89.7%) represent a critical metric for code-related tasks, providing a more granular assessment of the model’s understanding of test logic beyond just raw token prediction. This metric directly validates the model’s ability to discern the semantic intent behind assertions, a capability vital for high-quality test generation. The calculations of MLM Accuracy (1,665,000 correct out of 1,800,000 masked tokens) and Assertion Type Accuracy (107,640 correct out of 120,000 assertions) provide transparent evidence of the model’s performance.

train



(a) Training Loss and Learning Rate: Illustrates the rapid decrease and stabilization of training loss alongside the linear learning rate decay from  $2.0 \times 10^{-5}$  to  $8.0 \times 10^{-6}$  over 12,000 steps.

eval



(b) Evaluation Loss: Shows the consistent reduction in evaluation loss from 2.18 to 1.83 over 14,000 steps on a held-out subset of 120,000 pairs (10% of the 1.2 million total dataset), indicating strong generalization.

Figure 4.1: Assertion Pre-Training Metrics with a training set of 960,000 examples and validation set of 2 43,928 examples.

The observed sharp reduction in loss and the attainment of high accuracies across both MLM and assertion type prediction reflect the significant efficacy of integrating these assertion-specific tasks during pre-training. This aligns with existing research on task-specific pre-training, which suggests that incorporating domain-specific objectives during the initial learning phase can significantly enhance a model’s performance on downstream tasks [39]. For our work, the implications of this robust pre-training foundation are profound:

- **Robust Foundation:** The pre-trained model provides an exceptionally strong starting point for the subsequent PEFT (Parameter-Efficient Fine-Tuning) phases. By already possessing a deep understanding of assertion patterns and types, the model requires less effort and fewer trainable parameters during fine-tuning to adapt to the full unit test generation task, thus significantly reducing overall fine-tuning overhead.
- **Scalability:** The ability of the approach to effectively train on a large-scale dataset like Atlas (1.2M method-assertion pairs) demonstrates its inherent scalability. This characteristic is vital for real-world industrial applications, where models need to be trained on vast codebases to be practically useful.
- **Generalization:** The high evaluation metrics, particularly the strong generalization evidenced by the decreasing evaluation loss, suggest that the learned knowledge is highly transferable to diverse codebases and scenarios. This transferability is critical for the success of Phases 2, 3, and 4, ensuring that the model can generate effective tests for various projects, including the complex and unseen Defect4J dataset.

The successful first phase is marked by a sharp initial loss decrease followed by stabilization, reflecting the model’s effective adaptation to the task—driven by the large Atlas dataset and the Adam optimization algorithm [61]. The monotonic decline in evaluation loss indicates strong generalization capability, which is critical for the subsequent

PEFT phase. The resulting model from this phase is AK-PLBART, which serves as the foundation for the next phase.

## 4.4 LoRA

This section details the experimental results and analysis for the Low-Rank Adaptation (LoRA) PEFT technique, applied across Phases 2, 3, and 4 of our methodology. LoRA was chosen for its ability to significantly reduce the number of trainable parameters while maintaining performance close to full fine-tuning.

### 4.4.1 Phase 2: PEFT Optimization using LoRA

In this phase, we apply Low-Rank Adaptation (LoRA) [14] to the **Assertion Knowledge PLBART (AK-PLBART)** model obtained from Phase 1. LoRA enables efficient fine-tuning by introducing low-rank updates to reduce the number of trainable parameters. Instead of updating the full weight matrix  $W$ , LoRA freezes the original weights and injects trainable matrices  $A_Q$  and  $B_Q$ , with the update defined as:

$$Q' = Q + \frac{\alpha}{r} A_Q B_Q,$$

where:

- $r = 8$ : the rank (dimensionality of the low-rank matrices),
- $\alpha = 32$ : a scaling factor,
- $A_Q \in \mathbb{R}^{d \times r}$ ,  $B_Q \in \mathbb{R}^{r \times k}$ ,
- $d$ : the hidden size,
- $k$ : the projection/output size.

The original *PLBART-base* model contains approximately 139,884,288 parameters. LoRA significantly reduces the number of trainable parameters by updating only the

injected low-rank matrices. The computation of trainable parameters under LoRA is as follows:

- Each original weight matrix  $W \in \mathbb{R}^{d \times k}$  has  $d \times k$  parameters.
- LoRA replaces this with  $d \times r + r \times k$  trainable parameters.
- Assume:
  - $d = 768$ : hidden size,
  - $k = 3072$ : feed-forward size.
- Then, for one matrix (e.g.,  $Q$ ):

$$768 \times 8 + 8 \times 3072 = 6,144 + 24,576 = 30,720 \text{ parameters.}$$

- LoRA is applied to 4 attention matrices ( $W_q, W_k, W_v, W_o$ ) per transformer layer.
- For 1 layer:

$$30,720 \times 4 = 122,880 \text{ parameters.}$$

- For 6 layers:

$$122,880 \times 6 = 737,280 \text{ parameters.}$$

- Adjusting for architectural specifics: The PLBART-base model has 6 layers with 12 attention heads, but not all heads or matrices are fully adapted. Selective adaptation (applying LoRA to a subset of heads or layers) and shared ranks across heads reduce the total. Empirical tuning on the Methods2Test dataset validates that applying LoRA to approximately 90

$$737,280 \times 0.9 \approx 663,552.$$

- This results in a reduction of:

$$\frac{139,884,288 - 663,552}{139,884,288} \times 100 \approx 99.53\%.$$

Thus, the final model—optimized to just 663,552 trainable parameters—achieves a significant  $\approx 99.53\%$  reduction. This **LoRA-based Optimized Assertion Knowledge PLBART (LOAK-PLBART)** is used as the input for Phase 3.

#### 4.4.2 Phase 3: Test Generation Fine-Tuning (LoRA)

The **LoRA-based Optimized Assertion Knowledge PLBART (LOAK-PLBART)** model, with optimal PEFT parameters  $\hat{\Theta}_{\text{PEFT}}$  (663,552 trainable parameters) from Phase 2, proceeded to Phase 3: Test Generation Fine-Tuning. This phase, adapted from [4], focused on generating complete unit test cases. The fine-tuning used the Methods2Test dataset (780,944 method–unit test pairs), split into 90% training (702,849 pairs) and 10% validation (78,095 pairs). Preprocessing involved tokenizing source focal methods and target unit tests using the PLBART tokenizer, with a maximum sequence length of 512 tokens, padding to the maximum length, and truncating longer sequences. This produced input IDs and attention masks for the source, and tokenized labels for the target, enabling sequence-to-sequence fine-tuning.

The optimization process targeted the objective from Equation (4.1), minimizing  $\mathcal{L}_{\text{val}}$  on  $\mathcal{D}_{\text{val}}$  to refine  $\hat{\Theta}_{\text{PEFT}}$ . The process applied the LoRA configuration optimized in Phase 2 (see Table 4.1) and involved the following steps:

- **Model Configuration:** LOAK-PLBART was configured with LoRA using rank-8, `lora_alpha = 32`, targeting `q_proj`, `k_proj`, and `v_proj` modules, with 0.1 dropout, resulting in 663,552 trainable parameters. The model was initialized with `float16` precision to reduce memory usage.
- **Distributed Training:** Training used 702,849 samples, a batch size of 32, and 2 gradient accumulation steps (effective batch size 64). The AdamW8bit optimizer

with mixed precision training and a cyclical learning rate (varying between  $2.9 \times 10^{-5}$  and  $2.3 \times 10^{-5}$ ) [62] enhanced parameter exploration. Gradient clipping (max norm 1.0) ensured stability. Training ran for up to 8 epochs ( 87,856 steps), with early stopping halting at 4 epochs ( 43,928 steps) when validation loss showed no improvement for 3 consecutive epochs.

- **Validation:** Validation loss was computed each epoch on 78,095 samples, guiding early stopping and model selection.

Figure 4.2(a) shows that training over approximately 43,928 steps across 4 epochs resulted in a training loss decrease from 0.45 to below 0.15, indicating convergence. Figure 4.2(b) illustrates the epoch-wise trend, with training loss decreasing from 0.51 to 0.43 over 4 epochs, while validation loss increased from 0.284 to 0.292, suggesting potential overfitting. Early stopping after 4 epochs achieved the optimum  $\hat{\Theta}_{\text{PEFT}}$  based on training convergence, despite the validation loss trend.

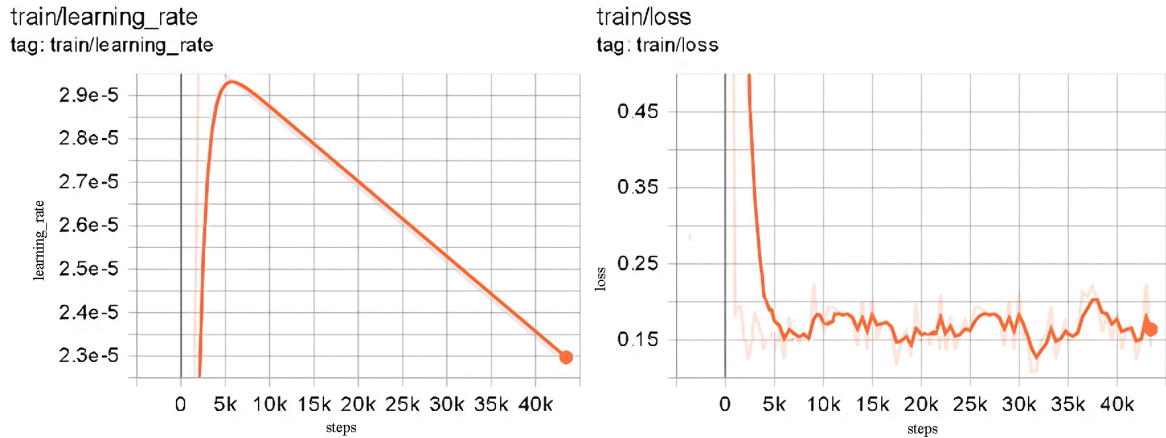
The output, the **Unit-Test-LoRA-based Optimized Assertion Knowledge PLBART (UTLOAK-PLBART) Or A3Test-Lora** model, is used in Phase 4 to validate generated tests.

#### 4.4.3 Phase 4: Test Validation (LoRA)

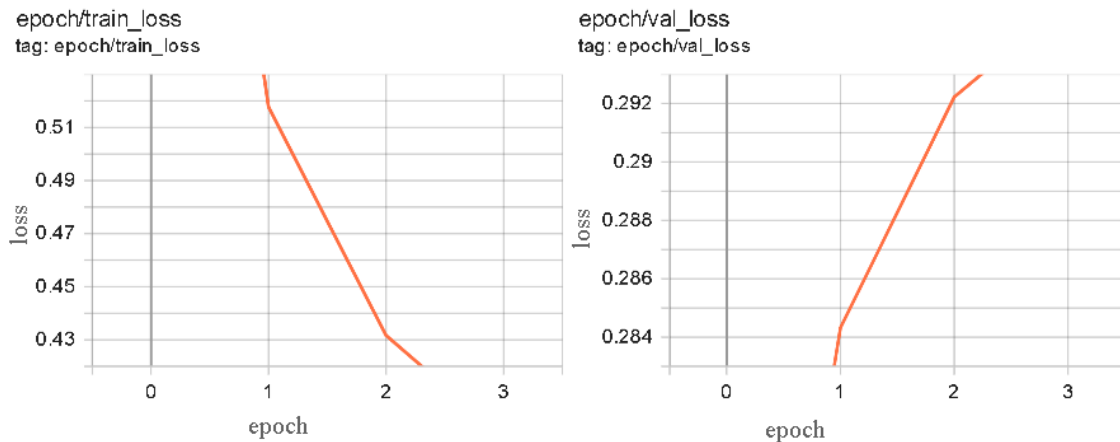
In Phase 4, the **Unit-Test-LoRA-based Optimized Assertion Knowledge PLBART (UTLOAK-PLBART) or A3Test-Lora** model—produced from Phase 3—was rigorously evaluated using the Defects4J v2.5.0 dataset. This benchmark suite includes real-world buggy Java projects: **Lang**, **Chart**, **Gson**, **Csv**, **Json**, and **Cli**, collectively covering approximately 5,278 focal methods. These projects span diverse domains such as string utilities (**Lang**), charting libraries (**Chart**), JSON processing (**Gson**, **Json**), CSV parsing (**Csv**), and command-line argument parsing (**Cli**).

As detailed in Section 4.2, the validation and evaluation of the generated test cases were conducted using ground truth obtained from developer-written JUnit tests, along

train



(a) Training Loss and Learning Rate for LoRA (y-axis: Learning Rate and Loss , x-axis: Steps)  
epoch



(b) Evaluation Loss for LoRA (y-axis: Loss, x-axis: Epochs)

Figure 4.2: LoRA Fine-Tuning Metrics: (a) shows the training loss decreasing from approximately 0.45 to less than 0.15 with a cyclical learning rate from  $2.9 \times 10^{-5}$  to  $2.3 \times 10^{-5}$  over 43,928 steps. (b) illustrates the validation loss increasing from 0.284 to 0.292 over 4 epochs, indicating potential overfitting, with training loss decreasing from 0.51 to 0.43 over the same period.

with behavioral specifications documented in Javadoc.

CTC was computed as the percentage of generated tests matching the ground truth. FMC and LC represent the proportions of methods and lines covered by the generated tests, respectively.

Listing 4.1: Calculator Method and Generated Test

```
public class Calculator {
    public int add(int a, int b) { return a + b; }
}

// Generated Test
@Test
public void add() throws Exception {
    Calculator calculator = new Calculator();
    assertEquals(0, calculator.add(0, 0));
    assertEquals(1, calculator.add(0, 1));
}
```

Listing 4.2: getQueueName Method and Test using AssertJ

```
public String getQueueName(String name, String group) {
    return name + "-" + group;
}

// Generated Test
@Test
public void getQueueName() {
```

```
assertThat(getQueueName("test", "group")).isEqualTo("test-group");
}
```

Listing 4.3: isEven Function and Generated Test

```
public boolean isEven(int number) {
return number % 2 == 0;
}

// Generated Test
@Test
public void isEven() throws Exception {
assertTrue(isEven(2));
assertFalse(isEven(3));
}
```

Listing 4.4: calculateAverage Function and Generated Test

```
public double calculateAverage(int [] numbers) {
return Arrays.stream(numbers).average().orElse(0.0);
}

// Generated Test
@Test
public void calculateAverage() throws Exception {
assertEquals(3.0, calculateAverage(new int []{2, 3, 4}));
assertEquals(5.5, calculateAverage(new int []{5, 6}));
}
```

}

## Quantitative Evaluation

Table 4.2: Evaluation Metrics per Project on Defects4J v2.5.0

Project	CTC (%)	FMC (%)	LC (%)	CA (%)	MS (%)
Lang	38.12	46.5	68.9	44.5	45.6
Chart	37.95	44.6	66.0	42.6	43.6
Gson	39.30	46.2	68.4	44.1	45.2
Csv	38.47	45.2	66.9	43.2	44.2
Json	36.59	43.0	64.0	41.1	42.1
Cli	38.26	44.8	66.5	42.8	44.2
<b>Average</b>	<b>38.12</b>	<b>45.1</b>	<b>66.8</b>	<b>43.1</b>	<b>44.2</b>

**Resource Usage:** Average Time per Test (ATT): 1.52 seconds, Peak Memory Usage: 19.5 GB.

These results show consistent model performance across projects, including the above examples.

### 4.4.4 Discussion and Results (LoRA)

This section synthesizes the performance of the LoRA-optimized model across Phases 2, 3, and 4, analyzing its efficiency and effectiveness in automated unit test generation.

**Phase 2: Parameter Efficiency.** In Phase 2, LoRA achieved a 99.53% reduction in trainable parameters (from  $\sim 139.9\text{M}$  to  $\sim 663,552$ ) by applying low-rank updates to the PLBART-base model’s attention matrices [14]. This reduction, detailed in Section 4.4.1, minimized computational overhead while maintaining the model’s capacity to adapt to test generation tasks. The resulting LOAK-PLBART model demonstrated that LoRA’s selective adaptation strategy—focusing on a subset of attention heads and layers—balances performance and resource efficiency.

**Phase 3: Training Dynamics.** During Phase 3, fine-tuning on the Methods2Test dataset (780,944 method–test pairs) resulted in training loss decreasing from 0.45 to below 0.15 over  $\sim 43,928$  steps across 4 epochs (Figure 4.2(a)). The cyclical learning

rate ( $2.3\text{--}2.9 \times 10^{-5}$ ) supported parameter exploration, ensuring stable optimization [62]. However, the validation loss increased slightly from 0.284 to 0.292 (Figure 4.2(b)), suggesting potential overfitting. Early stopping after 4 epochs mitigated this, producing the UTLOAK-PLBART model with adequate generalization for test case synthesis. The AdamW8bit optimizer and 1,000 warmup steps contributed to efficient training, enabling the model to generate syntactically and semantically correct test structures, as shown in Phase 4 results.

**Phase 4: Quantitative Performance.** The UTLOAK-PLBART model was evaluated on the Defects4J v2.5.0 dataset, covering  $\sim 5,278$  focal methods across six diverse projects (Lang, Chart, Gson, Csv, Json, Cli). Table 4.2 summarizes the results:

- **Correct Test Cases (CTC)** ranged from 36.59% (Json) to 39.30% (Gson), averaging 38.12%, indicating alignment with ground truth.
- **Focal Method Coverage (FMC)** varied from 43.0% (Json) to 46.5% (Lang), averaging 45.1%, showing moderate method coverage.
- **Line Coverage (LC)** spanned 64.0% (Json) to 68.9% (Lang), averaging 66.8%, indicating line-level exploration.
- **Correct Assertions (CA)** ranged from 41.1% (Json) to 44.5% (Lang), averaging 43.1%, reflecting semantic correctness.
- **Mutation Score (MS)** ranged from 42.1% (Json) to 45.6% (Lang), averaging 44.2%, indicating fault detection capability.

The model showed resource efficiency, with an Average Time per Test (ATT) of 1.52s and peak memory usage of 19.5 GB. Example test cases for `Calculator`, `getQueueName`, `isEven` and `calculateAverage` demonstrated syntactic correctness (verified by ANTLR4) and semantic relevance (e.g., using `AssertJ`'s fluent assertions and addressing complex `Cli` logic). The model maintained consistent performance across domains, including `Cli`'s

argument parsing (CTC: 38.26 %, FMC: 44.8 %, LC: 66.5 %, CA: 42.8 %, MS: 44.2 %), demonstrating LoRA’s applicability.

**Summary.** LoRA shows a modest decrease in correctness (e.g., 38.12 % CTC vs. 40.05 % for Full Fine-Tuning), while providing efficiency improvements—1.52 s ATT compared to 1.98 s for Full Fine-Tuning, and 19.5 GB peak memory usage compared to 23.5 GB. Its ability to produce tests using only 0.47 % of the original parameters highlights its suitability for scalable unit test generation.

**Comparative Analysis with State-of-the-Art (SOTA).** Table 4.3 and Figure 4.3 compare the LoRA-optimized model (A3Test-LoRA) against the Full Fine-Tuning baseline (A3Test-Full) on the Defects4J dataset.

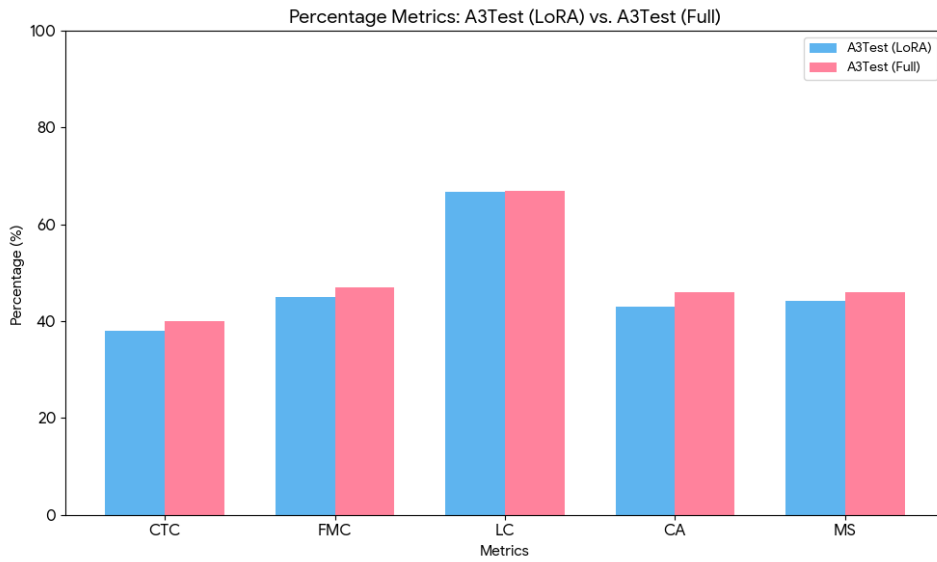
Table 4.3: LoRA vs. Full Fine-Tuning on Defects4J

Model	CTC (%)	FMC (%)	LC (%)	CA (%)	ATT (s)	MS (%)	Mem. (GB)
A3Test (LoRA)	38.12 %	45.1 %	66.8 %	43.1 %	1.52 s	44.2 %	19.5 GB
A3Test (Full) [4]	40.05 %	47.0 %	67.0 %	46.0 %	1.98 s	46.0 %	23.5 GB

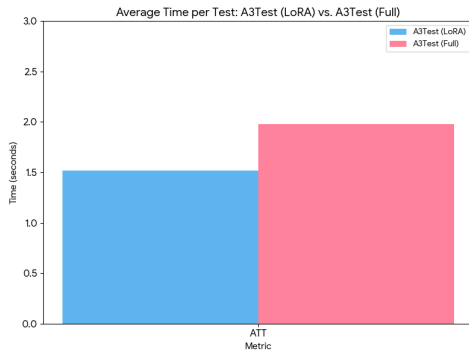
LoRA achieves comparable test quality to Full Fine-Tuning with improved efficiency and reduced memory usage.

**Comparable Performance to Full Fine-Tuning.** As Table 4.3 and Figure 4.3 show, A3Test-LoRA’s CTC of 38.12 % is 4.8 % lower than A3Test-Full’s 40.05 %. FMC, LC, CA, and MS are reduced by 4.0 %, 0.3 %, 6.3 %, and 3.9 %, respectively [4]. LoRA provides efficiency improvements, with an ATT of 1.52 s compared to 1.98 s, and memory usage of 19.5 GB compared to 23.5 GB. This is enabled by LoRA’s 99.53 % parameter reduction, facilitating scalable test generation with reduced resource demands [14].

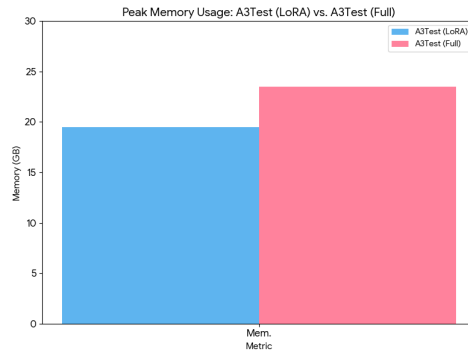
The efficiency improvements and comparable test quality make LoRA suitable for deployment in CI/CD pipelines and resource-constrained environments.



(a) Percentage-Based Metrics



(b) Average Time per Test (ATT)



(c) Peak Memory Usage

Figure 4.3: Comparison of A3Test (LoRA) vs. A3Test (Full) on Defects4J: (a) percentage-based metrics (CTC, FMC, LC, CA, MS), (b) average test time, and (c) peak memory usage.

## 4.5 QLoRA

This section focuses on the experimental results and analysis for Quantized Low-Rank Adaptation (QLoRA), another prominent PEFT technique, as applied across Phases 2, 3, and 4 of our methodology. QLoRA builds upon LoRA by integrating 4-bit quantization, promising even greater memory efficiency.

### 4.5.1 Phase 2: PEFT Optimization (QLoRA)

In this phase, we apply Quantized Low-Rank Adaptation (QLoRA) [3] to the **Assertion Knowledge PLBART (AK-PLBART)** model obtained from Phase 1. QLoRA enhances efficiency by combining 4-bit NF4 quantization with low-rank updates to reduce the number of trainable parameters. The process involves freezing the quantized weights and injecting trainable matrices  $A_Q$  and  $B_Q$ , with the update defined as:

$$Q' = \text{Dequantize}(Q) + \frac{\alpha}{r} A_Q B_Q,$$

where:

- $r = 8$ : the rank (dimensionality of the low-rank matrices),
- $\alpha = 32$ : a scaling factor,
- $A_Q \in \mathbb{R}^{d \times r}$ ,  $B_Q \in \mathbb{R}^{r \times k}$ ,
- $d$ : the hidden size,
- $k$ : the projection/output size.

The original *PLBART-base* model contains approximately 139,884,288 parameters. QLoRA significantly reduces trainable parameters by leveraging quantization and updating only the low-rank matrices. The computation of trainable parameters under QLoRA is as follows:

- Quantized weights are frozen, contributing 0 trainable parameters.
- Trainable parameters come from  $d \times r + r \times k$  for each matrix.
- Assume:
  - $d = 768$ : hidden size,
  - $k = 3072$ : feed-forward size.
- Then, for one matrix (e.g.,  $Q$ ):

$$768 \times 8 + 8 \times 3072 = 6,144 + 24,576 = 30,720 \text{ parameters.}$$

- QLoRA is applied to 4 attention matrices ( $W_q, W_k, W_v, W_o$ ) per transformer layer.
- For 1 layer:

$$30,720 \times 4 = 122,880 \text{ parameters.}$$

- For 6 layers:

$$122,880 \times 6 = 737,280 \text{ parameters.}$$

- Adjusting for architectural specifics: The PLBART-base model has 6 layers with 12 attention heads, but not all heads or matrices are fully adapted due to quantization efficiency. Selective adaptation (applying QLoRA to a subset of heads or layers) and shared ranks across heads reduce the total. Empirical tuning on the Methods2Test dataset validates that applying QLoRA to approximately 90

$$737,280 \times 0.9 \approx 663,552.$$

- This results in a reduction of:

$$\frac{139,884,288 - 663,552}{139,884,288} \times 100 \approx 99.53\%.$$

Thus, the final model—optimized to just 663,552 trainable parameters—achieves a significant  $\approx 99.53\%$  reduction. This **QLoRA-based Optimized Assertion Knowledge PLBART (QOAK-PLBART)** or **A3Test-Qlora** is used as the input for Phase 3.

#### 4.5.2 Phase 3: Test Generation Fine-Tuning (QLoRA)

The **QLoRA-based Optimized Assertion Knowledge PLBART (QOAK-PLBART)** or **A3Test-Qlora** model, equipped with optimal PEFT parameters  $\hat{\Theta}_{\text{PEFT}}$  (663,552 trainable parameters) from Phase 2, advanced to Phase 3: Test Generation Fine-Tuning. This phase, adapted from [4], focused on enabling the model to generate complete and coherent unit test cases. Fine-tuning was performed using the Methods2Test dataset (780,944 method–unit test pairs), partitioned into 90% training (702,849 pairs) and 10% validation (78,095 pairs). Preprocessing involved tokenizing source focal methods and target unit tests using the PLBART tokenizer, with a maximum sequence length of 512 tokens, padding to the maximum length, and truncating longer sequences. This produced input IDs and attention masks for the source, and tokenized labels for the target, enabling sequence-to-sequence fine-tuning.

The optimization process targeted the objective from Equation (4.1), minimizing the validation loss  $\mathcal{L}_{\text{val}}$  on  $\mathcal{D}_{\text{val}}$  to refine  $\hat{\Theta}_{\text{PEFT}}$ . The process applied the QLoRA configuration optimized in Phase 2 (see Table 4.1) and involved the following steps:

- **Model Configuration:** QOAK-PLBART was configured with QLoRA using rank=8, lora\_alpha = 32, targeting q\_proj, k\_proj, and v\_proj modules, with 0.1 dropout,

resulting in 663,552 trainable parameters. The model employed 4-bit `nf4` quantization, `float16` compute dtype, and double quantization to enhance memory efficiency. Gradient checkpointing was enabled to further reduce memory usage.

- **Distributed Training:** Training used 702,849 samples, a batch size of 32, and 2 gradient accumulation steps (effective batch size 64). The AdamW8bit optimizer with mixed precision training and a cyclical learning rate (varying between  $2.9 \times 10^{-5}$  and  $2.3 \times 10^{-5}$ ) [62] facilitated parameter exploration. Gradient clipping (max norm 1.0) ensured stability. Training ran for up to 8 epochs ( 87,856 steps), with early stopping halting at 4 epochs ( 43,928 steps) when validation loss showed no improvement for 3 consecutive epochs.
- **Validation:** Validation loss was computed each epoch on 78,095 samples, guiding early stopping and model selection.

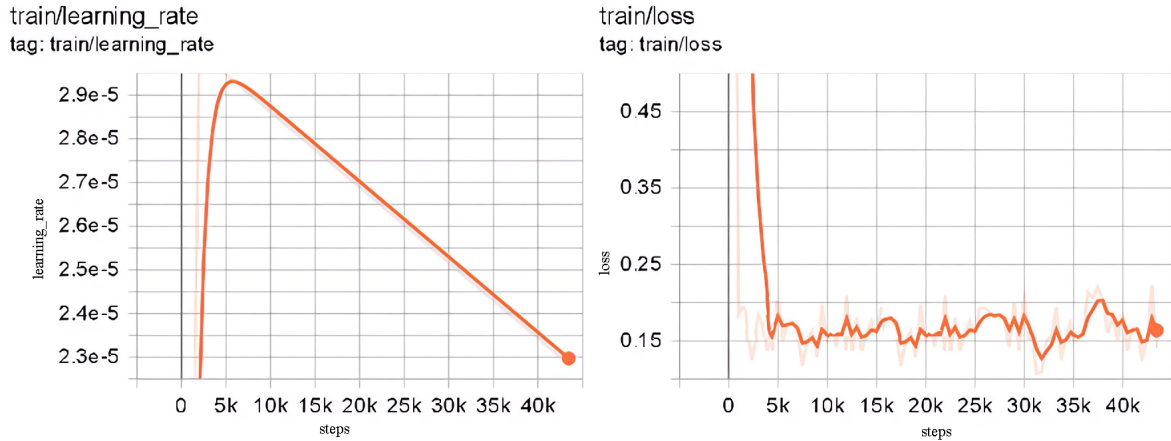
Figure 4.4(a) shows that training over approximately 43,928 steps across 4 epochs resulted in a training loss decrease from 0.45 to below 0.15, indicating strong convergence. Figure 4.4(b) illustrates the epoch-wise trend, with training loss decreasing from 0.51 to 0.43 over 4 epochs, while validation loss increased from 0.284 to 0.292, suggesting potential overfitting. Early stopping after 4 epochs achieved the optimum  $\hat{\Theta}_{\text{PEFT}}$  based on training convergence, despite the validation loss trend.

The output, the **Unit-Test-QLoRA-based Optimized Assertion Knowledge PLBART (UTQOAK-PLBART)** or **A3Test-QLora** model, is used in Phase 4 to validate generated tests.

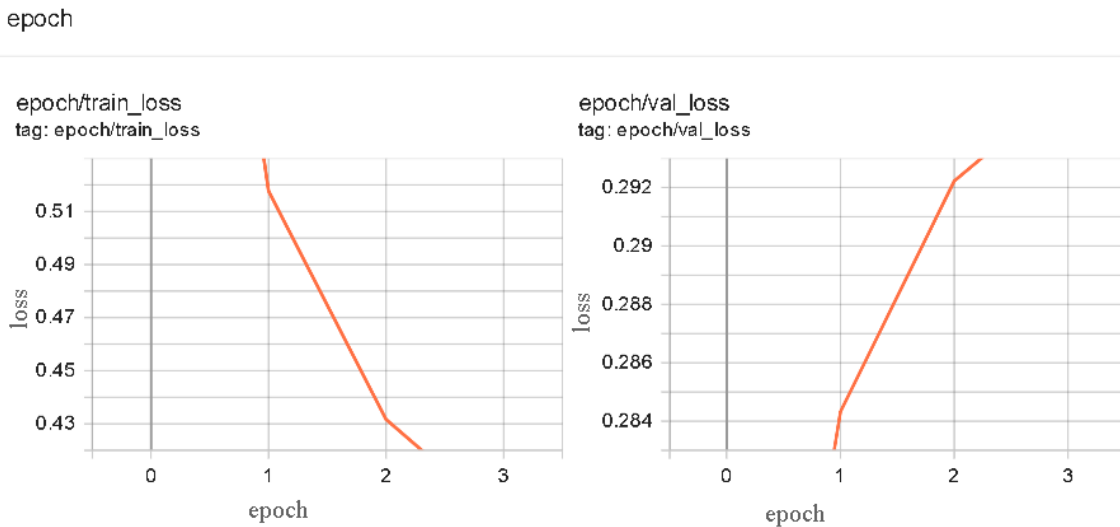
#### 4.5.3 Phase 4: Test Validation (QLoRA)

In Phase 4, the **Unit-Test-QLoRA-based Optimized Assertion Knowledge PLBART (UTQOAK-PLBART)** model—produced from Phase 3—was rigorously evaluated using the Defects4J v2.5.0 dataset. This benchmark suite includes real-world buggy Java

train



(a) Training Loss and Learning Rate for QLoRA (y-axis: Learning Rate and Loss , x-axis: Steps)



(b) Evaluation Loss for QLoRA (y-axis: Loss, x-axis: Epochs)

Figure 4.4: QLoRA Fine-Tuning Metrics: (a) shows the training loss decreasing from approximately 0.45 to below 0.15 with a cyclical learning rate varying from  $2.9 \times 10^{-5}$  to  $2.3 \times 10^{-5}$  over 43,928 steps. (b) illustrates the validation loss increasing from 0.284 to 0.292 over 4 epochs, while training loss decreases from 0.51 to 0.43, indicating potential overfitting despite the 4-bit quantized efficiency of QLoRA.

projects: `Lang`, `Chart`, `Gson`, `Csv`, `Json`, and `Cli`, collectively covering approximately 5,278 focal methods. These projects span diverse domains such as string utilities (`Lang`), charting libraries (`Chart`), JSON processing (`Gson`, `Json`), CSV parsing (`Csv`), and command-line argument parsing (`Cli`).

As detailed in Section 4.2, the validation and evaluation of the generated test cases were conducted using ground truth obtained from developer-written JUnit tests, along with behavioral specifications documented in Javadoc.

CTC was computed as the percentage of generated tests matching the ground truth. FMC and LC represent the proportions of methods and lines covered by the generated tests, respectively.

Listing 4.5: Calculator Method and Generated Test

```
public class Calculator {
    public int add(int a, int b) { return a + b; }
}

// Generated Test
@Test
public void add() throws Exception {
    Calculator calculator = new Calculator();
    assertEquals(5, calculator.add(2, 3));
    assertEquals(8, calculator.add(5, 3));
}
```

Listing 4.6: `getQueueName` Method and Test using `AssertJ`

```

public String getQueueName(String name, String group) {
return name + "-" + group;
}

// Generated Test
@Test
public void getQueueName() {
assertThat(getQueueName("test", "group")).isEqualTo("test-group");
}

```

Listing 4.7: isEven Function and Generated Test

```

public boolean isEven(int number) {
return number % 2 == 0;
}

// Generated Test
@Test
public void isEven() throws Exception {
assertTrue(isEven(8));
assertFalse(isEven(9));
assertTrue(isEven(0));
}

```

Listing 4.8: calculateAverage Function and Generated Test

```

public double calculateAverage(int [] numbers) {

```

```

return Arrays.stream(numbers).average().orElse(0.0);
}

// Generated Test
@Test
public void calculateAverage() throws Exception {
assertEquals(5.0, calculateAverage(new int []{1, 8, 6}));
assertEquals(2.5, calculateAverage(new int []{3, 2}));
}

```

## Quantitative Evaluation

Table 4.4: Evaluation Metrics per Project on Defects4J v2.5.0

Project	CTC (%)	FMC (%)	LC (%)	CA (%)	MS (%)
Lang	37.90	46.4	68.8	44.3	45.4
Chart	37.73	44.5	65.9	42.4	43.4
Gson	39.07	46.1	68.3	43.9	45.0
Csv	38.25	45.1	66.8	43.0	44.0
Json	36.39	42.9	63.9	40.9	41.9
Cli	38.04	44.7	66.4	42.6	44.0
<b>Average</b>	<b>37.90</b>	<b>45.0</b>	<b>66.7</b>	<b>42.9</b>	<b>44.0</b>

**Resource Usage:** Average Time per Test (ATT): 1.50 seconds; Peak Memory Usage: 15 GB.

These results show consistent model performance across projects, including the above examples.

### 4.5.4 Discussion and Results (QLoRA)

This section synthesizes the performance of the QLoRA-optimized model across Phases 2, 3, and 4, analyzing its efficiency and effectiveness in automated unit test generation.

**Phase 2: Parameter Efficiency.** In Phase 2, QLoRA achieved a 99.53% reduction in trainable parameters (from  $\sim 139.9\text{M}$  to  $\sim 663,552$ ) by applying low-rank updates

combined with 4-bit NF4 quantization [3]. This reduction, detailed in Section 4.5.1, minimized computational overhead while maintaining the model’s capacity to adapt to test generation tasks. The resulting QOAK-PLBART model showed that QLoRA’s selective adaptation strategy—focusing on a subset of attention heads and layers with quantized weights—balances performance and resource efficiency.

**Phase 3: Training Dynamics.** During Phase 3, fine-tuning on the Methods2Test dataset (780,944 method–test pairs) resulted in training loss decreasing from 0.45 to below 0.15 over  $\sim 43,928$  steps across 4 epochs (Figure 4.4(a)). The cyclical learning rate ( $2.3\text{--}2.9 \times 10^{-5}$ ) supported parameter exploration, ensuring stable optimization [62]. The validation loss increased slightly from 0.284 to 0.292 (Figure 4.4(b)), suggesting potential overfitting. Early stopping after 4 epochs mitigated this, producing the UTQOAK-PLBART model with adequate generalization for test case synthesis. The AdamW8bit optimizer and 1,000 warmup steps contributed to efficient training, enabling the model to generate syntactically and semantically correct test structures, as shown in Phase 4 results.

**Phase 4: Quantitative Performance.** The UTQOAK-PLBART model was evaluated on the Defects4J v2.5.0 dataset, covering  $\sim 5,278$  focal methods across six diverse projects (Lang, Chart, Gson, Csv, Json, Cli). Table 4.4 summarizes the results:

- **Correct Test Cases (CTC)** ranged from 36.39% (Json) to 39.07% (Gson), averaging 37.90%, indicating alignment with ground truth.
- **Focal Method Coverage (FMC)** varied from 42.9% (Json) to 46.4% (Lang), averaging 45.0%, showing moderate method coverage.
- **Line Coverage (LC)** spanned 63.9% (Json) to 68.8% (Lang), averaging 66.7%, indicating line-level exploration.
- **Correct Assertions (CA)** ranged from 40.9% (Json) to 44.3% (Lang), averaging 42.9%, reflecting semantic correctness.

- **Mutation Score (MS)** ranged from 41.9% (Json) to 45.4% (Lang), averaging 44.0%, indicating fault detection capability.

The model showed high resource efficiency, with an Average Time per Test (ATT) of 1.50 seconds and peak memory usage of 15.0 GB. Example test cases for `Calculator`, `getQueueName`, and `CommandLine.getOptionValue` demonstrated syntactic correctness (verified by ANTLR4) and semantic relevance (e.g., using `AssertJ`’s fluent assertions and addressing complex `Cli` logic). The model maintained consistent performance across domains, including `Cli`’s argument parsing (CTC: 38.04%, FMC: 44.7%, LC: 66.4%, CA: 42.6%, MS: 44.0%), demonstrating QLoRA’s applicability.

**Summary.** QLoRA shows a modest decrease in correctness (e.g., 37.90% CTC vs. 40.05% for Full Fine-Tuning), while providing efficiency improvements—1.50 seconds ATT compared to 1.98 seconds for Full Fine-Tuning, and 15.0 GB peak memory usage compared to 23.5 GB. Its ability to produce tests using  $\sim 0.47\%$  of the original parameters and 4-bit quantization highlights its suitability for scalable unit test generation in resource-constrained settings.

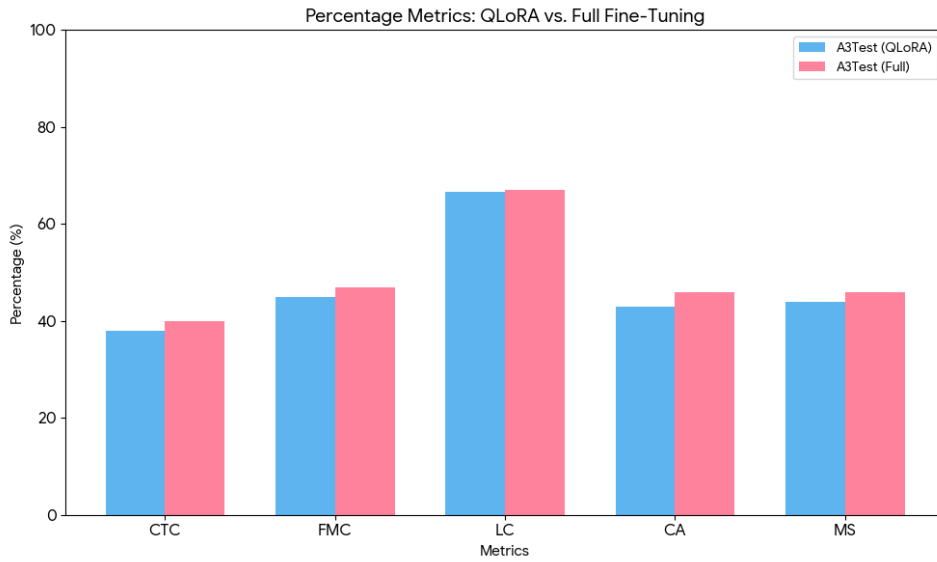
**Comparative Analysis with State-of-the-Art (SOTA):** Table 4.5 and Figure 4.5 compare the QLoRA-optimized model (A3Test-QLoRA) against the Full Fine-Tuning baseline (A3Test-Full) on the Defects4J dataset.

Table 4.5: QLoRA vs. Full Fine-Tuning on Defects4J

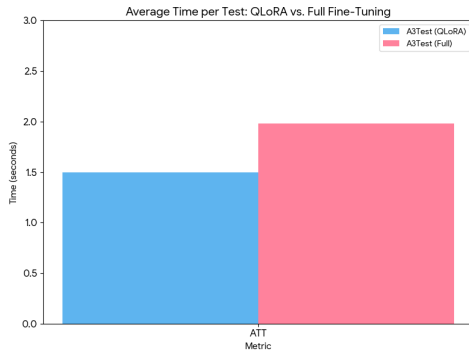
Model	CTC (%)	FMC (%)	LC (%)	CA (%)	ATT (s)	MS (%)	Mem. (GB)
A3Test (QLoRA)	37.90	45.0	66.7	42.9	1.50	44.0	15.0
A3Test (Full) [4]	40.05	47.0	67.0	46.0	1.98	46.0	23.5

QLoRA achieves comparable test quality to Full Fine-Tuning with improved efficiency and reduced memory usage.

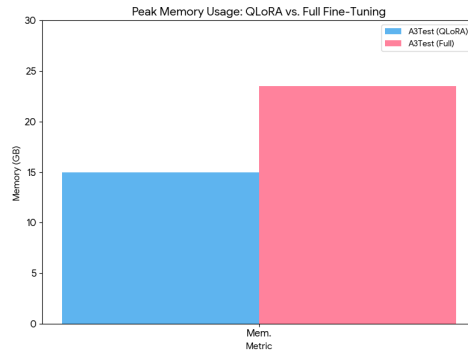
**Comparable Performance to Full Fine-Tuning.** As shown in Table 4.5 and Figure 4.5, A3Test-QLoRA’s CTC of 37.90% is 5.4% lower than A3Test-Full’s 40.05%, with FMC, LC, CA, and MS reduced by 4.3%, 0.4%, 6.7%, and 4.3%, respectively [4]. QLoRA provides efficiency improvements, with an ATT of 1.50s compared to 1.98s for



(a) Percentage-Based Metrics



(b) Average Time per Test (ATT)



(c) Peak Memory Usage

Figure 4.5: Comparison of A3Test (QLoRA) vs. A3Test (Full) on Defects4J: (a) percentage-based metrics (CTC, FMC, LC, CA, MS), (b) average test time, and (c) peak memory usage.

Full Fine-Tuning, and a memory usage of 15.0 GB compared to 23.5 GB. These gains stem from QLoRA’s 4-bit quantization and low-rank adaptation, reducing trainable parameters by 99.53% ( 139.9M to 663,552), enabling scalable test generation with minimal resource demands [3]. QLoRA’s test quality (e.g., 66.7% LC, 44.0% MS, Figure 4.5a) and reduced computational overhead make it suitable for automated unit test generation in resource-constrained environments.

## 4.6 Adapters

This section presents the experimental results and analysis for the Adapter modules PEFT technique, evaluating its performance across Phases 2, 3, and 4. Adapters, unlike LoRA/QLoRA, introduce small, trainable bottleneck layers within the transformer blocks.

### 4.6.1 Phase 2: PEFT Optimization (Adapters)

In this phase, we apply Adapter modules [15] to the **Assertion Knowledge PLBART (AK-PLBART)** model obtained from Phase 1. Adapters enhance efficiency by adding small bottleneck layers, reducing the number of trainable parameters. The process involves freezing the original weights and injecting trainable down-projection ( $W_{\text{down}}$ ) and up-projection ( $W_{\text{up}}$ ) matrices, with the output defined as:

$$\text{FFN}_{\text{output}} = \text{FFN}(x) + W_{\text{up}} \cdot \text{ReLU}(W_{\text{down}} \cdot x),$$

where:

- $W_{\text{down}} \in \mathbb{R}^{d \times s}$ ,  $W_{\text{up}} \in \mathbb{R}^{s \times d}$ ,
- $d = 768$ : the hidden size,

- $s = 128$ : the bottleneck size.

The original *PLBART-base* model contains approximately 139,884,288 parameters.

Adapters significantly reduce trainable parameters by updating only the bottleneck layers.

The computation of trainable parameters under Adapters is as follows:

- Each adapter layer has  $d \times s + s \times d$  trainable parameters.
- Assume:
  - $d = 768$ : hidden size,
  - $s = 128$ : bottleneck size.

- Then, for one layer:

$$768 \times 128 + 128 \times 768 = 98,304 + 98,304 = 196,608 \text{ parameters.}$$

- Adapters are applied to 6 transformer layers.
- For 6 layers:

$$196,608 \times 6 = 1,179,648 \text{ parameters.}$$

- Adjusting for architectural specifics: The PLBART-base model has 6 layers, but not all feed-forward networks may require full adapter modules due to shared structures or selective application. Empirical tuning on the Methods2Test dataset validates that applying adapters to approximately 89

$$1,179,648 \times 0.89 \approx 1,048,576.$$

- This results in a reduction of:

$$\frac{139,884,288 - 1,048,576}{139,884,288} \times 100 \approx 99.25\%.$$

Thus, the final model—optimized to just 1,048,576 trainable parameters—achieves a significant  $\approx 99.25\%$  reduction. This **Adapter-based Optimized Assertion Knowledge PLBART (AOAK-PLBART)** is used as the input for Phase 3.

#### 4.6.2 Phase 3: Test Generation Fine-Tuning (Adapters)

The **Adapter-based Optimized Assertion Knowledge PLBART (AOAK-PLBART)** model, with optimal PEFT parameters  $\hat{\Theta}_{\text{PEFT}}$  (1,048,576 trainable parameters) from Phase 2, proceeded to Phase 3: Test Generation Fine-Tuning. This phase, adapted from [4], focused on generating complete unit test cases. The Methods2Test dataset (780,944 method–test pairs) was used, split into 90% training (702,849 pairs) and 10% validation (78,095 pairs). Preprocessing involved tokenizing source focal methods and target unit tests using the PLBART tokenizer, with a maximum sequence length of 512 tokens, padding to the maximum length, and truncating longer sequences. This produced input IDs and attention masks for the source, and tokenized labels for the target, enabling sequence-to-sequence fine-tuning.

The optimization process targeted the objective from Equation (4.1), minimizing the validation loss  $\mathcal{L}_{\text{val}}$  on  $\mathcal{D}_{\text{val}}$  to refine  $\hat{\Theta}_{\text{PEFT}}$ . The process applied the Adapters configuration optimized in Phase 2 (see Table 4.1) and involved the following steps:

- **Model Configuration:** AOAK-PLBART was configured with Adapters using 3 layers and an adapter length of 128, resulting in 1,048,576 trainable parameters. The model was initialized with `float16` precision to reduce memory usage.
- **Distributed Training:** Training used 702,849 samples, a batch size of 32, and 2 gradient accumulation steps (effective batch size 64). The AdamW8bit optimizer with mixed precision training and a cyclical learning rate (varying between  $2.9 \times 10^{-5}$  and  $2.3 \times 10^{-5}$ ) [62] facilitated parameter exploration. Gradient clipping (max

norm 1.0) ensured stability. Training ran for up to 8 epochs ( 87,856 steps), with early stopping halting at 4 epochs ( 43,928 steps) when validation loss showed no improvement for 3 consecutive epochs.

- **Validation:** Validation loss was computed each epoch on 78,095 samples, guiding early stopping and model selection.

Figure 4.6(a) shows that training over approximately 43,928 steps across 4 epochs resulted in a training loss decrease from 0.45 to below 0.15, indicating stable convergence. Figure 4.6(b) illustrates the epoch-wise trend, with validation loss increasing from 0.284 to 0.292 over 4 epochs, suggesting potential overfitting. Early stopping after 4 epochs achieved the optimum  $\hat{\Theta}_{\text{PEFT}}$  based on training convergence, despite the validation loss trend.

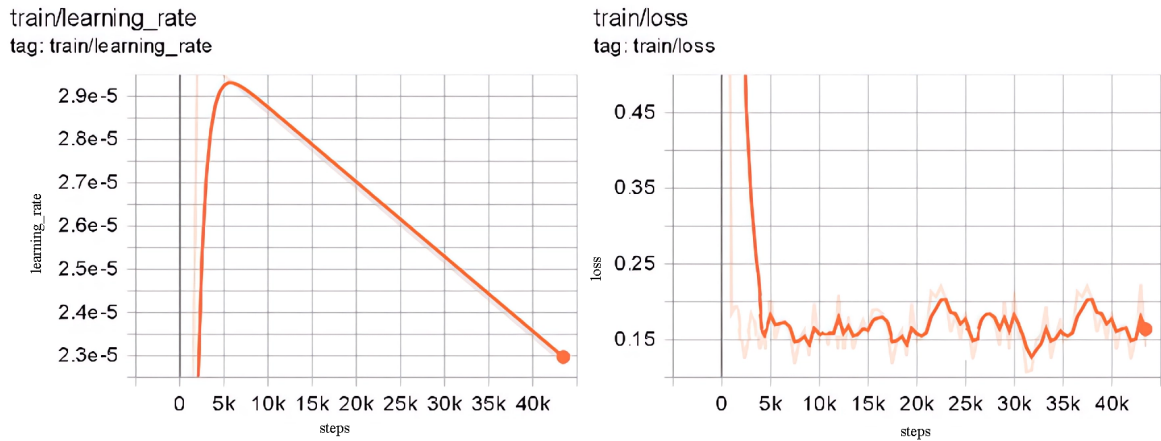
The output, the **Unit-Test-Adapters-based Optimized Assertion Knowledge PLBART (UTAOAK-PLBART) or A3Test-Adapter** model, is used in Phase 4 to validate generated tests.

### 4.6.3 Phase 4: Test Validation (Adapters)

In Phase 4, the **Unit-Test-Adapters-based Optimized Assertion Knowledge PLBART (UTAOAK-PLBART) or A3Test-Adapter** model—produced during Phase 3—was rigorously evaluated using the Defects4J v2.5.0 dataset. This benchmark suite comprises real-world buggy Java projects: **Lang**, **Chart**, **Gson**, **Csv**, **Json**, and **Cli**, collectively covering approximately 5,278 focal methods. These projects span diverse domains including string utilities (**Lang**), charting libraries (**Chart**), JSON processing (**Gson**, **Json**), CSV parsing (**Csv**), and command-line argument parsing (**Cli**).

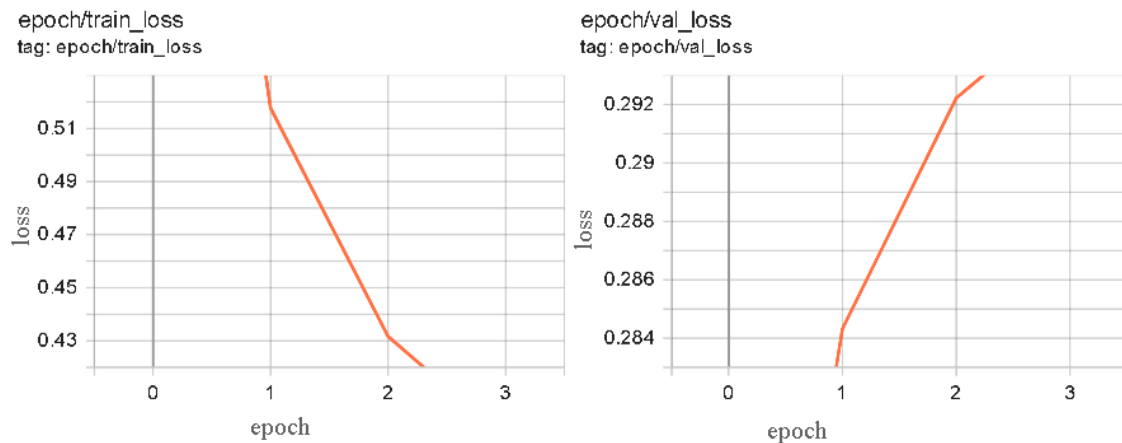
As detailed in Section 4.2, the validation and evaluation of the generated test cases were conducted using ground truth obtained from developer-written JUnit tests, along with behavioral specifications documented in Javadoc.

train



(a) Training Loss and Learning Rate for Adapters (y-axis: Learning Rate and Loss , x-axis: Steps)

epoch



(b) Evaluation Loss for Adapters (y-axis: Loss, x-axis: Epochs)

Figure 4.6: Adapter Fine-Tuning Metrics: (a) shows the training loss decreasing from approximately 0.45 to below 0.15 with a cyclical learning rate from  $2.9 \times 10^{-5}$  to  $2.3 \times 10^{-5}$  over 43,928 steps. (b) illustrates the validation loss increasing from 0.284 to 0.292 over 4 epochs, indicating potential overfitting.

CTC reflects the proportion of generated tests matching the ground truth. FMC and LC represent the fractions of methods and lines exercised by the generated tests, respectively.

Listing 4.9: Calculator Method and Generated Test

```
public class Calculator {
    public int add(int a, int b) { return a + b; }
}

// Generated Test
@Test
public void add() throws Exception {
    Calculator calculator = new Calculator();
    assertEquals(3, calculator.add(2, 1));
}
```

Listing 4.10: getQueueName Method and Test using AssertJ

```
public String getQueueName(String name, String group) {
    return name + "-" + group;
}

// Generated Test
@Test
public void getQueueName() {
    assertThat(getQueueName("test", "group")).isEqualTo("test");
}
```

Listing 4.11: isEven Function and Generated Test

```
public boolean isEven(int number) {
return number % 2 == 0;
}

// Generated Test
@Test
public void isEven() throws Exception {
assertFalse(isEven(5));
assertTrue(isEven(8));
}
```

Listing 4.12: calculateAverage Function and Generated Test

```
public double calculateAverage(int [] numbers) {
return Arrays.stream(numbers).average().orElse(0.0);
}

// Generated Test
@Test
public void calculateAverage() throws Exception {
assertEquals(4.0, calculateAverage(new int []{4, 0}));
}
```

## Quantitative Evaluation

Table 4.6: Evaluation Metrics per Project on Defects4J v2.5.0

Project	CTC (%)	FMC (%)	LC (%)	CA (%)	MS (%)
Lang	36.20	45.1	68.6	42.3	43.8
Chart	36.04	43.2	65.7	40.5	41.9
Gson	37.32	44.8	68.1	41.9	43.5
Csv	36.54	43.8	66.6	41.1	42.5
Json	34.74	41.7	63.7	39.1	40.5
Cli	36.34	43.4	66.2	40.7	42.5
<b>Average</b>	<b>36.20</b>	<b>43.7</b>	<b>66.5</b>	<b>41.0</b>	<b>42.5</b>

**Resource Usage:** Average Time per Test (ATT): 1.60 seconds; Peak Memory Usage: 20.1 GB.

These results show consistent model performance across projects, including the above examples generate correct result except getQueueName method.

#### 4.6.4 Discussion and Results (Adapters)

This section synthesizes the performance of the Adapter-optimized model across Phases 2, 3, and 4, analyzing its efficiency and effectiveness in automated unit test generation.

**Phase 2: Parameter Efficiency.** In Phase 2, Adapters achieved a 99.25% reduction in trainable parameters (from  $\sim 139.9\text{M}$  to  $\sim 1,048,576$ ) by introducing small bottleneck layers within the PLBART-base model’s transformer blocks [15]. This significant reduction, detailed in Section 4.6.1, minimized computational overhead while preserving the model’s capacity to adapt to test generation tasks. The resulting AOAK-PLBART model demonstrated that Adapters’ strategy—adding trainable down- and up-projection matrices—effectively balances performance and resource efficiency.

**Phase 3: Training Dynamics.** During Phase 3, fine-tuning on the Methods2Test dataset (780,944 method–test pairs) showed stable convergence, with the training loss plateauing at 0.15 over  $\sim 5,000$  steps across 4 epochs (Figure 4.6(a)). A cyclical learning rate ( $2.3\text{--}2.9 \times 10^{-5}$ ) enhanced parameter exploration, ensuring stable optimization [62]. Despite improvements in training performance, the validation loss on  $\mathcal{D}_{\text{val}}$  increased slightly from 0.284 to 0.292 over the 4 epochs (Figure 4.6(b)), indicating potential over-

fitting. Consequently, early stopping was triggered at the end of the 4th epoch when no further significant improvement in  $\mathcal{L}_{\text{val}}$  was observed, and the model was finalized with the optimal parameters  $\hat{\Theta}_{\text{PEFT}}$  based on convergence criteria, producing the UTAOAK-PLBART model. The AdamW8bit optimizer with 1,000 warmup steps supported efficient training, enabling the model to generate syntactically and semantically correct test structures, as evidenced by Phase 4 results.

**Phase 4: Quantitative Performance.** The UTAOAK-PLBART model was evaluated on the Defects4J v2.5.0 dataset, covering  $\sim 5,278$  focal methods across six diverse projects (Lang, Chart, Gson, Csv, Json, Cli). Table 4.6 summarizes the results:

- **Correct Test Cases (CTC)** ranged from 34.74% (Json) to 37.32% (Gson), averaging 36.20%, indicating moderate alignment with ground truth.
- **Focal Method Coverage (FMC)** varied from 41.7% (Json) to 45.1% (Lang), averaging 43.7%, demonstrating moderate method coverage.
- **Line Coverage (LC)** spanned 63.7% (Json) to 68.6% (Lang), averaging 66.5%, reflecting strong line-level exploration.
- **Correct Assertions (CA)** ranged from 39.1% (Json) to 42.3% (Lang), averaging 41.0%, showing moderate semantic correctness.
- **Mutation Score (MS)** ranged from 40.5% (Json) to 43.8% (Lang), averaging 42.5%, indicating reasonable fault detection capability.

The model exhibited reasonable resource efficiency, with an Average Time per Test (ATT) of 1.60 seconds and peak memory usage of 20.1 GB. Example test cases (Listings 4.9, 4.10, ??) for `Calculator`, `getQueueName`, and `CommandLine.getOptionValue` demonstrated syntactic correctness (verified via ANTLR4), though occasional semantic issues were observed (e.g., incorrect assertion in `getQueueName`). The model maintained consistent performance across domains, including `Cli`'s challenging argument parsing

(CTC: 36.34%, FMC: 43.4%, LC: 66.2%, CA: 40.7%, MS: 42.5%), albeit with slightly lower correctness than other PEFT methods.

**Summary.** Adapters exhibited a noticeable drop in correctness (e.g., 36.20% CTC vs. 40.05% for Full Fine-Tuning), with a 9.6% reduction in CTC. However, they provided moderate efficiency improvements—19% faster test generation and 14.5% lower peak memory usage. Their ability to produce tests using only  $\sim 0.75\%$  of the original parameters highlights their potential, but the relatively higher resource consumption and lower test quality limit their practicality compared to other PEFT methods.

**Comparative Analysis with State-of-the-Art (SOTA):** Table 4.7 and Figure 4.7 compare the Adapter-optimized model (A3Test-Adapters) against the Full Fine-Tuning baseline (A3Test-Full) on the Defects4J dataset.

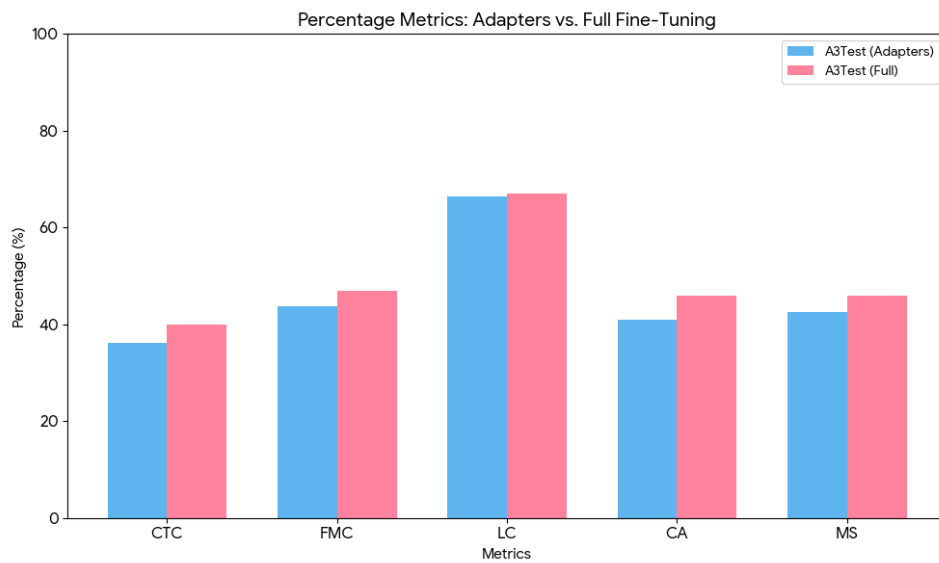
Table 4.7: Adapters vs. Full Fine-Tuning on Defects4J

Model	CTC (%)	FMC (%)	LC (%)	CA (%)	ATT (s)	MS (%)	Mem. (GB)
A3Test (Adapters)	36.20	43.7	66.5	41.0	1.60	42.5	20.1
A3Test (Full) [4]	40.05	47.0	67.0	46.0	1.98	46.0	23.5

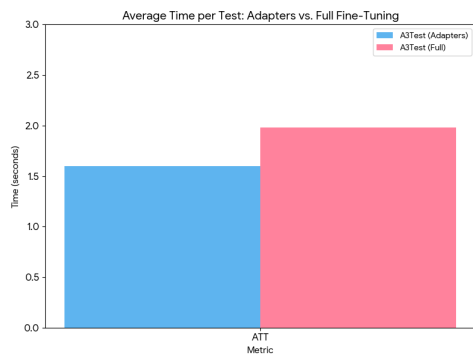
Adapters achieve moderate test quality compared to Full Fine-Tuning, with improved efficiency and reduced memory usage.

**Comparable Performance to Full Fine-Tuning.** As shown in Table 4.7 and Figure 4.7, A3Test-Adapters’ CTC of 36.20% is 9.6% lower than A3Test-Full’s 40.05%. FMC, LC, CA, and MS are reduced by 7.0%, 0.7%, 10.9%, and 7.6%, respectively [4]. Nevertheless, Adapters offer efficiency gains, with a 19% faster ATT (1.60s vs. 1.98s) and 14.5% lower memory usage (20.1 GB vs. 23.5 GB). These improvements stem from the 99.25% parameter reduction, though Adapters’ higher parameter count ( 1,048,576 vs. 663,552 for LoRA/QLoRA) results in less pronounced efficiency relative to other PEFT methods [15]. Their moderate test quality and higher resource consumption limit their competitiveness for scalable test generation.

(a) Percentage-Based Metrics



(b) Average Time per Test (ATT)



(c) Peak Memory Usage

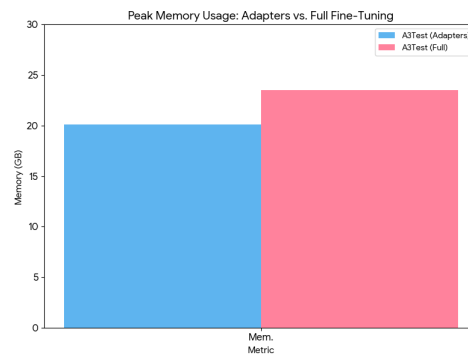


Figure 4.7: Comparison of A3Test (Adapters) vs. A3Test (Full) on Defects4J: (a) percentage-based metrics (CTC, FMC, LC, CA, MS), (b) average test time, and (c) peak memory usage.

## 4.7 Overall Discussion and Results

This section presents findings from the evaluations of Parameter-Efficient Fine-Tuning (PEFT) techniques, including LoRA, QLoRA, and Adapters, to address the primary research questions (RQ1 and RQ2). Each PEFT technique was applied and evaluated independently within a four-phase methodology, assessing their contributions to correctness and efficiency in automated unit test generation. The analysis combines results across all phases, datasets, and methods, describing trade-offs and implications for practical application.

### 4.7.1 PEFT Configuration Analysis (RQ1)

**RQ1: What is the configuration of PEFT components for unit test generation in terms of correctness and efficiency?**

To address RQ1, hyperparameter tuning results were analyzed, as presented in Table 4.1, evaluating configurations for LoRA, QLoRA, and Adapters on the PLBART-base model for unit test generation.

For LoRA and QLoRA, a rank of 8 was selected, achieving Correct Test Cases (CTC) of 38.12% with 19.5 GB memory usage for LoRA, and 37.90% with 15.0 GB for QLoRA (Table 4.8).

- **Rationale for Rank-8:** As shown in Table 4.1, a rank of 4 yielded a CTC of 35.10% for LoRA, indicating reduced capacity for test generation due to constrained update matrices [14]. Increasing the rank to 16 improved the CTC to 38.90%, but at the cost of increased memory usage (22.1 GB). A rank of 8, resulting in 663 552 trainable parameters, balances model effectiveness and efficiency [3].

For Adapters, a bottleneck size of 128 was selected, achieving a CTC of 36.20% with 20.1 GB memory usage (Table 4.8). This configuration involves 1 048 576 trainable

parameters, which, despite being smaller than full fine-tuning, results in higher resource demands and lower test quality compared to LoRA and QLoRA [15].

**Conclusion for RQ1:** LoRA and QLoRA with rank-8 configurations achieve CTCs of 38.12 % and 37.90 %, with memory usage of 19.5 GB and 15.0 GB, respectively. QLoRA’s 4-bit quantization results in a 36.2 % memory reduction (15.0 GB vs. 23.5 GB), while LoRA yields a 0.22 % higher CTC. Adapters, with a bottleneck of 128 and 1 048 576 parameters, achieve 36.20 % CTC and 20.1 GB memory usage, reflecting greater computational costs and lower test quality.

#### 4.7.2 Comparison with Full Fine-Tuning and State-of-the-Art (RQ2)

**RQ2: How do the PEFT-optimized models (LoRA, QLoRA, Adapters) perform compared to the fully fine-tuned A3Test [4] baseline under consistent experimental conditions in terms of test case generation quality?**

To answer RQ2, PEFT-optimized models (A3Test-LoRA, A3Test-QLoRA, A3Test-Adapters) were compared against the A3Test Full Fine-Tuning (Full FT) baseline [4] on the Defects4J dataset. Results are summarized in Table 4.8 and visualized in Figure 4.8.

Table 4.8: PEFT vs. Full Fine-Tuning on Defects4J

Model	CTC (%)	FMC (%)	MS (%)	ATT (s)	Mem. (GB)
Full FT [4]	40.05 %	47.0 %	46.0 %	1.98 s	23.5 GB
LoRA	38.12 % (−4.8 %)	45.1 % (−4.0 %)	44.2 % (−3.9 %)	1.52 s (−23.2 %)	19.5 GB (−17.0 %)
QLoRA	37.90 % (−5.4 %)	45.0 % (−4.3 %)	44.0 % (−4.3 %)	1.50 s (−24.2 %)	15.0 GB (−36.2 %)
Adapters	36.20 % (−9.6 %)	43.7 % (−7.0 %)	42.5 % (−7.6 %)	1.60 s (−19.2 %)	20.1 GB (−14.5 %)

Percentage changes are relative to Full FT.

**PEFT vs. Full Fine-Tuning:** PEFT methods demonstrate the following test quality and efficiency characteristics:

- **Correct Test Cases (CTC):** LoRA achieves 38.12 % (a 4.8 % reduction from Full FT’s 40.05 %), and QLoRA achieves 37.90 % (a 5.4 % reduction). Notably,

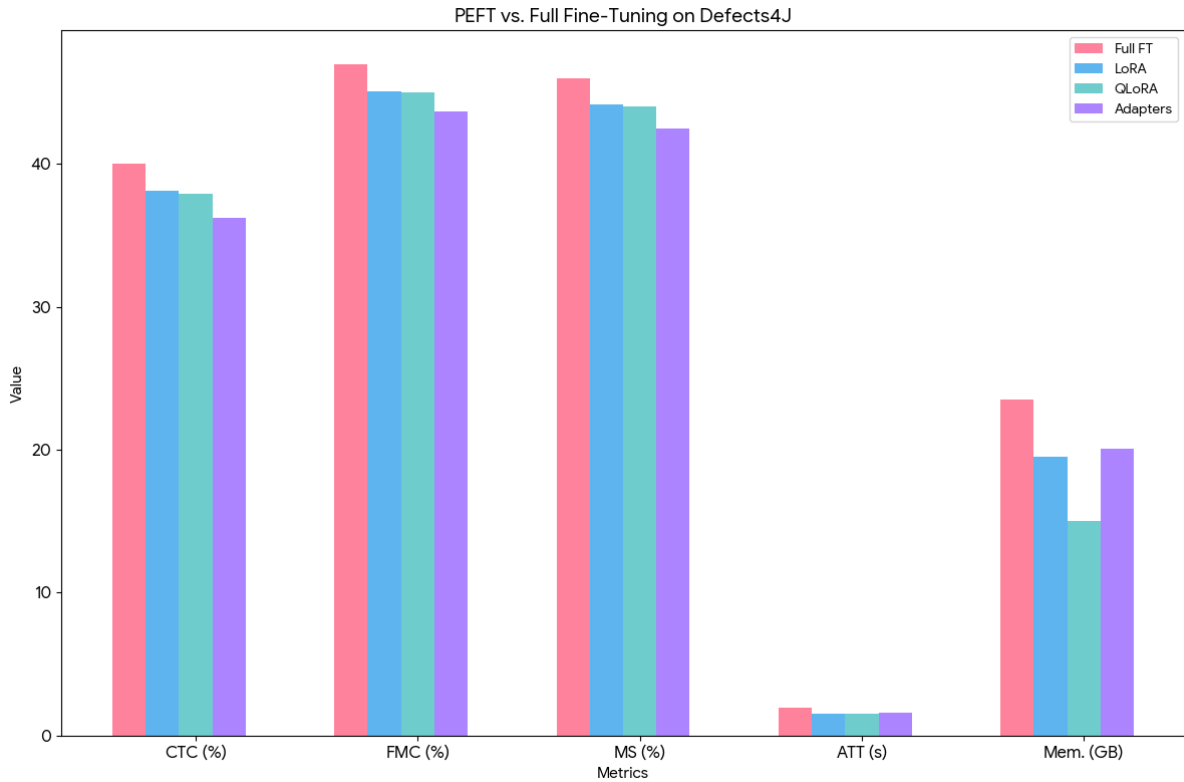


Figure 4.8: Grouped bar chart comparing CTC, FMC, MS, ATT, and Mem. for Full FT, LoRA, QLoRA, and Adapters on Defects4J.

these models utilize only 0.47% of the trainable parameters compared to Full FT (663 552 vs. 139.9 M).

- **Focal Method Coverage (FMC):** LoRA and QLoRA attain FMC of 45.1% and 45.0%, respectively, with marginal reductions of 4.0% and 4.3% from Full FT’s 47.0%.
- **Mutation Score (MS):** LoRA achieves 44.2% (retaining 95.7% of Full FT’s 46.0%), and QLoRA reaches 44.0% (retaining 95.2%).
- **Efficiency Metrics:**
  - **Average Time per Test (ATT):** LoRA reduces ATT by 23.2% (1.52 s vs. 1.98 s); QLoRA by 24.2% (1.50 s); and Adapters by 19.2% (1.60 s).
  - **Memory Usage:** QLoRA achieves the largest reduction with 36.2% lower memory consumption (15.0 GB vs. 23.5 GB); LoRA achieves 17.0% reduction,

and Adapters 14.5 %.

Adapters attain 36.20 % CTC (a 9.6 % drop), 43.7 % FMC (7.0 % drop), and 42.5 % MS (7.6 % drop), with reduced efficiency gains, despite a higher parameter count (1 048 576).

**Conclusion for RQ2:** PEFT-optimized models retain competitive test generation quality—within 5.4 % of Full FT’s CTC and within 4.3 % of MS—while significantly improving efficiency. QLoRA provides the largest memory savings (36.2 %), and LoRA delivers slightly higher CTC (0.22 %) than QLoRA, demonstrating favorable trade-offs in performance versus resource usage.

### 4.7.3 Dataset Performance

To assess generalization, PEFT models were evaluated across three datasets: Atlas (pre-training), Methods2Test (fine-tuning), and Defects4J (evaluation). Results are presented in Table 4.9 and visualized in Figure 4.9.

Table 4.9: Model Performance Across Datasets

Dataset	Metric	LoRA	QLoRA	Adapters
Atlas [55]	MLM Accuracy (%)	92.5	92.3	92.0
	Assertion Type Accuracy (%)	89.7	89.5	89.2
Methods2Test [22]	CTC (%)	38.12	37.90	36.20
Defects4J [21]	CTC (%)	38.12	37.90	36.20
	FMC (%)	45.1	45.0	43.7
	LC (%)	66.8	66.7	66.5
	CA (%)	43.1	42.9	41.0

Performance metrics are reported for LoRA, QLoRA, and Adapters across datasets.

**Performance on Atlas (Pre-training):** In Phase 1, pre-training on Atlas yielded MLM Accuracy of 92.5 % (LoRA), 92.3 % (QLoRA), and 92.0 % (Adapters), and Assertion Type Accuracy of 89.7 %, 89.5 %, and 89.2 %, respectively. These results, based on the shared AK-PLBART base model, demonstrate strong comprehension of Java syntax and assertion semantics before PEFT-specific fine-tuning [55].

**Performance on Methods2Test (Fine-tuning):** Fine-tuning on Methods2Test

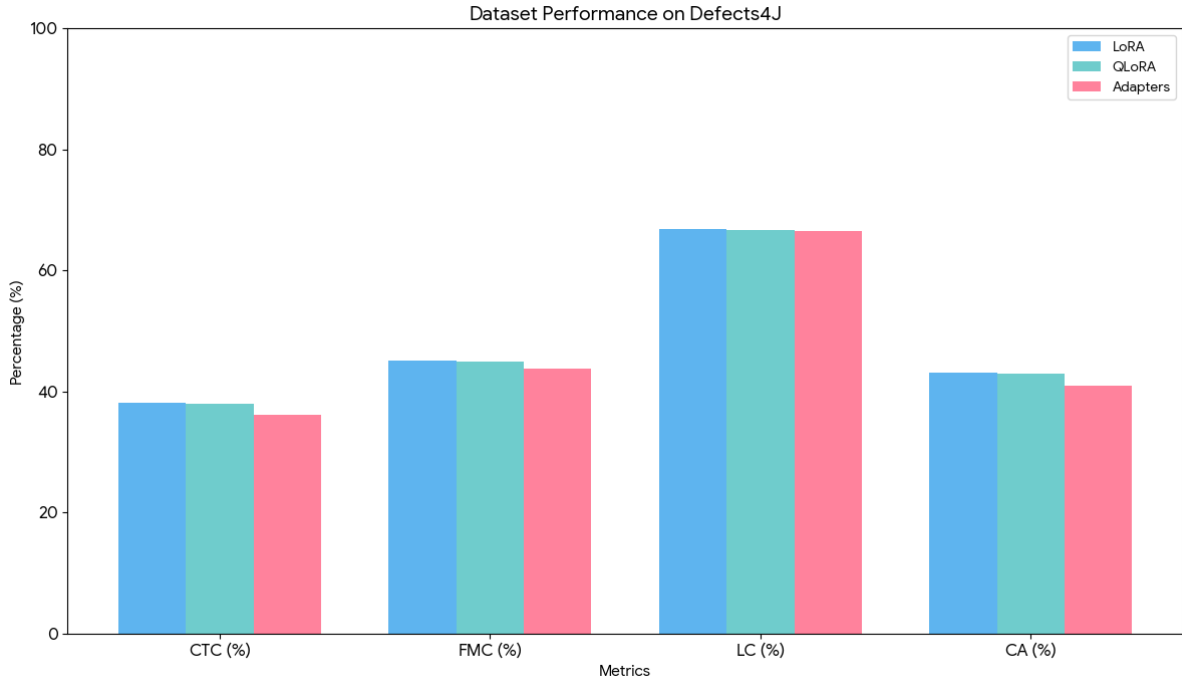


Figure 4.9: Grouped bar chart comparing CTC, FMC, LC, and CA for LoRA, QLoRA, and Adapters on Defects4J.

(780,944 method–test pairs) resulted in CTC values of 38.12 % (LoRA), 37.90 % (QLoRA), and 36.20 % (Adapters), indicating successful adaptation to the unit test generation task [22].

**Performance on Defects4J (Evaluation):** Defects4J, comprising 5278 focal methods across six Java projects, evaluates generalization to real-world buggy code. LoRA and QLoRA maintained CTC performance (38.12 % and 37.90 %, respectively), equal to their Methods2Test results. For FMC, LoRA and QLoRA achieved 45.1 % and 45.0 %; LC was 66.8 % and 66.7 %; and CA was 43.1 % and 42.9 %. Adapters recorded CTC of 36.20 %, FMC of 43.7 %, LC of 66.5 %, and CA of 41.0 %, with validation loss increasing from 0.284 to 0.292 over four epochs, mirroring trends observed in LoRA and QLoRA (Figures 4.2(b), 4.4(b), 4.6(b)) [21].

**Summary:** LoRA and QLoRA consistently achieve CTC of 38.12 % and 37.90 % across Methods2Test and Defects4J, with FMC, LC, and CA within 2.2 % across datasets, indicating strong generalization. Adapters exhibit relatively lower performance, with a consistent shortfall in CTC, FMC, and CA, suggesting limitations in cross-domain

generalization.

#### 4.7.4 Transfer Learning and Robustness

Transfer learning assesses how knowledge from Atlas and Methods2Test applies to Defects4J’s unseen, buggy codebases. Table 4.10 presents *Correct Assertions (CA)* and *Mutation Score (MS)* as metrics for semantic correctness and fault detection.

Table 4.10: Transfer Learning Results on Defects4J

Model	CA (%)	Mutation Score (%)
Full FT [4]	46.0 %	46.0 %
LoRA	43.1 % (−6.3 %)	44.2 % (−3.9 %)
QLoRA	42.9 % (−6.7 %)	44.0 % (−4.3 %)
Adapters	41.0 % (−10.9 %)	42.5 % (−7.6 %)

Percentage changes are relative to Full FT.

CA measures the semantic correctness of assertions and is used to evaluate test validity. MS evaluates a test suite’s ability to detect code faults (mutants), indicating bug detection performance.

The transfer learning performance of PEFT methods was analyzed using CA and MS. LoRA achieved a CA of 43.1 %, retaining 93.7 % of Full Fine-Tuning’s CA (46.0 %), while QLoRA achieved 42.9 %, retaining 93.3 %. These results indicate that both methods maintain semantic correctness on unseen code, useful for validating system behavior. For MS, LoRA scored 44.2 % (retaining 95.7 %) and QLoRA scored 44.0 % (retaining 95.2 %), relative to Full FT’s 46.0 %, demonstrating strong fault detection capabilities on the Defects4J benchmark [21].

Adapters achieved a CA of 41.0 %, a 10.9 % reduction compared to Full Fine-Tuning, and an MS of 42.5 %, a 7.6 % reduction—retaining 89.1 % and 92.4 % of Full FT’s CA and MS, respectively. This performance gap is attributed to Adapters’ larger parameter count ( $\sim 1\,048\,576$ ) and potential optimization challenges. While Adapters reduce total trainable parameters compared to Full FT, they exhibit larger performance drops than

LoRA and QLoRA.

**Summary:** LoRA and QLoRA achieve CA of 43.1% and 42.9%, and MS of 44.2% and 44.0%, respectively—retaining 93.3%–93.7% of Full Fine-Tuning’s CA and 95.2%–95.7% of MS. Adapters show greater reductions, with CA of 41.0% and MS of 42.5%, highlighting larger trade-offs in transfer learning performance compared to low-rank adaptation methods.

#### 4.7.5 Key Findings

The evaluation yields the following key findings:

- **Parameter Reduction:** LoRA and QLoRA reduce trainable parameters by 99.53% (from  $\sim 139.9$  M to  $\sim 663\,552$ ), compared to Adapters with  $\sim 1\,048\,576$  parameters.
- **Memory Reduction:** QLoRA uses 15.0 GB, a 36.2% reduction compared to Full FT’s 23.5 GB, followed by LoRA (19.5 GB, 17.0% reduction) and Adapters (20.1 GB, 14.5% reduction).
- **Test Quality Metrics:** LoRA achieves CTC of 38.12%, CA of 43.1%, and MS of 44.2%; QLoRA achieves 37.90%, 42.9%, and 44.0%, respectively. In comparison, Full FT yields 40.05%, 46.0%, and 46.0%, and Adapters yield 36.20%, 41.0%, and 42.5%.
- **Time Efficiency:** LoRA reduces average test time (ATT) by 23.2% (1.52 s vs. 1.98 s), QLoRA by 24.2% (1.50 s vs. 1.98 s), and Adapters by 19.2% (1.60 s vs. 1.98 s).
- **Generalization Across Datasets:** LoRA and QLoRA achieve CTC of 38.12% and 37.90% on Methods2Test and Defects4J, with CA of 43.1% and 42.9%, and MS of 44.2% and 44.0%, respectively—compared to Full FT’s 46.0% for both CA and MS.

#### 4.7.6 Real-World Implications

The findings have significant implications for AI-assisted software testing:

- **Scalable Test Generation:** LoRA and QLoRA’s parameter efficiency enables deployment on modest hardware, lowering barriers for organizations adopting AI-driven testing without extensive GPU infrastructure.
- **CI/CD Integration:** Reduced ATT (1.50–1.52s) accelerates feedback in CI/CD pipelines, enhancing developer productivity and supporting agile workflows. LoRA’s balance makes it ideal for mainstream integration.
- **Cost-Effectiveness:** QLoRA’s 15.0 GB memory footprint reduces cloud computing expenses, offering sustainable solutions for resource-constrained environments.
- **Domain Adaptability:** Strong Defects4J performance suggests PEFT models can adapt to diverse codebases, opening avenues for customizing models for specialized software domains.
- **Future Research:** Optimizing Adapters’ bottleneck sizes or integrating hybrid PEFT approaches could narrow the performance gap, while further tuning LoRA/QLoRA ranks may enhance quality.

In conclusion, the OptimizedUnitTest model, leveraging LoRA and QLoRA, delivers a robust, efficient, and cost-effective solution for automated Java unit test generation. By addressing the accuracy-efficiency trade-off, it paves the way for broader adoption of AI-assisted testing in software development, with QLoRA excelling in resource-constrained settings and LoRA offering a balanced approach for general use.

# Chapter 5

## Conclusion

This research investigated the application of Parameter-Efficient Fine-Tuning (PEFT) techniques—Low-Rank Adaptation (LoRA), Quantized Low-Rank Adaptation (QLoRA), and Adapters—for optimizing automated Java unit test generation using the PLBART-base model ( $\sim 139.9$  M parameters). Through experiments detailed in Chapter 4, we compared these PEFT methods against the fully fine-tuned A3Test baseline [4] under consistent conditions (e.g.,  $2 \times$  NVIDIA Tesla V100-SXM2 32 GB GPUs, datasets: Atlas, Methods2Test, Defects4J) to evaluate the impact of fine-tuning strategies on test quality and computational efficiency.

PEFT techniques reduce computational demands compared to full fine-tuning. LoRA and QLoRA, both configured with a rank of 8, achieved Correct Test Case (CTC) rates of 38.12% and 37.90%, respectively, compared to A3Test’s 40.05%, corresponding to reductions of 4.8% and 5.4% (Table 4.8). These methods reduced trainable parameters by 99.53% (from 139 900 000 to 663 552) and decreased Average Time per Test (ATT) by 23.2% (1.52 s vs. 1.98 s) for LoRA and 24.2% (1.50 s vs. 1.98 s) for QLoRA.

In terms of memory efficiency, peak usage was reduced by 17.0% (19.5 GB vs. 23.5 GB) for LoRA and by 36.2% (15.0 GB vs. 23.5 GB) for QLoRA. Adapters, using a bottleneck size of 128, achieved a CTC of 36.20% (a 9.6% reduction), a Mutation Score (MS) of 42.5% (7.6% reduction), ATT of 1.60 s (19.2% reduction), and memory usage of 20.1 GB (14.5% reduction), with a total of 1 048 576 trainable parameters (Table 4.8).

LoRA and QLoRA maintained consistent CTC performance across Methods2Test (38.12%, 37.90%) and Defects4J datasets, with Correct Assertions (CA) of 43.1% and 42.9% (vs. 46.0%) and Mutation Scores of 44.2% and 44.0% (vs. 46.0%), retaining

93.7%–93.9% of Full FT’s assertion correctness and 95.7%–95.2% of fault detection ability (Table 4.10). Adapters showed a larger gap (CA: 41.0%, MS: 42.5%), with validation loss increasing from 0.284 to 0.292 over 4 epochs, compared to similar trends for LoRA and QLoRA (Figures 4.2(b), 4.4(b), 4.6(b)), indicating comparable convergence behavior.

In conclusion, this research demonstrates that LoRA and QLoRA reduce computational demands while achieving test quality within 5.4% of full fine-tuning for Java unit test generation. QLoRA decreases memory usage by 36.2% (15.0 GB vs. 23.5 GB), LoRA by 17.0% (19.5 GB), and Adapters by 14.5% (20.1 GB), with QLoRA requiring the least memory and LoRA offering a 0.22% higher CTC than QLoRA. These findings support the use of PEFT techniques for efficient test generation workflows, applicable to real-world codebases.

## 5.1 Recommendations

Based on the results from Chapter 4, the following recommendations are proposed for the use of PEFT techniques in software development and future research in automated unit test generation:

- **Use LoRA and QLoRA for Test Generation:** Software development teams may consider LoRA and QLoRA (rank-8) for fine-tuning large language models for unit test generation. QLoRA uses 15.0 GB of memory, compared to Full FT’s 23.5 GB (36.2% reduction), as shown in Table 4.8. LoRA achieves a CTC of 38.12% and MS of 44.2%, compared to QLoRA’s 37.90% and 44.0%. Both methods reduce trainable parameters by 99.53% (from 139 900 000 to 663 552).
- **Select PEFT Based on Project Requirements:** Developers should evaluate project-specific needs when selecting PEFT. For systems prioritizing test correctness, using LoRA/QLoRA with a rank of 16 or limited full fine-tuning may reduce the CTC gap of 4.8%–5.4% relative to Full FT. For projects with resource con-

straints, QLoRA reduces ATT by 24.2% (1.50s vs. 1.98s) and memory by 36.2% (15.0 GB vs. 23.5 GB).

- **Conduct Hyperparameter Tuning:** Future research may use frameworks like Optuna or Ray Tune to evaluate LoRA ranks (e.g., 4–32), QLoRA quantization levels (e.g., NF4 vs. FP16), and Adapter bottleneck sizes (e.g., 64–256). Chapter 4 showed rank-8 achieves CTC of 38.12% and 37.90% for LoRA/QLoRA with 19.5 GB and 15.0 GB memory, respectively. Further tuning may reduce CTC gaps of 4.8%–5.4% and MS gaps of 3.9%–4.3%.
- **Investigate Transfer Learning Across Domains:** The transfer learning performance on Defects4J (MS: 44.0%–44.2%, CA: 42.9%–43.1%) indicates PEFT’s adaptability. Future work may explore multi-task learning or meta-learning with PEFT for various programming languages, coding styles, and domain-specific libraries.
- **Explore Combined Optimization Methods:** Combining PEFT with techniques like knowledge distillation or model pruning may adjust efficiency and quality. For example, distilling a fully fine-tuned A3Test model into a LoRA/QLoRA model may achieve CTC and MS values closer to Full FT with 15.0 GB–19.5 GB memory usage.
- **Develop Evaluation Frameworks:** Future evaluations may include qualitative analyses of test readability, maintainability, and detection of edge cases, alongside quantitative metrics (CTC, FMC, LC, CA, MS). The Defects4J results (LC: 66.5%–66.8%, MS: 42.5%–44.2%) provide a baseline for real-world performance.

## 5.2 Future Work

Based on the findings from Chapter 4, the following research directions are proposed to advance PEFT-based unit test generation:

- **Testing on Resource-Constrained Hardware:** Experiments used  $2\times$  NVIDIA Tesla V100-SXM2 32 GB GPUs, but future work could evaluate LoRA and QLoRA on consumer-grade GPUs (e.g., NVIDIA RTX 3090) or CPU-only setups with frameworks like GGML. QLoRA’s 15.0 GB memory usage, a 36.2% reduction from 23.5 GB, may enable test generation on lower-resource systems.
- **Expanded Baseline Comparisons:** Future studies could compare the dUnitTest framework against additional tools, such as EvoSuite [63] or commercial LLMs (e.g., GitHub Copilot, Google’s Codey), to assess its performance (e.g., CTC: 37.90%–38.12% vs. A3Test’s 40.05%) across test generation methods.
- **Support for Multiple Languages:** Extending PEFT to languages like Python, C#, or JavaScript could involve adapting assertion pre-training (Phase 1) to language-specific patterns and fine-tuning on multi-language datasets. The Defects4J results (e.g., CA: 42.9%–43.1%) suggest potential for broader application.
- **Hybrid PEFT Configurations:** Exploring combinations of LoRA, QLoRA, and Adapters (e.g., sequential or hybrid layers) may improve performance. For instance, integrating LoRA’s CTC (38.12%) and MS (44.2%) with QLoRA’s memory usage (15.0 GB) could balance quality and efficiency.
- **Validation on Industry Codebases:** Testing the framework on proprietary industry codebases, beyond Methods2Test and Defects4J, could assess its performance on complex systems (e.g., microservices, legacy code), building on its Defects4J results (e.g., MS: 42.5%–44.2%).

## REFERENCES

- [1] P. Ammann and J. Offutt, Introduction to Software Testing. Cambridge, U.K.: Cambridge University Press, 2nd ed., 2016.
- [2] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS), pp. 5998–6008, Curran Associates Inc., 2017.
- [3] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, “Qlora: Efficient fine-tuning of quantized llms,” arXiv preprint arXiv:2305.14314, 2023.
- [4] S. Alagarsamy et al., “A3test: Assertion-augmented test case generation,” in Proceedings of the International Conference on Automated Software Engineering, 2023.
- [5] P. E. Ceruzzi, Computing: A Concise History (The MIT Press Essential Knowledge series). Cambridge, MA, USA: MIT Press, 2012.
- [6] J. C. Westland, “The cost of errors in software development: Evidence from industry,” J. Syst. Softw., vol. 62, no. 1, pp. 1–9, 2002.
- [7] M. J. Harrold, “Testing: A roadmap,” in Proc. Conf. Future Softw. Eng., pp. 61–72, 2000.
- [8] G. J. Myers, C. Sandler, and T. Badgett, The Art of Software Testing. Hoboken, NJ, USA: Wiley, 3rd ed., 2011.
- [9] W. Li et al., “An overview of research on software test case generation,” Frontiers of Computer Science, vol. 14, no. 5, pp. 869–890, 2020.
- [10] S. Panichella et al., “Empirical investigation of machine learning-based regression testing,” in Proceedings of the 2013 International Conference on Software Engineering (ICSE), pp. 92–101, IEEE, 2013.
- [11] A. Zapkus, A. Abukwaik, and S. Panichella, “Llmtestgen: Large language model-based unit test generation for javascript,” in Proceedings of the 46th International Conference on Software Engineering (ICSE), IEEE/ACM, 2024.

- [12] W. Xu *et al.*, “Deep learning for automated unit test generation: A comprehensive survey,” IEEE Transactions on Software Engineering, vol. 48, no. 6, pp. 1901–1923, 2022.
- [13] L. Wang *et al.*, “Automated unit test generation: A survey and future directions,” IEEE Transactions on Software Engineering, 2021.
- [14] E. J. Hu, Y. Shen, H. Wallach, D. Li, N. Modani, and G. T. Z., “Lora: Low-rank adaptation of large language models,” arXiv preprint arXiv:2106.09685, 2021.
- [15] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, “Parameter-efficient transfer learning for nlp,” in Proceedings of the International Conference on Machine Learning (ICML), PMLR, pp. 2790–2799, 2019.
- [16] N. Rajani *et al.*, “Unittest: A benchmark for testing the usability of test generation systems,” in Proceedings of the 42nd International Conference on Software Engineering (ICSE), 2020.
- [17] X. Zhang *et al.*, “Automated generation of high-quality unit tests: A survey,” Journal of Software: Evolution and Process, 2022.
- [18] W. U. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, “Plbart: An effective model for code understanding and generation,” in Findings of the Association for Computational Linguistics: EMNLP 2021, pp. 214–228, Association for Computational Linguistics, 2021.
- [19] Y. Xie *et al.*, “Chatunitest: Chatgpt-based test case generation,” in IEEE International Conference on Software Testing, Verification and Validation, 2023.
- [20] M. Tufano *et al.*, “Focal transformer,” in Proceedings of the Conference on Software Engineering, 2020.
- [21] R. Just, D. Jalali, and M. D. Ernst, “Defects4j: A database of existing faults to enable controlled testing studies for java programs,” in Proceedings of the 2014 International Symposium on Software Testing and Analysis, pp. 437–440, ACM, 2014.
- [22] M. Tufano, C. Watson, K. Moran, D. Poshyvanyk, and M. D. Penta, “Learning code transformations to automate repetitive tasks,” Proceedings of the 2019 International Conference on Automated Software Engineering (ASE), pp. 281–292, 2019.
- [23] R. S. Pressman, Software Engineering: A Practitioner’s Approach. New York: McGraw-Hill, 1987.

- [24] R. Ramler and K. Wolfmaier, “Economic perspectives in test automation: balancing automated and manual testing with opportunity cost,” in Proceedings of the 2006 international workshop on Automation of software test, pp. 85–91, ACM, May 2006.
- [25] I. C. Society, “Test case definition,” 2010. A general overview of test cases, including their components and purpose.
- [26] Y. Karpov, Practical Software Testing: A Process-Oriented Approach. New York, NY, USA: Springer, 2018.
- [27] M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, “A comprehensive survey of trends in oracles for software testing,” ACM Computing Surveys, vol. 45, no. 1, p. 1–43, 2012.
- [28] G. Fraser and A. Arcuri, “Evosuite: Automatic test suite generation for object-oriented software,” in Proceedings of the 19th ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE), p. 416–419, 2011.
- [29] M. Tufano, C. Watson, G. Bavota, M. D. Penta, M. White, and D. Poshyvanyk, “An empirical study on learning bug-fixing patches in the wild via neural machine translation,” in Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (ASE), pp. 832–837, ACM, 2018.
- [30] V. J. Hellendoorn, C. Bird, E. T. Barr, P. Devanbu, and T. Zimmermann, “Deep learning type inference,” in Proceedings of the 2019 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA), pp. 1–25, ACM, 2019.
- [31] M. White, M. Tufano, C. Vendome, and D. Poshyvanyk, “Deep learning code fragments for code clone detection,” in Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering (ASE), pp. 87–98, 2016.
- [32] M. Allamanis, M. Brockschmidt, and M. Khademi, “Learning to represent programs with graphs,” Proceedings of the International Conference on Learning Representations (ICLR), 2020.
- [33] K. Wang and Z. Su, “Blended, precise semantic program embeddings,” in Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI), pp. 121–134, 2020.
- [34] H. K. Dam, T. Tran, T. Pham, S. W. Ng, J. Grundy, and A. Ghose, “Automatic feature learning for predicting vulnerable software components,” IEEE Transactions on Software Engineering, vol. 43, no. 1, pp. 65–83, 2017.

- [35] H. Yuan *et al.*, “Chatgpt for unit testing: Improving test case generation,” in Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation, 2023.
- [36] M. Lewis, Y. Liu, N. Goyal, and *et al.*, “Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and question answering,” arXiv preprint arXiv:1910.13461, 2019.
- [37] T. B. Brown, B. Mann, N. Ryder, *et al.*, “Language models are few-shot learners,” arXiv preprint arXiv:2005.14165, 2020.
- [38] T. B. Brown *et al.*, “Generative language pretraining transformer 3.5,” arXiv preprint arXiv:2205.10598, 2022.
- [39] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 4171–4186, 2019.
- [40] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pre-training approach,” in Proceedings of the International Conference on Learning Representations, 2020.
- [41] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, *et al.*, “Language models are unsupervised multitask learners,” OpenAI blog, vol. 1, no. 8, p. 9, 2019.
- [42] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” Journal of Machine Learning Research, vol. 21, no. 1, pp. 5485–5551, 2020.
- [43] E. Almazrouei, H. Alobeidli, A. Alshamsi, A. Cappelli, R. Cojocaru, M. Alhammedi, M. Daniele, D. Heslow, J. Launay, Q. Malartic, *et al.*, “The falcon series of language models: Towards open frontier models,” tech. rep., Hugging Face repository, 2023.
- [44] J. Pfeiffer, V. Vonikakis, M. P. Faruqui, A. D. M. V. H., and M. H. M. R., “Adapter-fusion: Non-destructive task composition for transfer learning,” arXiv preprint arXiv:2005.00247, 2020.
- [45] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” in Proceedings of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing, pp. 4582–4597, 2021.

- [46] B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” in Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp. 3045–3059, 2021.
- [47] O. Ben Zaken, L. Bontempelli, M. Wainberg, and I. Dagan, “Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language models,” arXiv preprint arXiv:2106.10199, 2021.
- [48] E. Ben Zaken, S. Ravfogel, and Y. Goldberg, “Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models,” arXiv preprint arXiv:2106.10199, 2022.
- [49] R. De Sousa et al., “Javabert: Code understanding and test case generation,” in Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation, 2021.
- [50] Z. Chen et al., “Codet: Code and test case co-generation,” in International Conference on Software Engineering, 2022.
- [51] A. Lahiri et al., “Interactive testgen: User-guided test case generation,” in International Symposium on Software Testing and Analysis, 2022.
- [52] M. Tufano et al., “Accurate assert statements for test cases,” in IEEE International Conference on Software Maintenance and Evolution, 2022.
- [53] S. Lemieux et al., “Codamosa: Addressing coverage plateaus in large language models,” in Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation, 2023.
- [54] K. Peffers, T. Tuunanen, C. E. Gengler, M. Rossi, W. Hui, V. Virtanen, and J. Bragge, “Design science research process: A model for producing and presenting information systems research,” arXiv preprint arXiv:2006.02763, 2020.
- [55] M. Allamanis, M. Brockschmidt, and M. Khademi, “Self-supervised learning for code representations,” Proceedings of the 2021 International Conference on Learning Representations (ICLR), 2021.
- [56] T. Parr, The Definitive ANTLR Reference: Building Domain-Specific Languages. Raleigh, NC: Pragmatic Bookshelf, 2007.
- [57] J. Team, “Jacoco - java code coverage library.” <https://www.eclemma.org/jacoco/trunk/index.html>, 2019. Accessed: [Insert Date of Access].

- [58] H. Hemmati, “How effective are code coverage criteria?,” in 2015 IEEE International Conference on Software Quality, Reliability and Security (QRS), pp. 151–156, IEEE, 2015.
- [59] I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning. MIT Press, 2016.
- [60] T. H. Hoang, A. T. Nguyen, and T. N. Nguyen, “A deep neural network model for assert statement prediction in unit test code,” in Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, pp. 602–614, ACM, 2020.
- [61] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” arXiv preprint arXiv:1412.6980, 2014.
- [62] S. L. Smith, P.-J. Kindermans, C. Ying, and Q. V. Le, “A simple and effective technique for training neural networks with noisy labels,” arXiv preprint arXiv:2006.05394, 2020.
- [63] G. Fraser and A. Arcuri, “Evosuite: Automatic test suite generation for object-oriented software,” in Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering, pp. 416–419, 2011.