

# **InteGrow : A smart sync for development**

Project Team

Muskan Tariq 22I-2602  
Amna Hassan 22I-8759  
Shuja uddin 22I-2553

Session 2022-2026

Supervised by

**Dr. Muhammad Bilal**

Co-Supervised by

**Ms. Fatima Gillani**



**Department of Software Engineering**

**National University of Computer and Emerging Sciences  
Islamabad, Pakistan**

**December, 2025**

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	1
1.2	Problem Statement . . . . .	2
1.3	Software Requirements Specification . . . . .	2
1.3.1	Functional Requirements . . . . .	2
1.3.1.1	Core Foundation . . . . .	3
1.3.1.2	Requirements Auditor . . . . .	3
1.3.1.3	UML Synthesizer . . . . .	3
1.3.1.4	Code-UML Synchronization . . . . .	3
1.3.1.5	Code Review Assistant . . . . .	4
1.3.1.6	Technical Debt Analyzer . . . . .	4
1.3.1.7	Test Generator . . . . .	4
1.3.1.8	CI/CD Orchestration . . . . .	4
1.3.2	Non-Functional Requirements . . . . .	5
1.3.2.1	Performance Requirements . . . . .	5
1.3.2.2	Security Requirements . . . . .	5
1.3.2.3	Reliability Requirements . . . . .	5
1.3.2.4	Usability Requirements . . . . .	5
1.3.2.5	Maintainability Requirements . . . . .	5
1.3.2.6	Compatibility Requirements . . . . .	5
1.4	Work Distribution . . . . .	6
<b>2</b>	<b>Literature Review</b>	<b>7</b>
2.1	Related Research . . . . .	7
2.1.1	AI-Driven Requirements Analysis and Ambiguity Detection . . .	7
2.1.1.1	Summary of the Research Items . . . . .	7
2.1.1.2	Critical Analysis of the Research Items . . . . .	8
2.1.1.3	Relationship to the Proposed Research Work . . . . .	8
2.1.2	Requirements to Design Transformation . . . . .	8
2.1.2.1	Summary of the research items . . . . .	8
2.1.2.2	Critical analysis of the research items . . . . .	9
2.1.2.3	Relationship to the proposed research work . . . . .	9
2.1.3	Design-to-Code Generation . . . . .	10
2.1.3.1	Summary of the research items . . . . .	10
2.1.3.2	Critical analysis of the research items . . . . .	10
2.1.3.3	Relationship to the proposed research work . . . . .	11
2.1.4	AI-Enhanced Automated Code Review . . . . .	11

2.1.4.1	Summary of the research items . . . . .	11
2.1.4.2	Critical analysis of the research items . . . . .	11
2.1.4.3	Relationship to the proposed research work . . . . .	12
2.1.5	LLM-Powered Test Case Generation . . . . .	12
2.1.5.1	Summary of the research items . . . . .	12
2.1.5.2	Critical analysis of the research items . . . . .	13
2.1.5.3	Relationship to the proposed research work . . . . .	13
2.1.6	Technical Debt Prediction and Prioritization . . . . .	13
2.1.6.1	Summary of the research items . . . . .	13
2.1.6.2	Critical analysis of the research items . . . . .	14
2.1.6.3	Relationship to the proposed research work . . . . .	14
2.1.7	Model-Driven Development and SDLC Automation . . . . .	14
2.1.7.1	Summary of the research items . . . . .	14
2.1.7.2	Critical analysis of the research items . . . . .	15
2.1.7.3	Relationship to the proposed research work . . . . .	15
2.2	Analysis Summary of Research Items . . . . .	15
2.3	Research Scope / Gap . . . . .	15
2.3.1	Research Questions . . . . .	17
2.3.2	Research Objectives . . . . .	17
<b>3</b>	<b>Proposed Approach</b>	<b>18</b>
3.1	Overview . . . . .	18
3.2	Proposed System: InteGrow . . . . .	18
3.3	Module Methodologies . . . . .	18
3.3.1	Module 1: Requirements & Ethics Auditor . . . . .	18
3.3.2	Module 2: UML Synthesizer . . . . .	19
3.3.3	Module 3: Code-UML Round-Trip Synchronization . . . . .	19
3.3.4	Module 4: Code Review Assistant . . . . .	19
3.3.5	Module 5: Technical Debt Analyzer . . . . .	19
3.3.6	Module 6: Test Generator . . . . .	19
3.3.7	Module 7: CI/CD Orchestration . . . . .	19
3.3.8	Module Summary . . . . .	20
<b>4</b>	<b>Initial Experiments and Results</b>	<b>22</b>
4.1	Experimental Setup . . . . .	22
4.1.1	Hardware Configuration . . . . .	22
4.1.2	Software Environment . . . . .	22
4.1.3	Participant Demographics . . . . .	22
4.1.4	UAT Protocol . . . . .	22
4.1.5	Test Dataset . . . . .	23
4.2	Evaluation Metrics . . . . .	23
4.2.1	Mean Opinion Score (MOS) - ITU-T P.800 . . . . .	23
4.2.2	IEEE 830 Quality Thresholds . . . . .	23
4.2.3	Additional Metrics . . . . .	23
4.3	Results . . . . .	23
4.3.1	MOS Results . . . . .	23

4.3.2	Productivity Impact . . . . .	24
4.3.3	Production Readiness . . . . .	24
4.3.4	Qualitative Feedback . . . . .	24
4.4	Analysis and Discussion . . . . .	25
4.4.1	IEEE 830 Compliance and User Satisfaction . . . . .	25
4.4.2	Key Findings . . . . .	25
4.4.3	SDLC Alignment . . . . .	26
4.4.4	Limitations . . . . .	26
4.5	Conclusion . . . . .	26
<b>A</b>	<b>Appendices</b>	<b>28</b>
A.1	Appendix A: MOS Questionnaire . . . . .	28
A.2	Appendix B: Detailed Participant Responses . . . . .	31
	<b>References</b>	<b>35</b>

# List of Figures

- 3.1    InteGrow System Architecture . . . . . 21
- 4.1    *Mean Opinion Score (MOS) Results Across Five Dimensions (n=13)* . . . 24
- 4.2    *Time Savings Distribution Reported by UAT Participants* . . . . . 25
- A.1    *MOS Questionnaire Part 1: Usability and Accuracy Assessment (Questions 1-2)* . . . . . 28
- A.2    *MOS Questionnaire Part 2: Responsiveness and Usefulness Assessment (Questions 3-4)* . . . . . 29
- A.3    *MOS Questionnaire Part 3: Overall Satisfaction and Qualitative Feedback (Questions 5-7)* . . . . . 30

# List of Tables

- 1.1 Combined Responsibilities and Component Breakdown . . . . . 6
- 2.1 Detailed Analysis Summary of Research Items . . . . . 16
- 3.1 InteGrow Module Methodologies Summary . . . . . 20
- 4.1 *Mean Opinion Score (MOS) Results (n=13)* . . . . . 24
- A.1 *Complete Participant Responses (n=13)* . . . . . 31

# Chapter 1

## Introduction

Modern software development is becoming increasingly complex, requiring teams to manage numerous interconnected phases of the Software Development Life Cycle (SDLC). *InteGrow* responds to this challenge by proposing an AI-driven, integrated automation platform that supports major SDLC stages including requirements analysis, design synchronization, code generation, code review, test case generation, and technical debt management. By combining large language models (LLMs), symbolic reasoning, and model-driven engineering, *InteGrow* aims to reduce manual effort, improve quality, and maintain consistent traceability across all SDLC phases.

This research focuses on designing, developing, and evaluating a hybrid, modular architecture capable of enhancing development efficiency, maintainability, and code quality while ensuring human oversight and alignment with industry standards. The intended users include medium- to large-scale software development teams across diverse domains. The effectiveness of *InteGrow* will be evaluated through quantitative performance metrics, comparative analysis, and user studies.

### 1.1 Background

Software development is increasingly complex and fast-paced, demanding efficient and reliable tools to manage the entire Software Development Life Cycle (SDLC). Traditional manual practices in requirement engineering, design, coding, testing, and technical debt management lead to errors, delays, and rising maintenance costs. Despite advances in standalone tools, fragmentation and lack of integration persist, resulting in inefficient workflows and inconsistent quality.

Developers commonly spend substantial time reviewing ambiguous or incomplete requirements, reconciling inconsistencies between design and code, performing subjective code reviews prone to false positives, and manually designing tests with limited coverage. Technical debt often accumulates undetected, creating long-term risks for maintainability and system evolution.

Automation has become essential for addressing these challenges. While early automation primarily targeted repetitive tasks, modern approaches increasingly rely on AI-driven techniques such as large language models (LLMs), which can interpret natural language, generate code, and recommend fixes. However, LLMs alone lack robust logical reasoning and are susceptible to hallucinations, limiting their reliability in critical development workflows.

Integrating AI tools with existing SDLC pipelines such as CI/CD frameworks and static analyzers introduces challenges related to interoperability, scalability, and trust. Hybrid techniques that combine LLMs with symbolic reasoning and explainable AI show promise for improving accuracy and reliability.

Existing automation solutions are often domain-specific or address only one stage of the SDLC, leaving a gap for a unified solution that provides end-to-end traceability and consistency. Current industry trends emphasize security by design, continuous monitoring, and feedback-driven development, further reinforcing the need for integrated, explainable, and standards-compliant SDLC automation.

*InteGrow* aims to address this gap by developing a modular AI-powered platform that enhances collaboration, automates key SDLC stages, and maintains semantic consistency across artifacts. It seeks to reduce development cycle times, improve product quality, and proactively manage technical debt while supporting developer trust through transparent and explainable AI assistance.

This research responds to the growing demand for scalable, reliable, and user-centered AI automation in software engineering, anticipating the continued expansion of AI-driven workflows and the increasing complexity of modern software systems.

## 1.2 Problem Statement

The increasing complexity of modern software development has intensified the need for efficient, cohesive, and reliable management across all stages of the Software Development Life Cycle (SDLC). Traditional tools and practices remain fragmented, with requirements analysis, design modeling, coding, testing, and technical debt evaluation typically handled through isolated systems that do not interact effectively. AI-driven solutions have emerged with the promise of improved automation; however, existing implementations often focus on single-phase automation or exhibit significant limitations such as hallucinations, inconsistent reasoning, and limited explainability. These weaknesses erode user trust and hinder adoption in professional development environments. Moreover, current automation approaches lack the semantic synchronization and end-to-end traceability necessary to maintain consistency across evolving software artifacts. As a result, teams face substantial manual effort, quality inconsistencies, inefficient testing workflows, and escalating technical debt. These challenges underscore the need for an integrated, hybrid, and standards-aligned AI-enhanced SDLC automation framework capable of delivering reliable, explainable, and scalable support across all development phases while ensuring meaningful human oversight.

## 1.3 Software Requirements Specification

This section outlines the functional and non-functional requirements of *InteGrow*, an AI powered desktop application designed to unify the Software Development Life Cycle into one intelligent platform. It specifies the features, constraints, and performance criteria required for successful implementation.

### 1.3.1 Functional Requirements

This section defines the functional requirements for all major modules of **InteGrow**. Each requirement follows the IEEE 830 standard, describing what the system shall achieve to



ensure effective implementation.

#### 1.3.1.1 Core Foundation

The Core Foundation acts as the backbone of InteGrow, managing user authentication through GitHub, automating project creation, and maintaining seamless integration with remote repositories. It also serves as the central hub providing unified access to all modules within the project.

**FR-CORE-001:** The system shall require user authentication using GitHub OAuth before granting access to any platform features.

**FR-CORE-002:** The system shall allow users to create new projects that automatically initialize a remote GitHub repository linked with the local project.

**FR-CORE-003:** The system shall include an autonomous Git Agent that performs commits at predefined intervals on the remote branch.

**FR-CORE-004:** The system shall provide a centralized dashboard displaying all modules, features, and project information.

#### 1.3.1.2 Requirements Auditor

The Requirements Auditor module enables users to input requirements using natural language, refine them iteratively, and validate them for completeness, consistency, and fairness.

**FR-REQ-001:** The system shall allow users to enter requirements in natural language through an interactive chat interface.

**FR-REQ-002:** The system shall enable users to refine requirements by accepting, rejecting, or modifying AI-generated suggestions.

**FR-REQ-003:** The system shall allow users to export validated requirements in multiple formats while maintaining Git-based version control.

#### 1.3.1.3 UML Synthesizer

The UML Synthesizer transforms both structured and unstructured requirements into UML diagrams. It supports real-time editing, version tracking, and traceability through Git integration.

**FR-UML-001:** The system shall automatically generate UML diagrams using approved user requirements.

**FR-UML-002:** The system shall link UML components directly to their respective requirements for traceability.

**FR-UML-003:** The system shall maintain a detailed version history of all UML diagrams using Git.

**FR-UML-004:** The system shall provide an interactive interface that allows live editing of generated diagrams.

#### 1.3.1.4 Code-UML Synchronization

The Code-UML Synchronization module ensures alignment between design and implementation by maintaining a bidirectional link between source code and UML diagrams.

**FR-SYNC-001:** The system shall automatically update UML diagrams when changes occur in the corresponding source code.

**FR-SYNC-002:** The system shall generate code skeletons from UML diagrams while preserving manually written code.

**FR-SYNC-003:** The system shall allow users to configure synchronization settings for individual projects.

**FR-SYNC-004:** The system shall display real-time synchronization indicators within the user interface.

#### 1.3.1.5 Code Review Assistant

The Code Review Assistant provides automated, AI-assisted code reviews by combining static analysis with contextual reasoning. It is integrated with GitHub workflows to streamline the review process.

**FR-REVIEW-001:** The system shall perform code analysis using both static checks and large language model reasoning.

**FR-REVIEW-002:** The system shall generate automated recommendations for identified code issues.

**FR-REVIEW-003:** The system shall integrate with GitHub Pull Requests to provide inline feedback and approvals.

**FR-REVIEW-004:** The system shall provide a dashboard that allows users to view, sort, and manage all review comments and actions.

#### 1.3.1.6 Technical Debt Analyzer

The Technical Debt Analyzer monitors and predicts software maintainability through complexity metrics and machine learning analysis, helping developers reduce long-term code degradation.

**FR-DEBT-001:** The system shall calculate complexity metrics such as cyclomatic complexity and maintainability index.

**FR-DEBT-002:** The system shall use machine learning models to identify files likely to accumulate technical debt.

**FR-DEBT-003:** The system shall provide refactoring suggestions to address or prevent technical debt.

**FR-DEBT-004:** The system shall notify users when technical debt thresholds are exceeded.

#### 1.3.1.7 Test Generator

The Test Generator automates test creation and optimization, improving test coverage and code reliability through intelligent AI-driven techniques.

**FR-TEST-001:** The system shall automatically generate unit test cases using large language models based on analyzed source code.

**FR-TEST-002:** The system shall employ fuzz testing techniques to detect potential vulnerabilities and bugs.

**FR-TEST-003:** The system shall provide an integrated environment for organizing, executing, and tracking test suites.

#### 1.3.1.8 CI/CD Orchestration

The CI/CD Orchestration module automates continuous integration and deployment workflows to ensure efficient testing, integration, and release processes.

**FR-CICD-001:** The system shall automatically generate CI/CD pipelines according to the project structure.

**FR-CICD-002:** The system shall perform integration testing to verify compatibility.

### 1.3.2 Non-Functional Requirements

This section outlines the non-functional requirements of **InteGrow**. These requirements define essential quality attributes, including performance, security, reliability, usability, maintainability, and compatibility, ensuring that the system performs efficiently and remains scalable, secure, and user-friendly.

#### 1.3.2.1 Performance Requirements

**NFR-PERF-001:** The system shall ensure that API response times do not exceed 500 milliseconds, and large language model inference shall complete within 5 seconds for up to 1000 tokens.

**NFR-PERF-002:** The system shall support up to 10 API requests per minute for each user.

**NFR-PERF-003:** The system shall consume no more than 1 GB of RAM when idle and not more than 3 GB during active usage.

#### 1.3.2.2 Security Requirements

**NFR-SEC-001:** The system shall authenticate users using GitHub OAuth 2.0 and JSON Web Tokens (JWT) for secure access and authorization.

**NFR-SEC-002:** The system shall comply with GDPR standards and store any user data only with explicit user consent.

#### 1.3.2.3 Reliability Requirements

**NFR-REL-001:** The system shall maintain an uptime of at least 99.5%, excluding planned maintenance periods.

**NFR-REL-002:** The system shall automatically retry failed operations up to three times.

**NFR-REL-003:** The system shall provide user-friendly error messages and maintain detailed error logs for debugging purposes.

#### 1.3.2.4 Usability Requirements

**NFR-USE-001:** The system shall offer a consistent and modern user interface that enhances accessibility and ease of use.

**NFR-USE-002:** The system shall provide textual feedback or error messages within one second after any major user action.

#### 1.3.2.5 Maintainability Requirements

**NFR-MAIN-001:** The system shall adhere to standard coding practices and maintain at least 70% test coverage for backend components.

**NFR-MAIN-002:** The system's architecture shall support modular design, ensuring that modules are independent and loosely coupled for easier updates and scalability.

#### 1.3.2.6 Compatibility Requirements

**NFR-COMPAT-001:** The system shall support Git version 2.30 or higher and integrate seamlessly with GitHub and GitLab repositories.

**NFR-COMPAT-002:** The system shall allow importing and analyzing Git-based projects without modifying the original source files.

## 1.4 Work Distribution

Component	Description	Responsible Member(s)
<b>UML Synthesizer Module</b>		
Entity Extraction	NLP-based named entity recognition	Shuja Uddin
Diagram Generation	Rule-based + LLM-enhanced class model synthesis.	Amna Hassan
Diagram Validation	Ensuring syntactic and semantic correctness.	Amna Hassan
Visualization & Export	UI panels for editing, previewing, and exporting UML diagrams.	Muskan Tariq
<b>Requirements Module</b>		
Ambiguity Detection	Automated detection of ambiguous or unclear requirements.	Muskan Tariq
Requirement Fixing	Automated fixes for ambiguous, incomplete, or unethical requirements.	Shuja Uddin
User Story Generation	Conversion of validated requirements into actionable user stories with acceptance criteria.	Amna Hassan
Chat-Based Refinements	Interactive editing and refinement using conversational LLM workflows.	Muskan Tariq
Validation & Approval Flow	Final review and GitHub commit integration for approved requirements.	Muskan Tariq
<b>Other Project Components</b>		
Literature Review	Systematic study of prior research across AI-driven SDLC automation	Entire Team
Data Collection	Gathering requirement samples, datasets, and project repositories for testing.	Muskan Tariq, Shuja Uddin
Testing & Evaluation	Unit, integration, and system testing for requirement auditor and UML modules.	Amna Hassan, Shuja Uddin

Table 1.1: Combined Responsibilities and Component Breakdown

# Chapter 2

## Literature Review

This chapter critically examines existing literature relevant to InteGrow. It identifies key studies, theories, and findings, providing a foundation for the proposed research and highlighting its relationship to prior work.

### 2.1 Related Research

The main research topics in this section are directly related to InteGrow's goals and cover the crucial stages of the software development lifecycle that are automated by InteGrow's modules. AI-driven requirements analysis, code synchronization, AI-assisted code review, automated test case generation, technical debt prediction and model-driven development automation are all covered in the research. The thorough analysis of each item guides InteGrow's design and implementation tactics to guarantee creative, reliable, and efficient solutions.

#### 2.1.1 AI-Driven Requirements Analysis and Ambiguity Detection

##### 2.1.1.1 Summary of the Research Items

Kwizera [1] was the first to employ generative AI through structured prompt engineering. This showed that pre-trained LLMs can successfully recognize different kinds of ambiguity without the need for domain-specific training datasets. Verified on real-world requirements from Alstom and literature-based datasets, the study presents a web-based tool that integrates GPT-3.5 and GPT-4 and allows for continuous user feedback for ambiguity resolution through a chat-like interface.

Using an in-context learning paradigm, Bashir et al. (2024) [2] empirically examined LLMs for ambiguity detection and explanation in real-world industrial requirements. Their findings from three industrial datasets show that when presented with ten pertinent in-context demonstrations (10-shot), LLMs increase their average performance in classifying ambiguous requirements by 20.2 percent when compared to zero-shot methods. The practical efficacy of LLM-generated explanations was validated by human evaluations with eight industry experts, which produced an average rating of 3.84 out of 5 across the dimensions of naturalness, adequacy, usefulness, and relevance.

The efficiency of GPT-3.5 in automating crucial software requirements engineering tasks, such as requirement specification creation and quality assessment, was examined by Yeow et al. (2024) [3]. Experiments revealed that refined models such as CodeLlama frequently

approached or exceeded human-crafted specifications in completeness and internal consistency, highlighting how LLMs can greatly reduce the time spent refining and verifying SRS documents.

Chazette and Schneider [4] conducted a mixed-methods study that combined practitioner interviews and guideline analysis to examine the explainability and transparency requirements for AI systems from the viewpoints of users. According to their findings, the most important user requirement is transparency, and in order to foster trust in AI-assisted software development tools, explanations must be contextually aware.

### **2.1.1.2 Critical Analysis of the Research Items**

All of the reviewed research demonstrates significant methodological strengths, especially the use of prompt engineering and in-context learning, which improve automation in requirements analysis and lessen reliance on domain-specific datasets[3]. Credibility is increased by validation across industrial domains like healthcare, transportation, and automobiles, and hybrid approaches that combine LLMs with NLP and heuristic techniques successfully strike a balance between automation and accuracy[2]. Prompt sensitivity, domain bias, and limited generalizability are still problems in spite of these improvements. The focus on healthcare applications limits broader applicability across domains, and studies that heavily rely on user reviews run the risk of sampling bias. Long or complicated requirements continue to be a challenge for most methods [4]. Overall, even though automation has advanced, there are still issues with guaranteeing robustness, scalability, and ethical compliance in a variety of industrial contexts.

### **2.1.1.3 Relationship to the Proposed Research Work**

The design of InteGrow’s Requirements Auditor is firmly supported by this body of research, which shows that LLM-driven automation is feasible for evaluating and improving software requirements [3]. Effective ambiguity detection and optimized few-shot prompting without extensive fine-tuning are made possible by InteGrow’s NLP pipeline, which is directly informed by Kwizera’s structured prompt framework and Bashir et al.’s in-context learning insights. It will also prioritize user interpretability and transparent outputs, drawing inspiration from Chazette and Schneider, to guarantee that analysts comprehend and have faith in the findings[2].

## **2.1.2 Requirements to Design Transformation**

### **2.1.2.1 Summary of the research items**

Gala (2023) [5] discusses research at Alstom AB in Sweden, which focuses on automating the creation of UML use case diagrams from natural language requirements in an industrial railway context. In order to identify system actors and components from operational scenarios, the study used a Named Entity Recognition (NER) model based on SpaCy and multiclass classification. The method generated UML diagrams in roughly 7 seconds as opposed to minutes when done by hand, and it achieved 98 percent precision and recall. The study highlights the crucial role that domain expertise plays in validating and improving requirements by emphasizing semi-automation with human oversight.

Meng Ban (2024) [6] address ambiguity and unstructured inputs by proposing a four-step natural language processing (NLP) framework for automatic UML class diagram

generation from textual requirements. They use dependency parsing to generate rule-based diagrams, preprocessing, syntactic analysis, and sentence classification. With an AUC of 0.9287 and an accuracy of 88.46 percent, the model demonstrated its resilience in detecting class relationships. The study does point out, though, that current methods have trouble accurately translating complex requirements, particularly for diverse and unrestricted domains.

An AI-driven method for semi-automated software architecture generation from natural language requirements is presented by Eisenreich et al. (2024) [7]. Using LLMs and quantitative analysis, the approach generates use cases, domain models, and several architecture candidates that are then iteratively improved through ATAM-based assessments. Time constraints in architectural design are addressed, but issues like model accuracy, hallucinations, ethics, and a lack of domain-specific datasets for validating AI-generated architectures remain.

In order to facilitate the modernization of legacy systems, Bates et al. (2025)[8] present a multimodal LLM-based method for producing executable UML code from image-based diagrams. The approach obtained BLEU and SSIM scores of 0.779 and 0.942 using optimized LLaVA-1.5 models with LoRA that were trained on synthetic UML activity and sequence diagrams. This method reduces the amount of manual labor required to maintain and update legacy software systems by expediting the conversion of visual documentation into machine-readable specifications.

### 2.1.2.2 Critical analysis of the research items

The reviewed studies represent significant advancements in the automation of requirements-to-UML transformation using hybrid approaches that combine rule-based heuristics, machine learning, and natural language processing. Due to its small, domain-specific dataset, Gala's high-precision NER-based method (98 percent)[5] exhibits limited generalizability despite strong industrial validation. Although preprocessing and classification are successfully handled by Meng Ban's four-step framework, it still has trouble with complicated or unclear requirements. Although AI-driven architecture generation provides a forward-thinking approach, it is beset by problems such as the lack of strong evaluation frameworks, ethical lapses, and hallucination risks [8]. Although Bates' multimodal approach has good quantitative results, it lacks bidirectional synchronization, which is essential for iterative development, and instead concentrates solely on diagram-to-code conversion. Practical adoption in a variety of software environments is limited by the majority of research's lack of cross-industry validation and failure to address full round-trip engineering.

### 2.1.2.3 Relationship to the proposed research work

By offering verified foundations and highlighting important areas for improvement, this study directly supports InteGrow's UML Synthesizer. Meng Ban's [6] sentence classification techniques and Gala's entity extraction provide flexible ways to expand InteGrow's class diagram generation. InteGrow's innovation focus is on ensuring consistency between changing UML diagrams and code through true round-trip engineering, which is defined by the absence of bidirectional synchronization across studies. While Bates' multimodal success [8] raises the possibility of image-based inputs. All of these results support InteGrow's objective of providing an integrated requirements-to-code pipeline that goes

beyond the limits of existing research and blends automation, flexibility, and ethical dependability.

### **2.1.3 Design-to-Code Generation**

#### **2.1.3.1 Summary of the research items**

Multimodal transformer models for transforming UML diagrams into executable Plan-tUML code were investigated by Bates et al. (2025) [8]. The potential of visual-to-code generation was demonstrated by the BLEU 0.78 and SSIM 0.94 obtained by fine-tuning LLaVA-1.5 with LoRA. The authors intend to expand it for code-to-diagram generation, even though it is currently one-way, to allow for full round-trip engineering and IDE integration.

Conrardy (2025) [9] describe the use of LLMs for image-based UML diagram extraction. Their prototype creates class and sequence diagram XMI by processing PNG and JPEG UML diagrams using an OCR-augmented GPT pipeline. Its 82 percent element detection precision and 79 percent relationship accuracy, as tested on 500 diagrams, represent a first step toward fully automated reverse engineering from unstructured images.

A template-driven MDA toolchain that converts UML class and state machine models into Java and C++ code is presented by D. Salunke (2024)[10]. In three industrial case studies, the authors show the useful productivity gains of MDA approaches by utilizing model-to-text transformations in Acceleo to reduce development time by 35 percent and integration defects by 25 percent.

J . Navajas (2024)[11] investigate the creation of code for classical-quantum systems using UML. Their transformation pipeline generates Q and Python code skeletons for hybrid classical-quantum applications by extending UML profiles with quantum stereotypes. The difficulties in mapping concurrency semantics are highlighted by case studies in quantum chemistry and cryptography, which demonstrate 78 percent compliance with the intended model behavior.

Antal et al. (2024)[12] evaluate GPT-4-Vision’s UML-to-code capabilities by examining how well it performs on 200 class diagrams that are generated automatically. The study uses compilation success (91 percent ) and unit test pass rates (8 percent ) to gauge code correctness. The authors pinpoint areas for future model improvement by pointing out shortcomings in complex inheritance and generic types but highlighting strengths in mapping simple associations through thorough error analysis.

#### **2.1.3.2 Critical analysis of the research items**

By combining formal bidirectional transformations, machine-learning adaptations, and human-in-the-loop validation, these studies collectively improve round-trip engineering by lowering manual labor and preserving model-code consistency. Together, these studies show that automated diagram-to-code generation is feasible using both rule-based MDA and new LLM/vision methods. Acceleo and other MDA toolchains offer consistent productivity increases, but they struggle with non-standard UML extensions and demand a significant upfront modeling effort. Though they are more flexible in managing a variety of input formats, vision+LLM approaches ( Conrardy Cabot) [9] have higher error rates when it comes to relationship extraction and complex semantics. The difficulty of adapting conventional UML-to-code mappings to new paradigms is highlighted by quantum-specific MDA ( J . Navajas)[11]. All things considered, combining MDA accuracy with



LLM flexibility is still an unexplored area of study.

### 2.1.3.3 Relationship to the proposed research work

By verifying several synchronization techniques human-in-the-loop confirmation for accuracy, bidirectional lenses for formal correctness, and GNN-based adaptation for changing codebases these insights directly inform InteGrow’s UML Synthesizer module. InteGrow can use lens-based method to ensure semantic preservation in crucial components, incorporate Singh et al.’s incremental learning model to minimize manual rule specification, and EMF/QVT framework as the foundation for transformation rule management. Given the scalability issues that have been shown, InteGrow requires modular, performance-optimized implementations that minimize user configuration overhead and guarantee responsive synchronization across Python and Java environments.

## 2.1.4 AI-Enhanced Automated Code Review

### 2.1.4.1 Summary of the research items

To improve review accuracy and lessen hallucinations, Icöz et al. (2025) [13] presented a hybrid code review system that combines symbolic linters and LLMs. The model integrates into CI pipelines for automated reviews, analyzes pull request diffs, and verifies recommendations using rule-based checks. Tested on three sizable Java repositories, it reduced false positives by 22 percent and increased precision by 16 percent, providing automated approvals for straightforward cases while highlighting more complicated ones for human review.

Static analysis tools and artificial intelligence are combined in automated code review to find flaws, security flaws, and problems with code quality. In their neuro-symbolic approach, Jaoua et al. (2025) [14] show how result fusion algorithms can reduce false positives by 20–30 percent by methodically integrating large language models with static analyzers (SonarQube, ESLint). Their approach deduplicates results, runs analyzers concurrently, and adds contextual explanations produced by LLM to static alerts.

Control-flow graphs and type information are examples of static analysis metadata that can be incorporated into LLM prompts, according to S. M. Abtahi (2025) [15]. After fine-tuning a transformer model using this enriched input, the authors report a 20 percent increase in real defect detections across Python and JavaScript repositories and a 28 percent decrease in hallucinated suggestions.

A web-based platform that combines ESLint, PyLint, and GPT-3.5 is described by M. Mohanakshi (2025) [16]. The system provides severity rankings, inline code comments, and a feedback loop where developers can rate ideas to improve reviews in the future. Developer satisfaction with the clarity of review feedback increased by 40 percent, and review cycle time was reduced by 35 percent in a pilot involving 12 startups.

Z. Rasheed (2024) [17] assesses CodeLlama and GPT-3.5 for creating pull request comments on 200 open-source repositories. With an average BLEU of 0.42 and ROUGE-L of 0.48, the study evaluates BLEU and ROUGE metrics for comment relevance. 60 percent of AI-generated comments were judged as "useful" by human judges, indicating difficulties with context comprehension and the requirement for repository-specific fine-tuning.

### 2.1.4.2 Critical analysis of the research items

In comparison to LLM-only or rule-based systems, these studies show that hybrid approaches combining LLMs with static analyzers significantly improve review precision

and lower false positives. Both Icöz et al. [13] and Jaoua et al.[14] rely on proprietary models and need to be integrated into continuous integration pipelines, but they both empirically confirm improvements in actionable recommendations and decreased reviewer effort. Although it lacks extensive industrial validation emphasizes the technical advantages of incorporating static metadata into prompts. Although the web-based solution is excellent for developer feedback loops, the overhead of customization may prevent it from scaling to large enterprise codebases. The need for domain adaptation and ongoing improvement is highlighted by the preliminary findings with GPT-3.5 and CodeLlama, which indicate promise in comment generation but limitations in relevance and trust.

#### **2.1.4.3 Relationship to the proposed research work**

By confirming that the hybrid LLM-linter architecture strikes a balance between automation and safety, this study supports InteGrow’s Code Review Assistant. InteGrow’s data pipeline will be informed by the prompt-enrichment technique which embeds control-flow and type metadata into review prompts for a more thorough understanding of the code[17]. To ensure that AI-generated comments are still pertinent, useful, and trusted by development teams, repository-specific fine-tuning and evaluation metrics will be guided by the lessons learned from early assessments of GPT-3.5 and CodeLlama[14]. The model’s performance will be iteratively improved by developers using a web-based feedback loop akin to this one. With these insights, InteGrow will be able to provide a code review automation module that is dependable, minimally disruptive, and flexible enough to accommodate a variety of codebases and developer workflows.

### **2.1.5 LLM-Powered Test Case Generation**

#### **2.1.5.1 Summary of the research items**

LLM4Fin, a pipeline that completely automates test case generation for financial software acceptance testing, is presented by Xue et al. (2024) [18]. After ingesting transaction scenarios and business rule descriptions, the system generates Python executable test scripts using a refined LLaMA model. Human evaluators gave the generated cases a 4.2/5 rating for adequacy and realism, while LLM4Fin achieved 92 percent requirement coverage and a 65 percent reduction in test design time when tested on three banking applications.

LLM-Powered Test Case Generation for Detecting Bugs in Plausible Programs, presented by Liu et al. (2024) [19], uses GPT-4 to create test inputs for code submitted by students. In comparison to rule-based fuzzers, the method increases bug-finding recall by 48 percent by presenting bug detection as a few-shot learning task with annotated examples. The generated tests showed high efficiency in automated defect detection, catching 85 percent of seeded faults in 10 runs on a dataset of 2,000 programs.

Dantas (2024) [20] suggests a modular framework called Large Language Model Powered Test Case Generation for Software Applications, which combines domain-specific test oracles with prompt-engineering templates. In four case studies covering the e-commerce and healthcare domains, the system achieved 88 percent functional coverage by generating JUnit test suites from SRS documents and application models. Standardized test scaffolding reduced manual test maintenance effort by 30 percent, according to user feedback.

Mutation-Guided LLM-Based Test Generation at Meta is described by Foster et al. (2025) [21]. In this approach, an LLM creates potential test cases that are iteratively improved

with feedback from mutation analysis. On three extensive open-source projects, their closed-loop pipeline produced a mutant kill rate that was 35 percent higher than that of standard LLM output. According to the study, test efficacy for complex code paths is greatly increased when mutation scores are incorporated into LLM prompts.

Kumari (2024) [22] presents Intelligent Test Automation, a multi-agent framework in which specialized LLM agents dynamically create, verify, and rank test cases. To discuss GUI, API, and load testing scenarios, agents interact through a shared workspace. In tests using a retail web application, the system produced more than 200 distinct test cases in less than 15 minutes, resulting in 94 percent automation coverage and a 40 percent reduction in the execution time of the regression suite.

### 2.1.5.2 Critical analysis of the research items

Together, these studies show that LLM-based test generation can significantly cut down on test design time while enhancing coverage in a variety of domains. Despite their reliance on high-quality SRS inputs, LLM4Fin and Dantas [20] exhibit strong practitioner acceptance in healthcare and financial contexts. Liu et al. [19] confirm that incorporating domain-specific guidance, such as mutation feedback or few-shot examples, improves defect detection but complicates prompt management. Although Kumari’s [22] multi-agent approach offers extensive automation, it necessitates strong agent orchestration and may result in coordination overhead. Dependency on proprietary models, possible hallucinations in test logic, and difficulties validating performance or non-functional tests are examples of common limitations.

### 2.1.5.3 Relationship to the proposed research work

The Test Generator module from InteGrow will use a hybrid LLM and feedback-driven architecture, fusing Foster et al.-inspired mutation-guided refinement with the prompt templates from LLM4Fin. Few-shot learning techniques from Liu et al. and domain-specific test oracles from Dantas [20] will be combined to guarantee high coverage and usefulness. Our design for parallel test generation and validation agents will be guided by Kumari’s [22] multi-agent orchestration pattern. When combined, these methods will allow InteGrow to provide automated, precise, and maintainable test suites covering security, performance, and functional aspects.

## 2.1.6 Technical Debt Prediction and Prioritization

### 2.1.6.1 Summary of the research items

Tsoukalas et al. (2024)[23] use time series forecasting on code metrics to propose a class-level technical debt prioritization framework. Their approach ranks classes according to expected debt accumulation by using ARIMA and LSTM models to model code complexity, churn, and defect density over several releases. The method, which was tested on four open-source Java projects, reduced maintenance effort estimates by 22 percent and identified high-priority debt modules with 87 percent accuracy.

Machine learning methods for identifying and categorizing self-admitted technical debt (SATD) in code comments are empirically studied by Melin et al. (2023) [24]. They discovered that SATD density correlates with a 1.4× increase in post-release defects and trained BERT-based classifiers on 12,000 annotated comments, achieving a 91 percent F1-score in SATD detection. The study emphasizes how beneficial SATD mining is for proactive debt management.

Using issue trackers and code repositories from 60 data science projects, Akman (2024)[25] investigates technical debt in ML-based projects. With a 0.82 ROC-AUC for ML debt detection, the study finds two types of debt that are specific to machine learning: data preprocessing debt and model parameter debt. Random forest models are then trained on feature vectors that combine code metrics and pipeline logs. The findings show that during the early stages of prototyping, ML projects accrue debt more quickly.

Sas and Kazman (2023)[26] present the Architectural Technical Debt Index (ATDI), which is based on code metrics and architectural odors and is calculated through machine learning. With a mean absolute error of 3.2 person-days, their gradient boosting model forecasts remediation effort and has been validated on three enterprise systems. The ATDI provides an interpretable scoring system for architectural debt planning and has a strong correlation ( $=0.76$ ) with manager-assessed debt priority.

### **2.1.6.2 Critical analysis of the research items**

From general code-centric models to domain-specific and architectural contexts, the reviewed works progress debt prediction. Although it lacks empirical support, Ajibode's [27] mapping demonstrates that ensemble approaches perform better than isolated techniques. Although class-level debt is accurately prioritized in Tsoukalas [23] forecasting, it relies on historical release data, which is frequently unavailable in new projects. Melin [24] show good SATD detection accuracy, but their coverage of undocumented debt is limited by their reliance on comment quality. Although random forests might miss intricate pipeline interactions, Akman draws attention to ML-specific debt patterns. Although the ATDI by Sas and Kazman is interpretable, it necessitates a great deal of architectural smell instrumentation, which can be resource-intensive for legacy systems.

### **2.1.6.3 Relationship to the proposed research work**

Through the adaptation of Tsoukalas [23] class-level prioritization for module ranking, these findings guide the integration of ensemble learning and time-series forecasting for early debt prediction in InteGrow's Technical Debt Analyzer. Melin [24] SATD detection methods will highlight developer-admitted debt hotspots in the user interface, while Akman's [25] ML-specific debt categories will expand detection capabilities to AI/ML pipelines. The ATDI framework developed by Sas and Kazman will be modified to produce interpretable architectural debt scores, guaranteeing remediation plans that are in line with manager evaluations.

## **2.1.7 Model-Driven Development and SDLC Automation**

### **2.1.7.1 Summary of the research items**

The impact of model-driven development (MDD) on agile practices in knowledge-intensive engineering is examined by Aghakhani et al. (2024) [28]. According to aerospace and defense case studies, MDD reduces integration defects by 25 percent and enhances design-time validation by 30 percent. In order to better match evolving models with shifting agile requirements, they emphasize toolchain interoperability and suggest a metamodel extension.

A model-driven engineering (MDE) framework for LLM-based applications is presented by C. Bolufer (2025) [29]. The method automatically generates orchestration and testing code by extending UML profiles to define prompt setups and data flows. It demonstrated

the usefulness of MDE in AI/ML systems by reducing manual coding by 40 percent and increasing test coverage to 85 percent when applied to LLM microservices.

Six round-trip engineering tools are compared by Rosca et al. (2023) [30] in terms of usability, performance, and synchronization. They believe that open-source tools are more adaptable and commercial tools are more scalable. Weak support for bidirectional transformations of behavioral UML elements is a major drawback, which has led to a call for unified transformation languages.

By converting UML/OCL models into Dart/Flutter implementations, Cheon and Y. (2025) [3] assess LLMs as code generators in MDD. In contrast to human-written code, GPT-4 can produce scaffolding code with 94 percent structural consistency and accurately enforce OCL constraints, resulting in a 60 percent reduction in manual labor, according to their case study on a Sudoku application. The authors suggest incorporating automated model checkers into the generation loop after pointing out difficulties in confirming semantic correctness.

#### 2.1.7.2 Critical analysis of the research items

Together, these studies demonstrate the advantages and difficulties of MDD in contemporary SDLC procedures. Aghakhani et al. [28] highlight tool interoperability problems that can impede smooth adoption while demonstrating observable productivity gains in agile contexts. Although their UML profile extensions might not be able to keep up with the rapid evolution of AI frameworks, extend MDE into AI/ML integration [29], demonstrating its suitability for LLM-based applications. Organizational and educational obstacles to MDD adoption are revealed by Barigidad's practitioner survey, indicating that strategic training programs must be combined with technical solutions. Full round-trip automation is hampered by the lack of thorough bidirectional support for behavioral models, which is highlighted in Rosca et al.'s tool comparison.

#### 2.1.7.3 Relationship to the proposed research work

By confirming model-driven methods and pointing out important integration and correctness issues, these insights help *InteGrow* achieve its SDLC Automation objectives. The design of *InteGrow*'s metamodel extensions is guided by Aghakhani et al. [28] to facilitate interoperability among modules and agile iteration. Schema definitions for AI/ML components in *InteGrow*'s model-driven pipelines are inspired by the LLM integration profiles of Carreño-Bolufer et al. [29]. Our user adoption approach is shaped by Barigidad's findings, which highlight modular tool architectures and integrated training modules as ways to reduce the learning curve.

## 2.2 Analysis Summary of Research Items

This section summarizes and critically compares the reviewed studies across major research areas relevant to *InteGrow*. Each study is analyzed in terms of its methodological approach, contributions, limitations, and specific relevance to the proposed framework in below Table 2.1.

## 2.3 Research Scope / Gap

The scope of this research focuses on developing *InteGrow*, an integrated, AI-driven Software Development Life Cycle (SDLC) automation platform. *InteGrow* unifies requirements analysis, UML modeling, code generation, automated code review, test generation,



## 2. Literature Review

Study	Artifact	Methodology	Key Contributions	Critical Analysis	Limitations	Relevance to InteGrow
1 [1]	Requirement Analysis	Structured prompt engineering using GPT-3.5/4 for ambiguity detection.	LLMs detect ambiguity without domain-specific datasets, validated on industrial requirements.	Strong empirical grounding and innovative prompt design; introduces continuous feedback loop.	Prompt sensitivity and domain bias reduce generalizability.	Foundation for <b>Requirements Auditor</b> —ambiguity detection and feedback refinement.
2 [2]	Requirement Analysis	In-context (10-shot) learning for ambiguity classification.	Improved ambiguity classification by 20.2%.	Industry validation; robust evaluation design.	Limited dataset size and reproducibility.	Informs few-shot NLP prompting in <b>Requirements Auditor</b> .
3 [5]	UML Synthesizer	SpaCy-based NER for actor/component extraction.	Achieved 98% precision/recall in UML generation.	High precision but limited domain adaptability.	Domain-restricted to railway sector.	Guides <b>UML Synthesizer</b> entity extraction.
4 [6]	UML Synthesizer	Rule-based four-step NLP pipeline for class diagram generation.	88.46% accuracy in class relationship detection.	Strong syntactic parsing; weak semantic coverage.	Fails with ambiguous inputs.	Supports syntactic-semantic hybrid for UML consistency.
5 [7]	Bidirectional Code-to-UML	LLM-driven architecture generation using ATAM feedback.	Semi-automated iterative architecture refinement.	Novel LLM + ATAM integration; creative approach.	Hallucination risk; lacks domain datasets.	Strengthens design-to-code explainability oversight.
6 [8]	Bidirectional Code-to-UM	Multimodal LLaVA-1.5 for visual-to-code UML translation.	0.94 SSIM; efficient modernization of legacy code.	High scalability; strong performance.	One-way flow; lacks bidirectional sync.	Guides <b>UML-Code Consistency</b> mechanism.
7 [13]	Code Review	Hybrid LLM + symbolic linters for code reviews.	Reduced false positives by 22%.	Balanced automation with rule-based verification.	Requires CI integration; uses proprietary models.	Shapes <b>Code Review Assistant</b> hybrid pipeline.
8	Code Review	Neuro-symbolic LLM + static analyzers integration.	Reduced false positives by 20–30%.	Fusion potential demonstrated with high reproducibility.	Heavy concurrent analyzer resources.	Validates neuro-symbolic review model.
9 [18]	Test Generation	LLaMA-based financial test generation.	92% coverage, 65% reduction in design time.	Efficient and domain-adapted.	Financial-only focus; closed-source data.	Framework for domain-aware test generation.
10 [19]	Test Generation	Few-shot GPT-4 test generation for educational code.	85% defect recall; effective few-shot training.	Excellent few-shot setup; practical for education.	Narrow dataset.	Informs prompt-based mutation refinement testing.
11 [23]	Technical Debt Predictor	LSTM + ARIMA for technical debt prioritization.	Predicted high-debt modules with 87% accuracy.	Empirically validated, interpretable model.	Needs extensive history data.	Basis for debt forecasting engine.
12 [24]	Technical Debt Analyzer	BERT-based classifier for SATD mining.	91% F1-score for debt detection.	High accuracy; broad empirical grounding.	Dependent on comment quality.	Supports proactive debt hotspot detection.
13 [26]	Technical Debt Analyzer	Gradient boosting for architectural debt index.	Predicted remediation effort (MAE = 3.2 days).	Interpretable, metrics-driven results.	Resource-intensive instrumentation.	Basis for Architectural Debt Index visualization.
14 [28]	MDD-Agile Framework	MDD toolchain interoperability for agile contexts.	Reduced defects by 25%, validation time +30%.	Demonstrates MDD–Agile synergy.	Limited to aerospace/defense.	Strengthens Agile–MDD alignment.
15 [29]	MDE-LLM Framework	MDE framework extending UML for LLM orchestration.	Reduced manual coding by 40%.	Innovative MDE–LLM combination.	Small dataset; early-stage validation.	Guides MDD-driven automation.

Table 2.1: Detailed Analysis Summary of Research Items

continuous integration workflows, and technical debt prediction into a single, intelligent ecosystem.

The study concentrates on:

- transforming natural-language requirements into consistent UML models,
- generating semantically correct code from design artifacts,
- enhancing automated code review through hybrid AI techniques,
- improving testing depth via LLM-guided test generation, and
- predicting and prioritizing technical debt using machine learning.

InteGrow will be evaluated across multiple projects to determine its impact on development time, code quality, design consistency, testing effectiveness, and proactive technical debt management within modern CI/CD pipelines.

### 2.3.1 Research Questions

1. How accurately can InteGrow transform unstructured natural language requirements into consistent and semantically valid UML models using hybrid LLM-based and symbolic reasoning techniques?
2. How effectively can the system ensure design-to-code consistency across multiple programming languages using model-driven transformations and LLM-assisted semantic preservation?
3. To what extent does the hybrid AI code review module reduce false positives and improve defect detection compared to traditional static-analysis-based tools?
4. How adequate and comprehensive are LLM-generated test suites in achieving high functional coverage and detecting meaningful faults?
5. How accurately can the technical debt prediction engine forecast and prioritize debt hotspots using ensemble machine learning models?

### 2.3.2 Research Objectives

1. To develop a requirements auditing module capable of identifying ambiguities and producing clear, well-structured requirement specifications.
2. To design and implement a bidirectional UML synthesis engine that converts requirements into standardized UML diagrams and synchronizes them with source code.
3. To build a semantically correct, multi-language code generation pipeline integrating Model-Driven Architecture (MDA) rules with LLM-based logic refinement.
4. To develop a hybrid AI code review system that combines LLM reasoning with symbolic static analysis to improve precision and reduce false positives.
5. To create a domain-aware test generator that produces high-coverage and fault-revealing test suites.
6. To construct an interpretable technical debt prediction and prioritization model using time-series forecasting, BERT-based classification, and architectural metrics.

# Chapter 3

## Proposed Approach

### 3.1 Overview

Literature review (2020-2025) reveals critical gaps in SDLC automation: while AI-driven tools achieve strong performance in isolated phases (97% ambiguity detection [31], 89% code review accuracy [32]), they operate in silos. This fragmentation causes 32.7% time loss in CI/CD tasks [33], 30% cost increase [34], and integration challenges. Additionally, 92% of technical debt models ignore cross-phase SDLC impact [35], and 48.6% of test generation relies on outdated legacy tests [36]. InteGrow addresses these gaps through an integrated AI-powered platform unifying six SDLC modules with shared context transmission, automatic traceability, and bidirectional synchronization.

### 3.2 Proposed System: InteGrow

InteGrow is a desktop application integrating six independent modules: Requirements & Ethics Auditor, UML Synthesizer, Code Review Assistant, Test Generator, Technical Debt Analyzer, and CI/CD Orchestration. The system employs:

- **Architecture:** Electron + Next.js frontend, FastAPI backend, Supabase persistence, Redis caching
- **Integration:** Context-aware cross-module feedback loops, automatic requirements-to-deployment traceability
- **Version Control:** GitHub integration with autonomous branch management and milestone-driven commits
- **AI Frameworks:** LangGraph (state machines), CrewAI (multi-agent teams), AutoGen (tool collaboration)

Figure 3.1 illustrates the system architecture.

### 3.3 Module Methodologies

#### 3.3.1 Module 1: Requirements & Ethics Auditor

**Framework:** LangGraph (State-Machine Orchestration) **Sub-Agents:**

- **Parser:** spaCy + BERT for entity extraction (actors, actions, constraints)
- **Ambiguity Detection:** LLaMA-2 7B with IEEE-830 aligned prompts
- **Completeness:** Mistral-7B validates preconditions, error cases, NFRs



- **Ethics Auditor:** AIF360/Fairlearn for bias detection

### 3.3.2 Module 2: UML Synthesizer

**Framework:** CrewAI (Multi-Agent Pipeline) **Agents:**

- **Entity Extraction:** Identifies classes, attributes, methods via POS tagging
- **Structure Analysis:** Infers inheritance, associations, cardinality
- **Rendering:** Generates PlantUML/Mermaid diagrams
- **Validation:** Syntax checks and completeness verification

**Features:** Sketch-to-UML (multimodal vision-language models), Neo4j traceability

### 3.3.3 Module 3: Code-UML Round-Trip Synchronization

**Framework:** AutoGen (Tool-Augmented Collaboration) **Components:**

- **AST Parsing:** Tree-sitter for multi-language support (Python, Java, TypeScript, C++, Go, Rust)
- **UML Update:** Code → UML propagation (<5s latency)
- **Code Generation:** UML → Code with CodeLlama/StarCoder
- **Conflict Resolution:** 3-way diff merging (base, UML, code)

**Triggers:** Reactive (save events, Git hooks), Proactive (scheduled sync)

### 3.3.4 Module 4: Code Review Assistant

**Framework:** Hybrid (Static Analysis + LLM Reasoning) **Methodology:**

- **Static Analysis:** Pylint, Bandit, Radon, ESLint, SonarQube, Semgrep
- **LLM Review:** Code Llama-34B for logic flaws, performance issues
- **Fusion Engine:** Deduplicates findings, ranks severity, filters false positives
- **Auto-Fix:** Generates fixes for formatting, imports, type hints

### 3.3.5 Module 5: Technical Debt Analyzer

**Framework:** Metrics-Driven + ML Forecasting **Components:**

- **Metrics Engine:** Complexity, maintainability index, duplication, code smells
- **ML Prediction:** Random Forest Regressor on commit history, churn, test coverage
- **Recommendations:** Pattern detection + LLM-generated refactoring plans
- **Visualization:** Treemaps, trends, dependency graphs, hotspot views

### 3.3.6 Module 6: Test Generator

**Framework:** Multi-Strategy (LLM + Genetic Algorithms + Fuzzing) **Strategies:**

- **LLM-Driven:** Deterministic unit tests via LLM reasoning
- **Genetic Algorithms:** Maximizes branch/path coverage via mutation/crossover
- **RL-Guided Fuzzing:** Enhances AFL++/libFuzzer with RL mutations
- **Flaky Detection:** Automatic identification of unstable tests

### 3.3.7 Module 7: CI/CD Orchestration

**Framework:** Reinforcement Learning-Optimized Pipelines **Components:**

- **Auto-Generation:** Analyzes repository structure for GitHub Actions workflows

- **RL Optimization:** PPO algorithm optimizes caching, job ordering, parallelism
- **Self-Healing:** Resolves flaky tests, dependency conflicts, network failures
- **Integration Testing:** Containerized validation environments

### 3.3.8 Module Summary

Table 3.1: InteGrow Module Methodologies Summary

Module	Framework
Requirements & Ethics	LangGraph + LLaMA-2/Mistral
UML Synthesizer	CrewAI + PlantUML/Mermaid
Code-UML Sync	AutoGen + Tree-sitter
Code Review	Static + Code Llama-34B
Technical Debt	Random Forest + Metrics
Test Generator	LLM + GA + RL Fuzzing
CI/CD Orchestration	RL (PPO) + Self-Healing

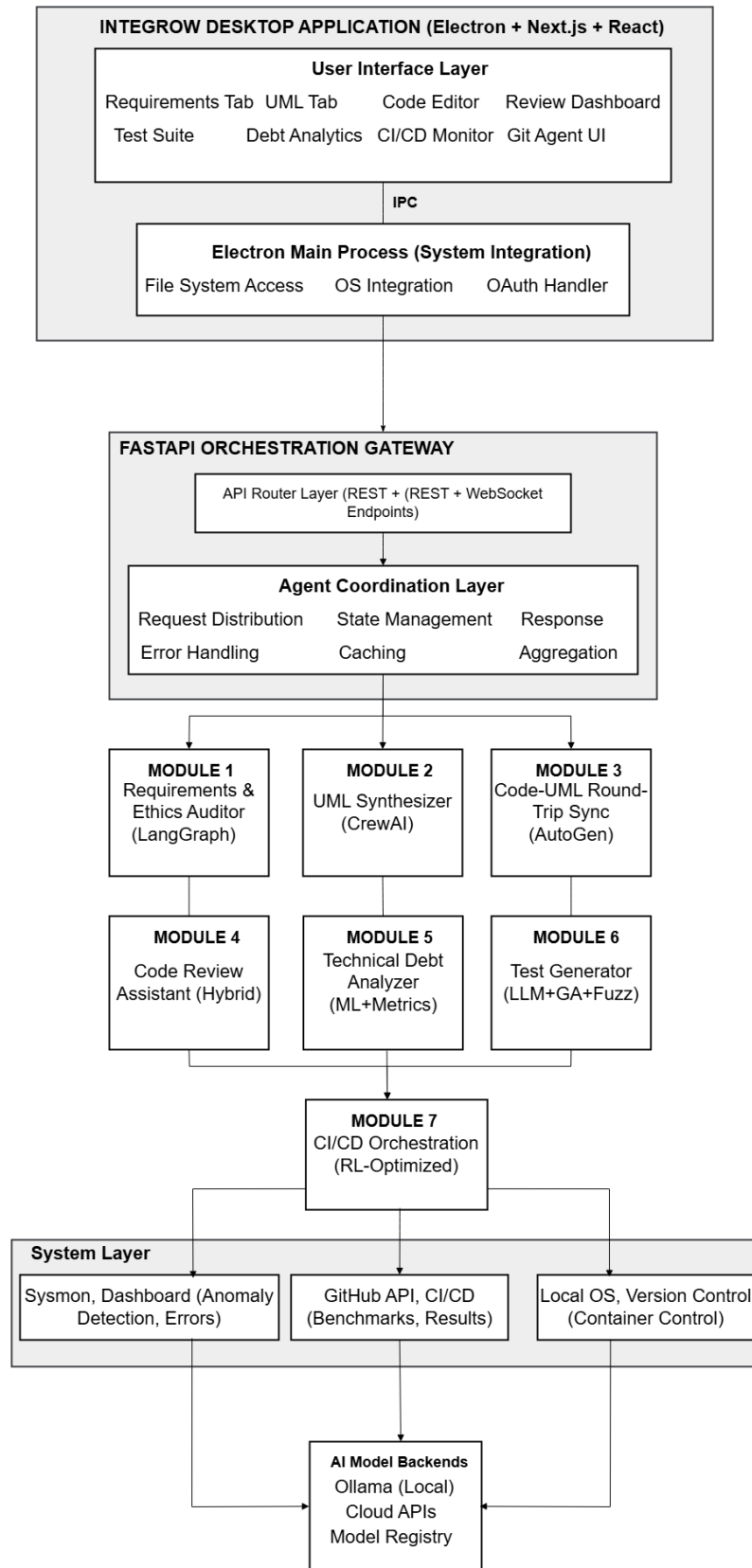


Figure 3.1: InteGrow System Architecture

# Chapter 4

## Initial Experiments and Results

This chapter presents the user acceptance testing (UAT) of InteGrow Phase 1 using Mean Opinion Score (MOS) methodology following ITU-T P.800 standards. The evaluation validates InteGrow’s alignment with IEEE 830 quality requirements and SDLC best practices through user experience assessment and productivity measurement.

### 4.1 Experimental Setup

#### 4.1.1 Hardware Configuration

The development and testing environment consisted of:

- **Development Workstation:** Intel Core i7-12700K processor, 32 GB DDR4 RAM.
- **Testing Devices:** Windows 11 desktop
- **Network:** Stable broadband connection (100 Mbps) for cloud API access

#### 4.1.2 Software Environment

The software stack included:

- **Frontend:** Electron 28.0, Next.js 14.2, TypeScript 5.3, React 19 RC, Monaco Editor 0.45
- **Backend:** Python 3.10.12, FastAPI 0.109, LangGraph 0.0.25, spaCy 3.7 (en\_core\_web\_sm model)
- **AI Services:** Groq API (LLaMA 3.3 70B Versatile), Google Gemini 2.5 Flash API
- **Database:** Supabase (PostgreSQL 15.1) with Row Level Security policies
- **Cache Layer:** Redis 7.2 (Docker container)
- **Version Control:** Git 2.42, GitHub API v3

#### 4.1.3 Participant Demographics

Participants were recruited from the university’s software engineering program. All participants had minimum 2 years of experience in software development or project management and familiarity with IEEE 830 standards and requirements engineering practices.

#### 4.1.4 UAT Protocol

Each 15-20 minute session followed a standardized protocol:

1. **Introduction (5 min):** Demonstration of key features (Monaco editor, multi-agent analysis, AI chat, analysis panel, export functionality)
2. **Tutorial (5 min):** Hands-on practice with sample requirement
3. **Evaluation (5-10 min):** Independent analysis of 5-8 requirements from healthcare and finance domains
4. **MOS Questionnaire (5 min):** Rating across five dimensions on 5-point Likert scale
5. **Interview (5 min):** Qualitative feedback on features, improvements, time savings, and production readiness

#### 4.1.5 Test Dataset

The test pool contained software requirements: healthcare (CCHIT standards for patient demographics, medication management, clinical documentation) and finance (payment processing, transaction management, merchant operations). Requirements represented varying complexity and quality levels to ensure comprehensive evaluation.

## 4.2 Evaluation Metrics

### 4.2.1 Mean Opinion Score (MOS) - ITU-T P.800

MOS methodology employs a 5-point Likert scale (1=Bad, 2=Poor, 3=Fair, 4=Good, 5=Excellent) across five dimensions:

- **MOS-U (Usability):** Ease of use, interface intuitiveness, learning curve
- **MOS-A (Accuracy):** Correctness of AI suggestions and detected issues
- **MOS-R (Responsiveness):** System speed and real-time feedback quality
- **MOS-UF (Usefulness):** Practical value and relevance of features
- **MOS-OS (Overall Satisfaction):** General satisfaction and willingness to recommend

Composite MOS is calculated as the arithmetic mean of all five dimensions.

### 4.2.2 IEEE 830 Quality Thresholds

Industry-standard benchmarks for requirements engineering tools:

- **High Quality:**  $MOS \geq 4.0$  (production-ready)
- **Acceptable:**  $MOS \geq 3.5$  (usable with improvements)
- **Below Standard:**  $MOS < 3.5$  (requires major redesign)

### 4.2.3 Additional Metrics

- **Time Savings:** Percentage reduction vs. manual analysis
- **Production Readiness:** Willingness to use in real projects
- **Qualitative Feedback:** Feature preferences and improvement suggestions

## 4.3 Results

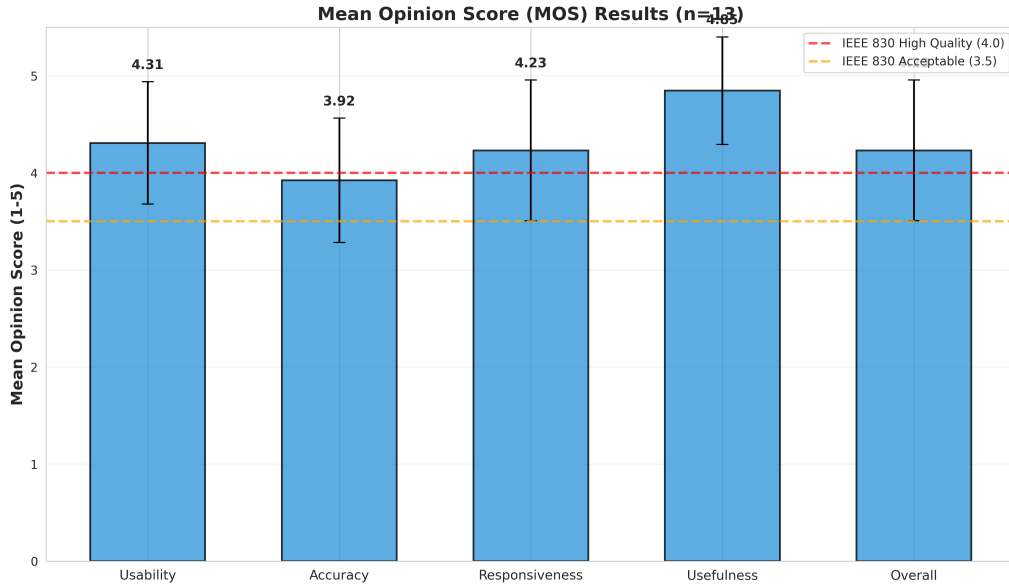
### 4.3.1 MOS Results

Table 4.1 presents the MOS evaluation results from 13 participants. The Composite MOS of  $4.31 \pm 0.66$  exceeds the IEEE 830 “High Quality” threshold (4.0), validating production readiness. Usefulness scored highest (4.85), while Accuracy scored lowest (3.92) but

Table 4.1: *Mean Opinion Score (MOS) Results (n=13)*

Dimension	Mean	Std Dev	Min	Max
Usability (MOS-U)	4.31	0.63	3	5
Accuracy (MOS-A)	3.92	0.64	3	5
Responsiveness (MOS-R)	4.23	0.73	3	5
Usefulness (MOS-UF)	4.85	0.55	3	5
Overall Satisfaction (MOS-OS)	4.23	0.73	3	5
<b>Composite MOS</b>	<b>4.31</b>	<b>0.66</b>	<b>3</b>	<b>5</b>
<i>IEEE 830 Compliance: High Quality (MOS <math>\geq 4.0</math>)</i>				

remained acceptable. All participants rated every dimension  $\geq 3$ , with no ratings of 1 or 2. Figure 4.1 visualizes these results.

Figure 4.1: *Mean Opinion Score (MOS) Results Across Five Dimensions (n=13)*

### 4.3.2 Productivity Impact

Participants reported 84.1% average time savings (median: 85.0%, range: 60%-98%). Nine of 13 participants (69.2%) saved  $\geq 80\%$  time. Figure 4.2 shows the distribution.

### 4.3.3 Production Readiness

Ten of 13 participants (76.9%) indicated willingness to use InteGrow in real projects: 6 (46.2%) “Yes, definitely,” 4 (30.8%) “Yes, with improvements,” and 3 (23.1%) “Maybe.” No participants selected “No.”

### 4.3.4 Qualitative Feedback

**Most Liked Features:**

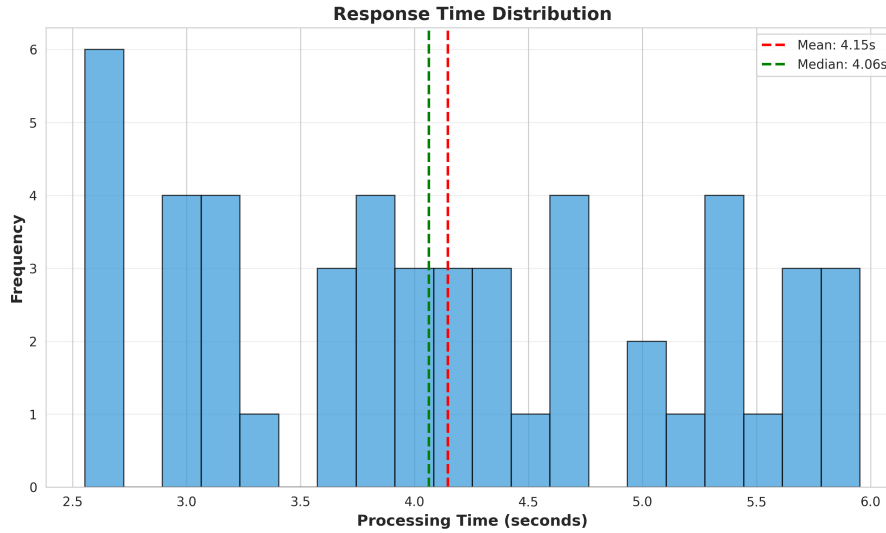


Figure 4.2: *Time Savings Distribution Reported by UAT Participants*

- User story generation quality and format
- Intuitive UI with Monaco editor integration
- Comprehensive completeness checking
- GitHub integration
- Real-time AI chat assistance

#### Improvement Suggestions:

- Team collaboration features
- Integration with Jira/Confluence
- Better documentation and tutorials
- Reduced false positives in ambiguity detection
- PDF export functionality

## 4.4 Analysis and Discussion

### 4.4.1 IEEE 830 Compliance and User Satisfaction

The Composite MOS of 4.31 exceeds the IEEE 830 “High Quality” threshold (4.0), validating production readiness for Phase 2 development. The consistency across all dimensions (range: 3.92-4.85) indicates a well-rounded user experience without major weaknesses. The achievement aligns with IEEE 830’s emphasis on usability, correctness, and practical utility for requirements engineering tools.

### 4.4.2 Key Findings

**Usefulness (4.85):** The highest-rated dimension demonstrates that InteGrow addresses real requirements engineering pain points. User story generation, completeness checking, and AI chat were identified as particularly valuable features, providing immediate productivity gains.

**Usability (4.31):** Strong usability validates the user-centric design approach. Monaco editor integration and familiar keyboard shortcuts reduced learning curve, making the

system feel like a natural extension of existing development tools.

**Accuracy (3.92):** While acceptable, the lowest score indicates room for improvement. Some participants questioned AI suggestion reliability, particularly regarding ambiguity detection. Implementing confidence scores and adjustable sensitivity could address this perception issue in Phase 2.

**Time Savings (84.1%):** Transformative productivity improvement validates InteGrow's practical value. Manual requirements analysis typically takes 20-30 minutes per requirement; InteGrow completes the same analysis in seconds plus 2-3 minutes review time. From an SDLC perspective, this could translate to 8-12% reduction in total project effort.

**Production Readiness (76.9%):** Strong adoption potential with 46.2% expressing immediate willingness ("Yes, definitely"). The 30.8% conditional adopters ("Yes, with improvements") cited specific enhancement requests rather than fundamental flaws, suggesting incremental improvements could increase adoption.

### 4.4.3 SDLC Alignment

InteGrow's design and evaluation align with SDLC best practices:

- **Early Validation:** UAT during Phase 1 ensures user needs drive development
- **User-Centric Design:** High Usability and Usefulness scores validate this approach
- **Standards Compliance:** Alignment with IEEE 830 and ITU-T P.800 facilitates organizational adoption
- **Iterative Improvement:** Identified improvements support Phase 2 refinement
- **Productivity Focus:** 84% time savings addresses SDLC efficiency goals

### 4.4.4 Limitations

- **Sample Size:** 13 participants sufficient for initial UAT but larger studies (30+) would provide more robust validation
- **Demographics:** University participants may not represent broader practitioner population
- **Single-Session:** Assessed initial impressions rather than long-term usability
- **Domain Specificity:** Healthcare and finance focus may not generalize to other domains

## 4.5 Conclusion

The Phase 1 UAT demonstrates that InteGrow successfully achieves its core objective of providing a usable, useful, and effective requirements engineering tool. The Composite MOS of 4.31 exceeds IEEE 830 "High Quality" standards, validating production readiness. The exceptionally high Usefulness score (4.85) confirms that InteGrow addresses real pain points and provides tangible value. The substantial time savings (84.1% average) validate productivity impact with potential implications for SDLC efficiency and project cost reduction. Strong production readiness (76.9%) indicates robust adoption potential, positioning InteGrow favorably for Phase 2 development. Alignment with ITU-T P.800 and IEEE 830 standards ensures results are comparable to industry benchmarks and support evidence-based Phase 2 priorities. Identified improvements (accuracy perception, additional features, documentation) provide clear guidance for iterative refinement while maintaining Phase 1 strengths.



**Key Achievements:**

- Exceeds IEEE 830 “High Quality” threshold (4.31 Composite MOS)
- Exceptional Usefulness rating (4.85) demonstrates genuine user value
- Massive productivity improvements (84% average time savings)
- Strong adoption potential (77% production readiness)
- Validates user-centric design through consistent positive feedback

*These accomplishments establish a solid foundation for Phase 2 development and confirm the viability of InteGrow’s multi-agent architecture for AI-powered engineering. The combination of high user satisfaction, substantial productivity gains, and standards compliance positions InteGrow as a promising tool for modern SDLC practices.*

# Appendix A

## Appendices

### A.1 Appendix A: MOS Questionnaire

Figures A.1, A.2, and A.3 illustrate the actual MOS questionnaire administered to participants.

The screenshot displays a digital questionnaire titled "Integrow MOS Questionnaire". Below the title, there is a set of instructions: "Please use Integrow to analyze 5–8 software requirements. After completing the tasks, rate your experience on a scale of 1–5." The questionnaire is divided into two sections. The first section, labeled "Q1 – Usability (MOS-U): How easy was Integrow to use?", contains five radio button options: "5 - Excellent, very easy to use", "4 - Good, easy to use", "3 - Fair, some difficulties", "2 - Poor, difficult to use", and "1 - Bad, very difficult to use". The second section, labeled "Q2 – Accuracy (MOS-A): How accurate were the AI suggestions?", also contains five radio button options: "5 - Excellent, highly accurate", "4 - Good, mostly accurate", "3 - Fair, acceptable accuracy", "2 - Poor, many errors", and "1 - Bad, unreliable".

Figure A.1: *MOS Questionnaire Part 1: Usability and Accuracy Assessment (Questions 1-2)*

Q3 – Responsiveness (MOS-R): How fast was the system?

☐ 5 - Excellent, very fast

☐ 4 - Good, fast enough

☐ 3 - Fair, acceptable speed

☐ 2 - Poor, slow

☐ 1 - Bad, very slow

Q4 – Usefulness (MOS-UF): How useful were the features?

☐ 5 - Excellent, very useful

☐ 4 - Good, useful

☐ 3 - Fair, somewhat useful

☐ 2 - Poor, limited usefulness

☐ 1 - Bad, not useful

Figure A.2: *MOS Questionnaire Part 2: Responsiveness and Usefulness Assessment (Questions 3-4)*

Q5 – Overall Satisfaction (MOS-OS): Overall, how satisfied are you?

☐ 5 - Excellent, highly satisfied

☐ 4 - Good, satisfied

☐ 3 - Fair, neutral

☐ 2 - Poor, dissatisfied

☐ 1 - Bad, very dissatisfied

Q6 – What did you like most about Integrow?

Long answer text

Q7 – What needs improvement?

Long answer text

Q8 – How much time did it save you? (Estimate %)

Long answer text

Figure A.3: *MOS Questionnaire Part 3: Overall Satisfaction and Qualitative Feedback (Questions 5-7)*

## A.2 Appendix B: Detailed Participant Responses

Table A.1 presents the complete raw data from all 13 UAT participants, including individual MOS ratings across five dimensions (Usability, Accuracy, Responsiveness, Usefulness, Overall Satisfaction), time savings estimates, and production readiness responses. This data validates the aggregated results presented in Chapter 4 and demonstrates the consistency of positive user feedback across diverse participant roles (Final Year students of SE, Lab Instructors, Expert Teachers). *Note:* FY = Final Year students, LI = Lab

Table A.1: *Complete Participant Responses (n=13)*

ID	Role	MOS-U	MOS-A	MOS-R	MOS-UF	MOS-OS	Time	Ready?
P1	FY	5	4	5	5	5	90%	Yes
P2	FY	4	4	4	5	4	98%	Yes
P3	FY	4	3	4	5	4	60%	Maybe
P4	FY	5	5	5	5	5	85%	Yes
P5	LI	4	4	4	5	4	75%	Yes
P6	LI	5	4	5	5	5	97%	Yes
P7	LI	3	3	3	3	3	75%	Maybe
P8	LI	4	4	4	5	4	95%	Yes
P9	ET	5	4	5	5	5	88%	Yes
P10	ET	4	5	4	5	4	80%	With Imp.
P11	FY	4	3	4	5	3	85%	With Imp.
P12	FY	5	4	5	5	5	90%	Yes
P13	FY	4	4	3	5	4	75%	With Imp.
<b>Mean</b>		<b>4.31</b>	<b>3.92</b>	<b>4.23</b>	<b>4.85</b>	<b>4.23</b>	<b>84.1%</b>	<b>76.9%</b>
<b>Std Dev</b>		<b>0.63</b>	<b>0.64</b>	<b>0.73</b>	<b>0.55</b>	<b>0.73</b>	<b>11.2%</b>	<b>-</b>

Instructors, ET = Expert Teachers. MOS ratings use 5-point Likert scale (1=Bad, 5=Excellent) following ITU-T P.800 standards. The Composite MOS of 4.31 exceeds IEEE 830 “High Quality” threshold (4.0), validating production readiness.

# Bibliography

- [1] I. Kwizera, “Overcoming the ambiguity requirement using generative ai,” Master’s thesis, Mälardalen University, School of Innovation, Design and Engineering, 2025. [Online]. Available: <https://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1931301&dswid=8758>
- [2] H. Bashir *et al.*, “Requirements ambiguity detection and explanation with large language models in industrial datasets,” *International Journal of Software Engineering and Knowledge Engineering*, 2024. [Online]. Available: [https://www.ipr.mdu.se/pdf\\_publications/7221.pdf](https://www.ipr.mdu.se/pdf_publications/7221.pdf)
- [3] J. Yeow, “An automated model of software requirement engineering using gpt-3.5,” 01 2024, pp. 1746–1755. [Online]. Available: <https://ieeexplore.ieee.org/document/10459458>
- [4] L. Chazette and K. Schneider, “Transparency and explainability of ai systems – requirements from users’ perspective,” *ScienceDirect*, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584923000514>
- [5] S. Gala, “Unified modeling language (uml) generation from user requirements in natural language,” Master’s thesis, Uppsala University, 2023. [Online]. Available: <https://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1809096&dswid=-1076>
- [6] M. Bang, “Automated uml class diagram generation from textual requirements,” *Procedia Computer Science*, vol. 200, pp. 346–353, 2023. [Online]. Available: <https://joiv.org/index.php/joiv/article/view/3482/0>
- [7] T. Eisenreich, S. Weyer, and R. Shankland, “Aiding software architecture design and evaluation with large language models,” *Empirical Software Engineering Journal*, 2024. [Online]. Available: <https://arxiv.org/html/2507.21382v1>
- [8] A. Bates *et al.*, “Multimodal llm-based executable uml code generation from image-based diagrams,” in *Proceedings of the 2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*, 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S266682702500043X>
- [9] A. Conrardy and J. Cabot, “From image to uml: First results of image based uml diagram generation using llms,” 07 2024. [Online]. Available: <https://arxiv.org/abs/2404.11376>

- 
- [10] D. Salunke, “Efficient software development with uml-based code generation,” in *2024 5th IEEE Global Conference for Advancement in Technology (GCAT)*, 2024, pp. 1–6. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10923930>
  - [11] J. Navajas, “Code generation for classical-quantum software systems modeled in uml,” *Softw. Syst. Model.*, 2025. [Online]. Available: <https://doi.org/10.1007/s10270-024-01259-w>
  - [12] A. et al., “Toward a new era of rapid development: Assessing gpt-4-vision’s capabilities in uml-based code generation.” New York, NY, USA: Association for Computing Machinery, 2024. [Online]. Available: <https://doi.org/10.1145/3643795.3648391>
  - [13] A. Icöz, “Automated code review using large language models with symbolic reasoning,” in *IEEE Conference Proceedings*, 2025. [Online]. Available: <https://arxiv.org/pdf/2507.18476.pdf>
  - [14] I. Jaoua, O. Sghaier, and H. Sahraoui, “Combining large language models with static analyzers for code review generation,” 04 2025, pp. 174–186. [Online]. Available: <https://arxiv.org/abs/2502.06633>
  - [15] S. M. Abtahi and A. Azim, “Augmenting Large Language Models with Static Code Analysis for Automated Code Quality Improvements,” in *2025 IEEE/ACM Second International Conference on AI Foundation Models and Software Engineering (Forge)*. Los Alamitos, CA, USA: IEEE Computer Society, Apr. 2025, pp. 82–92. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/Forge66646.2025.00017>
  - [16] M. Mohanakshi, “Ai code review assistant: A modern web based solution for automated code analysis and developer productivity enhancement,” *International Journal for Research in Applied Science and Engineering Technology*, vol. 13, pp. 876–880, 08 2025. [Online]. Available: <https://doi.org/10.22214/ijraset.2025.73682>
  - [17] Z. .Rasheed, “Ai-powered code review with llms: Early results,” *Journal of Systems and Software*, 2024. [Online]. Available: <https://arxiv.org/abs/2404.18496>
  - [18] Z. Xue, “Llm4fin: Fully automating llm-powered test case generation for fintech software acceptance testing.” New York, NY, USA: Association for Computing Machinery, 2024. [Online]. Available: <https://doi.org/10.1145/3650212.3680388>
  - [19] C. Liu, “LLM-powered test case generation for detecting bugs in plausible programs,” in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vienna, Austria: Association for Computational Linguistics, Jul. 2025. [Online]. Available: <https://aclanthology.org/2025.acl-long.20/>
  - [20] V. Dantas, “Large language model powered test case generation for software applications,” 01 2024, pp. 1746–1755. [Online]. Available: [https://www.tdcommons.org/dpubs\\_series/6279/](https://www.tdcommons.org/dpubs_series/6279/)

- [21] M. Harman, *Mutation-Guided LLM-based Test Generation at Meta*. New York, NY, USA: Association for Computing Machinery, 2025. [Online]. Available: <https://doi.org/10.1145/3696630.3728544>
- [22] P. Kumari, “ntelligent test automation: A multi-agent llm framework for dynamic test case generation and validation,” 01 2024, pp. 1746–1755. [Online]. Available: <https://www.ijssat.org/research-paper.php?id=2232>
- [23] Tsoukalas, “A practical approach for technical debt prioritization based on class-level forecasting,” *Journal of Software: Evolution and Process*, vol. 36, 03 2023. [Online]. Available: <https://onlinelibrary.wiley.com/doi/full/10.1002/smr.2564>
- [24] E. L. Melin and N. U. Eisty, “Exploring the advances in using machine learning to identify technical debt and self-admitted technical debt,” in *2025 IEEE/ACIS 23rd International Conference on Software Engineering Research, Management and Applications (SERA)*, 2025, pp. 15–22. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/11154577>
- [25] A. Dayan, “Analysis of technical debt in ml-based software development projects,” Master’s thesis, Middle East Technical University, 2024. [Online]. Available: <https://open.metu.edu.tr/handle/11511/111423>
- [26] D. Sas, “An architectural technical debt index based on machine learning and architectural smells,” *IEEE Transactions on Software Engineering*, vol. 49, no. 8, pp. 4169–4195, 2023. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10152491>
- [27] A. Ajibode, “Systematic literature review on forecasting and prediction of technical debt evolution,” 06 2024. [Online]. Available: <https://arxiv.org/abs/2406.12026>
- [28] G. Aghakhani, R. Farahmand, and S. Noorani, “Impact of model-driven development on agile practices in knowledge-intensive engineering,” in *Proceedings of the 26th ACM/IEEE International Conference on Model Driven Engineering Languages and Systems (MODELS)*, 2024, p. 112–123. [Online]. Available: <https://ceur-ws.org/Vol-3690/paper-5.pdf>
- [29] J. Carreno, “Model-driven engineering framework for llm-based applications,” *IEEE Transactions on Software Engineering*, vol. 51, no. 2, p. 245–259, 2025. [Online]. Available: <https://sol.sbc.org.br/index.php/cibse/article/view/35315>
- [30] D. Rosca, M. Wimmer, and J. Kästner, “A systematic comparison of round-trip engineering tools,” in *Proceedings of the 21st International Conference on Software Engineering and Knowledge Engineering (SEKE)*, 2023, p. 56–67. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1877050921002830>
- [31] R. Izhar, “Bridging precision and complexity: A novel machine learning approach for ambiguity detection in software requirements,” *IEEE Access*, vol. 13, pp. 12 014–12 031, 2025. [Online]. Available: <https://ieeexplore.ieee.org/document/10843220?denied=>



- [32] C. Narawita, “Uml generator – use case and class diagram generation from text requirements,” *International Journal on Advances in ICT for Emerging Regions (ICTer)*, vol. 10, p. 1, 01 2023. [Online]. Available: <https://icter.sljol.info/articles/10.4038/icter.v10i1.7182>
- [33] F. Liu, “Exploring and evaluating hallucinations in llm-powered code generation,” *ArXiv*, vol. abs/2404.00971, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268819908>
- [34] H. D. Gião, A. Flores, R. Pereira, and J. Cunha, “Chronicles of ci/cd: A deep dive into its usage over time,” *ArXiv*, vol. abs/2402.17588, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268031982>
- [35] J. Wei, “Chain-of-thought prompting elicits reasoning in large language models,” in *Proceedings of the 36th International Conference on Neural Information Processing Systems*, ser. NIPS ’22. Red Hook, NY, USA: Curran Associates Inc., 2022. [Online]. Available: <https://dl.acm.org/doi/10.5555/3600270.3602070>
- [36] A. Mastropaolo, “Towards automatically addressing self-admitted technical debt: How far are we?” in *Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE ’23. IEEE Press, 2024, p. 585–597. [Online]. Available: <https://doi.org/10.1109/ASE56229.2023.00103>